



Universiteit Antwerpen
| Faculteit Bedrijfswetenschappen
en Economie

Big Data

Hoe vinden we de weg in een massa gegevens?

Prof. David Martens

Introduction

- **About us**
 - PhD researchers at Faculty of Business & Economics
 - Research on Data Mining and Ethical AI

- **Economie Ontcijferd**
 - Clarify role of mathematics in Business & Economics

Big Data

- **Introduction to Big Data and Data Mining**
- **Decision Trees**
- **Artificial Neural Networks**
- **Ethics of Big Data**

Trends

"De marketeer van de toekomst is een data scientist met soft skills"



FINANCIAL TIMES
Data science is the big draw in business schools

Student demand for degrees in the subject soars as employers seek skilled analysts

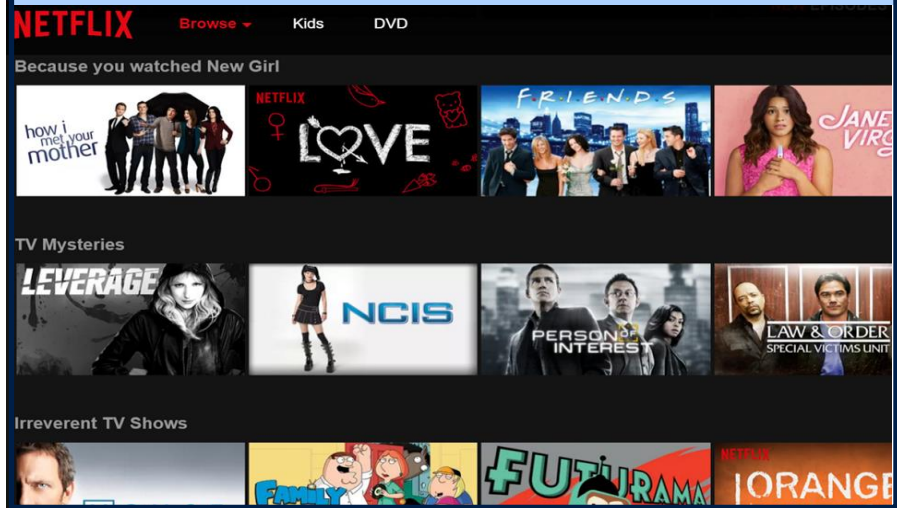
vacature.com

Data scientist, hét beroep van de toekomst

Data Mining

**“The automatic extraction of patterns
from large amounts of data”**

Recommender Systems



Targeted Advertising



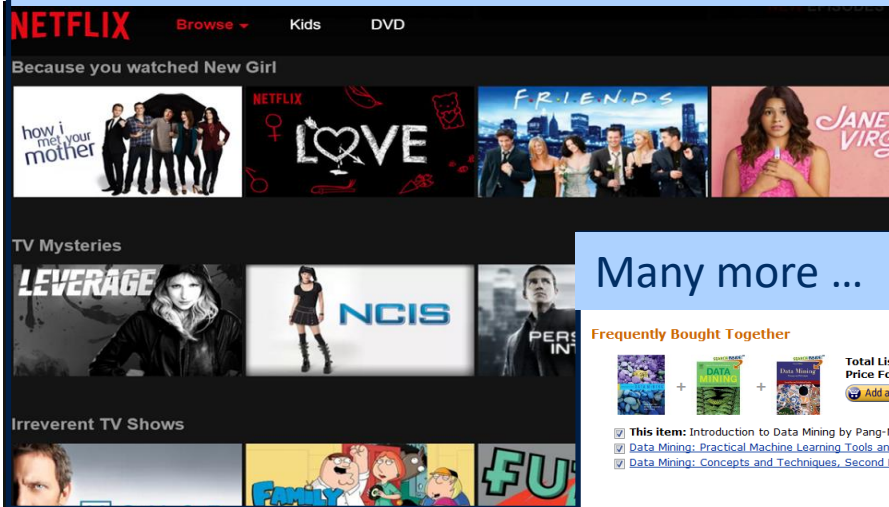
Fraud Detection



HR Analytics



Recommender Systems



Targeted Advertising



Many more ...

Frequently Bought Together



- This item: Introduction to Data Mining by Pang-Ning Tan
- Data Mining: Practical Machine Learning Tools and Techniques, Second Edition (Morgan Kaufmann Series in Data Management Systems) by Ian H. Witten
- Data Mining: Concepts and Techniques, Second Edition (The Morgan Kaufmann Series in Data Management Systems) by Micheline Kamber Jiawei Han



Fraud Detection



Data Mining

- Data mining: automatic extraction of knowledge from data
- Setting the scene with credit scoring example

Client	Income	Sex	Amount	Default
A	1.600	M	175.000	N
B	2.600	F	350.000	Y
C	3.280	M	50.000	N
D	950	M	120.000	Y
E	10.500	M	1.000.000	N
F	5.700	F	240.000	N
G	2.400	F	250.000	N

Data

Data Mining

Data mining technique

Classification Model
if income < 10.000 and Amount Loan > 100.000 and ... then default = yes

Pattern

Client	Income	Sex	Amount	Default
New client	2.000	F	500.000	Y

Classification: Classification Models

Mathematical classification models $f(x)$

$f(x) > 0.5 \Rightarrow$ customer = good

$f(x) \leq 0.5 \Rightarrow$ customer = bad

▪ Linear

- Linear, logistic regression; linear discriminant analysis
- Result: linear function of attributes
- $f(x) = 0.125 \text{ income} + 0.305 \text{ age} - 0.02 \text{ gender} + \dots - 3.1 \text{ amount loan} + 0.3$

▪ Non-linear

- Artificial Neural Networks, Support Vector Machines, RVM, ...
- Result: non-linear function of attributes
- $f(x) = 0.201 \text{ income}^2 \text{ age}^3 - 0.55 \text{ age}^3 - 5.21 \text{ gender income} + \dots + 3.6 \text{ gender}^2 \text{ amount loan}^2$

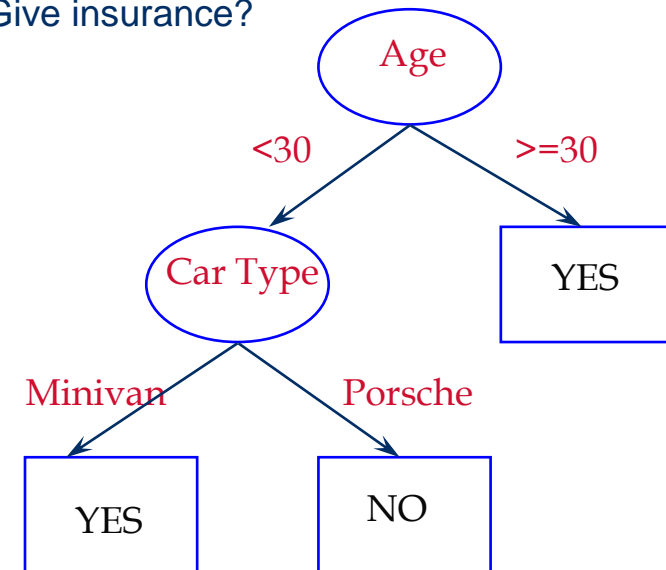
Classification: Classification Models

Rule-based classification models

- Decision Rules / Trees
 - C4.5, RIPPER, CN2, AntMiner+, ANN/SVM Rule extraction...
 - Result: set of rules or tree

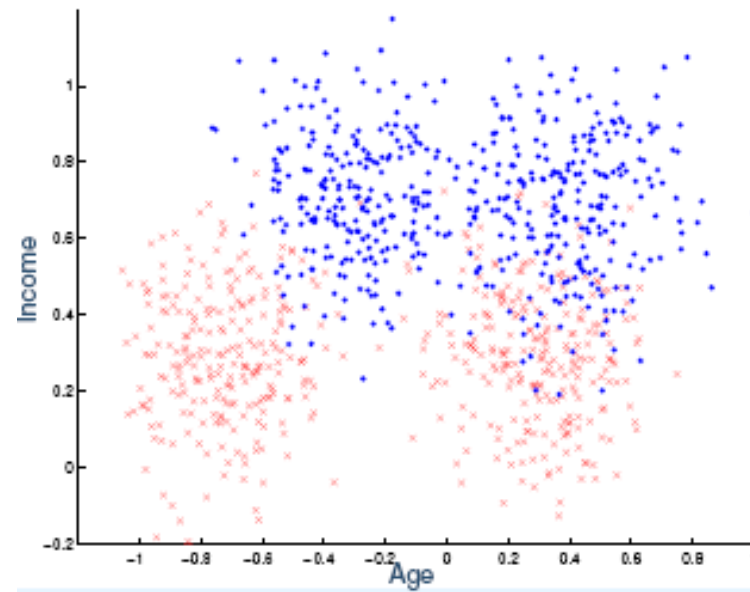
```
if (Checking Account < 200 DM and Duration > 15 m and
    Credit History = no credits taken and Savings Account < 1000 DM)
then class = bad
else if (Purpose = new car/repairs/education/others and
    Credit History = no credits taken/all credits paid back duly at this bank and
    Savings Account < 1000 DM)
then class = bad
else if (Checking Account < 0 DM and
    Purpose = furniture/domestic appliances/business and
    Credit History = no credits taken/all credits paid back duly at this bank and
    Savings Account < 500 DM)
then class = bad
else if (Checking Account < 0 DM and Duration > 15 m and
    Credit History = delay in paying off in the past and
    Savings Account < 500 DM)
then class = bad
else class = good
```

Give insurance?



Classification: Output Types

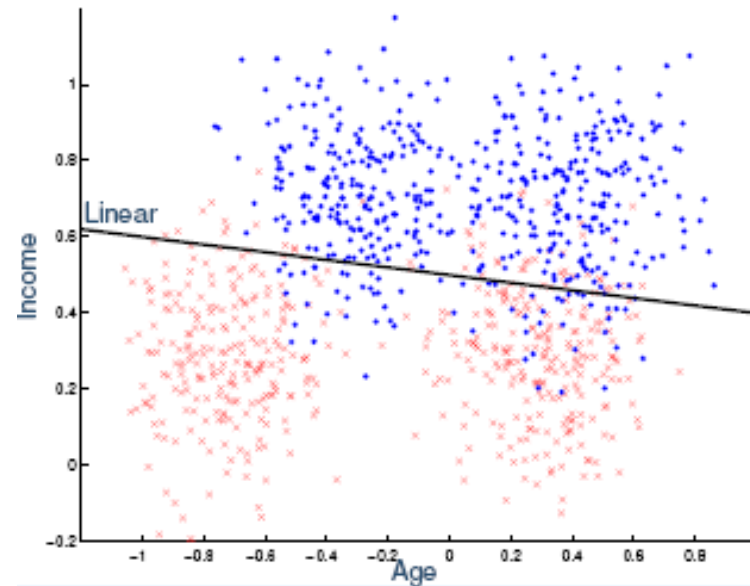
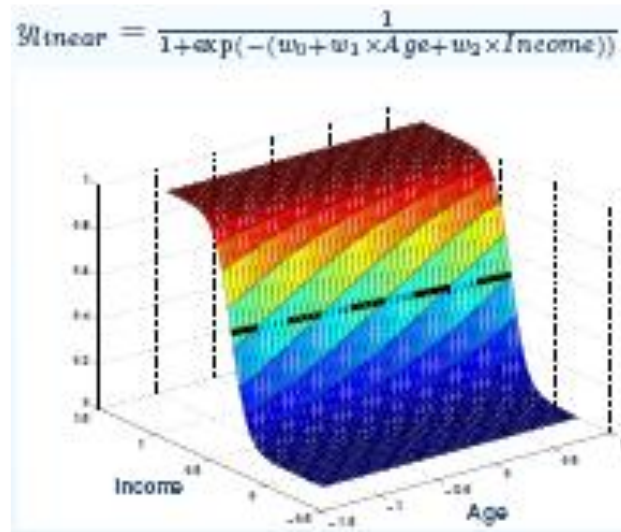
Different output types



