

## ▼ DonorsChoose

DonorsChoose.org receives hundreds of thousands of project proposals each year for classroom projects in need of funding. Right now, a large number of volunteers is needed to manually screen each submission before it's approved to be posted on the DonorsChoose.org website.

Next year, DonorsChoose.org expects to receive close to 500,000 project proposals. As a result, there are three main problems they need to solve:

- How to scale current manual processes and resources to screen 500,000 projects so that they can be posted as quickly and as efficiently as possible
- How to increase the consistency of project vetting across different volunteers to improve the experience for teachers
- How to focus volunteer time on the applications that need the most assistance

The goal of the competition is to predict whether or not a DonorsChoose.org project proposal submitted by a teacher will be approved, using the text of project descriptions as well as additional metadata about the project, teacher, and school. DonorsChoose.org can then use this information to identify projects most likely to need further review before approval.

## ▼ About the DonorsChoose Data Set

The train.csv data set provided by DonorsChoose contains the following features:

Feature	Description
project_id	A unique identifier for the proposed project. <b>Example:</b> p036502
project_title	Title of the project. <b>Examples:</b> <ul style="list-style-type: none"><li>• Art Will Make You Happy!</li><li>• First Grade Fun</li></ul>
project_grade_category	Grade level of students for which the project is targeted. One of the following enumerated values: <ul style="list-style-type: none"><li>• Grades PreK-2</li><li>• Grades 3-5</li><li>• Grades 6-8</li><li>• Grades 9-12</li></ul>
project_subject_categories	One or more (comma-separated) subject categories for the project from the following enumerated list of values: <ul style="list-style-type: none"><li>• Applied Learning</li><li>• Care &amp; Hunger</li><li>• Health &amp; Sports</li><li>• History &amp; Civics</li><li>• Literacy &amp; Language</li><li>• Math &amp; Science</li><li>• Music &amp; The Arts</li><li>• Special Needs</li><li>• Warmth</li></ul> <b>Examples:</b> <ul style="list-style-type: none"><li>• Music &amp; The Arts</li><li>• Literacy &amp; Language, Math &amp; Science</li></ul>
school_state	State where school is located ( <a href="#">Two-letter U.S. postal code</a> ). <b>Example:</b> WY
project_subject_subcategories	One or more (comma-separated) subject subcategories for the project. <b>Examples:</b> <ul style="list-style-type: none"><li>• Literacy</li><li>• Literature &amp; Writing, Social Sciences</li></ul>

Feature	Description
project_resource_summary	An explanation of the resources needed for the project. <b>Example:</b> <ul style="list-style-type: none"> <li>My students need hands on literacy materials to manage sensory needs!</li> </ul>
project_essay_1	First application essay*
project_essay_2	Second application essay*
project_essay_3	Third application essay*
project_essay_4	Fourth application essay*
project_submitted_datetime	Datetime when project application was submitted. <b>Example:</b> 2016-04-28 12:43:56.245
teacher_id	A unique identifier for the teacher of the proposed project. <b>Example:</b> bdf8baa8fedef6bfeec7ae4ff1c15c56
teacher_prefix	Teacher's title. One of the following enumerated values: <ul style="list-style-type: none"> <li>nan</li> <li>Dr.</li> <li>Mr.</li> <li>Mrs.</li> <li>Ms.</li> <li>Teacher.</li> </ul>
teacher_number_of_previously_posted_projects	Number of project applications previously submitted by the same teacher. <b>Example:</b> 2

\* See the section **Notes on the Essay Data** for more details about these features.

Additionally, the `resources.csv` data set provides more data about the resources required for each project. Each line in this file represents a resource required by a project:

Feature	Description
<code>id</code>	A <code>project_id</code> value from the <code>train.csv</code> file. <b>Example:</b> p036502
<code>description</code>	Description of the resource. <b>Example:</b> Tenor Saxophone Reeds, Box of 25
<code>quantity</code>	Quantity of the resource required. <b>Example:</b> 3
<code>price</code>	Price of the resource required. <b>Example:</b> 9.95

**Note:** Many projects require multiple resources. The `id` value corresponds to a `project_id` in `train.csv`, so you use it as a key to retrieve all resources needed for a project:

The data set contains the following label (the value you will attempt to predict):

Label	Description
-------	-------------

## ▼ Notes on the Essay Data

Prior to May 17, 2016, the prompts for the essays were as follows:

- \_\_project\_essay\_1:\_\_ "Introduce us to your classroom"
- \_\_project\_essay\_2:\_\_ "Tell us more about your students"
- \_\_project\_essay\_3:\_\_ "Describe how your students will use the materials you're requesting"
- \_\_project\_essay\_3:\_\_ "Close by sharing why your project will make a difference"

Starting on May 17, 2016, the number of essays was reduced from 4 to 2, and the prompts for the first 2 essays were changed to the following:

- \_\_project\_essay\_1:\_\_ "Describe your students: What makes your students special? Specific details about their background, your neighborhood, and your school are all helpful."
- \_\_project\_essay\_2:\_\_ "About your project: How will these materials make a difference in your students' learning and improve their school lives?"

For all projects with `project_submitted_datetime` of 2016-05-17 and later, the values of `project_essay_3` and `project_essay_4` will be NaN.

```
1 from google.colab import drive
2 drive.mount('/content/drive')

[Drive already mounted at /content/drive; to attempt to forcibly remount, call drive.mount("/content/drive", force_remount=True).]

1 %matplotlib inline
2 import warnings
3 warnings.filterwarnings("ignore")
4 import sqlite3
5 import pandas as pd
6 import numpy as np
7 import nltk
8 import string
9 import matplotlib.pyplot as plt
10 import seaborn as sns
11 from sklearn.feature_extraction.text import TfidfTransformer
12 from sklearn.feature_extraction.text import TfidfVectorizer
13 from sklearn.feature_extraction.text import CountVectorizer
14 from sklearn.metrics import confusion_matrix
15 from sklearn import metrics
16 from sklearn.metrics import roc_curve, auc
17 from nltk.stem.porter import PorterStemmer
18 import re
19 # Tutorial about Python regular expressions: https://pymotw.com/2/re/
20 import string
21 from nltk.corpus import stopwords
22 from nltk.stem import PorterStemmer
23 from nltk.stem.wordnet import WordNetLemmatizer
24 from gensim.models import Word2Vec
25 from gensim.models import KeyedVectors
26 import pickle
27 from tqdm import tqdm
28 import os
29 import chart_studio.plotly
30 # from plotly import plotly
31 import plotly.offline as offline
32 import plotly.graph_objs as go
33 offline.init_notebook_mode()
34 from collections import Counter
35 from scipy.sparse import hstack,vstack
36 from sklearn.model_selection import train_test_split
37 from sklearn.neighbors import KNeighborsClassifier
38 from sklearn.metrics import accuracy_score
39 from sklearn.model_selection import cross_val_score
40 from sklearn import model_selection
41 from sklearn.preprocessing import StandardScaler
42 from sklearn.model_selection import RandomizedSearchCV
43 #from sklearn.impute import SimpleImputer
44 from sklearn.datasets import load_digits
45 #from sklearn.feature_selection import SelectKBest, chi2
46 #from sklearn.model_selection import GridSearchCV
47 from sklearn.feature_selection import SelectKBest,f_classif
48 from prettytable import PrettyTable
49 from sklearn.naive_bayes import MultinomialNB
50 from sklearn.preprocessing import Normalizer
51 from sklearn.metrics import confusion_matrix
52 #import math
53 #from sklearn.linear_model import LogisticRegression
```

```

53 #from sklearn.linear_model import LogisticRegression
54 import nltk
55 from nltk.sentiment.vader import SentimentIntensityAnalyzer
56 #from sklearn.linear_model import SGDClassifier
57 from sklearn.tree import DecisionTreeClassifier
58 nltk.download('vader_lexicon')
59 import pdb
60 from sklearn.decomposition import TruncatedSVD
61
62 import graphviz
63 from sklearn import tree
64 from graphviz import Source
65 from sklearn.externals.six import StringIO
66 from IPython.display import Image
67 from sklearn.tree import export_graphviz
68 import pydotplus

```

## ▼ 1.1 Reading Data

```

1 Project_data = pd.read_csv('/content/drive/My Drive/Colab Notebooks/train_data.csv')
2 Resource_data = pd.read_csv('/content/drive/My Drive/Colab Notebooks/resources.csv')
3 print(Project_data.shape)
4 print(Resource_data.shape)

```

↳ (109248, 17)  
(1541272, 4)

```

1 # how to replace elements in list python: https://stackoverflow.com/a/2582163/4084039
2 cols = ['Date' if x=='project_submitted_datetime' else x for x in list(Project_data.columns)]
3 #sort datafram based on time pandas python: https://stackoverflow.com/a/49702492/4084039
4 Project_data['Date'] = pd.to_datetime(Project_data['project_submitted_datetime'])
5 Project_data.drop('project_submitted_datetime', axis=1, inplace=True)
6 Project_data.sort_values(by=['Date'], inplace=True)
7 # how to reorder columns pandas python: https://stackoverflow.com/a/13148611/4084039
8 Project_data = Project_data[cols]
9 Project_data.head(2)

```

↳

		Unnamed: 0	id	teacher_id	teacher_prefix	school_state	Date	project_grade_category	project_subject_categories	project_subject_subcategories	project_title	
55660	8393	p205479	2bf07ba08945e5d8b2a3f269b2b3cf5		Mrs.	CA	2016-04-27 00:27:36	Grades PreK-2	Math & Science	Applied Sciences, Health & Life Science	Engineering STEAM into the Primary Classroom	
76127	37728	p043609	3f60494c61921b3b43ab61bdde2904df		Ms.	UT	2016-04-27 00:31:25	Grades 3-5	Special Needs	Special Needs	Sensory Tools for Focus	

## ▼ 1.2 preprocessing of project\_subject\_categories

```

1 y = Project_data['project_is_approved'].values
2 Project_data.drop(['project_is_approved'], axis=1, inplace=True)
3 lnd = len(Project_data)

```

```

1 lpd = pd.read_csv('Project_data')
2 ys = np.zeros(lpd, dtype=np.int32)
3 X = Project_data

1 #Spliting the Dataset into three Train and Test
2 X_Train, X_Test, Y_Train, Y_Test = train_test_split(X, y, test_size=0.33, random_state=0, stratify=ys)
3 print('Shape of the X_Train data is {0} and Y_Train data is: {1}'.format(X_Train.shape,Y_Train.shape[0]))
4 print('Shape of the X_Test data is {0} and Y_Test data is : {1}'.format(X_Test.shape,Y_Test.shape[0]))

↳ Shape of the X_Train data is (73196, 16) and Y_Train data is: 73196
Shape of the X_Test data is (36052, 16) and Y_Test data is : 36052

1 # remove special characters from list of strings python: https://stackoverflow.com/a/47301924/4084039
2
3 # https://www.geeksforgeeks.org/removing-stop-words-nltk-python/
4 # https://stackoverflow.com/questions/23669024/how-to-strip-a-specific-word-from-a-string
5 # https://stackoverflow.com/questions/8270092/remove-all-whitespace-in-a-string-in-python
6
7 ****Train Data*****
8 categories = list(X_Train['project_subject_categories'].values)
9 cat_list = []
10 for i in categories:
11     temp = ""
12     # consider we have text like this "Math & Science, Warmth, Care & Hunger"
13     for j in i.split(','):# it will split it in three parts ["Math & Science", "Warmth", "Care & Hunger"]
14         if 'The' in j.split():# this will split each of the category based on space "Math & Science"=> "Math","&","Science"
15             j=j.replace('The','')# if we have the words "The" we are going to replace it with ''(i.e removing 'The')
16             j = j.replace(' ','')# we are placeing all the ' '(space) with ''(empty) ex:"Math & Science"=>"Math&Science"
17             temp+=j.strip()+" "# abc ".strip() will return "abc", remove the trailing spaces
18             temp = temp.replace('&','_')# we are replacing the & value into
19     cat_list.append(temp.strip())
20
21 X_Train['clean_categories'] = cat_list
22 X_Train.drop(['project_subject_categories'], axis=1, inplace=True)
23
24 from collections import Counter
25 my_counter = Counter()
26 for word in X_Train['clean_categories'].values:
27     my_counter.update(word.split())
28
29 cat_dict = dict(my_counter)
30 sorted_cat_dict_Train = dict(sorted(cat_dict.items(), key=lambda kv: kv[1]))
31 print(len(sorted_cat_dict_Train))
32
33 ****Test Data*****
34 categories = list(X_Test['project_subject_categories'].values)
35 cat_list = []
36 for i in categories:
37     temp = ""
38     # consider we have text like this "Math & Science, Warmth, Care & Hunger"
39     for j in i.split(','):# it will split it in three parts ["Math & Science", "Warmth", "Care & Hunger"]
40         if 'The' in j.split():# this will split each of the category based on space "Math & Science"=> "Math","&","Science"
41             j=j.replace('The','')# if we have the words "The" we are going to replace it with ''(i.e removing 'The')
42             j = j.replace(' ','')# we are placeing all the ' '(space) with ''(empty) ex:"Math & Science"=>"Math&Science"
43             temp+=j.strip()+" "# abc ".strip() will return "abc", remove the trailing spaces
44             temp = temp.replace('&','_')# we are replacing the & value into
45     cat_list.append(temp.strip())
46
47 X_Test['clean_categories'] = cat_list

```

```
48 X_Test.drop(['project_subject_categories'], axis=1, inplace=True)
```

49

↳ 9

## ▼ 1.3 preprocessing of project\_subject\_subcategories

```
1 # remove special characters from list of strings python: https://stackoverflow.com/a/47301924/4084039
2
3 # https://www.geeksforgeeks.org/removing-stop-words-nltk-python/
4 # https://stackoverflow.com/questions/23669024/how-to-strip-a-specific-word-from-a-string
5 # https://stackoverflow.com/questions/8270092/remove-all-whitespace-in-a-string-in-python
6 *****Train Data*****
7 sub_catogories = list(X_Train['project_subject_subcategories'].values)
8 sub_cat_list = []
9 for i in sub_catogories:
10    temp = ""
11    # consider we have text like this "Math & Science, Warmth, Care & Hunger"
12    for j in i.split(','):# it will split it in three parts ["Math & Science", "Warmth", "Care & Hunger"]
13        if 'The' in j.split():# this will split each of the category based on space "Math & Science"=> "Math","&","Science"
14            j=j.replace('The','')# if we have the words "The" we are going to replace it with ''(i.e removing 'The')
15            j = j.replace(' ','')# we are placeing all the ' '(space) with ''(empty) ex:"Math & Science"=>"Math&Science"
16            temp +=j.strip()+" "# abc ".strip() will return "abc", remove the trailing spaces
17            temp = temp.replace('&','_')
18    sub_cat_list.append(temp.strip())
19
20 X_Train['clean_subcategories'] = sub_cat_list
21 X_Train.drop(['project_subject_subcategories'], axis=1, inplace=True)
22
23 # count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039
24 my_counter = Counter()
25 for word in X_Train['clean_subcategories'].values:
26    my_counter.update(word.split())
27
28 sub_cat_dict = dict(my_counter)
29 sorted_sub_cat_dict_Train = dict(sorted(sub_cat_dict.items(), key=lambda kv: kv[1]))
30 print(len(sorted_sub_cat_dict_Train))
31
32 *****Test Data*****
33 sub_catogories = list(X_Test['project_subject_subcategories'].values)
34 sub_cat_list = []
35 for i in sub_catogories:
36    temp = ""
37    # consider we have text like this "Math & Science, Warmth, Care & Hunger"
38    for j in i.split(','):# it will split it in three parts ["Math & Science", "Warmth", "Care & Hunger"]
39        if 'The' in j.split():# this will split each of the category based on space "Math & Science"=> "Math","&","Science"
40            j=j.replace('The','')# if we have the words "The" we are going to replace it with ''(i.e removing 'The')
41            j = j.replace(' ','')# we are placeing all the ' '(space) with ''(empty) ex:"Math & Science"=>"Math&Science"
42            temp +=j.strip()+" "# abc ".strip() will return "abc", remove the trailing spaces
43            temp = temp.replace('&','_')
44    sub_cat_list.append(temp.strip())
45
46 X_Test['clean_subcategories'] = sub_cat_list
47 X_Test.drop(['project_subject_subcategories'], axis=1, inplace=True)
```

↳ 30

## ▼ 1.3 Text preprocessing

```
1 # merge two column text dataframe:  
2 X_Train["essay"] = X_Train["project_essay_1"].map(str) + \  
3             X_Train["project_essay_2"].map(str) + \  
4             X_Train["project_essay_3"].map(str) + \  
5             X_Train["project_essay_4"].map(str)  
6  
7 X_Test["essay"] = X_Test["project_essay_1"].map(str) + \  
8             X_Test["project_essay_2"].map(str) + \  
9             X_Test["project_essay_3"].map(str) + \  
10            X_Test["project_essay_4"].map(str)  
  
1 # https://stackoverflow.com/a/47091490/4084039  
2 import re  
3  
4 def decontracted(phrase):  
5     # specific  
6     phrase = re.sub(r"won't", "will not", phrase)  
7     phrase = re.sub(r"can't", "can not", phrase)  
8     # general  
9     phrase = re.sub(r"n't", " not", phrase)  
10    phrase = re.sub(r"\re", " are", phrase)  
11    phrase = re.sub(r"\s", " is", phrase)  
12    phrase = re.sub(r"\d", " would", phrase)  
13    phrase = re.sub(r"\ll", " will", phrase)  
14    phrase = re.sub(r"\t", " not", phrase)  
15    phrase = re.sub(r"\ve", " have", phrase)  
16    phrase = re.sub(r"\m", " am", phrase)  
17    return phrase  
  
1 # https://gist.github.com/sebleier/554280  
2 # we are removing the words from the stop words list: 'no', 'nor', 'not'  
3 stopwords = ['i', 'me', 'my', 'myself', 'we', 'our', 'ours', 'ourselves', 'you', "you're", "you've", \  
4     "you'll", "you'd", "your", 'yours', 'yourself', 'yourselves', 'he', 'him', 'his', 'himself', \  
5     'she', "she's", 'her', 'hers', 'herself', 'it', "it's", 'its', 'itself', 'they', 'them', 'their', \  
6     'theirs', 'themselves', 'what', 'which', 'who', 'whom', 'this', 'that', "that'll", 'these', 'those', \  
7     'am', 'is', 'are', 'was', 'were', 'be', 'been', 'being', 'have', 'has', 'had', 'having', 'do', 'does', \  
8     'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'or', 'because', 'as', 'until', 'while', 'of', \  
9     'at', 'by', 'for', 'with', 'about', 'against', 'between', 'into', 'through', 'during', 'before', 'after', \  
10    'above', 'below', 'to', 'from', 'up', 'down', 'in', 'out', 'on', 'off', 'over', 'under', 'again', 'further', \  
11    'then', 'once', 'here', 'there', 'when', 'where', 'why', 'how', 'all', 'any', 'both', 'each', 'few', 'more', \  
12    'most', 'other', 'some', 'such', 'only', 'own', 'same', 'so', 'than', 'too', 'very', \  
13    's', 't', 'can', 'will', 'just', 'don', "don't", 'should', "should've", 'now', 'd', 'll', 'm', 'o', 're', \  
14    've', 'y', 'ain', 'aren', "aren't", 'couldn', "couldn't", 'didn', "didn't", 'doesn', "doesn't", 'hadn', \  
15    "hadn't", 'hasn', "hasn't", 'haven', "haven't", 'isn', "isn't", 'ma', 'mightn', "mightn't", 'mustn', \  
16    "mustn't", 'needn', "needn't", 'shan', "shan't", 'shouldn', "shouldn't", 'wasn', "wasn't", 'weren', "weren't", \  
17    'won', "won't", 'wouldn', "wouldn't"]  
  
1 # Combining all the above stundents  
2 # tqdm is for printing the status bar  
3  
4 #-----PreProcessing of Essays in Train data set-----  
5 preprocessed_essays_Train = []  
6 for sentance in tqdm(X_Train['essay'].values):
```

```

7 sent = decontracted(sentance)
8 sent = sent.replace('\r', ' ')
9 sent = sent.replace('\'', ' ')
10 sent = sent.replace('\n', ' ')
11 sent = re.sub('[^A-Za-z0-9]+', ' ', sent)
12 # https://gist.github.com/sebleier/554280
13 sent = ' '.join(e for e in sent.split() if e.lower() not in stopwords)
14 preprocessed_essays_Train.append(sent.lower().strip())
15 # pdb.set_trace()
16
17 #-----PreProcessing of Essays in Test data set-----
18 preprocessed_essays_Test = []
19 for sentance in tqdm(X_Test['essay'].values):
20     sent = decontracted(sentance)
21     sent = sent.replace('\r', ' ')
22     sent = sent.replace('\'', ' ')
23     sent = sent.replace('\n', ' ')
24     sent = re.sub('[^A-Za-z0-9]+', ' ', sent)
25     # https://gist.github.com/sebleier/554280
26     sent = ' '.join(e for e in sent.split() if e.lower() not in stopwords)
27     preprocessed_essays_Test.append(sent.lower().strip())
28 # pdb.set_trace()

[?] 100%|██████████| 73196/73196 [00:40<00:00, 1798.75it/s]
[?] 100%|██████████| 36052/36052 [00:20<00:00, 1789.34it/s]

