

SCIENCE DES DONNÉES NUMÉRIQUES

TP3

Manel SELSANE ; Fatma GOUIAA ; Pépita KENESI-LAURENT ; Leo DECHAUMET

Novembre 2021

Table des matières

| | | |
|----------|----------------------------------|----------|
| 1 | Plus Proche Voisin | 2 |
| 2 | Classifieur Bayésien Naïf | 7 |

1 Plus Proche Voisin

1.

On crée la fonction PPV qui prend en entrée des données X et des étiquettes Y et qui renvoie une étiquette, pour chaque donnée, prédite à partir du plus proche voisin de cette donnée.

On évalue l'algorithme du Plus Proche Voisin avec la cross validation donc toutes les données sont données test.

Parameters :

X : numpy.ndarray

Une matrice de dimension 2 où chaque ligne correspond à une valeur

Y : array-like

Les labels

Returns :

numpy.array

Les predictions

```
1 def PPV(X,Y):
2     Ychapeau=[]
3     for e in X:
4         L=metrics.pairwise.euclidean_distances(X,e[np.newaxis])
5         L=L.reshape(1,-1)
6         L2=np.argsort(L)
7         Ychapeau.append(Y[L2[0],1])
8     return np.array(Ychapeau)
```

2.

On veut maintenant l'erreur de prediction, on calcule le nombre de valeur pour lesquelles X est different de Y.

On crée d'abord une fonction erreur qui renvoie l'erreur entre un jeu de donnée X et sa prédiction Y , puis on l'applique ensuite a PPV.

Parameters :

X : numpy.array

Le premier vecteur

Y : numpy.array

Le second vecteur

Returns :

float

Le pourcentage de valeurs différentes

```
1 def Erreur(X, Y):
2     N = X.size - sum(X == Y)
3     return N / X.size * 100
4
5 def ErreurPPV(X, Y):
6     return Erreur(PPV(X, Y), Y)
```

Un autre algorithme pour calculer directement l'erreur pour PPV :

```
1 def PPV3(Ychapeau, Y) :
2     s = 0
3     for i in range(0, len(Y)):
4         if Ychapeau[i] != Y[i]:
5             s=s+1
6     return s*100 / len(Y)
```

On teste maintenant nos fonctions sur les données d'iris :

```
PPV pour iris :  
[0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0  
0 0 0 0 0 0 0 0 0 0 0 0 0 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 2 1 2 1  
1 1 1 1 1 1 1 1 1 1 2 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 2 2 2 2 2 2 1 2 2 2 2  
2 2 2 2 2 2 2 2 1 2 2 2 2 2 2 2 2 2 2 2 2 2 1 2 2 2 2 2 2 2 2 2 2 2 2 2 2  
2 2]  
  
Erreur PPV pour iris  
4.0
```

4. On teste maintenant la fonction des K plus proches voisin de sklearn avec K=1 sur les données d'iris :

```
Ychapeau pour K=1:  
[0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0  
0 0 0 0 0 0 0 0 0 0 0 0 0 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1  
1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 2 2 2 2 2 2 2 2  
2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2  
2 2]
```

Erreur Ychapeau pour K=1 :
0.0

Erreur PPV :
4.0

Alors que dans notre algorithme, nous avons supprimé la possibilité de prendre comme PPV d'une donnée la donnée elle-même, ce qui explique la différence.

On teste ensuite l'algorithme KNC pour d'autres valeurs de K :

```

1 neigh2 = KNeighborsClassifier(n_neighbors=2)
2 neigh2.fit(X, Y)
3 Ychapeau2 = neigh2.predict(X)
4 print(Ychapeau2)
5 print(Erreur(Ychapeau2, Y))
6
7 neigh3 = KNeighborsClassifier(n_neighbors=3)
8 neigh3.fit(X, Y)
9 Ychapeau3 = neigh3.predict(X)
10 print(Ychapeau3)
11 print(Erreur(Ychapeau3, Y))

```

```
Ychapeau pour K=2:  
[0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0  
0 0 0 0 0 0 0 0 0 0 0 0 0 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1  
1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 2 2 2 2 2 2 1 2 2 2 2  
2 2 2 2 2 2 2 2 1 2 2 2 2 2 2 2 2 2 2 2 2 2 1 2 2 2 2 2 2 2 2 2 2 2  
2 2]
```

Erreur Ychapeau pour K=2 :
2.0

```
Ychapeau pour K=3:  
[0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0  
0 0 0 0 0 0 0 0 0 0 0 0 0 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 2 1 2 1  
1 1 1 1 1 1 1 1 1 1 2 1 1 1 1 1 1 1 1 1 1 1 1 1 1 2 2 2 2 2 2 1 2 2 2 2  
2 2 2 2 2 2 2 2 1 2 2 2 2 2 2 2 2 2 2 2 2 2 1 2 2 2 2 2 2 2 2 2 2 2 2 2  
2 2]
```

Erreur Ychapeau pour K=3 :
4.0

5.

