

Data Mining

TP1 - Sélection de modèles

Rémi MUSSARD - Thomas ROBERT

Question 1

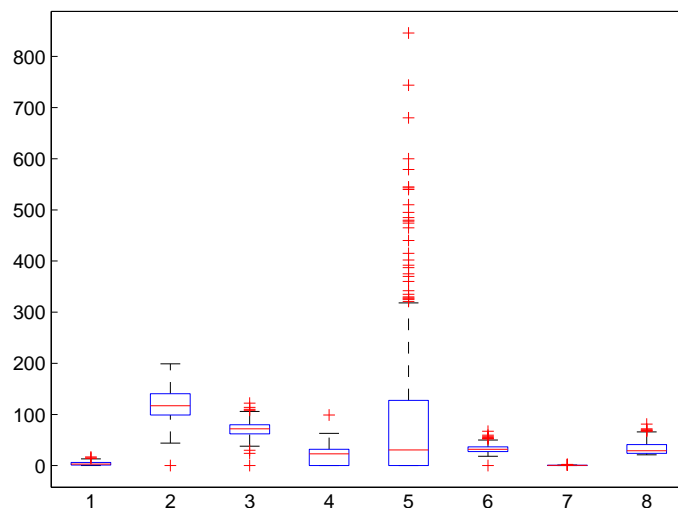
Effectuer une brève analyse statistique des données : moyenne, écart-type de chaque variable.

```
1 M = mean(X)
2 Med = median(X)
3 ET = std(X)
4 figure();
5 boxplot(X);
```

M =
3.8451 120.8945 69.1055 20.5365 79.7995 31.9926 0.4719 33.2409

Med =
3.0000 117.0000 72.0000 23.0000 30.5000 32.0000 0.3725 29.0000

ET =
3.3696 31.9726 19.3558 15.9522 115.2440 7.8842 0.3313 11.7602



Ces diverses valeurs permettent d'avoir une vue globale des données.

Question 2

Implémenter une méthode de k -ppv (k plus proche voisins) avec une distance euclidienne.

On utilise la fonction knn fournie qui renvoie les prédictions et les distances calculées à partir des données d'entrées.

Question 3

Séparer aléatoirement l'ensemble des données en un ensemble d'apprentissage et un ensemble de test en respectant au mieux la proportion des classes. L'ensemble de test ne sera utilisé qu'une fois. Nota : voir la fonction splitdata.m sur Moodle.

La fonction splitdata permet de séparer les données en respectant la proportion des classes. On sépare ainsi notre ensemble de données en deux ensembles de même taille.

```
1 [xapp, yapp, xtest, ytest] = splitdata(X, y, 0.5);
```

Question 4

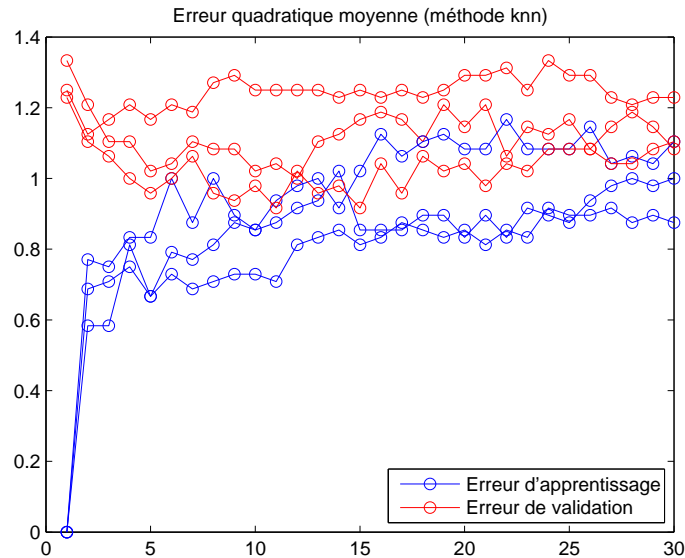
Séparer l'ensemble d'apprentissage en 2 ensembles : un autre ensemble d'apprentissage et un ensemble de validation. Tester votre méthode k-ppv sur l'ensemble d'ensemble d'apprentissage et l'ensemble de validation pour différentes valeurs de k ? Tracer une courbe de l'erreur d'apprentissage et une courbe de l'erreur de validation en fonction de k. Quelle est la valeur de k qui donne la plus faible erreur en validation ?

```
1 figure();
2
3 % on execute le code 3 fois pour voir son insatabilité
4 for i=1:3
5     % init de la matrice des erreurs pour chaque k de 1 à 30
6     kmax = 30;
7     errApp = zeros(kmax, 1);
8     errVal = zeros(kmax, 1);
9
10    % découpage en jeu d'apprentissage et de validation
11    [xapp2, yapp2, xval2, yval2] = splitdata(xapp, yapp, 0.5);
12
13    % pour chaque k
14    for k = 1:kmax
15
16        % prédiction
17        [ypredApp, Dist] = knn(xapp2, xapp2, yapp2, k);
18        [ypredVal, Dist] = knn(xval2, xapp2, yapp2, k);
19
20        % calcul d'erreur quadratique moyenne pour le k choisi
21        errApp(k) = mean((ypredApp - yapp2).^2);
22        errVal(k) = mean((ypredVal - yval2).^2);
23    end
24
25    % affichage des erreurs
26    plot(errApp, 'o-');
27    hold on;
28    plot(errVal, 'or-');
29    title('Erreur quadratique moyenne (méthode knn)');
30    leg = legend('Erreur d'apprentissage', 'Erreur de validation');
31    set(leg, 'Location', 'SouthEast');
32
33    % meilleure valeur de k en validation
34
35    [erreurMin, bestk] = min(errVal);
36    fprintf('Meilleur k par méthode knn : %i\n', bestk);
37 end
```