```

```

1 word_count_essay_Train = []
2 for a in tqdm(X_Train["essay"]) :
3     b = len(a.split())
4     word_count_essay_Train.append(b)
5
6 X_Train["word_count_essay_Train"] = word_count_essay_Train
7
8 word_count_essay_Test = []
9 for a in tqdm(X_Test["essay"]) :
10    b = len(a.split())
11    word_count_essay_Test.append(b)
12
13 X_Test["word_count_essay_Test"] = word_count_essay_Test

[?] 100%|██████████| 73196/73196 [00:01<00:00, 68071.92it/s]
[?] 100%|██████████| 36052/36052 [00:00<00:00, 67627.98it/s]

```

## 1.4 Preprocessing of project\_title

```

1 # Combining all the above stundents
2 # tqdm is for printing the status bar
3
4 #-----PreProcessing of Project Title in Train data set-----
5 preprocessed_titles_Train = []
6 for sentance in tqdm(X_Train['project_title'].values):
7     sent = decontracted(sentance)
8     sent = sent.replace('\r', ' ')
9     sent = sent.replace('\'', ' ')
10    sent = sent.replace('\n', ' ')
11    sent = re.sub('[^A-Za-z0-9]+', ' ', sent)
12    # https://gist.github.com/sebleier/554280

```

```

13     sent = ' '.join(e for e in sent.split() if e not in stopwords)
14     preprocessed_titles_Train.append(sent.lower().strip())
15 # pdb.set_trace()
16
17 #-----PreProcessing of Project Title in Test data set-----
18 preprocessed_titles_Test = []
19 for sentance in tqdm(X_Test['project_title'].values):
20     sent = decontracted(sentance)
21     sent = sent.replace('\\r', ' ')
22     sent = sent.replace('\\'', ' ')
23     sent = sent.replace('\\n', ' ')
24     sent = re.sub('[^A-Za-z0-9]+', ' ', sent)
25     # https://gist.github.com/sebleier/554280
26     sent = ' '.join(e for e in sent.split() if e not in stopwords)
27     preprocessed_titles_Test.append(sent.lower().strip())
28 # pdb.set_trace()

```

→ 100%|██████████| 73196/73196 [00:01<00:00, 41760.61it/s]  
100%|██████████| 36052/36052 [00:00<00:00, 41939.72it/s]

```

1 word_count_title_Train = []
2 for a in tqdm(X_Train["project_title"]) :
3     b = len(a.split())
4     word_count_title_Train.append(b)
5
6 X_Train["word_count_title_Train"] = word_count_title_Train
7
8
9 word_count_title_Test = []
10 for a in tqdm(X_Test["project_title"]) :
11     b = len(a.split())
12     word_count_title_Test.append(b)
13
14 X_Test["word_count_title_Test"] = word_count_title_Test

```

→ 100%|██████████| 73196/73196 [00:00<00:00, 849627.99it/s]  
100%|██████████| 36052/36052 [00:00<00:00, 874862.87it/s]

## ▼ 1.5 Preparing data for models

### ▼ 1.5.1 Vectorizing Categorical data

- <https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/handling-categorical-and-numerical-features/>

```

1 #-----Vectorizing categorical data for Train and Test-----
2
3 # we use count vectorizer to convert the values into one hot encoding
4 vectorizer_cat = CountVectorizer(vocabulary=list(sorted_cat_dict_Train.keys()), lowercase=False, binary=True)
5 vectorizer_cat.fit(X_Train['clean_categories'].values)
6 categories_one_hot_Train = vectorizer_cat.transform(X_Train['clean_categories'].values)
7 categories_one_hot_Test = vectorizer_cat.transform(X_Test['clean_categories'].values)
8 print(vectorizer_cat.get_feature_names())
9 print("-"*120)
10 print('Shape of Train dataset matrix after one hot encoding is: {}'.format(categories_one_hot_Train.shape))
11 print('Shape of Test dataset matrix after one hot encoding is: {}'.format(categories_one_hot_Test.shape))

```

```
↳ ['Warmth', 'Care_Hunger', 'History_Civics', 'Music_Arts', 'AppliedLearning', 'SpecialNeeds', 'Health_Sports', 'Math_Science', 'Literacy_Language']
```

```
-----  
Shape of Train dataset matrix after one hot encoding is: (73196, 9)  
Shape of Test dataset matrix after one hot encoding is: (36052, 9)
```

```
1 #-----Vectorizing Sub categorical data for Train and Test-----  
2  
3 # we use count vectorizer to convert the values into one hot encoding  
4 vectorizer_sub_cat = CountVectorizer(vocabulary=list(sorted_sub_cat_dict_Train.keys()), lowercase=False, binary=True)  
5 vectorizer_sub_cat.fit(X_Train['clean_categories'].values)  
6 sub_categories_one_hot_Train = vectorizer_sub_cat.transform(X_Train['clean_subcategories'].values)  
7 sub_categories_one_hot_Test = vectorizer_sub_cat.transform(X_Test['clean_subcategories'].values)  
8 print(vectorizer_sub_cat.get_feature_names())  
9 print("-"*120)  
10 print('Shape of Train dataset matrix after one hot encoding is: {}'.format(sub_categories_one_hot_Train.shape))  
11 print('Shape of Test dataset matrix after one hot encoding is: {}'.format(sub_categories_one_hot_Test.shape))
```

```
↳ ['Economics', 'CommunityService', 'FinancialLiteracy', 'ParentInvolvement', 'Civics_Government', 'Extracurricular', 'ForeignLanguages', 'Warmth', 'Care_Hunger', 'NutritionEducation', 'SocialScienc
```

```
-----  
Shape of Train dataset matrix after one hot encoding is: (73196, 30)  
Shape of Test dataset matrix after one hot encoding is: (36052, 30)
```

## School State

```
1 #-----Vectorizing categorical data of School state for Train dataset-----  
2  
3 school_catogories_Train = list(X_Train['school_state'].values)  
4 school_list_Train = []  
5 for sent in school_catogories_Train:  
6     school_list_Train.append(sent.lower().strip())  
7 X_Train['school_categories'] = school_list_Train  
8 X_Train.drop(['school_state'], axis=1, inplace=True)  
9  
10 # count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039  
11 my_counter_school_Train = Counter()  
12 for word in X_Train['school_categories'].values:  
13     my_counter_school_Train.update(word.split())  
14  
15 # dict sort by value python: https://stackoverflow.com/a/613218/4084039  
16 school_dict_Train = dict(my_counter_school_Train)  
17 sorted_school_dict_Train = dict(sorted(school_dict_Train.items(), key=lambda kv: kv[1]))  
18  
19 vectorizer_school = CountVectorizer(vocabulary=list(sorted_school_dict_Train.keys()), lowercase=False, binary=True)  
20 vectorizer_school.fit(X_Train['school_categories'].values)  
21 #print(vectorizer.get_feature_names())  
22  
23 school_one_hot_Train = vectorizer_school.transform(X_Train['school_categories'].values)  
24  
25 #-----Vectorizing categorical data of School state for Test dataset-----  
26  
27 school_catogories_Test = list(X_Test['school_state'].values)  
28 school_list_Test = []  
29 for sent in school_catogories_Test:  
30     school_list_Test.append(sent.lower().strip())  
31 X_Test['school_categories'] = school_list_Test  
32 X_Test.drop(['school_state'], axis=1, inplace=True)
```

```

33
34 # count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039
35 my_counter_school_Test = Counter()
36 for word in X_Test['school_categories'].values:
37     my_counter_school_Test.update(word.split())
38
39 # dict sort by value python: https://stackoverflow.com/a/613218/4084039
40 school_dict_Test = dict(my_counter_school_Test)
41 sorted_school_dict_Test = dict(sorted(school_dict_Test.items(), key=lambda kv: kv[1]))
42 school_one_hot_Test = vectorizer_school.transform(X_Test['school_categories'].values)
43 print("-"*120)
44 print('Shape of Train dataset matrix after one hot encoding is: {}'.format(school_one_hot_Train.shape))
45 print('Shape of Test dataset matrix after one hot encoding is: {}'.format(school_one_hot_Test.shape))

```

↳ -----  
 Shape of Train dataset matrix after one hot encoding is: (73196, 51)  
 Shape of Test dataset matrix after one hot encoding is: (36052, 51)

## Prefix

```

1 -----Vectorizing categorical data of Teacher Prefix for Train dataset-----
2
3 # remove special characters from list of strings python: https://stackoverflow.com/a/47301924/4084039
4 # https://www.geeksforgeeks.org/removing-stop-words-nltk-python/
5 # https://stackoverflow.com/questions/23669024/how-to-strip-a-specific-word-from-a-string
6 # https://stackoverflow.com/questions/8270092/remove-all-whitespace-in-a-string-in-python
7 prefix_catogories_Train = list(X_Train['teacher_prefix'].values)
8 prefix_list_Train = []
9 for sent in prefix_catogories_Train:
10    sent = re.sub('[^A-Za-z0-9]+', ' ', str(sent))
11    # https://gist.github.com/sebleier/554280
12    sent = ' '.join(e for e in sent.split())
13    prefix_list_Train.append(sent.lower().strip())
14 X_Train['prefix_catogories'] = prefix_list_Train
15 X_Train.drop(['teacher_prefix'], axis=1, inplace=True)
16
17 # count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039
18 my_counter_prefix_Train = Counter()
19 for word in X_Train['prefix_catogories'].values:
20     my_counter_prefix_Train.update(word.split())
21
22 # dict sort by value python: https://stackoverflow.com/a/613218/4084039
23 prefix_dict_Train = dict(my_counter_prefix_Train)
24 sorted_prefix_dict_Train = dict(sorted(prefix_dict_Train.items(), key=lambda kv: kv[1]))
25
26
27 vectorizer_prefix = CountVectorizer(vocabulary=list(sorted_prefix_dict_Train.keys()), lowercase=False, binary=True)
28 vectorizer_prefix.fit(X_Train['prefix_catogories'].values)
29 #print(vectorizer.get_feature_names())
30
31 prefix_one_hot_Train = vectorizer_prefix.transform(X_Train['prefix_catogories'].values)
32 #print("Shape of matrix after one hot encoding ",prefix_one_hot.shape)
33
34 -----Vectorizing categorical data of Teacher Prefix for Test dataset-----
35
36 prefix_catogories_Test = list(X_Test['teacher_prefix'].values)
37 prefix_list_Test = []
38 for sent in prefix_catogories_Test:

```

```

39 sent = re.sub('[^A-Za-z0-9]+', ' ', str(sent))
40 # https://gist.github.com/sebleier/554280
41 sent = ' '.join(e for e in sent.split())
42 prefix_list_Test.append(sent.lower().strip())
43 X_Test['prefix_catogories'] = prefix_list_Test
44 X_Test.drop(['teacher_prefix'], axis=1, inplace=True)
45
46 # count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039
47 my_counter_prefix_Test = Counter()
48 for word in X_Test['prefix_catogories'].values:
49     my_counter_prefix_Test.update(word.split())
50
51 # dict sort by value python: https://stackoverflow.com/a/613218/4084039
52 prefix_dict_Test = dict(my_counter_prefix_Test)
53 sorted_prefix_dict_Test = dict(sorted(prefix_dict_Test.items(), key=lambda kv: kv[1]))
54 prefix_one_hot_Test = vectorizer_prefix.transform(X_Test['prefix_catogories'].values)
55 print("-"*120)
56 print('Shape of Train dataset matrix after one hot encoding is: {}'.format(prefix_one_hot_Train.shape))
57 print('Shape of Test dataset matrix after one hot encoding is: {}'.format(prefix_one_hot_Test.shape))

```

→ -----  
 Shape of Train dataset matrix after one hot encoding is: (73196, 6)  
 Shape of Test dataset matrix after one hot encoding is: (36052, 6)

## project\_grade\_category

```

1 #-----Vectorizing categorical data of Project Grade for Train dataset-----
2
3 # remove special characters from list of strings python: https://stackoverflow.com/a/47301924/4084039
4 # https://www.geeksforgeeks.org/removing-stop-words-nltk-python/
5 # https://stackoverflow.com/questions/23669024/how-to-strip-a-specific-word-from-a-string
6 # https://stackoverflow.com/questions/8270092/remove-all-whitespace-in-a-string-in-python
7 grade_catogories_Train = list(X_Train['project_grade_category'].values)
8 grade_list_Train = []
9 for sent in grade_catogories_Train:
10    sent = sent.replace('-', '_')
11    sent = sent.replace(' ', '_')
12    # sent = re.sub('[^A-Za-z0-9]+', ' ', str(sent))
13    # https://gist.github.com/sebleier/554280
14    sent = ' '.join(e for e in sent.split())
15    grade_list_Train.append(sent.lower().strip())
16
17 # temp = temp.replace('-', '_')
18 X_Train['new_grade_category'] = grade_list_Train
19 X_Train.drop(['project_grade_category'], axis=1, inplace=True)
20
21 # count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039
22 my_counter_grade_Train = Counter()
23 for word in X_Train['new_grade_category'].values:
24     my_counter_grade_Train.update(word.split())
25
26 # dict sort by value python: https://stackoverflow.com/a/613218/4084039
27 grade_dict_Train = dict(my_counter_grade_Train)
28 sorted_grade_dict_Train = dict(sorted(grade_dict_Train.items(), key=lambda kv: kv[1]))
29
30 vectorizer_grade = CountVectorizer(vocabulary=list(sorted_grade_dict_Train.keys()), lowercase=False, binary=True)
31 vectorizer_grade.fit(X_Train['new_grade_category'].values)
32 #print(vectorizer.get_feature_names())
33

```

```

34 grade_one_hot_Train = vectorizer_grade.transform(X_Train['new_grade_category'].values)
35
36 #-----Vectorizing categorical data of Project Grade for Train dataset-----
37
38 grade_categories_Test = list(X_Test['project_grade_category'].values)
39 grade_list_Test = []
40 for sent in grade_categories_Test:
41     sent = sent.replace('-', '_')
42     sent = sent.replace(' ', '_')
43     # sent = re.sub('[^A-Za-z0-9]+', ' ', str(sent))
44     # https://gist.github.com/sebleier/554280
45     sent = ' '.join(e for e in sent.split())
46     grade_list_Test.append(sent.lower().strip())
47
48 # temp = temp.replace('-', '_')
49 X_Test['new_grade_category'] = grade_list_Test
50 X_Test.drop(['project_grade_category'], axis=1, inplace=True)
51
52 # count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039
53 my_counter_grade_Test = Counter()
54 for word in X_Test['new_grade_category'].values:
55     my_counter_grade_Test.update(word.split())
56
57 # dict sort by value python: https://stackoverflow.com/a/613218/4084039
58 grade_dict_Test = dict(my_counter_grade_Test)
59 sorted_grade_dict_Test = dict(sorted(grade_dict_Test.items(), key=lambda kv: kv[1]))
60
61 grade_one_hot_Test = vectorizer_grade.transform(X_Test['new_grade_category'].values)
62 print("-"*120)
63 print('Shape of Train dataset matrix after one hot encoding is: {}'.format(grade_one_hot_Train.shape))
64 print('Shape of Test dataset matrix after one hot encoding is: {}'.format(grade_one_hot_Test.shape))

```

↳ -----  
 Shape of Train dataset matrix after one hot encoding is: (73196, 4)  
 Shape of Test dataset matrix after one hot encoding is: (36052, 4)

## 1.5.2 Vectorizing Numerical features

```

1 price_data = Resource_data.groupby('id').agg({'price':'sum', 'quantity':'sum'}).reset_index()
2 X_Train = pd.merge(X_Train, price_data, on='id', how='left')
3 X_Test = pd.merge(X_Test, price_data, on='id', how='left')

1 price_norm = Normalizer(norm='l2', copy=False)
2 price_norm.fit(X_Train['price'].values.reshape(1,-1))
3
4 p=price_norm.transform(X_Train['price'].values.reshape(1,-1))
5 price_norm.transform(X_Test['price'].values.reshape(1,-1))
6 price_norm_Train = (X_Train['price'].values.reshape(-1,1))
7 price_norm_Test = (X_Test['price'].values.reshape(-1,1))
8 print("-"*120)
9 print('Shape of Train normalized price dataset matrix after one hot encoding is: {}'.format(price_norm_Train.shape))
10 print('Shape of Test normalized price dataset matrix after one hot encoding is: {}'.format(price_norm_Test.shape))

