

ALDA Fall 2018

HW 1

8/29/2018

HW1 contains 6 questions. Please read and follow the instructions.

- **DUE DATE FOR SUBMISSION: 9/13/2018 11: 45 PM**
 - **TOTAL NUMBER OF POINTS: 100+5** (5 bonus points if you follow all the instructions and 0 otherwise)
 - Make sure you clearly list each team member's **names and Unity IDs** at the top of your submission.
 - Your submission should be a **single zip file** containing a PDF of your answers, your code, and a readme file with running instructions. Please follow the naming convention for your zip file: H(homework group number)_HW(homework number), e.g. H1_HW1.
-

1. (13 points) [**Song Ju**] Classify the following attributes as binary, discrete, or continuous. Also classify them as nominal, ordinal, interval, or ratio. Some cases may have more than one interpretation, so briefly justify your answer if you think there may be some ambiguity.
 - (a) (1 point) Hair color (Black, Blonde, Red)
 - (b) (1 point) Level of agreement (yes, maybe, no)
 - (c) (1 point) Income earned in a week
 - (d) (1 point) Celsius temperature
 - (e) (1 point) Genotype (Bb, bb, BB, bB)
 - (f) (1 point) ISBN numbers for books.
 - (g) (1 point) Time in terms of AM or PM
 - (h) (1 point) Waiting number for restaurant
 - (i) (1 point) Years of work experience
 - (j) (1 point) Categorization of clothing (hat, shirt, pants, shoes)
 - (k) (1 point) Angles as measured in degrees between 0 and 360
 - (l) (1 point) Ratings of movies (G, PG, R)
 - (m) (1 point) Coat check number. (When you attend an event, you can often give your coat to someone who, in turn, gives you a numb that you can use to claim your coat when you leave.)

2. (10 points) [**Ruth Okoilu**] Data Transformation.

In natural language processing, we often use **term frequency** and **inverse document frequency** transformation (tf'_{ij}), defined by the following equation:

$$tf'_{ij} = tf_{ij} * \log \frac{m}{df_i} \quad (1)$$

where tf_{ij} is the term frequency of the i^{th} word (term) in the j^{th} document, m is the number of documents, and df_i is the number of documents in which the i^{th} term appears.

Alternatively, we can define (tf''_{ij}) as:

$$tf''_{ij} = tf_{ij} * \log \frac{\sum_{k=1}^m d_k}{\sum_{k=1}^{df_i} d_k} \quad (2)$$

where d_k is the length of a document k .

Assume the max term frequency tf_{ij} is p and answer the following questions.

- (a) (6 points) What are the maximum and minimum values of tf'_{ij} and tf''_{ij} respectively? Please specify what cases the max and min value achieves.
- (b) (4 points) Briefly explain the purpose for using tf'_{ij} and tf''_{ij} respectively in the context of natural language processing and also explain what is the main difference between tf'_{ij} and tf''_{ij} .
3. (8 points) [**Xi Yang**] Answer the following questions:
- (a) (4 points) A healthcare dataset contains 523,000 patients. Among these patients, 26,150 patients have albinism and the remaining 496,850 patients have normal skin. Suppose we will sample 1,000 patients from the dataset to conduct albinotic analysis, which sampling method should be selected to apply in this situation: simple random sampling or stratified sampling, and why? With the selected sampling method, how many albinotic and normal skin patients will be sampled, respectively?
- (b) (4 points) Consider the following scenario, a patient's systolic blood pressure (SBP) is *recorded* to be 250. When SBP is higher than 180, a patient is considered to have hypertensive crisis and need to seek the emergency care. For this given scenario, is the recorded data noise or outlier? And why? (no point will be given if you do not give a justification).
4. (15 points) [**Song Ju**] Write your code in Matlab, R or Python to perform the following tasks, please *report your outputs and key codes* in the document file, and also include your code file (end with .m, .r or .py) in the .zip file.
- (a) (1 point) Generate a 5*5 identity matrix A.
- (b) (1 point) Change all elements in the 2nd column of A to 3.
- (c) (1 point) Sum of all elements in the matrix (use a "for/while loop").
- (d) (1 point) Transpose the matrix A ($A = A^T$)

