

▼ SVD and Matrix Factorization on DonorsChoose

DonorsChoose.org receives hundreds of thousands of project proposals each year for classroom projects in need of funding. Right now 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
<code>project_id</code>	A unique identifier for the proposed project. Example: p036502
<code>project_title</code>	Title of the project. Examples: <ul style="list-style-type: none"> • Art Will Make You Happy! • First Grade Fun
<code>project_grade_category</code>	Grade level of students for which the project is targeted. One of the following enumerated values: <ul style="list-style-type: none"> • Grades PreK-2 • Grades 3-5 • Grades 6-8 • Grades 9-12

Feature	Description
<code>project_subject_categories</code>	<p>One or more (comma-separated) subject categories for the project from the following enumerated list of values:</p> <ul style="list-style-type: none"> • Applied Learning • Care & Hunger • Health & Sports • History & Civics • Literacy & Language • Math & Science • Music & The Arts • Special Needs • Warmth <p>Examples:</p> <ul style="list-style-type: none"> • Music & The Arts • Literacy & Language, Math & Science
<code>school_state</code>	<p>State where school is located (Two-letter U.S. postal code). Example: WY</p>
<code>project_subject_subcategories</code>	<p>One or more (comma-separated) subject subcategories for the project. Examples:</p> <ul style="list-style-type: none"> • Literacy • Literature & Writing, Social Sciences
<code>project_resource_summary</code>	<p>An explanation of the resources needed for the project. Example:</p> <ul style="list-style-type: none"> • My students need hands on literacy materials to manage sensory needs!
<code>project_essay_1</code>	First application essay*
<code>project_essay_2</code>	Second application essay*
<code>project_essay_3</code>	Third application essay*
<code>project_essay_4</code>	Fourth application essay*
<code>project_submitted_datetime</code>	Datetime when project application was submitted. Example: 2016-04-28 12:43:56.245
<code>teacher_id</code>	<p>A unique identifier for the teacher of the proposed project. Example: bdf8baa8fedef6bfeec7ae4ff1c15c56</p> <p>Teacher's title. One of the following enumerated values:</p> <ul style="list-style-type: none"> • nan • Dr. • Mr. • Mrs. • Ms. • Teacher.
<code>teacher_prefix</code>	
<code>teacher_number_of_previously_posted_projects</code>	Number of project applications previously submitted by the same teacher. Example: 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. Example: <code>p036502</code>
<code>description</code>	Description of the resource. Example: Tenor Saxophone Reeds, Box of 25
<code>quantity</code>	Quantity of the resource required. Example: 3
<code>price</code>	Price of the resource required. Example: 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 resources needed for a project:

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

Label	Description
<code>project_is_approved</code>	A binary flag indicating whether DonorsChoose approved the project. A value of 0 indicates the project was not approved, and a value of 1 indicates the project was approved.

▼ 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 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 b

```
1 ad&authuser=0&nonce=4jft1g4q59abm&user=03671904864144121079&hash=486rrlnvdfnbi5i5vdkq4ol6bnqk57be" -O "11_DonorsChoose_Tru
```

```
--2019-11-14 07:02:23-- https://doc-0g-5o-docs.googleusercontent.com/docs/securesc/gkcko5omvjm8g0ffrjouskdfeosq6vei/0hi
Resolving doc-0g-5o-docs.googleusercontent.com (doc-0g-5o-docs.googleusercontent.com)... 173.194.218.132, 2607:f8b0:400
Connecting to doc-0g-5o-docs.googleusercontent.com (doc-0g-5o-docs.googleusercontent.com)|173.194.218.132|:443... conne
HTTP request sent, awaiting response... 403 Forbidden
2019-11-14 07:02:23 ERROR 403: Forbidden.
```


```
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 import string
20 from nltk.corpus import stopwords
```

```


20 from nltk.corpus import stopwords
21 from nltk.stem import PorterStemmer
22 from nltk.stem.wordnet import WordNetLemmatizer
23 # from gensim.models import Word2Vec
24 # from gensim.models import KeyedVectors
25 import pickle
26 from tqdm import tqdm
27 import os
28 # from plotly import plotly
29 # import plotly.offline as offline
30 # import plotly.graph_objs as go
31 # offline.init_notebook_mode()
32 # from collections import Counter

```

```
1 ! wget --header="Host: doc-0k-5o-docs.googleusercontent.com" --header="User-Agent: Mozilla/5.0 (Windows NT 10.0; Win64; x86_64; rv:68.0) Gecko/20100101 Firefox/68.0" https://doc-0k-5o-docs.googleusercontent.com/docs/securesc/gkcko5omvjm8g0ffrjouskdfeosq6vei/cr...
```

 --2019-11-14 07:02:26-- <https://doc-0k-5o-docs.googleusercontent.com/docs/securesc/gkcko5omvjm8g0ffrjouskdfeosq6vei/cr...>
 Resolving doc-0k-5o-docs.googleusercontent.com (doc-0k-5o-docs.googleusercontent.com)... 173.194.218.132, 2607:f8b0:400d:20::...
 Connecting to doc-0k-5o-docs.googleusercontent.com (doc-0k-5o-docs.googleusercontent.com)|173.194.218.132|:443...
 HTTP request sent, awaiting response... 403 Forbidden
 2019-11-14 07:02:26 ERROR 403: Forbidden.

```
1 ! wget --header="Host: doc-0c-5o-docs.googleusercontent.com" --header="User-Agent: Mozilla/5.0 (Windows NT 10.0; Win64; x86_64; rv:68.0) Gecko/20100101 Firefox/68.0" https://doc-0c-5o-docs.googleusercontent.com/docs/securesc/gkcko5omvjm8g0ffrjouskdfeosq6vei/2g...
```

 --2019-11-14 07:02:26-- <https://doc-0c-5o-docs.googleusercontent.com/docs/securesc/gkcko5omvjm8g0ffrjouskdfeosq6vei/2g...>
 Resolving doc-0c-5o-docs.googleusercontent.com (doc-0c-5o-docs.googleusercontent.com)... 173.194.218.132, 2607:f8b0:400d:20::...
 Connecting to doc-0c-5o-docs.googleusercontent.com (doc-0c-5o-docs.googleusercontent.com)|173.194.218.132|:443...
 HTTP request sent, awaiting response... 403 Forbidden
 2019-11-14 07:02:26 ERROR 403: Forbidden.

```

1 project_data = pd.read_csv('train_data.csv')
2 resource_data = pd.read_csv('resources.csv')

```

```
1 project_data.head(5)
```

Unnamed:
0

id

teacher_id

teacher_prefix

school_state

project_submitted_datetime

project

0	160221	p253737	c90749f5d961ff158d4b4d1e7dc665fc	Mrs.	IN	2016-12-05 13:43:57
1	140945	p258326	897464ce9ddc600bced1151f324dd63a	Mr.	FL	2016-10-25 09:22:10
2	21895	p182444	3465aaf82da834c0582ebd0ef8040ca0	Ms.	AZ	2016-08-31 12:03:56
3	45	p246581	f3cb9bffbba169bef1a77b243e620b60	Mrs.	KY	2016-10-06 21:16:17
4	172407	p104768	be1f7507a41f8479dc06f047086a39ec	Mrs.	TX	2016-07-11 01:10:09

```

1 # We are removing the words from the stop words list: 'no', 'nor', 'not'
2 stopwords= ['i', 'me', 'my', 'myself', 'we', 'our', 'ours', 'ourselves', 'you', "you're", "you've", \
3             "you'll", "you'd", 'your', 'yours', 'yourself', 'yourselves', 'he', 'him', 'his', 'himself', \
4             'she', "she's", 'her', 'hers', 'herself', 'it', "it's", 'its', 'itself', 'they', 'them', 'their', \
5             'theirs', 'themselves', 'what', 'which', 'who', 'whom', 'this', 'that', "that'll", 'these', 'those', \
6             'am', 'is', 'are', 'was', 'were', 'be', 'been', 'being', 'have', 'has', 'had', 'having', 'do', 'does', \
7             'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'or', 'because', 'as', 'until', 'while', 'of', \
8             'at', 'by', 'for', 'with', 'about', 'against', 'between', 'into', 'through', 'during', 'before', 'after', \
9             'above', 'below', 'to', 'from', 'up', 'down', 'in', 'out', 'on', 'off', 'over', 'under', 'again', 'further', \
10            'then', 'once', 'here', 'there', 'when', 'where', 'why', 'how', 'all', 'any', 'both', 'each', 'few', 'more', \

```

```


11      'most', 'other', 'some', 'such', 'only', 'own', 'same', 'so', 'than', 'too', 'very', \
12      's', 't', 'can', 'will', 'just', 'don', "don't", 'should', "should've", 'now', 'd', 'll', 'm', 'o', 're', \
13      've', 'y', 'ain', 'aren', "aren't", 'couldn', "couldn't", 'didn', "didn't", 'doesn', "doesn't", 'hadn', \
14      "hadn't", 'hasn', "hasn't", 'haven', "haven't", 'isn', "isn't", 'ma', 'mightn', "mightn't", 'mustn', \
15      "mustn't", 'needn', "needn't", 'shan', "shan't", 'shouldn', "shouldn't", 'wasn', "wasn't", 'weren', "weren't".
16      'won', "won't", 'wouldn', "wouldn't"]

```

```

1 from tqdm import tqdm
2 clean_titles = []
3 for sent in tqdm(project_data["project_title"].values):
4     titles= " "
5     for w in sent.lower().split():
6         if w not in stopwords:
7             titles = titles + w + " "
8     clean_titles.append(titles.strip())
9
10 project_data["clean_titles"]=clean_titles
11 project_data.drop(["project_title"],axis =1 , inplace = True )
12
13
14
15
16


```

 100%|██████████| 109248/109248 [00:01<00:00, 94158.42it/s]

```

1 # Replace the NAN Values :
2 project_data = project_data.dropna(subset=["teacher_prefix"])
3 sujit = project_data[project_data["teacher_prefix"].isnull()]
4 sujit
5

```

 Unnamed: 0 id teacher_id teacher_prefix school_state project_submitted_datetime project_grade_category project_s

```
1 project_data = project_data[0 : 50000]
```

```
2 print(project_data.shape)
```


 (50000, 17)

```
1 Y=project_data["project_is_approved"]
2 len(Y)
```

 50000

```
1 # merge two column text dataframe:
2 project_data["essay"] = project_data["project_essay_1"].map(str) + \
3                             project_data["project_essay_2"].map(str) + \
4                             project_data["project_essay_3"].map(str) + \
5                             project_data["project_essay_4"].map(str)
```

```
1 project_data.head(2)
```

 Unnamed: 0 id teacher_id teacher_prefix school_state project_submitted_datetime projec

0	160221	p253737	c90749f5d961ff158d4b4d1e7dc665fc	Mrs.	IN	2016-12-05 13:43:57
1	140945	p258326	897464ce9ddc600bced1151f324dd63a	Mr.	FL	2016-10-25 09:22:10

```
1 # https://stackoverflow.com/a/47091490/4084039
```

```
2 import re
```



```


1 import re
2
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
9     # general
10    phrase = re.sub(r"n't", " not", phrase)
11    phrase = re.sub(r"\ 're", " are", phrase)
12    phrase = re.sub(r"\ 's", " is", phrase)
13    phrase = re.sub(r"\ 'd", " would", phrase)
14    phrase = re.sub(r"\ 'll", " will", phrase)
15    phrase = re.sub(r"\ 't", " not", phrase)
16    phrase = re.sub(r"\ 've", " have", phrase)
17    phrase = re.sub(r"\ 'm", " am", phrase)
18    return phrase

```

```

1 # Combining all the above students
2 from tqdm import tqdm
3 preprocessed_essays = []
4 # tqdm is for printing the status bar
5 for sentence in tqdm(project_data['essay'].values):
6     sent = decontracted(sentence)
7     sent = sent.replace('\r', ' ')
8     sent = sent.replace('\n', ' ')
9     sent = sent.replace('\n', ' ')
10    sent = re.sub('[^A-Za-z0-9]+', ' ', sent)
11    # https://gist.github.com/sebleier/554280
12    sent = ' '.join(e for e in sent.split() if e not in stopwords)
13    preprocessed_essays.append(sent.lower().strip())

```

 100% | ██████████ | 50000/50000 [00:26<00:00, 1898.74it/s]

```
1 project_data['preprocessed_essays'] = preprocessed_essays
```

```
1 project_data['combined'] = project_data['preprocessed_essays'] + ' ' + project_data['clean_titles']
```

```
1 project_data.head(2)
```



Unnamed:
0

id

teacher_id

teacher_prefix

school_state

project_submitted_datetime

project

0	160221	p253737	c90749f5d961ff158d4b4d1e7dc665fc	Mrs.	IN	2016-12-05 13:43:57
---	--------	---------	----------------------------------	------	----	---------------------

1	140945	p258326	897464ce9ddc600bced1151f324dd63a	Mr.	FL	2016-10-25 09:22:10
---	--------	---------	----------------------------------	-----	----	---------------------

```
1 categories = list(project_data['project_subject_categories'].values)
2 # remove special characters from list of strings python: https://stackoverflow.com/a/47301924/4084039
3
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 cat_list = []
8 for i in categories:
9     temp = ""
10    # consider we have text like this "Math & Science, Warmth, Care & Hunger"
11    for j in i.split(','): # it will split it in three parts ["Math & Science", "Warmth", "Care & Hunger"]
12        if 'The' in j.split(): # this will split each of the category based on space "Math & Science"=> "Math","&", "Scie
13            j=j.replace('The','') # if we have the words "The" we are going to replace it with ''(i.e removing 'The')
14            j = j.replace(' ','') # we are placing all the ' '(space) with ''(empty) ex:"Math & Science"=>"Math&Science"
15            temp+=j.strip()+" " # " abc ".strip() will return "abc", remove the trailing spaces
16            temp = temp.replace('&','_') # we are replacing the & value into
17    cat_list.append(temp.strip())
```

```

18
19 project_data['clean_categories'] = cat_list
20 project_data.drop(['project_subject_categories'], axis=1, inplace=True)
21
22 from collections import Counter
23 my_counter = Counter()
24 for word in project_data['clean_categories'].values:
25     my_counter.update(word.split())
26
27 cat_dict = dict(my_counter)
28 sorted_cat_dict = dict(sorted(cat_dict.items(), key=lambda kv: kv[1]))
29

1 sub_catogories = list(project_data['project_subject_subcategories'].values)
2 # remove special characters from list of strings python: https://stackoverflow.com/a/47301924/4084039
3
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
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 catogory based on space "Math & Science"=> "Math","&", "Scie
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 project_data['clean_subcategories'] = sub_cat_list
21 project_data.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 project_data['clean_subcategories'].values:
26     my_counter.update(word.split())

```

```
27
28 sub_cat_dict = dict(my_counter)
29 sorted_sub_cat_dict = dict(sorted(sub_cat_dict.items(), key=lambda kv: kv[1]))

1 price_data = resource_data.groupby('id').agg({'price':'sum', 'quantity':'sum'}).reset_index()
2 project_data = pd.merge(project_data, price_data, on='id', how='left')
```

► Computing Sentiment Scores

↳ 2 cells hidden

► Adding new feature Number of words in Essay and Title

↳ 6 cells hidden

► Split the Data-Matrix into Train,Test and Cross-validation

↳ 1 cell hidden

► Selecting the top-most features from essay and title .

1. The reasons behind the best 2000 features is , as the word corpus will have lot of words . so we will select the top 2000 words for the Matrix factoriz

↳ 7 cells hidden

► TruncatedSVD on calculated Co-Occurence matrix

↳ 4 cells hidden

► steps creating the word vector :

1. we have to ceate a dictionary of all the index with the key as the corosponding vectors .

↳ 5 cells hidden

► Vectorizing Textual Data

↳ 4 cells hidden

- ▼ Same has to be done with the cv data and the test data .

as in case we use the Glov vector for example :

here we have to give our own custom words for the list of words that we have to use for making the AVG word to vector

```
1. X_test_avg_w2v_vectors = []; # the avg-w2v for each sentence/review is stored in this list
2. for sentence in tqdm(X_test['clean_titles'].values): # for each review/sentence
3.     vector = np.zeros(300) # as word vectors are of zero length
4.     cnt_words = 0; # num of words with a valid vector in the sentence/review
5.     for word in sentence.split(): # for each word in a review/sentence
6.         if word in glove_words:
7.             vector += model[word]
8.         cnt_words += 1
9.     if cnt_words != 0:
10.        vector /= cnt_words
11. X_test_avg_w2v_vectors.append(vector)
```

```
1 #Word count for Cross validation data
2 cvfinal = calculating(xcv['combined'])
3 print("Number of words in CV data",len( cvfinal ))
4 tefinal = calculating(xtest['combined'])
5 print("Number of words in test data",len(tefinal))
```



```
100%|██████████| 11055/11055 [00:00<00:00, 14339.00it/s]
```

```
1 # Avg word vector of the cv data
2 xcvAW2V = avgw2v(cvfinal,svdFeat)
3 xcvAW2V=np.asarray(xcvAW2V)
4 print("shape of the AVG word Vector of CV data ",xcvAW2V.shape)
5 xteAW2V = np.asarray(avgw2v(tefinal,svdFeat))
6 print("shape of the AVG word Vector of test data ",xteAW2V.shape)
```

```
100%|██████████| 11055/11055 [00:20<00:00, 546.97it/s]
0%|██████████| 57/16500 [00:00<00:29, 563.47it/s]shape of the AVG word Vector of CV data (11055, 999)
100%|██████████| 16500/16500 [00:30<00:00, 548.78it/s]
shape of the AVG word Vector of test data (16500, 999)
```

► Vectorising the catogorial data

↳ 23 cells hidden

► Heatmap representaion of the confusion metrix

↳ 1 cell hidden

► Conclusion

↳ 2 cells hidden

