

Big data science Day 3

F. Legger - INFN Torino

<https://github.com/Course-bigDataAndML/MLCourse-2425>

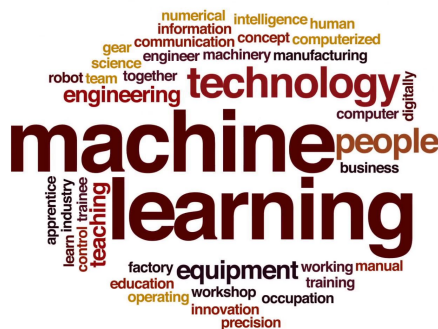
- So long, and thanks for all the images!
 - Taken freely from the web
 - Credits go to original creators



We learned

- Big data, analytics
- Distributed computing
- ML:
 - Supervised vs unsupervised
 - Preprocessing, feature engineering
 - Feature selection

Today



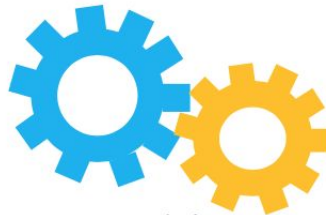
- Machine learning
 - Model architectures
 - Train model and evaluate performances

Remember



Machine learning involves
two mathematical entities

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



Model

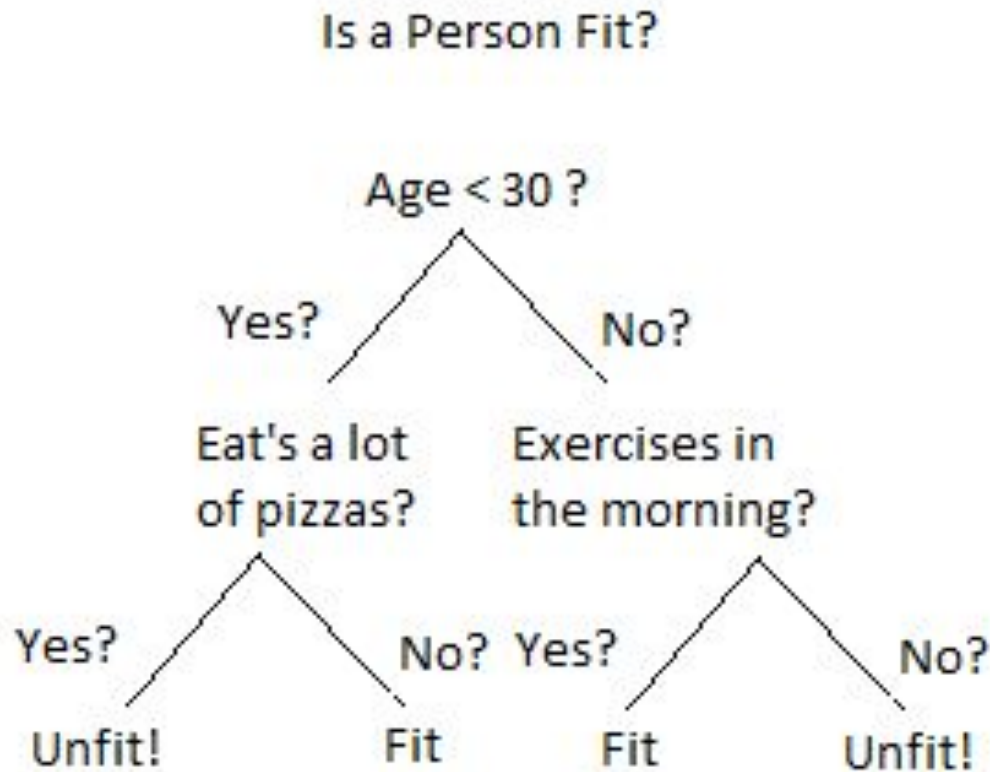


Features: a feature is a
representation of raw data

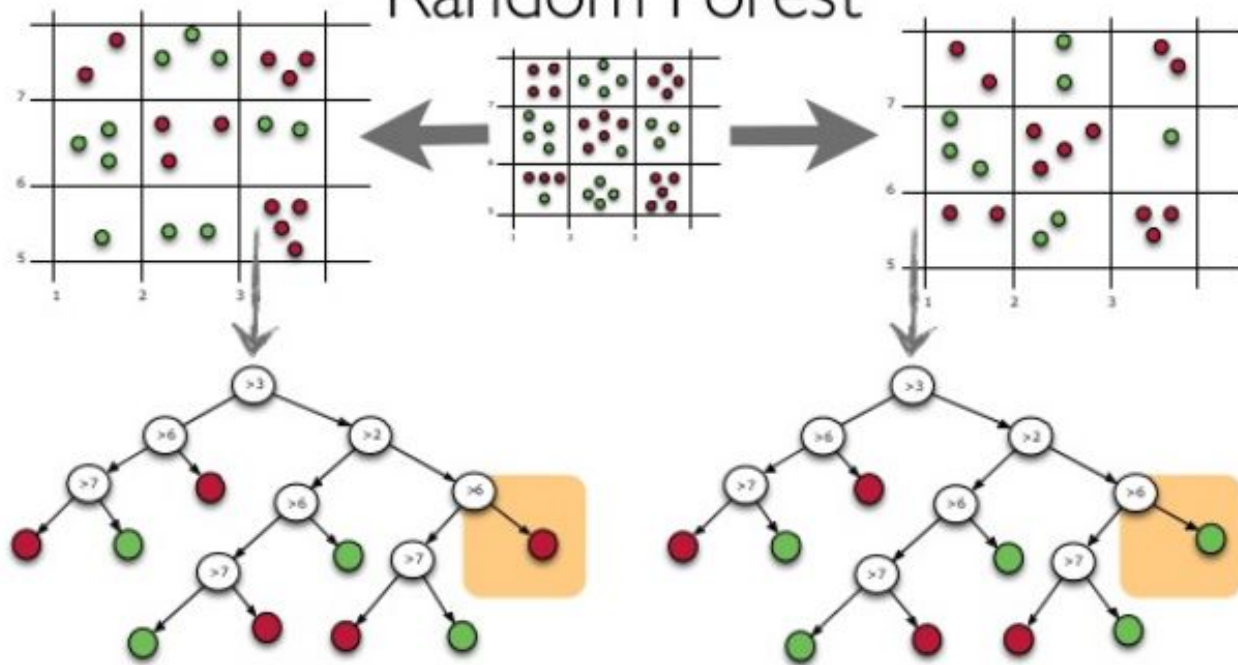
Supervised learning

Decision trees

- Single tree not very powerful
- Typically used in combinations
 - Random forest
 - Gradient Tree Boosting



