

ASSIGNMENT 2 - DATA ANALYTICS APPLICATION

UK SALARY PREDICTION MODEL BASED ON RANDOM FOREST REGRESSOR

GROUP NUMBER 10

MEMBER DETAILS:

- TEENU PRATHYUSH (GROUP COORDINATOR)

STUDENT ID:21262966

EMAIL ID:teenu.prathyush2@mail.dcu.ie

- NIRANJAN NIMBARTE

STUDENT ID:21264812

EMAIL ID:niranjn.nimbarte2@mail.dcu.ie

- TEENA SHARMA

STUDENT ID:21261593

EMAIL ID:teena.sharma5@mail.dcu.ie

- FIAZ ASHRAF

STUDENT ID:21263700

EMAIL ID:fiaz.ashraf2@mail.dcu.ie

- DIVYASHREE PARAMESHWARA

STUDENT ID:20211141

EMAIL ID:divyashree.parameshwara2@mail.dcu.ie

- ANISH PATIL

STUDENT ID:21260453

EMAIL ID:anish.patil4@mail.dcu.ie

LINK FOR THE DATASET

Dataset Link: <https://www.kaggle.com/c/job-salary-prediction/data>

1. INTRODUCTION AND MOTIVATION:

Salary plays a vital role in determining career progression as well as in determining occupational success. Additionally, the pay rate is also considered an influential motivator. With the increasing number of jobs available, it becomes even more crucial to determine how much a particular job can offer in monetary terms. Returns from investments in higher education have a strong correlation with salary outcomes and, hence, influence a student's decision in choosing a degree program or field of study.

Recent trends show that there has been a significant rise in the number of jobs available in the market in a wide range of different domains. It is possible to find a huge number of Job roles on e-Recruitment websites while applying for a job. Still, the underlying problem is that salary insight is often hidden or not readily available to job seekers. To find a solution to this problem, we propose an application that can be used to determine the salaries of the various roles in the dataset. This could be useful for gaining valuable insights. In our application, we analyze the United Kingdom job salary dataset from Kaggle and utilize the **Random Forest Regressor** to predict salaries. These predictions would be helpful for people living in the United Kingdom and those who live overseas who want to land a job in the U.K. Our web application is developed using the Flask framework and integrates machine learning as a key component for making predictions. Out of many possible models in machine learning, we have chosen the **Random Forest Regressor** model since it provided superior results.

2. RELATED WORK:

We can see different types of applications on the internet considering salary prediction. One of them is Morgan McKinley for salary prediction applications. This application predicts the salary for employees from lots of different disciplines which also can influence the prediction based on different cities in Ireland. There is also an application we found similar to our project which focuses on the challenge of predicting the salary offered by companies through job posts on the web. This application analyses job advertisements collected from Tecnoempleo, an e-Recruitment website specialized in information technology jobs in Spain. The classifiers implied on this application based on ensembles of decision trees (AB and RF) and based on voting ensembles (Vote and Vote3) achieve the best accuracy. The application's average accuracy is ≈ 0.84 , with the voting classifiers that lead to a slightly better median accuracy (≈ 0.841 for Vote and ≈ 0.85 for Vote3).

3. DESCRIPTION OF THE DATASET:

Link to the dataset: <https://www.kaggle.com/c/job-salary-prediction/data>

The dataset we used in our project was downloaded directly from Kaggle. The dataset consists of two parts - The training set and Test Set which is used for the Random Forest Regressor Model. We use the Training set for building our machine learning model.

In the training set, we have approximately 240,000 rows and 12 columns, but for our model we only require 7 columns, therefore, we removed the 5 unusable ones that are Id, FullDescription, LocationRaw, SalaryRaw, SourceName. The 7 columns that are used in our model are described below.

- **Title** column provides the designation of the particular employee as per the data.
- **LocationNormalized** is a normalization of the LocationRaw column that provides information about the city in which the company is established.
- **ContractType** column gives us information on whether the employee is working on a contract basis or working in a permanent job.
- **ContractTime** column gives us information on whether the employee is working full-time or part-time.
- **Category** column provides the field of the job for eg. Engineering, IT, etc
- **SalaryNormalized** is a normalization of the SalaryRaw column that provides the salary of the employee. This is also the target variable of our machine learning model where we will predict the salary for an employee looking for a job.
- **Company** column provides information about all the companies based in the UK.

4. DESCRIPTION OF DATA PROCESSING:

Link to the Pig Script: https://gitlab.com/computing.dcu.ie/prathyt2/ca675-assignment-2_group-10/-/blob/master/Code/Pigscript.pig

The dataset contains a large number of null values. It is necessary to clean the data before training it with the machine learning model. We have chosen Pig Technology on a Dataproc Cluster (in Google Cloud Platform) for cleaning the dataset. Using Pig, we have selected only the necessary columns and removed all the entries containing null values in the ContractTime and Company Columns. The ContractType column also contains many null values but based on the information present in the dataset we have normalized the missing values in this column.

The majority of the data processing/cleaning has been done on Pig. After cleaning the data on Pig we have used R programming to remove some of the messy characters in the title column. The rest of the columns are clean and do not require any processing.

5. TECHNOLOGIES USED:

Front-end: HTML , CSS , Bootstrap , jQuery, JavaScript.
 Back-end: Python
 Machine Learning Model: Random Forest Regressor
 Web Framework: Flask
 Libraries: scikit-learn, NumPy, pandas, flask-cors, pickle
 Other applications used: Microsoft Excel, R-Studio, Jupyter Notebook in GCP.

