



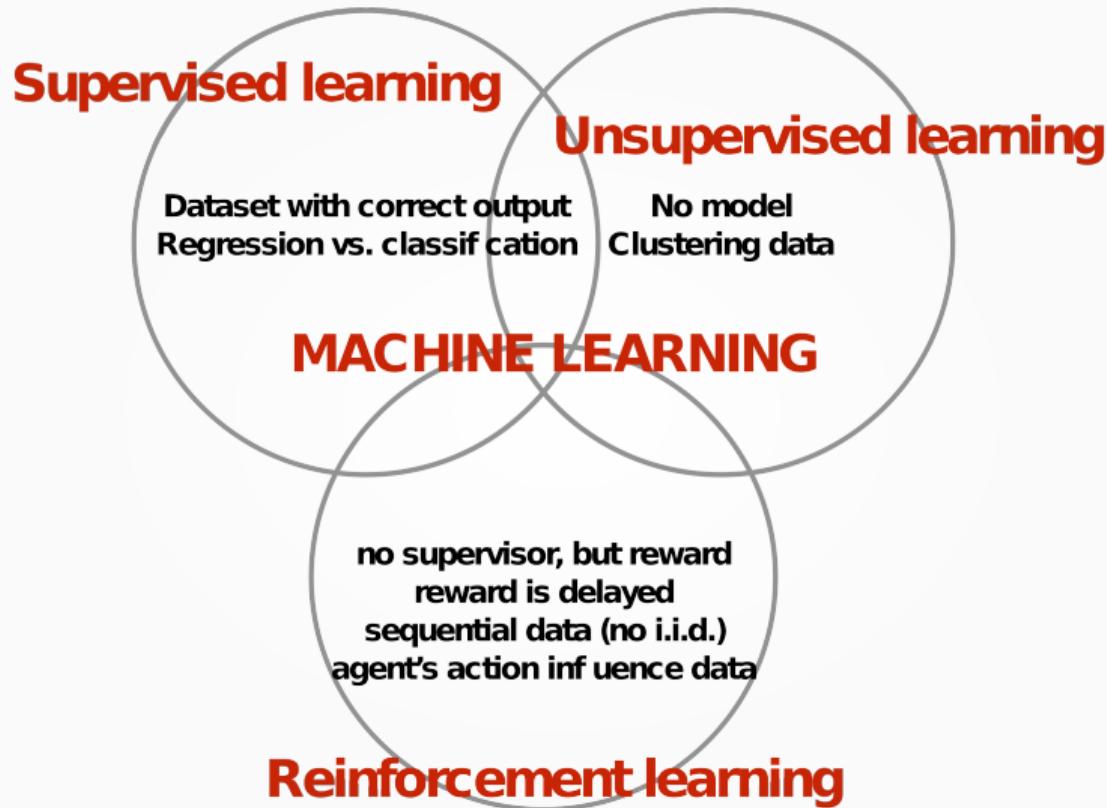
CenTuri Course 2022

Supervised Classification

Stefania Sarno

March 28, 2023

Branches of machine learning

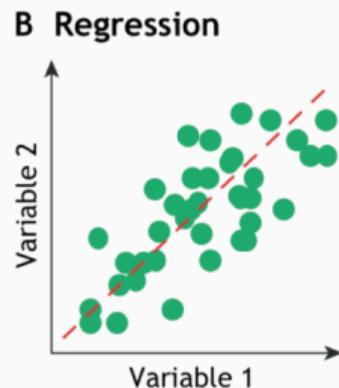
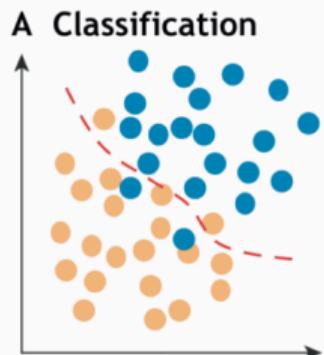


Some examples of machine learning problems

- Classify emails as spam or not Supervised Learning (Classification)
 - Predict the success of a movie Supervised Learning (Regression)
 - Diagnose from list of symptoms Supervised Learning (Classification)
 - Find groups in a social network Unsupervised Learning
 - Drive a car autonomously Reinforcement Learning
 - Defeat world champion of Go Reinforcement Learning

