

Classification: Part II

Statistical Analysis and Document Mining

Spring 2021

Vasilii Feofanov

Université Grenoble Alpes

vasilii.feofanov@univ-grenoble-alpes.fr

1 Classification and Artificial Intelligence

- 1.1 Biological Inspiration
- 1.2 Perceptron Algorithm
- 1.3 Further Developments

2 Classification: Practical Aspects

- 2.1 Performance Estimation
- 2.2 Hyperparameter Tuning
- 2.3 Metrics for the Class Imbalance

- AI is an active field of study these days that explores different ways to develop machines capable to learn and solve problems.

- AI is an active field of study these days that explores different ways to develop machines capable to learn and solve problems.
- A great introduction to the history of the artificial intelligence you can find in the following video.



Figure: https://www.youtube.com/watch?v=8FHBh_0mDsM

- One of the directions in AI is to mimic human intelligence: how human brain perceives information and make decisions.

- One of the directions in AI is to mimic human intelligence: how human brain perceives information and make decisions.
- Promising deep neural networks were originally inspired by biological neurons, electrically active cells.

- One of the directions in AI is to mimic human intelligence: how human brain perceives information and make decisions.
- Promising deep neural networks were originally inspired by biological neurons, electrically active cells.
- The brain has complex networks of connected to each other neurons that can receive, store or propagate information.

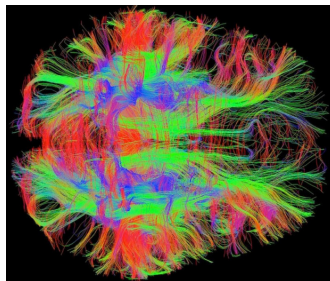
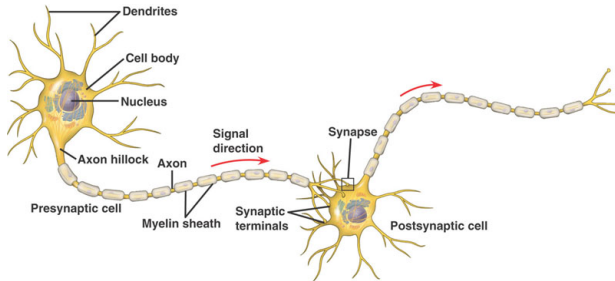
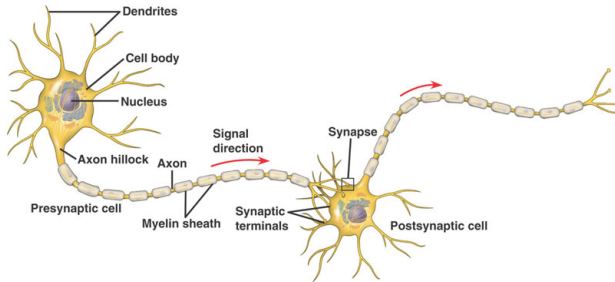


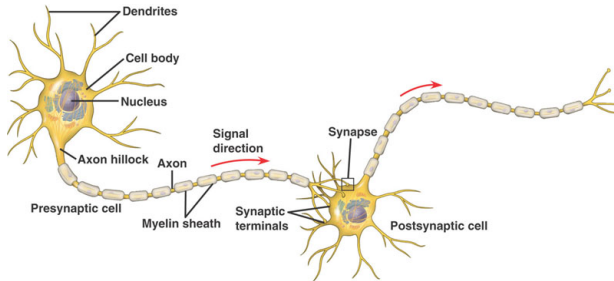
Figure: A map of nerve fibers (axons) in the human brain.



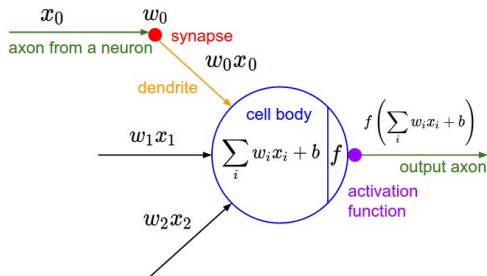
- One neuron sends a signal to another one through an axon to a synapse of a dendrite.