On modifie maintenant notre fonction PPV pour qu'elle prenne en entrée un nombre K de voisins . La classe prédite sera alors la classe majoritaire parmi les K voisins.

Parameters :

k : int

Le nombre de voisins à prendre en considération

X : numpy.ndarray

Une matrice de dimension 2 où chaque ligne correspond à une valeur

Y : array-like

Les labels

Returns :

numpy.array

Les predictions

```
1 def PPV_mod(k, X, Y):
2     Ychapeau = []
3     for i, e in enumerate(X):
4         L = metrics.pairwise.euclidean_distances(X, e[np.newaxis])
5         L = L.reshape(Y.size)
6         L2 = np.argsort(L)
7         G = [L2[j] for j in range(1, k+1)]
8         M = list(Y[G])
9     return np.array(Ychapeau)
```

2 Classifieur Bayesien Naïf

1.

On veut créer une fonction $\text{CBN}(X,Y)$ qui prend en entrée des données X et des étiquettes Y et qui renvoie une étiquette.

On commence par définir plusieurs fonctions afin de simplifier la démarche :

$\text{baricentre}(X, Y)$: Baricentre de toutes les classes :

Calcule le baricentre de toutes les classes comme étant la moyenne des points composant cette classe.

Parameters :

X : `numpy.ndarray`

Tableau à deux dimensions, chaque ligne est un point

Y : `numpy.array`

Les labels

Returns :

`numpy.array`

La liste des baricentres de chaque classe (l'ordre des classes est donné par la fonction `numpy.unique`).

$d(a,b)$:

Calcule la distance entre deux points.

Parameters : a : `numpy.array`

Un vecteur représentant le premier point

b : `numpy.array`

Un vecteur représentant le second point

Returns :

`float`

La distance euclidienne entre les deux points.

$P1(Y)$:

Pour chaque classe k dans Y , calcule la probabilité d'obtenir un élément de la classe k si on tire un élément au hasard.

Parameters :

Y : `numpy.array`

Un vecteur représentant les labels

Returns :

p_s : `numpy.array`

La liste des probabilités pour toutes les classes.

P2(x,k,baricentre) :

Calcule la probabilité d'avoir le point x sachant qu'il appartient à la classe k.

Parameters :

x : numpy.array

Un vecteur représentant point x

k : int

La classe à laquelle appartient le point x baricentre

baricentre : numpy.array Un vecteur representant les baricentres de toutes les classes.

Returns :

a : float

Un nombre entre 0 et 1 représentant la probabilité.

```
1  def baricentre(X, Y):
2      b_s = []
3      for k in np.unique(Y):
4          A = X[np.where(Y == k)]
5          b_s.append(np.mean(A, axis=0))
6      return np.array(b_s)
7
8
9  def d(a, b):
10     return np.linalg.norm(a - b)
11
12
13  def P1(Y):
14     p_s = []
15     for e in np.unique(Y):
16         p = sum(Y == e)
17         p_s.append(p/Y.size)
18     return p_s
19
20
21  def P2(x, k, baricentre):
22     sum_dist = np.sum(d(x, bar) for bar in baricentre)
23     a = 1 - (d(x, baricentre[k]) / sum_dist)
24     return a
```


On teste maintenant la fonction du Classifieur Bayésien Naïf de sklearn :

```
classificateur gaussien pour iris :  
[0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0  
0 0 0 0 0 0 0 0 0 0 0 0 0 1 1 2 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1  
1 1 1 2 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 2 2 2 2 2 1 2 2 2  
2 2 2 2 2 2 2 2 1 2 2 2 2 2 2 2 2 2 2 2 2 1 2 2 2 2 2 2 2 2 2 2 2  
2 2]  
  
Erreur classificateur gaussien pour iris  
4.0
```

10