Random Forest



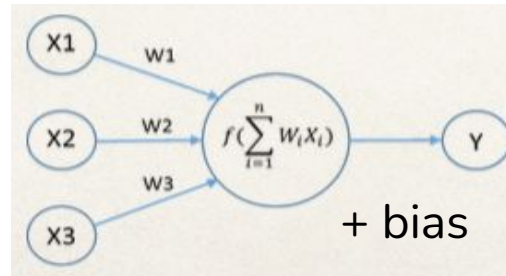
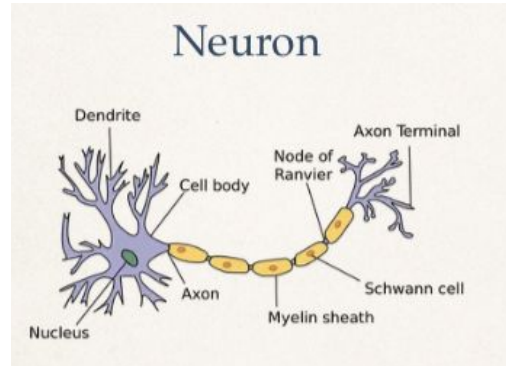
- Each tree sees part of the training sets and captures part of the information it contains

Ensembles

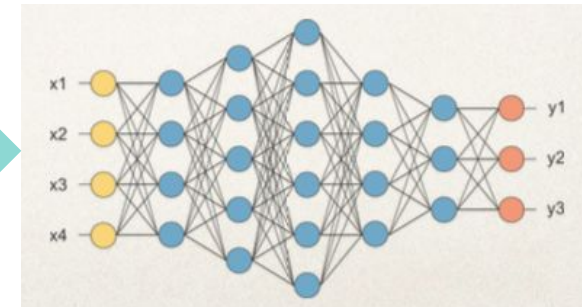
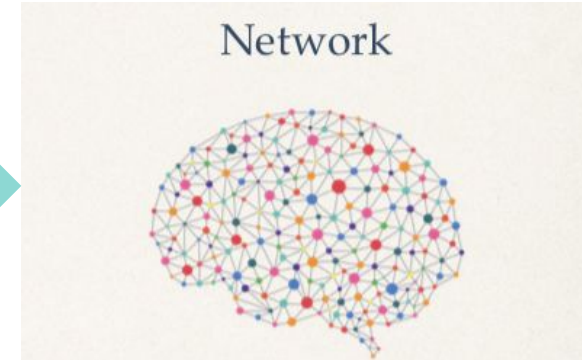
- **Bagging**
 - building multiple models (typically of the *same* type) from different subsamples of the training dataset
- **Boosting**
 - building multiple models (typically of the *same* type) each of which learns to fix the predictions errors of a prior model in the chain
- **Stacking**
 - building 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

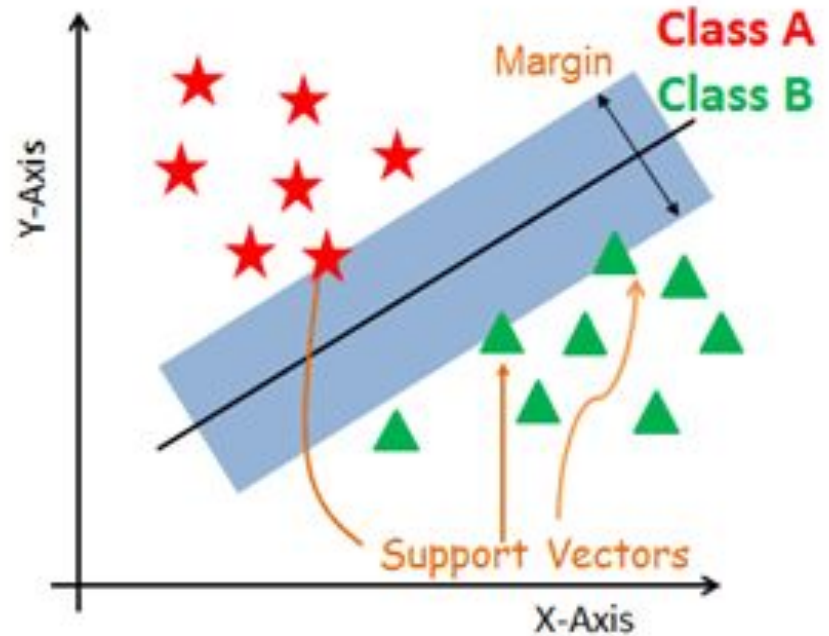
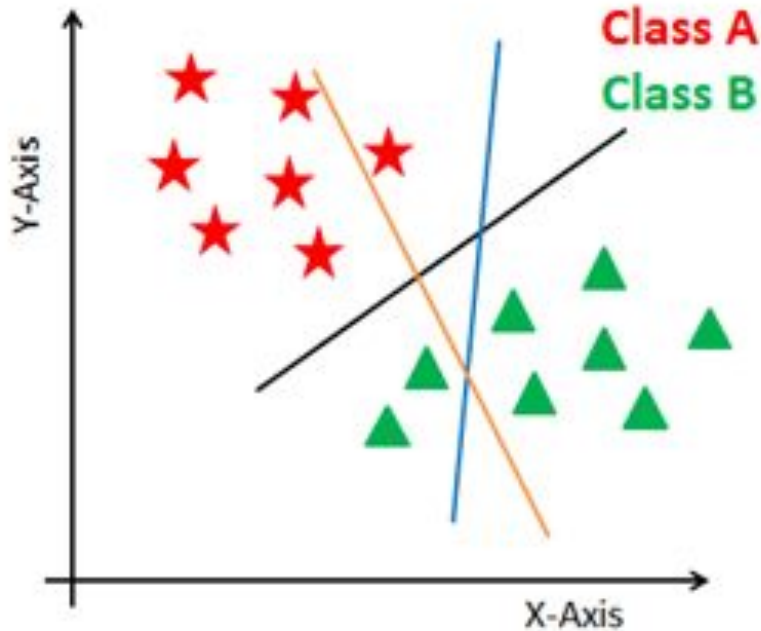
- Basic unit: **Neuron**
- A neuron takes inputs, does some math, and produces **one** output



More on NN tomorrow!



