

# NEURAL NETWORK AND FUZZY LOGIC

## IA-1 REPORT

### TOPIC: HOUSE PRICE PREDICTION



### GROUP MEMBERS

1813001 – ABHIRAM ANILKUMAR

1813009 – PRATHAMESH DAPHAL

1813031 – SHUBH MODY

## □ INTRODUCTION

In this project, we will develop and evaluate the performance and the predictive power of a model trained and tested on data collected from houses in the state of Tamil Nadu.

The data is collected from the 7 major cities in Tamil Nadu which have houses of different build type like housing, commercial and others and several other parameters which affect the sales price of these houses.

Once we get a good fit, we will use this model to predict the sales price of a house located in these regions.

A model like this would be very valuable for a real estate agent who could make use of the information provided in a daily basis.

You can find the complete project, documentation and dataset on:



## □ DATASET

The dataset is the Chennai\_house\_price\_prediction.csv file which is provided in the github repository. To import the dataset pandas and other libraries are imported for the entire preprocessing and graphical analysis, which is done using:

```

8 import pandas as pd
9 import numpy as np
10 import matplotlib.pyplot as plt
11 import seaborn as sn
12 import tensorflow as tf
13 from tensorflow.keras.models import Sequential
14 from tensorflow.keras.layers import Dense, Activation
15 from tensorflow.keras.optimizers import Adam
16 from tensorflow.keras.callbacks import EarlyStopping
17 from sklearn.model_selection import train_test_split
18 from sklearn.preprocessing import MinMaxScaler
19 from sklearn import metrics
20 from sklearn.preprocessing import LabelEncoder
21 import streamlit as st
22 from sklearn.linear_model import LinearRegression
23

```

The dataset has got 18 different features and 7109 different rows with the SALES\_PRICE being the target variable. This is checked using the df.shape command

The overview of the dataset is given below

AREA	INT_SQFT	T_MAINR	_BEDROOM	BATHROOM	N_ROOM	SALE_COND	ARK_FACI	BUILDTYPE	ILITY_AVA	STREET	MZZONE	_S_ROOM	_BATHROOM	_BEDROOM	S_OVERAL	COMMIS	SALES_PRICE
Karapakkam	1004	131	1	1	3	AbNormal	Yes	Commercial	AllPub	Paved	A	4	3.9	4.9	4.33	144400	7600000
Anna Nagar	1986	26	2	1	5	AbNormal	No	Commercial	AllPub	Gravel	RH	4.9	4.2	2.5	3.765	304049	21717770
Adyar	909	70	1	1	3	AbNormal	Yes	Commercial	ELO	Gravel	RL	4.1	3.8	2.2	3.09	92114	13159200
Velachery	1855	14	3	2	5	Family	No	Others	NoSewr	Paved	I	4.7	3.9	3.6	4.01	77042	9630290
Karapakkam	1226	84	1	1	3	AbNormal	Yes	Others	AllPub	Gravel	C	3	2.5	4.1	3.29	74063	7406250
Chrompet	1220	36	2	1	4	Partial	No	Commercial	NoSewr	No Access	RH	4.5	2.6	3.1	3.32	198316	12394750
Chrompet	1167	137	1	1	3	Partial	No	Other	AllPub	No Access	RL	3.6	2.1	2.5	2.67	33955	8488790

This also has a unique id for each row which can be dropped since it doesn't have any significance while predicting the model.

The features can be summarized as follows:

1. INT\_SQFT: It is the total square feet area of a particular house.
2. DIST\_MAINROAD: It is distance of a house from the main road.
3. N\_BEDROOM: Number of bedrooms in a house
4. N\_BATHROOM: Number of bathrooms in a house
5. N\_ROOM: Total number of rooms in a house which is the cumulative sum of all the different type of rooms in a house
6. SALE\_COND: The current condition in which the house is sold.
7. PARK\_FACIL: Used to identify whether a house has a parking facility or not
8. BUILDTYPE: The purpose for which the property is used, like housing or commercial or office etc.
9. UTIL\_AVAIL: The utilities available in the particular house or property.
10. STREET: The type of road leading to a particular property.
11. QS: Quality scores are given for bedrooms, bathrooms and for the total rooms in a house which can affect its sales price (QS\_ROOMS, QS\_BEDROOMS, QS\_BATHROOMS)
12. COMMISSION: The commission for each house which highly depends on the sales price.

## □ PREPROCESSING

In preprocessing we clean the data that is we fill or remove the missing values, change the data types and correct the errors in spelling. Also duplicate items should be removed since it would be redundant.

Using this command we get the total missing values in the dataset and the data type of each feature.

```
df.isnull().sum()
df.dtypes
```

These commands will give us output as:

<pre>In [8]: df.dtypes Out[8]: PRT_ID          object AREA            object INT_SQFT        int64 DIST_MAINROAD   int64 N_BEDROOM       float64 N_BATHROOM      float64 N_ROOM          int64 SALE_COND       object PARK_FACIL      object BUILDTYPE       object UTILITY_AVAIL   object STREET          object MZZONE          object QS_ROOMS        float64 QS_BATHROOM     float64 QS_BEDROOM      float64 QS_OVERALL      float64 COMMIS          int64 SALES_PRICE     int64 dtype: object</pre>	<pre>In [7]: df.isnull().sum() Out[7]: PRT_ID          0 AREA            0 INT_SQFT        0 DIST_MAINROAD   0 N_BEDROOM        1 N_BATHROOM       5 N_ROOM          0 SALE_COND        0 PARK_FACIL       0 BUILDTYPE       0 UTILITY_AVAIL    0 STREET          0 MZZONE          0 QS_ROOMS        0 QS_BATHROOM     0 QS_BEDROOM      0 QS_OVERALL      48 COMMIS          0 SALES_PRICE     0 dtype: int64</pre>
--	---

This shows that we have some missing values which is needed to be taken care of.

