

ÜK 259

ICT-Lösungen mit Machine
Learning entwickeln







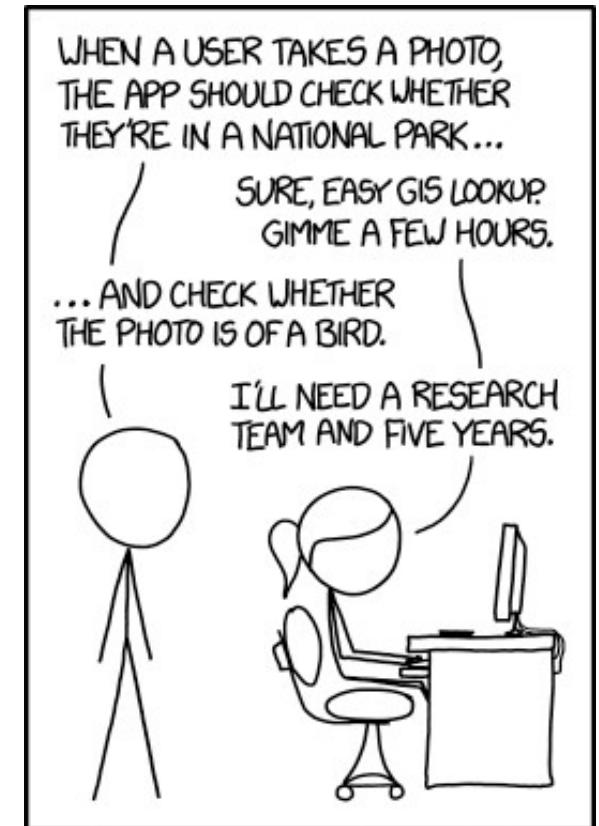




What is Machine Learning?

Discussion

- What is **Intelligence**?
- What is **Artificial Intelligence**?
- What's the difference between **machine learning** and conventional **data analysis**?



IN CS, IT CAN BE HARD TO EXPLAIN
THE DIFFERENCE BETWEEN THE EASY
AND THE VIRTUALLY IMPOSSIBLE.

What is Machine Learning?

- **What is Intelligence?**

“Ability to process information and use it to solve a problem”

- **What is Artificial Intelligence?**

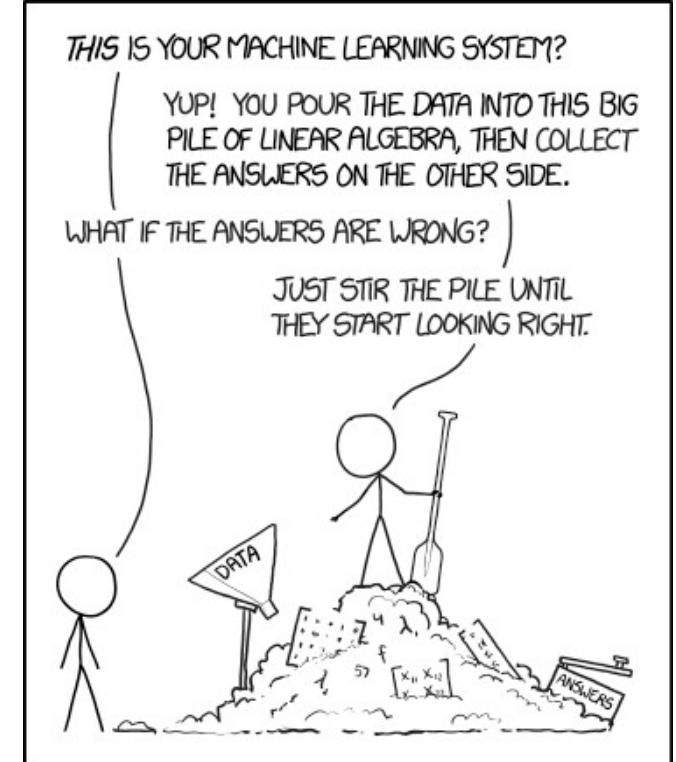
“Algorithms that imitate decision structures of humans to (somewhat) independently solve problems”



IN CS, IT CAN BE HARD TO EXPLAIN THE DIFFERENCE BETWEEN THE EASY AND THE VIRTUALLY IMPOSSIBLE.

Machine Learning is...

- “A lot of Maths”
- A set of algorithms that **improve automatically through trial and error** and by the use of data
- Extracting valuable information from data
- A method to solve a problem that we don’t know the solution for (yet)



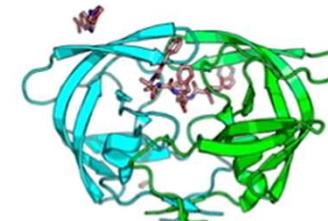
Use Cases



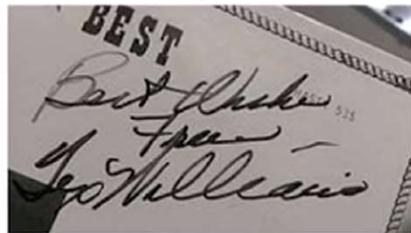
AlphaGo



Recommendation systems



Drug discovery



Character recognition

20 TWO SIGMA

Hedge fund stock predictions



Voice assistants



Assisted driving



Face detection/recognition

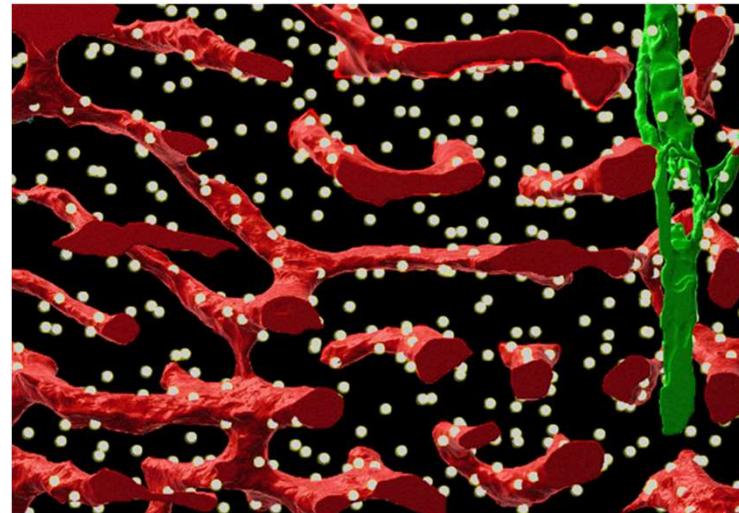
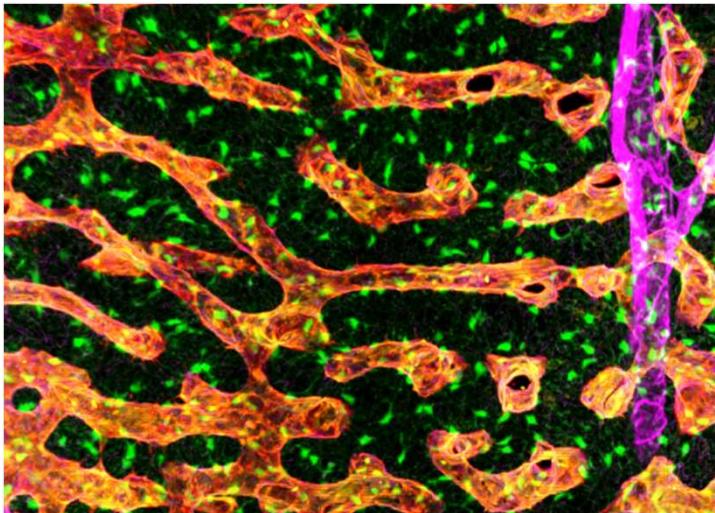


Cancer diagnosis

What can ML do?

Some examples

- DALL-E 2
- This Person does not exist: thispersondoesnotexist.com
- Github CoPilot: https://copilot.github.com/
- MINTIF: 3D microscopy segmentation



üK 259

Goals

- Understand how machines “learn”
- Get a feeling for common problems and workflows in data analysis
- Get an overview of current ML methods and their applications
- Gather first experiences in applying those methods on real world problems

üK 259

HANOKS

- 1.1 Kennt die verschiedenen Technologien im Machine Learning Umfeld und deren Anwendungsgebiete.
- 1.2 Kennt Lösungsvarianten und den Mehrwert für ICT-Lösungen im Vergleich zu bestehenden Lösungen
- 1.3 Kennt die Kategorien des Machine Learning und wählt aus diesen die geeignete Technologie für eine ICT-Lösung aus.
- 1.4 Kennt Modelle und Verfahren im Machine Learning Umfeld
- 2.1 Kennt die gesetzlichen Kriterien zur Bestimmung schützenswerter Daten.
- 2.2 Kennt Massnahmen zur Gewährleistung des Datenschutzes bei der Nutzung und Verarbeitung von schützenswerten Daten.
- 3.1 Kennt die Eigenschaften von Daten und die Vorgehensweise, zur Extraktion von Features für eine Datenanalyse.
- 3.2 Kennt die drei Datentypen und deren Eigenschaften.
- 3.3 Kennt die Möglichkeiten, Daten für eine weitere Verarbeitung aufzubereiten

üK 259

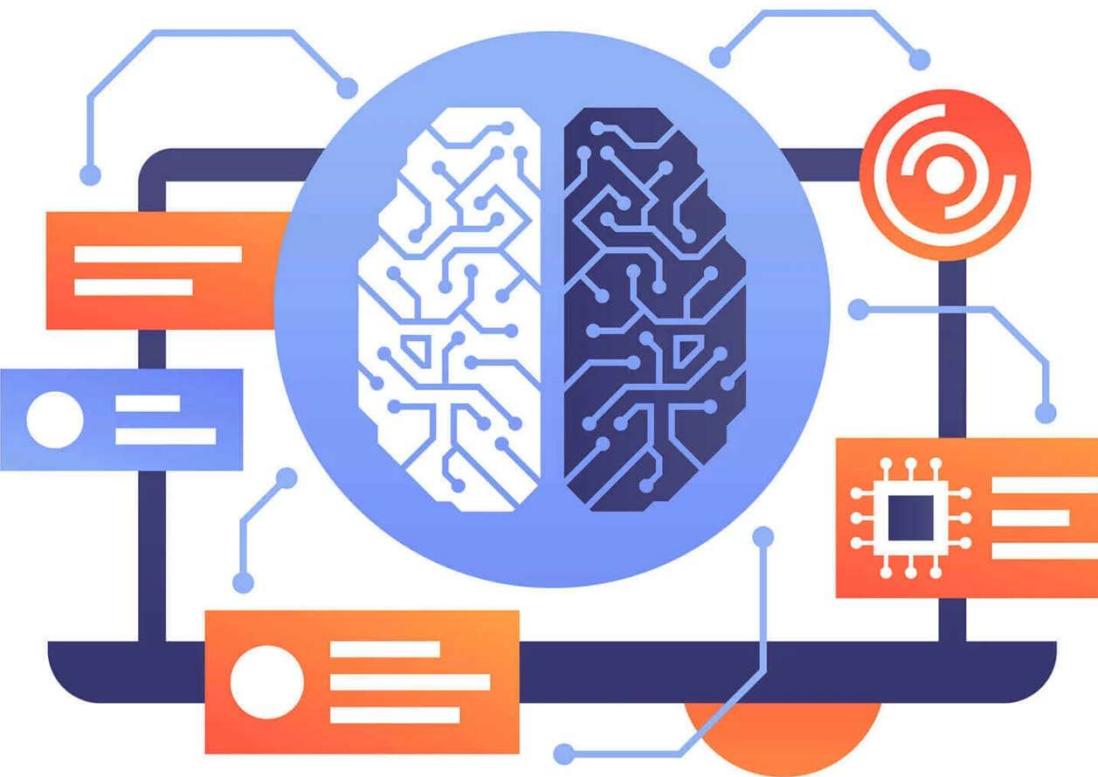
HANOKS

- 4.1 Kennt den Ablauf der Entwicklung einer ICT-Lösung mit Machine Learning gemäss den folgenden Schritten: Zieldefinition, Datenbeschaffung und Aufbereitung der Daten, Lernphase, Interpretation der Ergebnisse und produktivem Einsatz.
- 4.2 Kennt das Vorgehen zum Trainieren und Testen eines Modells
- 5.1 Kennt die Wahrheitsmatrix und deren Funktion.
- 5.2 Kennt die statistischen Gütekriterien zur Beurteilung der Wahrheitsmatrix.
- 6.1 Kennt die erforderlichen Zugriffsmechanismen und die benötigten Schnittstellen eines Machine Learning Dienstes.
- 6.2 Kennt die Komponenten und Dienste sowie das Vorgehen zur Konfiguration nach Vorgaben.
- 7.1 Kennt eine Programmierumgebung und deren Framework zur Entwicklung einer ICT-Lösung mit Machine Learning.
- 7.2 Kennt verschiedene Pipelines zur Entwicklung eines Machine Learning Modells.

ÜK 259

Topics

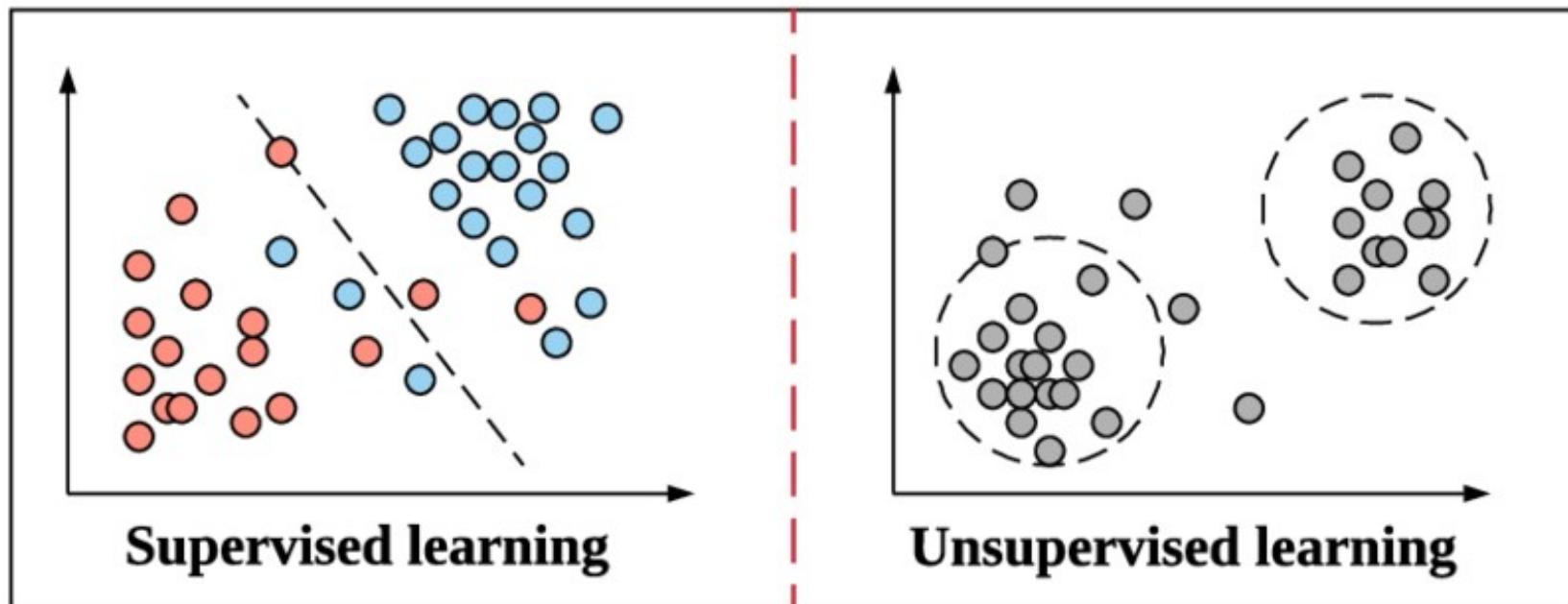
- Day 1: Data Handling & Visualization
- Day 2: Unsupervised Machine Learning
- Day 3: Supervised Machine Learning
- Day 4: Deep Learning & Neural Networks



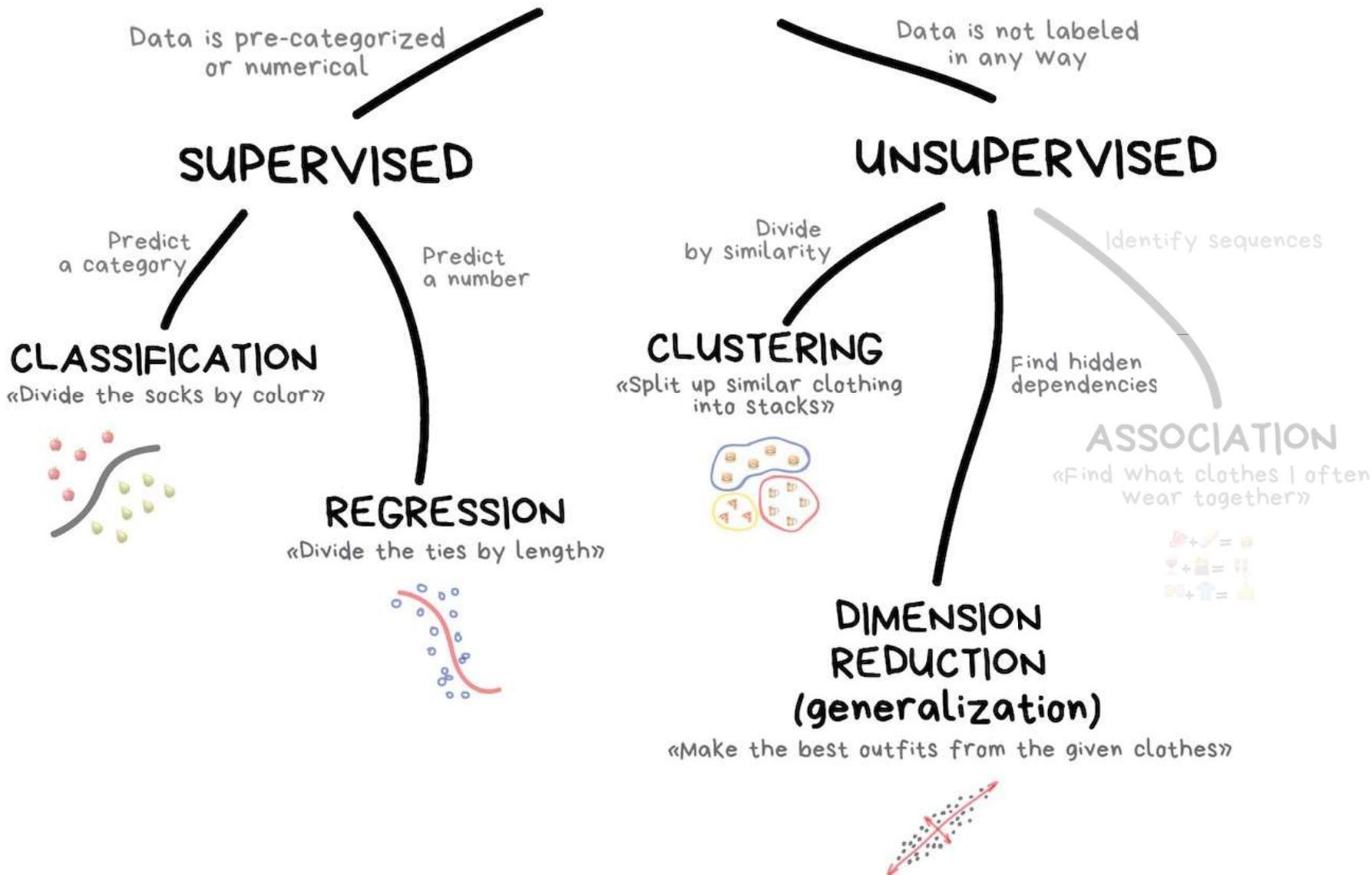
”