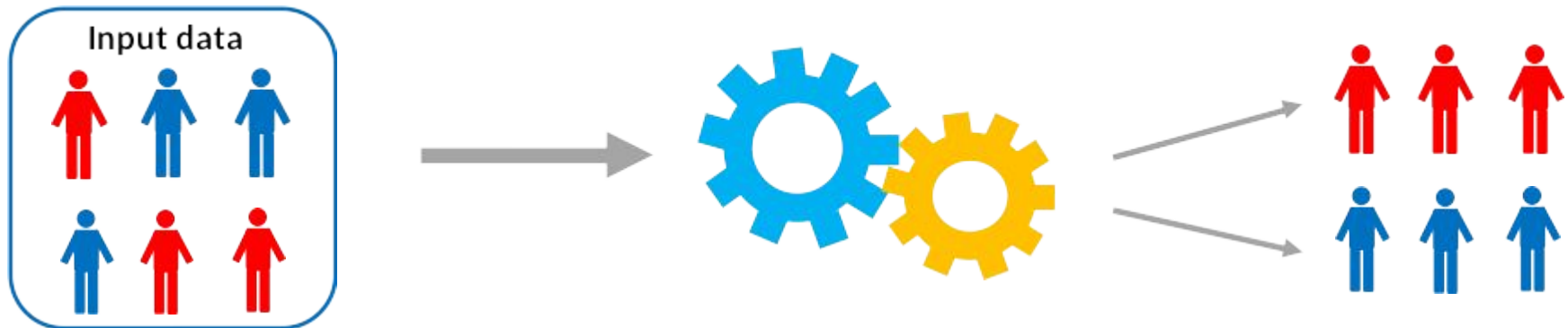
Support vector machines (SVG)



Unsupervised learning

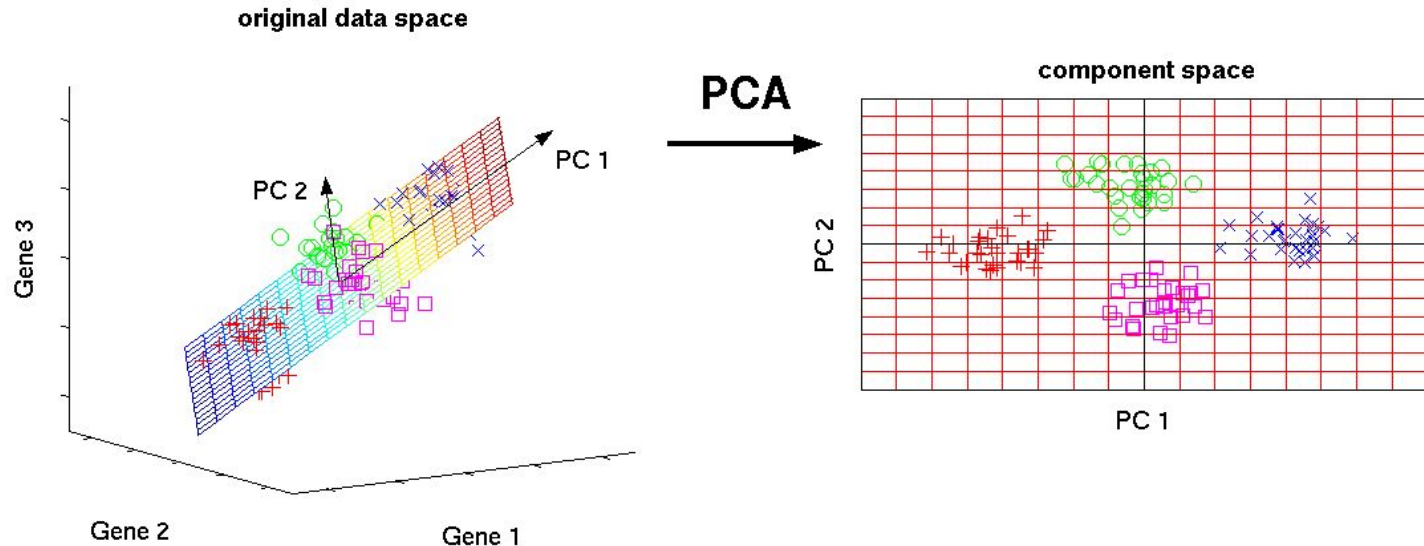
Challenges

- 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

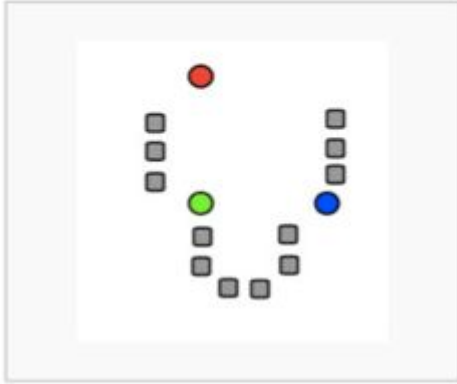


Principal Component Analysis (PCA)

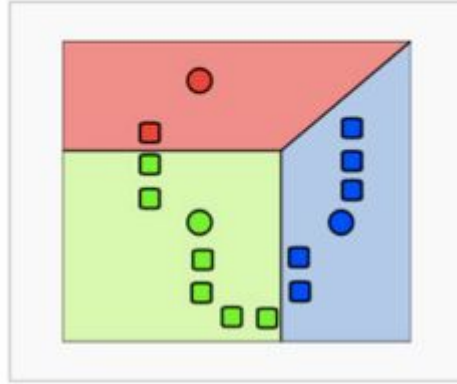
- Used for dimensionality reduction



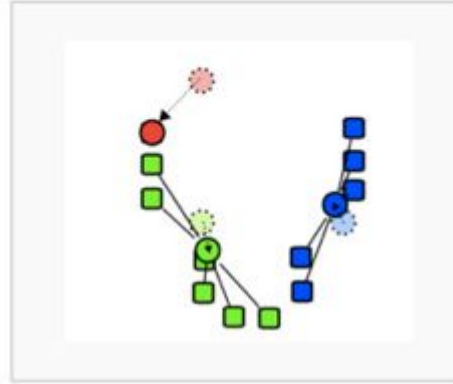
K-means clustering



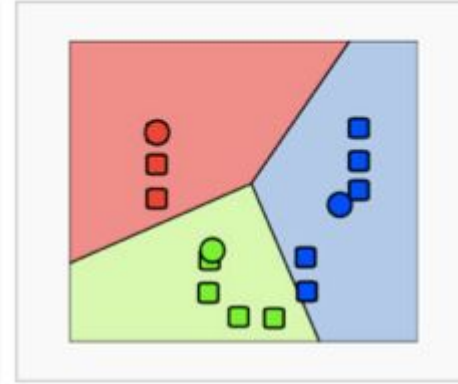
1. k initial "means" (in this case $k=3$) are randomly generated within the data domain (shown in color).



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.

