

Project rapid fire

Q.1 Can you explain your CDAC Project? Draw block diagram as appropriate.

→ My CDAC Project was titled 'Airbnb Data Analysis' & Price prediction' where me & my team project partner worked on a dataset to get some insight of the Airbnb listings & predict the price of new properties using this data.

The dataset is basically used for a kaggle competition. It consists of listings from 2008 till 2018. There were many missing ~~and~~ values in the dataset which is why a lot of cleaning was required. We performed feature selection by 3 different ways & selected the best method. Similarly, the choice of algorithm was done based on the performance metrics.

We mainly focused on RMSE & R² score values as for determining the performance.

Finally we used Random Forest and tuned it for best performance.

Further, ~~the~~ a web application is developed using flask framework and deployed using AWS EC2 instance.

Q.3 Why did you select this technology & framework for this project?

→ We have used many libraries in Python such as Pandas, matplotlib & Seaborn for multiple advantages.

These libraries have many advantages like easy implementation, readability, Openly available resources, the highly efficient data tools, platform independence and most importantly, the support of a large developer community.

We used Flask framework for creating a ^{small} web application. Flask does not require a lot of configuration, Hence it is suitable for small to medium projects.

After that, to deploy the application we used Amazon's AWS EC2 instance.

Q.5 Which advanced features have you used in your project?



~~I have~~

I have used some advanced data cleaning techniques in some features of the dataset like binning, imputing values according to another feature etc. The main aim during cleaning was to avoid the loss of data and make max use of the available usable data.

For feature selection I have used an external package called ~~featurez~~ which uses ~~&~~ ~~regboos~~ for feature selection.

Q.6 What was your role in your project & explain what you did in it?

→ I did this project with my project partner. In the initial phase, we ~~first~~ both searched for a suitable dataset. Then after getting one, ~~as~~ ~~the~~ the flow of the project was decided by me. It includes the steps ~~for~~ by step approach, the distribution of tasks and deadline which we had set for ourself.

Among the technical things, I was responsible for data cleaning and training. Other tasks like visualization & feature selection were done by my partner.

for, ~~for the~~ whatever part that I did, ~~I followed the approach of~~ I first looked for the standard methods to do that, then which python libraries can be used, how to use them etc. then I looked for how to apply it to our project.

Finally the documentation part - like collecting images, writing & editing the report was done by me.

Q.9 How will you deploy your project?

→ I am using Amazon AWS EC2 instance currently for deploying the model because it is free of cost. But for bigger storage and processing requirement, the instance can be scaled up.

Q.10 What are the limitations of your project?

→ I think, In our project, we did cleaning & visualization at a good level. But the accuracy of our model can be improved. And this can be done using deep learning algorithms. Our machines are old.

Another limitation which I faced is that I wanted to perform grid search for selecting optimum hyper parameters. But couldn't do it because of the configuration of our machines.

Q.11 What are the difficulties that you have faced in your project? How you overcame it?

→ The biggest difficulty which I found was handling the missing values in feature selection. Be my partner

During feature selection for better results I started to search on the different ways for feature selection. I came across an external package called featureuz which uses xgboost for feature selection.

Q.12 How will you improve performance of your project (memory related & response time)?

Q.15 mention sources of your project.

→ we used the dataset used for a Kaggle competition. It was in CSV format and included data from 2008 to 2019.

Q.16 Application of your project?

→ The main aim of the project was to build an application ~~which~~ which predicts the price of a property based on various features.

Property-owners who want to list their properties on Airbnb ~~can use this app~~ but are not sure about what price should they assign, can use this app.

Q.18 How much data have you allocated for training, validation & testing in your project?

→ Training : 80 %

Testing : 20 %

Not allocated data for validation as the data was small.

Q.21 Mention which steps you followed in your project.

→ Step ① Understanding the data.

From various sources, user experiences, host experiences, we understood the features which highly affect the price.

Step ② Visualization

To get better understanding of the data and to create some dynamic dashboards, visuals were created using PowerBI.

Step ③ Data pre-processing / Cleaning

Some features cannot be directly used for analytics. They have to be modified if they are of different type, contain null values or corrupted values.

It is the most crucial part, which was done keeping in mind the goal to avoid data loss as much as possible.

Step ④ Feature & Algorithm selection.

Out of the available features, only the few important ones are selected to improve accuracy. Further different algorithms are tested & the one with best results is selected.

Step ⑤ Training & hyperparameter tuning

On the selected algorithm, tuning is performed for obtaining the best hyperparameters.

Step ⑥ Create & deploy the application

A web application is created using flask framework in python. This app is deployed using Amazon EC2 instance on the cloud.

Q. 22 Which data cleaning methods used in your project?

→ For missing values, if no. of nulls is less then it can be replaced by mean, mode or median. But for a high no. of missing values various techniques have been used.

(i) Binning: Converting a numerical feature into a categorical one & assigning a new category for null values.

(ii) Imputing values according to some other independent variable

for ~~pre~~ processing the date columns,
first the date columns were converted to
~~date~~ no. of days till current date & then
by binning them into categories, missing
values were handled.

from the amenities feature, ~~and~~ unique
amenities were extracted & important
amenities were converted into new columns.