Meilleur k par méthode knn : 2

Meilleur k par méthode knn : 12

Meilleur k par méthode knn : 11



Avec cette méthode, on constate que la valeur de k trouvée est très dépendante du découpage aléatoire des données qui a été réalisé, entre ensemble de test et ensemble d'apprentissage.

Pour résoudre ce problème, on peut utiliser la méthode de validation croisée qui permet d'utiliser toute les données pour réaliser les tests et l'apprentissage de façon "rotative".

Question 5

Refaire l'expérience en utilisant une méthode de validation croisée sur les données d'apprentissage créées à la question 3. Quelle est la meilleure valeur de k ? Nota : voir sur Moodle la fonction *SeparateDataNfoldCV.m*.

On utilise la fonction *SeparateDataNfoldCV* pour découper l'ensemble d'apprentissage en plusieurs blocs afin d'y appliquer la méthode de validation croisée.

```

1 % constantes
2 kmax = 23;
3 Nfold = 20;
4 errApp = zeros(kmax, Nfold);
5 errVal = zeros(kmax, Nfold);
6
7 % pour chaque bloc
8 for NumFold = 1:Nfold
9     % séparation des données
10    [xapp2, yapp2, xval2, yval2] = SeparateDataNfoldCV(xapp, yapp, Nfold, NumFold);
11
12    % pour chaque k
13    for k = 1:kmax
14
15        % prédiction
16        [ypredApp, Dist] = knn(xapp2, xapp2, yapp2, k);
17        [ypredVal, Dist] = knn(xval2, xapp2, yapp2, k);
18
19        % calcul d'erreur quadratique moyenne pour le k choisi
20        errApp(k, NumFold) = mean((ypredApp - yapp2).^2);
21        errVal(k, NumFold) = mean((ypredVal - yval2).^2);
22    end
23 end
24
25 % erreur d'apprentissage et de validation moyenne pour chaque k
26 errAppK = mean(errApp');
27 errValK = mean(errVal');

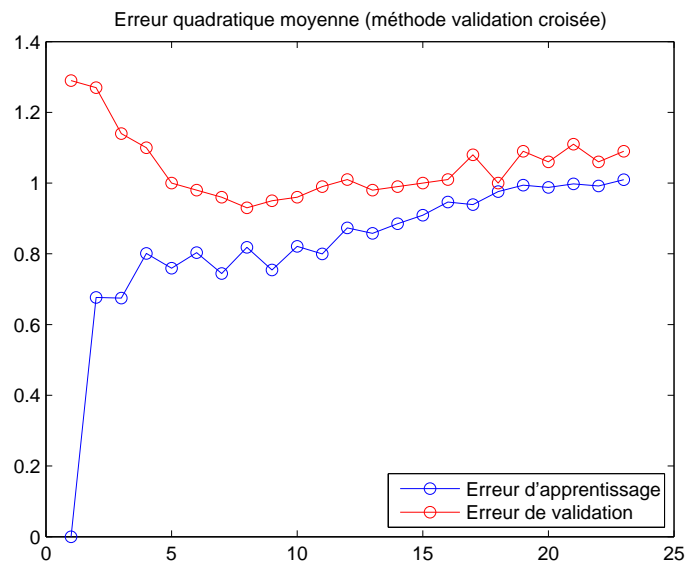
```

```

28
29 % affichage des erreurs
30 figure();
31 plot(errAppK, 'o-');
32 hold on;
33 plot(errValK, 'or-');
34 title('Erreur quadratique moyenne (méthode validation croisée)');
35 leg = legend('Erreur d\'apprentissage', 'Erreur de validation');
36 set(leg, 'Location', 'SouthEast');
37 hold off;
38
39 % meilleure valeur de k en validation
40 [erreurMin, bestk] = min(errValK);
41 fprintf('Meilleur k par méthode de validation croisée : %i\n', bestk);

```

Meilleur k par méthode de validation croisée : 8



On voit effectivement que la méthode de la validation croisée est plus stable.