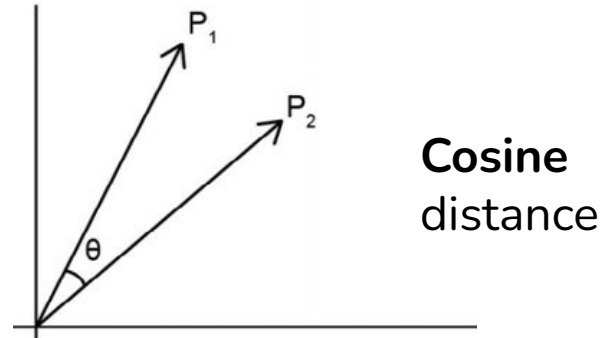
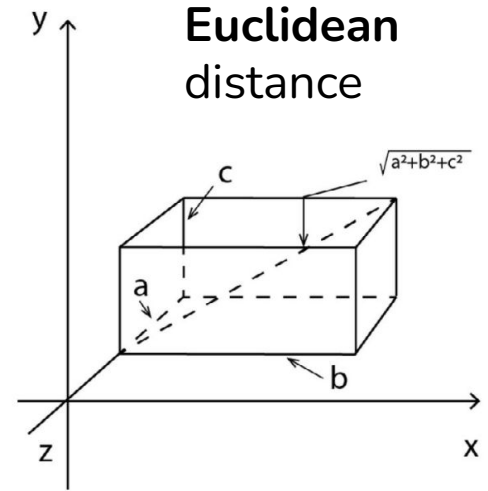
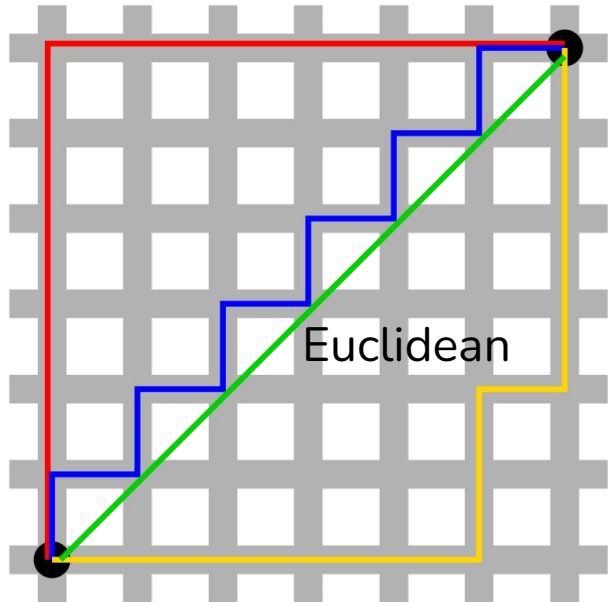


4. Steps 2 and 3 are repeated until convergence has been reached.

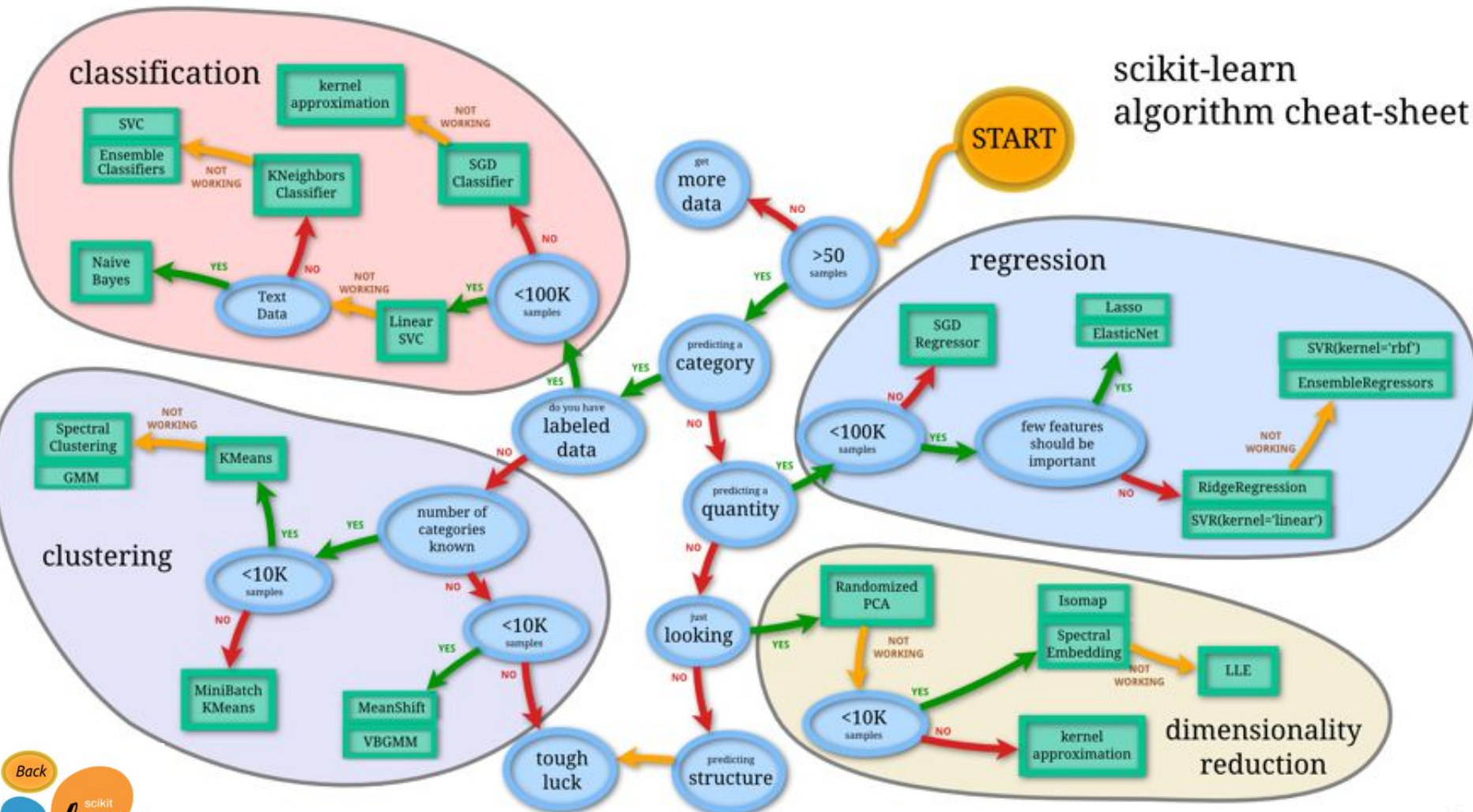
You must define **k**, the number of clusters, and which **distance** to use!