The web application is hosted on the Google AppEngine platform present in the GCP. Google Cloud SDK is used to manage the resources and deploy the application.

6. DEVELOPMENT OF THE APPLICATION PLATFORM:

The following steps have been taken in developing the application:

- The dataset is obtained from the Kaggle website.
- The dataset is then imported into the Google Cloud Staging Bucket present in GCP.
- We load the data into Pig and then process/clean the data.
- The processed data is then loaded into the Machine Learning model (Random Forest Regressor) for training.
- The trained machine learning model produces an accuracy score of 78% and is then saved as a pickle file (serializing into byte stream) for later use.
- The front-end of the application is developed using HTML, CSS, JavaScript, jQuery, and Bootstrap.
- A **main.py** file is then created to build the application using the front-end html file and integrate the trained machine learning model, which was previously saved as a pickle file.
- Flask framework is used for developing the application.
- The application is then hosted on the Google AppEngine platform in GCP using the Google Cloud SDK shell.
- The estimated salary is obtained by entering the values in the given form.

7. RELEVANT SCREENSHOTS:

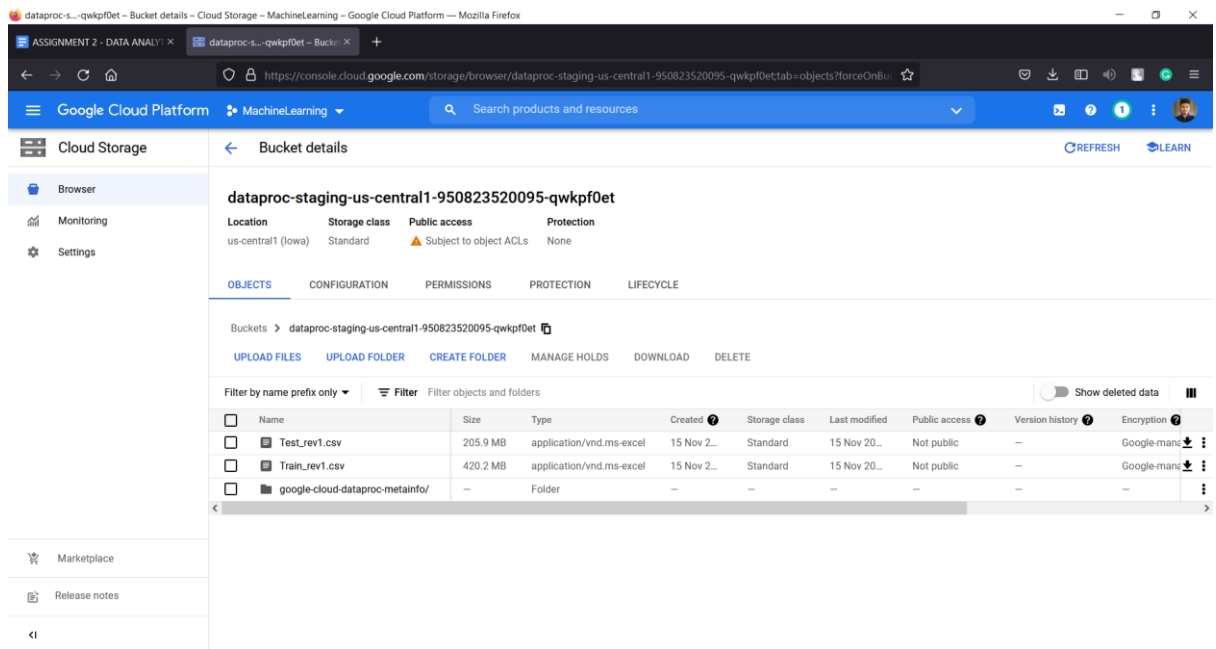


Figure 1: Cloud Storage - Google Cloud Staging Bucket

Figure 1 displays the datasets imported into the Google Cloud Staging Bucket for further processing.