Main Concepts of Machine Learning

Types of Machine Learning



CLASSICAL MACHINE LEARNING



JNG

Main Types of Machine Learning

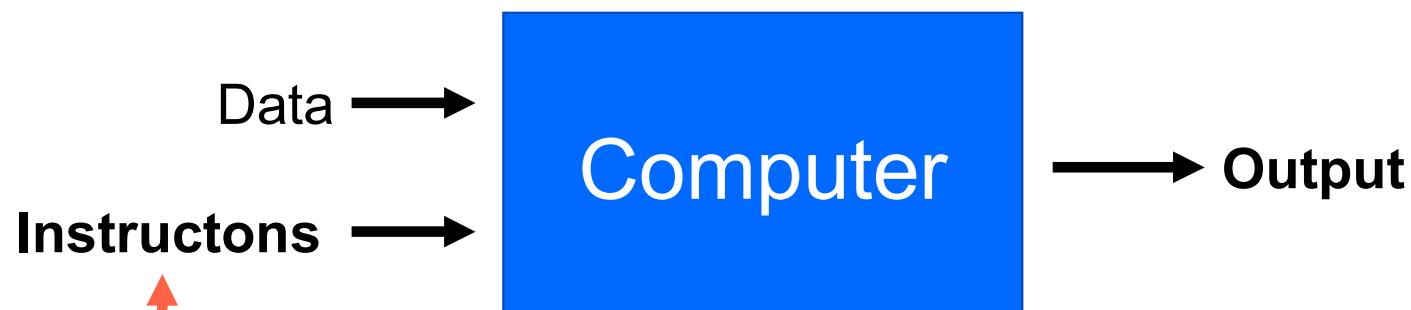
- **Unsupervised ML:** Find patterns & trends in the data, without any prior knowledge.
- **Supervised ML:** Given a dataset and the desired output, determine the best Algorithm to predict the output from the data.
- (Reinforcement Learning): Given a set of rules, develop a strategy to maximize a score

How do Neural Networks fit into this?

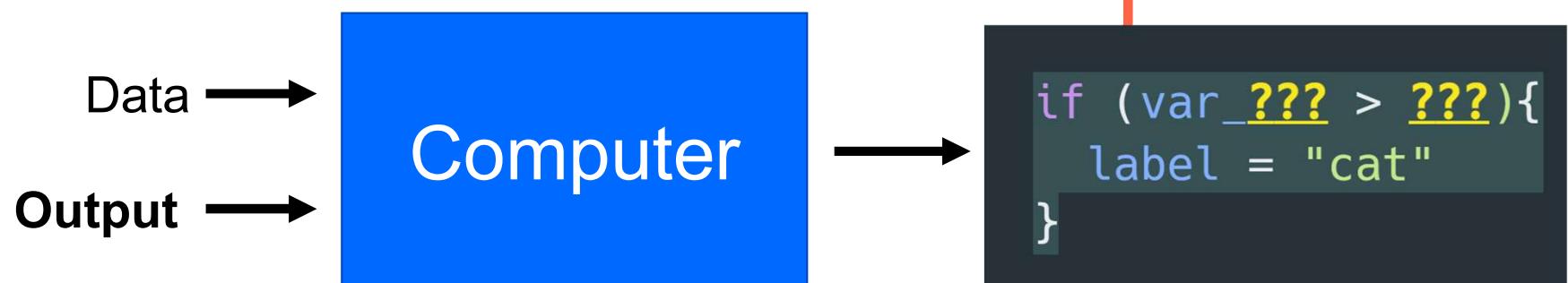
Neural Networks are a class of algorithms that can be used for supervised and unsupervised ML

Traditional Programming vs. ML

Traditional Programming



Machine Learning



What is a Model?

- A model is a set of instructions (Algorithm) used to come to a decision based on input data
 - It represents the general set of solutions / strategies used to solve the problem (for example a type of neural network)
 - A “trained” model is such an algorithm with tuned variables to solve a specific problem.



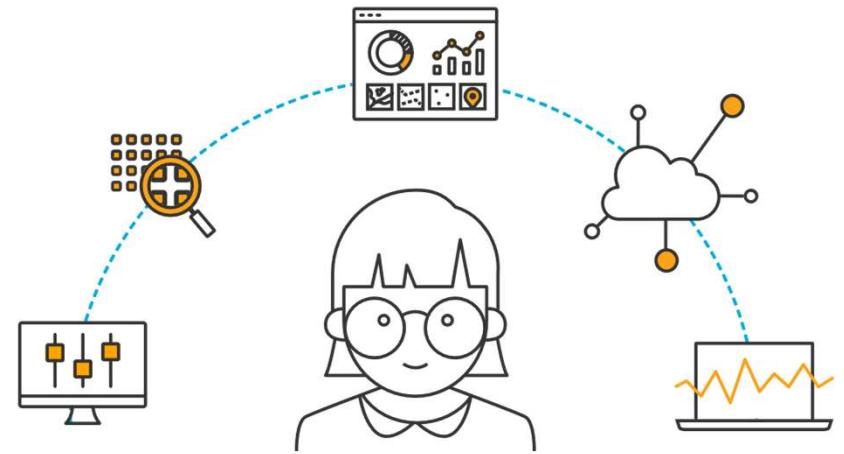
Tools for Machine Learning

- Python
- Anaconda
 - Numpy
 - Pandas
 - Matplotlib & Seaborn
 - Scikit Learn
 - Tensorflow
- Jupyter (Lab)



But First: Data Handling!

- Identifying types of data
- Data cleaning
- Standardization
- Data Augmentation
- Visualization





”

Data Handling

Data Handling

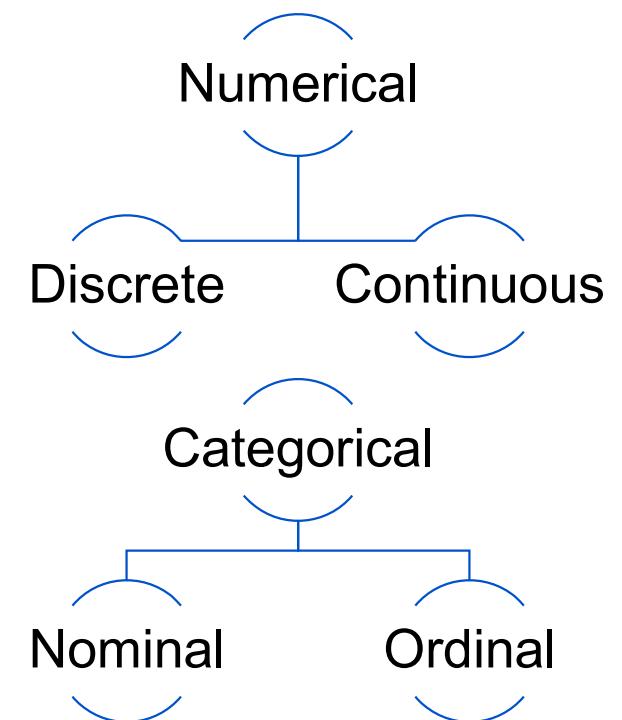
Common Tasks

- Identifying types of data
- Data cleaning
- Standardization
- Data Augmentation
- Visualization



Data Types

- **Numerical:** Numbers (duh)
 - Discrete or continuous
- **Ordinal:** Different states with a defined order
 - T-Shirt size: S < M < L
 - Low, medium, high
- **Nominal:** Multiple states without order
 - T-Shirt color
 - Gender



Data cleaning

Missing values

- Strategies for handling missing values:
 - Ignore (ಠ_ಠ)
 - Remove (lose statistical power)
 - Default values
(e.g. 0, may skew results)
 - Interpolate
(e.g. mean, max, may skew results)

BuildingArea	YearBuilt	CouncilArea
NaN	1981.0	NaN
133.0	1995.0	NaN
NaN	1997.0	NaN
157.0	1920.0	NaN
112.0	1920.0	NaN

Standardization

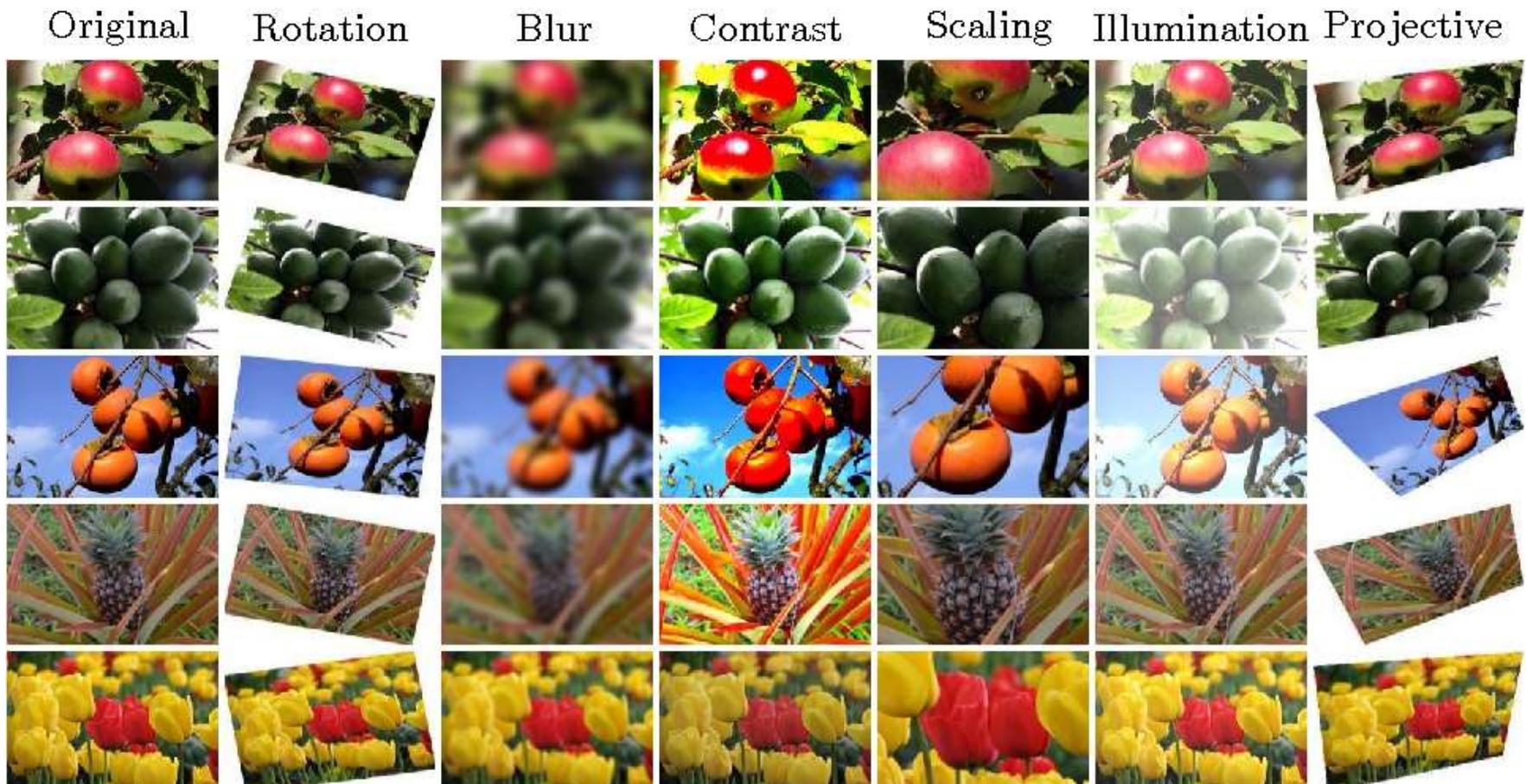
- **Problem:**
 - Features with a large scale are interpreted as having more weight (e.g. grams vs. kg).
 - Features with a large variance are interpreted as more informative.
- **Idea:** Scale features to same mean and variance:

$$\text{standard}(x_i) = \frac{x_i - \text{mean}(x)}{\text{stdev}(x)}$$

- Implemented by: `sklearn.preprocessing.StandardScaler`

Variance:
In general, how different are the values from the mean

Data Augmentation



Data Augmentation

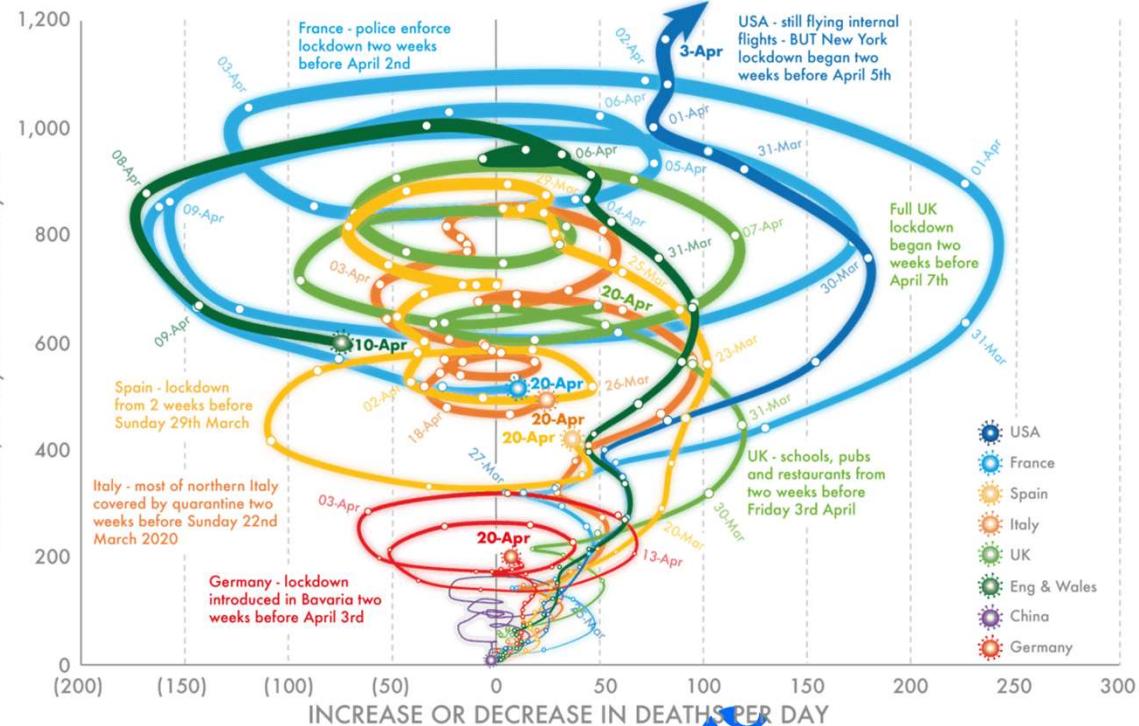
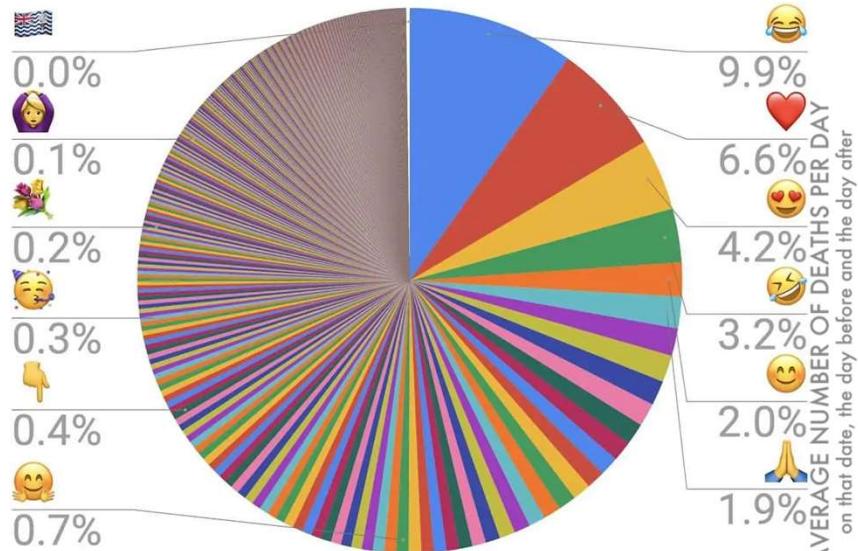
- **Idea: Modify data to augment the dataset**
- **Improving model prediction accuracy** by increasing generalizability and increasing the size of the training dataset.
 - *E.g. the model should still work with black and white images*

Generalizability:
Ability of a model to be applied to a wide variety of real world problems

Data-Visualization

Data-Visualization

How NOT to do it

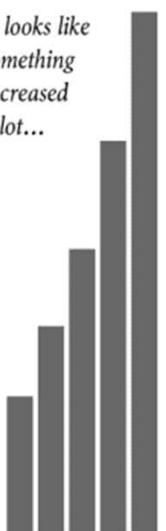


Data-Visualization

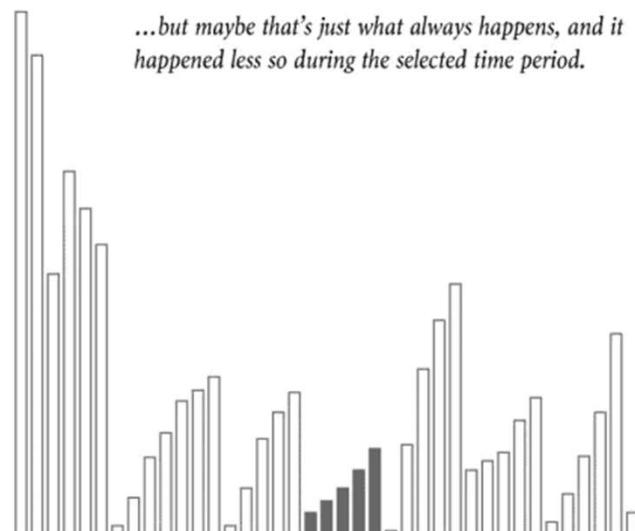
How NOT to do it

LIMITED SCOPE

It looks like something increased a lot...

