

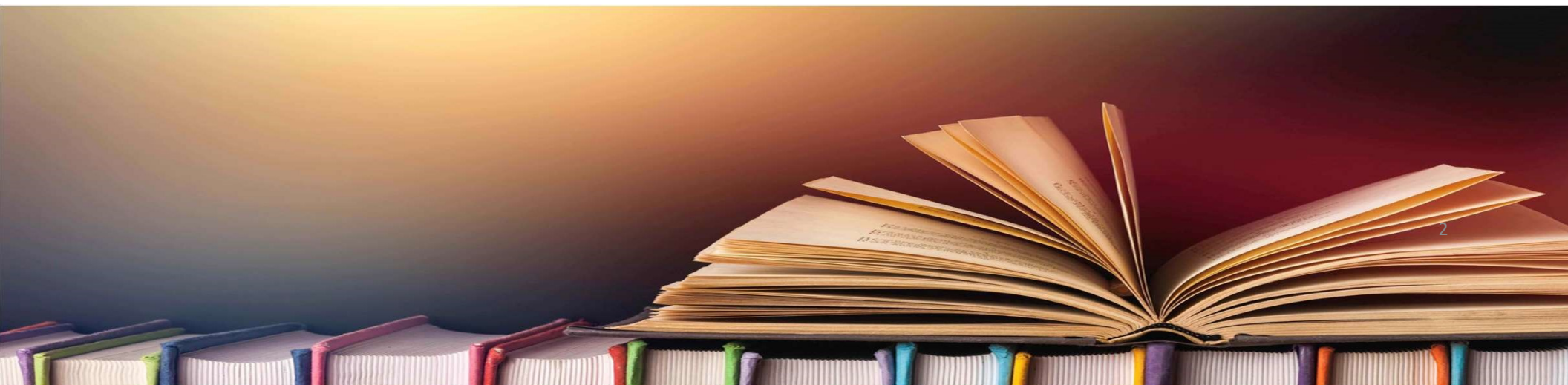
Module 1 : Machine Learning Review

Machine Learning Introduction



Bibliography

- Deep Learning book (Goodfellow, Bengio, Courville)
- Machine Learning @ Stanford (Prof Andrew Ng)
- Hands-On Machine Learning with Scikit-Learn & Tensorflow (Aurélien Géron)

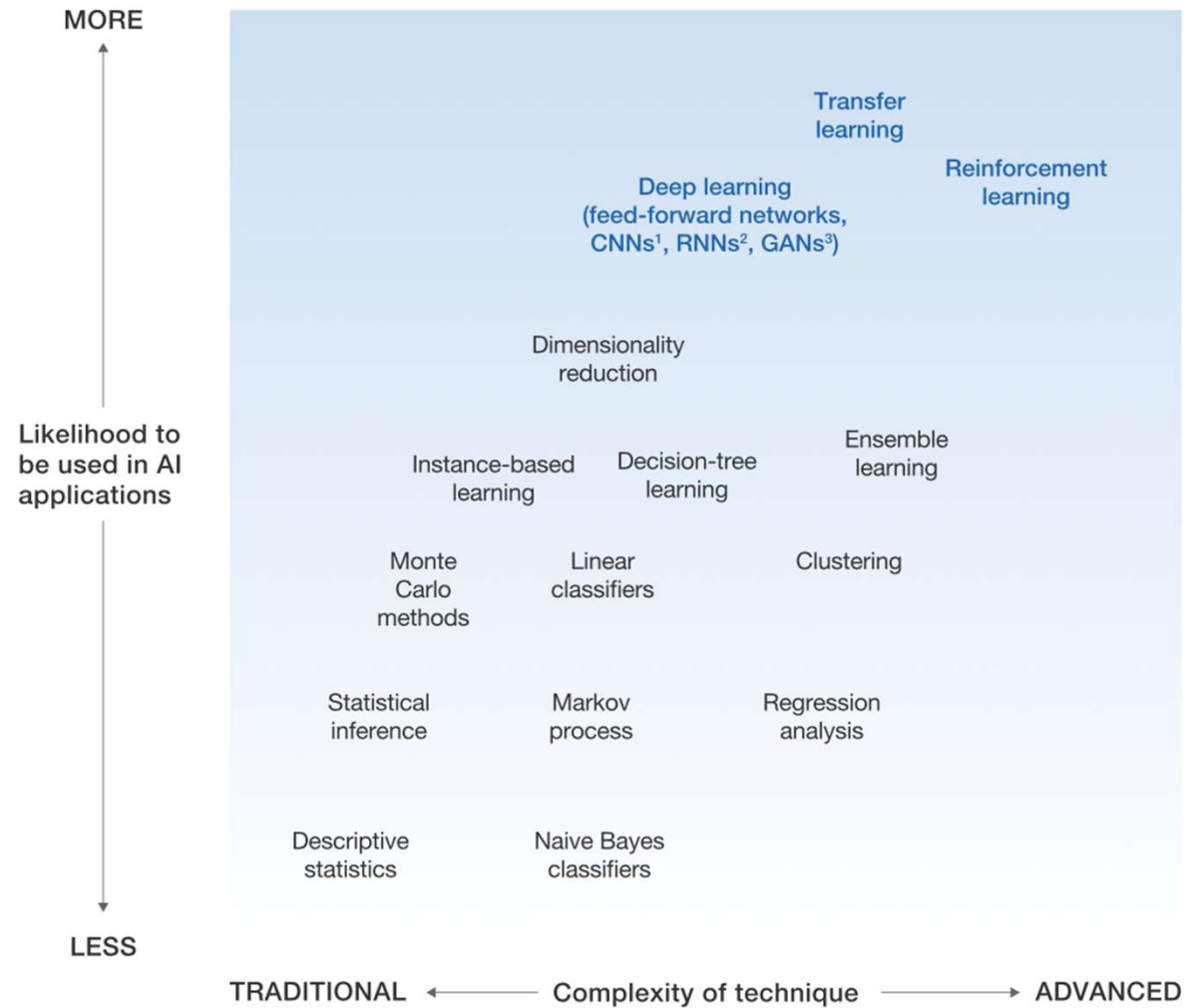




Learning Objectives

- What is Machine Learning ?
- Types of ML systems
- Main Challenges
- Data Preparation
- Optimization
- How to choose an algorithm ?
- Technical Details

Differences between ML and Statistical Modeling



THIS IS YOUR MACHINE LEARNING SYSTEM?

YUP! YOU POUR THE DATA INTO THIS BIG PILE OF LINEAR ALGEBRA, THEN COLLECT THE ANSWERS ON THE OTHER SIDE.

WHAT IF THE ANSWERS ARE WRONG?

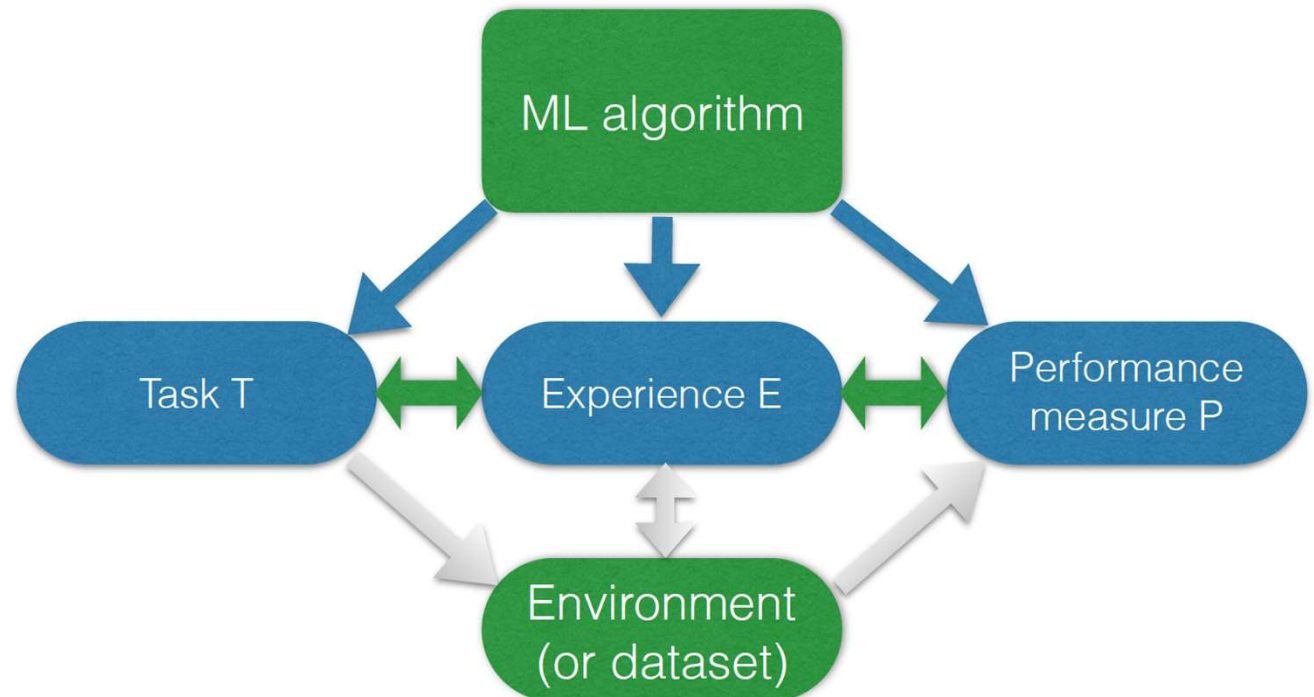
JUST STIR THE PILE UNTIL THEY START LOOKING RIGHT.



What is Machine Learning ?

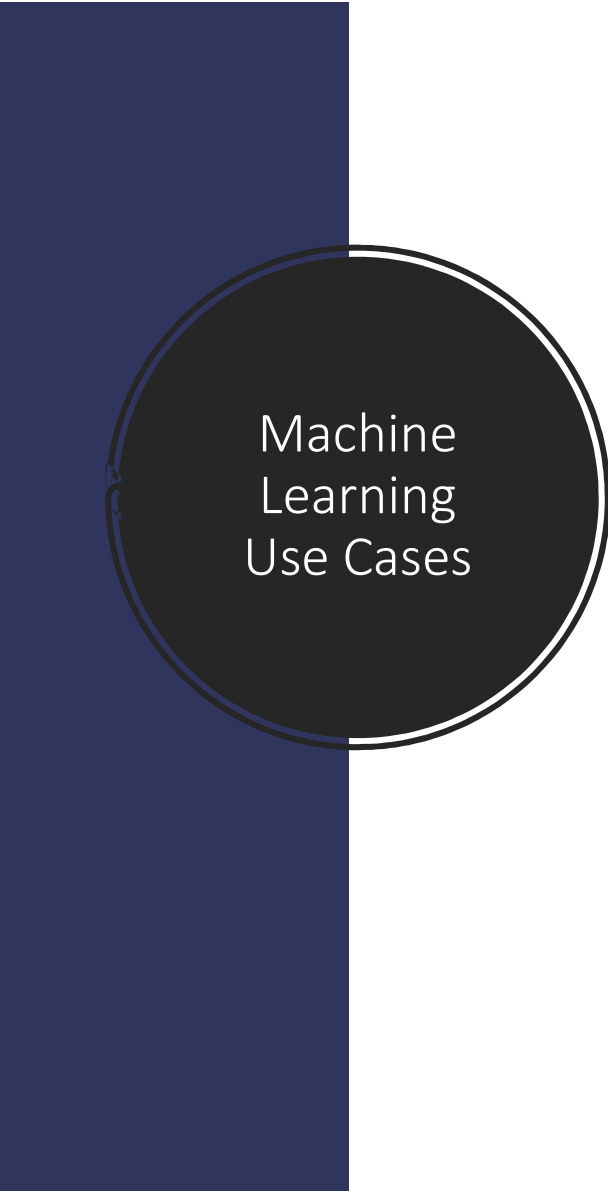
"Can machines do what we (as thinking entities) can do?" A. Turing

