



Machine Learning I & II

Hi! PARIS DataBootcamp 2024



Agenda of the course

Machine Learning I (09:00-10:15)

- Introduction to Machine Learning
- Data Preprocessing (+ demo)
- Model Training
- Linear Regression + Continuous Evaluation

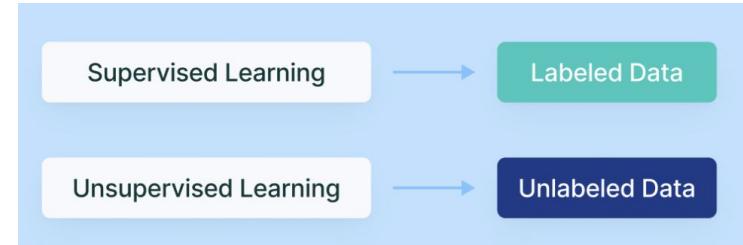
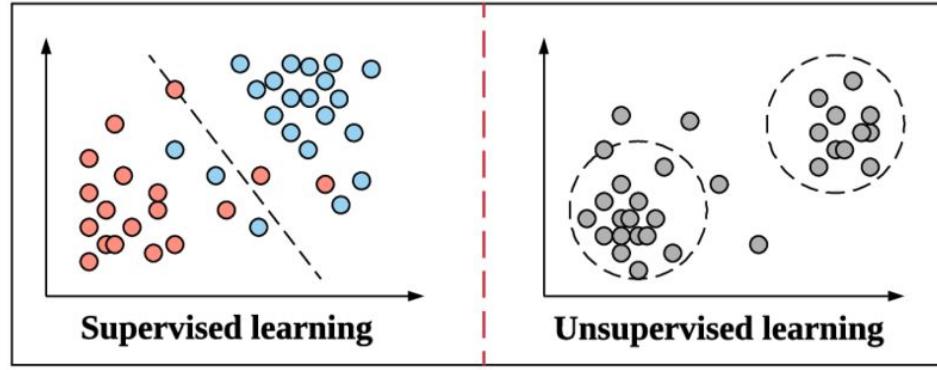
Machine Learning II (10:45-12:00PM)

- Logistic Regression + Categorical Evaluation (+ demo)
- KNN (+ demo)
- Decision Tree (+ demo)
- Improve the performance of your model (+ demo)

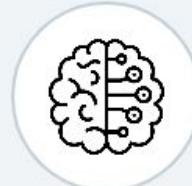




Introduction to Machine Learning



Machine Learning



Supervise Learning



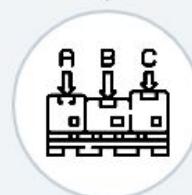
Unsupervised Learning



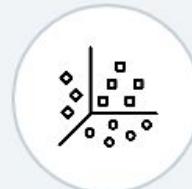
Regression



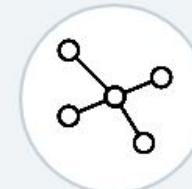
Classification

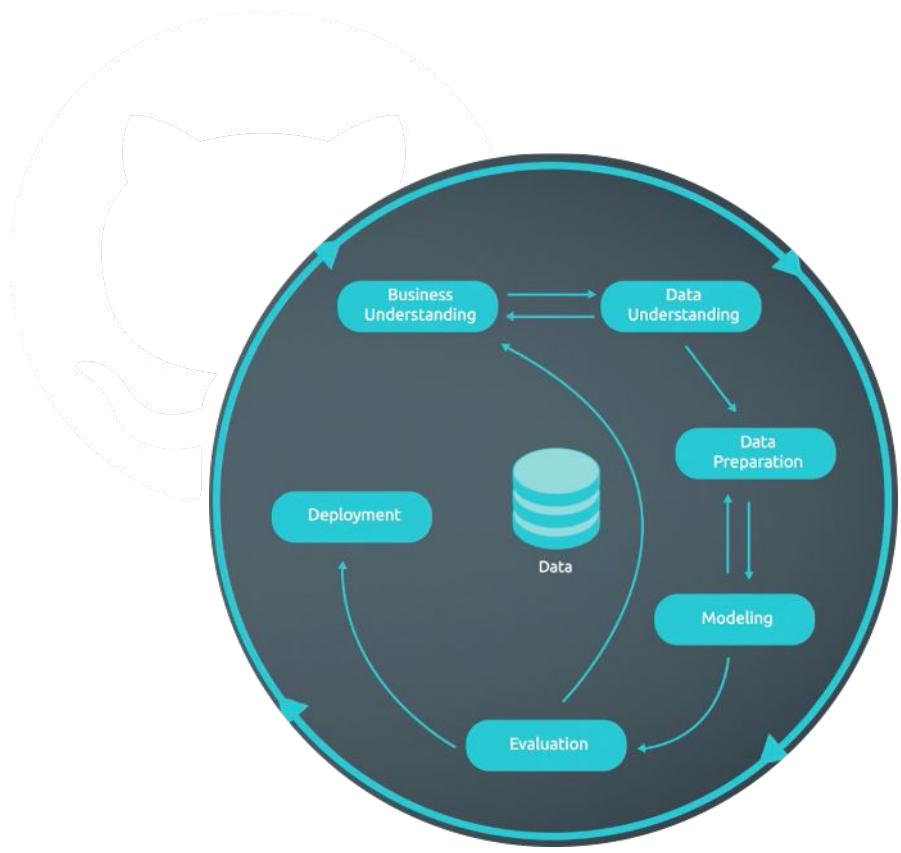


Clustering



Association





Feature vs Target variables



Feature variables: Attributes of the data used to make predictions.

- *Shouldn't include the desired output (target).*

Feature vs Target variables



Feature variables: Attributes of the data used to make predictions.

- *Shouldn't include the desired output (target).*



Target variable: Output variable that the model aims to predict.

- *Data is considered “labeled” when the target is available.*
- *Data is “unlabeled” when it isn’t.*



Let's give an example

Build a Machine Learning model that can **predict whether a passenger survived (or not) the titanic**

The dataset is
structured



Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare
0	3	Braund, Mr. Owen Harris	male	22.0	1	0	A/5 21171	7.2500
1	1	Cumings, Mrs. John Bradley (Florence Briggs Th... ...	female	38.0	1	0	PC 17599	71.2833
1	3	Heikkinen, Miss. Laina	female	26.0	0	0	STON/O2. 3101282	7.9250
1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	1	0	113803	53.1000

Let's give an example



Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare
0	3	Braund, Mr. Owen Harris	male	22.0	1	0	A/5 21171	7.2500
1	1	Cumings, Mrs. John Bradley (Florence Briggs Th...)	female	38.0	1	0	PC 17599	71.2833
1	3	Heikkinen, Miss. Laina	female	26.0	0	0	STON/O2. 3101282	7.9250
1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	1	0	113803	53.1000

Target

Feature variables

The target variable is available Labeled dataset

Let's give an example



Target

Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare
0	3	Braund, Mr. Owen Harris	male	22.0	1	0	A/5 21171	7.2500
1	1	Cumings, Mrs. John Bradley (Florence Briggs Th... ...	female	38.0	1	0	PC 17599	71.2833
1	3	Heikkinen, Miss. Laina	female	26.0	0	0	STON/O2. 3101282	7.9250
1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	1	0	113803	53.1000

Training data



Model training



Machine Learning
model

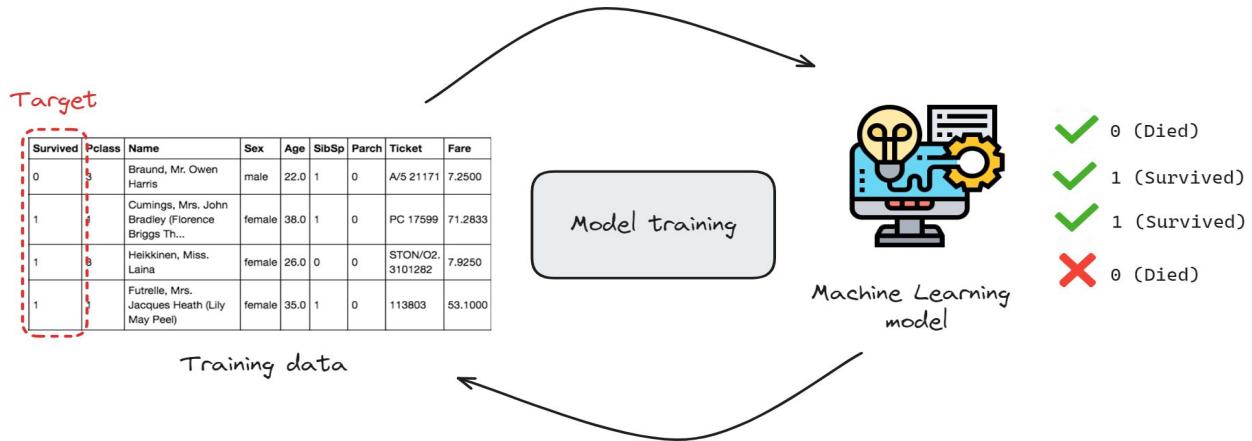
- ✓ 0 (Died)
- ✓ 1 (Survived)
- ✓ 1 (Survived)
- ✗ 0 (Died)



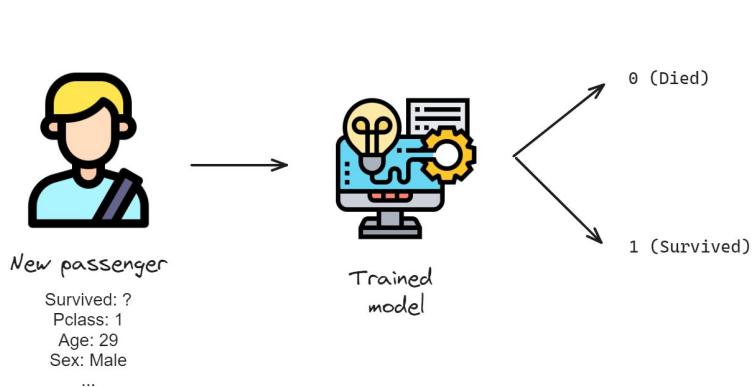
Let's give an example



Step 1: Training



Step 2: Final prediction



Continuous vs Categorical data



Continuous data: Numerical data that can take any value within a defined range

- Examples: *Weight, price, temperature*



Categorical data: Data with a finite number of possible values/categories

- Categories are ordinal or nominal
- Examples: Yes/No, 0/1,....

Continuous vs Categorical data



	Gender	Height	Weight	Index	Status
0	Male	174	96	4	Obesity
1	Male	189	87	2	Normal
2	Female	185	110	4	Obesity
3	Female	195	104	3	Overweight
4	Male	149	61	3	Overweight
5	Male	189	104	3	Overweight
6	Male	147	92	5	Extreme Obesity
7	Male	154	111	5	Extreme Obesity
8	Male	174	90	3	Overweight
9	Female	169	103	4	Obesity

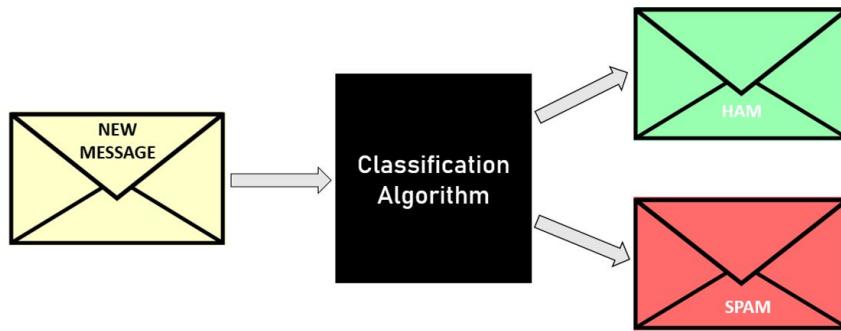
Height and Weight are continuous

Gender, Index and Status are categorical

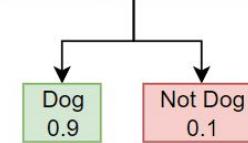


Supervised Learning

1. Classification: Task of predicting a categorical target



Spam classification (Spam, ham)

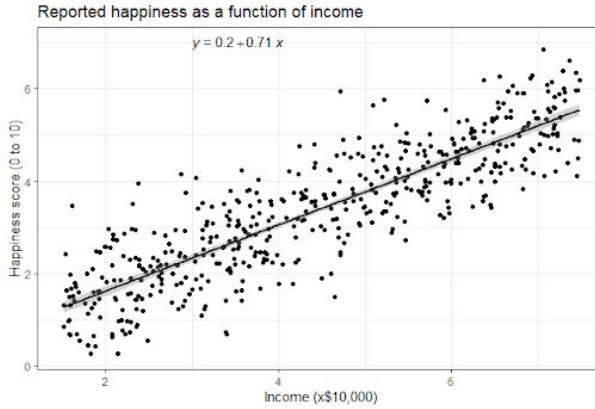


Object classification (Dog, not dog)

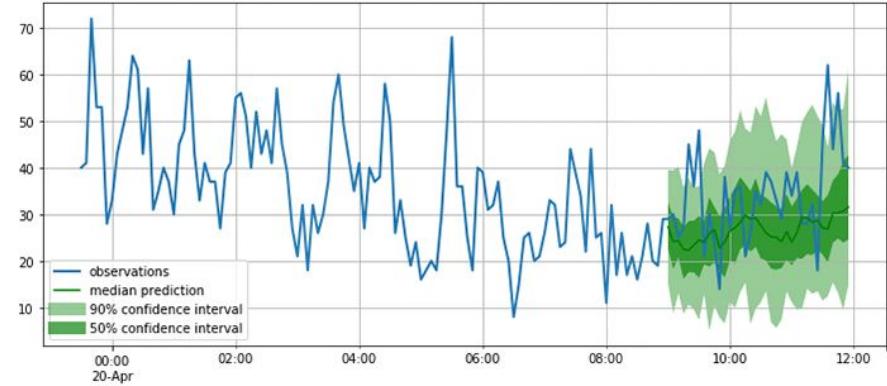


Supervised Learning

2. Regression: Task of predicting a **continuous target**



*Income prediction
(based on a happiness score)*

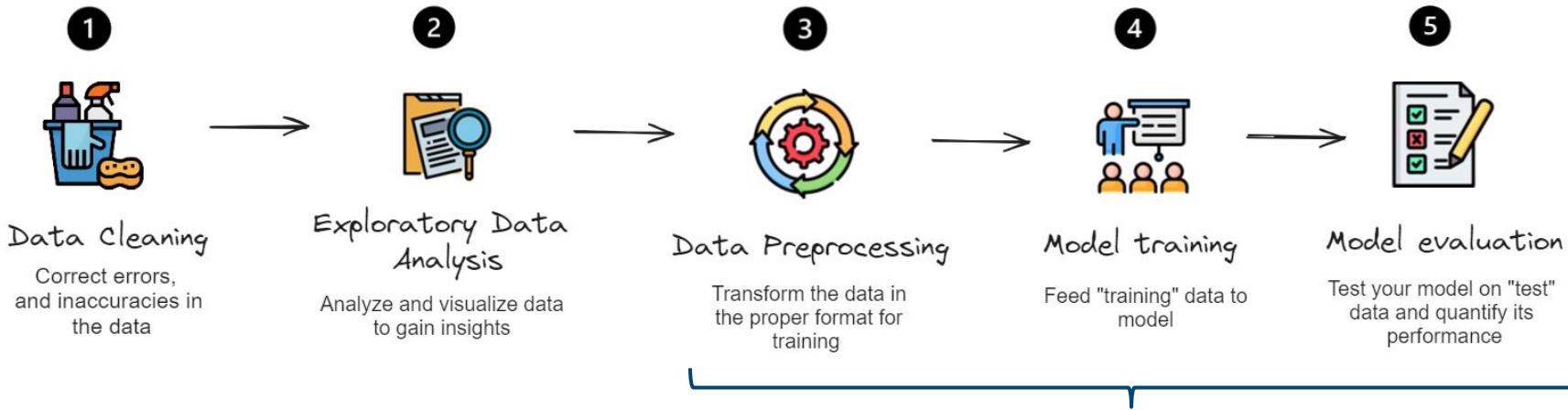


*Weather forecasting
(based on historical data)*

Machine Learning in practice



In practice, Machine Learning is more than building algorithms. Models need **clean and properly formated data** to generate good results.





Data Preprocessing

What is Data Preprocessing



Data Preprocessing is the process of transforming data it into a clean and suitable format for training.

Examples

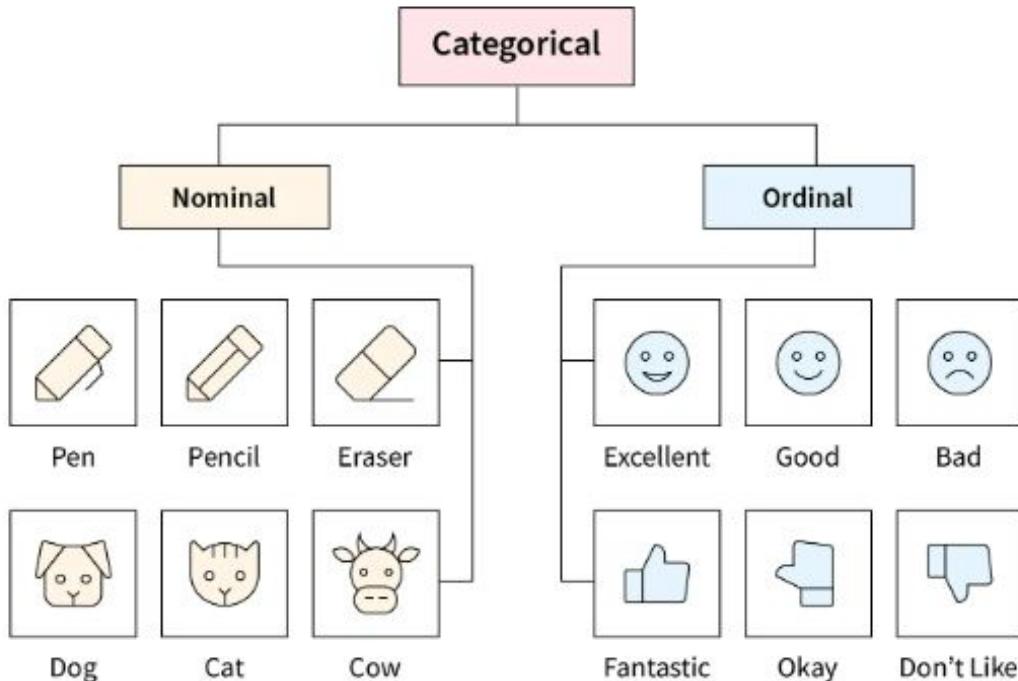
- **Categorical Encoding** (for categorical variables)
- **Feature Scaling** (for continuous variables)
- **Feature Selection** (not covered today)

Ordinal vs Nominal variables



- **Ordinal:** Categorical variable whose categories have an inherent order/hierarchy
- **Nominal (unordered):** Categorical variable whose categories have no hierarchy

Ordinal vs Nominal variables



Categorical encoding



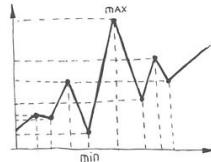
Categorical encoding means **converting categorical variables into numerical values** that can be understood and processed by a model.

Linear Algebra

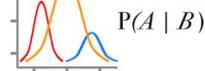
$$d = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

[matrix]

Graphs



Probability

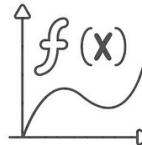


Statistics



Machine Learning

Calculus



Math behind Machine Learning models

Categorical encoding



1 Label Encoding

Replace categories by numerical values (0,1,2,...)

Each class gets attributed a numerical value within the same column

Grades
A
B
C
D
Fail



Grades	Encoded
A	4
B	3
C	2
D	1
Fail	0

Grades is an ordinal categorical variables



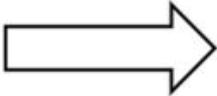
Categorical encoding

2 One-Hot Encoding

Create a new binary variable (0 or 1) for each class

Each class
has its own column with 0 or 1

Places
New York
Boston
Chicago
California
New Jersey



	New York	Boston	Chicago	California	New Jersey
New York	1	0	0	0	0
Boston	0	1	0	0	0
Chicago	0	0	1	0	0
California	0	0	0	1	0
New Jersey	0	0	0	0	1

Places is a nominal categorical variables

The variable is 1 if the place
is New York, else 0

Categorical encoding



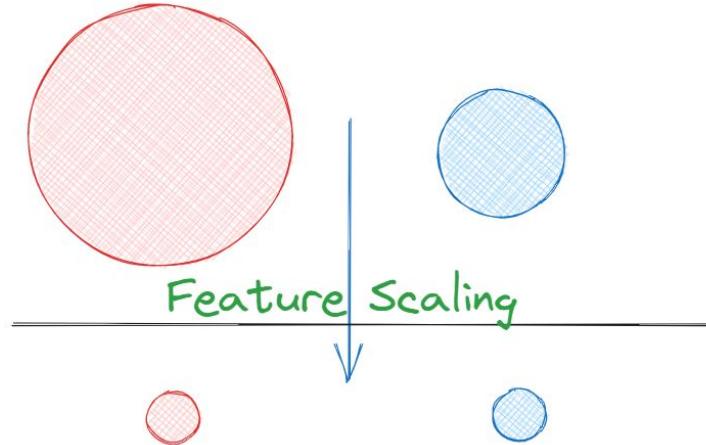
Important notions:

- Label Encoding should only be used on ordinal feature variables
- One-Hot Encoding can lead to a huge increase in the number of columns when the number of classes is large.