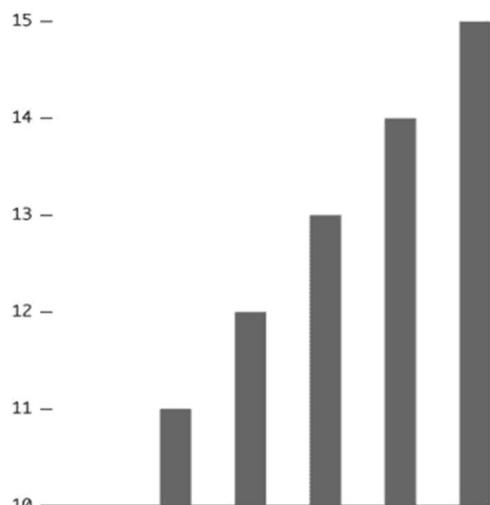


...but maybe that's just what always happens, and it happened less so during the selected time period.

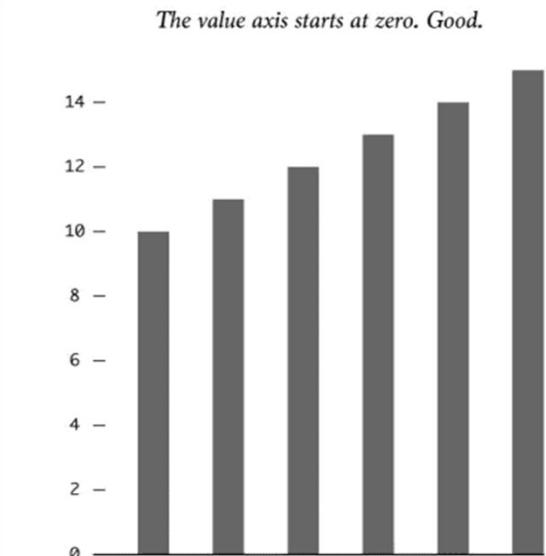


TRUNCATED AXIS

The value axis starts at ten. Liar, liar, pants on fire.



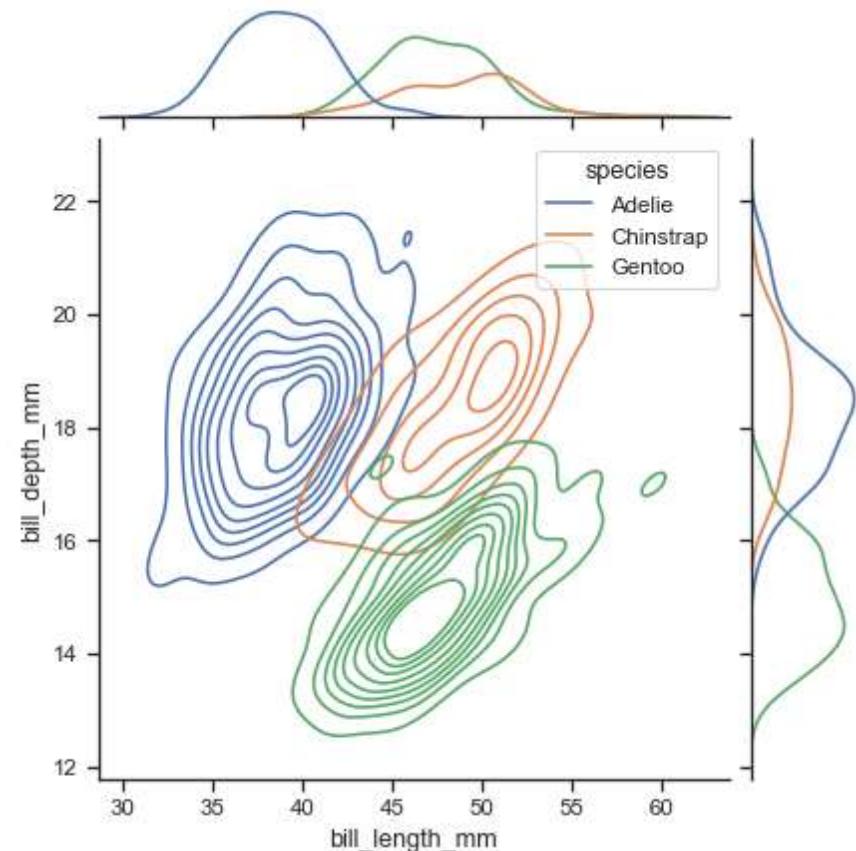
The value axis starts at zero. Good.



Data-Visualization

- **Visualizing your data should be the first and last thing you do!**
- Communicating results is a difficult but important part of Data Science
- The more complex the data the more important good and accurate data visualization becomes

<https://seaborn.pydata.org/examples/index.html>

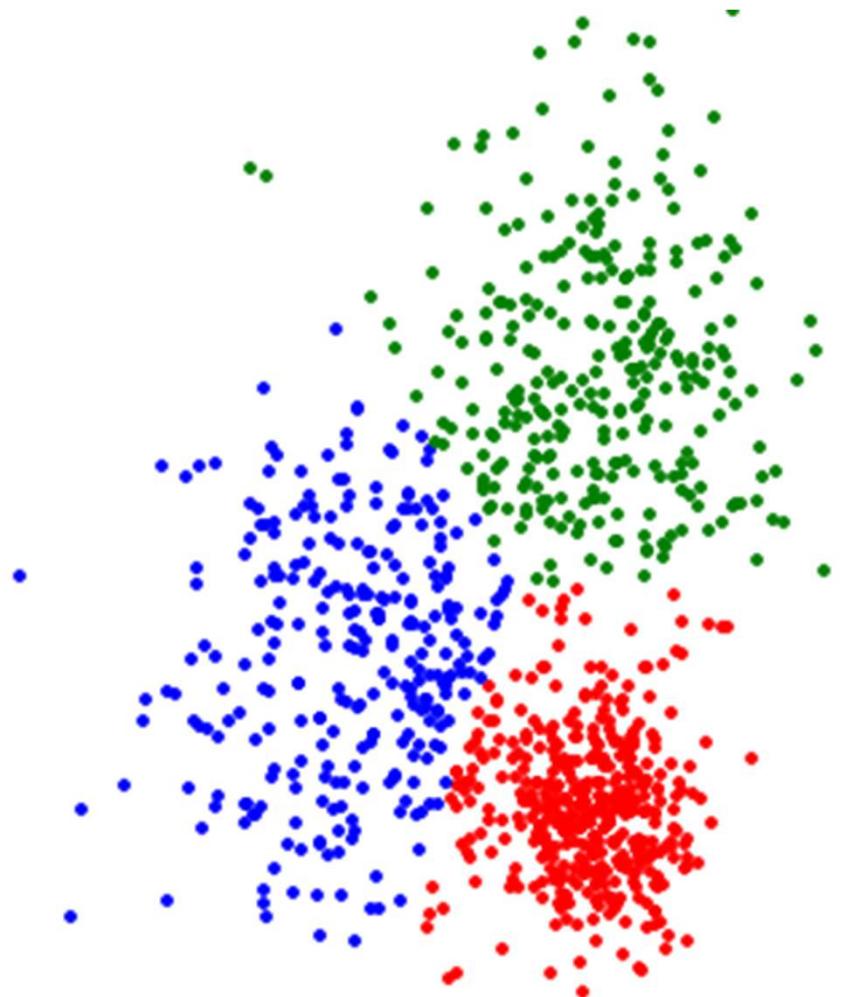


Hands-On

Part 2

Explore the Dataset “melb_data.csv” of the Melbourne Housing Market

1. Clean the dataset
2. Standardize the data
3. Think about how you would augment this dataset
4. Visualize and present an aspect of the dataset you find interesting

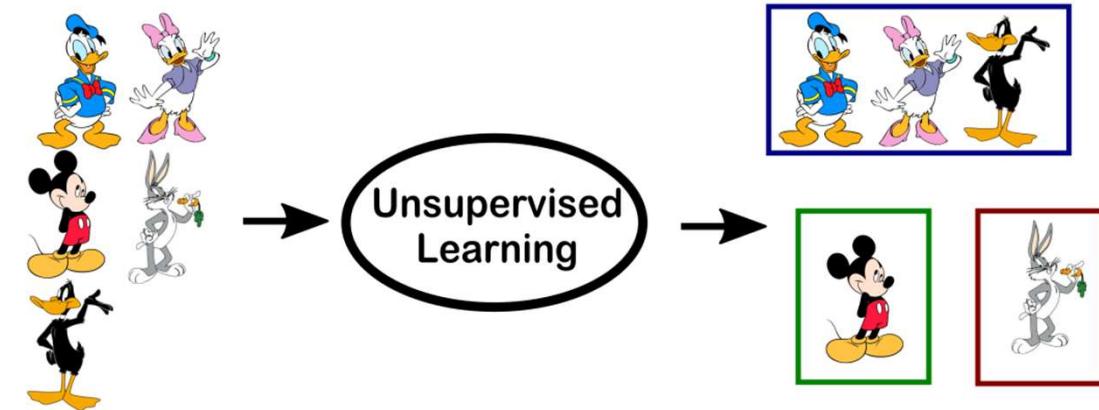


”

Unsupervised Machine Learning

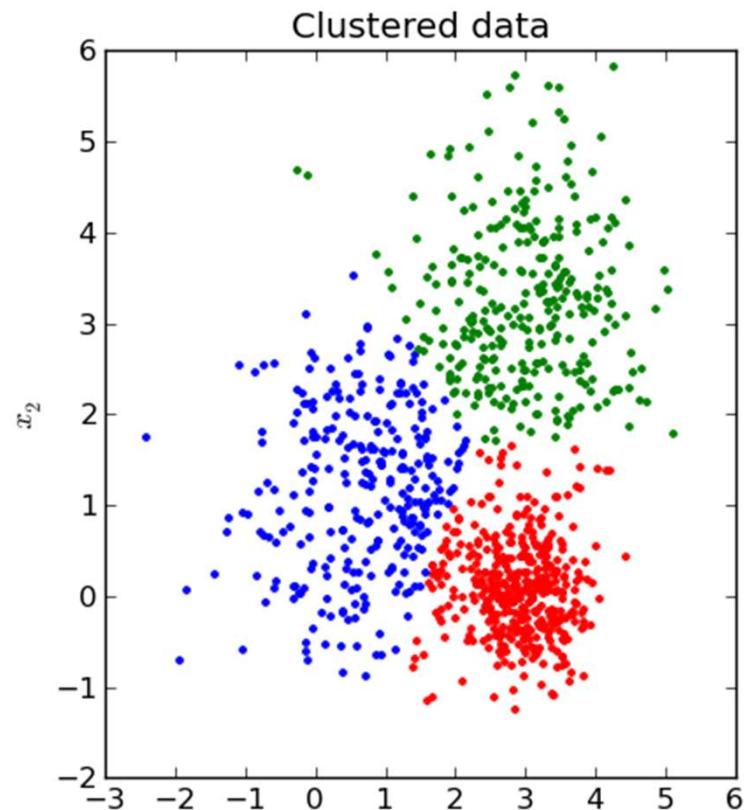
Unsupervised ML

- **Idea:** Find patterns & trends in the data, without any prior knowledge
- These patterns may give us new insights into our data
- **Main Types:**
 - Clustering
 - Dimensionality reduction



Clustering

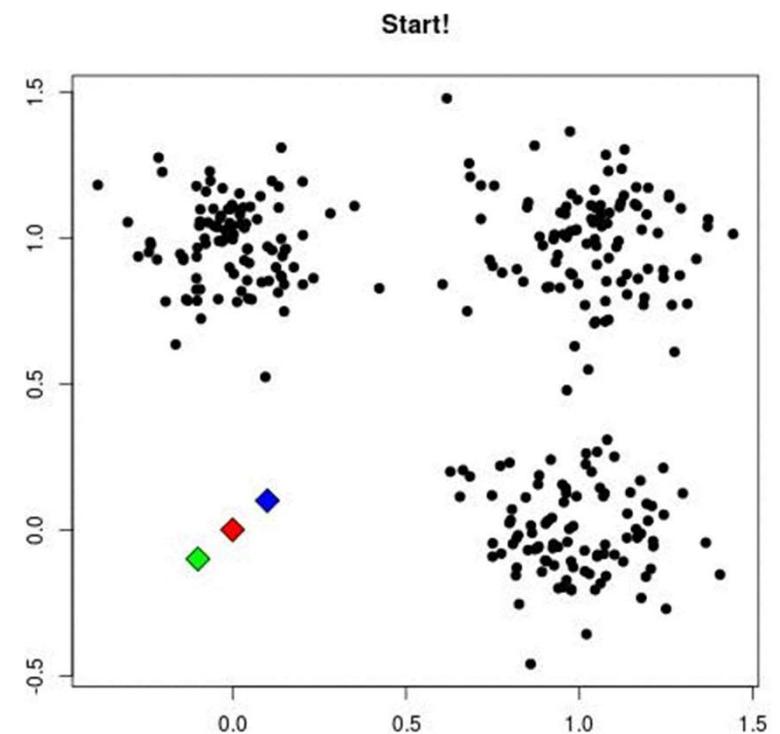
- Group datapoints into «close» groups
→ Works on some measure of similarity / distance
- **Applications:**
 - Customer segmentation
«what are the main groups in my customer base?»
 - **Recommender Systems**
«customers like you also bought...»
 - **Anomaly Detection**
«this does not look like the others»



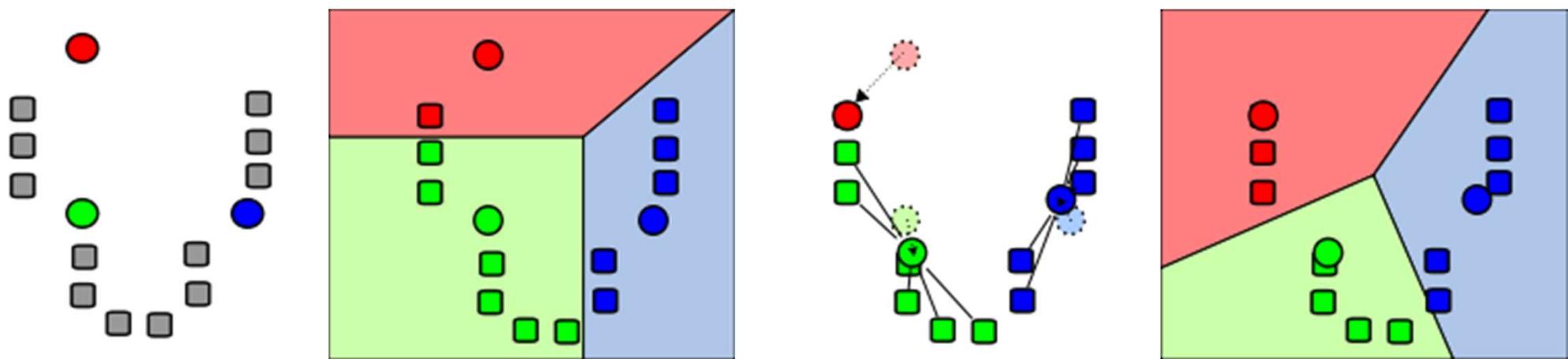
K-Means

Algorithm for Clustering

1. Initialize k random cluster-centers
2. do{
 1. Re-assign all points to closest cluster-center
 2. Recalculate cluster-centers
- } while #reassignments > 0