Definition



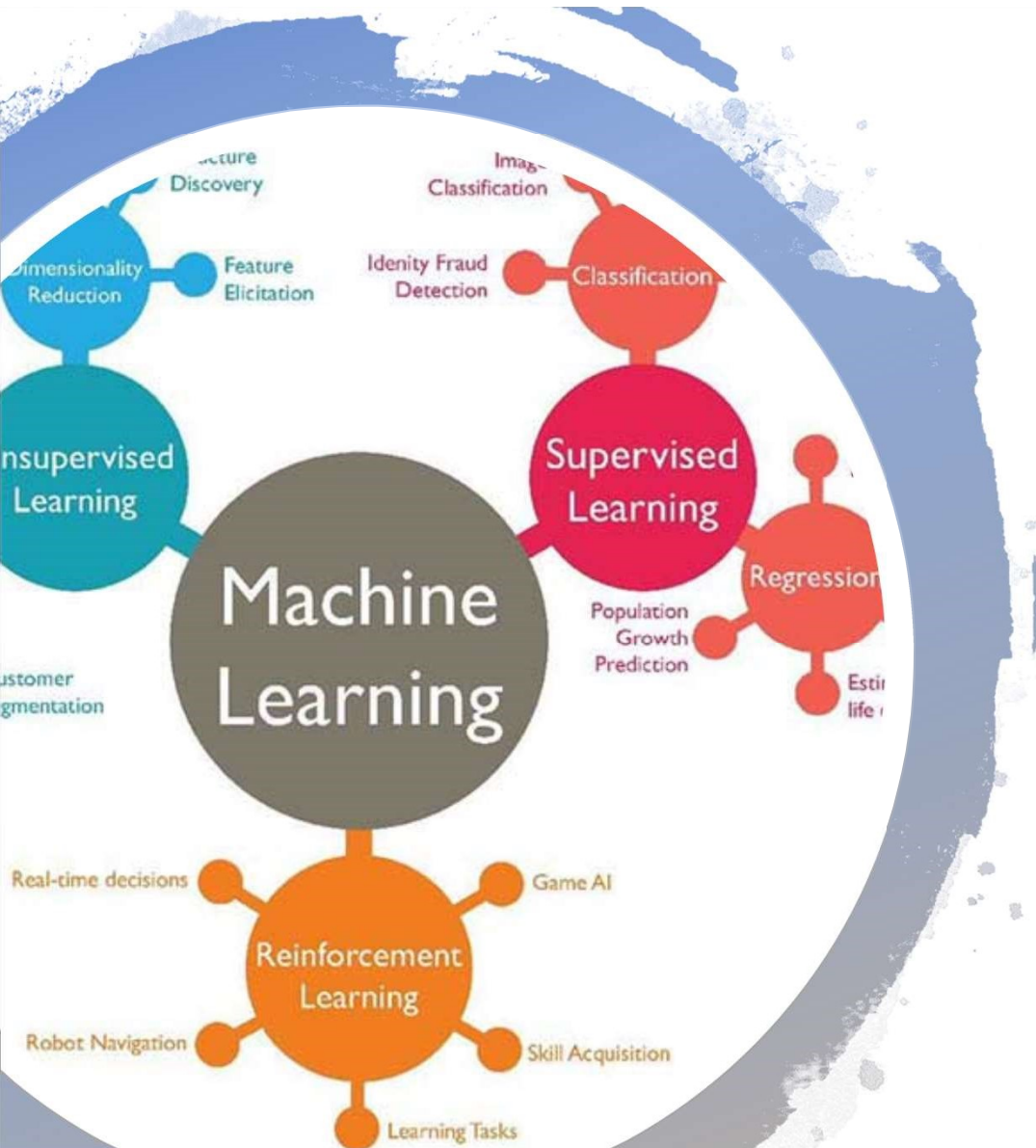
« A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P, if its performance at tasks in T, as measured by P, improves with experience E. »

Tom M. Mitchell (1997)



Machine Learning Use Cases

- Problems for which existing solutions require a **lot of hand-tuning** of **long lists of rules**
 - ML *simplifies* code and *performs better*
- **Complex** problems for which there is no good solution at all using a traditional approach
 - ML can *find a solution*
- **Fluctuating** environments
 - ML system can *adapt* to new data
- Getting insights about complex problems and **large amounts of data**



Types of Machine Learning Systems

Online learning

Supervised learning



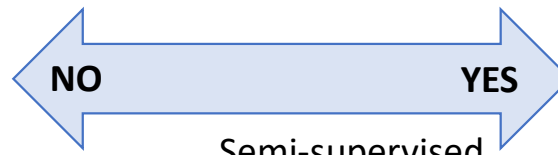
Machine Learning Systems

Batch learning

Unsupervised learning
instance-based learning

Model-based learning

Trained with **human supervision** (learning pillars) ?



Semi-supervised learning

Can learn **incrementally** on the fly ?



Compare new data points to **known data points** (or instead detect patterns) ?



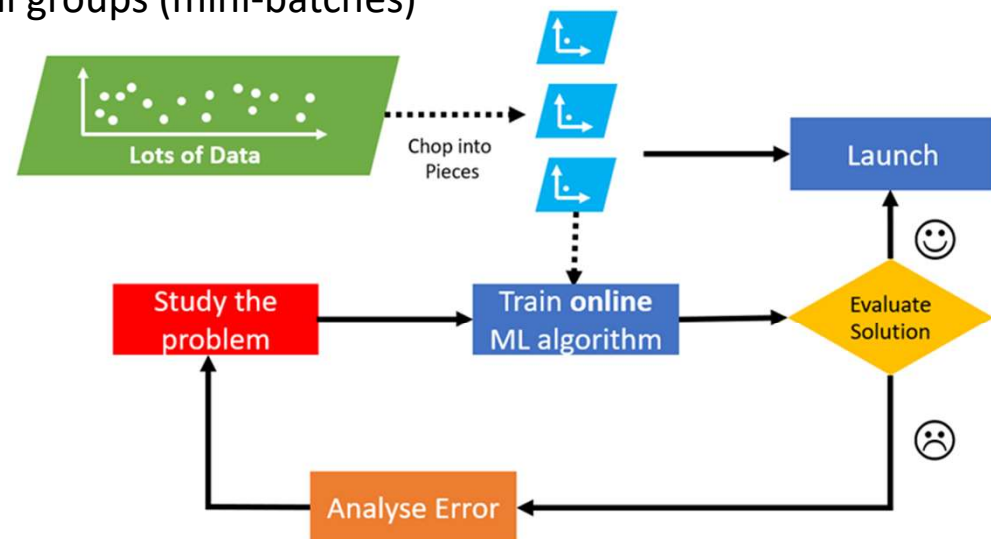


Batch Learning

- Not capable of learning incrementally
- Must be trained using all the available data
- With new data (new types), need to train a new version of the system from scratch on the full dataset
- Advantage :
- Drawbacks:

Online Learning

- Train the system **incrementally** by feeding data instances sequentially, system can learn about new data **on the fly**
 - Individually or in small groups (mini-batches)



- Important parameter : **learning rate**
 - *How fast the system should adapt to changing data*
 - **High learning rate** = adapt quickly, but forgets quickly
 - **Low learning rate** = system with more inertia
- **Challenge** : bad data will gradually kill the performance
 - Monitor your data (anomaly detection algorithms) and performance



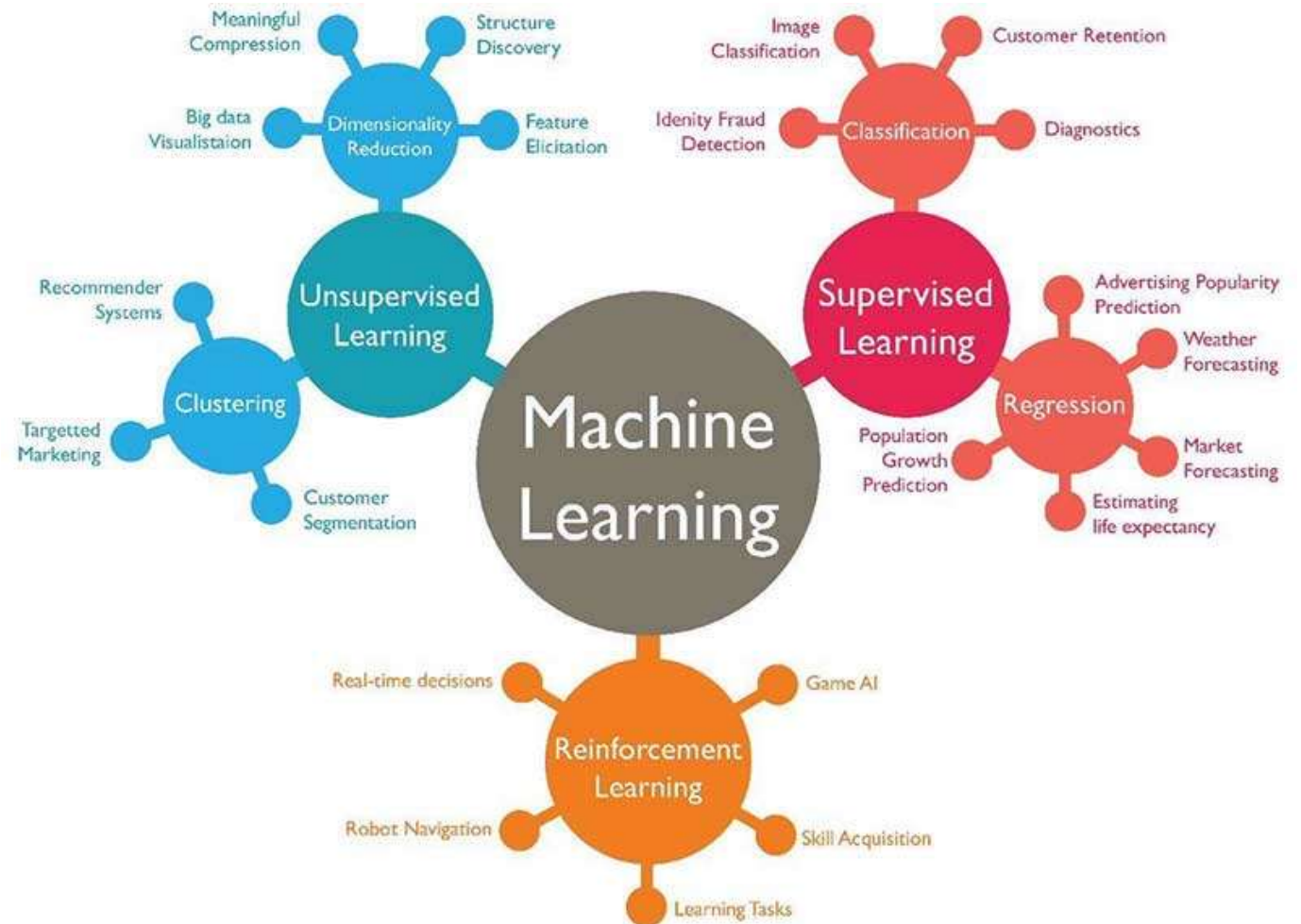
Generalization

- **Instance-based** Learning
 - System learns examples **by heart**
 - Generalizes to new cases using a **similarity measure**
 - Examples : k-nearest neighbors, decision trees
- **Model-based** learning
 - **Build a model** of the examples
 - Use the model to **make predictions**
 - Examples : Neural Networks,..



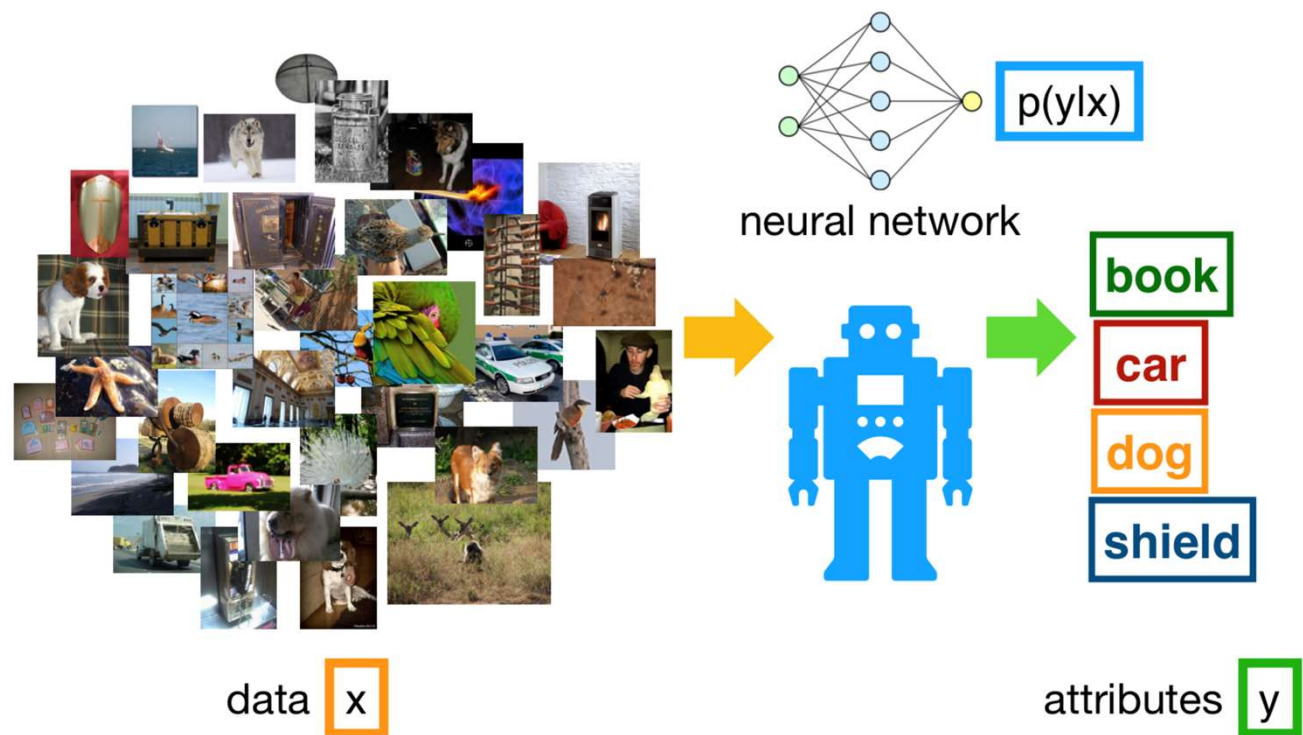
How does the number of parameters vary with the size of the training data ?

