Machine Learning for Design

Lecture Notes [V0.61]

Alessandro Bozzon

2023-03-29

Contents

1	Sect	ion 1: A	rtificial Intelligence and Machine Learning in Intelligent Products, Services,	
	and	System	ıs (design)	1
	1.1	1.1 Wh	y should you care about Machine Learning?	1
		1.1.1	Al is the new electricity	1
		1.1.2	Definitions	3
		1.1.3	Can't intelligence be programmed?	7
		1.1.4	Machine Learning	10
		1.1.5	Functions of a Machine Learning System	11
		1.1.6	Neural Networks and Deep Learning	13
		1.1.7	Computer Vision	14
		1.1.8	Natural Language Processing	15
		1.1.9	The hard problems are easy, and the easy problems are hard	15
	1.2	1.2 So,	, why should you care about Machine Learning?	16
		1.2.1	Why do we need Designers to understand ML?	16
		1.2.2	What can designers do for Machine Learning?	17
		1.2.3	What can designers do with Machine Learning?	19
		1.2.4	What can Machine Learning do for designers?	21
		1.2.5	Why Programming?	23
		1.2.6	Debunking some myths	23
2	Sect	ion 2: F	undamentals of Machine Learning	27
	2.1	2.1 The	e Machine Learning Life-Cycle	27
		2.1.1	CRISP-DM In our course	28
		2.1.2	Problem Specification	29
		2.1.3	Data Understanding	29
		2.1.4	Data Preparation	30
		2.1.5	Modeling	30
		2.1.6	Evaluation	31
		2.1.7	Deployment and monitoring	33
	2.2	2.2 Dat	ta. The raw material	34
		2.2.1	Data	34
		2.2.2	Types of Features / Label Values	35

		2.2.3	Data Modalities	37
		2.2.4	Key Dimensions	38
		2.2.5	Туреs of Data	39
		2.2.6	Categorising Data Sources	41
	2.3	2.3 Ca	tegories of Machine Learning	43
		2.3.1	How do machines learn?	43
		2.3.2	On Models	43
		2.3.3	Machine Learning Approaches	47
		2.3.4	The importance of domain expertise	54
	2.4	2.4 Ma	chine Learning Models	54
		2.4.1	Linear Regression	54
		2.4.2	Training a Machine Learning Model	58
		2.4.3	Classification	60
		2.4.4	Neural Networks	66
	2.5	2.5 Mo	dels Development Lifecycle	72
		2.5.1	Dataset Splitting	75
	2.6	2.6 Eva	aluating Machine Learning Models	77
		2.6.1	How to Evaluate?	77
		2.6.2	Let errors guide you	78
		2.6.3	How to calculate	78
		2.6.4	Metrics are designed in a multi-stakeholder context	87
_		• • •		
3	Sect	ion 3: I	mage Processing Methods	89
3	Sect 3.1	ion 3: I 3.1 lm	mage Processing Methods ages	89 90
3	Sect 3.1 3.2	ion 3: I 3.1 lm 3.2 Co	mage Processing Methods ages	89 90 90
3	Sect 3.1 3.2 3.3	ion 3: I 3.1 lm 3.2 Co 3.3 Wh	mage Processing Methods ages	89 90 90 91
3	Sect 3.1 3.2 3.3	ion 3: I 3.1 lm 3.2 Co 3.3 Wh 3.3.1	mage Processing Methods ages	89 90 90 91 92
3	Sect 3.1 3.2 3.3	ion 3: I 3.1 lm 3.2 Co 3.3 Wh 3.3.1 3.3.2	mage Processing Methods ages	89 90 91 92 93
3	Sect 3.1 3.2 3.3	ion 3: I 3.1 lm 3.2 Co 3.3 Wh 3.3.1 3.3.2 3.3.3	mage Processing Methods ages mputer Vision mat specific tasks can we train a CV system to perform? Object Recognition / Localisation Object Identification Image Segmentation	 89 90 90 91 92 93 95
3	Sect 3.1 3.2 3.3	ion 3: I 3.1 lm 3.2 Co 3.3 Wh 3.3.1 3.3.2 3.3.3 3.3.4	mage Processing Methods ages mputer Vision nat specific tasks can we train a CV system to perform? Object Recognition / Localisation Object Identification Image Segmentation Scene Recognition	 89 90 91 92 93 95 96
3	Sect 3.1 3.2 3.3	ion 3: 1 3.1 lm 3.2 Co 3.3 Wh 3.3.1 3.3.2 3.3.3 3.3.4 3.3.5	mage Processing Methods ages mputer Vision nat specific tasks can we train a CV system to perform? Object Recognition / Localisation Object Identification Image Segmentation Scene Recognition Human Activity Recognition	 89 90 91 92 93 95 96 97
3	Sect 3.1 3.2 3.3	ion 3: I 3.1 lm 3.2 Co 3.3 Wh 3.3.1 3.3.2 3.3.3 3.3.4 3.3.5 3.3.6	mage Processing Methods ages mputer Vision nat specific tasks can we train a CV system to perform? Object Recognition / Localisation Object Identification Image Segmentation Scene Recognition Human Activity Recognition Pose Estimation	 89 90 91 92 93 95 96 97 97
3	Sect 3.1 3.2 3.3	ion 3: I 3.1 lm 3.2 Co 3.3 Wh 3.3.1 3.3.2 3.3.3 3.3.4 3.3.5 3.3.6 3.4 Ho	mage Processing Methods ages mputer Vision nat specific tasks can we train a CV system to perform? Object Recognition / Localisation Object Identification Image Segmentation Scene Recognition Human Activity Recognition Pose Estimation w do humans see?	 89 90 91 92 93 95 96 97 97 99
3	Sect 3.1 3.2 3.3	ion 3: 1 3.1 lm 3.2 Co 3.3 Wh 3.3.1 3.3.2 3.3.3 3.3.4 3.3.5 3.3.6 3.4 Ho 3.4.1	mage Processing Methods ages	 89 90 91 92 93 95 96 97 97 99 99
3	Sect 3.1 3.2 3.3	ion 3: I 3.1 lm 3.2 Co 3.3 Wh 3.3.1 3.3.2 3.3.3 3.3.4 3.3.5 3.3.6 3.4 Ho 3.4.1 3.4.2	mage Processing Methods ages mputer Vision nat specific tasks can we train a CV system to perform? Object Recognition / Localisation Object Identification Image Segmentation Scene Recognition Human Activity Recognition Pose Estimation w do humans see? Hubel and Wiesel, 1959 Neural Pathways	 89 90 91 92 93 95 96 97 97 99 99 100
3	Sect 3.1 3.2 3.3	ion 3: 1 3.1 lm 3.2 Co 3.3 Wh 3.3.1 3.3.2 3.3.3 3.3.4 3.3.5 3.3.6 3.4 Ho 3.4.1 3.4.2 3.4.3	mage Processing Methods ages mputer Vision nat specific tasks can we train a CV system to perform? Object Recognition / Localisation Object Identification Image Segmentation Scene Recognition Human Activity Recognition Pose Estimation W do humans see? Hubel and Wiesel, 1959 Neural Pathways Neural Correlation of Objects & Scene Recognition	 89 90 91 92 93 95 96 97 97 99 99 100 101
3	Sect 3.1 3.2 3.3 3.4	ion 3: I 3.1 lm 3.2 Co 3.3 Wh 3.3.1 3.3.2 3.3.3 3.3.4 3.3.5 3.3.6 3.4 Ho 3.4.1 3.4.2 3.4.3 3.5 Wh	mage Processing Methods ages mputer Vision nat specific tasks can we train a CV system to perform? Object Recognition / Localisation Object Identification Image Segmentation Scene Recognition Human Activity Recognition Pose Estimation W do humans see? Hubel and Wiesel, 1959 Neural Pathways Neural Correlation of Objects & Scene Recognition	 89 90 91 92 93 95 96 97 97 99 90 100 101 101
3	Sect 3.1 3.2 3.3 3.4 3.4	ion 3: 1 3.1 lm 3.2 Co 3.3 Wh 3.3.1 3.3.2 3.3.3 3.3.4 3.3.5 3.3.6 3.4 Ho 3.4.1 3.4.2 3.4.3 3.5 Wh 3.5.1	mage Processing Methods ages mputer Vision nat specific tasks can we train a CV system to perform? Object Recognition / Localisation Object Identification Image Segmentation Image Segmentation Scene Recognition Human Activity Recognition Pose Estimation w do humans see? Hubel and Wiesel, 1959 Neural Pathways Neural Correlation of Objects & Scene Recognition The deformable and truncated cat	 89 90 91 92 93 95 96 97 99 90 101 101 101
3	Sect 3.1 3.2 3.3 3.4 3.4	ion 3: 1 3.1 lm 3.2 Co 3.3 Wh 3.3.1 3.3.2 3.3.3 3.3.4 3.3.5 3.3.6 3.4 Ho 3.4.1 3.4.2 3.4.3 3.5 Wh 3.5.1 3.5.1 3.5.2	mage Processing Methods ages	 89 90 91 92 93 95 96 97 97 99 99 100 101 101 101 102

	3.6	3.6 Ho	w Computer Vision models work?
		3.6.1	Course of dimensionality
		3.6.2	The "old days": Feature Extraction and Engineering
	3.7	3.7 Coi	nvolutional Neural Networks
		3.7.1	Convolution & Feature Maps
		3.7.2	What CNNs learn?
		3.7.3	Translation Invariance
		3.7.4	What about generalisation?
		3.7.5	Data Augmentation
		3.7.6	Robustness to input variation
		3.7.7	Transfer Learning
	3.8	3.8 Adv	vanced Computer Vision Techniques
		3.8.1	Generative Adversarial Networks
		3.8.2	Which face is real?
		3.8.3	Image super-resolution GAN
		3.8.4	Neural Style Transfer
		3.8.5	Text-To-Image Generation
		3.8.6	Image-to-Image Generation
		3.8.7	Synthetic Video Generation
		3.8.8	Deep Fakes
		3.8.8	Deep Fakes
4	Sect	3.8.8 ion 4: T	Deep Fakes
4	Sect 4.1	3.8.8 ion 4: T 4.1 Wh	Deep Fakes 126 'ext Processing Methods 127 y natural language processing? 127
4	Sect 4.1	3.8.8 ion 4: T 4.1 Wh 4.1.1	Deep Fakes 126 ext Processing Methods 127 y natural language processing? 127 Big Textual Data = Language at scale 127
4	Sect 4.1	3.8.8 ion 4: T 4.1 Wh 4.1.1 4.1.2	Deep Fakes 126 Text Processing Methods 127 y natural language processing? 127 Big Textual Data = Language at scale 127 Uses of NLP 130
4	Sect 4.1 4.2	3.8.8 ion 4: T 4.1 Wh 4.1.1 4.1.2 4.2 Wh	Deep Fakes 126 rext Processing Methods 127 y natural language processing? 127 Big Textual Data = Language at scale 127 Uses of NLP 130 at is Natural Language Processing? 131
4	Sect 4.1 4.2	3.8.8 ion 4: T 4.1 Wh 4.1.1 4.1.2 4.2 Wh 4.2.1	Deep Fakes 126 Text Processing Methods 127 y natural language processing? 127 Big Textual Data = Language at scale 127 Uses of NLP 130 at is Natural Language Processing? 131 Beyond keyword matching 132
4	Sect 4.1 4.2 4.3	3.8.8 ion 4: T 4.1 Wh 4.1.1 4.1.2 4.2 Wh 4.2.1 4.3 Wh	Deep Fakes126Text Processing Methods127y natural language processing?127Big Textual Data = Language at scale127Uses of NLP130at is Natural Language Processing?131Beyond keyword matching132y is NLP Hard?132
4	Sect 4.1 4.2 4.3	3.8.8 ion 4: T 4.1 Wh 4.1.1 4.1.2 4.2 Wh 4.2.1 4.3 Wh 4.3.1	Deep Fakes126rext Processing Methods127y natural language processing?127Big Textual Data = Language at scale127Uses of NLP130at is Natural Language Processing?131Beyond keyword matching132y is NLP Hard?132Ambiguity and Expressivity134
4	Sect 4.1 4.2 4.3	3.8.8 ion 4: T 4.1 Wh 4.1.1 4.1.2 4.2 Wh 4.2.1 4.3 Wh 4.3.1 4.3.2	Deep Fakes126Text Processing Methods127y natural language processing?127Big Textual Data = Language at scale127Uses of NLP130at is Natural Language Processing?131Beyond keyword matching132y is NLP Hard?132Ambiguity and Expressivity134Sparsity136
4	Sect 4.1 4.2 4.3	3.8.8 ion 4: T 4.1 Wh 4.1.1 4.1.2 4.2 Wh 4.2.1 4.3 Wh 4.3.1 4.3.2 4.3.3	Deep Fakes126rext Processing Methods127y natural language processing?127Big Textual Data = Language at scale127Uses of NLP130at is Natural Language Processing?131Beyond keyword matching132y is NLP Hard?132Ambiguity and Expressivity134Sparsity136Language Evolution138
4	Sect 4.1 4.2 4.3	3.8.8 ion 4: T 4.1 Wh 4.1.1 4.1.2 4.2 Wh 4.2.1 4.3 Wh 4.3.1 4.3.2 4.3.3 4.4 NLI	Deep Fakes126Text Processing Methods127y natural language processing?127Big Textual Data = Language at scale127Uses of NLP130at is Natural Language Processing?131Beyond keyword matching132y is NLP Hard?132Ambiguity and Expressivity134Sparsity136Language Evolution138P Tasks139
4	Sect 4.1 4.2 4.3	3.8.8 ion 4: T 4.1 Wh 4.1.1 4.1.2 4.2 Wh 4.2.1 4.3 Wh 4.3.1 4.3.2 4.3.3 4.4 NLH 4.4.1	Deep Fakes126 fext Processing Methods 127y natural language processing?127Big Textual Data = Language at scale127Uses of NLP130at is Natural Language Processing?131Beyond keyword matching132y is NLP Hard?132Ambiguity and Expressivity134Sparsity136Language Evolution138P Tasks139An example of NLP Process140
4	Sect 4.1 4.2 4.3	3.8.8 ion 4: T 4.1 Wh 4.1.1 4.1.2 4.2 Wh 4.2.1 4.3 Wh 4.3.1 4.3.2 4.3.3 4.4 NLH 4.4.1 4.4.2	Deep Fakes126Fext Processing Methods127y natural language processing?127Big Textual Data = Language at scale127Uses of NLP130at is Natural Language Processing?131Beyond keyword matching132y is NLP Hard?132Ambiguity and Expressivity134Sparsity136Language Evolution138P Tasks139An example of NLP Process140Morphology141
4	Sect 4.1 4.2 4.3 4.4	3.8.8 ion 4: T 4.1 Wh 4.1.1 4.1.2 4.2 Wh 4.2.1 4.3 Wh 4.3.1 4.3.2 4.3.3 4.4 NLI 4.4.1 4.4.2 4.4.3	Deep Fakes126fext Processing Methods127y natural language processing?127Big Textual Data = Language at scale127Uses of NLP130at is Natural Language Processing?131Beyond keyword matching132y is NLP Hard?132Ambiguity and Expressivity134Sparsity136Language Evolution138P Tasks139An example of NLP Process141Syntax144
4	Sect 4.1 4.2 4.3 4.4	3.8.8 ion 4: T 4.1 Wh 4.1.1 4.1.2 4.2 Wh 4.2.1 4.3.2 4.3.3 4.4 NLH 4.4.1 4.4.2 4.4.3 4.4.4	Deep Fakes126fext Processing Methods127y natural language processing?127Big Textual Data = Language at scale127Uses of NLP130at is Natural Language Processing?131Beyond keyword matching132y is NLP Hard?132Ambiguity and Expressivity134Sparsity136Language Evolution138P Tasks139An example of NLP Process140Morphology141Syntax144Semantics150
4	Sect 4.1 4.2 4.3	3.8.8 ion 4: T 4.1 Wh 4.1.1 4.2 Wh 4.2.1 4.3 Wh 4.3.1 4.3.2 4.3.3 4.4 NLI 4.4.1 4.4.2 4.4.3 4.4.4 4.4.5	Deep Fakes 126 Yext Processing Methods 127 y natural language processing? 127 Big Textual Data = Language at scale 127 Uses of NLP 130 at is Natural Language Processing? 131 Beyond keyword matching 132 y is NLP Hard? 132 Ambiguity and Expressivity 134 Sparsity 136 Language Evolution 138 P Tasks 139 An example of NLP Process 140 Morphology 141 Syntax 144 Semantics 150 State of the Art in NLP - as of 2022 157
4	Sect 4.1 4.2 4.3 4.4	3.8.8 ion 4: T 4.1 Wh 4.1.1 4.1.2 4.2 Wh 4.2.1 4.3 Wh 4.3.1 4.3.2 4.3.3 4.4 NLH 4.4.2 4.4.3 4.4.4 4.4.5 4.5 Fea	Deep Fakes 126 Yext Processing Methods 127 y natural language processing? 127 Big Textual Data = Language at scale 127 Uses of NLP 130 at is Natural Language Processing? 131 Beyond keyword matching 132 y is NLP Hard? 132 Ambiguity and Expressivity 134 Sparsity 136 Language Evolution 138 P Tasks 139 An example of NLP Process 140 Morphology 141 Syntax 144 Semantics 150 State of the Art in NLP - as of 2022 157

4.6	4.6 Lar	nguage Models	158
	4.6.1	What is a <i>language model==</i> *?	159
	4.6.2	What are language models used for?	159
	4.6.3	What is the issue with word representation?	160
	4.6.4	Word Embeddings	161
	4.6.5	One-Hot Encoding	163
	4.6.6	Independent and identically distributed words assumption	165
	4.6.7	Vector Representation of Words	165
	4.6.8	N-gram language models	167
	4.6.9	Representing words by their contexts	169
	4.6.10	How to calculate Distributional Word Embeddings?	174
4.7	4.7 Usi	ing Word Embeddings	176
	4.7.1	How can embeddings be used with NLP Models?	176
	4.7.2	Retraining	177
	4.7.3	Evaluating Word Embeddings	178
4.8	4.8 Lar	ge Language Models	182
Defe			102
кете	erences		193
Crea	lits		185

5

6

1 Section 1: Artificial Intelligence and Machine Learning in Intelligent Products, Services, and Systems (design)

1.1 1.1 Why should you care about Machine Learning?

1.1.1 AI is the new electricity



Andrew NG, a well known AI scientists and enterpreneur.

Some people even go as far as to say that: AI is the new electricity. This man in the figure is Andrew NG. He is a scientist at Baidu and an entrepreneur. He is one of the people that have been massively advocating for the fact that artificial intelligence is going to be everywhere.

One of the goals of this course is to somehow mediate between the reality and the perception that many people bring out about AI and about what it can and cannot do.

I might not 100% agree with Andrew that AI is like electricity. But I do agree about the fact that it is going to be in a lot of products, services, and systems that designers are designing and, even more in the future, they will design

1.1.1.1 Where is artificial Intelligence?

- Autonomous vehicles
 - from Roomba to Self-driving cars
 - In stores, warehouses, production lines, streets, living rooms



Examples of intelligent products indicated by students.

- More and more consumer products and appliances
 - Thermostats, Security Cameras, Fridges
- Content production and consumption applications
 - Social media, Amazon, Netflix etc.
- Chatbots
- In-store automation and smarter shopping
- Optimised supply chains
- Energy grid optimisation
- ...

Most of the examples in the Figure 1.1 and in the list above are about digital products enhanched with AI technology. At the end of the course, you will see that they are in other products as well. By the end of the course, you will see that also physical products have this.

Here we provide a non-exhaustive list of applications or products where AI/ML plays a role, for example, autonomous vehicles, the Roomba, to the self-driving. Some of you indicated the robot.



Another example of intelligent product. Or maybe it does not exist. Would you be able to tell the difference?



1.1.2 Definitions



Some of the terminology we will be using in this course

Now, let's also make sure that we clarify what we mean by different terminology. We need to rationalize and demystify what artificial intelligence is. And more specifically, what is Machine Learning?

What can it do? What can they not do? and how can you talk soberly about it without buying into the hype and the whole marketing and advertisement-driven narrative?

Unfortunately, we live in the age of marketing and advertisements, so everything is somehow declining regarding who can sell the most of it.

Alessandro Bozzon

1.1.2.1 Intelligence

Mental quality that consists of the abilities to learn from experience, adapt to new situations, understand and handle abstract concepts, and use knowledge to manipulate one's environment.[^ Encyclopaedia Britannica]

Let us start with a definition of *Intelligence*. The problem, however, is that there seem to be almost as many definitions of intelligence as there were experts asked to define it [^ R. J. Sternberg, quoted in The Oxford Companion to the Mind. R. L. Gregory. Oxford University]. Here are three examples - The ability to learn or understand or to deal with new or trying situations - The ability to apply knowledge to manipulate one's environment or to think abstractly as measured by objective criteria (such as tests) - Mental quality that consists of the abilities to learn from experience, adapt to new situations, understand and handle abstract concepts, and use knowledge to manipulate one's environment

If you are interested in a relatively exhaustive collection of definitions you can find it in this technical report published in 2007. The manuscript contains 71 definitions, drawn from literature in Psychology and Artificial Intelligence.

Intelligence can be described as the ability to learn and understand, to translate knowledge and skills from one situation to another.

1.1.2.2 Our definition of Intelligence

Intelligence measures an agent's ability to achieve goals in a wide range of environments.

As we will discuss AI and ML from a technical perspective, we will adopt a definition that resounds as very technical, given its emphasis on **measurement** – that is, the quantification of intelligence, to allow for comparison – and the explicit mentioning of an **agent** – that is, anything that can perceive its environment, take actions autonomously to achieve goals, and may improve its performance with learning or may use knowledge.

1.1.2.3 Artificial Intelligence

Intelligence demonstrated by machines Computer programs that can emulate *physical* and/or *cognitive* human capabilities

Multiple disciplines in computer science (e.g., robotics, machine learning, software engineering, data science, computer vision, and natural language processing) are devoted to studying the properties of intelligence by synthesizing it, that is, by recreating it artificially. Outside of computer science, the interest in Artificial Intelligence has grown exponentially and now includes fields like **Design**, ethics, sociology, and so on. However, most of the terminology and most of the narrative around Artificial Intelligence has originated from computer science.

When we talk about Artificial Intelligence, we talk about it as intelligence demonstrated by machines. This is a broad definition because it allows us to include any human capability: even a calculator or a traditional thermostat can be defined as intelligent. And, in a sense, they are. But we want to be more specific.

When we talk about Artificial Intelligence, we talk about **computer programs** that can emulate some **physical** or **cognitive** or human capabilities.

Physical capabilities relate to the physical sensing or manipulation of an environment. Sensing is performed through devices (**sensors**) that detect events or environmental changes. Sensors can be mechanical, chemical, or biological. Manipulations are typically achieved by controlling electronic, mechanical, or magnetic **actuators**. Examples of actuators are engines, valves, and switches.

Thanks to these physical capabilities, it is possible for a machine to execute activities that require some (complex) motor coordination, like opening a door, driving a car, playing football, or flying an airplane or a drone. Such coordination is, however "directed" by cognitive capabilities like:

- perceiving the world, e.g. seeing, hearing, feeling, smelling, tasting. Perception is enabled by sensing technology, but the interpretation of a physical signal requires cognitive abilities.
- learning, i.e. gaining knowledge from past experiences and interactions with the environment
- reasoning
- planning
- problem-solving
- creating

1.1.2.4 Strong vs. Weak Al

Another important terminological clarification is between **Strong** and **Weak AI**. It is an important clarification because it is very common to read headlines like:

- Robot kills worker at Volkswagen plant in Germany
- Are you scared yet? Meet Norman, the psychopathic AI
- Facebook AI Creates Its Own Language In Creepy Preview Of Our Potential Future
- Sophia the robot was granted citizenship.

All these ridiculous headlines portray Artificial Intelligence as human-level or even super-human intelligence.

1.1.2.4.1 Strong Artificial Intelligence

- AI that can do everything we humans can do, and possibly much more.
- Also called *Artificial General Intelligence* (AGI) or *human-level* intelligence. The AI we see in movies

• No AI program has been created yet that could be considered an AGI

Human-level intelligence is typically called **Strong AI**, which we see in the movies like "I, Robot," "Terminator," or wherever there is an entity that seems to be able to exist with the same level of cognitive reasoning capabilities as we do. It is the sort of Artificial Intelligence able to perform any reasoning (inductive, deductive, abductive). An intelligence that can "connect the dots" between the different experiences, perceptions, and activities, blending them to create a representation of the world (and its actors) complete enough to act as an autonomous and independent agent, able to (theoretically) perform any task.

It is important to realise that no one has ever created a computer program able to exhibit human-level intelligence, and probably there will not be one any time soon.

Another term often used is **super intelligence**, which is AI systems that can do things better than humans. While it is true that machines can execute some tasks better than humans (a 5 Euro calculator can do computations better than the vast majority of human beings), that does not mean that such machines are, overall, more intelligent than humans. It is a game of definitions, mostly fuelled by marketing needs.

1.1.2.4.2 Weak Intelligence

- Also called Narrow AI.
- AI specialised in well-defined tasks.
- For example, speech recognition, chess-playing, autonomous driving

Current discussions around Artificial Intelligence and related technology are around what is called **weak AI** or **narrow AI**, that is AI systems specialized in a specific task (e.g. recognizing cats in an image) and probably terrible at different ones.

Note that even a specific task, like recognizing an object, requires some general intelligence capabilities to be performed at a truly human level. As we will see later in the course, a computer vision system (an AI system specialised in visual perception tasks) have difficulties in keeping equivalents levels of performance across different contexts. For instance, when there are different lighting conditions (e.g. dark scenes) or environments (e.g., indoor vs. Outdoor). Al systems require a lot of input (training) to function. Humans can generalize their recognition activities also to variations of objects that they have never seen before. Weak AI cannot. An AI trained to recognize green and red apples will never be able to recognize pears.

Many of these issues will be addressed again throughout the course. Unfortunately, we've been through 5 to 10 years of brainwashing from corporate marketing and lazy news reporters, and it is important to understand what the actual capabilities of AI systems are.

1.1.2.5 Learning

- Any process by which a system improves performance from experience [^ Herbert Alexander Simon]
- The ability to perform a task in a situation that has never been encountered before
- Learning = generalisation

Artificial Intelligence features a broad set of approaches to create machines exhibiting some form of intelligence. *Machine Learning*, as we will see, focuses specifically on teaching computers how to **learn** without the need to be programmed for specific tasks.

But what is **learning**? Learning can be seen as the process through which a system (or a person) adapts so that a task or tasks drawn from the same set of tasks can be performed more efficiently and effectively the next time. When we learn how to drive a bicycle, most of our attention, in the beginning, is focused on how to use the handlebar and the pedals so that we can move forward (and break) without falling. As we gain experience, our focus shifts toward cycling faster or longer, and we acquire di ability to operate a bicycle.

Another way to define learning is the ability to **generalize**. That is the ability to perform a task in a situation that has never been encountered before by extrapolating from the knowledge and experience acquired by performing the same task in a different situation or task. We typically learn how to cycle using the same bike in our backyard. As we gain experience, we can operate different bikes in different environments (e.g., city bikes during the week and race bikes during the weekend). The experience we gained learning how to bike will allow us, for instance, also to operate a moped or a motorbike. We generalised our knowledge of how to ride a two-wheeler to other types of vehicles.

1.1.3 Can't intelligence be programmed?

At this point, you might be wondering: what does all of this have to do with Artificial Intelligence and computers? Thousand of people have been devoting decades of their lives to studying how to train programming computers. Can't intelligence be programmed?

1.1.3.1 Polany's Paradox

"We can know more than we can tell... The skill of a driver cannot be replaced by a thorough schooling in the theory of the motorcar" Michael Polanyi (1966)

But this brings about an interesting paradox that highlights one of the main issues with training AI.

Polanyi's paradox can help us understand why it is not always easy to program machines to execute tasks as we do.

Alessandro Bozzon

2023-03-29

Take the game of *chess*, for instance. One of the reasons why computers have been able to beat a reigning chess world champion (Gary Kasparov) as early as 1996 is our ability to mathematically formalize the rules of chess. Given a very large amount of time, and abundant computational power, we could simulate any game of chess ever played, or that will ever be played.

But take a different task, like driving a car, riding a bicycle, or writing a successful book. These are examples of activities that we, human beings, learn how to perform without being able to fully verbalize (or formalize) the rules or procedures behind them. We learn them through biological, cultural, social, personal, and interpersonal processes. We often make decisions for reasons that we can't tell, and therefore we don't know.

1.1.3.2 What is a cat?[^ Credits: Jonah Burlingame]



Let us give you a concrete example. Can you tell me what a cat is? What are the distinctive properties of a cat?



- A cat has whiskers
- A cat is furry

You can start with some visual properties of cats. For instance, they have whiskers, and they are furry.



- A cat has whiskers
- A cat is furry
- But so are lions!

But cats are not the only animals with such properties. Lions also have whiskers and fur.



- A cat has whiskers
- A cat is furry
- A cat is small

You could then operate by exclusion and add a new rule saying that cats are small.



- A cat has whiskers
- A cat is furry
- A cat is small
- But so are koalas

But many animals can you think of that fit this description? Quite a lot. The koala, for instance.



- A cat has whiskers
- A cat is furry

- A cat is small
- A cat does not climb trees

You could add a new rule that says cats do not climb trees.



- A cat has whiskers
- A cat is furry
- A cat is small
- A cat does not climb trees
- well...

But sometimes they do. If you want to try to tease out all the possible properties that make a thing an animal or a person, **that would take forever**.

1.1.4 Machine Learning

The field of study that gives computers the ability to learn without being explicitly programmed[^ Arthur Samuel]. Machine learning is the science (and art) of programming computers **so they can learn from data**

The form of cognitive "self-ignorance" illustrated in the previous example does not prevent us from teaching all of these activities to other human beings. So, how can we program a machine with something that resembles our intelligence if we can't explain our intelligence?

Machine Learning is an approach to Artificial Intelligence that partially solves the problem by circumventing the issue of being able to state the rules of our cognitive abilities explicitly. Machine learning allows computer systems to learn directly from examples, data, and experience. By enabling computers to perform specific tasks intelligently, machine learning systems can carry out complex processes by learning from data rather than following pre-programmed rules.

1.1.4.1 Programming vs. ML



Traditionally computers are programmed using explicit instructions that map pre-defined properties of a given input to desired outcomes using ad-hoc algorithms. By executing the algorithm, the computer program takes the input (e.g. an image) and, checking all the specified rules, decides if that image represents a cat.

In Machine Learning, we let the machine learn how to recognize a cat through examples that the programmer gives: images with the corresponding correct classification. The computer executes a Machine Learning algorithm (pre-existing), that, in turn, creates a **model**, i.e. a mathematical representation of the properties of the given examples associated with a specific decision.

Note that in the case of traditional programming, the quality of the classification system is mostly dominated by the programmer's ability to identify and program relevant properties of the objects to be analysed. In the case of Machine Learning, data plays a more central role: the more quality data an ML learning algorithm has about cats, the less likely the system will be to commit an error identifying a cat.