```

↳ -----  
 Shape of Train normalized price dataset matrix after one hot encoding is: (73196, 1)  
 Shape of Test normalized price dataset matrix after one hot encoding is: (36052, 1)

```
1 quantity_norm = Normalizer(norm='l2', copy=False)
2 quantity_norm.fit(X_Train['quantity'].values.reshape(1,-1))
3
4 quantity_norm.transform(X_Train['quantity'].values.reshape(1,-1))
5 quantity_norm.transform(X_Test['quantity'].values.reshape(1,-1))
6 quantity_norm_Train = quantity_norm.transform(X_Train['quantity'].values.reshape(-1,1))
7 quantity_norm_Test = quantity_norm.transform(X_Test['quantity'].values.reshape(-1,1))
8 print("-"*120)
9 print('Shape of Train normalized quantity dataset matrix after one hot encoding is: {}'.format(quantity_norm_Train.shape))
10 print('Shape of Test normalized quantity dataset matrix after one hot encoding is: {}'.format(quantity_norm_Test.shape))
```

↳ -----  
Shape of Train normalized quantity dataset matrix after one hot encoding is: (73196, 1)  
Shape of Test normalized quantity dataset matrix after one hot encoding is: (36052, 1)

```
1 teacher_prev_post_norm = Normalizer(norm='l2', copy=False)
2 teacher_prev_post_norm.fit(X_Train['teacher_number_of_previously_posted_projects'].values.reshape(1,-1))
3
4 teacher_prev_post_norm.transform(X_Train['teacher_number_of_previously_posted_projects'].values.reshape(1,-1))
5 teacher_prev_post_norm.transform(X_Test['teacher_number_of_previously_posted_projects'].values.reshape(1,-1))
6 teacher_prev_post_norm_Train = teacher_prev_post_norm.transform(X_Train['teacher_number_of_previously_posted_projects'].values.reshape(-1,1))
7 teacher_prev_post_norm_Test = teacher_prev_post_norm.transform(X_Test['teacher_number_of_previously_posted_projects'].values.reshape(-1,1))
8 print("-"*120)
9 print('Shape of Train normalized previously posted project dataset matrix after one hot encoding is: {}'.format(teacher_prev_post_norm_Train.shape))
10 print('Shape of Test normalized previously posted project dataset matrix after one hot encoding is: {}'.format(teacher_prev_post_norm_Test.shape))
```

↳ -----  
Shape of Train normalized previously posted project dataset matrix after one hot encoding is: (73196, 1)  
Shape of Test normalized previously posted project dataset matrix after one hot encoding is: (36052, 1)

```
1 title_norm = Normalizer(norm='l2', copy=False)
2 title_norm.fit(X_Train['word_count_title_Train'].values.reshape(1,-1))
3 title_norm.transform(X_Train['word_count_title_Train'].values.reshape(1,-1))
4 title_norm.transform(X_Test['word_count_title_Test'].values.reshape(1,-1))
5 word_count_title_Train = title_norm.transform(X_Train['word_count_title_Train'].values.reshape(-1,1))
6 word_count_title_Test = title_norm.transform(X_Test['word_count_title_Test'].values.reshape(-1,1))
7 print("-"*120)
8 print('Shape of Train normalized title dataset matrix after one hot encoding is: {}'.format(word_count_title_Train.shape))
9 print('Shape of Test normalized title dataset matrix after one hot encoding is: {}'.format(word_count_title_Test.shape))
```

↳ -----  
Shape of Train normalized title dataset matrix after one hot encoding is: (73196, 1)  
Shape of Test normalized title dataset matrix after one hot encoding is: (36052, 1)

```
1 essay_norm = Normalizer(norm='l2', copy=False)
2 essay_norm.fit(X_Train['word_count_essay_Train'].values.reshape(1,-1))
3 essay_norm.transform(X_Train['word_count_essay_Train'].values.reshape(1,-1))
4 essay_norm.transform(X_Test['word_count_essay_Test'].values.reshape(1,-1))
5 word_count_essay_Train = essay_norm.transform(X_Train['word_count_essay_Train'].values.reshape(-1,1))
6 word_count_essay_Test = essay_norm.transform(X_Test['word_count_essay_Test'].values.reshape(-1,1))
7 print("-"*120)
8 print('Shape of Train normalized title dataset matrix after one hot encoding is: {}'.format(word_count_essay_Train.shape))
9 print('Shape of Test normalized title dataset matrix after one hot encoding is: {}'.format(word_count_essay_Test.shape))
```

↳ -----  
Shape of Train normalized title dataset matrix after one hot encoding is: (73196, 1)  
Shape of Test normalized title dataset matrix after one hot encoding is: (36052, 1)

## ▼ 1.5.3 Vectorizing Text data

### ▼ 1.5.2.2 TFIDF vectorizer

```
1 from sklearn.feature_extraction.text import TfidfVectorizer
2 vectorizer_essays_tfidf = TfidfVectorizer(min_df=10)
3 text_tfidf_Train = vectorizer_essays_tfidf.fit_transform(preprocessed_essays_Train)
4 text_tfidf_Test = vectorizer_essays_tfidf.transform(preprocessed_essays_Test)
5 print("-"*120)
6 print("Applying TFIDF for Text Data")
7 print("-"*120)
8 print('Shape of Train dataset matrix after one hot encoding is: {}'.format(text_tfidf_Train.shape))
9 print('Shape of Test dataset matrix after one hot encoding is: {}'.format(text_tfidf_Test.shape))
```

→ -----  
Applying TFIDF for Text Data

```
-----  
Shape of Train dataset matrix after one hot encoding is: (73196, 14144)  
Shape of Test dataset matrix after one hot encoding is: (36052, 14144)
```

### TFIDF vectorizer for Project Title

```
1 vectorizer_titles_tfidf = TfidfVectorizer(min_df=10)
2 title_tfidf_Train = vectorizer_titles_tfidf.fit_transform(preprocessed_titles_Train)
3 title_tfidf_Test = vectorizer_titles_tfidf.transform(preprocessed_titles_Test)
4 print("-"*120)
5 print("Applying TFIDF for Project Title")
6 print("-"*120)
7 print('Shape of Train dataset matrix after one hot encoding is: {}'.format(title_tfidf_Train.shape))
8 print('Shape of Test dataset matrix after one hot encoding is: {}'.format(title_tfidf_Test.shape))
```

→ -----  
Applying TFIDF for Project Title

```
-----  
Shape of Train dataset matrix after one hot encoding is: (73196, 2631)  
Shape of Test dataset matrix after one hot encoding is: (36052, 2631)
```

```
1 # stronging variables into pickle files python: http://www.jessicayung.com/how-to-use-pickle-to-save-and-load-variables-in-python/
2 # make sure you have the glove_vectors file
3 #with open('glove_vectors', 'rb') as f:
4 #    model = pickle.load(f)
5 #    glove_words = set(model.keys())
6 with open('/content/drive/My Drive/Colab Notebooks/glove_vectors', 'rb') as f:
7     model = pickle.load(f)
8     glove_words = set(model.keys())
```

### ▼ 1.5.2.3 Using Pretrained Models: TFIDF weighted W2V

```
1 tfidf_model_essays = TfidfVectorizer()
2 tfidf_model_essays.fit(preprocessed_essays_Train)
3 # we are converting a dictionary with word as a key, and the idf as a value
4 dictionary = dict(zip(tfidf_model_essays.get_feature_names(), list(tfidf_model_essays.idf_)))
5 tfidf_words_essays = set(tfidf_model_essays.get_feature_names())
```

```

1 # average Word2Vec
2 # compute average word2vec for each review.
3 tfidf_w2v_vectors_Train = [] # the avg-w2v for each sentence/review is stored in this list
4 for sentence in tqdm(preprocessed_essays_Train): # for each review/sentence
5     vector = np.zeros(300) # as word vectors are of zero length
6     tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
7     for word in sentence.split(): # for each word in a review/sentence
8         if (word in glove_words) and (word in tfidf_words_essays):
9             vec = model[word] # getting the vector for each word
10            # here we are multiplying idf value(dictionary[word]) and the tf value((sentence.count(word)/len(sentence.split())))
11            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tfidf value for each word
12            vector += (vec * tf_idf) # calculating tfidf weighted w2v
13            tf_idf_weight += tf_idf
14    if tf_idf_weight != 0:
15        vector /= tf_idf_weight
16    tfidf_w2v_vectors_Train.append(vector)
17 #-----
18 tfidf_w2v_vectors_Test = [] # the avg-w2v for each sentence/review is stored in this list
19 for sentence in tqdm(preprocessed_essays_Test): # for each review/sentence
20     vector = np.zeros(300) # as word vectors are of zero length
21     tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
22     for word in sentence.split(): # for each word in a review/sentence
23         if (word in glove_words) and (word in tfidf_words_essays):
24             vec = model[word] # getting the vector for each word
25             # here we are multiplying idf value(dictionary[word]) and the tf value((sentence.count(word)/len(sentence.split())))
26             tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tfidf value for each word
27             vector += (vec * tf_idf) # calculating tfidf weighted w2v
28             tf_idf_weight += tf_idf
29    if tf_idf_weight != 0:
30        vector /= tf_idf_weight
31    tfidf_w2v_vectors_Test.append(vector)
32
33 print(len(tfidf_w2v_vectors_Test))
34 print(len(tfidf_w2v_vectors_Test[0]))
35

```

⌚ 100% | ██████████ | 73196/73196 [02:04<00:00, 585.70it/s]  
100% | ██████████ | 36052/36052 [01:01<00:00, 582.80it/s] 36052  
300

## Using Pretrained Models: TFIDF weighted W2V on project\_title

```

1 # Similarly you can vectorize for title also
2 tfidf_model_title = TfidfVectorizer()
3 tfidf_model_title.fit(preprocessed_titles_Train)
4 # we are converting a dictionary with word as a key, and the idf as a value
5 dictionary = dict(zip(tfidf_model_title.get_feature_names(), list(tfidf_model_title.idf_)))
6 tfidf_words_title = set(tfidf_model_title.get_feature_names())
7
8 # compute tfidf word2vec for each title.
9 tfidf_w2v_vectors_title_Train = [] # the avg-w2v for each sentence/review is stored in this list
10 for sentence in tqdm(preprocessed_titles_Train): # for each review/sentence
11     vector_title = np.zeros(300) # as word vectors are of zero length
12     tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
13     for word in sentence.split(): # for each word in a review/sentence
14         if (word in glove_words) and (word in tfidf_words_title):
15             vec = model[word] # getting the vector for each word
16             # here we are multiplying idf value(dictionary[word]) and the tf value((sentence.count(word)/len(sentence.split())))
17             tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tfidf value for each word
18             vector_title += (vec * tf_idf) # calculating tfidf weighted w2v
19             tf_idf_weight += tf_idf
20     if tf_idf_weight != 0:
21         vector_title /= tf_idf_weight
22     tfidf_w2v_vectors_title_Train.append(vector_title)
23
24 print(len(tfidf_w2v_vectors_title_Train))
25 print(len(tfidf_w2v_vectors_title_Train[0]))
26

```

```

15     vec = model[word] # getting the vector for each word
16     # here we are multiplying idf value(dictionary[word]) and the tf value((sentence.count(word)/len(sentence.split())))
17     tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tfidf value for each word
18     vector_title += (vector_title * tf_idf) # calculating tfidf weighted w2v
19     tf_idf_weight += tf_idf
20 if tf_idf_weight != 0:
21     vector_title /= tf_idf_weight
22 tfidf_w2v_vectors_title_Train.append(vector_title)
23
24 #-----
25
26
27 tfidf_w2v_vectors_title_Test = [] # the avg-w2v for each sentence/review is stored in this list
28 for sentence in tqdm(preprocessed_titles_Test): # for each review/sentence
29     vector_title = np.zeros(300) # as word vectors are of zero length
30     tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
31     for word in sentence.split(): # for each word in a review/sentence
32         if (word in glove_words) and (word in tfidf_words_title):
33             vec = model[word] # getting the vector for each word
34             # here we are multiplying idf value(dictionary[word]) and the tf value((sentence.count(word)/len(sentence.split())))
35             tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tfidf value for each word
36             vector_title += (vector_title * tf_idf) # calculating tfidf weighted w2v
37             tf_idf_weight += tf_idf
38 if tf_idf_weight != 0:
39     vector_title /= tf_idf_weight
40 tfidf_w2v_vectors_title_Test.append(vector_title)
41
42 print(len(tfidf_w2v_vectors_title_Test))
43 print(len(tfidf_w2v_vectors_title_Test[0]))
44
45
```

⌚ 100% |██████████| 73196/73196 [00:02<00:00, 31772.34it/s]  
⌚ 100% |██████████| 36052/36052 [00:01<00:00, 31641.90it/s] 36052  
300

## Calculating the sentiment score's of each of the essay

```

1
2 sid = SentimentIntensityAnalyzer()
3
4 essays = X_Train['essay']
5 essays_sentiment_TR_P = []
6 essays_sentiment_TR_N = []
7 essays_sentiment_TR_NE = []
8 essays_sentiment_TR_C = []
9 for essay in tqdm(essays):
10     res = sid.polarity_scores(essay)
11     essays_sentiment_TR_P.append(res['pos'])
12     essays_sentiment_TR_N.append(res['neg'])
13     essays_sentiment_TR_NE.append(res['neu'])
14     essays_sentiment_TR_C.append(res['compound'])
15 X_Train['sentiment_essay_TR_P'] = essays_sentiment_TR_P
16 X_Train['sentiment_essay_TR_N'] = essays_sentiment_TR_N
17 X_Train['sentiment_essay_TR_NE'] = essays_sentiment_TR_NE
18 X_Train['sentiment_essay_TR_C'] = essays_sentiment_TR_C
19
20
```

```

21 essays = X_Test['essay']
22 essays_sentiment_TS_P = []
23 essays_sentiment_TS_N = []
24 essays_sentiment_TS_NE = []
25 essays_sentiment_TS_C = []
26 for essay in tqdm(essays):
27     res = sid.polarity_scores(essay)
28     essays_sentiment_TS_P.append(res['pos'])
29     essays_sentiment_TS_N.append(res['neg'])
30     essays_sentiment_TS_NE.append(res['neu'])
31     essays_sentiment_TS_C.append(res['compound'])
32 X_Test['sentiment_essay_TS_P'] = essays_sentiment_TS_P
33 X_Test['sentiment_essay_TS_N'] = essays_sentiment_TS_N
34 X_Test['sentiment_essay_TS_NE'] = essays_sentiment_TS_NE
35 X_Test['sentiment_essay_TS_C'] = essays_sentiment_TS_C
36
37 sentiment_norm_P = Normalizer(norm='l2', copy=False)
38 sentiment_norm_N = Normalizer(norm='l2', copy=False)
39 sentiment_norm_NE = Normalizer(norm='l2', copy=False)
40 sentiment_norm_C = Normalizer(norm='l2', copy=False)
41
42
43 sentiment_norm_P.fit(X_Train['sentiment_essay_TR_P'].values.reshape(1,-1))
44 sentiment_norm_N.fit(X_Train['sentiment_essay_TR_N'].values.reshape(1,-1))
45 sentiment_norm_NE.fit(X_Train['sentiment_essay_TR_NE'].values.reshape(1,-1))
46 sentiment_norm_C.fit(X_Train['sentiment_essay_TR_C'].values.reshape(1,-1))
47
48 sentiment_Train_P = sentiment_norm_P.transform(X_Train['sentiment_essay_TR_P'].values.reshape(1,-1))
49 sentiment_Test_P = sentiment_norm_P.transform(X_Test['sentiment_essay_TS_P'].values.reshape(1,-1))
50 sentiment_Train_N = sentiment_norm_N.transform(X_Train['sentiment_essay_TR_N'].values.reshape(1,-1))
51 sentiment_Test_N = sentiment_norm_N.transform(X_Test['sentiment_essay_TS_N'].values.reshape(1,-1))
52 sentiment_Train_NE = sentiment_norm_NE.transform(X_Train['sentiment_essay_TR_NE'].values.reshape(1,-1))
53 sentiment_Test_NE = sentiment_norm_NE.transform(X_Test['sentiment_essay_TS_NE'].values.reshape(1,-1))
54 sentiment_Train_C = sentiment_norm_C.transform(X_Train['sentiment_essay_TR_C'].values.reshape(1,-1))
55 sentiment_Test_C = sentiment_norm_C.transform(X_Test['sentiment_essay_TS_C'].values.reshape(1,-1))
56
57 sentiment_Train_P = (X_Train['sentiment_essay_TR_P'].values.reshape(-1,1))
58 sentiment_Test_P = (X_Test['sentiment_essay_TS_P'].values.reshape(-1,1))
59 sentiment_Train_N = (X_Train['sentiment_essay_TR_N'].values.reshape(-1,1))
60 sentiment_Test_N = (X_Test['sentiment_essay_TS_N'].values.reshape(-1,1))
61 sentiment_Train_NE = (X_Train['sentiment_essay_TR_NE'].values.reshape(-1,1))
62 sentiment_Test_NE = (X_Test['sentiment_essay_TS_NE'].values.reshape(-1,1))
63 sentiment_Train_C = (X_Train['sentiment_essay_TR_C'].values.reshape(-1,1))
64 sentiment_Test_C = (X_Test['sentiment_essay_TS_C'].values.reshape(-1,1))
65
66
67 print("Shape of sentiment Train matrix after one hot encoding ",sentiment_Train_P.shape)
68 print("Shape of sentiment Test matrix after one hot encoding ",sentiment_Test_P.shape)
69 print("Shape of sentiment Train matrix after one hot encoding ",sentiment_Train_N.shape)
70 print("Shape of sentiment Test matrix after one hot encoding ",sentiment_Test_N.shape)
71 print("Shape of sentiment Train matrix after one hot encoding ",sentiment_Train_NE.shape)
72 print("Shape of sentiment Test matrix after one hot encoding ",sentiment_Test_NE.shape)
73 print("Shape of sentiment Train matrix after one hot encoding ",sentiment_Train_C.shape)
74 print("Shape of sentiment Test matrix after one hot encoding ",sentiment_Test_C.shape)

```

```
100% | 73196/73196 [03:08<00:00, 387.92it/s]
100% | 36052/36052 [01:33<00:00, 385.50it/s]Shape of sentiment Train matrix after one hot encodig (73196, 1)
Shape of sentiment Test matrix after one hot encodig (36052, 1)
Shape of sentiment Train matrix after one hot encodig (73196, 1)
Shape of sentiment Test matrix after one hot encodig (36052, 1)
Shape of sentiment Train matrix after one hot encodig (73196, 1)
Shape of sentiment Test matrix after one hot encodig (36052, 1)
Shape of sentiment Train matrix after one hot encodig (73196, 1)
Shape of sentiment Test matrix after one hot encodig (36052, 1)
```