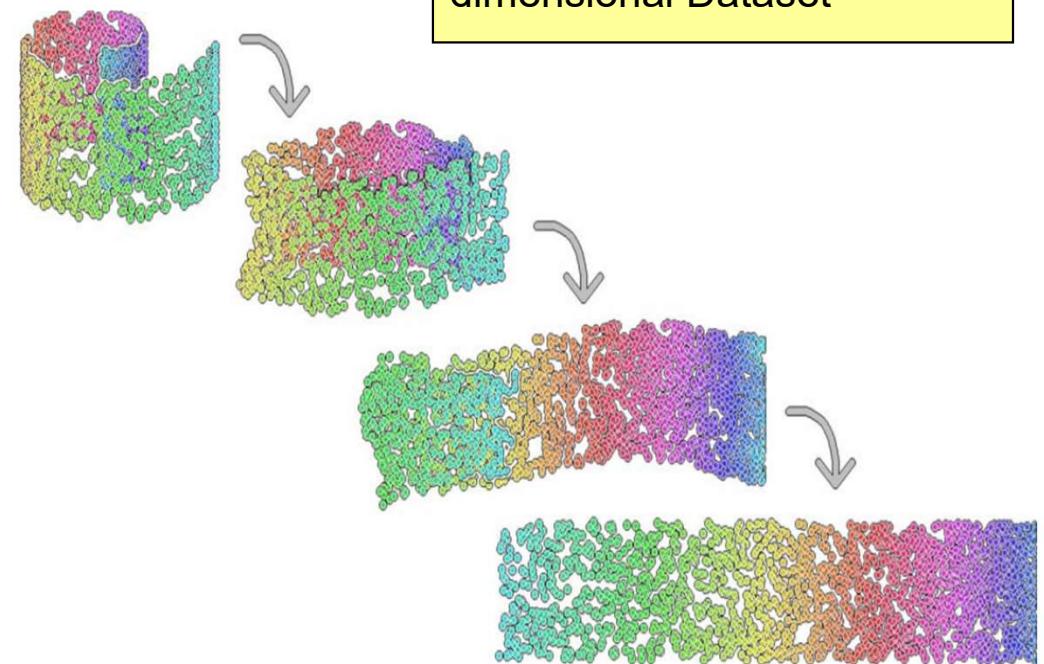


K-Means



Dimensionality Reduction

- **Idea:** represent a high-dimensional dataset in lower dimensions, while preserving local structures
- **Uses:**
 - **Data visualisation**
«how do I visualize a 10-D dataset?!»
 - **Denoising**
«real world variance vs. measurement-error»



Dimension in Data-Science / Mathematics:

One axis or column of a dataset. E.g. a Dataset with 10 Columns is a 10 dimensional Dataset

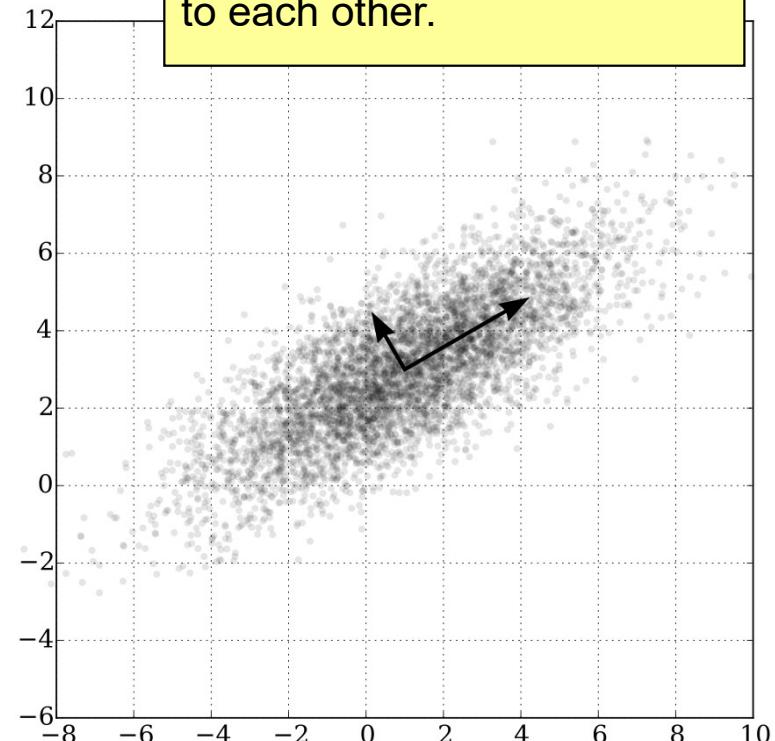
Principal Component Analysis

For dimensionality reduction

- **Idea:** Find the **principal components** that best describe the variations in the data
 - “Intuitive” explanations of PCA:
 - Shift the coordinate system such that you can discard one or more axis without loosing much information
→ PC is the main axis of variance
 - Combine multiple columns of the dataset into one in the optimal way
- Implemented in `sklearn.decomposition.PCA`

Principal Component:

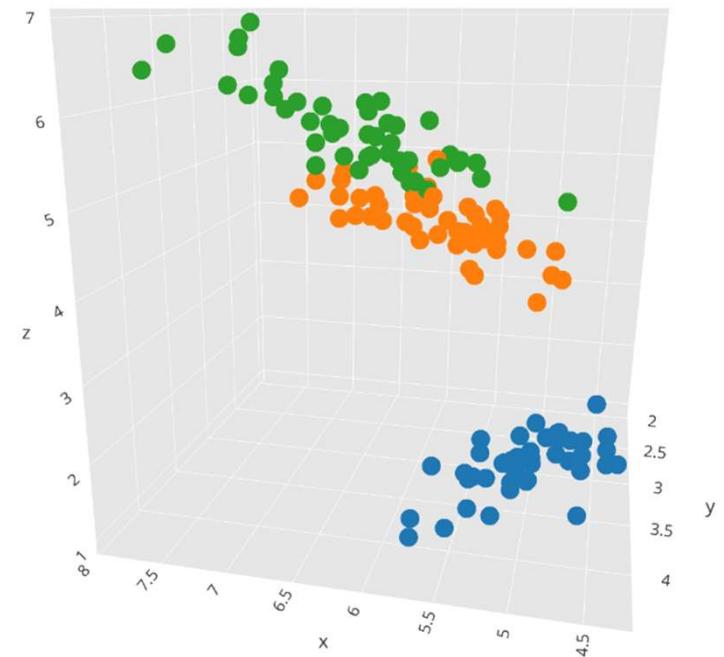
Main Axis of Variance.
Always perpendicular (90°)
to each other.



Additional Information Unsupervised Learning

- **kMeans:** <https://www.youtube.com/watch?v=mfqmoUN-Cuw>
- **PCA:** https://www.youtube.com/watch?v=_UVHneBUBW0
 - Math: <https://www.youtube.com/watch?v=PFDu9oVAE-g>

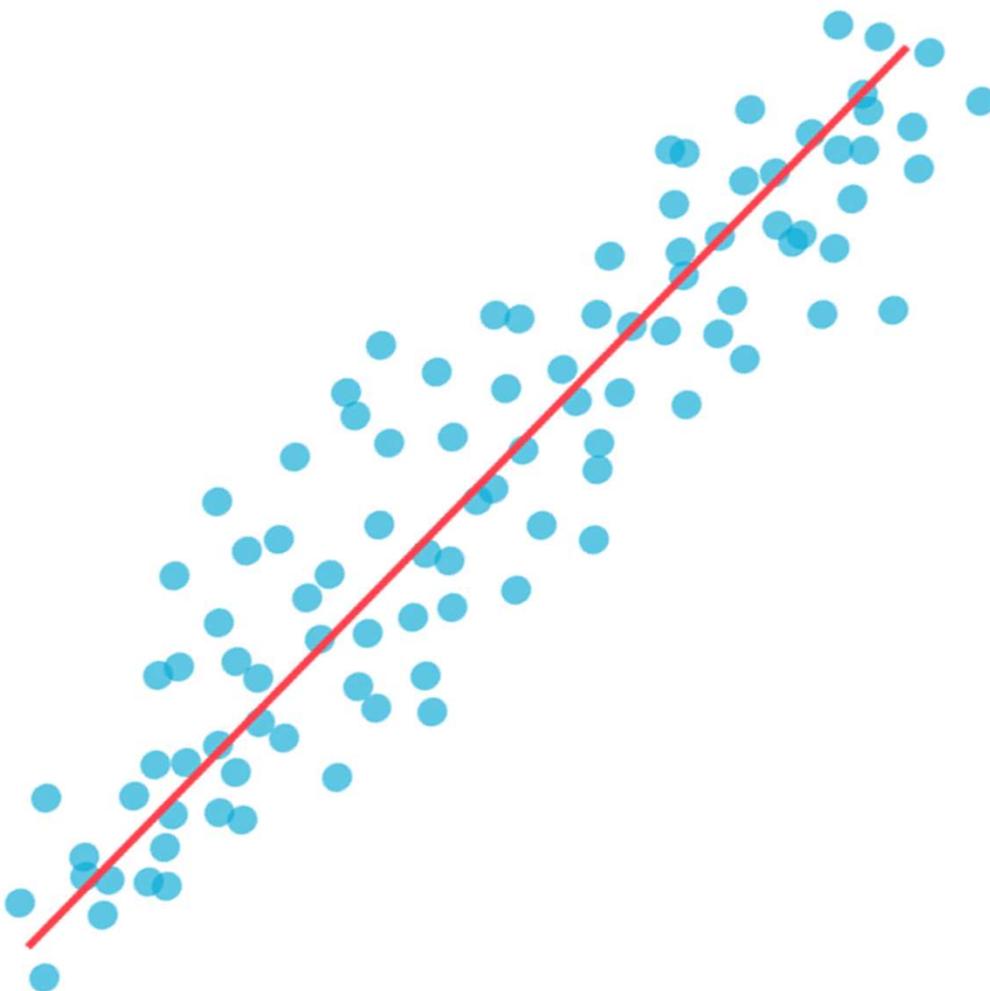
curse of dimensionality



Hands-On

Part 3

1. Implement the K-Means Algorithm for a set of random 2D datapoints (use the ‘`sklearn make_blobs`’ function to get a random dataset with underlying clusters
 - Visualize your results (Bonus: can you animate the graph to show each iteration of the algorithm?)
 - How could you improve the initialization-step to reduce strange results?
2. Think about how you could use your implementation to categorize a new (previously unknown) datapoint.
 - Bonus: Implement your idea and visualize the result

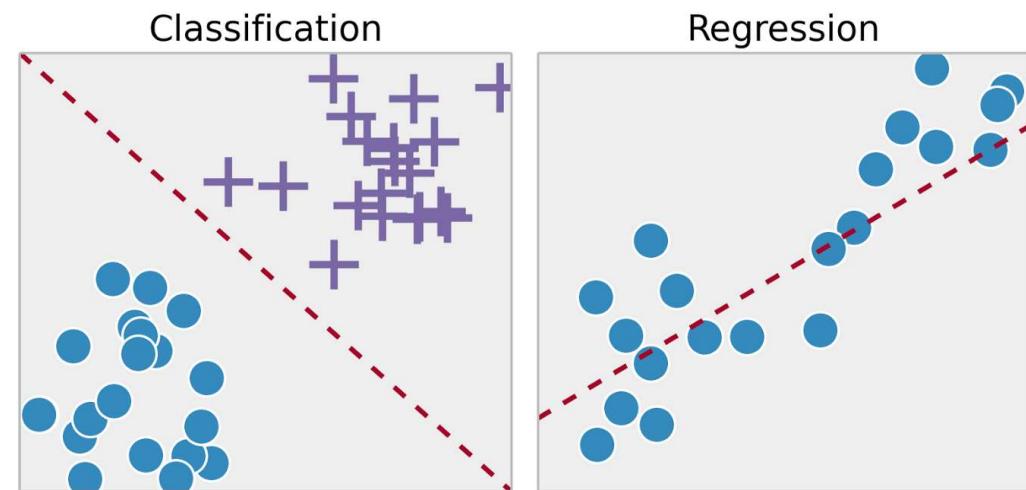


”

Supervised Machine Learning

Supervised ML

- **Idea:** given a dataset and its corresponding desired output, determine the best algorithm & parameters to predict the output from the data
- **Use cases:**
 - Classification
«which ad should I show this customer»
 - Regression / Prediction
- **Common Methods:**
 - Decision Trees
 - Support Vector Machines
 - Neural Networks

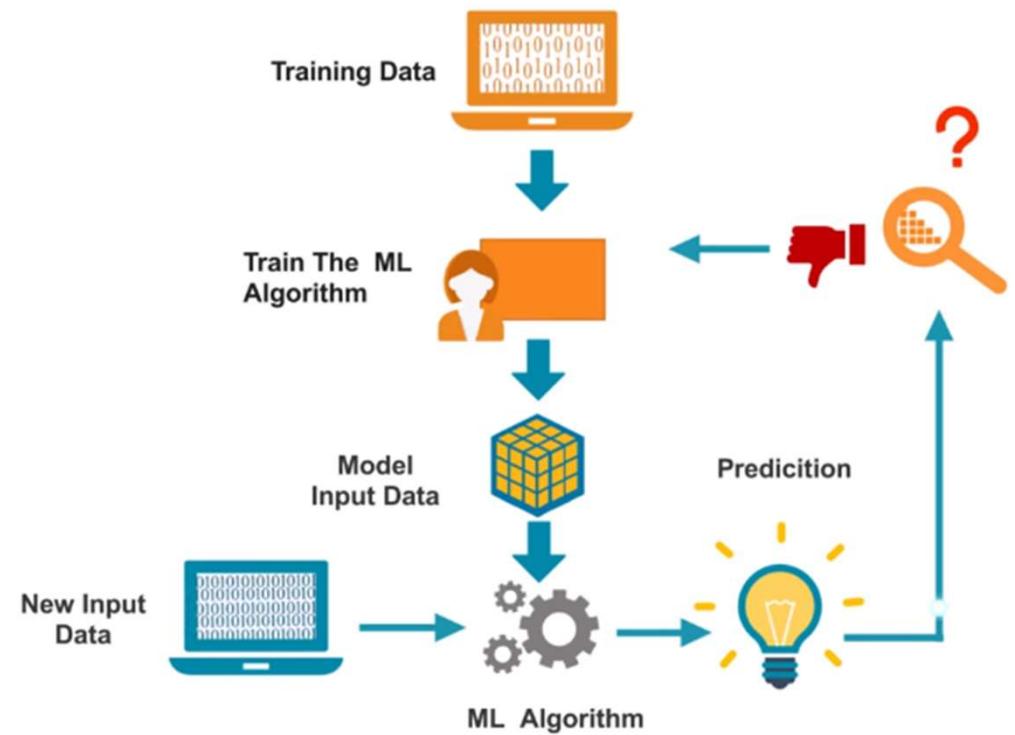


How does “supervision” work?

- **Intuition:** 1 Algorithm **builds models**, 1 Algorithm **scores the models** and chooses the best.
→ The builder Algorithm tweaks the best model and repeats
- **Building:** Depends on specific method, generally tweaking of model parameters.
- **Scoring:** Calculate a predefined **cost function / score**
 - E.g. Sum of Squared Errors (SSE)
- How Machines Learn [CGPGrey]

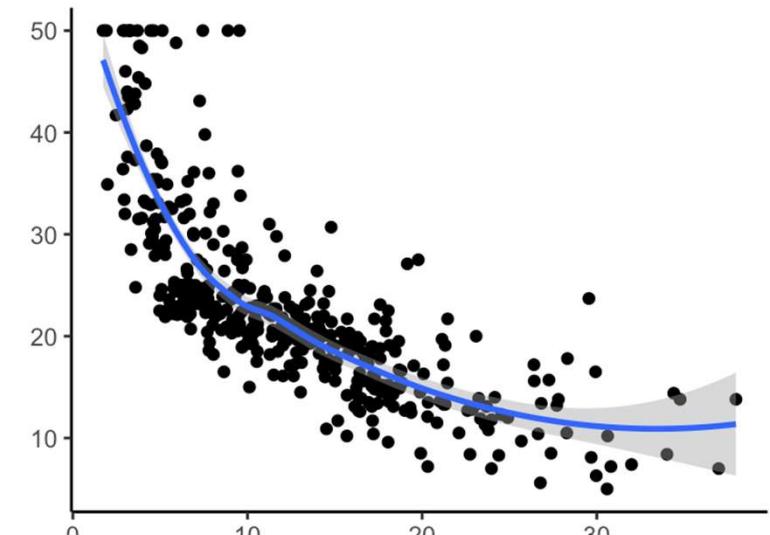
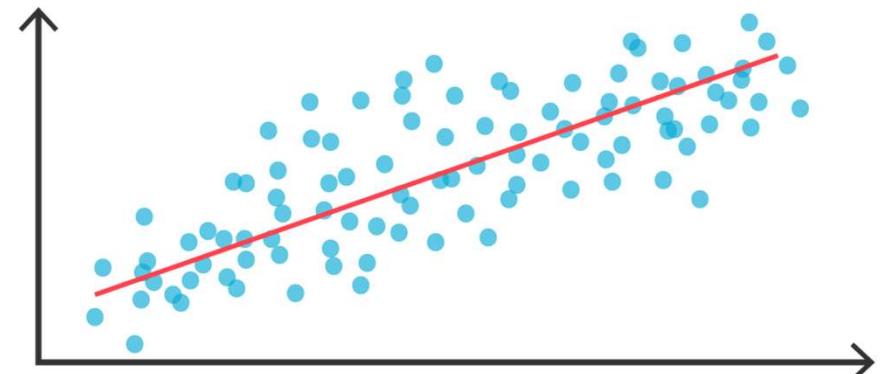
Typical Supervised ML Workflow