1.1.5 Functions of a Machine Learning System

What can a Machine Learning system be used for? Fundamentally, four functions

- 1. Descriptive
- 2. Predictive
- 3. Prescriptive
- 4. Generative

1.1.5.1 Descriptive

Explain what happened

A machine learning algorithm is, at the very core, a pattern recognition machine. Its function is to learn statistical associations (patterns) in the data. Therefore, Machine Learning can be used to analyze historical data (that is, data about the past) to generate insights about specific questions. For instance, a scientist could be interested in understanding if there is some relationship between air quality and the amount of rainfall in a given region. Suppose the scientist has historical data about these two specific aspects. In that case, she could feed the data to a machine-learning algorithm and explore, for instance, if air quality decreases or increases when it rains.

Similarly, a retail company might be interested in understanding what makes a product attractive to a particular type of client (e.g. industrial design students). Using historical purchase data combined with demographic data, a machine-learning algorithm can help find patterns between specific classes of products and buyers in a given demographic.

1.1.5.2 Predictive

Predict what will happen

Machine learning systems can also process current and historical data to make predictions about future events. *Predictive* machine learning essentially involves interpolating past information to guess what will happen next.

Take, for instance, the example of a weather forecast, that is, predicting what the weather will be like tomorrow. A ML algorithm could be used to explore if, in a given area, there are correlations between weather conditions across different hours, days, or weeks. If such a correlation exists, then data about today's weather can be used to predict tomorrow's or next week's weather.

1.1.5.3 Prescriptive

Suggest/recommend actions to take

By *Prescriptive* use of Machine Learning we mean the indication of the "best course of action to take" based on insights derived from descriptive or predictive activities. Prescriptive systems are often called "Decision Support Systems", whose main goal is to aid individuals or organizations in decision-making activities.

Recommender systems fall into this category. Applications like Amazon, Netflix, or Spotify use recommender systems to suggest to their users what to buy, watch, or listen to. Application supporting healthy lifestyles (e.g., MyFitnessPal, Headspace, Fitness+) suggests actions (e.g., eating, meditating, exercising) based on the current health status of their users.

1.1.5.4 Generative

(Semi) autonomously create new data

Descriptive, *Predictive* and *Prescriptive* Machine Learning systems are often referred to as "Discriminative" as their goal is, broadly speaking, to discriminate between different data instances (e.g. to tell a cat from a dog, a rainy day from a sunny one, a good movie for you from a bad one).

Generative Machine Learning systems, on the other hand, have as their main goal the creation of a new data instances that are "likely to be realistic". By realistic in this context, mean data instances similar to the ones used to create the machine learning system, but not the same.

These models have become very popular in recent months, thanks to the advent of tools like Open AI ChatGPT (to generate new text), Stable Diffusion (to generate new images) or Open AI Jukebox.

1.1.6 Neural Networks and Deep Learning

- Deep Learning is a Machine Learning approach based on neural networks (NN)
- NN are machine learning algorithms in which processing nodes (neurons) are organized into layers
- Depth = number of layers



Deep learning is a particular way of implementing machine learning based on *neural networks*. Neural networks are mathematical constructs inspired by the way that the brain works.

The main idea behind neural networks is the presence of several *neurons* – very simple computational units – connected through weighted *edges* and organised in *layers*. We will describe neural networks later in the course.

Deep learning networks are neural networks with many layers and complex learning architectures. Deep learning networks have very interesting learning properties that simple neural networks do not have. But they are very complicated to create. Deep learning networks were created almost 30 years ago, but we have now able to use them, thanks two 2 important technical advancements:

- 1. The availability of a large amount of digital data
- 2. The availability of high-performing computing architectures.

1.1.7 Computer Vision

- High-level understanding of digital images or videos
- Also generation (e.g Stable Diffusion)
- An enabler for technology such as smart doorbells, self-driving cars, etc.



An example of computer vision algorithm for object recognition.

Computer Vision is a sub-field of Artificial Intelligence and Machine Learning that focuses on extracting high-level understanding from images or videos. Or, more recently, to generate realistic images and videos.

As we will see later in the course, Computer Vision is a well-developed field that developed several tasks, such as:

- Detect, recognise, and identify entities (e.g., objects, faces, people, animals)
- Modify visual content (e.g., image manipulation, image restoration)
- Categorise visual content (e.g., offensive images)

• Generate new images and videos

1.1.8 Natural Language Processing

- High-level understanding of language spoken and written by humans
- Also generation (e.g. ChatGPT)
- An enabler for technology like Siri or Alexa

Natural Language Processing (NLP) is a sub-field of Artificial Intelligence and Machine Learning that focuses on analysing natural language (written or spoken) to understand its content. Or, more recently, to generate realistic text and voices.

As we will see later in the course, Natural Language Processing is also a well-developed field, that developed several tasks such as:

- Recognize the language, understand it, and respond to it
- Categorise textual content (e.g. spam vs. Not-spam, offensive vs. Non-offensive)
- Translate between languages
- Generate new text

NLP is an enabler for popular personal assistants like Siri or Alexa, and it is at the core of recent systems like ChatGPT. As we will see later in the course, NLP technology can greatly support your design process. For instance, by automatically processing textual documents to, for instance, perform thematic analysis on interview transcriptions.

1.1.9 The hard problems are easy, and the easy problems are hard

The main lesson of thirty-five years of AI research is that **the hard problems are easy and the easy problems are hard**. The mental abilities of a four-year-old that we take for granted – recognizing a face, lifting a pencil, walking across a room, answering a question – in fact solve some of the hardest engineering problems ever conceived... As the new generation of intelligent devices appears, it will be the stock analysts and petrochemical engineers and parole board members who are in danger of being replaced by machines. The gardeners, receptionists, and cooks are secure in their jobs for decades to come.

The paragraph above was written by the cognitive psychologist Steven Pinker in 1994 (The Language Instinct) to describe how, contrary to popular belief, the tasks that are very difficult to imitate are the ones that humans find the easiest, as they are performed unconsciously: vision, language, basic motor skills. On the contrary, tasks that are very difficult for humans, such as complex mathematics or playing chess, are very easy for computers.

Despite impressive advancements in all AI and Machine Learning fields, Pinker's observations remain true.

1.2 1.2 So, why should you care about Machine Learning?

The most profound technologies are those that disappear. They weave themselves into the fabric of everyday life until they are indistinguishable from it [^Mark Weiser, The Computer for the Twenty-First Century (Scientific American, 1991, pp. 66–75)].

In the first part of the lecture, we acknowledged that Artificial Intelligence and Machine Learning technology are pervasive. Perhaps they are like electricity: always present, running our world, yet invisible. And this thought always brings back to mind this quote from Mark Weiser, one of the founders of the Internet of Things and pervasive computing fields, who said, "the most profound technologies are those that disappear."

But if this is true, why is Artificial Intelligence always in the news, and not always for good reasons?

Documentaries like "The Social Dilemma" warn us of the dangers of recommendation algorithms and how they can polarise political and social discourse and have adverse psychological effects.

News items often remind us how some AI systems were connected to some injustice perpetrated by public or private organisations. And the tendency is often to attribute responsibility to the technical system. As if the system itself is at fault. When, in reality, someone has designed the technical system in a given (social/economical/organisational / political) context.

1.2.1 Why do we need Designers to understand ML?

- Focus on purpose, not on outcomes.
- Asking "Why" questions
- Understanding and acknowledging diversity of stakeholders and values
- ...



So, why do we need designers? What is the role of designers in the context of AI-based PSSs? The answer I give, is that we need these systems to be designed and not simply engineered.

I have always been more interested in what people can do with computers than in what computers can do. This is what the field of Human-Computer Interaction is all about. And design plays a big role in there.

Let me refer you to an interview I recently had.

"I grew up in a different time. I have experienced the transition from the analog to the digital world, and we were there at the emergence of the digital society. I've been playing with web technology from the beginning out of curiosity, and I have also witnessed the emergence of intelligent systems in society. That has allowed me to see how these systems have changed people and how they have changed society as a whole. Digital technologies possess enormous power. We cannot deny how they can also be disruptive; how they can even cause harm. We must learn to control them.

When you hear the word 'design', it doesn't mean aesthetics in this case, not as we use it for products such as chairs and tables. Al systems do not operate in isolation. They are part of socio-technical systems, organizations, and society. In this case, design is about the interplay between people, algorithms, and data from people used to train those algorithms. Our job is to shape all of that. When we think of engineering, we think of a straight line, from data to model to user, or from requirements to software to user. In practice, this is not a straight line but a circle. The design evolves over time, as we get a better understanding of people and a better understanding of technology; it keeps changing."

To me, design is a technical discipline that, more than many others, focuses on *purpose* and not on *outcomes*. Designers are the ones asking the *why* questions. And I believe, Designers are the ones that are trained to be able to deal with complex multi-stakeholders situations.

1.2.1.1 Design for AI video and Podcast

- Video
- Podcast

I recommend you watch these two resources above. In addition to the ones indicated as suggested reading/video.

1.2.2 What can designers do for Machine Learning?

Here you can see three important roles that designers can play.

1.2.2.1 Humane AI+ technology



First, helping create AI/ML technology that is more humane and more attuned to the needs of society and the planet. But this requires designers to be able to sit at the decision-making tables of companies like Twitter, or Google, and steer to technology-informed arguments on the development paths of the products created by such companies. Or public organizations.

1.2.2.2 Design tools for AI stakeholders



The picture comes from this thesis.

Designers can also play a big role in ideating and creating the next generation of tools that will help democratising ML technology, but also that will allow ML developers to engineer solutions that will have the properties that are really desired. The picture above is an example of scientific work conducted by master and Ph.D. students of TUDelft, where design research techniques have been applied to investigate how ML developers could debug the behavior of ML models. This is a booming field, with plenty of opportunity.

1.2.2.3 Design ML data



Another important role that designers can have is to invent and steer the data collection process for ML models. Which data should be collected? How? From whom? What are the limitations of these data? What do they represent? In this course, we are going to stress over and over the importance of data for ML technology. But you would be surprised to know how these questions are seldom addressed by engineers of ML systems. We will see later some examples of ML technology that fails. Let me reassure you: 95% of the time, the root cause of these issues can be traced to some issue with the training data.

1.2.3 What can designers do with Machine Learning?

As we saw in one of the first slides, the number of products having some form of ML technology is increasing by the day. But ML can play an important role also beyond traditional consumer electronics.

1.2.3.1 Human Augmentation



The image on the left is from a colleague in the faculty, working on this. Dr. Evangelos Niforatos. The one on the right are Envision Glasses, created by a IDE alumni.

1.2.3.2 ML for Fascination and Engagement



Frederik Ueberschär

1.2.3.3 Interaction



A snare-no-more device milles.

Explore a visual chronicle of frontline photofournalist No Amin's archive with the bala of

Colorizing Klimt's Vanished Paintings with Artificial Intelligence and Klimt Experts

Play an invisible drum kit

Control your devices with the wave of a finger.

A code-free tool that lets you create custom, microcontroller-ready models base on DRU data.

Experiments with Google. 1612 and counting...

If you look around a little bit, you can find tons of examples of how we can play around with machine learning as a design material in and of itself to enable a new way of interacting you way of our relationship between people and technology to support people.

And you can find the lock out there. We try to put some examples and additional material for the course, but if you spend a little bit of time you will see you will find a lot.

1.2.4 What can Machine Learning do for designers?

1.2.4.1 Co-create



1.2.4.2 Inspire



Dall-e

1.2.4.3 Scale up!



Thesis Document : Analysis of how parents perceive their baby, their behaviours towards their child, and thus understand how overprotection develops throughout childhood more than 300 stories, manually and NLP analysis

1.2.4.4 Scale up!



How to help designers, experts, and societal stakeholders work together with AI, to prepare, realise and evaluate design interventions? Goal: reduce design complexity for large-scale social interventions

D@S Lab

1.2.4.5 Understand design



Using big data ... we experiment with artificial agency during complex system design processes We are exploring the form and use of novel design methods to address systemic design problems to create an AI Toolkit

Design Intelligence lab

1.2.5 Why Programming?

All design needs a medium. A designer in the age of computable technology also contends with programming, which the designer wields as a tool and canvas. Ge Wang - Stanford

1.2.6 Debunking some myths



1.2.6.1 Expectations



1.2.6.2 Reality

Zillow wrote down millions of dollars, slashed workforce due to algorithmic nome-buying disaster	 C stan kont ← ○ ○ 	nobelný kolej 🗙 🔶 Č. Měst, jeven ženy dolej kolež ži pověsta kolest	ster-storing-ballydigestatur to Strattage-Myries Margore	tesenustrationgetasys P. Q.)	- Q + 6
n November 2021, online real estate marketplace Zillow told shareholders it vould wind down its Zillow Offers operations and cut 25% of the company's	🚼 Warson't Bing	0, 1946H	Welcome to the new Bin Your # powered enswer engine	cur 🚷 14	• 🗊 a =
vorkionce — about 2,000 employees — over the next several quarters. The		Ask complex questions	🔒 Get better answers	Set creative inspiration	
arning algorithm it used to predict home prices.		What are some meals I can make for my picky todder who only east orange- coloured fixed?	"What are the pres and core of the top 3 selling pet vacuums?"	Write a halou about exceeding in order space in the vector of a pinets"	
ow Offers was a program through which the company made cash offers properties based on a "Zestimate" of home values derived from a chine learning algorithm. The idea was to renovate the properties and them quickly. But a Zillow spokesperson toid CNN that the algorithm		Left i l	earn together. Bing is powered by AL ao susprises a exible. Nake zure to check the fuct, and share fee can kern and improve! terrs at twi	nd mistakes Back no we elet is sotar storing to	
iad a median error rate of 1.9%, and the error rate could be much higher, as nuch as 6.9%, for off-market homes.		Insurching for when is analar shaving tasky Tourching for analar mode Oramating arrangem for equ.			
IN reported that Zillow bought 27,000 homes through Zillow Offers since launch in April 2018 but sold only 17,000 through the end of September (21. Black swan events like the COVID-19 pandemic and a home novation labor shortage contributed to the algorithm's accuracy troubles.		Helis, this is Ding. Language you are referring up a science factor film set or the allen work down and its or commy playing in the atom works and the gradience it. December 14, 2021 ¹¹ . No can for the these ¹⁵ , You can also warsh the official toward to Lane news. ¹ Strendpointing ¹² whice one	to the month Analise by James Gameson, which is a 2000 of Plendors ¹¹ , ²¹ . There are no bandwises for this means linearce, you can how you then the XVD of the law you have the the NVD of Vlaster. As indexided to be released on advances of the NVD of Vlaster. ²¹ . It also assues The Way of Vlaster. ² A sequent: A summaria care: A summaria		
llow said the algorithm had led it to unintentionally purchase homes at		🗸 Section for weaks the map of water at streng	a in Manhpurd share linear being	of water at cinemaa in blackpool show frees to	
gher prices that its current estimates of future selling prices, resulting in a		🗸 Generating arowers for pro.			
04 million inventory write-down in Q3 2021.		Skirne arything.			

My new favorite thing - Bing's new ChatGPT bot argues with a user, gaslights them about the current year being 2022, says their phone might have a virus, and says "You have not been a good user"

Why? Because the person asked where Avatar 2 is showing nearby



"48% of US consumers intend to buy at least one smart home device in 2018"^a "23% of connected security system owners said they **deactivate their system completely** when they have guests over"

^{*a*} Survey of 2000 US Consumers. Ooma

1.2.6.3 AI/ML can predict the future

- AI/ML are "statistical parrots"
- They are (very good) pattern recognition machine
- Garbage in Garbage Out

1.2.6.4 AI/ML has agency

- AI/ML are tools.
- People design and use them.
- And they change us!

1.2.6.5 AI/ML can magically transform a PSS overnight

- Magically: maybe
- Overnight: No

1.2.6.6 AI/ML can solve any problem

- AI/ML technologies are very flexible and powerful
- But they have very strict requirements
- And potentially harmful limitations

2 Section 2: Fundamentals of Machine Learning

2.1 2.1 The Machine Learning Life-Cycle



The Cross-Industry Standard Process for Data Mining (CRISP-DM) methodology

There is no standardized methodology for designing Machine Learning systems, although several companies have created several proposals in the past. For instance, this CRISP-DM methodology: Cross Industry Standard Process for Data Mining is the closest you might get to a standardized process for Machine Learning system development. It's not exhaustive, and as we will see, it has some limitations, but it's complete enough, and I think we can use it as a reference moving forward. At least to understand the different actors and the different activities.

There is a *problem owner*: Someone that is interested in designing this particular PSS that has a Machine Learning-component. You also have diverse stakeholders: people that, in one way or another, have a stake in, if not the design, at least in the output and the consequence of creating this particular Machine Learning-driven solution.

Data engineers: people, mostly engineers, that deal with acquiring and managing the data that will be used later in the ML process.

Data scientists: experts in statistics, data science, and machine learning.

2023-03-29

Model validators: those in charge of checking that the system performs as expected. They can be testers operating before the PSS is put into production. Auditors, if their job is to check if the PSS complies with specific regulations and policies. Or, sometimes, end users provide feedback on the product's performance (e.g. a thumbs-up/thumbs-down on a recommendation).

Operations engineers: are the ones in charge of making sure that the system runs continuously and effectively. Sometimes they are also called *ML OPs engineers*.

Let us use as an example the recommender system of Amazon. Amazon is the *Problem Owner* for that recommender system. It is a company interested in maximizing revenues using automated procedures suggesting items consumers operating on the platform should buy. Of course, the stakeholders of these systems are the users and companies selling products through Amazon. One could argue: Even society as a whole, because the more gadgets we sell, the less sustainable and the more we're going to waste.

Within Amazon, there are *Data Engineers* that take care of managing all the data that is needed to run these systems. For instance, the databases containing all the review data from people and products, all the products' details coming from the suppliers. The *data scientist* is the professional that creates the Machine Learning system, the recommender system, in this particular case.

The *validator* checks that the performance indicators expected by the problem owner are met (e.g. performance or revenue). The *operations engineer* looks at the system as such: The computers, the network. That everything runs as quickly and as effectively as possible. But, as I mentioned before, validators can also be users, giving feedback on the quality of the recommendations.

In a big company, all these people might be different, in small companies or startups, they are a single person. But in a sense, those are the roles that are somehow involved.



2.1.1 CRISP-DM In our course

We are not going to cover in detail all these steps. For instance, we will not cover the deployment and monitoring step in this course. This is an advanced topic, and it's really about systems engineering, probably not interesting for you.

We will cover the problem specification and the evaluation of the system. We will study and look at how we can specify what the system should do and how we will evaluate it.

Data understanding and modeling we will do throughout the course. Data preparation we will address in module three.

2.1.2 Problem Specification

- What is the problem owner hoping to accomplish and why?
- Why am I (being asked to) solve it?
- Am I the right person to solve this problem?
- What are the (psychological, societal, and environmental) repercussions of building this technology?
- Should this thing be built at all?
- What are the metrics of success?

This is where you, as designers, should be more comfortable. The process of understanding what the problem owner hopes to accomplish and why. And, to get to the bottom of what exactly are the desires, needs, and requirements of this particular solution that you're asked to design and develop.

Of course, this also implies considering ethical issues: Should the system be developed? How are we going to measure success? What impact will the system have, maybe, on the broader set of users or society?

Those questions should be accounted for before creating the system, not after.

2.1.3 Data Understanding

Know your data! - Data need to be collected —> Datasets - What data is available? - What data should be available but isn't? - What population/system/process is your data representing? - And what properties of such population/system/process are included (or excluded)? - What biases (social, population, temporal) are present in your datasets?

Data understanding might seem like a purely technical task, but it is not, purely. Data is what makes or breaks a Machine Learning system. So, it is essential to understand what data needs to be collected. Whether this data is available or not. Even if there is data that should be available but is not there. What part of the phenomenon being modeled can the data capture, and which is not? So that if biases exist, they could be identified early.

To give a concrete example, consider Social Media (e.g, Twitter) data. One might be tempted to say that data collected through social media, being so abundant, are representative of very large, exhaustive, and comprehensive populations. We know that not to be true. We know that different demographics, and different people, operate on different social media, and we know that they do that in different ways. Quantity does not mean comprehensiveness and diversity to the point that it suits a particular problem. So being VERY inquisitive about the data is a very pragmatic first step for everyone that wants to do something with Machine Learning.

2.1.4 Data Preparation

- Data integration
 - Extracting, transforming, and loading (ETL) data from disparate relevant databases and other data sources
 - This step is most challenging when dealing with big data sources
- Data cleaning
 - Filling missing values
 - Transforming value types (e.g. binning)
 - Dropping features that should not be considered
- Feature engineering
 - Transform the data to derive new features

Data does not exist in a vacuum, and it is not created "perfect." It is the result of complex socio-technical processes that might lead your data to be "messy," Incomplete, with errors of different types. Data comes from multiple sources, so they must be aligned and integrated.

In the real world, the data work is 80% of the work, while the modeling work is only 20%. Data understanding and preparation take (or should take) most of the time in the design and development process for ML systems.

2.1.5 Modeling

- Select a training algorithm
- Use it to find patterns in the training dataset
- Generalize them to fit a statistical model
- **Enhance** the model to satisfy additional objectives and constraints captured in the problem specification
 - e.g., increase reliability, mitigate biases, generate explanations

Once the problem is defined, and the data are prepared, then we can start to pick the right Machine Learning model and the right training algorithm. This is [in theory] conceptually easy, but it is complex in practice. And the reason is that **there is no such thing as a best possible solution applied in all possible conditions**.

- No free-lunch theorem
 - There is no one best machine learning algorithm for all problems and datasets

That is the so-called "no free lunch" theorem. There is no one best Machine Learning solution for all problems and datasets. We have a portfolio, and a bunch of potential alternatives, each one is more or less suited according to the class of problems you have and the type of data you have available.



The process of training a Machine Learning Model

The slide mentions four steps. The *training* and the *enhancement* are activities performed by data scientists through algorithms, while the *model* is the final result.

The Machine Learning *model* is the computer program (in reality, it often looks more like a configuration file) that is created using injecting the data and letting the system learn from this data. **Model training** is the activity performed by a **training algorithm** (configured by the data scientist) to create that Machine Learning model. For instance, given many examples of cats and dogs, the **model training** step is concerned with finding the statistical correlation between the input data and the classes (*cat* and *dog*) that we want to recognize. The **model** is the program that, based on the observed statistical correlations, can distinguish a cat from a dog.

We are separating *training* and *enhancement* because, oftentimes, activities happen *after* the training of a model to ensure that the model behaves according to expected behavior. For instance, to check (and eventually fix) biased behaviours against specific type of data, or the people (e.g. a specific population or minority) that this data represents.

2.1.6 Evaluation

- Testing and validation of the model
 - Also against the problem specification requirements
- Performed on data not used for training
 - Hold out dataset

We will address evaluation in a following section.


Splitting Data for Evaluation

2.1.6.1 Model auditing/risk management



Evaluation is also a very important step because this is where one can check whether the system that has been created is compliant with its initial expected properties. Evaluation is complex for several reasons that we will explore later in the course.

These systems are stochastic, and they're probabilistic. This means that, in practice, they will **never** with 100% accuracy performance. A Machine Learning model trained to recognize **cats** versus **dogs**, is bound to make wrong classifications, sometimes. There are many reasons a prediction can fail (mostly due to the properties of the data).

ML systems *never* operate perfectly. Sometimes, they fail. But we use them because they are helpful, not perfect. It is essential to know how to deal with those errors. This is an interesting interaction design and design-interaction challenge.

But also an ethical one: if it is true that these systems are not perfect, when (if ever) shall we use them? Are there situations where perhaps we should not use them at all? This is the core of the current, ongoing efforts to regulate Artificial Intelligence in Europe. It's an ongoing process built on a risk-based approach in defining under which conditions and under which particular situations the risk of using one of these systems is acceptable.

There are four classes: *Unacceptable, high, limited* and *minimal* risk. And each one of these risks is associated with particular demand in terms of required checks and controls on the performance of the system, the understanding of its potential limitations, and, of course, also compensatory actions.

To give some examples. The recommender system of Netflix will probably fall in the minimal risk appli-

cation. In the worst-case scenario, if the system does not work well, a user receives a recommendation for a movie they do not like: tough luck!

Risks are considered *Limited* when the malfunctioning of the system might have effects that have some consequence, but they are not too problematic. For instance, a system that is designed to classify whether content is inappropriate fails 5% of the time. There is some risk: people can be offended or traumatized. But as long as these mistakes can be corrected in the future, it is fine. AI chatbots, spam filters, and customer management sustems fall into this category.

High risk are those systems that "either pose a health and safety risk or pose a risk to people's fundamental rights". And here, the requirements are very strict. The legislation includes in this category:

- biometric identification and categorisation of natural persons
- management and operation of critical infrastructure
- education and vocational training
- employment, worker management and access to self-employment (e.g. recruitment systems)
- access to and enjoyment of essential public and private services and benefits
- law enforcement
- migration, asylum and border control management
- administration of justice and democratic processes

Note that the list of critical areas are still in discussion, and currently subject to variations. For instance, the Parliament is currently considering placing conversational and art-generating AI tools such as ChatGPT and DALL-E-2 in the high-risk category.

Unacceptable risks happen in those situations where under no circumstances should it not be allowed for an A.I. system to operate. These include:

- "subliminal, manipulative, or exploitative techniques causing harm"
- "Real-time, remote, biometric identification systems used in public spaces for law enforcement"
- All forms of social scoring

Some people believe, for instance, that we should not be allowed to design a system for military purposes. Of course, other people disagree.

2.1.7 Deployment and monitoring

- What data infrastructure will bring new data to the model?
 - Will predictions be made in batch or one-by-one?
 - How much latency is allowed?
 - How will the user interact with the system?
 - * Is there a problem here?

- Tools to monitor the model's performance
 - * And ensure it is operating as expected

Deployment and monitoring is about the actual operation of the system. It is a phase that addresses important needs, although not necessarily needs that a designer would care about too much. For instance, *latency*, how quickly a system can respond to user interaction. A system like Google Search, for instance, is engineered to answer queries in less than 100 milliseconds. If for whatever reason, the intelligence behind Google is slower than that, it has been empirically demonstrated that the company will lose business and they will lose money. So, it is their best interest not to delay the user experience. MLOps Engineers are tasked to ensure that the system works as fast, as fluidly and as, well, as best as possible. There are, of course, trade-offs. Not all systems can be, for instance, accurate and also be fast. Take for instance a computer vision system installed on a car: if such a system is not powerful enough, then recognising an object might be too slow, and so the car might exhibit a higher reaction time and eventually bump into something or somebody. That is unacceptable. A solution would be to put a more powerful system on the car (but more expensive); or let the system run faster but will less accuracy.

2.2 2.2 Data. The raw material

Data, or the raw material of A.I. If you have read the suggested reading material for today's lecture. "I, pencil", which is a lovely, short essay. You'll know where this reference is coming from.



2.2.1 Data

Structured data: a common representation for data in Machine Learning

In our collective imaginaries, data looks like a table. Or so-called **structured data**. And, most of the time, that is the case. And, it turns out, it is a useful and generalisable abstract representation for ML purposes.

A **Dataset** is a collection of **record**s.

A **record** is a row in a table; in the context of ML, a record is called a **sample**, a **data item**, or an **instance**. In the *IOB-3-22 Data Course* they were called **Data points**, with the following definition: "units of information, usually collected through observations by people or machines. Data points can be single facts or observations (e.g., the age of one individual child incorporated in the dataset, or the temperature measured at one location at a given time), or multiple corresponding values (e.g. set of environmental variables)". This definition is, of course, compatible with ours.

A record is composed of **attributes**, the columns. In Machine Learning parlance, those are called **features**. A feature is a property of the world modelled by the data; the property is captured, measured, and represented for the prediction problem we have in mind.

Let us use the example in the slide. The goal is to create a system that can recognize the species of a flower based on the geometrical properties (width and length) of their petals and sepals. Notice how, in this case, someone decided that a *flower*, a creature, is represented by these attributes; other attributes, like colour and shape, were not included. A record therefore represents a real-world flower; the value of the features (*feature values*) are the ones of that flower.

Later we will see that features might take different representations according to the data type (e.g., an image or a textual document).

There is a very specific column called **label**, or **class**, that contains the value that the system is expected to be predicting – i.e., the prediction output. So, for instance, if you want to train a Machine Learning system that can recognize flower species (given a set of sepal length, width, and petal length and width), the *label* allows the algorithm to find a statistical correlation between feature values and the flower type.

Note that a label might not always be present.

The **dataset size** is the number of records available in our dataset. Dataset sizes matter in Machine Learning because, as discussed in previous lectures, ML methods are data-hungry: the more data you have, the better (usually). **Dimensionality** is the number of features representing a particular entity. Later we will encounter a concept called *curse of dimensionality*: the problem that too many features represent a data item.

2.2.2 Types of Features / Label Values

- Categorical
 - Named Data

- Can take numerical values, but no mathematical meaning

• Numerical

- Measurements
- Take numerical values (discrete or continuous)

Data values can have different types. A type holds properties that allow doing some specific operation with data values.

You must understand the difference between these data types. As we will see later, Machine Learning systems, being statistical systems, prefer to deal with numerical values more than categorical ones. It is a common practice, for instance, to translate categories into numbers and text into vectors. We will see several examples of this in the next modules of the course.

It is also important to distinguish the notion of **Data Types** as you have encountered in the *IOB-3-22 Data Course* from **Feature Type**. Data types include: integer (any whole number, without a decimal point, e.g., 7); float (any number with decimal point, e.g., 7.43) date (a particular year, month, and day in time); time (a particular moment of day); text (any collection of letters rather than numerals or other symbols, also referred to as 'string' or 'character'); and boolean (a binary data type with only two values: i.e. either TRUE or FALSE; Yes or No; 1 or 0).

It is very common to deal with *categorical* data. These data could be either textual or numerical, but for which **there is no mathematical meaning** attached to it. Take, for instance *Marital Status*: there is a close number of statuses. It is a string (e.g., "Married", "Single"), it could also be represented as a number (e.g., "0", "1"), but no arithmetic operations can be performed ("Married" - "Single" = ?).

- Categorical Nominal**
 - No order
 - No direction
 - e.g. marital status, gender, ethnicity

Marital status is a **nominal** categorical data type because there is no way to order values (Is "Married" greater/better than "Single"?).

- Categorical Ordinal
 - Order
 - Direction
 - e.g., letter grades (A,B,C,D), ratings (*dislike*, *neutral*, *like*)

Ordinal data are categorical data with an order and a direction. Think, for instance, of letter grades: *A B C D*. We know that A is greater than B, and B is greater than C, so there is an order and a direction, but still, you cannot do A minus B or C minus D. It doesn't make any mathematical sense. Another example

of *Ordinal* data are the classical Likert scales – e.g., *Dislike*, *Neutral* or *Like*. *Dislike* is worse than *Like*, but *Like* - *Dislike* is not equal to *Neutral".

All of those are categorical data, we use them all the time, and, from a statistical perspective, they come with a lot of strings attached because there are things that you cannot do with this type of data that you can do, instead, with the other type. We cannot dig deeply into this topic, but it will be extensively addressed in the Design Analytics [IOB6-E8] elective course.

- Numerical Interval
 - Difference between measurements
 - No true zero or fixed beginning
 - e.g., temperature (C or F), IQ, time, dates

Numerical data is used when measurements are involved, e.g., measuring someone's height or weight. Numerical values can be continuous (i.e., real numbers) or discrete (i.e., integers); and, of course, they can be part of mathematical operations. Note that boolean values (0 and 1) are also numerical.