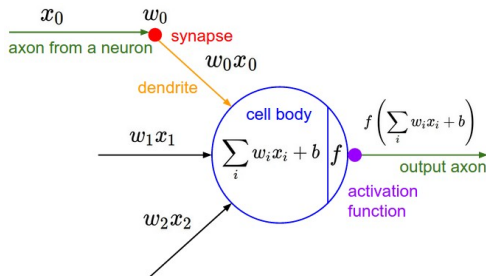
- One neuron sends a signal to another one through an axon to a synapse of a dendrite.
- The synapse processes the information and can regularize its amplitude and its frequency.



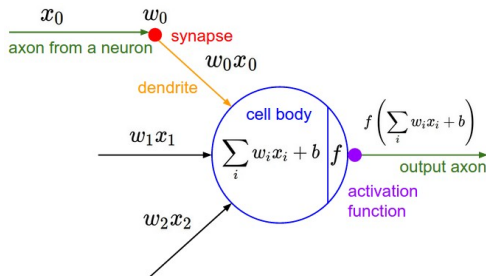
- One neuron sends a signal to another one through an axon to a synapse of a dendrite.
- The synapse processes the information and can regularize its amplitude and its frequency.
- The cell body receives the processed information via the dendrite. The neuron decides then whether to fire it further.



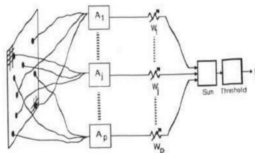
- In 1943, a simplified mathematical model of the neuron was proposed.



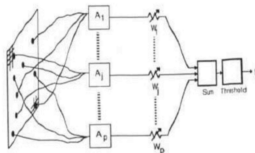
- In 1943, a simplified mathematical model of the neuron was proposed.
- Neurons $0, 1, 2 \dots$ send information $x_0, x_1, x_2 \dots$. The synapse determines the strength (weight) of information.



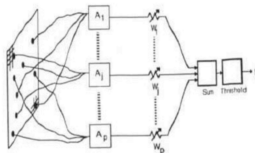
- In 1943, a simplified mathematical model of the neuron was proposed.
- Neurons $0, 1, 2 \dots$ send information $x_0, x_1, x_2 \dots$. The synapse determines the strength (weight) of information.
- The cell body sums up all information and decides via activation function whether to propagate it further.



- Using the neuron model, Rosenblatt implemented an algorithm in hardware for image recognition of geometric figures.



- Using the neuron model, Rosenblatt implemented an algorithm in hardware for image recognition of geometric figures.
- An array of 20x20 photocells are randomly connected to form *features* (input neurons).



- Using the neuron model, Rosenblatt implemented an algorithm in hardware for image recognition of geometric figures.
- An array of 20x20 photocells are randomly connected to form *features* (input neurons).
- The weights represent the resistance of electric motors and are updated during learning by purely electric process.

- **Summation:** From an image we extract features $\mathbf{x} = (x_1, \dots, x_d)$ and linearly combine them inside the decision neuron.

$$h_{\mathbf{w}}(\mathbf{x}) := \langle \mathbf{w}, \mathbf{x} \rangle + w_0 = \sum_{j=1}^d w_j x_j + w_0.$$

¹We consider the binary classification framework, i.e. $\mathcal{Y} = \{-1, +1\}$.

- **Summation:** From an image we extract features $\mathbf{x} = (x_1, \dots, x_d)$ and linearly combine them inside the decision neuron.