Feature scaling



Feature scaling means **scaling numerical variables to a similar range of values**



Why is it useful ?

- Improve a model's performance
- Decrease training time
- Not useful for every model
(usually distance-based ones)



Feature scaling

1 Normalization (MinMax Scaling)

Scale values to a specific range, usually between 0 and 1

Age	Normalized Age
44	0.80952381
27	0
30	0.142857143
38	0.523809524
40	0.619047619

Salary	Normalized Salary
73000	0.838709677
47000	0
53000	0.193548387
62000	0.483870968
57000	0.322580645

$$x_{new} = \frac{x - x_{min}}{x_{max} - x_{min}}$$

The normalized values of **Age** and **Salary** are between 0 and 1



Feature scaling

2 Standardisation (Standard Scaler)

Scale a variable to a normal distribution, with a 0 mean and 1 standard deviation

input	standardized
0	-1.46385
1	-0.87831
2	-0.29277
3	0.29277
4	0.87831
5	1.46385

$$x_{new} = \frac{x - \mu}{\sigma}$$

The standardized input follows a normal distribution

Normalization vs Standardisation



1. Normalization

- Better for variables without a normal distribution
- More sensitive to outliers

2. Standardisation

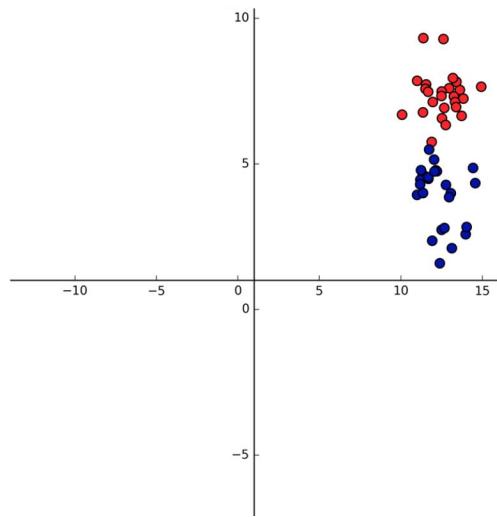
- Better for variables that have close to a normal distribution
- Less sensitive to outliers

The best way to choose is to try both and compare results

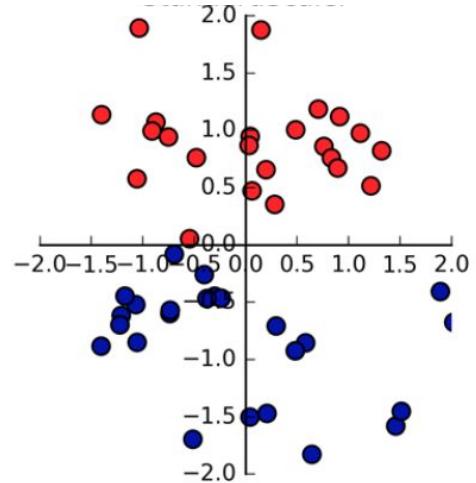
Normalization vs Standardisation



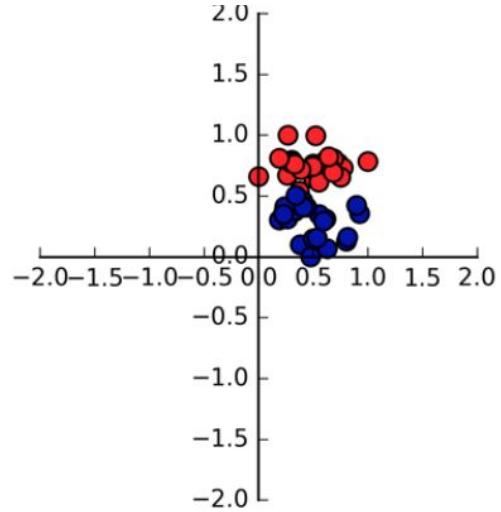
No scaling



Standardisation



Normalization



Data Preprocessing with Scikit-learn



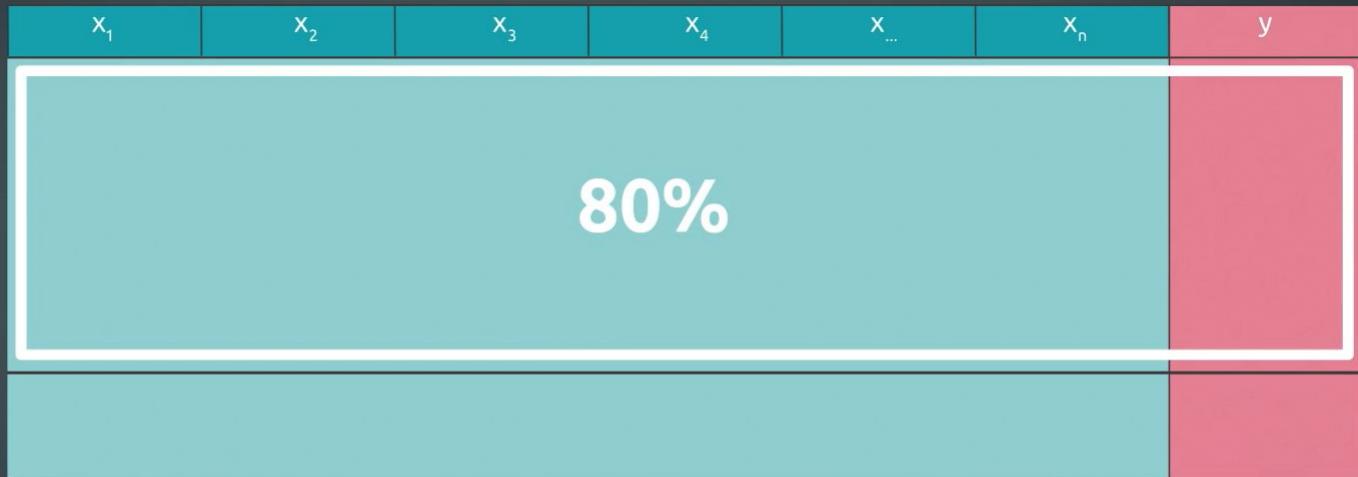
Demo of how to use scikit-learn for data preprocessing

<https://scikit-learn.org/stable/modules/preprocessing.html>



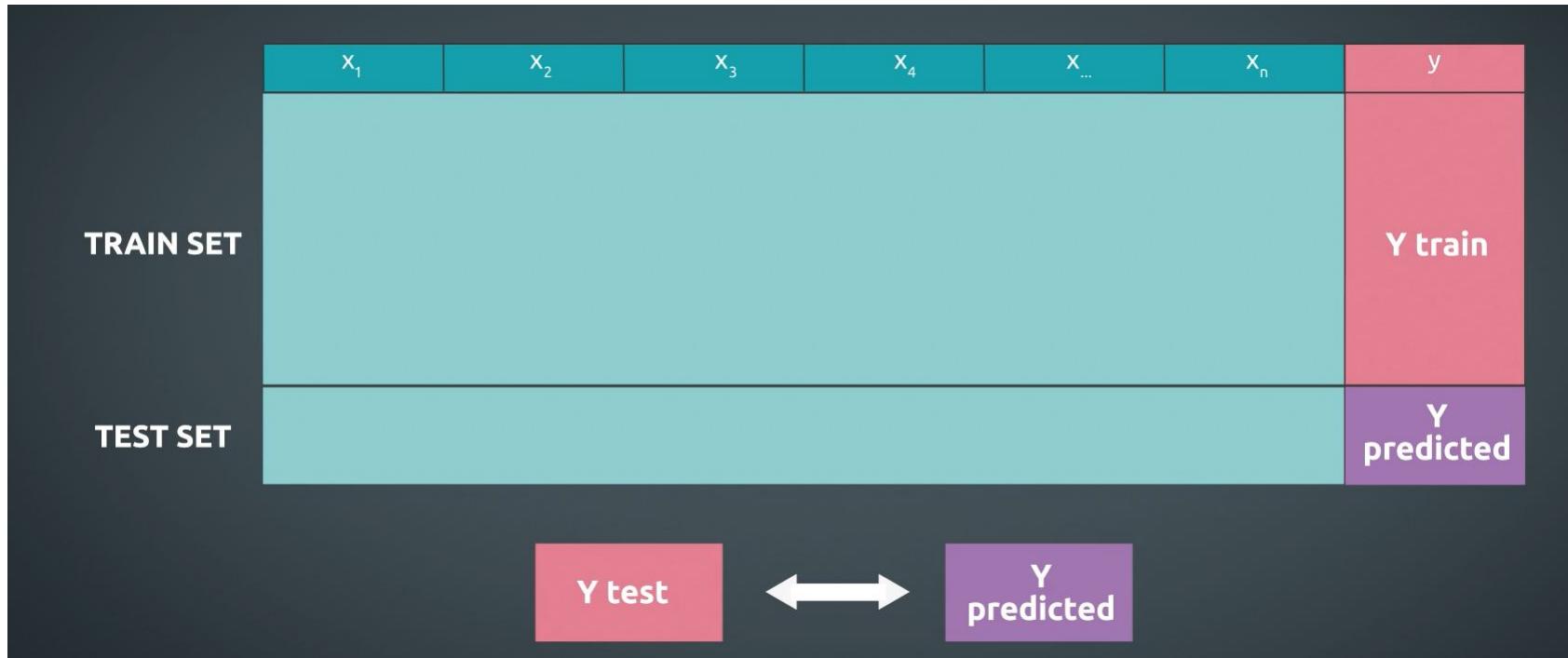
Model training

TRAIN SET

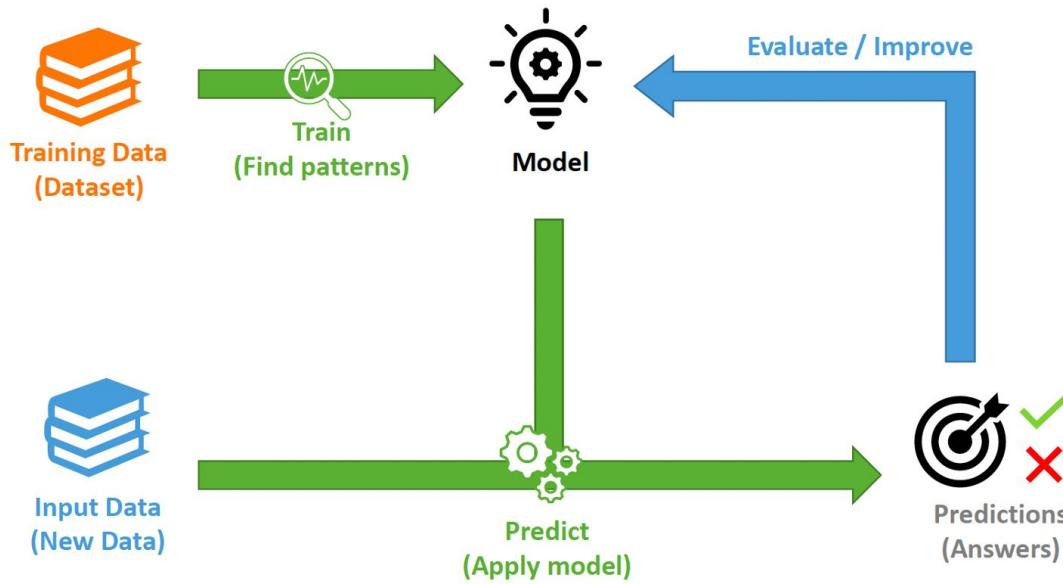


	x_1	x_2	x_3	x_4	x_{\dots}	x_n	y
TRAIN SET							
TEST SET			20%				

	x_1	x_2	x_3	x_4	x_{\dots}	x_n	y
TRAIN SET				X train			Y train
TEST SET				X test			Y test



Machine Learning Process



Supervised Learning models



Most supervised learning models have **variants for Regression and Classification** (but not all).

Models	Classification	Regression
Linear Regression	No	Yes
Logistic Regression	Yes	No
K nearest neighbors	Yes	Yes
Decision Trees	Yes	Yes
Random Forest	Yes	Yes
Boosting	Yes	Yes



Linear Regression

Linear Regression

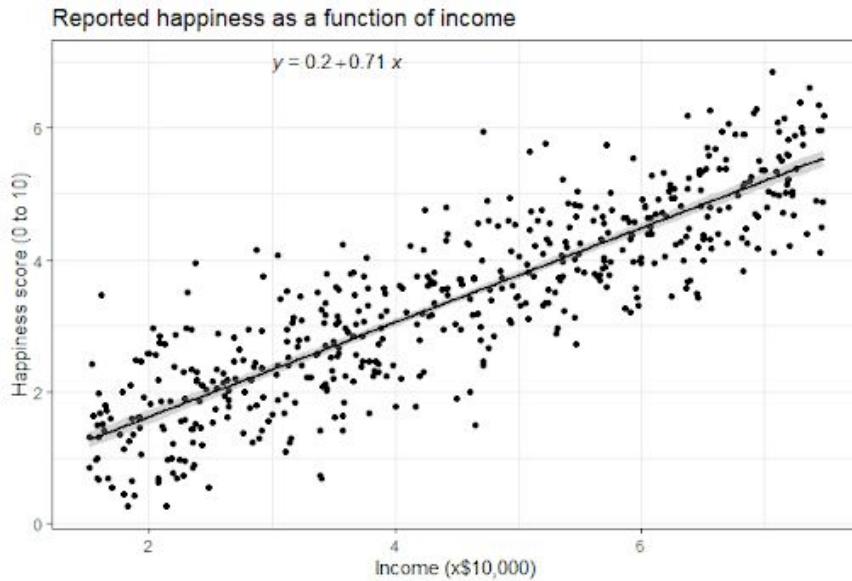


Find the best **linear equation** (or line) that describes the relationship between the features and target variable

Linear Regression

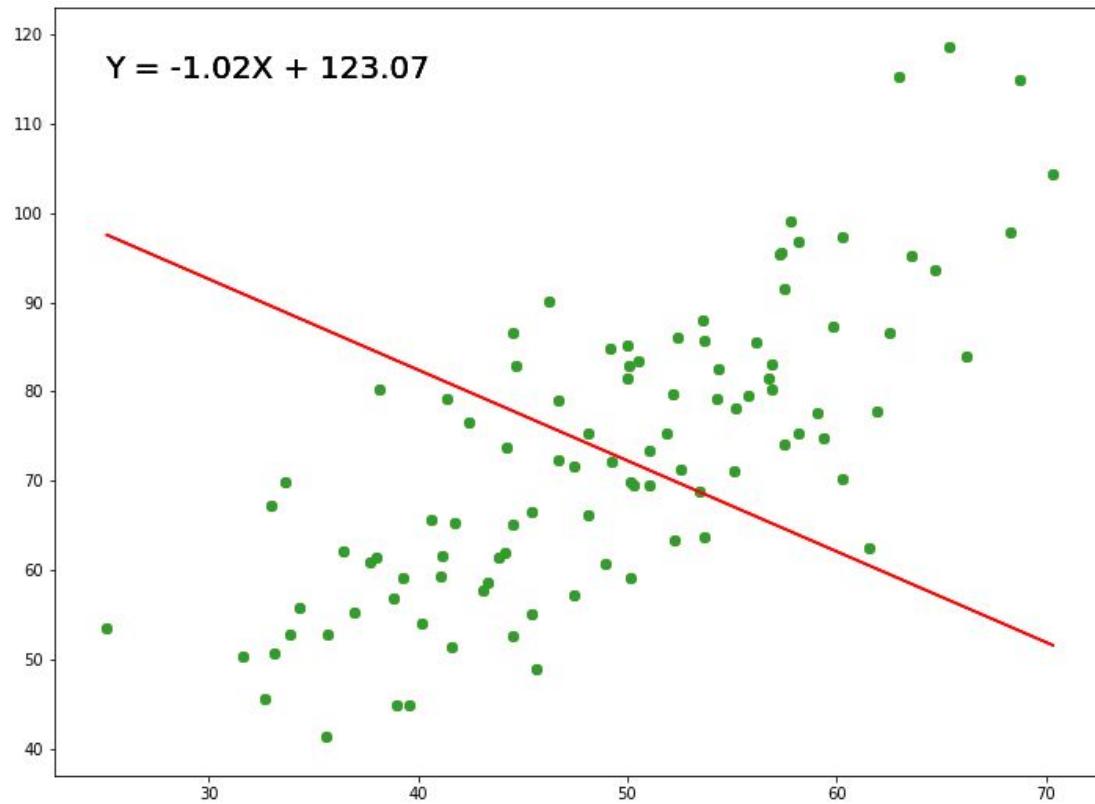


Find the best **linear equation** (or line) that describes the relationship between the features and target variable



Feature (X): Income

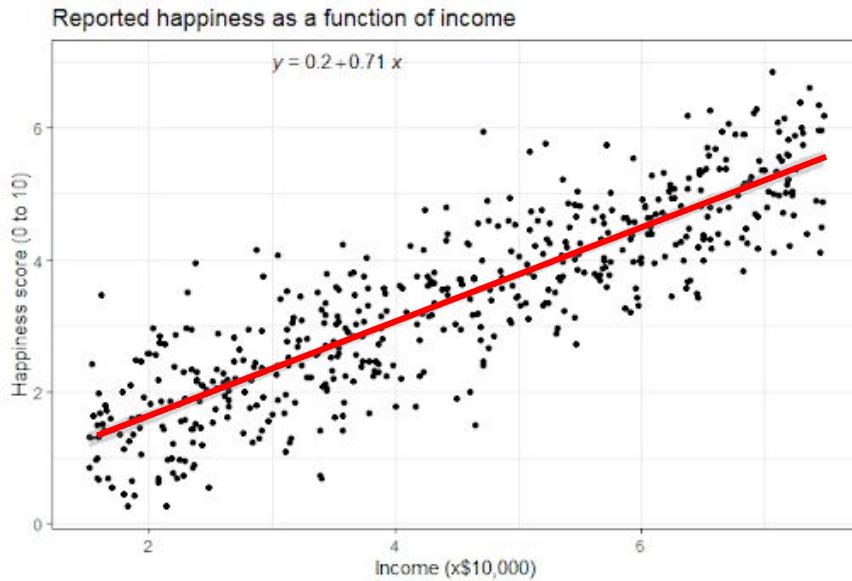
Target (Y): Happiness score (0 to 10)



Linear Regression



Find the best **linear equation** (or line) that describes the relationship between the features and target variable



Linear equation

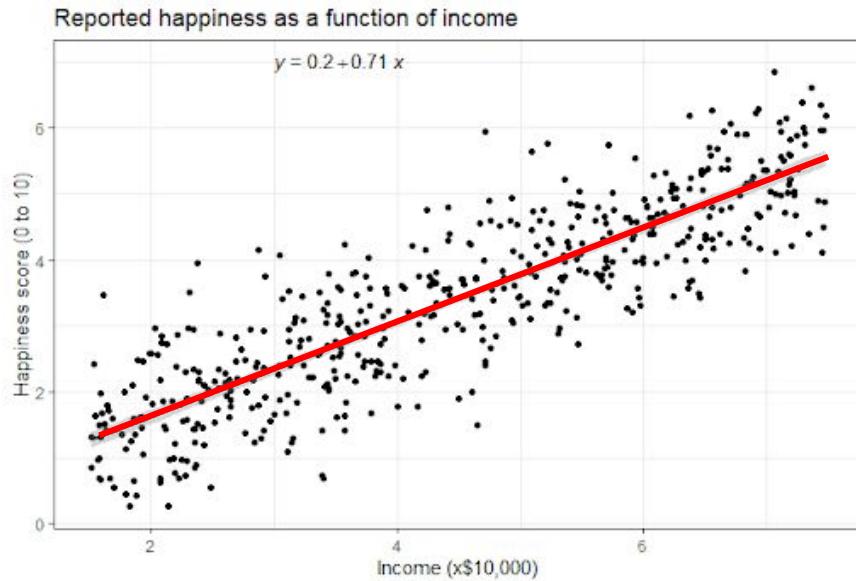
$$Y = a + X * b$$

$$Y = a + X_1 * b_1 + X_2 * b_2 + \dots + X_n * b_n$$

Linear Regression



Find the best **linear equation** (or line) that describes the relationship between the features and target variable