Classification: Output Types

Different output types

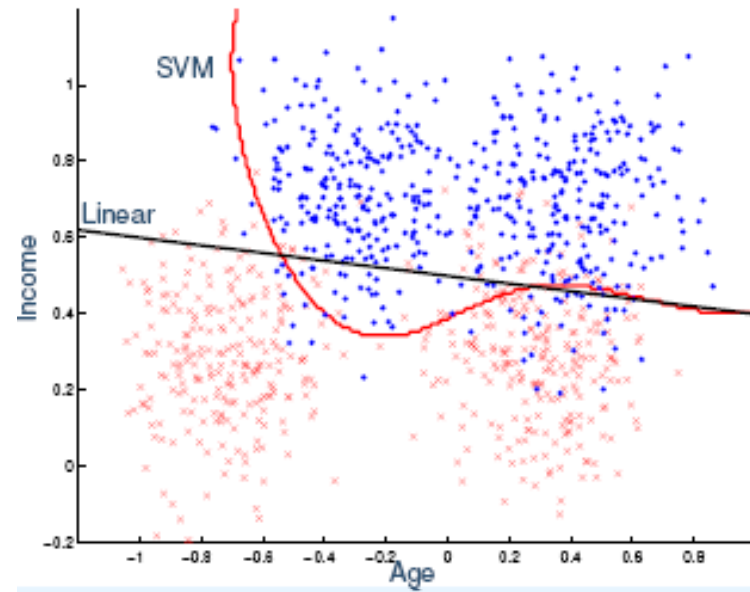
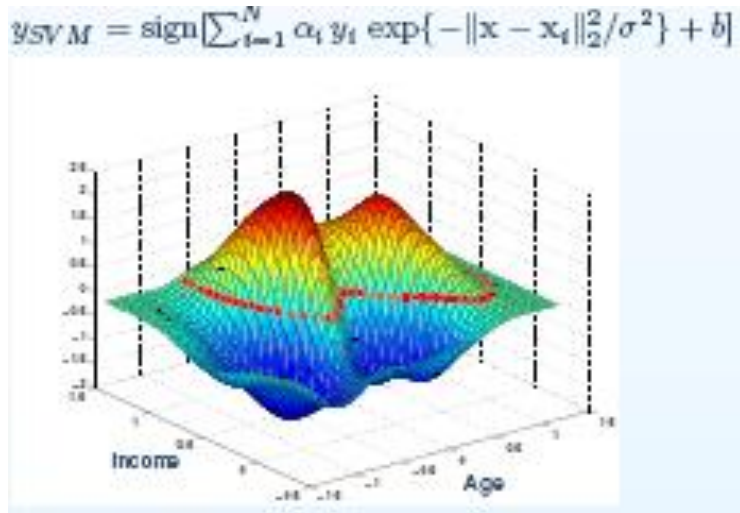
- Linear



Classification: Output Types

Different output types

- Linear
- Non-linear

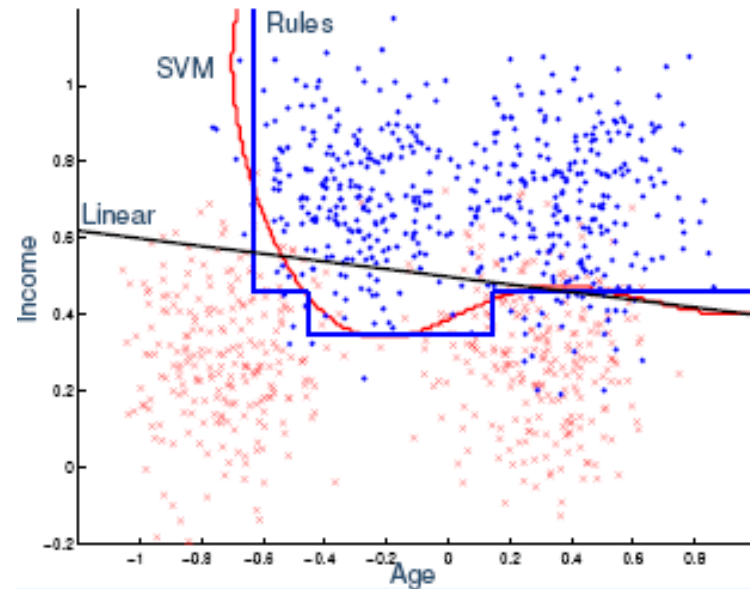


Classification: Output Types

Different output types

- Linear
- Non-linear
- Rule-based

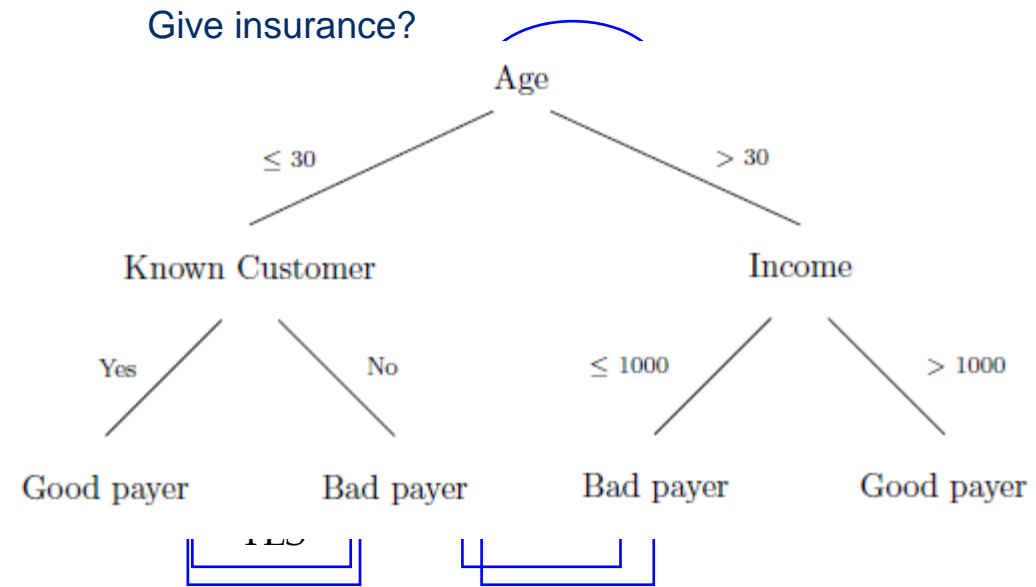
```
If age < 24 then Bad  
Else If Income < 2000 then Bad  
Else Good
```



Big Data

- Introduction to Big Data and Data Mining
- Decision Trees
- Artificial Neural Networks
- Ethics of Big Data

Decision Trees



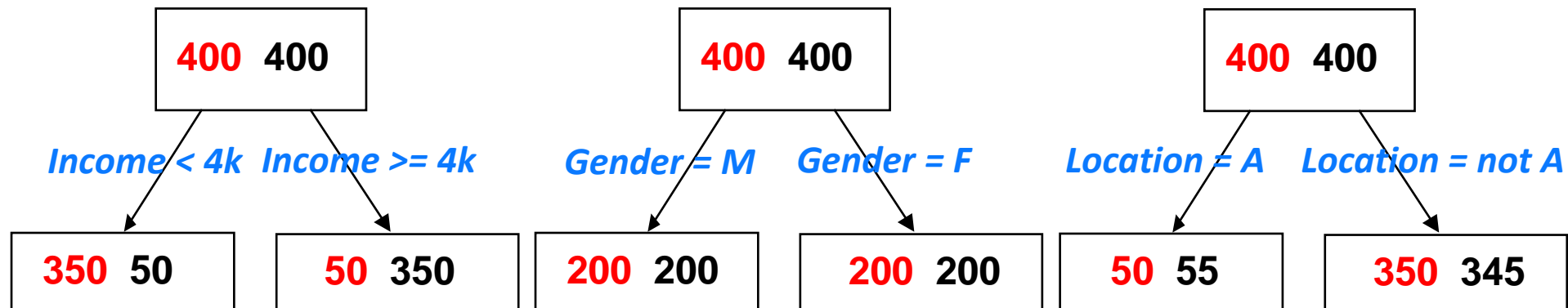
Finding informative variables from the data

- **Fundamental notion of data science:**
 - Finding and selecting **informative** variables
 - What is information? Reduces uncertainty about something.
 - *“So, if an old pirate gives me information about where a treasure is hidden ...”*



Finding informative variables from the data

- Predicting credit default
- Suppose your bank has this data
 - 800 customers, 400 known **good** ones, 400 known **bad** ones
 - 3 variables: income, gender and location
- Which variable is “most informative”?



Finding informative variables from the data

- Measuring the uncertainty
- “Technically, we would like the resulting groups to be as pure as possible” –

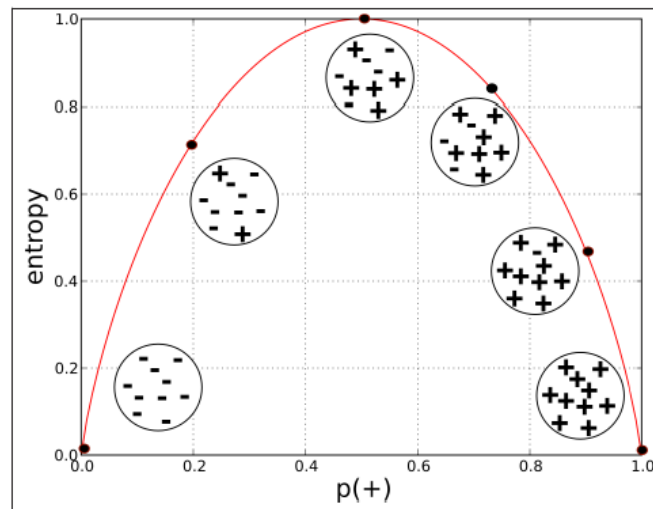


Figure 3-3. Entropy of a two-class set as a function of $p(+)$.