1. Define Goal
2. Get Data
3. Prepare Data
4. Create & Train A Model
5. Evaluate & Improve
6. Make Predictions



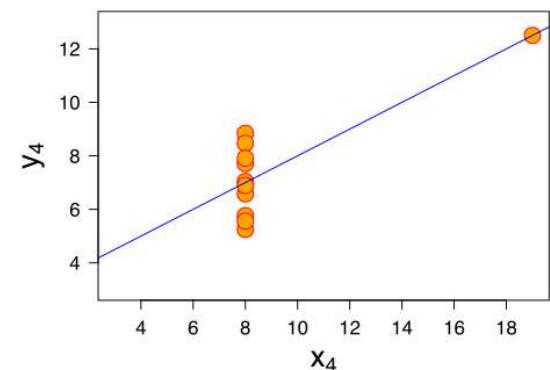
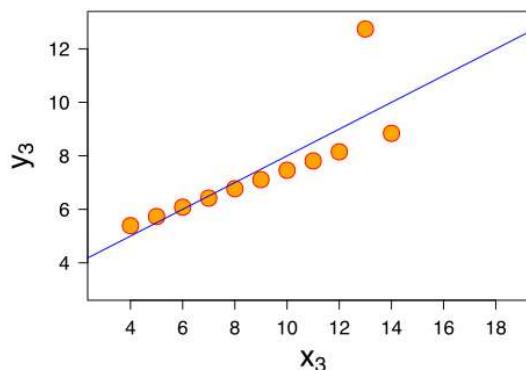
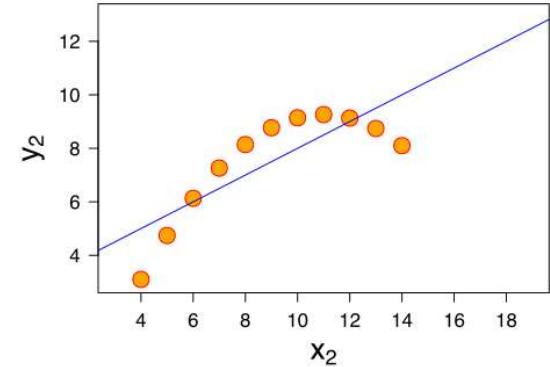
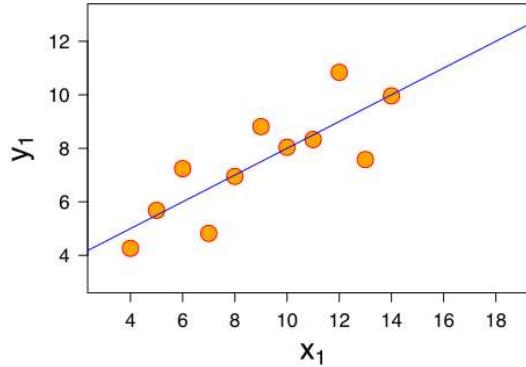
Regression

- **Idea:** Given a set of input variables, predict a numerical output variable
- **Use cases:**
 - Prediction in Marketing, Medicine, Finance etc.
- **Method:**
 - Find line that minimizes error between prediction and real values



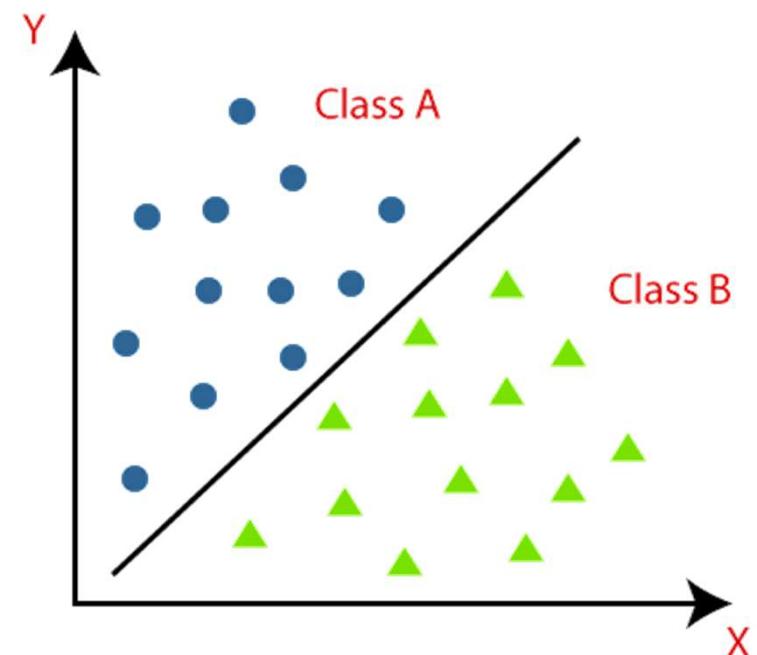
Regression Problems

- Can be susceptible to outliers (this can be overcome)
- Can be susceptible to over- / underfitting



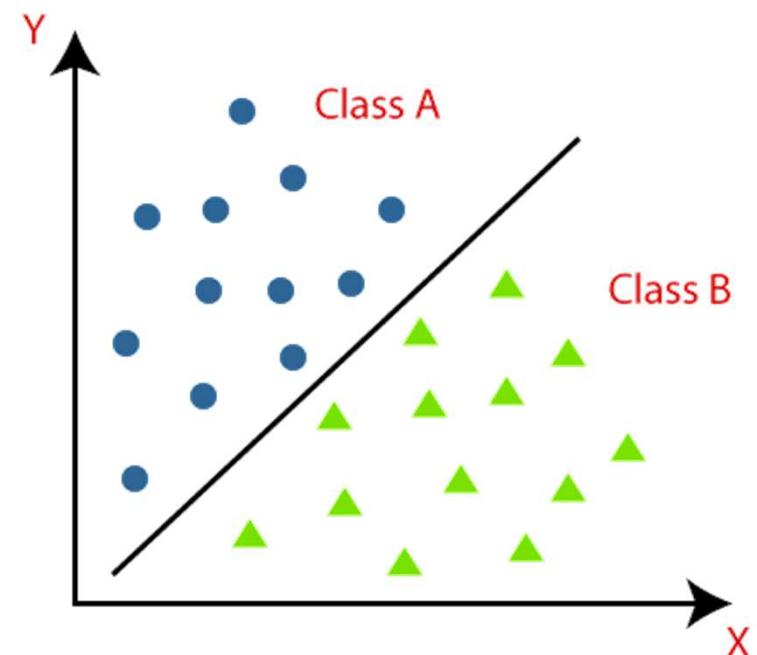
Classification

- **Idea:** given a set of input variables, predict a class for each datapoint
- **Use cases:**
 - Image-/ Speech-recognition
 - Medical prognosis
 - Customer Segmentation
 - Spam detection



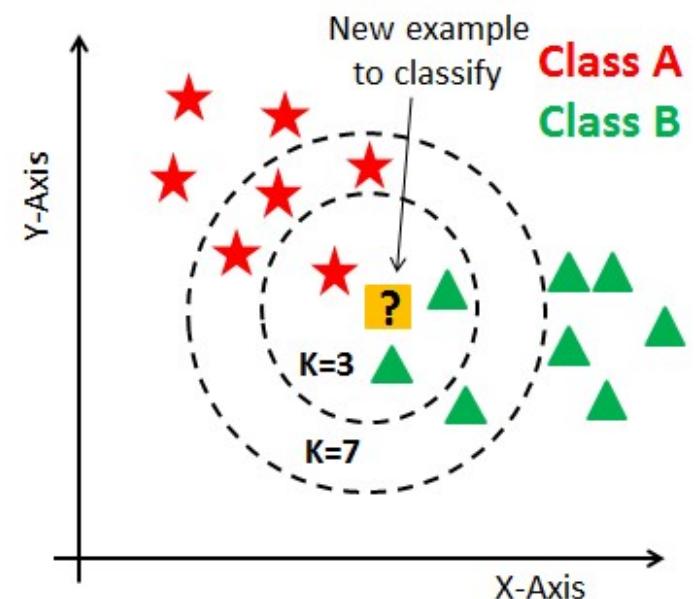
Classification Method

- In general, the aim is to minimize the number and severity of wrong assignments
- **Selection of classification Algorithms:**
 - K-nearest neighbours
 - Decision Trees
 - Support Vector Machines
 - Neural Networks



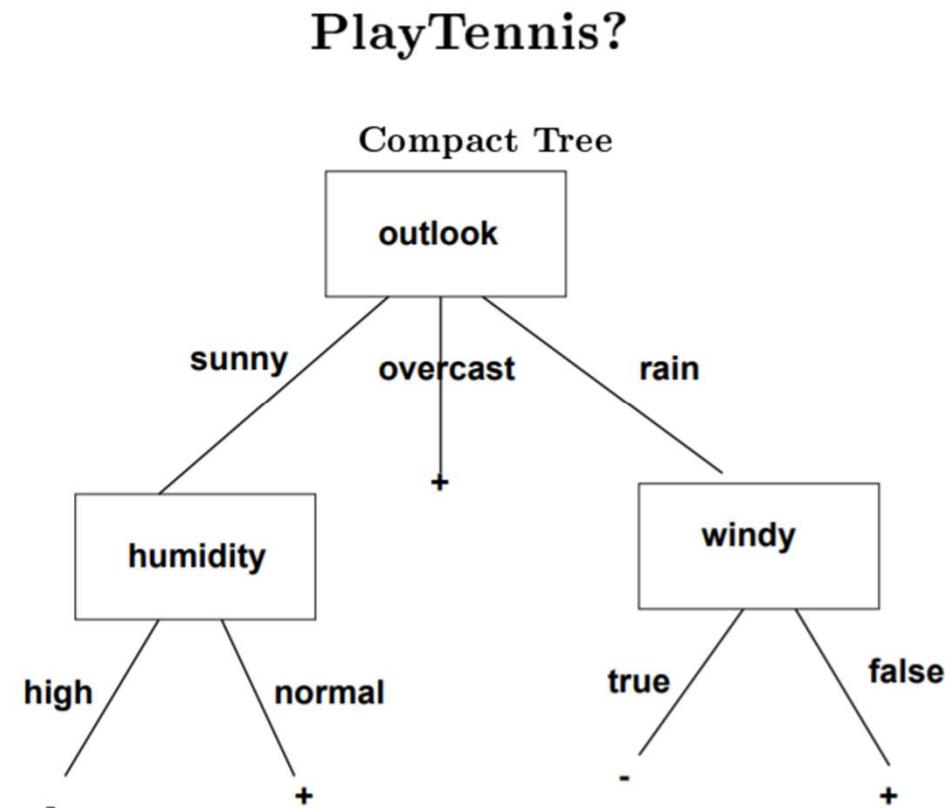
K-nearest-neighbour

- **Idea:** classify a point as the majority class of its k nearest neighbours.
- **Pros:** no training, simple, easy to incorporate more data
- **Con:** cannot handle very large, very high dimensional or imbalanced datasets. Sensitive to outliers



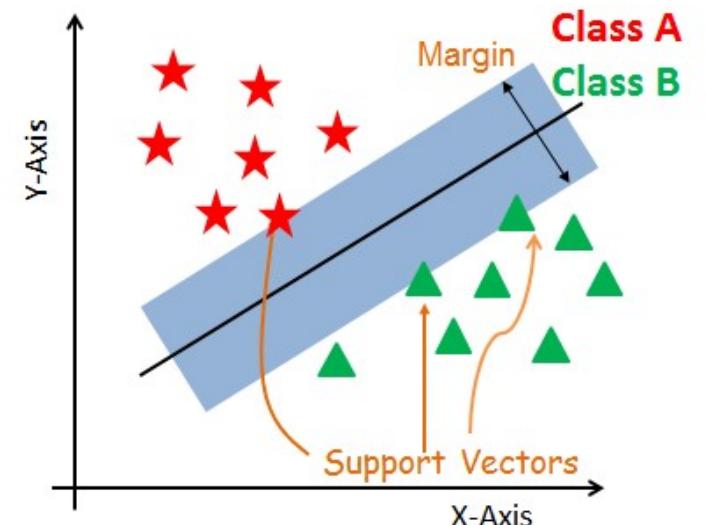
Decision Trees

- **Idea:** find a combination of binary decisions that best classify all datapoints into output classes
- **Pros:** less preprocessing required
Easy to interpret and may give additional insights
- **Con:** unstable results → ensemble methods (random forest)



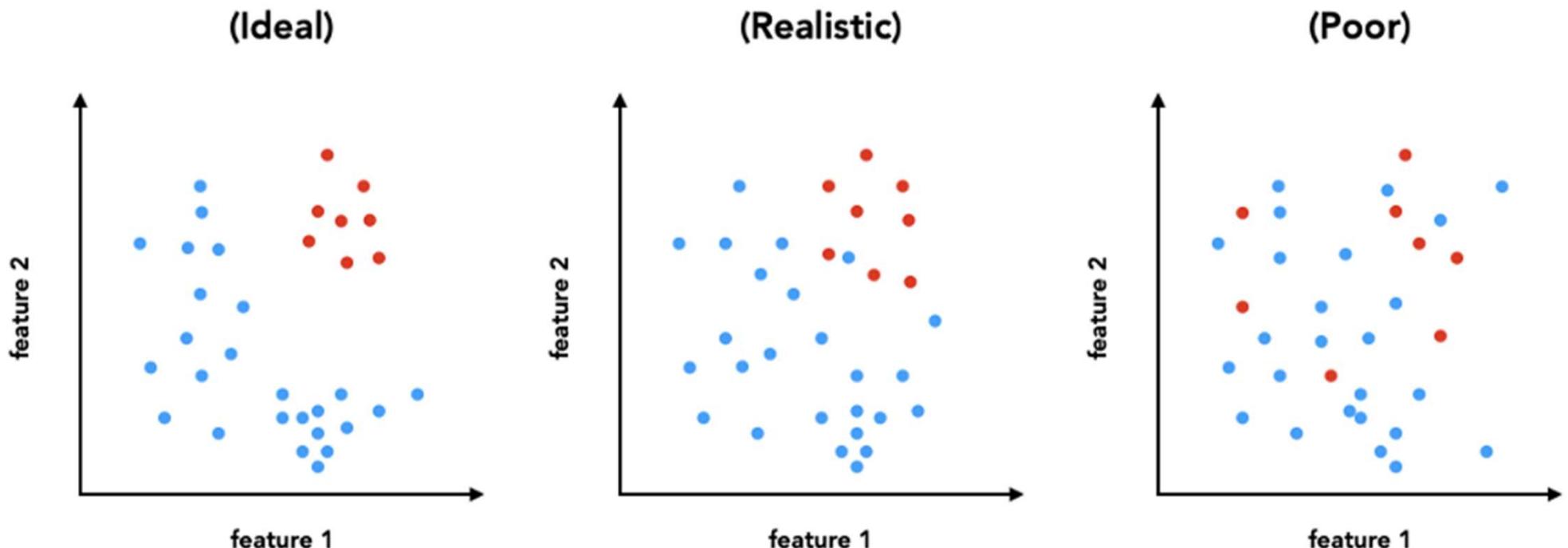
Support Vector Machines

- **Idea:** find the hyperplane (Line in 2D) that best separates the classes (largest margin)
- **Pros:** handles high dimensional data well
- **Cons:** problems with noisy / overlapping data.



Problems

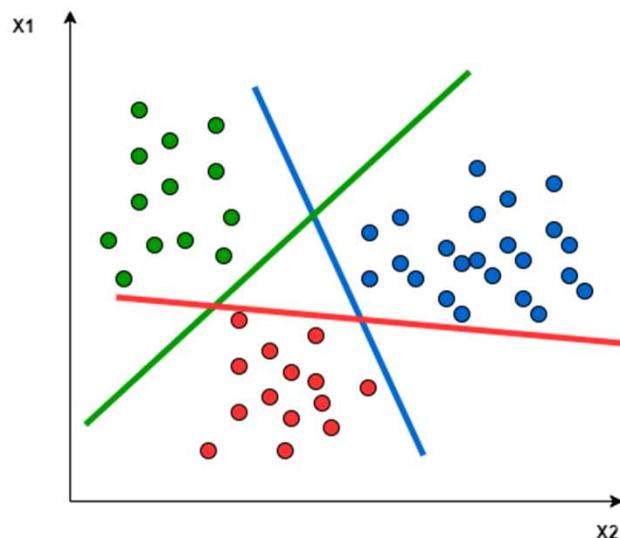
Poor class separation



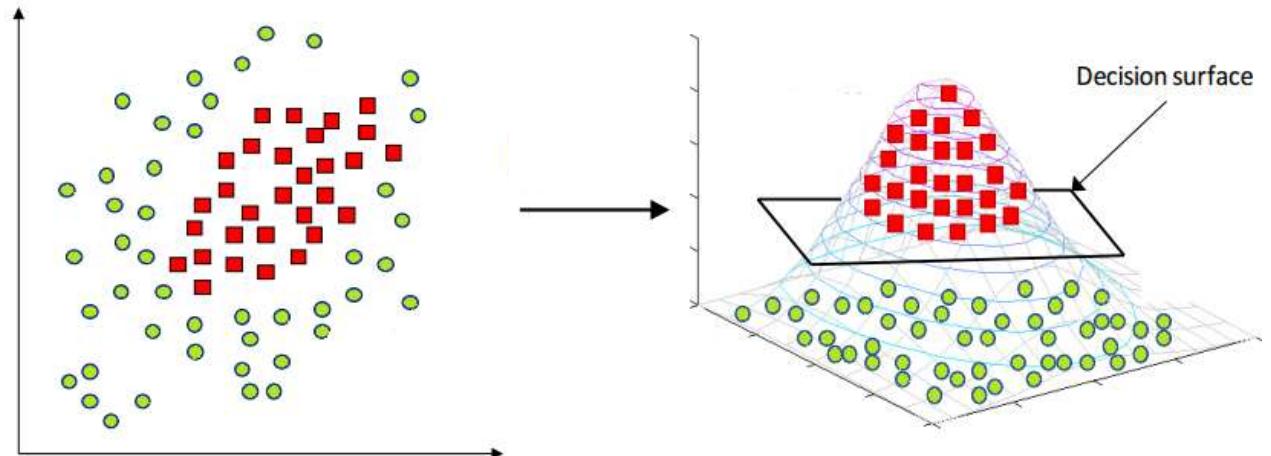
Problems

Advanced

Multi-Class Classification



Feature Transformation

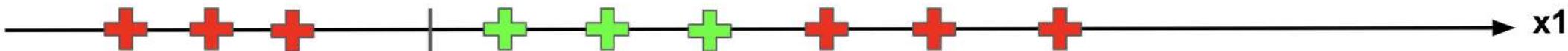


Evaluating Model Performance

Training & Test Data

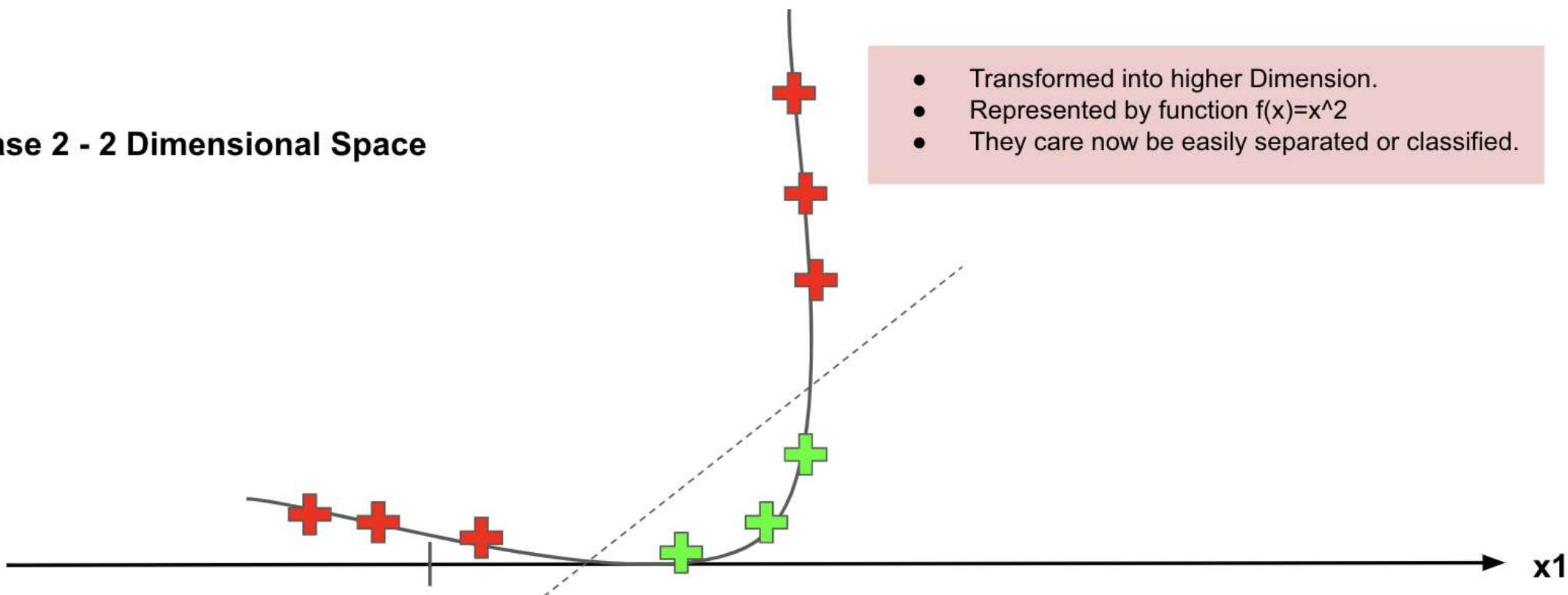
Case 1 - 1 Dimensional Space

- Points in 1 Dimension Plan.
- Represented by function $f(x)=x$
- They cannot be separated or classified.



