# Big data science Day 3

F. Legger - INFN Torino
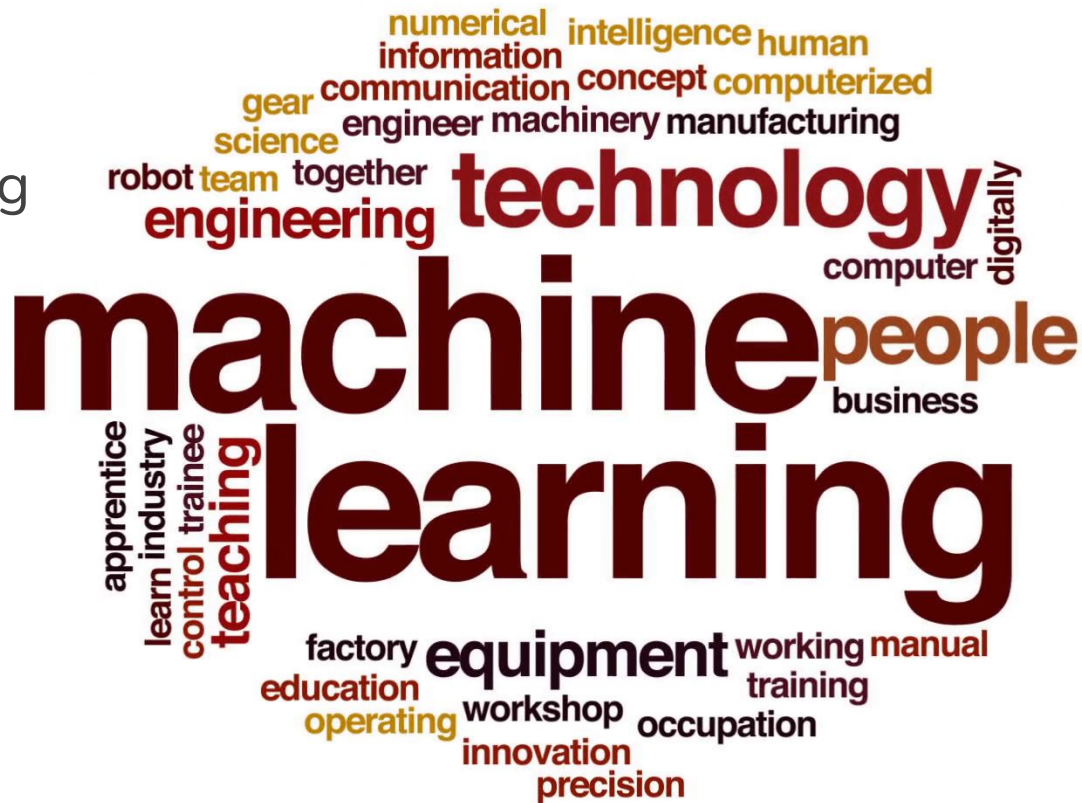https://github.com/Course-bigDataAndML/MLCourse-2324

# Yesterday

- Big data, analytics
- Distributed computing
- ML: Feature engineering

# Today

- **Machine learning**
  - Architectures
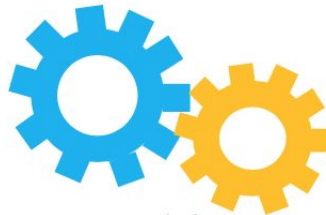  - Train model and evaluate

# Remember:

Machine learning involves **two mathematical entities**

**Model**: a mathematical model describes the relationship between different aspects of the data

**Features**: a feature is a representation of raw data

# Feature engineering

Raw Data

Feature Vector

Feature Engineering

0 : {
  house_info : {
  num_rooms: 6
  num_bedrooms: 3
  street_name: "Shorebird Way"
  num_basement_rooms: -1
  ...
  }
}

[
  6.0,
  1.0,
  0.0,
  0.0,
  0.0,
  9.321,
  -2.20,
  1.01,
  0.0,
  ...,
]

Raw data doesn't come to us as feature vectors.

Process of creating features from raw data is **feature engineering**.

4

# Example: supervised classification



## Ingredients

- **Inputs: X**, is a matrix of size **n** (number of samples) x **m** (number of features)
- **Features: X**, transformed inputs, matrix **n**x**m**
- **Labels: y,** vector size **n**

# Recipe: supervised classification
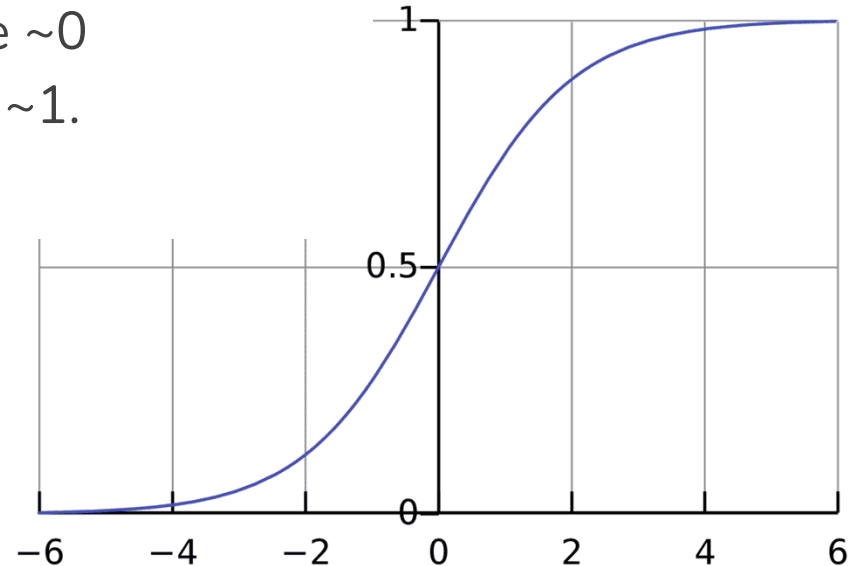
- For each input vector $X_i$ predict $z_i$, i=1...n

- **Predictions:**
  - $z_i = \varphi(W^T X_i)$ yields (0,1)

  - **Weights: W** (matrix) contains the model parameters
  - **Activation function: $\varphi$** (step function, sigmoid)

