**Ain Shams University**

**Faculty of Computer and information science**

**Machine learning**

**Team ID: SC_7**

**Project name: Hotel Rating Prediction**

**Project Description:**

Predicting the reviewer score on a hotel depending on some features.

| Name | Section | ID |
|---|---|---|
| احمد اسامة عبدالرزاق حسن (team leader) | 1 | 20201700020 |
| االاء مصطفى صادق على | 2 | 20201700137 |
| شهد اشرف احمد عبداللاه | 4 | 20201700407 |
| شهد محمد كامل مصطفي | extra | 20201700409 |
| احمد اديب معضماني | 1 | 20201701044 |
| مصعب ثابت محمد عبد الآخر | 7 | 20201700848 |

## 1- The Preprocessing:

We worked on every column.

- **(lat & lng -> address)**
  After trying to find a pattern in Hotel Address column we decided to take the latitude and longitude and connect the model to google maps to return an object that has the city, country to put them in columns to the next stage.

- **(tags)**
  Splitting on the "," and dropping the extra punctuation marks then putting each of the room type, trip type, nights spent and the people in a new column to later then encode.

- **(Trip Type)**

  After splitting the "tags" column into different categories, we found that the "trip type" category had only two values: "Leisure trip" and "Business trip". Some rows were empty, so we performed logistic regression to classify the values based on the provided tags into the aforementioned two classes.

- **(Review_Date)**
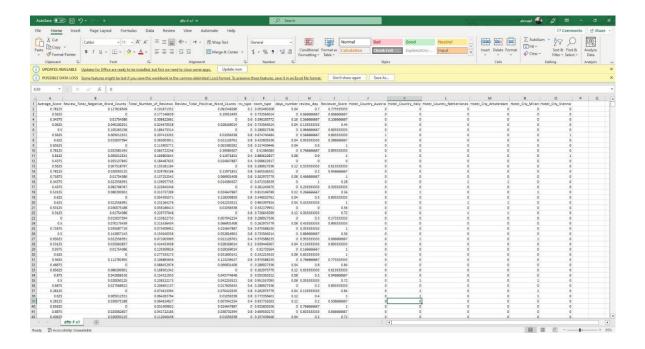  Turning it to date time then putting day, month and year then putting each into a column to encode.

After that dropping the nulls, we moved to the next part.

**Encoding:**

Using Label and one hot Encoding to turn all our data to numerical to use them properly.

**Scaling:**

Re scaling the data from 0 to 1 so that all columns become on one scale.

## 2- Feature Selection (Correlation):

our data is numerical and categorical so after research we found that the best methods to handle them is Pearson and ANOVA correlation.

- ANOVA Correlation was applied on room type and year of the review as they are categorical and we want the output to be numerical.
- As for Pearson we tried it on the rest of the columns as they are numerical and we want the output to be numerical.

## 3- Models:

We tried more than one model to choose the best of them according to accuracy.

| Model name | MSE | Accuracy |
|---|---|---|
| Linear | 0.03357117405195145 | 0.2826225059200652 |
| Random Forest | 0.030295117059987713 | 0.3526280872476978 |
| SVR | 0.034258421636806695 | 0.2679368130851052 |