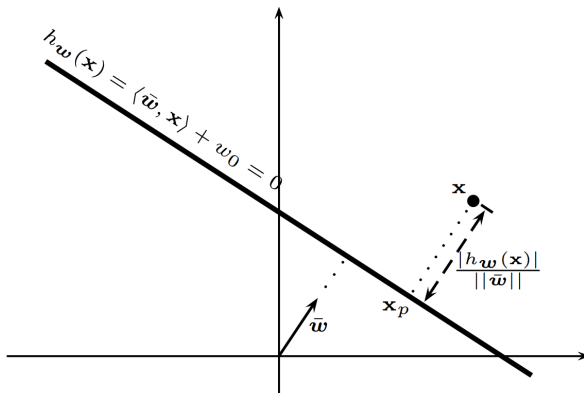
$$h_{\mathbf{w}}(\mathbf{x}) := \langle \mathbf{w}, \mathbf{x} \rangle + w_0 = \sum_{j=1}^d w_j x_j + w_0.$$

- **Activation:** We compare the sign of the true label y and the output $h_{\mathbf{w}}(\mathbf{x})$. If there is no mistake ($y h_{\mathbf{w}}(\mathbf{x}) \geq 0$), the activation is 1, otherwise, it is -1 ¹.

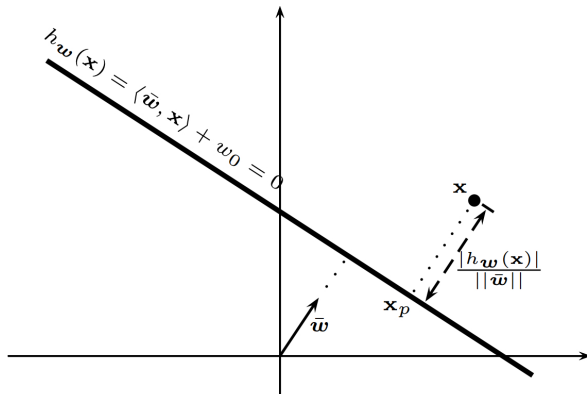
$$a(\mathbf{x}) = \begin{cases} +1, & \text{if } y h_{\mathbf{w}}(\mathbf{x}) \geq 0; \\ -1, & \text{if } y h_{\mathbf{w}}(\mathbf{x}) < 0. \end{cases}$$

¹We consider the binary classification framework, i.e. $\mathcal{Y} = \{-1, +1\}$.

We want to find parameters $(w_0, \bar{\mathbf{w}})$ such that the distance between the misclassified examples and the decision boundary is minimised.



By this, we push the boundary away from correctly classified examples that will have high margins. In this case, the margin $\frac{|h_{\mathbf{w}}(\mathbf{x})|}{\|\mathbf{w}\|}$ represents the confidence score.



- We would like to minimise the *perceptron error*:

$$\hat{L}(\mathbf{w}, w_0) = - \sum_{i' \in \mathcal{I}} y_{i'} (\langle \mathbf{w}, \mathbf{x}_{i'} \rangle + w_0).$$

- We would like to minimise the *perceptron error*:

$$\hat{L}(\mathbf{w}, w_0) = - \sum_{i' \in \mathcal{I}} y_{i'} (\langle \mathbf{w}, \mathbf{x}_{i'} \rangle + w_0).$$

- If we take derivatives with respect to the weights:

$$\frac{\partial \hat{L}(\mathbf{w})}{\partial w_0} = - \sum_{i' \in \mathcal{I}} y_{i'},$$

$$\nabla \hat{L}(\mathbf{w}) = - \sum_{i' \in \mathcal{I}} y_{i'} \mathbf{x}_{i'}.$$

- We would like to minimise the *perceptron error*:

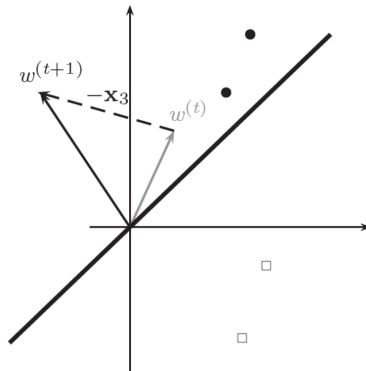
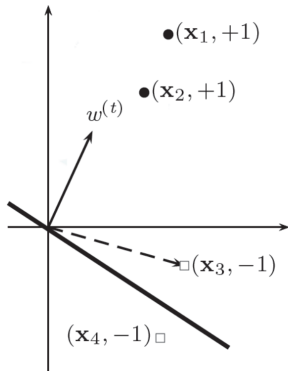
$$\hat{L}(\mathbf{w}, w_0) = - \sum_{i' \in \mathcal{I}} y_{i'} (\langle \mathbf{w}, \mathbf{x}_{i'} \rangle + w_0).$$