Filling out the missing values:

```
df['N_BEDROOM'].mode()
df['N_BEDROOM'].fillna(value=(df['N_BEDROOM'].mode()[0]), inplace=True)
```

The missing values of the bedroom is filled with the mode of the bedroom feature

```
for i in range(0, len(df)):
    if pd.isnull(df['N_BATHROOM'])[i]==True:
        if (df['N_BEDROOM'][i]==1.0):
            df['N_BATHROOM'][i] = 1.0
        else:
            df['N_BATHROOM'][i] = 2.0
```

If there is only one bedroom then only one bathroom is required otherwise at least two bathrooms are filled in the missing value space.

```
for i in range(0, len(df)):
    if pd.isnull(df['QS_OVERALL'])[i]==True:
        df['QS_OVERALL'][i] = (df['QS_ROOMS'][i] + df['QS_BEDROOM'][i] + df['QS_BATHROOM'][i])/3
```

QS-OVERALL value is approximated as the average of the QS\_ROOMS, QS\_BEDROOM, QS\_BATHROOM.

Therefore now the total missing values after preprocessing is:

```
In [13]: df.isnull().sum()
Out[13]:
AREA                0
INT_SQFT             0
DIST_MAINROAD       0
N_BEDROOM           0
N_BATHROOM          0
N_ROOM              0
SALE_COND            0
PARK_FACIL          0
BUILDTYPE           0
UTILITY_AVAIL       0
STREET              0
MZZONE              0
QS_ROOMS            0
QS_BATHROOM         0
QS_BEDROOM          0
QS_OVERALL          0
COMMIS              0
SALES_PRICE         0
dtype: int64
```

Duplicates in the dataset is removed by the command `df.drop_duplicates()`

The spelling error correction is done below using the following code. Spelling errors can increase the unique values of each feature and create imbalance of data due to which the accuracy of the model is very much affected.

```
CorrectNames = ['AREA', 'SALE_COND', 'PARK_FACIL', 'BUILDTYPE', 'UTILITY_AVAIL', 'STREET', 'MZZONE']

for i in CorrectNames:
    print('*****Value Counts*****')
    print(df[i].value_counts())
    print('Unique values are:', df[i].nunique())

df['AREA'].replace({'Chrompt': 'Chrompet', 'Chormpet': 'Chrompet', 'Chrmpt': 'Chrompet',
                    'TNagar': 'T Nagar', 'Ana Nagar': 'Anna Nagar', 'Ann Nagar': 'Anna Nagar',
                    'Karapakam': 'Karapakkam', 'Velchery': 'Velachery', 'Adyr': 'Adyar',
                    'KKNagar': 'KK Nagar'}, inplace=True)
df['AREA'].value_counts()

df['BUILDTYPE'].replace({'Comercial': 'Commercial', 'Other': 'Others'}, inplace=True)

df['SALE_COND'].replace({'Adj Land': 'AdjLand', 'Ab Normal': 'AbNormal', 'Partiall': 'Partial',
                        'PartialL': 'Partial'}, inplace=True)
df['SALE_COND'].value_counts()

df['PARK_FACIL'].replace({'Noo': 'No'}, inplace=True)
df['PARK_FACIL'].value_counts()

df['UTILITY_AVAIL'].replace({'All Pub': 'AllPub'}, inplace=True)
df['UTILITY_AVAIL'].value_counts()

df['STREET'].replace({'Pavd': 'Paved', 'NoAccess': 'No Access'}, inplace=True)
df['STREET'].value_counts()
```

.replace function is used which consists of a dictionary which replaces the incorrect word with correct word.

The data types change using the .astype function which is implemented below:

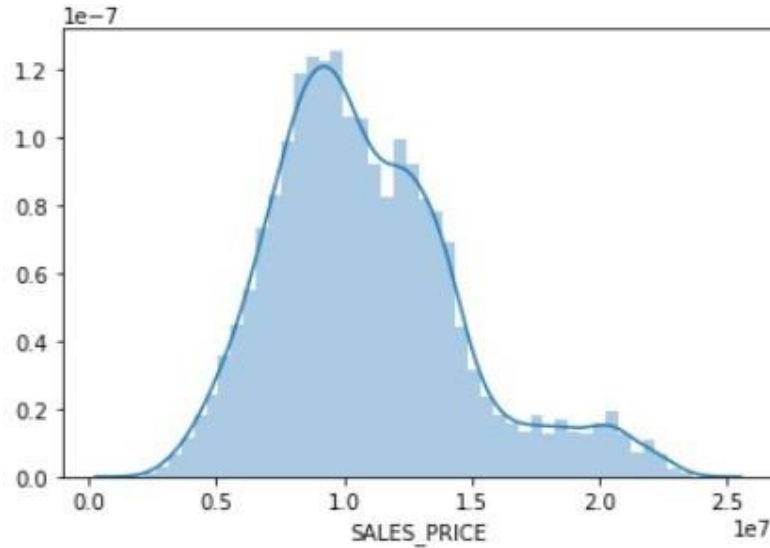
```
df = df.astype({'N_BATHROOM': 'int64', 'N_BEDROOM': 'int64'})
```

The data types of N\_BEDROOM and N\_BATHROOM are changed from float to integers.

