

# Analyse de données - TP 5

ISEP – 24 Novembre 2020

Instructions : Déposez votre rapport au format pdf sur Moodle

## Bibliothèques

Ce TP nécessite l'installation des librairies suivantes : pandas, matplotlib, numpy, scipy, sklearn et pydotplus.

## A Données du titanic

### A.1 Analyse préliminaire

1. Ouvrez les fichiers "titanic\_train.csv" et "titanic\_test.csv"
  - Pourquoi y a-t-il 2 fichiers ? Comparez-les. Sont-ils identiques ?
  - A l'aide de la fonction *head()* et de la propriété *dtype*, décrivez les différents attributs du jeu de données *titanic\_train*.
  - En utilisant les fonctions *isnull()* et *sum()* que pouvez-vous dire des données manquantes des différents attributs ?
2. A partir du fichier *titanic\_train*, affichez les informations suivantes :
  - Affichez un histogramme pour l'âge des passagers. Justifiez votre gestion des valeurs manquantes.
  - Affichez le pourcentage de personnes décédées.
  - Affichez le pourcentage de personnes en 1ère, 2nde et 3ème classe. Affichez un diagramme en camembert.
  - Affichez les pourcentages d'hommes et de femmes à bord et le taux de survie dans les 2 catégories.
  - En sachant que l'âge adulte était à 18 ans, estimez les pourcentages de survivants chez les enfants et les adultes ? Justifiez votre gestion des données manquantes.
  - Estimez les taux de survie pour toutes les combinaisons adulte/enfant, homme/femme.
  - Affichez les taux de survie en 1ère, 2nd et 3ème classe. Expliquez à nouveau votre gestion des données manquantes.
3. A partir des réponses précédentes, que pouvez-vous dire sur le respect de la politique "les femmes et les enfants d'abord" lors du naufrage du Titanic ?
4. Affichez une matrice de corrélation des différents attributs. Commentez sur d'éventuelles corrélations remarquables.

5. Dans les 2 jeux des données (training et test), créez un attribut "Fare2" dont la valeur est 1 si le prix du ticket était inférieur à 10, 2 en dessous de 20, 3 en dessous de 30, et 4 sinon. En utilisant les bons outils mathématiques, évaluez la force du lien entre votre nouvel attribut "Fare2" et la classe dans laquelle les passagers voyageaient. Affichez ensuite le taux de survivants pour chaque valeur de "Fare2". Commentez et concluez sur l'influence du prix des tickets et les chances de survie.

## A.2 Prédiction avec Naive Bayes

1. Si vous ne l'aviez pas déjà fait, rajoutez un attribut "Child" dans les deux jeux de données avec un valeur de "1" pour les passagers de moins de 18 ans, et 0 sinon. N'oubliez pas de préciser votre politique pour les données manquantes en indiquant les avantages et les inconvénients de votre choix.
2. Utilisez l'algorithme Naive Bayes pour prédire les taux de survie des passagers en fonction de leur sexe et de s'ils étaient adultes ou non. Confirmez vos résultats sur les jeux de données d'apprentissage et de test. Expliquez et commentez les différentes métriques de la fonction *classification\_report*. Un exemple de code est disponible ci-dessous.

```

1  # needed imports
2  from sklearn import metrics
3  from sklearn.naive_bayes import GaussianNB
4
5  #classifier choice
6  gnbModel=GaussianNB()
7  #choice of the training set, considered attributes and variable to \
   predict
8  gnbModel.fit(df_train[['Child', 'Sex_cat']], df_train['Survived']) \
   #Sex_cat stands for transformed categorical array
9
10 #expected results are stored in a separate vector
11 expected =df_train['Survived']
12
13 #predictions on the training set
14 predicted = gnbModel.predict(df_train[['Child', 'Sex_cat']])
15 #displaying relevant metrics
16 print(metrics.classification_report(expected, predicted))
17
18 #same when applying the model to the test set
19 expected =df_test['Survived']
20 predicted = gnbModel.predict(df_test[['Child', 'Sex_cat']])
21 print(metrics.classification_report(expected, predicted))

```

3. Appliquez à nouveau l'algorithme Naive Bayes avec le sexe, adulte ou non, et l'attribut Fare2. Commentez les résultats.

## A.3 Prédiction à partir d'un arbre de décision

1. Entraînez un arbre de décision pour prédire si les passagers vont survivre ou non. Pour cela, utilisez tous les attributs qui vous semblent pertinents, y compris ceux que vous avez créés. Vous pouvez essayer plusieurs combinaisons. Commentez les résultats.

Remarque 1 : Tous les classifieurs de sklearn fonctionnent de la même façon et ont une syntaxe similaire. La fonction *fit()* permet d'entraîner l'algorithme sur le jeu de données d'entraînement. La fonction *predict()* permet de calculer la classe la plus probable à partir des attributs de nouvelles données. Le code ci-dessous est un exemple de comment utiliser un arbre de décision et afficher le résultat.

Remarque 2 : Si vous rencontrez des difficultés avec la librairie *graphviz*. Utilisez plutôt la fonction *plot\_tree* de *sklearn*.

```
1 # needed imports
2 import pydotplus
3 import collections
4 from IPython.display import Image
5 from sklearn.tree import export_graphviz
6 from sklearn import tree
7
8 clf = tree.DecisionTreeClassifier()
9 clf = clf.fit(df_train[['Pclass', 'Fare2', 'Child', 'Sex_cat']], df_train['Survived\
    ']) #Sex_cat stands for transformed categorical array
10
11 dot_data = tree.export_graphviz(clf, out_file=None, feature_names=\
    data_feature_names, class_names=True, filled=True, rounded=True, precision\
    =0)
12 graph = pydotplus.graph_from_dot_data(dot_data)
13 Image(graph.create_png()) # to plot a tree
14 Image(graph.write_png('./filename.png')) # to save the plot
```