

TP 3 : analyses statistiques

Exercice I : Le test de student

1. Construire un vecteur (nommé X) de 100 valeurs dont chaque élément est issu d'une loi normale de moyenne nulle et de variance unitaire. Construire un vecteur (nommé Y) de 50 valeurs dont chaque élément est issu d'une loi normale centrée en 1 et de variance unitaire.
2. Réaliser un test de student pour tester l'égalité des moyennes
3. Quel type d'objet est renvoyé par la fonction utilisée en question 2.
4. Récupérer la p-valeur du test. Qu'en concluez-vous ?

Exercice II : La régression linéaire simple

Dans le package datasets est disponible le jeu de données cars. Ce jeu de données est composé de deux variables la première indiquant la vitesse de voitures et la seconde le temps de freinage.

1. Charger le jeu de donnée 'cars'
2. Tracer le graphe bivarié (abscisses=vitesse, ordonnées= distance). Que constate t'on ?
3. Construire un modèle linéaire qui analyse la distance de freinage comme fonction de la vitesse. Stiquer le résultat du modèle dans un objet nommé result.lm.
4. Ajouter au graphe bivarié construit à la question 1, la droite de régression résultant du modèle de la question 2.
5. Donner les valeurs prédites par le modèle pour chacune des observations.
6. Calculer la différence entre les valeurs prédites par le modèle et les valeurs observées. Comparer ce résultat aux résidus disponible dans l'objet residuals de la liste des outputs
7. Construire et visualiser des intervalles de prédiction et de confiance.

Exercice III : l'analyse en composantes principales

Pour illustrer l'Analyse en Composantes Principales, nous allons, une nouvelle fois, utiliser le jeu de données bordeaux_R.txt.

1. Charger le jeu de données bordeaux_R.

2. Réaliser une analyse en composantes principales (à partir des 4 variables température, soleil, chaleur, pluie) sur les variables centrées réduites.
3. Afficher sur un même graphique les variables et les individus (biplot) sur les deux premiers axes principaux.
4. Afficher les qualités des vins comme labels des individus.
5. Interpréter le résultat de l'analyse.
6. Donner le pourcentage de variance expliquée par les deux premières composantes

Exercice IV : la régression logistique

1. Récupérer les deux premières composantes principales de l'exercice précédent et stocker les dans un objet nommé X .
2. Construire une variable binaire (nommée y) prenant la valeur 1 si le vin est « bon » et prenant la valeur 0 si le vin est « moyen » ou « médiocre »
3. En utilisant la fonction glm() construire un modèle logistique de Y sur X.
4. Récupérer les probabilités prédites par le modèle
5. En déduire les décisions du modèle pour chacun des individus (On seuillera les probabilités d'appartenance à un seuil fixé à 0.5).
6. Renvoyer le tableau de contingence. En déduire le taux de bonne classification en apprentissage.
7. On s'intéresse maintenant au taux de bonne classification en test. Pour ce faire sélectionner aléatoirement 22 individus (avec la fonction 'sample') pour construire votre modèle et le restant pour le tester.
8. Comparer les taux tableaux de contingence d'apprentissage et de test. En déduire les taux d'erreur en apprentissage et en test.

Exercice V : La methode des k-means

1. Récupérer l'objet nommé X avec les scores de l'analyse en composantes principales de l'exercice précédent et tracer le graphe bivarié du premier score sur le deuxième ; colorer les points selon leur qualité (« bon », « moyen », « médiocre »).
2. Construire une partition de 3 groupes de ce nuage de points via l'algorithme des kmeans.
3. Lier chacune des observations à son centroïde par un segment et compter visuellement le taux de mal classés.