## ▼ Assignment 8: DT

### 1. Apply Decision Tree Classifier(`DecisionTreeClassifier`) on these feature sets

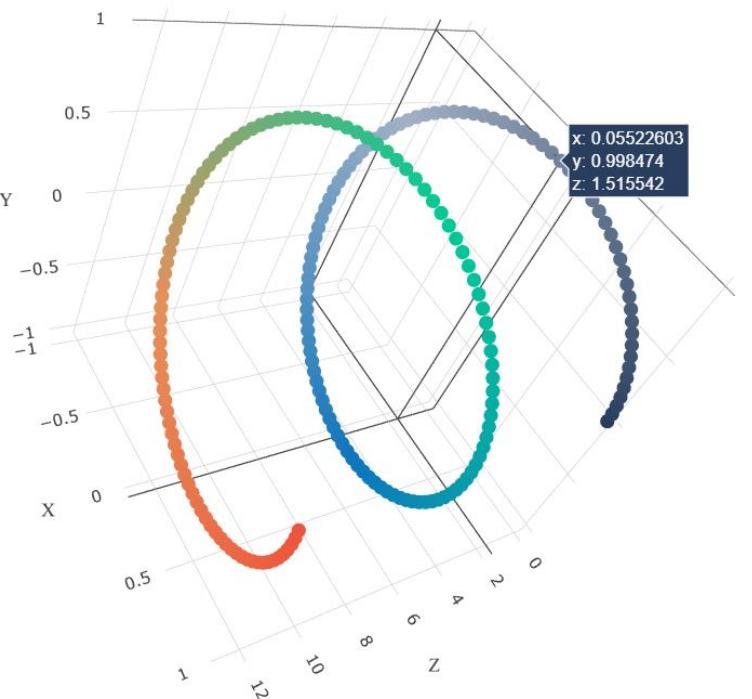
- **Set 1:** categorical, numerical features + `project_title(TFIDF)`+ `preprocessed_eassay (TFIDF)`
- **Set 2:** categorical, numerical features + `project_title(TFIDF W2V)`+ `preprocessed_eassay (TFIDF W2V)`

### 2. The hyper parameter tuning (**best `depth` in range [1, 5, 10, 50]**, and the best `min\_samples\_split` in range [5, 10, 100, 500])

- Find the best hyper parameter which will give the maximum [AUC](#) value
- find the best hyper parameter using k-fold cross validation(use gridsearch cv or randomsearch cv)/simple cross validation data(you can write your own for loops refer sample solution)

### 3. Representation of results

- You need to plot the performance of model both on train data and cross validation data for each hyper parameter, like shown in the



figure

with X-axis as `min_sample_split`, Y-axis as `max_depth`,

and Z-axis as **AUC Score** , we have given the notebook which explains how to plot this 3d plot, you can find it in the same drive

`3d_scatter_plot.ipynb`

or

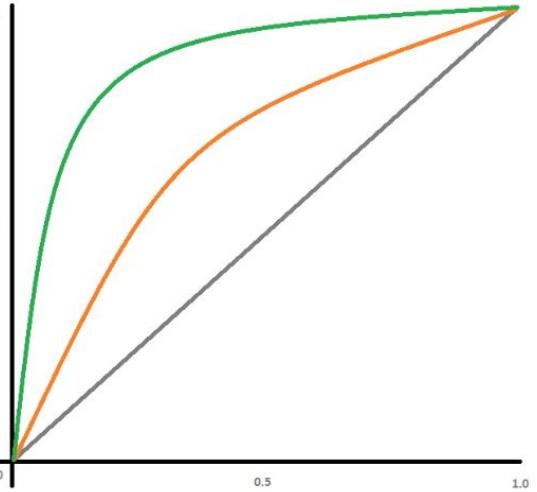
- You need to plot the performance of model both on train data and cross validation data for each hyper parameter, like shown in the



figure [seaborn heat maps](#) with rows as **n\_estimators**, columns as **max\_depth**, and values

inside the cell representing **AUC Score**

- You choose either of the plotting techniques out of 3d plot or heat map
- Once after you found the best hyper parameter, you need to train your model with it, and find the AUC on test data and plot the ROC



curve on both train and test.

- Along with plotting ROC curve, you need to print the [confusion matrix](#) with predicted and original labels of test data points

	Predicted: NO	Predicted: YES
Actual: NO	TN = ??	FP = ??
Actual: YES	FN = ??	TP = ??

- Once after you plot the confusion matrix with the test data, get all the `false positive data points`

- Plot the WordCloud(<https://www.geeksforgeeks.org/generating-word-cloud-python/>) with the words of essay text of these `false positive data points`
- Plot the box plot with the `price` of these `false positive data points`
- Plot the pdf with the `teacher\_number\_of\_previously\_posted\_projects` of these `false positive data points`

**4. Task 2:** For this task consider set-1 features. Select all the features which are having non-zero feature importance. You can get the feature importance using `feature\_importances\_` (<https://scikit-learn.org/stable/modules/generated/sklearn.tree.DecisionTreeClassifier.html>), discard the all other remaining features and then apply any of the model of you choice i.e. (Decision tree, Logistic Regression, Linear SVM), you need to do hyperparameter tuning corresponding to the model you selected and procedure in step 2 and step 3

Note: when you want to find the feature importance make sure you don't use max\_depth parameter keep it None.

5. You need to summarize the results at the end of the notebook, summarize it in the table format

Vectorizer	Model	Hyper parameter	AUC
BOW	Brute	7	0.78
TFIDF	Brute	12	0.79
W2V	Brute	10	0.78
TFIDFW2V	Brute	6	0.78

#### ▼ 1.5.4 Merging all the above features

- we need to merge all the numerical vectors i.e categorical, text, numerical vectors

```
1 # merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
2 TFIDF_Train = hstack((categories_one_hot_Train,sub_categories_one_hot_Train,school_one_hot_Train,grade_one_hot_Train,prefix_one_hot_Train,text_tfidf_Train,t
3 print(TFIDF_Train.shape)
4 TFIDF_W2V_Train = hstack((categories_one_hot_Train,sub_categories_one_hot_Train,school_one_hot_Train,grade_one_hot_Train,prefix_one_hot_Train,tfidf_w2v_vect
5 print(TFIDF_W2V_Train.shape)

⇒ (73196, 16884)
(73196, 709)
```

```
1 TFIDF_Test = hstack((categories_one_hot_Test,sub_categories_one_hot_Test,school_one_hot_Test,grade_one_hot_Test,prefix_one_hot_Test,text_tfidf_Test,title_tf
2 print(TFIDF_Test.shape)
3 TFIDF_W2V_Test = hstack((categories_one_hot_Test,sub_categories_one_hot_Test,school_one_hot_Test,grade_one_hot_Test,prefix_one_hot_Test,tfidf_w2v_vectors_Te
4 print(TFIDF_W2V_Test.shape)

⇒ (36052, 16884)
(36052, 709)
```

#### ▼ Applying Decision Tree on TFIDF, SET 1

```
1 %%time
2 Dec_Tree = DecisionTreeClassifier(class_weight = 'balanced')
3 parameters = {'max_depth': [1, 5, 10, 50], 'min_samples_split': [5, 10, 100, 500]}
4 Dec_Tree_C = RandomizedSearchCV(Dec_Tree, parameters, cv=3, scoring='roc_auc', return_train_score=True, n_iter=30)
5 Dec_Tree_C.fit(TFIDF_Train, Y_Train)
6 print(Dec_Tree_C.best_estimator_)
7 #print(Dec_Tree_C.cv_results_)

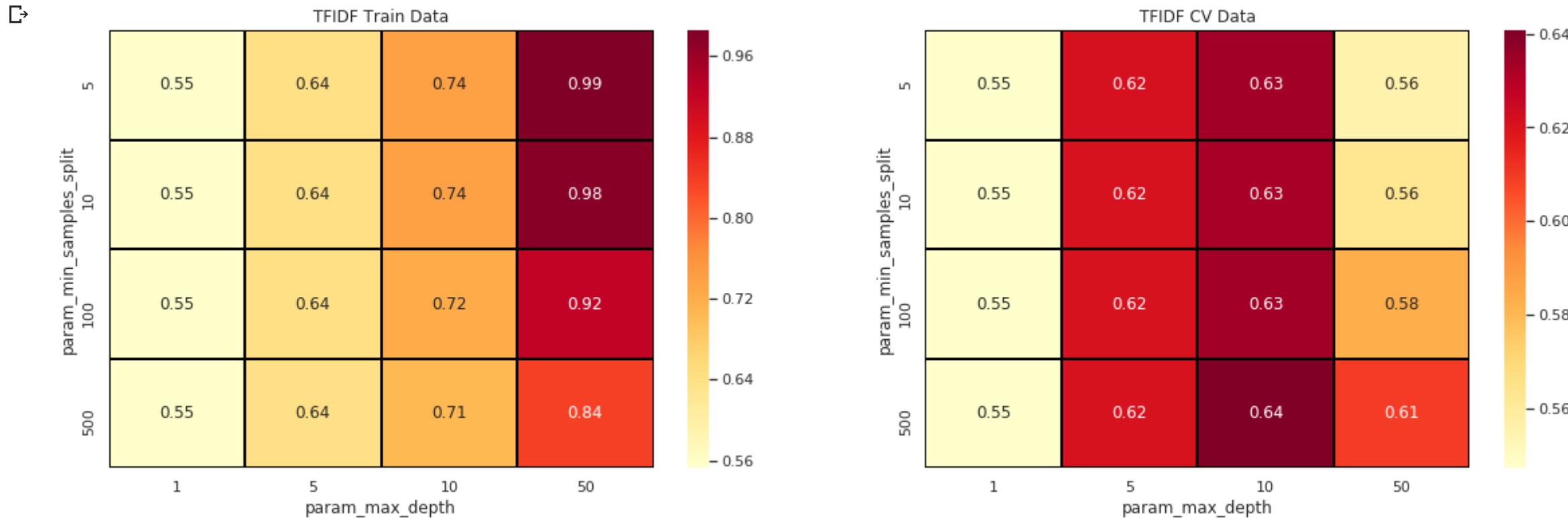
⇒ DecisionTreeClassifier(class_weight='balanced', criterion='gini', max_depth=10,
                           max_features=None, max_leaf_nodes=None,
                           min_impurity_decrease=0.0, min_impurity_split=None,
                           min_samples_leaf=1, min_samples_split=500,
                           min_weight_fraction_leaf=0.0, presort=False,
                           random_state=None, splitter='best')
CPU times: user 14min 51s, sys: 30.3 ms, total: 14min 52s
Wall time: 14min 52s
```

#### ▼ TFIDF Heatmap

```

1 sns.set()
2 Dec_Tree_CV = pd.DataFrame(Dec_Tree_C.cv_results_).groupby(['param_min_samples_split', 'param_max_depth']).max().unstack()[['mean_test_score','mean_train_sc
3 fig, ax = plt.subplots(1,2, figsize=(20,6))
4 sns.heatmap(Dec_Tree_CV.mean_train_score, annot = True,cbar=True, fmt='%.2g', cmap='YlOrRd', linewidths=1, linecolor='black', ax=ax[0])
5 sns.heatmap(Dec_Tree_CV.mean_test_score, annot = True, fmt='%.2g', cmap='YlOrRd', linewidths=1, linecolor='black', ax=ax[1])
6 ax[0].set_title('TFIDF Train Data')
7 ax[1].set_title('TFIDF CV Data')
8 plt.show()

```



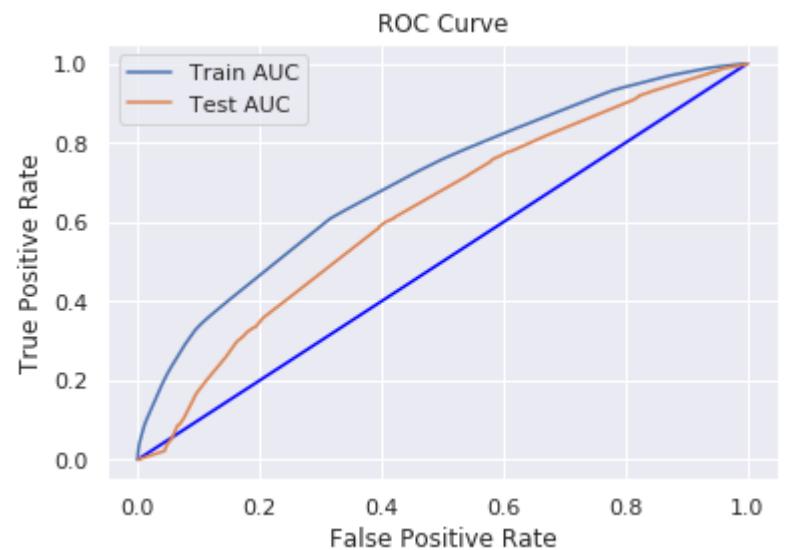
```

1 #Fitting Model to Hyper-Parameter Curve
2 # https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc\_curve.html#sklearn.metrics.roc\_curve
3
4 TFIDF_MD=(Dec_Tree_C.best_params_['max_depth'])
5 TFIDF_MSS=(Dec_Tree_C.best_params_['min_samples_split'])
6 Dec_Tree_C_OP=DecisionTreeClassifier(class_weight = 'balanced',max_depth=TFIDF_MD,min_samples_split=TFIDF_MSS)
7 Dec_Tree_C_OP.fit(TFIDF_Train, Y_Train)
8
9 Y_Train_PR = Dec_Tree_C_OP.predict_proba(TFIDF_Train) [:,1]
10 Y_Test_PR = Dec_Tree_C_OP.predict_proba(TFIDF_Test) [:,1]
11
12 fpr_Train, tpr_Train, thresholds_Train = roc_curve(Y_Train, Y_Train_PR)
13 fpr_Test, tpr_Test, thresholds_Test = roc_curve(Y_Test, Y_Test_PR)
14
15 TFIDF_AUC = auc(fpr_Test, tpr_Test)
16 plt.plot([0,1],[0,1], 'k-', color='blue')
17 plt.plot(fpr_Train, tpr_Train, label="Train AUC")
18 plt.plot(fpr_Test, tpr_Test, label="Test AUC")
19 plt.legend()
20 plt.xlabel("False Positive Rate")
21 plt.ylabel("True Positive Rate")
22 plt.title("ROC Curve")
23 plt.show()
24 print("-"*120)
25 print("Optimal value of AUC Train =",auc(fpr_Train, tpr_Train))
26 print("Optimal value of AUC Test =",auc(fpr_Test, tpr_Test))
27 pred1 = Dec_Tree_C_OP.predict(TFIDF_Train)
28 pred2 = Dec_Tree_C_OP.predict(TFIDF_Test)

```

```
28 pred2 = Dec_Tree_C_0P.predict(TFIDF_test)
```

➡

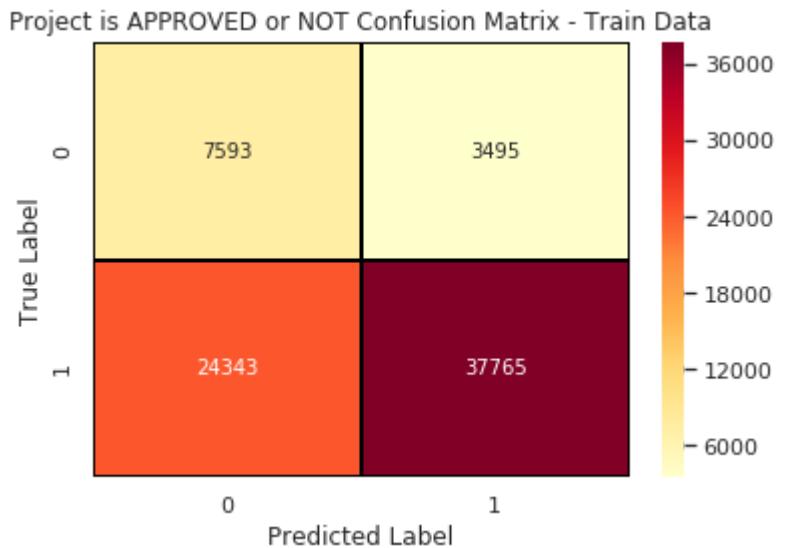


Optimal value of AUC Train = 0.7007161862927224  
Optimal value of AUC Test = 0.6218123067835468

## ▼ Confusion Matrix

```
1 https://seaborn.pydata.org/generated/seaborn.heatmap.html
2 https://getaravind.com/blog/confusion-matrix-seaborn-heatmap/
3 https://stackoverflow.com/questions/19233771/sklearn-plot-confusion-matrix-with-labels
4 %matplotlib inline
5 Train = confusion_matrix(Y_Train, pred1)
6 sns.heatmap(Train, annot=True, cbar=True, fmt='d', cmap='YlOrRd', linewidths=1, linecolor='black')
7 plt.ylabel('True Label')
8 plt.xlabel('Predicted Label')
9 plt.title('Project is APPROVED or NOT Confusion Matrix - Train Data')
```

➡ Text(0.5, 1, 'Project is APPROVED or NOT Confusion Matrix - Train Data')



## OBSERVATION:

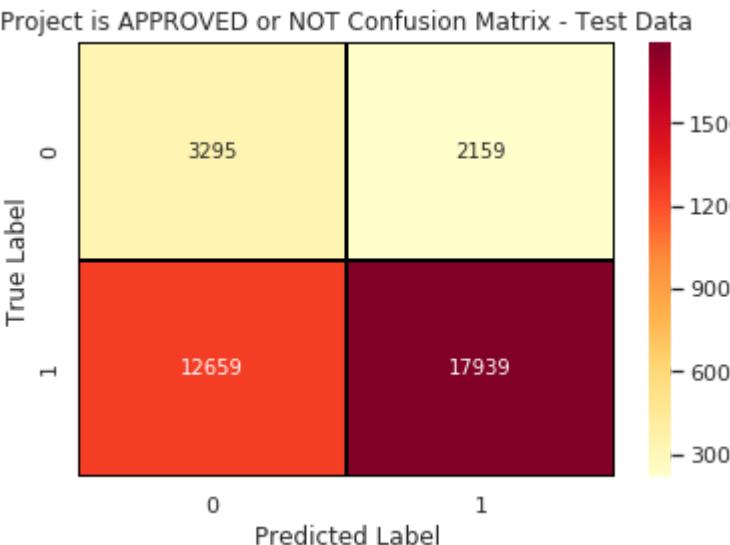
True Negative = 7593; False Negative = 24343; True Positive = 37765; False Positive = 3495  
Accuracy (Overall, how often is the classifier correct) = 0.62  
Precision(When it predicts yes, how often is it correct) = 0.92  
Misclassification (Overall, how often is it wrong) = 0.39

```

1 #https://seaborn.pydata.org/generated/seaborn.heatmap.html
2 #https://getaravind.com/blog/confusion-matrix-seaborn-heatmap/
3 #https://stackoverflow.com/questions/19233771/sklearn-plot-confusion-matrix-with-labels
4 Test = confusion_matrix(Y_Test, pred2)
5 sns.heatmap(Test, annot=True, cbar=True, fmt='d', cmap='YlOrRd', linewidths=1, linecolor='black')
6 plt.ylabel('True Label')
7 plt.xlabel('Predicted Label')
8 plt.title('Project is APPROVED or NOT Confusion Matrix - Test Data')

```

↳ Text(0.5, 1, 'Project is APPROVED or NOT Confusion Matrix - Test Data')



### OBSERVATION:

True Negative = 3295; False Negative = 12659; True Positive = 17939; False Positive = 2159