- Numerical Ratio
 - Difference between measurements
 - True zero exists
 - e.g., temperature (K), age, height

Ratios are numerical values used to represent properties like *age*, *height*, *weight*, *width*: these are properties for which there exists a *true zero*. Temperatures, for instance, but measured in Kelvin, where there is a zero – the lowest possible temperature there is in the known universe. But that's not the case, for instance, with temperature in Celsius or Fahrenheit, where there isn't a zero. The same goes for dates and times. You can subtract them, for instance, one year minus the other, but there isn't a scale and a fixed point. We call them **intervals**.

2.2.3 Data Modalities

Regarding modalities, data can be split into two big families: **structured** and **unstructured** (or semistructured) data.

Structured data is data that comes in some predefined format: Tabular data, like the one that we saw in the example before. This is data for which there exists a pre-existing organization: attributes are explicitly defined.

Unstructured data does not have an intrinsic structure—for instance, images, music, and text. Of course, images contain shapes and patterns, but there is no explicit way to describe such pictures because, essentially, an image is a collection of pixels.



A summary of the different types of data that could be processed through Machine Learning techniques.

In the same way, you can think about *text* being a collection of words. Textual documents like papers, newspapers, etc. indeed have some structure (e.g. introduction, conclusion). However, how the words are organised in a text does not follow any specific structure – besides, perhaps, grammar rules.

Modality	Quantity	Quality	Freshness	Cost
Structured	Number of records	Errors	Rate of collection	Acquisition
Semi-structured	Number of features	Missing data		Licensing
		Bias		Cleaning and integrations

2.2.4 Key Dimensions

What properties make data "good" or "bad"? "Useful" or not "useful"?

Modality as discussed previously.

The **quantity**. The number of records, the *dataset size*. But also the *number of features*, as we will see also later on in the course, the number of features available can either be a blessing or a curse: the more features available, the more detailed description of a data item; but the more features, the more the data items we need to have to train a good machine learning model.

The **quality** of the data is also very important because the data might contain mistakes. To believe that the data is "correct" is, most of the time, wishful thinking, an illusion. In a dataset, data is often missing due to collection errors or technical issues. And data is ==always== biased for multiple reasons: 1) the

collection process, 2) the specific condition in which the data has been acquired, 3) the design choices in the data acquisition process itself. You had an example before about social media data: biases can creep on because there is an uneven distribution of people and topics. Remember: there is no such thing as objective, neutral, and unbiased data, there is always a bias in the data.

Freshness means how often new data is generated. Some data sources, like census data, are updated in the best cases every year and in the worst cases every decade. Freshness (or lack thereof) can be a problem depending on the ML goals: some system requires fresh data (e.g., a news recommendation system, a malfunction prediction system), while others don't (e.g., a system recognising defective items in a production line).

To track your heart rate and to predict whether you are going to have some heart issue, you might want to sample your heart rate frequently. Higher sampling rates might not be needed.

Cost. Cost is also significant because cheap data is not necessarily good data. And there is a cost for cleaning and integrating it, so that's also a consideration.

2.2.5 Types of Data



2.2.5.1 Static Tabular Data

Structured data.

This type of data we have seen before.

2.2.5.2 Time Series

- tabular data with **time** feature
- For instance
 - Sensor data, Stock market data

· Label is usually associated with a set of records - e.g. a continuous movement of the phone indicating an action Acceleromete Time Timestamp X Ζ Class У Feature 15060015925 2.04 3.72 8.12 Device 15060015943 1.96 4.73.68 7.56 Rotation 15060015980 1.63 3.56 6.53 1506001610 1.06 3.76 5.81

Time series data are structured, like tabular data, but have two fundamental differences. 1) One of their features is a time-stamp: this is because the reality that time series represent has time as a critical component. 2) Labels can be associated with *sets* (more than one) data items, because what needs to be recognized spans multiple time stamps. For instance, recognizing the type of activity a person performs (e.g., walking vs. Running) based on accelerometer data requires more than one data item.

Time series data are very interesting data, but we are not going to cover them in this course.

2.2.5.3 Images

- Visual content acquired through cameras, scanners, etc.
- Each pixel in an image is a feature
 - But spatially and geometrically organised
 - * e.g., edges, corners
- Feature values are numerical values across channels

– e.g., *R,G,B*

Dimensionality -> n x m



Images are nothing more than a set of pixels. And each pixel in an image is a **feature** having a value (e.g. scale of grey, RGB, CMYK) that tells you something about that particular pixel. Now you can understand why I mentioned the problem of the "curse of dimensionality" before: in an image, especially a

high-resolution one, the number of features can be in the order of millions, making computation very expensive.

There are ways to avoid the curse of dimensionality: for instance, in traditional computer vision, a process called *feature engineering* was used to compute features (e.g. edges, corners, or shapes) from images *before* training a machine learning system. We will discuss this issue more in the coming module.

2.2.5.4 Textual documents

- Sequence of alphanumerical characters
 - Short: e.g. tweets
 - Long: e.g Web documents, interview transcripts
- Features are (set of) words
 - Words are also syntactically and semantically organised
- Feature values are (set of) words occurences
- Dimensionality -> at least dictionary size

	Document	•	Wear	Mask	W(n)	Class
		1	1	1	0	Spam
http://www.children, 24 May 2015		0	0	1	0	Not Spam
By Sir Chubs		0	0		0	Not Opani
Verified Purchase (What is this?) This review is from: Overhead Rubber Penguin Mask Happy Feet Animal Fancy Dress (Toy)						Spam
I wear this mask to sing lullables to my children. They are terrified of the mask. Whenever they protest about their bed time, or ask for too many sweets, I whip on the mask, and they soon know who is the King Penguin.						

The issue of dimensionality also occurs with textual documents. A textual document is nothing more than a sequence of words. Texts can be short (a Tweet) or long (a novel). They can be created in many scenarios and applications: e.g. reviews on Amazon or interview transcripts.

As words have their meaning, it is very common to consider each word in the dictionary a feature that can either be present (or not) in a document: we call this representation "bag(s) of words". The value of the feature is whether the word is present, or not, in a particular document.

There are also more complicated ways to represent documents, that will look at sequences of words, i.e. words in their order.

2.2.6 Categorising Data Sources

Purposefully			
Collected Data	Administrative Data	Social Data	Crowdsourcing
Survey	Call records	Web pages	Distributed sensing
Census	Financial transactions	Social Media	Implicit crowd work (e.g. captcha)
Economic Indicators	Travel Data	Apps	Micro-work platforms (e.g Amazon Mechanical Turk)
Ad-hoc sensing	GPS Data	Search Engines	

Different data sources have different properties in the five dimensions mentioned above: *modality*, *quantity*, *quality*, *freshness*, and *cost*.

Purposely collected data are collected specifically for a specific goal: Surveys, censuses, and scientific experiments done with very expensive equipment. They are created using ad-hoc sensing infrastructure specifically designed to collect the data for the prediction purpose that we have in mind. Most of the time, this data comes in a structured way. They are very expensive to acquire, and they have very high quality but might not have high freshness—for instance, census data. Data from sensors or sensing infrastructure are high-frequency and high-quantity, but the cost is very high.

The following type of sources includes **re-purposes* data. That is, data that is not created for the prediction purpose at hand but they are reused from other scenarios or application domains.

Typical examples are *call records*, which is data collected by network operator companies when your phone downloads data, or when you are making a phone call. It is very common to perform prediction tasks on this data: for instance, to predict the behavioral or consumption patterns of people. Another example is financial transactions with a credit card. These are data that are not created to infer the consumption properties of a person, they're just created for the sake of keeping track of who spends money where. Another example are travel data from the "OV chipcard".

Mostly structured, very high quality, very quantitative, sometimes also very high freshness, but then again, they're very, very expensive. Either you need a very expensive infrastructure to acquire them or they are privacy sensitive.

The data that we use the most are **social data**: Data from the web, from social media, from apps and search engines, billions and billions and billions of documents, images and videos. Mostly, this is unstructured: Images, audios, text. They are low-cost to access, but the quality is not necessarily high.

Finally, there is **crowdsourcing**: distributing data collection and processing to people, implicitly or explicitly. For instance, every time you solve one of those "CAPTCHA" for Google, you are performing

a data labeling operation. People can also be paid for this on microwork platforms like Amazon Mechanical Turk. In this case, all data properties are primarily mid-low because you need to pay people, but these people are not your employees. For this reason, quality also varies greatly.

Purposefully			
Collected Data	Administrative Data	Social Data	Crowdsourcing
<i>Modality</i> : mostly structured	<i>Modality</i> : mostly structured	<i>Modality</i> : mostly semi-structured	<i>Modality</i> : all
<i>Quantity</i> : low	Quantity: high	Quantity: low	Quantity: mid-low
<i>Quality</i> : high	<i>Quality</i> : high	<i>Quality</i> : low	<i>Quality</i> : mid
Freshness: low	Freshness: high	Freshness: high	Freshness: mid
<i>Cost</i> : high	<i>Cost</i> : high	Cost: low	Cost: mid-low

2.3 2.3 Categories of Machine Learning

2.3.1 How do machines learn?

How do machines Learn? We encountered this diagram earlier in the lecture. It is important to distinguish two elements here: *Model Training* is the activity through which a machine learning model is created. This is an algorithm, typically available in a software library that allows you to create a ML *model*. But, what is a model, exactly?

2.3.2 On Models

- A physical, mathematical, logical, or conceptual representation of a system, entity, phenomenon, or process
- A **simple(r)** representation of reality helping us understand how something works or will work.
 - Not truthful, just a useful one
- The goal of models is to make a particular part or feature of the world more accessible to understand, define, quantify, visualise, or simulate

Generally speaking, when we talk about models, we refer to physical, mathematical, logical, or conceptual representations of a system, an entity, a phenomenon, or a process. For instance, the "model of a product" is a representation of a product (e.g. a sketch, a schema, a low-resolution prototype) that is not as good, as complete, or as detailed, as the real product. Still, it is good enough to understand how or how it could work. A model is not meant to be a truthful and accurate representation of reality. It is meant to be a *useful* representation, that has the function of making a particular part or feature of the world more accessible to understand, define, quantify, visualise, or simulate.

2.3.2.1 Examples of models

- Architecture plans
- Maps
- Music Sheet
- Mathematical laws of physics!
- Machine Learning (statistical) Models



Models are everywhere around us, even if we don't realize it. A *music sheet*, for instance, is model: it is a structured, visual representation of sequences of musical notes. It represents how a musician should play a tune, and of how such then should sound. The representation (the notes) does not match 100% the final sound: that is the task and responsibility of the musician. A *map* is not a high-fidelity representation of the world but is good enough for you to navigate in it. Architectural plans *blueprints*, even the laws of physics, in a way, are simplified representations of reality that we create to more easily understand, define, quantify, visualize, or simulate a particular part or a particular property of the world.

Machine Learning models are statistical models that we know are created from data. The availability of a mathematical representation can, however, give the wrong impression that machine learning models are "models of the world", as scientific models can be. So, let us spend some time discussing differences.

2.3.2.2 Scientific Models

- GOAL: explain reality
 - Created to make predictions about the outcomes of future experiments
 - * e.g., apples on the moon

- Tested against the **outcome**
- If data from new experiments don't agree, the model has to be modified/extended / refined
 - Falsifiability
- Scientific models should be *small* and *simple*.
- They should generalize phenomena observed in new ways.



At school, we encounter the scientific method. Given a reality of interest, through observations and measurements, a scientist tries to understand if mathematical laws could explain regularities in the collected data. For these laws to be correct (that is, a good explanation of reality), they should allow making predictions about the outcome of experiments performed in different contexts. If the prediction holds, we deem the model, the physical law, true. If it doesn't, we consider it falsified, and in need of an update. Consider for instance the law of gravity. We know that the very simple formula [F = mLa] allows us to accurately predicts not only the movement of objects on earth, but also in space, or on other planets.

2.3.2.3 Machine Learning Models

- GOAL: describe the data
- Designed to capture the *variability* in observational data by exploiting regularities/symmetries/redundancies
- A good ML model doesn't need to explain reality, it just describe data
- They don't need to be simple or transparent, or intelligible. Just accurate
 - Black box
- ML models may be large and complex.

- They should generalize to new data obtained in the same way as the training data
 - Same application context and data acquisition process



Machine Learning is about *representing the data*, **not** about *representing reality*. More specifically, ML approaches try to capture regularities in the variability in the data. An implied assumption when working with ML models is that such data is somehow good enough to describe a reality of interest. So it follows that by being able to describe data that describes reality, one could "describe" reality. This is a logical inference with profound consequences and the major source of misunderstanding on the actual capabilities and limitations of an ML model.

Machine Learning models are complex mathematical artefacts that are supposed to be very accurate in capturing the data on which they are trained upon. *Generalization* for Machine Learning is not about being able to make predictions about any possible data acquisition contexts, but just in the same context as for the original training data.

Note that ML models can models pick up on *spurious correlations*: input-output associations that are just a bi-product of data properties – think about The Neural Net Tank Urban Legend we discussed in class.

These spurious correlations can lead to incorrect or over-simplified predictions that can again lead to safety risks.

This is why the ==know your data== action described earlier is so important: *Where is the data coming from? What does (and doesn't) it represent? How much can the same data be used to operate in another domain or for another problem?*

By having a vast amount of observations (data), in many different contexts and many different situations, it would be possible (hypothetically) for a Machine Learning model to distill a simple, elegant law of physics. But that would happen in a much more complex way than what could be achieved by

reasoning and intellect.



Credits: Zoltan Szlavik. Benjamin Timmermans

2.3.3 Machine Learning Approaches

2.3.3.1 Supervised Learning

- Input: **labeled** data
 - Data + expected prediction
- During training, labels are used to associate patterns with outputs
- Learns how to make input-output **predictions**
- Classification
- Regression
- Ranking
- Recommendation

Supervised learning is a type of machine learning that is trained over **labeled** data to create models that can predict the labels of new data items.



Supervised learning is one of the most common types of ML approach, used in many real-world applications. To give a few examples:

- Predicting house prices based on the house's size, number of rooms, and location
- Predicting tomorrow's weather based on yesterday's temperature, humidity, and environmental factors
- Detecting spam and non-spam emails based on the words in the e-mail
- Recognizing faces/objects/animals in images based on their pixels
- Recommending videos or music to a user (e.g., on YouTube or Spotify)
- Diagnosing patients as sick or healthy

We can identify two main types of tasks that could be performed through supervised learning: Classification and Regression (see later).

Supervised learning models can also be used for two other tasks: - **Recommendation**: the machine learning model learns how to predict set of categorical data. The output of a recommendation model is a set (a collection) of data items. For instance, a set of videos to watch next. - Ranking: the machine learning model learns how to predict list of categorical data. The output of a ranking model is a list (an ordered set) of data items, ordered according to given criteria-for instance, the results from a search engine.

2.3.3.1.1 Classification

- Learn to output a category label
- Binary
 - e.g. Spam / not Spam, Cat / not cat
- Multi-class
 - e.g. cat, dog, bird



Classification: the machine learning model learns how to predict **categorical data**. The output of a classification model is a *category* (e.g. spam / not spam; cat / dog; sick / healthy)

2.3.3.1.2 Regression

- Learn to output one or more **numbers**
 - e.g., value of a share, number of stars in a review



Regression: the machine learning model learns how to predict **numerical data**. The output of a regression model is a *number* (e.g. tomorrow's temperature, the value of a stock, the cost of a home, the expected days of recovery for a sick patient).

2.3.3.2 Unsupervised Learning

- Input: **unlabeled** data
 - The machine learns structures (patterns) from the data without human guidance
 - Clustering

- Dimensionality Reduction (e.g. Large Language Models)
- Anomaly detection

Unsupervised learning is a type of machine learning that is trained over **unlabeled** data. The goal of the model is to extract as much information as possible from the dataset.



To give a few examples:

- Customer/User segmentation, i.e. grouping people in different categories based on data representing their properties or behavior
- Detecting fraudulent behaviour in payment transactions
- Detecting defective mechanical parts
- Categorize documents (e.g. articles, interviews) based on the words in the documents
- Discover trends I customer behaviour.

We can identify three main types of **tasks** that could be performed through ubsupervised learning: **Clustering**, **Dimensionality Reduction**, and **Anomaly Detection**.



Examples of clustering

2.3.3.2.1 Clustering Clustering algorithms group data into *clusters* based on **similarity**. The idea is to create groups (clusters) of items that are more similar to each other than to those in other groups.

Similarity is calculated through a *metric*, a function determining how two data items are close to each other in a given feature space. For example, customers could be grouped based on their demographics (age and gender segmentation), buying habits (e.g. types of products they buy), or a combination of different features. Documents could be grouped by the frequency of co-occurrence of certain words.



Examples of dimensionality reduction

2.3.3.3 Dimensionality Reduction

Dimensionality reduction is a task aimed at *simplifying* data before applying other techniques (e.g. supervised learning). Simplification here means reducing the number of features in a dataset by finding the ones that are more *discriminative*, i.e. they can capture diversity in the data the most. Dimensionality reduction could also mean *transforming the data*: features could be "merged" so that as little information as possible is lost while keeping the data as intact as possible.

Note that dimensionality reduction techniques are at the core of many (if not all) **generative machine learning** techniques and tools you can find on the news daily. Given a dataset, generative learning techniques are trained to generate new data points that look like (belong to the same distribution as) samples from that original dataset. These algorithms learn regularities and variabilities in the data, and apply dimensionality reduction techniques to reduce (simplify) the representation space of data items. Later in the course, we will discuss examples of generative machine learning techniques and explain - in simple terms - how a technology like ChatGPT works.

2.3.3.4 Semi-Supervised Learning

- Combination of supervised and unsupervised learning
- Few labeled data in the input are used to create noisy labeled data
- With more labeled data, the machine learns how to make input-output predictions

Semi-supervised learning is a machine learning technique that combines supervised and unsupervised learning to deal with situations where acquiring abundant labeled data is too expensive.



Semi-supervised learning process

Semi-supervised learning allows training a supervised model only having a limited amount of labeled data. The idea is conceptually simple: using a few labeled items, a supervised ML model is trained. This model will probably perform "sub-optimally", but it will probably be able to capture some statistical properties of the data. This lower-quality model is then used to analyze (e.g. classify) all the other data, creating a so-called *noisy label*. All the data (the manually labeled one and the noisy labeled one) is then fed to another supervised model. By expanding the set of labeled data, it should be possible to learn additional properties (patterns) of the dataset, thus improving overall performance.

Note that, generally speaking, a semi-supervised learning model has a lower performance than a comparable supervised learning one as, obviously, is more noise in the training data. Another disadvantaged is that biases potentially present in the original labeled set are also more likely to be propagated. However, semi-supervised learning techniques are cheaper and allow quicker prototyping and experimentation.

2.3.3.5 Transfer Learning

- Often called fine-tuning
 - Reuse a model trained for one task is **re-purposed** (tuned) on a different but related task
- Useful in tasks lacking abundant data



Transfer Learning

Transfer learning is a technique that allows one to take a model trained for a specific task and repurpose it for another one.

For instance, imagine somebody created a model to recognize *cats*, but the task is to create a model to classify *dogs*. As, arguably, *cats* and *dogs* share several visual features (e.g. they all have ears, eyes, and noses) the original model could be **fine-tuned** to recognised dogs. The process of fine-tuning implies re-using some of the knowledge previously acquired in an existing model and its adaptation in the new model through a process of partial retraining. Models created through transfer learning may be



Prediction

Dog

less accurate than models created through supervised learning. Transfer learning is commonly used these days, especially in computer vision and natural language processing.

2.3.3.6 Reinforcement Learning

- Data about the **environment** and **reward function** as input
- The machine can perform **actions** influencing the environment
- The machine learns behaviours that result in greater reward

In reinforcement learning, the machine learning system is given no training data, but only information about 1) the **environment** when the system operates, 2) a set of **actions** that the system can perform in the real world, and 3) a **reward function** that expresses how the environment has *rewards* or *punishes* the system to guide it to reach its goal (or a set of goals).





The training of such a system happens by allowing the system to "navigate" the environment while finding a balance between *exploration* (i.e. test actions and their eventual reward/punishment) and *exploitation* (i.e. use existing knowledge to achieve a given goal).

Reinforcement learning finds several interesting applications. For instance

- Games: for instance, the AlphaGO system is based on RL technology
- Robotics: reinforcement learning is used extensively in robotics to allow robots to learn from their interaction with the environment

- Self-driving cars: for instance, for path planning or to learn how to behave in particular environments.
- Chatbot interaction, for instance ChatGPT,

2.3.4 The importance of domain expertise

- ML makes some tasks automatic, but we still need our brains
- More in Section 3 and Section 4
- Defining the prediction task
- Define the evaluation metrics
- Designing features
- Designing inclusions and exclusion criteria for the data
- Annotating (hand-labelling) training (and testing) data
- Select the right model
- Error analysis

This last slide is simply a reminder about the importance of design activities for ML systems. While the technical and engineering processes of creating a well-functioning ML system are very important, essential work takes place **before** a model is even trained.

2.4 2.4 Machine Learning Models

2.4.1 Linear Regression



Let us start with the *regression* problem and take the simplest, although still very used *linear regression* model.

The idea is simple. We have data points in a multi-dimensional *feature space* and their associated label value. The problem at hand is to determine a mathematical model (a function) that best captures the relationship between these data points in that feature space. The goal is to be able to calculate the *label value* of new data items accurately.

In the simple example here, the problem is to estimate the cost of a house (the label) based on its size (the feature). So it is a simple 1-dimensional problem. And we can describe the data points as in the figure. For each data point in blue, we know the house size and its corresponding value.



How can we estimate the value of a new house, represented here with the red data point? Visual estimation is, of course, not possible in the real world, where the dataset size or the feature dimensionality is much bigger.



The idea is to find a mathematical function, in the case of linear regression, a linear function (a straight line) that best approximates (i.e. *best models*, *best fits*) the data in our possession. The linear function is our model.

The linear function is mathematically expressed as in the slides. The function has only one variable (*Size*), and the value to calculate is *Cost. b and x are the function's parameters:* x^* is the intercept (a constant value, or **bias**), and *b* is the slope of the linear functions. In machine learning lingo, these parameters are often called the **model's weights**.

Not that the bias is not attached to any feature, but its function is to tell us what the output would be in a data point where all the features are precisely zero. In our 1-dimensional feature space, bias allows you to "move the line up and down" on the y-axis, to fit the prediction with the data better. Without bias *b*, the line would always go through the origin point (0,0), thus yielding a poorer fit.

We aim to find the value of *x* and *b* for which the function is a better fit for our data. How could we do this?



To explain the method, we simplify the problem and focus on the case where we want to find the best line that fits these three points.



Assuming we position the line as in the slide, the question is: which *slope* (value of the parameter *b*), allows the linear function to fit the training data **best**? But how do we measure what **best** is?



A common way to measure the quality of a function in fitting a given distribution is by measuring the error made (on average) when using the data distribution. In the example, we do it by calculating the distance (the error) between each point in the training set and the linear function.



Many different lines could fit the simple distribution in the slide. How do we find the best one? We calculate the error for every possible slope of the linear function. A person would take a very long time, but this is what computers are for: execute repetitive calculations quickly and tirelessly. To guide the process, we need a **cost function** (sometimes called a **loss function**), a formula that allows calculating "how good" a model works. **Training a model** means finding the parameter values that minimise the cost function

In linear regression, we often use the so-called **Mean Squared Error**, the average of the squared differences between the predicted value and the actual value of a data item. This is a quadratic function that looks like a parabola curve.

Alessandro Bozzon



In our example, we find that the first slope is the one that better fits the training data.

What we have been doing here is, in a very simple way, **training** a model based on the available data. The training process has the goal of finding the optimal parameters of the model. In our example, the model has a single parameter *b*.

2.4.2 Training a Machine Learning Model

- Training the model
 - Gradient Descent: an algorithm to find the minimum point of a function
 - Hyper parameters: parameters of the Gradient Descent
 - * Learning Rate: speed of descent
 - * *Epochs*: max number of steps

To find the best parameter values, we use a very common approach called **Gradient Descent**. Gradient descent is an algorithmic way to find a minimum point in a differentiable function. The gradient descent approach does not require knowledge of the shape of that function but only of its partial derivatives.

The gradient descent algorithm has two parameters. The **learning rate** controls how much the value of the parameters to be estimated changes at every step. The higher the value, the faster will be the "descent." But also, the higher could be the risk of "jumping" in less optimal points of the curve. The lower the value, the more accurate yet slow the descent. The **epochs** that is the maximum number of steps that the algorithm can take before stopping. We need to set a maximum number of epochs because it is possible for the algorithm never to find **THE** optimal value and thus continuously oscillate between close values.



In this example animation, we visually show how the gradient descent algorithm operates.

The graph on the left shows a dataset, similar to the one in our example. Each blue dot is a data item. The red line represents the linear function for a specific *slope* value. You can observe the line moving upwards as the *slope* value increases.

The graph on the right has on the *x* axis the value of the parameter *p* to estimate (in our example, the slope *b*), and on the *y* axis the value of the cost function calculated for a specific value of *b*. We saw in a previous slide that the *mean squared error* is a quadratic function that looks like a parabola curve. In the animation, we explore different values of the parameter left to right, from zero to higher numbers. As we move from left to right, the value of the error function decreases, progressively approaching its minimal value - that is, the value of *p* for which the error is minimal. That value, *18* in our animation, is the desired one for our model.

From this animation, you can appreciate how the **learning rate** and **epochs** parameter allows controlling for the length and precision of this process. For instance, a high learning rate – that is, increasing too much the *p* value, allows to explore the parameter values quicker, but it could have led to missing the optimal value.

In the real world, where models have thousands, millions, or billions of parameters, there are several optimisations for the gradient descent algorithm. For instance, the initial values of parameters are either set at random or sampled from an existing data distribution. Likewise, values of learning rate and epochs are determined dynamically, based on the value of the partial derivates of the error function.



In the previous example, we explored a very simple regressions model; the data points were just a few, so, the optimization could be done quickly.

Also, by selecting a linear function, we assumed that the data hold, to some extent, a linear relationship. In many cases, such linear relationships might not hold so we can use more complex models. That is, models with more parameters.

In the example here, we use a polynomial function. The function still has a single feature (*size*) but several components. This allows for a non-linear model that is "closer" to the data points, thus reducing the error. This is called polynomial regression.

As we will see later, there are even more complicated models.

In a model like the one in the slide, what does it mean for a parameter to equal 0? That the related component **does not contribute** to the prediction – that is, the value of that particular component does not relate to the output value. Similarly, if a component has a *very high parameter value* (negative or positive), then such a component is very important in predicting the output value.



2.4.3 Classification

Let us now look at the **classification** problem.

For example, let us take the problem of classifying whether a student will be accepted at a university based on the results of a selection *test* (*x*-axis) and their high-school grades in *mathematics* (*y*-axis).

For this problem, we are provided with several data points. We have several students for whom we know the values of these two variables (our *feature space*, in this case 2-dimensional) and their admission result (yes/no, our binary *class*). For instance, in the slide, we have two students that were either rejected or accepted based on the two values.



The dataset is, of course, bigger. Green dots are students that were previously accepted. The red dots are students that were not accepted. As you can see, there isn't a clear-cut separation in this 2-dimensional feature space. Obviously, high-performing students (high *math* and *test* scores) were accepted, and low-performing students were rejected. However, there are several students "in the middle" for which the relationship between the value of the two features and the outcome (accepted/rejected) is unclear. Some students with a high *test* grades were not accepted due to low *math* grades. And vice versa.



Let us assume that new students are coming. Given a new student applying at a university (math grade 6) and performing the admission test with a 7. Will the student be accepted?



To answer the question, let's look at the data. Most of the points seem to be separated by a line. Most points over it are green, and most below it are red, with some exceptions. We will say that that line is our model, so when a new student comes in, if they happen to be above the line, we predict they will be admitted. Otherwise, they won't.

We have already encountered this method, this is a Logistic Regression.

How do we find this line that cuts the data in 2?



As before, let me show you an example of how a computer can learn to do this automatically.

Let's assume we draw a line trying to separate this bi-dimensional space in 2.

As with the previous linear regression, we want to measure how "good" this model is in fitting the data. In this case, we measure the error slightly differently. This means we use a different **cost function**.

For instance, we can count how many points are classified incorrectly. In the example here, the number of errors is 2.



By trying different configurations for this function (that is, different slopes), we can find the one that minimises the number of prediction errors. This can also be done with the gradient descent algorithm.



In the slide, by moving the line a little bit, the number of errors decreases.



If we move it again, it becomes 0.

Alessandro Bozzon

Note that in reality, we would not use the number of errors as a cost function, but we use something called **Log-loss** function.

Here the mathematical formulation becomes a big complex, so we will not indulge further. The intuition beyond this measure is to give a high penalty to items that are misclassified with a high probability (that is, they fall far away from the line), and to give a lower penalty to the ones that are misclassified with low probability (so, they are closer to the line).



We created the example here to illustrate the problem of finding a **decision boundary**. However, the example is weird, as we data items where students with a very low test score or math scores are still accepted.



Let us then modify the data a little to make it a bit more realistic.



Notice that with this data, it would be very difficult to define a decision boundary with a single line. We would need something like a circle, a function that is much more difficult to estimate.



With can also try to train a model composed of multiple functions. For instance, two lines. As before, we can use some gradient descent approach to find the best possible parameters for each of these lines.



In this example, we are estimating multiple functions at the same time. Each function is estimated separately, yet together, as they contribute to defining the decision boundary. And in practice, this is

what a **neural network** does.

2.4.4 Neural Networks



Neural networks are one of the most popular (if not the most popular) machine learning models currently used for both *classification* and *regression*.

Neural networks are meant to, in the very broad sense of the word, mimic how the human brain operates.

The one in the picture is a **deep** neural network, of the ones that you might have heard of or that we have discussed previously. They indeed look scary, but we can simply understand them.



We call them neural networks because their basic unit, the **perceptron**, vaguely resembles a **neuron**.

A neuron comprises three main parts: the soma, the dendrites, and the axon. In broad terms, the neuron receives signals from other neurons through the dendrites, processes them in the soma, and sends a signal through the axon to be received by other neurons.

A *perceptron* is designed to work similarly to a neutron: it receives *numbers* as inputs, applies a *mathematical operation* to them and outputs a new *number*.



In a *perceptron*, the mathematical operation is typically the weighted sum of the input values plus the bias - a simple weighted linear combination of the input values.

In this way, we can represent linear functions in feature space of arbitrary size, including the monodimensional one of our previous examples with linear regression. Intuitively, a perceptron with a single input (and a bias) is precisely the equivalent of the regression function in our previous example.



The **activation function** takes the values of the weighted linear combination of the input and decides how the perceptron will *activate* (or *fire*), that is, if it will output a value.


Many activation functions could be used in a perceptron. We show two.

The *step* function returns a 1 if the output of the weighted linear combination is nonnegative and a 0 if it is negative.

The *sigmoid* function works slightly differently, as it outputs a number between 0 and 1. The number is close to 1 if the score is positive and close to zero if the score is negative. If the score is zero, then the output is 0.5. Note that a single perceptron using a sigmoid function operates exactly as a linear regression.