Learning Pillars



Supervised Learning

- Prediction of an output **y** given an input **x**



Regression



\$82000



\$55500



???

Predict results within *a continuous output*

Classification



Dog: 96%

Cat: 29%

Duck: 2%

Bird: 0%



Dog: 36%

Cat: 94%

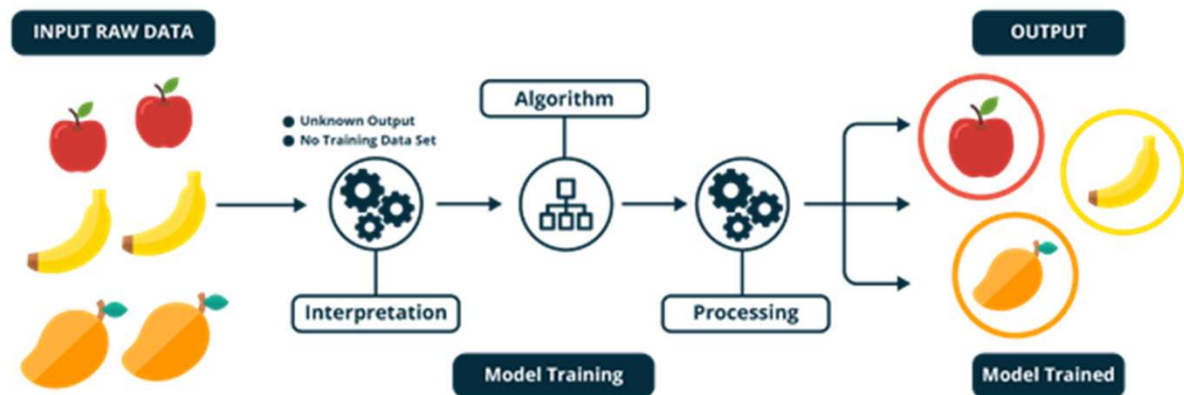
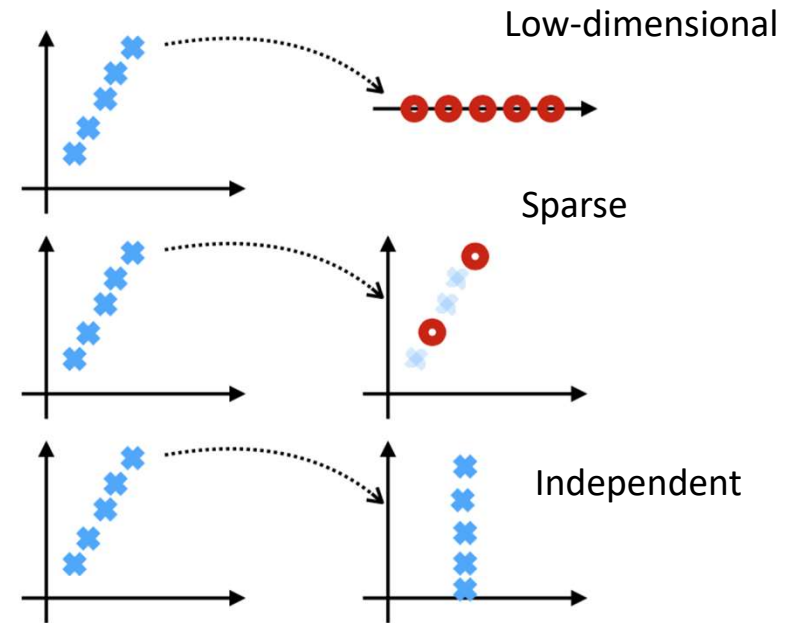
Duck: 2%

Bird: 1%

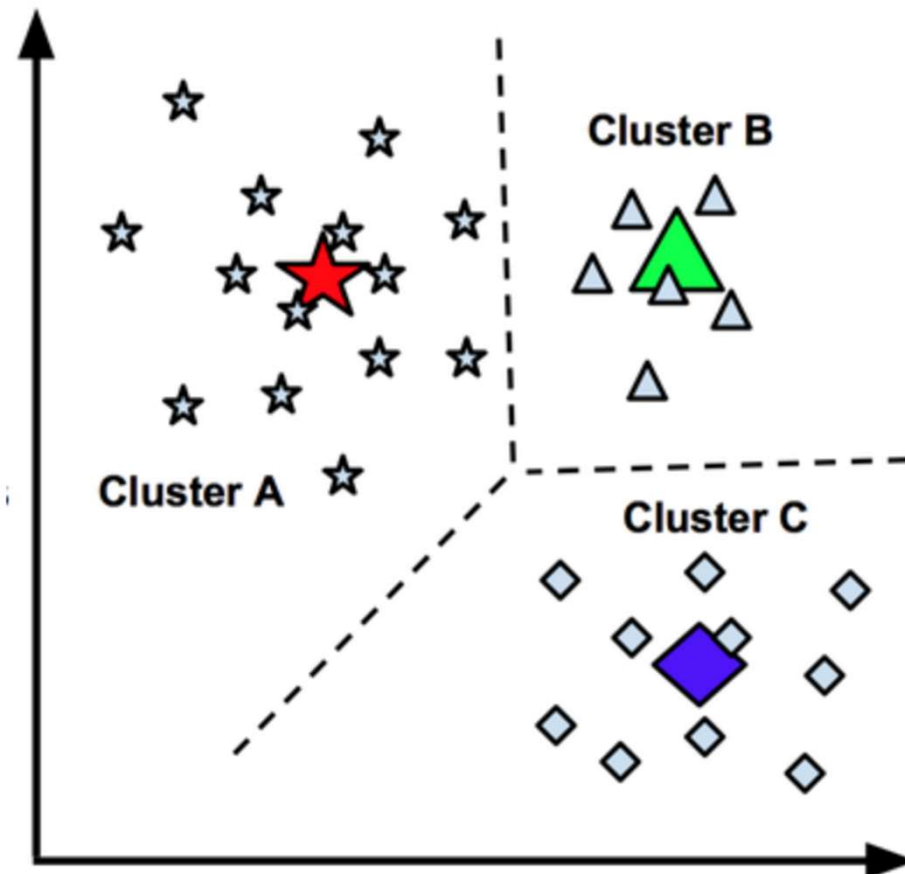
categorize new inputs as belonging to one of a set of categories
→ Predict results within a *discrete output (categories)*

Unsupervised Learning

- Find a *suitable data representation*
 - Preserving all task-relevant information
 - Simpler than the original data and easier to use



Clustering



Create a **set of categories**, for which individual data instances have a set of **common or similar characteristics**.

Representation Learning

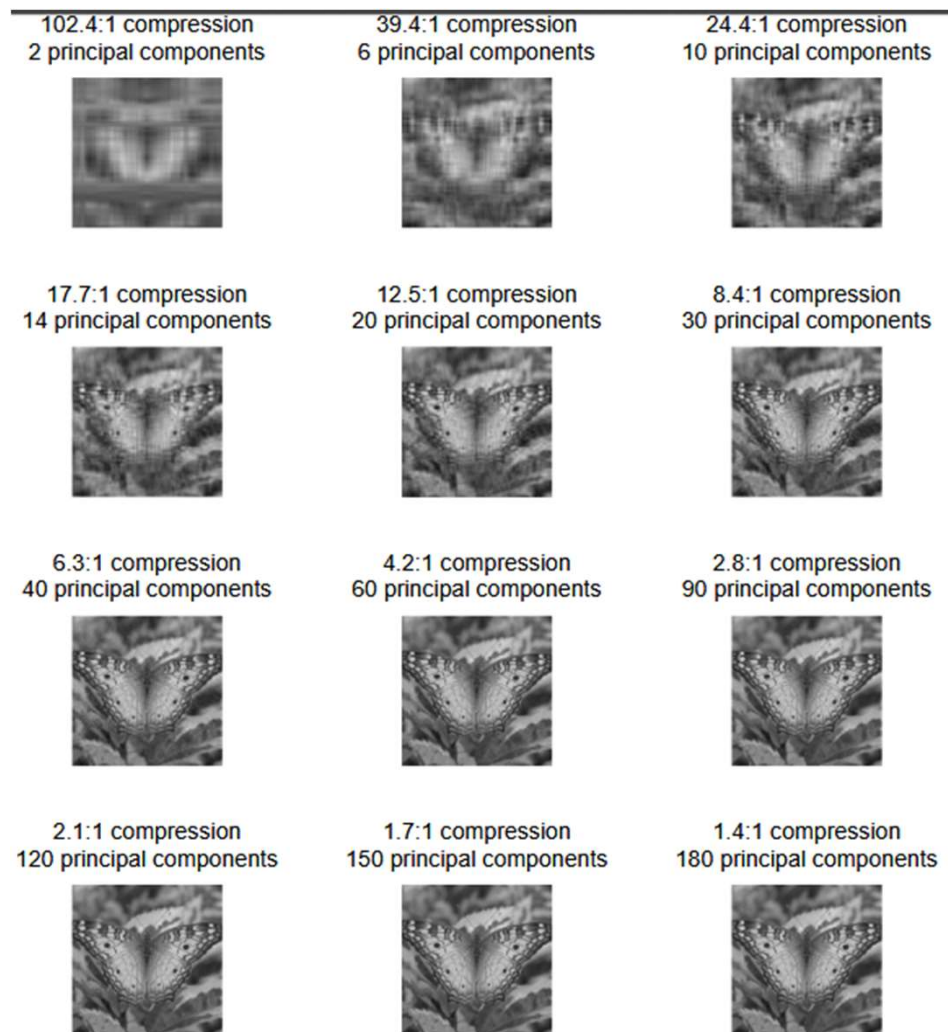


Figure 10: The visual effect of retaining principal components

More
nuances...



Supervised Classification



Semi-supervised Learning



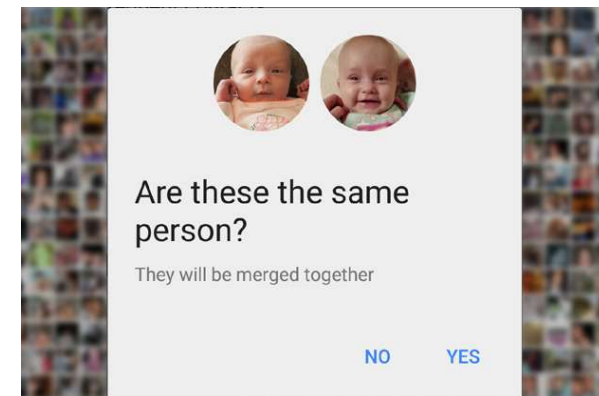
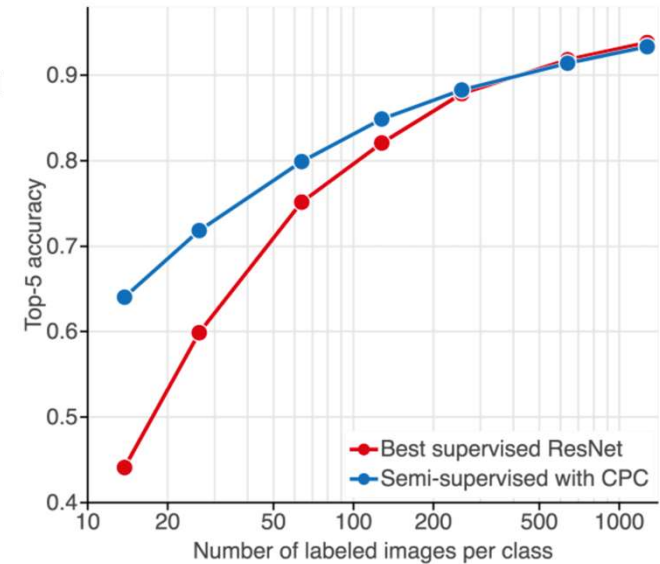
Transfer Learning