K-means: distances

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



scikit-learn algorithm cheat-sheet



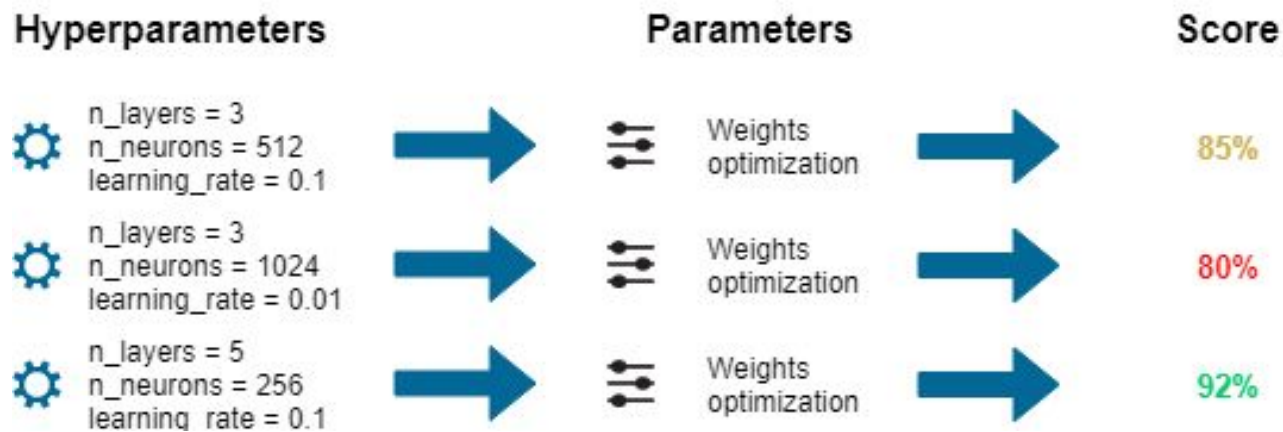
	TYPE	NAME	DESCRIPTION	ADVANTAGES	DISADVANTAGES
Linear		Linear regression	The “best fit” line through all data points. Predictions are numerical.	Easy to understand – you clearly see what the biggest drivers of the model are.	<ul style="list-style-type: none"> ✗ Sometimes too simple to capture complex relationships between variables. ✗ Tendency for the model to “overfit”.
		Logistic regression	The adaptation of linear regression to problems of classification (e.g., yes/no questions, groups, etc.)	Also easy to understand.	<ul style="list-style-type: none"> ✗ Sometimes too simple to capture complex relationships between variables. ✗ 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.	<ul style="list-style-type: none"> ✗ Not often used on its own for prediction because it's also often too simple and not powerful enough for complex data.
		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.	<ul style="list-style-type: none"> ✗ Can be slow to output predictions relative to other algorithms. ✗ Not easy to understand predictions.
		Gradient Boosting	Uses even weaker decision trees, that are increasingly focused on “hard” examples .	High-performing.	<ul style="list-style-type: none"> ✗ A small change in the feature set or training set can create radical changes in the model. ✗ 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.	<ul style="list-style-type: none"> ✗ Very, very slow to train, because they have so many layers. Require a lot of power. ✗ Almost impossible to understand predictions.

*All models are wrong, but
some are useful*

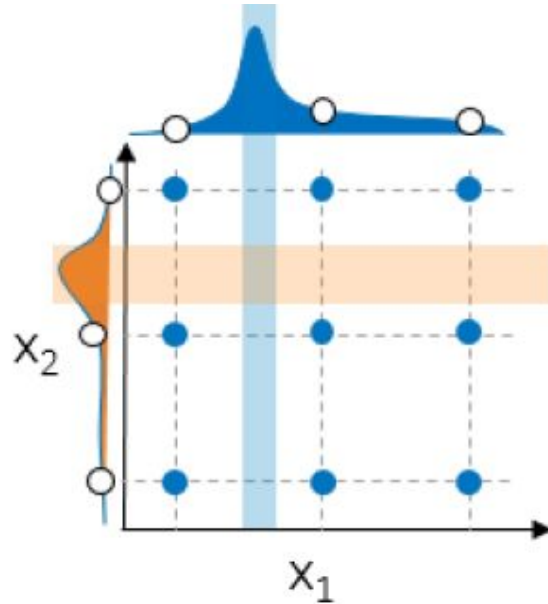
George Box

Hyperparameters vs parameters

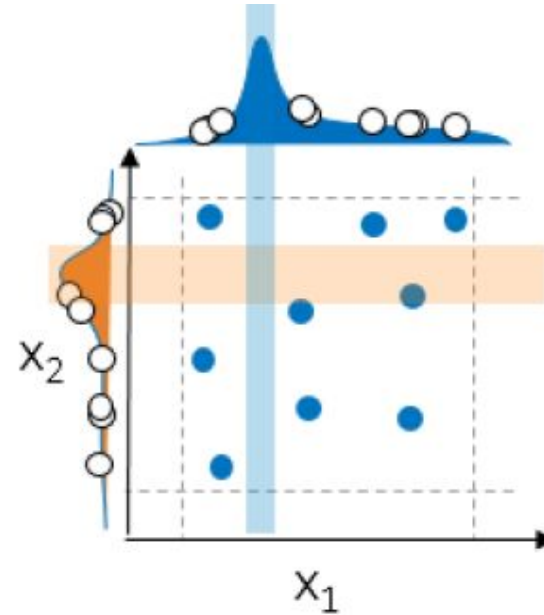
- Model **parameters** are learned during **training** when we optimize the loss function -> *weights*
- **Hyperparameters** are not model parameters and they **cannot** be directly trained from the data -> *model architecture architecture*



Hyperparameter tuning



(a) Standard Grid Search



(b) Random Search

Modeling Algorithm

Tune

hyperparameters (tuning options)

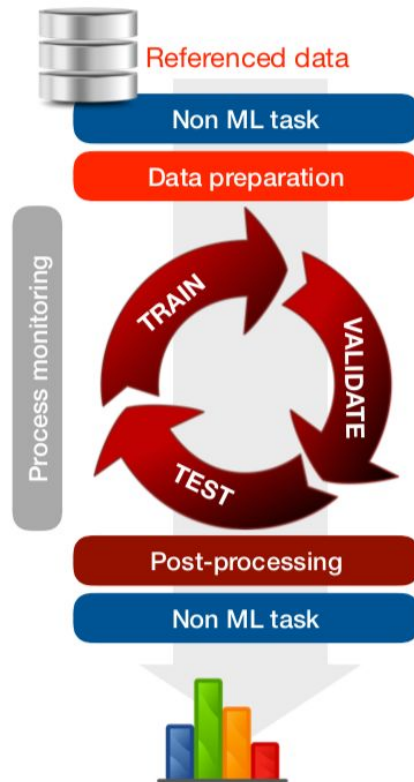
- Polynomial order, penalty parameter, ...
- Network configuration, solver options, ...
- Max tree depth, splitting criterion, ...