Case 2 - 2 Dimensional Space

- Transformed into higher Dimension.
- Represented by function $f(x)=x^2$
- They can now be easily separated or classified.



Hands-On

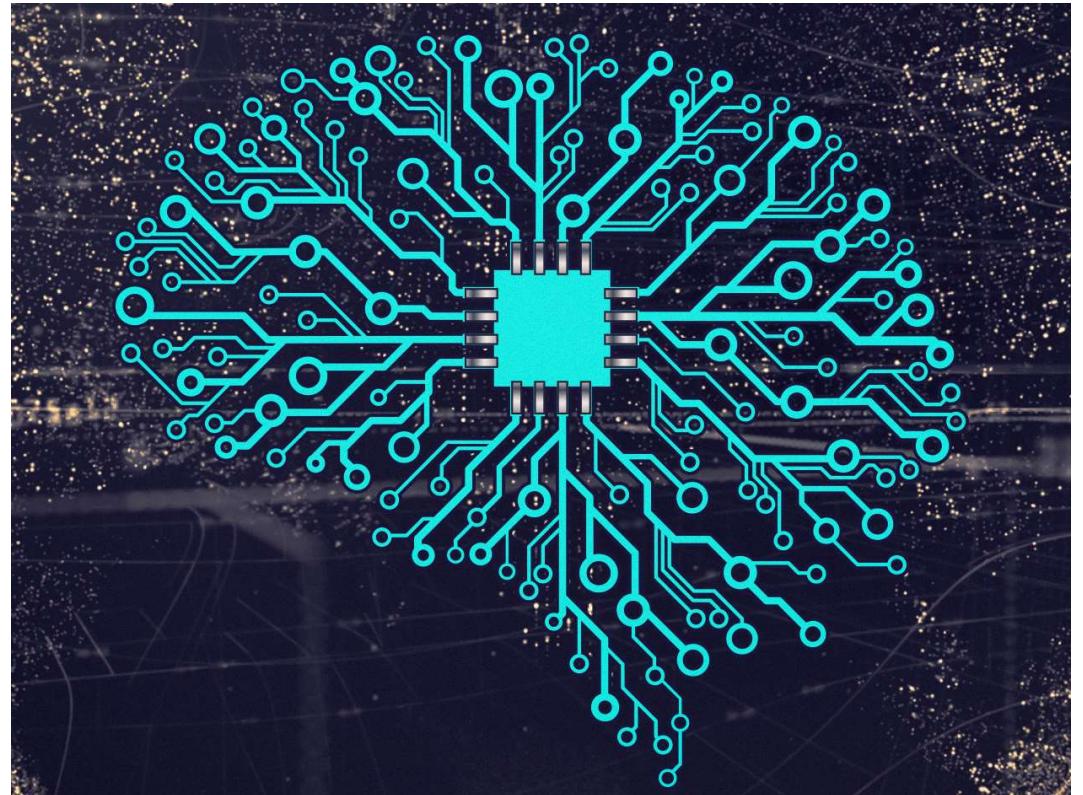
Part 4

1. Implement the KNN Algorithm for a set of 2D datapoints (use the ‘`sklearn make_blobs`’ function to get a random dataset with underlying clusters)
2. Use `seaborn.lmplot` to perform a linear regression on the tips dataset
 - `tips = sns.load_dataset("tips")`
3. Use `sklearn.svm.SVC` to **train** a SVM classifier and plot the data with your trained classifier.
 - `cancer = datasets.load_breast_cancer()`
 - Use your trained classifier to predict the classes of new datapoints.

Neural Networks

What are they not?

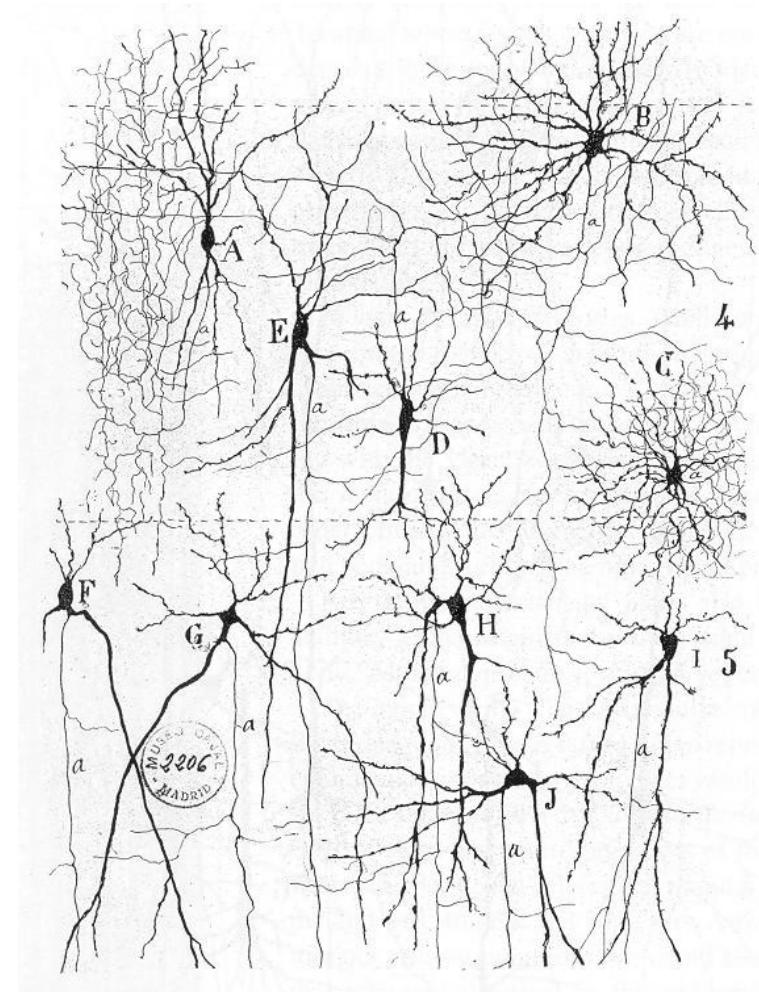
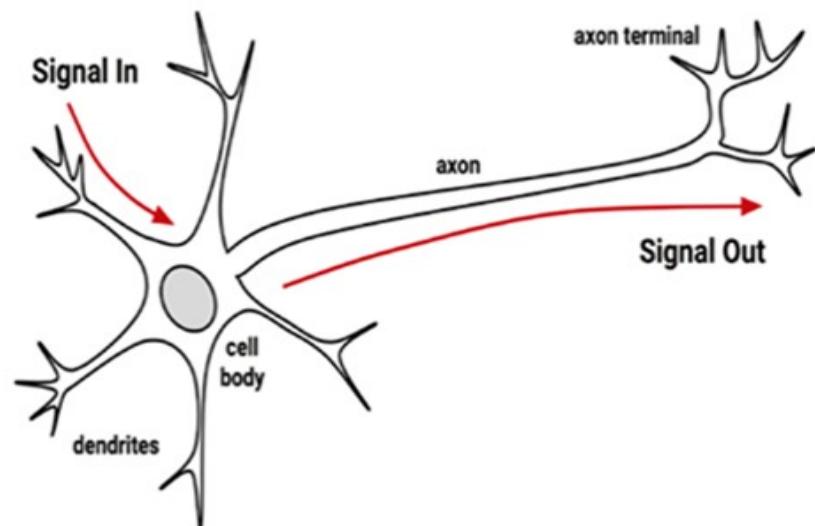
- Artificial Brains
- Super Computers
- Flexible ?
- “Intelligent” ?



Neural Networks

What are they?

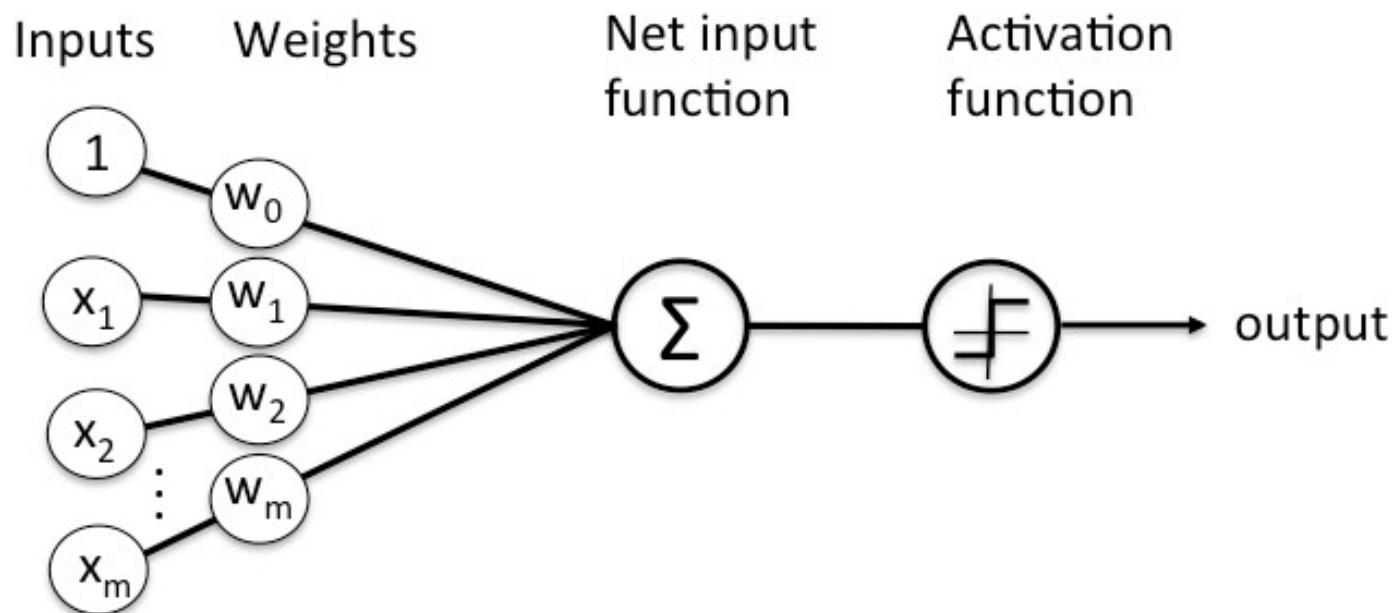
- Inspired by biological “neural networks” aka brains.



Neural Networks

What are they?

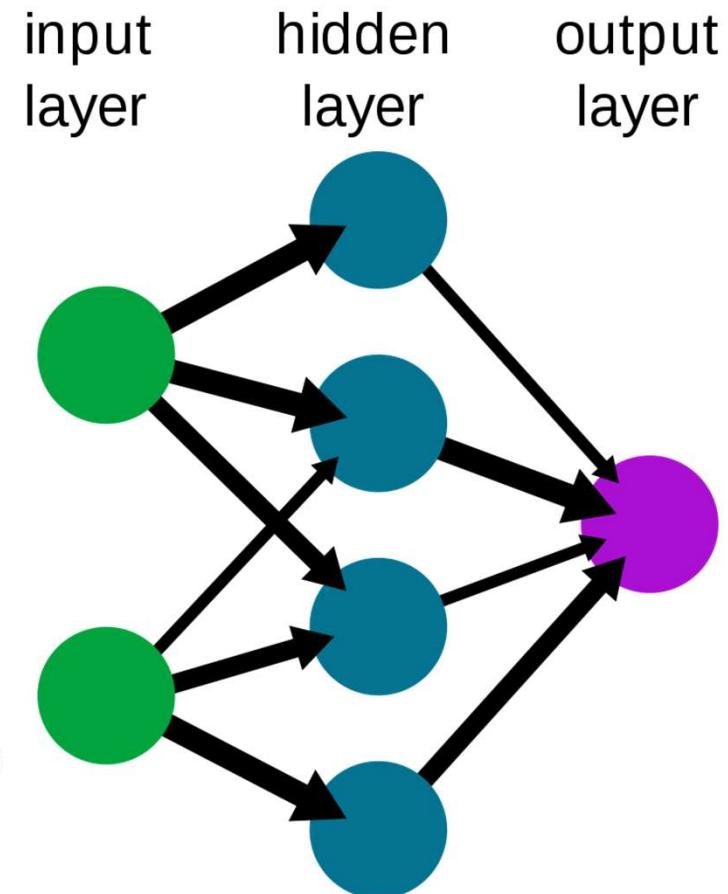
- A network of logical modules or **Nodes**



Neural Networks

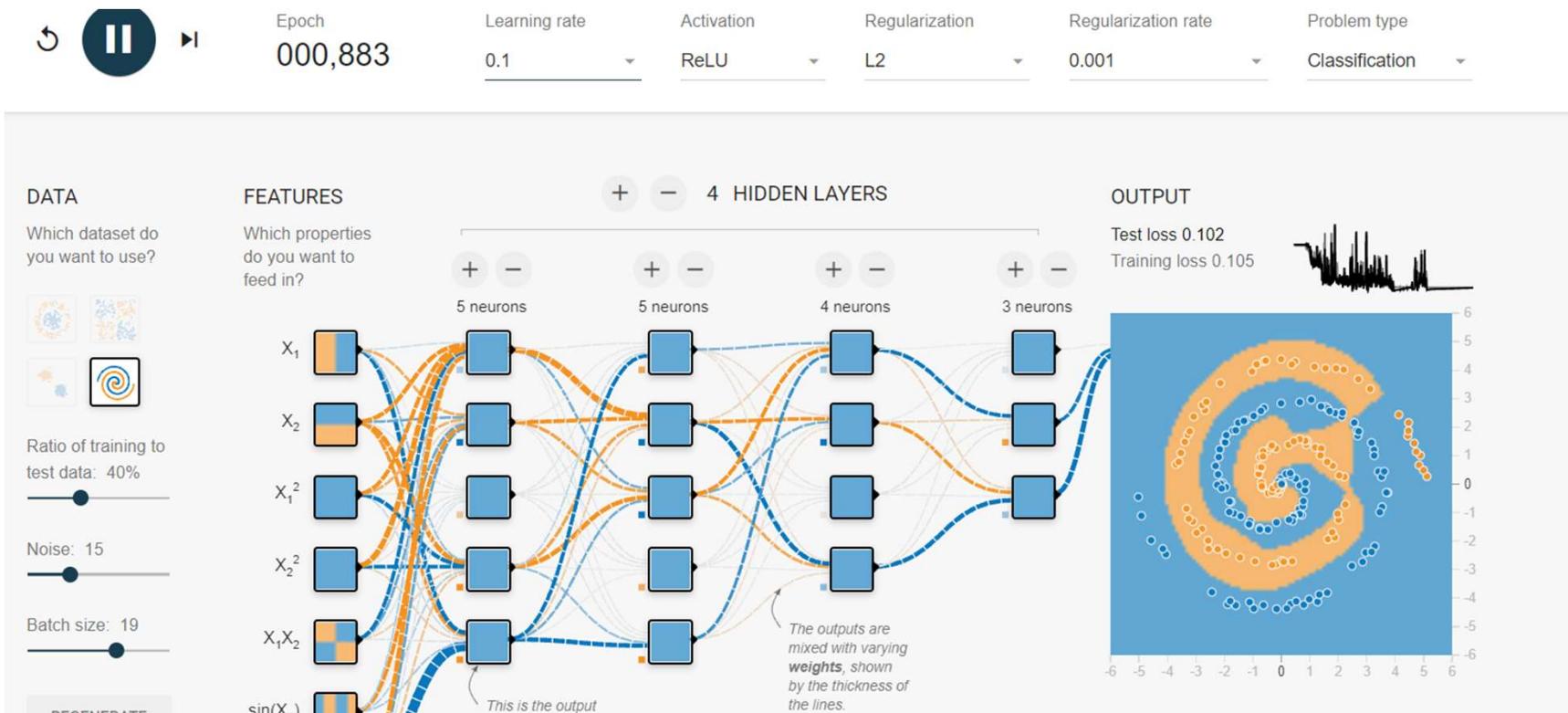
What are they?

