# Training Piscine datascience - 3

## TO DO

*Summary:*   *Today, you will see some TODO.*

*Version: 1.00*

# Contents

# Chapter I

# General rules

- Your project must be realized in a virtual machine.

- Your virtual machine must have all the necessary software to complete your project. These softwares must be configured and installed.

- You can choose the operating system to use for your virtual machine.

- You must be able to use your virtual machine from a cluster computer.

- You must use a shared folder between your virtual machine and your host machine.

- During your evaluations you will use this folder to share with your repository.

- Your functions should not quit unexpectedly (segmentation fault, bus error, double free, etc) apart from undefined behaviors. If this happens, your project will be considered non functional and will receive a `0` during the evaluation.

- We encourage you to create test programs for your project even though this work **won't have to be submitted and won't be graded**. It will give you a chance to easily test your work and your peers' work. You will find those tests especially useful during your defence. Indeed, during defence, you are free to use your tests and/or the tests of the peer you are evaluating.

- Submit your work to your assigned git repository. Only the work in the git repository will be graded. If Deepthought is assigned to grade your work, it will be done after your peer-evaluations. If an error happens in any section of your work during Deepthought's grading, the evaluation will stop.

# Chapter II

# Specific instructions of the day

Module de data Scientist, ils utilisent souvent comme techno le python, Jupyter Notebook ...

A vous de trouver les outils qui vous conviennent ce module est en language libre.

Le role du data Scientist est de predire "l'avenir" avec des modele d'aprentissage automatique sur des donnes passees, il doit etre force de proposition pour expliquer l'interet possitif a la mise en place de ses modeles, creer des outils d'aide a la prise de decision
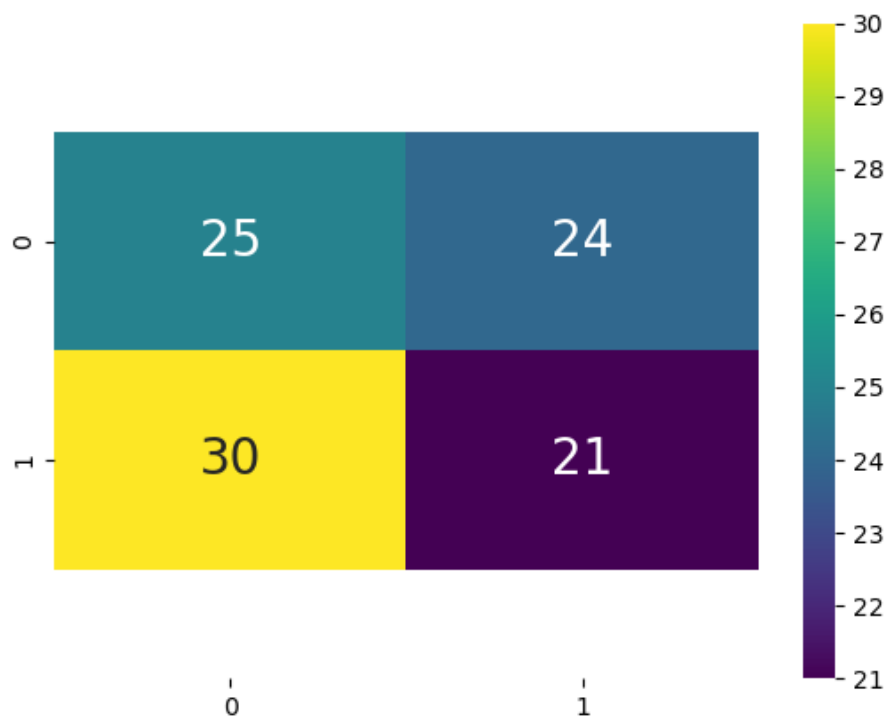
# Chapter III

# Exercise 00

| | Exercise 00 |
|---|---|
| | Exercice 00 : TO DO |
| Turn-in directory : *ex00/* | |
| Files to turn in : `TO DO.*` | |
| Allowed functions : `None` | |

- faire un prog qui prend en arg 1er fichier avec les predictions, 2nd fichier les truth, puis print

- faire une Confusion matrix

```
       precision   recall f1-score     total

  0       0.51      0.45     0.48        55
  1       0.41      0.47     0.44        45

accuracy                     0.46       100

[[25 24]
 [30 21]]
```

⚠️ Si cette exo est faux, on ne passe pas a la suite, c'est finito

ℹ️ soyez sur d'avoir bien compris cette exo car cela va etre verifie

# Chapter IV

# Exercise  01

| | Exercise  01 |
|---|---|
| | Exercice  01 : |
| Turn-in directory : *ex*01/ | |
| Files to turn in : `TO DO.*` | |
| Allowed functions : `None` | |

- Faire une Heatmap pour voir le Correlation Coefficient entre les donees

# Chapter V

# Exercise  02

| | Exercise  02 |
|---|---|
| | Exercice  02 : |
| Turn-in directory : *ex*02/ | |
| Files to turn in : `TO DO.*` | |
| Allowed functions : `None` | |

- additioner les variances pour voir combien sont utiles dans 90% du modele

# Chapter VI

# Exercise  03

| | Exercise  03 |
|---|---|
| | Exercice  03 : Feature Selection |
| Turn-in directory : *ex03/* | |
| Files to turn in : `TO DO.*` | |
| Allowed functions : `None` | |

les donnee semble trop Multicollinearity, vous aller devoir faire un modele de Detecting Multicollinearity

Il en existe un certain nombre de modele de detection d'importance de variables (Lasso, Backward Elimination, Step Forward Selection, ...)

Mais ici vous aller devoir utiliser le Variance Inflation Factor (VIF)

- Le VIF de vos features doit etre inferieur a 5

- pas de hard coding (on change les data dans l'eval)

# Chapter VII

# Exercise 04

|  | Exercise 04 |
|---|---|
|  | Exercice 04 : TO DO |
| Turn-in directory : *ex04/* | |
| Files to turn in : `TO DO.*` | |
| Allowed functions : `None` | |

- faire un modele de Decision Tree Classifier ou Random Forest Classifier (l'aisser de choix pour que le stud cherche ?)

- faire un prog qui prend en arg le fichier de train en 1er arg et le fichier de test en 2nd arg et qui ecrit un fichier avec les reponses dedans (comme dans dslr)

- afficher l'arbre dans un graphique

## Decision tree trained on all Knights features

X[27] <= 0.489
gini = 0.471
samples = 381
value = [145, 236]

True / False

X[20] <= 0.345
gini = 0.148
samples = 249
value = [20, 229]

X[13] <= 0.013
gini = 0.1
samples = 132
value = [125, 7]

X[13] <= 0.053
gini = 0.073
samples = 236
value = [9, 227]

X[8] <= 0.231
gini = 0.26
samples = 13
value = [11, 2]

gini = 0.0
samples = 3
value = [0, 3]

X[6] <= 0.206
gini = 0.06
samples = 129
value = [125, 4]

X[27] <= 0.467
gini = 0.036
samples = 219
value = [4, 215]

X[19] <= 0.098
gini = 0.415
samples = 17
value = [5, 12]

gini = 0.0
samples = 2
value = [0, 2]

gini = 0.0
samples = 11
value = [11, 0]

X[28] <= 0.282
gini = 0.49
samples = 7
value = [4, 3]

X[0] <= 0.159
gini = 0.016
samples = 122
value = [121, 1]

X[19] <= 0.014
gini = 0.019
samples = 213
value = [2, 211]

X[28] <= 0.226
gini = 0.444
samples = 6
value = [2, 4]

X[17] <= 0.23
gini = 0.494
samples = 9
value = [5, 4]

gini = 0.0
samples = 8
value = [0, 8]

X[12] <= 0.113
gini = 0.375
samples = 4
value = [1, 3]

gini = 0.0
samples = 3
value = [3, 0]

gini = 0.0
samples = 1
value = [0, 1]

gini = 0.0
samples = 121
value = [121, 0]

X[27] <= 0.326
gini = 0.245
samples = 7
value = [1, 6]

X[21] <= 0.553
gini = 0.01
samples = 206
value = [1, 205]

gini = 0.0
samples = 2
value = [2, 0]

gini = 0.0
samples = 4
value = [0, 4]

gini = 0.0
samples = 3
value = [0, 3]

X[21] <= 0.227
gini = 0.278
samples = 6
value = [5, 1]

gini = 0.0
samples = 3
value = [0, 3]

gini = 0.0
samples = 1
value = [1, 0]

gini = 0.0
samples = 6
value = [0, 6]

gini = 0.0
samples = 1
value = [1, 0]

gini = 0.0
samples = 189
value = [0, 189]

X[21] <= 0.581
gini = 0.111
samples = 17
value = [1, 16]

gini = 0.0
samples = 1
value = [0, 1]

gini = 0.0
samples = 5
value = [5, 0]

gini = 0.0
samples = 1
value = [1, 0]

gini = 0.0
samples = 16
value = [0, 16]
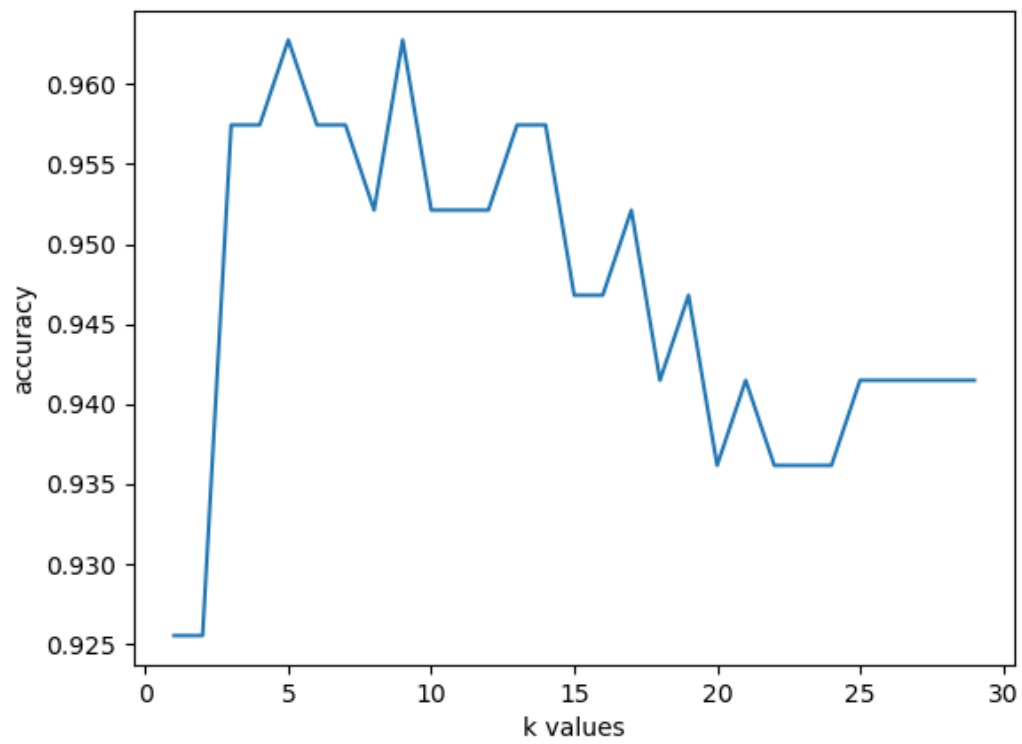
# Chapter VIII

# Exercise  05

| | Exercise  05 |
|---|---|
| | Exercice  05 : TO DO |
| Turn-in directory : *ex05/* | |
| Files to turn in : `TO DO.*` | |
| Allowed functions : `None` | |

- faire un KNN qui calcul le % de precision en fonction du nb de k value

- afficher le graph

- faire un prog qui prend en arg le fichier de train en 1er arg et le fichier de test en 2nd arg et qui ecrit un fichier avec les reponses dedans (comme dans dslr)

# Chapter IX

# Submission and peer-evaluation

Turn in your assignment in your `Git` repository as usual. Only the work inside your repository will be evaluated during the defense. Don't hesitate to double check the names of your folders and files to ensure they are correct.

> The evaluation process will happen on the computer of the evaluated group.