Model

Train

model parameters

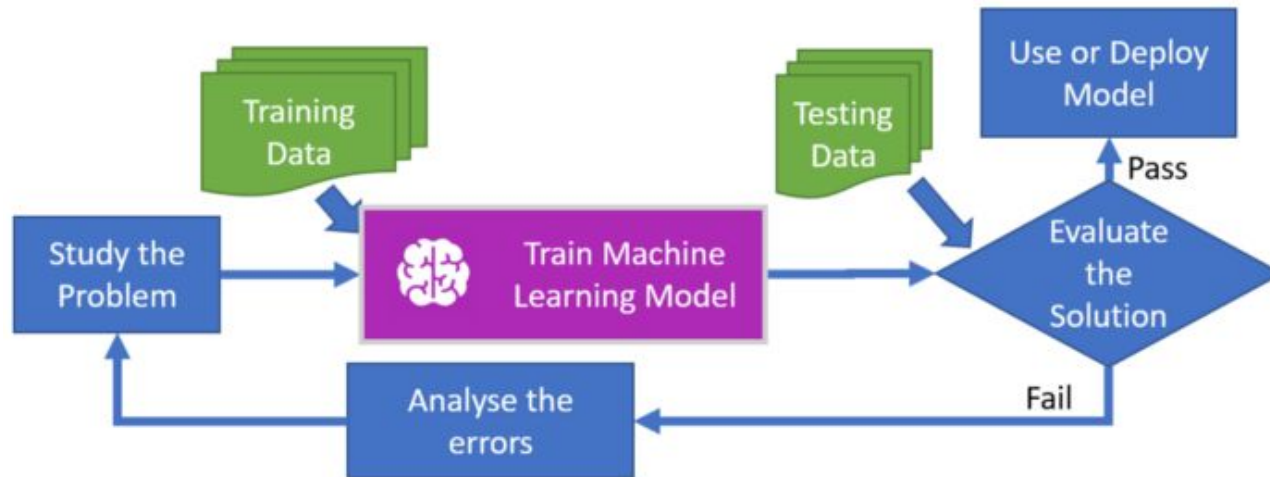
- Regression coefficients
- Neural net weights
- Tree splitting rules
- ...



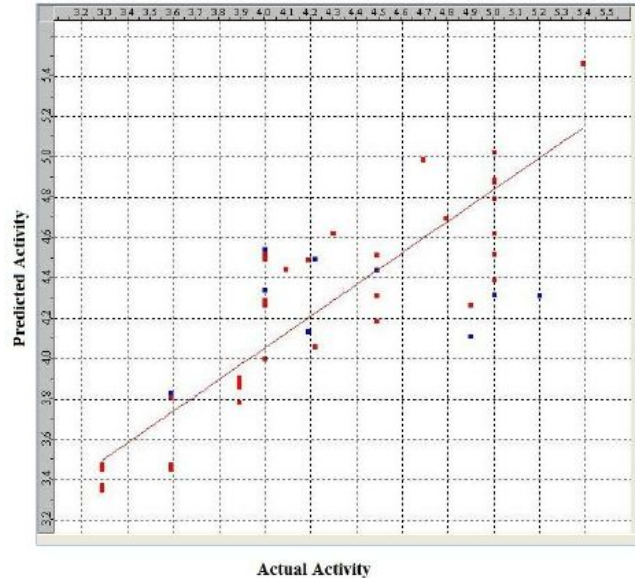
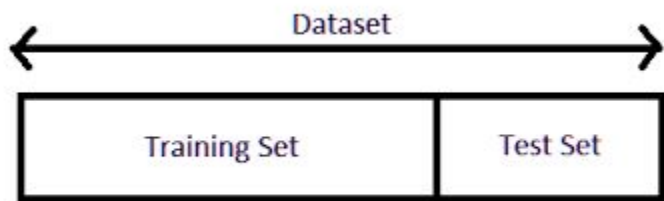
The machine learning pipeline

- Input data is typically split into:
 - Training dataset
 - Testing dataset

Remember the ML pipeline from Day 1!

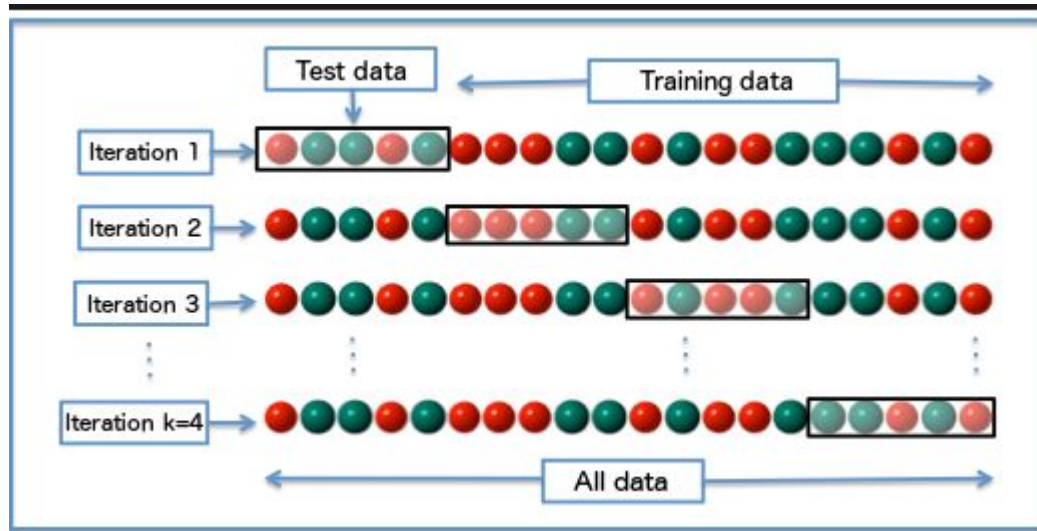


Training and test set



Cross-validation

- Train/test split
- K-folds cross validation



Classification metrics

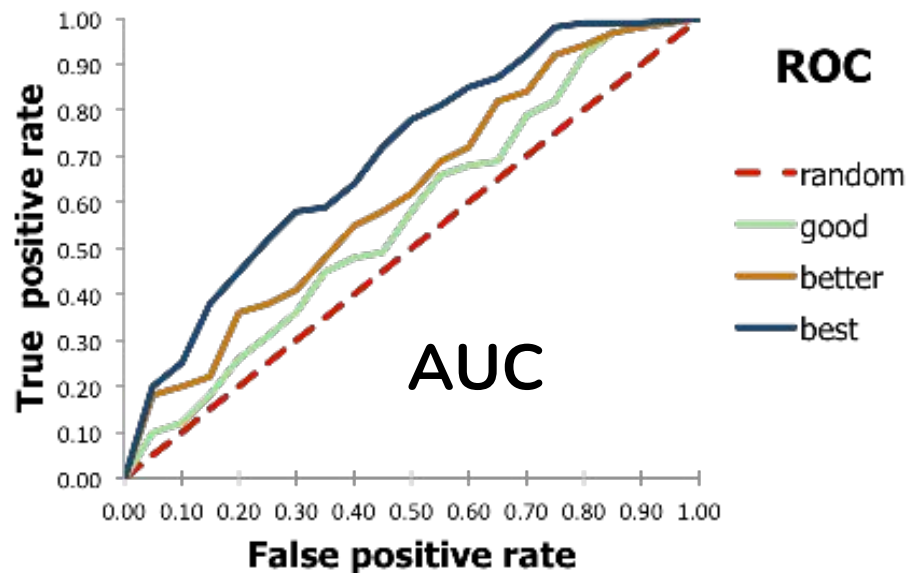
- TPR: True positive rate
- FPR: False positive rate
- TNR: True negative rate
- FNR: False negative rate