Why “**Entropy**”? The story goes that Shannon didn't know what to call his new information measure, so he asked **von Neumann**, who said ‘You should call it **entropy** ... [since] ... no one knows what entropy really is, so in a debate you will always have the advantage’ ([Tribus 1971](#))

Decision Trees: Impurity

Measuring the uncertainty

- **Impurity $I(n)$ of node n**

- at maximum when observations are distributed evenly over all classes
- at minimum when all observations belong to a single class

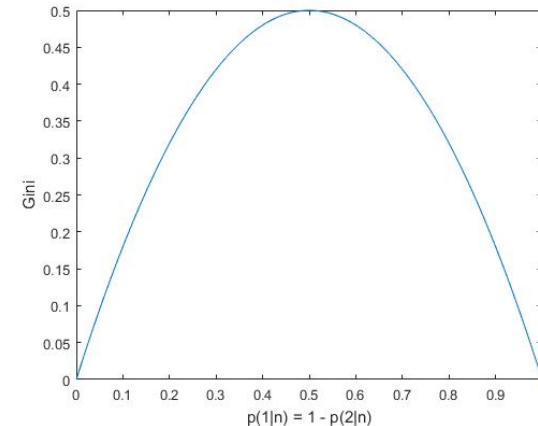
- **Two popular measures:**

1. **Entropy** measure

$$I(t) = -p(1|n) \log(p(1|n)) - p(2|n) \log(p(2|n))$$

2. **Gini** index of diversity

$$I(t) = 2 p(1|n) p(2|n)$$



Decision Trees: Impurity

Measuring the uncertainty

Gini index of diversity

$$I(t) = 2 p(1|n) p(2|n)$$

All customers

400	400
B	G

$$p(1|n) = 400 / 800$$

$$p(2|n) = 400 / 800$$

$$I(n) = 2 \times 0,5 \times 0,5 = 0,5$$

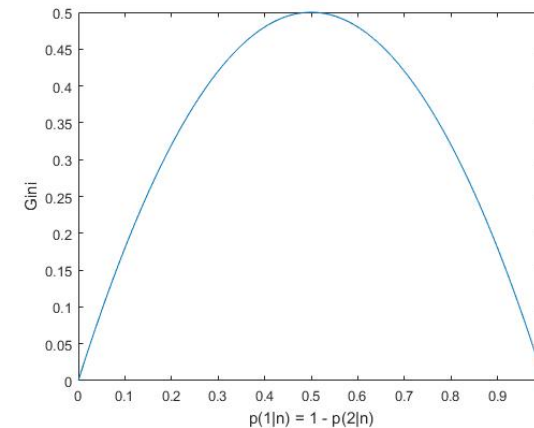
Customers with high income

100	300
B	G

$$p(1|n) =$$

$$p(2|n) =$$

$$I(n) =$$



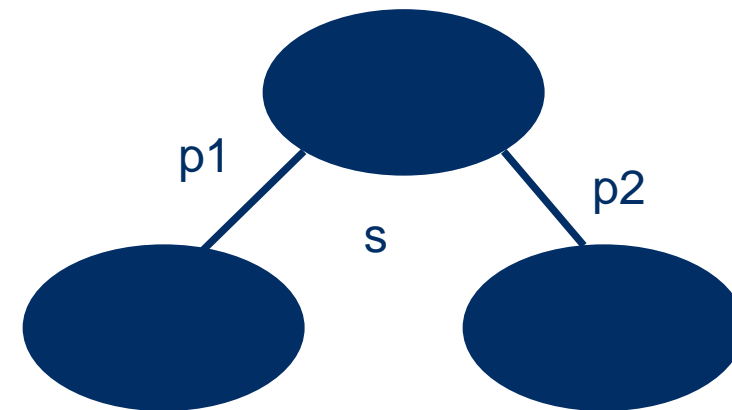
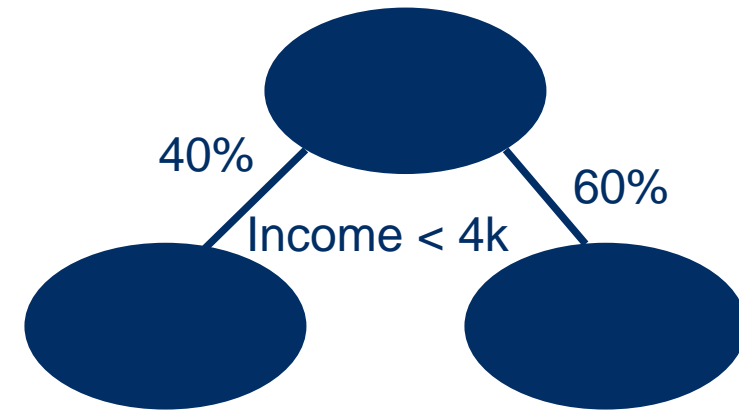
What is the uncertainty about the class to predict, for this group (as measured by Gini)? Indicate the point on the graph.

Decision Trees

Measuring the reduction in uncertainty

- Consider candidate split s of node n
- Notations
 - p_1 : proportion of the data in n that ends up in n_1
 - $I(n)$: the impurity of node n
- The goodness of the split is weighted:
mean decrease in impurity

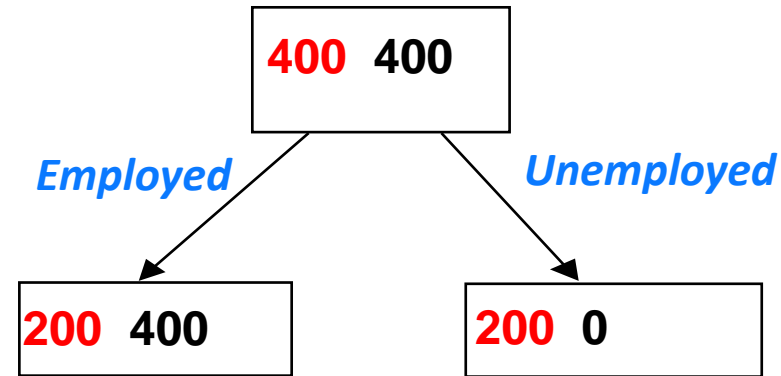
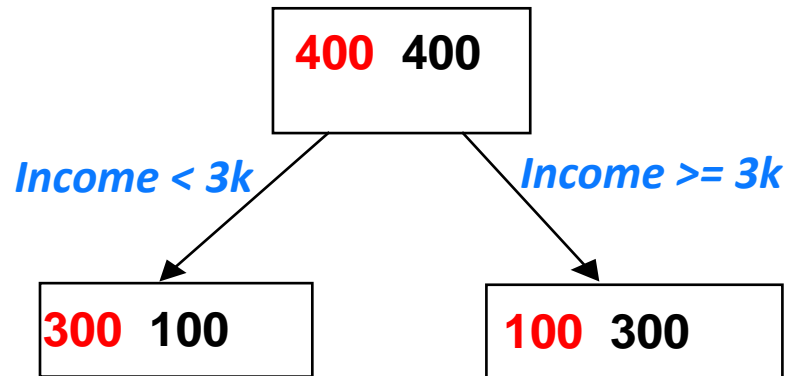
$$\Delta I(s,n) = I(n) - p_1 I(n_1) - p_2 I(n_2)$$



Decision Trees: Impurity

Notations:

- $I(n)$: impurity at node n
- $p(1|n)$: probability of being class 1 at node n
- $p(2|n)$: probability of being class 2 at node n



$I(n)$ = impurity of top node
 $I(n1)$ = impurity of left node
 $I(n2)$ = impurity of right node

$$\text{Gini: } I(n) = 2 \times p(1|n) \times p(2|n)$$

$$\begin{aligned}
 I(n) &= 2 \times (1/2) \times (1/2) = 0,5 \\
 I(n1) &= 2 \times (3/4) \times (1/4) = 0,375 \\
 I(n2) &= 2 \times (1/4) \times (3/4) = 0,375
 \end{aligned}$$

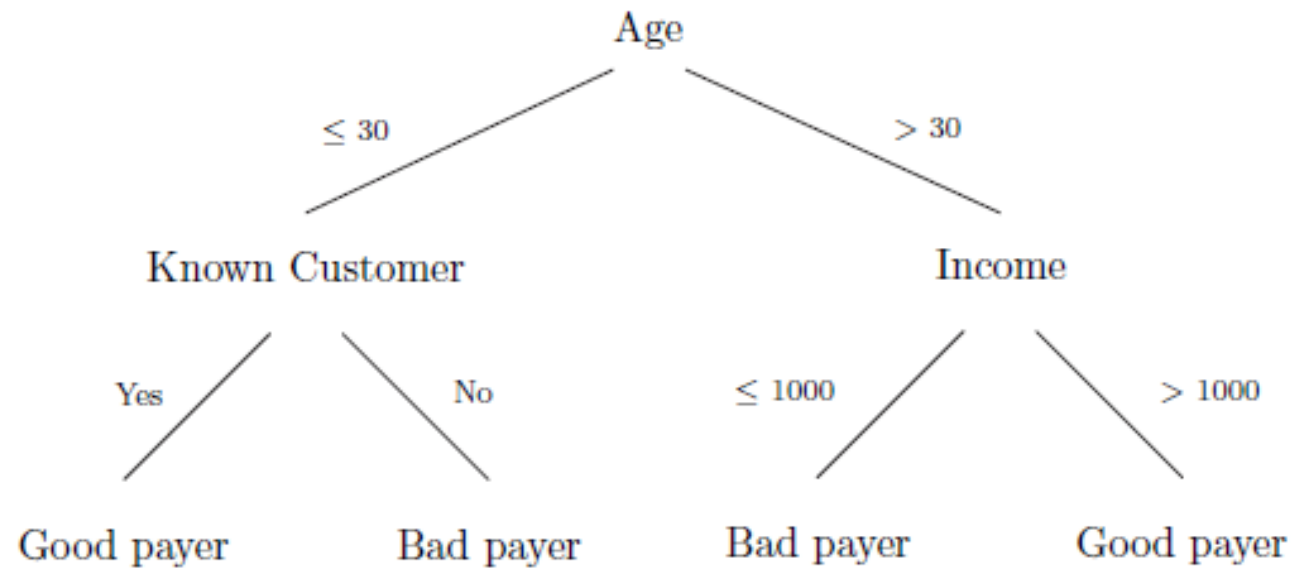
$$\begin{aligned}
 \Delta I(n) &= I(n) - (400/800) \times I(n1) - (400/800) \times I(n2) \\
 &= 0,5 - (1/2) \times 0,375 - (1/2) \times 0,375 \\
 &= \mathbf{0,125}
 \end{aligned}$$

$$\begin{aligned}
 I(n) &= \\
 I(n1) &= \\
 I(n2) &=
 \end{aligned}$$

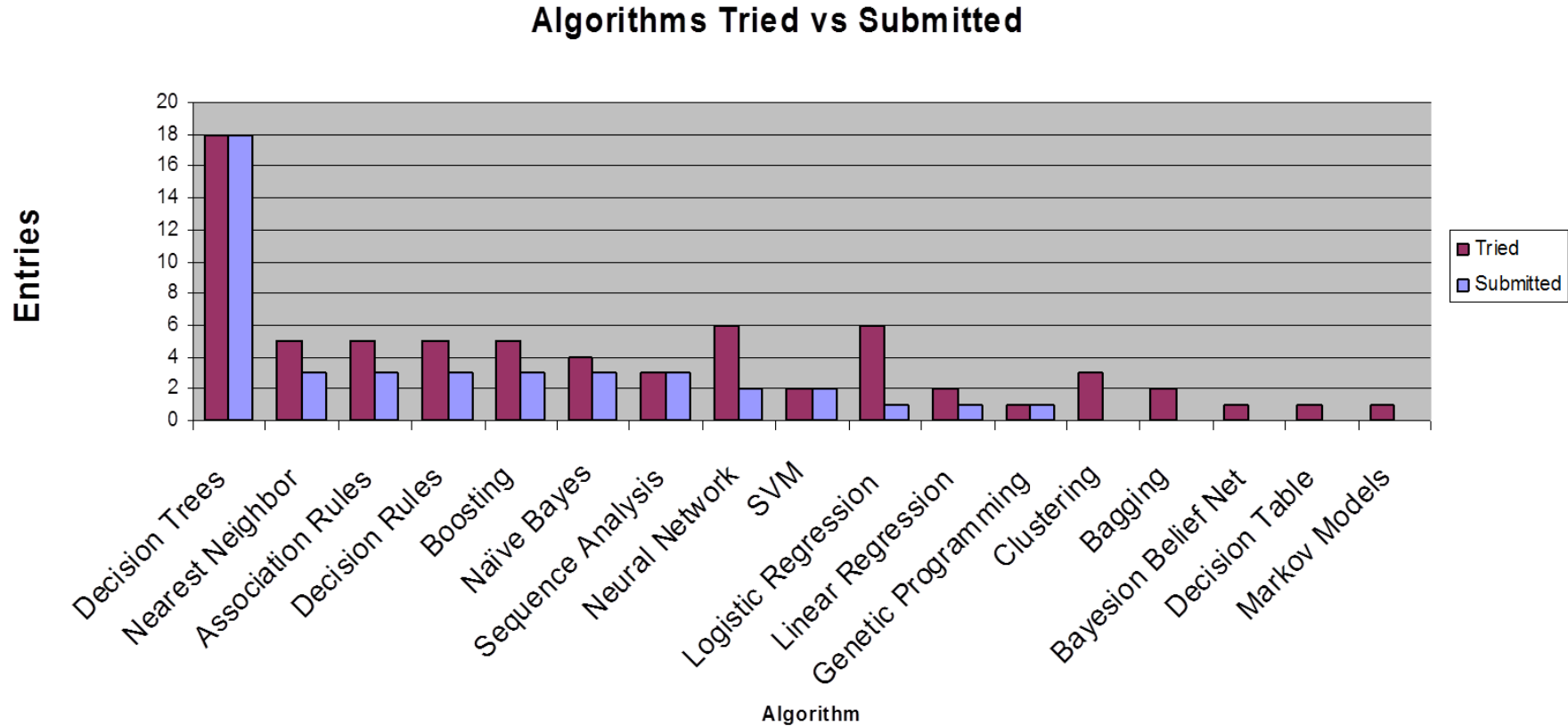
$$\begin{aligned}
 \Delta I(n) &= \\
 &= \\
 &= \mathbf{0,166}
 \end{aligned}$$