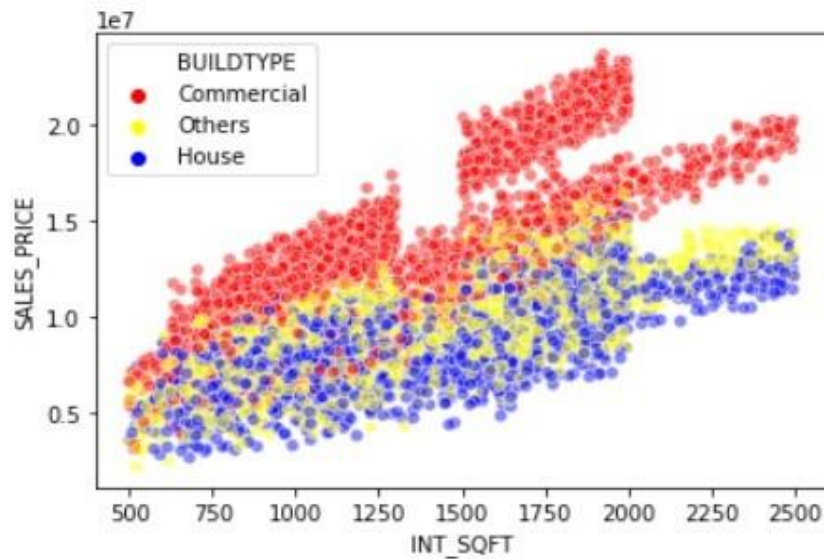


### Graphical analysis:

The following are some of the graphical analysis made from the dataset using the seaborn and matplotlib libraries.

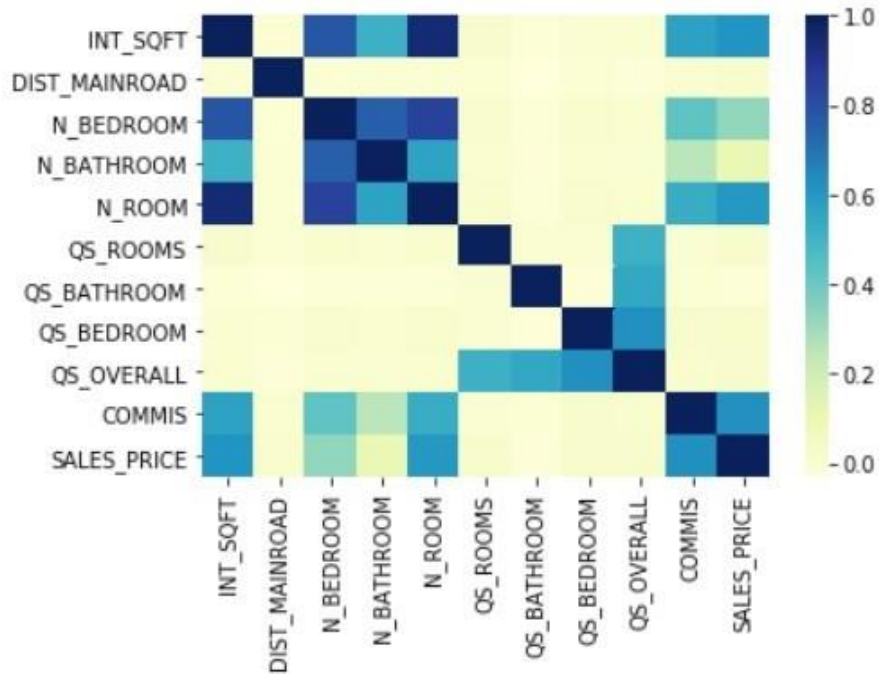


This following is a histogram with density plot which shows that the sales price of most of the house falls in the 1 crore mark.