Supervised learning

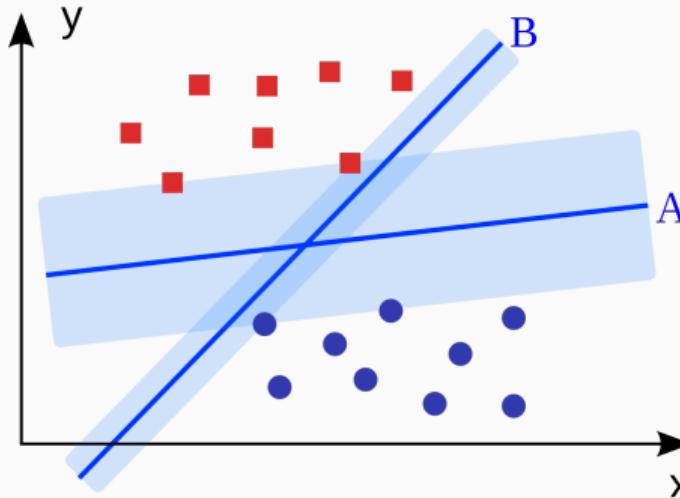
- Dataset with correct output
- Two main tasks: Regression and Classification



Classification

Support vector machine (SVM): intuition

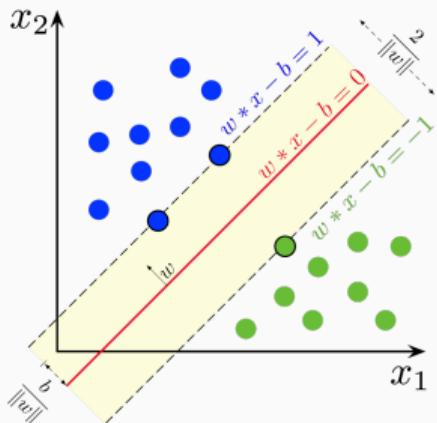
- Multiple hyperplanes could correctly classify binary data
- Can we find the hyperplane that best separates our two classes?



The SVM finds the hyperplane maximizing the margin between our classes

SVM: formal equations

- Our training data are $\{(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(N)}, y^{(N)})\}$ with $\mathbf{x}^i \in \mathbb{R}^n$ and $y^i = \pm 1$



- We look for an hyperplane: $\mathbf{w}^T \mathbf{x} - b = 0$
- When $y^{(i)} = 1$ the points lie on or above the boundary and when $y^{(i)} = -1$ on or below it:

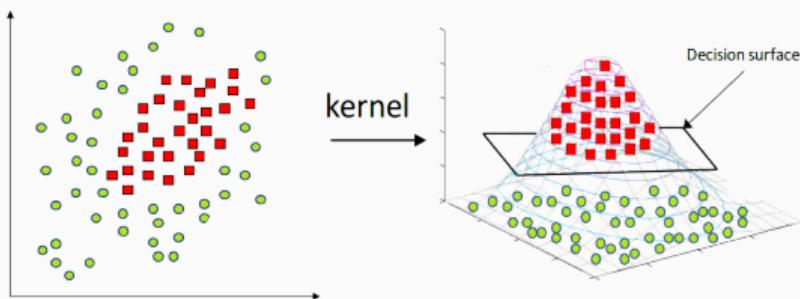
$$\begin{cases} \mathbf{w}^T \mathbf{x}^{(i)} - b \geq 1, & \text{if } y^{(i)} = 1 \\ \mathbf{w}^T \mathbf{x}^{(i)} - b \leq -1, & \text{if } y^{(i)} = -1 \end{cases}$$

- The margin's width is $2/\|\mathbf{w}\|$

We want to minimize $\|\mathbf{w}\|$ under the constraint $y^{(i)} (\mathbf{w}^T \mathbf{x}^{(i)} - b) \geq 1, \forall i$

SVM: the kernel trick

What happen when the data are not linearly separable?