Accuracy (Overall, how often is the classifier correct) = 0.59

Precision(When it predicts yes, how often is it correct) = 0.90

Misclassification (Overall, how often is it wrong) = 0.42

### WordCloud

```

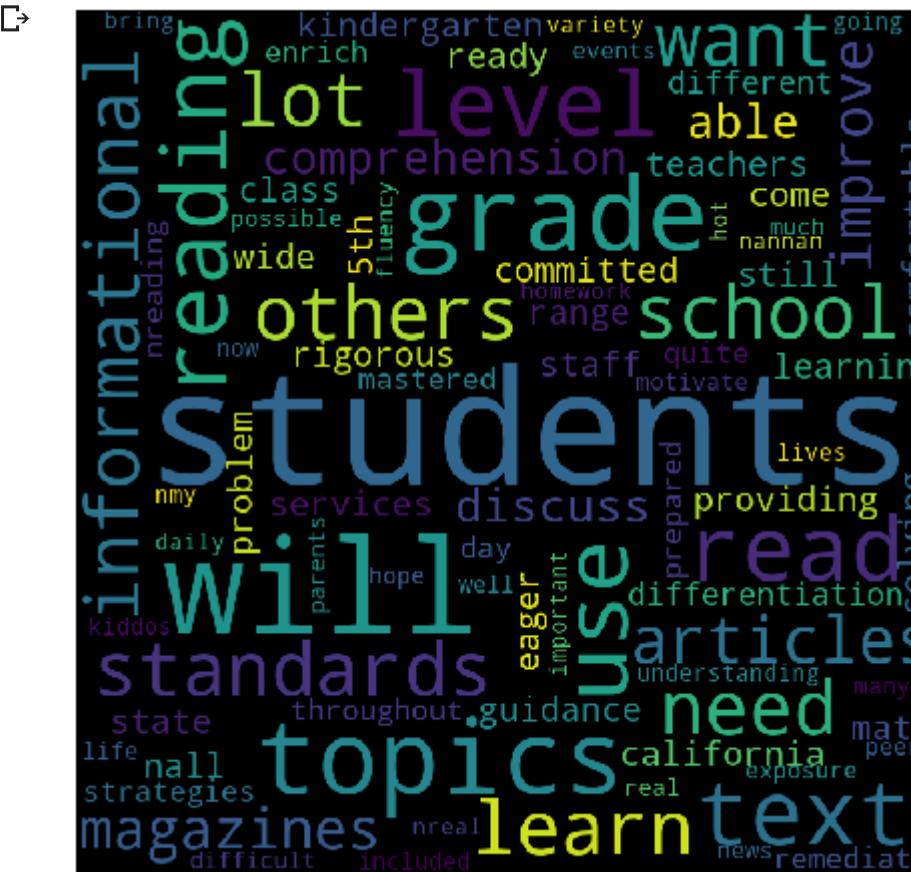
1 #https://www.geeksforgeeks.org/generating-word-cloud-python/
2 fpi = []
3 for i in range(len(Y_Test)) :
4     if (Y_Test[i] == 0) & (pred2[i] == 1) :
5         fpi.append(i)
6
7 X_Test_ES = pd.DataFrame(X_Test['essay'])
8
9 fp_essay1 = []
10 for i in fpi :
11     #pdb.set_trace()
12     fp_essay1.append(X_Test_ES['essay'][i])
13 # Word cloud of essay
14 from wordcloud import WordCloud, STOPWORDS
15 comment_words = ' '
16 stopwords = set(STOPWORDS)
17 for val in fp_essay1 :
18     val = str(val)
19     tokens = val.split()
20 for i in range(len(tokens)):
21     tokens[i] = tokens[i].lower()
22

```

```

22 for words in tokens :
23     comment_words = comment_words + words + ' '
24
25 wordcloud = WordCloud(width = 800, height = 800, background_color ='black', stopwords = stopwords,
26 min_font_size = 10).generate(comment_words)
27
28 plt.figure(figsize = (6, 6), facecolor = None)
29 plt.imshow(wordcloud)
30 plt.axis("off")
31 plt.tight_layout(pad = 0)
32 plt.show()

```

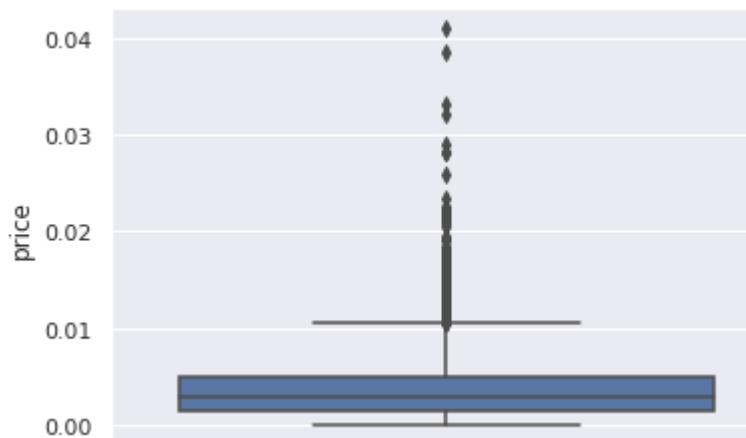


```

1 cols = X_Test.columns
2 X_Test_fp = pd.DataFrame(columns=cols)
3 for i in fpi :
4     X_Test_fp = X_Test_fp.append(X_Test.filter(items=[i], axis=0))
5 sns.boxplot(y='price', data=X_Test_fp)

```

```
↳ <matplotlib.axes._subplots.AxesSubplot at 0x7f4679451748>
```

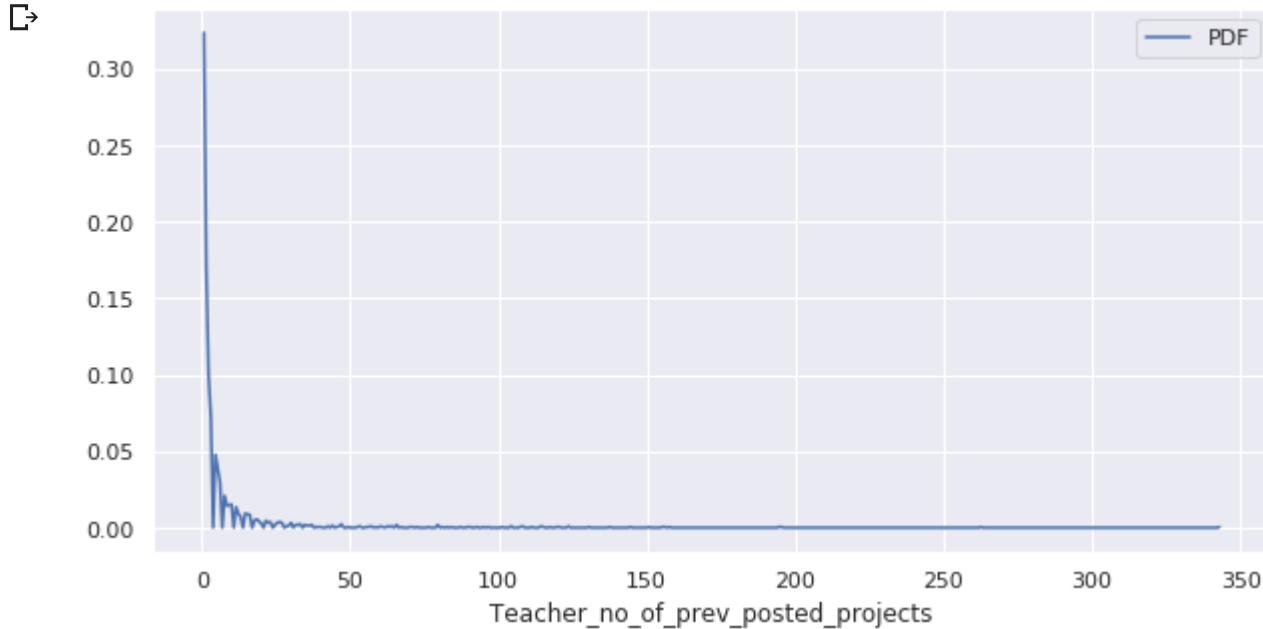


```

1 #PDF (FP ,teacher_number_of_previously_posted_projects)
2 plt.figure(figsize=(10,5))

```

```
3 counts, bin_edges = np.histogram(X_Test_fp[ 'teacher_number_of_previously_posted_projects' ],bins='auto', density=True)
4 pdf = counts/sum(counts)
5 pdfP, = plt.plot(bin_edges[1:], pdf)
6 plt.legend([pdfP], [ "PDF"])
7 plt.xlabel('Teacher_no_of_prev_posted_projects')
8 plt.show()
```



```
1 #Feature aggregation
2 feat_1=vectorizer_cat.get_feature_names()
3 feat_2=vectorizer_sub_cat.get_feature_names()
4 feat_3=vectorizer_school.get_feature_names()
5 feat_4=vectorizer_prefix.get_feature_names()
6 feat_5=vectorizer_grade.get_feature_names()
7 feat_ES=vectorizer_essays_tfidf.get_feature_names()
8 feat_TT=vectorizer_titles_tfidf.get_feature_names()
9
10 feat_net_tfidf = feat_1 + feat_2 + feat_3 + feat_4 + feat_5 + feat_ES + feat_TT
11 # p is price, q is quantity, t is teacher previous year projects
12 feat_net_tfidf.append('price')
13 feat_net_tfidf.append('quantity')
14 feat_net_tfidf.append('prev_proposed_projects')
15 feat_net_tfidf.append("title_word_count")
16 feat_net_tfidf.append("essay_word_count")
17 feat_net_tfidf.append('sentiment_Test_P')
18 feat_net_tfidf.append('sentiment_Test_N')
19 feat_net_tfidf.append('sentiment_Test_NE')
20 feat_net_tfidf.append('sentiment_Test_C')
21
22 print(len(feat net tfidf))
```

```
[4]: 16884

1 # https://medium.com/@rnbrown/creating-and-visualizing-decision-trees-with-python-f8e8fa394176
2 Dec_Tree_C_OP=DecisionTreeClassifier(class_weight = 'balanced',max_depth=TFIDF_MD,min_samples_split=TFIDF_MSS)
3 Dec_Tree_C_OP.fit(TFIDF_Train, Y_Train)
4 dot_data = tree.export_graphviz(Dec_Tree_C_OP, out_file=None, feature_names=feat_net_tfidf)
5 graph = graphviz.Source(dot_data)
6 graph.format = 'png'
7 graph.render("Tfidf decision tree",view = True)
8 dot_data = StringIO()
```

```

9 export_graphviz(Dec_Tree_C_OP, out_file=dot_data, filled=True, rounded=True, special_characters=True, feature_names=feat_net_tfidf, rotate=True)
10 graph = pydotplus.graph_from_dot_data(dot_data.getvalue())
11 Image(graph.create_png())

```

## ▼ Applying Decision Tree on TFIDF\_W2V, SET 2

```

1 %%time
2 Dec_Tree = DecisionTreeClassifier(class_weight = 'balanced')
3 parameters = {'max_depth': [1, 5, 10, 50], 'min_samples_split': [5, 10, 100, 500]}
4 Dec_Tree_C = RandomizedSearchCV(Dec_Tree, parameters, cv=3, scoring='roc_auc', return_train_score=True, n_iter=30)
5 Dec_Tree_C.fit(TFIDF_W2V_Train, Y_Train)
6 print(Dec_Tree_C.best_estimator_)
7 #print(Dec_Tree_C.cv_results_)

[→ DecisionTreeClassifier(class_weight='balanced', criterion='gini', max_depth=5,
                           max_features=None, max_leaf_nodes=None,
                           min_impurity_decrease=0.0, min_impurity_split=None,
                           min_samples_leaf=1, min_samples_split=500,
                           min_weight_fraction_leaf=0.0, presort=False,
                           random_state=None, splitter='best')
CPU times: user 25min 46s, sys: 409 ms, total: 25min 46s
Wall time: 25min 46s

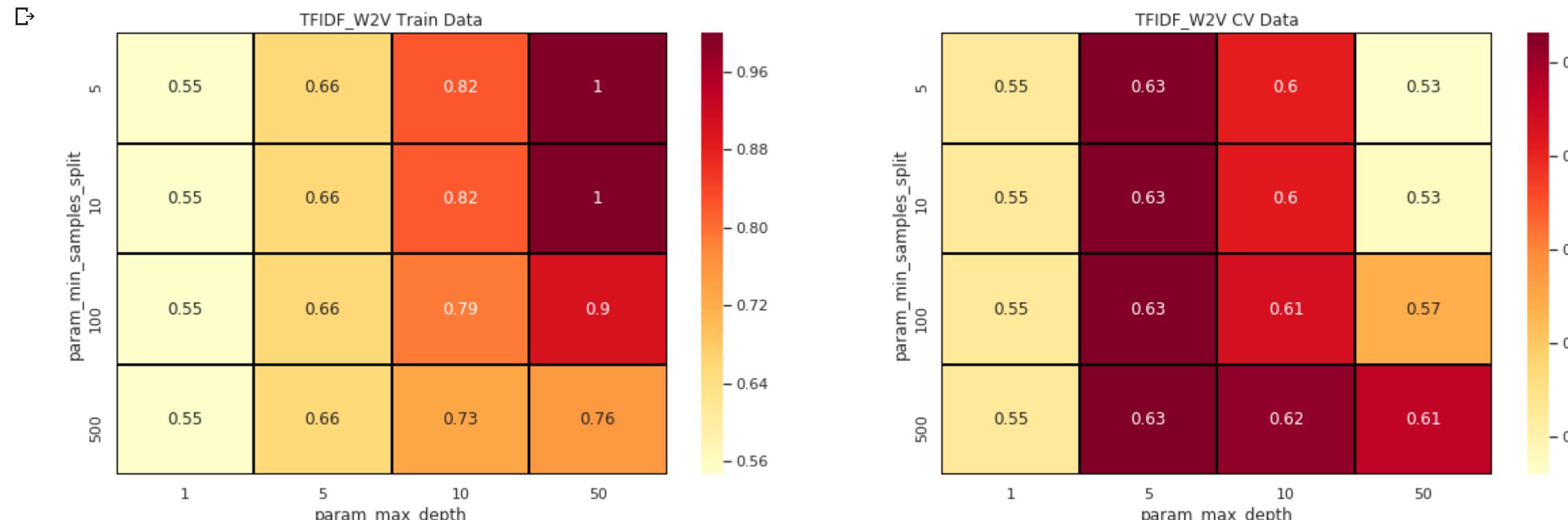
```

## ▼ TFIDF W2V Heatmap

```

1 sns.set()
2 Dec_Tree_CV = pd.DataFrame(Dec_Tree_C.cv_results_).groupby(['param_min_samples_split', 'param_max_depth']).max().unstack()[['mean_test_score', 'mean_train_score']]
3 fig, ax = plt.subplots(1, 2, figsize=(20, 6))
4 sns.heatmap(Dec_Tree_CV.mean_train_score, annot = True, cbar=True, fmt='.2g', cmap='YlOrRd', linewidths=1, linecolor='black', ax=ax[0])
5 sns.heatmap(Dec_Tree_CV.mean_test_score, annot = True, fmt='.2g', cmap='YlOrRd', linewidths=1, linecolor='black', ax=ax[1])
6 ax[0].set_title('TFIDF_W2V Train Data')
7 ax[1].set_title('TFIDF_W2V CV Data')
8 plt.show()

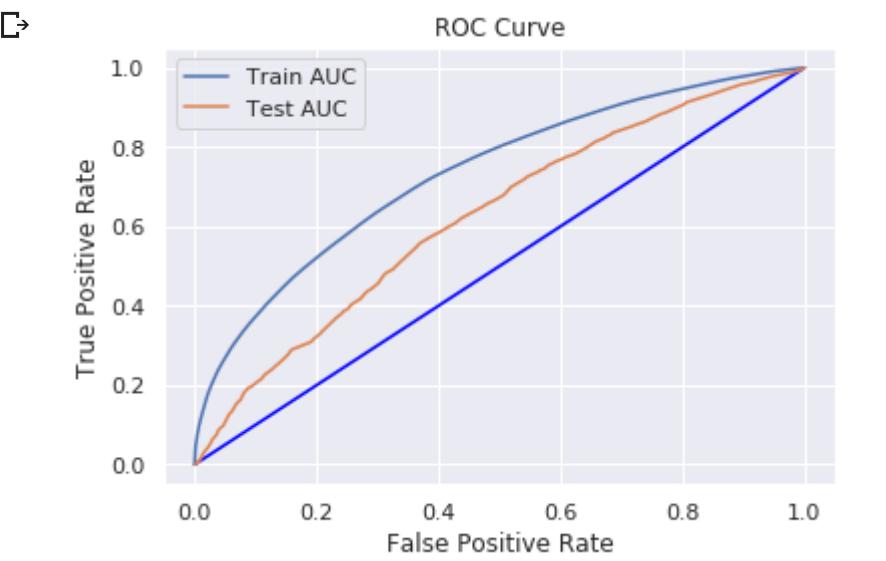
```



```

1 #fitting model to hyper-parameter curve
2 # https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc_curve.html#sklearn.metrics.roc_curve
3
4 TFIDF_W2V_MD=(Dec_Tree_C.best_params_['max_depth'])
5 TFIDF_W2V_MSS=(Dec_Tree_C.best_params_['min_samples_split'])
6 Dec_Tree_C_OP=DecisionTreeClassifier(class_weight = 'balanced',max_depth=TFIDF_MD,min_samples_split=TFIDF_MSS)
7 Dec_Tree_C_OP.fit(TFIDF_W2V_Train, Y_Train)
8
9 Y_Train_PR = Dec_Tree_C_OP.predict_proba(TFIDF_W2V_Train) [:,1]
10 Y_Test_PR = Dec_Tree_C_OP.predict_proba(TFIDF_W2V_Test) [:,1]
11
12 fpr_Train, tpr_Train, thresholds_Train = roc_curve(Y_Train, Y_Train_PR)
13 fpr_Test, tpr_Test, thresholds_Test = roc_curve(Y_Test, Y_Test_PR)
14
15 TFIDF_W2V_AUC = auc(fpr_Test, tpr_Test)
16 plt.plot([0,1],[0,1],'k-',color='blue')
17 plt.plot(fpr_Train, tpr_Train, label="Train AUC")
18 plt.plot(fpr_Test, tpr_Test, label="Test AUC")
19 plt.legend()
20 plt.xlabel("False Positive Rate")
21 plt.ylabel("True Positive Rate")
22 plt.title("ROC Curve")
23 plt.show()
24 print("-"*120)
25 print("Optimal value of AUC Train =",auc(fpr_Train, tpr_Train))
26 print("Optimal value of AUC Test =",auc(fpr_Test, tpr_Test))
27 pred3 = Dec_Tree_C_OP.predict(TFIDF_W2V_Train)
28 pred4 = Dec_Tree_C_OP.predict(TFIDF_W2V_Test)

```



Optimal value of AUC Train = 0.7335238586109045  
Optimal value of AUC Test = 0.6240460475988554

