

# DATA MINING

## ASSIGNMENT 2

Ali Gholami

Department of Computer Engineering & Information Technology  
Amirkabir University of Technology

<http://ceit.aut.ac.ir/~aligholamee>  
[aligholamee@aut.ac.ir](mailto:aligholamee@aut.ac.ir)

### Abstract

In this assignment, several paramount concepts of *Data Analysis* will be explained. we'll discuss the importance of metrics in the first theoretical problem. A quick review on the *Apriori* algorithm for the *Association Rule Mining* will be explained also. We'll also show how *Weka* can be used for *Association Rule Mining*. Furthermore, The effectiveness of *Normalization* concept is proposed. Finally, an *Statistical* point of view will help us to demonstrate and rationalize the relationship between the *Performance* of the *Learning Algorithm* and the amount of *Data* available. A chief section of this assignment is dedicated to solve the *Titanic* problem, which is a great practice of data mining concepts in production. We'll use *Python* programming language and three main libraries; *Scikit-Learn*, *Pandas* and *Numpy* to tackle this problem. The Python implementation of the Titanic problem is provided on a *Jupyter Notebook* attached with this report.

**Keywords.** *Apriori, Association Rule Mining, Normalization, Generalization, Preprocessing, Feature Engineering, Scikit-Learn, Pandas, Numpy, Python 3.5.*

## 1 Data Preprocessing

In this section, we'll be looking at our training data from different aspects. First, we need to get a quick intuition of how data looks like, how is that distributed and what to do with that! To do this, we'll be using some functions as described below.

---

```
separate_output('Training Data Types')  
print(train_data.dtypes)
```

---

In this part, we have printed the data types of our training set. Note that *separate-output* is a self-defined function to make thing more clear in the terminal. Now, its time for some statistics. To get a full understanding of how our numerical data is distributed, we use the following code.

---

```
separate_output('Statistical Information')  
print(train_data.describe())
```

---

The result of this part of code will be some statistical parameters such as: *variance*, *mean*, *max*, *min*, *counts*. These can be useful in the future to make decisions about *data normalization*. Another amazing feature that *Pandas* has provided for us is the ability to separately describe each column in the dataset. As an example, the first column contains 686 missing values. Use the following code to believe this fact.

---

```
separate_output('Counts Values on a Column')
print(train_data['col_1'].value_counts())
```

---

This column may not be seem much useful at the first glance, but we keep it since there are some good values for that column which might make it useful while we go further in the classification task. Some of these columns are completely useless. Let's find them. The following function will return a dictionary consisting of number of missing values of each column.

---

```
def compute_nans(df):

    nans_dict = {}

    for col in df:
        nan_col_counter = 0
        for row in df[col]:
            if(row == '?'):
                nan_col_counter += 1
        nans_dict[str(col)] = [nan_col_counter]

    return nans_dict
```

---