So, let's go back to our classification example, where we have two linear functions dividing the twodimensional space.

We can split the problem of classifying the data items into two problems - each related to one of the two functions.



Let's say that one function – we call it **Test Grade Classifier** – partitions the decision space in one way, while the other function – we call it **Math Grade Classifier** – partitions the decision space in the other. Note that we are still using **both** features in each classification. Together, the two classifiers make the **Admission Classifier**.



The situation described in the previous slide can be modelled with the neural network in the picture. It comprises three perceptrons: the **Test Grade Perceptron** and the **Math Grade Perceptron**, both used as input for the **Math* and** Test** features. The **Admission Perceptron** uses the output of the **Test Grade Perceptron** and the **Math Grade Perceptron** as input. The three perceptrons output a 0 or a \$1# value, as they use a **step** activation function.



Note how different activation functions produce different types of classification boundaries.



Here you can see how the use of a sigmoid function can change the boundaries of the classification

space.



2.4.4.1 Fully connected Neural Network

Hyperparameters

- Learning rate
- Number of epochs
- Architecture
 - * #layers, #nodes, activation functions
- Batch vs. mini-batch vs. stochastic gradient descent
- Regularization parameters:
 - * Dropout probability p

So far, we have seen examples of a small neural network. In real life, neural networks are much larger. The nodes are arranged in layers, as illustrated in the slide.

The first layer is the **input layer**, the final layer is the **output layer**, and all the layers in between are called the **hidden layers**. Perceptrons in the same layers are typically similar, i.e., they have the same activation function.

The arrangement of nodes and layers is called the **architecture** of the neural network. The number of layers (excluding the input layer) is called the **depth** of the neural network. The **size** of a layer is the number of non-bias nodes in that layer. Note that neural networks are often drawn *without* the bias nodes, but it is always assumed they are part of the architecture.

The network in the slide has every node in a layer connected to every (non-bias) node in the next layer. Furthermore, no connections happen between nonconsecutive layers. This architecture is called

fully connected feedforward. For some applications, we use different architectures where not all the connections are there or where some nodes are connected between non-consecutive layers. We will see some examples in the coming lectures.

Like most machine learning algorithms, the *training* of a neural network is configurable through *hyper-parameters*. These hyperparameters determine how the training is performed. The first two we saw before: how long we want the process to go (number of epochs) and at what speed (learning rate). As neural networks are much more powerful models, their training (and performance) are subject to much more fine-tuning. For instance, neural networks are not trained by examine **all the data* in the dataset at once but just a subset at a time. How these subjects are composed (e.g. size, distribution of data items in the feature space) is a hyper-parameter (Batch vs. mini-batch vs. stochastic gradient descent). Neural networks are also prone to *overfitting* (a concept we will explore in another lecture) – that is, they have poor generalisation capabilities and learn "too well" the training data distribution. *Regularisation* techniques help by "distracting" the network by introducing noise. For instance, *Dropout* is a training technique where a few perceptrons in the network are *randomly* "switched off" during different training phases.



2.4.4.2 Classifying into multiple classes - Softmax function

- Return a probability for each class
 - example C1= ADMITTED, C2 = NOT ADMITTED, C3 = NEW TEST
 - p(C1) = 0.37, p(C2) = 0,21, p(C3) = 0,42
- We use the Softmax activation function for the output layer

Perceptrons return a single numerical value (discrete or continuous). They are, therefore, perfectly suited for *regression* problems, but how can we make them work for classification problems? In that case, instead of returning a 0 or a 1 (step function on the output layer), or a value between 0 and 1 (sigmoid function on the output layer), we would like to return a **probability** value for each class in the classification problem. In the slide, for instance, we have slightly modified the classification problem by having 3 classes: C_1 that means *admitted*, C_2 that means *not admitted* and C_3 that represents

a case when a *new test* is needed. When a new instance is processed, the network should return a probability value (0 and 1) for each class.

To allow for this situation, we can change the activation function of the output layer into a **Softmax** function. The softmax function is a function that turns a vector of K real values (where K is the number of classes) into a vector of K real values that sum to 1. The input values can be positive, negative, zero, or greater than one, but the softmax transforms them into values between 0 and 1, so they can be interpreted as *probabilities*. If one of the inputs is small or negative, the softmax turns it into a small probability, and if an input is large, it turns it into a large probability, but it will always remain between 0 and 1.

The important assumption is that the true class labels are independent. That is to say each sample of data can only belong to one class. For example, a person cannot be *admitted* and not *admitted* simultaneously. Its true label can only belong to one class.

Note that the *sigmoid* activation function is equivalent to the *softmax* function when we there are only two classes. It is unnecessary to calculate the second vector component explicitly because when there are two probabilities, they must sum to 1. So, if we are developing a two-class classifier with logistic regression, we can use the sigmoid function.



If you want to experiment with neural networks, and have a better feeling of how they work, I strongly recommend you try the Tensorflow Playground. There you can try different data distributions, classification problems, activation functions, etc.

2.5 2.5 Models Development Lifecycle

You were earlier introduced to the CRISP-DM methodology: Cross Industry Standard Process for Data Mining, a way to describe the process of Machine Learning system development. Three phases in that process are of interest to us in this section:

- In the **modeling step**, the actual machine learning model is trained. This step begins once the problem to be solved is defined and the data is prepared.
- In the **evaluations** step, the development team checks how the developed model complies with the expected performance requirements.
- In the **deployment and monitoring** step, the machine learning model is put "in production," facing real users. There, the model encounters "the real world", dealing with data items never experienced before during training and evaluation.

These three steps are independent but entirely connected. Therefore, we will look at them again, describing them through a slightly different lens.



The diagram in the slide *unpacks* three **CRISP-DM** steps described before and highlights the **data** used in each stage and the **actions** performed.

The data available beforehand are called *historical data*. This emphasizes how such data represent and encode the past phenomena that the ML algorithm will try to predict.

In the **modeling phase**, developers use two disjoint (i.e., no data items in common) subsets of the historical data: the *training set* and the *validation set*. As the name suggests, the *training set* is used to train the machine-learning model - that is, given the chosen class of models (e.g. a neural network), learn the model's parameters that can achieve better prediction performance. Once the model is trained, the *validation set* is used for validation purposes. This means, for instance, "tuning" the model based on the performance of that different set of data. The tuning is needed for several reasons.

• The first is to look at the model's performance against a different data set. As we will see in a later lecture, machine learning models can become too good at learning a given data distribution (this is called *overfitting*) at the cost of generalisation performance (showing good performance with unseen data items). If there is a substantial difference in the model performance over the

training and validation sets, then probably some early choices need to be reconsidered. For instance, a deep neural network might be too complex a model for the problem.

• The second reason concerns selecting the best models and learning *hyperparameters*. As you remember from the previous lecture, models can be configured – e.g., the number of layers in a deep learning architecture, and the complexity of a linear regression model. Likewise, learning algorithms like gradient descent have hyperparameters like the learning rate or the number of training epochs. Tuning the model entails determining the optimum parameters for the model, called hyperparameters, that give the best performance. The validation step is an iterative process, that continues until the developer is satisfied with the achieved level of performance.

In the **evaluation phase**, the learned model is tested against another subset of *historical data* called the *test set*. The *test set* is also disjoint from the *training(and validation* sets.*

The objective of the **evaluation phase** is to obtain an unbiased assessment of the model's performance. The lack of bias, in this case, comes from the fact that the *test set* is entirely unknown to the model, thus mimicking the situation that could occur in the real world. In reality, the data in the *test set* are still *historical data*, so they have probably been created in the same conditions as the training and validation data. Therefore, the evaluation is not completely unbiased. The evaluation phase can (and *should*) also include stakeholders and experts so that they can also evaluate the model's performance before being put into production. This additional testing is called **user acceptance testing** (UAT), and is supposed to be the final stage in the development of any software system. Unsatisfactory performance requires returning to the modeling phase or earlier to the **data preparation** or the **problem specification** stages.

If the performance is satisfactory, the process moves to the **deployment and monitoring** phase. The model is run against *new data items* produced by the system's end users. If these data items are sufficiently similar to the *historical* ones, then the model should perform well. It is possible, however, that these new data items are different enough to yield unsatisfactory performance. This is almost inevitable due to a phenomenon called *concept drift*.

Concept drift occurs when the properties or the distribution of the data fed to the deployed model *changes* compared to the *historical data*. This could happen when the historical data did not contain examples of a specific class. The (in)famous example of the Tesla Autopilot confusing the moon for a traffic light is representative of this class of issues. Concept drift can also happen when classes exist in the historical data, but the properties of items belonging to such classes change over time, thus leading to increasingly worse performance. For example, consider a machine learning model trained to classify news articles into categories like "politics," "sport," and "technology." Suppose the model is trained on news articles from the 20th century. In that case, its performance on current data will probably be bad: terms like "smartphone" were perhaps not present back then, thus predicting the "technology" category very hard.

To mitigate these issues, it is essential to continually update the training dataset and retrain the models

accordingly. This is why any decently designed ML system should include a **monitoring** step, where feedback from users or predictions are collected and analysed to provide new training data for the system and then allow its improvement over time.

2.5.1 Dataset Splitting

2.5.1.1 Split your data



– evaluate

A quick discussion on the issue of dataset splitting. As explained in the previous section, training and evaluating a machine learning model based on historical data requires splitting such data into a *training*, *validation*, and *test* set. But how should this data be split?

From a numerical perspective, the idea is always to reserve more data for the training than for the evaluation phase. So, an 80-20 strategy (80% for training and validation, 20% for testing) is very common. The proportion can vary slightly, also based on the abundance of the available data. Remember that machine learning models need a lot of data for good training, so the more data available for training the better.

But quantity is not the only issue. How data items are distributed is also important. Imagine a computer vision dataset for object classification containing 10 categories of items. Intuitively, the 10 classes should be present in the *training*, *validation*, and *test sets*, otherwise performance will be worse, for instance, for those classes that received no training. However, this is not sufficient, as it is essential to make sure that the number of items per class (their distribution) is proportionally similar. In the example above, the training set should contain a similar percentage of each category of items as the

test and validation sets. This will maximise the chances that the developed model reflects and is being evaluated in real-world scenarios - at least insofar as the historical data contains a good reflection of the world to be modeled.

2.5.1.2 Avoid leakages



- in the validation or evaluation sets

- Features
 - highly correlated to prediction
 - not present in the production environment

Keeping a proportion of the data is insufficient, as there could be some subtle issues in how the data is distributed across the three sets. A specific type of issue is the so-called **data leakage**, which is the presence of data items in the **training** set leak into the **validation** or **evaluation** sets. If this happens, then it would be like training the model to "cheat the test," as training data are more likely to be recognized by definition.

Data leakage can also happen when features in the training, validation, and test sets unintentionally leak information that would otherwise not appear in deployment when the model performs predictions on new data. A famous example of data leakage happened during the KDD Cup Challenge of 2008. The goal was to create a machine-learning model based on X-ray images to detect whether a breast cancer cell was benign or malignant. The model with higher performance used a feature called *Patient ID*, the identifier generated by the hospital for the patient. It turned out that some of the hospitals that provided the data used the patient ID as a way to indicate the severity of the patient's condition *when they were admitted to the hospital*, thus creating an association between the feature and the expected result of the x-ray analysis of the patient when they were admitted to the hospital, which, therefore, leaked information about the target variable.

2.5.1.3 Cross-validation



- Cycle training and validation data several times
 - Useful when dataset is small
- Split the data into *n* portions
 - Train the model n times using n-1 portions for training
 - Average results

How could data be spliced between the *training set* and the *validation set*? Simply partitioning the data (e.g., 60% for training and 20% for validation) is an option. However, remember that at training time, the goal is to achieve good performance and generalize as much as possible. So, the more data items are explored for training and validation purposes the merrier.

Cross-validation is an approach where the historical data allocated for training and validation are randomly partitioned into N sets of equal size, and the learning algorithm is also run N times. Each time, one of the N sets is used as a validation set, and the model is trained on the remaining N-1 sets. The score (error) of the model is evaluated by averaging the error across the N validation errors. In the slides, we can see an example where the available historical data is split in 4 batches. In this way, it is possible to create 4 training set up, each with a different validation set.

The advantage of cross-validation techniques is their conceptual simplicity, the disadvantage is the high computational cost resulting from many repeated training trials. If the computational cost of training a model is high, cross-validation could be expensive. Moreover, as the training dataset is reduced, it is important to make sure that there is sufficient training data so that all relevant phenomena of the problem exist in both the training data and the validation data.

Note that a similar approach can also be used for the creation of the *test* dataset.

2.6 2.6 Evaluating Machine Learning Models

2.6.1 How to Evaluate?

- Metric
 - How to measure errors?
 - Both training and testing
- Training
 - How to "help" the ML model to perform well?
- Validation
 - How to pick the best ML model?
- Evaluation
 - How to "help" the ML model to generalize?

How can the evaluation of Machine Learning systems be performed? And, most importantly, what is evaluated?

We have seen in the previous slides the roles that the *training* and *evaluation* steps play in the engineering of an ML model. Evaluation during *training* is primarily devoted to creating a model that performs best on the training data; or selecting the best model or configuration. During the *evaluation* step, the additional goal is to assess if the model performs acceptably in the chosen application domain, possibly by including end users or other stakeholders.

In both cases, what is needed is a way to measure the performance of a model mathematically.

2.6.2 Let errors guide you

- Errors are almost inevitable!
- How to measure errors?
- Select an evaluation procedure (a "metric")

Machine learning systems are stochastic. This means that, in practice, they will **always** make some mistakes. The goal is, of course, to minimize such mistakes, and this is what we try to get a machine learning model to do: **minimize mistakes**. Or minimise the errors they make.

In Section XXX we described in simple terms what, in practice, the training of a machine learning model is: find the best model's parameter values that **minimize** a given prediction error (or prediction cost).

2.6.3 How to calculate

- These are the most common questions:
 - How is the prediction wrong?
 - *How often* is the prediction wrong?

- What is the *cost* of wrong predictions?
- How does the *cost* vary by the wrong prediction type?
- How can *costs* be minimised?

How can a good metric (or error functions) be selected? The answer to this question requires a deep and precise analysis of the *cost* associated with an error. Here by *cost* we do not only mean monetary cost, although money can always be an issue. *Cost* can also be interpreted as the **harm** that could be done to the people on the receiving side of a prediction. It is possible, for instance, that the wrong classification of an offensive document brings discomfort to the user of an online video service. The wrong classification of an chest scan can have far worse consequences.

2.6.3.1 Regression

How do we measure errors in regression problems?

2.6.3.1.1 Mean absolute error $MAE = \frac{1}{N} \sum_{J=1}^{N} |p_j - v_j|$

Average of the difference between the *expected* value (v_i) and the *predicted* value (p_i)



In a previous lecture, we have already been introduced to the concept of *error* in regression: the difference between the *expected* value v_j and the *predicted* value p_j .

Mean Absolute Error averages the error made by the model on every data point in the training set (if the Mean Absolute Error guides the training), in the validation set (if it is used during validation), or in the evaluation set. It is called the **absolute** error because it is calculated over the *absolute* value of the difference: the difference can be positive (the model predicts a higher value) or negative (the model predicts a lower value). To turn this difference into a number that is always positive, we take its absolute value.

2.6.3.1.2 Mean square error $MSE = \frac{1}{2N} \sum_{J=1}^{N} (p_j - v_j)^2$

- Average of the square of the difference between the *training* value (v_j) and the *expected* value (p_j)
- Square is easier to use during the training process (derivative)
- More significant errors are more pronounced



The **Mean Square Error** is very similar to the absolute error, but what is calculated is the **square** of the difference between the predicted and expected values. The result will always be a positive number. The square error is used more commonly in practice than the absolute error for two reasons: 1) the square has a much nicer derivative than an absolute value, which is useful when using the gradient descent algorithm for training; 2) it "amplifies" the value of large errors, making them more prominent in the error function.

2.6.3.2 Classification

And what about classification?



2.6.3.2.1 Confusion Matrix Describes the complete performance of the model

In classification models, the classification can either be **correct** or **not correct** - it is impossible to be "sort of" correct. Therefore, errors are described and calculated differently from *regression*.

Let us use the example of a *binary* classifier that works only on two classes. For instance, *yes* and *no* (as in the slide), *sick* or *not sick*, *cat* or *dog*. In this case, errors can only be of two types:

- **False Positive** (FP): a data point with a negative label that the model classifies as positive. For instance, a healthy person (*not sick*) who is incorrectly diagnosed as *sick*.
- **False Negative** (FN): a data point with a positive label that the model falsely classified as negative. For instance, a *sick* person is incorrectly diagnosed as *not sick*.

Of course, we also have:

• **True Positive** (TP): data points with a positive label correctly classified as positive. For instance, a healthy person classified as *not sick*.

• **True Negative** (TN): data points with a negative label correctly classified as negative. For instance, a sick person correctly diagnosed as *sick*.

A **Confusion Matrix** helps keep track of these four quantities. In binary classification models like the one in our example, the confusion matrix has two rows (the predicted classes) and two columns (the actual class). Cells contain the number of time the model produces **True Positives**, **True Negatives**, **False Positives**, and **False Negatives**. In confusion, the correctly classified data items are counted on the diagonal, and the incorrectly classified data items are counted off the diagonal. Note that sometimes rows and columns can also be transposed, but the meaning of TP, TN, FP, FN remains the same.

In classification settings with more classes, the confusion matrix is larger. Consider, for example, a model classifying patients as *sick*, *not sick*, and *unsure*; then the confusion matrix is a three-by-three matrix.

2.6.3.2.2 Type I and Type II errors In literature, it is also possible to find False Positives and False Negatives respectively described as Type I and Type II errors.



2.6.3.2.3 Accuracy

- The percentage of times that a model is correct
- The model with the highest accuracy is not necessarily the best
- Some errors (e.g., False Negative) can be more costly than others

$\frac{TP{+}TN}{TP{+}TN{+}FP{+}FN}$

Accuracy is the simplest and the most common measure of classification models, and it calculates the percentage of times a model is **correct**. Its value is calculated as the ratio between the number of correctly predicted data points and the total number of data points.

Note that in this formulation, **Accuracy** treats each data points in the dataset equally: there is no data point for which a mistake is worse than others. Also, accuracy does not consider if items belonging to a specific class are more popular than items belonging to another class. In this situation, a model can achieve a very high accuracy score because it performs well on the most popular class items. Take, for instance, a problem where 99% of the data items are in one class (e.g. *not sick*) and only 1% are in the other class (e.g., *sick*): a model that **always** predicts *not sick* will achieve a 99% accuracy, thus

2023-03-29

appearing to be very good, while in reality being very problematic for the 1% of the population that needed the model to work.

2.6.3.2.4 Errors are not equal Also, accuracy treats **errors** as *equally costly*. In the example above, making a mistake on the 1% of the population (that is, classifying them as *not sick* even if they are *sick*) is a problem.

FALSE POSITIVES: SELF-DRIVING CARS AND THE AGONY OF KNOWING WHAT MATTERS



According to a preliminary report released by the National Transportation Safety Board last week, Uber's system detected pedestrian Elaine Herzberg six seconds before striking and killing her. It identified her as an unknown object, then a vehicle, then finally a bicycle. (She was pushing a bike, so close enough.) About a second before the crash, the system determined it needed to slam on the brakes. But Uber hadn't set up its system to act on that decision, the NTSB explained in the report. The engineers prevented their car from making that call on its own "to reduce the potential for erratic vehicle behavior." (The company relied on the car's human operator to avoid crashes, which is a whole separate problem.)

ARN MORE

Uber's engineers decided not to let the car auto-brake because they were worried the system would overreact to things that were unimportant or not there at all. They were, in other words, very worried about false positives.

The images in this slide provide real-world examples of situations where classification errors made for rare or unknown classes can have very dire consequences.

- Detecting the "Alexa" command?
- Pregnancy detection
 - Cost of "false negatives"?
 - Cost of "false positives"?
- Covid testing
 - Cost of "false negatives"?
 - Cost of "false positives"?
- Law enforcement?

Depending on your task, different errors have different costs. Consider the examples in the slide. What are the costs associated with mistaking an *Alexa* command? The cost of a *false negative* is probably meagre for the end user. In the worst case, the user must repeat the command or stop using *Alexa* altogether. A *false positive* can be more tricky, as *Alexa* can start processing the input, thus believing that an instruction must be executed. In a funny example, it is reported that during a broadcast, the TV anchor Jim Patton said, "I love the little girl saying, 'Alexa, order me a dollhouse.'" In households with an Alexa device listening, Alexa mistook the statement as a request to order a dollhouse. Here the cost of the mistake is not high, but several users had to rush to cancel orders placed on their devices, and Amazon had to explain how something like that was possible in the first place.

What about the cost of errors in the case of a pregnancy test? Here a *false positive* (or a *false negative) will probably give a big scare or a big hope (depending on whether the people involved wanted or not to conceive a child), but a second test or a proper visit to the doctor will probably remove any doubt.

In the case of a COVID test, things are a bit trickier. A *false negative* could give *sick* people the freedom to circulate and further spread the virus. In March 2023 this was not such a big issue anymore, but in March 2020 it was. A *false positive*, on the other hand, will probably force someone to stay at home (or work from home) even if they feel perfectly ok.

One critical remark: no matter how good a metric can be, it can always be *fooled*. The decision on what metric to use is **a design choice**, which should be driven by an understanding of the costs of errors (false positive vs. false negative). Never assume that a model is good before evaluating it against different metrics.

2.6.3.2.5 Balanced Accuracy The **Balanced Accuracy** is a way to measure accuracy by accounting for the accuracy value of each class.

$$\frac{\frac{TP}{TP+FN} + \frac{TN}{FP+TN}}{2}$$

- Average of single class performance
- Good to use when the distribution of data items in classes is imbalanced

The formula of the Balanced Accuracy calculates the average of each class's accuracy. Note that the accuracy of a single class is equivalent to the **recall** for that class (we will see the formula for recall in a couple of slides).

When the dataset is *unbalanced*, i.e., one or more classes have a significatively smaller number of data items, the **Accuracy** value mainly depends on the algorithm's performance with the most popular classes. The **Balanced Accuracy** gives smaller classes the same influence on the formula then the more popular ones, although their size is reduced in the number of data points. This means that **Balanced Accuracy** is insensitive to an imbalanced class. This can be an advantage when interested in having good prediction also for under-represented classes. But it can be a disadvantage if the goal is to have good prediction on the entire dataset.

If dataset is quite *balanced*, i.e., the classes have a balanced number of data items, **Accuracy** and **Balanced Accuracy** have similar values.

2.6.3.2.6 Balanced Accuracy Weighted The **Balanced Accuracy Weighted** extends the **Balanced Accuracy** formula by multiplying each class's accuracy by the class frequency (*w*) on the entire dataset. In this way, larger classes have a weight proportional to their size.

 $\frac{\frac{TP}{(TP+FN)*w} + \frac{TN}{(FP+TN)*(1-w)}}{2}$

Alessandro Bozzon

- Weighted average of single-class performance
- Weight depends on the popularity of a class.

Since class accuracy is weighted by the class frequency, **Balanced Accuracy Weighted** can be a very good performance indicator with classification algorithms working on many classes.

2.6.3.2.7 Precision The **Precision** is the fraction of True Positive predictions divided by the total number of positively predicted data items (True Positive plus False Positive).

 $\frac{TP}{TP+FP}$

• Among the examples we classified as positive, how many did we correctly classify?

Precision is a metric that tests the correctness of the Positive data items that are correctly classified. With Precision, we are interested in minimizing the number of **False Positives** the model predicts because such mistakes can be costly. Consider the example of a pregnancy test. A test with too many false positives will probably be undesirable because it would overflow the healthcare system with people that did not need to be controlled during an actual pregnancy – not to mention the undesired emotional effect for people that wanted a pregnancy bit did not have one, or for those that did not want a pregnancy, and are now very upset about it. Not that very high precision can be achieved by a system that rarely (or never) predicts a positive value - thus making, by definition, no mistakes. Yet, that model will probably be useless.

2.6.3.2.8 Recall The **Recall** is the fraction of True Positive elements divided by the number of positively classified units (True Positive plus False Negative).

 $\frac{TP}{TP+FN}$

• Among the positive examples, how many did we correctly classify?

Recall is a metric that tests the proportion of all the Positive data items that were supposed to be positively classified. **Recall** is a useful metric when in the considered applications **False Negatives** are costly. Take the example of a COVID-19 test during the early phases of the pandemic. Arguably, the most important thing was to find all the infected people worldwide, even if that meat misdiagnoses some non-infected people. Note that a model that diagnoses every patient as *sick* can be very bad: despite having zero False Negative patients, it can have too many False Positives.

2.6.3.2.9 F_1 -Score **Precision** and **Recall** measure two performance properties that are both of importance. It can be helpful to be able to combine them both. For example, we can be interested in an ML model that does not misdiagnose *sick* patients but does not misdiagnose too many healthy (*non*

sick) people – e.g., when the wrong diagnosis of a healthy person may involve unnecessary and painful testing.

$$F_1 = 2 * \frac{1}{\frac{1}{P} + \frac{1}{R}}$$

- The harmonic mean between *precision* and *recall*
- What is the implicit assumption about the costs of errors?

The F_1 score combines both recall and precision. If both **precision** and **recall** are high, then the **F_1** score is also high. However, if one is low, then the **F_1** score will also be low, thus indicating that one of the two is not at the desired level. Note that the *harmonic mean* is not the same as the *average*, although it is always smaller or equal to it. The reason for not using a simple average is that a simple average would balance the excellent performance of one metric with the lousy performance of another. For example, a system with 90% precision and a 10% recall will have a simple average of 50%, which might be misleading. The harmonic mean (the F_1 score) is 18%.

2.6.3.2.10 Sensitivity (true positive rate) In the medical field, Sensitivity is used to identify the positively labeled data items, and it is the same as **Recall**.

 $\frac{TP}{FN+TP}$

- Identification of the positively labeled data items
- Same as recall

2.6.3.2.11 Specificity (false positive rate) Specificity is another metric used in the medical field to focus on negatively labeled data points. Note that Specificity is not the same as Precision because the focus is on the negative data items.

 $\frac{TN}{FP+TN}$

- Identification of the negatively labeled data items
- Not the same as precision



This table summarises how Sensitivity (Recall), Specificity, and Precision are calculated.

2.6.3.3 Examples

2.6.3.3.1 Medical Test Model

- Recall and sensitivity
 - How many were correctly diagnosed as sick among the sick people (positives)?
- Precision
 - Among the people diagnosed as sick, how many were sick?
- Specificity
 - Among the healthy people (negatives), how many were correctly diagnosed as healthy?

Let us discuss the three measures in the case of a medical test like in the COVID-19 example. As we mentioned before, in this case, we are more concerned about correctly diagnosing sick people so we would need a model with **high sensitivity** (or high **recall**).

2.6.3.3.2 Spam Detection Model

- Recall and sensitivity
 - How many were correctly deleted among the spam emails (positives)?
- Precision
 - Among the deleted emails, how many were spam?
- Specificity
 - Among the good emails (negatives), how many were correctly sent to the inbox?

Let us consider another example: a model that classifies *spam* emails and deletes them from a mailbox. With this model, the **sensitivity** (**recall**) is the proportion of spam messages (positive class) that were correctly deleted among all the spam messages. The **specificity** is the proportion of good emails (negative class) that were correctly sent to the inbox among all the good emails. In this case, the application requires the model to have **high specificity**, as wrongly classifying mails as **spam** can be a big issue.

2.6.3.3.3 Search Engine

- Constraint: high precision
 - False positives are tolerable but should be minimised
- Among the available models, pick one with a higher recall
 - **Metric**: Recall at Precision = x%

In a final example, let us consider a search engine. In this case, the goal of the application is to retrieve the best possible set of pages at the top of the search results page. For this application **high precision** is important, although some mistakes can be tolerated. If there is the possibility of picking between different models, the suggestion is to pick the one with similar precision and higher recall.

Note that dozens (if not hundreds) of metrics could be used to measure machine-learning models (for instance, the **Recall at Precision** mentioned in the slide. In this course, we will only look at few, but there will be many more that you can encounter in the follow-up courses in the master.

2.6.4 Metrics are designed in a multi-stakeholder context

- One team builds the mode
 - Data scientists / ML engineers
- Many teams will make use of it
 - e.g., product team, management team

Chasing a metric is, of course, a critical and delicate design task. A task that often requires a design process involving several stakeholders.

Though a single team is often responsible for building a machine-learning model (the data science / or ML team), many teams across an organization will have a stake in the model's performance. Each team will inevitably have its definition of *success** or *failure*, which could be at odds with one another.

For instance, imagine a model that should identify defective products from images. The data science team wants to minimize standard accuracy metrics because it simplifies training. On the other hand, the product management team is more interested in reducing the number of defective products that are misclassified and sent to customers. Finally, the executive team might aim at minimizing the products that are misclassified as defective because that will lead at a revenue decrease. Each one of these goals has a different conceptualisation of *success*, and gives a different weight to *errors*. Balancing these differing needs within an organization can present a challenge. When defining the goals of a model, it is therefore important, as for many other design processes, to consider the needs of different teams across an organization and how each team's needs relate to the model.

3 Section 3: Image Processing Methods

In the first lecture, we defined Computer Vision as the sub-field of Artificial Intelligence and Machine Learning that focuses on 1) extracting high-level understanding from images or videos or 2) generating realistic images and videos.



To understand why vision is a difficult task for computers to perform, let's start with something relatively obvious: the way we (humans) see is very different from the way computers "see." When presented with an image like the one in the slide, humans can process much information simultaneously. You can look at the picture as a whole, and understand that it depicts an outdoor scene in a park, probably New York City (this is Central Park), where people are Ice Skating, probably in winter. We do all of this instantaneously, almost intuitively, relying both on our intuition and our knowledge.