- **Cost function == loss function == prediction error**, function of the model parameters **W**
- Aim: find weights W that <u>minimize cost function</u>
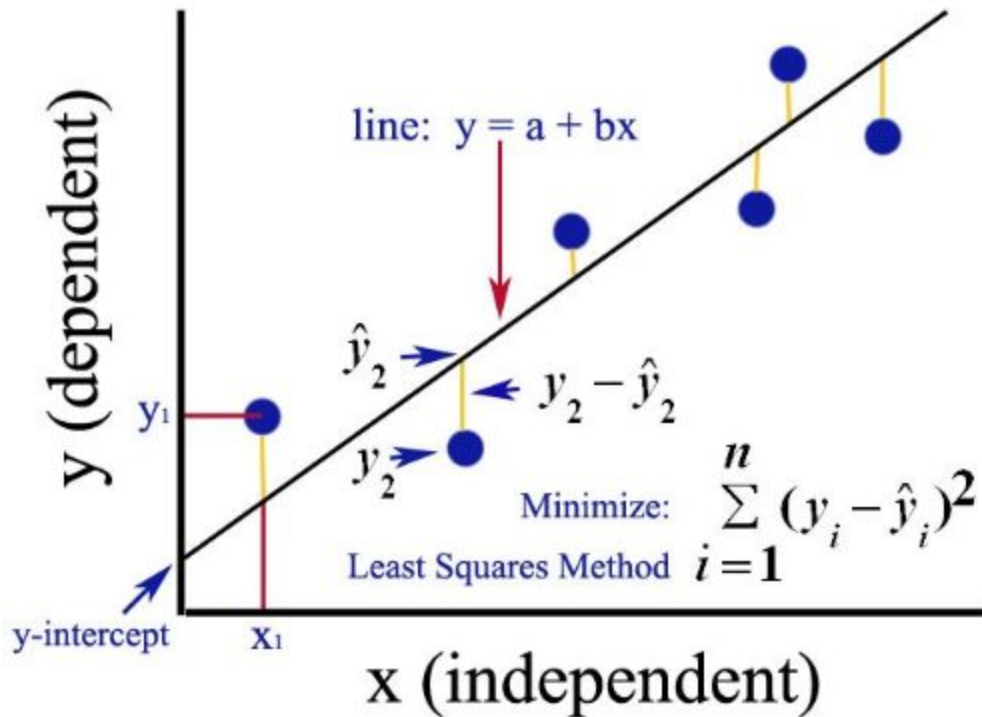
# Activation function

- Turns unbounded output into a known range/shape
- For example, **sigmoid** function only outputs numbers in the range (0, 1)
  - big negative numbers become ~0
  - big positive numbers become ~1.

$$S(x) = \frac{1}{1 + e^{-x}} = \frac{e^x}{e^x + 1}.$$
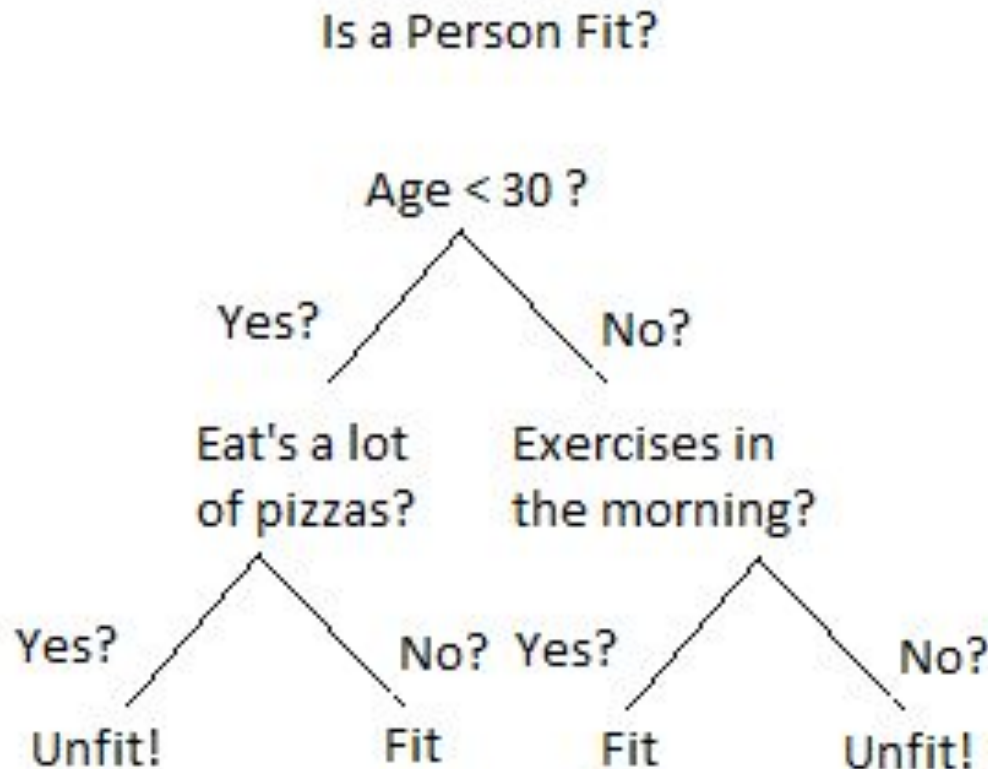
# Another example, linear regression

- Inputs (features): $x_i$
- Labels: $y_i$
- Model: **y = a + bx**
- Weight+bias (parameters to be found): **a, b**
- Cost function: **Mean Square Error (MSE)**
- No **activation function:** problem is linear

line: $y = a + bx$

$\hat{y}_2$

$y_2 - \hat{y}_2$

$y_1$

$y_2$

Minimize: $\sum\limits_{i=1}^{n} (y_i - \hat{y}_i)^2$

Least Squares Method

y-intercept    $x_1$

y (dependent)

x (independent)

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^{n} (Y_i - \hat{Y}_i)^2.$$

# Supervised learning

# Decision trees: supervised classification



Is a Person Fit?

Age < 30 ?

Yes? / No?

Eat's a lot of pizzas?    Exercises in the morning?

Yes? / No?    Yes? / No?

Unfit!    Fit    Fit    Unfit!

Typically used in combinations (Random forest, Gradient Tree Boosting)

# Random Forest

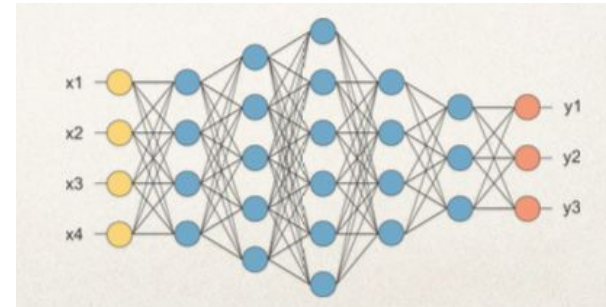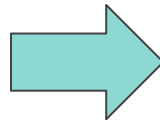- Each tree sees part of the training sets and captures part of the information it contains
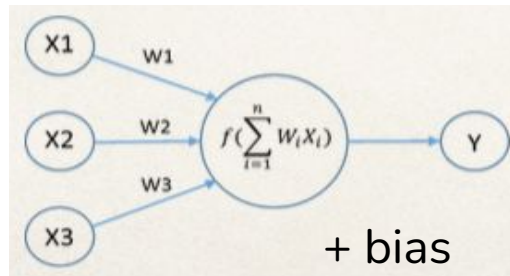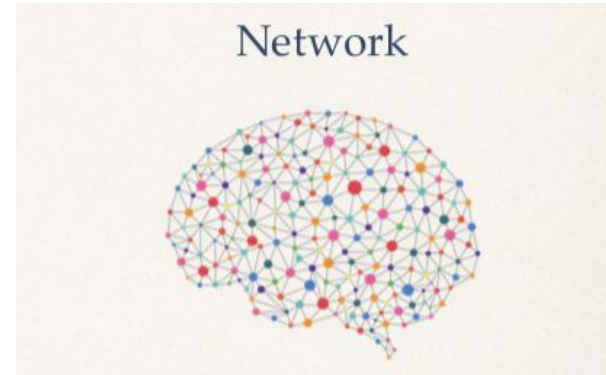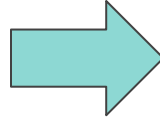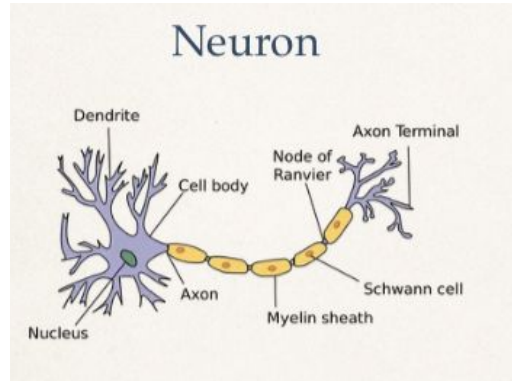
# Ensembles