Decision Trees

Decision Tree algorithm: recursively repartitioning the data



Commonly Used Induction Algorithms



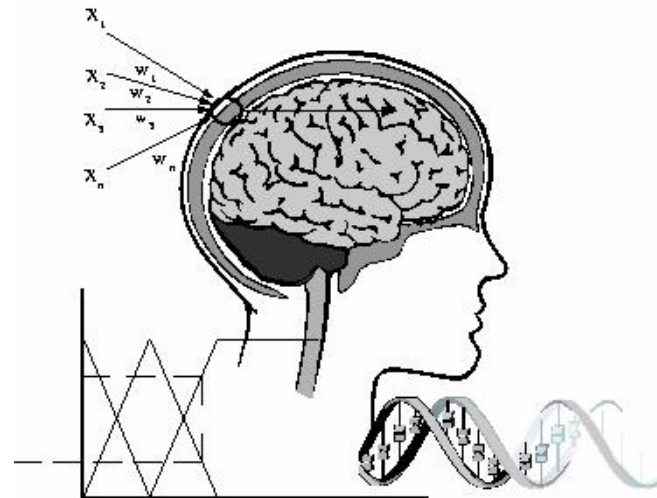
Post-mortem analysis of a popular data mining competition
Thanks to Carla Brodley & Ron Kohavi

Big Data

- **Introduction to Big Data and Data Mining**
- **Decision Trees**
- **Artificial Neural Networks**
- **Ethics of Big Data**