## ▼ Confusion Matrix

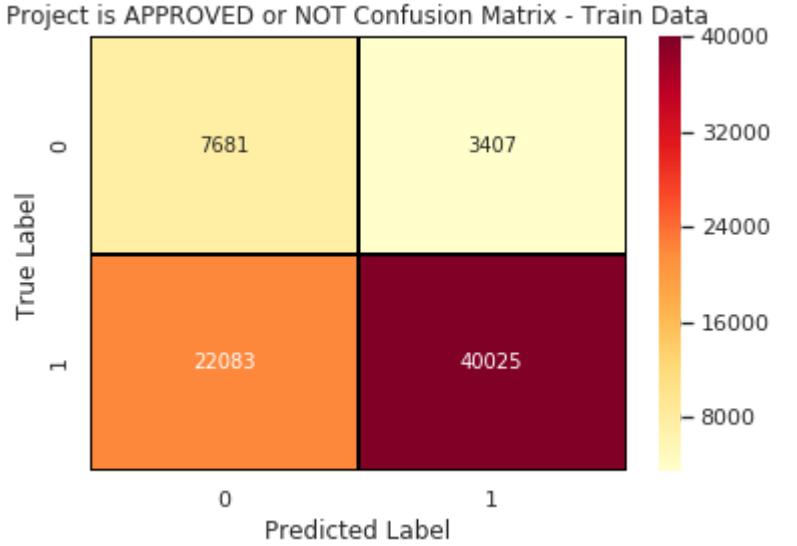
```

1 #https://seaborn.pydata.org/generated/seaborn.heatmap.html
2 #https://getaravind.com/blog/confusion-matrix-seaborn-heatmap/
3 #https://stackoverflow.com/questions/19233771/sklearn-plot-confusion-matrix-with-labels
4 %matplotlib inline
5 Train = confusion_matrix(Y_Train, pred3)
6 sns.heatmap(Train, annot=True, cbar=True, fmt='d', cmap='YlOrRd', linewidths=1, linecolor='black')
7 plt.ylabel('True Label')
8 plt.xlabel('Predicted Label')

```

```
9 plt.title('Project is APPROVED or NOT Confusion Matrix - Train Data')
```

```
↳ Text(0.5, 1, 'Project is APPROVED or NOT Confusion Matrix - Train Data')
```



#### OBSERVATION:

True Negative = 7681; False Negative = 22083; True Positive = 40025; False Positive = 3407

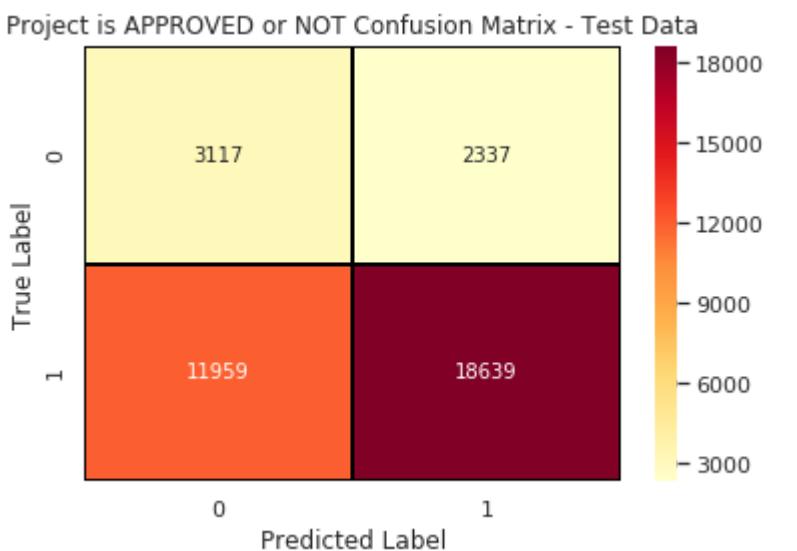
Accuracy (Overall, how often is the classifier correct) = 0.66

Precision(When it predicts yes, how often is it correct) = 0.93

Misclassification (Overall, how often is it wrong) = 0.35

```
1 #https://seaborn.pydata.org/generated/seaborn.heatmap.html
2 #https://getaravind.com/blog/confusion-matrix-seaborn-heatmap/
3 #https://stackoverflow.com/questions/19233771/sklearn-plot-confusion-matrix-with-labels
4 Test = confusion_matrix(Y_Test, pred4)
5 sns.heatmap(Test, annot=True, cbar=True, fmt='d', cmap='YlOrRd', linewidths=1, linecolor='black')
6 plt.ylabel('True Label')
7 plt.xlabel('Predicted Label')
8 plt.title('Project is APPROVED or NOT Confusion Matrix - Test Data')
```

```
↳ Text(0.5, 1, 'Project is APPROVED or NOT Confusion Matrix - Test Data')
```



#### OBSERVATION:

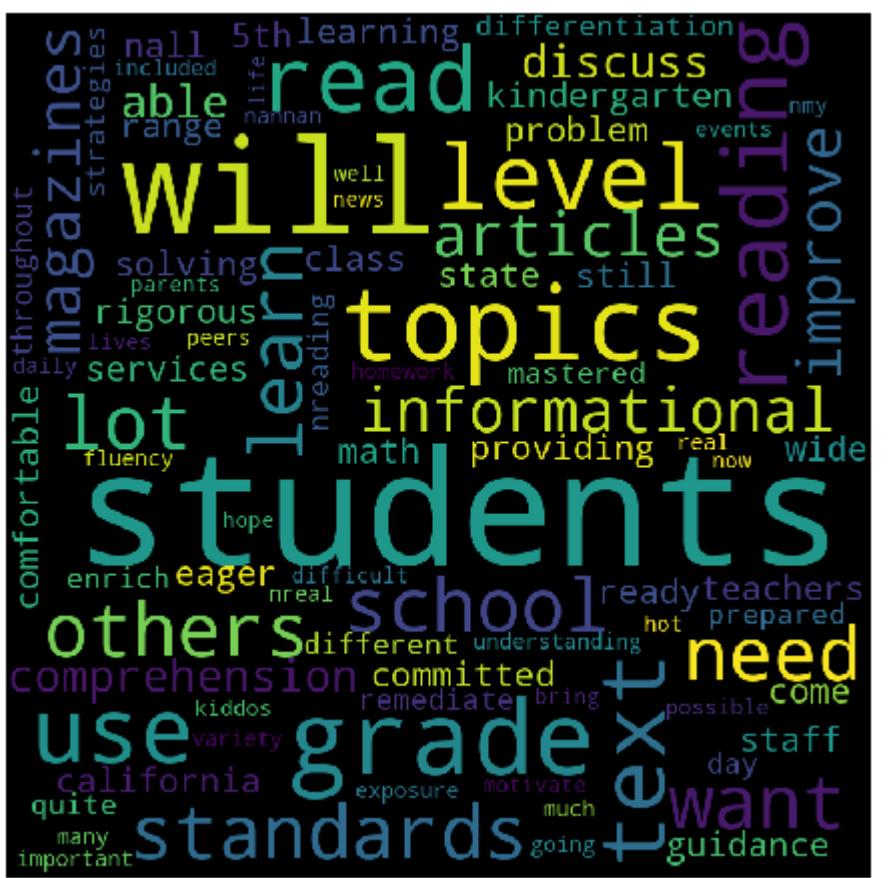
True Negative = 3117; False Negative = 11959; True Positive = 18639; False Positive = 2337

Accuracy (Overall, how often is the classifier correct) = 0.61

## ▼ WordCloud

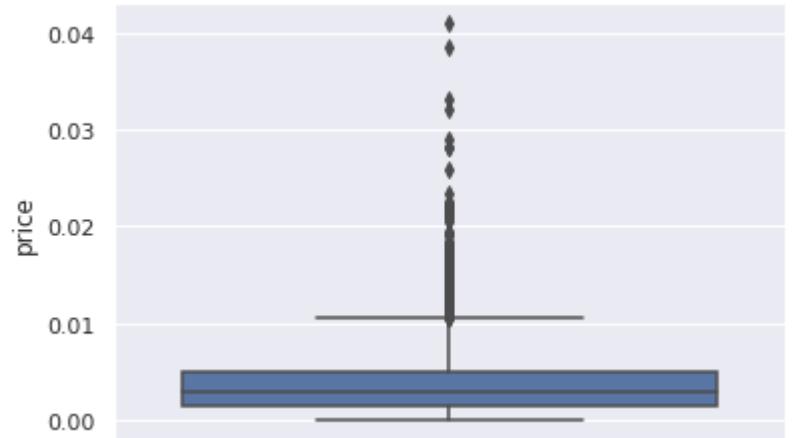
```
1 #https://www.geeksforgeeks.org/generating-word-cloud-python/
2 fpi = []
3 for i in range(len(Y_Test)) :
4     if (Y_Test[i] == 0) & (pred2[i] == 1) :
5         fpi.append(i)
6
7 X_Test_ES = pd.DataFrame(X_Test['essay'])
8
9 fp_essay1 = []
10 for i in fpi :
11     #pdb.set_trace()
12     fp_essay1.append(X_Test_ES['essay'][i])
13 # Word cloud of essay
14 from wordcloud import WordCloud, STOPWORDS
15 comment_words = ' '
16 stopwords = set(STOPWORDS)
17 for val in fp_essay1 :
18     val = str(val)
19     tokens = val.split()
20 for i in range(len(tokens)):
21     tokens[i] = tokens[i].lower()
22 for words in tokens :
23     comment_words = comment_words + words + ' '
24
25 wordcloud = WordCloud(width = 800, height = 800, background_color ='black', stopwords = stopwords,
26 min_font_size = 10).generate(comment_words)
27
28 plt.figure(figsize = (6, 6), facecolor = None)
29 plt.imshow(wordcloud)
30 plt.axis("off")
31 plt.tight_layout(pad = 0)
32 plt.show()
```





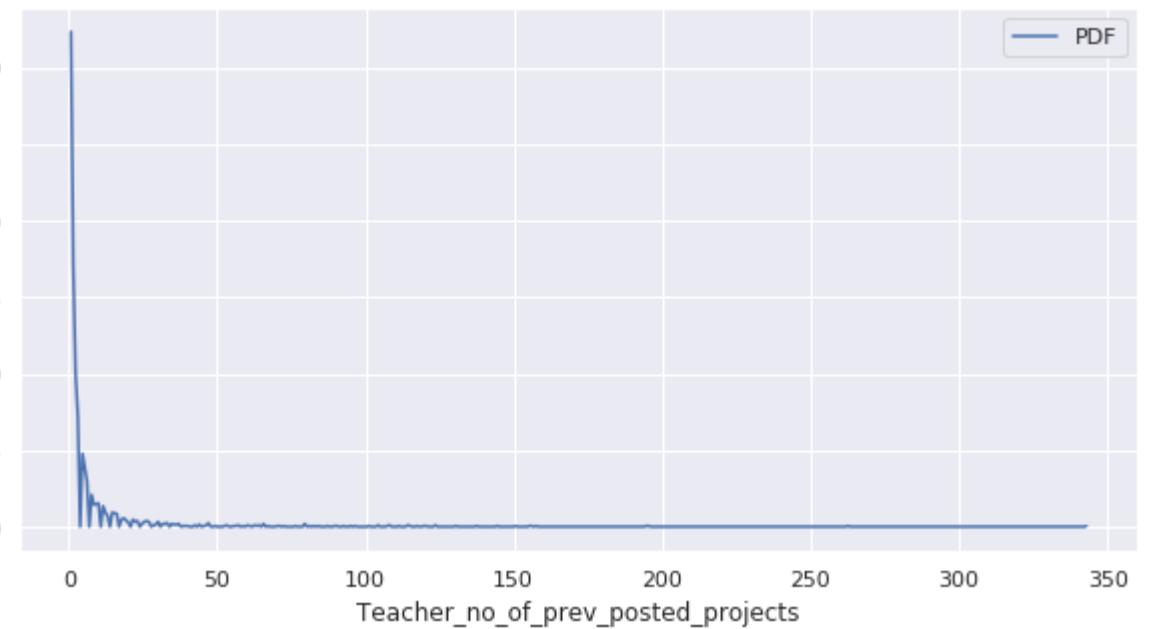
```
1 cols = X_Test.columns
2 X_Test_fp = pd.DataFrame(columns=cols)
3 for i in fpi :
4     X_Test_fp = X_Test_fp.append(X_Test.filter(items=[i], axis=0))
5 sns.boxplot(y='price', data=X_Test_fp)
```

→ <matplotlib.axes.\_subplots.AxesSubplot at 0x7f46c816ccf8>



```
1 #PDF (FP ,teacher_number_of_previously_posted_projects)
2 plt.figure(figsize=(10,5))
3 counts, bin_edges = np.histogram(X_Test_fp[ 'teacher_number_of_previously_posted_projects' ],bins='auto', density=True)
4 pdf = counts/sum(counts)
5 pdfP, = plt.plot(bin_edges[1:], pdf)
6 plt.legend([pdfP], ["PDF"])
7 plt.xlabel('Teacher_no_of_prev_posted_projects')
8 plt.show()
```

→



#### ▼ [Task-2] Applying Decision Tree on Non-zero features importance of Set-1

```

1 #https://www.geeksforgeeks.org/numpy-nonzero-in-python/
2 # Coverting the Non-zero important features into a new dataset
3
4 Dec_Tree_C1=DecisionTreeClassifier(class_weight = 'balanced',max_depth=None,min_samples_split=TFIDF_MSS)
5 Dec_Tree_C1.fit(TFIDF_Train, Y_Train)
6 Fe_Im=Dec_Tree_C1.feature_importances_
7 NZFI = np.nonzero(Fe_Im)
8 df = pd.DataFrame(TFIDF_Train.toarray())
9
10 ls=[]
11 NZFI_TR= pd.DataFrame()
12 NZFI1_TR=list(NZFI)
13 for i in NZFI1_TR:
14     ls.append(i)
15 for j in ls:
16     for k in j:
17         m=int(k)
18         fn=feat_net_tfidf[m]
19         t_ls=list(df[m])
20         NZFI_TR[fn] = t_ls
21
22 df_ts = pd.DataFrame(TFIDF_Test.toarray())
23 ls_ts=[]
24 NZFI_TS= pd.DataFrame()
25
26 NZFI1_TS=list(NZFI1_TR)
27 for i in NZFI1_TS:
28     ls_ts.append(i)
29 for j in ls_ts:
30     for k in j:
31         m=int(k)
32         fn=feat_net_tfidf[m]
33         t_ls_ts=list(df_ts[m])
34         NZFI_TS[fn] = t_ls_ts
35 print("-"*120)
36 print("Shape of the Non Zero features important Train Dataset:",NZFI_TR.shape)

```

```
37 print("Shape of the Non Zero features important Test Dataset:",NZFI_TS.shape)
```

```
38
```

```
↳ -----  
Shape of the Non Zero features important Train Dataset: (73196, 1077)  
Shape of the Non Zero features important Test Dataset: (36052, 1077)
```

```
1 # Non-zero important features of Set-1  
2 for j in ls:  
3     for k in j:  
4         m=int(k)  
5         fn=feat_net_tfidf[m]  
6         print(fn)
```

```
↳
```

Literacy\_Language  
PerformingArts  
wv  
ar  
mn  
ga  
ny  
mr  
grades\_6\_8  
100  
11th  
12  
14  
15  
17  
20  
2006  
2016  
2nd  
35  
3d  
3rd  
4th  
54  
70  
700  
80  
83  
abilities  
able  
abstract  
academics  
acceptance  
accomplish  
account  
accuracy  
accurate  
acquire  
across  
active  
actively  
activities  
activity  
actual  
adapted  
adapters  
add  
added  
adding  
additional  
admired  
adulthood  
adults  
affluent  
age  
agency  
aid  
air  
alaska  
alive  
allow  
allowed  
allowing  
alphabet  
already  
also

also  
always  
amazed  
amazing  
analysis  
another  
answering  
appeal  
appear  
apple  
appreciated  
approaches  
apps  
argue  
arts  
asking  
assess  
assigned  
assignments  
atoms  
attend  
attitude  
attracts  
available  
babysit  
background  
backpacks  
bad  
bags  
balance  
ball  
balls  
band  
bands  
based  
basic  
basketball  
bathroom  
batteries  
beach  
bean  
beating  
beautiful  
became  
become  
becoming  
beethoven  
beg  
behavior  
behavioral  
behind  
belong  
best  
better  
beyond  
big  
binders  
bins  
biology  
black  
blessed  
blessing  
blocks  
bloomington  
board  
boards  
boats

boats  
boiling  
book  
books  
boring  
borrow  
bottles  
bought  
bouncy  
box  
boy  
braille  
brave  
break  
breakout  
breakoutedu  
bridging  
broken  
budget  
builders  
bundle  
business  
camera  
capable  
cardio  
cards  
care  
caring  
carpet  
cart  
catalyst  
cause  
center  
centered  
centers  
certified  
chair  
chairs  
challenges  
chance  
channel  
character  
charts  
chatting  
check  
child  
choice  
chorus  
chromebook  
chromebooks  
circumstances  
class  
classroom  
clean  
cleaning  
clear  
clever  
close  
closer  
clothes  
coach  
coding  
collaborating  
collaborative  
collaboratively  
colony  
color

color  
colored  
colorful  
colors  
come  
comes  
comfort  
communicate  
community  
comparing  
comparison  
compassionate  
competition  
compile  
complete  
component  
computer  
concentration  
concept  
concerts  
conditions  
confidence  
confident  
confused  
connections  
consequently  
considering  
consists  
constantly  
construct  
constructing  
contagious  
content  
continents  
continually  
continue  
continuing  
contribution  
control  
conversations  
core  
could  
counters  
cozy  
create  
creative  
creativity  
criteria  
culturally  
culture  
curious  
current  
currently  
customs  
cut  
cycle  
cycles  
daily  
dance  
day  
deaf  
dealing  
decided  
deficits  
demands  
depending  
describe

describe  
deserves  
deserving  
designs  
desk  
desks  
desperately  
determine  
detroit  
devastating  
develop  
diagnosis  
different  
dig  
direct  
disabilities  
discipline  
discover  
discoveries  
disorder  
display  
district  
diverse  
dollars  
donated  
donation  
door  
drag  
drills  
drinks  
drum  
dvd  
dynamic  
early  
ease  
easel  
easier  
easily  
economic  
economically  
economy  
ecosystem  
ecosystems  
education  
educators  
effective  
eighth  
electricity  
electronically  
electronics  
elements  
eliminating  
emphasize  
encouragement  
ended  
endurance  
energetic  
enforce  
engagement  
engaging  
engineers  
enjoying  
enough  
enriching  
enrichment  
enroll  
enrollment

environment  
ensure  
enthusiasm  
enthusiastic  
entire  
environment  
equipment  
especially  
essentials  
esteem  
ethnic  
event  
events  
every  
everywhere  
excellent  
excited  
existing  
exit  
expand  
expect  
expecting  
experience  
experiences  
experiment  
experiments  
explicit  
exploration  
exploring  
expose  
exposure  
expression  
extension  
eye  
fabulous  
face  
faced  
fail  
failures  
fairy  
families  
family  
fascinating  
feel  
feels  
fellow  
fidget  
fill  
filled  
financial  
fine  
finish  
fire  
firm  
fish  
fitness  
flash  
flood  
floor  
fluently  
folder  
follow  
following  
foods  
foreign  
forever  
fossils

fossils  
foundational  
fourth  
fractions  
freshman  
frog  
frost  
frustration  
fueled  
fun  
functioning  
fund  
funded  
furniture  
gain  
galaxies  
games  
gather  
gears  
gets  
giant  
gifted  
girls  
give  
gives  
giving  
glad  
go  
goal  
going  
good  
graders  
grammar  
graph  
graphic  
great  
greater  
greatest  
grow  
habitats  
half  
hamper  
hand  
handy  
hang  
hard  
hardworking  
harlem  
harry  
harsh  
headphones  
hear  
heard  
hearing  
heart  
help  
helped  
higher  
highlighters  
highlights  
highly  
hispanic  
historical  
history  
home  
homeless  
homes

homes  
homework  
honors  
hooked  
hoops  
hope  
hopefully  
hoping  
housekeeping  
houses  
humanity  
humble  
hundred  
idea  
identification  
identity  
imaginings  
imagine  
implement  
implementation  
implications  
important  
improve  
improvement  
include  
incredible  
independent  
india  
individualized  
industries  
ink  
inquiry  
inspiration  
inspirational  
instant  
instead  
instill  
instrument  
instruments  
intellectual  
interact  
interactive  
interest  
interested  
intervention  
introduce  
introduced  
invention  
investigative  
investment  
involvement  
involves  
ipad  
ipads  
ipod  
items  
jackson  
keep  
kid  
kids  
kinetic  
kit  
kits  
lab  
laminate  
laminating  
land

land  
language  
languages  
laptop  
laptops  
large  
largely  
larger  
later  
leadership  
learn  
learning  
learns  
lego  
legos  
less  
lessons  
let  
letter  
leveled  
levels  
lexile  
licenses  
lifestyle  
lighting  
likes  
limited  
lines  
list  
listen  
listening  
literature  
little  
living  
logic  
long  
longer  
look  
looking  
loose  
love  
lovers  
loving  
low  
lunch  
made  
magna  
magnetic  
magnifying  
majority  
make  
makers  
makes  
man  
manipulative  
many  
markers  
master  
masterpieces  
mat  
material  
materials  
mats  
max  
maze  
means  
measurement

measurable  
mechanism  
meet  
members  
mention  
mentors  
mexico  
middle  
military  
mindfulness  
mini  
minimal  
minimum  
mistakes  
moment  
momentum  
monkey  
mornings  
motivated  
mouse  
movement  
multiple  
muscles  
music  
names  
nannan  
nc  
near  
necessary  
need  
needs  
neighborhoods  
neighbors  
new  
newly  
next  
nice  
nights  
notebooks  
novels  
numbers  
nurturing  
nutrition  
nv  
objects  
observing  
obstacles  
occasions  
odd  
odds  
offer  
often  
ohio  
old  
online  
opened  
opening  
opportunities  
opportunity  
order  
orderly  
outcomes  
outlet  
overall  
ownership  
pack  
pads

pads  
pain  
paintings  
pan  
parent  
park  
participants  
participate  
partnerships  
passionate  
patricia  
pbl  
pe  
peers  
pen  
pencils  
pens  
per  
percent  
perform  
period  
periods  
permanent  
permission  
perseverance  
personal  
philosophies  
physical  
pick  
picked  
pictures  
piece  
pillows  
pitch  
place  
plant  
platform  
play  
point  
points  
pollination  
poor  
portable  
pot  
practice  
pre  
preparing  
present  
presentation  
presentations  
pressure  
pretend  
prevent  
prevents  
prezi  
price  
priceless  
primarily  
primary  
print  
printer  
printing  
prior  
problems  
procedures  
process  
professional

professional  
proficiency  
program  
programmers  
progress  
progresses  
project  
projector  
promoting  
prompts  
properly  
protect  
protection  
proud  
provide  
provided  
provides  
providing  
public  
put  
puzzle  
quest  
quote  
race  
range  
rapid  
rarely  
rays  
reaches  
read  
reader  
readers  
reading  
ready  
realize  
receive  
receiving  
recent  
recess  
recognized  
recommendations  
record  
reduced  
refurbished  
relatable  
relevant  
rely  
remember  
request  
requesting  
research  
resources  
responsibilities  
rest  
restaurant  
review  
revolution  
rhyming  
rid  
rise  
risk  
robots  
rocking  
roles  
room  
rope  
ropes

ropes  
rotating  
rotation  
rounded  
rug  
ruined  
rural  
safe  
safely  
saver  
school  
schooling  
schools  
science  
scientist  
scientists  
score  
scores  
seating  
seats  
second  
sedentary  
seeing  
seek  
seeking  
selected  
senior  
sensory  
services  
set  
sets  
several  
shaped  
share  
shared  
sharing  
sharpener  
shed  
sheets  
show  
showing  
shows  
shy  
side  
sign  
silly  
similar  
since  
sing  
singing  
single  
sitting  
sixth  
sizes  
skills  
skin  
small  
smaller  
smoothies  
snacks  
soccer  
social  
socializing  
society  
soft  
software  
solution

solution  
something  
sort  
sound  
space  
spanish  
spanning  
sparking  
special  
specially  
specific  
specifically  
speech  
spell  
spelling  
spends  
spoon  
st  
stamina  
standard  
standards  
star  
start  
started  
starter  
stations  
statuses  
stomach  
stools  
storage  
store  
strategies  
strategy  
stream  
streamline  
streets  
strengths  
stronger  
structured  
structures  
struggle  
struggles  
students  
subscriptions  
suburb  
success  
suitable  
supplies  
support  
supporting  
supportive  
sweet  
systems  
table  
tables  
tablets  
take  
takes  
taking  
talented  
tales  
talkative  
talked  
tasks  
teach  
teacher  
teachers

teachers  
teaching  
team  
technological  
technology  
tell  
temporary  
terms  
test  
tests  
text  
texts  
thief  
thinking  
thirsty  
thousand  
three  
tight  
tissue  
title  
titles  
together  
told  
tomorrow  
tools  
totem  
touch  
toward  
towers  
track  
traditional  
trampoline  
translates  
trauma  
traveling  
treasured  
tree  
tremendously  
tried  
trimester  
trouble  
true  
try  
tucked  
twice  
types  
typically  
unable  
uniforms  
unique  
unit  
units  
us  
use  
used  
using  
variety  
vary  
varying  
verbal  
videos  
virtually  
vision  
visit  
visualize  
visuals  
vital

vital  
vocabulary  
wait  
walk  
want  
wanted  
waste  
water  
waters  
web  
webcam  
website  
weeks  
well  
wellness  
went  
whatever  
wheel  
whole  
wide  
within  
wobble  
wonderful  
wonderfully  
words  
work  
worked  
workforce  
working  
works  
world  
worry  
worrying  
would  
wrap  
writing  
year  
yearning  
yoga  
young  
younger  
zero  
2016  
2017  
ahead  
basketballs  
better  
big  
book  
books  
bookworms  
carpet  
centers  
chromebooks  
classic  
classroom  
create  
discovery  
dramatic  
drumming  
ed  
enhance  
explore  
first  
fit  
flood  
great

```
g,cat  
grow  
health  
hidden  
imaginative  
ipads  
keeps  
kindergarten  
kindle  
leaders  
learners  
library  
love  
manipulatives  
move  
my  
natural  
need  
needed  
not  
organization  
place  
playground  
preschool  
printer  
readers  
recess  
scientist  
see  
sit  
space  
special  
stand  
super  
teach  
teaching  
time  
under  
up  
us  
want  
we  
wiggle  
without  
writing  
price  
prev_proposed_projects  
sentiment_Test_P  
sentiment_Test_N  
sentiment_Test_NE  
sentiment_Test_C
```

```
1 %%time  
2 Dec_Tree = DecisionTreeClassifier(class_weight = 'balanced')  
3 parameters = {'max_depth': [1, 5, 10, 50], 'min_samples_split': [5, 10, 100, 500]}  
4 Dec_Tree_C = RandomizedSearchCV(Dec_Tree, parameters, cv=3, scoring='roc_auc',return_train_score=True,n_iter=30)  
5 Dec_Tree_C.fit(NZFI_TR, Y_Train)  
6 print(Dec_Tree_C.best_estimator_)  
7 #print(Dec_Tree_C.cv_results_)
```



```

DecisionTreeClassifier(class_weight='balanced', criterion='gini', max_depth=10,
                      max_features=None, max_leaf_nodes=None,
                      min_impurity_decrease=0.0, min_impurity_split=None,
                      min_samples_leaf=1, min_samples_split=500,
                      min_weight_fraction_leaf=0.0, presort=False,
                      random_state=None, splitter='best')
CPU times: user 6min 9s, sys: 71.7 ms, total: 6min 9s
Wall time: 6min 9s

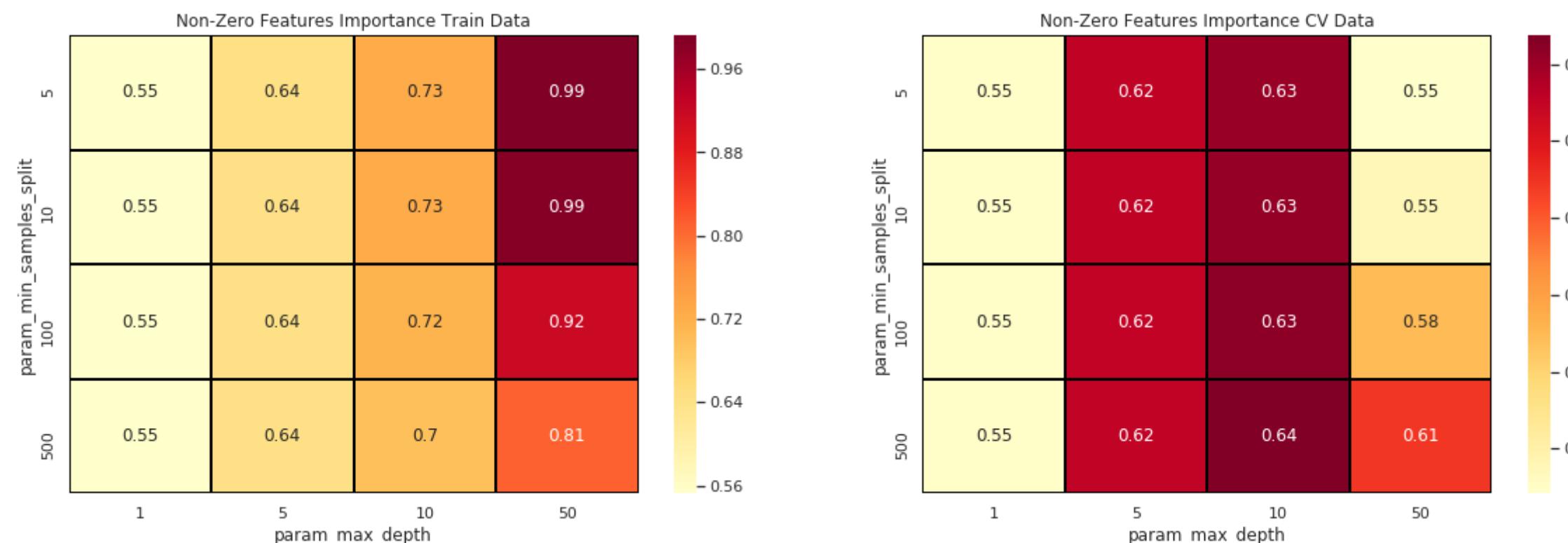
```

## ▼ Non-zero features importance Heatmap

```

1 sns.set()
2 Dec_Tree_CV = pd.DataFrame(Dec_Tree_C.cv_results_).groupby(['param_min_samples_split', 'param_max_depth']).max().unstack()[['mean_test_score', 'mean_train_score']]
3 fig, ax = plt.subplots(1,2, figsize=(20,6))
4 sns.heatmap(Dec_Tree_CV.mean_train_score, annot = True,cbar=True, fmt='.2g', cmap='YlOrRd', linewidths=1, linecolor='black', ax=ax[0])
5 sns.heatmap(Dec_Tree_CV.mean_test_score, annot = True, fmt='.2g', cmap='YlOrRd', linewidths=1, linecolor='black', ax=ax[1])
6 ax[0].set_title('Non-Zero Features Importance Train Data')
7 ax[1].set_title('Non-Zero Features Importance CV Data')
8 plt.show()

```



```

1 #Fitting Model to Hyper-Parameter Curve
2 # https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc\_curve.html#sklearn.metrics.roc\_curve
3
4 NZIF_MD=(Dec_Tree_C.best_params_['max_depth'])
5 NZIF_MSS=(Dec_Tree_C.best_params_['min_samples_split'])
6 Dec_Tree_C_OP=DecisionTreeClassifier(class_weight = 'balanced',max_depth=NZIF_MD,min_samples_split=NZIF_MSS)
7 # for visualization
8 Dec_Tree_C_OP.fit(NZFI_TR, Y_Train)
9 #https://scikit-learn.org/stable/modules/generated/sklearn.linear\_model.SGDClassifier.html#sklearn.linear\_model.SGDClassifier.decision\_function
10 Y_Train_PR = Dec_Tree_C_OP.predict_proba(NZFI_TR) [:,1]
11 Y_Test_PR = Dec_Tree_C_OP.predict_proba(NZFI_TS) [:,1]
12
13 fpr_Train, tpr_Train, thresholds_Train = roc_curve(Y_Train, Y_Train_PR)
14 fpr_Test, tpr_Test, thresholds_Test = roc_curve(Y_Test, Y_Test_PR)
15

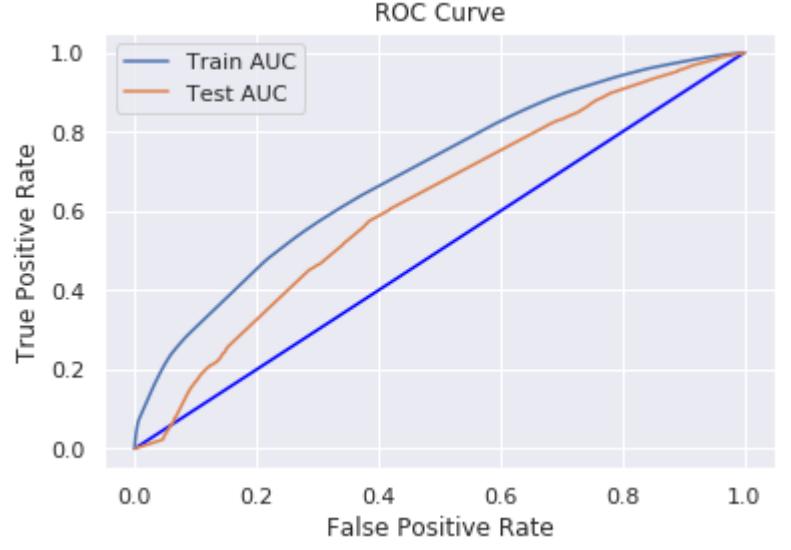
```

```

16 NZIF_AUC = auc(fpr_Test, tpr_Test)
17 plt.plot([0,1],[0,1],'k-',color='blue')
18 plt.plot(fpr_Train, tpr_Train, label="Train AUC")
19 plt.plot(fpr_Test, tpr_Test, label="Test AUC")
20 plt.legend()
21 plt.xlabel("False Positive Rate")
22 plt.ylabel("True Positive Rate")
23 plt.title("ROC Curve")
24 plt.show()
25 print("-"*120)
26 print("Optimal value of AUC Train =",auc(fpr_Train, tpr_Train))
27 print("Optimal value of AUC Test =",auc(fpr_Test, tpr_Test))
28 pred5 = Dec_Tree_C_OP.predict(NZFI_TR)
29 pred6 = Dec_Tree_C_OP.predict(NZFI_TS)

```

↳




---

Optimal value of AUC Train = 0.6938484698394856  
Optimal value of AUC Test = 0.6155285512428185

## ▼ Confusion Matrix

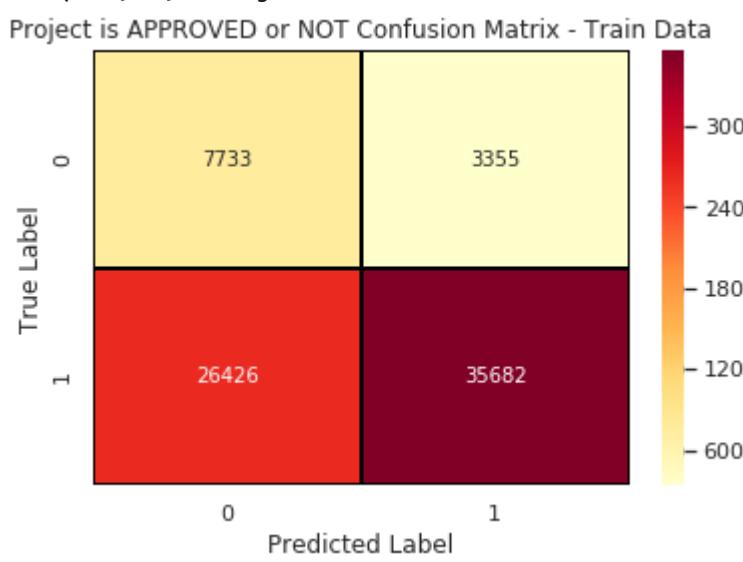
```

1 https://seaborn.pydata.org/generated/seaborn.heatmap.html
2 https://getaravind.com/blog/confusion-matrix-seaborn-heatmap/
3 https://stackoverflow.com/questions/19233771/sklearn-plot-confusion-matrix-with-labels
4 %matplotlib inline
5 Train = confusion_matrix(Y_Train, pred5)
6 sns.heatmap(Train, annot=True, cbar=True, fmt='d', cmap='YlOrRd', linewidths=1, linecolor='black')
7 plt.ylabel('True Label')
8 plt.xlabel('Predicted Label')
9 plt.title('Project is APPROVED or NOT Confusion Matrix - Train Data')

```

↳

Text(0.5, 1, 'Project is APPROVED or NOT Confusion Matrix - Train Data')



#### OBSERVATION:

True Negative = 7733; False Negative = 26426; True Positive = 35682; False Positive = 3355