- A network of logical modules or **Nodes**, generally arranged in multiple (hidden) layers with an input and output layer
- The **type, arrangement and connections** between the modules lead to complex **emergent patterns**.
- A Neural network is **trained by tuning the connections and parameters** of each node



Neural Networks

Try it out yourselves

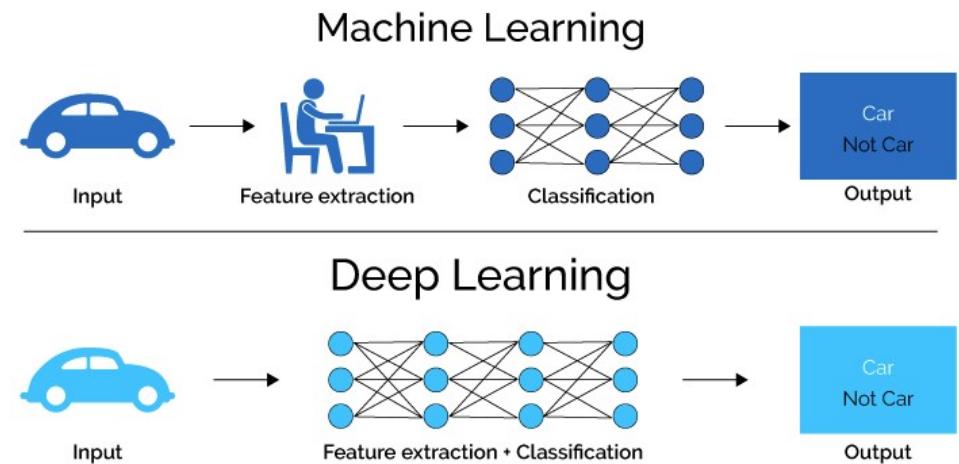


<https://playground.tensorflow.org/>



Deep Neural Nets

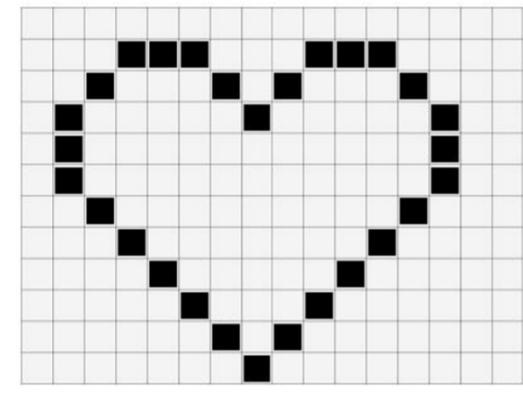
- A subset of machine learning algorithms, which extract their own training features from “raw” data
 - Requires less human intervention
 - Requires much more data



Convolutional Neural Net (CNN)

Example of Deep Neural Nets

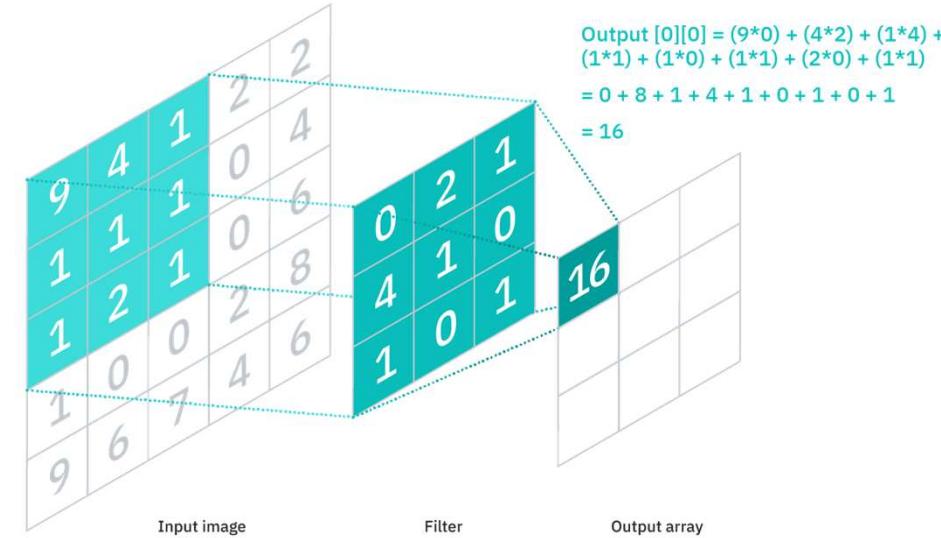
- **Idea:** instead of datapoints use whole matrices (such as images) as input.


$$\begin{bmatrix} 0000000000000000 \\ 0001110001110000 \\ 0010001010001000 \\ 0100000100000100 \\ 0100000000000100 \\ 0100000000000100 \\ 0100000000000100 \\ 0010000000001000 \\ 0001000000010000 \\ 0000100000100000 \\ 0000001010000000 \\ 0000000100000000 \end{bmatrix}$$

Convolutional Neural Net (CNN)

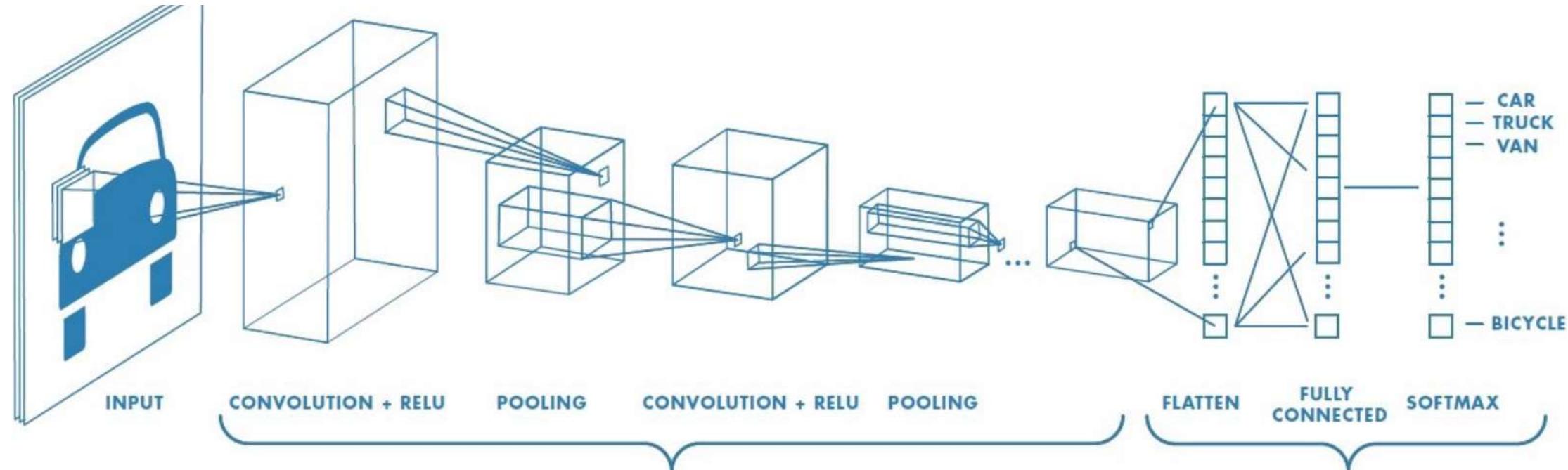
Example of Deep Neural Nets

- **Idea:** instead of datapoints use whole arrays/matrices (such as images) as input.
- Use **Convolution** to extract features from matrix
- **Convolution:** application of a filter to an input. The output is used in subsequent layers
- <https://poloclub.github.io/cnn-explainer/>



Convolutional Neural Net (CNN)

Example of Deep Neural Nets



Neural Networks

Pros & Cons

- **Pros:**

- Extremely flexible
- Handle complex big datasets (the more data the better)
- Once trained, predictions are fast

- **Cons:**

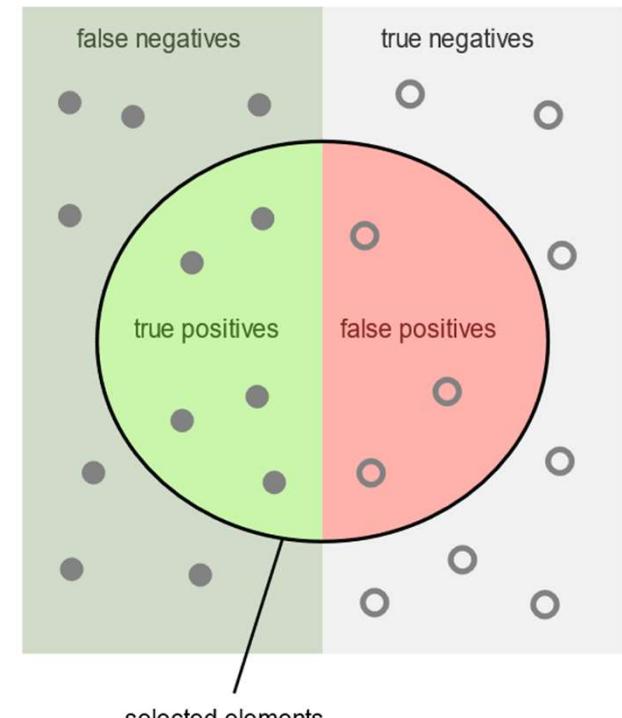
- Model is a Black Box
- Expensive training, inefficient with CPUs



Metrics

What are they & why should we care?

- **Problem:** We want to quantify the quality & performance of our model
- **Idea:** Predict with the model and compare the model predictions to the true values.
- These metrics can be used to make automatic decisions (e.g. in training, model-selection)



Confusion matrix / error matrix

For Binary Classifier

		Predicted condition	
		Positive (PP)	Negative (PN)
Total population $= P + N$			
Actual condition	Positive (P)	True positive (TP), hit	False negative (FN), type II error, miss, underestimation
	Negative (N)	False positive (FP), type I error, false alarm, overestimation	True negative (TN), correct rejection

Derived Metrics

Example of binary classifier (positive vs. negative)

$$\bullet \text{Recall} = \frac{TP}{TP+FN}$$

- Probability of detecting positive cases

$$\bullet \text{Precision} = \frac{TP}{TP+F}$$

- Probability of positive prediction being correct

$$\bullet \text{Specificity} = \frac{TN}{TN+F}$$

- Probability of correct negative prediction

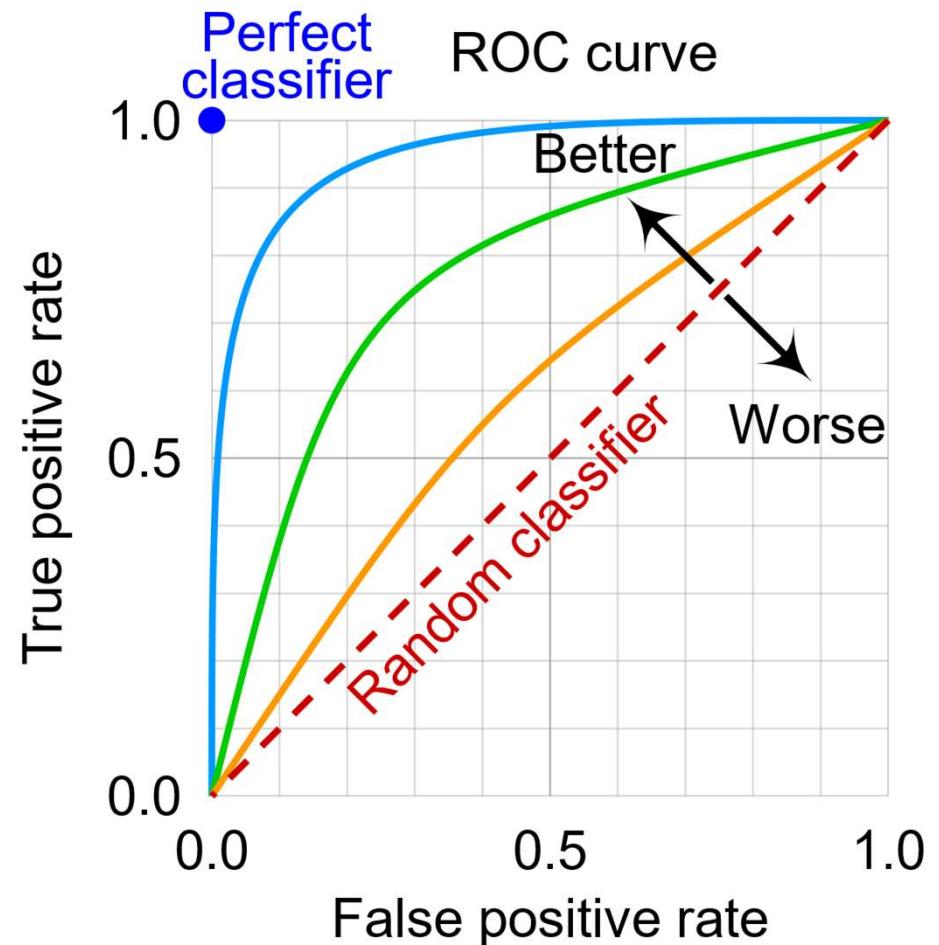
$$\bullet \text{Accuracy} = \frac{TP + TN}{TP+TN+FP+F}$$

- Probability of predicting correctly

		Predicted condition	
		Positive (PP)	Negative (PN)
		Total population = P + N	
Actual condition	Positive (P)	True positive (TP), hit	False negative (FN), type II error, miss, underestimation
	Negative (N)	False positive (FP), type I error, false alarm, overestimation	True negative (TN), correct rejection

ROC - Curve

- Diagnostic ability of a binary classifier system as its discrimination threshold is varied.
- Better classifiers have a larger area under the curve (AUC)



Common Problems in ML Applications

- Missing Values
- Noise
- Not enough data
- Overfitting

Challenges in ML applications

Noise

- Plot viel noise wenig noise (clustering)
- Plot varianz (nur Mean=> dann mit varianz)
- Measurement uncertainty
- Wrong Labels
- Errors

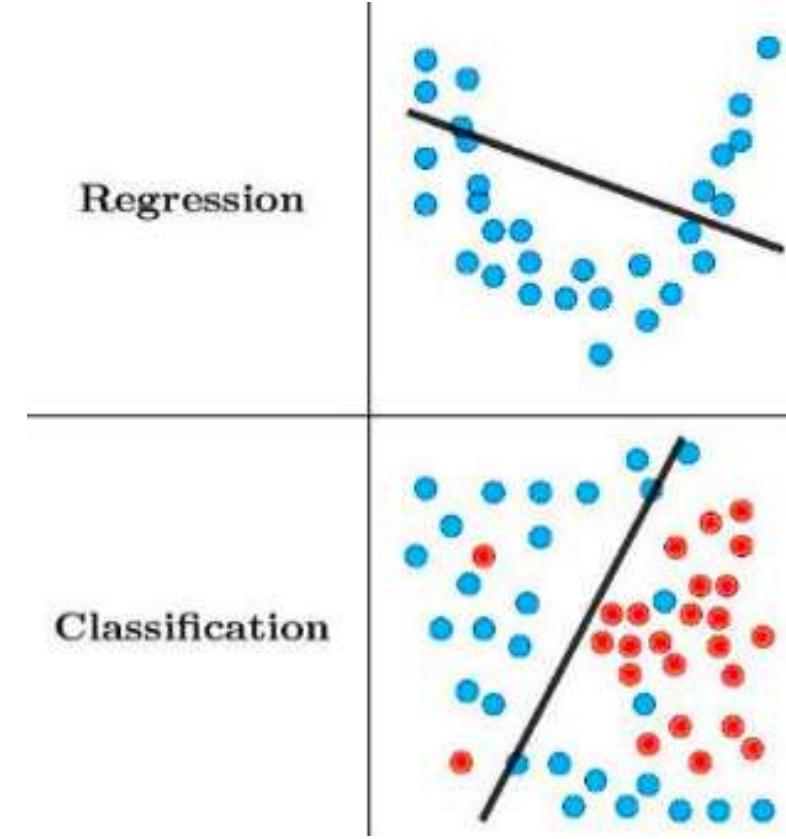
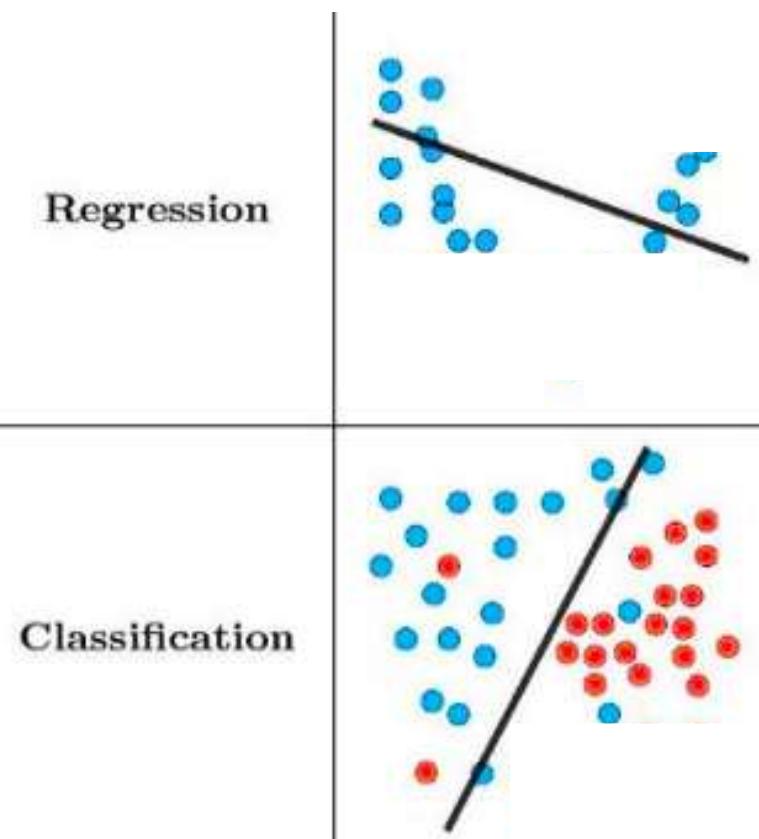
Challenges in ML applications

