

CMSC 435 Assignment 2

Fall 2019

(individual work; 10 pts total)

This assignment asks you to implement and evaluate four algorithms for the imputation of missing values using two provided datasets. You will also evaluate and compare quality of the imputed values by comparing them with the corresponding values in the “complete” dataset.

Datasets

There are three datasets: the original dataset without missing values and two derived datasets where the missing values were introduced at two different amounts:

- *dataset_complete.csv* file is the complete dataset. It includes 9 features (including the class/predicted features) and 8795 objects.
- *dataset_missing05.csv* and *dataset_missing20.csv* files include the same dataset with 5% and 20% of missing values, respectively.

You will impute the missing values in each of the latter two datasets and compare these imputed values to the true/correct values that are available in the *dataset_complete.csv* file to evaluate and compare quality of different imputation algorithms. The three files are in the comma separated value (CSV) format. The first line in each file defines the names of features and the remaining lines include the values of the corresponding 8795 objects. The first 8 features are numeric and continuous with values in $[0, 1]$ interval. The last, class feature is symbolic and binary with values Yes and No. The missing values are represented by ?. There are no missing values in the class feature.

Algorithms for missing data imputation

You will implement four algorithms for the imputation of missing values and apply each of them on the corresponding two datasets that have missing values: *dataset_missing05.csv* and *dataset_missing20.csv*.

Algorithm 1. Mean imputation

Missing value for a specific feature and object is imputed with the mean value computed using the complete values of this feature.

Example

	F1	F2	F3	Class
Object 1	0.40256	0.14970	?	No
Object 2	0.41139	0.30140	?	Yes
Object 3	0.24752	0.32148	0.11169	No
Object 4	0.24609	?	0.13986	Yes
Object 5	?	0.58306	0.08910	No

To impute the missing value for feature F3 from object 1, we compute the mean of all complete values of F3: $mean = (0.11169 + 0.13986 + 0.08910) / 3 = 0.11355$.

	F1	F2	F3	Class
Object 1	0.40256	0.14970	0.11355	No
Object 2	0.41139	0.30140	?	Yes
Object 3	0.24752	0.32148	0.11169	No
Object 4	0.24609	?	0.13986	Yes
Object 5	?	0.58306	0.08910	No

The imputed values **must not** be used to compute the means. This means that all missing values for a given feature are imputed with the same mean value.

	F1	F2	F3	Class
Object 1	0.40256	0.14970	0.11355	No
Object 2	0.41139	0.30140	0.11355	Yes
Object 3	0.24752	0.32148	0.11169	No
Object 4	0.24609	?	0.13986	Yes
Object 5	?	0.58306	0.08910	No

Algorithm 2. Conditional mean imputation

Missing value for a specific feature and object is imputed with the mean value computed using the complete values of this feature for objects that satisfy a condition defined by the class feature. For instance, a missing value for an object 1 for which class = No is imputed based on the mean value computed using all objects for which class = No.

Example

	F1	F2	F3	Class
Object 1	0.40256	0.14970	?	No
Object 2	0.41139	0.30140	?	Yes
Object 3	0.24752	0.32148	0.11169	No
Object 4	0.24609	?	0.13986	Yes
Object 5	?	0.58306	0.08910	No

For object 1 for which class = No, the missing value for the feature F3 is imputed as

$$mean_{No} = (0.11169 + 0.0891) / 2 = 0.10039$$

For object 2 for which class = Yes, the missing value for the feature F3 is imputed as

$$mean_{Yes} = 0.13986 / 1 = 0.13986$$

The two imputed values are

	F1	F2	F3	Class
Object 1	0.40256	0.14970	0.10039	No
Object 2	0.41139	0.30140	0.13986	Yes
Object 3	0.24752	0.32148	0.11169	No
Object 4	0.24609	?	0.13986	Yes
Object 5	?	0.58306	0.08910	No

Algorithm 3. Hot deck imputation

Missing values for features that have missing values in a given object are imputed with the values for the same features copied from another, the most similar object. First, similarity of a given object that has missing values with every other object in the dataset is computed using the Euclidean distance. The object with the smallest distance is assumed to be the most similar and its values are used for the imputation. If that object is missing some of the values that should be imputed then the second most similar object is used to impute the remaining missing values, and so on. In other words, you should use the first complete value that you find by screening objects by their increasing values of the distance.

Given two objects $\mathbf{x} = \{x_1, x_2, \dots, x_i, \dots, x_n\}$ and $\mathbf{y} = \{y_1, y_2, \dots, y_i, \dots, y_n\}$, the Euclidean distance is

$$\text{calculated as } d(\mathbf{x}, \mathbf{y}) = \frac{\sqrt{\sum_{i=1}^m (x_i - y_i)^2}}{m} \text{ where } x_i \text{ and } y_i \text{ are values of feature } i \text{ for objects } \mathbf{x} \text{ and } \mathbf{y},$$

respectively, n is the number of features (excluding the class feature), and $m \leq n$ is the number of features that do not have missing values in either of the two objects. In other words, the distance is computed using the m features that have complete values. Note that you **must not** use the class feature in the calculation of the distance.