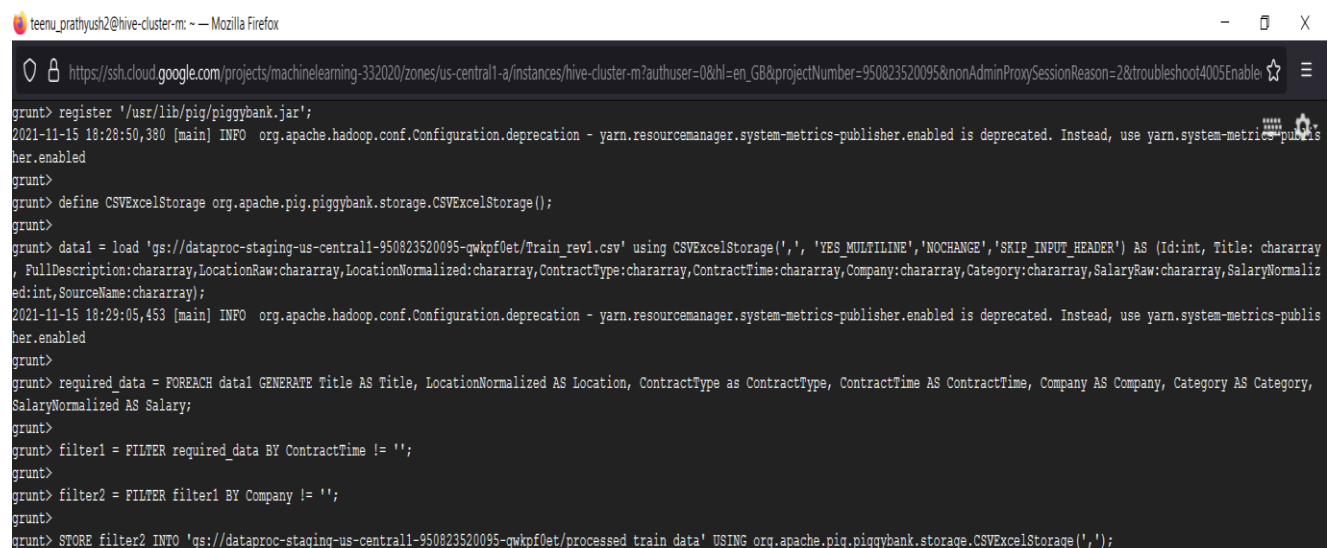


Figure 2: Pig Commands

Figure 2 displays the commands run in the Pig shell for processing/cleaning the data

```
teenu.prathyush2@hive-cluster-m: ~ -- Mozilla Firefox
https://ssh.cloud.google.com/projects/machinelearning-332020/zones/us-central1-a/instances/hive-cluster-m/authtoken=0&hl=en_GB&projectNumber=950823520095&nonAdminProxySessionReason=2&troubleshoot4005enable

2021-11-15 18:30:36,777 [main] INFO org.apache.hadoop.mapred.ClientServiceDelegate - Application state is completed. FinalApplicationStatus=SUCCEEDED. Redirecting to job history server
2021-11-15 18:30:37,422 [main] INFO org.apache.hadoop.yarn.client.RMProxy - Connecting to ResourceManager at hive-cluster-m/10.128.0.3:8032
2021-11-15 18:30:37,423 [main] INFO org.apache.hadoop.yarn.client.AHSProxy - Connecting to Application History server at hive-cluster-m/10.128.0.3:10200
2021-11-15 18:30:37,428 [main] INFO org.apache.hadoop.mapred.ClientServiceDelegate - Application state is completed. FinalApplicationStatus=SUCCEEDED. Redirecting to job history server
2021-11-15 18:30:37,449 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.reduce.tasks is deprecated. Instead, use mapreduce.job.reduces
2021-11-15 18:30:37,451 [main] INFO org.apache.hadoop.yarn.client.RMProxy - Connecting to ResourceManager at hive-cluster-m/10.128.0.3:8032
2021-11-15 18:30:37,451 [main] INFO org.apache.hadoop.yarn.client.AHSProxy - Connecting to Application History server at hive-cluster-m/10.128.0.3:10200
2021-11-15 18:30:37,456 [main] INFO org.apache.hadoop.mapred.ClientServiceDelegate - Application state is completed. FinalApplicationStatus=SUCCEEDED. Redirecting to job history server
2021-11-15 18:30:37,523 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - 100% complete
2021-11-15 18:30:37,527 [main] INFO org.apache.pig.tools.pigstats.mapreduce.SimplePigStats - Script Statistics:

HadoopVersion  PigVersion  UserId  StartedAt  FinishedAt  Features
2.9.2  0.17.0  teenu.prathyush2  2021-11-15 18:30:04  2021-11-15 18:30:37  FILTER

Success!

Job State (time in seconds):
JobId  Maps  Reduces  MaxMapTime  MinMapTime  AvgMapTime  MedianMapTime  MaxReduceTime  MinReduceTime  AvgReduceTime  MedianReduceTime  Alias  Feature  Outputs
job_163699988049_0001  4  0  16  8  14  15  0  0  0  0  data,filter,required_data  MAP_ONLY  gs://dataproc-staging-us-central1-9508235200

Input(s):
Successfully read 244769 records (1588 bytes) from: "gs://dataproc-staging-us-central1-950823520095-qwkp0et/Train_rev1.csv"

Output(s):
Successfully stored 161888 records in: "gs://dataproc-staging-us-central1-950823520095-qwkp0et/processed_train_data"

Counters:
Total records written : 161888
Total bytes written : 0
Spillable Memory Manager spill count : 0
Total bags proactively spilled: 0
Total records proactively spilled: 0

Job DAG:
job_163699988049_0001

2021-11-15 18:30:37,530 [main] INFO org.apache.hadoop.yarn.client.RMProxy - Connecting to ResourceManager at hive-cluster-m/10.128.0.3:8032
2021-11-15 18:30:37,531 [main] INFO org.apache.hadoop.yarn.client.AHSProxy - Connecting to Application History server at hive-cluster-m/10.128.0.3:10200
2021-11-15 18:30:37,536 [main] INFO org.apache.hadoop.mapred.ClientServiceDelegate - Application state is completed. FinalApplicationStatus=SUCCEEDED. Redirecting to job history server
2021-11-15 18:30:37,568 [main] INFO org.apache.hadoop.yarn.client.RMProxy - Connecting to ResourceManager at hive-cluster-m/10.128.0.3:8032
2021-11-15 18:30:37,569 [main] INFO org.apache.hadoop.yarn.client.AHSProxy - Connecting to Application History server at hive-cluster-m/10.128.0.3:10200
2021-11-15 18:30:37,574 [main] INFO org.apache.hadoop.mapred.ClientServiceDelegate - Application state is completed. FinalApplicationStatus=SUCCEEDED. Redirecting to job history server
2021-11-15 18:30:37,591 [main] INFO org.apache.hadoop.yarn.client.RMProxy - Connecting to ResourceManager at hive-cluster-m/10.128.0.3:8032
2021-11-15 18:30:37,591 [main] INFO org.apache.hadoop.yarn.client.AHSProxy - Connecting to Application History server at hive-cluster-m/10.128.0.3:10200
2021-11-15 18:30:37,604 [main] INFO org.apache.hadoop.mapred.ClientServiceDelegate - Application state is completed. FinalApplicationStatus=SUCCEEDED. Redirecting to job history server
2021-11-15 18:30:37,627 [main] WARN org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Encountered Warning FIELD_DISCARDED_TYPE_CONVERSION_FAILED 1 time(s).
2021-11-15 18:30:37,627 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!

grant;
```

