

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
<code>project_id</code>		A unique identifier for the proposed project. Example: p036502
<code>project_title</code>	<ul style="list-style-type: none">••	Title of the project. Examples: <code>Art Will Make You Happy!</code> <code>First Grade Fun</code>
<code>project_grade_category</code>	<ul style="list-style-type: none">••••	Grade level of students for which the project is targeted. One of the following enumerated values: <code>Grades PreK-2</code> <code>Grades 3-5</code> <code>Grades 6-8</code> <code>Grades 9-12</code>
<code>project_subject_categories</code>	<ul style="list-style-type: none">•••••••••	One or more (comma-separated) subject categories for the project from the following enumerated list of values: <code>Applied Learning</code> <code>Care & Hunger</code> <code>Health & Sports</code> <code>History & Civics</code> <code>Literacy & Language</code> <code>Math & Science</code> <code>Music & The Arts</code> <code>Special Needs</code> <code>Warmth</code> Examples: <ul style="list-style-type: none">• <code>Music & The Arts</code>• <code>Literacy & Language, Math & Science</code>
<code>school_state</code>		State where school is located (Two-letter U.S. postal code). Example: WY
<code>project_subject_subcategories</code>	<ul style="list-style-type: none">••	One or more (comma-separated) subject subcategories for the project. Examples: <code>Literacy</code> <code>Literature & Writing, Social Sciences</code>
<code>project_resource_summary</code>	<ul style="list-style-type: none">•	An explanation of the resources needed for the project. Example: <code>My students need hands on literacy materials to manage sensory needs!</code>
<code>project_essay_1</code>		First application essay*
<code>project_essay_2</code>		Second application essay*
<code>project_essay_3</code>		Third application essay*

Feature	Description
project_essay_4	Fourth application essay
project_submitted_datetime	Datetime when project application was submitted. Example: 2016-04-28 12:43:56.245
teacher_id	A unique identifier for the teacher of the proposed project. Example: bdf8baa8fedef6bfeec7ae4ff1c15c56
teacher_prefix	Teacher's title. One of the following enumerated values: <ul style="list-style-type: none"> nan Dr. Mr. Mrs. Ms. Teacher.
teacher_number_of_previously_posted_projects	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
id	A <code>project_id</code> value from the <code>train.csv</code> file. Example: p036502
description	Description of the resource. Example: Tenor Saxophone Reeds, Box of 25
quantity	Quantity of the resource required. Example: 3
price	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 all resources needed for a project:

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

Label	Description
project_is_approved	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 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.

In [1]:

```
%matplotlib inline
import warnings
warnings.filterwarnings("ignore")

import sqlite3
import pandas as pd
import numpy as np
import nltk
import string
import matplotlib.pyplot as plt
```

```

import seaborn as sns
from sklearn.feature_extraction.text import TfidfTransformer
from sklearn.feature_extraction.text import TfidfVectorizer

from sklearn.feature_extraction.text import CountVectorizer
from sklearn.metrics import confusion_matrix
from sklearn import metrics
from sklearn.metrics import roc_curve, auc
from nltk.stem.porter import PorterStemmer

import re
# Tutorial about Python regular expressions: https://pymotw.com/2/re/
import string
from nltk.corpus import stopwords
from nltk.stem import PorterStemmer
from nltk.stem.wordnet import WordNetLemmatizer

from gensim.models import Word2Vec
from gensim.models import KeyedVectors
import pickle

import time
from tqdm import tqdm
import os
import pickle

from chart_studio import plotly
import plotly.offline as offline
import plotly.graph_objs as go
offline.init_notebook_mode()
from collections import Counter

```

1.1 Reading Data

In [2]:

```

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

```

In [3]:

```

print("Number of data points in train data", project_data.shape)
print('-'*50)
print("The attributes of data :", project_data.columns.values)

```

Number of data points in train data (109248, 17)

The attributes of data : ['Unnamed: 0' 'id' 'teacher_id' 'teacher_prefix' 'school_state' 'project_submitted_datetime' 'project_grade_category' 'project_subject_categories' 'project_subject_subcategories' 'project_title' 'project_essay_1' 'project_essay_2' 'project_essay_3' 'project_essay_4' 'project_resource_summary' 'teacher_number_of_previously_posted_projects' 'project_is_approved']

In [4]:

```

# how to replace elements in list python: https://stackoverflow.com/a/2582163/4084039
cols = ['Date' if x=='project_submitted_datetime' else x for x in list(project_data.columns)]

#sort dataframe based on time pandas python: https://stackoverflow.com/a/49702492/4084039
project_data['Date'] = pd.to_datetime(project_data['project_submitted_datetime'])
project_data.drop('project_submitted_datetime', axis=1, inplace=True)
project_data.sort_values(by=['Date'], inplace=True)

# how to reorder columns pandas python: https://stackoverflow.com/a/13148611/4084039
project_data = project_data[cols]

project_data.head(2)

```

Out[4]:

Unnamed: 0	id	teacher_id	teacher_prefix	school_state	Date	project_grade_category	project_
55660	8393 p205479	2bf07ba08945e5d8b2a3f269b2b3cfe5	Mrs.	CA	2016-04-27 00:27:36	Grades PreK-2	
76127	37728 p043609	3f60494c61921b3b43ab61bdde2904df	Ms.	UT	2016-04-27 00:31:25	Grades 3-5	

In [5]:

```
print("Number of data points in train data", resource_data.shape)
print(resource_data.columns.values)
print(resource_data.head(2))

# https://stackoverflow.com/questions/22407798/how-to-reset-a-dataframes-indexes-for-all-groups-in-one-step
price_data = resource_data.groupby('id').agg({'quantity':'sum', 'price':'sum'}).reset_index()

# Join two data frames
project_data = pd.merge(project_data, price_data, on='id', how='left')
project_data.head(5)
```

```
Number of data points in train data (1541272, 4)
['id' 'description' 'quantity' 'price']
   id                        description  quantity  \
0  p233245  LC652 - Lakeshore Double-Space Mobile Drying Rack      1
1  p069063      Bouncy Bands for Desks (Blue support pipes)      3

   price
0  149.00
1   14.95
```

Out[5]:

Unnamed: 0	id	teacher_id	teacher_prefix	school_state	Date	project_grade_category	project_subject
0	8393 p205479	2bf07ba08945e5d8b2a3f269b2b3cfe5	Mrs.	CA	2016-04-27 00:27:36	Grades PreK-2	
1	37728 p043609	3f60494c61921b3b43ab61bdde2904df	Ms.	UT	2016-04-27 00:31:25	Grades 3-5	
2	74477 p189804	4a97f3a390bfe21b99cf5e2b81981c73	Mrs.	CA	2016-04-27 00:46:53	Grades PreK-2	Litera
3	100660 p234804	cbc0e38f522143b86d372f8b43d4cff3	Mrs.	GA	2016-04-27 00:53:00	Grades PreK-2	A
4	33679 p137682	06f6e62e17de34fcf81020c77549e1d5	Mrs.	WA	2016-04-27 01:05:25	Grades 3-5	Litera

1.2 preprocessing of project_subject_categories

In [6]:

```
categories = list(project_data['project_subject_categories'].values)
# remove special characters from list of strings python:
```

<https://stackoverflow.com/a/47301924/4084039>

```
# https://www.geeksforgeeks.org/removing-stop-words-nltk-python/
# https://stackoverflow.com/questions/23669024/how-to-strip-a-specific-word-from-a-string
# https://stackoverflow.com/questions/8270092/remove-all-whitespace-in-a-string-in-python
cat_list = []
for i in categories:
    temp = ""
    # consider we have text like this "Math & Science, Warmth, Care & Hunger"
    for j in i.split(','): # it will split it in three parts ["Math & Science", "Warmth", "Care & Hunger"]
        if 'The' in j.split(): # this will split each of the category based on space "Math & Science" => "Math", "&", "Science"
            j = j.replace('The', '') # if we have the words "The" we are going to replace it with '' (i.e. removing 'The')
            j = j.replace(' ', '') # we are placing all the ' ' (space) with '' (empty) ex: "Math & Science" => "Math&Science"
            temp += j.strip() + " " # " abc ".strip() will return "abc", remove the trailing spaces
            temp = temp.replace('&', '_') # we are replacing the & value into
    cat_list.append(temp.strip())

project_data['clean_categories'] = cat_list
project_data.drop(['project_subject_categories'], axis=1, inplace=True)

from collections import Counter
my_counter = Counter()
for word in project_data['clean_categories'].values:
    my_counter.update(word.split())

cat_dict = dict(my_counter)
sorted_cat_dict = dict(sorted(cat_dict.items(), key=lambda kv: kv[1]))
```

1.3 preprocessing of project_subject_subcategories

In [7]:

```
sub_categories = list(project_data['project_subject_subcategories'].values)
# remove special characters from list of strings python:
https://stackoverflow.com/a/47301924/4084039

# https://www.geeksforgeeks.org/removing-stop-words-nltk-python/
# https://stackoverflow.com/questions/23669024/how-to-strip-a-specific-word-from-a-string
# https://stackoverflow.com/questions/8270092/remove-all-whitespace-in-a-string-in-python

sub_cat_list = []
for i in sub_categories:
    temp = ""
    # consider we have text like this "Math & Science, Warmth, Care & Hunger"
    for j in i.split(','): # it will split it in three parts ["Math & Science", "Warmth", "Care & Hunger"]
        if 'The' in j.split(): # this will split each of the category based on space "Math & Science" => "Math", "&", "Science"
            j = j.replace('The', '') # if we have the words "The" we are going to replace it with '' (i.e. removing 'The')
            j = j.replace(' ', '') # we are placing all the ' ' (space) with '' (empty) ex: "Math & Science" => "Math&Science"
            temp += j.strip() + " " # " abc ".strip() will return "abc", remove the trailing spaces
            temp = temp.replace('&', '_')
    sub_cat_list.append(temp.strip())

project_data['clean_subcategories'] = sub_cat_list
project_data.drop(['project_subject_subcategories'], axis=1, inplace=True)

# count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039
my_counter = Counter()
for word in project_data['clean_subcategories'].values:
    my_counter.update(word.split())

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

1.3 Text preprocessing

In [8]:

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

In [9]:

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

Out[9]:

Unnamed: 0	id	teacher_id	teacher_prefix	school_state	Date	project_grade_category	project_title
0	8393 p205479	2bf07ba08945e5d8b2a3f269b2b3cfe5	Mrs.	CA	2016-04-27 00:27:36	Grades PreK-2	Engineering STEAM into the Primary Classroom
1	37728 p043609	3f60494c61921b3b43ab61bdde2904df	Ms.	UT	2016-04-27 00:31:25	Grades 3-5	Sensory Tools for Focus

In [10]:

```
#### 1.4.2.3 Using Pretrained Models: TFIDF weighted W2V
```

In [11]:

```
# printing some random reviews
print(project_data['essay'].values[0])
print("="*50)
print(project_data['essay'].values[150])
print("="*50)
print(project_data['essay'].values[1000])
print("="*50)
print(project_data['essay'].values[20000])
print("="*50)
print(project_data['essay'].values[99999])
print("="*50)
```

I have been fortunate enough to use the Fairy Tale STEM kits in my classroom as well as the STEM journals, which my students really enjoyed. I would love to implement more of the Lakeshore STEM kits in my classroom for the next school year as they provide excellent and engaging STEM lessons. My students come from a variety of backgrounds, including language and socioeconomic status. Many of them don't have a lot of experience in science and engineering and these kits give me the materials to provide these exciting opportunities for my students. Each month I try to do several science or STEM/STEAM projects. I would use the kits and robot to help guide my science instruction in engaging and meaningful ways. I can adapt the kits to my current language arts pacing guide where we already teach some of the material in the kits like tall tales (Paul Bunyan) or Johnny Appleseed. The following units will be taught in the next school year where I will implement these kits: magnets, motion, sink vs. float, robots. I often get to these units and don't know if I am teaching the right way or using the right materials. The kits will give me additional ideas, strategies, and lessons to prepare my students in science. It is challenging to develop high quality science activities. These kits give me the materials I need to provide my students with science activities that will go along with the curriculum in my classroom. Although I have some things (like magnets) in my classroom, I don't know how to use them effectively. The kits will provide me with the right amount of materials and show me how to use them in an appropriate way.

I teach high school English to students with learning and behavioral disabilities. My students all vary in their ability level. However, the ultimate goal is to increase all students' literacy levels. This includes their reading, writing, and communication levels. I teach a really dynamic group of students. However, my students face a lot of challenges. My students all live in poverty and in a dangerous neighborhood. Despite these challenges, I have students who have the desire to def

a dangerous neighborhood. Despite these challenges, I have students who have the desire to defeat these challenges. My students all have learning disabilities and currently all are performing below grade level. My students are visual learners and will benefit from a classroom that fulfills their preferred learning style. The materials I am requesting will allow my students to be prepared for the classroom with the necessary supplies. Too often I am challenged with students who come to school unprepared for class due to economic challenges. I want my students to be able to focus on learning and not how they will be able to get school supplies. The supplies will last all year. Students will be able to complete written assignments and maintain a classroom journal. The chart paper will be used to make learning more visual in class and to create posters to aid students in their learning. The students have access to a classroom printer. The toner will be used to print student work that is completed on the classroom Chromebooks. I want to try and remove all barriers for the students learning and create opportunities for learning. One of the biggest barriers is the students not having the resources to get pens, paper, and folders. My students will be able to increase their literacy skills because of this project.

"Life moves pretty fast. If you don't stop and look around once in awhile, you could miss it." from the movie, Ferris Bueller's Day Off. Think back...what do you remember about your grandparents? How amazing would it be to be able to flip through a book to see a day in their lives? My second graders are voracious readers! They love to read both fiction and nonfiction books. Their favorite characters include Pete the Cat, Fly Guy, Piggie and Elephant, and Mercy Watson. They also love to read about insects, space and plants. My students are hungry bookworms! My students are eager to learn and read about the world around them. My kids love to be at school and are like little sponges absorbing everything around them. Their parents work long hours and usually do not see their children. My students are usually cared for by their grandparents or a family friend. Most of my students do not have someone who speaks English at home. Thus it is difficult for my students to acquire language. Now think forward... wouldn't it mean a lot to your kids, nieces or nephews or grandchildren, to be able to see a day in your life today 30 years from now? Memories are so precious to us and being able to share these memories with future generations will be a rewarding experience. As part of our social studies curriculum, students will be learning about changes over time. Students will be studying photos to learn about how their community has changed over time. In particular, we will look at photos to study how the land, buildings, clothing, and schools have changed over time. As a culminating activity, my students will capture a slice of their history and preserve it through scrap booking. Key important events in their young lives will be documented with the date, location, and names. Students will be using photos from home and from school to create their second grade memories. Their scrap books will preserve their unique stories for future generations to enjoy. Your donation to this project will provide my second graders with an opportunity to learn about social studies in a fun and creative manner. Through their scrapbooks, children will share their story with others and have a historical document for the rest of their lives.

"A person's a person, no matter how small." (Dr. Seuss) I teach the smallest students with the biggest enthusiasm for learning. My students learn in many different ways using all of our senses and multiple intelligences. I use a wide range of techniques to help all my students succeed. Students in my class come from a variety of different backgrounds which makes for wonderful sharing of experiences and cultures, including Native Americans. Our school is a caring community of successful learners which can be seen through collaborative student project based learning in and out of the classroom. Kindergarteners in my class love to work with hands-on materials and have many different opportunities to practice a skill before it is mastered. Having the social skills to work cooperatively with friends is a crucial aspect of the kindergarten curriculum. Montana is the perfect place to learn about agriculture and nutrition. My students love to role play in our pretend kitchen in the early childhood classroom. I have had several kids ask me, "Can we try cooking with REAL food?" I will take their idea and create "Common Core Cooking Lessons" where we learn important math and writing concepts while cooking delicious healthy food for snack time. My students will have a grounded appreciation for the work that went into making the food and knowledge of where the ingredients came from as well as how it's healthy for their bodies. This project would expand our learning of nutrition and agricultural cooking recipes by having us peel our own apples to make homemade applesauce, make our own bread, and mix up healthy plants from our classroom garden in the spring. We will also create our own cookbooks to be printed and shared with families. Students will gain math and literature skills as well as a life long enjoyment for healthy cooking.

My classroom consists of twenty-two amazing sixth graders from different cultures and backgrounds. They are a social bunch who enjoy working in partners and working with groups. They are hard-working and eager to head to middle school next year. My job is to get them ready to make this transition and make it as smooth as possible. In order to do this, my students need to come to school every day and feel safe and ready to learn. Because they are getting ready to head to middle school, I give them lots of choice- choice on where to sit and work, the order to complete assignments, choice of projects, etc. Part of the students feeling safe is the ability for them to come into a welcoming, encouraging environment. My room is colorful and the atmosphere is casual. I want them to take ownership of the classroom because we ALL share it together. Because my time with them is limited, I want to ensure they get the most of this time and enjoy it to the best of their abilities. Currently, we have twenty-two desks of differing sizes, yet the desks are similar to the ones the students will use in middle school. We also have a kidney table with crates for seating. I allow my students to choose their own spots while they are working independently or in groups. More often than not, most of them move out of their desks and onto the crates. Believe it or not, this has proven to be more successful than making them stay at their desks! It is because of this that I am looking toward the "Flexible Seating" option for my classroom. The students look forward to their work time so they can move around the room. I would like to get rid of the

look forward to their work time so they can move around the room. I would like to get rid of the constricting desks and move toward more "fun" seating options. I am requesting various seating so my students have more options to sit. Currently, I have a stool and a papasan chair I inherited from the previous sixth-grade teacher as well as five milk crate seats I made, but I would like to give them more options and reduce the competition for the "good seats". I am also requesting two rugs as not only more seating options but to make the classroom more welcoming and appealing. In order for my students to be able to write and complete work without desks, I am requesting a class set of clipboards. Finally, due to curriculum that requires groups to work together, I am requesting tables that we can fold up when we are not using them to leave more room for our flexible seating options.

I know that with more seating options, they will be that much more excited about coming to school! Thank you for your support in making my classroom one students will remember forever!

nannan

In [12]:

```
# https://stackoverflow.com/a/47091490/4084039
import re

def decontracted(phrase):
    # specific
    phrase = re.sub(r"won't", "will not", phrase)
    phrase = re.sub(r"can't", "can not", phrase)

    # general
    phrase = re.sub(r"n't", " not", phrase)
    phrase = re.sub(r"'re", " are", phrase)
    phrase = re.sub(r"'s", " is", phrase)
    phrase = re.sub(r"'d", " would", phrase)
    phrase = re.sub(r"'ll", " will", phrase)
    phrase = re.sub(r"'t", " not", phrase)
    phrase = re.sub(r"'ve", " have", phrase)
    phrase = re.sub(r"'m", " am", phrase)
    return phrase
```

In [13]:

```
sent = decontracted(project_data['essay'].values[20000])
print(sent)
print("="*50)
```

"A person is a person, no matter how small." (Dr.Seuss) I teach the smallest students with the biggest enthusiasm for learning. My students learn in many different ways using all of our senses and multiple intelligences. I use a wide range of techniques to help all my students succeed. Students in my class come from a variety of different backgrounds which makes for wonderful sharing of experiences and cultures, including Native Americans. Our school is a caring community of successful learners which can be seen through collaborative student project based learning in and out of the classroom. Kindergarteners in my class love to work with hands-on materials and have many different opportunities to practice a skill before it is mastered. Having the social skills to work cooperatively with friends is a crucial aspect of the kindergarten curriculum. Montana is the perfect place to learn about agriculture and nutrition. My students love to role play in our pretend kitchen in the early childhood classroom. I have had several kids ask me, "Can we try cooking with REAL food?" I will take their idea and create "Common Core Cooking Lessons" where we learn important math and writing concepts while cooking delicious healthy food for snack time. My students will have a grounded appreciation for the work that went into making the food and knowledge of where the ingredients came from as well as how it is healthy for their bodies. This project would expand our learning of nutrition and agricultural cooking recipes by having us peel our own apples to make homemade applesauce, make our own bread, and mix up healthy plants from our classroom garden in the spring. We will also create our own cookbooks to be printed and shared with families. Students will gain math and literature skills as well as a life long enjoyment for healthy cooking.

nannan

In [14]:

```
# \r \n \t remove from string python: http://texthandler.com/info/remove-line-breaks-python/
sent = sent.replace('\r', ' ')
sent = sent.replace('\n', ' ')
sent = sent.replace('\t', ' ')
print(sent)
```

A person is a person, no matter how small. (Dr.Seuss) I teach the smallest students with the biggest enthusiasm for learning. My students learn in many different ways using all of our senses and

multiple intelligences. I use a wide range of techniques to help all my students succeed. Students in my class come from a variety of different backgrounds which makes for wonderful sharing of experiences and cultures, including Native Americans. Our school is a caring community of successful learners which can be seen through collaborative student project based learning in and out of the classroom. Kindergarteners in my class love to work with hands-on materials and have many different opportunities to practice a skill before it is mastered. Having the social skills to work cooperatively with friends is a crucial aspect of the kindergarten curriculum. Montana is the perfect place to learn about agriculture and nutrition. My students love to role play in our pretend kitchen in the early childhood classroom. I have had several kids ask me, Can we try cooking with REAL food? I will take their idea and create Common Core Cooking Lessons where we learn important math and writing concepts while cooking delicious healthy food for snack time. My students will have a grounded appreciation for the work that went into making the food and knowledge of where the ingredients came from as well as how it is healthy for their bodies. This project would expand our learning of nutrition and agricultural cooking recipes by having us peel our own apples to make homemade applesauce, make our own bread, and mix up healthy plants from our classroom garden in the spring. We will also create our own cookbooks to be printed and shared with families. Students will gain math and literature skills as well as a life long enjoyment for healthy cooking.annan

In [15]:

```
#remove spacial character: https://stackoverflow.com/a/5843547/4084039
sent = re.sub('[^A-Za-z0-9]+', ' ', sent)
print(sent)
```

A person is a person no matter how small Dr Seuss I teach the smallest students with the biggest enthusiasm for learning My students learn in many different ways using all of our senses and multiple intelligences I use a wide range of techniques to help all my students succeed Students in my class come from a variety of different backgrounds which makes for wonderful sharing of experiences and cultures including Native Americans Our school is a caring community of successful learners which can be seen through collaborative student project based learning in and out of the classroom Kindergarteners in my class love to work with hands on materials and have many different opportunities to practice a skill before it is mastered Having the social skills to work cooperatively with friends is a crucial aspect of the kindergarten curriculum Montana is the perfect place to learn about agriculture and nutrition My students love to role play in our pretend kitchen in the early childhood classroom I have had several kids ask me Can we try cooking with REAL food I will take their idea and create Common Core Cooking Lessons where we learn important math and writing concepts while cooking delicious healthy food for snack time My students will have a grounded appreciation for the work that went into making the food and knowledge of where the ingredients came from as well as how it is healthy for their bodies This project would expand our learning of nutrition and agricultural cooking recipes by having us peel our own apples to make homemade applesauce make our own bread and mix up healthy plants from our classroom garden in the spring We will also create our own cookbooks to be printed and shared with families Students will gain math and literature skills as well as a life long enjoyment for healthy cooking.annan

In [16]:

```
# https://gist.github.com/sebleier/554280
# we are removing the words from the stop words list: 'no', 'nor', 'not'
stopwords= ['i', 'me', 'my', 'myself', 'we', 'our', 'ours', 'ourselves', 'you', "you're", "you've",
\
            "you'll", "you'd", 'your', 'yours', 'yourself', 'yourselves', 'he', 'him', 'his',
'himself', \
            'she', "she's", 'her', 'hers', 'herself', 'it', "it's", 'its', 'itself', 'they', 'them',
'their',\
            'theirs', 'themselves', 'what', 'which', 'who', 'whom', 'this', 'that', "that'll",
'these', 'those', \
            'am', 'is', 'are', 'was', 'were', 'be', 'been', 'being', 'have', 'has', 'had', 'having',
'do', 'does', \
            'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'or', 'because', 'as', 'until', '
while', 'of', \
            'at', 'by', 'for', 'with', 'about', 'against', 'between', 'into', 'through', 'during',
'before', 'after',\
            'above', 'below', 'to', 'from', 'up', 'down', 'in', 'out', 'on', 'off', 'over', 'under'
, 'again', 'further',\
            'then', 'once', 'here', 'there', 'when', 'where', 'why', 'how', 'all', 'any', 'both', 'e
ach', 'few', 'more',\
            'most', 'other', 'some', 'such', 'only', 'own', 'same', 'so', 'than', 'too', 'very', \
            's', 't', 'can', 'will', 'just', 'don', "don't", 'should', "should've", 'now', 'd', 'll'
, 'm', 'o', 're', \
            've', 'y', 'ain', 'aren', "aren't", 'couldn', "couldn't", 'didn', "didn't", 'doesn', "dc
esn't", 'hadn',\
            "hadn't", 'hasn', "hasn't", 'haven', "haven't", 'isn', "isn't", 'ma', 'mightn',
"mightn't", 'mustn', "mustn't", 'shan', "shan't", 'won', "won't", 'wouldn', "wouldn't"]
```

