



FACULTÉ DES SCIENCES APPLIQUÉES  
ELEN062-1 : INTRODUCTION TO MACHINE LEARNING

---

## Classification Algorithms

---

*Teachers :*

Louis Wehenkel  
Pierre Geurts

*Group :*

LAMBERMONT Romain  
LOUIS Arthur

23 octobre 2022

# 1 Decision trees

## 1.1 Impact of the tree's depth on the decision boundary

- a) The lower the depth of the tree, the lower the decision boundary will be. Indeed, when the depth is equal to 1, the boundary is a simple line and as the depth increases, the boundary becomes more complex, fitting the data better. This phenomenon can be seen in the following figure 1.

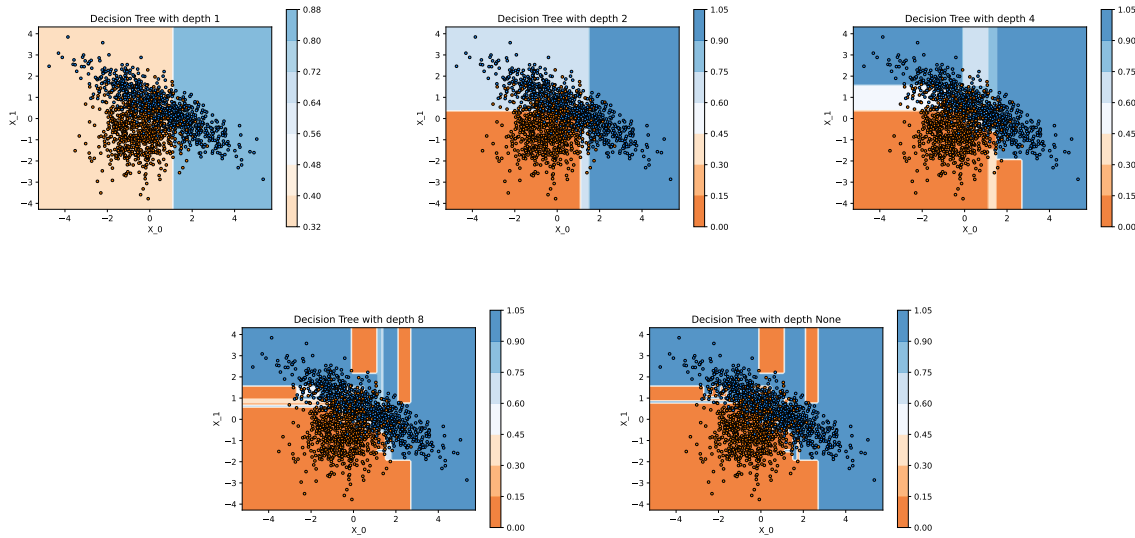


FIGURE 1 – Decision trees with different depths

- b) On one hand, when the depth is 1, the separation between areas is a simple line. The model is too complex to be fit in this way so the depth 1 is clearly underfitting the problem. For the depth 2, the model is more complex and fits the data better but it is still underfitting the problem. On the other hand, when the depth is 8, the model is too complex and the decision boundary is overfitting the problem. When the hyperparameter `depth` is set to `None`, the nodes are expanded until all leaves are pure or until all leaves contain less than `min_samples_split` samples, which is the minimum number of samples required to split an internal node. The best depth is 4 because it is the one that fits the problem the best. As demonstrated in the next point, the depth 4 is the one that gives the best accuracy.
- c) As previously said, when the depth becomes large, the model overfits the data, leading to a better confidence on the training set but a lower confidence on the test set.

## 1.2 Tests accuracy report

We report the results obtained for the different depths in the following table. The models were tested on 5 iterations of the dataset. We can clearly see that the best depth is 4, as it gives the best accuracy while keeping the standard deviation low. When the depth gets past

4, the model overfits the data leading to decreased accuracy and increased standard deviation (the model is fitting the noise in the data).

max_depth	Mean	Standard Deviation
1	0.68	0.031
2	0.80	0.021
4	0.83	0.014
8	0.81	0.016
None	0.77	0.02

TABLE 1 – Decision trees results

## 2 K-nearest neighbors

### 2.1 Impact of the number of neighbors on the decision boundary

a) This is the graphics with the 6 parameters.