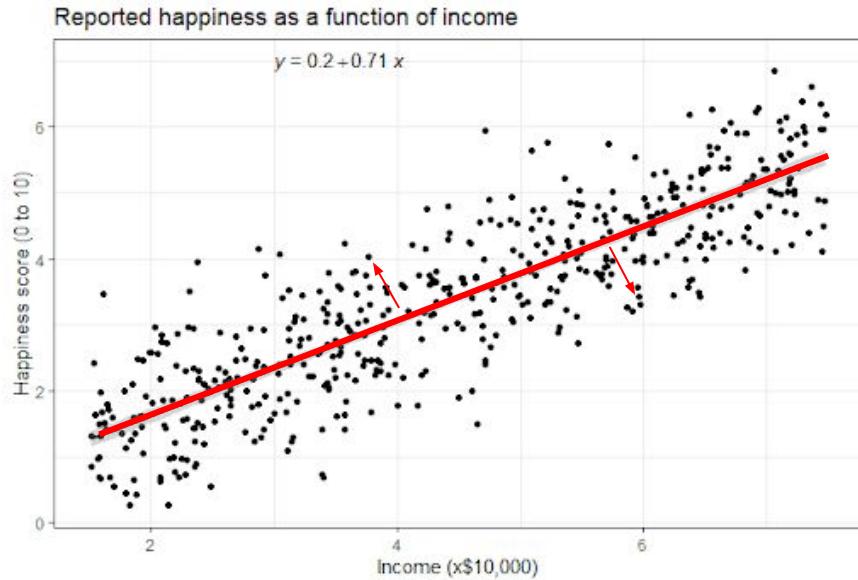
Linear equation

*Happiness Score = a + Income * b*

Linear Regression



Find the best **linear equation** (or line) that describes the relationship between the features and target variable



Find optimal parameters by
minimizing distances with the
estimated line

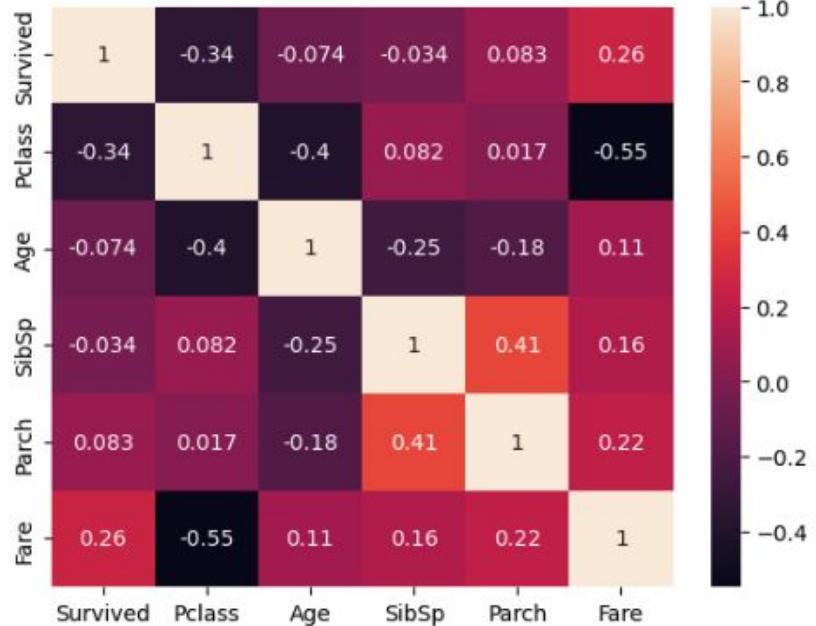
Linear Regression - Pros

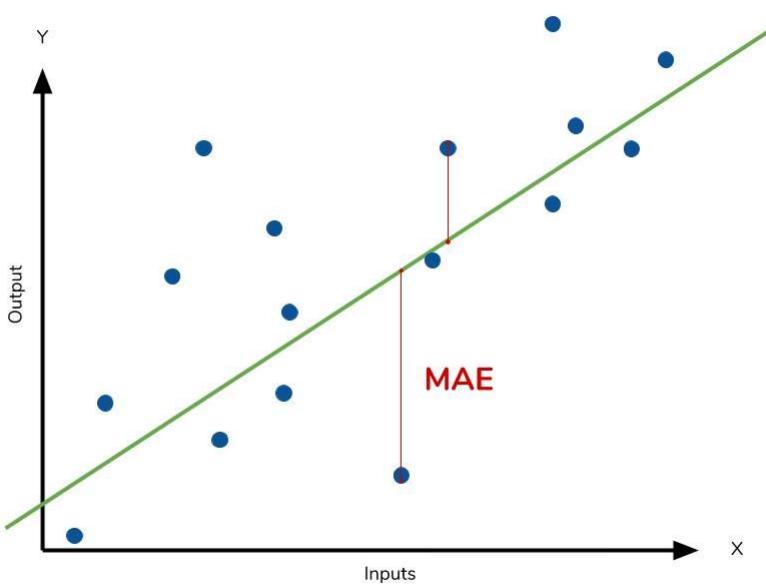
- **Simple model :** The Linear regression model is the simplest equation using which the relationship between the multiple predictor variables and predicted variable can be expressed.
- **Computationally efficient :** The modeling speed of Linear regression is fast as it does not require complicated calculations and runs predictions fast when the amount of data is large.
- **Interpretability of the Output:** The ability of Linear regression to determine the relative influence of one or more predictor variables to the predicted value when the predictors are independent of each other is one of the key reasons of the popularity of Linear regression. The model derived using this method can express the what change in the predictor variable causes what change in the predicted or target variable.

```
call:  
lm(formula = sales ~ ., data = adv_training)  
  
Residuals:  
    Min      1Q  Median      3Q     Max  
-8.6331 -0.8971  0.2283  1.1971  2.9630  
  
Coefficients:                                         P_class < 0.05 to be used for prediction  
              Estimate Std. Error t value Pr(>|t|)  
(Intercept) 2.707724  0.354788  7.632 2.05e-12 ***  
TV          0.046521  0.001561 29.801 < 2e-16 ***  
radio        0.188857  0.009423 20.041 < 2e-16 ***  
newspaper    0.001619  0.006255   0.259    0.796  
---  
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1  
  
Residual standard error: 1.677 on 158 degrees of freedom  
Multiple R-squared:  0.8967,    Adjusted R-squared:  0.8955  
F-statistic: 457.1 on 3 and 158 DF,  p-value: < 2.2e-16
```

R : % of variations that could be explained by the variation of variables.

Correlation between numerical variables





Method of Least Squares

$$\sum e_i^2 = \sum (Y_i - \hat{Y}_i)^2$$

random error for X_i
 $e_i = Y_i - \hat{Y}_i$

observed value for Y_i
predicted value for Y_i

rvalue 1 = 80k

rvalue 2 = 20k

pred1 = 70k

pred2 = 60k

$$(80k - 70k) + (20k - 60k) = -30k$$

$$| 80k - 70k | + | 20k - 60k | = 50k \text{ mistakes}$$

$$(80k - 70k)^2 + (20k - 60k)^2 = 100 + 1600 = 1700 \text{ mistakes}$$

Evaluation for Regression models



Mean squared error:

Measure the distance between the predicted value and the true value

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

Mean Error Squared

- Average errors of every predicted observation/point

Evaluation for Regression models



Root mean squared error:

Compute the square root of the mean squared error

$$RMSE = \sqrt{\frac{\sum_{i=1}^n (\hat{y}_i - y_i)^2}{n}}$$

**Mean squared
error**

- More intuitive as it is the same unit as the original data

$$MSE = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

Mean Square Error formula

```
from sklearn.metrics import mean_squared_error
import math
print(mean_squared_error(Y_test, Y_predicted))
print(math.sqrt(mean_squared_error(Y_test, Y_predicted)))
# MSE: 2017904593.23
# RMSE: 44921.092965684235
```

$$MAE = \frac{1}{N} \sum_{i=1}^N |y_i - \hat{y}_i|$$

Mean Absolute Error formula

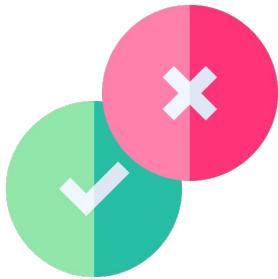
```
from sklearn.metrics import mean_absolute_error
print(mean_absolute_error(Y_test, Y_predicted))
#MAE: 26745.1109986
```

$$MAE = \frac{1}{n} \sum_{i=1}^n |\hat{y}_i - y_i|$$

\hat{y}_i = *predicted value*

y_i = *actual value*

n = *# of observations*



Linear Regression

Pros

- Simple method
- Good interpretation
- Easy to implement

Cons

- Assumes linear relationship between dependent and independent variables, which is incorrect in most cases
- Sensitive to outliers
- If the number of observations are less, it leads to over fitting, it starts considering noise.

Pseudo Code

```
import sklearn
```

- 1) >X= ["nbr bed", "level",
"efficiency", "radius closest metro"]
and y = "price"
- 2) > X_train, y_train, X_test, y_test =
train_test_split(X,y, split_ratio=0.2)
- 3) >lr = LinearRegression()
- 4) >lr.fit(X_train,y_train) training phase
- 5) >pred = lr.predict(X_test) test phase
- 6) > Evaluation phase: (pred-y_test)



Logistic Regression - Binary

Logistic Regression

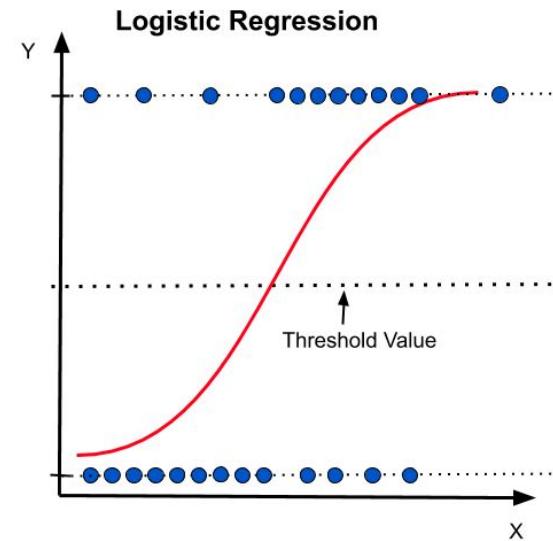
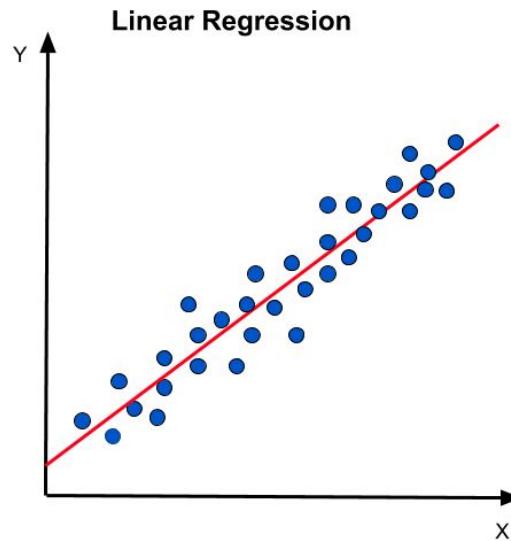


Estimate the probability that an **observation belongs to two possible classes** (binary classification)

Logistic Regression



Estimate the probability that an **observation belongs to two possible classes**



Logistic Regression



Logistic Regression is an **adaptation of Linear Regression** for binary classification tasks

$$Y = a + X * b \quad + \quad \sigma(z) = \frac{1}{1 + \exp(-z)}$$

Linear regression *Sigmoid function*

Logistic Regression



Logistic Regression is an **adaptation of Linear Regression** for binary classification tasks

$$Y = a + X * b \quad + \quad \sigma(z) = \frac{1}{1 + \exp(-z)} \quad = \quad P(Y = 0)$$

Linear regression *Sigmoid function* *Probability of the target being 0*

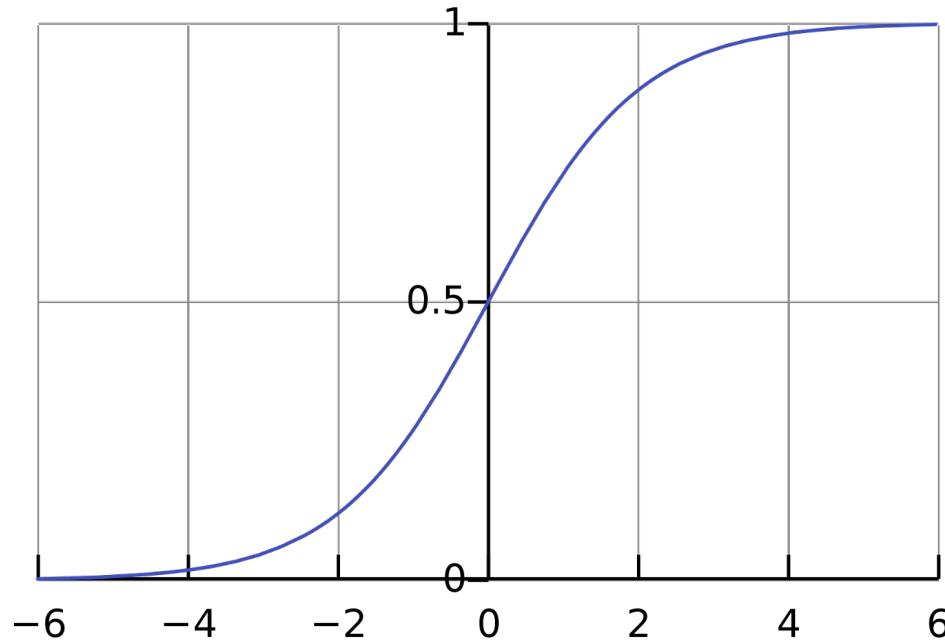
Logistic Regression



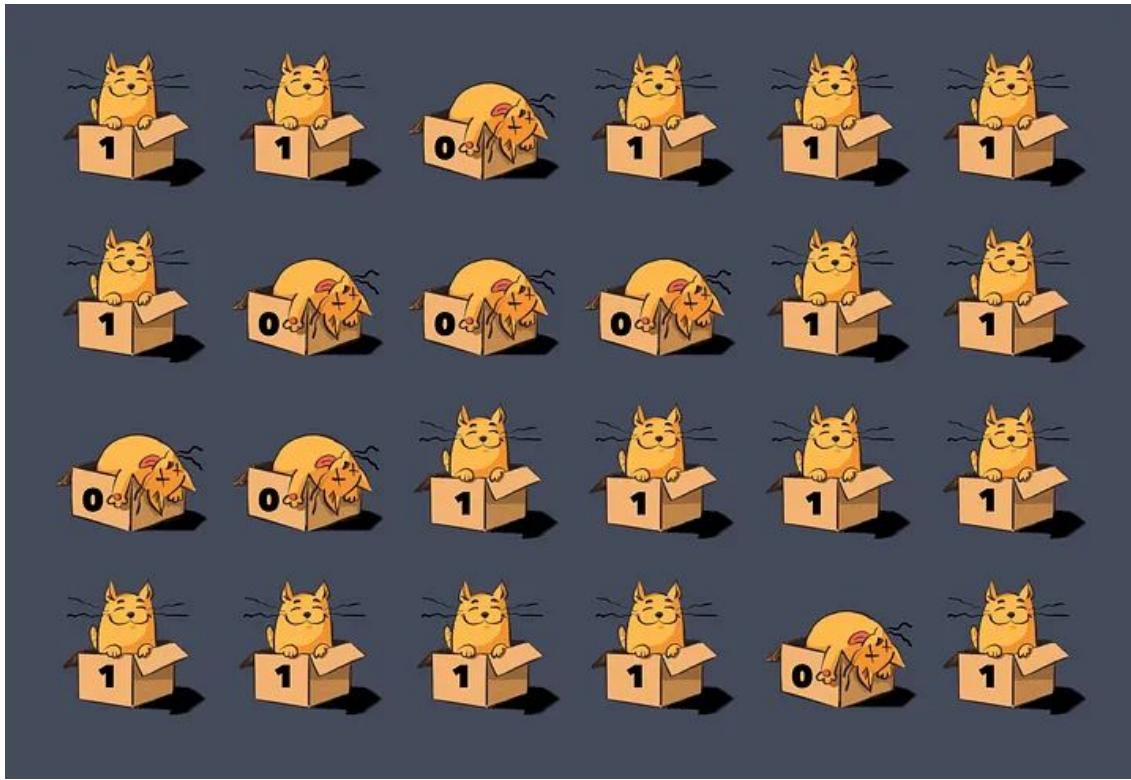
Final predictions are made by setting a threshold value (usually 0.5)

$P(Y = 0) \geq 0.5 \rightarrow \text{The model predicts 0}$

$P(Y = 0) < 0.5 \rightarrow \text{The model predicts 1}$



$$f(x) = \frac{1}{1 + e^{-x}}$$



✓ TRUE POSITIVE

True label



Prediction



✗ FALSE POSITIVE

True label



✗ FALSE NEGATIVE

True label



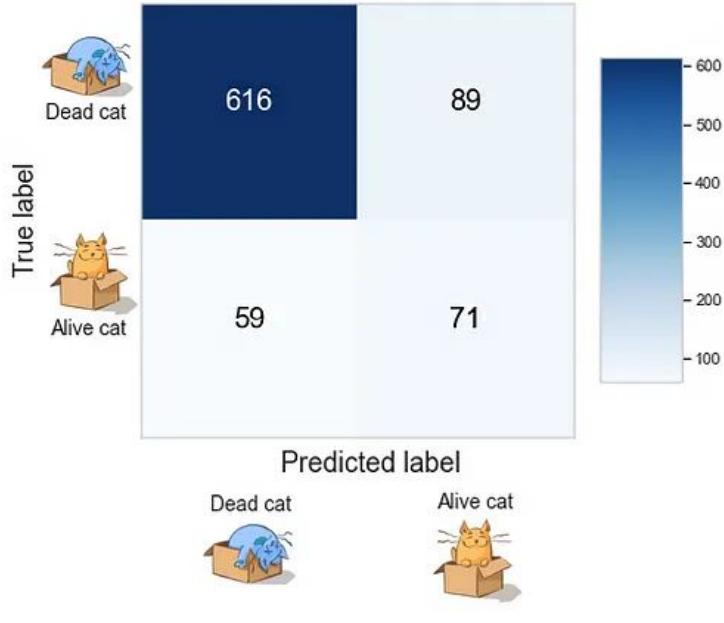
Prediction



✓ TRUE NEGATIVE

True label

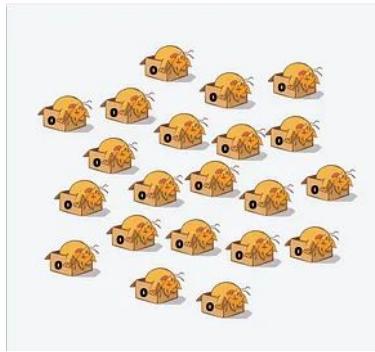




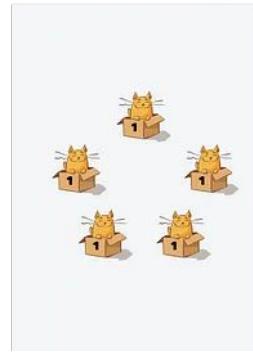
$$\frac{(71(\text{TP}) + 616(\text{TN}))}{(71+616+89+59)} = \frac{687}{834} = 82\%$$

$$\text{Accuracy Score} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{TN} + \text{FP} + \text{FN}}$$

Dead - 95



Alive - 5



$$\text{Accuracy Score} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{TN} + \text{FP} + \text{FN}} = \frac{0 + 95}{0 + 95 + 0 + 5} = 0.95$$

Evaluation for Classification models



F1-score

Combine precision and recall into a single metric (harmonic mean)

$$F1\ score = \frac{2 \times Precision \times Recall}{Precision + Recall}$$

Evaluation for Classification models



Example of Precision and Recall

Let's take the same example

		Predicted label	
		Positive (1)	Negative (0)
True label	Positive (1)	5 True positive	4 False negative
	Negative (0)	10 False positive	100 True negative

Evaluation for Classification models



Example of Precision and Recall

		Predicted label	
		Positive (1)	Negative (0)
True label	Positive (1)	5 True positive	4 False negative
	Negative (0)	10 False positive	100 True negative

$$Precision = \frac{True\ Positive}{True\ Positive + False\ Positive} = \frac{5}{15} = 0.33$$

When the model predicts positive (1), it is correct 33% of the time

Evaluation for Classification models



Example of Precision and Recall

		Predicted label	
		Positive (1)	Negative (0)
True label	Positive (1)	5 True positive	4 False negative
	Negative (0)	10 False positive	100 True negative

$$Recall = \frac{True\ Positive}{True\ Positive + False\ Negative} = \frac{5}{9} = 0.55$$

The model correctly identifies 55% of all positive classes

Evaluation for Classification models



Example of F1 score

Precision = 0.33

Recall = 0.55



$$F1\ score = \frac{2 \times Precision \times Recall}{Precision + Recall} = 0.41$$