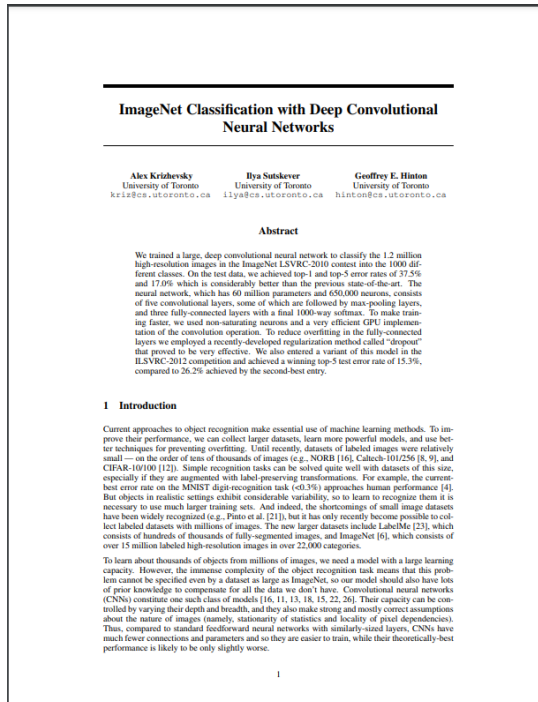
Artificial Neural Networks

- **Non-linear models**
- **Mimic human brain**
- **Good performance**



Neural Networks

AlexNet



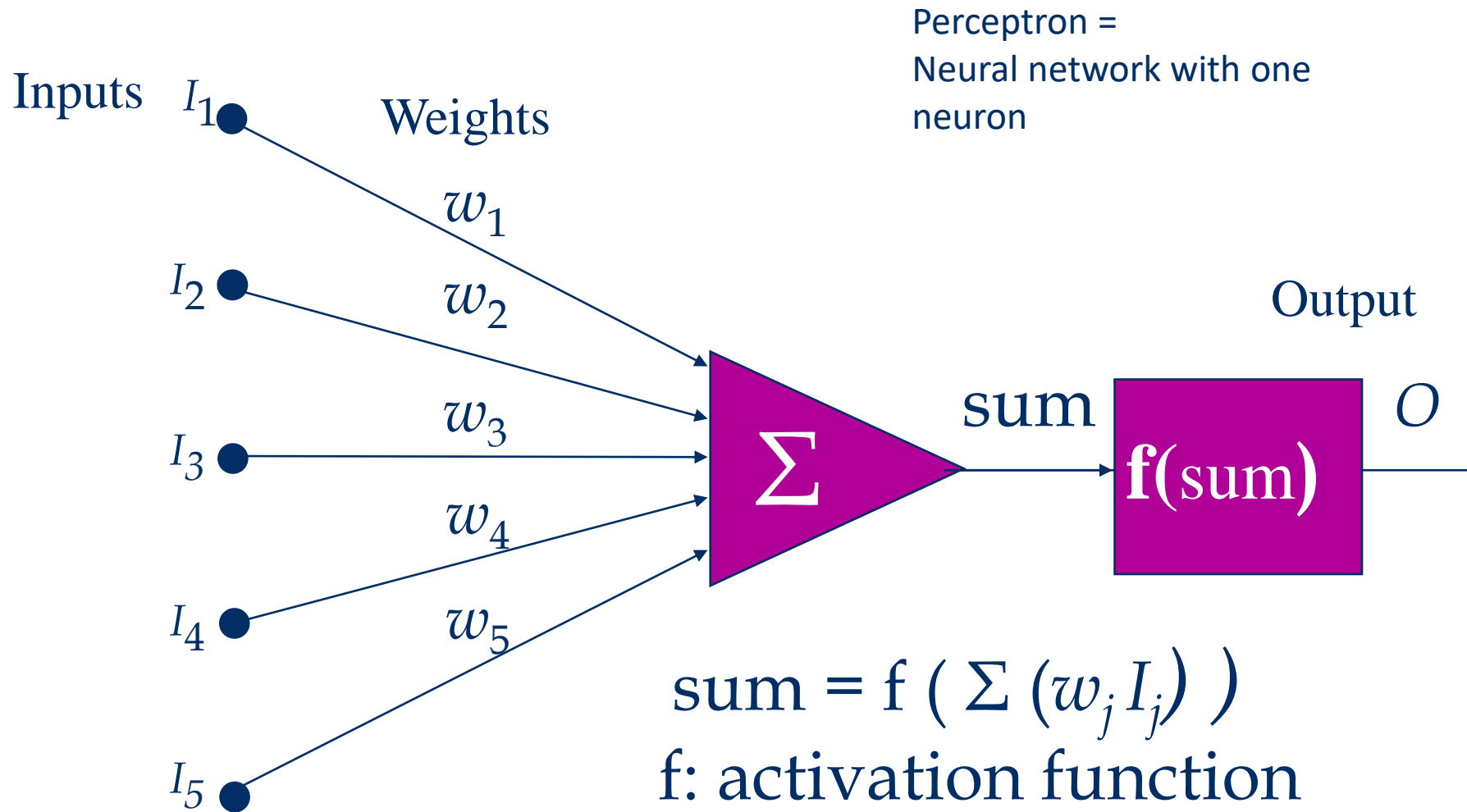
AlphaGo



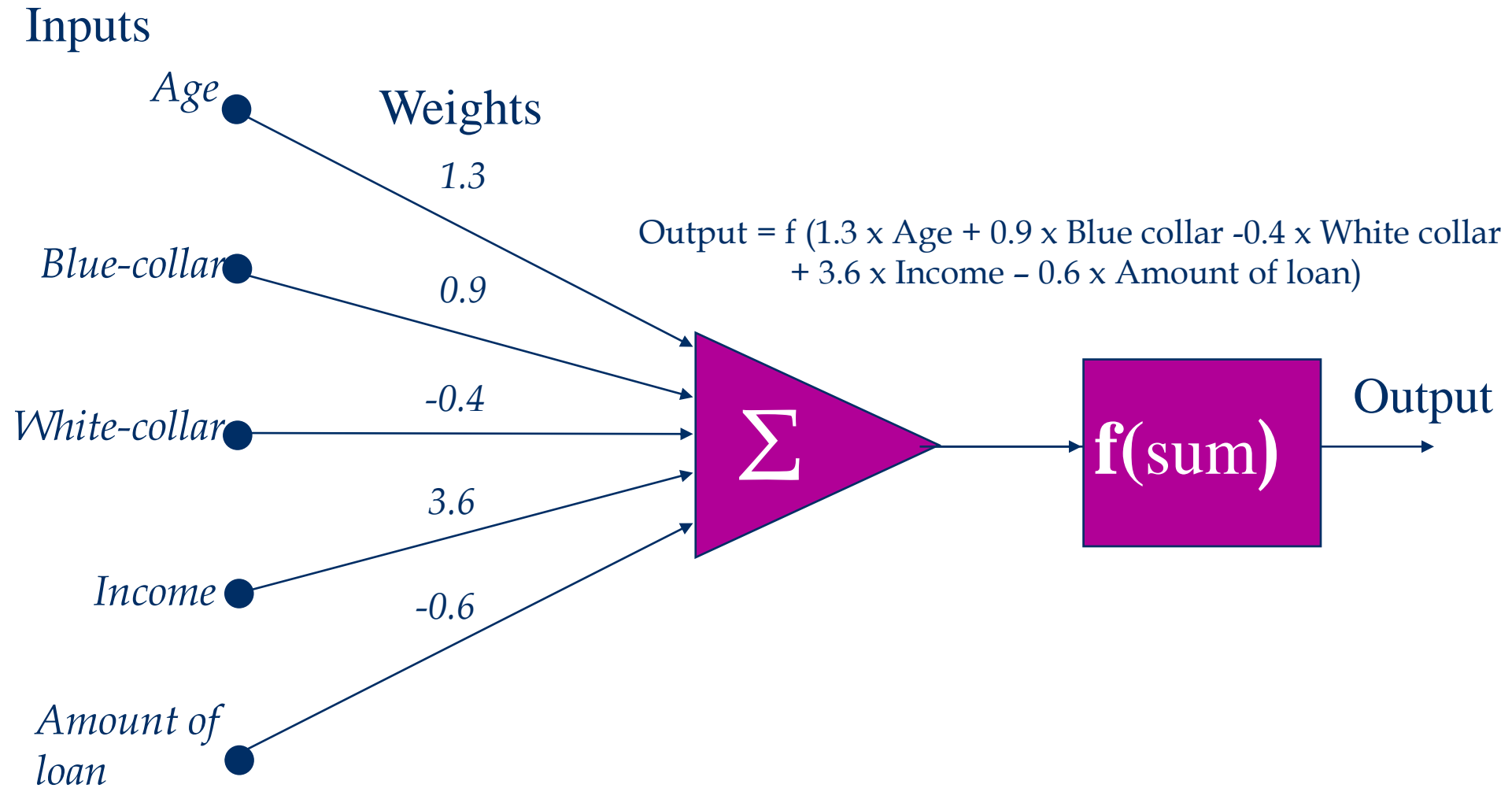
ChatGPT



The neuron model

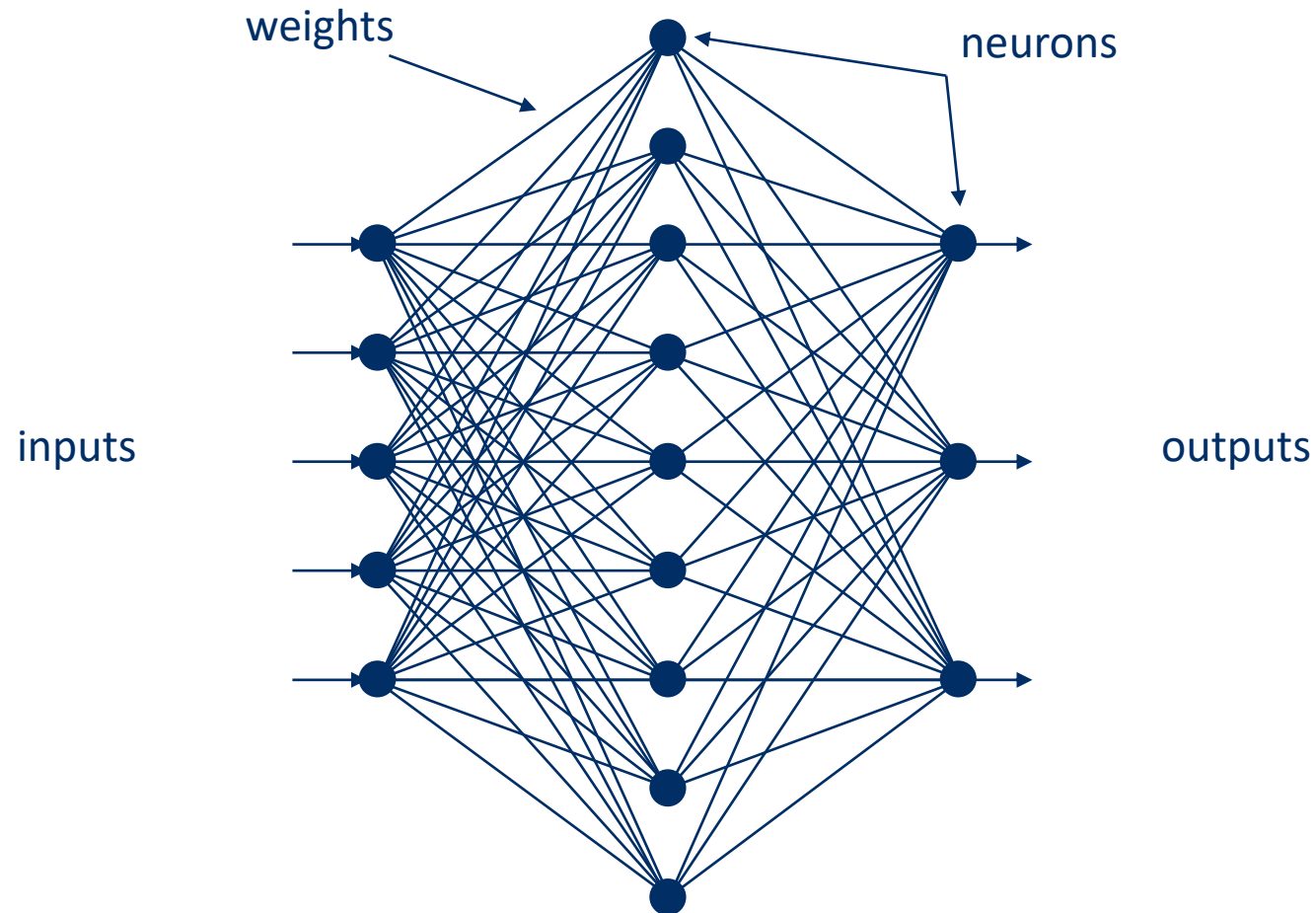


The neuron model



The Multi Layer Perceptron (MLP)

- Organise neurons into layers



Deep Learning

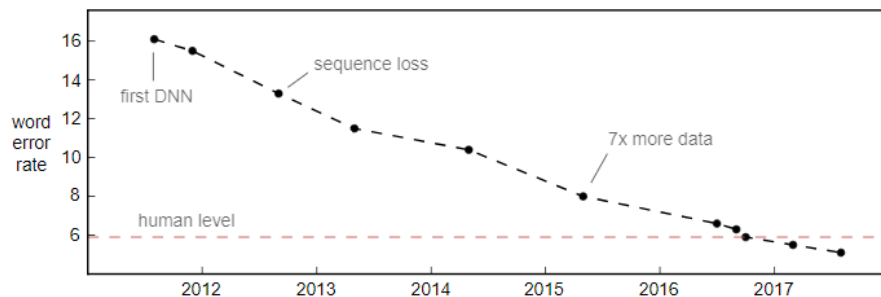
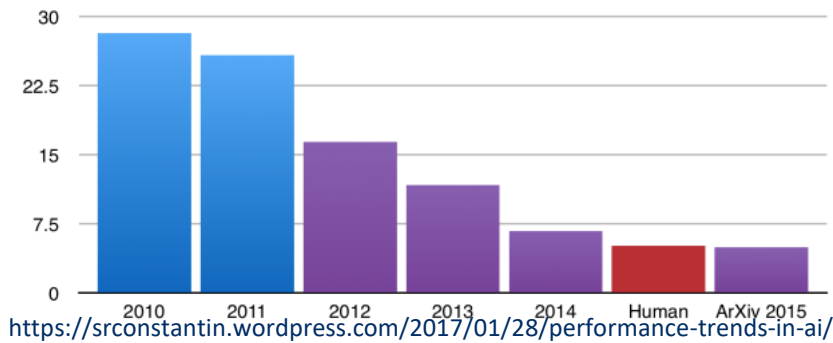
- Neural network with **many layers** (deep)
- Used primarily for image, voice, text
- Automatically learns shapes, without need of supervision!



Deep Learning

- Superhuman results: more accurate than a human

ILSVRC top-5 error on ImageNet

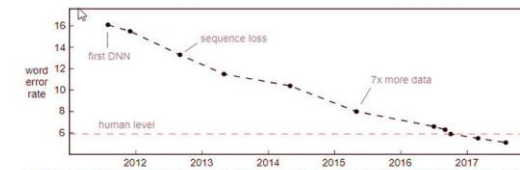


<https://awni.github.io/speech-recognition/>

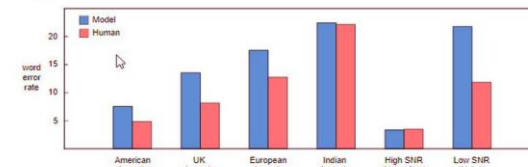
But: Challenges remain

Tesla that crashed in Autopilot mode sped up before hitting truck - police

Police report says Model S accelerated for 3.5 seconds prior to collision with stopped firetruck that left two injured



Improvements in word error rate over time on the Switchboard conversational speech recognition benchmark. The test set was collected in 2000. It consists of 40 phone conversations between two random native English speakers.



Comparison of human transcribers to Baidu's Deep Speech 2 model on various types of speech. Notice the humans are worse at transcribing the non-American accents. This is probably due to an American bias in the transcriber pool. I would expect transcribers native to a given region to have much lower error rates for that region's accents.

ChatGPT



? Help D Personal ☰

Playground

Load a preset... ▾

Save View code Share ... ⚙️

Why are neural networks better for natural language processing compared to tree-based models? 🗣️

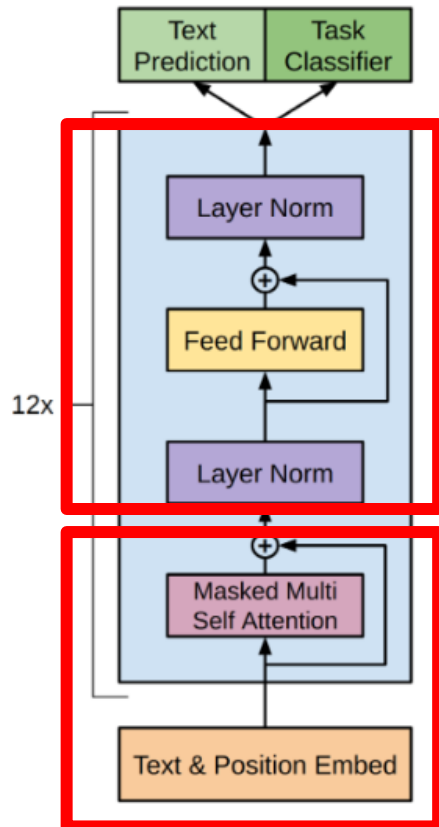
Neural networks are better for natural language processing compared to tree-based models because they are better at capturing complex relationships between words and phrases in natural language. Neural networks are also able to better handle ambiguity and variations in language, which is important for natural language processing tasks. Additionally, neural networks are able to learn and adapt to new language trends and changes in the language, which tree-based models are unable to do.



Submit ↶ ↷ 🗣️ 🍏 102

- Natural Language Processing (NLP)
- Generative model
- Trained on massive dataset (Wikipedia, Books, Reddit, Stackoverflow,..)
- Main components are neural networks

ChatGPT



Neural Network

Convert words to vectors

1. Meaning
2. Position
3. Attention

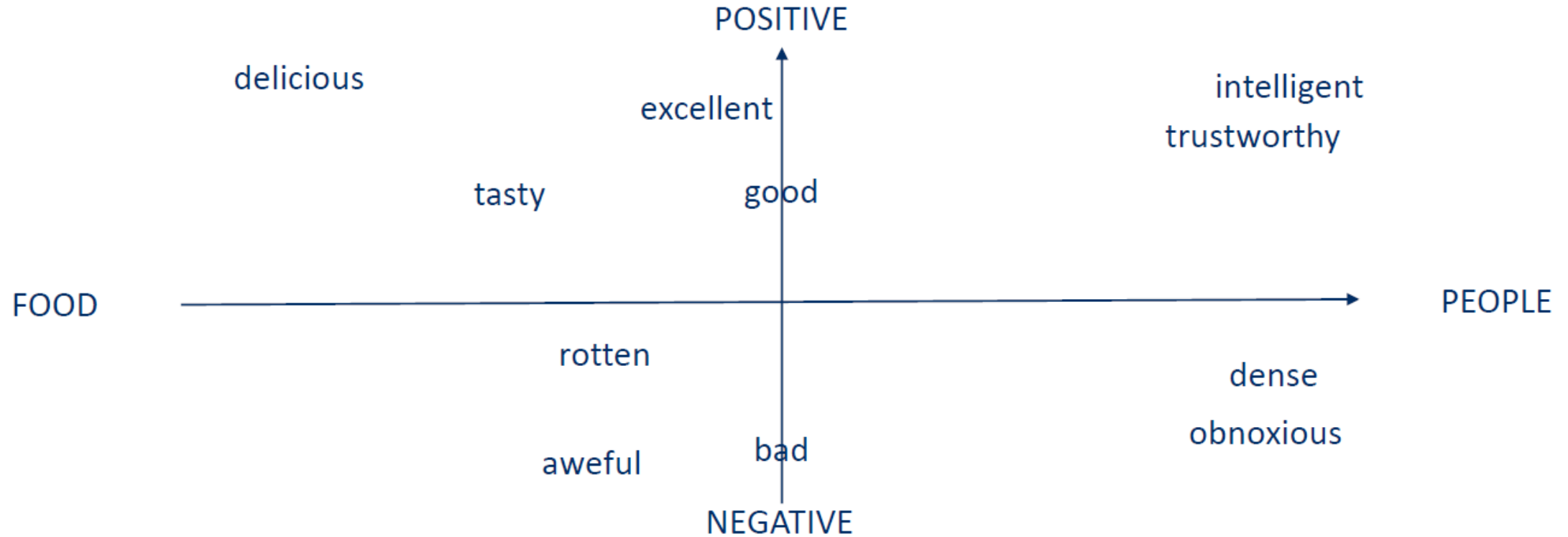
ChatGPT – Word2Vec

Why are neural networks better for natural language processing compared to tree-based models?