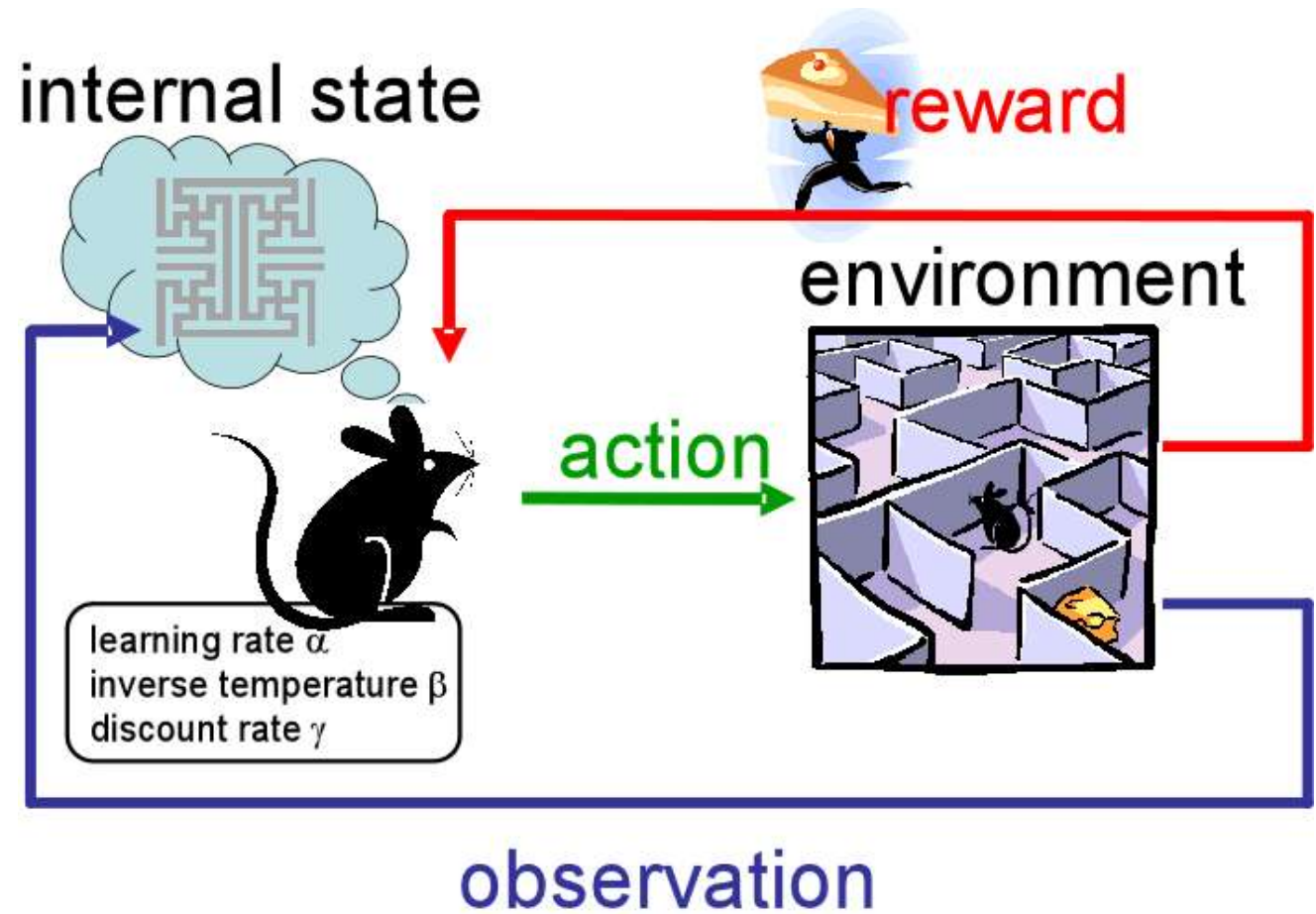
Self-taught Learning

Semi-supervised Learning

- Partially labelled training data
 - Lots of unlabeled data and a little bit of labelled data
- Why bother ?
- Most semi-supervised learning algorithms are combinations of unsupervised and supervised algorithms
 - Example : Google Photos



Reinforcement Learning





Two-Minute Papers

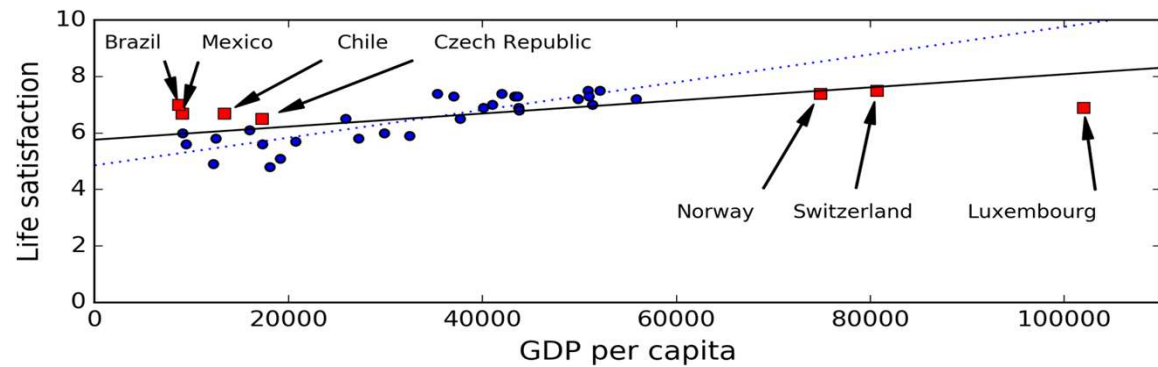


Main Challenges

- Bad data
- Bad model

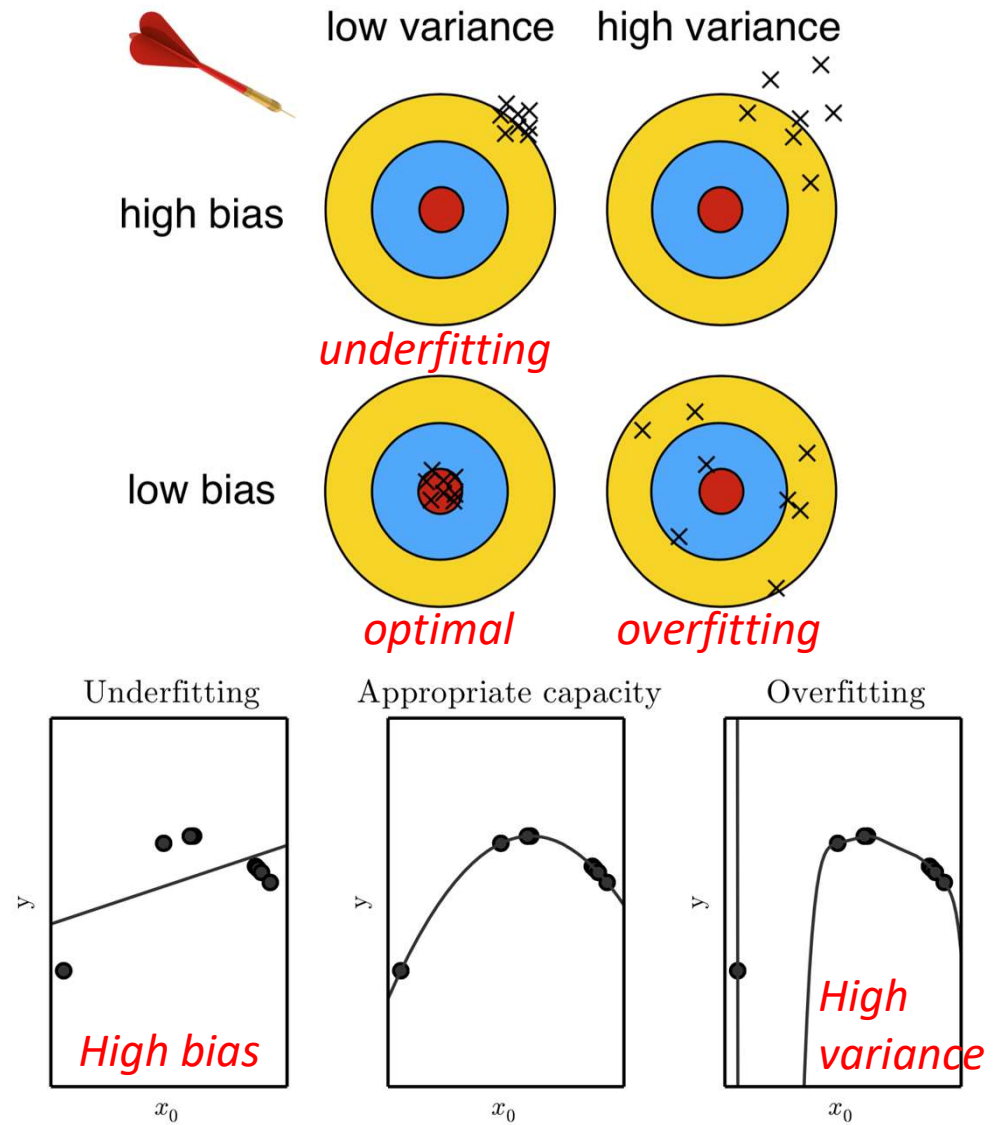
Bad Data

- **Insufficient** training data
 - Try to get more data, data augmentation,...
- **Non-representative** training data

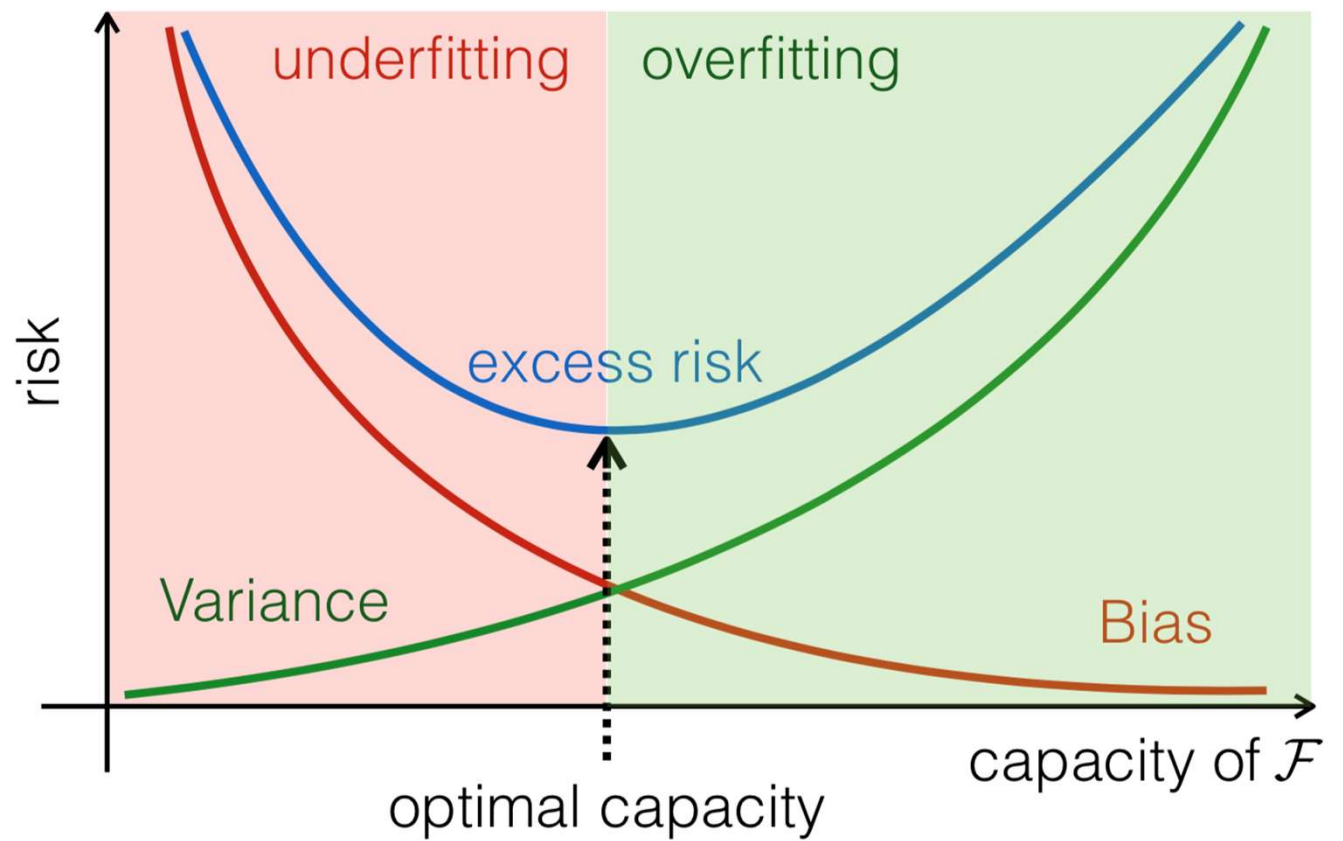


- **Poor quality** data
 - Data cleaning
- **Irrelevant** Features
 - **Feature selection** (select most useful features)
 - **Feature extraction** (combine features to produce a more useful one)