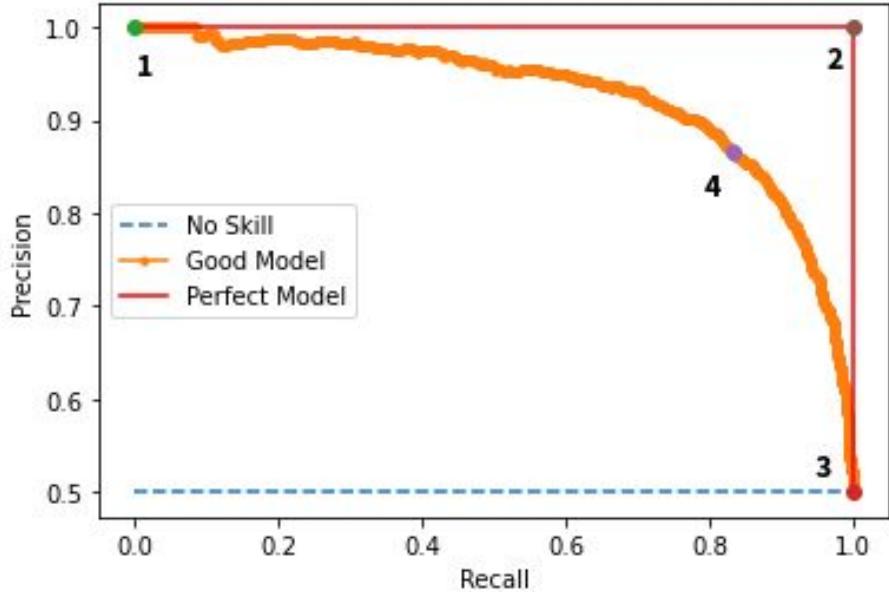
The F1 score is significantly lower than accuracy
(41% vs 88%)

The model has a hard time predicting class 1 correctly

Evaluation for Classification models



Precision and Recall Curve



A « good model » has its curve pointing to the **top right of the graph**

- Tradeoff between Precision and Recall

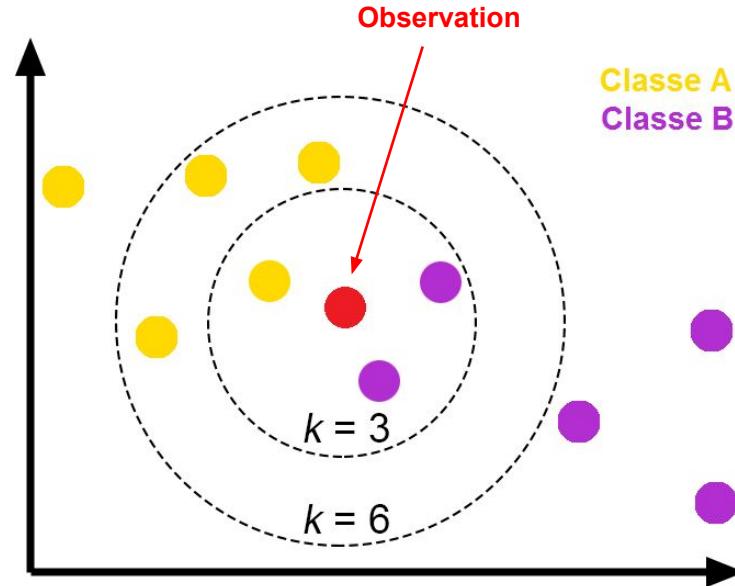


KNN - Classification more than 2 classes

K nearest neighbors (KNN)



Predict the output of an observation using **the data points that are closest to it** (in its neighborhood)



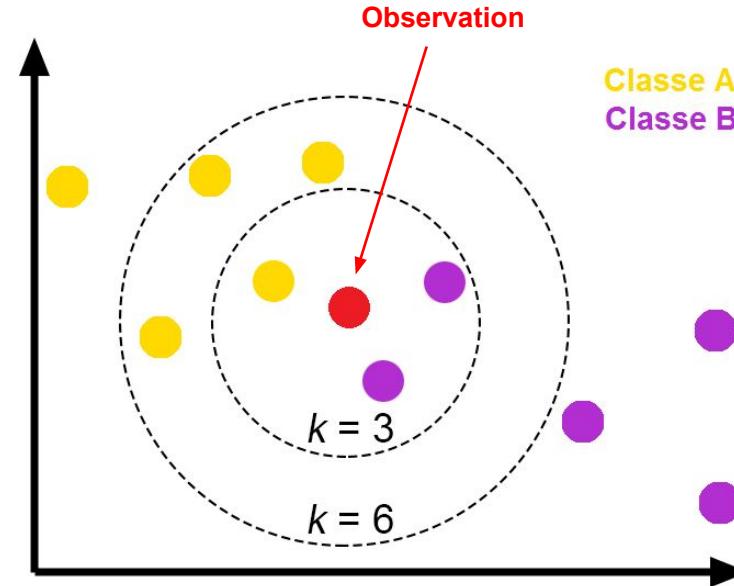
K nearest neighbors (KNN)



Predict the output of an observation using **the data points that are closest to it** (in its neighborhood)

k indicates how many points belong to a neighborhood

- “hyperparameter” that must be chosen

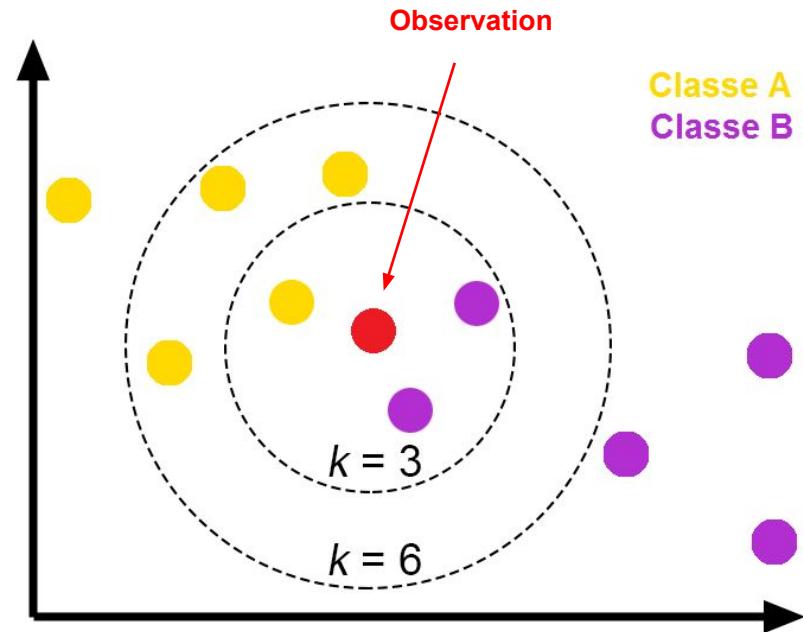


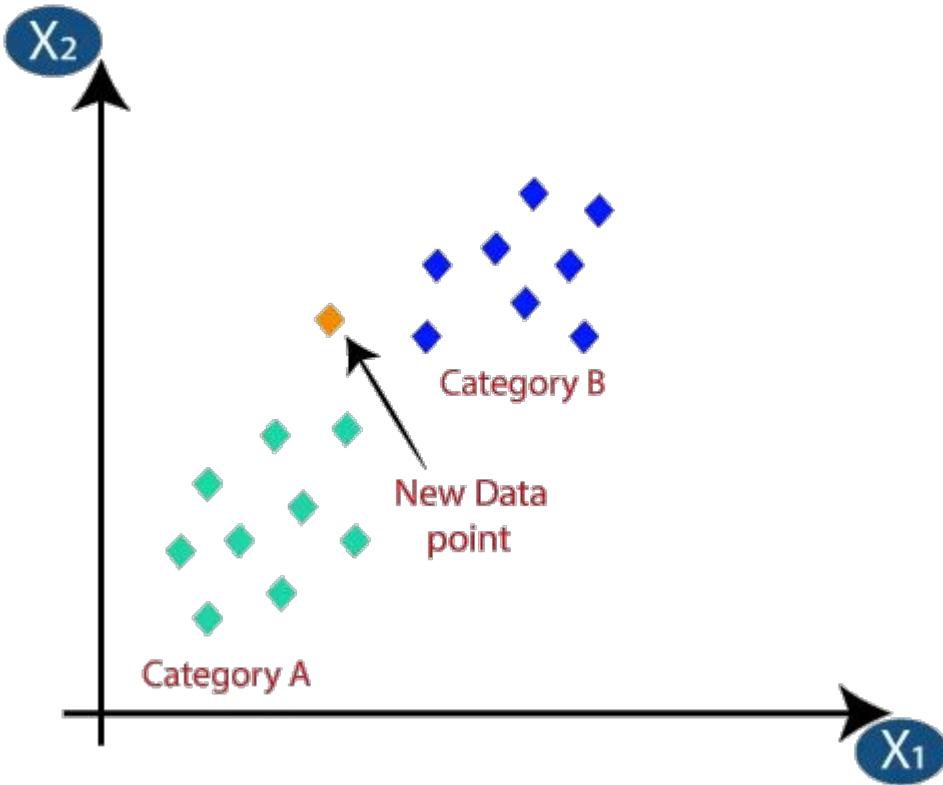
K nearest neighbors (KNN)

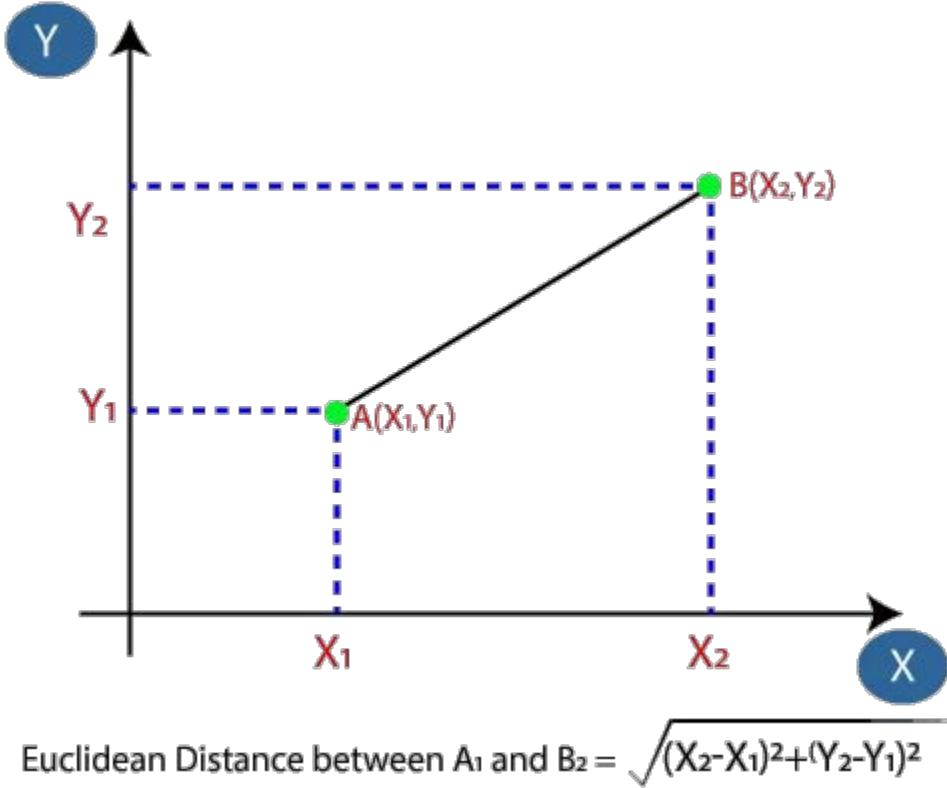


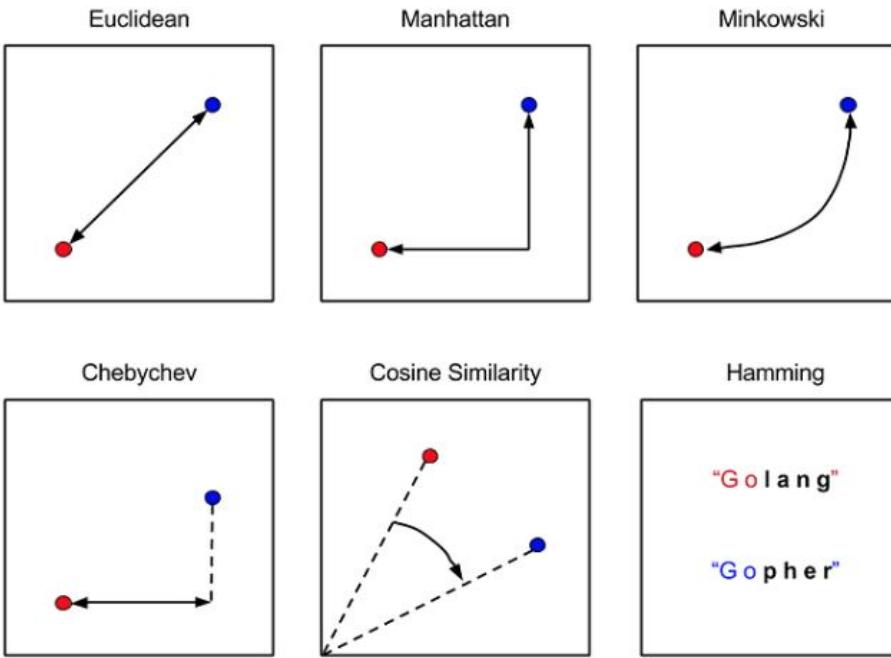
Final prediction:

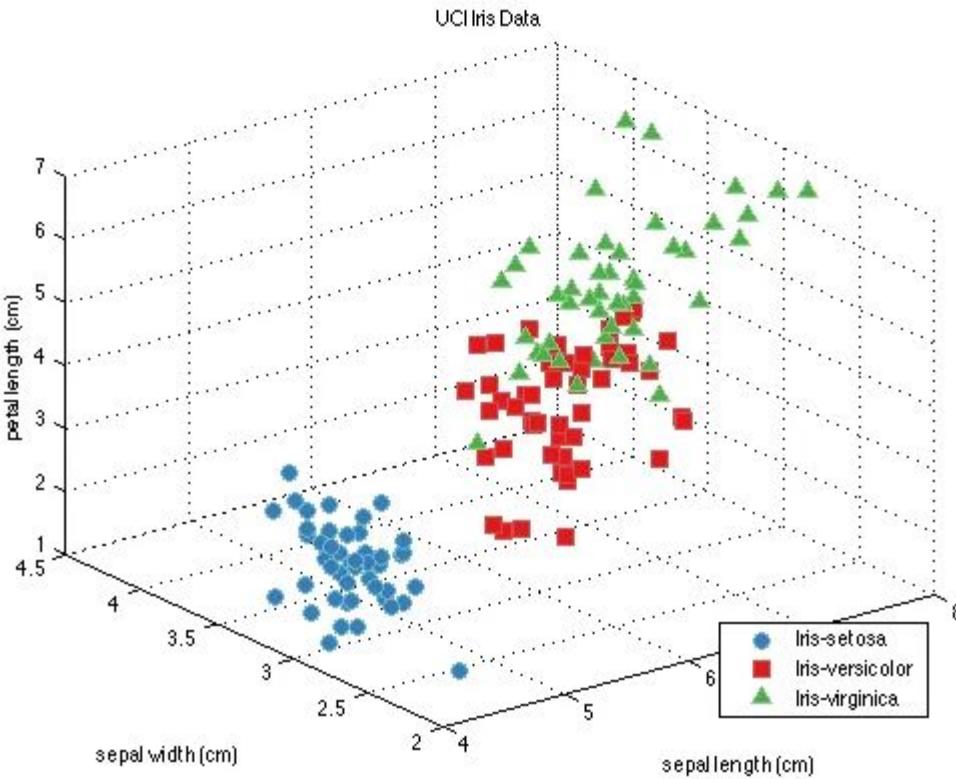
- **Regression**: Compute the mean value of neighbor points
- **Classification**: Assign the most commonly found class in the neighborhood

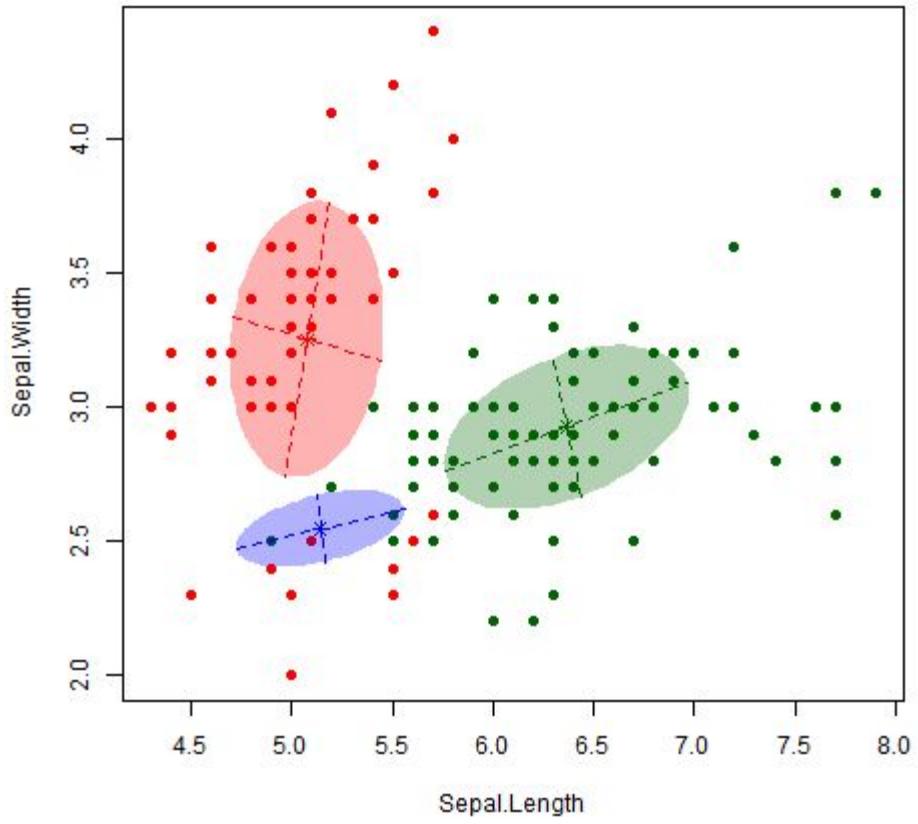












Advantages of KNN Algorithm:

- It is simple to implement.
- It is robust to the noisy training data
- It can be more effective if the training data is large.

Disadvantages of KNN Algorithm:

- Always needs to determine the value of K which may be complex sometimes.
- The computation cost is high because of calculating the distance between the data points for all the training samples.

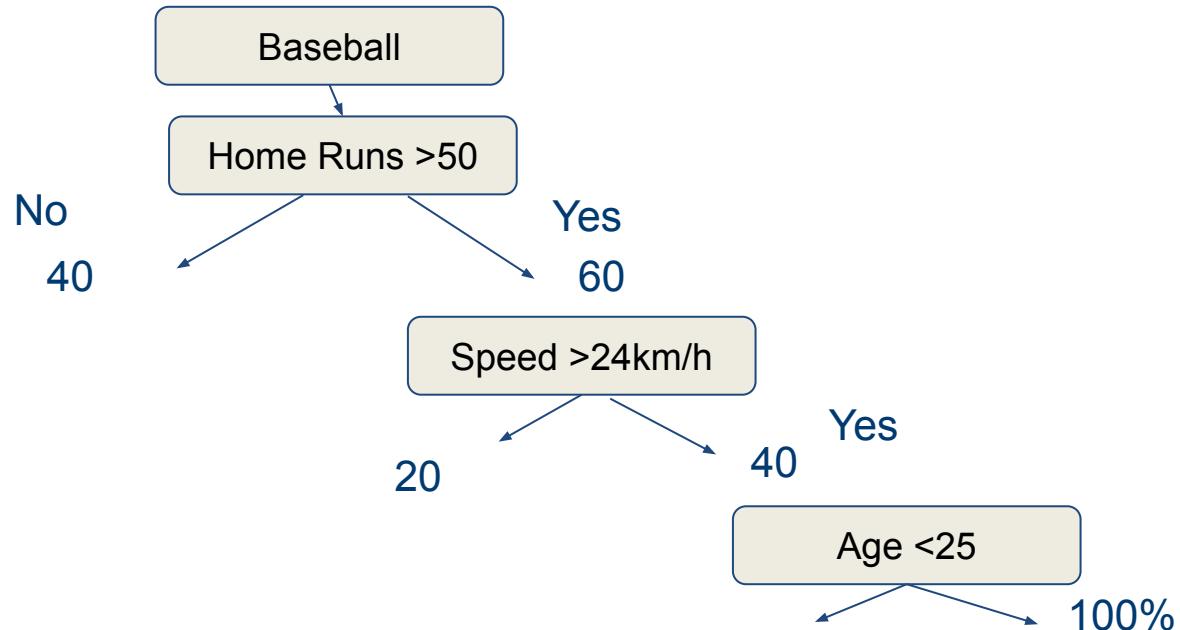


Decision Tree

Decision trees



Build a decision tree that makes predictions by **splitting data into homogenous groups** (group points with similar characteristics)

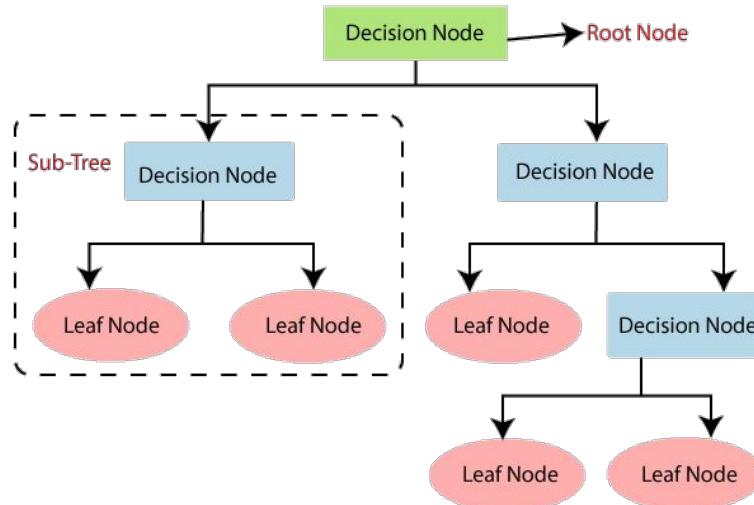


Decision trees



Build a decision tree that makes predictions by **splitting data into homogenous groups** (group points with similar characteristics)

- The model splits data using **binary decision rules**

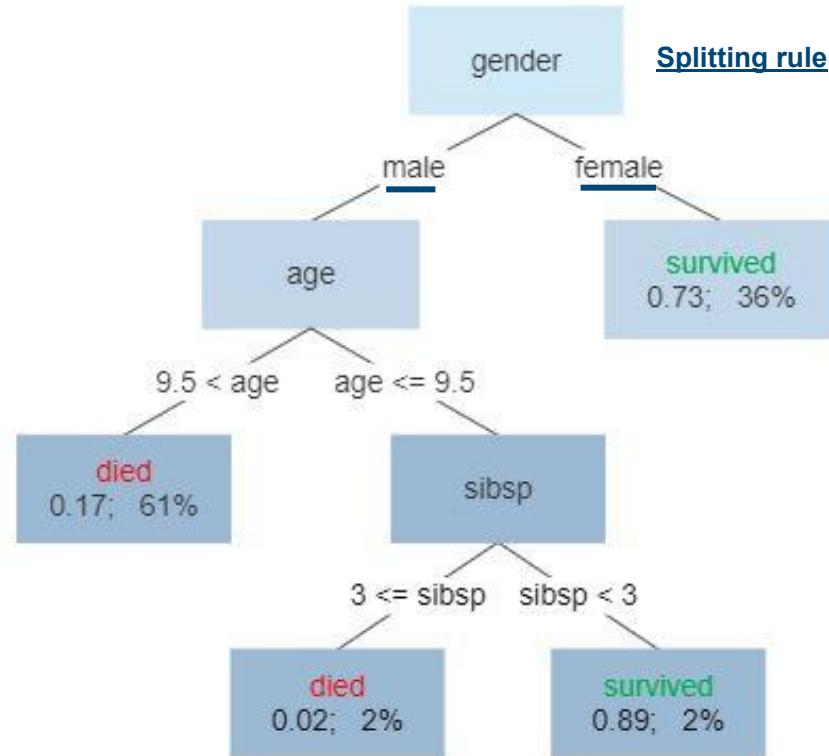


Decision trees



Example of a decision rule:

What gender is the passenger ?