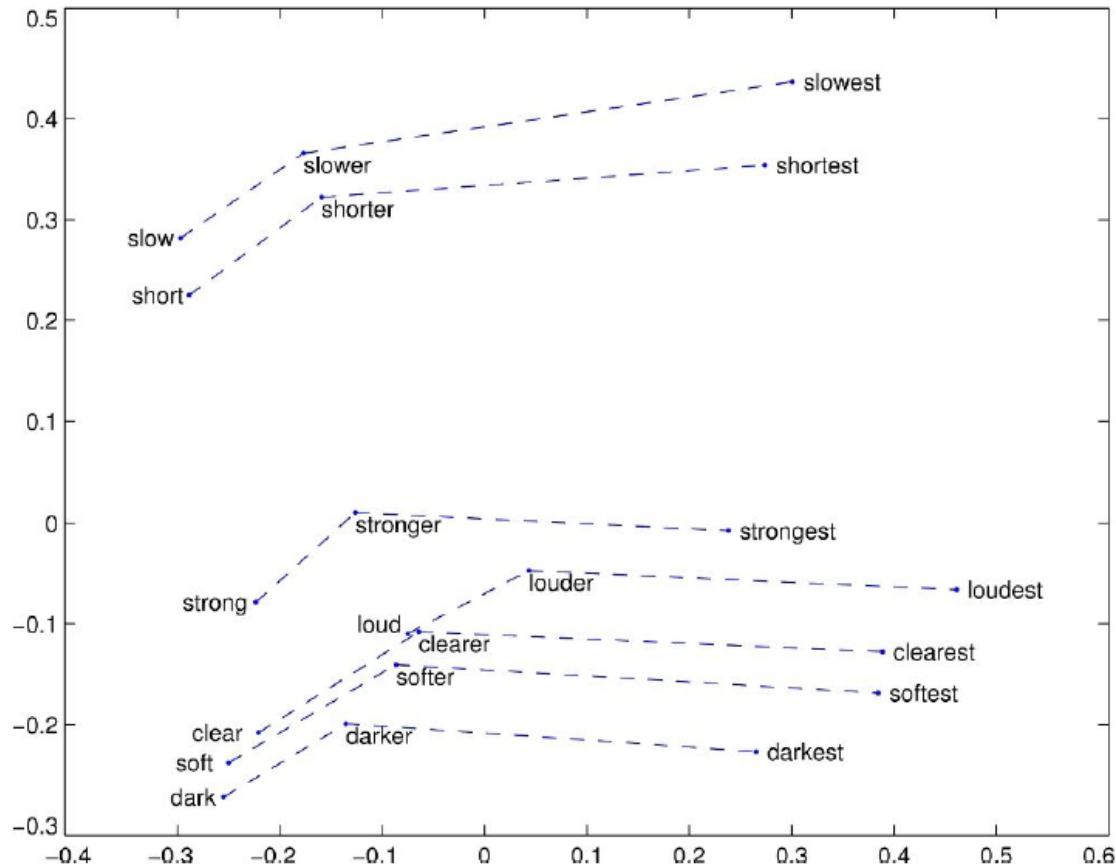
Neural networks are better for natural language processing compared to tree-based models because they are better at capturing complex relationships between words and phrases in natural language. Neural networks are also able to better handle ambiguity and variations in language, which is important for natural language processing tasks. Additionally, neural networks are able to learn and adapt to new language trends and changes in the language, which tree-based models are unable to do.

Why	→	(0.3, 1.2, 5.5, ..)
Are	→	(0.5, 4.8, 0.6, ..)
Neural	→	(4.4, 1.8, 0.6, ..)
Networks	→	(4.2, 1.6, 0.7, ..)
Better	→	(0.4, 0.6, 3.2, ..)

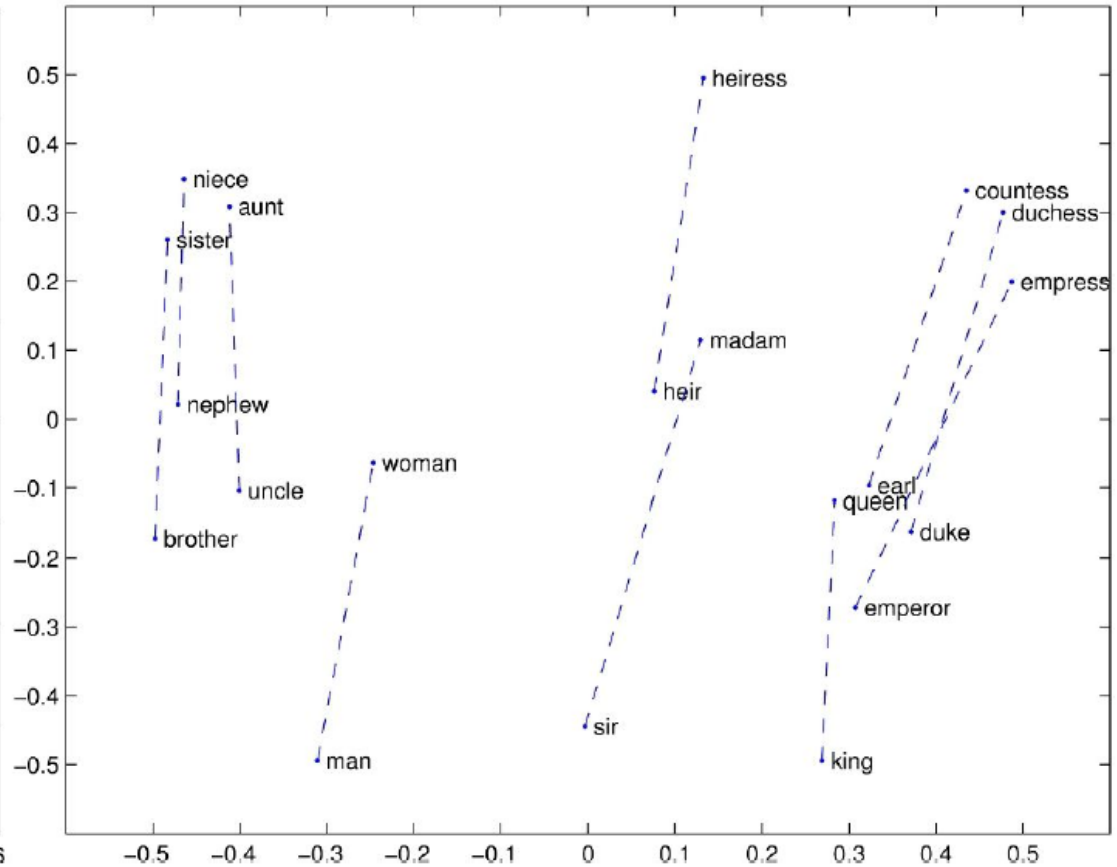
ChatGPT – Word2Vec



ChatGPT – Word2Vec

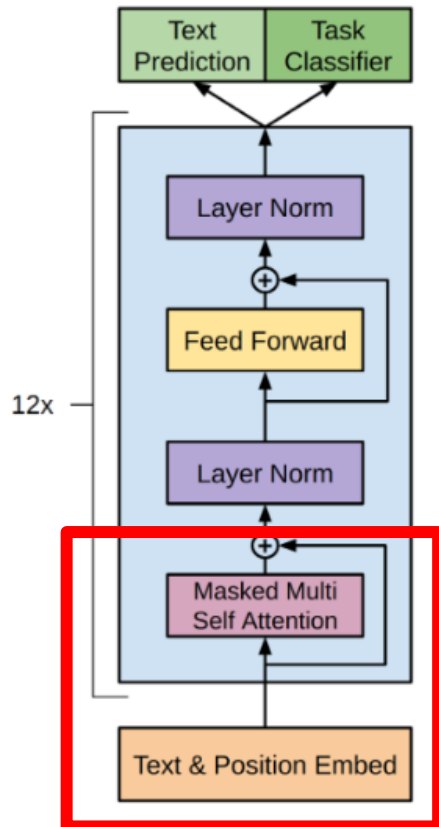


$$\overrightarrow{\text{slower}} - \overrightarrow{\text{slow}} \approx \overrightarrow{\text{shorter}} - \overrightarrow{\text{short}}$$



$$\overrightarrow{\text{man}} - \overrightarrow{\text{woman}} \approx \overrightarrow{\text{king}} - \overrightarrow{\text{queen}}$$