Limited Data (Underfitting)

- **Curse of dimensionality:** With increasing dimensionality of the data, much more datapoints are needed to reliably fit a model
 - Additionally; the complexer the model, the more data you need!
-
- You need **LOTS** of data to train state of the art ML models
- Obtaining, storing & handling this data is a key challenge of “Big Data” applications

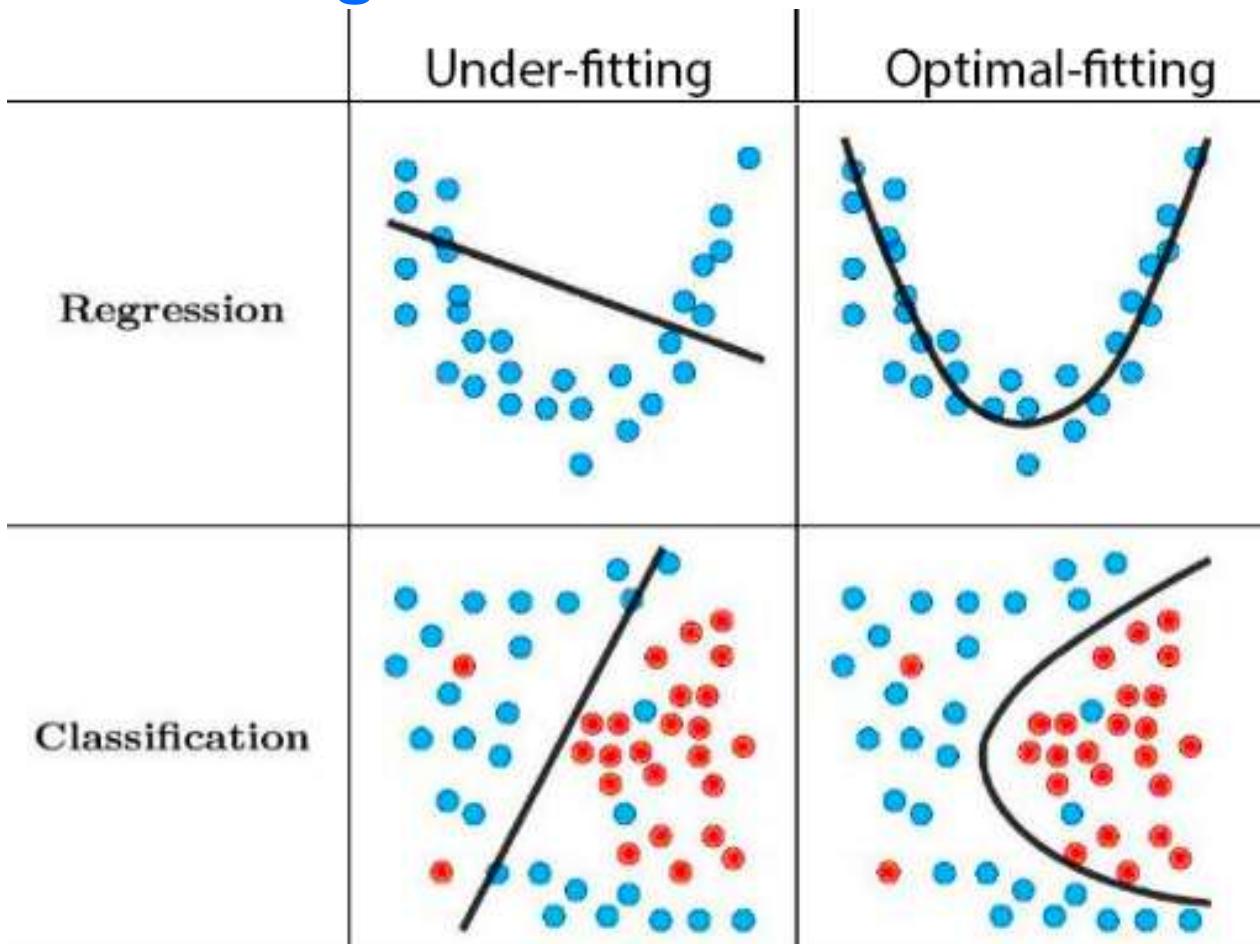
Challenges in ML applications → LinReg

Underfitting



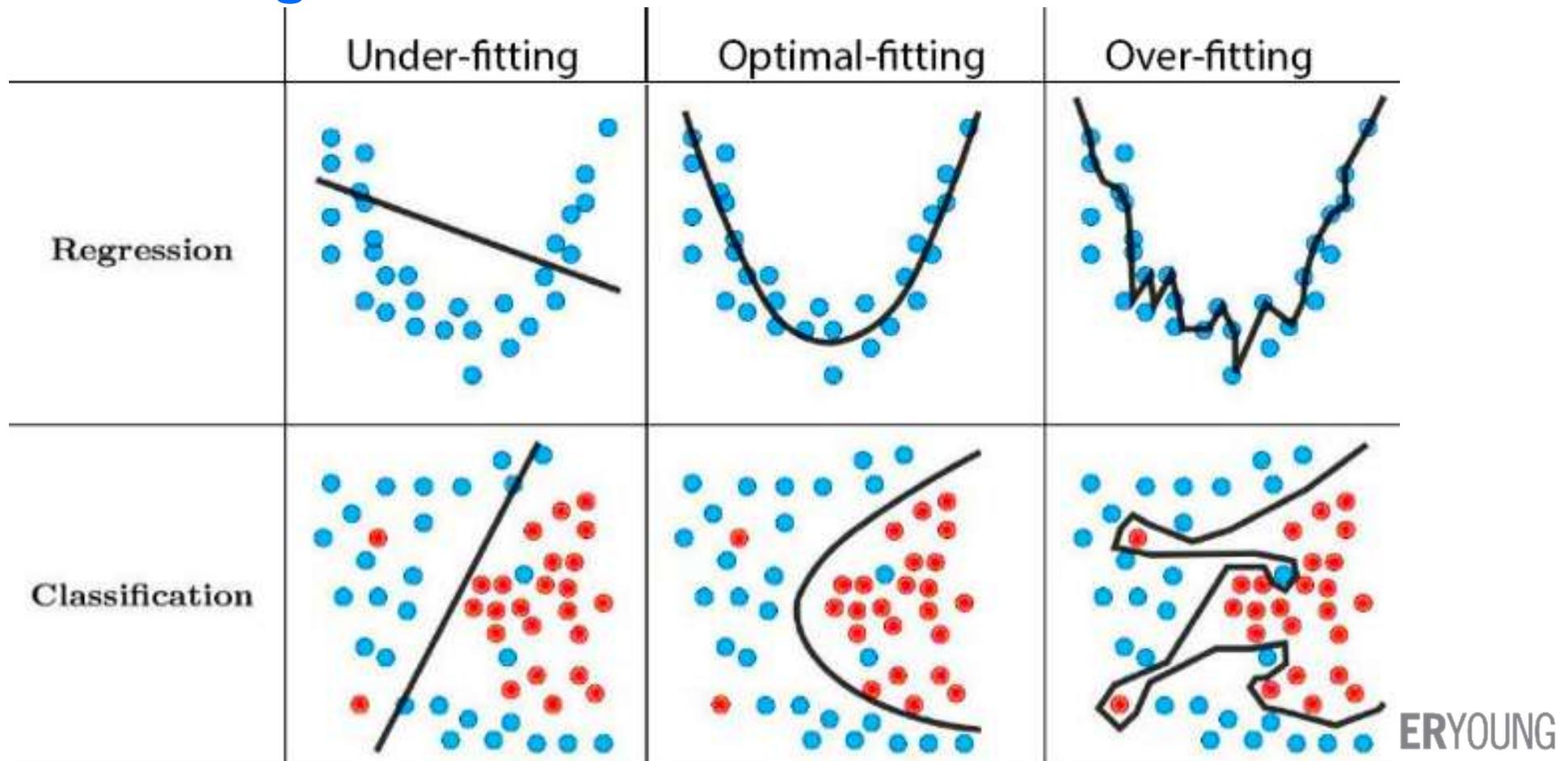
Challenges in ML applications

Overfitting



Challenges in ML applications

Overfitting



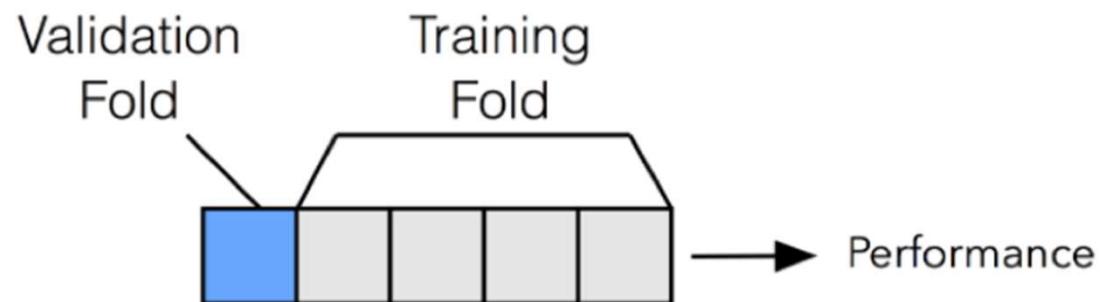
Challenges in ML applications

Overfitting

- Creating a model that fits the training-data **too well!**
 - I.e. it learns patterns that aren't there (Noise)
- **Fixes:**
 - **Removing noise** from the dataset (e.g. by dimensionality reduction)
 - Manually **limiting the allowable complexity** of the model
 - **Punishing overcomplicated models** by including a complexity metric in the cost function
 - Using test sets or crossvalidation to evaluate your model

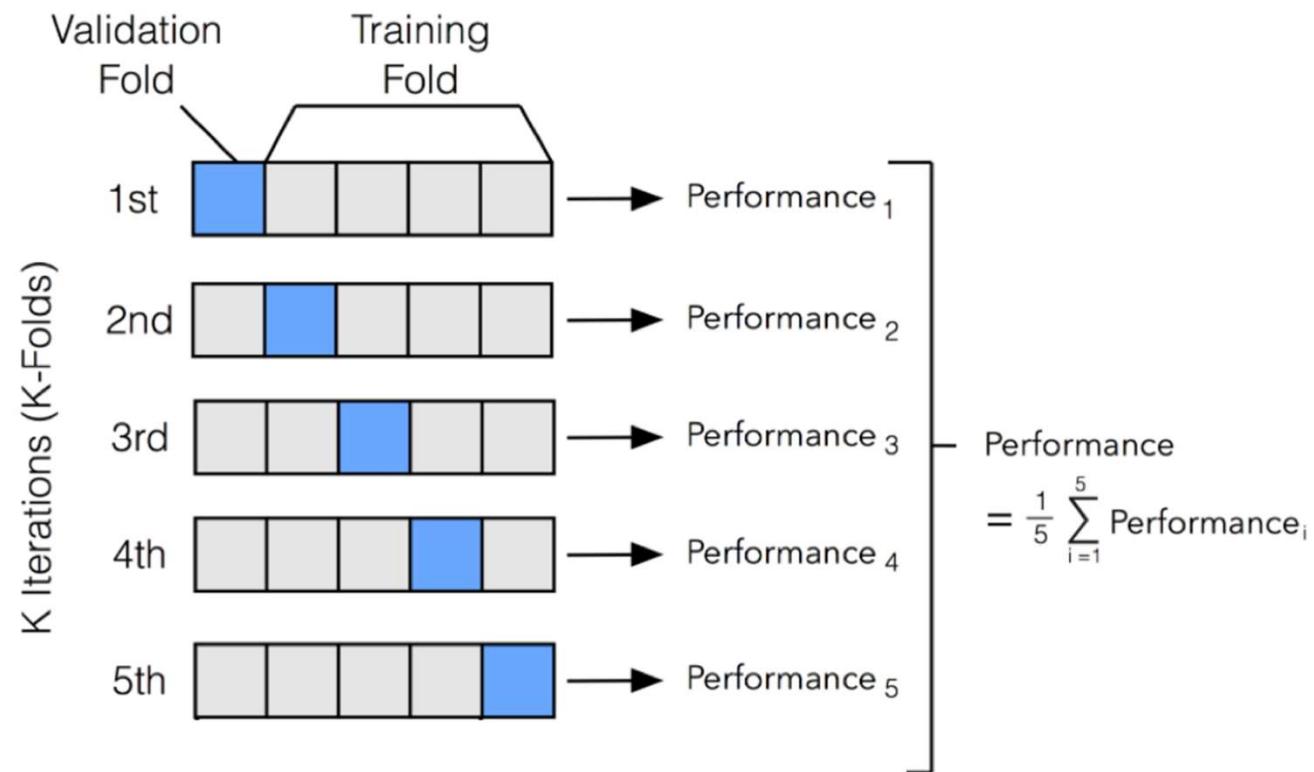
Cross Validation

- **Problem:** selecting a model based on training scores leads to overfitting
- **Idea:** Keep a subset of the data to validate the model. Choose the best model based on performance on the validation set
- **This data is never used in training!**



K-fold Cross Validation

- **Problem:** we loose valuable data to the test set
 - **Idea:** train and evaluate multiple times, with different splits, average results
- Every datapoint is used



Problems to solve with machine learning



Classification



Clustering



Regression

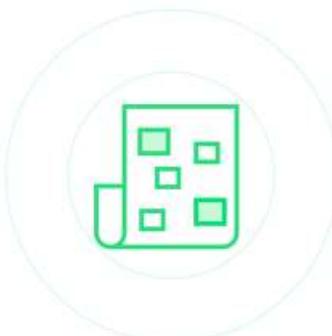


Dimensionality
Reduction

When is traditional computing better than machine learning?



Not enough data



Noisy Data



No time & money



Simple problem
to solve

Personendaten im Kontext von Big Data

- «Alle Angaben, die sich auf eine bestimmte oder bestimmbare Person beziehen» (Art. 3 Bst. a DSG)
- **Problem:** De-Anonymisierung bzw. Re-Identifikation
 - <https://amiunique.org/fp>
 - Set von Einkäufe mit gleicher Kreditkarte
 - Sets von Webrequests von gleicher IP Adresse

Personendaten im Kontext von Big Data

- «Besonders schützenswerte Daten»:
 - Weitreichende Einschränkungen für Gebrauch solcher Daten
- Quantität, Qualität und geplanter Gebrauch spielen eine Rolle

Beispiele

- Mit Fitness-Armbändern werden das Bewegungsverhalten (Anzahl Schritte), der Puls und das Schlafverhalten (Bewegungsaktivität während der Nacht) gemessen. Zudem kann der Armbandträger Informationen zu eingenommenen Mahlzeiten und zu seiner Befindlichkeit selbst erfassen. Während jede dieser Informationen bei isolierter Betrachtung kaum Hinweise auf den Gesundheitszustand einer Person zulässt, liefert die Kombination und Analyse dieser Daten über einen längeren Zeitraum Erkenntnisse über die Entwicklung des Pulses, die Zeitdauer, bis sich der Puls nach einer sportlichen Aktivität (hoher Puls, grosse Anzahl Schritte in kurzer Dauer) normalisiert hat, mögliche Schlafstörungen und deren Einfluss auf die Befindlichkeit oder den Hang zu Adipositas (stark fett- oder zuckerhaltige Ernährung und wenig Bewegung).
- In Loyalitätsprogrammen werden Transaktionsdaten aller Einkäufe erhoben, was deren spätere Auswertung im Rahmen von Warenkorbanalysen ermöglicht. Der Warenkorb eines isolierten Einkaufs (z.B. laktosefreie Milch, Wein, Pommes Chips, Zigaretten) hat dabei kaum Aussagekraft. Die Analyse der Einkäufe über einen längeren Zeitraum lässt aber möglicherweise Schlüsse auf den Konsum von Suchtmitteln, auf Laktoseintoleranz oder auf Hang zu Fettleibigkeit zu.
- Der Besuch einer einzelnen Website durch eine Person oder das Aufschalten eines einzelnen Eintrags in einem Blog lassen kaum Rückschlüsse auf deren Gesundheit zu. Anders ist dies möglicherweise dann zu beurteilen, wenn das Surfverhalten oder die «Posts» einer Person über einen längeren Zeitraum erfasst und analysiert werden: Die häufige Verwendung bestimmter Suchbegriffe und der Besuch der als Treffer angezeigten Websites (z.B. Informationsseiten über eine spezifische Krankheit, Websites von Selbsthilfegruppen wie Weight Watchers oder Anonyme Alkoholiker) können Hinweise darauf geben, dass die betreffende Person an einer bestimmten Krankheit oder Sucht leidet

Schutzmaßnahmen

1. Vertraulichkeit

- Authentifizierung und Authorisierung
- **Trennungsgebot:** gewährleisten, dass zu unterschiedlichen Zwecken erhobene Daten getrennt verarbeitet werden können

2. Integrität

- Weitergabe-, Verarbeitungs- und Dokumentationskontrolle
- Resilienz von Systemen/ Diensten

3. Regelmäßige Evaluierung der Maßnahmen

Project Work

- Choose a Dataset and a question you want to answer. Use one of the following techniques to answer your question:
 - Classification
 - Clustering
 - Prediction
- Choose a fitting ML method and metric to solve your problem
- Incorporate techniques such as data pre-processing, cross-validation and data visualization in your solution.
- **Present your solutions**

Reinforcement Learning

Further Reading / Watching

- Basics:
 - <https://www.youtube.com/watch?v=h0e2HAPTGF4>
 - https://www.youtube.com/playlist?list=PL8dPuuaLjXtO65LeD2p4_Sb5XQ51par_b
- More Math:
https://www.youtube.com/playlist?list=PLZHQQObOWTQDNU6R1_67000Dx_ZCJB-3pi
- Industry Trends and Opinions
<https://towardsdatascience.com/>
- More In-Depth Tutorial in GoogleColab
<https://developers.google.com/machine-learning/crash-course>