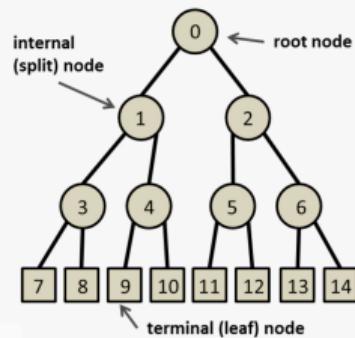


- Proper feature mapping can make non-linear to linear separable
- We do not need the feature space transformation, we only need the kernels
- Popular kernels:
- Polynomial: $k(x, x') = (\mathbf{x}^T \mathbf{x}' + c)^d$
- Gaussian: $k(x, x') = \exp(-\|\mathbf{x} - \mathbf{x}'\|^2 / 2\sigma)$

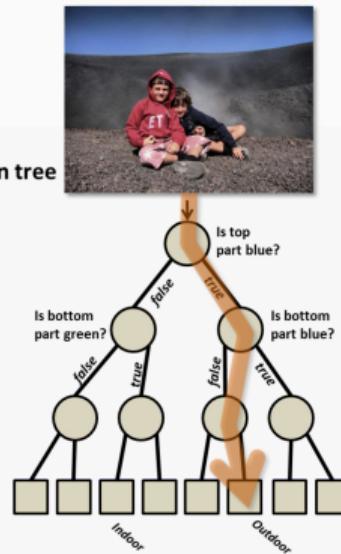
Decision Trees

Image classification example

A general tree structure



A decision tree



[Criminisi et al, 2011]

How do we build a tree

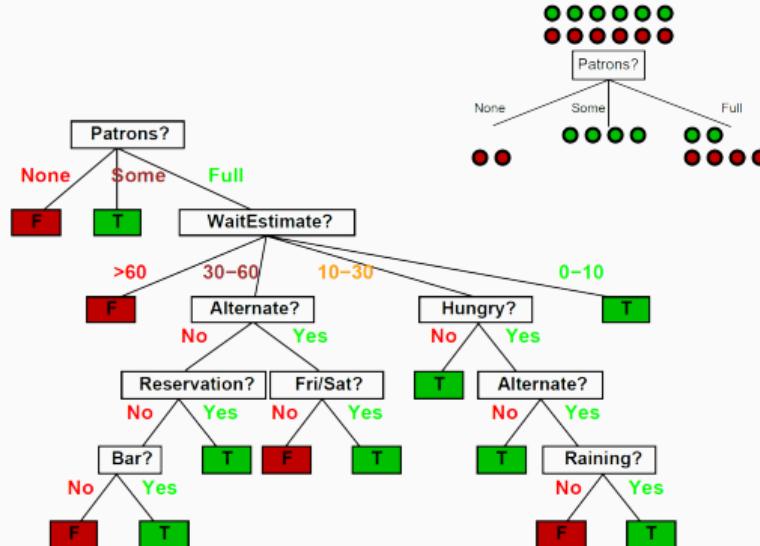
Building a node from data

Example	Input Attributes										Goal WillWait
	Alt	Bar	Fri	Hun	Pat	Price	Rain	Res	Type	Est	
x_1	Yes	No	No	Yes	Some	\$\$\$	No	Yes	French	0–10	$y_1 = \text{Yes}$
x_2	Yes	No	No	Yes	Full	\$	No	No	Thai	30–60	$y_2 = \text{No}$
x_3	No	Yes	No	No	Some	\$	No	No	Burger	0–10	$y_3 = \text{Yes}$
x_4	Yes	No	Yes	Yes	Full	\$	Yes	No	Thai	10–30	$y_4 = \text{Yes}$
x_5	Yes	No	Yes	No	Full	\$\$\$	No	Yes	French	>60	$y_5 = \text{No}$
x_6	No	Yes	No	Yes	Some	\$\$	Yes	Yes	Italian	0–10	$y_6 = \text{Yes}$
x_7	No	Yes	No	No	None	\$	Yes	No	Burger	0–10	$y_7 = \text{No}$
x_8	No	No	No	Yes	Some	\$\$	Yes	Yes	Thai	0–10	$y_8 = \text{Yes}$
x_9	No	Yes	Yes	No	Full	\$	Yes	No	Burger	>60	$y_9 = \text{No}$
x_{10}	Yes	Yes	Yes	Yes	Full	\$\$\$	No	Yes	Italian	10–30	$y_{10} = \text{No}$
x_{11}	No	No	No	No	None	\$	No	No	Thai	0–10	$y_{11} = \text{No}$
x_{12}	Yes	Yes	Yes	Yes	Full	\$	No	No	Burger	30–60	$y_{12} = \text{Yes}$