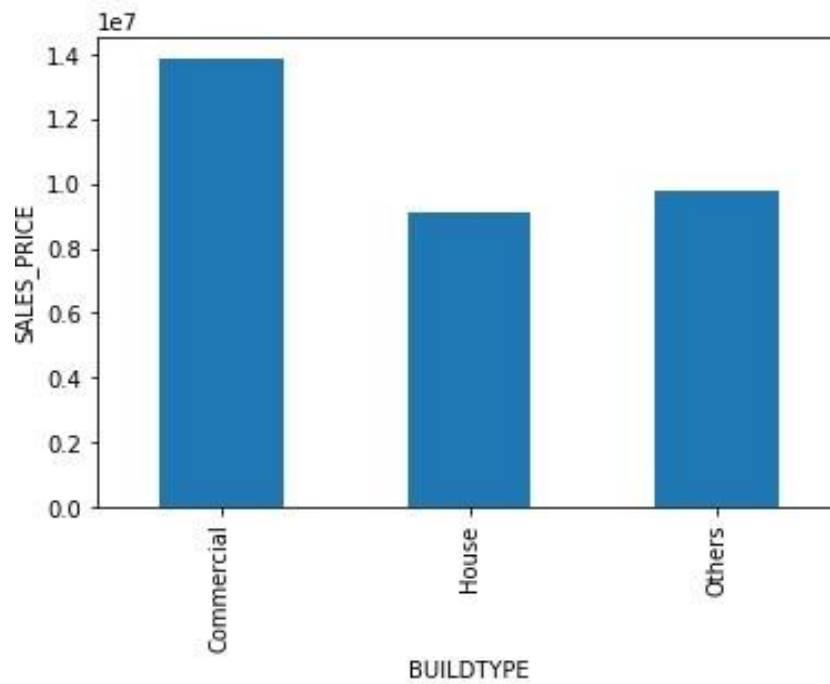


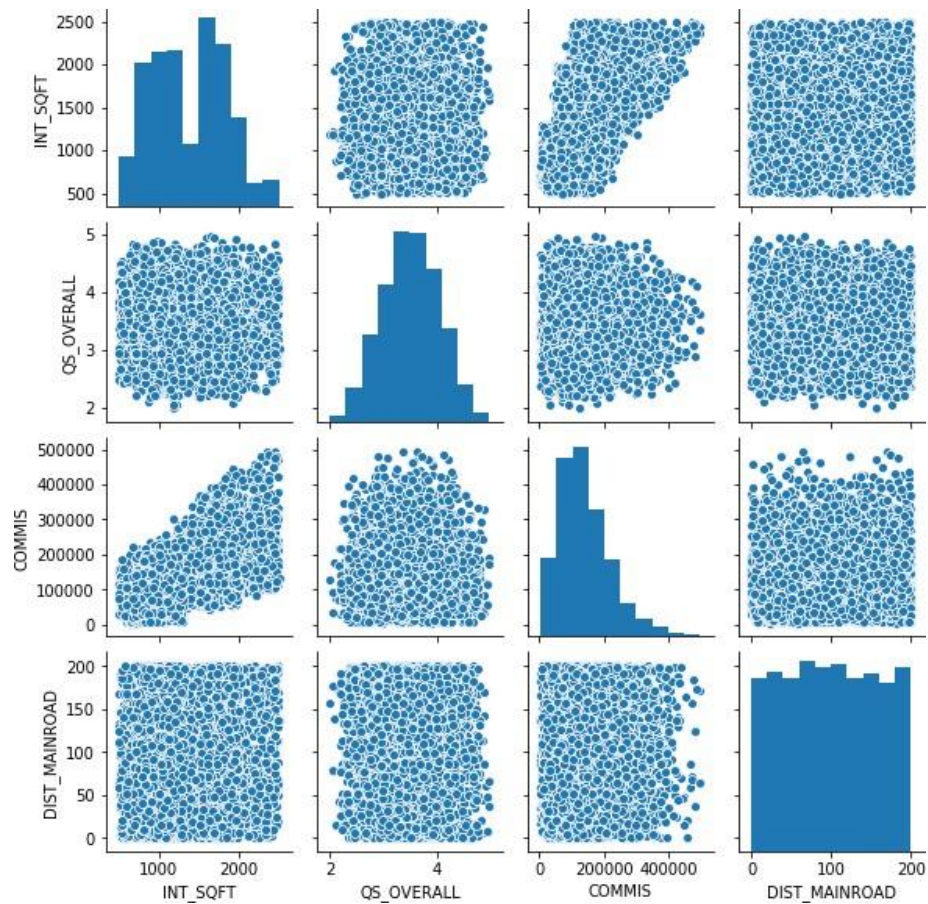
This scatter plot shows the relation between square feet and sales price of the house with the hue as buildtype.





Heatmap show the correlation between various features of the dataset.





Following is a pair plot which shows relationship between various continuous variables

## Label encoding:

Label Encoding refers to converting the labels into numeric form so as to convert it into the machine-readable form. Machine learning algorithms can then decide in a better way on how those labels must be operated. It is an important pre-processing step for the structured dataset in supervised learning

Code :

```
enc=LabelEncoder()
df.iloc[:,0]=enc.fit_transform(df.iloc[:,0])
df.iloc[:,6]=enc.fit_transform(df.iloc[:,6])
df.iloc[:,7]=enc.fit_transform(df.iloc[:,7])
df.iloc[:,8]=enc.fit_transform(df.iloc[:,8])
df.iloc[:,9]=enc.fit_transform(df.iloc[:,9])
df.iloc[:,10]=enc.fit_transform(df.iloc[:,10])
df.iloc[:,11]=enc.fit_transform(df.iloc[:,11])
```

## An overview off the dataset after one-hot encoding and normalization

Index	AREA	INT_SQFT	T_MAINRC	_BEDROO	BATHROC	N_ROOM	ALE_CONI	ARK_FACI	BUILDTYPEI	ILITY_AVF	STREET	MZZONE
0	4	1004	131	1	1	3	0	1	0	0	2	0
1	1	1986	26	2	1	5	0	0	0	0	0	3
2	0	909	70	1	1	3	0	1	0	1	0	4
3	6	1855	14	3	2	5	2	0	2	3	2	2
4	4	1226	84	1	1	3	0	1	2	0	0	1
5	2	1220	36	2	1	4	4	0	0	2	1	3
6	2	1167	137	1	1	3	4	0	2	0	1	4
7	6	1847	176	3	2	5	2	0	0	0	0	5
8	2	771	175	1	1	2	1	0	2	3	2	5
9	6	1635	74	2	1	4	0	0	2	1	1	2
10	2	1203	78	2	1	4	1	1	0	0	1	5
11	2	1054	143	1	1	3	4	0	2	3	0	5