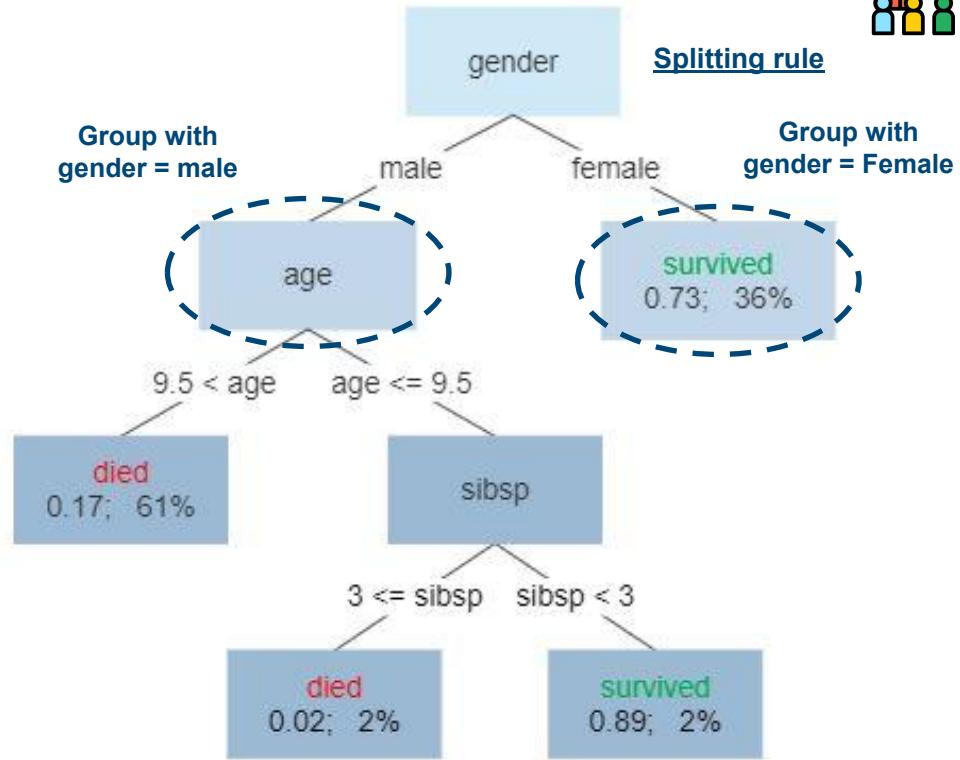


Decision trees



Example of a decision rule:

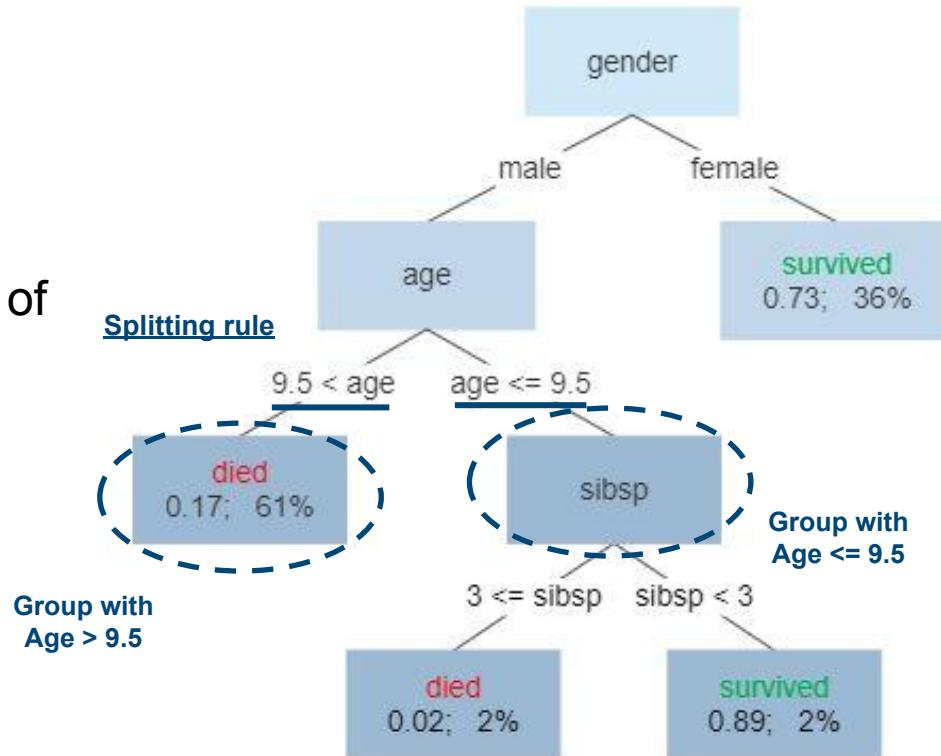
- Gender = Male: left split
- Gender = Female: right split



Decision trees



Continue till you reach the end of the tree (leaf node)



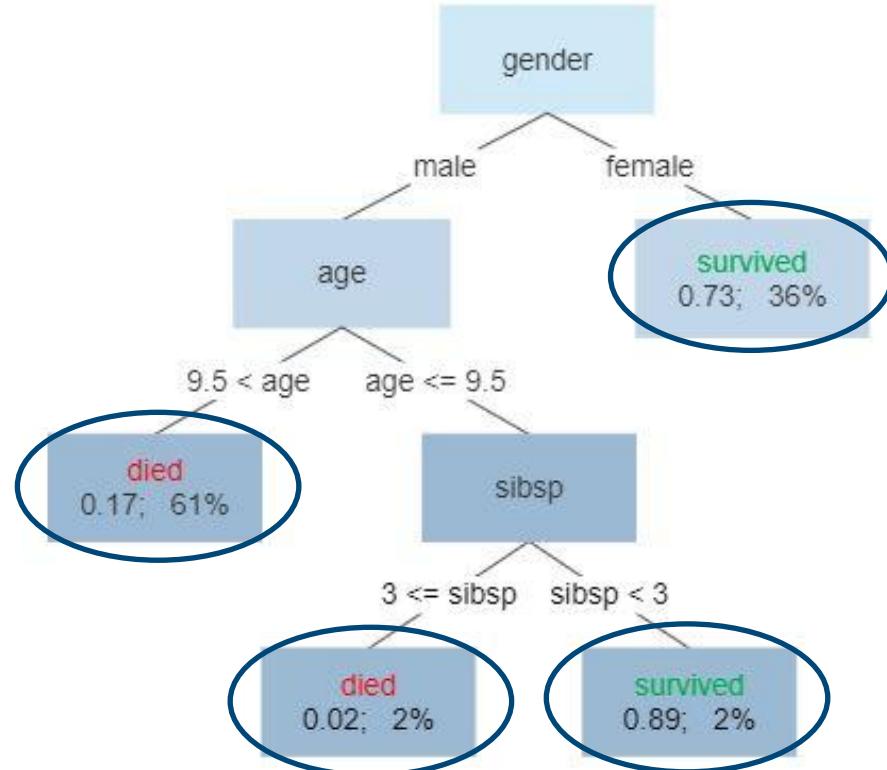
Decision trees



How does the tree stop splitting ?

- Set a max tree depth
- Set a minimum number of points in a leaf/final node

“hyperparameters” that must be chosen

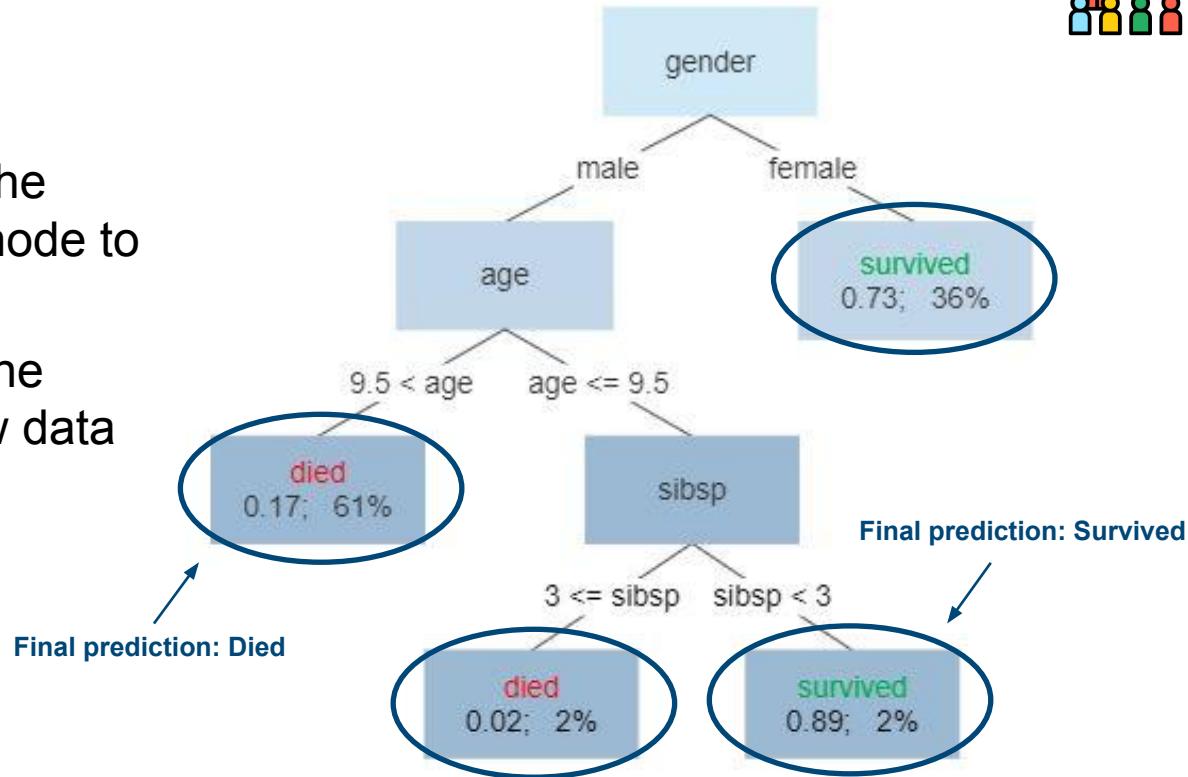


Decision trees

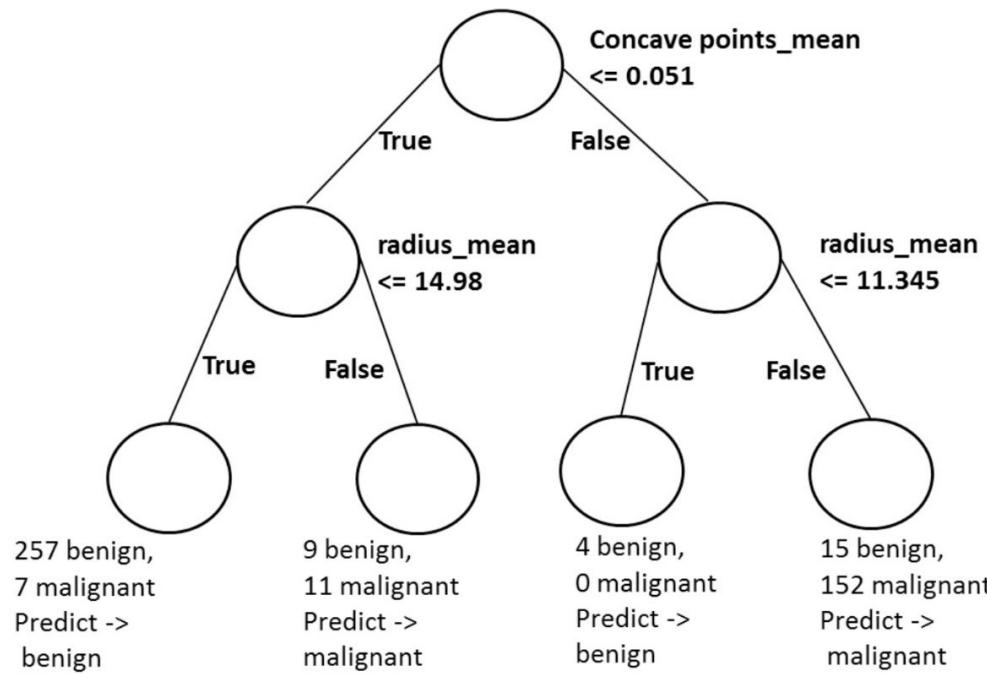


Final prediction

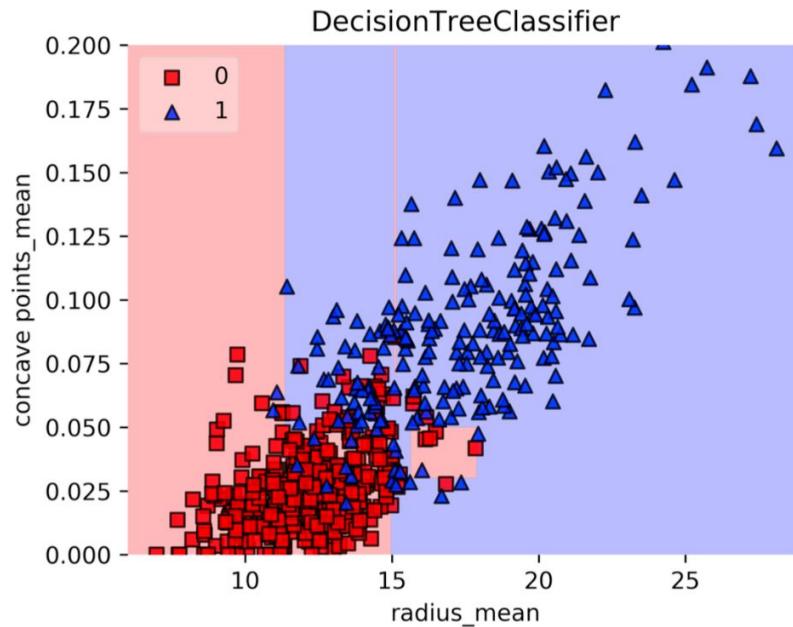
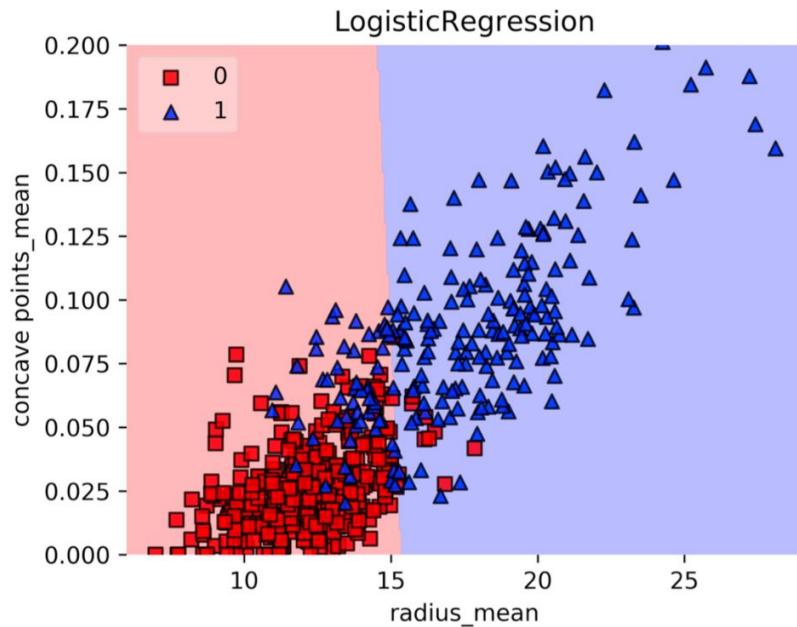
- **Regression:** Compute the average value in a leaf node to make a prediction
- **Classification:** Assign the majority class to the new data point



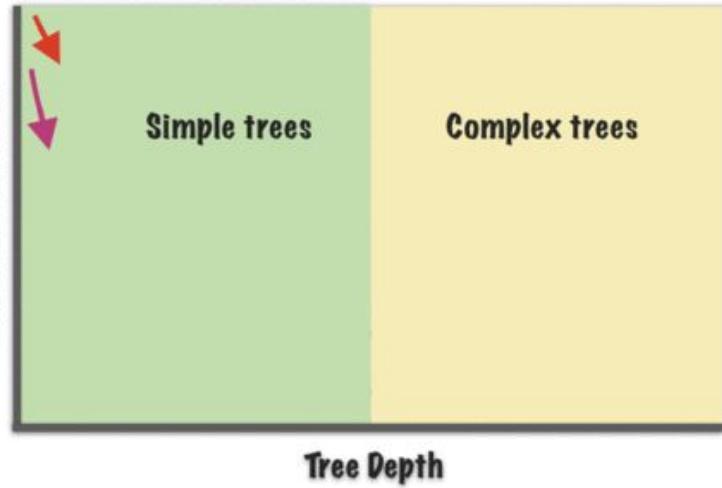
Decision-tree Diagram



Decision Regions: CART vs. Linear Model



Classification Error



- True error
- Training error



Random Forest



Build multiple decorrelated trees than aggregate final results to make a prediction on a new data point

Random Forest



Build multiple decorrelated trees than aggregate final results to make a prediction on a new data point