- If we take derivatives with respect to the weights:

$$\frac{\partial \hat{L}(\mathbf{w})}{\partial w_0} = - \sum_{i' \in \mathcal{I}} y_{i'},$$
$$\nabla \hat{L}(\mathbf{w}) = - \sum_{i' \in \mathcal{I}} y_{i'} \mathbf{x}_{i'}.$$

- Minimisation of the perceptron error is done by stochastic gradient descent algorithm:

$$\text{if } y(\langle \mathbf{w}, \mathbf{x} \rangle + w_0) < 0, \text{ then } \begin{pmatrix} w_0 \\ \mathbf{w} \end{pmatrix} \leftarrow \begin{pmatrix} w_0 \\ \mathbf{w} \end{pmatrix} + \eta \begin{pmatrix} y \\ y\mathbf{x} \end{pmatrix}.$$



Algorithm Perceptron

Input: Training set $S = \{\mathbf{x}_i, y_i\}_{i=1}^n$.

Max number of iterations T .

Initialisation: Weights $\mathbf{w}^{(0)} \leftarrow \mathbf{0}, w_0 \leftarrow 0$;

Counter $t \leftarrow 0$, learning rate $\eta > 0$.

repeat

Choose randomly an example $(\mathbf{x}^{(t)}, y^{(t)}) \in S$

if $y \cdot (\langle \mathbf{w}^{(t)}, \mathbf{x}^{(t)} \rangle + w_0) < 0$ **then**

$$w_0^{(t+1)} \leftarrow w_0^{(t)} + \eta y^{(t)}$$

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} + \eta y^{(t)} \mathbf{x}^{(t)}$$

end if

$t \leftarrow t + 1$

until $t > T$

Output: $(w_0^{(T)}, \mathbf{w}^{(T)})$.

- What is the time complexity during training phase?

- What is the time complexity during training phase?
- What is the time complexity to predict a label for new x ?

- What is the time complexity during training phase?
- What is the time complexity to predict a label for new \mathbf{x} ?
- How to tune learning rate η ?

- What is the time complexity during training phase?
- What is the time complexity to predict a label for new \mathbf{x} ?
- How to tune learning rate η ?
- Does the algorithm converge to an exact solution?

- What is the time complexity during training phase?
- What is the time complexity to predict a label for new \mathbf{x} ?
- How to tune learning rate η ?
- Does the algorithm converge to an exact solution?
- Under what condition the algorithm does not converge?

- Large success of the perceptron were followed by public disappointment and a sharp decline of AI funding.
- Minsky and Papert showed that the perceptron is not able to approximate the XOR operation, a non-linear problem.

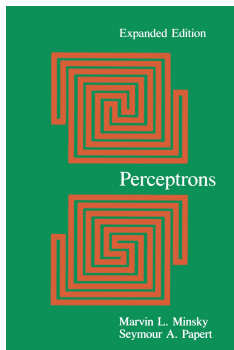


Figure: (Minsky & Papert, 1969)



Figure: Minsky and Papert in 1971.