## □ DEVELOPING THE MODEL

Splitting and Normalizing the data:

```

from sklearn.cross_validation import train_test_split
X = df.drop('SALES_PRICE',axis=1).values
y = df['SALES_PRICE'].values
X_train,X_test,y_train,y_test=train_test_split(X,y,test_size=0.2,random_state=101)
scaler = MinMaxScaler()
X_train = scaler.fit_transform(X_train)
X_test=scaler.fit_transform(X_test)

```

Data after Normalizing:

	0	1	2	3	4	5	6
0	0.666667	0.437719	0.19	0.333333	1	0.5	0.25
1	0.333333	0.206103	0.355	0	0	0.25	0.75
2	0.166667	0.743372	0.545	0.333333	0	0.75	1
3	1	0.715358	0.7	0.666667	1	0.75	0.25
4	0	0.3997	0.31	0.333333	1	0.5	0.75
5	0.333333	0.312156	0.12	0	0	0.25	0.25
6	0.333333	0.247624	0.205	0	0	0.25	0.5
7	0.333333	0.376688	0.675	0.333333	0	0.5	0.5
8	0.166667	0.71986	0.65	0.333333	0	0.75	0.75
9	0.833333	0.63932	0.625	0	0	0.5	0.25

Our Deep Learning Model:

```
model=Sequential() #initialising ANN
model.add(Dense(17,activation='relu'))
model.add(Dense(17,activation='relu'))
model.add(Dense(17,activation='relu'))
model.add(Dense(17,activation='relu'))
model.add(Dense(17,activation='relu'))
model.add(Dense(17,activation='relu'))
model.add(Dense(17,activation='relu'))
model.add(Dense(17,activation='relu'))
model.add(Dense(17,activation='relu'))
model.add(Dense(17,activation='relu'))
model.add(Dense(17,activation='relu'))

model.add(Dense(1))

model.compile(optimizer='adam',
              loss='mae',
              )
model.fit(x=X_train,y=y_train,
          validation_data=(X_test,y_test),
          batch_size=128,epochs=200,
          )
model.summary()
```

Output:

## NNFL IA-1

Console 2/A		
45/45	[=====]	- 0s 4ms/step - loss: 323007.8539 - val_loss: 321186.0312
Epoch 199/200		
45/45	[=====]	- 0s 4ms/step - loss: 313749.6746 - val_loss: 345241.0000
Epoch 200/200		
45/45	[=====]	- 0s 6ms/step - loss: 318047.4321 - val_loss: 325624.0938
Model: "sequential"		
Layer (type)	Output Shape	Param #
=====		
dense (Dense)	(None, 17)	306
dense_1 (Dense)	(None, 17)	306
dense_2 (Dense)	(None, 17)	306
dense_3 (Dense)	(None, 17)	306
dense_4 (Dense)	(None, 17)	306
dense_5 (Dense)	(None, 17)	306
dense_6 (Dense)	(None, 17)	306
dense_7 (Dense)	(None, 17)	306
dense_8 (Dense)	(None, 17)	306
dense_9 (Dense)	(None, 17)	306
dense_10 (Dense)	(None, 17)	306
dense_11 (Dense)	(None, 1)	18
=====		
Total params: 3,384		
Trainable params: 3,384		
Non-trainable params: 0		

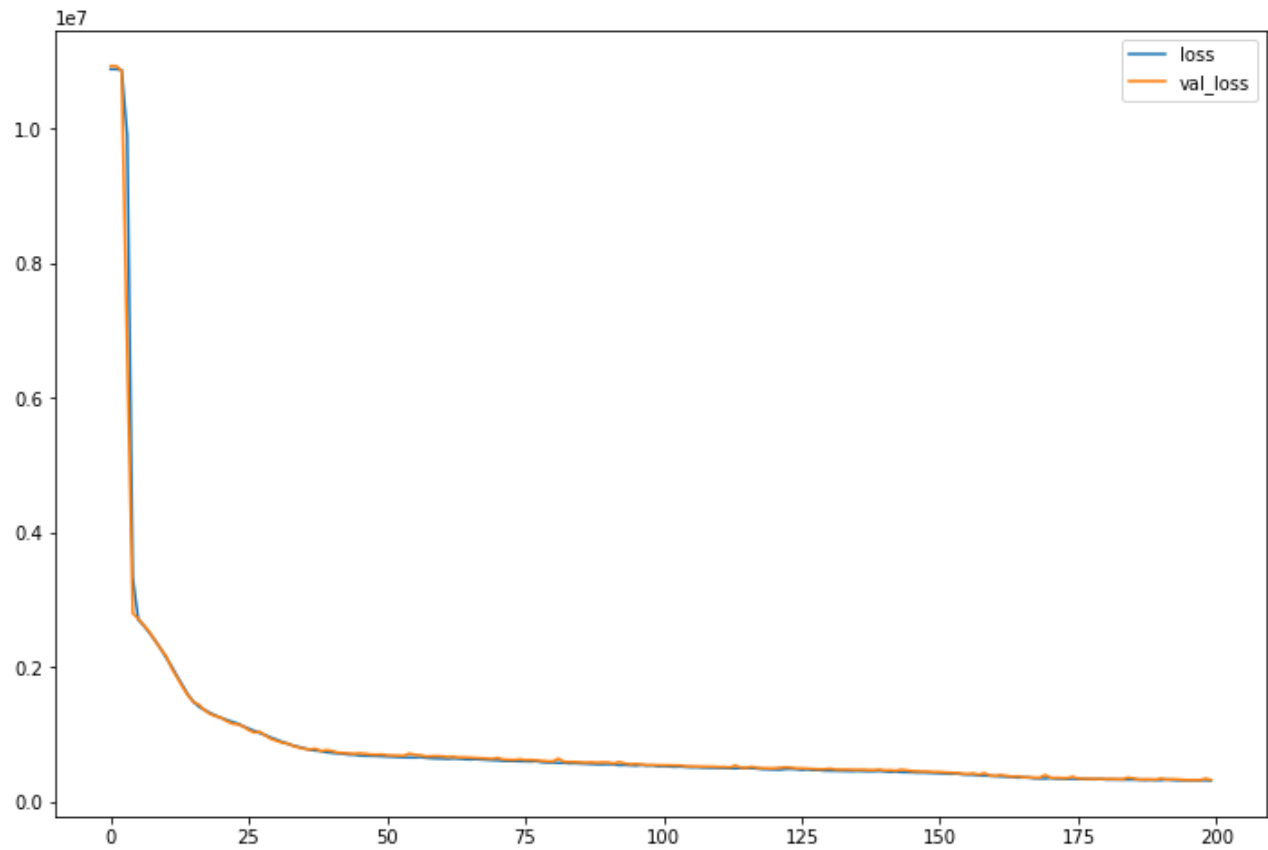
MAE on the test data set and the Accuracy of the model:

Trainable params: 3,384	
Non-trainable params: 0	
<hr/>	
MAE: 325624.11251758086	
MSE: 167282056707.2385	
RMSE: 409001.2918160999	
VarScore: 0.9892304917510264	

So we have got almost 99% accuracy on our model.



Plot of training and Validation loss:



Interpretation of Results:

```
11/10/0  
print(df['SALES_PRICE'].describe())
```

```
In [3]: runcell(8, 'C:/Users/Prathamesh/Desktop/NNFL_IA/NNFL_final.py')
count    7.109000e+03
mean     1.089491e+07
std      3.768603e+06
min      2.156875e+06
25%      8.272100e+06
50%      1.033505e+07
75%      1.299390e+07
max      2.366734e+07
Name: SALES_PRICE, dtype: float64

In [4]:
```

---

From statistics, the mean price of a house in these regions of TamilNadu is around 1 crore. Therefore mean absolute error of about 3 lakhs means that we are bound to overshoot or undershoot by around 0 – 4 lakhs. This means we have an error of about 0% - 4% when compared to the target variable.

## DEPLOYING model aa WEB APPLICATION:

We have used Streamlit library of python for deploying model on web.

```
7
8 import pandas as pd
9 import numpy as np
10 import matplotlib.pyplot as plt
11 import seaborn as sn
12 import tensorflow as tf
13 from tensorflow.keras.models import Sequential
14 from tensorflow.keras.layers import Dense, Activation
15 from tensorflow.keras.optimizers import Adam
16 from tensorflow.keras.callbacks import EarlyStopping
17 from sklearn.model_selection import train_test_split
18 from sklearn.preprocessing import MinMaxScaler
19 from sklearn import metrics
20 from sklearn.preprocessing import LabelEncoder
21 import streamlit as st
22 from sklearn.linear_model import LinearRegression
23
```

---

Code for the same can be referred through the uploaded python file on github.

HOW THE INTERFACE LOOKS LIKE??

Command for running code:

```
Anaconda Prompt (anaconda3) - streamlit run NNFL_final.py

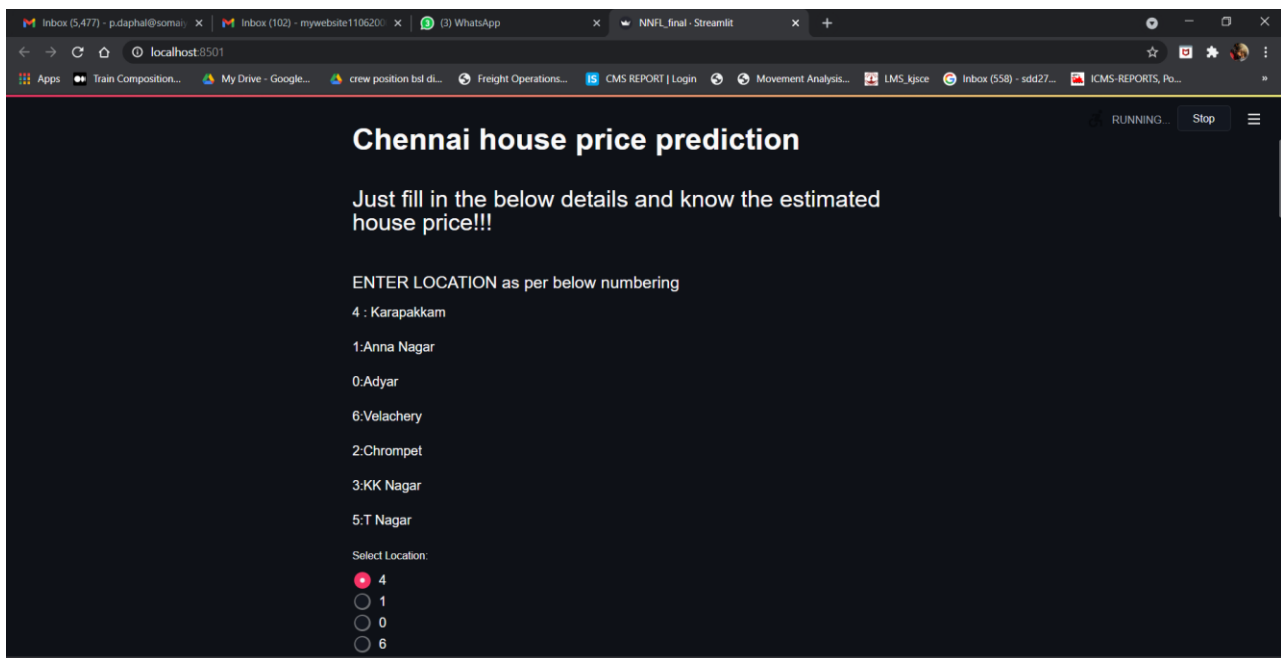
(base) C:\Users\Prathamesh>cd desktop

(base) C:\Users\Prathamesh\Desktop>cd NNFL_IA

(base) C:\Users\Prathamesh\Desktop\NNFL_IA>streamlit run NNFL_final.py

You can now view your Streamlit app in your browser.

Local URL: http://localhost:8501
Network URL: http://192.168.1.102:8501
```