Model : Overfitting and underfitting



Model :
Bias,
variance and
capacity





Cross-validation

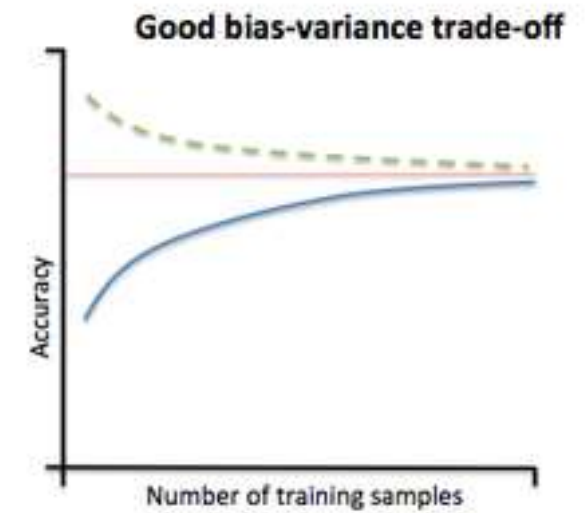
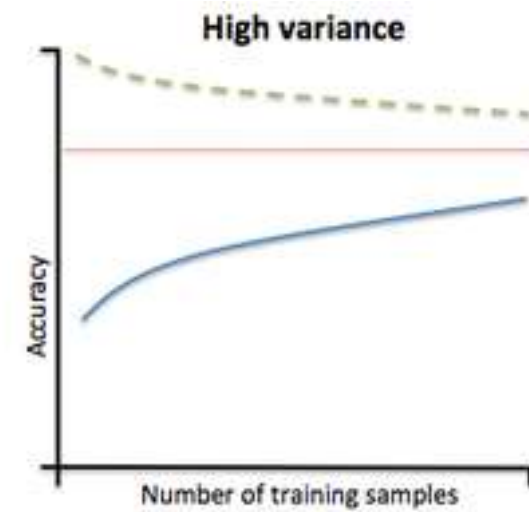
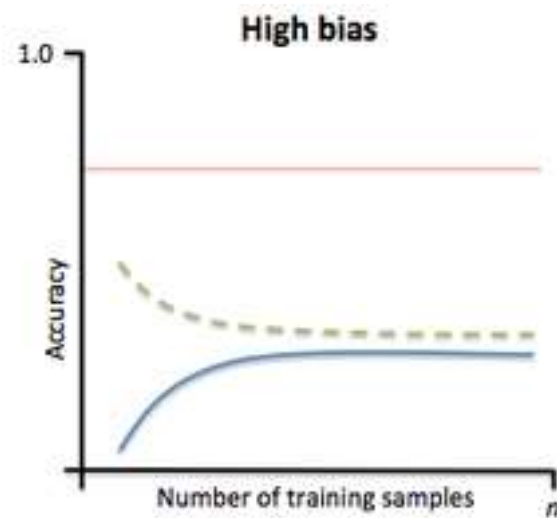
- **Validation technique** that assesses how the forecasts of a model generalize on unseen data
- A **score** can be computed on the forecasts across multiple test sets

Cross-
validation :
Holdout

- Separate the data into 2 or 3 sets
 - Training set for training
 - Validation set to find the best parameters
 - (Test set to estimate the performance in case of competing models)
- Separation depends on size of the dataset

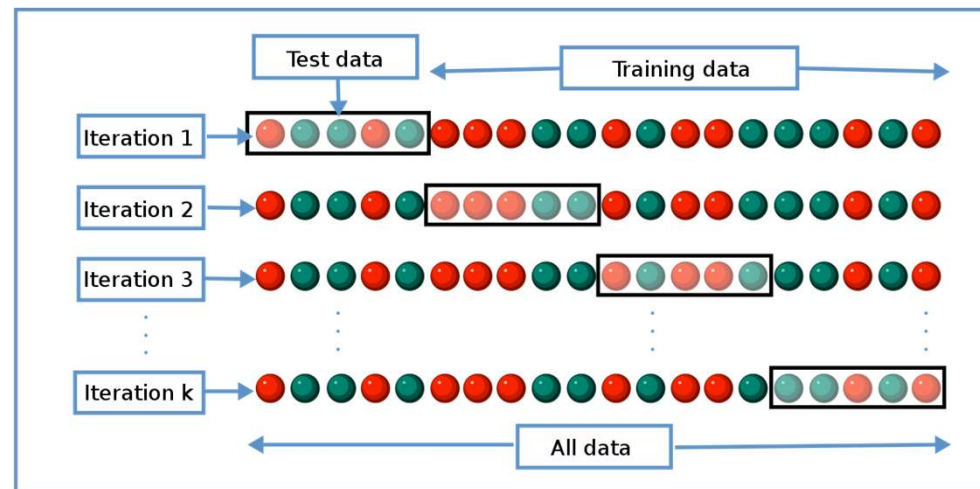


Bias-Variance Trade-off



Cross-validation : k-fold

- To avoid “wasting” too much training data
- test on **multiple splits** so we can get a better idea on how the model will perform on unseen data



- **In principle better than using the holdout method** because the holdout method score is dependent on how the data is split into train and test sets



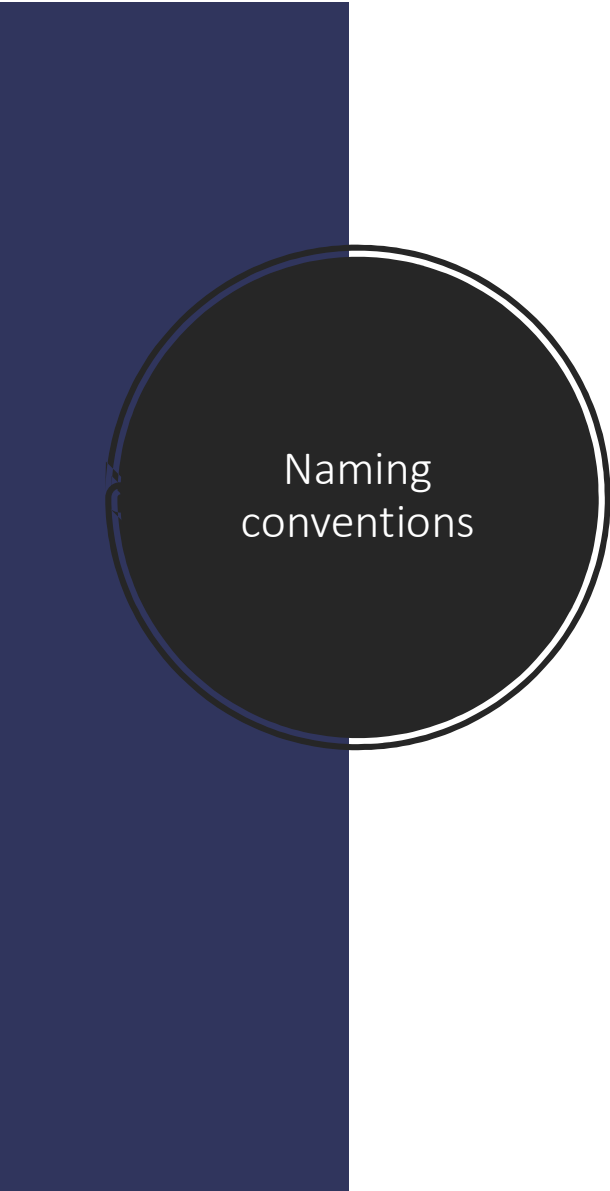
With which type of input data shouldn't it be used ?

Summary of cross- validation techniques

Name	Description
Leave-one-out	Use one observation as the test set, and the remaining observations as the train set. Out of N observations, there are N train-test splits.
Leave-p-out	Use p observations as the test set, and the remaining observations as the train set. Out of N observations, there are $\binom{N}{p}$ train-test splits.
Holdout	Randomly assign observations to the train and test sets. Caveat: There is a single run, and the result may not be representative.
K-fold	Split the N observations into K regular folds and apply a leave-one-out CV on the folds. The data may or may not be shuffled prior to forming the folds.
Combinatorial K-fold	Split the N observations into K regular folds and apply a leave-p-out CV on the folds. The data may or may not be shuffled prior to forming the folds.
Monte Carlo	Train sets are generated from random subsampling of the rows of (X, y) (random sampling without replacement)

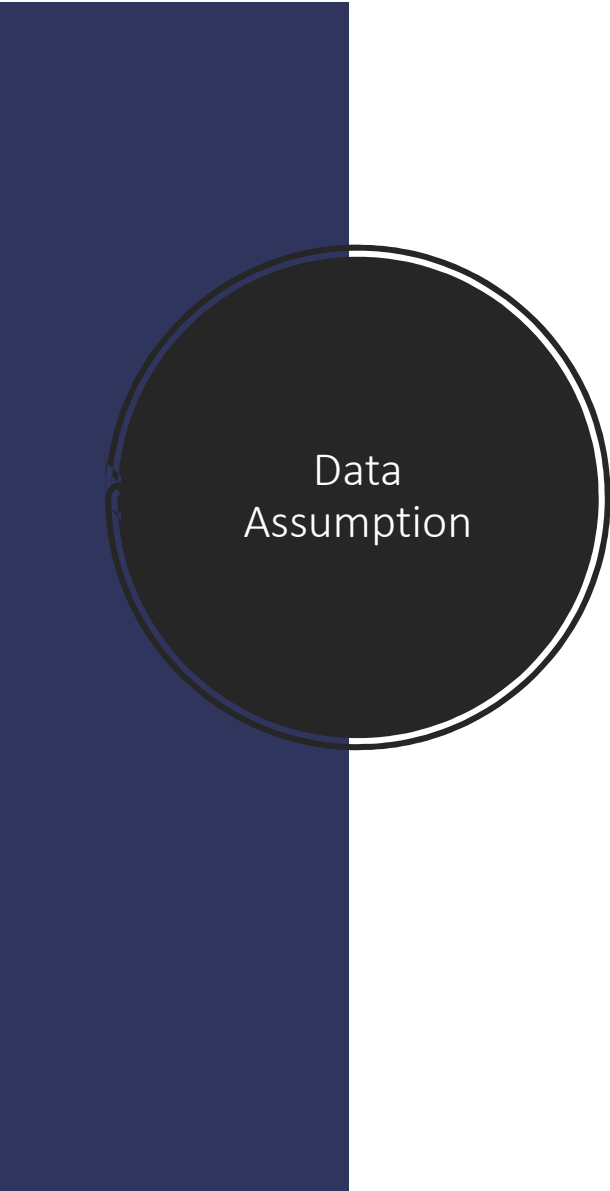


Data Preparation



Naming conventions

- **Data Cleaning** : identify and fix errors in the data prior to modeling (outliers, missing values, ...)
- **Feature Engineering** : create new input variables from raw data
- **Data Wrangling** : more general or colloquial term for data preparation that might include some data cleaning and feature engineering.