- (e) (2 points) Calculate sum of the 3rd row, and the diagonal in the matrix A.
 - (f) (1 point) Generate a 5*5 matrix B following Gaussian Distribution with mean 5 and variance 3.
 - (g) (2 points) From B, using matrix operations to get a new matrix C such that, the first row of C is equal to the first row of B times the second row of B, the second row of C is equal to the sum of the 3rd and 4th row of B minus the 5th row of B.
 - (h) (2 points) From C, using one matrix operation to get a new matrix D such that, the first column of D is equal to the first column of C times 2, the second column of D is equal to the second column of C times 3 and so on.
 - (i) (2 points) $X = [2, 4, 6, 8]^T$, $Y = [6, 5, 4, 3]^T$, $Z = [1, 3, 5, 7]^T$. Compute the covariance matrix of X, Y and Z.
 - (j) (2 points) Verify the equation: $\bar{x^2} = (\bar{x})^2 + \sigma^2(x)$, using $x = [2, 4, 6, 8, 10, 12, 14, 16, 18, 20]^T$. $\sigma(x)$ is the standard deviation.
5. (33 points) [**Ruth Okoilu**] For this exercise, use the provided 'seeds.csv' file, which contains a list of 210 data instances. The examined group comprised kernels belonging to three different varieties of wheat: Kama, Rosa and Canadian, 70 elements each, randomly selected for the experiment. (Source: <https://archive.ics.uci.edu/ml/datasets/seeds>) There are 8 columns representing: 1) area A, 2) perimeter P, 3) compactness, 4) length of kernel, 5) width of kernel, 6) asymmetry coefficient, 7) length of kernel, and 8) groove Class (Type of wheat). For the purpose of this exercise, you consider two features, 'area.A' and 'kernel.width' (columns 1 & 5) of the provided 'seeds.csv' dataset. Write your codes in Matlab, R or Python to perform the following tasks, please *report your outputs and key codes* in the document file, and also include your code file (end with .m, .r or .py) in the .zip file.
- (a) (3 points) Load the file and read 'area.A' and 'kernel.width' columns and save them as the original *raw dataset*. Apply normalization (transformed data $z \in [0, 1]$) to the raw dataset to get the *normalized dataset* and apply the standardization to the raw dataset to get the *standardized dataset*. Show the range of the two features in each dataset.
 - (b) (30 points) Perform the following operations on the raw, normalized and standardized datasets respectively.
 - i. (3 points) Make a 2D plot of the values and label the axes (area_A should be x-axis and kernel.width should be y-axis). Compare the three plots.
 - ii. (3 points) Compute the mean of area.A and kernel.width values. Consider this point as P.
 - iii. (9 points) Compute the distance between P and the 210 data points using the following distance measures: 1) Euclidean distance, 2) Mahalanobis distance, 3) City block metric, 4) Minkowski metric (for $r=3$), 5) Chebyshev distance, 6) Cosine distance and 7) Canberra distance.
 - iv. (3 points) For each distance measure, identify the 10 points from the dataset that are the closest to the point P from (ii). (You are allowed to use any package functions to calculate the distances.)

- v. (6 points) Create plots, one for each distance measure. Place an 'X' for P and mark the 10 closest points. To mark them, you could place a circle or draw the line between these closest neighbors and the points 'X'. Make sure the points can be uniquely identified.
 - vi. (3 points) Verify if the set of points is the same across all the distance measures. If there is any big difference, briefly explain why it is.
 - vii. (3 points) Reason about your results and state the importance of data transformation in the dataset.
6. (21 points) [**Xi Yang**] In this question, please summarize and explore data in the provided file "hw1q6_data.csv", which comes from the Pima Indians Diabetes Database (<https://archive.ics.uci.edu/ml/machine-learning-databases/pima-indians-diabetes/>). In this data file, each row indicates the data for a patient. The first 6 columns are features for patients, and the last column "Class" indicates if a patient has diabetes: 1 (diabetic) or 0 (nondiabetic). The specific meaning for each feature is as follows:
- 1. *Glucose*: Plasma glucose concentration a 2 hours in an oral glucose tolerance test.
 - 2. *BloodPressure*: Diastolic blood pressure (mm Hg).
 - 3. *SkinThickness*: Triceps skinfold thickness (mm).
 - 4. *BMI*: Body mass index (weight in kg/(height in m)²).
 - 5. *DiabetesPedigreeFunction*: Diabetes pedigree function.
 - 6. *Age*: (years).

Write code in Matlab, R or Python to perform the following tasks. Please *report your outputs and key codes* in the document file, and also include your code file (end with .m, .r or .py) in the .zip file.

- (a) (1 point) How many diabetic and nondiabetic patients are in the dataset?
- (b) (2 points) There are missing values in the features which are marked as 0. What is the missing rate (%) for each feature?
- (c) (4 points) Specify two methods for missing data handling and discuss their respective advantages and disadvantages.

Remove the patients (rows) in dataset with missing values, then answer the following questions based on the **remaining data**:

- (d) (1 point) How many diabetic and nondiabetic patients are in the remaining data?
- (e) (3 points) Compute the mean, median, standard deviation, range, 25th percentiles, 50th percentiles, 75th percentiles for each feature.
- (f) (4 points) Create histogram plot using 10 bins for the two features *BloodPressure* and *DiabetesPedigreeFunction*, respectively.
- (g) (6 points) Quantile-quantile plot can be used for comparing the distribution of data against the normal distribution. Create quantile-quantile plot for the two features *BloodPressure* and *DiabetesPedigreeFunction*, respectively. Give a brief analysis for the two plots.