## NNFL IA-1

2  
3  
5

Select the area

500 981 2500

Selected: 981

Distance from main road

0 39 200

Selected: 0

Number of rooms

2

Number of bedrooms

1

Number of bathrooms

1

RUNNING... Stop

Enter the sale condition as per below numbering

0:AbNormal

2:Family

4:Partial

1:AdjLand

3:Normal

Sale Condition:

0  
2  
4  
1  
3

Enter the requirement of Parking facility

1:YES

0:NO

Select Parking facility:

1  
0

RUNNING... Stop

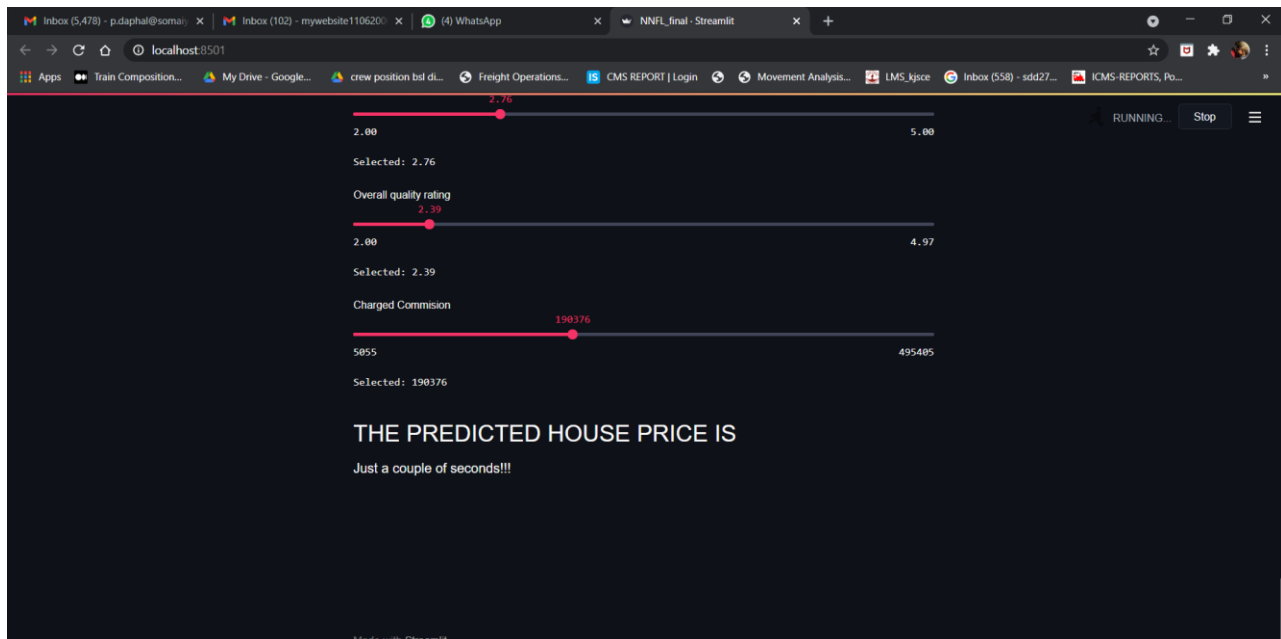
In this way we take inputs from user...