Data
Assumption

Independent and Identically Distributed (IID)

- 1) Come from the same distribution

$$p_{x^{(i)}}(x) = p_{x^{(j)}}(x)$$

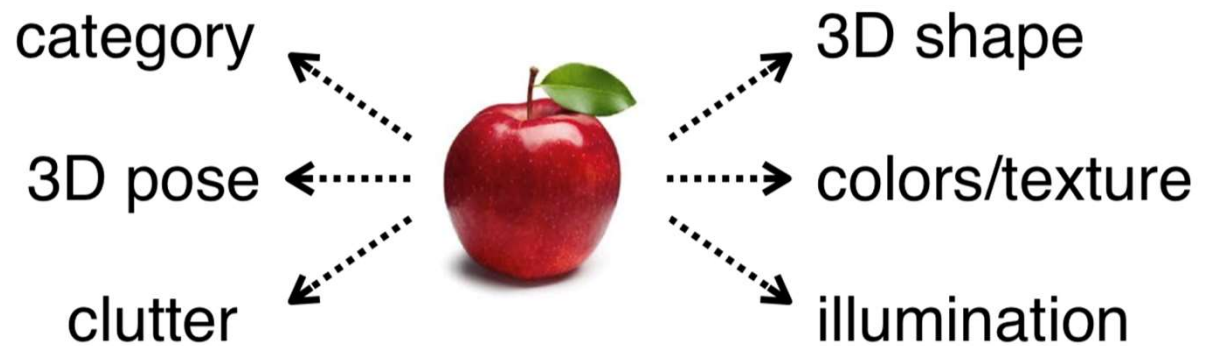
- 2) Are *independent*


$$p\left(x^{(1)}, \dots, x^{(m)}\right) = \prod_{i=1}^m p\left(x^{(i)}\right)$$

Data often encoded into **more focused relevant information** (features or internal representation) to simplify the decision

Data
Preparation

$$\text{data} \xrightarrow{\quad} x \rightarrow \phi(x) \xleftarrow{\quad} \text{feature}$$





Data Pre-Processing

- Convert labels into **one-hot encoded vectors** to avoid biases

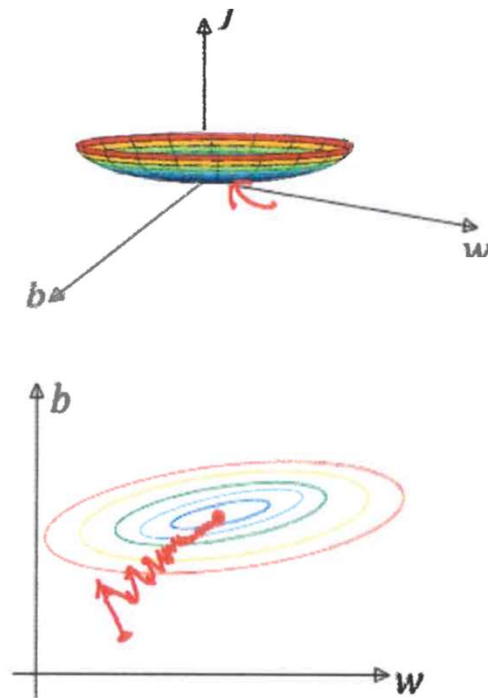
$[1\ 2\ 3\ 0\ 2] \rightarrow [0100\ 0010\ 0001\ 1000\ 0010]$

- **Shuffling** to ensure independent samples and an unbiased estimate
- Remove **under-represented** categories

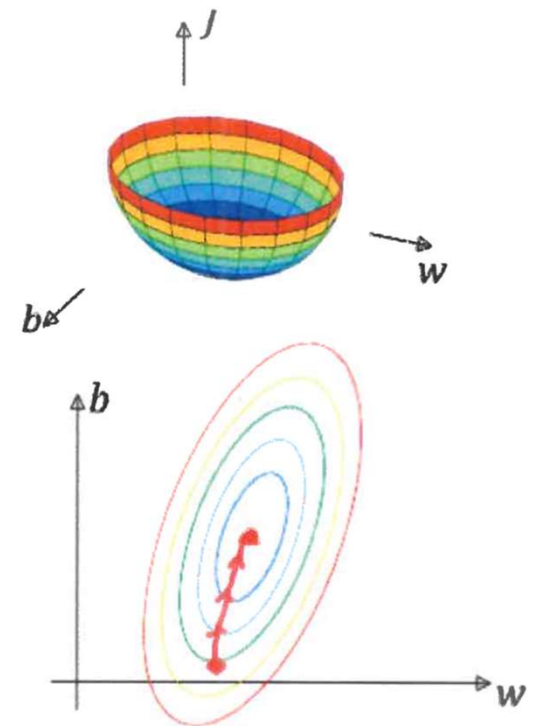
Standardization

- for faster convergence and more reliable predictions

Unnormalized :



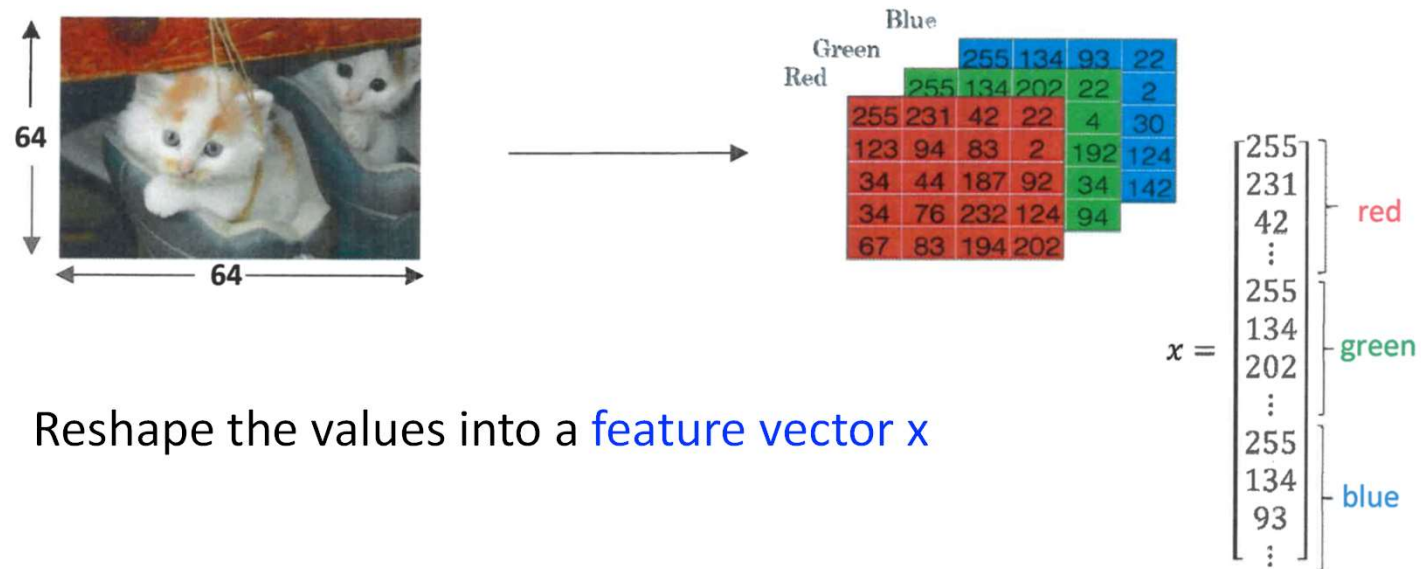
Normalized :



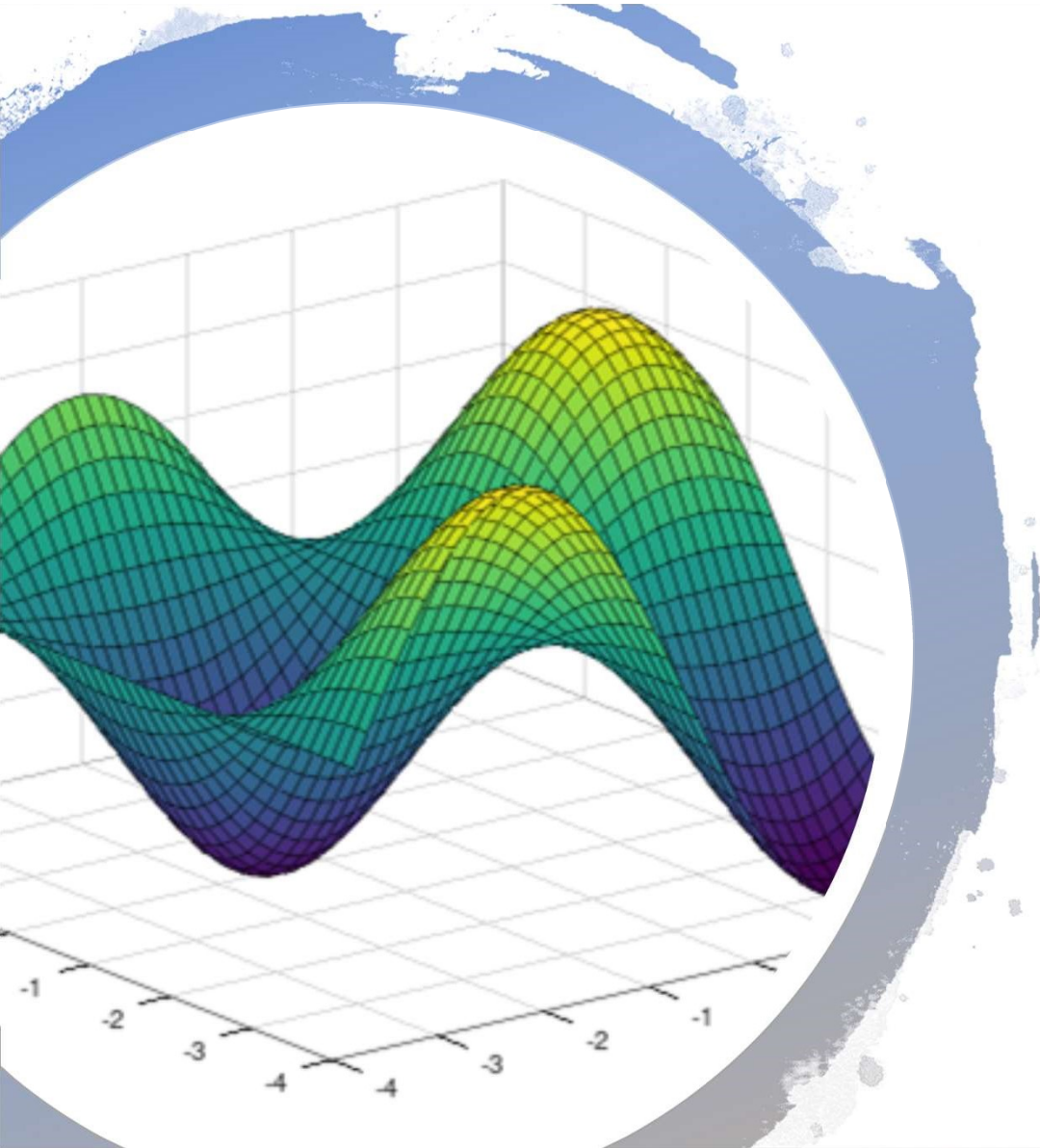
Reshaping into
feature vector

Example with binary classification :

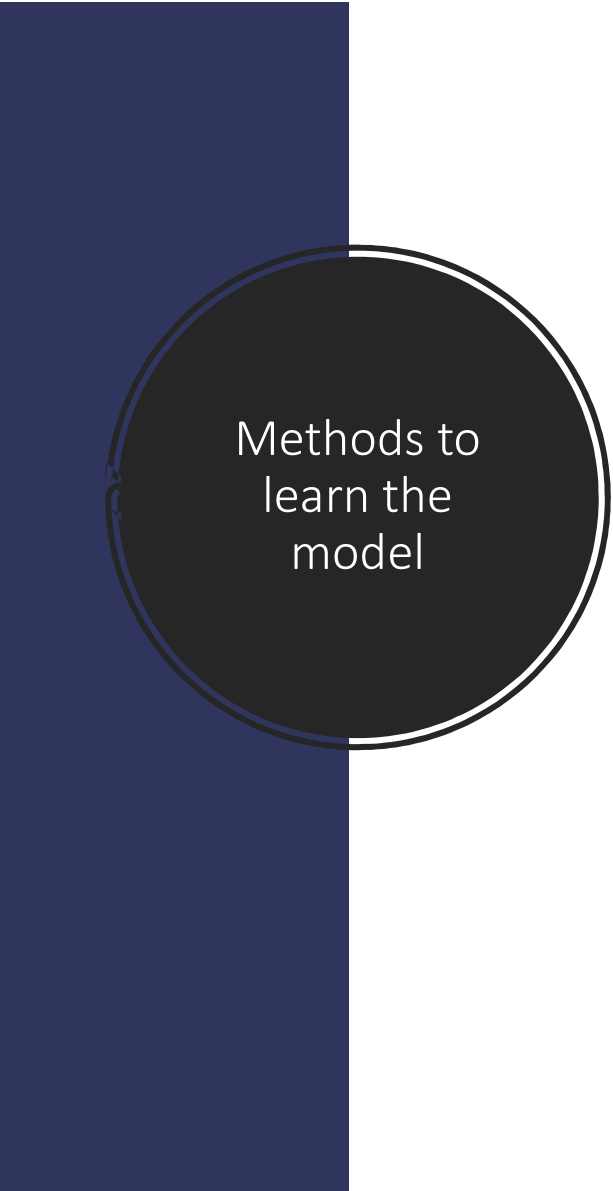
- Image stored in 3 matrices (red, green, blue color channels)
 - Similar size as the image



- Reshape the values into a **feature vector** x
- Normalization : divide by 255 (more convenient for data)



Optimization




Methods to
learn the
model

1) Simple linear regression

- use statistics (means, std, correlations and covariance)
- Need for all the data

