

# Homework 1 Clustering and Regression

## Instructions

Answer the questions and upload your answers to courseville. Answers can be in Thai or English. Answers can be either typed or handwritten and scanned. the assignment is divided into several small tasks. Each task is weighted equally (marked with **T**). For this assignment, each task is awarded 0.9 points. There are also optional tasks (marked with **OT**) counts for 0.5 points each.

## Fun with matrix algebra

Prove the following statements. All of them can be solved by first expanding out the matrix notation as a combination of their elements, and then use the definitions of trace and matrix derivatives to help finish the proof. For example, the  $(i, j)$  element of  $Y = AB$  is  $Y_{i,j} = \sum_m A_{i,m} B_{m,j}$ .

**T1.**  $\nabla_A \text{tr} AB = B^T$

**T2.**  $\nabla_{A^T} f(A) = (\nabla_A f(A))^T$

**T3.**  $\nabla_A \text{tr} ABA^T C = CAB + C^T AB^T$

Hint: Try first solving the easier equation of  $\nabla_A \text{tr} BAC = (CB)^T = B^T C^T$

## Hello Clustering

Recall from lecture that K-means has two main steps: the points assignment step, and the mean update step. After the initialization of the centroids, we assign each data point to a centroid. Then, each centroids are updated by re-estimating the means.

Concretely, if we are given N data points,  $x_1, x_2, \dots, x_N$ , and we would like to form K clusters. We do the following;

1. **Initialization:** Pick K random data points as K centroid locations  $c_1, c_2, \dots, c_K$ .

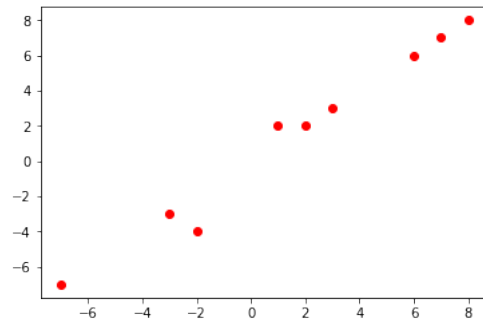
2. **Assign:** For each data point k, find the closest centroid. Assign that data point to the centroid. The distance used is typically Euclidean distance.

3. **Update:** For each centroid, calculate the mean from the data points assigned to it.

4. **Repeat:** repeat step 2 and 3 until the centroids stop changing (convergence).

Given the following data points in x-y coordinates (2 dimensional)

| x  | y  |
|----|----|
| 1  | 2  |
| 3  | 3  |
| 2  | 2  |
| 8  | 8  |
| 6  | 6  |
| 7  | 7  |
| -3 | -3 |
| -2 | -4 |
| -7 | -7 |



**T4.** If the starting points are (3,3), (2,2), and (-3,3). Describe each assign and update step. What are the points assigned? What are the updated centroids? You may do this calculation by hand or write a program to do it.

**T5.** If the starting points are (-3,-3), (2,2), and (-7,-7), what happens?

**T6.** Between the two starting set of points in the previous two questions, which one do you think is better? How would you measure the ‘goodness’ quality of a set of starting points?

In general, it is important to try different sets of starting points when doing k-means.

**OT1.** What would be the best K for this question? Describe your reasoning.

My heart will go on<sup>1</sup>



In this part of the exercise we will work on the Titanic dataset provided by Kaggle. The Titanic dataset contains information of the passengers boarding the Titanic on its final voyage. We will work on predicting whether a given passenger will survive the trip.

Let's launch Jupyter and start coding!

We start by importing the data using Pandas

```
train_url = "http://s3.amazonaws.com/assets.datacamp.com/course/Kaggle/train.csv"
train = pd.read_csv(train_url) #training set
```

```
test_url = "http://s3.amazonaws.com/assets.datacamp.com/course/Kaggle/test.csv"
test = pd.read_csv(test_url) #test set
```

Both `train` and `test` are dataframes. Use the function `train.head()` and `train.tail()` to explore the data. What do you see?

Use the function `describe()` to get a better understanding of the data. You can read the meaning of the data fields at <https://www.kaggle.com/c/titanic/data>

Looking at the data, you will notice a lot of missing values. For example, some age is NaN. This is normal for real world data to have some missing values. There are several ways to handle missing values. The simplest is to throw away any rows that have missing values. However, this usually reduce the amount of training data you have. Another method is to guess what the missing value should be. The simplest guess is to use the Median or Mode of the data. For this exercise we will proceed with this.

**T7.** What is the median age of the training set? You can easily modify the age in the dataframe by

---

<sup>1</sup>Many parts of this exercise are adapted from Kaggle Python Tutorial on Machine Learning

```
train["Age"] = train["Age"].fillna(train["Age"].median())
```

Note that you need to modify the code above a bit to fill with `mode()` because `mode()` returns a series rather than a single value.

**T8.** Some fields like ‘Embarked’ are categorical. They need to be converted to numbers first. We will represent S with 0, C with 1, and Q with 2. What is the mode of Embarked? Fill the missing values with the mode. You can set the value of Embarked easily with the following command.

```
train.loc[train["Embarked"] == "S", "Embarked"] = 0
```

Do the same for Sex.

**T9.** Write a logistic regression classifier using matrix inversion as learned in class. Use PClass, Sex, Age, and Embarked as input features. You can extract the features from Pandas to Numpy by

```
data = np.array(train[["PClass", "Sex", "Age", "Embarked"]].values)
```

Check the datatype of each values in data, does it make sense? You can force data to be the any datatype by using the command

```
data = np.array(train[["PClass", "Sex", "Age", "Embarked"]].values, dtype = float)
```

**T10.** Re-do the above process by using gradient descent instead. Are the weights learned from the two methods similar? Report the Mean Squared Errors (MSE) of the difference between the two weights.

When you evaluate the trained model on the test set, you will need to make a final decision. Since logistic regression outputs a score between 0 and 1, you will need to decide whether a score of 0.3 (or any other number) means the passenger survive or not. For now, we will say if the score is greater than or equal to 0.5, the passenger survives. If the score is lower than 0.5 the passenger will be dead. This process is often called ‘Thresholding.’ We will talk more about this process later in class.

To evaluate your results, we will use Kaggle. Kaggle is a website that hosts many machine learning competitions. Many companies put up their data as a problem for anyone to participate. If you are looking for a task for your course project, Kaggle might be a good place to start. You will need to make sure that your output is in line with the submission requirements of Kaggle: a csv file with exactly 418 entries and two columns: PassengerId and Survived. Then, use the code provided to make a new data frame using `DataFrame()`, and create a csv file using `to_csv()` method from Pandas.

To submit your prediction, you must first sign-up for an account on [Kaggle.com](https://www.kaggle.com/). Click participate to the competition at <https://www.kaggle.com/c/titanic/> then submit your csv file for the score.

The output file should have two columns: the passengerId and a 0,1 decision (0 for dead, 1 for survive). As shown below:

```
PassengerId,Survived
892,0
893,1
894,0
```

**T11.** Submit a screenshot of your submission (with the scores). Upload your code to courseville.

**OT2.** Try adding some higher order features to your training  $(x_1^2, x_1x_2, \dots)$ . Does this model has better accuracy on the training set? How does it perform on the test set?

**OT3.** What happens if you reduce the amount of features to just Sex and Age?