And this is the DL model running at backend:

## NNFL IA-1

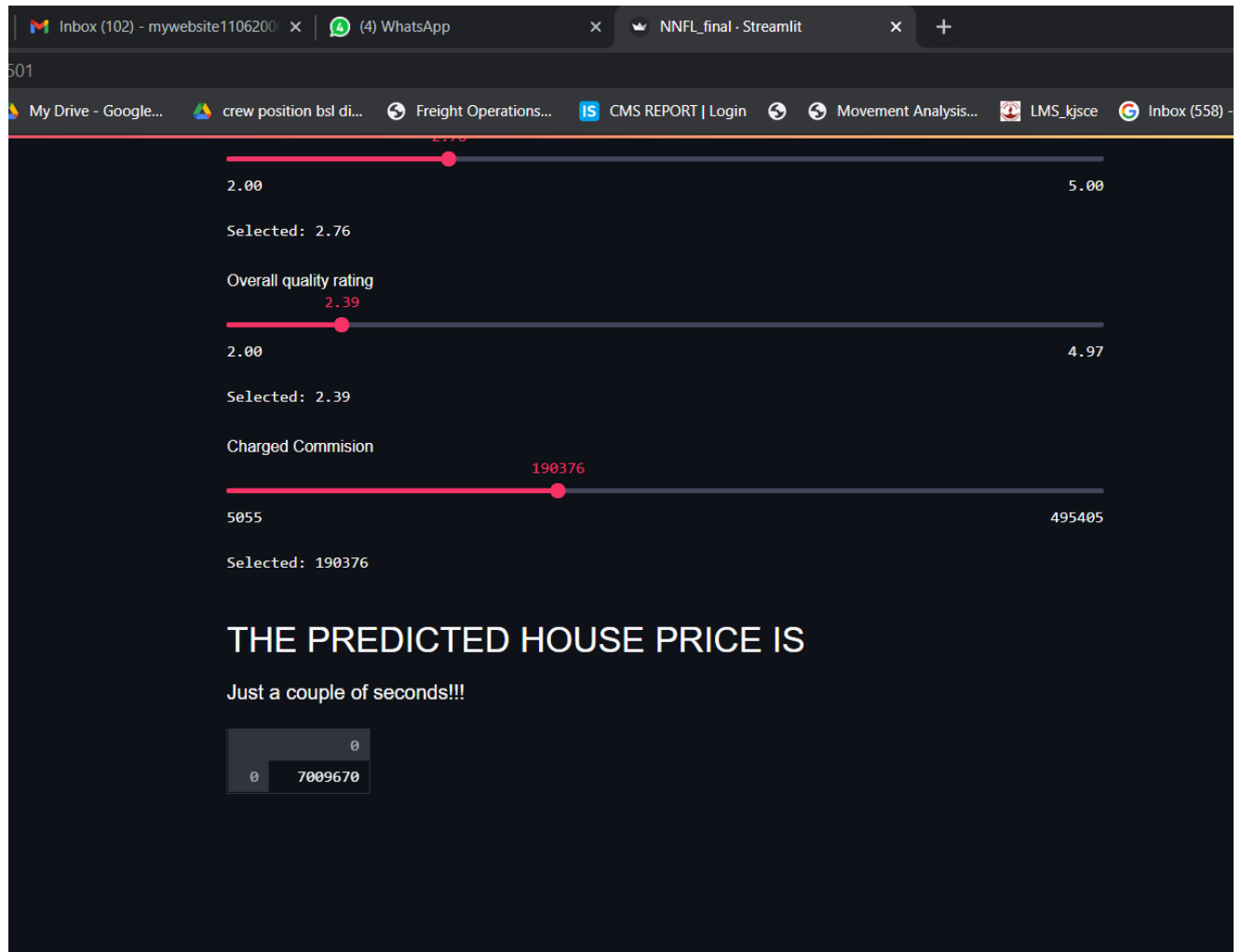
```
Anaconda Prompt (anaconda3) - streamlit run NNFL_final.py
Epoch 85/200
45/45 [=====] - 0s 5ms/step - loss: 575918.7480 - val_loss: 602741.3750
Epoch 86/200
45/45 [=====] - 0s 4ms/step - loss: 565147.1332 - val_loss: 596361.3750
Epoch 87/200
45/45 [=====] - 0s 5ms/step - loss: 572807.2649 - val_loss: 592287.2500
Epoch 88/200
45/45 [=====] - 0s 6ms/step - loss: 563589.8913 - val_loss: 587236.6250
Epoch 89/200
45/45 [=====] - 0s 3ms/step - loss: 560799.7215 - val_loss: 587456.3750
Epoch 90/200
45/45 [=====] - 0s 5ms/step - loss: 561443.1848 - val_loss: 580403.7500
Epoch 91/200
45/45 [=====] - 0s 3ms/step - loss: 554110.6753 - val_loss: 576878.1875
Epoch 92/200
45/45 [=====] - 0s 4ms/step - loss: 553975.9022 - val_loss: 574521.4375
Epoch 93/200
45/45 [=====] - 0s 7ms/step - loss: 549834.4565 - val_loss: 571509.9375
Epoch 94/200
45/45 [=====] - 0s 6ms/step - loss: 550161.5557 - val_loss: 566839.7500
Epoch 95/200
45/45 [=====] - 0s 6ms/step - loss: 557469.4429 - val_loss: 574953.5000
Epoch 96/200
45/45 [=====] - 0s 4ms/step - loss: 546190.5815 - val_loss: 567444.0000
Epoch 97/200
45/45 [=====] - 0s 6ms/step - loss: 544921.4457 - val_loss: 559030.3750
Epoch 98/200
45/45 [=====] - 0s 4ms/step - loss: 540614.6834 - val_loss: 554367.2500
Epoch 99/200
1/45 [.....] - ETA: 0s - loss: 609462.5625
```

It takes some time for results to appear because model is running at the backend.



Predicted OUTPUT value:

## NNFL IA-1



Thankyou!!!