- **Bagging**
  - building multiple models (typically of the *same* type) from different subsamples of the training dataset
- **Boosting**
  - building multiple models (typically of the *same* type) each of which learns to fix the predictions errors of a prior model in the chain
- **Stacking**
  - building multiple models (typically of *different* types) and a supervisor model that learns how to best combine the predictions of the primary model
- **Weighting|Blending**
  - combine multiple models into single prediction using different weight functions
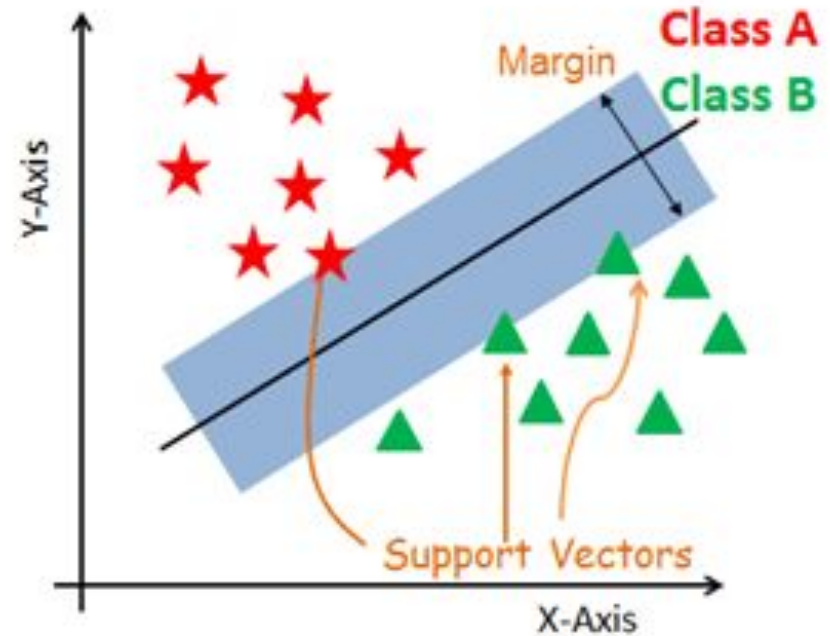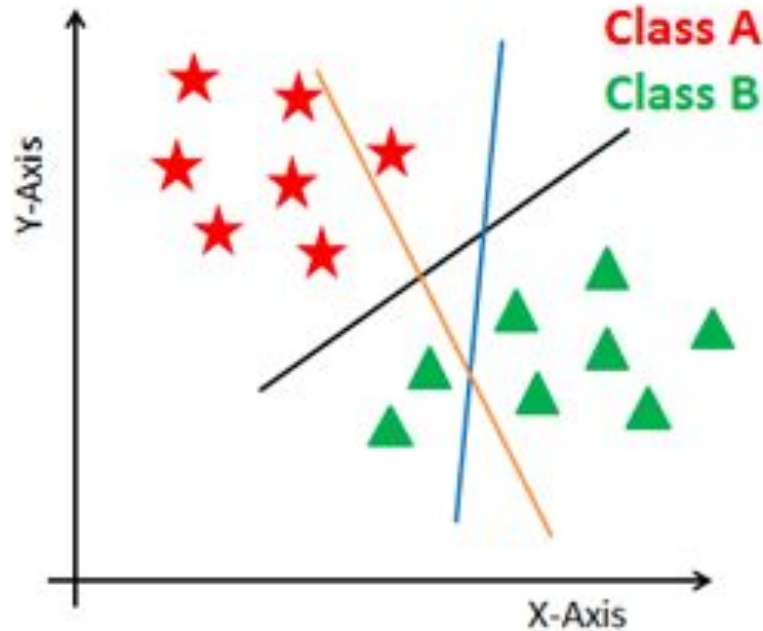
# Neural networks: supervised classification

- Basic unit: **Neuron.** A neuron takes inputs, does some math with them, and produces **one** output
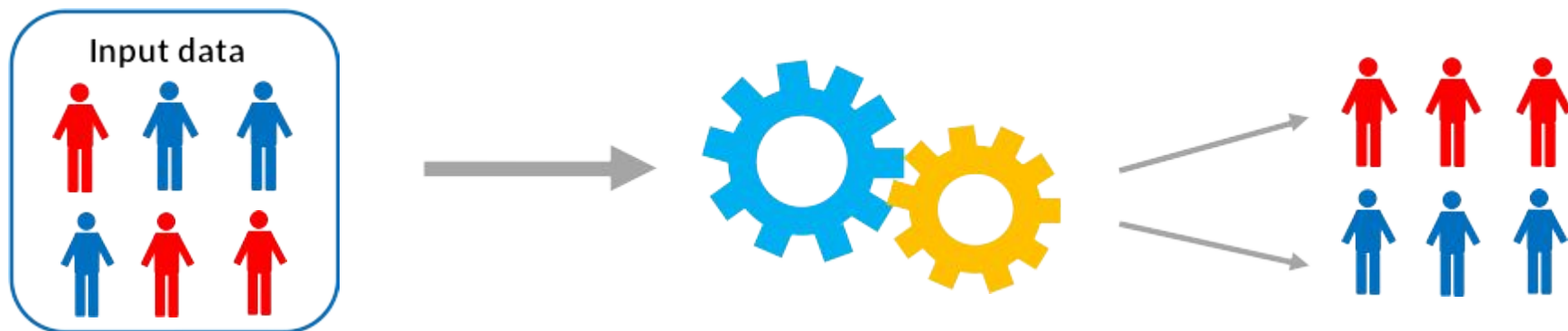
*More on NN tomorrow!*