```
"mustn't", 'needn', "needn't", 'shan', "shan't", 'shouldn', "shouldn't", 'wasn',  
"wasn't", 'weren', "weren't", \  
'won', "won't", 'wouldn', "wouldn't"]
```

In [17]:

```
# Create function that will filter sentence  
def filterSentence(sentence):  
    sent = decontracted(sentence)  
    sent = sent.replace('\r', ' ')  
    sent = sent.replace('\n', ' ')  
    sent = sent.replace('\n', ' ')  
    sent = re.sub('[^A-Za-z0-9]+', ' ', sent)  
    sent = sent.lower()  
    # https://gist.github.com/sebleier/554280  
    sent = ' '.join(e for e in sent.split() if e not in stopwords)  
    return sent.strip()
```

In [18]:

```
# Combining all the above students  
from tqdm import tqdm  
preprocessed_essays = []  
# tqdm is for printing the status bar  
for sentence in tqdm(project_data['essay'].values):  
    preprocessed_essays.append(filterSentence(sentence))
```

```
100%|██████████| 109248/109248 [00:59<00:00, 1821.58it/s]
```

In [19]:

```
# after preprocessing  
preprocessed_essays[20000]
```

Out[19]:

```
'person person no matter small dr seuss teach smallest students biggest enthusiasm learning  
students learn many different ways using senses multiple intelligences use wide range techniques h  
elp students succeed students class come variety different backgrounds makes wonderful sharing exp  
periences cultures including native americans school caring community successful learners seen coll  
aborative student project based learning classroom kindergarteners class love work hands materials  
many different opportunities practice skill mastered social skills work cooperatively friends cruc  
ial aspect kindergarten curriculum montana perfect place learn agriculture nutrition students love  
role play pretend kitchen early childhood classroom several kids ask try cooking real food take id  
ea create common core cooking lessons learn important math writing concepts cooking delicious heal  
thy food snack time students grounded appreciation work went making food knowledge ingredients cam  
e well healthy bodies project would expand learning nutrition agricultural cooking recipes us peel  
apples make homemade applesauce make bread mix healthy plants classroom garden spring also create  
cookbooks printed shared families students gain math literature skills well life long enjoyment he  
althy cooking nannan'
```

1.4 Preprocessing of `project_title`

In [20]:

```
# similarly you can preprocess the titles also  
# Combining all the above students  
from tqdm import tqdm  
preprocessed_titles = []  
# tqdm is for printing the status bar  
for sentence in tqdm(project_data['project_title'].values):  
    preprocessed_titles.append(filterSentence(sentence))
```

```
100%|██████████| 109248/109248 [00:02<00:00, 42371.51it/s]
```

In [21]:

```
# after preprocessing
```

```
print(preprocessed_titles[20000])
```

health nutritional cooking kindergarten

In [22]:

```
# similarly you can preprocess the project_resource_summary also
# Combining all the above students
from tqdm import tqdm
preprocessed_resource_summary = []
# tqdm is for printing the status bar
for sentence in tqdm(project_data['project_resource_summary'].values):
    preprocessed_resource_summary.append(filterSentence(sentence))
```

```
100%|██████████| 109248/109248 [00:06<00:00, 16231.19it/s]
```

In [23]:

```
# after preprocessing
print(preprocessed_resource_summary[20000])
```

students need cooking supplies help us healthy learn nutrition mixer apple spiralizer kitchen
tools nutrition kit kid friendly healthy literature ink make cookbooks

In [24]:

```
# Preprocess teacher_prefix
from tqdm import tqdm
preprocessed_teacher_prefix = []
# tqdm is for printing the status bar
for teacher_prefix in tqdm(project_data['teacher_prefix'].values):
    teacher_prefix = str(teacher_prefix)
    clean_teacher_prefix = decontracted(teacher_prefix)
    clean_teacher_prefix = clean_teacher_prefix.replace('\\r', ' ')
    clean_teacher_prefix = clean_teacher_prefix.replace('\\\"', ' ')
    clean_teacher_prefix = clean_teacher_prefix.replace('\\n', ' ')
    clean_teacher_prefix = re.sub('[^A-Za-z0-9]+', ' ', clean_teacher_prefix)
    clean_teacher_prefix = clean_teacher_prefix.lower()
    if clean_teacher_prefix in stopwords:
        continue
    preprocessed_teacher_prefix.append(clean_teacher_prefix.strip())
```

```
100%|██████████| 109248/109248 [00:01<00:00, 61923.06it/s]
```

In [25]:

```
preprocessed_teacher_prefix[0:10]
```

Out[25]:

```
['mrs', 'ms', 'mrs', 'mrs', 'mrs', 'mrs', 'mrs', 'ms', 'ms', 'mrs']
```

In [26]:

```
# Preprocess project_grade_category
from tqdm import tqdm
preprocessed_project_grade_category = []
# tqdm is for printing the status bar
for project_grade_category in tqdm(project_data['project_grade_category'].values):
    project_grade_category = str(project_grade_category)
    clean_project_grade_category = decontracted(project_grade_category)
    clean_project_grade_category = clean_project_grade_category.replace('\\r', ' ')
    clean_project_grade_category = clean_project_grade_category.replace('\\\"', ' ')
    clean_project_grade_category = clean_project_grade_category.replace('\\n', ' ')
    clean_project_grade_category = re.sub('[^A-Za-z0-9]+', ' ', clean_project_grade_category)
    clean_project_grade_category = clean_project_grade_category.lower()
    if clean_project_grade_category in stopwords:
        continue
    clean_project_grade_category = clean_project_grade_category.strip()
# substitute any words in the problem because we are trusting the preprocessed feature
```

```
# whitespace are creating problems because we are treating this as categorical feature
preprocessed_project_grade_category.append(clean_project_grade_category.replace(' ', '_'))
```

```
100%|██████████| 109248/109248 [00:03<00:00, 34930.81it/s]
```

In [27]:

```
preprocessed_project_grade_category[0:10]
```

Out[27]:

```
['grades_prek_2',
 'grades_3_5',
 'grades_prek_2',
 'grades_prek_2',
 'grades_3_5',
 'grades_3_5',
 'grades_3_5',
 'grades_3_5',
 'grades_prek_2',
 'grades_3_5']
```

In [28]:

```
# Replace original columns with preprocessed column values
project_data['clean_essays'] = preprocessed_essays
project_data['clean_titles'] = preprocessed_titles
project_data['project_resource_summary'] = preprocessed_resource_summary
project_data['teacher_prefix'] = preprocessed_teacher_prefix
project_data['project_grade_category'] = preprocessed_project_grade_category
# Drop essays column
project_data.drop(['project_essay_1'], axis=1, inplace=True)
project_data.drop(['project_essay_2'], axis=1, inplace=True)
project_data.drop(['project_essay_3'], axis=1, inplace=True)
project_data.drop(['project_essay_4'], axis=1, inplace=True)
```

In [29]:

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

Out[29]:

Unnamed: 0	id	teacher_id	teacher_prefix	school_state	Date	project_grade_category	project_title
0	8393 p205479	2bf07ba08945e5d8b2a3f269b2b3cfe5	mrs	CA	2016-04-27 00:27:36	grades_prek_2	Engineering STEAM into the Primary Classroom
1	37728 p043609	3f60494c61921b3b43ab61bdde2904df	ms	UT	2016-04-27 00:31:25	grades_3_5	Sensory Tools for Focus
2	74477 p189804	4a97f3a390bfe21b99cf5e2b81981c73	mrs	CA	2016-04-27 00:46:53	grades_prek_2	Mobile Learning with a Mobile Listening Center
3	100660 p234804	cbc0e38f522143b86d372f8b43d4cff3	mrs	GA	2016-04-27 00:53:00	grades_prek_2	Flexible Seating for Flexible Learning
4	33679 p137682	06f6e62e17de34fcf81020c77549e1d5	mrs	WA	2016-04-27 01:05:25	grades_3_5	Going Deep: The Art of Inner Thinking!

In [30]:

```
project_data.tail(5)
```

Out[30]:

Unnamed: 0		id	teacher_id	teacher_prefix	school_state	Date	project_grade_category	project_title
109243	45036	p194916	29cf137e5a40b0f141d9fd7898303a5c	mrs	HI	2017-04-30 23:11:45	grades_9_12	Nature of Science
109244	12610	p162971	22fee80f2078c694c2d244d3ecb1c390	ms	NM	2017-04-30 23:23:24	grades_prek_2	Organic Chemistry
109245	179833	p096829	c8c81a73e29ae3bdd4140be8ad0bea00	mrs	IL	2017-04-30 23:25:42	grades_3_5	Biology Agriculture Sustainability
109246	13791	p184393	65545a295267ad9df99f26f25c978fd0	mrs	HI	2017-04-30 23:27:07	grades_9_12	Mathematics
109247	124250	p028318	1fff5a88945be8b2c728c6a85c31930f	mrs	CA	2017-04-30 23:45:08	grades_prek_2	Nature of Science

In [31]:

```
print(set(preprocessed_project_grade_category))
```

```
{'grades_prek_2', 'grades_3_5', 'grades_9_12', 'grades_6_8'}
```

In [32]:

```
project_data['teacher_prefix'] = project_data['teacher_prefix'].fillna('null')
```

In [33]:

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

Out[33]:

Unnamed: 0		id	teacher_id	teacher_prefix	school_state	Date	project_grade_category	project_title
0	8393	p205479	2bf07ba08945e5d8b2a3f269b2b3cfe5	mrs	CA	2016-04-27 00:27:36	grades_prek_2	Engineering STEAM into the Primary Classroom
1	37728	p043609	3f60494c61921b3b43ab61bdde2904df	ms	UT	2016-04-27 00:31:25	grades_3_5	Sensory Tools for Focus

Unnamed:

id

teacher id teacher prefix school state

Date project grade category project title

1.5 Preparing data for models

In [34]:

```
project_data.columns
```

Out[34]:

```
Index(['Unnamed: 0', 'id', 'teacher_id', 'teacher_prefix', 'school_state',  
      'Date', 'project_grade_category', 'project_title',  
      'project_resource_summary',  
      'teacher_number_of_previously_posted_projects', 'project_is_approved',  
      'quantity', 'price', 'clean_categories', 'clean_subcategories', 'essay',  
      'clean_essays', 'clean_titles'],  
      dtype='object')
```

we are going to consider

- school_state : categorical data
- clean_categories : categorical data
- clean_subcategories : categorical data
- project_grade_category : categorical data
- teacher_prefix : categorical data
- project_title : text data
- text : text data
- project_resource_summary: text data (optinal)
- quantity : numerical (optinal)
- teacher_number_of_previously_posted_projects : numerical
- price : numerical

In [35]:

```
print(project_data.shape)  
  
# I am taking 5% of data points for my analysis  
project_data = project_data.sample(frac=0.05)  
  
print(project_data.shape)
```

(109248, 18)

(5462, 18)

In [36]:

```
# Assigning data  
y = project_data['project_is_approved'].values  
project_data.drop(['project_is_approved'], axis=1, inplace=True)  
X = project_data  
project_data.shape
```

Out[36]:

(5462, 17)

In [37]:

```
X_train = X  
y_train = y
```

In [38]:

```
print('Train Data Set', X_train.shape, y_train.shape)
```

```
print('Train Data Set', X_train.shape, y_train.shape)
print('*'*100)
```

Train Data Set (5462, 17) (5462,)

1.5.1 Vectorizing Categorical data

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

In [39]:

```
# One hot encoding of Categorical Feature
# - school_state : categorical data
vectorizer = CountVectorizer()
vectorizer.fit(X_train['school_state'].values)# Fit has to happen only on train data

X_train_school_state_ohe = vectorizer.transform(X_train['school_state'].values)

school_state_features = vectorizer.get_feature_names()

print(X_train_school_state_ohe.shape, y_train.shape)
print(vectorizer.get_feature_names())
print('*'*100)
```

(5462, 51) (5462,)
['ak', 'al', 'ar', 'az', 'ca', 'co', 'ct', 'dc', 'de', 'fl', 'ga', 'hi', 'ia', 'id', 'il', 'in', 'k
s', 'ky', 'la', 'ma', 'md', 'me', 'mi', 'mn', 'mo', 'ms', 'mt', 'nc', 'nd', 'ne', 'nh', 'nj', 'nm',
'nv', 'ny', 'oh', 'ok', 'or', 'pa', 'ri', 'sc', 'sd', 'tn', 'tx', 'ut', 'va', 'vt', 'wa', 'wi', 'wv',
'wy']

In [40]:

```
# One hot encoding of Categorical Feature
# - clean_categories : categorical data
vectorizer = CountVectorizer()
vectorizer.fit(X_train['clean_categories'].values)# Fit has to happen only on train data