Figure 3: Pig-processed data output

Figure 3 - The processed data is stored back into the Google Cloud Staging bucket using the STORE command. After removing the null values we are left with 160,000 rows (approx.) out of 240,000 rows(approx.) in total. The same procedure is repeated for the test dataset.

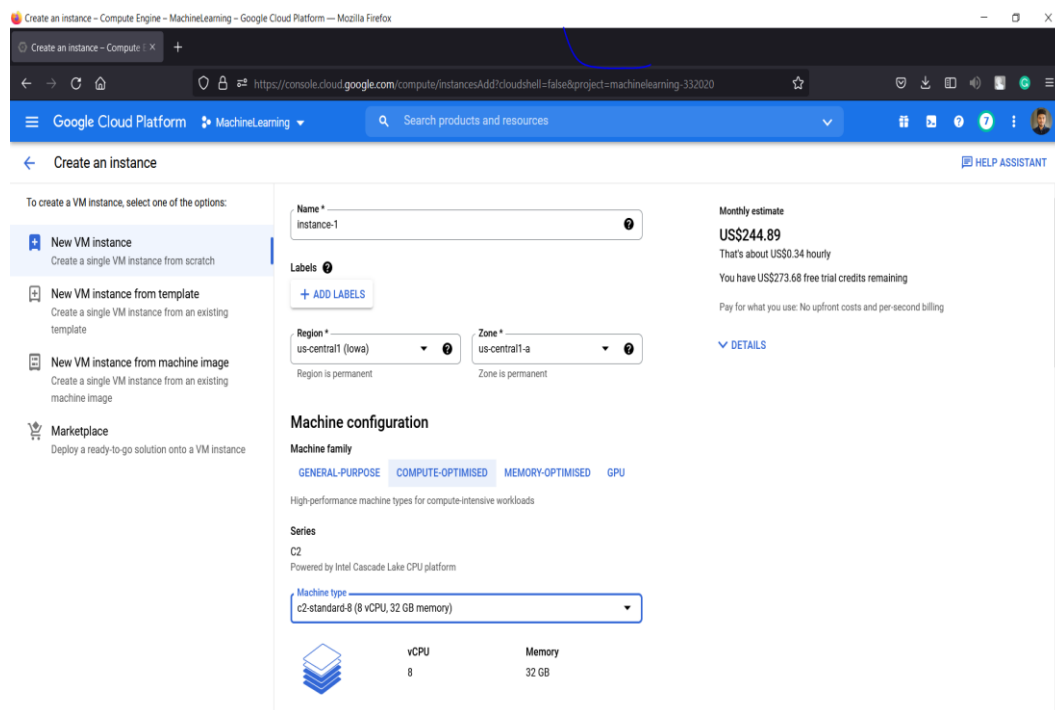


Figure 4: Creating an instance on Google Cloud Platform

Figure 4 displays the process of creating an instance for using Jupyter Notebook on the Cloud.

```

teenu_prathyush2@instance-1:~$ wget http://repo.continuum.io/archive/Anaconda3-4.0.0-Linux-x86_64.sh
--2021-11-14 22:36:25-- http://repo.continuum.io/archive/Anaconda3-4.0.0-Linux-x86_64.sh
Resolving repo.continuum.io (repo.continuum.io)... 104.18.200.79, 104.18.201.79, 2606:4700::6812:c84f, ...
Connecting to repo.continuum.io (repo.continuum.io)|104.18.200.79|:80... connected.
HTTP request sent, awaiting response... 301 Moved Permanently
Location: https://repo.anaconda.com/archive/Anaconda3-4.0.0-Linux-x86_64.sh [following]
--2021-11-14 22:36:25-- https://repo.anaconda.com/archive/Anaconda3-4.0.0-Linux-x86_64.sh
Resolving repo.anaconda.com (repo.anaconda.com)... 104.16.131.3, 104.16.130.3, 2606:4700::6810:8203, ...
Connecting to repo.anaconda.com (repo.anaconda.com)|104.16.131.3|:443... connected.
HTTP request sent, awaiting response... 200 OK
Length: 417798602 (398M) [application/x-sh]
Saving to: 'Anaconda3-4.0.0-Linux-x86_64.sh'

Anaconda3-4.0.0-Linux-x86_64 100%[=====>] 398.44M  149MB/s  in 2.7s

2021-11-14 22:36:28 (149 MB/s) - 'Anaconda3-4.0.0-Linux-x86_64.sh' saved [417798602/417798602]

teenu_prathyush2@instance-1:~$