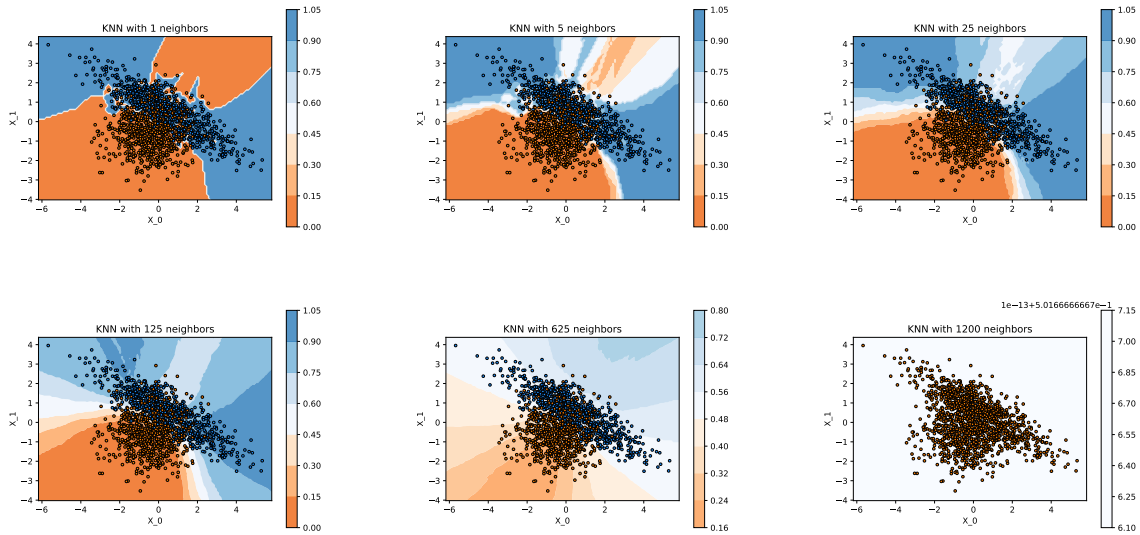


FIGURE 2 – K-nearest neighbors with different  $n\_neighbors$  values

b) The different boundaries are represented in figure 2. The lower the number of neighbors, the more the decision boundary is affected by the noise in the data (it starts to overfit). When the number of neighbors increases, e.g for 5, 25 and 125, there's less overfitting and the decision boundary is more stable. When the number of neighbors gets even greater, e.g 625, the model starts to underfit the data, as it is not complex enough to fit the data. And in the worst case, when the number of neighbors is equal to the number of samples, the model simply makes the same prediction for every sample, which is the mean of the training set, giving the worst accuracy possible.

## 2.2 Tests accuracy report

Reporting to the table, we can see that the best number of neighbors is 125. Indeed, it gives good accuracy while keeping the standard deviation low. When the number of neighbors is too low, the model overfits the data and when it is too high, the model underfits the data.

n_neighbors	Mean	Standard Deviation
1	0.78	0.024
5	0.81	0.013
25	0.84	0.013
125	0.85	0.021
625	0.81	0.010
1200	0.48	0.014

TABLE 2 – Decision trees results

## 3 Quadratic/Linear discriminant analysis

### 3.1 Decision boundaries

In the 2 models, the decision boundaries is the curve where :

$$P(y = 1|x) = P(y = 0|x) \quad (1)$$

$$\frac{f_0(x)\pi_0}{\sum_{l=0}^1 f_l(x)\pi_l} = \frac{f_1(x)\pi_1}{\sum_{l=0}^1 f_l(x)\pi_l} \quad (2)$$

$$(3)$$

With :

$$f_k(x) = \frac{1}{2\pi\sqrt{|\Sigma_k|}} \exp\left(-\frac{1}{2}(x - \mu_k)^T \Sigma_k^{-1} (x - \mu_k)\right) \quad (4)$$

Dividing both sides by  $P(y = 1|x)$  and taking the logarithm, we get the following equation :

$$\log \frac{P(y = 0|x)}{P(y = 1|x)} = \log \frac{f_0(x)\pi_0}{f_1(x)\pi_1} = 0 \quad (5)$$

$$\log \frac{f_0(x)}{f_1(x)} + \log \frac{\pi_0}{\pi_1} = 0 \quad (6)$$

The first term of the line 6 change between the Linear and Quadratic form. In the linear case, we have  $|\Sigma_0| = |\Sigma_1|$  so in the division, we can simplify the 2 discriminant and the quadratic form disappears. In the quadratic case, we have  $|\Sigma_0| \neq |\Sigma_1|$  so we have to keep the quadratic form.