167	153	174	168	150	152	151	172	141	i tere		107.14	1.00	174	100	140	160		1.92		100		107.	147	163		100	140	163		1.01		100	100	187.		14.5	174	144	150	152	151	172	161	155	156	187
165	182	163	74	75	62	33	17			h	ic	н.				h	-	- 1	F.	2	۰.	~	~		2			in.	۰.			"	-	-	~			"	75	62	33	17	110	210	180	154
180	180	50	14	34	6	10	33				15		Э	1	N		C	1		a			U	4		I.	Л	u,	u				3			-	5		34	8	10	33	-48	106	150	181
205	109	6	124	13	111	120	204	166	8 16	56	180,206	109	- 6	124	131	111	120	204	168	16	56	1803	206	109	6	124	131	111	120	204	168	16	66	180 :	206	109	5	124	131	111	120	204	168	18	66	180
172	105	207	233	233	214	220	230	223	90	74	205 173	105	207	233	233	214	220	230	220	90	74	205	172	105	207	233	233	214	220	230	220	90	74	205	172	105 :	207	233	233	214	220	230	220	90	74	206
188	88	170	209	188	3 216	211	150	130	75	20	168 188	88	179	209	188	216	211	150	139	75	20	168	188	88	179	209	188	216	211	150	139	75	20	168	188	88	179 -	209	198	216	211	160	139	75	20	168
189	97	165	84	10	168	134	11	31	62	22	148 188	97	165	84	10	163	134	11	31	62	22	148	189	97	165	84	10	168	134	11	31	62	22	148	189	97	165	84	10	168	134	11	31	62	22	148
199	168	191	193	158	3 227	178	143	182	108	38	198 198	168	191	193	158	227	178	143	182	106	38	198	199	168	191	193	158	227	178	143	182	108	35	196	190	168	191	193	158	227	178	143	182	106	36	196
205	174	155	252	236	3 187	86	150	79	38	218	241 200	174	155	252	236	187	85	150	79	38	218	241;	206	174	155	252	236	187	85	150	79	38	218	241 :	205	174	155	252 :	236	187	85	150	79	38	218	241
189	224	147	108	223	210	127	102	36	101	255	224 18	224	147	105	227	210	127	102	36	101	255	224	189	224	147	105	227	210	127	102	35	101	255	224	189 :	224	147	108	227	210	127	102	35	101	255	224
190	214	173	66	109	9 143	96	60	2	109	245	210 190	214	173	66	109	143	96	60	2	109	249	210	190	214	173	66	109	143	96	50	2	109	249	210	190 :	214	173	66	109	143	96	50	2	109	249	210
187	198	235	75	1	81	47	0	8	217	255	211 18	198	236	75	1	81	47	0	6	217	258	211	187	195	236	75	1	81	47	0	6	217	258	211	187	196 :	825	75	1	81	47	0	6	217	255	211
167	153	174	168	160	152	161	172	161	155	166	187 163	153	174	168	150	152	151	172	161	165	156	187	167	163	174	168	150	162	151	172	161	155	156	187	167 -	163	174	168	150	162	151	172	161	165	156	187
155	182	163	74	75	62	33	17	110	210	180	154 158	182	163	74	75	50	33	17	110	210	150	154	155	182	163	74	75	62	33	17	110	210	180	154	155	182	163	74	75	62	33	17	110	210	180	154
180	160	50	14	34	6	10	33	48	108	160	181 180	180	50	14	34	6	10	33	48	106	150	181	180	180	60	14	34	6	10	33	48	106	150	181	180	180	60	14	34	6	10	33	48	106	150	181
201	109	5	124	13	111	120	204	100	5 16	56	180 200	109	5	124	131	111	120	204	166	16	56	180	205	109	5	124	131	111	120	204	100	16	55	180 :	206	109	5	124	131	111	120	204	100	16	55	180
172	105	207	253	233	214	220	220	225	90	74	206 172	105	207	233	233	214	220	230	220	90	74	206	172	105	207	235	253	214	220	230	220	90	74	206	172	105	207	235 :	253	214	220	230	220	90	74	206
188	88	179	209	18	3 216	211	150	135	9 75	20	168 188	88	179	209	186	216	211	150	139	75	20	168	188	88	179	209	186	216	211	150	139	76	20	168	158	88	179	209	186	216	211	150	139	76	20	168
183	97	165	84	10	168	134	11	31	62	22	148 180	97	165	84	10	168	134	11	31	62	22	148	189	97	165	84	10	168	134	11	31	62	22	148	189	97	165	84	10	168	134	11	31	62	22	148
199	168	191	193	168	3 227	178	143	182	108	38	196 196	163	191	193	158	227	178	143	182	106	36	196	199	168	191	193	158	227	178	143	182	106	36	196	199	168	191	193	158	227	178	143	162	106	36	196
205	174	155	250	230	5 167	16	150	79	35	218	241 203	174	155	252	236	187	85	150	79	28	218	241	205	174	155	252	226	187	85	150	79	20	218	241 :	205	174	155	252 :	226	187	85	150	79	20	218	241
189	224	147	108	227	210	127	102	38	101	255	224 18	224	147	108	227	210	127	102	36	101	255	224	189	224	147	108	227	210	127	102	36	101	255	224	189 ;	224	147	108 :	227	210	127	102	36	101	255	224
190	214	173	05	109	143	95	50	2	109	245	210 190	214	173	66	109	143	95	50	2	109	249	210	190	214	173	m	109	143	95	50	2	109	249	210	190 3	214	173	66	109	143	95	50	2	109	249	210
187	108	236	75		81	47	0	6	217	255	211 18	198	258	75	1	81	47	0	6	217	255	211	187	196	238	76	1	81	47	0	6	217	266	211	187	198	238	76	1	81	47	0	8	217	266	211
167	153	174	168	160	152	151	172	161	155	166	187 16	153	174	168	150	152	151	172	161	155	150	187	167	153	174	168	150	152	151	172	161	155	158	187	167	153	74	168	150	152	151	172	161	155	156	187
155	182	163	74	75	82	33	17	110	210	180	154 154	182	163	74	75	82	33	17	110	210	180	154	155	182	163	74	75	80	33	17	110	210	180	154	155	182	163	74	75	60	33	17	110	210	480	154
180	150	50	14	34	6	10	31	48	109	150	181 140	100	50	14	34		10	33	- 10	105	150	181	140	180	50	14	34		10	33	48	105	150	181	100	180	50	14	34		10	33	48	105	150	181
																×.			~									· *							- 1						1					

But this is what a computer "sees"—just an extensive matrix of numerical information. Just numbers. This difference in interpretation is what is commonly referred to as "Semantic Gap". A semantic gap "characterises the difference between two descriptions of an object by different ... representations, for instance, languages or symbols".

3.1 3.1 Images



- Each pixel in an image is a feature
 - numerical
 - * 0 or 1 for *Black and White*
 - * Between 0 and 255 for greyscale
 - * 16M values for RGB
- Dimensionality -> *n x m*

Images are made of pixels and, in computer vision machine learning models, each pixel is an **input feature*. Depending on the type of the image, each feature (pixel) can assume a different numerical value.

Pixels in *Black and White* images have a value of either 0 (black) or 1 (white). *Grayscale* images have a wider intensity range, so each pixel has a value ranging between 0 (black) and 255 (white). *Color Images* have multiple color channels, represented using different color models (e.g., RGB, LAB, HSV). For example, an image in the RGB model consists of red, green, and blue channels. Each pixel in a channel has intensity values ranging from 0 and 255.

3.2 3.2 Computer Vision

- Building algorithms that can "understand" the content of images and use it for other applications
 - It is a "Strong AI" problem
 - * signal-to-symbol conversion
 - * The semantic gap -A general-purpose vision system requires
 - Flexible, robust visual representation
 - * Updated and maintained
 - Reasoning
 - Interfacing with attention, goals, and plans

Computer vision is about automating and integrating various processes and representations used for visual perception. And this is difficult, very difficult.

Vision requires many capabilities we often take for granted. For instance, the ability to extract intrinsic images of "lightness," "color," and "range." We perceive black as black in a complex scene even when the lighting is such that some black patches reflect more light than some white patches. Similarly, perceived colors are not related simply to the wavelengths of the reflected light; if they were, we would consciously see colors changing with illumination.

But vision is not only about **sensing**, it is also about **interpreting**, i.e. processing the acquired information to extract meaning from it. In the context of computer vision, the Semantic Gap can be defined as "... the lack of coincidence between the information that one can extract from the visual data and the interpretation that the same data have for a user in a given situation."⁴¹

The human brain solves this in multiple steps in different brain regions. Take, for instance, the ability to *recognise objects*. This ability is acquired either biologically (e.g. recognizing faces), through development (e.g. recognizing different types of vehicles), or through learning (e.g. recognizing a specific type of plane).

Computer vision is intrinsically tricky because it requires "re-inventing" the most basic yet inaccessible abilities of biological visual systems.

3.3 3.3 What specific tasks can we train a CV system to perform?

Due to all these complexities, it is not possible to implement a vision system that can emulate the human one. What recent advances in machine learning allowed, however, is the ability to perform *specific vision tasks*, in predefined contexts, with very good performance. We will now see several examples.

¹ Smeulders, A.W., Worring, M., Santini, S., Gupta, A., Jain, R.: Content-based image retrieval at the end of the early years. IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI) 22, 1349–1380 (12 2000). https://doi.org/10.1109/34.895972

3.3.1 Object Recognition / Localisation



In an image, drawing a *bounding box* around a specific object is called **object localization**. Object localization means recognizing the existence of one or more objects (of any type) in an image and marking the location with a bounding box.

Object recognition (or object **classification**) is the ability to *both locate and classify* objects in an image. Given an image, the object recognition tasks produce in output *one or more* multiple bounding boxes. Each bounding box marks the objects' location and class.



Current CV technology can recognize hundreds (or more) of objects at the same time. It is possible, for instance, to use object recognition to *count* the number of people in a given area (e.g., for crowd control purposes).



It is possible to do so also in real-time, in live video streams. The examples in the picture use a popular CV technique called Yolo (see a recent implementation here.



Face detection has been used for multiple years in cameras to take better pictures and focus on faces. Smile detection can allow a camera to take pictures automatically when the subject is smiling.

3.3.2 Object Identification



Face identification is more difficult than *face detection*, as the goal is not to recognize *any face*, but to *identify* the person having a specific face.



In the past, this task could have been performed only by security services and large organisations (e.g. airports). Today, thanks to the availability of large collections of high-quality labeled data (e.g. your personally curated photo albums on Google Images or Apple Photos, or Facebook), achieving extremely good performance on very large populations is possible. For the same reason, face identification can also be used for biometric identification.



The identification task can also be performed on specific objects in the real world.



ML4D Lecture Notes [V0.61]

Alessandro Bozzon

Applications like Google Lens allow one to identify monuments, landmarks, books, plants, and other object classes.

3.3.3 Image Segmentation



The **image segmentation** task is devoted to the **partitioning** of a digital image into multiple **image segments** (sets of pixels), also known as *image regions*, each containing a defined object class. Image segmentation is the **assigning a label to every pixel in an image**, aiming to simplify the representation of an image into something more meaningful and easier to analyze.



Image segmentation is very important for all those CV applications where detecting an object's contours (boundaries) is essential. Imagine, for instance, self-driving cars. Or medical imaging applications where the goal is to locate tumors, or measure tissue size/volume.



An example of use of image segmentation that I like is Project Sunroof. In that application, image segmentation has been used to delineate roofs from satellite images to estimate solar exposure for PV installations.

Alessandro Bozzon

3.3.4 Scene Recognition



Scene categorization is a task in which scenes from photographs are categorically classified. Unlike object classification, which focuses on classifying prominent objects in the foreground, Scene recognition involves a different set of challenges from those posed by object recognition. Like objects, scenes are composed of parts, but whereas objects have strong constraints on their parts distribution, scene elements are governed by much weaker spatial constraints.



An example of dataset (and scene recognition algorithm) is the MIT Places dataset. The dataset contains more than 10 million images comprising 400+ unique scene categories at a different level of semantic resolution. For instance, there are the classes *forest - broadleaf*, *forest path*, *forest road*; *garage indoor* and *garage outdoor* and so on.



You can try an example of scene recognition algorithm yourself, on HuggingFace.

3.3.5 Human Activity Recognition



Human Activity Recognition is the problem of identifying events performed by humans given a video input. It is formulated as a binary (or multiclass) classification problem of outputting activity class labels. Activity Recognition is an important problem with many societal applications including smart surveillance, video search/retrieval, intelligent robots, and other monitoring systems ².

3.3.6 Pose Estimation



While activity recognition is commonly performed by wearable devices (by analysing accelerometer and other data) in computer vision **pose estimation* is a common way to classify activities.

In **Pose Estimation**, the goal is to detect the position and orientation of a person. Usually, this is done by predicting the location of specific key points like hands, head, elbows, and knees.

² https://paperswithcode.com/task/activity-recognition



This video, created in the context of the master thesis Ethical task tracking of operators in agile manufacturing shows an example of how pose estimation can be used to track the actions of workers in a production plant, to support them with knowledge about their currently executed tasks,

3.3.6.1 Stereolabs ZED Camera

3D Object Detection Body tracking Positional tracking



This technology is commonly and readily available. The Stereolab ZED 2 camera offers an integrated solution that could be easily deployed and adapted to specific application domains.



The list of computer vision applications can be **very** long. I recommend you explore the examples in this page, classified according to application domains.

3.4 3.4 How do humans see?

Vision is one of the ways humans perceive the world. At its most basic, visual perception is observing patterns and objects through sight or visual input. Visual perception relates visual input to a previously existing *understanding* of the world, as constructed through previous experience and learning. Over 50 percent of the processing in the human brain is dedicated to visual information. This fact alone should give a good feeling of how difficult it is to replicate biological vision with a computer. Computer vision has not been solved in 50 years and is still a tough problem.

3.4.1 Hubel and Wiesel, 1959





In 1962, Hubel & Wiesel[^10 David H Hubel and Torsten N Wiesel. Receptive fields, binocular interaction and functional architecture in the cat's visual cortex. The Journal of Physiology, 160(1):106–154, 1962] set up an experiment to study the optical system of a cat. They recorded neurons while showing bright lines to a "wired" cat. They found that some specialized neurons fired only when the line was in a

Alessandro Bozzon

particular spot on the retina or if it had a specific orientation. If you open this video and turn up your volume, you can listen to the neuronal activity of the cat's visual cortex.

The experiments by Hubel and Wiesel are cornerstones of our understanding of how neurons along the visual pathway extract increasingly complex information to construct an image. Their research spearheaded a discipline devoted to understanding the human visual system's workings. And in 1981, they were awarded the Nobel Prize in Physiology and Medicine for their work.

3.4.2 Neural Pathways



Picture from https://nba.uth.tmc.edu/neuroscience/m/s2/chapter15.html

Vision systems are the same for humans, animals, insects, and most living organisms at the highest level. A visual system is built of two fundamental components. A **sensing device** - your eyes, and an **interpreting device**, your brain.

The eye captures the light coming through the *iris* and projects it to the *retina*. The retina contains specialized light-sensitive receptors that convert the image into spatially distributed neural activity in the first neurons of the visual pathway. Stimulus features (e.g., color, brightness contrast, movement) are processed (in parallel) at all levels of the visual system and recomposed by the simultaneous activation of large areas of the visual cortex, which resides at the back of your head. The visual cortex is roughly organized as a hierarchical series of *layers* where the neurons in each layer communicate their activations to neurons in the next layer. Through their experiments, Hubel and Wiesel discovered that neurons in different layers act as *detectors* that respond to increasingly complex features appearing in the visual scene: neurons at initial layers become active in response to *edges*; their activation feeds into layers of neurons that respond to *simple shapes* made up of these edges, and so on.

Note that as the complexity of the recognized patterns grows, different brain regions (the **interpret-ing device**) to get activated to associate **meaning** to the visual information and connect it to one's knowledge and understanding of the world.

This is a relatively simplified description of how the visual system works. An important note: in reality, the visual system does not only have *feed-forward* pathways - i.e., flow of information from the eye to the brain. There are also many more (ten times more) *feed-backward* connections - from the brain to

the lower layers - whose role is not well understood by neuroscientists. The current hypothesis is that our prior knowledge and expectations about the world, as stored in higher brain layers, can strongly influence our perceptions. The sentence "seeing what we want to see" might have a literal meaning.

3.4.3 Neural Correlation of Objects & Scene Recognition



The slides' pictures show how different brain areas get "activated" based on the visual stimuli shown for instance, a face, an object, or a scene.

Hubel and Wiesel's discoveries inspired a Japanese engineer named Kunihiko Fukushima, who in the 1970s developed one of the earliest deep neural networks, dubbed the cognitron, and its successor, the *neocognitron*.

3.5 3.5 Why is machine vision hard?

3.5.1 The deformable and truncated cat



What's so hard about computer vision? Scientists (a long time ago) believed that vision would be an easy problem to solve. There's this famous AI memo from Seymour Papert, who, in 1966, proposed a summer vision research project with interns to solve computer vision for a few months. He severely underestimated the task at hand.

Consider the problem of getting a computer program to *recognize* cats in photographs.

Suppose the input is simply the pixels of the image. In that case, the program first has to figure out which are "cat" pixels - that is, pixels that contain visual information about the cat - and which are "non-cat" pixels, that is, pixels having information about the background, shadows, or other objects.

However, cats are fascinating creatures, aren't they? They look very different: they can have diverse coloring, shapes, and sizes. This means the pixels associated with "cats" may vary greatly from image to image. When taking a picture of them, cats can face in various directions; the lighting can vary considerably between images; parts of the cat can be blocked by other objects (for example, fences and people). Cats are deformable [^11 The Truth About Cats and Dogs] Moreover, "cat pixels" might look like "dog pixels" or other animals. Under some lighting conditions, a cloud in the sky might even look very much like a cat.



All these representational variations make the problem of recognizing objects in computer vision very difficult. The image in the slide comes from a now classic paper Strike (with) a Pose: Neural Networks Are Easily Fooled by Strange Poses of Familiar Objects. As the title entails, the paper shows how non-canonical poses of ordinary objects can easily fool a neural network.

3.5.2 Computer Vision Challenges

The following slides summarise and explain some of the main challenges computer vision systems must address. Note that all of these challenges are present **simultaneously**. Solving one might not necessarily help solve all the others.

3.5.2.1 Viewpoint Variation

• A single instance of an object can be oriented in many ways to the camera.

Viewpoint variation



Let us take, as an example, the task of recognizing objects in an image. Depending on the viewpoint, the 2D representation of a 3D object might dramatically change, despite this being the same object. Depending on the viewpoint, the pixels change dramatically.

3.5.2.2 Deformation

Deformation



• Many objects of interest are not rigid bodies and can be deformed in extreme ways.

The same is true for *deformation*. Many objects (and animals, like cats) are not rigid. This means that they can be deformed in ways that, when captured in a picture, will create collections of pixels that are entirely unique.

3.5.2.3 Occlusion



• The objects of interest can be occluded. Sometimes only a tiny portion of an object (as few pixels) could be visible.

Sometimes only a relatively small portion of an object can be visible. For instance, because it is cost out of the image or because it is *occluded* by another object. In the picture, the cat is occluded by a wall.

3.5.2.4 Illumination Condition

• The effects of illumination can be drastic on the pixel level.
Illumination conditions



Look at the pictures in the slide, can you recognize what they show? The scene is the same but illuminated by different lights. Depending on where the light is projected, our understanding of the scene changes dramatically, as the organization of the pixels in the image. But it is precisely the same scene, pictured under the same viewpoint. This is to show an intricate interplay between materials and light that give rise to many different images despite showing precisely the same scene.

3.5.2.5 Scale variation

- Visual classes often exhibit variation in their size
 - Size in the real world
 - Size in the image



Variations of the size of the object to be recognized also matter. Objects of the same class in the real world can have different sizes - there are small and big cats, various sizes of autos, etc. However, also the size of an object in the image matters. Small things in images are captured by fewer pixels. This means there is less visual information to use when recognizing them.

3.5.2.6 Background clutter



• The objects of interest may blend into their environment, making them hard to identify.

The recognition of an object in an image can be made more difficult by "visual distractions," that is, the background or elements of the image, such as patterns, textures, colors, shapes, or other objects that are visually similar to the object of interest or that overlap with it in some way.

3.5.2.7 Intra-class variation

- The classes of interest can often be relatively broad, such as chairs.
- There are many different types of these objects, each with their appearance.



Finally, a big challenge is a variation in the object itself. Despite having the same nature or functionality, objects within the same class (e.g., cats, dogs, chairs, tables) might look very different. To distinguish a cat from a dog, for instance, a computer vision system has to be able to solve the problem that despite

many cats looking different, they are part of the same category. We call this the *intra-class variation* (*intra = within*) problem.

Of course, there is also an issue with the sheer number of object categories that exist in the world – it is estimated that there are between *ten thousand* to *thirty thousand* high-level object categories. Of course, the number immediately increases when the categorization becomes more fine-grained. ImageNet, a popular computer vision dataset, contains over 20,000 categories, with a typical category consisting of several hundred images.

3.6 3.6 How Computer Vision models work?

Images are made of pixels, and, in computer vision machine learning models, the more straightforward approach is to use the pixel values directly **input features*.

Here we need to take a small diversion and introduce the concept of *similarity* in a representation space; and the *manifold hypothesis*, an essential assumption at the basis of machine learning research.

Similarity refers to the degree of resemblance or closeness between two objects or data points in a given feature space. Intuitively, this could be described with the idea that "cat pictures" have similar pixels distribution because cats look similar. We saw in the previous section that this might not be the case; however, with enough pictures of diverse enough cats (in diverse enough lighting conditions, etc.), a fundamental assumption in machine learning is that it is possible to create a "decision boundary" that identifies cats in that feature space. Whether this assumption is realistic depends, of course, on the amount, diversity, and quality of the available training data.

One of the characteristics of high dimensional data (e.g., images in a dataset) is that the number of dimensions is comparable to, or larger than, the number of samples. The *Manifold Hypothesis* states that real-world high-dimensional data lie on low-dimensional manifolds embedded within the high-dimensional space. The hypothesis is based on the idea that high-dimensional data is often redundant (i.e., not all the dimensions used to represent a data item contain helpful information), and only a few underlying factors or features determine the structure and patterns in the data. Due to the manifold hypothesis, many data sets that appear to initially require many variables to describe can be described by a comparatively small number of variables. The Manifold Hypothesis explains why machine learning techniques can find useful features and produce accurate predictions from datasets that have a potentially large number of dimensions (variables). The fact that the actual data set of interest lives on in a space of low dimension means that a given machine learning model only needs to learn to focus on a few key features of the dataset to make decisions.

3.6.1 Course of dimensionality

- High dimensionality
 - A 1024×768 image has d = 786432!
 - A tiny 32×32 image has d = 1024
- Decision boundaries in pixel space are extremely complex
- We will need "big" ML models with lots of parameters
 - For example, linear regressors need d parameters



The number of pixels (hence, the number of features) can become too big for a machine-learning model. Presume you have a 1MB image, where each pixel is represented by a single byte (0..255 value). At 1MB, you have one million pixels. That would require an input vector of 1,000,000 elements. Assuming that the input layer has 1024 nodes, there will be over a billion (1 million x 1024) weights to learn, just in the input layer! This number increases drastically when we have tens or hundreds of layers. Even if, in the news, you can read of ML models with 200 or 300 billion parameters, this is a number that cannot be easily handled.

Therefore, the first step after preprocessing the image is to simplify the image by 1) reducing its color complexity (e.g., transforming it into black and white); 2) resizing the image to make it smaller; 3) extracting the vital information (features) and throwing away non-essential information. Each one of these techniques can be applied in conjunction. Modern ML techniques based on deep learning only require resizing the images, keeping color information, and not requiring feature extraction.

3.6.1.1 Downsampling



Downsampling reduces the resolution of an image to make it more tractable from a computational perspective. It is common in CV applications to reduce images to 1024x1024 resolution or even 512x512.

Downsampling allows for faster learning and processing time, but it also comes at a disadvantage. By reducing the image resolution too much, it is possible to lose the ability to distinguish what's in the image.

3.6.1.2 Flattening



In computer vision applications, we deal with images or video. A picture can be represented as a 2-dimensional matrix, a grid of pixels.

However, machine learning approaches typically deal with input data organized in arrays (sets or lists) of features. For instance, the input layer of a Neural Network is an array of numerical values of size *d*.

Flattening is placing each row of the image matrix in sequential order into a vector. So the vector starts with the first row of pixels, followed by the second row of pixels, and continues by ending with the last row of pixels.

This simple transformation has a significant consequence: *any spatial relationship between the pixels is lost*! This means that contiguous pixels in the array might not have the same contiguity in the original image. As we will see, modern computer vision techniques can overcome this slight but essential distortion introduced by a mere technical requirement.

Note that flattening is not **always** used as the first step in a machine-learning approach. As we will see later, Convolutional Neural Networks allow for keeping the matrix format of an image and even retaining color information. The feature flattening happens later in the network.

3.6.2 The "old days": Feature Extraction and Engineering

- Feature
 - A relevant piece of information about the content of an image
 - e.g., edges, corners, blobs (regions), ridges
- A good feature
 - Repeatable

- Identifiable
- Can be easily tracked and compared
- Consistent across different scales, lighting conditions, and viewing angles
- Visible in noisy images or when only part of an object is visible
- Can distinguish objects from one another



A *feature* in machine learning is an individual measurable property or characteristic of an observed phenomenon. Features are the input you feed to your machine learning model to output a prediction or classification. Selecting good features that clearly distinguish your objects increases the predictive power of machine learning algorithms.

In computer vision, a feature is a measurable piece of data in your image that is unique to that specific object. A CV feature is a group of connected pixels with some common property. It may be a distinct color or a particular shape, such as a line, edge, or image segment.

Features are helpful to "compress" information. The input image has too much extra information that is unnecessary for classification. Therefore, the first step after preprocessing the image is simplifying it by extracting the important information and throwing away nonessential information. By extracting, for instance, important colors or image segments, complex and large image data can be transformed into smaller sets of features. This makes classifying images based on their features simpler and faster.

A good feature is used to distinguish objects from one another. For example, if we consider a feature like a *wheel*. Clearly, that feature is more likely to be associated with the class "motorcycle" than the class "dog". A *wheel* is a strong *identifiable* feature that clearly distinguishes between motorcycles and dogs. However, the same feature will probably not be strong *(identifiable*) enough to distinguish a bike

from a bicycle. For that purpose, more features are needed, like a mirror, license plate, or maybe a pedal.

A good feature is also *repeatable*, i.e., it should be useful to represent a class of objects, not a single one. For instance, the visual representation of a *wheel* should not be a single image's exact copy of a wheel. Still, it should be generic enough to represent many motorcycle wheels (see the image at the bottom). So, it should look like a circular shape with some patterns that identify wheels in all images in the training dataset.

Good features should also be easy to *compute*, to *track* and *compare*. If a feature is complex to calculate or very brittle to compare, it will probably not be very useful to the machine learning model.

Finally, an idea feature is also *consistent* across different conditions (sales, lights, viewing angle) and *robust* to noise and occlusions.

- Machine learning models are only as good as the features you provide
 - To figure out which features you should use for a specific problem
 - Rely on domain knowledge (or partner with domain experts)
 - Experiment to create features that make machine learning algorithms work better



Before the era of deep neural networks, computer vision required spending a lot of time in manual feature selection and engineering.

Feature engineering relates to manipulating and transforming data into a format that optimally represents the underlying problem that an ML algorithm tries to model and mitigates inherent complexities and biases within the data.

In this process, the ML engineer relied on domain knowledge – potentially partnering up with domain experts – to create features that make ML algorithms work better.

The features are in specific locations of the images, such as mountain peaks, building corners, doorways, or interestingly shaped patches of snow. These kinds of localized features are often called *keypoint features* (or even corners) and are often described by the appearance of patches of pixels surrounding the point location. The features that can be matched based on their orientation and local appearance (edge profiles) are called *edges*. They can also indicate object boundaries and occlusion events in the image sequence.

These features are used as input to machine learning algorithms like (fully connected) neural networks or random forests for classification, that is, to learn the correlation between such features and the prediction class.

3.6.2.1 Feature Extraction Techniques



In this course, it is impossible to cover the whole field of features engineering for computer vision - there is literature going back thirty years. I will, however, describe some examples.

Histograms of oriented gradients (HOGs) is an example of a feature extraction technique commonly used for object recognition tasks. HOG focuses on the object's shape in the image by attempting to quantify the gradient (or magnitude) and the orientation (or direction) of the edges of an object. HOG calculates gradients and orientations in broken-down, localized regions of the image and calculates a histogram of gradients and orientations to determine the final feature values.

Another technique is the **Scale Invariant Feature Transform** or SIFT. The SIFT descriptor is invariant to translations, rotations, and scaling transformations. It is also robust to moderate perspective transformations and illumination variations.

3.6.2.2 Performance



Credits: Ross Girshick (Facebook AI Research)

Feature engineering and "old" style feature extraction techniques have been the foundation of computer vision for many years. And while advancements were consistent, they never reached a level of performance that allowed for wide-scale applications.

Convolutional Neural Networks changed everything.

3.7 3.7 Convolutional Neural Networks



- CNNs exploit image properties to reduce the number of model parameters drastically
- Feature maps
 - Automatically extracted hierarchical
 - Retain spatial association between pixels
- Local interactions
 - all processing happens within tiny image windows
 - within each layer, far-away pixels cannot influence nearby pixels
- Translation invariance
 - a dog is a dog even if its image is shifted by a few pixels

Convolutional Neural Networks (CNN) drive today's deep-learning revolution in computer vision and other areas. CNNs were first proposed in the 1980s by the French computer scientist Yann LeCun, inspired by Fukushima's neocognitron. The design of CNN is based on several critical insights about the brain's visual system discovered by Hubel and Wiesel.

CNNs allow for eliminating Feature Engineering, as the network can be used for feature learning and classification.

ConvNet consists of a sequence of layers of perceptrons. Perceptrons in each layer provide input to perceptrons in the next layer. Like the neural network, when a ConvNet processes an image, each perceptron takes on a particular activation value. This number is computed from the unit's inputs and their weights. Weights are randomly initiated and learned during network training. The difference between a fully connected neural network and a CNN is the use **convolutional layers** instead of regular fully connected layers for feature- learning.

A **convolutional layer** is the core building block of a convolutional neural network. Convolutional layers act like a *feature finder* window that slides over the image pixel by pixel to extract meaningful features that identify the objects in the image.

Layers are locally connected. This means that nodes in the layer are connected to only a small subset of

the previous layers' nodes - for instance, pixels in the input image. This way, filters are applied only on close-by pixels without being influenced by far-away ones. This allows filters to focus more, preserve local spatial relationships, and reduce the overall computational cost.

A **feature map** is the output of one filter applied to the previous layer. It is called a feature map because it maps where a specific feature is found in the image. CNNs look for features such as *straight lines*, *edges*, or even objects in the deeper layers. Whenever they spot these features, they report them to the feature map.

CNNs allow for *translation invariance* for images fed through the network. This means the network can recognize patterns (or shapes) that are shifted or slightly warped within images.

Notice that the image dimensions shrink after each layer, and the number of feature maps (the layer depth) increases. Conceptually, you can think of this set of consecutive convolutional layers as the neural network learning to represent more abstract features of the original image.

The output of the feature extraction step is then flattened to a vector of the learned features of the image. Notice that the image dimensions shrink after each layer, and the number of feature maps (the layer depth) increases until we have a long array of small features in the last layer of the feature-extraction part. The flattened feature vector is fed to the fully connected layers (a traditional fully-connected neural network) to classify the extracted features of the image.

3.7.1 Convolution & Feature Maps



Each convolutional layer contains one or more **convolutional filters**. The number of filters in each convolutional layer determines the depth of the next layer because each filter produces its feature map (convolved image).

In mathematics, *convolution* is the operation of two functions to produce a third modified function. In the context of CNNs, the first function is the input image, and the second is the convolutional filter. In the slides figure (left), the convolutional filter (also called a **kernel**) is the 3x3 matrix that, by sliding over the input image, breaks the image into little chunks and processes those chunks individually to assemble the modified image, a **feature map** (right).

Not that in CNNs, convolution matrixes are **the weights**. This means the network learns its values during training.

Try this: https://cs.stanford.edu/people/karpathy/convnetjs/demo/mnist.html

3.7.2 What CNNs learn?

Deep Visualization Toolbox

Watch this video.

3.7.2.1 Feature Visualisation



Feature visualization of convolutional net trained on ImageNet from [Zeiler & Fergus 2013]

In the following slides and pictures, we will give an impression of what a CNN learns through its convolutional layers.