[AI book of of Stuart Russel and Peter Norvig]

How do we build a tree: A learned tree

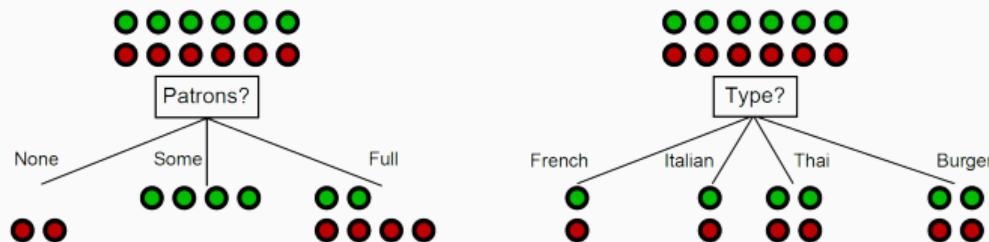
A learned tree



How do we build a tree

Which split is better?

Ideally we want to separate negative and positive examples



[AI book of of Stuart Russel and Peter Norvig]

How do we build a tree

We use the information gain as criterion:

- Shannon Entropy

$$H = - \sum_i p_i \log_2(p_i)$$

- Expected Entropy (for a feature F with K values)

$$EH(F) = - \sum_{i=1}^K \frac{n_i}{N} H_i$$

- Information Gain $I(F)$

$$I(F) = H(Data) - EH(F)$$

How do we build a tree

The patron vs type example



- Entropy data:

$$H(Data) = -2 \cdot [1/2 \cdot \log_2(1/2)] = 1$$

- Entropy left split:

$$\begin{aligned} EH(L) = & -[2/12 \cdot 1 \cdot \log_2(1) + 4/12 \cdot 1 \cdot \log_2(1) + \\ & + 6/12 \cdot (2/6 \cdot \log_2(2/6) + 4/6 \cdot \log_2(4/6))] \approx 0.46 \end{aligned}$$

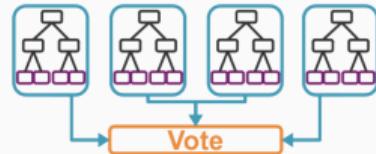
- Entropy right split:

$$EH(R) = -(2/12 + 2/12 + 4/12 + 4/12) \cdot [1/2 \cdot \log_2(1/2)] = 1$$

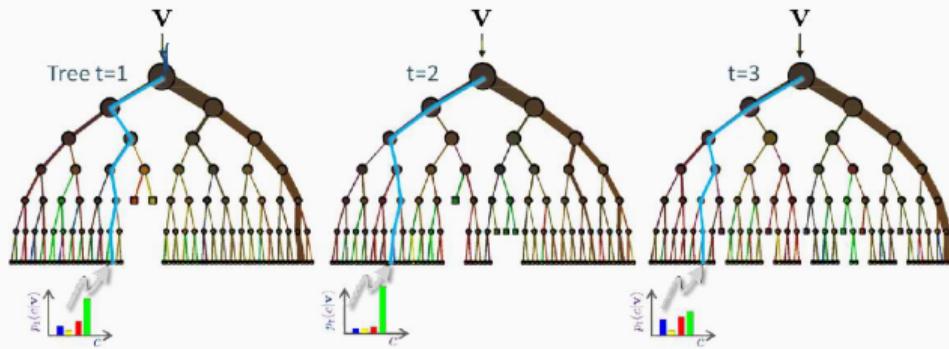
The information gain is clearly higher in the left split!

Random forest: intuition and example

A multitude of decision trees are generated at training time. The output of the random forest is the class selected by the forest (average or majority vote).



Example: We trained a forest of $T = 3$ trees and we want to classify a new point v .



The new point is classified as green!

Random forest algorithm

Given a training set $(\mathbf{x}_1, \dots, \mathbf{x}_N)$ with responses (y_1, \dots, y_N) :

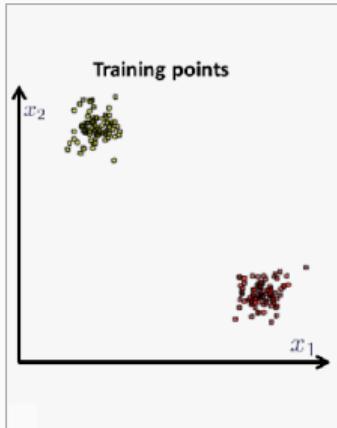
1. For $b = 1$ to B :
 - (a) Draw a bootstrap sample (with replacement) Z of size N from the training data.
 - (b) Grow a random-forest tree T_b to the bootstrapped data, by recursively repeating the following steps for each terminal node of the tree, until the minimum node size n_{\min} is reached.
 - i. Select m variables at random from the p variables.
 - ii. Pick the best variable/split-point among the m .
 - iii. Split the node into two daughter nodes.
2. Output the ensemble of trees $\{T_b\}_1^B$.

To make a prediction at a new point \mathbf{x} :

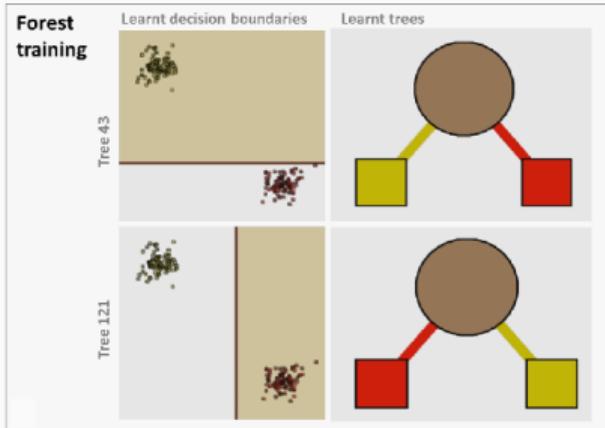
- Let $\hat{C}_b(x)$ be the class prediction of the b -th random-forest tree.
- Then $\hat{C}_{rf}^B(x) = \text{majority vote or average of } \{\hat{C}_b(x)\}_1^B$

[From the book of Hastie, Friedman and Tibshirani]

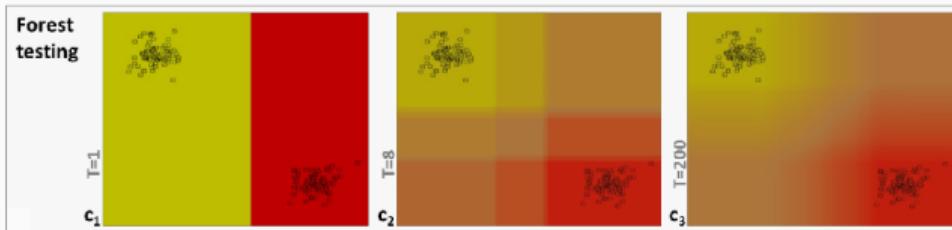
Effects of size



(a)



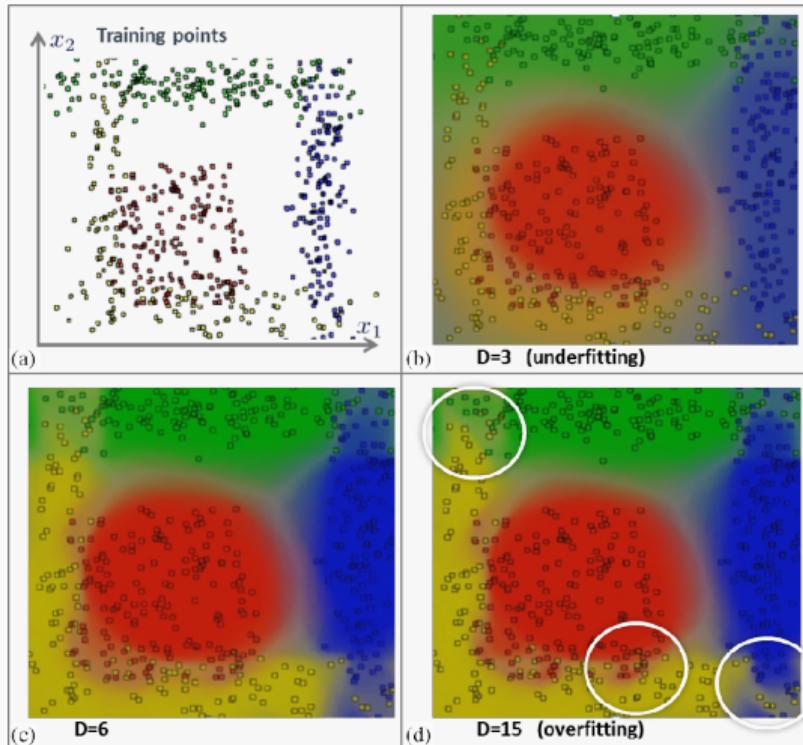
(b)



(c)

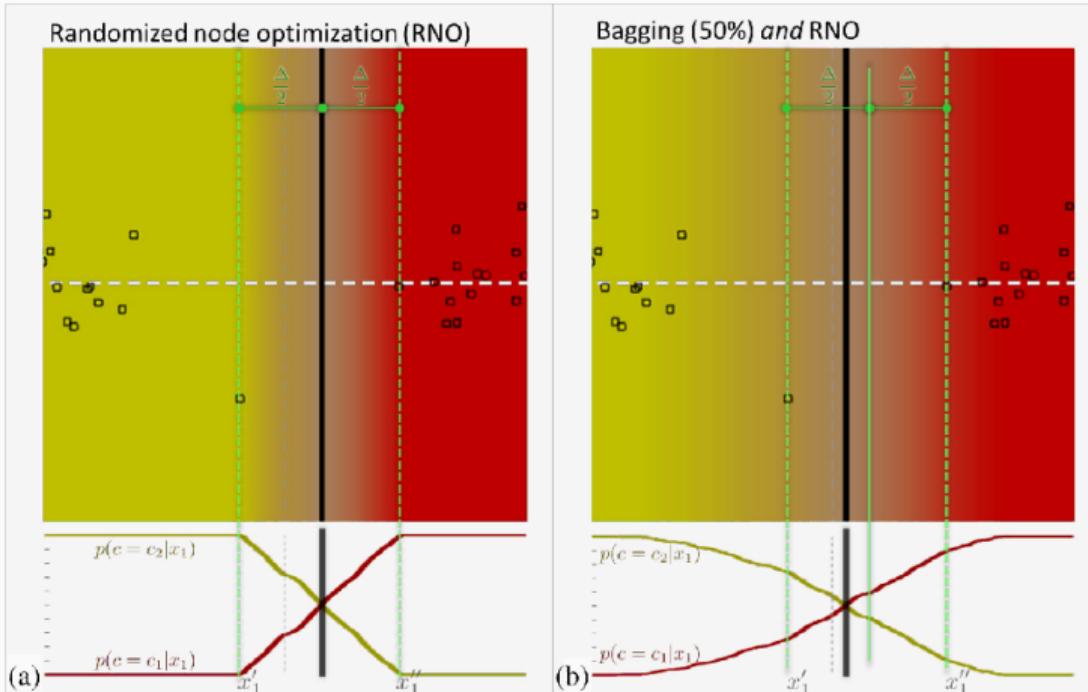
[Criminisi et al, 2011]

Effects of depth



[Criminisi et al, 2011]

Effects of bagging



[Criminisi et al, 2011]

Model evaluation and selection

Measure of performance for a classification model

Confusion matrix

Example

		Predicted Class	
		Positive	Negative
Actual Class	Positive	True Positive (TP)	False Negative (FN) <i>Type II Error</i>
	Negative	False Positive (FP) <i>Type I Error</i>	True Negative (TN)

		Predicted Class	
		Spam	Non-Spam
Actual Class	Spam	TP=45	FN=20
	Non-Spam	FP=5	TN=30

Sensitivity/True positive rate

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

Precision/Positive predicted value

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

Specificity/True negative rate

$$TNR = \frac{TN}{FP + TN}$$

F₁-score

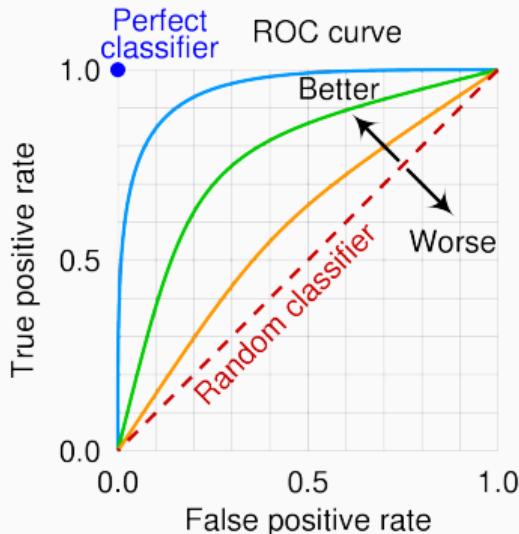
$$F_1 = 2 \frac{PPV \cdot TPR}{PPV + TPR}$$

The receiver operating characteristic (ROC) curve

How to construct the ROC curve:

1. Getting model predictions.
2. Calculate the TPR and FPR.
3. Plot TPR and FPR for every cut-off.

		PREDICTED VALUE	
		Positive	Negative
ACTUAL VALUE	Positive	TP	FN
	Negative	FP	TN



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

$$TNR = \frac{TN}{FP + TN}$$

Maximizing the area under the curve allows to choose the best model

Model evaluation and selection: training/test sets

Model evaluation

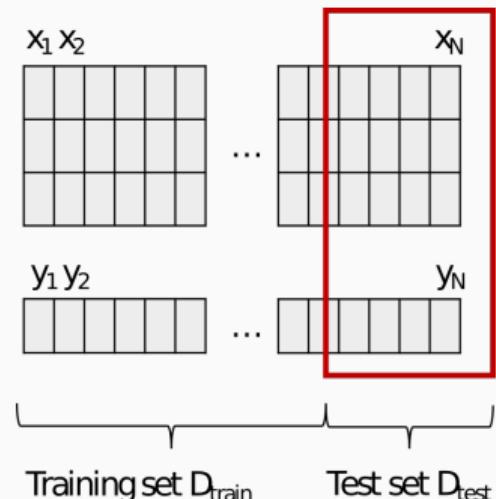
We evaluate a model h by computing the empirical risk ER over the training set D_{train} of size N_{train}

$$ER = \frac{1}{N_{\text{train}}} \sum_{i \in D_{\text{train}}} J(h(x_i), y_i)$$

Model selection

For K models h_1, h_2, \dots, h_K we select the best as:

$$h^* = \operatorname{argmin}_{k=1, \dots, K} \frac{1}{N_{\text{test}}} \sum_{i \in D_{\text{test}}} J(h_k(x_i), y_i)$$



Training/test set error and appropriate fitting

	Underfitting	Just right	Overfitting
Symptoms	<ul style="list-style-type: none">• High training error• Training error close to test error• High bias	<ul style="list-style-type: none">• Training error slightly lower than test error	<ul style="list-style-type: none">• Very low training error• Training error much lower than test error• High variance
Regression illustration			
Classification illustration			
Deep learning illustration			
Possible remedies	<ul style="list-style-type: none">• Complexify model• Add more features• Train longer		<ul style="list-style-type: none">• Perform regularization• Get more data

Parameters vs Hyperparameters

In machine learning models there are **two types of "parameters"**:

1. **Model Parameters:** learned from data automatically via the optimization procedure
2. **Model Hyperparameters:** set manually or tuned and used in processes to help estimate model parameters

