

AI for Biotechnology

Exercise 2

Prof. Dr. Dominik Grimm

Bioinformatics Research Lab

TUM Campus Straubing for Biotechnology and Sustainability

Exercise E2.1

The euclidean distance between two points $\mathbf{x}^{[a]} \in \mathbb{R}^{1 \times m}$ and $\mathbf{x}^{[b]} \in \mathbb{R}^{1 \times m}$ is defined as follows:

$$d(\mathbf{x}^{[a]}, \mathbf{x}^{[b]}) = \sqrt{\sum_{i=1}^m (\mathbf{x}_i^{[a]} - \mathbf{x}_i^{[b]})^2} \quad (1)$$

To implement Equation 1 in Python a **for** loop is needed. This could be computationally expensive when m is large (several thousand to millions of features). The dot-product between two arbitrary vectors $\mathbf{a} \in \mathbb{R}^{n \times 1}$ and $\mathbf{b} \in \mathbb{R}^{n \times 1}$ is defined as follows:

$$\mathbf{a}^\top \mathbf{b} = (a_1, \dots, a_n) \cdot \begin{pmatrix} b_1 \\ \vdots \\ b_n \end{pmatrix} = \sum_{i=1}^n a_i b_i \quad (2)$$

- a) Show that you can reformulate Equation 1 as a dot product.

Solution:

Let $\mathbf{c} \in \mathbb{R}^{m \times 1}$ be a column-vector. Then we can write:

$$\mathbf{c} = \mathbf{x}^{[a]\top} - \mathbf{x}^{[b]\top}$$

Thus, the dot product $\mathbf{c}^\top \mathbf{c} = (\mathbf{x}^{[a]\top} - \mathbf{x}^{[b]\top})^\top (\mathbf{x}^{[a]\top} - \mathbf{x}^{[b]\top})$ is equivalent to $\sum_{i=1}^m (\mathbf{x}_i^{[a]} - \mathbf{x}_i^{[b]})^2$. Thus we can rewrite the Euclidean distance as:

$$d(\mathbf{x}^{[a]}, \mathbf{x}^{[b]}) = \sqrt{\mathbf{c}^\top \mathbf{c}}$$

This expression can be written in Python in a single line and no **for** loop is needed.

- b) Implement both versions in Python (`euclidean_distance_naive` and `euclidean_distance_dot`) and test your implementations with the following code:

```
1 import numpy as np
2 import time
3
4 def euclidean_distance_naive(a,b):
5     dist = 0
6     for i in range(a.shape[1]):
7         dist += (a[0,i] - b[0,i])**2
8     dist = np.sqrt(dist)
9     return dist
10
11 def euclidean_distance_dot(a,b):
```

```

12     dist = np.sqrt(((a.T-b.T).T).dot(a.T-b.T))
13     return dist[0,0]
14
15 #Simple Test with row vectors of size 1 x 2
16 a = np.array([[3,5]])
17 b = np.array([[6,9]])
18 print(euclidean_distance_naive(a,b))
19 print(euclidean_distance_dot(a,b))
20 print()
21
22 #Test with random row vectors of size 1 x 10 Million
23 #Random numbers with seed
24 np.random.seed(42)
25 m = 10**7 #10 Million Features
26 a = np.random.random((1,m))
27 b = np.random.random((1,m))
28
29 #stop time of computation
30 start = time.process_time()
31 #compute distance
32 print(euclidean_distance_naive(a,b))
33 delta = (time.process_time()-start)
34 print("Computation Time Naive: %f s" % delta)
35
36 start = time.process_time()
37 print(euclidean_distance_dot(a,b))
38 delta = (time.process_time()-start)
39 print("Computation Time Fast: %f s" % delta)

```

Exercise E2.2

We use the `scikit-learn` package to simulate some toy data as illustrated in Figure 1:

```

1 %matplotlib inline
2 import sklearn.datasets as datasets
3 import pylab as pl
4 #Simulate Toy Dataset
5 X, y = datasets.make_circles(n_samples=150, shuffle=True,
6                               noise=0.2, random_state=42,
7                               factor=0.1)

```

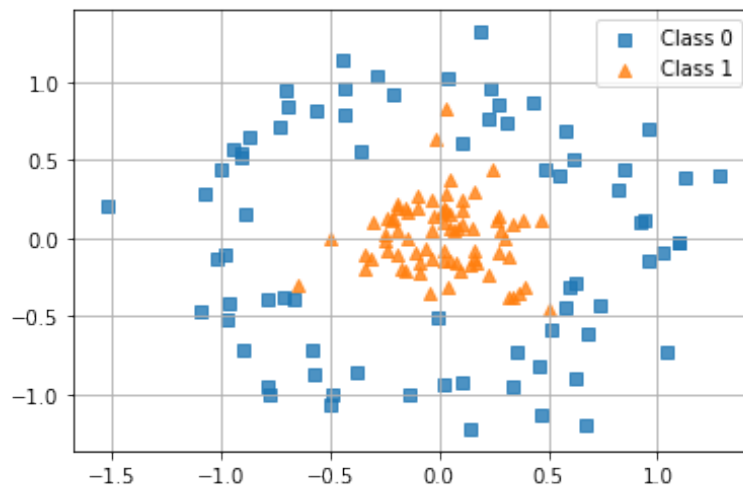


Figure 1: Toy data simulated with `scikit-learn`

- a) Write a Python snippet to determine how many samples, features and classes the dataset has. Return also the number of samples for each of the classes (hint: have a look at the NumPy Method `unique`).

```
1 #Extract Information from Toy Dataset
2 print("Total #Samples:\t\t" + str(X.shape[0]))
3 print("Number of Features:\t" + str(X.shape[1]))
4 labels, counts = np.unique(y, return_counts=True)
5 for i in range(len(labels)):
6     print("#Samples class " + str(labels[i]) + ":\t" + str(counts[i]))
```

Output:

```
Total #Samples: 150
Number of Features: 2
#Samples Class 0: 75
#Samples Class 1: 75
```

- b) Reproduce the plot shown in Figure 1 using the `Matplotlib` library (hint: have a look at the `marker` argument of the scatter function using the matplotlib documentation).

```
1 #Plot the two features for each class
2 X_class0 = X[y==0,:] #Samples for class 0
3 X_class1 = X[y==1,:] #Samples for class 1
4 pl.scatter(X_class0[:,0], X_class0[:,1],
5           alpha=0.8, label="Class 0", marker="s")
6 pl.scatter(X_class1[:,0], X_class1[:,1],
7           alpha=0.8, label="Class 1", marker="^")
8 pl.grid()
9 pl.legend()
```

- c) Add the two query points $q_1 = (-0.8, 0.3)$ and $q_2 = (-0.1, 0.1)$ to the scatter plot.

```
1 q1 = np.array([-0.8, 0.3])
2 q2 = np.array([-0.1, 0.1])
3 pl.scatter(q1[0], q1[1], alpha=1, label="Q1", marker="D")
4 pl.scatter(q2[0], q2[1], alpha=1, label="Q2", marker="P")
```

```
5 pl.grid()
6 pl.legend()
```

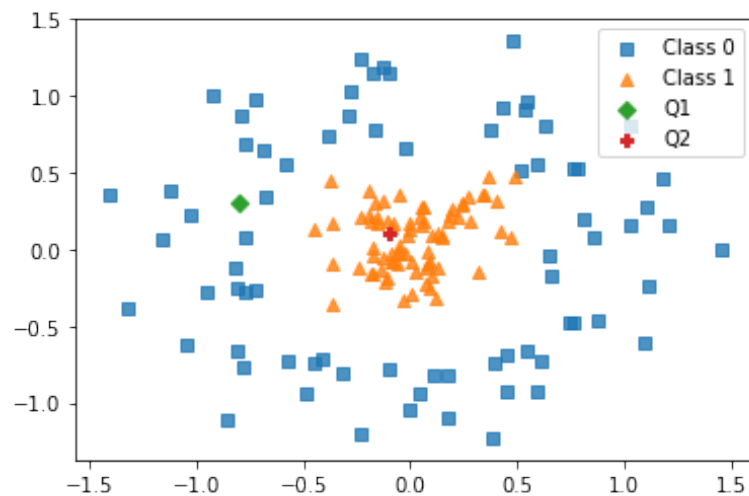


Figure 2: Toy data including query points q_1 and q_2

Exercise E2.3

Implement a function to classify a given query point q using the k-Nearest-Neighbor algorithm from scratch. Do not use the algorithms implemented in `scikit-learn`. The function should have the name `knn_predict` and should accept the arguments, `X`, `y`, `query_point` and `k`. Test your implemented function using the two query points q_1 and q_2 from the previous exercise for $k \in \{1, 10, 50, 100\}$. How are the query points classified using your implementation? (Hint: the NumPy functions `unique`, `argsort`, `argmax` and `norm` could be useful. However, there are several ways how to implement this function. You can also compare your implementation with the results from `scikit-learn`.)

```
1 #Implementation without for loops
2 def knn_predict(X,y,query_point,k=1):
3     diff = X-query_point
4     dist = np.linalg.norm(diff,axis=1)
5     knn = np.argsort(dist)[:k]
6     (label, counts) = np.unique(y[knn],return_counts=True)
7     prediction = np.argmax(counts)
8     return label[prediction]
9
10 #Alternative implementation
11 def knn_predict(X,y,query_point,k=1):
12     distances = []
13     #compute all distances
14     for i in range(X.shape[0]):
15         distance = euclidean_distance_dot(X[i,:],query_point)
16         distances.append(distance)
17     #sort distances in descending order
18     indices_closest_points = np.argsort(distances)
19     #use slicing to get the indices for the first k closest neighbors
20     knn_indices = indices_closest_points[:k]
21     #majority vote
22     class0_count = (y[knn_indices]==0).sum()
```

```
23     class1_count = (y[knn_indices]==1).sum()
24     if class0_count>class1_count:
25         return 0
26     else:
27         return 1
```

Output:

```
1  q1 = np.array([-0.8,0.3])
2  print("q1 for k1: ", knn_predict_fast(X,y,q1,k=1))
3  print("q1 for k10: ", knn_predict_fast(X,y,q1,k=10))
4  print("q1 for k50: ", knn_predict_fast(X,y,q1,k=50))
5  print("q1 for k100: ", knn_predict_fast(X,y,q1,k=100))
6  q2 = np.array([-0.1,0.1])
7  print("q2 for k1: ", knn_predict_fast(X,y,q2,k=1))
8  print("q2 for k10: ", knn_predict_fast(X,y,q2,k=10))
9  print("q2 for k50: ", knn_predict_fast(X,y,q2,k=50))
10 print("q2 for k100: ", knn_predict_fast(X,y,q2,k=100))
```

```
1  q1 for k1:  0
2  q1 for k10:  0
3  q1 for k50:  1
4  q1 for k100:  1
5  q2 for k1:  1
6  q2 for k10:  1
7  q2 for k50:  1
8  q2 for k100:  1
```