- Prevent overfitting with regular Decision Trees

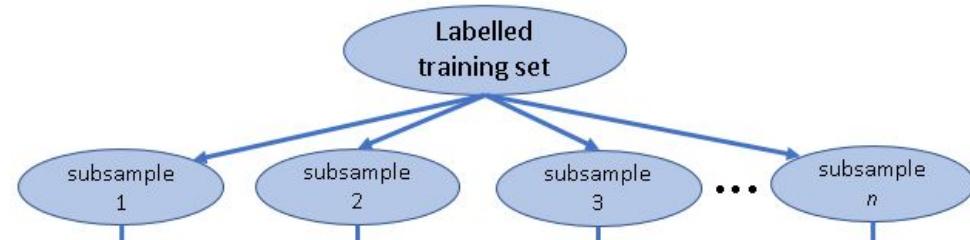
Random Forest



Step 1: Generate
bootstrap samples



Bootstrap
sampling

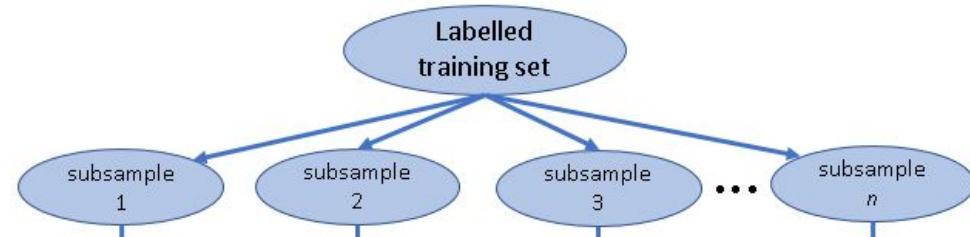


Random Forest

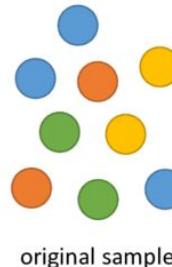


Step 1: Generate bootstrap samples

→ Bootstrap sampling



- ☐ Bootstrapping means resampling multiple times the original data



1 2 3 4 5 6 7 8 9 bootstrap sample 1

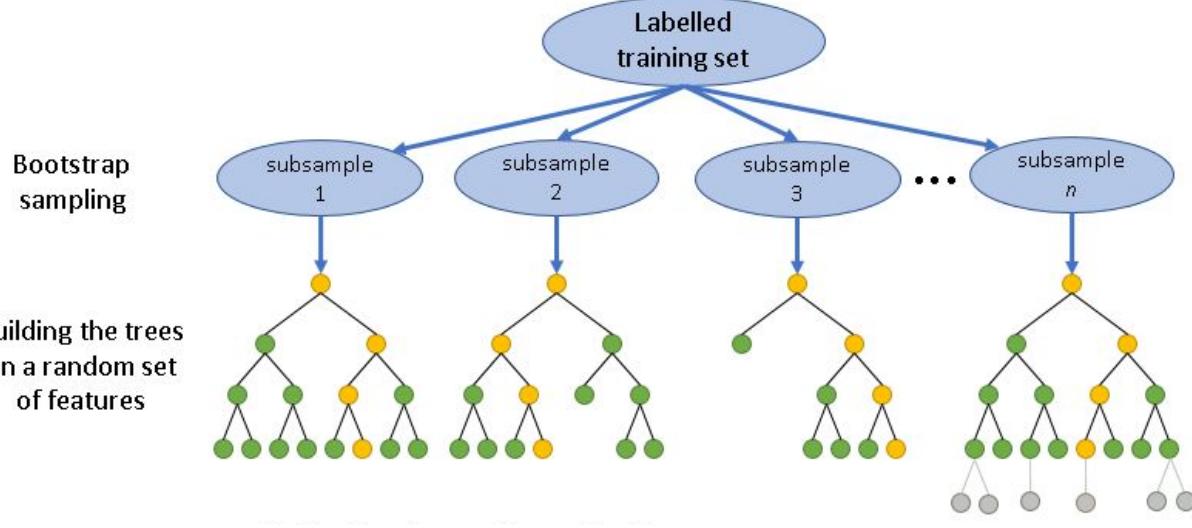
1 2 3 4 5 6 7 8 9 bootstrap sample 2

1 2 3 4 5 6 7 8 9 bootstrap sample 3

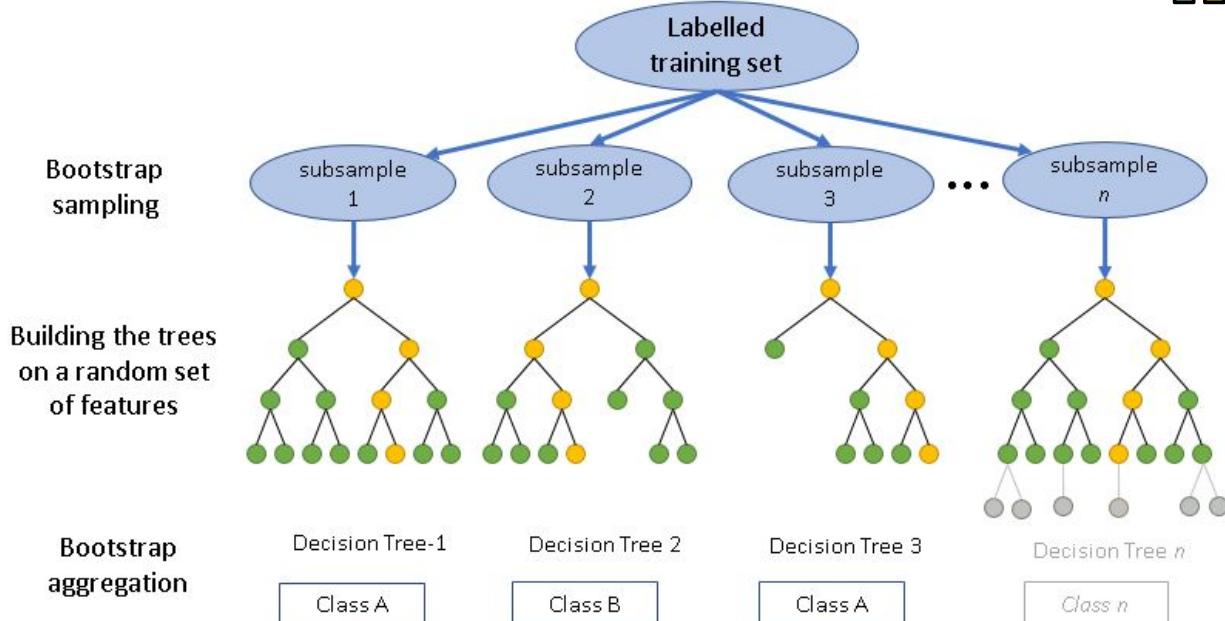
Random Forest



Step 2: Build trees using the bootstrap samples

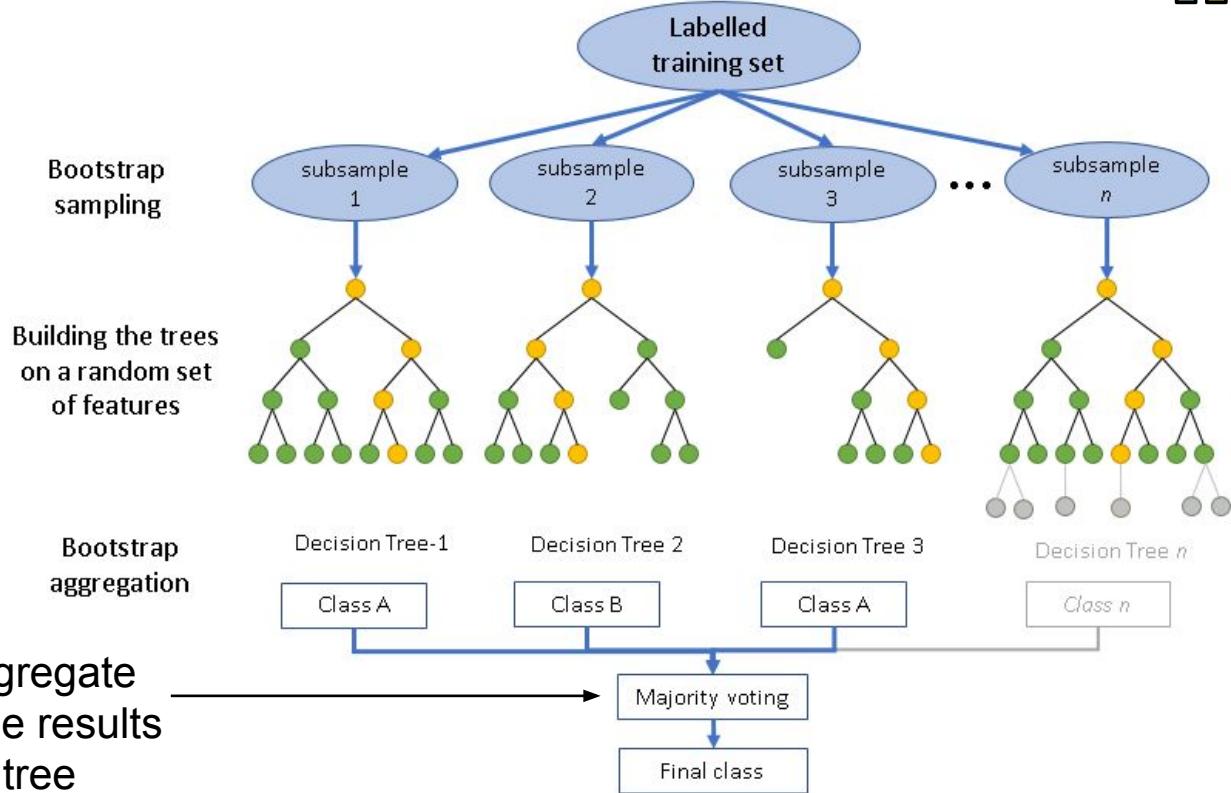


Random Forest



Step 3: Each tree makes its own prediction

Random Forest

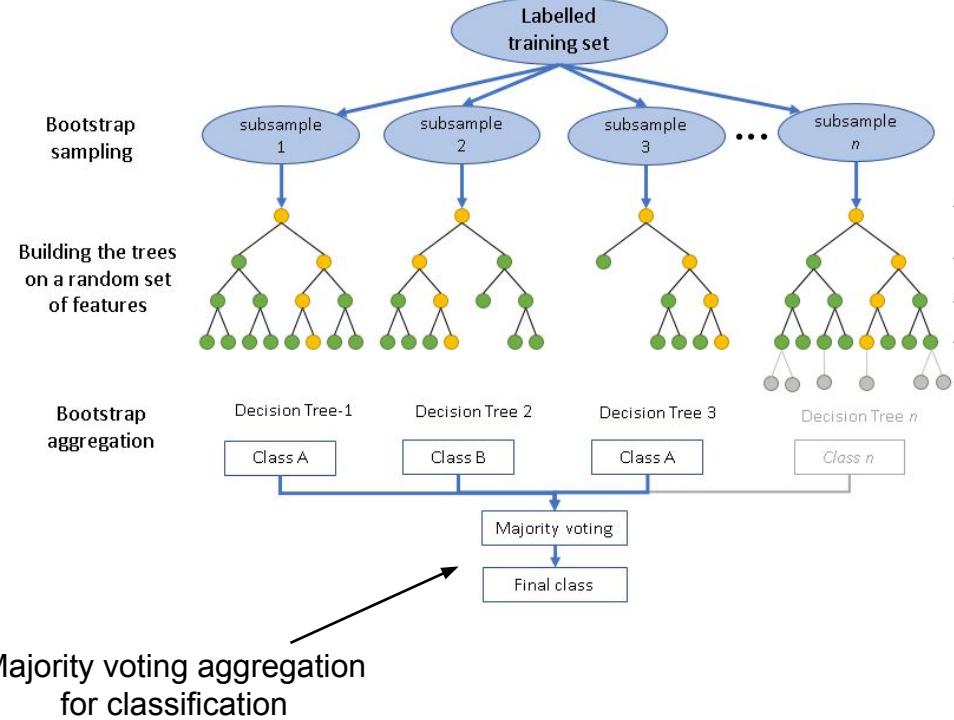


Random Forest



Final prediction

- **Regression:** Compute the average value of each tree predictions
- **Classification:** Apply a majority vote between each class predicted



Summary of Supervised models

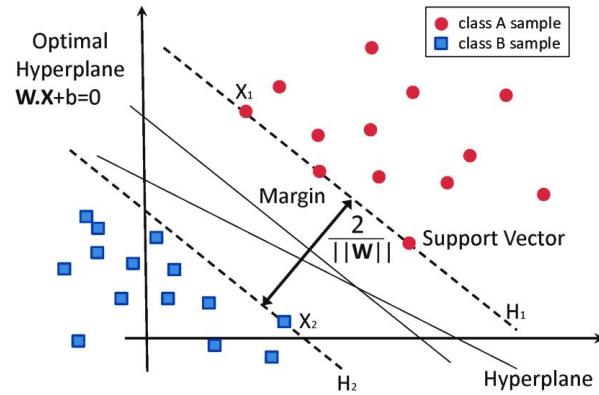


Models	Complexity	Hyperparameters
Linear Regression	Low	None
Logistic Regression	Low	<ul style="list-style-type: none">C: regularization strength
K nearest neighbors	Low	<ul style="list-style-type: none">n_neighbors: number of neighbors (k)
Decision Trees	Medium	<ul style="list-style-type: none">max_depth: depth of the treemin_samples_leaf: min number of points in a leaf node)
Random Forest	High	<ul style="list-style-type: none">Same as Decision Treen_estimators: Number of treesmax_features: Number of features for building trees
Boosting	High	<ul style="list-style-type: none">Same as Decision Tree + n_estimatorslearning rate: Scale each trees contributionSubsample: Fraction of samples used to build individual trees

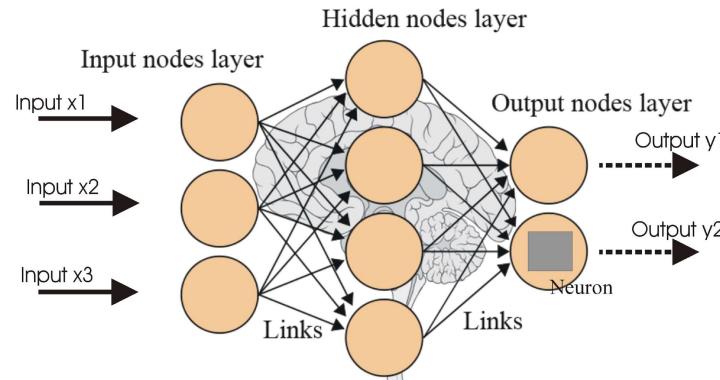
Other supervised models (classification)



- **Support Vector Machine (SVM):** Find the line that best separates classes by maximizing the margin between them.
- **Neural networks (MLP):** Layers of “neurons” that process input data by adjusting weights, mimicking the way the human brain works.



Support Vector Machine (SVM)



Neural Network

Model training with Scikit-learn



Demo of how to use scikit-learn for model training

https://scikit-learn.org/stable/supervised_learning.html



Model evaluation

Boosting



Build multiple trees sequentially by adapting from previous errors, instead of building them independently

Boosting



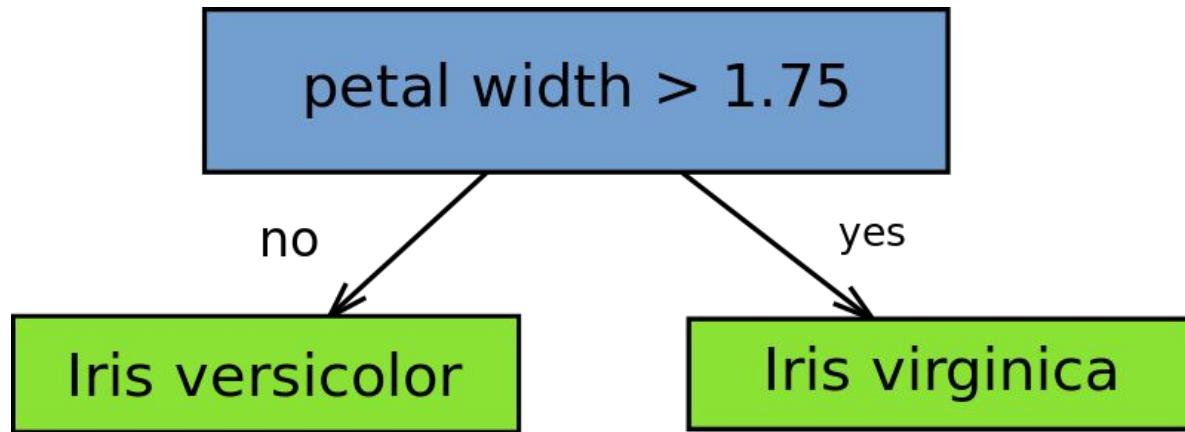
Build multiple trees sequentially by adapting from previous errors, instead of building them independently

- Any **weak learner** can be used, not only decision trees

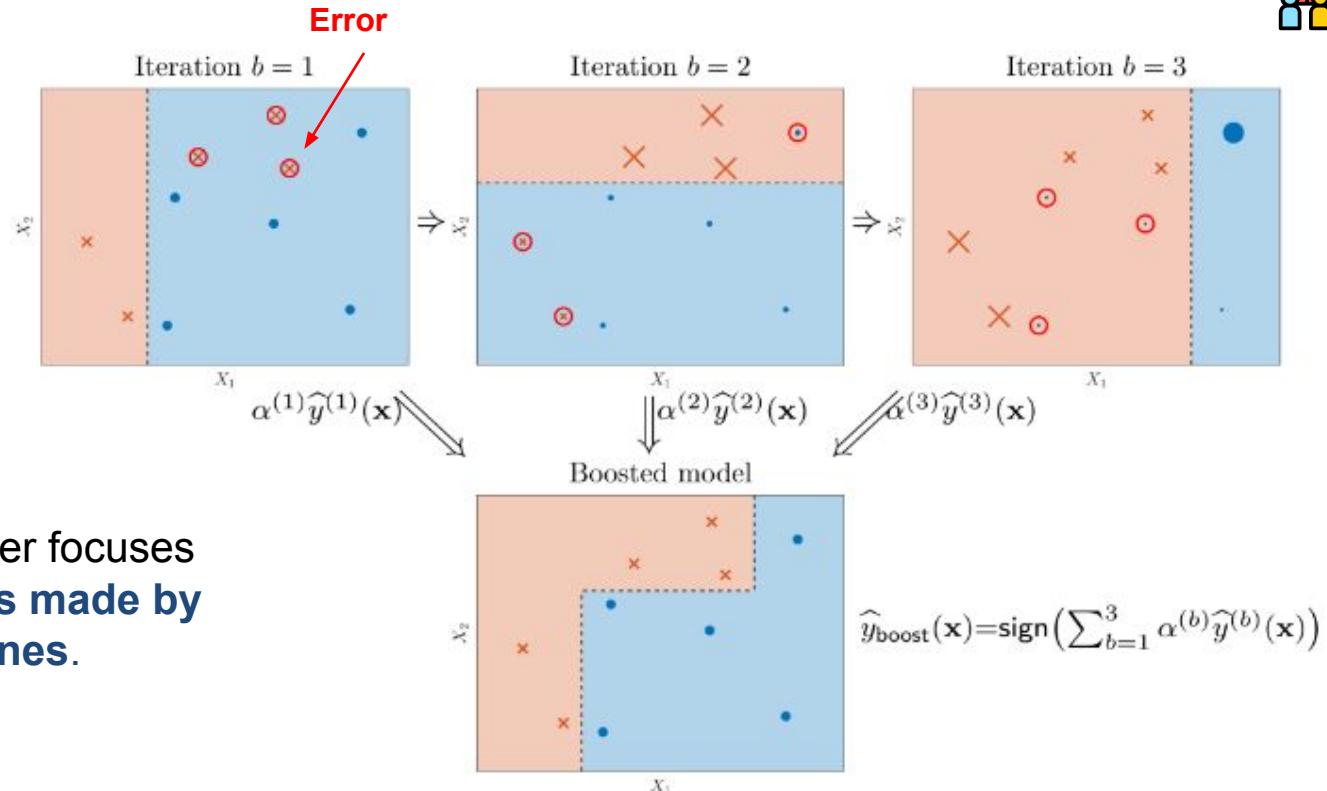
Boosting



Weak learners are models that perform slightly better than a random guess (with a 50% accuracy)



Boosting

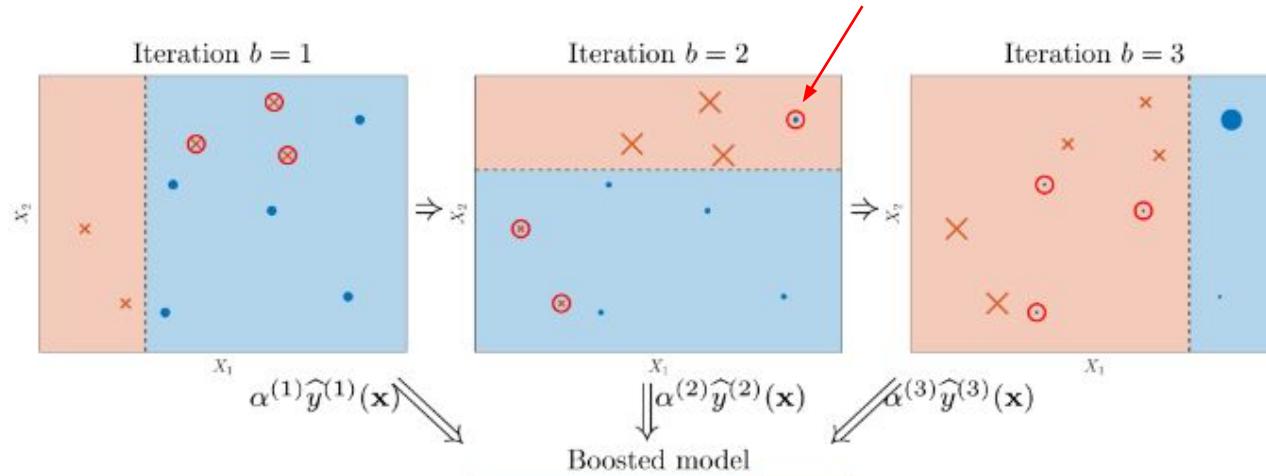


Each new learner focuses
on the **mistakes made by**
the previous ones.