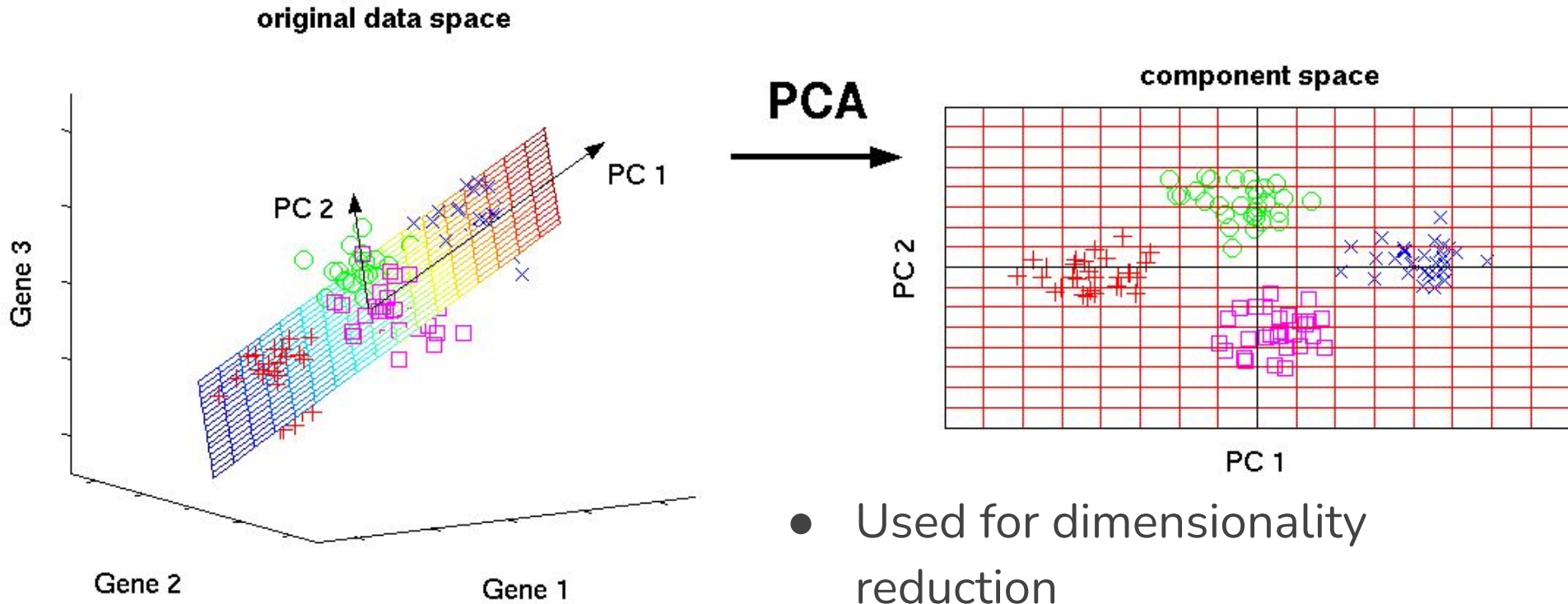# Support vector machines (SVG), supervised classification
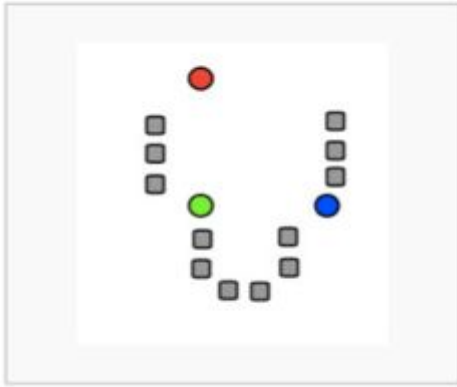
# Unsupervised learning

# Challenges

- No label (ground truth) in input dataset
- The system must have the ability to recognize patterns in the data without explicitly being told what patterns to identify
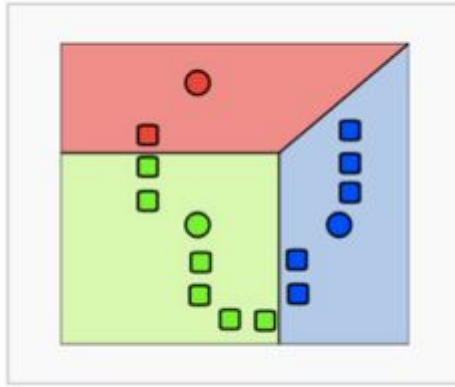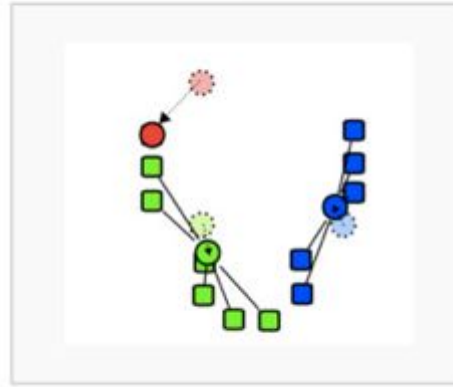
# Principal Component Analysis (PCA), unsupervised



- Used for dimensionality reduction
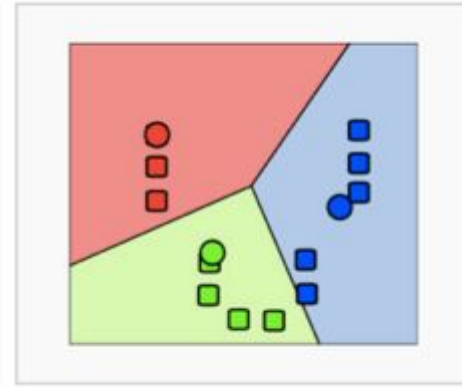
# K-means clustering, unsupervised



1. *k* initial "means" (in this case *k*=3) are randomly generated within the data domain (shown in color).

2. *k* clusters are created by associating every observation with the nearest mean. The partitions here represent the Voronoi diagram generated by the means.