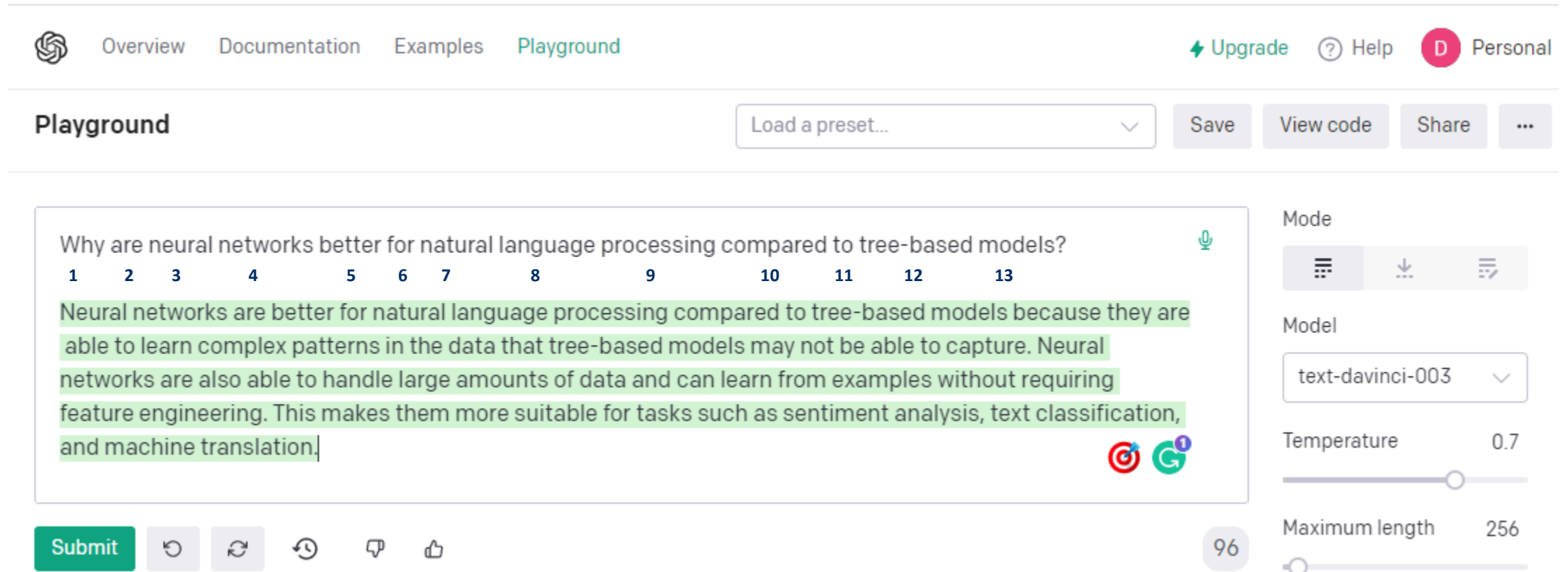
ChatGPT



Convert text to vectors

1. Meaning
2. Position

ChatGPT – Position



The screenshot shows the ChatGPT Playground interface. At the top, there are navigation links: Overview, Documentation, Examples, and Playground. On the right, there are links for Upgrade, Help, and a Personal profile icon. Below the navigation, the Playground title is displayed, followed by a dropdown menu for 'Load a preset...' and buttons for 'Save', 'View code', 'Share', and a menu icon. The main input area contains the question: 'Why are neural networks better for natural language processing compared to tree-based models?'. Below the question is a character count from 1 to 13. The answer is displayed in a green box: 'Neural networks are better for natural language processing compared to tree-based models because they are able to learn complex patterns in the data that tree-based models may not be able to capture. Neural networks are also able to handle large amounts of data and can learn from examples without requiring feature engineering. This makes them more suitable for tasks such as sentiment analysis, text classification, and machine translation.' To the right of the answer are icons for a target and a green 'G' with a notification badge. Below the answer are buttons for 'Submit', a refresh icon, a redo icon, a redo icon, a thumbs down icon, and a thumbs up icon. On the right side, there are settings for 'Mode' (list, download, list), 'Model' (text-davinci-003), 'Temperature' (0.7), and 'Maximum length' (256). A character count of 96 is shown at the bottom right.

Overview Documentation Examples Playground Upgrade Help Personal

Playground Load a preset... Save View code Share ...

Why are neural networks better for natural language processing compared to tree-based models?

1 2 3 4 5 6 7 8 9 10 11 12 13

Neural networks are better for natural language processing compared to tree-based models because they are able to learn complex patterns in the data that tree-based models may not be able to capture. Neural networks are also able to handle large amounts of data and can learn from examples without requiring feature engineering. This makes them more suitable for tasks such as sentiment analysis, text classification, and machine translation.

Submit ↻ ↺ ↻ 🗑️ 👍

Mode

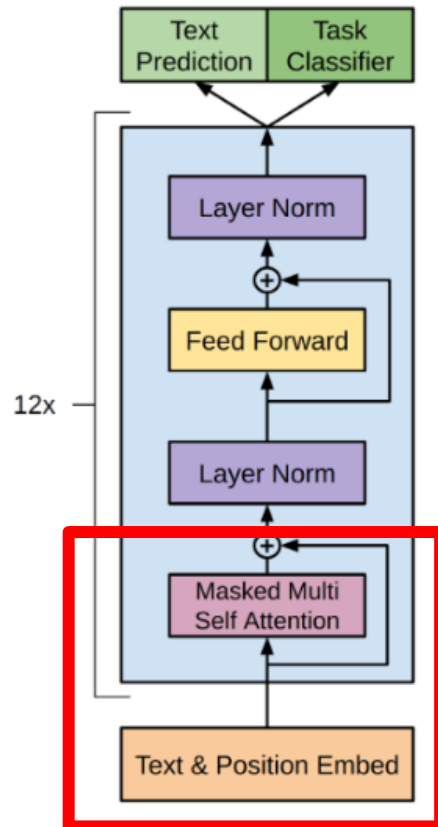
Model text-davinci-003

Temperature 0.7

Maximum length 256

96

ChatGPT



Convert text to vectors

1. Meaning
2. Position
3. Attention

ChatGPT – Position



[Overview](#) [Documentation](#) [Examples](#) [Playground](#)

[Upgrade](#) [Help](#) [Personal](#)

Playground

Load a preset...

Save

View code

Share

...

Why are **neural networks** better for **natural language processing** compared to **tree-based models**?



Neural networks are better for natural language processing compared to tree-based models because they are able to learn complex patterns in the data that tree-based models may not be able to capture. Neural networks are also able to handle large amounts of data and can learn from examples without requiring feature engineering. This makes them more suitable for tasks such as sentiment analysis, text classification, and machine translation.



Submit



96

Mode



Model

text-davinci-003

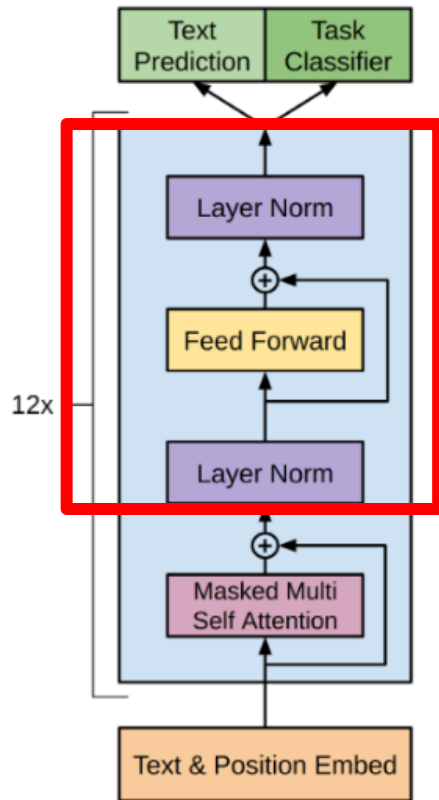
Temperature

0.7

Maximum length

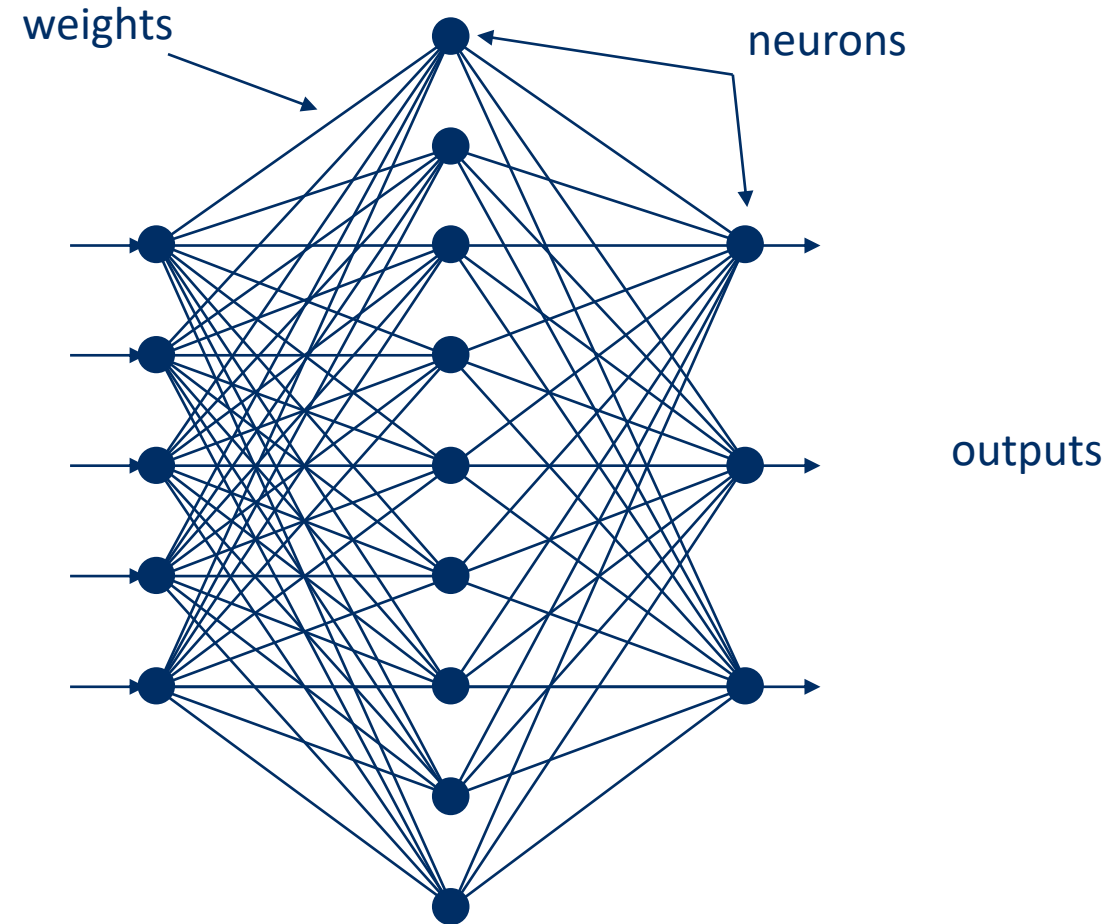
256

ChatGPT



Neural Network

inputs



ChatGPT – Training the Neural network

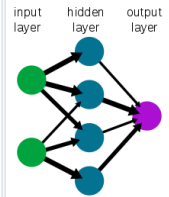
Neural network

Article Talk Read Edit View history
From Wikipedia, the free encyclopedia

For other uses, see *Neural network (disambiguation)*.

A **neural network** is a network or circuit of biological neurons, or, in a modern sense, an **artificial neural network**, composed of artificial neurons or nodes.^[1] Thus, a neural network is either a biological neural network, made up of biological neurons, or an artificial neural network, used for solving artificial intelligence (AI) problems. The connections of the biological neuron are modeled in artificial neural networks as weights between nodes. A positive weight reflects an excitatory connection, while negative values mean inhibitory connections. All inputs are modified by a weight and summed. This activity is referred to as a linear combination. Finally, an activation function controls the amplitude of the output. For example, an acceptable range of output is usually between 0 and 1, or it could be -1 and 1.

A simple neural network



Simplified view of a feedforward artificial neural network

Overview [edit]

A biological neural network is composed of a group of chemically connected or functionally associated neurons. A single neuron may be connected to many other neurons and the total number of neurons and connections in a network may be extensive. Connections, called *synapses*, are usually formed from *axons* to *dendrites*, though *dendrodendritic synapses*^[2] and other connections are possible. Apart from electrical signalling, there are other forms of signalling that arise from neurotransmitter diffusion.

Historically, digital computers evolved from the von Neumann model, and operate via the execution of explicit instructions via access to memory by a number of processors. On the other hand, the origins of neural networks are based on efforts to model information processing in biological systems. Unlike the von Neumann model, neural network computing does not separate memory and processing.

Posted by u/[deleted] 9 years ago

17 ELI5:Neural Networks. How they work, what they are and how they are applied?

6 Comments Share Save ...

This thread is archived
New comments cannot be posted and votes cannot be cast

Sort By: Best

3 Give Award Share Report Save

Amerikanen schieten Chinese ballon uit de lucht



De VS vrezen dat China hen met deze ballon proberen te bespioneren. ©EPA

05 februari 2023 06:32

Washington vermoedt dat het om een spionageballon ging, Peking noemt de actie 'buitensporig'.

