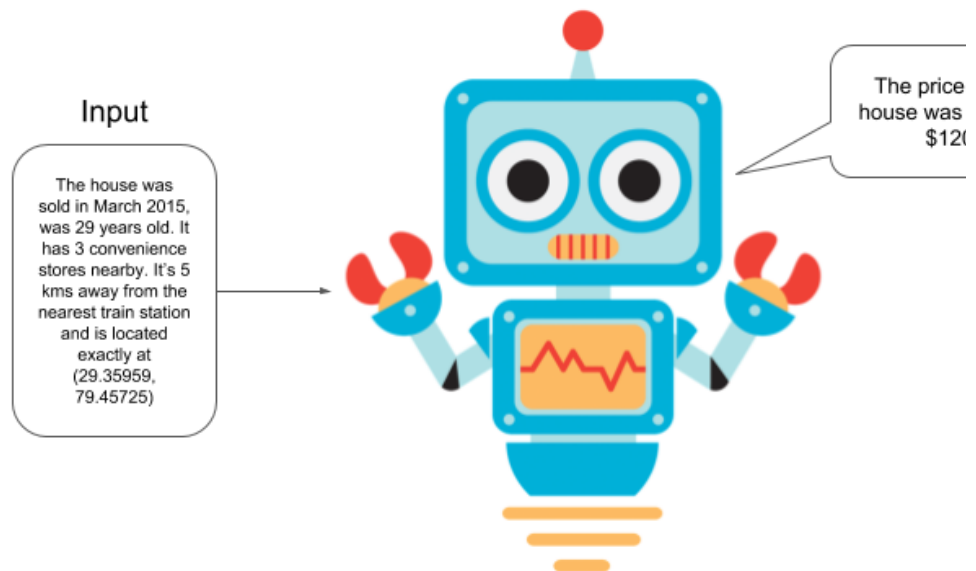


## Task 1: Introduction

For this project, we are going to work on evaluating price of houses given the f

1. Year of sale of the house
2. The age of the house at the time of sale
3. Distance from city center
4. Number of stores in the locality
5. The latitude
6. The longitude



Note: This notebook uses python 3 and these packages: tensorflow , pandas , matplotlib , sklearn .

### 1.1: Importing Libraries & Helper Functions

First of all, we will need to import some libraries and helper functions. This includes some utility functions that I've written to save time.

```
In [1]: import pandas as pd
import matplotlib.pyplot as plt
import tensorflow as tf

from utils import *
from sklearn.model_selection import train_test_split
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Dense, Dropout
from tensorflow.keras.callbacks import EarlyStopping, LambdaCallback

%matplotlib inline
tf.logging.set_verbosity(tf.logging.ERROR)

print('Libraries imported.')

Libraries imported.
```

## Task 2: Importing the Data

### 2.1: Importing the Data

The dataset is saved in a `data.csv` file. We will use `pandas` to take a look at son

```
In [2]: df= pd.read_csv('data.csv', names= column_names)
df.head()
```

	serial	date	age	distance	stores	latitude	longitude	price
0	0	2009	21	9	6	84	121	14264
1	1	2007	4	2	3	86	121	12032
2	2	2016	18	3	7	90	120	13560
3	3	2002	13	2	2	80	128	12029
4	4	2014	25	5	8	81	122	14157

### 2.2: Check Missing Data

It's a good practice to check if the data has any missing values. In real world data is common and must be taken care of before any data pre-processing or model training.

```
In [3]: df.isna().sum()

serial      0
date        0
age         0
distance    0
stores      0
latitude    0
longitude   0
price       0
dtype: int64
```

## Task 3: Data Normalization

### 3.1: Data Normalization

We can make it easier for optimization algorithms to find minimas by normalizing training a model.

```
In [4]: df = df.iloc[:, 1:] #1 onwards column, all rows
df_norm = (df - df.mean()) / df.std() #columnwise mean & std
df_norm.head()
```

	date	age	distance	stores	latitude	longitude	price
0	0.015978	0.181384	1.257002	0.345224	-0.307212	-1.260799	0.350088
1	-0.350485	-1.319118	-0.930610	-0.609312	0.325301	-1.260799	-1.836486
2	1.298598	-0.083410	-0.618094	0.663402	1.590328	-1.576456	-0.339584
3	-1.266643	-0.524735	-0.930610	-0.927491	-1.572238	0.948803	-1.839425
4	0.932135	0.534444	0.006938	0.981581	-1.255981	-0.945141	0.245266

### 3.2: Convert Label Value

Because we are using normalized values for the labels, we will get the predicted values from the trained model in the same distribution. So, we need to convert the predicted values back to the original distribution if we want predicted prices.

```
In [5]: y_mean = df['price'].mean()
y_std = df['price'].std()

def convert_label_value(pred):
    return int(pred * y_std + y_mean)

print(convert_label_value(0.350088))

14263
```

## Task 4: Create Training and Test Sets

## 4.1: Select Features

Make sure to remove the column **price** from the list of features as it is the label used as a feature.

```
In [6]: X= df_norm.iloc[:, :6]
X.head()
```

	date	age	distance	stores	latitude	longitude
0	0.015978	0.181384	1.257002	0.345224	-0.307212	-1.260799
1	-0.350485	-1.319118	-0.930610	-0.609312	0.325301	-1.260799
2	1.298598	-0.083410	-0.618094	0.663402	1.590328	-1.576456
3	-1.266643	-0.524735	-0.930610	-0.927491	-1.572238	0.948803
4	0.932135	0.534444	0.006938	0.981581	-1.255981	-0.945141

## 4.2: Select Labels

```
In [8]: y= df_norm.iloc[:, -1]
y.head()

0    0.350088
1   -1.836486
2   -0.339584
3   -1.839425
4    0.245266
Name: price, dtype: float64
```

## 4.3: Feature and Label Values

We will need to extract just the numeric values for the features and labels as the model will expect just numeric values as input.

```
In [11]: X_arr= X.values
         y_arr=y.values

print('Features array shape: ',X_arr.shape)
print('Labels array shape: ',y_arr.shape)

Features array shape: (5000, 6)
Labels array shape: (5000,)

array([[ 0.01597778,  0.18138426,  1.25700164,  0.34522379, -0.30721158,
        -1.26079862],
       [-0.35048517, -1.31911814, -0.93060999, -0.60931203,  0.32530146,
        -1.26079862],
       [ 1.29859812, -0.08341028, -0.61809404,  0.66340239,  1.59032754,
        -1.57645598],
       ...,
       [ 1.4818296 , -1.14258845,  1.56951759,  0.02704518,  1.59032754,
         0.00183081],
       [ 0.19920926,  1.59362182, -0.61809404,  0.02704518, -1.25598114,
         0.94880289],
       [ 1.66506107, -0.87779391, -1.24312594,  1.2997596 ,  1.59032754,
         0.63314553]])
```

## 4.4: Train and Test Split

We will keep some part of the data aside as a **test** set. The model will not use it for training and it will be used only for checking the performance of the model in the trained states. This way, we can make sure that we are going in the right direction for training.

```
In [12]: X_train,X_test,y_train,y_test= train_test_split(X_arr, y_arr, test_size=0.05, random_state=42)
print('Training Set: ',X_train.shape, y_train.shape)
print('Test Set: ',X_test.shape, y_test.shape)

Training Set: (4750, 6) (4750,)
Test Set: (250, 6) (250,)
```

## Task 5: Create the Model

### 5.1: Create the Model

Let's write a function that returns an untrained model of a certain architecture.

```
In [13]: def get_model():
          model=Sequential([
              Dense(10, input_shape=(6,), activation='relu'),
              Dense(20, activation='relu'),
              Dense(5, activation='relu'),
              Dense(1)
          ])#since regression prob no need of activation in last layer, rectified lin
          model.compile(loss='mse',
                        optimizer='adam')
          #loss func. that optimization algorithm tries to minimize
          return model
```

```
get_model().summary()
```

Layer (type)	Output Shape	Param #
dense (Dense)	(None, 10)	70
dense_1 (Dense)	(None, 20)	220
dense_2 (Dense)	(None, 5)	105
dense_3 (Dense)	(None, 1)	6
Total params: 401		
Trainable params: 401		
Non-trainable params: 0		

## Task 6: Model Training

### 6.1: Model Training

We can use an `EarlyStopping` callback from Keras to stop the model training if it stops decreasing for a few epochs.

```
In [14]: es_cb= EarlyStopping(monitor='val_loss', patience=5)
#es_cb will moniot validation loss and wait for 5 epochs before it decides to s

model= get_model()
preds_on_untrained= model.predict(X_test)

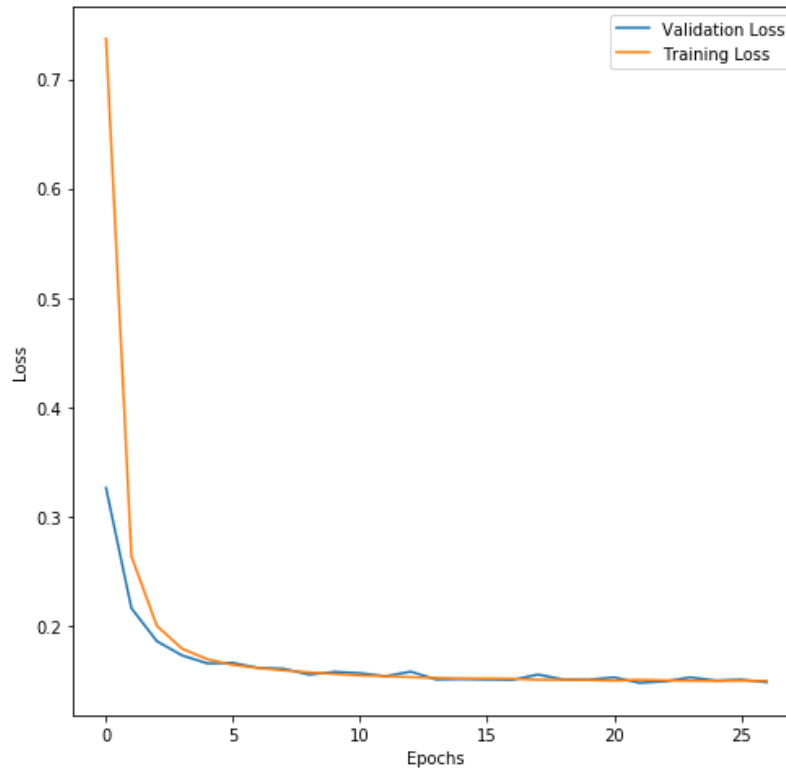
history= model.fit(X_train, y_train, validation_data= (X_test, y_test),
                  epochs=100,
                  callbacks=[es_cb])
```

```
Train on 4750 samples, validate on 250 samples
Epoch 1/100
4750/4750 [=====] - 1s 160us/sample - loss: 0.7370 - val_loss: 0.3265
Epoch 2/100
4750/4750 [=====] - 0s 59us/sample - loss: 0.2640 - val_loss: 0.2165
Epoch 3/100
4750/4750 [=====] - 0s 66us/sample - loss: 0.2004 - val_loss: 0.1864
Epoch 4/100
4750/4750 [=====] - 0s 61us/sample - loss: 0.1793 - val_loss: 0.1733
Epoch 5/100
4750/4750 [=====] - 0s 37us/sample - loss: 0.1698 - val_loss: 0.1660
Epoch 6/100
4750/4750 [=====] - 0s 39us/sample - loss: 0.1646 - val_loss: 0.1665
Epoch 7/100
4750/4750 [=====] - 0s 40us/sample - loss: 0.1618 - val_loss: 0.1619
Epoch 8/100
4750/4750 [=====] - 0s 40us/sample - loss: 0.1596 - val_loss: 0.1612
Epoch 9/100
4750/4750 [=====] - 0s 40us/sample - loss: 0.1578 - val_loss: 0.1557
Epoch 10/100
4750/4750 [=====] - 0s 42us/sample - loss: 0.1564 - val_loss: 0.1583
Epoch 11/100
4750/4750 [=====] - 0s 41us/sample - loss: 0.1551 - val_loss: 0.1571
Epoch 12/100
4750/4750 [=====] - 0s 43us/sample - loss: 0.1543 - val_loss: 0.1544
Epoch 13/100
4750/4750 [=====] - 0s 41us/sample - loss: 0.1535 - val_loss: 0.1585
Epoch 14/100
4750/4750 [=====] - 0s 46us/sample - loss: 0.1528 - val_loss: 0.1513
Epoch 15/100
4750/4750 [=====] - 0s 61us/sample - loss: 0.1523 - val_loss: 0.1517
Epoch 16/100
4750/4750 [=====] - 0s 41us/sample - loss: 0.1522 - val_loss: 0.1514
Epoch 17/100
4750/4750 [=====] - 0s 43us/sample - loss: 0.1520 - val_loss: 0.1511
Epoch 18/100
4750/4750 [=====] - 0s 70us/sample - loss: 0.1511 - val_loss: 0.1559
Epoch 19/100
4750/4750 [=====] - 0s 49us/sample - loss: 0.1512 - val_loss: 0.1514
Epoch 20/100
4750/4750 [=====] - 0s 43us/sample - loss: 0.1509 - val_loss: 0.1512
Epoch 21/100
4750/4750 [=====] - 0s 59us/sample - loss: 0.1503 - val_loss: 0.1534
Epoch 22/100
4750/4750 [=====] - 0s 52us/sample - loss: 0.1511 - val_loss: 0.1482
Epoch 23/100
4750/4750 [=====] - 0s 53us/sample - loss: 0.1505 - val_loss: 0.1496
Epoch 24/100
4750/4750 [=====] - 0s 70us/sample - loss: 0.1503 - val_loss: 0.1533
Epoch 25/100
4750/4750 [=====] - 0s 44us/sample - loss: 0.1500 - val_loss: 0.1504
Epoch 26/100
4750/4750 [=====] - 0s 76us/sample - loss: 0.1503 - val_loss: 0.1513
Epoch 27/100
```

## 6.2: Plot Training and Validation Loss

Let's use the `plot_loss` helper function to take a look training and validation lo

```
In [15]: plot_loss(history)
```



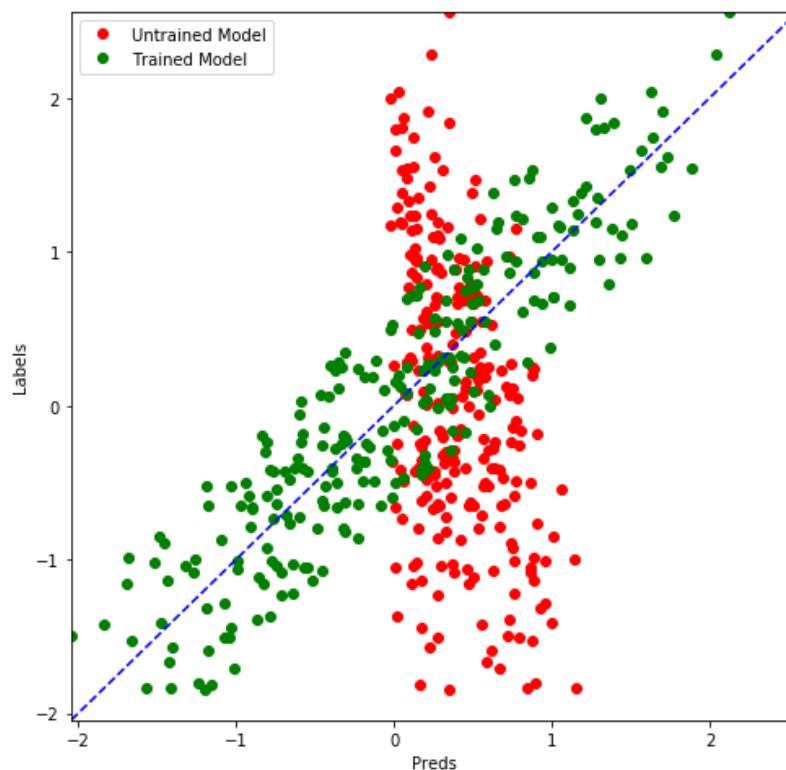
## Task 7: Predictions

### 7.1: Plot Raw Predictions

Let's use the `compare_predictions` helper function to compare predictions from was untrained and when it was trained.



```
In [16]: preds_on_trained= model.predict(X_test)
compare_predictions(preds_on_untrained, preds_on_trained, y_test)
#blue line is the original label line
```



## 7.2: Plot Price Predictions

The plot for price predictions and raw predictions will look the same with just o  
x and y axis scale is changed.

```
In [17]: price_untrained= [convert_label_value(y) for y in preds_on_untrained]
price_trained= [convert_label_value(y) for y in preds_on_trained]
price_test= [convert_label_value(y) for y in y_test]

compare_predictions(price_untrained, price_trained, price_test)
```