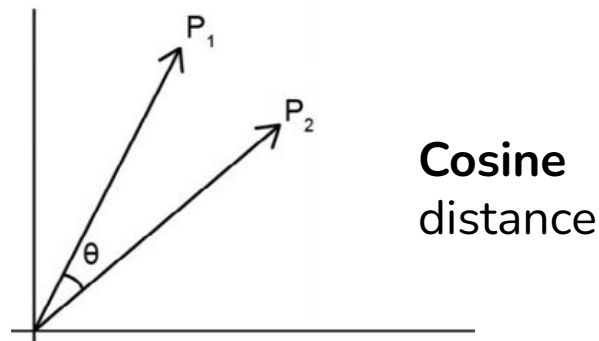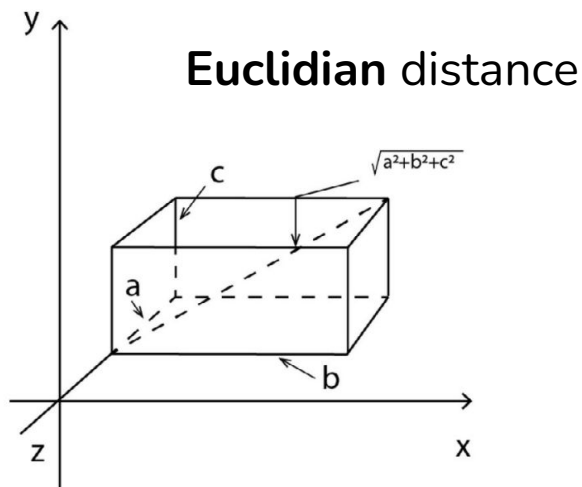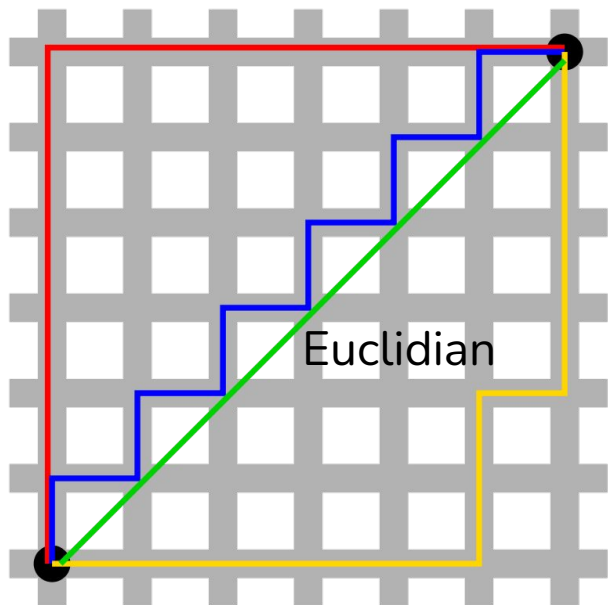3. The centroid of each of the *k* clusters becomes the new mean.

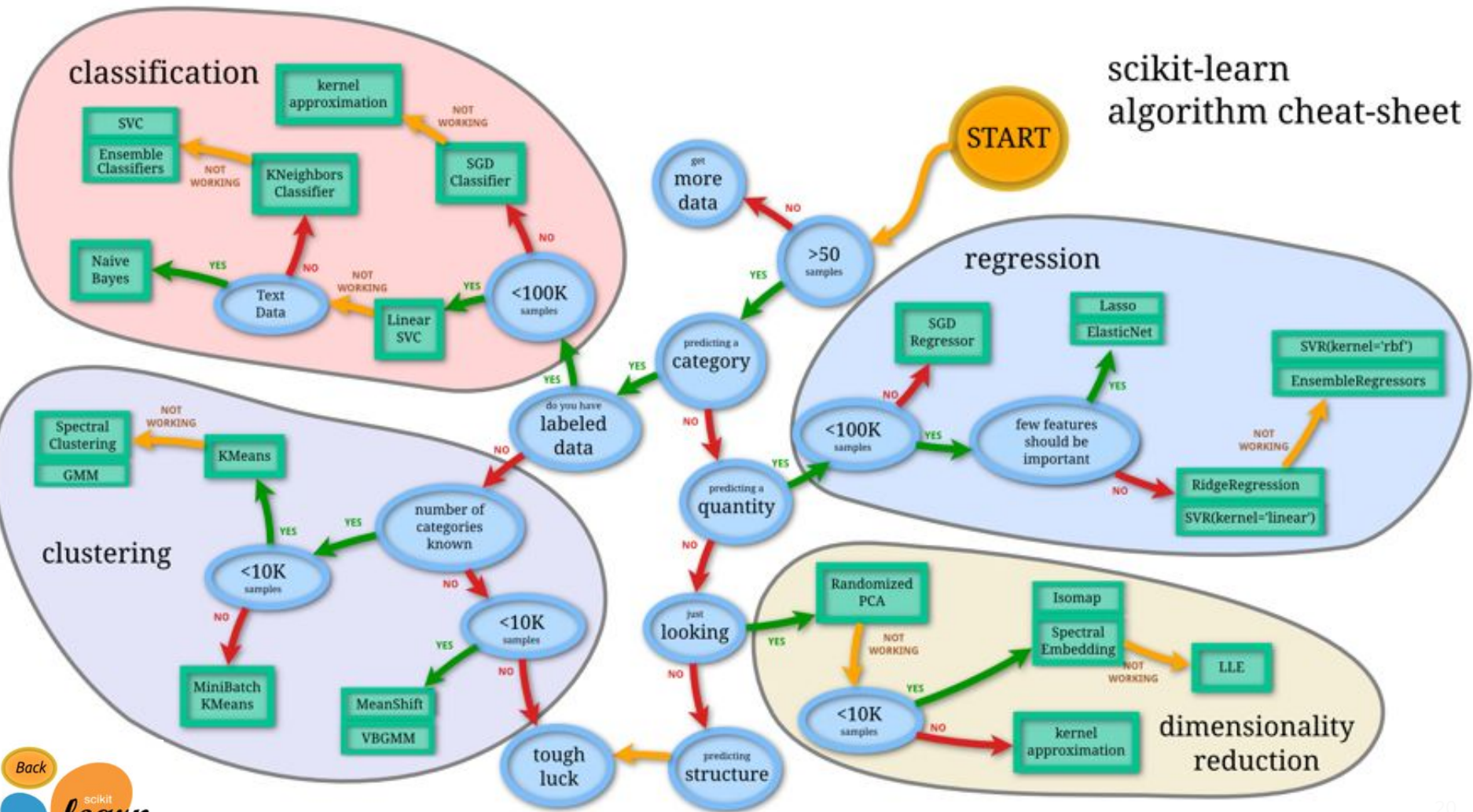4. Steps 2 and 3 are repeated until convergence has been reached.

You must define **k, the number of clusters,** and which **distance** to use!

