```
In [ ]:
from google.colab import drive
drive.mount('/content/drive')
Mounted at /content/drive
In [ ]:
import os
os.chdir("/content/drive/My Drive/Classroom/projects/Mercari")
!ls -1
total 7772458
-rw----- 1 root root
                           151 Nov 19 17:35 akarshan.1711@gmail.com CS1.gdoc
-rw----- 1 root root
                        192263 Jan 2 21:08 'Copy of HptTfidf2.ipynb'
                         151 Dec 16 13:22 EDA+FE.gdoc
-rw----- 1 root root
-rw----- 1 root root
                         2441752 Dec 20 16:29 EDA.ipynb
-rw----- 1 root root
                         14393 Dec 27 21:06 FE+prep+modelling.ipynb
-rw----- 1 root root
                          30163 Dec 29 18:34 HptBrnandImpute.v1.0.ipynb
-rw----- 1 root root
                         249493 Jan 2 20:56 HptTfidf2.ipynb
                         192396 Jan 2 21:54
-rw----- 1 root root
                                              HptTfidf.ipynb
-rw----- 1 root root 117131678 Jan 1 12:07
                                              lgbt2.csv
-rw----- 1 root root 68399264 Jan 1 02:00 lgbt3.csv
-rw----- 1 root root 927353 Dec 28 15:17
                                              mercari_mainV2.ipynb
-rw----- 1 root root
                        380928 Jan 2 21:54 Mercari_to3.db
-rw----- 1 root root
                         77824 Jan 2 14:07 Mercari_to4.db
-rw----- 1 root root
                        249856 Jan 2 20:56 Mercari_to5.db
                       196608 Jan 2 21:08 Mercari_to6.db
-rw----- 1 root root
-rw----- 1 root root 11853944 Dec 30 21:08 price log2.pickle
-rw----- 1 root root 11853944 Dec 31 07:52 price log.pickle
-rw----- 1 root root
                          23640 Jan 2 21:57 Stack.ipynb
-rw----- 1 root root 308669128 Dec 10 2019 test stg2.tsv.zip
-rw----- 1 root root 3474387330 Dec 30 21:08 tfidf2.pickle
-rw----- 1 root root 3623909034 Dec 30 20:50 tfidf.pickle
-rw----- 1 root root 337809843 Nov 11 2017 train.tsv
-rw----- 1 root root
                            272 Jan 2 19:42 Untitled
In [ ]:
#importing modules/libraries
import pandas as pd
import numpy as np
import scipy
import seaborn as sns
import matplotlib.pyplot as plt
import gc
import sys
import os
import psutil
# from scipy.stats import randint as sp randint
# from scipy.stats import uniform as sp uniform
from tqdm.notebook import tqdm
# from collections import Counter
# from collections import defaultdict
import re
import random
# from random import sample
# from bs4 import BeautifulSoup
import pickle
import inspect
import time
import sklearn
from sklearn.feature_extraction.text import TfidfVectorizer,CountVectorizer
from sklearn.model selection import train test split
from sklearn.preprocessing import StandardScaler, LabelBinarizer
```

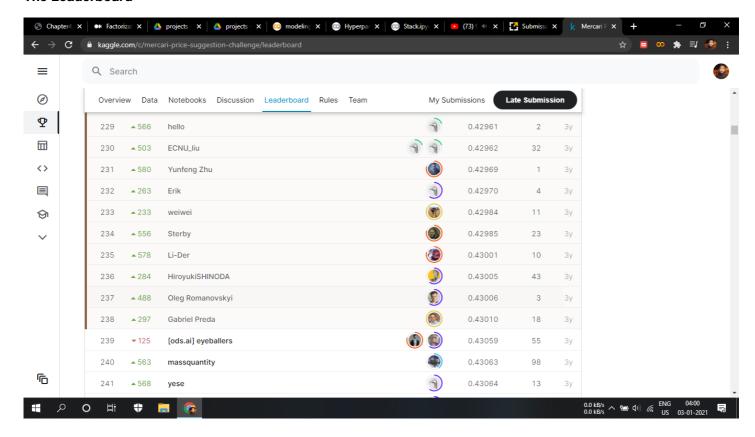
```
from sklearn.model_selection import RandomizedSearchCV
from sklearn.model_selection import GridSearchCV
from sklearn.metrics import mean squared error
import lightgbm as lgb
from sklearn.linear model import Lasso,Ridge
# import string
# # import emoji
# # from wordcloud import WordCloud
# import nltk
# nltk.download("stopwords")
# # nltk.download("brown")
# # nltk.download("names")
# # nltk.download('punkt')
# nltk.download('wordnet')
# # nltk.download('averaged perceptron_tagger')
# # nltk.download('universal_tagset')
# # from nltk.tokenize import word tokenize
# from nltk.corpus import stopwords
# from nltk.stem.wordnet import WordNetLemmatizer
# # from nltk.stem.porter import PorterStemmer
import warnings
warnings.filterwarnings("ignore")
In [ ]:
# defining root mean square error over Log transformed y test data
# (as linear models homoscedasticity can be kept in check for better prediction)
# and hence an effective Root Mean Square Log Error
def error(y test, predictions):
  return np.sqrt(mean squared error( y test, predictions ))
In [ ]:
# making a submmissin file compatible with kaggels submission format
def submission(model, file name):
  test pr = pd.read csv('test stg2.tsv.zip', sep='\t', usecols = ['test id'] )
  with open('tfidf.pickle','rb') as f:
    data = pickle.load(f)
  # loading only the test data that does not have a target values/ test data from kaggle
  data = data[1481661:,:]
  # making prediciton on them and saving them with tets id
  test pr['price'] = np.expm1(model.predict(data))
  test pr.to csv(file name+'.csv')
In [ ]:
tr len = 1185329# demarkation of cv data(0.8 percent)
whole tr = 1481661# whole train data
In [ ]:
with open('tfidf.pickle','rb') as f:
  df=pickle.load(f)
with open('price log.pickle','rb') as f:
  y=pickle.load(f)
In [ ]:
df = df[:whole tr] # only taking train and cv data
gc.collect()
Out[]:
```

376

```
In [ ]:
X train, X test, Y train, Y test = train test split(df,y, train size = round(0.8*df.shap
e[0]))
In [ ]:
model= Ridge(alpha=4.5, max iter=10000, tol=0.0005, solver='auto', random state=34)
model.fit(X_train, Y_train)
op rdg1 = model.predict(X test)
In [ ]:
error(op rdg1,Y test)
Out[]:
0.4417180108419465
In [ ]:
# this submission helped me achieve 0.4488 rmsle on kaggle
submission(modle, 'ridge')
In [ ]:
gc.collect()
lgbt = lgb.LGBMRegressor(boosting type='gbdt', objective='regression',
                                               random state=11, n jobs=4, subsample for bin=81920,
                                               learning rate=0.19989416216581643, num leaves=137, ma
x depth=40, n estimators=2194)
lgbt.fit(X train, Y train)
op lgb1 = lgbt.predict(X test)
In [ ]:
# this submission helped me achieve 0.4209 rmsle on kaggle
error(op lgb1,Y test)
Out[]:
0.4209546895125615
In [ ]:
submission(lgbt, 'lgbt2')
 📀 Chapterd: X | 🖦 Factoriza: X | 🚵 projects - X | 🚵 projects - X | 🚳 projects - X | 🔞 Myperpair X | 🔞 Stackipy: X | 📵 (73) 🖽 X | 🚰 Submissi X
 \leftarrow \rightarrow \mathbf{C} ^{\circ} kaggle.com/c/mercari-price-suggestion-challenge/submissions
           Q Search
  \equiv
  0
            Overview Data Notebooks Discussion Leaderboard Rules Team
                                                                    My Submissions
                                                                                 Late Submission
  Ψ
            All Successful Selected
  Submission and Description
                                                          Private Score
                                                                      Public Score
                                                                                Use for Final Score
  <>
                                                                       0.42990
            mockrey
                                                            0.43016
  (version 14/14)
            a day ago by Akarshan kumar
  0
            From "mockrey" Notebook
                                                            0.44904
                                                                       0.44881
            (version 13/14)
            2 days ago by Akarshan kumai
            From "mockrey" Notebook
                                                            0.44976
                                                                       0.44951
            mockrey
            (version 9/14)
             4 days ago by Akarshan kuma
            From "mockrey" Notebook
            mockrev
                                                            0.48692
                                                                       0.48613
            (version 8/14)
```



#### The Leaderboard



## In [ ]:

#### Out[]:

LGBMRegressor(boosting\_type='gbdt', class\_weight=None, colsample\_bytree=1.0, importance\_type='split', learning\_rate=0.2, max\_depth=40, min\_child\_samples=20, min\_child\_weight=0.001, min\_split\_gain=0.0, n\_estimators=2200, n\_jobs=4, num\_leaves=140, objective=None, random\_state=11, reg\_alpha=0.0, reg\_lambda=0.0, silent=True, subsample=1.0, subsample for bin=196892, subsample freq=0)

#### In [ ]:

```
# second best lgbm model
op_lgb1 = lgbt.predict(X_test)
error(op_lgb1,Y_test)
```

#### Out[]:

0.4210177447417679

### In [ ]:

# In [ ]:

```
op_lgb1 = lgbt.predict(X_test)
error(op_lgb1,Y_test)
```