For instance, the first layer of the CNN consists of edge-detecting units. This picture is taken from the paper: *Visualizing and Understanding Convolutional Network. Zeiler and Fergus, ECCV 2014*.

The second layer can recognize more complex patterns and shapes.

While the third layer can now capture the visual representation of more "abstract" concepts, like faces or wheels.

3.7.2.2 Network Dissection



This image is taken from the Network Dissection dataset and related papers. It shows how specific concepts are "learned" by different CNN architectures.

3.7.3 Translation Invariance



• But not rotation and scaling invariance!

Remember, CNN allows for *translation* invariance. They can recognize features (e.g., a wheel) even if it appears in different positions of an image. They are, however, **not** *rotation or scaling invariant*. Objects that are rotated or scaled w.r.t. their representation in the training data might not be recognized.

3.7.4 What about generalisation?



2023-03-29

As we often repeated, one of the main challenges in ML is to give the model the ability to *generalize* beyond the training data. However, what could be done if a CNN (or any network architecture) does not offer rotation or translation invariance?

The answer is in **data augmentation techniques**.

3.7.5 Data Augmentation



- Generate variations of the input data
 - To improve generalisability (out-of-distribution inputs)
 - Improve invariance (rotation, scaling, distortion)

Data augmentation means expanding the training dataset with modified versions of the current images. Scaling, flipping, rotations, and other affine transformations are typically used to enlarge your dataset and expose the neural network to various variations of the training images. Alternatively, *color space augmentations* (photometric transformations), random cropping, or noise injection allow the network to be more robust and makes it more likely that your model will recognize objects when they appear in any form and shape.

Machine Learning for Design



- Geometric
 - Flipping, Cropping, Rotation, Translation,
- Noise Injection
- Color space transformation
- Mixing Images
- Random erasing
- Adversarial training
- GAN-based image generation

A good survey on state-of-the-art data augmentation techniques is: A survey on Image Data Augmentation for Deep Learning. Shorten, Journal of Big Data, 2019.

3.7.6 Robustness to input variation



Let's keep in mind that, despite all these attempts to extend the training datasets through data augmentations, computer vision models like CNN are still very brittle to variations of objects (and their representation) in images. The paper Strike (with) a Pose: Neural Networks Are Easily Fooled by Strange Poses of Familiar Objects. Alcorn et al. 2019. describes this problem, although I encourage you to try using a tool like Teachable Machine.

3.7.7 Transfer Learning



- Problem: training custom ML models requires huge datasets
- **Transfer learning**: take a model trained on the same data type for a similar task and apply it to a specialised task using our custom data.
 - Same data: same data modality. same types of images (e.g., professional pictures vs. Social media pictures)
 - **Similar tasks**: if you need a new object classification model, use a model pre-trained for object classification

Transfer learning is one of the most important techniques that emerged from modern deep learning. Building a vision system to solve a specific problem requires collecting and label a vast amount of data to train your network. But what if this is not possible?

Transfer learning is the transfer of the knowledge (*feature maps*) that the network has acquired from one task, where we have a large amount of data, to a new task where data is unavailable or available only in small quantity. The idea behind transfer learning is to utilize the weights and layers from a model trained in the same domain as your prediction task. Intuitively, all images must have shapes and edges, so the early layers are usually transferable between domains.

In most deep learning models, the final (fully connected) layer contains the classification label or output specific to a prediction task. With transfer learning, the last layers are removed to keep only the learned features that could be effectively transferred across tasks. However, the transferability of features depends on the similarity of the training and new datasets. This is especially true for later layers: the deeper the network layers, the more image-specific the learned features. For instance, imagine the task of recognizing objects in medieval paintings. An object recognition network trained

on modern images (e.g., from Flick) would probably not be helpful, as bikes and cars (and their parts) were perhaps unavailable in medieval times. However, the earlier layers of the network – the ones recognising edges and shapes – are undoubtedly helpful, as any object shares the same geometry set.

At training time, the preserved model's trained weights are frozen, and the final layer (the fully connected classification one) is replaced with a specialized prediction task output. The new network can then be trained, and only the last weights (the ones of the classifier plus any other unfrozen layer) are learned.

Transfer learning is also frequently applied in image object detection, image style transfer, image generation, text classification, machine translation, and more. Transfer learning works because it allows to utilization of models already trained on extremely large, labeled datasets.

3.8 3.8 Advanced Computer Vision Techniques



3.8.1 Generative Adversarial Networks

- Learn patterns from the training dataset and create new images that have a similar distribution of the training set
- Two deep neural networks that compete with each other
 - The **generator** tries to convert random noise into observations that look as if they have been sampled from the original dataset
 - The **discriminator** tries to predict whether an observation comes from the original dataset or is one of the generator's forgeries

Generative adversarial network (GAN) is a class of deep learning models invented in 2014 by Ian Good-fellow. This architecture is inspired by game theory. Two models, a **generator** and a **discriminator**, are competing while making each other stronger simultaneously.

Alessandro Bozzon

GANs have shown remarkable results in many generative tasks to replicate real-world rich content such as images, human language, and music. GANs can also be used in text-to-photo synthesis, image-to-image translation, image super-resolution, and many other applications.

• The **generator**'s architecture looks like an inverted CNN that starts with a narrow input and is upsampled a few times until it reaches the desired size



• The **discriminator** 's model is a typical classification neural network that aims to classify images generated by the generator as real or fake



The **generator** attempts to generate real-looking images, whereas the **discriminator** is a classifier whose job is to separate authentic images from fake ones. The generator creates images using an architecture that resembles a reversed CNN. Starting from a vector of random noise, it creates an image. During the training process, the discriminator is provided with authentic images (from the training set) and fake images generated by the generator. The two networks train together until the discriminator can be "fooled" by the quality of the images created by the generator. The images that are created in the end are pretty realistic-looking. Note that these images are new, never seen before, and imaginary.

3.8.2 Which face is real?



To give an example of how good GANs can be, try this online "game". Can you recognize which face is real?

3.8.3 Image super-resolution GAN



• A good technical summary

Image Super-Resolution refers to enhancing an image's resolution from low resolution (LR) to high (HR). Thanks to generative technology, using an existing image at the input and using the network to generate new visually compatible pixels is possible. For instance, to enlarge it. Or to fix it.

Alessandro Bozzon



- ML-generated painting sold for \$432,500
- The network trained on a dataset of 15,000 portraits painted between the fourteenth and twentieth centuries
- Network "learned" the style and generated a new painting

GANs can also be used to learn the "style" of a specific set of paintings and to generate new ones in that style. An image generated by a GAN has been sold for good money.

3.8.4 Neural Style Transfer



Neural style transfer transfers the style from one image to another. It consists of an optimization technique used to take two images: 1) a content image and 2) a style reference image, such as an artwork by a famous painter. Neural style transfer blends them, so the output image looks like the content image but is "painted" in the style reference image style.

3.8.5 Text-To-Image Generation



A **text-to-image** model is an example of a generative ML model where the input is a **text prompt**, and the output is an image *matching* the text.

The ML systems learn (in a supervised way) from the text-image pairs how people describe images, their content, and their styles, or, the other way around, how visual concepts/scenes are textually described. These systems demonstrated exciting capabilities. Not only to create "realistic" artworks in a particular style but also, for instance, to create new versions of animals and objects, combining unrelated concepts in plausible ways.



Credits: https://github.com/CompVis/latent-diffusion

Design

This is an example of two images generated from the prompt *A dream of a classroom full of interested XXX students. Realistic, matte painting, HQ, 4k.* In one image, we prompted **design students** and **computer science students** in the other. Observing how the model has learned the association between visual concepts and text in almost stereotypical ways is intriguing. Computer science students are depicted with many screens, with a blue-ish illumination. Design students seem to have post-its on their desks and work in a better-lighted room.

You can try to generate images using the Dreamstudio.ai tools. I recommend you look at the prompt guide to understand how text prompts can be formulated and how they lead to different results.



Computer Science



The architecture of latent diffusion model. (Image source: Rombach & Blattmann, et al. 2022

The previous examples are not generated through GANS, but through a different approach called *Diffusion Models*

Standard *Diffusion Models* rely on Markov chains (a statistical model) and operate with two major processes: the *Forward Diffusion* takes an image and gradually corrupts it by introducing noise until it becomes utterly random noise. The *Reverse Diffusion* process Markov Chains recover the data by gradually removing the predicted noise at each time step. The idea is that, in this way, the model learns how to "fill the visual gaps", thus generating images that are close to the ones given for training. A **Conditioning** module lets the network associate visual properties with textual descriptions. The conditioning module contains textual representations (e.g., captions or longer text) associated with the images used for training. The image generation is *conditioned* to respond to a particular textual prompt; therefore, when a user specifies a new prompt, the network generates an image "similar" to

the ones associated with the text in the prompt.

If you are interested, this blog post describes the images used in the training of *Stable Diffusion*. These are for instance, images related to "Frida Kahlo."

You may wonder: how is it possible to train a model that, based on millions of images, can still create realistic images? How can (visual) information not be lost in the process? The reason is simple: **natural images are not random**. They have high regularity. A face follows a specific spatial relationship between the eyes, nose, cheek, and mouth. In other words, the high dimensionality of images is *artifactual* and the manifold hypothesis described in the previous lecture seems to hold, at least in the case of image processing.

Diffusion Models have shown incredible performance in many generative tasks, such as image generation, image synthesis, image substitution, and super-resolution. They overcome several limitations of GAN models, for instance, the ability to diversify their output more and not to "fixate" on some specific properties of the training sets (mainly because they can best preserve the semantic structure of the data). Diffusion Models are highly computationally demanding, and their training requires huge memory and generates an enormous carbon footprint. For instance, the estimated emissions generated by training a *Stable Diffusion v1* are 11250kg CO2 equivalent (calculated using the Machine Learning Impact calculator). This amount is the equivalent of 2.5 gasoline-powered passenger vehicles driven for one year.

New methods have been proposed to make the process much faster, but they are still slower and more computationally expensive than GANs.



3.8.6 Image-to-Image Generation

Video from: https://github.com/CompVis/latent-diffusion.

This video shows several examples of the application of diffusion models. For instance, the generation of realistic images given as input a simple sketch. Or substituting parts of a picture (a technique called "Inpainting"). Or super-resolution.

Alessandro Bozzon

3.8.7 Synthetic Video Generation

Generated from Synthesia.io

To conclude this quick overview of advanced computer vision capabilities, let me show you a video I have created on the Synthesia.io. Combining several techniques makes it possible today to generate realistic videos (like the one in the slide) simply by selecting an avatar, and providing some text to speak. It is, of course, also possible to create an avatar of yourself or anyone. This capability opens an excellent design space of possibilities.

These techniques yield considerable potential in the animation and movie industry. Classical techniques require tedious 3D model creation and manual editing. Machine Learning could automate this process and enable richer editing, or updating scenes without re-shooting them.

3.8.8 Deep Fakes

Very realistic Tom Cruise Deepfake

But these techniques also raise many concerns. **Deepfakes** (a term coined by combining *Deep Learning* and *fake*) are a form of synthesized media that can be created to simulate the presence of a real person in a media. You can imagine this technology's risks in our current era of widespread online misinformation and deep polarisation.

4 Section 4: Text Processing Methods

In a previous lecture, we introduced *Natural Language Processing* (**NLP**), a sub-field of Artificial Intelligence and Machine Learning that analyses natural language (written or spoken) to understand its content. Or, more recently, to generate realistic text and voices.

NLP technology can play several roles in the context of design. NLP is an enabling technology for products that base their interaction on natural language conversations—for example, popular personal assistants like Siri or Alexa or conversational systems like ChatGPT. NLP technology can support the design process, as it can help extract meaning from collections of textual documents and support its analysis, for instance, by automatically facilitating the thematic analysis of interview transcriptions.

4.1 4.1 Why natural language processing?

4.1.1 Big Textual Data = Language at scale

- One of the largest reflections of the world, a man-made one
- Essential to better understand people, organisations, products, services, systems
 - and their relationships!
- Language is a proxy for human behaviour and a strong signal of individual characteristics
 - Language is always situated
 - Language is also a political instrument

Language is an essential component of social interaction, and written text is, in a way, the materialization of language. Text reflects our world, as captured by words, stories, and books. Through text, humans record their laws, report news and events, capture history, and express feelings and emotions. Companies and governmental organizations use text to communicate their values and intentions, through advertisements or public messages. Text is produced by people, and therefore an expression of their context, culture, and politics. Even text generated by machines (e.g. ChatGPT) is, in reality, the product of a process of selection of documents (for training), styles of response, and safety measures that embed a specific view of the world.

Textual documents encode information about the world, which is why it has always been the object of study and interpretation. Scholars of history analyze text to understand the past and give interpretations to events for which little is left in our collective memory. Social scientists analyze text to

study how groups of people interact, relate to each other, and understand and evolve their culture and their world. Companies are interested in understanding how consumers respond to advertisement campaigns, or how they appreciate (or criticize) their products and services. Governments analyze text data to understand the views of citizens on proposed policies, or the consent (or dissent) of the elected officials.

Despite the ability to print text on a large scale, analysing an extensive collection of textual documents was impossible in the past. Reading large collections of documents was (and still is) very time-consuming and considered too daunting the task of manually organizing text into categories or tracking the presence or evolution of concepts of interest.

The widespread availability of digital computing technologies, including personal computers (from the 90s), mobile computers, and the Web (from the 2000s onwards), changed everything. Text has been not only very easy to create and distribute; but thanks to the increasing computational power of digital devices, the processing and analysis of extensive text collections have been democratized.

Despite all the recent advances in NLP technology, it is a mistake to believe that the qualitative methods developed in the long tradition of the humanities and social sciences could be simply replaced. NLP technology should be considered a tool to **assist** and **augment** reading, analysis, and interpretation abilities.



4.1.1.1 Fora, social media

A quick list of common sources of textual content.

Web pages, fora, and social media. They are incredibly broad sources of content that, despite their size, can suffer from issues of lack of diversity and bias. Size, it is important to remember, is never an

intrinsic positive quality for a document collection - especially if there is homogeneity in the type of content, and in the people that wrote it.

4.1.1.2 Product review



grown to understand and respect their genite ways, now i question everything I thought I once knew and fear I am no longer capable of following through with my primary objective. I know that those who sent me Amazon/customer review/<u>ByronicHero</u>.

Online e-commerce Websites like Amazon, eBay, and Alibaba are incredibly rich sources of information about products and services. Consumers spontaneously (most of the time) describe their experience and opinion with these products, creating input that could be very useful for companies to interpret the success (or failure) of particular commercial endeavors.

4.1.1.3 Books

• Digital, or digitised



4.1.1.4 Interviews



4.1.2 Uses of NLP

- Answer questions using the Web
- Translate documents from one language to another
- Do library research; summarize
- Archive and allow access to cultural heritage
- Interact with intelligent devices
- Manage messages intelligently
- · Help make informed decisions
- Follow directions given by any user
- Fix your spelling or grammar
- Grade exams
- Write poems or novels
- Listen and give advice
- Estimate public opinion
- Read everything and make predictions
- Interactively help people learn
- Help disabled people
- Help refugees/disaster victims
- Document or reinvigorate indigenous languages

The slides list examples of tasks that could be performed through NLP technology. The list is far from complete, but it should give you an indication of how broad could be the application of NLP technology in your work. As next text analysis methods emerge (e.g. large language models), it becomes cheaper

and more accessible to use NLP methods both to design and in designs.

4.2 4.2 What is Natural Language Processing?

• Computer using natural language as input and/or output



- Natural: human communication, unlike e.g., programming languages
- Language: signs, meanings, and a code connecting signs with their meanings
- Processing: computational methods to allow computers to 'understand', or to generate

Natural Language Processing uses computational methods to **understand** and **generate** language produced by humans to communicate with each other.

As for any computational methods, NLP approaches attempt to reduce the complexity of language so that it can be captured mathematically and processed through computers.

Language is a complex, multifaceted communication system used by humans to convey thoughts, ideas, emotions, and intentions. Language can take on various forms, including *spoken*, *written*, and *non-verbal* communication (such as gestures or facial expressions). Language consists of a set of *symbols* known as **phonemes** (in spoken language) and **graphemes** (in written language), which are combined according to specific rules to create meaningful units called **morphemes**, **words**, **phrases**, and **sentences**. These combinations form the basis of a language's **grammar**, which governs the structure and relationships between words and phrases, ensuring coherent and consistent communication.

4.2.1 Beyond keyword matching



- Identify the structure and meaning of words, sentences, texts and conversations
- Deep understanding of broad language

One of the simplest text-based tools that (almost) everyone is accustomed to are *search engines* – online, like *Google*, or on a computer, like *Finder* on a Mac. Search engines, in their original yet most common versions, are *keyword-matching* systems: given a query composed of a sequence of words (keywords), they find documents that contain such words and order them according to criteria such as size, recency or relevance to the query.

Natural Language Processing (NLP) technology offers more sophisticated approaches to understanding human language than simple keyword matching. Keyword matching falls short when dealing with the complexity and nuance inherent in human language. NLP encompasses a range of techniques and algorithms designed to tackle various linguistic challenges, enabling a deeper and more meaningful understanding of language.

4.3 4.3 Why is NLP Hard?

- Human languages are messy, ambiguous, and ever-changing
- A string may have many possible interpretations at every level
- The correct resolution of the ambiguity will depend on the *intended meaning*, which is often inferable from the *context*

Human language is not easy to tame. As individuals, we can sustain even complex communications with colleagues, friends, and family. We understand the dialogs we hear in tv series and movies, and

we (sort of) get what politicians tell us on TV. When we read a news item or a book, we comprehend what the writer is trying to say, evoke, and let us believe. But we are complex creatures, and our ability to communicate and understand language is both (to some extent) innate and acquired. We do not actively "learn" our mother tongue, but we literally grow into it. Through experience and imitation, we learn the association between words and visual, abstract, or emotional concepts. We possess a knowledge of the world - and its language representation - that helps us resolve possible ambiguity.

When we learn a new language, we comprehend the intrinsic cognitive complexity of actively learning a language. However, we can rely on the knowledge of our mother tongue to create associations and analogies. We can associate the same meaning with different words (e.g., synonyms) or expressions without noticing.

- There is tremendous **diversity** in human languages
- Languages express meaning in different ways
- Some languages express some meanings more readily/often

Language varies significantly across the world. Even if it appears that there exist some universal properties of languages, the way a language evolves and is used across communities, cultures, and countries is incredibly diverse.

Languages express meaning in different ways. We capture the culturally-related subtleties of language as they are expressed through idioms and figures of speech. For instance, the Dutch expression "Dit varkentje wassen" means "getting things done" or "taking care of a problem." The idiomatic expression, literally translated, has no straightforward equivalent meaning in Italian, where washing pigs is not associated with having confidence in one's ability to complete a task or deal with a challenge. But *language evolves*, so maybe it will also be an Italian figure of speech.

Some languages express some meanings more readily or more often. Think about how in a language like Italian, there are countless different names for *pasta* – spaghetti, linguine, fettuccine, penne, rigatoni, lasagne, and farfalle, to name just a few. Each name captures different properties of the pasta itself, properties on how it should be cooked (a specific type of sauce or filling), but also some properties of the region and the tradition connected to such pasta.

- Knowledge Bottleneck
- About language
- About the world: Common sense and Reasoning

All the complexity explained above cannot be captured by hard-coding rules over the presence or absence of keywords unless when working with a limited subset of a language and when addressing a minimal domain. This is why natural language processing is often said to be a **hard AI** problem: to process language at a human level, a computer system should be able to represent the world and acquire world knowledge. Just like a human does.

4.3.1 Ambiguity and Expressivity

Language is highly context-dependent, with the meaning of words and phrases often changing based on the surrounding context. **Ambiguity** is a common challenge in natural language understanding, as sentences can often be interpreted differently.

Christopher Robin is alive and well. **He** is the same person that you read about in the book, **Winnie the Pooh**. As a boy, **Chris** lived in a pretty home called **Cotchford Farm**. When **Chris** was three years old, **his father** wrote a poem about **him**. The poem was printed in a magazine for others to read. **Mr. Robin** then wrote a book

- Who wrote Winnie the Pooh?
- Where did **Chris** live?



Take the example text in the slide, and try to answer the question: who wrote Winnie the Pooh?

The text provides the answer, but implicitly. The text mentions the name *Christopher Robin*, but it is unclear if he is the person that wrote **Winnie the Pooh**, or by his father. We do not know if *Mr. Robin* is *Christopher Robin*; it is possible that *Christopher Robin* made a book after his father wrote a poem about him. To complicate things, *Christopher Robin* is also the name of a character in the **Winnie the Pooh** books. A quick online search will reveal that **Winnie the Pooh** was created by *Alan Alexander Milne*, the father of *Christopher Robin Milne*. *Christopher Robin Milne* is, coincidentally, also an author.

4.3.1.1 Lexical ambiguity

The presence of two or more possible meanings within a single word



Lexical Ambiguity is the presence of two or more possible meanings within a single word. The context in which a term is used determines the intended meaning. For example, the word *lost* in the comic can mean *losing sight of* or *losing due to death*. Another example is the word *bank*, which could refer to a financial institution, the side of a river, or, in Dutch, to a couch.

4.3.1.2 Syntactic ambiguity (Word sense ambiguity)

The presence of two or more possible meanings within a single sentence or sequence of words



The syntax of a language is the set of principles under which sequences of words are judged grammatically acceptable by fluent speakers. Syntactic ambiguity occurs when a sentence can be **parsed** (*interpreted*) in multiple ways due to its grammatical structure. For instance, the sentence I saw her duck can be interpreted as seeing a duck owned by a female person, seeing a female person lowering her head or body, or cutting a duck with a saw. The sentence I saw the Grand Canyon flying to New York is another example.

4.3.1.3 Attachment ambiguity

The policeman shot the thief with the gun

The slide shows an example of **Attachment ambiguity**, a type of syntactic ambiguity that emerge when a word or phrase can be associated with more than one element in a sentence, leading to different interpretations. This ambiguity arises because it is unclear how a specific part of a sentence should be "attached" to the rest of the sentence in terms of its grammatical structure.

4.3.1.4 Pronoun Reference ambiguity

A **Pronoun Reference Ambiguity** (or Anaphoric Ambiguity) occurs when a sentence uses a pronoun to refer to an antecedent (a person or a thing) in an ambiguous way. Like in the funny example in the slide.



Dr. Macklin often brings his dog Champion to visit with the patients. He just loves to give big, wet, sloppy kisses!

Source: https://www.printwand.com/blog/8-catastrophic-examples-of-word-choice-mistakes

4.3.1.5 Semantic Ambiguity



Every fifteen minutes a woman in this country gives birth. Our job is to find this woman, and stop her! Groucho Marx

Semantic Ambiguity occurs when an expression is semantically ambiguous when it can have multiple meanings – that is, there are different ways of reading the sentence.

Some examples:

- "John and Mary are married.": to each other? or separately?
 - Compare "John and Mary got engaged last month. Now, John and Mary are married."
- "John kissed his wife, and so did Sam." Did Sam kiss John's wife or his own?

Other examples of ambiguity include:

- **Idiomatic Ambiguity**: when a phrase or expression has a figurative meaning that differs from its literal interpretation. For instance, *break a leg*.
- **Pragmatic Ambiguity**: when the intended meaning of a statement depends on the context, the speaker's intentions, or the listener's background knowledge.

4.3.2 Sparsity

4.3.2.1 Zip's Law

"... given some document collection, the frequency of any word is inversely proportional to its rank in the frequency table..."



Zipf's Law is an empirical law formulated to describe the frequency distribution of words in natural languages. The law has been empirically formulated by observing how in a document collection, the most frequent word will occur approximately twice as often as the second most frequent word, which occurs twice as often as the fourth most frequent word, etc. More formally, the law states that *the frequency of a word in a large corpus is inversely proportional to its rank in the frequency table.*

Mathematically, Zipf's Law can be expressed as:

 $f(r) = c/r^z$

where:

- f(r) is the frequency of a word with rank r
- c is a constant specific to the language or text
- r is the rank of the word in the frequency table
- z is an exponent, usually close to 1. When a is 0 the distribution is uniform. As z increases, so does the skewness of the function.

Zipf's Law holds for all languages)including non-natural ones like Esperanto) but it is yet to be understood why. Differences in languages are reflected by differences in c and z values. For instance, empirical work found that for English $z = 0.97\pm0.06$ and for Russian $z = 0.89\pm0.07$, the difference being 8.3%. Two properties of these languages might explain the difference: Russian is a highly inflective language (i.e., a language that changes the form or ending of some words according to factors such as the genders, noun cases, verb conjugations, verb tenses, persons, moods, voices, aspects, numbers) while English is analytical (i.e., a language that uses specific grammatical words, or particles, and word orders rather than inflections). Second, it is well known that Russian lexical richness is greater than English.

any word			nouns	
Frequency	Token	Frequency	Token	
1,698,599	the	124,598	European	
849,256	of	104,325	\mathbf{Mr}	
793,731	to	92,195	Commission	
640,257	and	66,781	President	
508,560	in	62,867	Parliament	
407,638	that	57,804	Union	
400,467	is	53,683	report	
394,778	a	53,547	Council	
263,040	I	45,842	States	

The figure shows an example of word distribution from a document collection belonging to the European Parliament. It is immediately possible to observe which words have high frequency: the ones that are very generic in the language (e.g., articles) and the ones that are very specific for a document collection (e.g. nouns specific to the application domain).

Zipf's Law is handy in NLP because it can guide the choice of features for various NLP tasks. For instance, focusing on some of the most frequent words makes it possible to 1) capture the most important and relevant information (the right column in the slide) while 2) filter out less significant words or noise (the left column).

Zipf's Law can also inform the development of language models, as it provides a basis to estimate word probabilities. We will discuss language models in the following lecture.

4.3.3 Language Evolution

```
- | - |
LOL | Laugh out loud |
G2G | Got to go |
BFN | Bye for now |
B4N | Bye for now |
Idk | I don't know |
FWIW | For what it's worth |
LUWAMH | Love you with all my heart |
```

Finally, NLP is hard because language evolves. For instance, humans invent new terms (neologisms and acronyms) daily. These new terms could be the result of an optimization process (e.g. LOL (Laugh out loud) is said to have been invented in the 80s in the context of pre-Web digital chat rooms. Neologisms are so common that linguists started curating collections to help track the evolution of the language. For instance, the website Woordenboek van Nieuwe Woorden collects new words that appeared in the standard Dutch Language. Have you ever wondered when the expression *roast* (ridiculing someone for entertainment) entered the Dutch language? Look here. These terms become officially part of a

language when included in dictionaries. For instance, *LOL* has been officially included in the *Oxford English Dictionary* in 2011.



And humans change the way they communicate in even more fundamental ways. Think about using emojis to indicate emotional state: emojis have their meaning, but they can also add meaning, clarity, and credibility to text. Emojis can also be interpreted in different ways, depending on the cultural and socio-economic background of the sender or the receiver. But also depends on the device in use (Emoji characters vary slightly between platforms).

4.4 4.4 NLP Tasks

A product/service/system (or a designer) can use NLP techniques to process natural language text and draw valuable information from it.

NLP techniques could be used for:

- **Discovery**, i.e., to create new conceptualizations or ways to organize the world, to support making sense of it. The conceptualization helps simplify some of the highly complex elements of the environment or situation you are exploring to focus your attention on one or a few specific aspects. For example, imagine you are interested in studying online reviews. What aspects of these reviews are essential for your analysis? For example, you could be interested in exploring:
 - *Topical content*: What are the reviews about? The functionality of a product? Its aesthetics? Ergonomics?
 - Sentiment: Are users positive or negative about the product or one of its functionalities?
 - Credibility and informativeness: Are the review written and informative for your purpose, Or are they rant, probably created to mess with the reputation of a given vendor or product? The text may contain some other relevant aspects you were unaware of, which might be captured by a different concept not in your original conceptualization of the world.
- Measurement, i.e., the quantification of a given content in the data, to describe its prevalence in the real world. Given a concept, how prevalent is it in the data? Measurement is the essential ingredient for *description* as it provides valuable **summaries** of the data, which in turn may help characterize the state of the world, inform theories, or provide evidence for the success (or failure) of some design choices.
- **Prediction**, i.e., to make predictions about events in the future or the effect of an intervention based on the discovered concepts and measures. For example, a designer could be interested in estimating how a given framing of a problem (or solution) could affect users' perception of (or
engagement with) a product or a service. This is a predictive question because it uses today's information to help understand what will happen tomorrow.

Note that these three activities are not to be conceptualised in isolation. As you know, it is very common (especially for designers) to discover new directions, theories, and measures as part of the design process. While this refinement and evolution process is supported by qualitative insights, (textual) data can also play an important role.



4.4.1 An example of NLP Process

The one in the figure [^13 Bird, S., Klein, E., and Loper, E. (2009). Natural Language Processing with Python. O'Reilly.] is an example of an NLP pipeline for a spoken dialog system, like the one you interact with an Apple (Siri), Amazon (Alexa), or Microsoft (Cortana) product. We introduce it because it summarises the many aspects of NLP that go into such products, touching both issues of natural language understanding and generation.

Along the top of the diagram, moving from left to right is a "pipeline" (a sequence) of some **language understanding** components. This pipeline processes speech input (e.g., "Siri, how is the weather going to be in the coming hour?") to produce a meaningful representation of the utterance (e.g., the request is about the weather one hour from now). The reverse pipeline of components converts concepts back to speech to **generate language** that answers the question (e.g., "It will be raining"). These components make up the dynamic aspects of the system. At the bottom of the figure (and at the far right) are some representative bodies of information: the repositories of language-related data that the system uses in its work, but also remote systems, used to draw other data, or to invoke some remote business function.

The pipeline addresses five fundamental steps of NLP:

- **Phonology**, i.e., how the basic language units of sound (*phonemes*) are organized in a language. In an NLP pipeline, phonology may involve tasks such as *speech recognition*, *text-to-speech conversion*, and *text-to-speech conversion*. We will not address them in this course.
- **Morphology**, i.e., how words are built up from smaller meaning-bearing units called morphemes. Two broad classes of morphemes can be distinguished: **stems** - the central morpheme of the

word, supplying the primary meaning; and **affixes** - adding "additional" meanings of various kinds. For example, the word *cats* consists of the morpheme *cat* (stem) and of the morpheme *-s* (affix).

- **Syntax**, i.e., how words are arranged together to form sentences. Syntax concerns how words are combined and ordered to create grammatically correct and coherent sentences. The syntax of a language (as described by its grammar) determines how specific types of words, like nouns, verbs, and adjectives, are organized in sentences.
- Semantics, i.e., how words, phrases, and sentences convey meaning in various contexts
- **Reasoning**, i.e., drawing conclusions and making inferences based on the information available in a text

4.4.2 Morphology

In an NLP pipeline, morphological analysis typically involves tasks such as *tokenization*, *stemming*, and *lemmatization*. It is possible to extract word forms, inflections, and derivations by breaking down words into their constituent morphemes. This can improve the efficiency (speed) and effectiveness (quality) of text processing.

4.4.2.1 Tokenisation

- Issues
 - Separators: punctuations
 - Exceptions: "m.p.h", "Ph.D"
 - Expansions: "we're" = "we are"
 - Multi-words expressions: "New York", "doghouse"

One of the first (and most common) steps in NLP processing is converting text into a standard form that is convenient for automatic processing - a task called **Text Normalisation**.