Example

	F1	F2	F3	Class
Object 1	0.40256	0.14970	?	No
Object 2	0.41139	0.30140	?	Yes
Object 3	0.24752	0.32148	0.11169	No
Object 4	0.24609	?	0.13986	Yes
Object 5	?	0.58306	0.08910	No

To impute missing value for feature F3 from object 1, we compute distances to every other object

$$d(obj1, obj2) = \sqrt{(0.40256 - 0.41139)^2 + (0.1497 - 0.3014)^2/2} = 0.0760$$

$$d(obj1, obj3) = \sqrt{(0.40256 - 0.24752)^2 + (0.1497 - 0.32148)^2/2} = 0.1157$$

$$d(obj1, obj4) = \sqrt{(0.40256 - 0.24609)^2/1} = 0.1565$$

$$d(obj1, obj5) = \sqrt{(0.1497 - 0.58306)^2/1} = 0.4334$$

Since object 2 that is the most similar to object 1 has a missing value for the feature F3, the second nearest object 3 is used and the missing value is imputed as follows

	F1	F2	F3	Class
Object 1	0.40256	0.14970	0.11169	No
Object 2	0.41139	0.30140	?	Yes
Object 3	0.24752	0.32148	0.11169	No
Object 4	0.24609	?	0.13986	Yes
Object 5	?	0.58306	0.08910	No

The imputed values **must not** be used to compute the distances. In other words, all missing values for each feature are imputed based on the distances that use the dataset before the imputation. This ensures that the errors inherent in the imputed values are not propagated to compute the imputation.

Algorithm 4. Conditional hot deck imputation

Missing values for features with missing values in a given object are imputed with the values for the same features copied from another, most similar object that satisfies a condition defined by the class feature. For instance, a missing value for an object 1 for which class = No is imputed based on similarity only to the objects for which class = No. The calculation of the similarities follows the unconditional hot deck imputation.

Example

	F1	F2	F3	Class
Object 1	0.40256	0.14970	?	No
Object 2	0.41139	0.30140	?	Yes
Object 3	0.24752	0.32148	0.11169	No
Object 4	0.24609	?	0.13986	Yes
Object 5	?	0.58306	0.08910	No

To impute missing value for feature F3 from object 1, we first compute the two distances to the objects that share the same value of the class feature

$$d(obj1, obj3) = \sqrt{(0.40256 - 0.24752)^2 + (0.1497 - 0.32148)^2/2} = 0.1157$$

$$d(obj1, obj5) = \sqrt{(0.1497 - 0.58306)^2/1} = 0.4334$$

The missing value is imputed with the value for the same feature F3 from the closest object 3 as follows.

	F1	F2	F3	Class
Object 1	0.40256	0.14970	0.11169	No
Object 2	0.41139	0.30140	?	Yes
Object 3	0.24752	0.32148	0.11169	No
Object 4	0.24609	?	0.13986	Yes
Object 5	?	0.58306	0.08910	No

Like for the unconditional hot deck imputation, the imputed values **must not** be used to compute the distances.

Calculation of the imputation error

You will use the two datasets that were imputed with the four methods to calculate the corresponding eight imputation errors. You will evaluate quality of these imputations based on the Mean Absolute

Error (MAE) between the imputed values and the corresponding complete values that are available in the *dataset_complete.csv* file. This dataset should be used only to calculate MAE values, not to perform the imputations. The MAE values should be used to judge and compare the quality of each imputation.

Given the imputed values $\mathbf{x} = \{x_1, x_2, \dots, x_i, \dots, x_N\}$ computed from a dataset that has missing values and the corresponding complete values $\mathbf{t} = \{t_1, t_2, \dots, t_i, \dots, t_N\}$ in the complete dataset, MAE is defined as

$$MAE = \frac{1}{N} \sum_{i=1}^N |x_i - t_i|$$

where N is the total number of missing values, x_i is a the imputed value in the dataset that has missing values, x_i and t_i are values for the same object and same feature in the two datasets, and $|\cdot|$ denotes the absolute value.

