

TP1 KNN

Entrée []:



```
1  # imports mnist et KNN
2  from sklearn.datasets import fetch_mldata
3  from sklearn import datasets
4  from sklearn.model_selection import train_test_split
5  from sklearn import neighbors
6  import numpy as np
7
8  #imports mesure de performance
9  import matplotlib.pyplot as plt
10 from sklearn import metrics
11 import time
12
13 mnist=datasets.fetch_mldata('MNIST original')
14
15
```

Entrée []:



```
1  #QUESTION 1
2
3  #imprime la structure et un aperçu des data
4  print(mnist)
5  #imprime un aperçu des data
6  #print(mnist.data)
7
8  # imrime la cible ???
9  #print(mnist.target)
10
11 #Longueur de data
12 #len(mnist.data) #70000
13 print(mnist.data.shape) #(70000, 784)
14 print(mnist.target.shape) #(70000,)
15
16 #mnist.data[0]
```

Entrée []:



```
1  #QUESTION 2
2
3  #Afficher une image
4  images = mnist.data.reshape((-1,28,28))
5  plt.imshow(images[0],cmap=plt.cm.gray_r,interpolation="nearest")
6  plt.show()
7  print("classe : ", mnist.target[0])
8
9
```

Entrée []:



```
1  #####
2  #          EXERCICE 2
3  #####
4
5
6  #ECHANTILLONS DE 10000
7  indices = np.random.randint(70000, size=10000)
8  data = mnist.data[indices]
9  target = mnist.target[indices]
10 #pour séparer le dataset en 2 échantillons de trainint et de test
11 xtrain, xtest, ytrain, ytest =train_test_split(data, target, train_size=0.8)
12
13
```

Entrée []:



```
1  clf=neighbors.KNeighborsClassifier(10)
2  clf.fit(xtrain,ytrain)
3  predict = clf.predict(xtest)
4  proba=clf.predict_proba(xtest)
5  score = clf.score(xtest,ytest)
6
7  #Résultat 1er test
8  print("3eme image : prédiction ",predict[3], "reel : ", ytest[3])
9  print("Pour K=10 score test : ", score*100, "%")
10
```

```
1  #VARIATION des data prises et de K avec Kfold cross validation
2  #ATTENTION, ce n'est pas une utilisation adéquate de K-fold,
3  # mais c'est un choix pour voir l'effet du changement de K sur des données dif
4  # voir plus bas pour une utilisation de K-fold correcte
5  ResScore = []
6  ResPrecision = []
7  ResRecall = []
8  ResLoss = []
9  ResTimeTraining = []
10 ResTimePrediction = []
11
12 from sklearn.model_selection import KFold
13 k_fold = KFold( n_splits=15, shuffle=True)
14 print(k_fold)
15 k=1
16 for train_indices, test_indices in k_fold.split(data, target):
17     k+=1
18     clf=neighbors.KNeighborsClassifier(k)
19     # Algo
20     startTrain =time.time()
21     clf.fit(data[train_indices],target[train_indices])
22     endTrain = time.time()
23     startpred= time.time()
24     clf.predict(data[test_indices])
25     endpred = time.time()
26     # Metrics
27     predict = clf.predict(data[test_indices])
28     score = clf.score(data[test_indices],target[test_indices])
29     precision = metrics.precision_score(target[test_indices], predict, avera
30     recall = metrics.recall_score(target[test_indices], predict, average ='mac
31     loss01 = metrics.zero_one_loss(target[test_indices], predict)
32     timetrain = endTrain - startTrain
33     timePred = endpred - startpred
34     # Append
35     ResScore.append(score)
36     ResPrecision.append(precision)
37     ResRecall.append(recall)
38     ResLoss.append(loss01)
39     ResTimeTraining.append(timetrain)
40     ResTimePrediction.append(timePred)
41     print("Pour K = ", k, "score = ", clf.score(data[test_indices],target[test
42
43
```

Entrée []:



```
1 # pour garder une trace des résultats et ne pas runner encore une fois
2 print(ResScore)
3 print(ResPrecision)
4 print(ResRecall)
5 print(ResLoss)
6 print(ResTimeTraining)
7 print(ResTimePrediction)
8
```