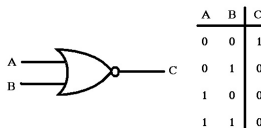
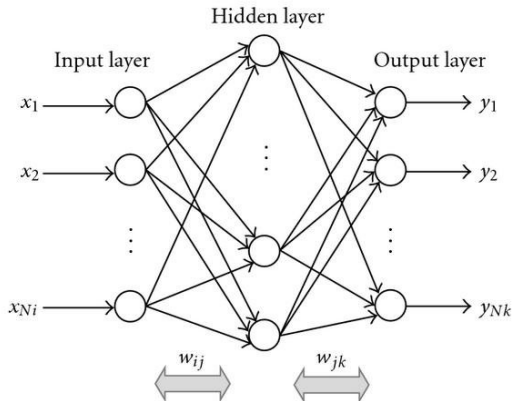
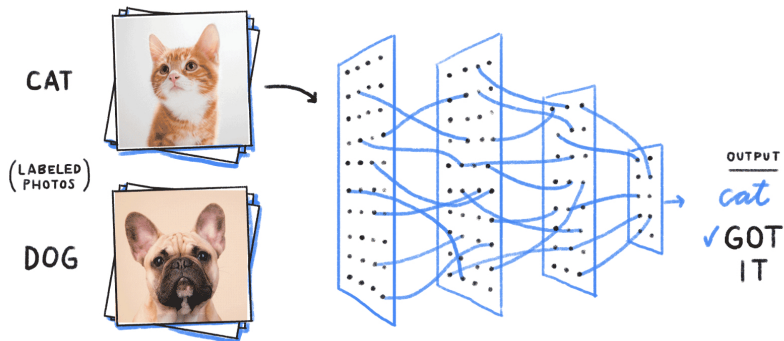


Figure: The XOR operation.

- In the 1980's, the perceptron returned as *multi-layer*: a hidden layer of neurons connects the input and the output layers.
- The weights are updated layer by layer propagating errors back through the whole network (*backpropagation*).



- Nowadays, with large computational resources, the neural networks has become *deep* with a large number of hidden layers of different functionality.
- The goal is not only prediction itself, but also learning new appropriate data representation for better prediction.



1 Classification and Artificial Intelligence

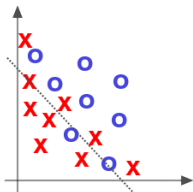
- 1.1 Biological Inspiration
- 1.2 Perceptron Algorithm
- 1.3 Further Developments

2 Classification: Practical Aspects

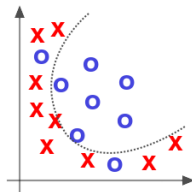
- 2.1 Performance Estimation
- 2.2 Hyperparameter Tuning
- 2.3 Metrics for the Class Imbalance

- How to select a learning model with the best accuracy score² on *unseen examples*?
- How to tune hyperparameters of an algorithm?
- Are there metrics that take into account imbalance in classes? Importance of one class over another one?

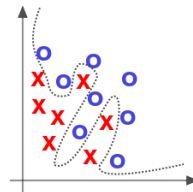
²The accuracy score is the proportion of correctly predicted examples.



(a) Underfitting

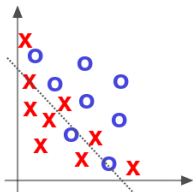


(b) Appropriate

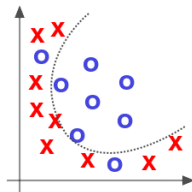


(c) Overfitting

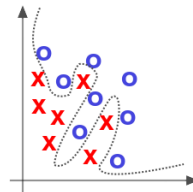
- **Underfitting:** The classifier does not fit data and has a large error value.



(a) Underfitting

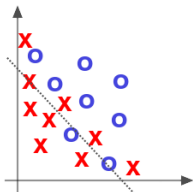


(b) Appropriate

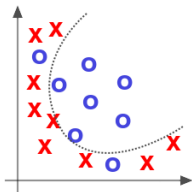


(c) Overfitting

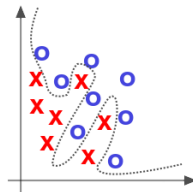
- **Underfitting:** The classifier does not fit data and has a large error value.
- **Overfitting:** Being perfect on training data, the classifier has poor generalization making lots of mistakes on new examples.



(a) Underfitting



(b) Appropriate



(c) Overfitting

- **Underfitting:** The classifier does not fit data and has a large error value.
- **Overfitting:** Being perfect on training data, the classifier has poor generalization making lots of mistakes on new examples.
- **What we want:** The classifier does not take into account noise in the training data and approximates well the true boundary.