2) Ordinary Least Squares

- Minimize the sum of the squared residuals
- Treats the data as a matrix and uses linear algebra to estimate the optimal values for the coefficients
- Need for all the data
- Need for enough memory



Loss and Cost Functions

- **Loss function** $L(\hat{y}^{(i)}, y^{(i)})$, also called error function, measures **how different** the prediction $\hat{y} = f(x)$ and the desired output y are

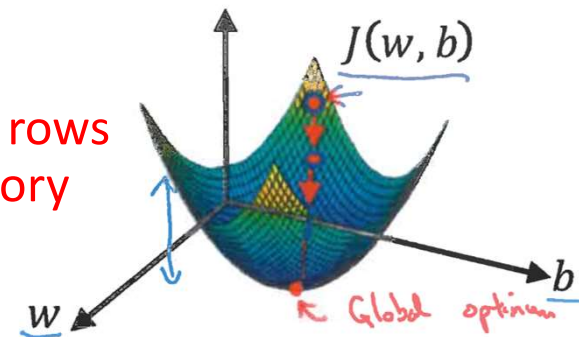
- **Cost function** $J(w, b)$ is the average of the loss function on the **entire training set**

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)})$$

- Goal of the optimization is to find the **parameters** $\theta = (w, b)$ that minimize the cost function

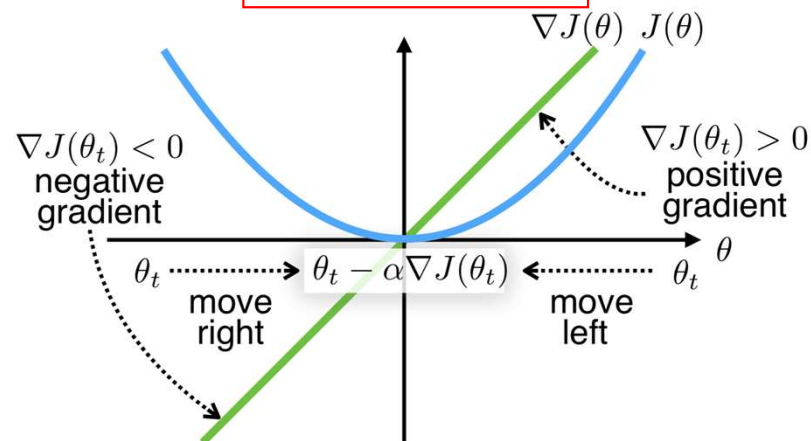
Gradient Descent

- **iterative** method to find the parameters $\theta = (w, b)$ that minimize $J(\theta)$
- Useful with very large datasets (either rows or columns) that may not fit into memory
- Learning rate α



gradient descent

$$\theta_{t+1} = \theta_t - \alpha \nabla J(\theta_t)$$

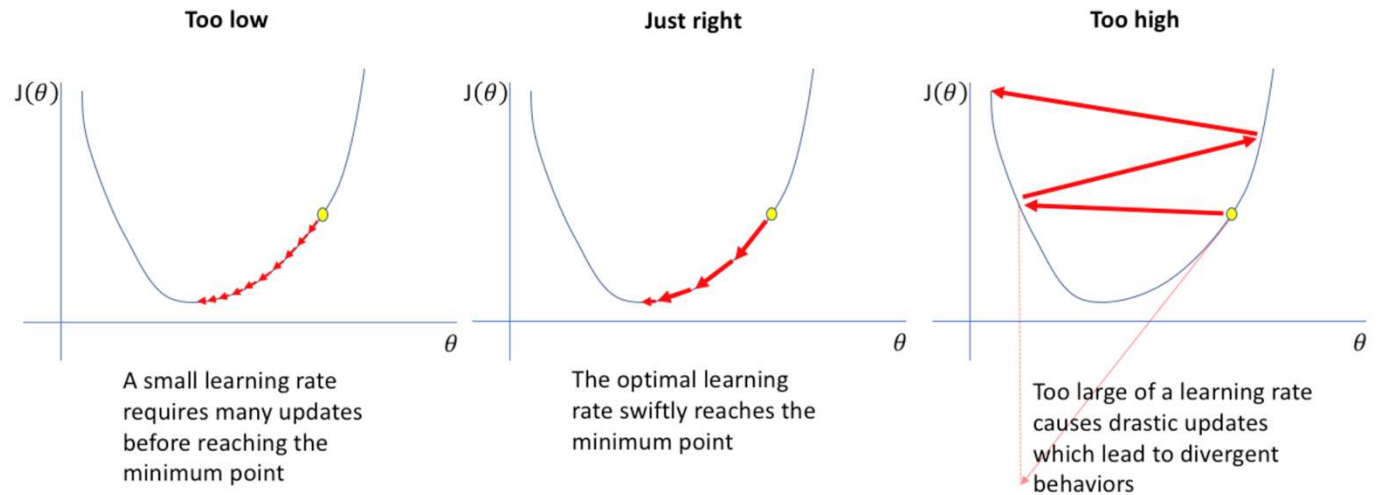


$$\nabla J(w) = \frac{dJ(w, b)}{dw}$$

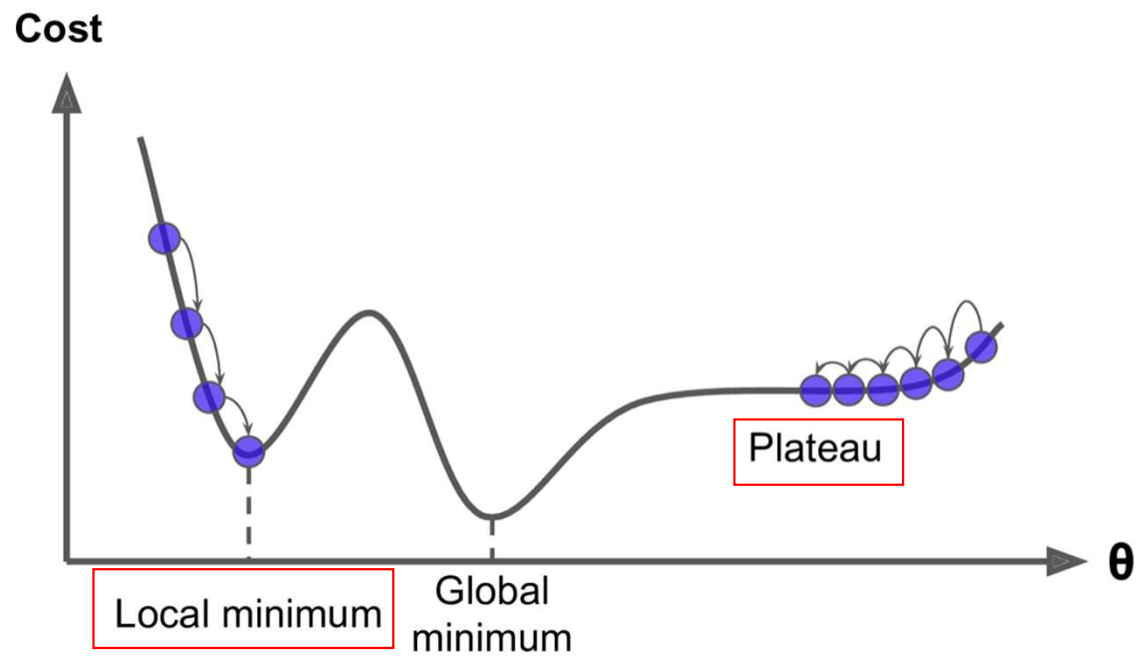
$$\nabla J(b) = \frac{dJ(w, b)}{db}$$

learning rate
 α

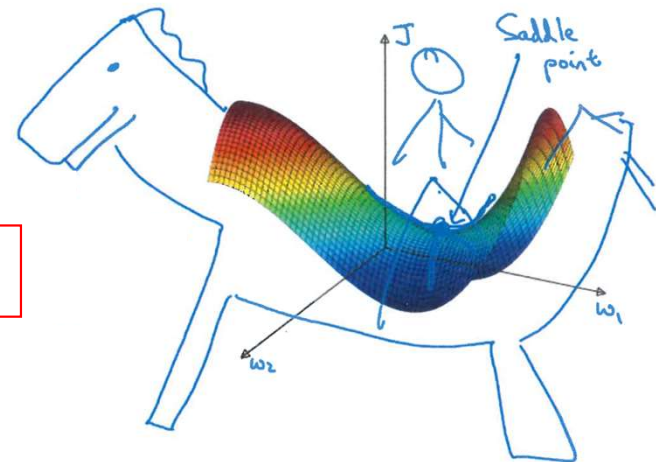
- Has a significant impact on the model performance, while being **one of the most difficult parameters** to set



Optimization Pitfalls



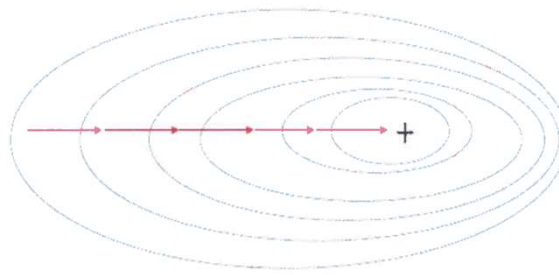
Saddle point



Mini-batch
Gradient Descent
and Stochastic
Gradient Descent

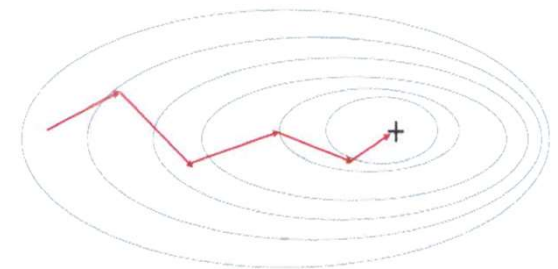
- At each iteration :

Gradient descent (GD)



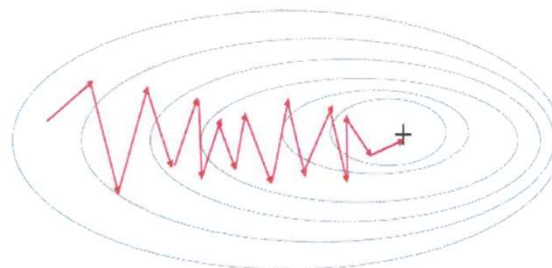
the *whole* training set

Mini-batch gradient
descent



a *batch* of samples

Stochastic gradient
descent (SGD)




1 sample

Batch size choice *typically*
32,64,128,256,512



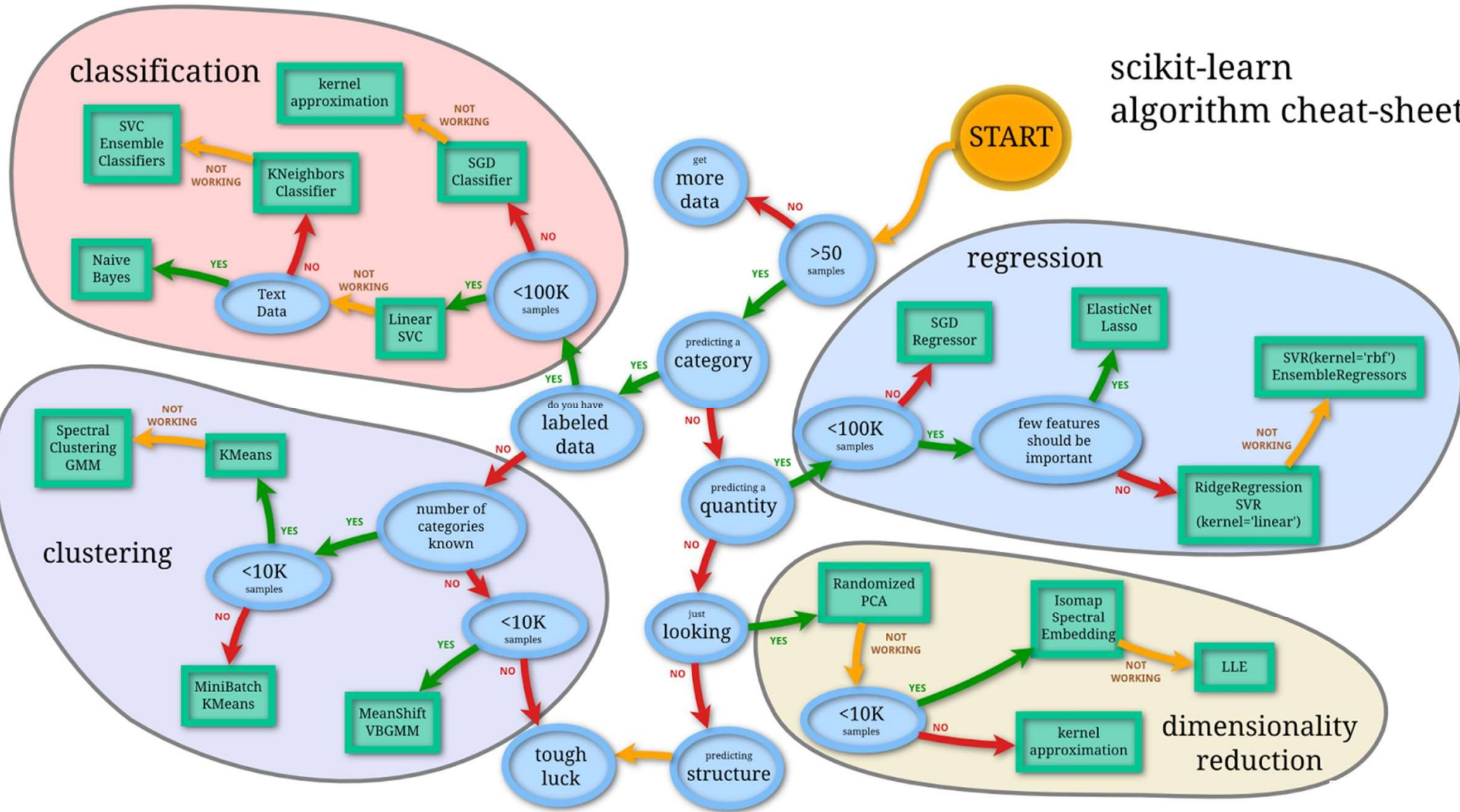
How to find its way in
the algorithm jungle ?