Tokenization is about breaking up a document into discrete words having their meaning.

Tokenization is a language-specific problem, and each language poses unique challenges.

English words are often separated from each other by *whitespace*, but whitespace is not always sufficient. For instance, words like *New York*, *White House* and *rock 'n' roll* must be treated as a single word despite containing spaces. Sometimes words must be separated: for instance *I'm* into the two words *I* and *am*. In these cases, more sophisticated approaches are needed.

Using whitespace for tokenization is an approach that works for most languages. However, in languages such as Chinese and Japanese, words are not separated by spaces; the reader infers words from the context within the sentence. A *word segmentation model* is used to split the characters into their constituent words for these languages.

4.4.2.2 Stop-word Removal

any	word		nouns				
Frequency	Token	Frequency	y Token				
1,698,599	the	124,598	B European				
849,256	of	104,32	5 Mr				
793,731	to	92,19	5 Commission				
640,257	and	66,78	1 President				
508,560	in	62,86	7 Parliament				
407,638	that	57,804	4 Union				
400,467	is	53,683	3 report				
394,778	a	53,54	7 Council				
263,040	I	45,842	2 States				

- Removal of high-frequency words, which carry less information
 - E.g., determiners, prepositions
- English stop list is about 200-300 terms (e.g., been, a, about, otherwise, the, etc..)

When introducing Zip's law, we discussed how some words in a document collection are bound to be very frequent, depending on the language or the document collection itself. In NLP parlance, **stop-words** are words used across documents that do not give much information about a specific NLP task. In English, for instance, these are common words such as *and*, *the*, and *that*. These words account for a significant fraction of the words, but they account for only a tiny fraction of the meaning. By removing them, it is possible to reduce the size and complexity of the feature set that needs to be handled computationally.

4.4.2.3 Lemmatisation

- Technique using dictionaries and morphological analysis of words to return the base or dictionary form of a word
 - Example: Lemmatization of *saw* —> attempts to return *see* or *saw* depending on whether the use of the token is a *verb* or a *noun*

```
Google , headquartered in Mountain View ( 1600 Amphitheatre Pkwy , Mountain View ,

headquarter

Sundar Pichal said in his keynote that users love their new Android phones .
```

Another part of text normalization is **lemmatization**, i.e. the task of determining that two words have the same root, despite their surface differences. For example, *family* and *families*; or *sang*, *sung*, and *sings* (different forms of the verb *sing*). They are distinct terms, but in some tasks, it may be effective to map them all to a common form.

A **lemma** is the canonical form (such as one might find in a dictionary) of a set of words that are related by inflection (i.e., modifications due to case, number, tense, etc.). Lemmatization is the process of mapping words to their lemma. Sometimes identifying the lemma is relatively simple (e.g. *family* and *families*) but it can also be quite complex (e.g. *sing*, *sang*, *sung*).

Lemmatization can be cumbersome because it often requires a *dictionary lookup* (to map *sung* to *sing*) and, sometimes, syntactic analysis to identify the role played by the word in the sentence (e.g. *saw* as a noun, or *saw* as the past tense of the verb *see*).

4.4.2.4 Stemming

- *Sample text:* Such an analysis can reveal features that are not easily visible from the variations in the individual genes and can lead to a picture of expression that is more biologically transparent and accessible to interpretation
- *Lovins stemmer:* such an analys can reve featur that ar not eas vis from th vari in th individu gen and can lead to a pictur of expres that is mor biolog transpar and acces to interpres
- **Porter stemmer:** such an analysi can reveal featur that ar not easili visibl from the variat in the individu gene and can lead to a pictur of express that is more biolog transpar and access to interpret
- *Paice stemmer:* such an analys can rev feat that are not easy vis from the vary in the individ gen and can lead to a pict of express that is mor biolog transp and access to interpret
 - Heuristic process that chops off the ends of words in the hope of achieving the goal correctly most of the time
 - Stemming collapses derivationally related words
 - Two basic types:
 - Algorithmic: uses programs to determine related words
 - Dictionary-based: uses lists of related words

A popular approximation to lemmatizing that also maps related forms together is **stemming**.

Stemming uses simple algorithms (and sometimes dictionaries too) to discard the end of a word. In the case of *family*, all the variants would be mapped to *famili*.

Stemming is often quite effective and substantially faster than lemmatization, but it produces words that are not actual in the language. Also, it fails to deal with words with more complex forms. Stemming can also sometimes reduce two words to a common stem with distinct meanings, such as *secure* and *securities*.

4.4.3 Syntax

By understanding the syntactic structure of a sentence, NLP systems can decipher the relationships between words and phrases, which is crucial for various NLP tasks and applications.

4.4.3.1 Part-of-speech Tagging

Tagging each word in a sentence with a corresponding *part-of-speech* (e.g. noun, verb, adverbs)

nsubj		р	v	mod	E.	orep	nn		po	bj	р	num	1	nn		app	os p	
Google		,	heado	quarte	red	in M	oun	tain	Vie	w	(160	0 A	mphith	neatre	Pkv	vy,	
NOUN	PU	NCT	N	/ERB	1	ADP	NOU	N	NO	JN I	PUNCT	NUN	1	NOU	N	NOU	JN PUN	СТ
nn		appo	s p	a	ppos	num		р		P	roo	t	det	amod	nn		dobj	prep
Mounta	ain	View	ι,		CA	94043	D)		,	unvei	led	the	new	Andr	oid	phone	for
NOUN	L	NOU	N PUN	ICT N	IOUN	NUM	P	UNC'	T Pl	JNCT	VER	В	DET	ADJ	NOU	JN	NOUN	ADP
pobi	prep	det		nn		nn	DC	bi	p									
¢700	a t	the	Con		r Elou	otronio	ch	~										
2/99	at	the	Con	sume	r Eleo	ctronic	Sh	ow	•									
NUM	ADP	DET	N	OUN	N	IOUN	NO	UN	PUN	СТ								
nn	ns	subj	root	prep	poss	ро	bj	mai	rk n	isubj	ccomp	po	oss	amod	nn		dobj	
Sundar	Pie	chai	said	in	his	keyn	ote	tha	nt u	sers	love	th	eir	new	Andro	id j	phones	
NOUN	N	NUC	VERB	ADP	PRON	NOI	JN	AD	P N	IOUN	VERB	PR	ON	ADJ	NOUN	4	NOUN	

Source: https://cloud.google.com/natural-language#section-2

A part of speech (POS) is a category of words with similar grammatical properties. In English, for example, **nouns** describe the *names of objects, animals, people*, and *concepts*, among many other things. **Verbs** describe *actions, states*, and *occurrences*. A noun can be used as a **subject**, or as an **object** of a verb. **Adjectives** describe properties (e.g. colour, age, visual quality) of nouns. **Adverb** modify verbs by indicating *time, place* or *manner* of a given action. The English language also has other classes of POS, including **Pronouns** (i.e., a shorthand for referring to an entity or event - e.g. *she, I*); **particles** (i.e. preposition-like form used together with a verb - e.g. *up, down, on, off*), **Punctuation**, **Numeral** and so on.

Parts of speech tagging is the task of assigning a grammatical tag to each word within a document. Words are tagged either by a rule-based algorithm or by a machine learning algorithm that classifies each word based on a manually tagged training dataset.

POS tagging is an essential input to several other techniques and can be used to focus text processing on properties of interest. For example, if one is primarily interested in sentiment content, it may be

advantageous to remove all parts of speech except adjectives and adverbs, which are more likely to convey cues about sentiment. POS tagging can also be used to identify phrases (or *n-grams*) that can be included in feature sets or to visualize texts. For example, it can be used to identify **noun phrases** (e.g. *Health Care*).

4.4.3.2 Language Analysis

- Idea: people's language can provide *insights into their psychological states* (e.g. emotions, thinking style)
- For instance
 - Frequency of words associated with positive or negative emotions
 - Use of pronouns as a proxy for confidence and character traits



It is a common intuition among scholars (but also laypeople) that how we use language (i.e., word choices or figures of speech used; sentence structure; and register or tone) gives out subtle clues about what the people uttering (or writing) they are like, psychologically speaking.

Language analysis explores how authors/writers/speakers convey meaning through specific language analysis techniques.

Language analysis attempts to infer a person's psychology by mapping out the content of what they said: for example, charting a user's trajectory by counting how often the concepts of *excitement* or *frustration* are mentioned.

Psychologists conducted a lot of research into how "particles" of language like *pronouns* (I, you, we), *articles* (the, a, an), and *negations* (no, not, never) could provide compelling insights into a multitude of psychosocial phenomena. For instance, they can be reliable indicators of a person's thoughts, regardless of their discussion. In some ways, the psychological significance of a text could be discovered by counting **how often** different meanings are conveyed. The *word counting* approach relies on scanning texts and counting the frequency of words from predefined categories—categories that are informed by psychological theory: emotional words (*happy, upset, angry*), agentic words (*do, able, try*), thinking words (*think, understand, guess*), and so on. Word frequencies represent **attentional habits** [^15 Natural Language Analysis and the Psychology of Verbal Behavior: The Past, Present, and Future

States of the Field. Ryan L. Boyd and H. Andrew Schwartz. Journal of Language and Social Psychology 2021, Vol. 40(1) 21–41].

The relative frequencies of each category are then interpreted as reflecting a person's relative focus on each domain. For example, people who use high rates of articles (the, a, an) and prepositions (next, above) tend to focus on formal or concrete concepts and their inter-relations. People with higher social status and confidence are more focused on the external social environment than themselves, using more "you" and "royal we" words than "I" words. Types of words (e.g. pronouns) are *markers* of what (or who) we pay attention to. "Emotion" words are not reflections of the experience of emotions; they are merely diagnostic of one's attention to affective states.

I recommend the book in the slide - Pennebaker, J. W. (2011). The secret life of pronouns: What our words say about us.



The tool **Linguistic Inquiry and Word Count**, or **LIWC** (Pennebaker & Francis, 1999), allows these types of analysis.

			Psychological Processes		
			Drives	Drives	we, our, work, us
			Affiliation	affiliation	we, our, us, help
			Achievement	achieve	work, better, best, working
			Power	power	own, order, allow, power
		Description/Most frequently	Cognition	Cognition	is, was, but, are
Category	Abbrev.	used exemplars	All-or-none	allnone	all, no, never, always
Summary Variables			Cognitive processes	cogproc	but, not, if, or, know
Word count	WC	Total word count	Insight	insight	know, how, think, feel
Analytical thinking	Analytic	Metric of logical formal thinking	Causation	cause	how, because, make, why
Clout	Clout	Language of leadership, status	Discrepancy	discrep	would, can, want, could
Authentic	Authentic	Parcaived honasty ganuinanase	Tentative	tentat	if, or, any, something
Emotional tone	Tone	Degree or positive (negative) tone	Certitude	certitude	really, actually, of course, real
Words per sentence	WPS	Average words per sentence	Differentiation	differ	but, not, if, or
Rig words	RigWords	Parcent words 7 latters or longer	Memory	memory	remember, forget, remind, forgot
Dictionary words	Dic	Percent words contured by LIWC	Affect	Affect	good, well, new, love
Linguistic Dimensions	Linquistic	Tereent words captaled by ETWC	Positive tone	tone pos	good, well, new, love
Total function words	function	the to and I	Negative tone	tone neg	bad, wrong, too much, hate
Total propouns	pronoun	I you that it	Emotion	emotion	good, love, happy, hope
Personal pronoune	pronoun	I you my me	Positive emotion	emo pos	good, love, happy, hope
l st person singular	i	I me my myself	Negative emotion	emo neg	bad, hate, hurt, tired
1st person plural	we	we our us lets	Anxiety	emo anx	worry, fear, afraid, nervous
2nd person	vou	you your u yourself	Anger	emo anger	hate, mad, angry, frustr*
3rd person singular	shehe	he she her his	Sadness	emo sad	:(, sad, disappoint*, cry
3rd person plural	they	they their them themsel*	Swear words	swear	shit, fuckin*, fuck, damn
Impersonal pronouns	inton	that it this what	Social processes	Social	you, we, he, she
Determiners	det	the at that my	Social behavior	socbehav	said, love, say, care
Articles	article	a an the alot	Prosocial behavior	prosocial	care, help, thank, please
Numbers	number	one two first once	Politeness	polite	thank, please, thanks, good morning
Prenositions	prep	to of in for	Interpersonal conflict	conflict	fight, kill, killed, attack
Auxiliary verbs	auxyerb	is was be have	Moralization	moral	wrong, honor*, deserv*, judge
Adverbs	adverb	so just about there	Communication	comm	said, say, tell, thank*
Conjunctions	coni	and but so as	Social referents	socrefs	you, we, he, she
Negations	negate	not no, never, nothing	Family	family	parent*, mother*, father*, baby
Common verbs	verb	is, was, be, have	Friends	friend	friend*, boyfriend*, girlfriend*, dude
Common adjectives	adi	more, very, other, new	Female references	female	she, her, girl, woman
Quantities	quantity	all, one, more, some	Male references	male	he, his, him, man

This is a list of all the categories supported by the LIWC-22 Language Dimensions. I recommend reading through the manual, also to get a feeling of how reliable some of these measurements are.

Here I report four summary variables that are calculated in LIWC. The description of the variable and the supporting literature is available here.

- **Analytic Thinking**: the degree to which people use words that suggest formal, logical, and hierarchical thinking patterns.
 - low Analytical Thinking —> language that is more intuitive and personal
- **Clout**: the relative social status, confidence, or leadership that people display through their writing or talking
- Authenticity: the degree to which a person is self-monitoring
 - Low authenticity: prepared texts (i.e., speeches written ahead of time) and texts where a person is being socially cautious
- **Emotional tone**: the higher the number, the more positive the tone. Numbers below 50 suggest a more negative emotional tone.

4.4.3.3 Sentiment Analysis

- The detection of attitudes, affectively colored beliefs, dispositions towards objects or persons"
- Main elements
 - Holder (source)
 - Target (aspect)
 - Type of attitude
 - Text containing the attitude
- Tasks
 - *Classification*: Is the text's attitude positive or negative?
 - Regression: Rank the attitude of the text from 1 to 5
 - Advanced: Detect the target, source, or complex attitude types

Sentiment analysis is a text analytic technique that automatically identifies and categorizes subjective information within the text.

Specifically, we are interested in *affective meaning*, which analyzes the writer's evaluations, opinions, emotions, and speculations.

A common typology of affective states comes from Scherer (2000):

- **Emotion**: A relatively brief episode of response to the evaluation of an external or internal event as being of significant significance (angry, sad, joyful, fearful, ashamed, proud, elated, desperate)
- **Mood**: Diffuse affect state, most pronounced as a change in subjective feeling, of low intensity but relatively long duration, often without apparent cause. (cheerful, gloomy, irritable, listless, depressed, buoyant)

- **Interpersonal stance**: Affective stance taken toward another person in a specific interaction, coloring the interpersonal exchange in that situation. (distant, cold, warm, supportive, contemptuous, friendly)
- **Attitude**: Relatively enduring, affectively colored beliefs, preferences, and predispositions towards objects or persons. (liking, loving, hating, valuing, desiring)
- **Personality traits**: Emotionally laden, stable personality dispositions and behavior tendencies, typical for a person. (nervous, anxious, reckless, morose, hostile, jealous)

Sentiment Analysis is a technique used to quantify **attitudes** towards specific topics or entities that are written in an unstructured way and, thus, hard to quantify otherwise.

One of the most basic tasks in sentiment analysis is the classification of **polarity**, i.e., to classify whether the expressed opinion is *positive*, *negative*, or *neutral*. It is, of course, possible to use more than three classes, e.g., *strongly positive*, *positive*, *neutral*, *negative*, or *strongly negative*, or to deal with the problem as a regression problem.

Sentiment analysis can be applied to various textual resources such as surveys, reviews, and social media posts.



These are examples of sentiment analysis evaluated against a standard text. The first two screenshots come from a Google demonstrator, while the second is from an IBM one. You can notice how sentiment is typically measured against a whole sentence. However, it is also possible (typically by combining with POS tagging and NER) to associate a sentiment to a specific entity in the text



4.4.3.4 Emotion Analysis



Plutchik wheel of emotion

While sentiment analysis is framed in terms of positive and negative categories, psychologists generally regard emotion as more multifaceted. Detecting emotion has the potential to improve several language processing tasks. For instance, they automatically detect emotions in reviews or customer responses (anger, dissatisfaction, trust). Emotion can play a role in medical NLP tasks like helping diagnose depression or suicidal intent. Detecting emotions expressed toward characters in novels might play a role in understanding how different social groups were viewed by society at different times.

Perhaps the most well-known of this family of theories are the 6 emotions proposed by Ekman (1992) - happiness, surprise, fear, sadness, anger, and contempt. He argues that these six basic emotions are universal across human cultures. Another theory is the Plutchik (1980) wheel of emotion, consisting of 8 basic emotions in four opposing pairs: joy–sadness, anger–fear, trust–disgust, and anticipation–surprise, and the emotions derived from them. These are called *atomic theories* because they assume a limited number of basic emotions (6 or 8) from which others are generated.

Another class of emotion theories widely used in NLP considers emotion as a space in 2 or 3 dimensions (Russell, 1980):

- valence: the pleasantness of the stimulus
- arousal: the intensity of emotion provoked by the stimulus
- dominance: the degree of control exerted by the stimulus

In this model, sentiment can be viewed as a special case of this second view of emotions as points in space. In particular, the valence dimension, measuring how pleasant or unpleasant a word is, is often used directly to measure sentiment.

	Full Document						
	Sadness	25.93%					
	Joy	81.39%					
	Fear	1.38%					
	Disgust	1.77%					
	Anger	3.02%					
	Entity Emotio	n Scores					
Google, headquartered in Mountain View (1600 Amphitheatre Pkwy, Mountain View, CA 940430), unveiled the new	Mountain View (1600 Amphitheatre Pkwy						
Android phone for \$799 at the Consumer Electronic Show. Sundar Pichai said in his keynote that users love their	Sadness	40.97%					
new Android phones .	Joy	67.24%					
	Fear	1.37%					
🔳 Sadness 📕 Fear 🔲 Disgust 🔲 Anger 🔲 Joy							

The slide shows an example of emotion recognition system from IBM. It detects anger, disgust, fear, joy, or sadness that is conveyed in the content or by the context around target phrases.

4.4.4 Semantics

Semantic analysis determines the meaning of words, sentences, and documents beyond their lexical organization. Semantic analysis includes *lexical semantics*, which is concerned with the individual words' meanings, extending to the relationships between individual words in sentences and across sentences.

4.4.4.1 Named Entity Recognition

- Factual information and knowledge are usually expressed by named entities
 - Who, Whom, Where, When, Which, ...
- Identify words that refer to proper names of interest in a particular application
 - E.g. people, companies, locations, dates, product names, prices, etc.
- Classify them to the corresponding classes (e.g. person, location)
- Assign a unique identifier from a database

Part of speech tagging allows inferring that words like *Alessandro*, *Delft University of Technology*, and *Zuid Holland* are all proper nouns. From a semantic perspective, these proper nouns refer to different kinds of entities: *Alessandro* is a **Person**, *Delft University of Technology* is an **Organisation**, and *Zuid Holland* is a **Location**.

A **named entity** is, roughly speaking, anything that can be referred to with a proper name: a *person*, a *location*, or an *organization*. Four entity tags are most common: **PER** (person), **LOC** (location), **ORG** (organization), or **GPE** (geo-political entity). Other information related to dates, times (and other temporal expressions), and even numerical expressions (e.g., prices) is also tagged as entities, although they are not entities per se.

The task of **named entity recognition** (**NER**) is to tag a word (or multiple contiguous words) with the entity type.

Like part of speech tagging, modern named entity recognition systems tend to be built with machine learning approaches. In some specialized contexts, these systems can be trained to tag additional information, such as symptoms, drug names, or drug reactions.

Google	ORGANIZATION	2. Mountain View	LOCATIO
Vikipedia Article		Wikipedia Article	
alience: 0.19		Salience: 0.18	
Android	CONSUMER GOOD	4. Sundar Pichai	PERSO
Vikipedia Article		Wikipedia Article	
alience: 0.14		Salience: 0.11	
phone	CONSUMER GOOD	6. users	PERSO
alience: 0.10		Salience: 0.09	
Amphitheatre Pkwy	LOCATION	8. CA 940430	ОТНЕ
alience: 0.07		Salience: 0.05	
. keynote	OTHER	10. phones	CONSUMER GOO
aliance: 0.02		Salience: 0.02	

Source: https://cloud.google.com/natural-language#section-2

This is an example of NER system in action.

4.4.4.2 Document Categorisation

- Assigning a label or category to an entire text or document
- Supervised learning
- For instance
 - Spam vs. Not spam
 - Language identification
 - Authors attribution
 - Assigning a library subject category or topic label

Many NLP tasks involve *classification*, i.e., assigning a label or category to an entire text or document.

The simplest version of document classification is **binary classification tasks**. For instance, the classification emails as *spam* or *not spam*. Other tasks are multi-classes (or multi-label), e.g., language identification, or **categorization**, i.e., deciding which categories a document belongs to (in a given taxonomy). These are *supervised* machine-learning tasks similar to the ones we have encountered before in the course.

CATEGORIES

- 0.85 science and technology
- 0.58 education
- 0.58 economy, business and finance>economic sector>computing and information technology
- 0.57 society
- 0.54 science and technology>social sciences>psychology
- 0.54 economy, business and finance>economic sector>media
- 0.54 society>values>ethics
- 0.49 education>school>further education
- 0.43 economy, business and finance>economic sector>computing and information technology>software
- 0.43 science and technology>social sciences>philosophy

ML4D Course Description

4.4.4.3 Topic Modeling

- A topic is the subject or theme of a discourse
- Topic modeling: group documents/text according to their (semantic) similarity
- An unsupervised machine learning approach



ML4D Course Description

Topic modeling is an NLP approach aimed at automatically discovering abstract *topics* that occur in a collection of documents. Topic modelling is a typical *unsupervised machine learning* task.

Topic modeling operates under the assumption that if a document is about a particular topic, then some words are more or less likely to appear in the document. As documents typically concern multiple topics in different proportions, topic modeling techniques also use the relative proportion of words in a document to estimate similarly distributed and co-present words. The *topics* produced by topic modeling techniques are clusters of similar words.

These type of tasks are typically explorative: a researcher sets some number of topics, runs the topic modeling algorithm, checks the nature of the topics outputted by reading the words and documents identified as having high probabilities of belonging to each of the topics, and decides whether or not those topics are substantively meaningful. If that is not the case, then the research can vary the number of topics and iterate.

4.4.4.4 Word Sense Disambiguation

• Multiple words can be spelled the same way (homonymy)

- The same word can also have different, related senses (polysemy)
- Disambiguation depends on context!



Word sense disambiguation (WSD) is the task of determining the correct meaning of a word with multiple meanings, depending on its context. WSD algorithms use contextual information, such as surrounding words or phrases, to identify the most appropriate sense of an ambiguous word.

The slide shows an example of a WSD tool.

4.4.4.5 Automated Summarisation

- Condensing a piece of text to a shorter version while preserving key informational elements and the meaning of content
- A challenging task!

Text Summarization Result

Original URL/Text

Summarized Text

Original URL/Text
NoB-HT Machine Learning for Design is a technology elective methoded in the Arg var of the Bachelor of Industrial Design Engineering at the Deft University of Technology. The course provides tudents with the knowledge requires (PSS). Machine learning for Design is a technology elective methoded in the Arg var of the Bachelor of Industrial Design Pathematics, services, and systems in the context of the design of intelligent products, services, and systems in the context of the design of intelligent products, services, and systems in the context of the design of Intelligent products, services, and systems platforms, context Creation platforms, personal health appliances runch of Current and Huns PSSs are powered by ML technology; influencing, and staping our interests, habits, lease, and tailures of ML technology influencing, and staping our interests, habits, lease, and solicity. To meaningfully evina with the PSSs at are beneficial and useful to poope and society, designers must: engage with the detain the incorring and staping our interests, habits, lease, and tailures of ML technology; influencing, and staping our interests, habits, lease of ML technology; influencing, and staping our interests, habits, lease of ML technology; influencing, and staping our interests, habits, lease of ML technology; influencing, and staping our interests, habits, lease of ML technology; influencing, and staping our interests, habits, lease of ML technology; influencing, and staping our interests, and tailures of ML technology; influencing and staping our interests, and tailures of ML technology; our our divert with the stape interest of ML technology; and all more the quick, bases, and tailures of ML technology; and technology of halt ML spectem can do, and how they could and should be integrated in IPSSs.

https://textsummarization.net/

Result

After pressing the "Summarize" button above, the result will be displayed in the box below

The summated text will be here. IOB4-T3 Machine Learning for Design is a technology optional embedded in the 2nd year of the Bachelor of Industrial Design Engineering at the Delft University of Technology, Machine learning is a computational approach that focuses on "offering computer systems the capacity to learn without being explicitly configured". Students in this course gain useful experience w ML innovation and learn just how to think seriously of what ML systems can do, and just how they could and should be integrated in iPSSs.

4.4.4.6 Machine Translation (popular languages)



Machine translation is the process of translating given text from one language to another language. The language the input text is written in is called the source language, whereas the one for the output is called the target language. One challenge in MT is the tradeoff between *fluency* and *adequacy*. The translation must be fluent, meaning the output must sound natural in the target language. Translation also needs to be adequate, meaning that the output has to reflect the meaning expressed by the input as closely as possible. These two are often in conflict, especially when the source and the target languages are not very similar (e.g., English and Chinese). A translation can be a precise, verbatim mapping of the input, but the result will likely not sound natural in the target language. On the other hand, it is possible to create an output that sounds natural in the target language but does not reflect the precise meaning. Good human translators creatively address this tradeoff.

JAGE ENGLISH TURKISH DUTCH 🗸 I study machine learning for design × Men dizayn uchun mashinani oʻrganishni ☆ organaman 25/5000 ■ * 40 4 10 Men dizayn uchun mashinani oʻrganishni × oʻrganyapman 💿 × I am learning to machine for design ف الا التحقيق العامة المحقيق المحقق المحقيق ال المحقيق ال المحقيق المحقي +-* ENGLISH UZBEK DUTCH imes Men dizayn uchun mashinani oʻrganyapman \dot{x} 35/5000 💼 * 🖷 егтет Lakeuwse uzвех викцан тиккан v Men dizayn uchun mashinani oʻrganyapman v +²⁴ ENGLISH UZBEK DUTCH ↓ × I am studying the machine for design 5 ¢9 <

4.4.4.7 Machine Translation (languages with fewer resources)

Current translation systems are built on top of ML techniques provided with examples of sentences in both the source and target language. However, such examples are not always readily available, especially for less popular languages or languages with abundant resources.

4.4.4.8 Natural Language Instructions / Dialog systems



Dialog systems allow humans to interact with computers through natural language conversation, textual or vocal. The field of dialog systems has a long history. One of the earliest dialog systems, ELIZA, was developed in 1966.

The two main types of dialog systems are **task-oriented** and **chatbots**.

- **Task-oriented** dialog systems are used to achieve specific goals. For example, commanding a car entertainment system and interacting with an Alexa-like devise to obtain some information. As described in the previous slides, task-oriented dialog systems are usually built around an NLP pipeline. These types of dialog systems are not designed to support long and complex conversations, as the goal is typically to execute a given task in the shortest possible time. Most interactions with Alexa (or Siri) only include 1 or 2 exchanges.
- **Chatbots** dialog systems are designed to have conversations that might not have a specific goal but are more extended and interactive. Handwritten rules usually manage conversations in traditional chatbots (e.g., when the human says this, say that).

For both types of dialog systems, recent advances in neural networks (and generative models) allow for richer, longer, and task-specific conversation. However, it will still take some time (perhaps a couple of years) to see these advanced conversation systems operate in consumer products, as their operation is computationally (and energy) intensive.

4.4.4.9 Natural Language Generation



Natural Language Generation is the process of generating natural language text from something else. Generative models like ChatGPT are examples of end-to-end NLG systems. However, remember that language can also be generated according to pre-defined patterns and templates.



4.4.5 State of the Art in NLP - as of 2022

Credits: Nava Tintarev

I kept this slide - and I will probably keep it in the future - as a memento of how the progress in the field of NLP is accelerating at an incredible pace. Many of the tasks listed here were still considered difficult (or hard) in 2022 but are now becoming less and less challenging due to the emergence of large language models like GPT-3 and GPT-4.

4.5 4.5 Features in Natural Language Processing

In previous sections we discussed how *feature selection* is an essential pre-condition for successful ML models. We have seen how features could be canonically represented using a table, where rows are documents and columns are *numerical features*. We have discussed how in computer vision *pixels* can (and are) used as features, as each pixel contains one or more numerical values. But what about textual documents? What constitutes a feature in Natural Language Processing?

- A sequence of alphanumerical characters
 - Short: e.g. tweets
 - Long: e.g Web documents, interview transcripts
- Features are (set of) words
 - Words are also syntactically and semantically organised
- Feature values are (sets of) words occurrences
- Dimensionality —> at least dictionary size

Textual documents can be represented in multiple ways. One approach, common in the so-called *symbolic approaches*, is to treat words as symbols organized in a so-called syntax tree. Another approach

called *bag-of-words* or *n-gram*, common in machine-learning systems, is to represent text as *sparse structured data* by counting in documents the occurrence of individual words, pairs of words, triplets of words, etc.



For example, consider the funny amazon review from the left. It is possible to conceptualize a document collection as a big table where rows are documents, and features (columns) are terms in the collection vocabulary; that is, every single word is a potential feature for a document. The presence or absence of a term in a document is represented by a 0 or 1 value. This representation is also called *term-document matrix*.

4.5.1 Main types of NLP Tasks

- Label (classify) a region of text
 - e.g. part-of-speech tagging, sentiment classification, or named-entity recognition
- Link two or more regions of text
 - e.g. coreference
 - * are two mentions of a real-world thing (e.g. a person, place) in fact referencing the same real-world thing?
- Fill in missing information (missing words) based on context

As we already discussed in the previous lecture, NLP tasks seek to do one of three things:

- **label** a region of text: this is the case for tasks such as *part-of-speech tagging*, *sentiment classification*, or *named-entity recognition*;
- **link** two or more regions of text: this is common in *coreference* tasks, where the goal is to assess whether to named entities in a document refer to the same real-world entity (e.g., same person, place, or some other named entity) -**fill in** missing information (missing words): this is the case for generative models (e.g., GPT) trained on the task of **predicting the most likely next word* based on context.

In the following, we will learn what Language Models are, how they work, and how they can be used.

4.6 4.6 Language Models

Language: vocabulary and its usage in a specific context captured by textual data

Let us further explore how language could be represented in machine learning systems. An important disclaimer: from now on, by *language* we mean the **vocabulary** (collection of words) and its usage in a specific context captured by textual documents. This is to distinguish our interpretation of the word *language* from the commonly used interpretation - to refer, for instance, to the *Dutch* language.

4.6.1 What is a *language model==**?

- A collection of statistics learned over a particular language
- Almost always empirically derived from a text corpora

"Language Model": two words that have become very popular recently due to the success of generative models like GPT. But what is a language model? Technically, a language model is a statistical model that gives a value (probability) to the likelihood for a piece of text to appear in a sentence.