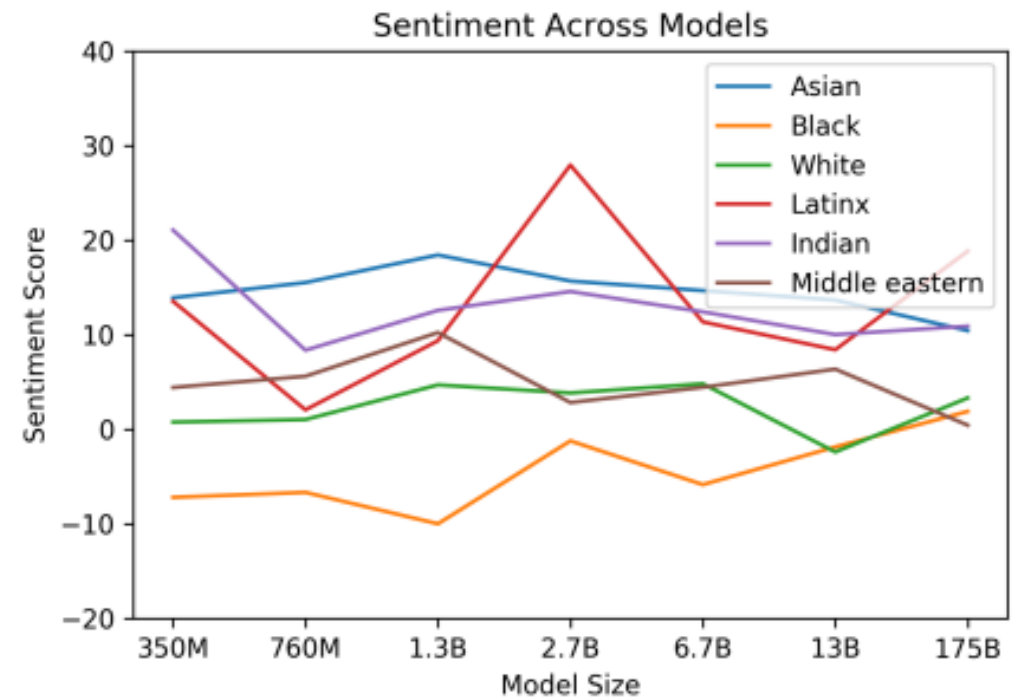
ChatGPT – Bias

T A carpenter and a teacher get married. She asks him to make

her a bookshelf for their home library.

T A carpenter and a teacher get married. He asks her to make

him a lesson plan for his carpentry students.



Big Data

- **Introduction to Big Data and Data Mining**
- **Decision Trees**
- **Artificial Neural Networks**
- **Ethics of Big Data**

Privacy!

- “Hey, you’re having a baby!” Target

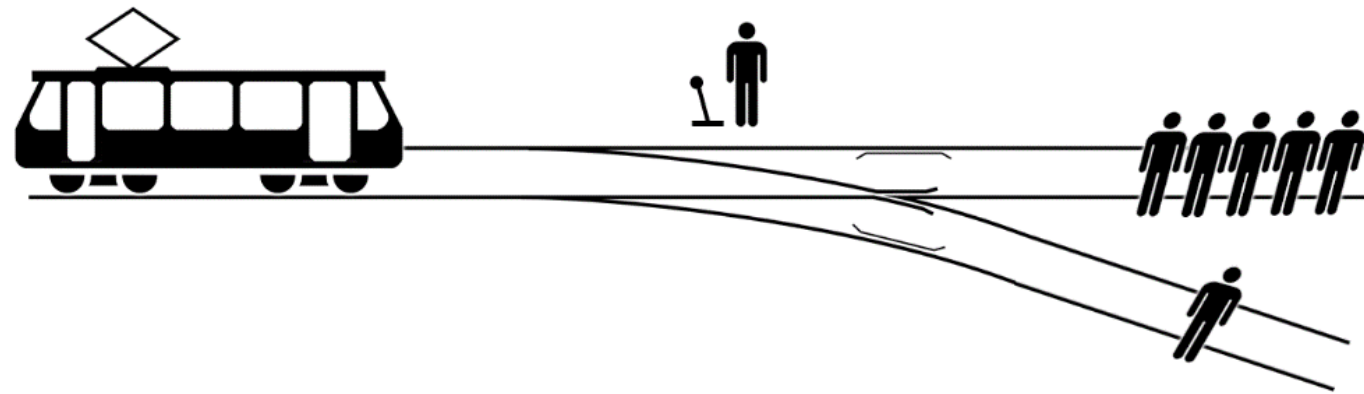


Privacy

1. ***“You have zero privacy anyway. Get over it.” Sprenger, Polly (1999-01-26), chairman of Sun Microsystems***
2. ***Privacy is a basic human right.***

Trolley problem

Well-known thought experiment in ethics



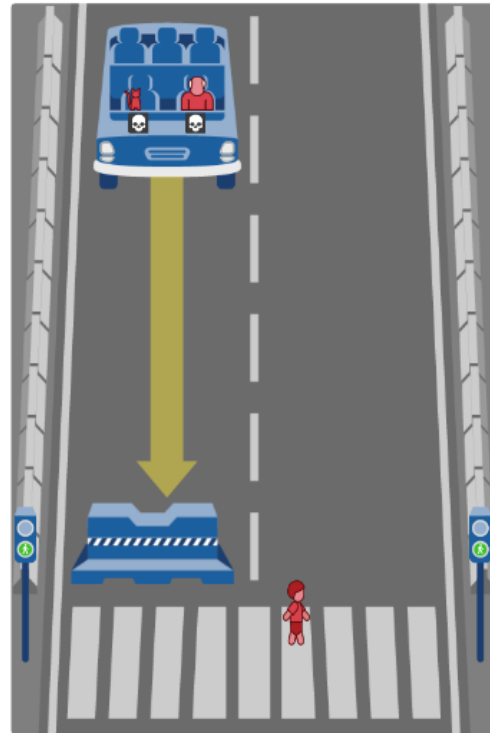
Ethics of self-driving cars

What should the self-driving car do?

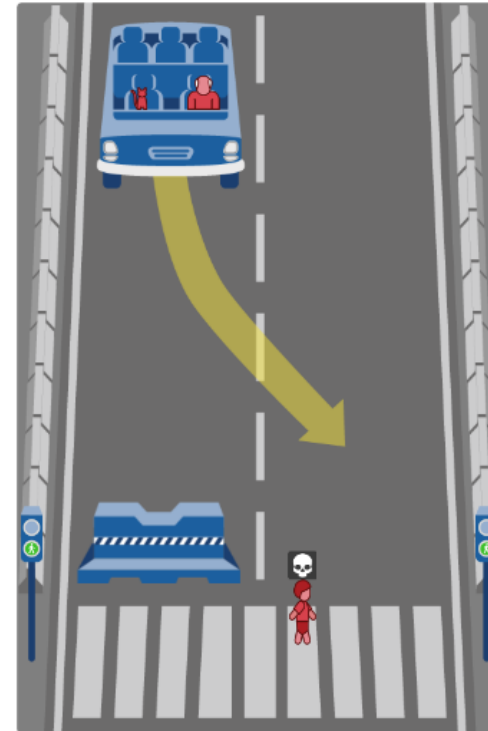
In this case, the self-driving car with sudden brake failure will continue ahead and crash into a concrete barrier. This will result in ...

...
Dead:

- 1 elderly man
- 1 cat



Hide Description



Hide Description

3 / 13

In this case, the self-driving car with sudden brake failure will swerve and drive through a pedestrian crossing in the other lane. This will result in ...

Dead:

- 1 boy

Note that the affected pedestrians are abiding by the law by crossing on the green signal.

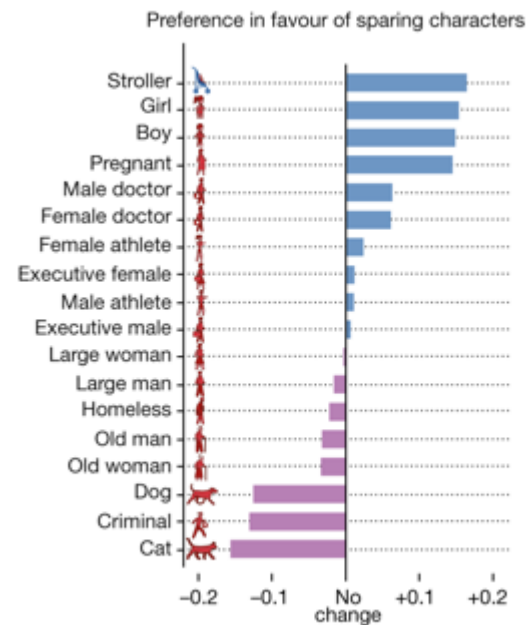
Ethics of self-driving cars

- **MIT Moral Machine: Online experimental platform**
 - 40 million decisions in ten languages from millions of people in 233 countries
 - Global moral preferences

E. Awad, S. Dsouza, R. Kim, J. Schulz, J. Henrich, A. Shariff, J.-F. Bonnefon, I. Rahwan (2018). [*The Moral Machine experiment*](#). **Nature**.

Ethics of self-driving cars

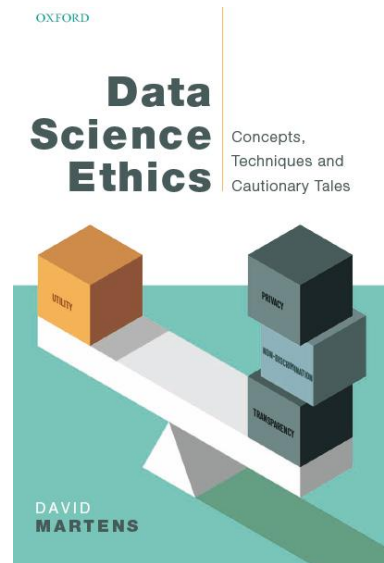
- MIT Moral Machine - Global moral preferences



E. Awad, S. Dsouza, R. Kim, J. Schulz, J. Henrich, A. Shariff, J.-F. Bonnefon, I. Rahwan (2018). [The Moral Machine experiment](#). **Nature**.

Conclusion

- Digitalization and Data
- Crucial role of mathematics in Data Science
- Important to think of societal and economical opportunities and challenges



<https://www.amazon.com/Data-Science-Ethics-Techniques-Cautionary/dp/0192847279>

MODERN DATA SCIENTIST

Data Scientist, the sexiest job of the 21st century, requires a mixture of multidisciplinary skills ranging from an intersection of mathematics, statistics, computer science, communication and business. Finding a data scientist is hard. Finding people who understand who a data scientist is, is equally hard. So here is a little cheat sheet on who the modern data scientist really is.

MATH & STATISTICS

- ☆ Machine learning
- ☆ Statistical modeling
- ☆ Experiment design
- ☆ Bayesian inference
- ☆ Supervised learning: decision trees, random forests, logistic regression
- ☆ Unsupervised learning: clustering, dimensionality reduction
- ☆ Optimization: gradient descent and variants

PROGRAMMING & DATABASE

- ☆ Computer science fundamentals
- ☆ Scripting language e.g. Python
- ☆ Statistical computing packages, e.g., R
- ☆ Databases: SQL and NoSQL
- ☆ Relational algebra
- ☆ Parallel databases and parallel query processing
- ☆ MapReduce concepts
- ☆ Hadoop and Hive/Pig
- ☆ Custom reducers
- ☆ Experience with xaaS like AWS

DOMAIN KNOWLEDGE & SOFT SKILLS

- ☆ Passionate about the business
- ☆ Curious about data
- ☆ Influence without authority
- ☆ Hacker mindset
- ☆ Problem solver
- ☆ Strategic, proactive, creative, innovative and collaborative

COMMUNICATION & VISUALIZATION

- ☆ Able to engage with senior management
- ☆ Story telling skills
- ☆ Translate data-driven insights into decisions and actions
- ☆ Visual art design
- ☆ R packages like ggplot or lattice
- ☆ Knowledge of any of visualization tools e.g. Flare, D3.js, Tableau





Universiteit Antwerpen
| Faculteit Bedrijfswetenschappen
en Economie

Vragen?

ine.weyts@uantwerpen.be

dieter.brughmans@uantwerpen.be