X_train_clean_categories_ohe = vectorizer.transform(X_train['clean_categories'].values)

clean_categories_features = vectorizer.get_feature_names()

print(X_train_clean_categories_ohe.shape, y_train.shape)
print(vectorizer.get_feature_names())
print('*'*100)
```

(5462, 9) (5462,)
['appliedlearning', 'care_hunger', 'health_sports', 'history_civics', 'literacy_language',
'math_science', 'music_arts', 'specialneeds', 'warmth']

In [41]:

```
# One hot encoding of Categorical Feature
# - clean_subcategories : categorical data
vectorizer = CountVectorizer()
vectorizer.fit(X_train['clean_subcategories'].values)# Fit has to happen only on train data

X_train_clean_subcategories_ohe = vectorizer.transform(X_train['clean_subcategories'].values)

clean_subcategories_features = vectorizer.get_feature_names()

print(X_train_clean_subcategories_ohe.shape, y_train.shape)
print(vectorizer.get_feature_names())
print('*'*100)
```

```
(5462, 30) (5462,)
['appliedsciences', 'care_hunger', 'charactereducation', 'civics_government',
'college_careerprep', 'communityservice', 'earlydevelopment', 'economics', 'environmentalscience',
'esl', 'extracurricular', 'financialliteracy', 'foreignlanguages', 'gym_fitness',
'health_lifescience', 'health_wellness', 'history_geography', 'literacy', 'literature_writing', 'mathematics', 'music', 'nutritioneducation', 'other', 'parentinvolvement', 'performingarts', 'socialsciences', 'specialneeds', 'teamsports', 'visualarts', 'warmth']
*****
```

In [42]:

```
print(X_train['project_grade_category'])
# One hot encoding of Categorical Feature
# - project_grade_category : categorical data
# Convert one hot encoding for project grade category
vectorizer = CountVectorizer()
vectorizer.fit(X_train['project_grade_category'].values)# Fit has to happen only on train data

X_train_project_grade_category_ohe = vectorizer.transform(X_train['project_grade_category'].values
)

project_grade_category_features = vectorizer.get_feature_names()

print(X_train_project_grade_category_ohe.shape, y_train.shape)
print(vectorizer.get_feature_names())
print('*'*100)
```

```
39389      grades_6_8
41647      grades_3_5
25354      grades_prek_2
21471      grades_3_5
85606      grades_prek_2
67572      grades_6_8
56900      grades_prek_2
92610      grades_9_12
41691      grades_prek_2
60640      grades_3_5
22504      grades_3_5
60983      grades_prek_2
22025      grades_3_5
50206      grades_6_8
75482      grades_3_5
73229      grades_prek_2
19887      grades_3_5
91931      grades_prek_2
106348     grades_3_5
50600      grades_9_12
37412      grades_3_5
37574      grades_6_8
77818      grades_3_5
28593      grades_prek_2
72789      grades_prek_2
7232       grades_3_5
89920      grades_3_5
3949       grades_6_8
67510      grades_3_5
18418      grades_prek_2
...
491        grades_prek_2
89590      grades_prek_2
558        grades_9_12
51732      grades_6_8
18269      grades_6_8
42951      grades_3_5
37507      grades_3_5
59057      grades_6_8
27183      grades_9_12
104750     grades_3_5
55559      grades_6_8
17246      grades_3_5
77082      grades_3_5
2665       grades_3_5
2612       grades_prek_2
98424      grades_3_5
```



```

50121      grades_3_5
79389      grades_prek_2
71555      grades_3_5
59055      grades_6_8
61752      grades_6_8
41854      grades_prek_2
86417      grades_prek_2
81927      grades_6_8
28667      grades_prek_2
70908      grades_9_12
72808      grades_prek_2
94460      grades_3_5
36383      grades_prek_2
103896     grades_6_8
6142      grades_prek_2
Name: project_grade_category, Length: 5462, dtype: object
(5462, 4) (5462,)
['grades_3_5', 'grades_6_8', 'grades_9_12', 'grades_prek_2']
*****

```



In [43]:

```
print(X_train_project_grade_category_ohe.toarray())
```

```

[[0 1 0 0]
 [1 0 0 0]
 [0 0 0 1]
 ...
 [0 0 0 1]
 [0 1 0 0]
 [0 0 0 1]]

```

In [44]:

```

# One hot encoding of Categorical Feature
# - teacher_prefix : categorical data
print(X_train['teacher_prefix'])
vectorizer = CountVectorizer()
vectorizer.fit(X_train['teacher_prefix'].values) # Fit has to happen only on train data

X_train_teacher_prefix_ohe = vectorizer.transform(X_train['teacher_prefix'].values)

teacher_prefix_features = vectorizer.get_feature_names()

print(X_train_teacher_prefix_ohe.shape, y_train.shape)
print(vectorizer.get_feature_names())
print('*'*100)

```

```

39389      mrs
41647      mrs
25354      ms
21471      mrs
85606      mrs
67572      ms
56900      ms
92610      mr
41691      mrs
60640      mrs
22504      mrs
60983      mrs
22025      mrs
50206      mr
75482      teacher
73229      mrs
19887      mrs
91931      mrs
106348     ms
50600      mrs
37412      mrs
37574      ms
77818      ms
28593      mrs
72789      ms
50000

```

```

7232      ms
89920     mrs
3949      mrs
67510     ms
18418     mrs
...
491       ms
89590     mrs
558       mrs
51732     mrs
18269     ms
42951     mrs
37507     mrs
59057     ms
27183     mrs
104750    ms
55559     mrs
17246     mrs
77082     mr
2665      ms
2612      mrs
98424     mr
79389     ms
71555     mrs
59055     mrs
61752     mrs
41854     mrs
86417     mrs
81927     mr
28667     mrs
70908     mr
72808     mrs
94460     mrs
36383     mrs
103896    mr
6142      ms
Name: teacher_prefix, Length: 5462, dtype: object
(5462, 4) (5462,)
['mr', 'mrs', 'ms', 'teacher']
*****

```

In [45]:

```
print(X_train_teacher_prefix_ohe.toarray())
```

```

[[0 1 0 0]
 [0 1 0 0]
 [0 0 1 0]
 ...
 [0 1 0 0]
 [1 0 0 0]
 [0 0 1 0]]

```

1.5.2 Vectorizing Text data

1.5.2.1 Bag of words

In [46]:

```

# - project_title : text data
print(X_train.shape, y_train.shape)

print("***100)

# We are considering only the words which appeared in at least 10 documents (rows or projects).
vectorizer = CountVectorizer(min_df=10, ngram_range=(1,4), max_features=5000)
vectorizer.fit(X_train['clean_titles'].values) # fit has to happen only on train data

# we use the fitted CountVectorizer to convert the text to vector
X_train_title_bow = vectorizer.transform(X_train['clean_titles'].values)

```

```
clean_titles_bow_features = vectorizer.get_feature_names()
```

```
print("After vectorizations")
print(X_train_title_bow.shape, y_train.shape)
# print(vectorizer.get_feature_names())
print("*"*100)
```

```
(5462, 17) (5462,)
*****
```

```
After vectorizations
(5462, 407) (5462,)
*****
```



In [47]:

```
# - text : text data
print(X_train.shape, y_train.shape)

print("*"*100)

# We are considering only the words which appeared in at least 10 documents(rows or projects).
vectorizer = CountVectorizer(min_df=10,ngram_range=(1,4), max_features=5000)
vectorizer.fit(X_train['clean_essays'].values) # fit has to happen only on train data

# we use the fitted CountVectorizer to convert the text to vector
X_train_essay_bow = vectorizer.transform(X_train['clean_essays'].values)

easy_bow_features = vectorizer.get_feature_names()

print("After vectorizations")
print(X_train_essay_bow.shape, y_train.shape)
# print(vectorizer.get_feature_names())
print("*"*100)
```

```
(5462, 17) (5462,)
*****
```

```
After vectorizations
(5462, 5000) (5462,)
*****
```



In [48]:

```
# - project_resource_summary: text data (optinal)
print(X_train.shape, y_train.shape)

print("*"*100)

# We are considering only the words which appeared in at least 10 documents(rows or projects).
vectorizer = CountVectorizer(min_df=10,ngram_range=(1,4), max_features=5000)
vectorizer.fit(X_train['project_resource_summary'].values) # fit has to happen only on train data

# we use the fitted CountVectorizer to convert the text to vector
X_train_project_resource_summary_bow = vectorizer.transform(X_train['project_resource_summary'].values)

project_resource_summary_bow_features = vectorizer.get_feature_names()

print("After vectorizations")
print(X_train_project_resource_summary_bow.shape, y_train.shape)
# print(vectorizer.get_feature_names())
print("*"*100)
```

```
(5462, 17) (5462,)
*****
```

```
After vectorizations
(5462, 1560) (5462,)
*****
```

1.5.2.2 TFIDF vectorizer

In [49]:

```
# - project_title : text data
print(X_train.shape, y_train.shape)

print("*"*100)

from sklearn.feature_extraction.text import TfidfVectorizer
# We are considering only the words which appeared in at least 10 documents(rows or projects).
vectorizer = TfidfVectorizer(min_df=10)
vectorizer.fit(X_train['clean_titles'].values) # fit has to happen only on train data

# we use the fitted CountVectorizer to convert the text to vector
X_train_title_tfidf = vectorizer.transform(X_train['clean_titles'].values)

clean_titles_tfidf_features = vectorizer.get_feature_names()

print("After vectorizations")
print(X_train_title_tfidf.shape, y_train.shape)
print("*"*100)
```

```
(5462, 17) (5462,)
*****

After vectorizations
(5462, 359) (5462,)
*****
```

In [50]:

```
# - text : text data
print(X_train.shape, y_train.shape)

print("*"*100)

from sklearn.feature_extraction.text import TfidfVectorizer
# We are considering only the words which appeared in at least 10 documents(rows or projects).
vectorizer = TfidfVectorizer(min_df=10)#, ngram_range=(2,2), max_features=5000
vectorizer.fit(X_train['clean_essays'].values) # fit has to happen only on train data

# we use the fitted CountVectorizer to convert the text to vector
X_train_essay_tfidf = vectorizer.transform(X_train['clean_essays'].values)

easy_tfidf_features = vectorizer.get_feature_names()

print("After vectorizations")
print(X_train_essay_tfidf.shape, y_train.shape)
print("*"*100)
```

```
(5462, 17) (5462,)
*****

After vectorizations
(5462, 4486) (5462,)
*****
```

In [51]:

```
# - project_resource_summary: text data (optinal)
print(X_train.shape, y_train.shape)

print("*"*100)

from sklearn.feature_extraction.text import TfidfVectorizer
# We are considering only the words which appeared in at least 10 documents(rows or projects).
vectorizer = TfidfVectorizer(min_df=10)
vectorizer.fit(X_train['project_resource_summary'].values) # fit has to happen only on train data
```

```

vectorizer.fit(X_train[ project_resource_summary ].values, # fit has to happen only on train data

# we use the fitted CountVectorizer to convert the text to vector
X_train_project_resource_summary_tfidf = vectorizer.transform(X_train['project_resource_summary'].
values)

project_resource_summary_tfidf_features = vectorizer.get_feature_names()

print("After vectorizations")
print(X_train_project_resource_summary_tfidf.shape, y_train.shape)
print("***100)

```

```

(5462, 17) (5462,)
*****

After vectorizations
(5462, 921) (5462,)
*****

```



1.5.2.3 Using Pretrained Models: Avg W2V

In [52]:

```

'''
# Reading glove vectors in python: https://stackoverflow.com/a/38230349/4084039
def loadGloveModel(gloveFile):
    print ("Loading Glove Model")
    f = open(gloveFile,'r', encoding="utf8")
    model = {}
    for line in tqdm(f):
        splitLine = line.split()
        word = splitLine[0]
        embedding = np.array([float(val) for val in splitLine[1:]])
        model[word] = embedding
    print ("Done.",len(model)," words loaded!")
    return model
model = loadGloveModel('glove.42B.300d.txt')

# =====
Output:

Loading Glove Model
1917495it [06:32, 4879.69it/s]
Done. 1917495 words loaded!