Example

Incomplete dataset

Object 1
Object 2
Object 3
Object 4
Object 5

F1	F2	F3	F4	Class
0.40256	0.14970	0.16870	?	No
0.41139	0.30140	0.47033	?	Yes
0.24752	0.32148	0.41167	0.11169	No
0.24609	?	?	0.13986	Yes
?	0.58306	0.52568	0.08910	No

Dataset where values were imputed using the unconditional mean imputation

Object 1
Object 2
Object 3
Object 4
Object 5

F1	F2	F3	F4	Class
0.40256	0.14970	0.16870	0.11355	No
0.41139	0.30140	0.47033	0.11355	Yes
0.24752	0.32148	0.41167	0.11169	No
0.24609	0.33891	0.39409	0.13986	Yes
0.32689	0.58306	0.52568	0.08910	No

Complete dataset

Object 1
Object 2
Object 3
Object 4
Object 5

F1	F2	F3	F4	Class
0.40256	0.14970	0.16870	0	No
0.41139	0.30140	0.47033	0.14175	Yes
0.24752	0.32148	0.41167	0.11169	No
0.24609	0.21359	0.24071	0.13986	Yes
0.70541	0.58306	0.52568	0.08910	No

Given the above imputation, the MAE is calculated as follows.

$$MAE = \frac{1}{5} (|0.11355 - 0| + |0.11355 - 0.14175| + |0.33891 - 0.21359| + |0.39409 - 0.24071| + |0.32689 - 0.70541|) = 0.1598$$

The MAE values must be computed with precision of **four digits** after the decimal point.

The imputed values must be computed with precision of **five digits** after the decimal point.

Implementation

Your code must perform imputation, display the eight values of MAE on the screen and save the eight imputed datasets in the csv format. The imputed datasets should be named as follows:

Vnumber_missing05_imputed_mean.csv

Vnumber_missing05_imputed_mean_conditional.csv

Vnumber_missing05_imputed_hd.csv

Vnumber_missing05_imputed_hd_conditional.csv

Vnumber_missing20_imputed_mean.csv

Vnumber_missing20_imputed_mean_conditional.csv

Vnumber_missing20_imputed_hd.csv

Vnumber_missing20_imputed_hd_conditional.csv

where *Vnumber* is your V number, e.g., *V12345678_missing01_imputed_mean.csv*

The MAE values should be displayed on the screen in the following format

```
MAE_05_mean = 0.1234
MAE_05_mean_conditional = 0.5678
MAE_05_hd = 0.1234
MAE_05_hd_conditional = 0.5678
MAE_20_mean = 0.1234
MAE_20_mean_conditional = 0.5678
MAE_20_hd = 0.1234
MAE_20_hd_conditional = 0.5678
```

You must use Java or Python3 to implement all computations including loading the datasets from the csv files, coding the four imputation methods, calculation of the MAE values, printing the MAE values on the screen, and saving of the eight imputed datasets. You may use multiple classes and functions, but they must be included in **a single source code (.java or .py) file**. This java/python file must successfully compile and produce the above mentioned outputs. Make sure to use appropriate data types to ensure the required precision of results.

If you use Java then make sure that the program can be run using the following commands:

```
javac a2.java
java a2
```

This means that your main class should be named a2 and you should not specify a package in your code.

If you use Python3 then make sure that the program can be run using the following command:

```
python a2.py
```

In both cases the program should expect that the three input csv files are located in the same working directory from which the code is run.

Deliverables

1. Java or Python3 source code in a single .java or .py file. The file **must be named** a2.java or a2.py.
2. A pdf document (**named** with your last name, e.g., Smith.pdf) with answers to the following four questions
 - 2a. What are the MAE values for the eight results? You should copy the output from the screen.
 - 2b. Which of the considered four imputation methods is the least accurate for the *dataset_missing05.csv* dataset? Briefly explain **why** this method is worse than the other three methods.
 - 2c. Are the results for the same algorithm on the two datasets the same/worse/better (e.g. mean on *dataset_missing05* vs. mean on *dataset_missing20*)? Briefly explain **why**.
 - 2d. Which of the two unconditional methods (mean vs. hot deck) is faster, i.e., requires fewer computations? Briefly explain **why**. Give computational complexity of both methods as a function of the number of objects n , and use it to support your explanation.

Notes

- Include your **name, class number and title** at the top of the pdf file. Use separate and **clearly marked and numbered paragraphs** for answers to each of the four questions. We will deduct points if you do not follow these instructions.
- Copy the submission email to yourself to have a proof for the time of your submission.

- Do not procrastinate and start early – this assignment requires a substantial amount of time and effort. Late submissions will be subject to deductions: 15% in first 12 hours and 30% for between 12 and 48 hours. We will not accept submissions that are over 48 hours late.
- We will check if the source code runs correctly, validate the results on the screen and in the files, and mark the answers to the four questions.
- We will **deduct** points if the files names and/or the outputs on the screen do not follow the above defined format.
- We will check for **plagiarism**. Make sure to write your own code and provide your own answers.

Due Date

Your assignment must be received by 12:45pm Eastern Time (beginning of the lecture) on September 24 (Tuesday), 2019. Send the two files (remember to name them as explained in the deliverables section) in a single email to the class TA, Mr. Sina Ghadermarzi, at ghadermarzis@mymail.vcu.edu. The email should be entitled “CMSC 435 assignment 2 submission”.