# K-means: Distances

**Euclidian** distance

**Taxicab** or **Manhattan** distance: sum of the projections along all axis



Euclidian



**Cosine** distance

scikit-learn algorithm cheat-sheet

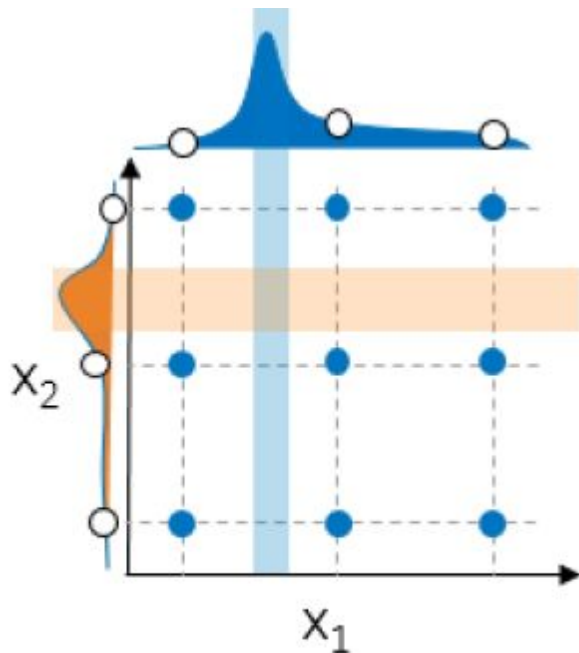| TYPE | NAME | DESCRIPTION | ADVANTAGES | DISADVANTAGES |
|---|---|---|---|---|
| Linear | Linear regression | The "best fit" line through all data points. Predictions are numerical. | Easy to understand — you clearly see what the biggest drivers of the model are. | **X** Sometimes too simple to capture complex relationships between variables. <br> **X** Tendency for the model to "overfit". |
| Linear | Logistic regression | The adaptation of **linear regression** to problems of classification (e.g., yes/no questions, groups, etc.) | Also easy to understand. | **X** Sometimes too simple to capture complex relationships between variables. <br> **X** Tendency for the model to "overfit". |
| Tree-based | Decision tree | A graph that uses a **branching method** to match all possible outcomes of a decision. | Easy to understand and implement. | **X** Not often used on its own for prediction because it's also often too simple and not powerful enough for complex data. |
| Tree-based | Random Forest | Takes the average of many decision trees, each of which is made with a sample of the data. Each tree is weaker than a full decision tree, but **by combining them we get better overall performance.** | A sort of "wisdom of the crowd". Tends to result in very high quality models. Fast to train. | **X** Can be slow to output predictions relative to other algorithms. <br> **X** Not easy to understand predictions. |
| Tree-based | Gradient Boosting | Uses even weaker decision trees, that are increasingly **focused on "hard" examples.** | High-performing. | **X** A small change in the feature set or training set can create radical changes in the model. <br> **X** Not easy to understand predictions. |
| Neural networks | Neural networks | Mimics the behavior of the brain. Neural networks are interconnected neurons that pass messages to each other. Deep learning uses several **layers of neural networks put one after the other.** | Can handle extremely complex tasks - no other algorithm comes close in image recognition. | **X** Very, very slow to train, because they have so many layers. Require a lot of power. <br> **X** Almost impossible to understand predictions. |

*All models are wrong, but some are useful (George Box)*
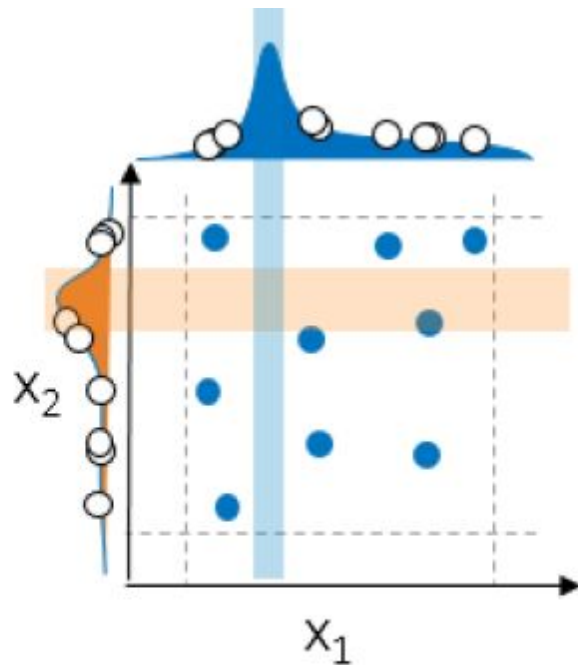
# Hyperparameters vs parameters

- Model **parameters** are learned during **training** when we optimize a loss function -> *weights*
- **Hyperparameters** are not model parameters and they **cannot** be directly trained from the data -> *model architecture*

| Hyperparameters | | Parameters | | Score |
|---|---|---|---|---|
| n_layers = 3<br>n_neurons = 512<br>learning_rate = 0.1 | ➡ | Weights optimization | ➡ | 85% |
| n_layers = 3<br>n_neurons = 1024<br>learning_rate = 0.01 | ➡ | Weights optimization | ➡ | 80% |
| n_layers = 5<br>n_neurons = 256<br>learning  rate = 0.1 | ➡ | Weights optimization | ➡ | 92% |

# Hyperparameter tuning



(a) Standard Grid Search

(b) Random Search

# Modeling Algorithm

## *Tune*

### hyperparameters (tuning options)

- Polynomial order, penalty parameter, ...
- Network configuration, solver options, ...
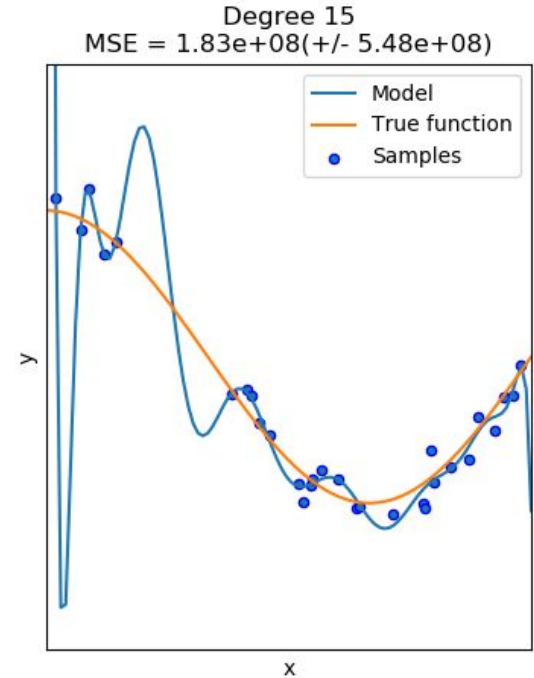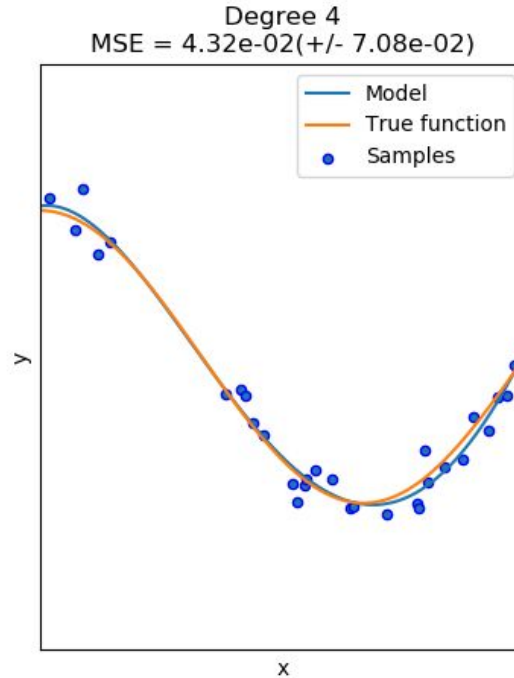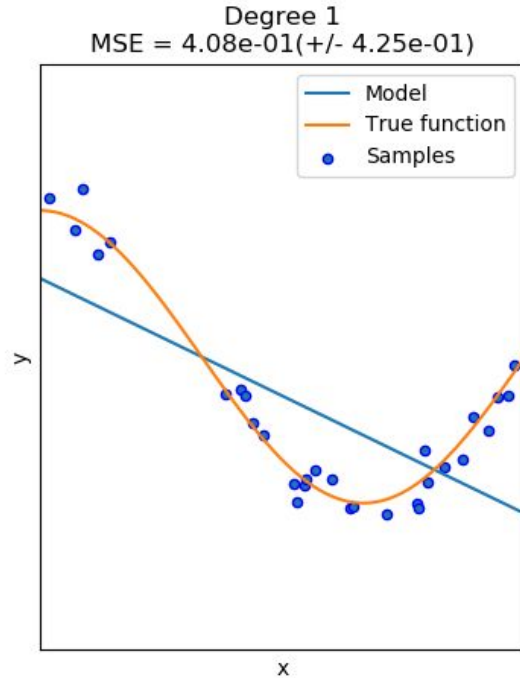- Max tree depth, splitting criterion, ...

## Model

### *Train*

#### model parameters

- Regression coefficients
- Neural net weights
- Tree splitting rules
- ...

# Overfitting / underfitting



Degree 1
MSE = 4.08e-01(+/- 4.25e-01)

Degree 4
MSE = 4.32e-02(+/- 7.08e-02)

Degree 15
MSE = 1.83e+08(+/- 5.48e+08)

**Underfitting**
Model doesn't have enough
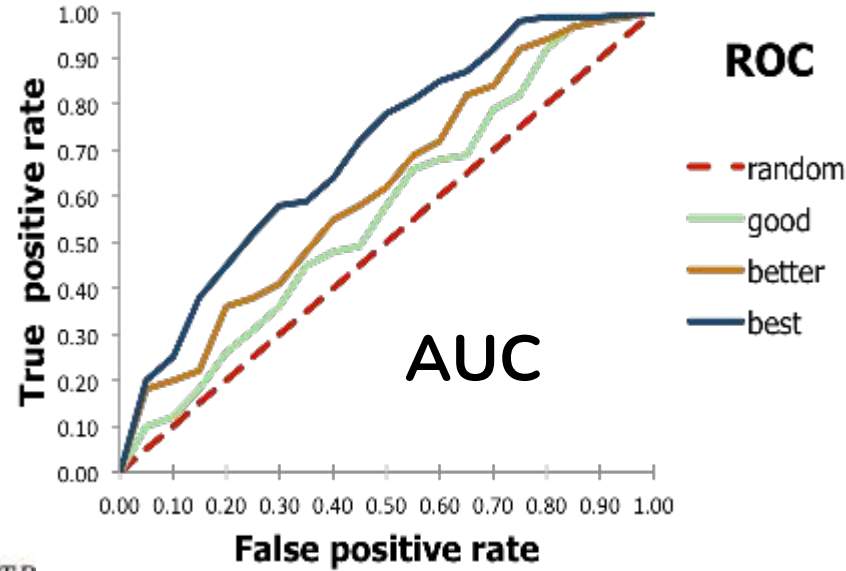(hyper-)parameters to describe data

**Overfitting**
Model has too many
(hyper-)parameters
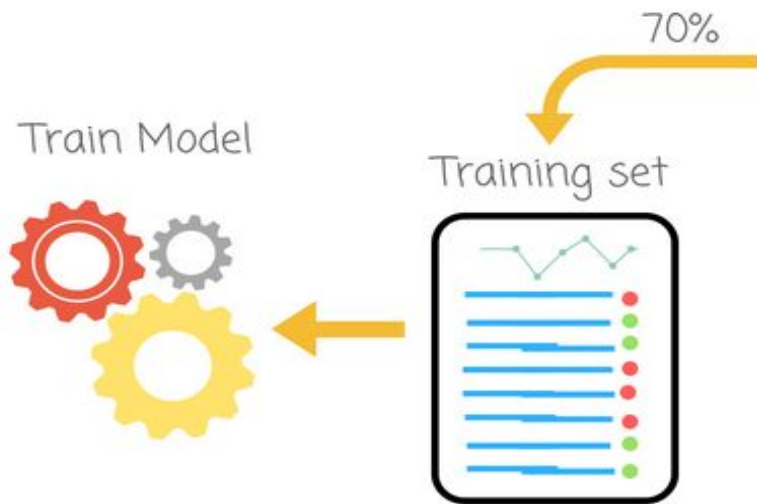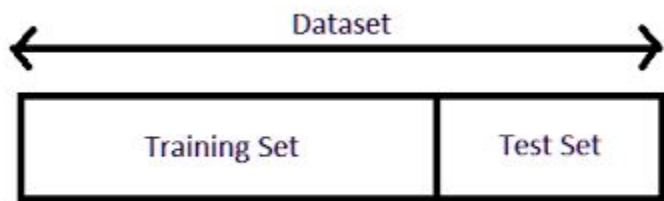
# Classification metrics

- **ROC**: Receiver Operating Characteristics
- **AUC**: Area under the curve
- **TPR**: True positive rate
- **FPR**: False positive rate
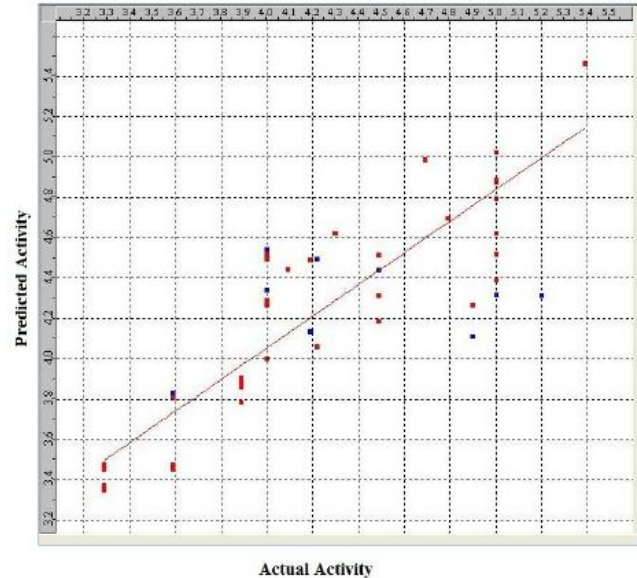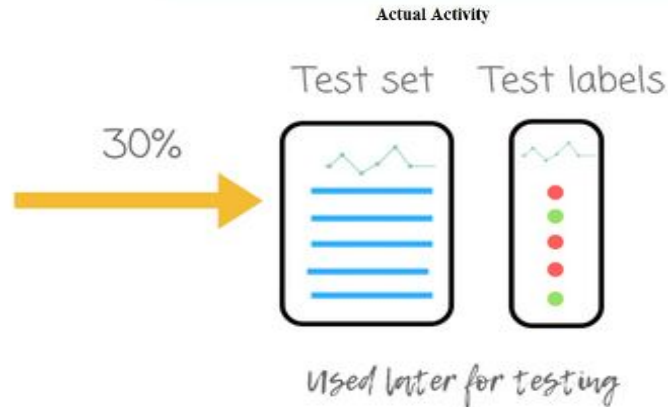- **TNR/FNR**: True/False negative rate