# =====

words = []
for i in preproced_texts:
    words.extend(i.split(' '))

for i in preproced_titles:
    words.extend(i.split(' '))
print("all the words in the coupus", len(words))
words = set(words)
print("the unique words in the coupus", len(words))

inter_words = set(model.keys()).intersection(words)
print("The number of words that are present in both glove vectors and our coupus", \
      len(inter_words), "(", np.round(len(inter_words)/len(words)*100, 3), "%) ")

words_courpus = {}
words_glove = set(model.keys())
for i in words:
    if i in words_glove:
        words_courpus[i] = model[i]
print("word 2 vec length", len(words_courpus))

# stronging variables into pickle files python: http://www.jessicayung.com/how-to-use-pickle-to-sa
ve-and-load-variables-in-python/

import pickle
with open('glove_vectors'. 'wb') as f:

```

```

with open('glove_vectors', 'wb') as f:
    pickle.dump(words_courpus, f)

```

```
'''
```

Out[52]:

```

'\n# Reading glove vectors in python: https://stackoverflow.com/a/38230349/4084039\ndef
loadGloveModel(gloveFile):\n    print ("Loading Glove Model")\n    f = open(gloveFile,\nencoding="utf8")\n    model = {}\n    for line in tqdm(f):\n        splitLine = line.split()\nword = splitLine[0]\n        embedding = np.array([float(val) for val in splitLine[1:]])\n        model[word] = embedding\n    print ("Done.",len(model)," words loaded!")\n    return model\nmodel =
loadGloveModel('glove.42B.300d.txt')\n\n# =====\nOutput:\n\nLoading G
love Model\n1917495it [06:32, 4879.69it/s]\nDone. 1917495 words loaded!\n\n#
=====
\n\nwords = []\nfor i in preproced_texts:\n    words.extend(i.split('\n\n#
\n\nwords_glove = {}
\nwords_courpus = {}
\nwords_glove =
set(model.keys())\nfor i in words:\n    if i in words_glove:\n        words_courpus[i] = model[i]\n
print("word 2 vec length", len(words_courpus))\n\n\n# stronging variables into pickle files python
: http://www.jessicayung.com/how-to-use-pickle-to-save-and-load-variables-in-python/\n\nimport pic
kle\nwith open('glove_vectors', 'wb') as f:\n    pickle.dump(words_courpus, f)\n\n\n'

```

In [53]:

```

# stronging variables into pickle files python: http://www.jessicayung.com/how-to-use-pickle-to-sa
ve-and-load-variables-in-python/
# make sure you have the glove_vectors file
with open('glove_vectors', 'rb') as f:
    model = pickle.load(f)
    glove_words = set(model.keys())

```

In [54]:

```

# average Word2Vec for train text
# compute average word2vec for each review.
avg_w2v_vectors_text_train = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_train['clean_essays'].values): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    cnt_words = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if word in glove_words:
            vector += model[word]
            cnt_words += 1
    if cnt_words != 0:
        vector /= cnt_words
    avg_w2v_vectors_text_train.append(vector)

print(len(avg_w2v_vectors_text_train))
print(len(avg_w2v_vectors_text_train[0]))

```

```
100%|██████████| 5462/5462 [00:01<00:00, 3047.29it/s]
```

5462

300

In [55]:

```

# Similarly you can vectorize for title also
# average Word2Vec
# compute average word2vec for each review.
avg_w2v_vectors_title_train = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_train['clean_titles']): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    cnt_words = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if word in glove_words:
            vector += model[word]
            cnt_words += 1

```

```

        cnt_words += 1
    if cnt_words != 0:
        vector /= cnt_words
    avg_w2v_vectors_title_train.append(vector)

print(len(avg_w2v_vectors_title_train))
print(len(avg_w2v_vectors_title_train[0]))

```

```
100%|██████████| 5462/5462 [00:00<00:00, 72933.74it/s]
```

```
5462
300
```

In [56]:

```

# Similarly you can vectorize for project_resource_summary also
# average Word2Vec
# compute average word2vec for each review.
avg_w2v_vectors_project_resource_summary_train = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_train['project_resource_summary']): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    cnt_words = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if word in glove_words:
            vector += model[word]
            cnt_words += 1
    if cnt_words != 0:
        vector /= cnt_words
    avg_w2v_vectors_project_resource_summary_train.append(vector)

print(len(avg_w2v_vectors_project_resource_summary_train))
print(len(avg_w2v_vectors_project_resource_summary_train[0]))

```

```
100%|██████████| 5462/5462 [00:00<00:00, 32580.38it/s]
```

```
5462
300
```

1.5.2.3 Using Pretrained Models: TFIDF weighted W2V

In [57]:

```

# S = ["abc def pqr", "def def def abc", "pqr pqr def"]
tfidf_model = TfidfVectorizer()
tfidf_model.fit(X_train['clean_essays'])
# we are converting a dictionary with word as a key, and the idf as a value
dictionary = dict(zip(tfidf_model.get_feature_names(), list(tfidf_model.idf_)))
tfidf_words = set(tfidf_model.get_feature_names())

```

In [58]:

```

# average Word2Vec
# compute average word2vec for each review.
tfidf_w2v_vectors_text_train = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_train['clean_essays']): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if (word in glove_words) and (word in tfidf_words):
            vec = model[word] # getting the vector for each word
            # here we are multiplying idf value(dictionary[word]) and the tf
            value((sentence.count(word)/len(sentence.split())))
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tf
            idf value for each word
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf

```

```

    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    tfidf_w2v_vectors_text_train.append(vector)

print(len(tfidf_w2v_vectors_text_train))
print(len(tfidf_w2v_vectors_text_train[0]))

```

100%|██████████| 5462/5462 [00:12<00:00, 436.74it/s]

5462
300

In [59]:

```

# S = ["abc def pqr", "def def def abc", "pqr pqr def"]
tfidf_model = TfidfVectorizer()
tfidf_model.fit(X_train['clean_titles'])
# we are converting a dictionary with word as a key, and the idf as a value
dictionary = dict(zip(tfidf_model.get_feature_names(), list(tfidf_model.idf_)))
tfidf_words = set(tfidf_model.get_feature_names())

```

In [60]:

```

# Similarly you can vectorize for title also
# average Word2Vec
# compute average word2vec for each review.
tfidf_w2v_vectors_title_train = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_train['clean_titles']): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if (word in glove_words) and (word in tfidf_words):
            vec = model[word] # getting the vector for each word
            # here we are multiplying idf value(dictionary[word]) and the tf
            value((sentence.count(word)/len(sentence.split())))
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tf
            idf value for each word
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    tfidf_w2v_vectors_title_train.append(vector)

print(len(tfidf_w2v_vectors_title_train))
print(len(tfidf_w2v_vectors_title_train[0]))

```

100%|██████████| 5462/5462 [00:00<00:00, 27854.56it/s]

5462
300

In [61]:

```

# S = ["abc def pqr", "def def def abc", "pqr pqr def"]
tfidf_model = TfidfVectorizer()
tfidf_model.fit(X_train['project_resource_summary'])
# we are converting a dictionary with word as a key, and the idf as a value
dictionary = dict(zip(tfidf_model.get_feature_names(), list(tfidf_model.idf_)))
tfidf_words = set(tfidf_model.get_feature_names())

```

In [62]:

```

# Similarly you can vectorize for title also
# average Word2Vec
# compute average word2vec for each review.
tfidf_w2v_vectors_project_resource_summary_train = []; # the avg-w2v for each sentence/review is s
tored in this list

```



```

for sentence in tqdm(x_train['project_resource_summary']): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if (word in glove_words) and (word in tfidf_words):
            vec = model[word] # getting the vector for each word
            # here we are multiplying idf value(dictionary[word]) and the tf
            value((sentence.count(word)/len(sentence.split())))
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tf
            idf value for each word
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    tfidf_w2v_vectors_project_resource_summary_train.append(vector)

print(len(tfidf_w2v_vectors_project_resource_summary_train))
print(len(tfidf_w2v_vectors_project_resource_summary_train[0]))

```

100%|██████████| 5462/5462 [00:00<00:00, 8989.50it/s]

5462
300

1.5.3 Vectorizing Numerical features

In [63]:

```

# You no need to perform standardization/normalization on numerical data,
# because you will classify data by using gini impurity in decision tree classifier.
# - quantity : numerical (optional)

```

```
X_train_quantity_norm = X_train['quantity'].values.reshape(-1,1)
```

```

print("After vectorizations")
print(X_train_quantity_norm.shape, y_train.shape)
print("=="*100)

```

After vectorizations
(5462, 1) (5462,)



In [64]:

```

# You no need to perform standardization/normalization on numerical data,
# because you will classify data by using gini impurity in decision tree classifier.
# One hot encoding of numerical feature
# - teacher_number_of_previously_posted_projects : numerical

```

```
X_train_teacher_number_of_previously_posted_projects_norm =
X_train['teacher_number_of_previously_posted_projects'].values.reshape(-1,1)
```

```

print("After vectorizations")
print(X_train_teacher_number_of_previously_posted_projects_norm.shape, y_train.shape)
print("=="*100)

```

After vectorizations
(5462, 1) (5462,)



In [65]:

```

# You no need to perform standardization/normalization on numerical data,
# because you will classify data by using gini impurity in decision tree classifier.
# - price : numerical

```

```
X_train_price_norm = X_train['price'].values.reshape(-1,1)
```

```
print("After vectorizations")
print(X_train_price_norm.shape, y_train.shape)
print("=="*100)
```

```
After vectorizations
(5462, 1) (5462,)
```



1.5.4 Merging all the above features

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

In [66]:

```
# print(categories_one_hot.shape)
# print(sub_categories_one_hot.shape)
# print(text_bow.shape)
# print(price_standardized.shape)
print('Categorical Features')
print('=='*100)
print(X_train_school_state_ohe.shape, y_train.shape)
print('=='*100)
print(X_train_clean_categories_ohe.shape, y_train.shape)
print('=='*100)
print(X_train_clean_subcategories_ohe.shape, y_train.shape)
print('=='*100)
print(X_train_project_grade_category_ohe.shape, y_train.shape)
print('=='*100)
print(X_train_teacher_prefix_ohe.shape, y_train.shape)
print('=='*100)
print('Text Encoding Features')
print('=='*100)
print(X_train_title_bow.shape, y_train.shape)
print('=='*100)
print(X_train_essay_bow.shape, y_train.shape)
print('=='*100)
print(X_train_project_resource_summary_bow.shape, y_train.shape)
print('=='*100)
print(X_train_title_tfidf.shape, y_train.shape)
print('=='*100)
print(X_train_essay_tfidf.shape, y_train.shape)
print('=='*100)
print(X_train_project_resource_summary_tfidf.shape, y_train.shape)
print('=='*100)
print(len(avg_w2v_vectors_text_train))
print(len(avg_w2v_vectors_text_train[0]))
print('=='*100)
print(len(avg_w2v_vectors_title_train))
print(len(avg_w2v_vectors_title_train[0]))
print('=='*100)
print(len(avg_w2v_vectors_project_resource_summary_train))
print(len(avg_w2v_vectors_project_resource_summary_train[0]))
print('=='*100)
print(len(tfidf_w2v_vectors_text_train))
print(len(tfidf_w2v_vectors_text_train[0]))
print('=='*100)
print(len(tfidf_w2v_vectors_title_train))
print(len(tfidf_w2v_vectors_title_train[0]))
print('=='*100)
print(len(tfidf_w2v_vectors_project_resource_summary_train))
print(len(tfidf_w2v_vectors_project_resource_summary_train[0]))
print('=='*100)
print('Numerical Features')
print('=='*100)
print(X_train_quantity_norm.shape, y_train.shape)
print('=='*100)
print(X_train_teacher_number_of_previously_posted_projects_norm.shape, y_train.shape)
print('=='*100)
print(X_train_price_norm.shape, y_train.shape)
```

Categorical Features

(5462, 51) (5462,)

(5462, 9) (5462,)

(5462, 30) (5462,)

(5462, 4) (5462,)

(5462, 4) (5462,)

Text Encoding Features

(5462, 407) (5462,)

(5462, 5000) (5462,)

(5462, 1560) (5462,)

(5462, 359) (5462,)

(5462, 4486) (5462,)

(5462, 921) (5462,)

5462
300

5462
300

5462
300

5462
300

5462
300

5462
300

Numerical Features

(5462, 1) (5462,)

(5462, 1) (5462,)

(5462, 1) (5462,)



```

# merge the sparse matrices together, each column is a row of 1s and 0s
from scipy.sparse import hstack
# with the same hstack function we are concatenating a sparse matrix and a dense matrix :)
# X = hstack((categories_one_hot, sub_categories_one_hot, text_bow, price_standardized))
# X.shape

X_train_real = X_train

X_train = hstack((X_train_school_state_ohe, X_train_clean_categories_ohe,
X_train_clean_subcategories_ohe, X_train_project_grade_category_ohe, X_train_teacher_prefix_ohe, X
_train_title_bow, X_train_essay_bow, X_train_project_resource_summary_bow, X_train_title_tfidf,
X_train_essay_tfidf, X_train_project_resource_summary_tfidf, avg_w2v_vectors_text_train,
avg_w2v_vectors_title_train, avg_w2v_vectors_project_resource_summary_train,
tfidf_w2v_vectors_text_train, tfidf_w2v_vectors_title_train,
tfidf_w2v_vectors_project_resource_summary_train, X_train_quantity_norm,
X_train_teacher_number_of_previously_posted_projects_norm, X_train_price_norm)).tocsr()

print(X_train_real.shape)
print(X_train.shape)

```

```

(5462, 17)
(5462, 14634)

```

Computing Sentiment Scores

In [68]:

```

import nltk
from nltk.sentiment.vader import SentimentIntensityAnalyzer

# import nltk
# nltk.download('vader_lexicon')

sid = SentimentIntensityAnalyzer()

for_sentiment = 'a person is a person no matter how small dr seuss i teach the smallest students w
ith the biggest enthusiasm \
for learning my students learn in many different ways using all of our senses and multiple intelli
gences i use a wide range\
of techniques to help all my students succeed students in my class come from a variety of differen
t backgrounds which makes\
for wonderful sharing of experiences and cultures including native americans our school is a carin
g community of successful \
learners which can be seen through collaborative student project based learning in and out of the
classroom kindergarteners \
in my class love to work with hands on materials and have many different opportunities to practice
a skill before it is\
mastered having the social skills to work cooperatively with friends is a crucial aspect of the ki
ndergarten curriculum\
montana is the perfect place to learn about agriculture and nutrition my students love to role pla
y in our pretend kitchen\
in the early childhood classroom i have had several kids ask me can we try cooking with real food
i will take their idea \
and create common core cooking lessons where we learn important math and writing concepts while co
oking delicious healthy \
food for snack time my students will have a grounded appreciation for the work that went into maki
ng the food and knowledge \
of where the ingredients came from as well as how it is healthy for their bodies this project woul
d expand our learning of \
nutrition and agricultural cooking recipes by having us peel our own apples to make homemade apple
sauce make our own bread \
and mix up healthy plants from our classroom garden in the spring we will also create our own cook
books to be printed and \
shared with families students will gain math and literature skills as well as a life long enjoymen
t for healthy cooking \
nannan'
ss = sid.polarity_scores(for_sentiment)

for k in ss:
    print('{0}: {1}, '.format(k, ss[k]), end='')

# we can use these 4 things as features/attributes (neg, neu, pos, compound)
# neg: 0.0, neu: 0.753, pos: 0.247, compound: 0.93

```

```

neg: 0.01. neu: 0.745. pos: 0.245. compound: 0.9975.

```

neg: 0.01, neu: 0.10, pos: 0.20, compound: 0.69,

Assignment 10: Clustering

- **step 1:** Choose any vectorizer (data matrix) that you have worked in any of the assignments, and got the best AUC value.
- **step 2:** Choose any of the [feature selection/reduction algorithms](#) ex: selectkbest features, pretrained word vectors, model based feature selection etc and reduce the number of features to 5k features
- **step 3:** Apply all three kmeans, Agglomerative clustering, DBSCAN
 - **K-Means Clustering:**
 - Find the best 'k' using the elbow-knee method (plot k vs inertia_)
 - **Agglomerative Clustering:**
 - Apply [agglomerative algorithm](#) and try a different number of clusters like 2,5 etc.
 - You can take less data points (as this is very computationally expensive one) to perform hierarchical clustering because they do take a considerable amount of time to run.
 - **DBSCAN Clustering:**
 - Find the best 'eps' using the [elbow-knee method](#).
 - You can take a smaller sample size for this as well.
- **step 4:** Summarize each cluster by manually observing few points from each cluster.
- **step 5:** You need to plot the word cloud with essay text for each cluster for each of algorithms mentioned in **step 3**.

2. Clustering

2.1 Choose the best data matrix on which you got the best AUC

In [69]:

```
%%time
# I would prefer Navie Bayes BOW. It gives me 0.70 AUC

# Please write all the code with proper documentation

# Prepare data for BOW
X_train_bow = hstack((X_train_school_state_one, X_train_clean_categories_one,
X_train_clean_subcategories_one, X_train_project_grade_category_one, X_train_teacher_prefix_one, X
_train_title_bow, X_train_essay_bow, X_train_project_resource_summary_bow, X_train_quantity_norm,
X_train_teacher_number_of_previously_posted_projects_norm, X_train_price_norm)).tocsr()

print(X_train_bow.shape, y_train.shape)

(5462, 7068) (5462,)
CPU times: user 53.2 ms, sys: 34.3 ms, total: 87.4 ms
Wall time: 66.5 ms
```

Note: I already completed steps 2.2 & 2.3 previously, So I didn't copy code in below cells.

2.2 Make Data Model Ready: encoding numerical, categorical features

In [70]:

```
# I already computed in above cells.
```

2.3 Make Data Model Ready: encoding eassay, and project_title

In [71]:

```
# I already computed in above cells.
```

2.4 Dimensionality Reduction on the selected features

In [72]:

```
# Train a multinomial naive bayes to get important features
from sklearn.naive_bayes import MultinomialNB

nb = MultinomialNB(alpha=1, class_prior=[0.5,0.5])
nb.fit(X_train_bow, y_train)
```

Out[72]:

```
MultinomialNB(alpha=1, class_prior=[0.5, 0.5], fit_prior=True)
```

In [73]:

```
# Select 5K important features
from sklearn.feature_selection import SelectFromModel

model = SelectFromModel(nb, prefit=True, max_features=5000)
X_train_bow_5k = model.transform(X_train_bow)
X_train_bow_5k.shape
```

Out[73]:

```
(5462, 4107)
```

In [74]:

```
type(X_train_bow_5k)

X_train_bow_5k
```

Out[74]:

```
<5462x4107 sparse matrix of type '<class 'numpy.float64'>'
  with 127171 stored elements in Compressed Sparse Row format>
```

2.5 Apply Kmeans

In [75]:

```
from sklearn.cluster import KMeans

# define clusters
clusters = [2, 5, 10, 15, 20]
inertia = []

# apply kmeans and collect inertia
for i in tqdm(clusters):
    Kmean = KMeans(n_clusters=i)
    Kmean.fit(X_train_bow_5k)

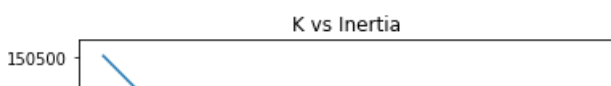
    inertia.append(Kmean.inertia_)
```

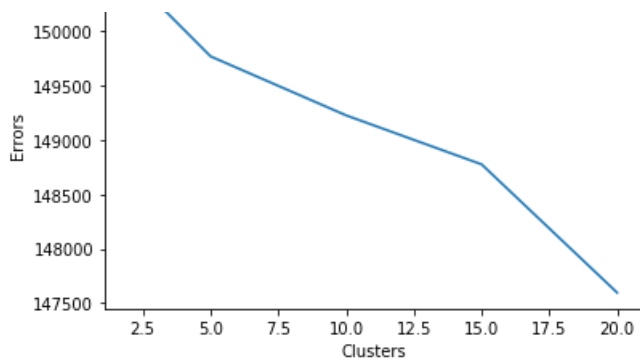
100%|██████████| 5/5 [01:19<00:00, 15.81s/it]

In [76]:

```
# plot k vs inertia_

plt.plot(clusters, inertia)
plt.xlabel('Clusters')
plt.ylabel('Errors')
plt.title('K vs Inertia')
plt.show()
```





In [77]:

```
Kmean = KMeans(n_clusters=15)
Kmean.fit(X_train_bow_5k)
```

Out[77]:

```
KMeans(algorithm='auto', copy_x=True, init='k-means++', max_iter=300,
        n_clusters=15, n_init=10, n_jobs=None, precompute_distances='auto',
        random_state=None, tol=0.0001, verbose=0)
```

In [78]:

```
Kmean.n_clusters
```

Out[78]:

```
15
```

In [79]:

```
Kmean.labels_
```

Out[79]:

```
array([0, 0, 0, ..., 0, 0, 0], dtype=int32)
```

In [80]:

```
# Collect data points from each cluster
clusters_set = {i: np.where(Kmean.labels_ == i)[0] for i in range(Kmean.n_clusters)}

clusters_set
```

Out[80]:

```
{0: array([ 0,  1,  2, ..., 5459, 5460, 5461]),
 1: array([377]),
 2: array([ 3, 45, 51, 66, 79, 135, 136, 192, 197, 225, 238,
          247, 254, 285, 290, 311, 317, 319, 362, 401, 430, 446,
          466, 500, 501, 539, 592, 613, 679, 683, 688, 699, 712,
          718, 731, 740, 743, 749, 754, 781, 807, 818, 834, 838,
          855, 857, 859, 875, 879, 980, 1087, 1120, 1130, 1134, 1143,
          1194, 1199, 1252, 1257, 1263, 1267, 1282, 1312, 1339, 1355, 1376,
          1386, 1407, 1423, 1462, 1465, 1501, 1512, 1540, 1580, 1597, 1598,
          1608, 1619, 1622, 1642, 1651, 1668, 1680, 1684, 1717, 1733, 1744,
          1800, 1812, 1848, 1856, 1870, 1891, 1912, 1914, 1947, 1948, 1952,
          1955, 1982, 2011, 2087, 2088, 2124, 2143, 2180, 2195, 2207, 2315,
          2323, 2328, 2340, 2345, 2378, 2391, 2402, 2407, 2440, 2446, 2453,
          2493, 2498, 2514, 2534, 2578, 2581, 2594, 2600, 2605, 2732, 2735,
          2736, 2763, 2787, 2789, 2815, 2847, 2866, 2871, 2887, 2896, 2897,
          2921, 2925, 2934, 2949, 2959, 2994, 2996, 3053, 3059, 3072, 3081,
          3082, 3099, 3109, 3147, 3160, 3169, 3173, 3201, 3235, 3257, 3285,
          3343, 3350, 3476, 3515, 3554, 3558, 3580, 3642, 3670, 3680, 3717,
          3732, 3733, 3746, 3748, 3764, 3768, 3772, 3778, 3797, 3799, 3866,
          3868, 3881, 3886, 3904, 3912, 3918, 3936, 3943, 3957, 3984, 3987,
          3999, 4036, 4039, 4052, 4062, 4069, 4108, 4131, 4132, 4133, 4143,
          4151, 4161, 4167, 4172, 4201, 4213, 4226, 4235, 4253, 4288, 4311,
```