		<u>True class</u>			
		p	n		
<u>Hypothesized class</u>	Y	True Positives	False Positives	$fp\ rate = \frac{FP}{N}$	$tp\ rate = \frac{TP}{P}$
	N	False Negatives	True Negatives	$precision = \frac{TP}{TP+FP}$	$recall = \frac{TP}{P}$
Column totals:		P	N	$accuracy = \frac{TP+TN}{P+N}$	$F\text{-measure} = \frac{2}{1/precision + 1/recall}$

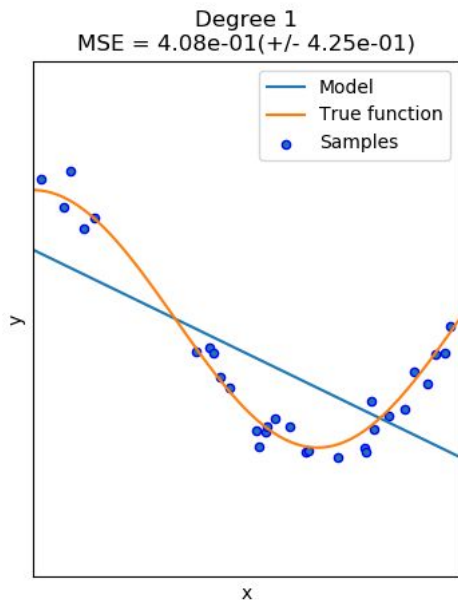
Confusion
matrix

Classification metrics

- **ROC**: Receiver Operating Characteristics
- **AUC**: Area under the curve

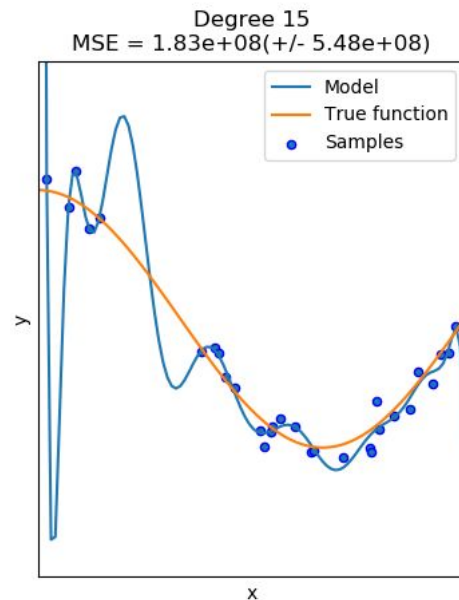
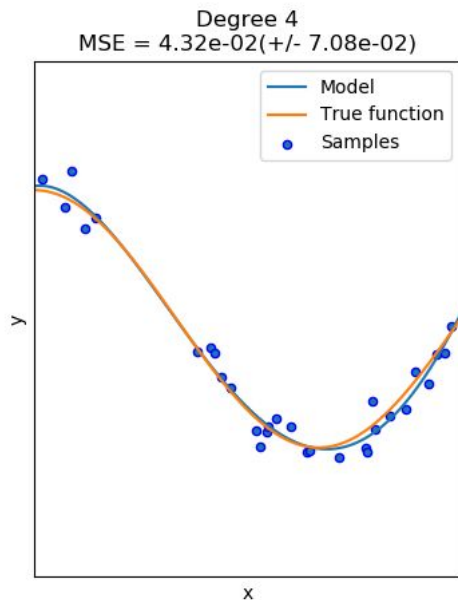


Overfitting / underfitting



Underfitting

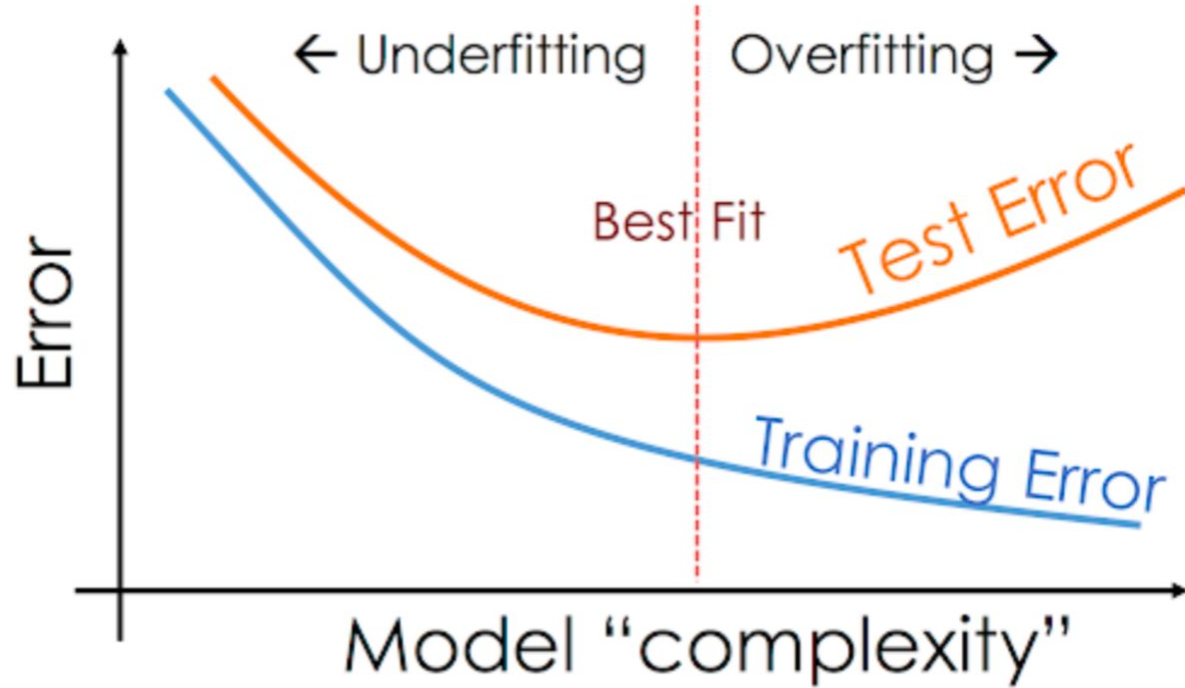
Model doesn't have enough (hyper-)parameters to describe data