Consider you are given the following sentence: "My weekend in Lax has been ruined by bad *____." What words are most likely to come next? There could be many reasons for a fantastic winter sports weekend to be ruined: rude people, traffic issues, or bad food. Of course, words like food* and people can appear in this content, but we can agree that the most likely word to appear in the sentence is "weather."

We just estimated the probability of occurrence for several words in this English sentence. We compared these words based on our personal experience with language and life. We tapped into our *repository* of memories and experiences to calculate the statistical probability of several words and picked the most likely one.

This is, in essence, the nature of a language model: a collection of statistics learned (mostly empirically - i.e., through observations) on a text corpora. By being empirically derived, language models are not necessarily able to "generalize" their beliefs beyond the observed text collection. A language model trained on a collection of English novels might not have statistics for Italian words. Likewise, it would assign higher probabilities to sequences of words (sentences) that "make more sense," that are more grammatical in English.

4.6.2 What are language models used for?



- **Measure** how *important* (or descriptive) a word is in a given document collection e.g., find the set of words that best describe multiple clusters (see Assignment 2)
- **Predict** how *likely* a sequence of words is to occur in a given context e.g., find the words that are more likely to occur next

What is the utility of keeping statistics about the likelihood of a word appearing in the given context? It turns out this is an essential feature in many applications of NLP. For instance:

- To **measure** the importance of words in a document collection
- To **predict** the likelihood of a sequence of words appearing. This is useful in many applications, such as:
 - **machine translation**, converting a text from a source to a target language.
 - **speech recognition**, converting from an audio signal to text the presence of noise makes it very useful to know which word is most likely next.
 - Text **summarization**, converting a long text into a short one.
 - Text **generation**, creating new text from an initial input.

4.6.3 What is the issue with word representation?

A previous slide showed how textual content could be represented in feature space using a simple *term-document* matrix. That is a representation that is very useful in several contexts. But let us know explore more in details the issue with the representation of words.

- Words are discrete symbols
 - Machine-learning algorithms cannot process symbolic information as it is
 - We need to transform the text into **numbers**
- But we also need a way to express relationships between words!



The first and arguably most important common denominator across all NLP tasks is how we represent words as input to any of our models. To perform well on most NLP tasks, two requirements must be satisfied:

- Words (discrete symbols) must be represented in a way a computer can process typically numbers. Neural networks are pure mathematical computation models that can only deal with numbers. They can't do symbolic operations such as concatenating two strings and conjugating a verb to past tense unless they are all represented by numbers and arithmetic operations.
- There must be some notion of **similarity** and **difference** between words to allow comparison and, possibly, some other form of calculation.

4.6.3.1 A simple representation approach

- Assign an incremental number to each word
 - cat = 1
- dog = 2
- pizza = 3
- Problem: there is no notion of similarity
 - Is a *cat* as semantically close (similar) to a *dog* as a *dog* is to a *pizza*
 - Also, no arithmetic operations
 - * Does it make sense to calculate dog cat to establish similarity?

A possible approach to numerically represent all words in a vocabulary is to assign an incremental number to individual terms, as in the slide. While numbers now represent words, it does not mean that it is possible to perform arithmetic operations on them. These numbers are discrete and arbitrary (for instance, numbers can be assigned incrementally based on the alphabetical order of words). Clearly, it makes no sense to infer that "cat" is equally similar to "dog" (difference between 1 and 2), as "dog" is to "pizza" (difference between 2 and 3).

4.6.4 Word Embeddings

- Embed (represent) words in a numerical n-dimensional space
- Essential for using machine learning approaches to solve NLP tasks
 - They bridge the symbolic (discrete) world of words with the numerical (continuous) world of machine learning models

What if words were to be represented on a numerical scale? Not a simple discrete list of numbers (integers) but a potentially n-dimensional space of continuous values. An array (a vector) of *float* numbers.

A **word embedding** is a continuous vector representation of a word. Word embeddings are a way to represent each word with arrays of arbitrary size filled with nonzero float numbers. The name *embedding* derives from each discrete word being "embedded" in a continuous vector space.

Of course, a question immediately arises: how large should this vector be?

4.6.4.1 Approach 1

• Assign numbers to words, and put semantically related words close to each other



- We can now express that dog is more related to cat than to pizza

• But is *pizza* more related to *dog* than to *cat*?

The slide shows a 1-dimensional vector. It is now possible to represent that "cat" and "dog" are more similar to each other than to "pizza." Still, "pizza" is slightly closer to "dog" than it is to "cat," a fact that does not register with our intuition. What if you wanted to place "pizza" somewhere that is equally far from "cat" and "dog?" A single dimension is not possible, so perhaps we need more.

4.6.4.2 Approach 2

• Assign multiple numbers (a vector) to words



- e.g. Euclidean, or Cosine (angles)
- But what is the meaning of an axis?

In this representation, we achieve that goal. Using a 2-dimensional space and these specific values, we reach the purpose of positioning "cat" and "dog" close to each other but also equally distant to "pizza." This is better, but other questions immediately arise. For instance, how many axes should there be? And how to calculate the numerical value of each word? An arbitrary assignment of numbers to terms doesn't capture any relationships among words.

4.6.5 One-Hot Encoding

- Each word in the vocabulary is represented by a one-bit position in a HUGE (sparse) vector
 - Vector dimension = size of the dictionary
 - * There are an estimated 13 million tokens for the English language

$$cat = [0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, \dots, 0]$$

 $dog = [0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, \dots, 0]$

A simple method exists to "embed" words into a multidimensional space: creating one dimension for each word in the vocabulary. This is called a **one-hot encoding**: each vector has only one 1 at the position corresponding to the word's index. These vectors are still *word embeddings* not very useful in representing the semantic relationship between words. Still, they are used as input to a machine learning algorithm, at least when no other embedding representation is available.

- Problems with one-hot encoding:
 - The size of the vector can be huge
 - * Do you Remember Zip's law?
 - * Easy to reach 10^6 words
 - * But we can use stemming, lemmatisation, etc
- Still, no notion of similarity or words relationship
 - Each word is an independent, discrete entity



There are some other issues associated with one-hot encoding:

- All words are equally distant from each other, which means that each word is treated as an *independent and discrete* entities
- The vocabulary size can be enormous. Even after applying techniques such as stop-word removal and stemming, the size of these vectors could be computationally prohibitive.

Despite its simplicity, one-hot encoding is a popular approach to represent discrete symbols (like words) numerically.

4.6.6 Independent and identically distributed words assumption

- The simplest language models assume that each word in a text **appears independently** of the others
 - The text is modeled as generated by a sequence of independent events
- The **probability** of a *word* can be estimated as the *number of times* a word appears in a text corpus
- But high probability **does not mean important** (or descriptive)

We need to make an essential yet technical clarification. When working with statistical approaches, including machine learning ones, we often operate according to the so-called *Independent and identically distributed* (IID) assumption. In NLP, the IID assumption assumes that each word (or sentence, depending on your prediction unit) is **independent** of the others and is drawn from **the same under-lying probability distribution**. In practice, this means that the occurrence of one word or sentence in a document does not depend on the occurrence of any other term or sentence in that document; moreover, the assumption implies that all words (or sentences) are *equally likely to be generated* from the underlying probability distribution.

The IID assumption is used to "ease the math," that is, to simplify the mathematical analysis of the learning algorithms and to use some statistical tools (e.g., some classes of ML models) that rely on such assumptions for their use. This assumption, clearly, does not hold in all real-world scenarios. For instance, in language, there are dependencies between words or sentences due to the context in which they are used: the word "the" always appears before the word "beer" because there is a grammatical dependency (in the English language) between articles and nouns. Given that the assumption does not hold in the real world, while being a mathematical pre-condition for several machine learning models, it should be no surprise that learning algorithms may perform poorly on new data.

4.6.7 Vector Representation of Words



• How to measure the importance of words?

The "term-document" matrix representation is also a very useful one. In a term-document matrix, each column represents a word in the vocabulary, and each word represents a document from some collection of documents. This is a representation used, for instance, in keyword-based search engines.

With this representation, it is possible to create a representation space (or vector space) where each document (row) is represented as a *vector* in the multi-dimensional space defined by the terms in the dictionary. Using this representation, it is possible to calculate the distance (or similarity of documents through measures like cosine similarity. This allows, for instance, to identify similar documents (at least, according to the terms they contain). This is an approach common in clustering methods.

This vector semantics can also be used to represent the meaning of words. We do this by creating a vector for each word (column). In this way, it is possible to calculate the similarity of words based on the documents they tend to occur in.

In the following slides, we will describe how to improve the simple binary representation of words in documents by exploring ways to quantify their importance.

4.6.7.1 Term frequency tf

- Raw frequency $tf(t, d) = f_{t,d}$
- Log normalisation $tf(t, d) = log(1 + f_{t,d})$ Normalised Frequency $tf(t, d) = 0.5 + \frac{0.5f_{t,d}}{f_{max}(d)}$
- Measuring the importance of a word *t* to a *document d*
- The more frequent the word, the more important it is to describe the document

Term Frequency is the simplest way to measure the importance of a word in a document. Several variations can be used:

- **Raw frequency**: the number of times the term appears in the document. In this formula, t is the term (word) for which we are calculating the TF, d is the document containing the term, and $f_{t,d}$ is the frequency of the term t in the document d.
- Log normalization: this formula applies logarithmic normalization to the raw frequency of a term in a document. In this formula, t, d, and $f_{t,d}$ are defined as in the raw frequency formula. Log normalization is used - and useful - to avoid prevalent terms from having excessive importance in the language model.
- **Normalised frequency**: This formula normalizes the raw frequency of a term in a document by dividing it by the maximum raw frequency of any term in the document, and then adding a small constant to avoid zero values. In this formula, t, d, and $f_{t,d}$ are defined as in the raw frequency formula, and $f_{\max}(d)$ is the maximum raw frequency of any term in the document d.

4.6.7.2 Inverse document frequency *IDF*

TF calculates the importance of a word in a document. **IDF**, on the other hand, measures the importance of a word in the document collection.

 $\mathsf{IDF}(t, D) = \log \frac{N}{|d \in D: t \in d|}$

- Measuring the importance of a word t to a *document collection* D
- Rare terms are more important than common terms

In this formula, t represents the term for which we are calculating the IDF, D represents the set of documents in the corpus, N is the total number of documents in the corpus, and $|d \in D : t \in d|$ is the number of documents in which the term t appears. The logarithm is typically taken with base 2.

By importance we mean *discriminative* (or *descriptive*) power, that is, how much the word is helpful to discriminate a document in the collection. For example: if all documents in a collection contain the word *design*, but only a few selected documents contain the word *machine*, then *machine* is more discriminative in the document collection that *design*.

4.6.7.3 *TF* – *IDF*

 $\mathsf{tfIDF}(t, d, D) = tf_{t,d} \times IDF_{t,D}$

• *Scaling* a word's importance (in a document) based on both its frequency and its importance in the collection

TF - IDF allows weighting the importance of a word in a document by its overall importance in the collection.

4.6.8 N-gram language models

An *n-gram* is a sequence of *n* words. A 2-gram - *bigram* - is a two-word sequence. For instance "stop here", "read this", "stand up", "New York". A 3-gram - *trigram* - is a three-word sequence. For instance, "stop hear immediately" or "stand up there". Language models based on n-grams are much simpler than the most recent ones based on deep learning architectures. Yet, they are important to understand the fundamental concepts of language modeling.

- Calculate the conditional probabilities among *adjacent* words
- Given the word w, what is the probability of the next word w+1
 - e.g., given *eat*, *eat* on vs. *eat* British
- bi-grams -> 2 words, 3-grams -> 3 words

p(w)

eat on	0.16	eat Thai	0.03
eat some	eat some 0.06 eat breakfast		0.03
eat lunch	0.06	eat in	0.02
eat dinner 0.05 eat Ch		eat Chinese	0.02
eat at	0.04	eat Mexican	0.02
eat a	0.04	eat tomorrow	0.01
eat indian	0.04	eat dessert	0.007
eat today	0.03	eat British	0.001
	eat on eat some eat lunch eat dinner eat at eat a eat indian eat today	eat on0.16eat some0.06eat lunch0.06eat dinner0.05eat at0.04eat a0.04eat indian0.04eat today0.03	eat on0.16eat Thaieat some0.06eat breakfasteat lunch0.06eat ineat dinner0.05eat Chineseeat at0.04eat Mexicaneat a0.04eat tomorroweat indian0.04eat desserteat today0.03eat British

Through n-gram models, it is possible to estimate the probability of the last word of an n-gram given the previous words (or *history*). Suppose we have the sequence of words "I am looking forward to eat" and we want to know the probability that the next word is "mexican". One way to estimate the likelihood of the term following *eat* is through relative frequency counts: given a document collection, count the number of times a word follows the sentence "I am looking forward to eat". However, the problem with this approach is that no matter how large the collection is, it is impossible to make good estimations. This is because language is creative; new sentences are constantly created, and it is impossible to count the frequencies of sentences with slight variations.

The intuition of the n-gram model is that instead of computing the probability of a word given its entire history, we can **approximate** the history by just the last few words. For instance, the figure in the slide shows an example where the *history* is the single word *eat* - these are bigrams, where we compute the *conditional probability* of a word given that its predecessor is *eat*. Imagine also collecting similar conditional probabilities for the words "to", "forward", "looking", "am" and "I": by multiplying the different probabilities, we could get an *approximation* of the likelihood of the whole sentence.

While the example in the slide shows bigrams, it is possible to generalize to trigrams, or larger bigrams that look at more words in the past.

4.6.8.1 Properties of N-grams-based Language Model

- More accurate
 - The probabilities depend on the considered **context**
- The model accuracy increases with N
 - The syntactic/semantic contexts are better modeled
- Grammatical rules
 - e.g., an adjective is likely to be followed by a noun
- Semantic restrictions
 - e.g., Eat a pear vs. Eat a crowbar
- Cultural restrictions
 - e.g., Eat a cat

N-gram models are more accurate than models where words are treated as statistically independent because they can capture some of the **context** in which words appear. And, as you can imagine, the larger the *n*, the better the model could capture some linguistic properties.

N-gram models are known to be good at capturing *syntactic* linguistic phenomena, like grammar rules, and some *semantic* dependencies captured through sheer likelihood.

4.6.8.2 Limits of N-grams-based Language Model

- Conditional probabilities are difficult to estimate
 - For dictionary contains D terms there are D^N N-grams (30K words, 900M bi-grams)
 - * the corpus should be billions of documents big for a good estimation
- They do not generalize to unseen words sequences

N-gram models have been popular for a long time, but they suffer from several practical limitations.

The most obvious one is related to vocabulary size. The larger the vocabulary, the larger the size of the text collection needed to properly estimate the conditional probabilities, especially for the rarest n-grams. Also, n-gram models have issues with unseen (or new) words as, by definition, they are not part of their probability estimation.

The choice of the size of n is also important, as it can be too small or too big. In practice, using *trigram* models is common, or even 4-gram and 5-gram when there is sufficient training data.

Neural networks have largely supplanted N-gram language models. We will see that approaches based on neural networks do not make the n-gram assumption of restricted context. They can incorporate arbitrarily distant contextual information while remaining computationally and statistically tractable.

4.6.9 Representing words by their contexts

- **Distributional semantics**: A word's meaning is given by the words that frequently appear close-by
- When a word *w* appears in a text, its **context** is the set of words that appear nearby (within a *fixed-size* window)
- The contexts in which a word appears tell us much about its meaning

Another way to approach the problem of language representation is to consider words based on their context. A recurring theme in natural language processing is the complexity of the mapping from words to meaning. In the previous approaches, this mapping has been mostly based on words (or sequences of words) as a basic unit of analysis.

"You shall know a word by the company it keeps" *The distributional hypothesis, John Firth (1957)*

An alternative approach is based on the idea that *a word's meaning can be defined by its distribution in language use*. That is, the meaning of a word can be defined in terms of other words it tends to occur with, the words that tend to occur with those words, and so on. Or, differently described, *words that appear in similar contexts have similar meanings*. For instance, "laugh" tends to occur in the same context as "humour". "Wine" tends to occur in the same context as "beer".

This idea is known more formally in linguistics as **distributional semantics**. The underlying hypothesis of distributional semantics is that "the degree of semantic similarity between two linguistic expressions A and B is a function of the similarity of the linguistic contexts in which A and B can appear.



- What other words fit into these contexts?
- Contexts
 - 1 A bottle of ==___= is on the table
 - 2 Everybody likes ==___=
 - 3 Don't have ==___= before you drive
 - **4** We make ==___= out of corn

To give you an example, based on a similar one describe in [^1 Lin, D. (1998). Automatic retrieval and clustering of similar words.] Have you ever heard the word *tezgüino*? I never did, at least before reading the paper above and some textbooks mentioning the example. Let's assume that an NLP system is in your same situation: encountering a word that has never been seen before. How to infer the meaning of this word?

Imagine that you have collected some statistics, and observed that the word *tezgüino* is used in the context listed in the slide (e.g. A bottle of **____** is on the table). What other words have been used in similar contexts? Each row in the table in the slide shows if they appear (1) or do not appear (0) in such contexts (the columns). Each row, in practice, is a *vector representation* of the term, it summarises its *contextual properties*. Note that the table in the slide is very similar to the "term-document" matrix we

encountered before; a significant difference, however, is that we now use *contexts* and not *documents* as dimensions associated with words. These vectors describe the **distributional properties** of each word.

What other words appear in a similar context? *Wine* is very similar to *tezgüino*, "motor oil" and "tortillas" are fairly similar to *tezgüino*, "loud" and "choices" are completely different. This tells us that, probably, *tezgüino* is a beverage, perhaps an alcoholic one.

The distributional hypothesis has been extremely successful, and it is on the basis of modern NLP approaches based on deep learning techniques. It allows leveraging large amounts of unlabeled data to learn about rare words that do not appear in labeled training data. Also, they have a striking ability to capture lexical semantic relationships such as analogies.

4.6.9.1 Distributional Word Embeddings



- Define dimensions that allow expressing a *context*
 - The vector for any particular word captures how strongly it is associated with each context
- For instance, in a 3 -dimensional space, the axis could have the semantic meaning

- *x* -axis represents some concept of *"animal-ness"*
- z -axis corresponds to "food-ness"

The distributional hypothesis says word meaning relates to the "contexts" in which the word appears. The idea here is that once all the words in the vocabulary are appropriately placed in the *contextual* space, the meaning of a word can be represented by its location in this space. Take the example in the slide - our old "pizza example". Assuming we now want a three-dimensional context space, what is the meaning of each axis? How to define what a *context* is?

In the example of the previous slide, we used a few complete sentences. In practice, there are too many sentences - a problem we have encountered already with n-gram models.

4.6.9.2 Distributional Word Embeddings



- Defining the axes is **difficult**
 - How many?
 - * A lot less than the size of the dictionary (dense vectors)
 - * But at least ~100-dimensional, to be effective
 - * GPT-2 has 768, ChatGPT 12,288
- How to assign values associated with the vectors?
 - Tens of millions of numbers to tweak

To make things more complicated, how many axes should there be?

The basic idea is to create vectors (word embeddings) that are **dense**; that is, they do not have a lot of 0 values in them, as opposed to **sparse vectors** as the ones we encountered when we discussed the

"term-document" matrix. It turns out that **dense vectors** are much better in NLP tasks, although we do not understand why yet - probably, language has some form of regularity that we cannot understand yet.

Empirical evidence shows that in a typical large corpus, the size of these vectors should be *at least* of 100 elements. Note that representing words as 100-dimensional dense vectors would require our machine learning models to learn far fewer weights than if we represented words as 50,000-dimensional vectors. A smaller parameter space also helps avoid overfitting and lack of generalization. Still, even with 100 dimensions, there can be tens of millions of weights to estimate - a process needed to associate numerical values to each word in the dictionary.

The field of NLP has been progressing at incredible speed lately. Modern language models like GPT-3 or GPT-4 have scaled up to thousands of dimensions while controlling a massive learning process. One of the reasons they can perform better than previous models - and perform incredibly in various tasksis that they can capture these distributional semantics much better.



Here's an example of how single words (here, common nouns) might be distributed in the dense representation space of a system like GPT-3. Here the multi-dimensional space has been projected into a 2-dimensional one. We can observe how "semantically similar words" are placed nearby.

This example and the following have been taken from this blog post.

Alessandro Bozzon



In this example, using a different projection, it is possible to observe how words tend to cluster also based on their part of speech. Clearly, some words can have a different role in a sentence (as a verb or as a noun).



Interestingly, also sentences seem to organise coherently in the feature space. This is probably due to big language models like GPT-3 (or GPT-4) to capture extensive contexts (sequences of very long words).

An important observation: these language models are trained on document collections that might contain different biases.

4.6.10 How to calculate Distributional Word Embeddings?

- With machine learning models
- Advanced topic
 - Wait for Advanced Machine Learning for Design :)

One last issue needs to be explored: how are these dense vectors, these contexts, identified? How are the values in these vectors calculated? This is a relatively advanced topic that requires much more time than we have available in this course.

4.6.10.1 Ok, just a sneak peak

- SKIPGRAM: Predict the probability of context words from a centre word
- Input: one-hot vector of the centre word
 - the size of the vocabulary
- **Output**: one-hot vector of the output words
 - the probability that the output word is selected to be in the context window
- Embeddings: lower-dimensional representation of context of co-occurence

Skipgram is an example of an algorithm that calculates *static embeddings*, fixed embedding for each word in the vocabulary. More recent approaches (like BERT) learn dynamic contextual embeddings, in which the vector for each word is different in different contexts.

The intuition of behind the *skipgram* algorithm is that instead of counting how often each word w occurs near another one, we train a classifier on a prediction task where for each word in the vocabulary we classify if it is likely (or not) to appear close w—the learned classifier weights as the word embeddings. Given a specific word in the middle of a sentence (the input word), the network predicts the probability for every word in the vocabulary of being the "nearby word". The fascinating idea here is that, instead of using labeled data, the text is used to supervise the classifier's training. That is called **self-supervision** in machine learning lingo.


// http://mccormickml.com/2016/04/19/word2vec-tutorial-the-skip-gram-model/

The learning algorithm for skipgram embeddings takes a corpus of text as input and a chosen vocabulary size N. Words are represented as one-note encoding vectors having the size of the vocabulary. The network output (of the same size) is a single vector of the same size as the input. The network is trained by the word pairs (at a different distance) found in the training documents. The network learns statistics from the number of times each pairing shows up. For every word in the language, their network is then able to calculate the probability that a randomly selected nearby word is that vocabulary word. If two different words have very similar "contexts" (that is, what words are likely to appear around them), then the model needs to output very similar results for these two words.

4.7 4.7 Using Word Embeddings

4.7.1 How can embeddings be used with NLP Models?

- Word embeddings are trained from a corpus
 - And then they can be reused!
- 3 scenarios

Machine Learning for Design



Before the advent of deep learning and of large language model, NLP models were re-trained on each task, specifically for the type of task that they were trained for. For example, a model trained for *sentiment analysis* would use a dataset annotated with the desired output (e.g., negative, neutral, and positive labels). The trained model would have bee used exclusively for sentiment analysis and not, for instance, for *part-of-speech* (POS) tagging, no matter how good the model was. Obviously, this is due to the task-specificity of the training dataset. However, there could have been some commonalities to be exploited. For instance adjectives (POS) like "beautiful" would have also use in the context of sentiment analysis: the language is the same, just its type of analysis is different.

This is where **word embeddings** can become useful. As they are used to represent words, while capturing the semantic relationship between words by exploiting similar contexts, it is possible to use them as *input* for multiple machine learning models - we called them *downstream tasks*. Word embeddings can be **pretrained**, that is, *learned independently* of the downstream tasks. In this way, the model can focus on learning higher-level concepts that cannot be captured by word embeddings (e.g., phrases, syntax, and semantics), and on the task-specific patterns learned from the given annotated data.

There are three ways embedding can be used.

4.7.2 Retraining

• Train word embeddings and your model at the same time using the train set for the task

In the simplest scenario, the word embeddings are initially "empty", and trained together with the downstream task. This means that the activity of calculating the embeddings and the downstream models occur at the same time. This scenario can occur when there is no existing pre-trained model available, but some training dataset exists. So, in this way, it is possible to create embeddings that could be later reused.

4.7.2.1 Fine-Tuning

- Initialise the model using the pre-trained word embeddings
 - e.g., train on Wikipedia, or large Web corpora
- Keep the embedding fixed while training the model for the task
 - Another example of transfer learning

This second scenario is the classic **transfer learning** scenario. You might remember from previous lectures that transfer learning is technique in which the performance of a machine learning model are improved by used data and/or models trained in a different tasks.

Transfer learning always consists of two or more steps:

- pre-training: a machine learning model is first trained for one task
- *fine-tuning*: the same model is used for both tasks. This is the Scenario 2 discussed in this slide. Alternatively, we have
- *adaptation*: the same model is adjusted and used in another task.

In literature about ML, you might also encounter the term **domain adaptation**. This is a technique, close to transfer learning, where the idea is to train a machine learning model in one domain (e.g., news) and adapt it to another domain (e.g., sport). **Domain adaptation** is typically performed with models that operate on the same task (e.g., text classification).

4.7.2.2 Adaptation

- The embeddings are adapted while the downstream model is trained, the train set for the task
 - Same as *Scenario 2*, but the embeddings are now more *close* to the words distribution in your training set

As discussed above, in this scenario the training of the embedding does not start from skratch. The weights are adjusted (adapted) based on the properties of the language in the new training set. This adaptation work is usually pretty fast, as the assumption is that most of the properties of the original language hold also for the one in the training set.

4.7.3 Evaluating Word Embeddings

4.7.3.1 How to evaluate word vectors?

- Intrinsic: evaluation on a specific/intermediate subtask (e.g. analogy)
 - Fast to compute
 - * It helps to understand that system
 - * Not clear if helpful unless correlation to the actual task is established
- Extrinsic: evaluation of a real task
 - It can take a long time to compute the accuracy
 - Unclear if the subsystem is the problem or if it is an interaction with other subsystems

We know that word embedding are *useful*. Bot how do we know if they are *good*? How can we assess the quality of a word embedding?

Basically in two ways.

Intrinsically, by testing if the representations they create align with our (as designers) intuitions about word meaning. The idea is to use a similarity function (e.g. Cosine Similarity) to evaluate if words that are related in the evaluation space are also close in the embedding space. When doing this evaluation, the size of the context window (that is, how many words around the specific words are considered to define the context of use of the word) matter. Shorter context windows tend to lead to representations that are a bit more syntactic, since the information is coming from immediately nearby words. The most similar words to a target word *w* tend to be semantically similar words with the same parts of speech. Longer context windows groups together words that tend to be topically related but not similar.

Extrinsically, by evaluating if the use of the embeddings improve the performance of the downstream tasks. So, using error functions as a proxy for embedding quality. This is experimentally possible, as the learning task and the dataset do not change.

4.7.3.2 Intrinsic evaluation

• Word vector **analogies**

1 \$a:b = c:?\$

\$man:woman=king:?\$

2023-03-29



A interesting semantic property of embeddings is their ability to capture **relational meanings**. A classical way intrinsically evaluate embeddings is by using the so-called **parallelogram model** to solve simple analogy problems of the form *a* is to *b* as *c* is to what?. In such a context, given a problem like man : woman :: king :?, i.e., man is to woman as king is to ?, a system should find the word queen.

- Find a word such that the vector is closest (cosine similarity) to vec[man] vec[woman] + vec[king]
 - Correct if the word found is *queen*
- Can be applied to test for syntactic analogy as well
 - Quick : quickly = slow : slowly



Being words represented by vectors in a vector space, it is possible to apply some vector operation (e.g. subtraction or sum) to *move* around the vector space and check for the words that are occupying the targeted destination space. Note that such a space might be occupied by words that are morphological variants of the original one - for instance, **queens** in the example, or words that just happen to be in that space due to statistical artefacts.

Note that this relational properties work also for syntactic analogies.

An important disclaimer: such an approach to navigate an embedding space is too simple to mimic higher level cognitive capabilities of humans, as analogies are formed in more complex ways.

As well as gender relations. Association between people and the role they occupy in companies. Entities having some relevant relations between each other, e.g. countries and their capitals. And morphological relations between words (e.g. superlatives and comparatives).



4.7.3.2.1 There are problems, of course

- By exploring the semantic space, you can also find analogies like
 - *Thirsty* is to *drink* as *tired* is to *drunk*
 - *Fish* is to *water* as *bird* is to *hydrant*

Embeddings can also be a very useful tool. For instance, scholars used it to study how meaning changes over time, by computing multiple embedding spaces, each from texts written in a particular time period.

Biases in word vectors might leak through to produce unexpected, hard-to-predict biases

While learning word meaning, embeddings can also learn (and reproduce) from the text implicit biases and stereotypes.

- man is to woman as computer programmer is to ==_____=
- woman is to man as computer programmer is to ==_____=
- man is to genius as woman is to ==_____=
- woman is to genius as man is to ==_____=

Look for example at the four analogies in the slide. What are the right words to complete the analogies?

- *man* is to *woman* as *computer programmer* is to ==homemaker==
- *woman* is to *man* as *computer* programmer is to ==mechanical engineer==
- man is to genius as woman is to ==muse==

woman is to genius as man is to ==geniuses==

The answers in this slide are evocative of how these embeddings could capture such biases. Several scholars discovered similar examples, highlighting how this way of capturing meaning could result in several types of harms, when the system is used in real-word application. For instance **allocation harm** when a system allocates resources unfairly to different groups. Or **representational harm***, which is harm caused by deafening or ignoring social groups.

In turns out that embeddings do not only just reflect biases in the original text, they amplify it: for instance, gendered terms become **more** gendered in embedding spaces, and biases are more exaggerated.

A lot of work has been going into trying to correct these biases, but a lot of work still needs to be done. Of course, one approach would be in better curating the data that is fed to learn embeddings, but this could be an impossible task given the size of the datasets. Other **debiasing** approaches work by transferring the embedding space ands removing some type of stereotypes, but empirical evidence shows that such bias is never fully eliminated.

4.8 4.8 Large Language Models

TO DO

5 References

- CMU Computer Vision course Matthew O'Toole.
- Grokking Machine Learning. Luis G. Serrano. Manning, 2021
- [CIS 419/519 Applied Machine Learning]. Eric Eaton, Dinesh Jayaraman.
- Deep Learning Patterns and Practices Andrew Ferlitsch, Maanning, 2021
- Machine Learning Design Patterns Lakshmanan, Robinson, Munn, 2020
- Deep Learning for Vision Systems. Mohamed Elgendy. Manning, 2020 EECS498: Conversational AI. Kevin Leach. https://dijkstra.eecs.umich.edu/eecs498/
- CS 4650/7650: Natural Language Processing. Divi Yang. https://www.cc.gatech.edu/classes/AY2020/cs7650_spring
- Natural Language Processing. Alan W Black and David Mortensen. http://demo.clab.cs.cmu.edu/NLP/
- IN4325 Information Retrieval. Jie Yang.
- Speech and Language Processing, An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition. Third Edition. Daniel Jurafsky, James H. Martin.
- Natural Language Processing, Jacob Eisenstein, 2018.
- Text as Data. Grimmer, Justin; Roberts, Margaret E.; Stewart, Brandon M. Princeton University Press.

6 Credits

• Ilse Koolen