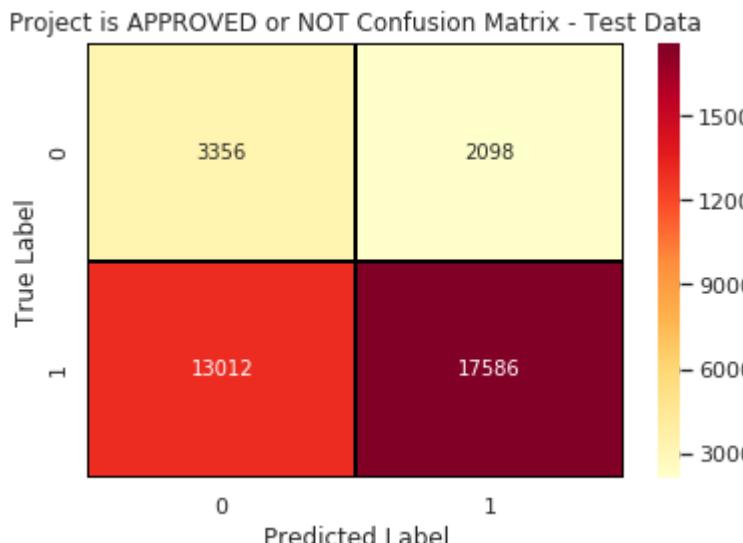
Accuracy (Overall, how often is the classifier correct) = 0.60

Precision(When it predicts yes, how often is it correct) = 0.92

Misclassification (Overall, how often is it wrong) = 0.41

```
1 #https://seaborn.pydata.org/generated/seaborn.heatmap.html
2 #https://getaravind.com/blog/confusion-matrix-seaborn-heatmap/
3 #https://stackoverflow.com/questions/19233771/sklearn-plot-confusion-matrix-with-labels
4 Test = confusion_matrix(Y_Test, pred6)
5 sns.heatmap(Test, annot=True, cbar=True, fmt='d', cmap='YlOrRd', linewidths=1, linecolor='black')
6 plt.ylabel('True Label')
7 plt.xlabel('Predicted Label')
8 plt.title('Project is APPROVED or NOT Confusion Matrix - Test Data')
```

→ Text(0.5, 1, 'Project is APPROVED or NOT Confusion Matrix - Test Data')



#### OBSERVATION:

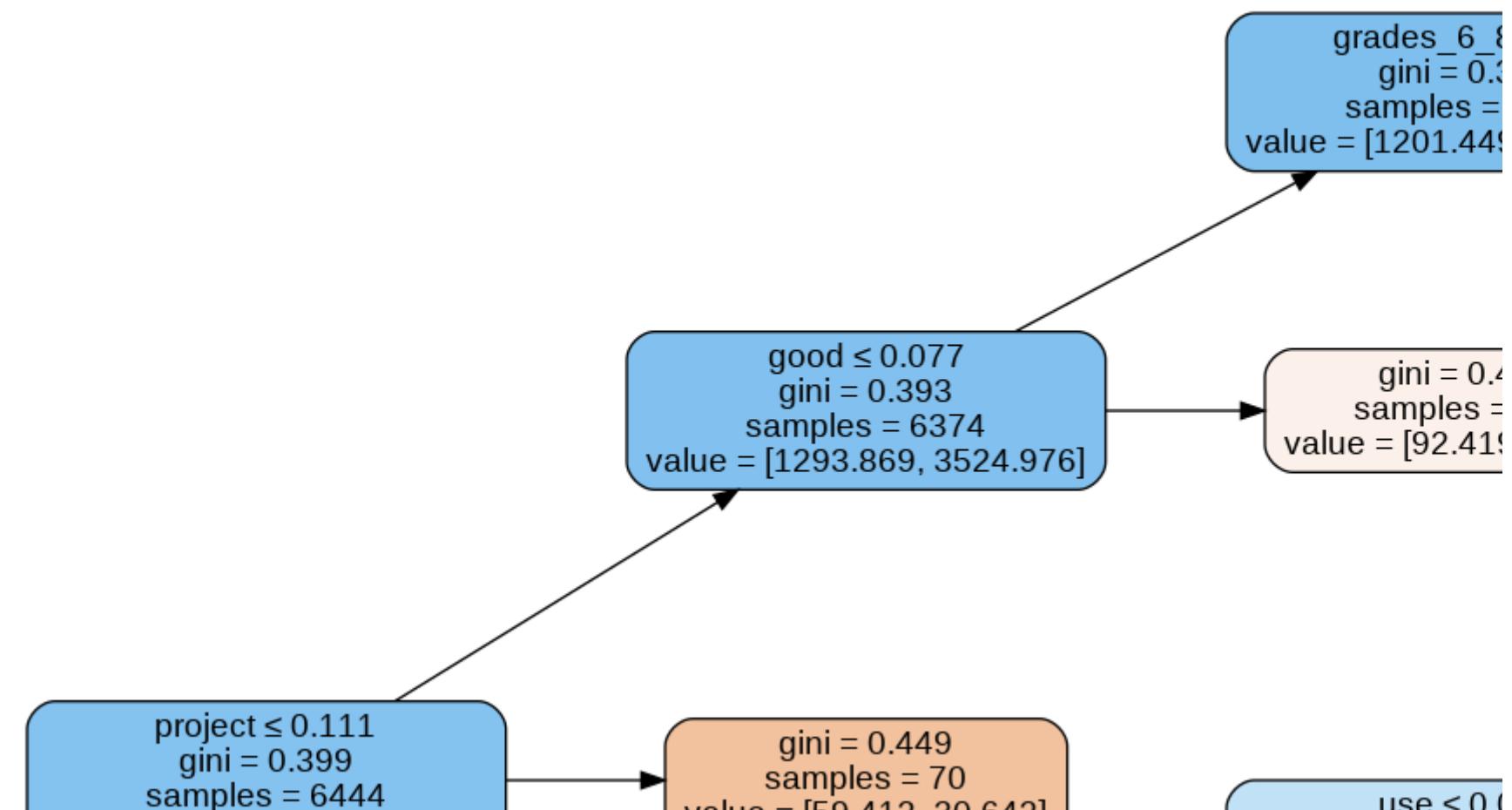
True Negative = 3356; False Negative = 13012; True Positive = 17586; False Positive = 2098

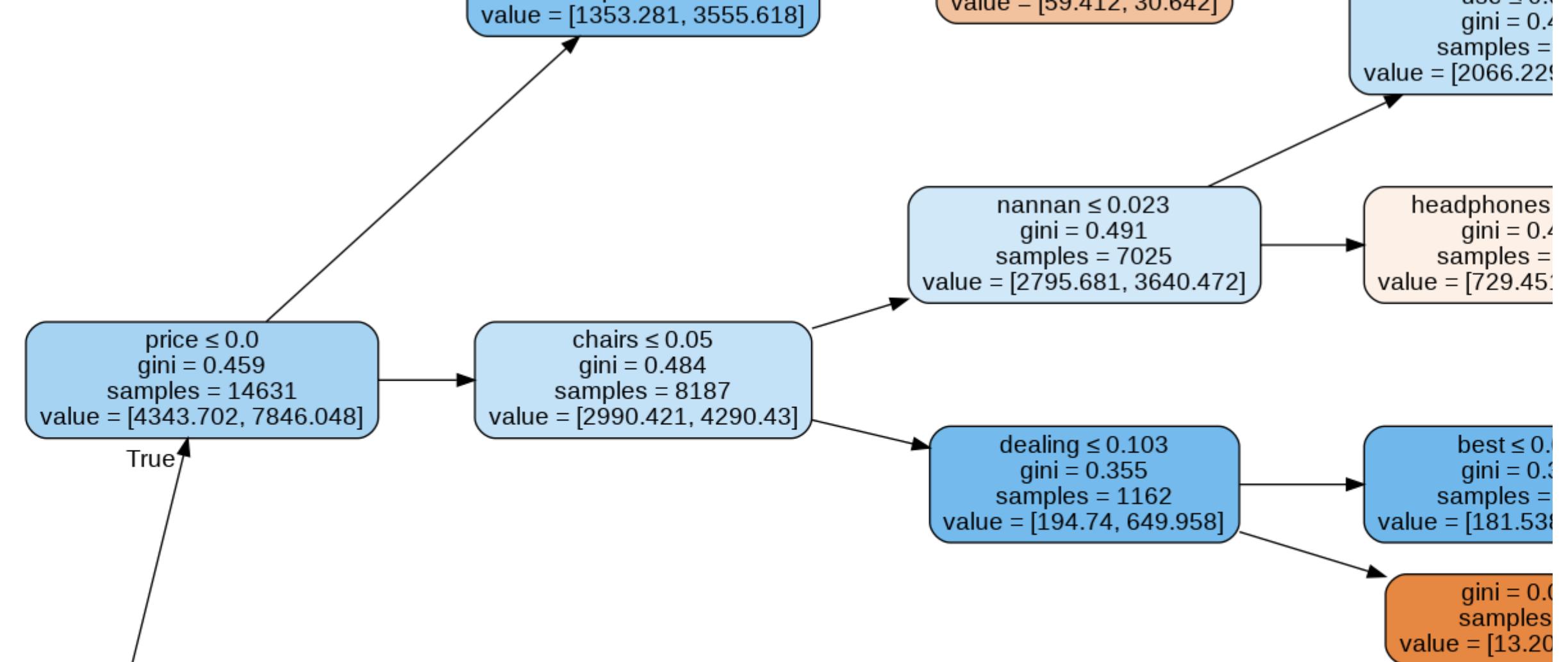
Accuracy (Overall, how often is the classifier correct) = 0.59

Precision(When it predicts yes, how often is it correct) =0.90

```
1 # https://medium.com/@rnbrown/creating-and-visualizing-decision-trees-with-python-f8e8fa394176
2 Dec_Tree_C_OP=DecisionTreeClassifier(class_weight = 'balanced',max_depth=TFIDF_MD,min_samples_split=TFIDF_MSS)
3 Dec_Tree_C_OP.fit(NZFI_TR, Y_Train)
4 graph.format = 'png'
5 dot_data = StringIO()
6 export_graphviz(Dec_Tree_C_OP, out_file=dot_data, filled=True, rounded=True, special_characters=True, feature_names=NZFI_TR.columns,rotate=True)
7 graph = pydotplus.graph_from_dot_data(dot_data.getvalue())
8 Image(graph.create_png())
```

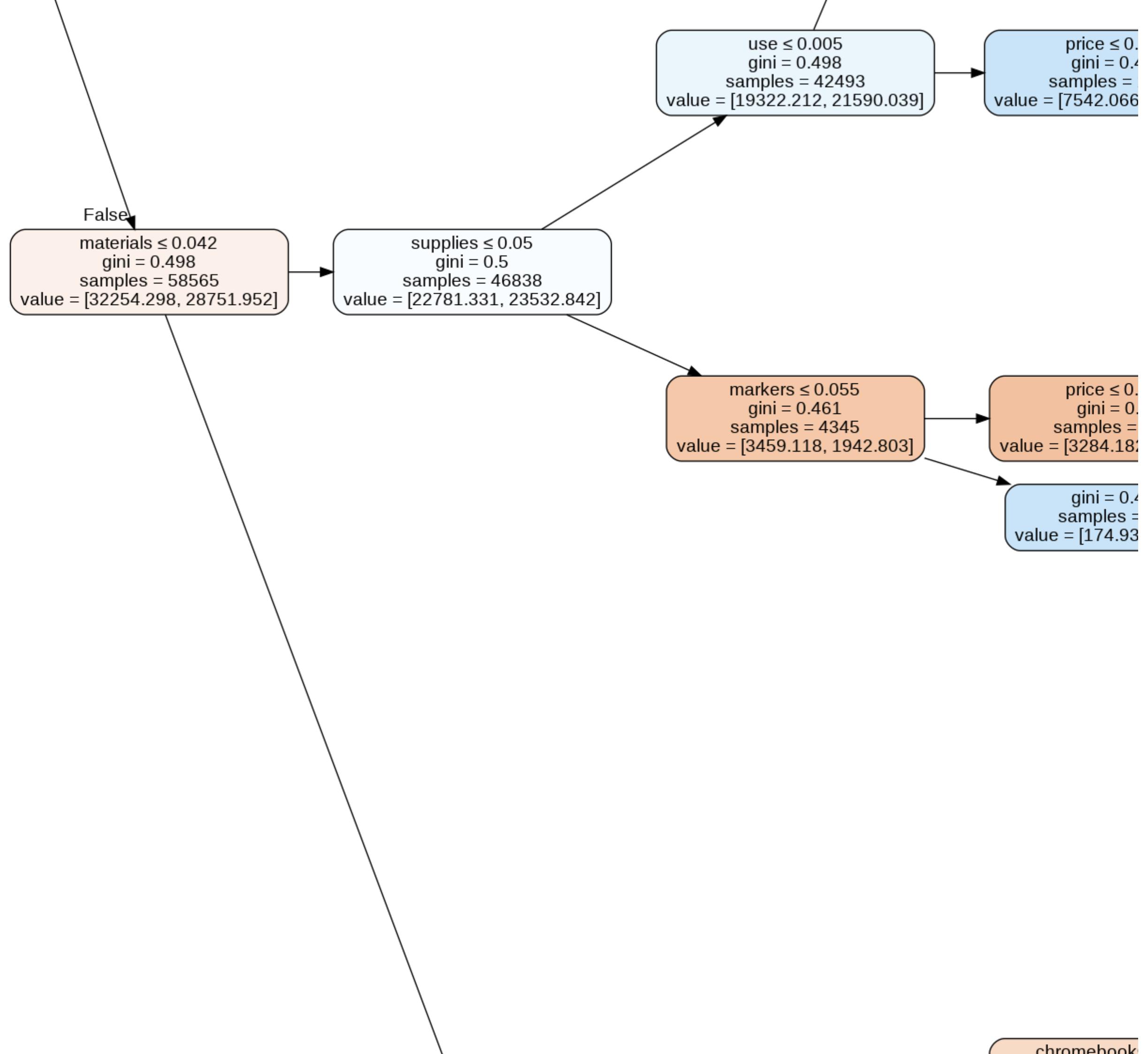
⇨

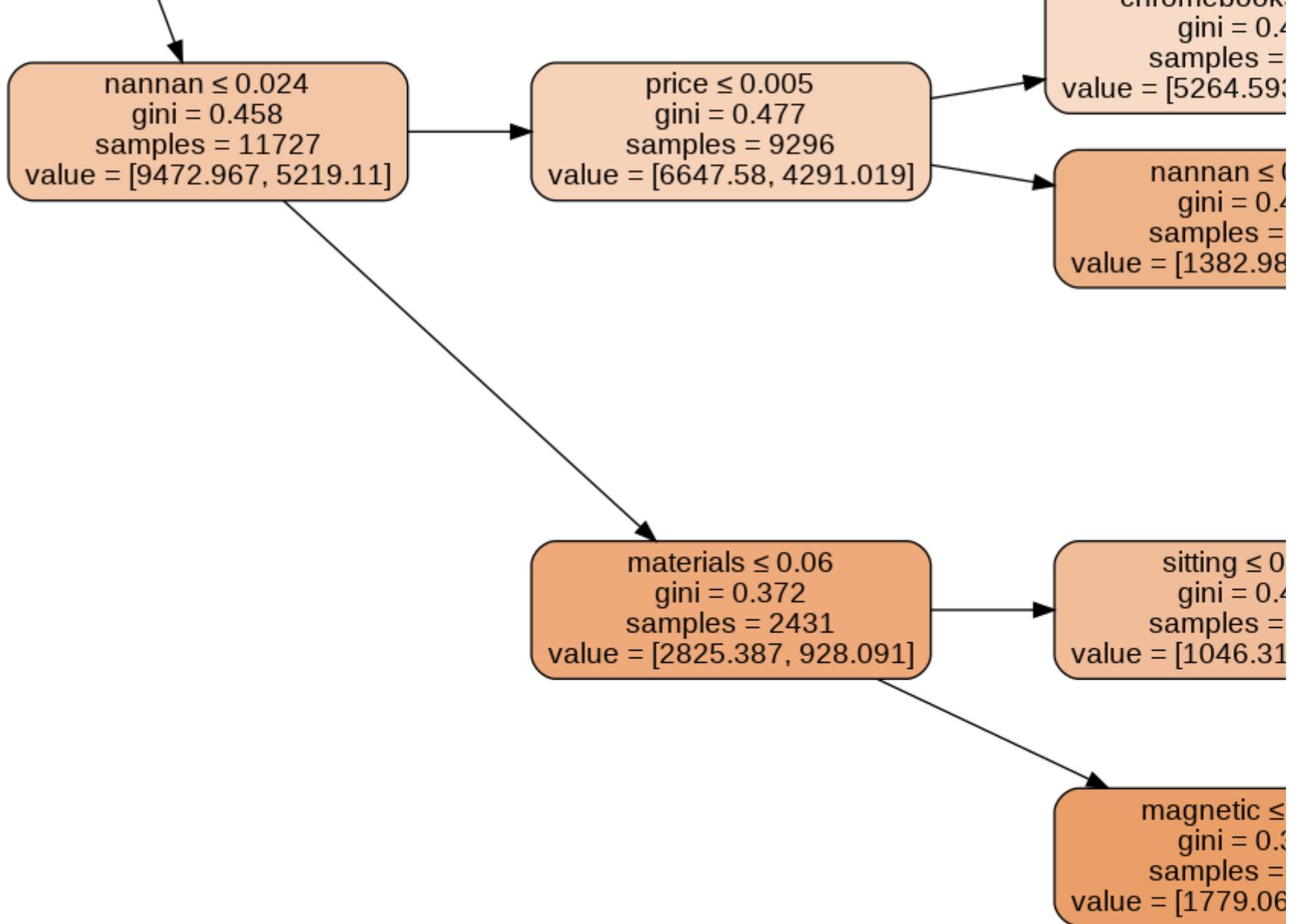




price  $\leq$  0.001  
gini = 0.5  
samples = 73196  
value = [36598.0, 36598.0]

stools  $\leq$  0  
gini = 0  
samples = 11780  
value = [11780.140, 11780.140]





### 3. Conclusions

```

1 # Please compare all your models using Prettytable library
2
3 pt = PrettyTable()
4 pt.field_names= ("S.No", "Vectorizer", "Model", "Max Depth", "Min Samples Split", "AUC")
5 pt.add_row(["1", "TFIDF", "DecisionTree", TFIDF_MD, TFIDF_MSS, TFIDF_AUC])
6 pt.add_row(["2", "TFIDF_W2V", "DecisionTree", TFIDF_W2V_MD, TFIDF_W2V_MSS, TFIDF_W2V_AUC])
7 pt.add_row(["3", "[Task-2]-Non Zero Features Importance", "DecisionTree", NZIF_MD, NZIF_MSS, NZIF_AUC])
8 print(pt)

```

S.No	Vectorizer	Model	Max Depth	Min Samples Split	AUC
1	TFIDF	DecisionTree	10	500	0.6218123067835468
2	TFIDF_W2V	DecisionTree	5	500	0.6240460475988554
3	[Task-2]-Non Zero Features Importance	DecisionTree	10	500	0.6155285512428185

#### SUMMARY:

1. Compare to KNN and SVM "**DecisionTree**" had less space and time complexity.
2. Overall test accuracy score is average.
3. There is no much difference in the AUC score for Non Zero important fetaures and TFIDF. Both are look similar only.
4. "**TFIDF W2V**" performing very well in accuracy and misclassification scores. When compare to other vectorizer's.