Given the data sample $S = \{\mathbf{x}_i, y_i\}_{i=1}^n$, are the following estimates of the misclassification error good (in terms of bias/variance)?

- Compute error on the same examples the model was trained?

Given the data sample $S = \{\mathbf{x}_i, y_i\}_{i=1}^n$, are the following estimates of the misclassification error good (in terms of bias/variance)?

- Compute error on the same examples the model was trained?
 - No! Among various models we will choose the one that is optimistically biased. This is exactly *overfitting*.
- Split the data into two parts: one is used for training, error is evaluated on the examples of the second one?

Given the data sample $S = \{\mathbf{x}_i, y_i\}_{i=1}^n$, are the following estimates of the misclassification error good (in terms of bias/variance)?

- Compute error on the same examples the model was trained?
 - No! Among various models we will choose the one that is optimistically biased. This is exactly *overfitting*.
- Split the data into two parts: one is used for training, error is evaluated on the examples of the second one?
 - This is appropriate. The method is widely used when the number of examples n is large. However, when n is small, we cherish every example we have due to the risk of *underfitting*.
- Data is divided by k folds, one fold is for evaluation and $k - 1$ rest are for training; the errors are averaged over k rounds?

Given the data sample $S = \{\mathbf{x}_i, y_i\}_{i=1}^n$, are the following estimates of the misclassification error good (in terms of bias/variance)?

- Compute error on the same examples the model was trained?
 - No! Among various models we will choose the one that is optimistically biased. This is exactly *overfitting*.
- Split the data into two parts: one is used for training, error is evaluated on the examples of the second one?
 - This is appropriate. The method is widely used when the number of examples n is large. However, when n is small, we cherish every example we have due to the risk of *underfitting*.
- Data is divided by k folds, one fold is for evaluation and $k - 1$ rest are for training; the errors are averaged over k rounds?
 - *Cross-validation* (CV) is the most popular approach. It reduces not only the bias, but also the data variability used for training.

- The smaller k , the less number of examples we use for training. If k is large, we train too many classifiers.
- For medium/large sample sizes, usual choices of k are 5 or 10.
- If the sample size n is really small, the popular paradigm is to take $k = n$, a.k.a. the *leave-one-out* (LOO) cross-validation.

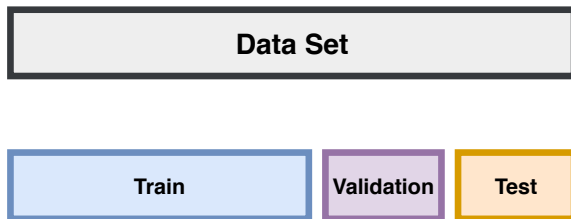
Data Set							
Round 1	Ex. 1	Ex. 2	Ex. 3	Ex. 4	Ex. 5	Ex. 6	Ex. 7
Round 2	Ex. 1	Ex. 2	Ex. 3	Ex. 4	Ex. 5	Ex. 6	Ex. 7
Round 3	Ex. 1	Ex. 2	Ex. 3	Ex. 4	Ex. 5	Ex. 6	Ex. 7
Round 4	Ex. 1	Ex. 2	Ex. 3	Ex. 4	Ex. 5	Ex. 6	Ex. 7
Round 5	Ex. 1	Ex. 2	Ex. 3	Ex. 4	Ex. 5	Ex. 6	Ex. 7
Round 6	Ex. 1	Ex. 2	Ex. 3	Ex. 4	Ex. 5	Ex. 6	Ex. 7
Round 7	Ex. 1	Ex. 2	Ex. 3	Ex. 4	Ex. 5	Ex. 6	Ex. 7

Figure: LOO for a data set with 7 examples (blue = train, orange = test).