**ROC**



- - - random
— good
— better
— best

**AUC**

**Confusion matrix**

$$\text{fp rate} = \frac{FP}{N} \qquad \text{tp rate} = \frac{TP}{P}$$

$$\text{precision} = \frac{TP}{TP+FP} \qquad \text{recall} = \frac{TP}{P}$$

$$\text{accuracy} = \frac{TP+TN}{P+N}$$

$$\text{F-measure} = \frac{2}{1/\text{precision}+1/\text{recall}}$$

True class

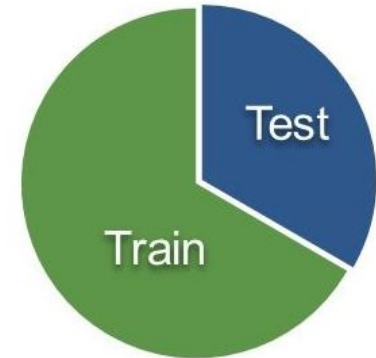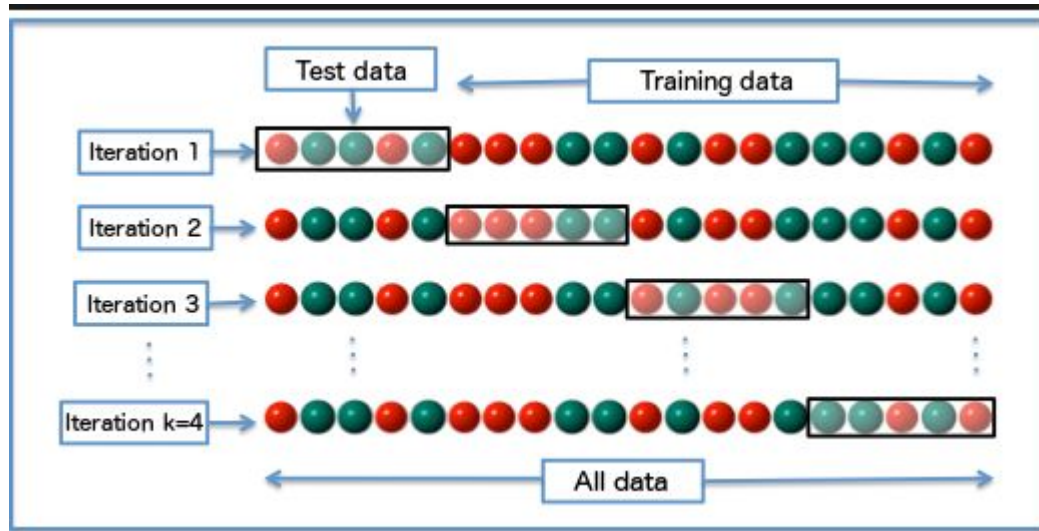|  |  | p | n |
|---|---|---|---|
| Hypothesized class | Y | True Positives | False Positives |
|  | N | False Negatives | True Negatives |
| Column totals: |  | P | N |

# Training and test set

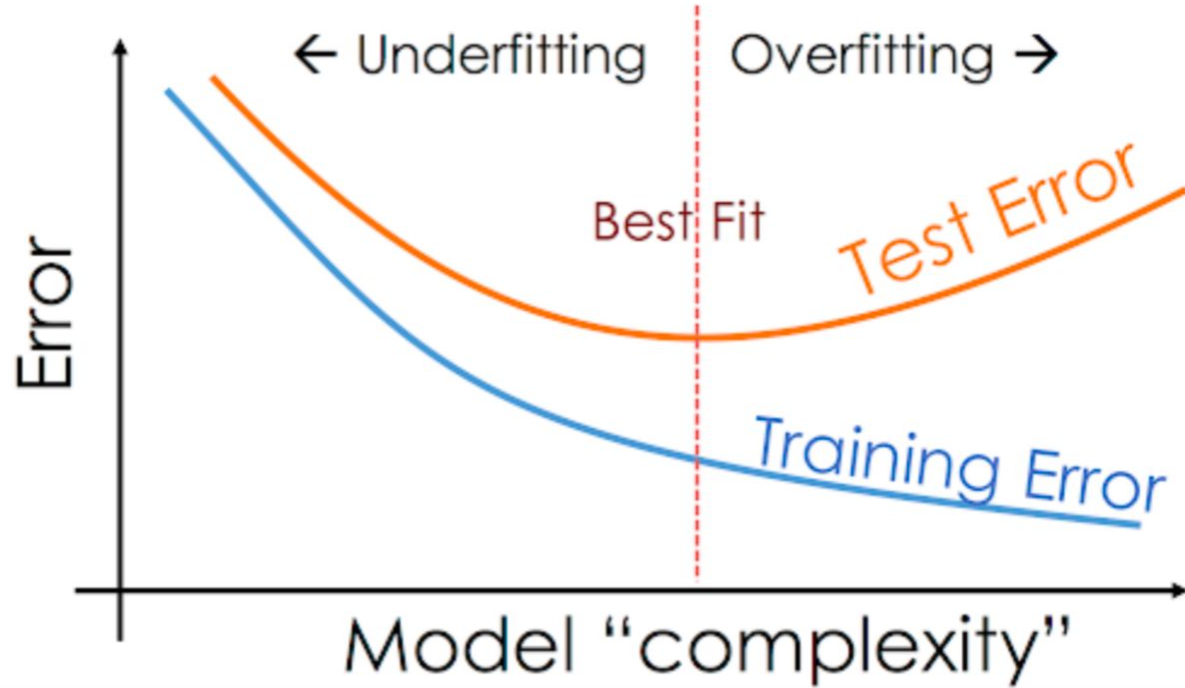# Cross-validation: is you model robust?

- Train/test split
- K-folds cross validation

Loss, or could be any other metrics of interest

Underfitting ← | Overfitting →

Best Fit

Test Error

Training Error

Error

Model "complexity"

Or training epochs (number of iterations)

# Uncertainties

- **Covariate shift/domain shift**: change in the distribution of the input variables between the training and the test datasets
- **How to treat:**
    - **Propagate all errors** in your input variables to the output model
    - Use **cross-validation** and get an estimate of the variance of the model performances when changing hyperparameters
    - Compare model performances when **dropping/adding** features
    - If possible, compare your model results with another architecture
    - Ideally, measure error with independent model-data comparison (such as tag&probe method) just like any other classifier
- [Reference](Reference)