```

Figure 5: Installing Anaconda on the created instance

Figure 5 displays the installation process of Anaconda, which is required for using Jupyter Notebook.

```

teenu_prathyush2@instance-1:~$ jupyter notebook --generate-config
Writing default config to: /home/teenu_prathyush2/.jupyter/jupyter_notebook_config.py
teenu_prathyush2@instance-1:~$ sudo nano /home/teenu_prathyush2/.jupyter/jupyter_notebook_config.py
teenu_prathyush2@instance-1:~$ jupyter-notebook --no-browser --port=5000
[I 22:47:59.513 NotebookApp] Writing notebook server cookie secret to /run/user/1000/jupyter/notebook_cookie_secret
[W 22:47:59.567 NotebookApp] WARNING: The notebook server is listening on all IP addresses and not using encryption. This is not recommended.
[W 22:47:59.567 NotebookApp] WARNING: The notebook server is listening on all IP addresses and not using authentication. This is highly insecure and not recommended.
[I 22:47:59.571 NotebookApp] Serving notebooks from local directory: /home/teenu_prathyush2
[I 22:47:59.571 NotebookApp] 0 active kernels
[I 22:47:59.571 NotebookApp] The Jupyter Notebook is running at: http://[all ip addresses on your system]:5000/
[I 22:47:59.571 NotebookApp] Use Control-C to stop this server and shut down all kernels (twice to skip confirmation).


```

Figure 6: Launching Jupyter Notebook

Figure 6 displays the process of launching Jupyter Notebook on Google Cloud Platform which will be used to create our machine learning model.

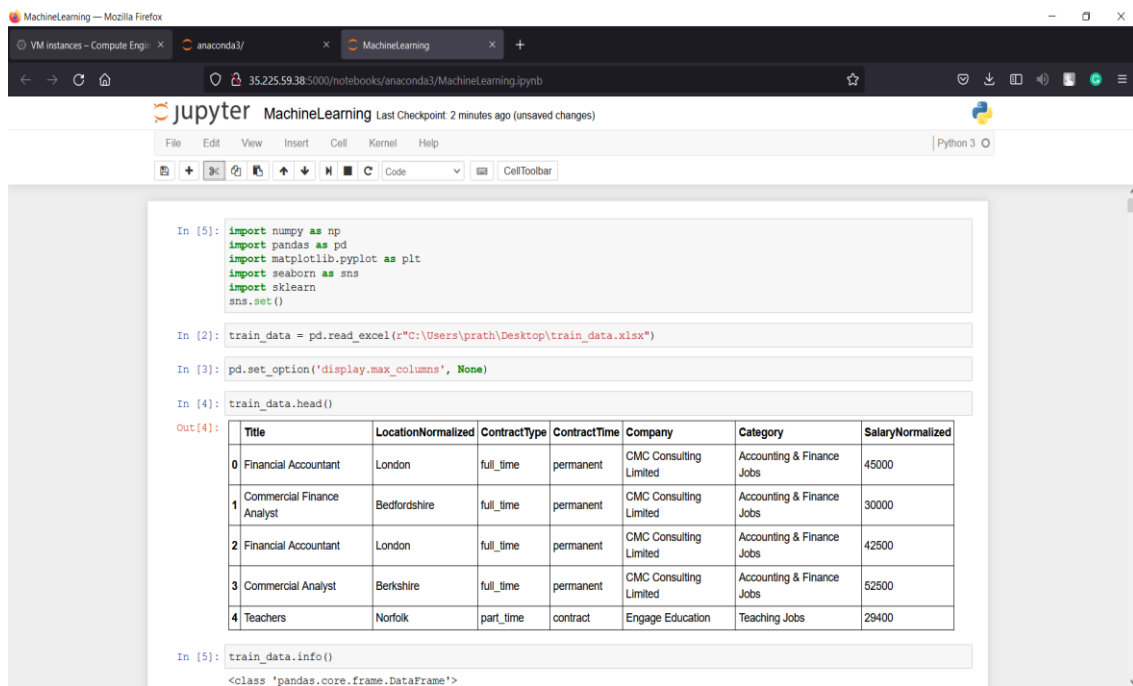


Figure 7: Creating the Machine Learning model

Figure 7 displays the creation of the machine learning model on Jupyter Notebook in Google Cloud Platform.