```
4334, 4344, 4347, 4349, 4404, 4411, 4417, 4442, 4455, 4461, 4505,
4529, 4532, 4540, 4551, 4568, 4578, 4590, 4612, 4616, 4663, 4666,
4674, 4687, 4714, 4732, 4744, 4751, 4789, 4808, 4816, 4838, 4878,
4879, 4890, 4893, 4902, 4903, 4912, 4967, 5013, 5029, 5032, 5046,
5053, 5062, 5082, 5107, 5110, 5122, 5128, 5156, 5158, 5220, 5253,
5262, 5271, 5303, 5332, 5362, 5366, 5372, 5387, 5394, 5401, 5441,
5454, 5455]],
3: array([ 46, 245, 253, 1669, 2064, 3068, 3378, 4450, 4733, 5238]),
4: array([1817, 2184, 2255, 2635, 4365, 4587, 4717, 5170]),
5: array([ 28, 69, 70, 145, 219, 289, 379, 447, 523, 537, 538,
560, 566, 567, 580, 717, 732, 767, 888, 986, 1071, 1183,
1190, 1214, 1226, 1235, 1258, 1275, 1311, 1331, 1372, 1648, 1808,
2024, 2061, 2139, 2150, 2160, 2199, 2257, 2275, 2335, 2343, 2363,
2431, 2518, 2571, 2672, 2676, 2689, 2718, 2771, 2773, 2798, 2950,
2954, 3220, 3437, 3495, 3597, 3743, 3826, 3981, 4054, 4092, 4159,
4272, 4312, 4336, 4388, 4401, 4431, 4515, 4561, 4676, 4868, 4885,
4991, 5075, 5115, 5208, 5211, 5234, 5239, 5244, 5252, 5268, 5281,
5343, 5354, 5377, 5390, 5456]),
6: array([5337]),
7: array([3222]),
8: array([ 460, 2739]),
9: array([ 261, 517, 1506, 2646, 2751, 3304, 3885]),
10: array([1242]),
11: array([ 806, 927, 1028, 1739, 1843, 2189, 2592, 2627, 3089, 3844, 5236]),
12: array([ 17, 186, 257, 260, 414, 547, 816, 851, 861, 894, 983,
1052, 1109, 1139, 1210, 2073, 2362, 2375, 2392, 2488, 2495, 2781,
2930, 2931, 3055, 3066, 3251, 3657, 3720, 3839, 3842, 3861, 3976,
4035, 4114, 4142, 4413, 4437, 4647, 4745, 4785, 4842, 4869, 4939,
5057, 5081, 5099, 5412, 5438, 5458]),
13: array([ 9, 12, 16, 23, 33, 44, 58, 72, 84, 87, 90,
92, 102, 112, 114, 119, 121, 123, 125, 126, 134, 142,
143, 162, 170, 183, 207, 216, 232, 235, 248, 258, 259,
266, 267, 280, 297, 298, 303, 316, 324, 328, 335, 355,
364, 367, 392, 394, 400, 407, 408, 423, 458, 475, 485,
487, 520, 521, 526, 527, 544, 559, 576, 589, 599, 609,
620, 632, 637, 639, 670, 680, 685, 694, 697, 702, 720,
724, 728, 737, 744, 746, 787, 790, 795, 798, 821, 832,
833, 844, 864, 877, 880, 881, 897, 901, 909, 913, 918,
920, 921, 924, 932, 961, 964, 969, 977, 1011, 1015, 1019,
1025, 1047, 1048, 1065, 1089, 1095, 1097, 1102, 1113, 1117, 1128,
1133, 1135, 1136, 1160, 1165, 1166, 1177, 1179, 1202, 1219, 1222,
1232, 1240, 1259, 1288, 1291, 1315, 1320, 1330, 1366, 1378, 1384,
1387, 1404, 1406, 1414, 1425, 1426, 1439, 1457, 1463, 1464, 1480,
1508, 1517, 1539, 1545, 1566, 1567, 1573, 1574, 1578, 1581, 1583,
1604, 1613, 1632, 1639, 1663, 1671, 1673, 1688, 1693, 1703, 1732,
1776, 1786, 1794, 1796, 1797, 1805, 1821, 1825, 1852, 1867, 1868,
1877, 1895, 1896, 1916, 1965, 1976, 1994, 1998, 2004, 2007, 2008,
2028, 2033, 2040, 2050, 2051, 2055, 2062, 2063, 2072, 2079, 2085,
2090, 2098, 2115, 2118, 2123, 2164, 2168, 2185, 2191, 2198, 2219,
2236, 2239, 2250, 2277, 2296, 2302, 2313, 2316, 2325, 2336, 2365,
2366, 2370, 2389, 2406, 2416, 2421, 2434, 2442, 2454, 2458, 2465,
2481, 2497, 2508, 2531, 2538, 2549, 2560, 2582, 2585, 2587, 2597,
2610, 2613, 2632, 2637, 2643, 2652, 2655, 2657, 2658, 2691, 2695,
2706, 2715, 2728, 2729, 2755, 2761, 2807, 2821, 2827, 2843, 2844,
2849, 2850, 2861, 2886, 2898, 2936, 2942, 2961, 2969, 2971, 2976,
2983, 2991, 3012, 3035, 3036, 3040, 3042, 3050, 3058, 3074, 3078,
3084, 3095, 3107, 3125, 3134, 3140, 3143, 3154, 3172, 3175, 3181,
3186, 3191, 3218, 3234, 3238, 3259, 3269, 3271, 3274, 3278, 3281,
3288, 3301, 3310, 3311, 3319, 3321, 3327, 3339, 3349, 3357, 3363,
3364, 3373, 3376, 3379, 3380, 3386, 3387, 3389, 3394, 3414, 3419,
3421, 3436, 3444, 3446, 3461, 3465, 3487, 3489, 3526, 3537, 3565,
3571, 3583, 3589, 3592, 3610, 3616, 3640, 3647, 3654, 3663, 3686,
3690, 3693, 3710, 3714, 3724, 3730, 3735, 3737, 3739, 3742, 3769,
3793, 3804, 3814, 3819, 3827, 3863, 3864, 3865, 3873, 3883, 3884,
3901, 3907, 3910, 3913, 3931, 3935, 3947, 3962, 3970, 3993, 4022,
4027, 4040, 4047, 4058, 4060, 4067, 4086, 4093, 4098, 4100, 4105,
4118, 4125, 4140, 4147, 4150, 4160, 4165, 4166, 4178, 4182, 4190,
4205, 4206, 4215, 4223, 4239, 4240, 4246, 4262, 4278, 4279, 4306,
4318, 4324, 4330, 4337, 4342, 4343, 4356, 4359, 4362, 4374, 4405,
4427, 4434, 4452, 4453, 4458, 4465, 4468, 4479, 4480, 4487, 4489,
4497, 4504, 4512, 4518, 4521, 4538, 4555, 4556, 4569, 4573, 4591,
4598, 4613, 4622, 4627, 4642, 4646, 4650, 4664, 4681, 4691, 4698,
4722, 4734, 4749, 4757, 4764, 4779, 4781, 4784, 4794, 4797, 4801,
4802, 4806, 4829, 4831, 4856, 4864, 4892, 4898, 4904, 4910, 4919,
4920, 4925, 4932, 4936, 4940, 4942, 4958, 4980, 4982, 4986, 4987,
5008, 5021, 5024, 5040, 5045, 5054, 5061, 5071, 5093, 5102, 5123,
5125, 5127, 5130, 5142, 5152, 5162, 5163, 5166, 5191, 5215, 5216,
```



```

5222, 5226, 5227, 5235, 5240, 5245, 5250, 5263, 5295, 5298, 5348,
5353, 5355, 5359, 5361, 5368, 5392, 5395, 5398, 5402, 5403, 5409,
5415, 5421, 5422, 5432, 5436, 5443, 5451]),
14: array([ 57, 139, 621, 727, 1261, 1271, 1274, 1324, 1488, 1549, 1591,
1603, 1753, 1804, 1866, 1956, 2725, 3111, 3139, 3265, 3393, 3785,
3792, 3926, 4031, 4120, 4126, 4163, 4302, 4758, 4848, 5457])})

```

In [81]:

```
X_train_real.columns
```

Out[81]:

```

Index(['Unnamed: 0', 'id', 'teacher_id', 'teacher_prefix', 'school_state',
      'Date', 'project_grade_category', 'project_title',
      'project_resource_summary',
      'teacher_number_of_previously_posted_projects', 'quantity', 'price',
      'clean_categories', 'clean_subcategories', 'essay', 'clean_essays',
      'clean_titles'],
      dtype='object')

```

In [82]:

```

# Collect easy text for word cloud
cluster_text_essay = dict()
for i in clusters_set:
    for j in clusters_set[i]:
        cluster_text_essay[i] = cluster_text_essay.get(i, '') + X_train_real.iloc[i]['clean_essays']
]

```

In [83]:

```

print(len(cluster_text_essay[0]))
len(cluster_text_essay)

```

6899200

Out[83]:

15

In [84]:

```

%%time
# create world cloud
from wordcloud import WordCloud, STOPWORDS

for i in tqdm(cluster_text_essay):
    wordcloud = WordCloud(width = 800, height = 800,
                          background_color = 'white',
                          min_font_size = 10).generate(cluster_text_essay[i])

    # plot the WordCloud image
    plt.figure(figsize = (6, 6), facecolor = None)
    plt.imshow(wordcloud)
    plt.axis("off")
    plt.title(f'Cluster {i}')
    plt.tight_layout(pad = 0)

plt.show()

```

100%|██████████| 15/15 [00:17<00:00, 1.13s/it]

Cluster 0





Cluster 4



Cluster 5



kids kids currently send charter diverse interested extracurricular
 spend compete love asking working hard
 school budget civic service expensive school offered expected
 students purchased

Cluster 6

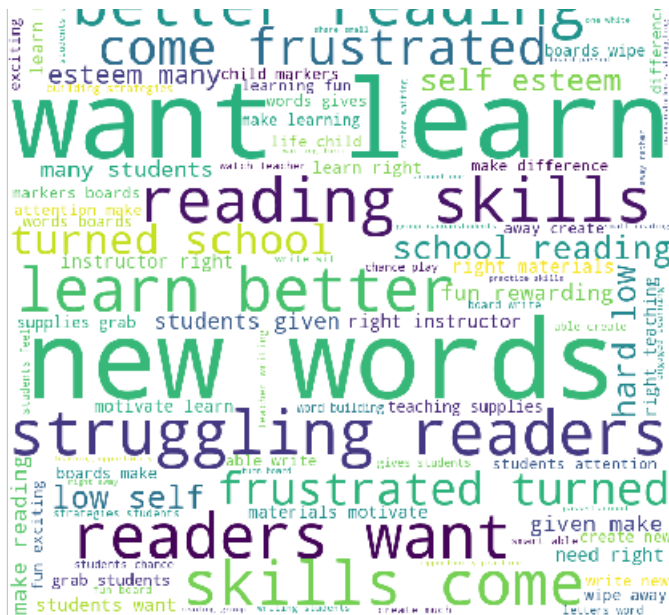
learning playing puppet love building class play dough students words center
 poetry science teach magic kit sweet one
 really energetic art another wet ny
 chef pieces kinder garden tinkertoys group diverse
 day games pumpkin going school new exploring
 car beloved success school play time
 extremely dynamic dancing zoo enthusiastic singing
 creating four visiting sight popular poems choice
 books fairy funny children math need
 lotsa cooks math need
 vocabulary pot facts writing fabulous adore stories brooklyn

Cluster 7

important balls students average material four kickball many
 activity equipment male physical activities simple full
 working playing energy well
 help see
 need basic wellness education sport play
 level project use growth heart
 currently monitors female needed want love
 classes budgets required improvement covering essential mixture
 opt school fitness basketball footballs immediate lowered

Cluster 8

low socioeconomic creative inviting provide
 want basis eye create critical gathering
 early lessons love receives daily
 independent classroom library new
 learners responsible seat soft creativity
 thinking learn elementary books opportunity part work
 students paired mastering throughout cozy excited practicing react
 environment place look process forward go model able located attend
 kids school readers foster
 workshop



```
CPU times: user 16.7 s, sys: 1.5 s, total: 18.2 s
Wall time: 18.3 s
```

In [85]:

```
from prettytable import PrettyTable
from collections import Counter

table = PrettyTable()
table.field_names = ["Cluster No", "No. of words in cluster", "Most frequent words"]

for key, val in cluster_text_essay.items():
    freq_dict = dict(sorted(list(Counter(val.split()).items()), key=lambda x: x[1], reverse=True))
    table.add_row([key, len(val.split()), ".".join(list(freq_dict.keys())[:10])])

print(table)
```

Cluster No	No. of words in cluster	Most frequent words
0	893201	students, materials, come, families, class, feel, small, teaching, years, love
1	119	students, printer, make, want, also, classroom, able, share, learning, school
2	46081	students, skills, year, learning, school, help, great, quality, early, providing
3	2091	students, class, well, classroom, instruction, assist, success, technology, plan, math
4	1193	students, table, would, amazing, reading, not, day, work, group, learn
5	10231	kids, students, get, seek, love, not, academics, extracurricular, interested, school
6	122	students, words, would, love, building, class, play, dough, center, kindergarten
7	136	class, students, basic, play, fitness, level, need, balls, important, high
8	287	reading, students, school, classroom, create, year, workshop, learning, writing, low
9	995	students, seating, day, group, classroom, learning, options, enhance, student, choose
10	104	students, reading, writing, new, year, spot, make, school, would, comfortable
11	1519	students, class, learn, work, interactive, whiteboard, classroom, learning, ready, share
12	7801	students, learning, classroom, best, work, environment, learn, options, seating, challenges
13	60157	students, technology, use, school, economic, free, also, kids, many, not

```
|      14      |      3361      |  
students,reading,right,learn,make,fun,board,want,skills,boards      |  
+-----+-----+-----+-----+-----+-----+-----+-----+-----+  
-----+  
◀────────────────────────────────────────────────────────────────────────────────▶
```

2.6 Apply AgglomerativeClustering

In [86]:

```
%%time  
from sklearn.cluster import AgglomerativeClustering  
from sklearn.metrics import silhouette_score  
  
X_train_bow_5k_aggl = X_train_bow_5k.toarray()  
  
clusters = [2, 3, 5, 7]  
scores = []  
for i in tqdm(clusters):  
    print(f'Starting Cluster #{i}')  
    aggl_cluster = AgglomerativeClustering(n_clusters=i)  
    aggl_cluster.fit(X_train_bow_5k_aggl)  
  
    score = silhouette_score(X_train_bow_5k_aggl, aggl_cluster.labels_, random_state=42)  
    scores.append(score)
```