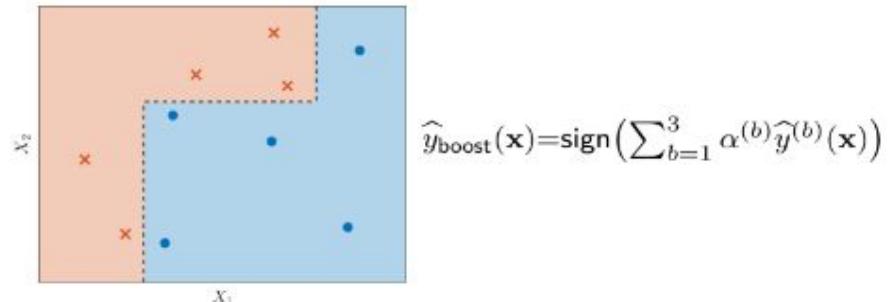
Boosting



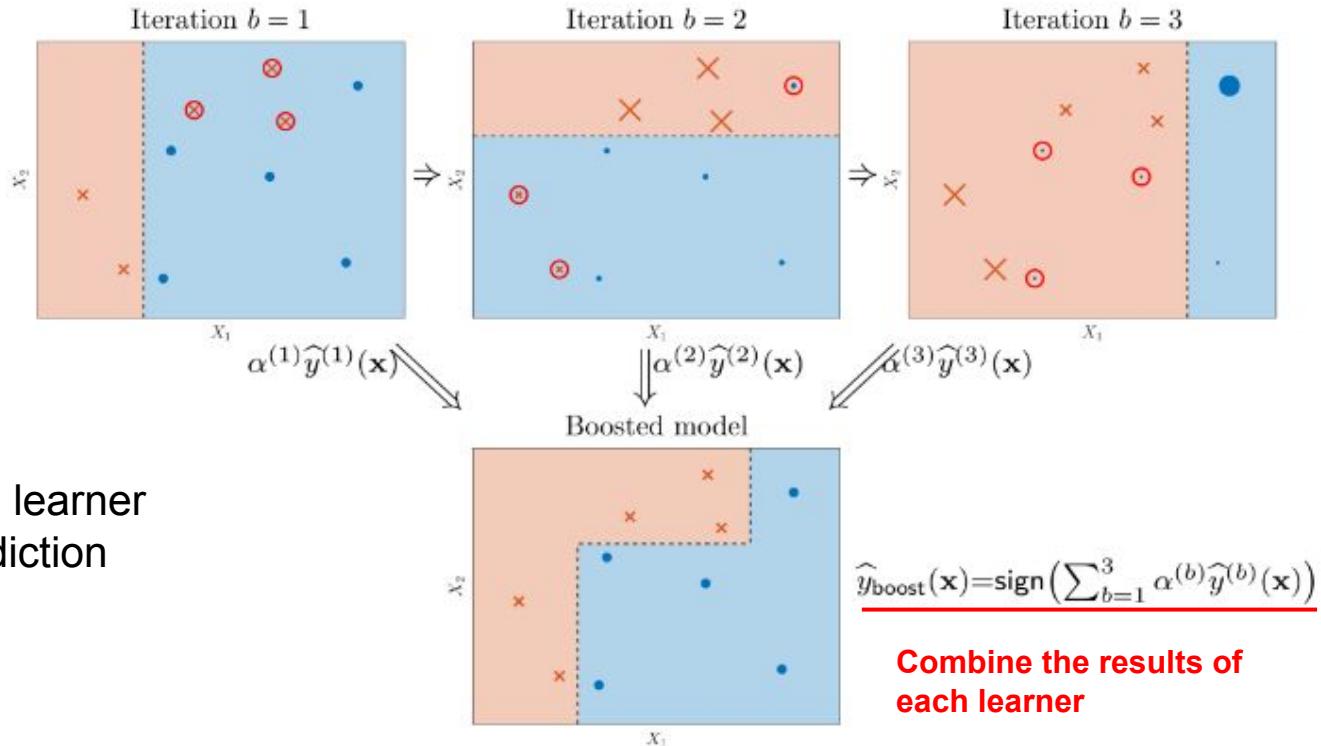
Focus on the previous errors



Each new learner focuses on the **mistakes made by the previous ones**.



Boosting



Different types of Boosting



- **AdaBoost**: Adjust the weights of training data
- **Gradient Boosting**: Correct residual errors
- **Extreme Gradient Boosting (XGBoost)**: Fast implementation of Gradient Boosting

XGBoost

Evaluation for Regression models



Important information:

- A lower MSE and RMSE signals a better performance
- MSE and RMSE are **not bounded between 0 and 1**
- Should only be used to compare models trained on the same data

Evaluating a Machine Learning model

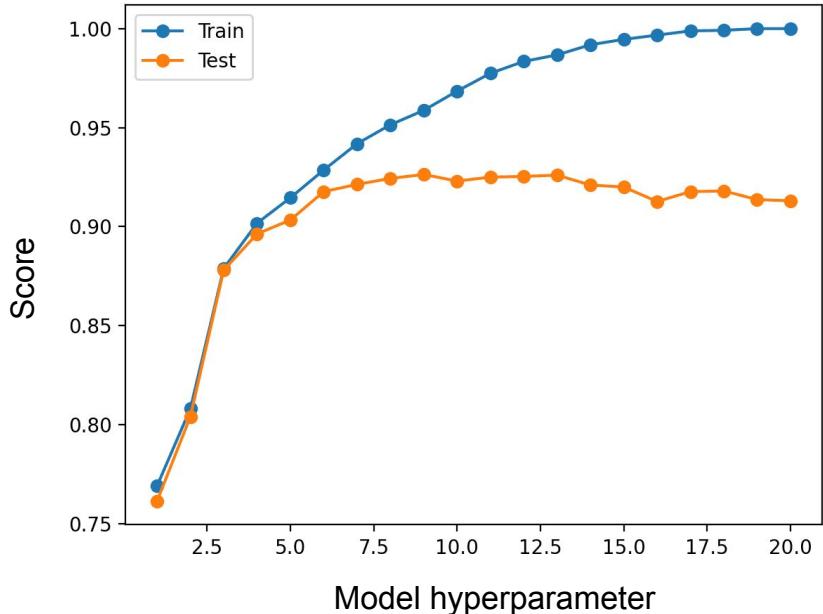


- Evaluation in Machine Learning is performed by **comparing the predicted values of a model with true values**

Evaluating a Machine Learning model



- Evaluation is done on the test data. This ensure the model is **reliable** on data it hasn't seen during training.

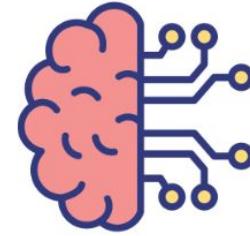
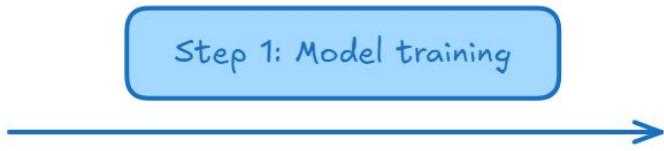


Evaluating a Machine Learning model



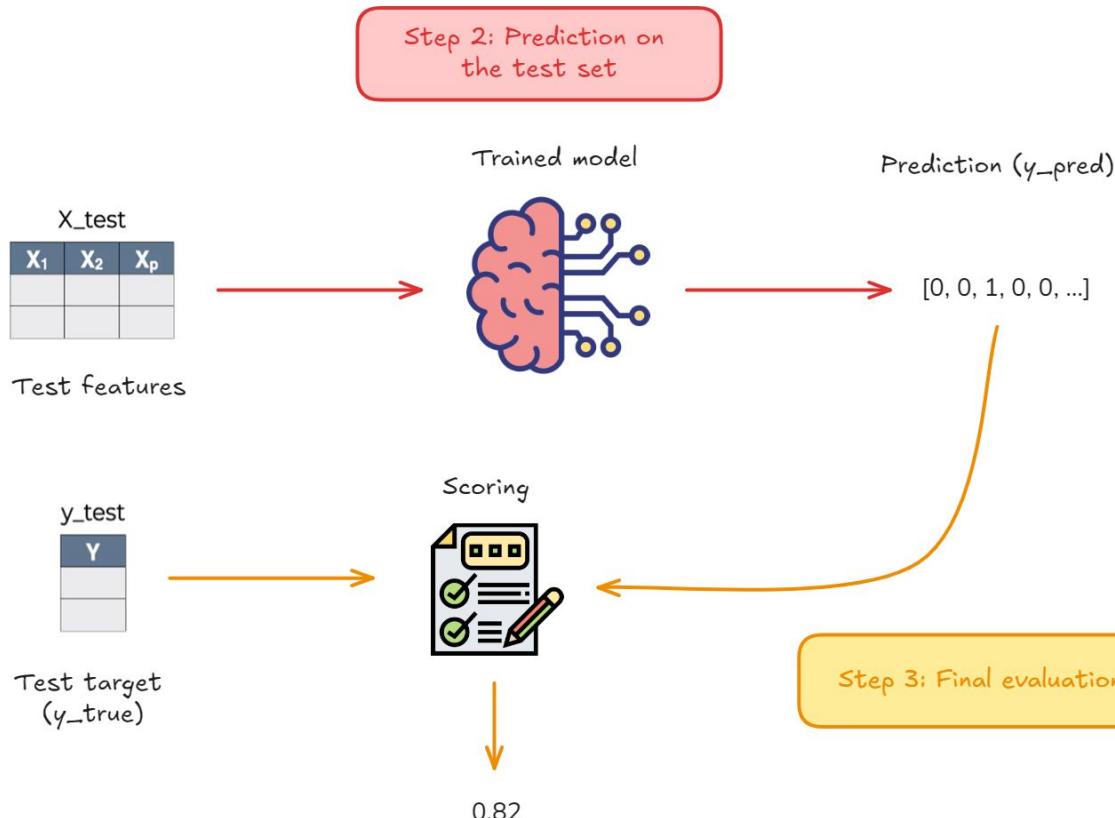
X_train			y_train
X ₁	X ₂	X _p	Y

Training features and target
(train split)



Machine Learning model

Evaluating a Machine Learning model



Evaluation for Classification models



Accuracy

Ratio of correct predictions to the total number of predictions

$$Accuracy = \frac{Nbr\ of\ correct\ predictions}{Total\ nbr\ of\ predictions}$$

- Accuracy is always between 0 and 1
- Reflects the **overall performance of a model**

Evaluation for Classification models



Example of Accuracy

Our classification model predicts two classes (0 and 1)

Let's test its performance on **109 new observations**

- **105 are correct predictions**
- 14 are incorrect predictions

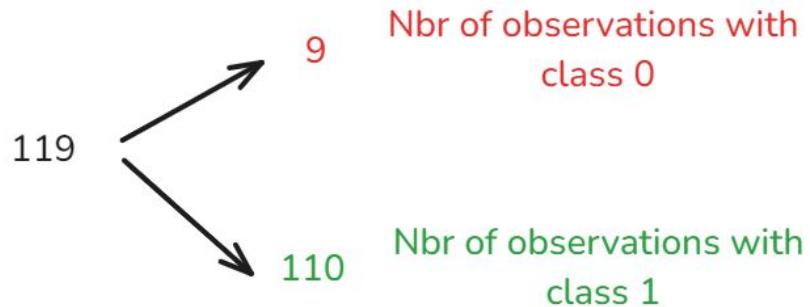
Accuracy = $105/119 = 0.88$ (88%)

Evaluation for Classification models



Example of Accuracy

Now let's study its performance separately on each predicted class



Evaluation for Classification models



Example of Accuracy

Now let's study its performance separately on each predicted class



Accuracy can be misleading if the target classes are not in the same proportion

Evaluation for Classification models



Confusion Matrix

Summarize the performance of a classifier for each predicted class

		Predicted label	
		Positive (1)	Negative (0)
True label	Positive (1)	5 True positive	4 False negative
	Negative (0)	10 False positive	100 True negative

Evaluation for Classification models



Confusion Matrix

Summarize the performance of a classifier for each predicted class

		Predicted label	
		Positive (1)	Negative (0)
True label	Positive (1)	5 True positive	4 False negative
	Negative (0)	10 False positive	100 True negative

Nbr of correct predictions for 1

Nbr of correct predictions for 0

True positive

False negative

False positive

True negative

Evaluation for Classification models



Confusion Matrix

Summarize the performance of a classifier for each predicted class

		Predicted label	
		Positive (1)	Negative (0)
True label	Positive (1)	5 True positive	4 False negative
	Negative (0)	10 False positive	100 True negative

Nbr of false predictions for 0

Nbr of false predictions for 1

Evaluation for Classification models



Confusion Matrix

Summarize the performance of a classifier for each predicted class

		Predicted label	
		Positive (1)	Negative (0)
True label	Positive (1)	5 True positive	4 False negative
	Negative (0)	10 False positive	100 True negative

→ $\frac{5}{5 + 4} = 0.55$ Class 1 accuracy

→ $\frac{10}{10 + 100} = 0.90$ Class 0 accuracy

Evaluation for Classification models



Precision

Proportion of predicted positives that were actually correct

$$\frac{TP}{TP + FP}$$

Recall

Proportion of true positives that were correctly identified

$$\frac{TP}{TP + FN}$$

Learn more about Precision and Recall

<https://developers.google.com/machine-learning/crash-course/classification/precision-and-recall>

Supervised Learning



Regression



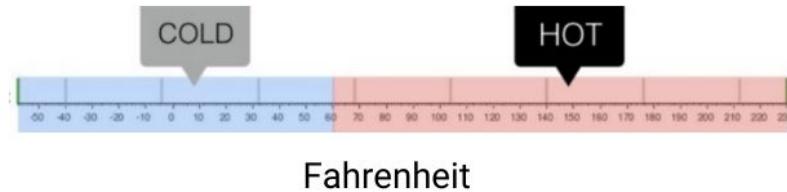
What will be the temperature tomorrow?



Classification



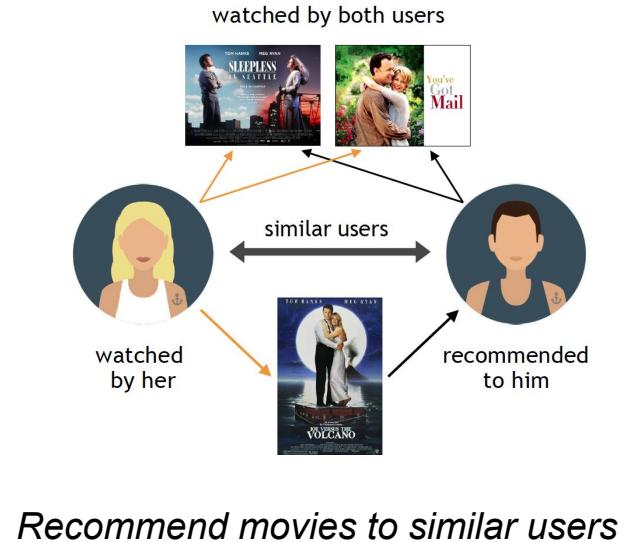
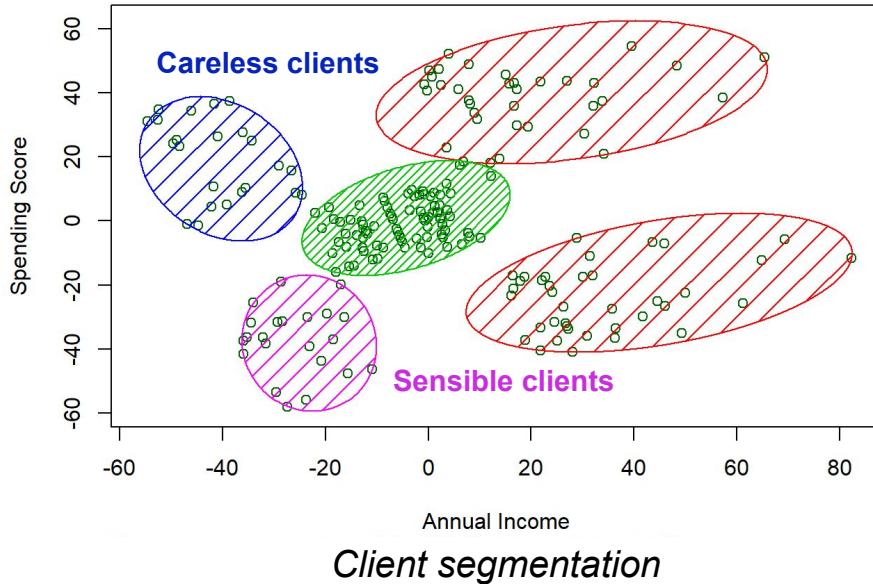
Will it be hot or cold tomorrow?





Unsupervised Learning

1. Clustering: Group data together based on distances/similarities



Summary of evaluation methods



Metric	Type	Characteristics	Bounded	Goal
Accuracy	Classification	Overall performance of the model	Yes [0,1]	Maximise
Confusion Matrix	Classification	Study the performance on all predicted classes	x	x
Precision, Recall and F1 score	Classification	Focus on the performance on the positive class (1)	Yes [0,1]	Maximise

Summary of evaluation methods



Metric	Type	Characteristics	Bounded	Goal
Accuracy	Classification	Overall performance of the model	Yes [0,1]	Maximise
Confusion Matrix	Classification	Study the performance on all predicted classes	x	x
Precision, Recall and F1 score	Classification	Focus on the performance on the positive class (1)	Yes [0,1]	Maximise

For regression models, metrics such as the **Mean Squared Error** or the **R2 score** are used instead

Model Evaluation with Scikit-learn



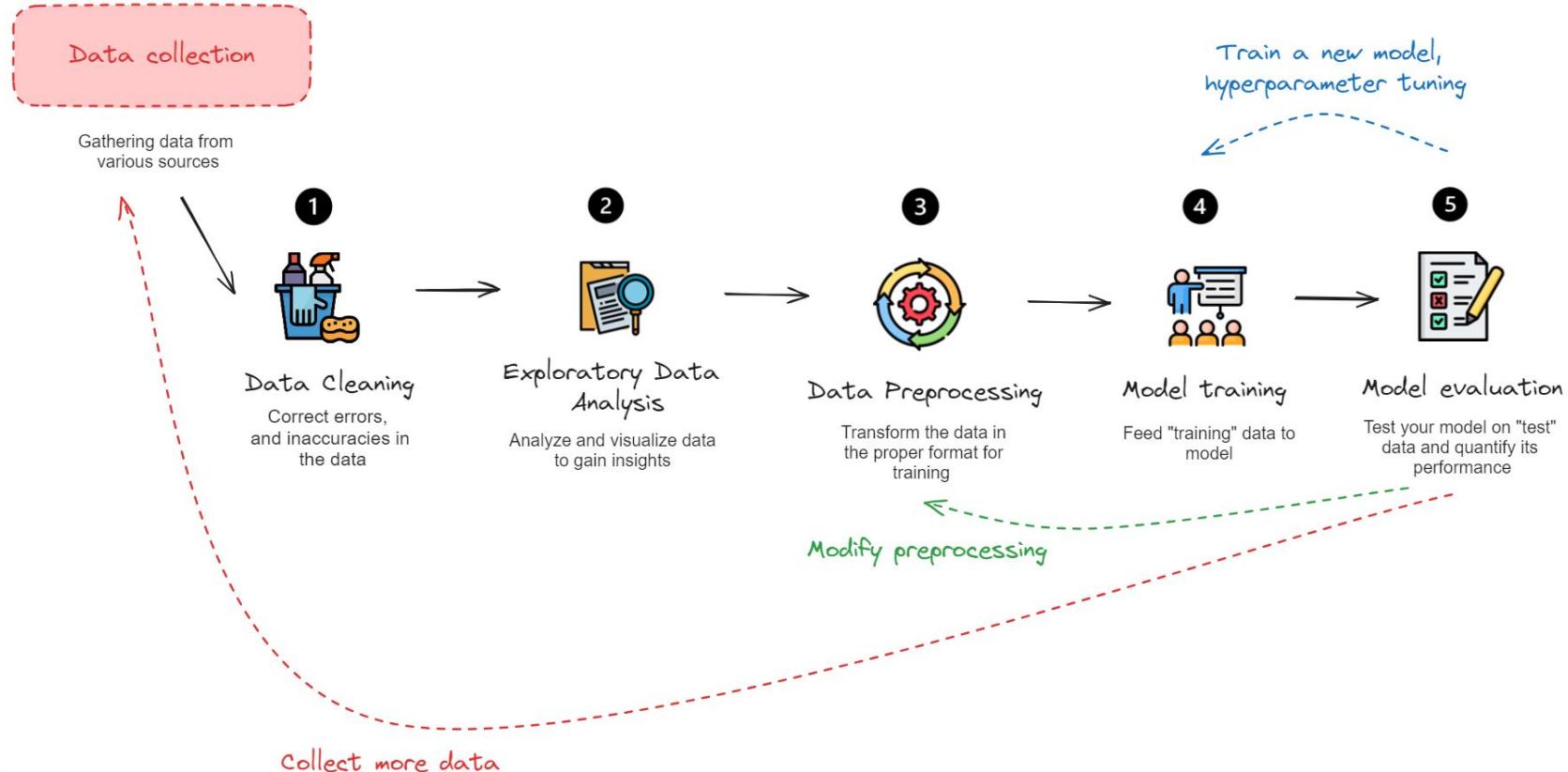
Demo of how to use scikit-learn for model evaluation

https://scikit-learn.org/stable/modules/model_evaluation.html



Improve the performance of a model

Improve the performance of a model

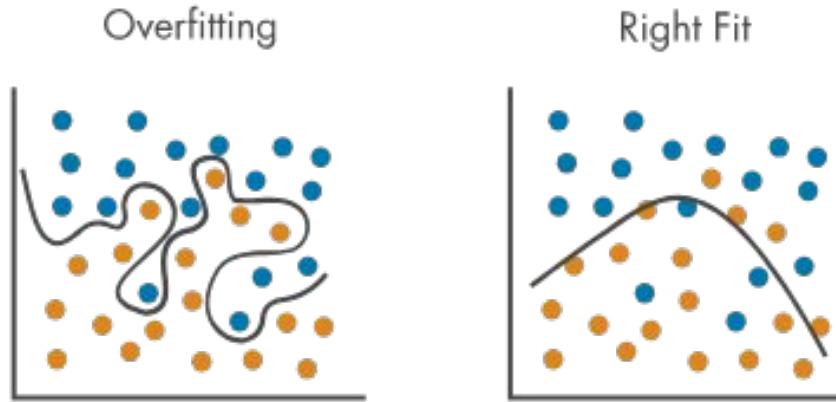


Common reasons for low performances



- **Overfitting:** The model has learned too closely on the training data and can't generalize to new data (test set)

This can lead to unstable predictions

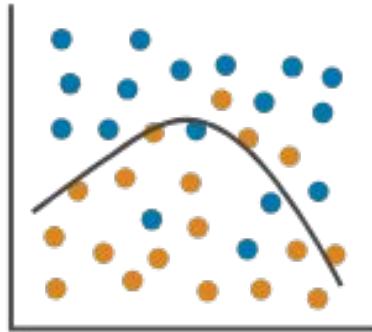


Common reasons for low performances

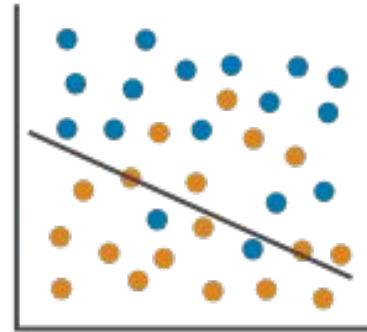


- **Underfitting:** The model hyperparameters you selected for training aren't optimal for your task

Right Fit



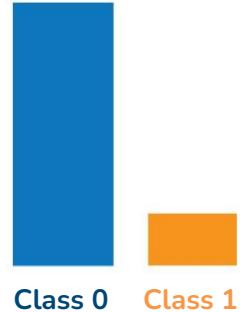
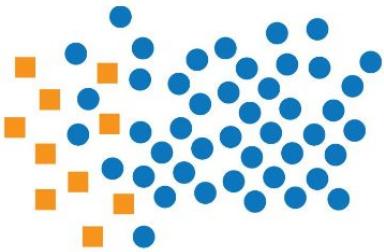
Underfitting



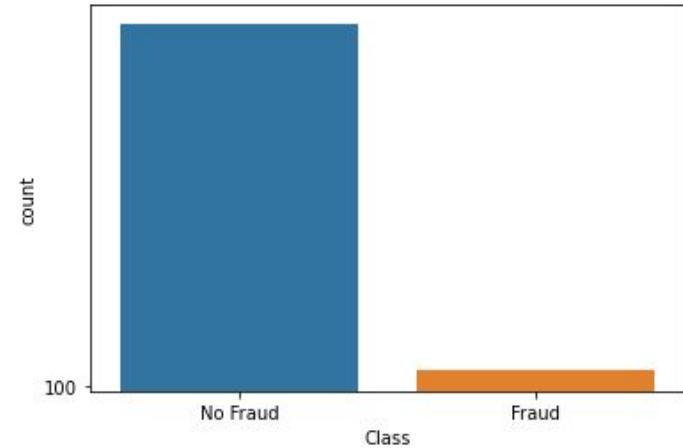
Common reasons for low performances



- **Unbalanced data:** The classes to predict in the target variable aren't in the same proportion.



Example of imbalanced data
with fraud detection



Reduce overfitting



K-Fold cross validation

Evaluate a model on K number of train/test splits (usually K=5)



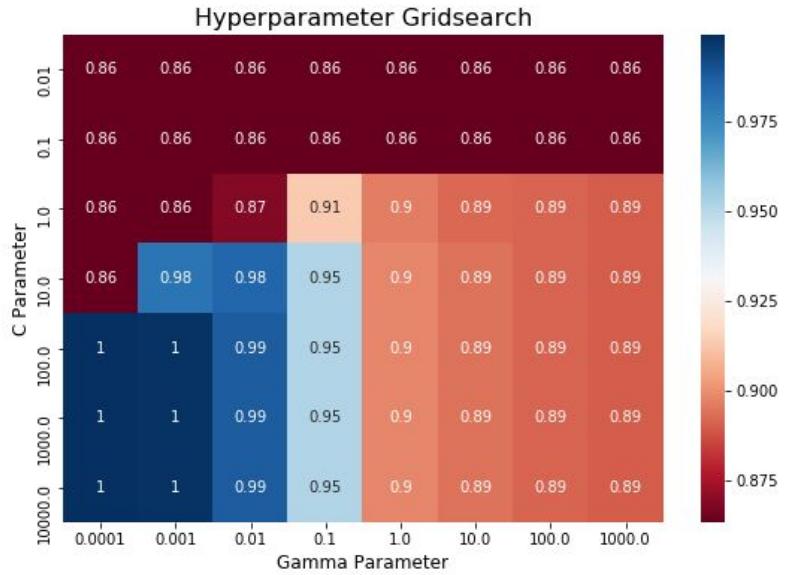
$$\text{Overall score: } \frac{1}{5} \sum_{i=1}^5 Score_i$$

Hyperparameter tuning



Grid Search

Find the optimal hyperparameters of a model by testing every possible combination



Improve the performance on unbalanced data



Most models have a hyperparameter called `class_weight` to help with unbalanced data.

`class_weight` : *dict or 'balanced', default=None*

Weights associated with classes in the form `{class_label: weight}`. If not given, all classes are supposed to have weight one.

The “balanced” mode uses the values of y to automatically adjust weights inversely proportional to class frequencies in the input data as `n_samples / (n_classes * np.bincount(y))`.

Note that these weights will be multiplied with `sample_weight` (passed through the `fit` method) if `sample_weight` is specified.

Improve the performance on unbalanced data



Most models have a hyperparameter called `class_weight` to help with unbalanced data.

`class_weight` : *dict or 'balanced', default=None*

Weights associated with classes in the form `{class_label: weight}`. If not given, all classes are supposed to have weight one.

The “balanced” mode uses the values of y to automatically adjust weights inversely proportional to class frequencies in the input data as `n_samples / (n_classes * np.bincount(y))`.

Note that these weights will be multiplied with `sample_weight` (passed through the `fit` method) if `sample_weight` is specified.

When equal to « **`balanced`** », errors on the minority class will be penalized more heavily during training

Thank you for listening ! 

Do you have any questions ?

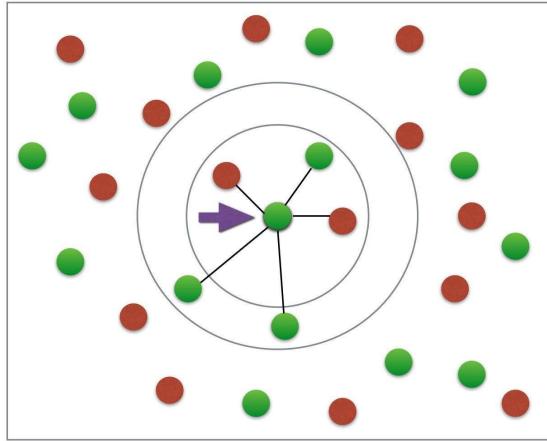


Annexe

How does scaling improve performances ?



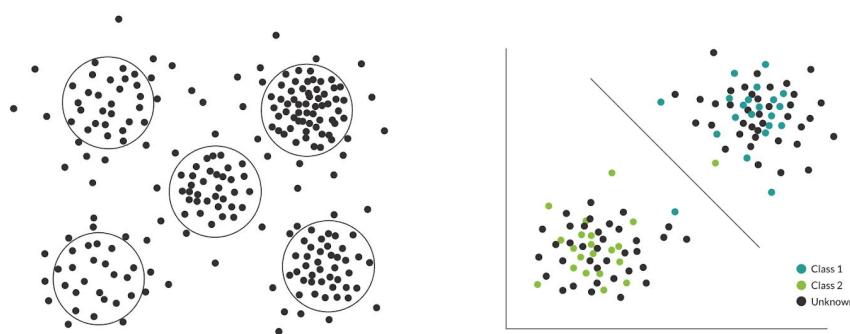
- **Distance-based models**: Models that rely on distances to measure similarity between points and make predictions
- Numerical variables with a large range will contribute more to the decrease in distance ☐ **false sense of importance**



How to choose a model to train ?



- **What is the business objective, the end goal ?**
 - Predict known values or discover unknown patterns ? (*Supervised vs Unsupervised Learning*)
 - Predict a continuous or categorical value ? (*Regression vs Classification*)



How to choose a model to train ?



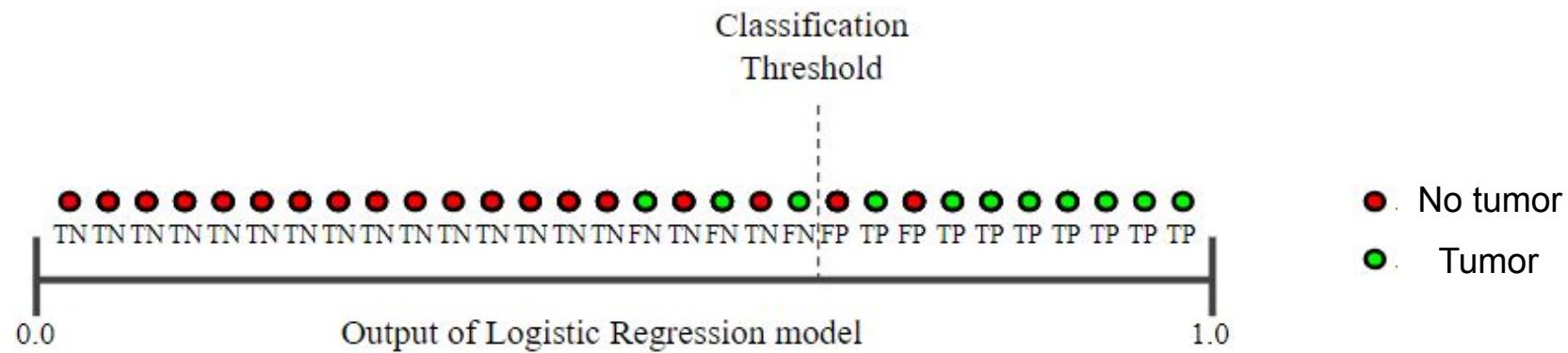
- Size of the training data ?
- Accuracy vs interpretability of the output ?
- Access to powerful computing ressources ?

Evaluation for Classification models



Tradeoff between Precision and Recall

Improving Precision tends to reduce Recall (and vice versa)



Evaluation for Regression models



R2 score:

Measure the **goodness of fit** of a regression model (how well it explained the variations in the target variable Y)

$$\text{R2 Squared} = 1 - \frac{\text{SSr}}{\text{SSm}}$$

SSr = Squared sum error of regression line

SSm = Squared sum error of mean line

Evaluation for Regression models



R2 score:

Measure the **goodness of fit** of a regression model (how well it explained the variations in the target variable Y)

- R2 is **bounded between 0 and 1** (not like MSE and RMSE)
- A higher R2 signals a better model

Summary of evaluation methods



Metric	Type	Characteristics	Bounded	Goal
Accuracy	Classification	Overall performance of the model	Yes [0,1]	Maximise
Confusion Matrix	Classification	Study the performance on all predicted classes	x	x
Precision, Recall and F1 score	Classification	Focus on the performance on the positive class (1)	Yes [0,1]	Maximise
Mean squared error, Root mean squared error	Regression	Measure the distance between predictions and true values	No	Minimise
R2 score	Regression	Evaluate how well the features explain the variations of the target	Yes [0,1]	Maximise

What is Machine Learning (ML) ?



“Machine Learning is the field of study that gives computers the ability to learn without explicitly being programmed.”

Arthur Samuel, Pioneer of AI in the 1950's

What is Machine Learning (ML) ?



Algorithms learn how to make decisions by analyzing large amounts of data, that can be stored in different formats

What is Machine Learning (ML) ?

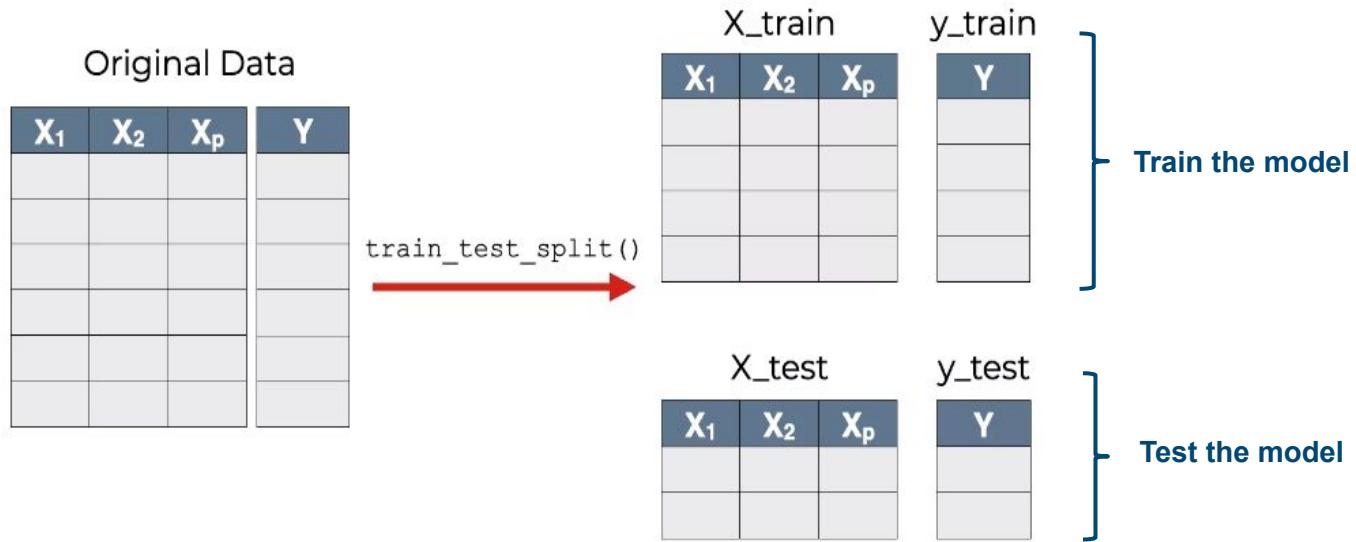


Model training: A model is given many examples and improves its results by adjusting its parameters

Training vs test data



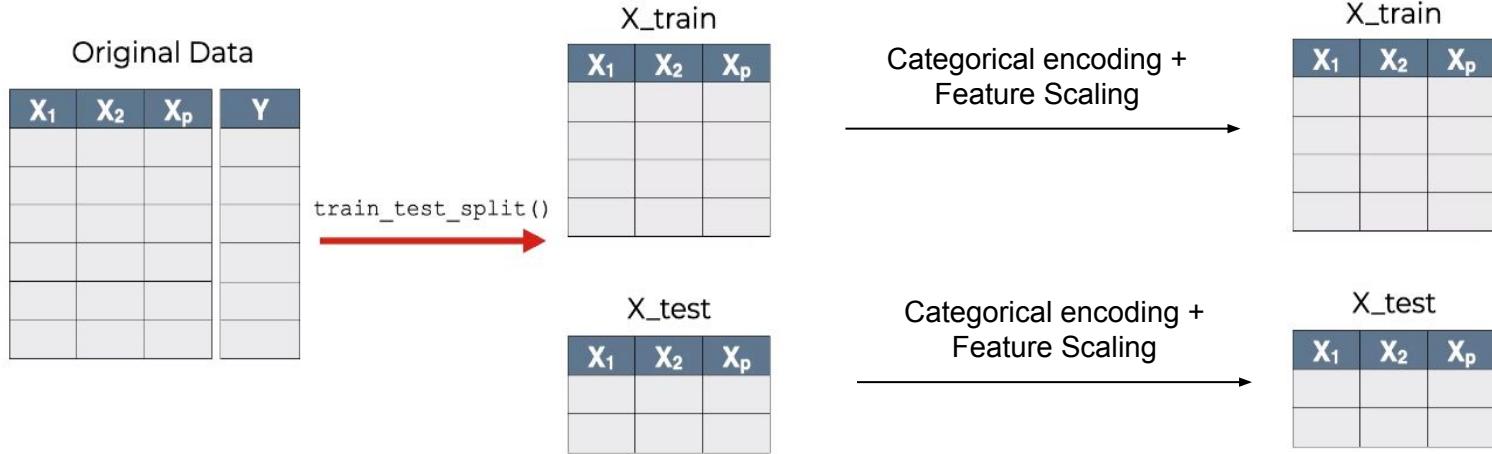
A dataset is split into a **training** and a **test set** to evaluate a model on new/unseen data



Training vs test data



Data preprocessing should be applied separately to each set to prevent information in the training data to “leak” to the test data

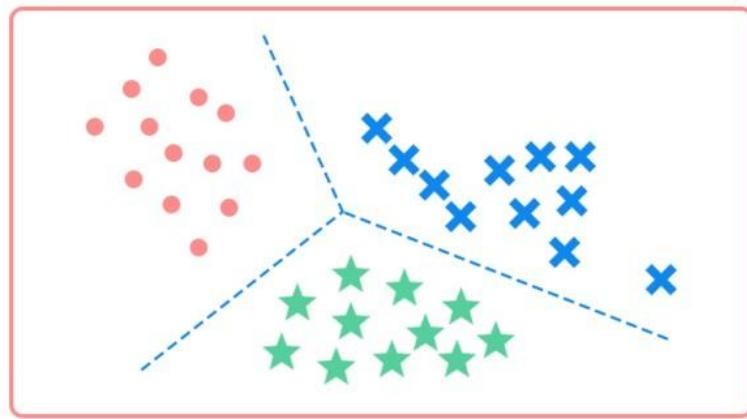


Data preprocessing is only applied to the feature variables (X) !

Main types of Machine Learning

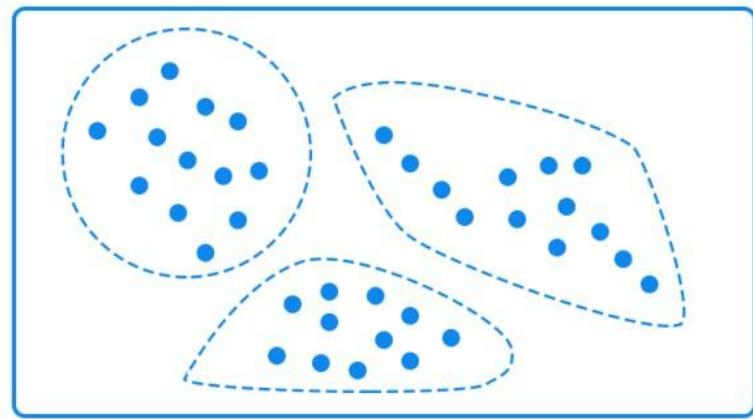


Supervised Learning



- Knowledge of the values to predict
- Model learns using **labeled data** (examples)

Unsupervised Learning

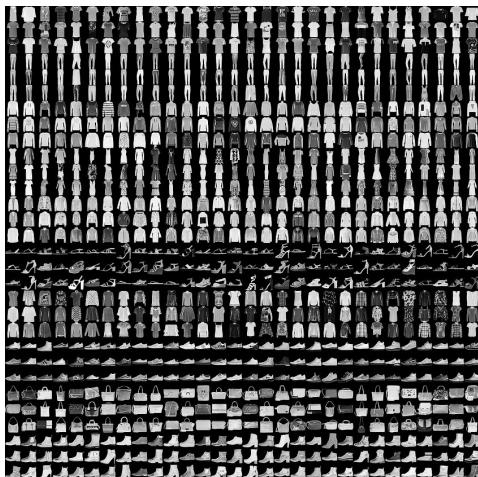


- No prior knowledge on what to predict
- Let the model discover groupings / patterns on its own
- **unlabeled data**

Unsupervised Learning



2. Dimension reduction: Represent data in a lower number of features to simplify visualization and interpretability



Fashion MNIST Dataset



Random Search



Try randomly chosen combinations of hyperparameters to find well performing (but not necessary optimal) ones