No Free Lunch Theorem

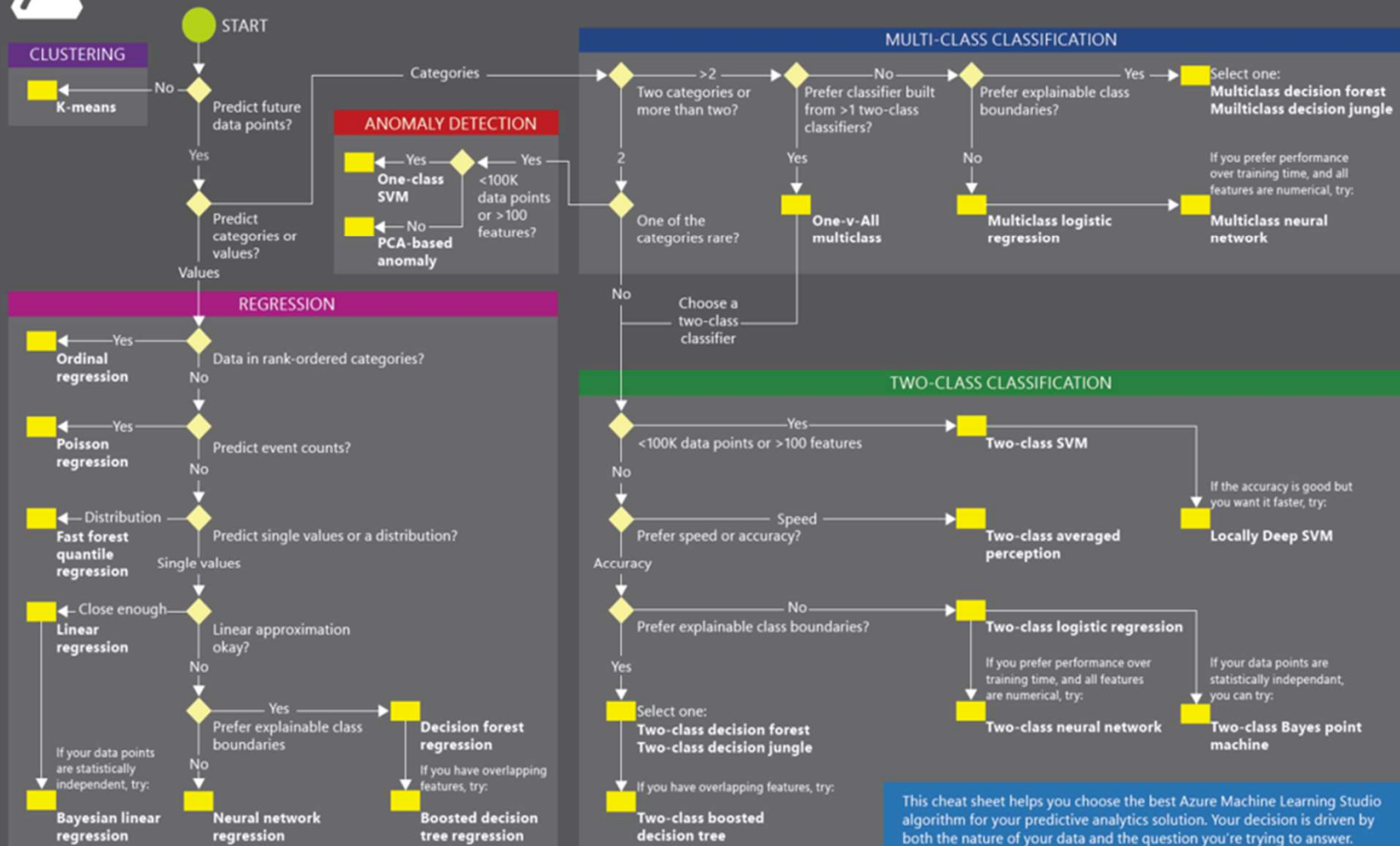
- If you make **no assumption about the data**, there is no reason to prefer one model over any other
- No model that is **a priori guaranteed to work better**
- In practice, **impossible to test all the models**, so make some **reasonable assumptions** about the data and evaluate only a few models

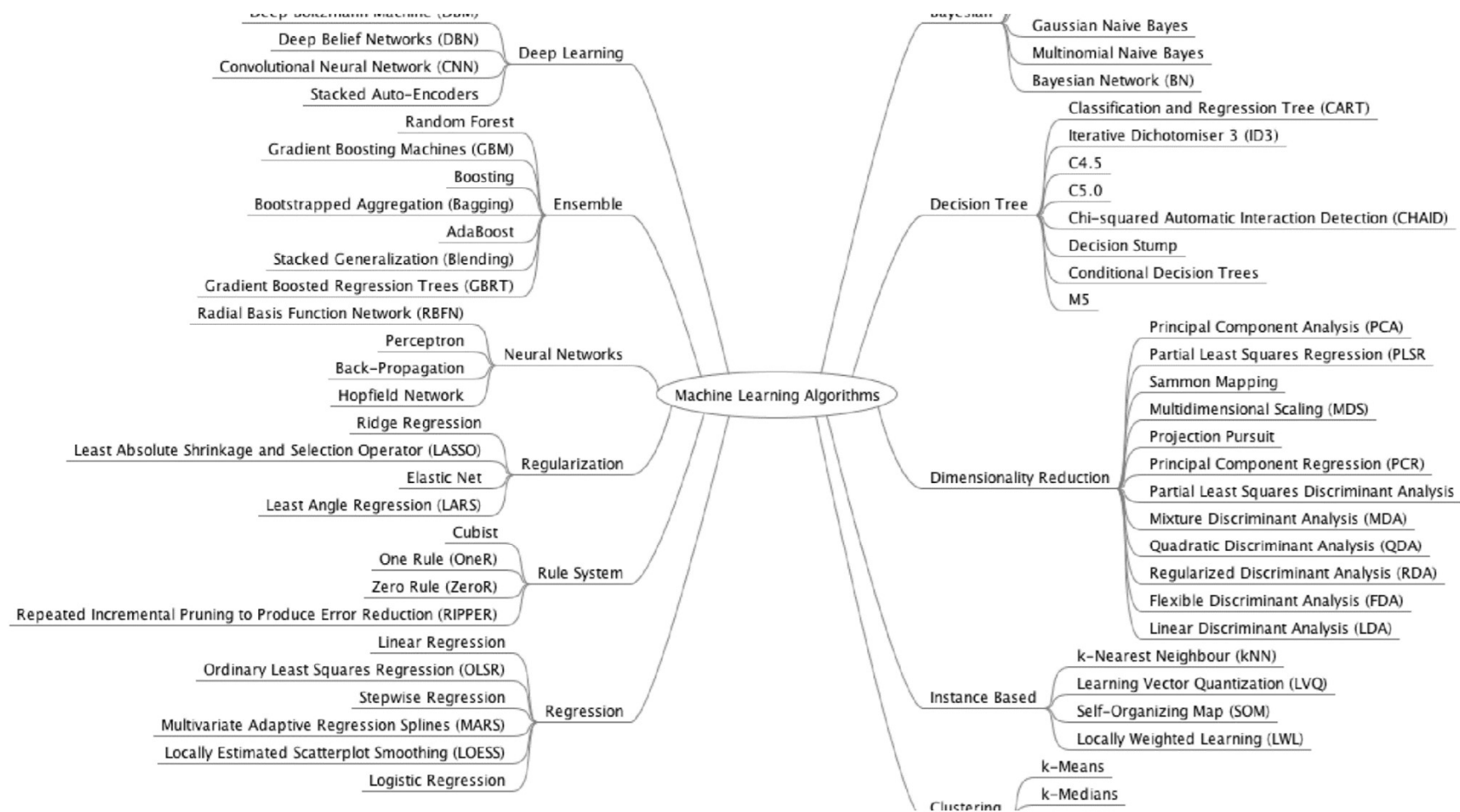
scikit-learn algorithm cheat-sheet





Microsoft Azure Machine Learning: Algorithm Cheat Sheet





The Lesson

The lesson is always to fit a simple model first, and then only adopt a more complicated ML model if the extra predictive accuracy (value) it provides is worth it.





Exercise

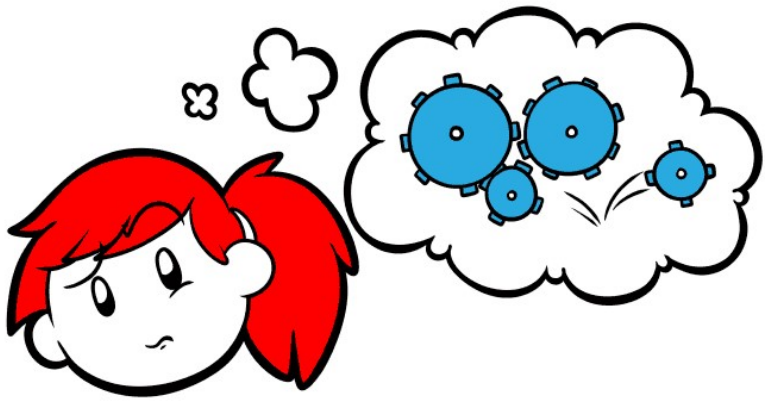
- Create your own summary (graph) for :
 - Regression
 - Classification



Quiz

<https://b.socrative.com/login/student/>

Room : CONTI6128

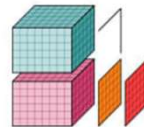


Technical Details

Tools

data structure & analysis

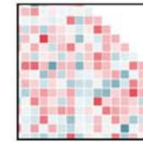
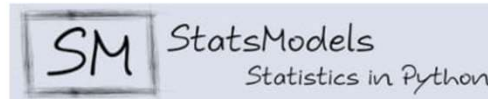
pandas
 $y_{it} = \beta' x_{it} + \mu_i + \epsilon_{it}$



xarray



scikit-image
image processing in python



Seaborn

main package for scientific computing in Python



matplotlib



Python-based ecosystem of open-source software for mathematics, science, and engineering.



IP[y]:
IPython

famous library to plot graphs in Python

interactive coding environments embedded in a webpage

provides simple and efficient tools for data mining and data analysis

h5py : common package to interact with a dataset that is stored on an H5 file



Real Data

- <http://archive.ics.uci.edu/lm/>
- <https://www.Kaggle.com/datasets>
- <http://aws.amazon.com/fr/dataset/>



Models

- <https://github.com/tensorflow/models>

Notebooks (1)

Use Anaconda Navigator

New environment with tensorflow **v.2.1.0:**

conda create -n tf tensorflow

conda activate tf

Add pip **v.20.1.1** :

conda install pip

Notebooks (2)

Add libraries :

conda install matplotlib

pip install sklearn

pip install seaborn

pip install ipywidgets

jupyter nbextension enable --py widgetsnbextension

Notebooks (3)

- **Define git command :**
- install git : <https://git-scm.com/download/win> (64-bit Git for Windows setup)
- Modifying PATH on Windows 10:
- In the Start Menu or taskbar search, search for "environment variable".
- Select "Edit the system environment variables".
- Click the "Environment Variables" button at the bottom.
- Double-click the "Path" entry under "System variables".
- With the "New" button in the PATH editor, add C:\Program Files\Git\bin\ and C:\Program Files\Git\cmd\ to the end of the list.
- Close and re-open your console.