Entrée []:



```
1 fig, axarr = plt.subplots(6, sharex=True, figsize=(10,10))
2 axarr[0].plot(range(15), ResScore)
3 axarr[0].set_title('KNN for 10000 training data, for k 1..16, with k_fold')
4 axarr[0].set_ylabel('Score')
5 axarr[1].plot(range(15), ResPrecision)
6 axarr[1].set_ylabel('Precision')
7 axarr[2].plot(range(15), ResRecall)
8 axarr[2].set_ylabel('Recall')
9 axarr[3].plot(range(15), ResLoss)
10 axarr[3].set_ylabel('Zero-to-one Loss')
11 axarr[4].plot(range(15), ResTimeTraining)
12 axarr[4].set_ylabel('Training Time in sec')
13 axarr[5].plot(range(15), ResTimePrediction)
14 axarr[5].set_ylabel('Prediction Time in sec')
15
16 plt.show()
17
```

Entrée []:



```
1 # VARIATION DE K sur Les même xtrain, ytrain
2 #AUTRE METHODE + Simple
3 ResScore1 = []
4 ResPrecision1 = []
5 ResRecall1 = []
6 ResLoss1 = []
7 ResTimeTraining1 = []
8 ResTimePrediction1 = []
9
10 for k in range(2,16):
11     clf=neighbors.KNeighborsClassifier(k)
12     startTrain =time.time()
13     clf.fit(xtrain,ytrain)
14     endTrain = time.time()
15     startpred= time.time()
16     clf.predict(xtest)
17     endpred = time.time()
18     # Metrics
19     predict = clf.predict(xtest)
20     score = clf.score(xtest,ytest)
21     precision = metrics.precision_score(ytest, predict, average='macro')
22     recall = metrics.recall_score(ytest, predict, average = 'macro')
23     loss01 = metrics.zero_one_loss(ytest, predict)
24     timetrain = endTrain - startTrain
25     timePred = endpred - startpred
26     # Append
27     ResScore1.append(score)
28     ResPrecision1.append(precision)
29     ResRecall1.append(recall)
30     ResLoss1.append(loss01)
31     ResTimeTraining1.append(timetrain)
32     ResTimePrediction1.append(timePred)
33
34     print("Pour K = ", k, "score = ", clf.score(xtest,ytest)*100, "%")
35
```

Entrée []:



```
1 # pour garder une trace des résultats et ne pas runner encore une fois
2 print(ResScore1)
3 print(ResPrecision1)
4 print(ResRecall1)
5 print(ResLoss1)
6 print(ResTimeTraining1)
7 print(ResTimePrediction1)
8
```

Entrée []:



```
1
2 ResScore1 =[0,0,0.9355, 0.948, 0.9465, 0.947, 0.945, 0.946, 0.9435, 0.9415, 0.
3 ResPrecision1 =[0,0,0.9369785302063333, 0.9483535681114812, 0.9478015367975028
4 ResRecall1 =[0,0,0.9337392168963335, 0.9462140532025879, 0.9447111859424187, 0
5 ResLoss1 =[0,0,0.0645, 0.052000000000000046, 0.05349999999999999, 0.0530000000
6 ResTimeTraining1 =[0,0,0.5166172981262207, 0.49174928665161133, 0.474731922149
7 ResTimePrediction1 =[0,0,22.78105330467224, 24.293938875198364, 24.77076554298
```

Entrée []:



```
1 fig, axarr = plt.subplots(6, sharex=True, figsize=(10,10))
2 axarr[0].plot(range(16), ResScore1)
3 axarr[0].set_title('KNN for 10000 training data, for k 1..16')
4 axarr[0].set_ylabel('Score')
5 axarr[0].set_ylim([0.93,0.95])
6 axarr[1].plot(range(16), ResPrecision1)
7 axarr[1].set_ylabel('Precision')
8 axarr[1].set_ylim([0.93,0.95])
9 axarr[2].plot(range(16), ResRecall1)
10 axarr[2].set_ylabel('Recall')
11 axarr[2].set_ylim([0.92,0.95])
12 axarr[3].plot(range(16), ResLoss1)
13 axarr[3].set_ylabel('Zero-to-one Loss')
14 axarr[3].set_ylim([0.05,0.07])
15 axarr[4].plot(range(16), ResTimeTraining1)
16 axarr[4].set_ylabel('Training Time in sec')
17 axarr[4].set_ylim([0.46,0.68])
18 axarr[5].plot(range(16), ResTimePrediction1)
19 axarr[5].set_ylabel('Prediction Time in sec')
20 axarr[5].set_ylim([22,25])
21 plt.xlim(left=2)
22
23 plt.show()
```

```
1  #VARIATION des data prises avec Kfold cross validation
2  #Ici, l'utilisation de K-fold CV est OK
3  ResScore4 = []
4  ResPrecision4 = []
5  ResRecall4 = []
6  ResLoss4 = []
7  ResTimeTraining4 = []
8  ResTimePrediction4 = []
9
10 from sklearn.model_selection import KFold
11 k_fold = KFold( n_splits=15, shuffle=True)
12 print(k_fold)
13 clf=neighbors.KNeighborsClassifier(3)
14 for train_indices,test_indices in k_fold.split(data,target):
15     # Algo
16     startTrain =time.time()
17     clf.fit(data[train_indices],target[train_indices])
18     endTrain = time.time()
19     startpred= time.time()
20     clf.predict(data[test_indices])
21     endpred = time.time()
22     # Metrics
23     predict = clf.predict(data[test_indices])
24     score = clf.score(data[test_indices],target[test_indices])
25     precision = metrics.precision_score(target[test_indices], predict, average='macro')
26     recall = metrics.recall_score(target[test_indices], predict, average='macro')
27     loss01 = metrics.zero_one_loss(target[test_indices], predict)
28     timetrain = endTrain - startTrain
29     timePred = endpred - startpred
30     # Append
31     ResScore4.append(score)
32     ResPrecision4.append(precision)
33     ResRecall4.append(recall)
34     ResLoss4.append(loss01)
35     ResTimeTraining4.append(timetrain)
36     ResTimePrediction4.append(timePred)
37     print("Pour K = 3 score = ", clf.score(data[test_indices],target[test_indices]))
38
39
40
41 # pour garder une trace des résultats et ne pas runner encore une fois
42 print(ResScore4)
43 print(ResPrecision4)
44 print(ResRecall4)
45 print(ResLoss4)
46 print(ResTimeTraining4)
47 print(ResTimePrediction4)
48
```

```
1
2 fig, axarr = plt.subplots(6, sharex=True, figsize=(10,10))
3 axarr[0].plot(range(15), ResScore4)
4 axarr[0].set_title('KNN for 10000 training data, 15-fold Cross Validation')
5 axarr[0].set_ylabel('Score')
6 axarr[1].plot(range(15), ResPrecision4)
7 axarr[1].set_ylabel('Precision')
8 axarr[2].plot(range(15), ResRecall4)
9 axarr[2].set_ylabel('Recall')
10 axarr[3].plot(range(15), ResLoss4)
11 axarr[3].set_ylabel('Zero-to-one Loss')
12 axarr[4].plot(range(15), ResTimeTraining4)
13 axarr[4].set_ylabel('Training Time in sec')
14 axarr[5].plot(range(15), ResTimePrediction4)
15 axarr[5].set_ylabel('Prediction Time in sec')
16
17 plt.show()
```



```
1  # VARIATION de la taille des échantillons
2  ResScore2 = []
3  ResPrecision2 = []
4  ResRecall2 = []
5  ResLoss2 = []
6  ResTimeTraining2 = []
7  ResTimePrediction2 = []
8
9
10 for k in range(1,10):
11     clf=neighbors.KNeighborsClassifier(3)
12     x_train, x_test, y_train, y_test =train_test_split(data, target, train_size=0.2)
13     startTrain =time.time()
14     clf.fit(x_train,y_train)
15     endTrain = time.time()
16     startpred= time.time()
17     clf.predict(x_test)
18     endpred = time.time()
19     # Metrics
20     predict = clf.predict(xtest)
21     score = clf.score(xtest,ytest)
22     precision = metrics.precision_score(ytest, predict, average='macro')
23     recall = metrics.recall_score(ytest, predict, average = 'macro')
24     loss01 = metrics.zero_one_loss(ytest, predict)
25     timetrain = endTrain - startTrain
26     timePred = endpred - startpred
27     # Append
28     ResScore2.append(score)
29     ResPrecision2.append(precision)
30     ResRecall2.append(recall)
31     ResLoss2.append(loss01)
32     ResTimeTraining2.append(timetrain)
33     ResTimePrediction2.append(timePred)
34     print("Pour un train/test = ", (k*10), "%, score = ", clf.score(x_test,y_test))
35
36 #Pour un train/test = 10 %, score = 81.44444444444444 %
37 #Pour un train/test = 20 %, score = 86.575 %
38 #Pour un train/test = 30 %, score = 88.31428571428572 %
39 #Pour un train/test = 40 %, score = 90.0 %
40 #Pour un train/test = 50 %, score = 91.12 %
41 #Pour un train/test = 60 %, score = 91.05 %
42 #Pour un train/test = 70 %, score = 89.8 %
43 #Pour un train/test = 80 %, score = 91.2 %
44 #Pour un train/test = 90 %, score = 93.2 %
```

Entrée []:



```
1 # pour garder une trace des résultats et ne pas runner encore une fois
2 print(ResScore2)
3 print(ResPrecision2)
4 print(ResRecall2)
5 print(ResLoss2)
6 print(ResTimeTraining2)
7 print(ResTimePrediction2)
8
9 #[0.8565, 0.8985, 0.91, 0.924, 0.935, 0.939, 0.941, 0.949, 0.953]
10 #[0.8754876643697613, 0.9069422336372133, 0.9176068152884023, 0.92849192544389
11 #[0.8507976794260326, 0.8943515897570011, 0.9066121362228408, 0.92128075126968
12 #[0.14349999999999996, 0.10150000000000003, 0.08999999999999997, 0.07599999999
13 #[0.02293848991394043, 0.06283235549926758, 0.10571765899658203, 0.17253923416
14 #[12.847708702087402, 22.783091068267822, 26.58786702156067, 34.31918358802795
```

Entrée []:



```
1 fig, axarr = plt.subplots(6, sharex=True, figsize=(10,10))
2 axarr[0].plot(range(9), ResScore2)
3 axarr[0].set_title('KNN for 15000 training data, changind size of train/test s
4 axarr[0].set_ylabel('Score')
5 axarr[1].plot(range(9), ResPrecision2)
6 axarr[1].set_ylabel('Precision')
7 axarr[2].plot(range(9), ResRecall2)
8 axarr[2].set_ylabel('Recall')
9 axarr[3].plot(range(9), ResLoss2)
10 axarr[3].set_ylabel('Zero-to-one Loss')
11 axarr[4].scatter(range(9), ResTimeTraining2)
12 axarr[4].set_ylabel('Training Time in sec')
13 axarr[5].scatter(range(9), ResTimePrediction2)
14 axarr[5].set_ylabel('Prediction Time in sec')
15
16 plt.show()
```

```

1  #VARIATION de la taille des échantillons
2  ResScore3 = []
3  ResPrecision3 = []
4  ResRecall3 = []
5  ResLoss3 = []
6  ResTimeTraining3 = []
7  ResTimePrediction3 = []
8
9  for k in range(1,6):
10     indices = np.random.randint(70000, size=5000*k)
11     data = mnist.data[indices]
12     target = mnist.target[indices]
13     #pour séparer le dataset en 2 échantillons de trainint et de test
14     x_train, x_test, y_train, y_test = train_test_split(data, target, train_size=0.8)
15     clf = neighbors.KNeighborsClassifier(3)
16     startTrain = time.time()
17     clf.fit(x_train, y_train)
18     endTrain = time.time()
19     startpred = time.time()
20     clf.predict(x_test)
21     endpred = time.time()
22     # Metrics
23     predict = clf.predict(x_test)
24     score = clf.score(x_test, y_test)
25     precision = metrics.precision_score(y_test, predict, average='macro')
26     recall = metrics.recall_score(y_test, predict, average='macro')
27     loss01 = metrics.zero_one_loss(y_test, predict)
28     timetrain = endTrain - startTrain
29     timePred = endpred - startpred
30     # Append
31     ResScore3.append(score)
32     ResPrecision3.append(precision)
33     ResRecall3.append(recall)
34     ResLoss3.append(loss01)
35     ResTimeTraining3.append(timetrain)
36     ResTimePrediction3.append(timePred)
37     print("Pour un échantillon de taille = ", (k*5000), ", score = ", clf.score(x_test, y_test))
38
39 #Pour un échantillon de taille = 5000 ,score = 92.5 %
40 #Pour un échantillon de taille = 10000 ,score = 93.60000000000001 %
41 #Pour un échantillon de taille = 15000 ,score = 94.63333333333334 %
42 #Pour un échantillon de taille = 20000 ,score = 94.69999999999999 %
43 #Pour un échantillon de taille = 25000 ,score = 95.54 %

```

Entrée []:



```
1 # pour garder une trace des résultats et ne pas runner encore une fois
2 print(ResScore3)
3 print(ResPrecision3)
4 print(ResRecall3)
5 print(ResLoss3)
6 print(ResTimeTraining3)
7 print(ResTimePrediction3)
8 #[0.9185, 0.9375, 0.947, 0.9545, 0.9565]
9 #[0.9239451509342509, 0.940184968052854, 0.94913040865386, 0.956196591830037,
10 #[0.9156434687814252, 0.9357598432718817, 0.9457375638031905, 0.95337192960119
11 #[0.08150000000000002, 0.0625, 0.05300000000000005, 0.04549999999999985, 0.04
12 #[0.1795210838317871, 0.6402866840362549, 0.894573450088501, 1.822124719619751
13 #[6.629263162612915, 24.836551904678345, 47.89191031455994, 92.35796546936035,
```

Entrée []:



```
1 fig, axarr = plt.subplots(6, sharex=True, figsize=(10,10))
2 axarr[0].scatter(range(5), ResScore3)
3 axarr[0].set_title('KNN for training data between 5000 and 25000, for K=3')
4 axarr[0].set_ylabel('Score')
5 axarr[1].scatter(range(5), ResPrecision3)
6 axarr[1].set_ylabel('Precision')
7 axarr[2].scatter(range(5), ResRecall3)
8 axarr[2].set_ylabel('Recall')
9 axarr[3].scatter(range(5), ResLoss3)
10 axarr[3].set_ylabel('Zero-to-one Loss')
11 axarr[4].scatter(range(5), ResTimeTraining3)
12 axarr[4].set_ylabel('Training Time in sec')
13 axarr[5].scatter(range(5), ResTimePrediction3)
14 axarr[5].set_ylabel('Prediction Time in sec')
15
16 plt.show()
```

Entrée []:



```
1 # VARIATION de la distance
2 # Par défaut, la distance est la distance euclidienne (minkowskiDistance avec
3 # see here for more details : https://scikit-learn.org/stable/modules/generate
4
5 print("Avec la distance Euclidienne, le score = ", clf.score(xtest,ytest)*100
6
7 # test ManhattanDistance "manhattan"
8 manha=neighbors.KNeighborsClassifier(3, metric='manhattan')
9 manha.fit(xtrain,ytrain)
10 manha.predict(xtest)
11 print("Avec la distance de Manhattan, le score = ", manha.score(xtest,ytest)*1
12
13 #test ChebyshevDistance "chebyshev"
14 cheby=neighbors.KNeighborsClassifier(3, metric='chebyshev')
15 cheby.fit(xtrain,ytrain)
16 cheby.predict(xtest)
17 print("Avec la distance de Chebyshev, le score = ", cheby.score(xtest,ytest)*1
18
19 #test Hamming distance
20 ham=neighbors.KNeighborsClassifier(3, metric='hamming')
21 ham.fit(xtrain,ytrain)
22 ham.predict(xtest)
23 print("Avec la distance de Hamming, le score = ", ham.score(xtest,ytest)*100,
24
25 #test minkowski p=3,4,5
26 mink3=neighbors.KNeighborsClassifier(3, p=3, metric='minkowski')
27 mink3.fit(xtrain,ytrain)
28 mink3.predict(xtest)
29 print("Avec la distance de Minkowski p=3 , le score = ", mink3.score(xtest,yte
30
31 mink4=neighbors.KNeighborsClassifier(3, p=4, metric='minkowski')
32 mink4.fit(xtrain,ytrain)
33 mink4.predict(xtest)
34 print("Avec la distance de Minkowski p=4, le score = ", mink4.score(xtest,ytes
35
36 mink5=neighbors.KNeighborsClassifier(3, p=5, metric='minkowski')
37 mink5.fit(xtrain,ytrain)
38 mink5.predict(xtest)
39 print("Avec la distance de Minkowski p=5, le score = ", mink5.score(xtest,ytes
40
```

