

# Assignment #3

Monday, September 25, 2023 4:17 PM



11260658

## Data Mining

### Assignment Three

10 points

*Due: Thursday, October 5 @ 11:50 PM*

Modify Tutorial\_4\_Data\_Preprocessing.ipynb to do a few steps with the attached university salary and faculty size data set (aaup.csv).

Here is a description of the variables in the dataset.

Univ\_id: id number  
Univ\_name: Name of institution  
State: 2 letter state code  
Type: (I, IIA, or IIB)  
fp\_sal: Average salary - full professors  
ac\_sal: Average salary - associate professors  
at\_sal: Average salary - assistant professors  
to\_sal: Average salary - all ranks  
fp\_com: Average compensation - full professors  
ac\_com: Average compensation - associate professors  
at\_com: Average compensation - assistant professors  
to\_com: Average compensation - all ranks  
fp\_#: Number of full professors  
ac\_#: Number of associate professors  
at\_#: Number of assistant professors  
in\_#: Number of instructors  
to\_#: Number of faculty - all ranks

#### Tasks

Replicate the preprocessing steps applied to the breast cancer example as guided below:

1. Input the data into a Pandas dataframe; create the data columns of your choice; print the number of observations and attributes.
2. Recode the missing values to NaN. This dataset uses \*. Print the counts of missing values across the attributes.
3. How do you handle missing values in this dataset? Explain your selection. (put your answer in the box below)

The missing values in the dataset are not considered and therefore dropped. By doing so, the attributes we are analyzing provide some numeric attribute that would be helpful in finding meaningful data within the dataset. I think using median did not feel appropriate for NaN did not feel appropriate in this circumstance.

4. Explore for outliers. Apply the boxplot display. Are there any outliers? (put responses in the box below)

Yes, there are outliers present in the dataset across all attributes available. The outliers are beyond the 'max' within a boxplot.

5. Are there any duplicate records? No, there are not any duplicate records present in the dataset.

6. Can you aggregate the institutions within each state using the grouping operation from Pandas? So you should end up with ~50 observations. Which statistics are you aggregating on? (describe in the box below)

Yes, you can aggregate the number of institutions within each state. After doing so, I computed the mean to find the average salary or compensation, in each of the respective attributes that were selected.

7. Explore some sampling from the original data set, not the aggregate. What did you find to be the best? Why? (put answer in the box below)

Because there are no duplicates within our data, it would be not of our concern to consider sampling with replacement. The best sampling method would be sampling without replacement, but randomly selecting a certain number of rows of data to be shown. The dataset contains 1100 rows of data, which is not considered large in the machine learning realm, as well as that the dataset can be computationally ran on our device without much trouble.

8. Pick a salary column to discretize and pick a count to discretize. Why did you choose your type of discretization? (put in the box below)

I decided to apply two discretization methods to the 'fp\_sal' and 'fp\_#' columns: equal width and equal frequency. I used equal width discretization to create bins with consistent salary ranges, showing uniformity in the width of each category. I also experimented with equal frequency discretization to distribute data points more evenly across bins, making it a valuable for identifying groups with balanced representation. Equal width prioritizes uniformity and equal frequency makes sure that each category contains a similar number of data points. My decision to use one method over the other depended on whether I wanted consistent range widths or balanced representation in my analysis. If I had to choose, I would use equal width discretization as it provides a simpler and more intuitive way to understand salary distributions by creating consistent bins with fixed ranges.

**Submission:** Submit both the .ipynb file and the word document.

Good luck!