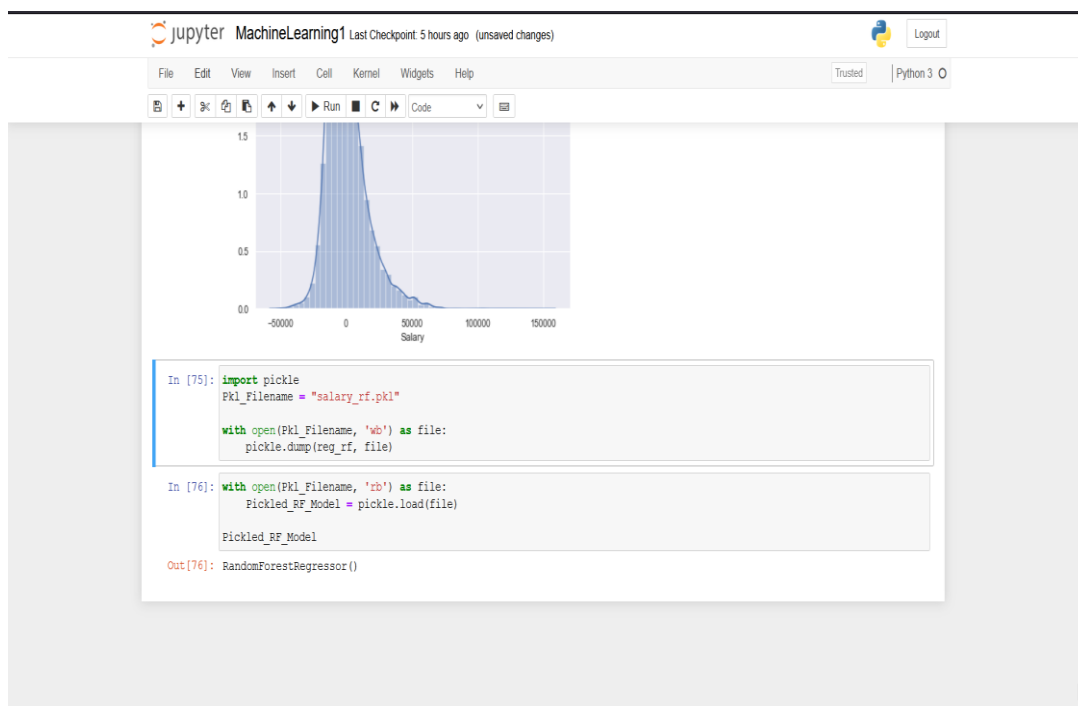


Figure 8: Serializing the machine learning model into a Pickle file

Figure 8 - The trained random forest regression model is saved as a pickle file for later use


```
main - Notepad
File Edit Format View Help
from flask import Flask, request, render_template
from flask_cors import cross_origin
import sklearn
import pickle
import pandas as pd

app = Flask(__name__)
model = pickle.load(open("salary_rf.pkl", "rb"))

@app.route("/")
@cross_origin()
def home():
    return render_template("home.html")

@app.route("/predict", methods = ["GET", "POST"])
@cross_origin()
def predict():
    if request.method == "POST":
        Title=request.form['Title']
```

Figure 9: Loading the pickle file into main.py

Figure 9 displays the code snippet of **main.py** which is used to load the machine learning model (pickle file) and also implements the flask framework for developing the web application.