Entrée []:



1

Entrée []:



```
1 # MESURE avec njobs
2 import time
3 job1=neighbors.KNeighborsClassifier(3, n_jobs=1)
4 job=neighbors.KNeighborsClassifier(3, n_jobs=-1)
5
```

Entrée []:



```
1 # MESURE avec njobs
2 # TODO :A FAIRE TOURNER !!!!!
3 #import timeit
4 #job1=neighbors.KNeighborsClassifier(10, n_jobs=1)
5 #job=neighbors.KNeighborsClassifier(10, n_jobs=-1)
6 # for help https://docs.python.org/2/library/timeit.html
7
8 def trainjob1():
9     starttime= time.time()
10    job1.fit(xtrain,ytrain)
11    endtime = time.time()
12    print("temps entraînement 1 job : ", endtime - starttime , "sec.")
13
14 def predjob1():
15     starttime= time.time()
16     job1.predict(xtest)
17     endtime = time.time()
18     job1.score(xtest,ytest)
19     print("temps prediction 1 job : ", endtime - starttime , "sec.")
20
21 def trainjob():
22     starttime= time.time()
23     job.fit(xtrain,ytrain)
24     endtime = time.time()
25     print("temps entraînement -1 job : ", endtime - starttime , "sec.")
26
27
28 def predjob():
29     starttime= time.time()
30     job.predict(xtest)
31     endtime = time.time()
32     job.score(xtest,ytest)
33     print("temps prediction -1 job : ", endtime - starttime , "sec.")
34
35 trainjob1()
36 predjob1()
37 trainjob()
38 predjob()
39
```

Entrée []:



```
1 # quels sont les avantages de KNN ??
2 # -optimalité ?
3 # - temps de calcul ?
4 # - passage à l'échelle ?
5
6
```

Entrée []:



```
1 #Meilleur KNN
2 # K=3
3 # Toutes les ressources (njobs=-1)
4 # distance euclidienne par défaut
5
6 #ECHANTILLONS DE 10000
7 indices = np.random.randint(70000, size=20000)
8 data = mnist.data[indices]
9 target = mnist.target[indices]
10 #pour séparer le dataset en 2 échantillons de trainint et de test
11 xtrain, xtest, ytrain, ytest =train_test_split(data, target, train_size=0.8)
12
13
14 #Classifier
15 clf=neighbors.KNeighborsClassifier(3,n_jobs=-1)
16 # temps de training et prédiction mesuré
17 startTrain =time.time()
18 clf.fit(xtrain,ytrain)
19 endTrain = time.time()
20 startpred= time.time()
21 clf.predict(xtest)
22 endpred = time.time()
23 # Metrics (score, precision, recall, zero-to-one loss, temps de training et d
24 predict = clf.predict(xtest)
25 score = clf.score(xtest,ytest)
26 precision = metrics.precision_score(ytest, predict, average='macro')
27 recall = metrics.recall_score(ytest, predict, average = 'macro')
28 loss01 = metrics.zero_one_loss(ytest, predict)
29 timetrain = endTrain - startTrain
30 timePred = endpred - startpred
```

Entrée []:



```
1 print("Ce modèle de KNN, K=3, à un score de ", score*100, "%.")
2 print("4eme image : prédiction ",predict[3], "reel : ", ytest[3])
3 print ("ce KNN à une précision de", precision*100, "%.")
4 print ("ce KNN à un recall de",recall*100, "%.")
5 print ("ce KNN à un zero-one_loss de",recall*100, "%.")
6 print("    temps apprentissage : ", timetrain, "sec , temps prediction = ", ti
7
```

Entrée []:



1	
---	--