```
0%|          | 0/4 [00:00<?, ?it/s]
```

Starting Cluster #2

```
25%|██        | 1/4 [01:07<03:22, 67.40s/it]
```

Starting Cluster #3

```
50%|██████    | 2/4 [02:14<02:14, 67.33s/it]
```

Starting Cluster #5

```
75%|██████████| 3/4 [03:22<01:07, 67.50s/it]
```

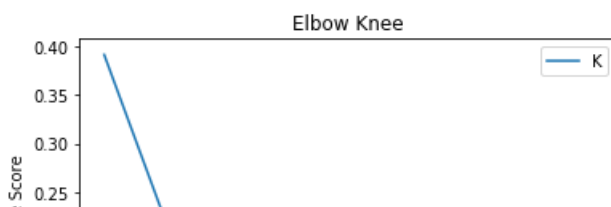
Starting Cluster #7

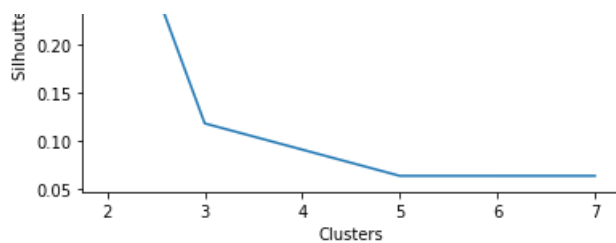
```
100%|██████████| 4/4 [04:30<00:00, 67.78s/it]
```

CPU times: user 4min 34s, sys: 980 ms, total: 4min 35s
Wall time: 4min 30s

In [87]:

```
# plot k vs Silhoutte Score  
  
plt.plot(clusters, scores)  
plt.xlabel('Clusters')  
plt.ylabel('Silhoutte Score')  
plt.title('Elbow Knee')  
plt.legend('Knee')  
plt.show()
```





In [88]:

```
aggl_cluster = AgglomerativeClustering(n_clusters=3)
aggl_cluster.fit(X_train_bow_5k_aggl)
```

Out[88]:

```
AgglomerativeClustering(affinity='euclidean', compute_full_tree='auto',
                        connectivity=None, distance_threshold=None,
                        linkage='ward', memory=None, n_clusters=3,
                        pooling_func='deprecated')
```

In [89]:

```
aggl_cluster.n_clusters
```

Out[89]:

3

In [90]:

```
aggl_cluster.labels_
```

Out[90]:

```
array([2, 2, 2, ..., 2, 2, 2])
```

In [91]:

```
# Collect data points from each cluster
clusters_set = {i: np.where(aggl_cluster.labels_ == i)[0] for i in range(aggl_cluster.n_clusters)}

clusters_set

# Collect easy text for word cloud
cluster_text_essay = dict()
for i in clusters_set:
    for j in clusters_set[i]:
        cluster_text_essay[i] = cluster_text_essay.get(i, '') + X_train_real.iloc[j]['clean_essays']
]

print(len(cluster_text_essay[0]))
len(cluster_text_essay)
```

87808

Out[91]:

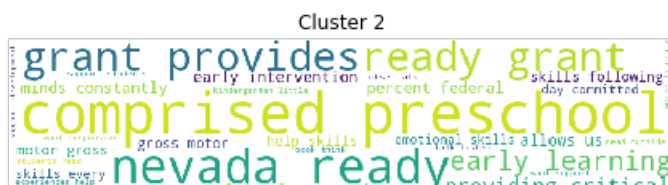
3

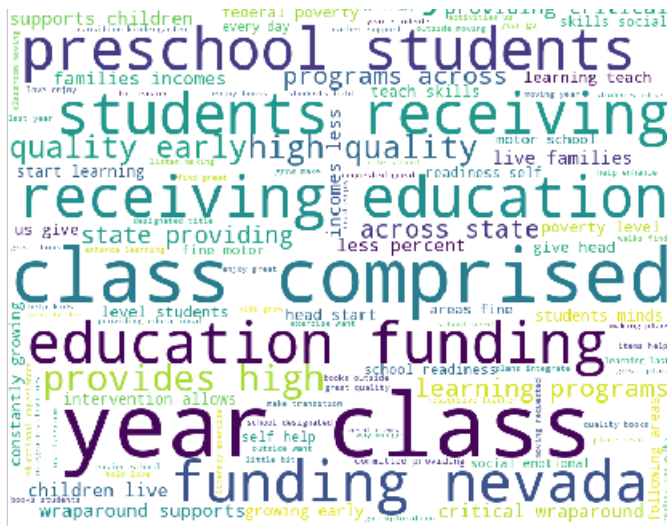
In [92]:

```
%%time
# create world cloud
from wordcloud import WordCloud, STOPWORDS

for i in tqdm(cluster_text_essay):
    wordcloud = WordCloud(width = 800, height = 800,
                          background_color = 'white',
```

```
100%|██████████| 3/3 [00:04<00:00, 1.58s/it]
```





CPU times: user 4.58 s, sys: 378 ms, total: 4.96 s
Wall time: 4.98 s

In [93]:

```
from prettytable import PrettyTable
from collections import Counter

table = PrettyTable()
table.field_names = ["Cluster No", "No. of words in cluster", "Most frequent words"]

for key, val in cluster_text_essay.items():
    freq_dict = dict(sorted(list(Counter(val.split()).items()), key=lambda x: x[1], reverse=True))
    table.add_row([key, len(val.split()), ",".join(list(freq_dict.keys())[:10])])

print (table)
```

```
+-----+-----+-----+-----+
+-----+
| Cluster No | No. of words in cluster | Most frequent words |
+-----+-----+-----+-----+
+-----+
| 0 | 11369 | students,materials,come,families,class,feel,small,teaching,years,love |
| 1 | 3777 | students,printer,make,want,also,classroom,able,share,learning,school |
| 2 | 859841 | students,skills,year,learning,school,help,great,quality,early,providing |
+-----+-----+-----+-----+
+-----+
```

2.7 Apply DBSCAN

In [94]:

```
%%time
# citation https://www.kaggle.com/rohit0812/clustering-on-donors-choose-data-set
from sklearn.neighbors import KDTree

X_train_bow_5k_dbscan = X_train_bow_5k.todense()[:5000]

algo_title = 'DBSCAN Clustering'
minPts = 8
tree = KDTree(X_train_bow_5k_dbscan)

idx = 0
epss = []
for x_i in tqdm(X_train_bow_5k_dbscan):
    epss.append(tree.query(X_train_bow_5k_dbscan[idx], return_distance=True, k=minPts)[0][0][-1])
    idx += 1
epss.sort()
```

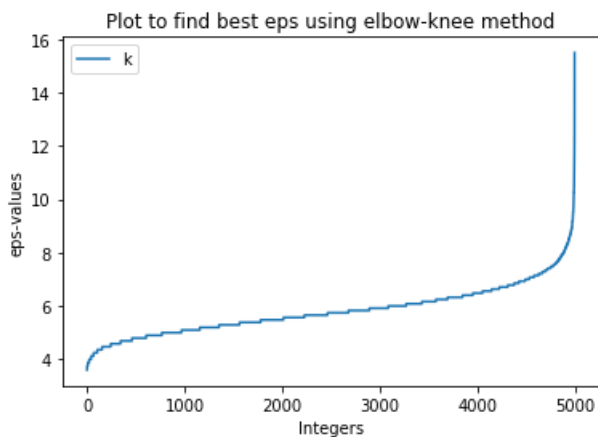
```
plt.plot(range(0,5000), epss[:5000])
plt.title("Plot to find best eps using elbow-knee method")
plt.xlabel('Integers')
plt.ylabel('eps-values')
plt.legend('kneee')
```

100%|██████████| 5000/5000 [02:40<00:00, 31.19it/s]

CPU times: user 2min 35s, sys: 1.25 s, total: 2min 36s
Wall time: 2min 41s

Out[94]:

<matplotlib.legend.Legend at 0x1a298e3080>



In [95]:

```
def getCorupusDict(essay_hot_info, y_pred):
    one_hot_featr, one_hot_enc = essay_hot_info
    one_hot_enc_cols = one_hot_enc.shape[1]
    corpus_dict = {}
    i = 0
    for each_x in tqdm(y_pred):
        if each_x not in corpus_dict: corpus_dict[each_x] = ''
        for j in range(one_hot_enc_cols):
            if one_hot_enc[i][j] >= 0.5:
                corpus_dict[each_x] = "%s %s"%(corpus_dict[each_x], one_hot_featr[j])
        i += 1

    return corpus_dict
```

In [96]:

```
from sklearn.cluster import DBSCAN

cluster = DBSCAN(eps=7, min_samples=minPts).fit(X_train_bow_5k_dbscan)
corpus_dict = getCorupusDict((X_train_bow_5k, X_train_bow_5k.toarray()), cluster.labels_)
print ("number of clusters gotten:", len(corpus_dict))
```

100%|██████████| 5000/5000 [12:55<00:00, 3.16it/s]

number of clusters gotten: 3

In [97]:

```
cluster.labels_
```

Out[97]:

```
array([0, 0, 0, ..., 0, 0, 0])
```

In [98]:

```
# Collect data points from each cluster
clusters_set = {i: np.where(cluster.labels_ == i)[0] for i in range(len(corpus_dict))}

clusters_set

# Collect easy text for word cloud
cluster_text_essay = dict()
for i in clusters_set:
    for j in clusters_set[i]:
        cluster_text_essay[i] = cluster_text_essay.get(i, '') + X_train_real.iloc[j]['clean_essays']
]

print(len(cluster_text_essay[0]))
len(cluster_text_essay)
```

7360192

Out[98]:

2

In [99]:

```
%%time
# create world cloud
from wordcloud import WordCloud, STOPWORDS

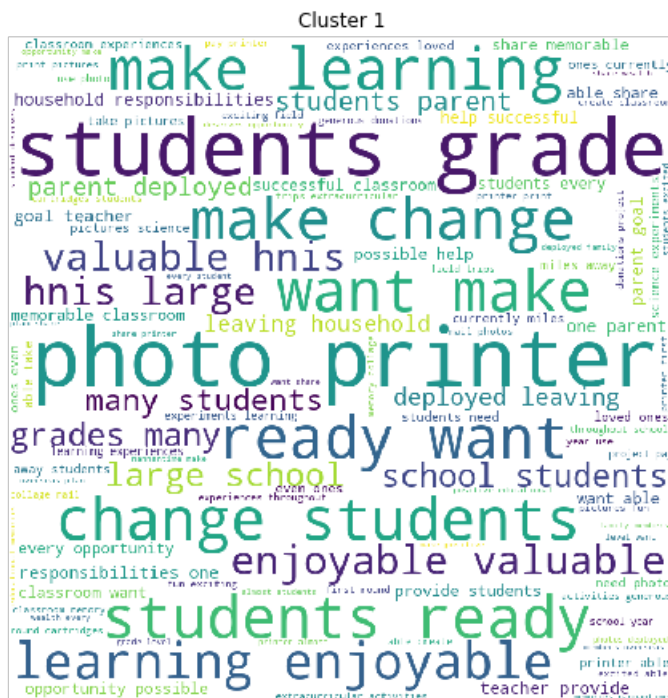
for i in tqdm(cluster_text_essay):
    wordcloud = WordCloud(width = 800, height = 800,
                           background_color = 'white',
                           min_font_size = 10).generate(cluster_text_essay[i])

    # plot the WordCloud image
    plt.figure(figsize = (6, 6), facecolor = None)
    plt.imshow(wordcloud)
    plt.axis("off")
    plt.title(f'Cluster {i}')
    plt.tight_layout(pad = 0)

plt.show()
```

100% |██████████| 2/2 [00:04<00:00, 2.74s/it]





CPU times: user 4.37 s, sys: 362 ms, total: 4.74 s
Wall time: 4.81 s

In [101]:

```
from prettytable import PrettyTable
from collections import Counter

table = PrettyTable()
table.field_names = ["Cluster No", "No. of words in cluster", "Most frequent words"]

for key, val in cluster_text_essay.items():
    freq_dict = dict(sorted(list(Counter(val.split()).items()), key=lambda x: x[1], reverse=True))
    table.add_row([key, len(val.split()), ",".join(list(freq_dict.keys())[:10])])

print(table)
```

```
+-----+-----+-----+
-----+
| Cluster No | No. of words in cluster | Most frequent words
|
+-----+-----+-----+
-----+
| 0 | 952883 | students,materials,come,families,class,feel,small,teaching,years,love |
| 1 | 1771 | students,printer,make,want,also,classroom,able,share,learning,school |
+-----+-----+-----+
-----+
```

Conclusion

Select multinomial naive bayes with BOW for best AUC

Select 5K best features using SelectFromModel

Kmeans

1. Train the model and plot K vs Inertia Plot

2. Select K=15 as best clusters using knee-elbow method

3. I found 1 cluster is very dense with 893201 words followed by 3 medium sized clusters

4. Plot word cloud for every cluster and show top words from each cluster

AgglomerativeClustering

1. Train the model and plot Silhouette Score

2. Select 3 best clusters

3. I found 1 cluster is very dense

4. Plot word cloud for every cluster and show top words from each cluster

DBSCAN

1. Use KDtree to find best eps

2. I found two clusters with 0 and -1 label

3. Plot word cloud for every cluster and show top words from each cluster

In []: