

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:

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

- Show that you can reformulate Equation 1 as a dot product.
- 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     #Write down your code here
6
7 def euclidean_distance_dot(a,b):
8     #Write down your code here
9
10 #Simple Test with row vectors of size 1 x 2
11 a = np.array([[3,5]])
12 b = np.array([[6,9]])
13 print(euclidean_distance_naive(a,b))
14 print(euclidean_distance_dot(a,b))
15 print()
16
17 #Test with random row vectors of size 1 x 10 Million
18 #Random numbers with seed
19 np.random.seed(42)
20 m = 10**7 #10 Million Features
21 a = np.random.random((1,m))
22 b = np.random.random((1,m))
23
24 #stop time of computation
25 start = time.process_time()
26 #compute distance

```

```

27 print(euclidean_distance_naive(a,b))
28 delta = (time.process_time()-start)
29 print("Computation Time Naive: %f s" % delta)
30
31 start = time.process_time()
32 print(euclidean_distance_dot(a,b))
33 delta = (time.process_time()-start)
34 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)

```

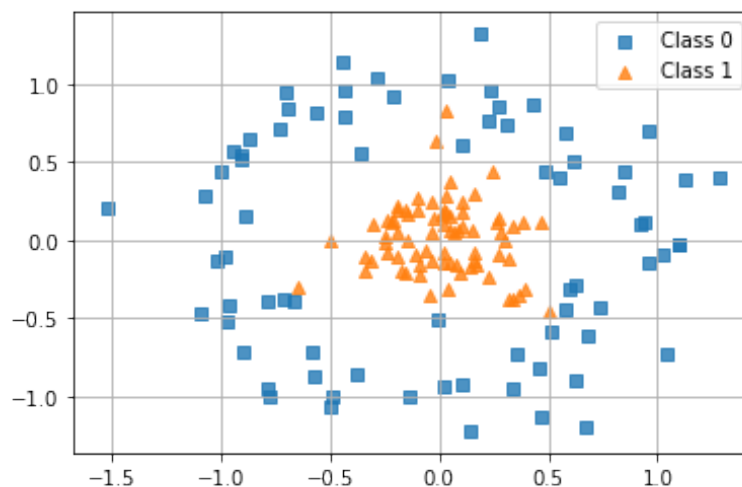


Figure 1: Toy data simulated with `scikit-learn`

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

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`.)