Linear case :

$$\log \frac{\pi_0}{\pi_1} - \frac{1}{2}(\mu_0 + \mu_1)^T \Sigma^{-1} (\mu_0 - \mu_1) + x^T \Sigma^{-1} (\mu_0 - \mu_1) = 0 \quad (7)$$

Quadratic case :

$$\log \frac{\pi_0}{\pi_1} - \frac{1}{2}(\mu_0 + \mu_1)^T \Sigma^{-1}(\mu_0 - \mu_1) + x^T \Sigma^{-1}(\mu_0 - \mu_1) + \log \frac{|\Sigma_1|}{|\Sigma_0|} = 0 \quad (8)$$

### 3.2 Quadratic/Linear discriminant analysis

As we can see, LDA is a special case of QDA where  $\Sigma_0 = \Sigma_1 = \Sigma$ . This is why the decision boundary is linear in LDA and quadratic in QDA. It is also the reason why the curve of iso probability is parallel line in LDA. We can also see that QDA is much more confident in general than LDA.

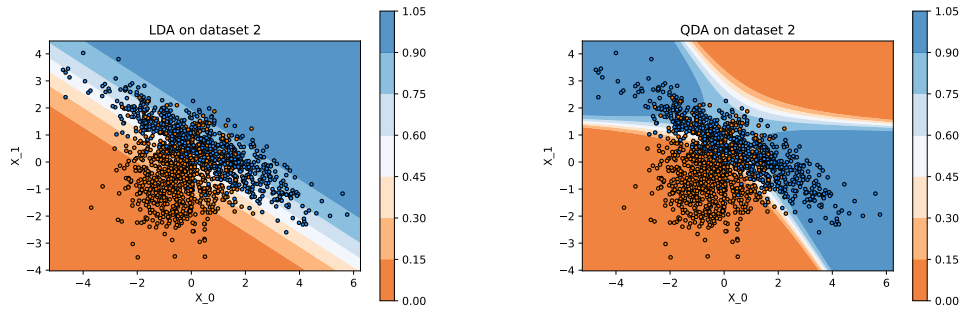


FIGURE 3 – Quadratic/Linear discriminant analysis

### 3.3 Tests accuracy report

We can see that QDA is better than LDA in general. But the results are really close for dataset 1 because this dataset is created using circular gaussian distribution. So the covariance matrix is diagonal and the decision boundary is linear.

LDA	QDA	Mean	Standard Deviation
LDA1		0.85	0.068
QDA1		0.86	0.057
LDA2		0.72	0.053
QDA2		0.76	0.062

TABLE 3 – Decision trees results

## 4 Comparison

### 4.1 Tuning of hyperparameters

To tune our classifiers hyperparameters, we can use the cross-validation method. This method consists in splitting the dataset into  $k$  folds. Then, we train the model on  $k - 1$  folds

and test it on the remaining fold. We repeat this process  $k$  times, each time using a different fold as the test set. Finally, we average the results to get the final accuracy. We used different values of the hyperparameters and we chose the one that gave the best accuracy. We implement this method in the `tune.py` file using the `GridSearchCV` class from the `sklearn.model_selection` module. We used the 5-fold cross-validation method and the accuracy as the scoring parameter.

## 4.2 Results of our implementation

After finding the best hyperparameters, we trained the model using it and we tested it on 5 iterations of the dataset. Here are the results we obtained :

	Mean	Standard Deviation
KNN_1_tuned (neighbors = 25)	0.91	0.012
DT_1_tuned (depth = 8)	0.89	0.018
KNN_2_tuned (neighbors = 125)	0.83	0.019
DT_2_tuned (depth = 4)	0.82	0.02

TABLE 4 – Tuned hyperparameters results

## 4.3 Comparison with the QDA/LDA results

A quick review of the 2 tables 3 and 4 show that an Algorithms with good parameters can be better than QDA/LDA even if they are not as complex as them. This point is the proof that it is really important to tune the hyperparameters of the algorithms to get the best results and it can't be done randomly.