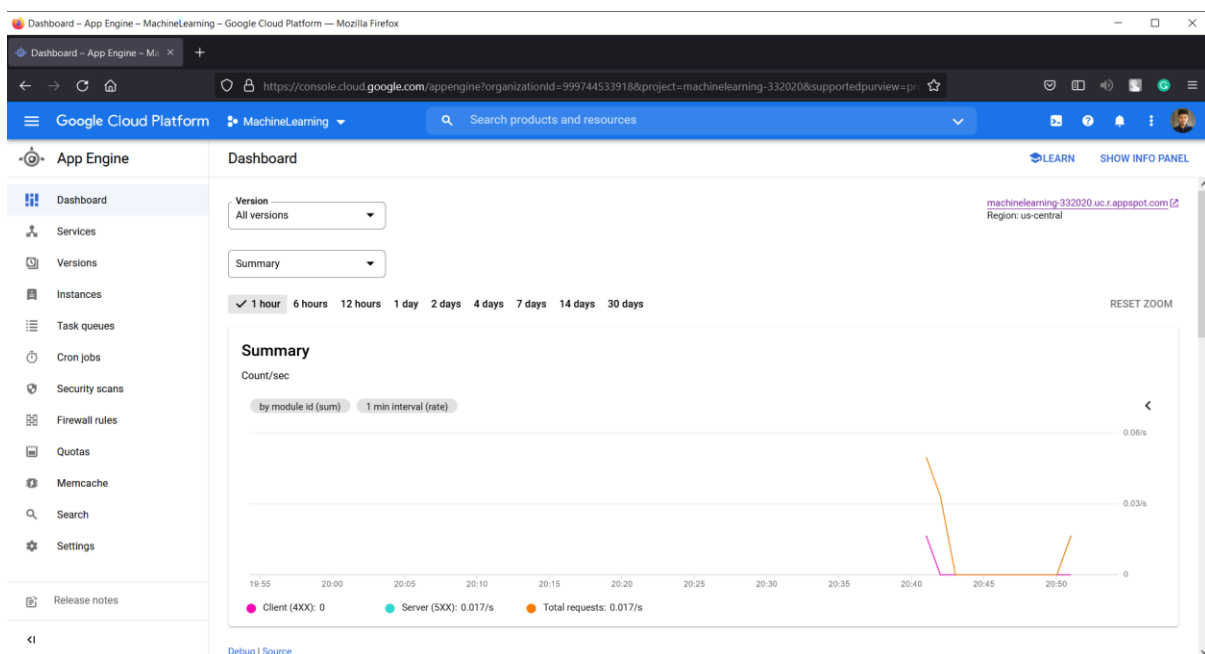


Figure 10: AppEngine Dashboard

Figure 10 displays the dashboard of the AppEngine platform present in GCP.

```
Google Cloud SDK Shell
Welcome to the Google Cloud SDK! Run "gcloud -h" to get the list of available commands.
...
C:\Users\prath\AppData\Local\Google\Cloud SDK>cd ..
C:\Users\prath\AppData\Local\Google>cd ..
C:\Users\prath\AppData\Local>cd ..
C:\Users\prath>cd ..
C:\Users\prath>cd Desktop
C:\Users\prath\Desktop>cd SalaryPredictor
C:\Users\prath\Desktop\SalaryPredictor>gcloud init
Welcome! This command will take you through the configuration of gcloud.

Settings from your current configuration [exit] are:
accessibility:
  screen_reader: 'False'
compute:
  region: us-central1
  zone: us-central1-b
core:
  account: teenu.prathyush2@mail.dcu.ie
  disable_usage_reporting: 'True'
  project: machinelearning-332020

Pick configuration to use:
[1] Re-initialize this configuration [exit] with new settings
[2] Create a new configuration
[3] Switch to and re-initialize existing configuration: [default]
Please enter your numeric choice: 1

Your current configuration has been set to: [exit]

You can skip diagnostics next time by using the following flag:
  gcloud init --skip-diagnostics

Network diagnostic detects and fixes local network connection issues.
Checking network connection...done.
Reachability Check passed.
Network diagnostic passed (1/1 checks passed).

Choose the account you would like to use to perform operations for this configuration:
[1] teenu.prathyush2@mail.dcu.ie
[2] Log in with a new account
Please enter your numeric choice: 1
```

Figure 11(a): Commands to host the web application

Figure 11(a) displays the commands run on Google SDK to host the web application using the AppEngine service of GCP.

```
Google Cloud SDK Shell
You are logged in as: [teenu.prathyush2@mail.dcu.ie].

Pick cloud project to use:
[1] machinelearning-332020
[2] Create a new project
Please enter numeric choice or text value (must exactly match list item): 1

Your current project has been set to: [machinelearning-332020].

Do you want to configure a default Compute Region and Zone? (Y/n)? y

Which Google Compute Engine zone would you like to use as project default?
If you do not specify a zone via a command line flag while working with Compute Engine resources, the default is assumed.
[1] us-east1-b
[2] us-east1-c
[3] us-east1-d
[4] us-east4-c
[5] us-east4-b
[6] us-east4-a
[7] us-central1-c
[8] us-central1-a
[9] us-central1-f
[10] us-central1-b
[11] us-west1-b
[12] us-west1-c
[13] us-west1-a
[14] europe-west4-a
[15] europe-west4-b
[16] europe-west4-c
[17] europe-west1-b
[18] europe-west1-d
[19] europe-west1-c
[20] europe-west3-c
[21] europe-west3-a
[22] europe-west3-b
[23] europe-west2-c
[24] europe-west2-b
[25] europe-west2-a
[26] asia-east1-b
[27] asia-east1-a
[28] asia-east1-c
[29] asia-southeast1-b
[30] asia-southeast1-a
[31] asia-southeast1-c
[32] asia-northeast1-b
[33] asia-northeast1-c
[34] asia-northeast1-a
[35] asia-south1-c
[36] asia-south1-b
[37] asia-south1-a
```

Figure 11(b): Commands to the host the web application

Figure 11(b) is a continuation of the previous step. Here, we configure the zone in which the project will be used.

```
Google Cloud SDK Shell
Your project default Compute Engine region has been set to [us-central1].
You can change it by running [gcloud config set compute/region NAME].

Your Google Cloud SDK is configured and ready to use!

* Commands that require authentication will use teenu.prathyush2@mail.dcu.ie by default
* Commands will reference project 'machinelearning-332020' by default
* Compute Engine commands will use region 'us-central1' by default
* Compute Engine commands will use zone 'us-central1-b' by default

Run 'gcloud help config' to learn how to change individual settings

This gcloud configuration is called [exit]. You can create additional configurations if you work with multiple accounts and/or projects.
Run 'gcloud topic configurations' to learn more.

Some things to try next:

* Run 'gcloud --help' to see the Cloud Platform services you can interact with. And run 'gcloud help COMMAND' to get help on any gcloud command.
* Run 'gcloud topic --help' to learn about advanced features of the SDK like arg files and output formatting

C:\Users\prath\Desktop\SalaryPredictor>gcloud app deploy app.yaml --project machinelearning-332020
Services to deploy:

descriptor:      [C:\Users\prath\Desktop\SalaryPredictor\app.yaml]
source:          [C:\Users\prath\Desktop\SalaryPredictor]
target project:  [machinelearning-332020]
target service:  [default]
target version:  [20211115t205924]
target url:      [https://machinelearning-332020.uc.r.appspot.com]
target service account: [App Engine default service account]

Do you want to continue (Y/n)? y

Beginning deployment of service [default]...
#=====#
# Uploading 1 file to Google Cloud Storage      =#
#=====#
File upload done.
Updating service [default]...done.
Setting traffic split for service [default]...done.
Deployed service [default] to [https://machinelearning-332020.uc.r.appspot.com]

You can stream logs from the command line by running:
$ gcloud app logs tail -s default

To view your application in the web browser run:
$ gcloud app browse

C:\Users\prath\Desktop\SalaryPredictor>
```

Figure 11(c): Commands to the host the web application

Figure 11(c) displays the progress of deployment of the web application using the AppEngine service of GCP.