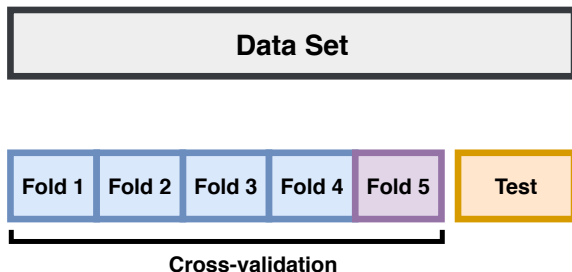
- The CV can be also used to find the best hyperparameter (k for knn, η for perceptron, etc.).
- We define a grid of values (e.g. $\eta \in \{0.01, 0.1, 1\}$) and compute the CV score for each parameter value.
- Then, we select a model with the best CV-score.

- The CV can be also used to find the best hyperparameter (k for knn, η for perceptron, etc.).
- We define a grid of values (e.g. $\eta \in \{0.01, 0.1, 1\}$) and compute the CV score for each parameter value.
- Then, we select a model with the best CV-score.

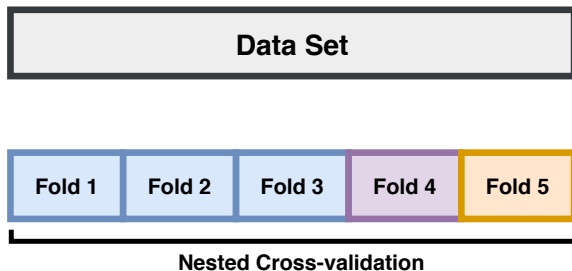
However, there is still a chance of overfitting. If you wanna be on the safe side, one of the following strategy can be applied.



- *Strategy 1:* Split into three parts: learn models on Train, choose the best hyperparameter on Validation, estimate the performance on Test.



- *Strategy 2*: Split into two parts: The first part is used to perform CV. By the cv-score we choose the best hyperparameter, which is used then to learn a final model using all examples from the first part. Finally, the performance is estimated on Test.



- *Strategy 3:* We average the effect of choosing the best hyperparameter: at each round, $k - 2$ folds are used for learning, orange fold is for estimating performance of the hyperparameter chosen on the purple fold. The final cv-score is the average performance over k orange folds.

- To see the difference between the predicted and the true outputs in more detail, the confusion matrix can be used.
- Its (i, j) -entry is the number of examples that were assigned to a class j given the true label is i .

		Prediction	
		-1	+1
Real	-1	65	5
	+1	7	123

- To see the difference between the predicted and the true outputs in more detail, the confusion matrix can be used.
- Its (i, j) -entry is the number of examples that were assigned to a class j given the true label is i .
- In the binary classification, the entries are named as follows: True Negative (TN), False Positive (FP), False Negative (FN), True Positive (TP).

	Prediction		
	-1	+1	
Real	-1	TN	FP
	+1	FN	TP

- The class imbalance is an often situation in real applications. For instance, the number of patients without a disease (-1) can be much larger than with it (+1).
- In this case, the accuracy score might be not the best choice to evaluate the prediction quality.

- The class imbalance is an often situation in real applications. For instance, the number of patients without a disease (-1) can be much larger than with it (+1).
- In this case, the accuracy score might be not the best choice to evaluate the prediction quality.
- Example: Classifier 2 is better in terms of accuracy. However, we would like to choose Classifier 1, because it predicts less frequently "no disease" given a diseased patient (small FN).

		Prediction	
		-1	+1
Real	-1	150	20
	+1	3	27

(a) Classifier 1.

		Prediction	
		-1	+1
Real	-1	170	0
	+1	20	10

(b) Classifier 2.

Figure: Confusion matrices of two classifiers.

- If we want to minimize FN, the *sensitivity* (a.k.a. recall) score can be used to choose the best model:

$$\text{SEN} = \frac{\text{TP}}{\text{TP} + \text{FN}}.$$

- If we want to minimize FP, the *specificity* score can be used:

$$\text{SPE} = \frac{\text{TN}}{\text{TN} + \text{FP}}.$$

- Finally, with the *balanced accuracy* we minimize both FP and FN discarding the effect of the class distribution:

$$\text{BACC} = \frac{\text{SPE} + \text{SEN}}{2}.$$

To see a simple example of classification in R, please go to Chamilo:

Lectures → Classification Part II → R script with a classification example