

Lab Course Machine Learning

Exercise 7

Prof. Dr. Dr. Lars Schmidt-Thieme,
Mofassir ul Islam Arif
Information Systems and Machine Learning Lab
University of Hildesheim
Submission: 14.12.2019 LearnWeb 3115

December 7, 2019

Instructions

Please following these instructions for solving and submitting the exercise sheet.

1. You should submit two things a) [python scripts\(zipped\)](#) / [jupyter notebook](#) and b) [a pdf document](#).
2. In the pdf document you will explain your approach (i.e. how you solved a given problem), and present your results in form of graphs and tables.
3. The submission should be made before the deadline, only through learnweb.
4. **Unless explicitly mentioned, you are not allowed to use scikit, sklearn or any other library for solve any part. All implementations must be done yourself.**

This lab has 2 components. 1st is this sheet worth 10 points. 2nd is the lecture that will take place in the lab. You need to be present in the lab in order to get those 5 points

1 Exercise Sheet 7a

Datasets

- 1. Classification Datasets: You can use one of the two datasets (or optionally, both datasets).
 - (a) Iris dataset D1: Target attribute class:Iris Setosa, Iris Versicolour, Iris Virginica
<https://archive.ics.uci.edu/ml/datasets/Iris>
 - (b) Wine Quality called D2: (use winequality-red.csv)
<http://archive.ics.uci.edu/ml/datasets/Wine+Quality>
- Note: Dataset D2 can also be used for a regression problem.

```

1 predict-knn-reg( $q \in \mathbb{R}^M, \mathcal{D}^{\text{train}} := \{(x_1, y_1), \dots, (x_N, y_N)\} \in \mathbb{R}^M \times \mathbb{R}, K \in \mathbb{N}, d$ ):
2   allocate array  $D$  of size  $N$ 
3   for  $n := 1 : N$ :
4      $D_n := d(q, x_n)$ 
5    $C := \text{argmin-k}(D, K)$ 
6    $\hat{y} := \frac{1}{K} \sum_{k=1}^K y_{C_k}$ 
7   return  $\hat{y}$ 

```

Figure 1: KNN for regression

Exercise 1: Implement K-Nearest Neighbor (KNN)

Part A:(10 Points) Your task is to implement KNN algorithm. To implement KNN you have to

- Split data into a train and a test split (70% and 30% respectively).
- Implement a similarity (or a distance) measure. To begin with you can implement the Euclidean Distance
- Implement a function that returns top K Nearest Neighbors for a given query (data point).
- You should provide the prediction for a given query (for a classification task you can use majority voting and for a regression you can use mean).
- Measure the quality of your prediction. [Hint: You have to choose a quality criterion according to the task you are solving i.e. a regression or a classification task **Defend your choice**].

Part B: (5 Points): **Determine Optimal Value of K in KNN algorithm.** In this exercise you have to provide the optimal value of K for given datasets.

- 1. How you can choose value of K for KNN. Give a criterion to choose an optimal value of K
- 2. Implement the criterion for choosing the optimal value of K.
- 3. Experimentally, give evidence that your chosen value is better than other values of K. [Hint: run your experiment with different values of K and plot the error measure for each value].

2 Algorithms

2.1 ANNEX

- You can use numpy
- You can use pandas to read and processing data

```

1 predict-knn-class( $q \in \mathbb{R}^M$ ,  $\mathcal{D}^{\text{train}} := \{(x_1, y_1), \dots, (x_N, y_N)\} \in \mathbb{R}^M \times \mathcal{Y}$ ,  $K \in \mathbb{N}$ ,  $d$ ):
2   allocate array  $D$  of size  $N$ 
3   for  $n := 1 : N$ :
4      $D_n := d(q, x_n)$ 
5    $C := \text{argmin-k}(D, K)$ 
6   allocate array  $\hat{p}$  of size  $|\mathcal{Y}|$ 
7   for  $k := 1 : K$ :
8      $\hat{p}_{C_k} := \hat{p}_{C_k} + 1/K$ 
9   return  $\hat{p}$ 

```

Figure 2: KNN for classification

```

1 argmin-k( $x \in \mathbb{R}^N$ ,  $K \in \mathbb{N}$ ) :
2   allocate array  $T$  of size  $K$ 
3   for  $n = 1 : \min(K, N)$ :
4     insert-bottomk( $T_{1:n}, n, \pi_x, 1$ )
5   for  $n = K + 1 : N$ :
6     if  $x_n < x_{T_K}$ :
7       insert-bottomk( $T, n, \pi_x, 0$ )
8   return  $T$ 
9
10 insert-bottomk( $T \in \mathcal{X}^K$ ,  $n \in \mathcal{X}$ ,  $\pi : \mathcal{X} \rightarrow \mathbb{R}$ ,  $s \in \mathbb{N}$ ) :
11    $k := \text{find-sorted}(T_{1:K-s}, n, \pi)$ 
12   for  $l := K : k + 1$  decreasing:
13      $T_l := T_{l-1}$ 
14    $T_{k+1} := n$ 

```

Note: $\pi_x(n) := x_n$ comparison by x -values. Here, $\mathcal{X} := \mathbb{N}$.

Figure 3: KNN argmin

```

1 find-sorted-linear( $x \in \mathcal{X}^K$ ,  $z \in \mathcal{X}$ ,  $\pi : \mathcal{X} \rightarrow \mathbb{R}$ ) :
2    $k := K$ 
3   while  $k > 0$  and  $\pi(z) < \pi(x_k)$ :
4      $k := k - 1$ 
5   return  $k$ 

```

- requires
 - x is sorted (increasingly w.r.t. π)
- returns smallest index k with $\pi(x_k) \leq \pi(z)$
 - 0, if $\pi(z) < \pi(x_1)$

Note: Esp. for larger K it is better to use binary search.

Figure 4: KNN sorted

- You can use matplotlib for plotting.
- You should not use any machine learning library for solving the problem i.e. scikit-learn etc. If you use them you will not get any points for the task.