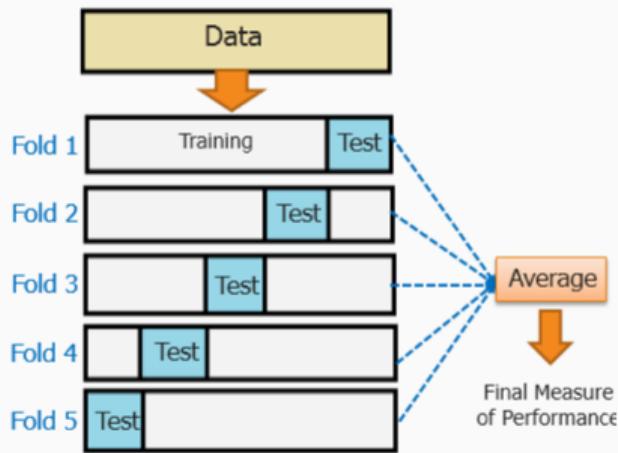
Example Hyperparameters of random forest:

- depth of the trees
- size of the forest

K-fold cross validation

K-fold cross validation:

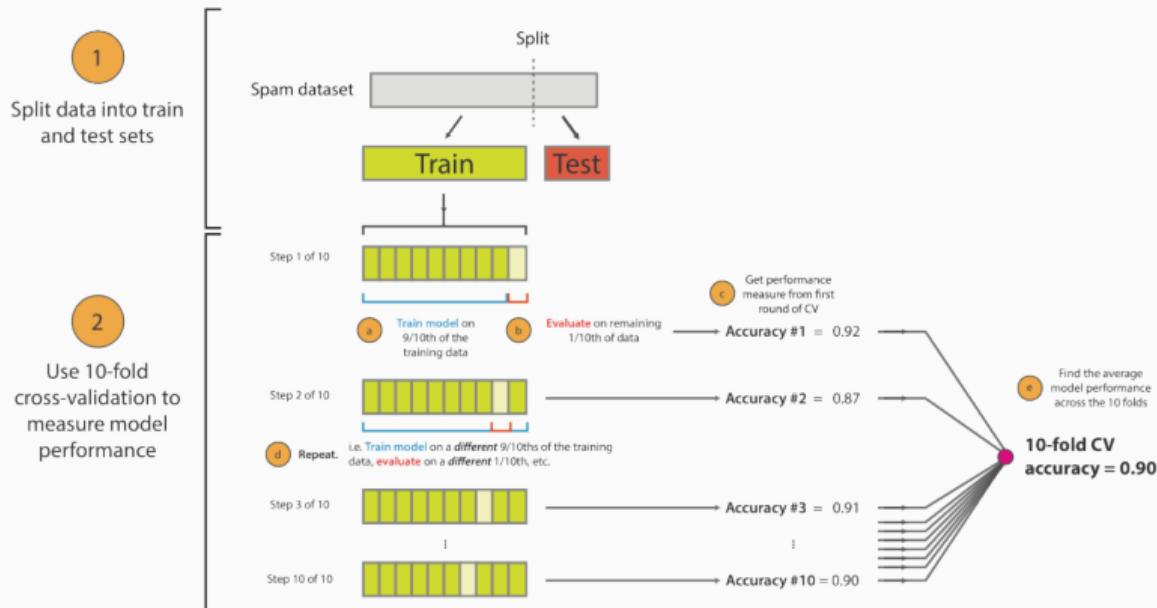
1. Partition the data into K subsets of similar sizes D_1, \dots, D_K
2. For $i = 1, \dots, K$: train on D_1, \dots, D_{i-1}, D_K
3. At each iteration of the loop evaluate on D_i



(A) Way 1

- Apply k-fold validation on entire data set
- Calculate average of each validation score
- Using it as a final score

K-fold cross validation



(B) Way 2

- Split data into train set, test set
- Apply K-fold validation on train set
- Find best estimator from average of each validation score
- Calculate final score on test set with best estimator

Cross validation

Both ways can be correct depending on what the purpose of the procedure is.

- **(A)** uses cross validation for **validation** (or rather, verification), that is, to estimate generalization error.
- **(B)** uses cross validation for **optimizing some hyperparameters** (e.g. model complexity) and then tests the optimized model with the "test" data.

Other useful remarks:

- The term "**cross validation**" refers only to this particular scheme of splitting data (drawing **without replacement** and calculating a pre-determined number of surrogate models so that each case is used for testing exactly once).
- The **K-fold** procedure is called **leave-one-out when $K = N$** , with N indicating the number of data.
- **Bootstrap** can be alternatively used for validation: similar strategy, except that samples are randomly drawn **with replacement** from the dataset.