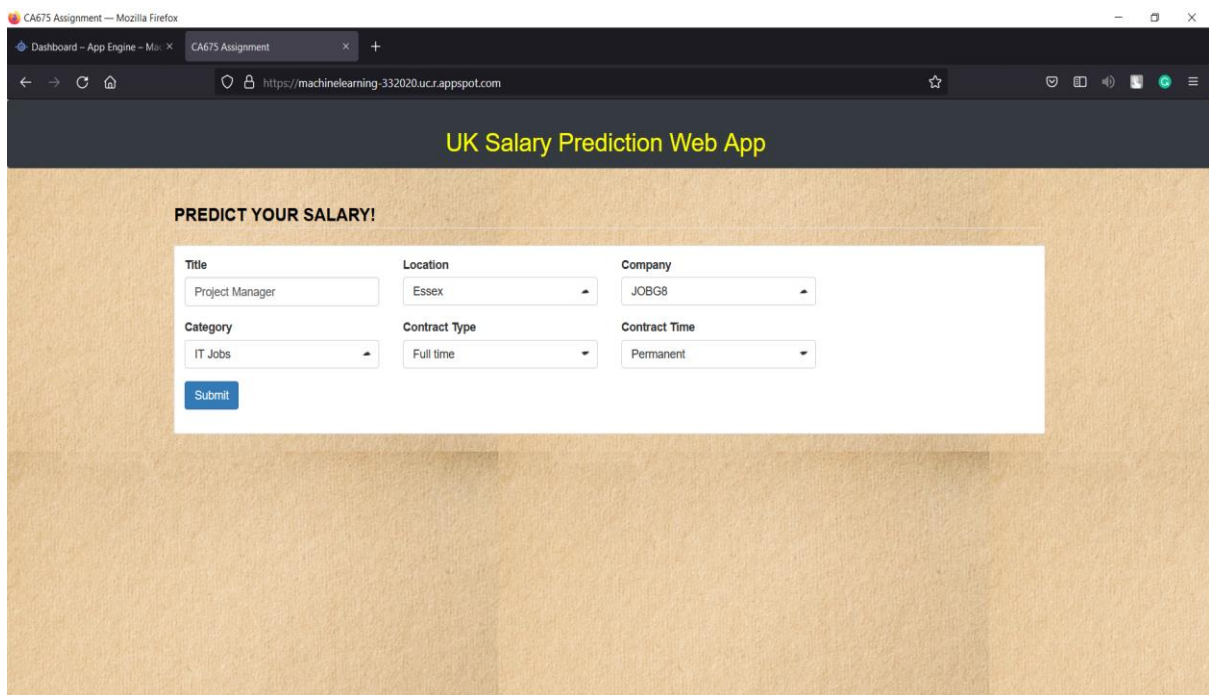


Figure 12: Deployment of the web application

Figure 12 displays the deployment of the web application using AppEngine Service of the Google Cloud Platform, where a user can enter values in the form to get the estimated salary.

UK Salary Prediction Web App

PREDICT YOUR SALARY!

Title:

Location:

Company:

Category:

Contract Type:

Contract Time:

Estimated Salary is €43570.54

Figure 13: Predicting the estimated salary

Figure 13 displays the estimated salary for the values entered in the form (refer figure 12).

8. CHALLENGES AND LESSONS LEARNED:

The challenges and lessons learned in building the web application are described below:

- **Acquiring the dataset:** Not many datasets are available containing big data for this particular project. We have taken the UK job salary prediction dataset from Kaggle as we found it to be the most suitable dataset for implementing our project.
- **Training the machine learning model with a suitable algorithm:** Having trained the model with different machine learning algorithms like Support Vector Machines, Decision trees, etc., we found Random Forest Regressor to be the best performing model.
- **Designing the web application:** The design of the web application is simple, elegant, and responsive. However, with more time we could have improved the design. For example, the **Title** field in the web application form takes a user input instead of a list of values. Another design feature we could implement in the future is that when a user enters a value in one of the fields, the values in related fields should be auto-populated.

In the future, we could enhance this model to include more features and an improved front-end design. We could also develop a mobile application for predicting salaries.

9. LINK TO THE WALKTHROUGH VIDEO:

<https://www.youtube.com/watch?v=6KYlb2Kp4Jc&t=1s>

10. TEAM ROLES AND TASKS:

- **Data Collection and Processing:** Fiaz Ashraf and Divyashree Parameshwara - SATISFACTORY
- **Front-end Setup:** Teena Sharma, Anish Patil - SATISFACTORY
- **Back-end Setup:** Teenu Prathyush, Niranjana Nimbarte - SATISFACTORY
- **WebApp Deployment:** Teenu Prathyush - SATISFACTORY
- **Report:** All the members of the group contributed in writing the report.

LIST OF FIGURES

Fig No.	Figure Name	Page No.
1.	Cloud Storage - Google Cloud Staging Bucket	5
2.	Pig Commands	5
3.	Pig - processed data output	6
4.	Creating an instance on Google Cloud Platform	6
5.	Installing Anaconda on the created instance	7
6.	Launching Jupyter Notebook	7
7.	Creating the Machine Learning model	8
8.	Serializing the machine learning model into a Pickle file	8
9.	Loading the pickle file into main.py	9
10.	AppEngine Dashboard	9
11(a).	Commands to host the web application	10
11(b).	Commands to host the web application	10
11(c).	Commands to host the web application	11
12.	Deployment of the web application	11
13.	Predicting the estimated salary	12