Overfitting

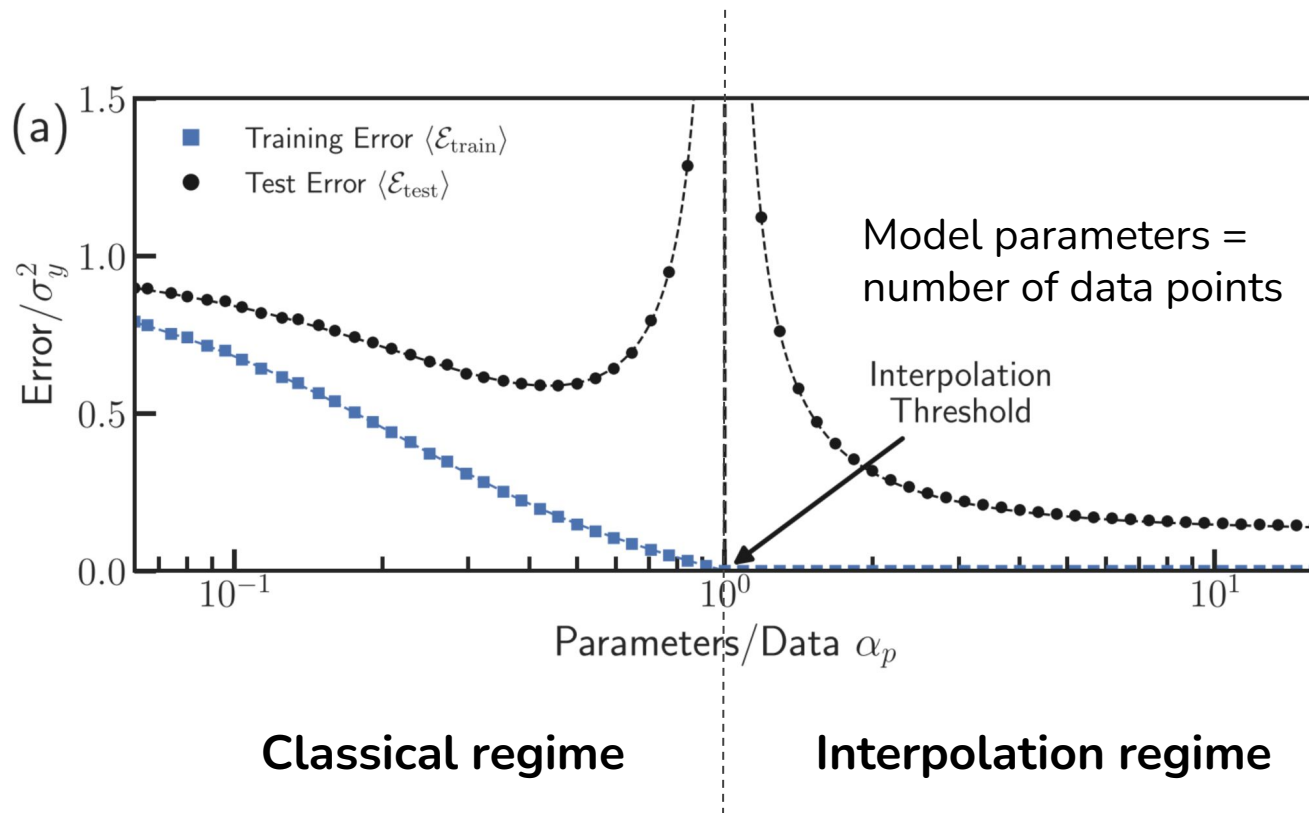
Model has too many (hyper-)parameters

Loss, or could
be any other
metrics of
interest

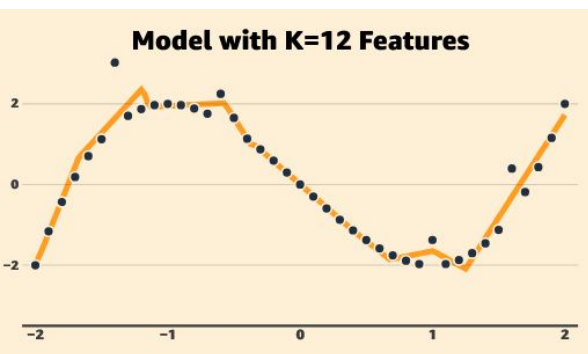
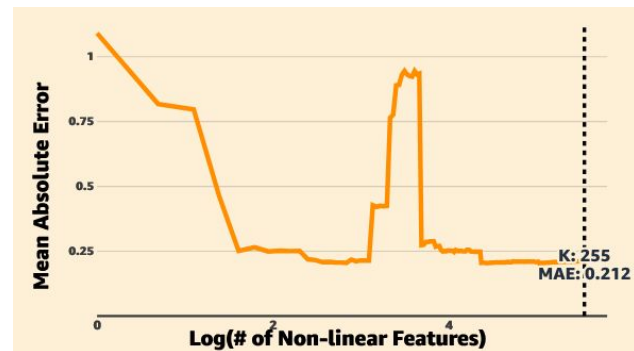


Or training epochs (number of iterations)

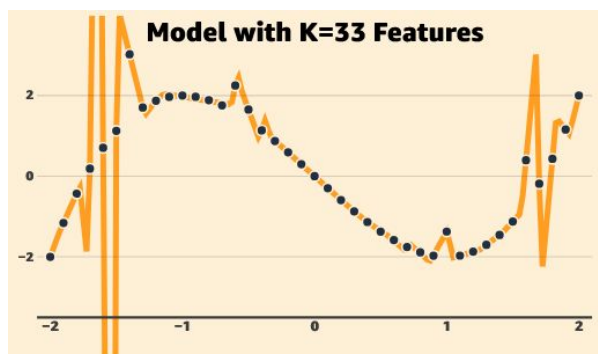
Double descent



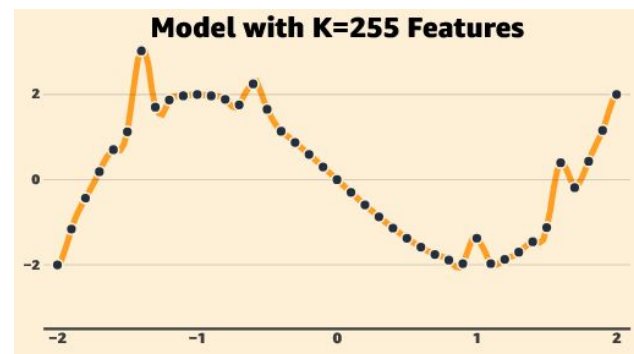
Double descent: explained?



Classical regime



Interpolation threshold



Interpolation regime

Uncertainties

- **Covariate shift/domain shift:** change in the distribution of the input variables between the training and the test datasets
- **How to mitigate:**
 - **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](#)