

# Choosing features and metrics for nearest neighbor search

When exploring a large set of documents -- such as Wikipedia, news articles, StackOverflow, etc. -- it can be useful to get a list of related material. To find relevant documents you typically

- Decide on a notion of similarity
- Find the documents that are most similar

In the assignment you will

- Gain intuition for different notions of similarity and practice finding similar documents.
- Explore the tradeoffs with representing documents using raw word counts and TF-IDF
- Explore the behavior of different distance metrics by looking at the Wikipedia pages most similar to President Obama's page.

## If you are using GraphLab Create

An IPython Notebook has been provided below to you for this assignment. This notebook contains the instructions, quiz questions and partially-completed code for you to use as well as some cells to test your code.

- Download the Wikipedia people dataset in SFrame format: [people\\_wiki.gl.zip](#)
- Download the companion IPython notebook: [0\\_nearest-neighbors-features-and-metrics\\_blank.ipynb](#)
- Save both of these files in the same directory (where you are calling IPython notebook from) and unzip the data file.

**Open the companion IPython notebook and follow the instructions in the notebook. The instructions below do not apply to users of GraphLab Create.**

# If you are not using GraphLab Create

It is possible to complete this assignment without using GraphLab Create. The instructions below are geared towards Python users, but you are free to adapt them to your specific environment.

**Disclaimer.** We have tested the assessment using the standard Python installation (with access to scikit-learn). However, the assessment may not be compatible with other tools (e.g. Matlab, R).

## Download the dataset

- Download the Wikipedia people dataset in SFrame format: [people\\_wiki.gl.zip](#). (Those experimenting with other tools, get [people\\_wiki.csv.zip](#) instead.)
- Download the mapping between words and integer indices: [people\\_wiki\\_map\\_index\\_to\\_word.gl.zip](#) (or alternatively, [people\\_wiki\\_map\\_index\\_to\\_word.json.zip](#))
- Download the pre-processed set of word counts: [people\\_wiki\\_word\\_count.npz](#)
- Download the pre-processed set of TF-IDF scores: [people\\_wiki\\_tf\\_idf.npz](#)

## Import packages

```
1 import sframe          # see below for install instruction
2 import matplotlib.pyplot as plt # plotting
3 import numpy as np      # dense matrices
4 from scipy.sparse import csr_matrix # sparse matrices
5 %matplotlib inline
```

**About SFrame.** SFrame is a dataframe library Dato has released free-of-charge. Its source code is available [here](#). You may install SFrame via pip:

```
1 pip install --upgrade sframe
```

## Load in the dataset

We will be using the same dataset of Wikipedia pages that we used in the Machine Learning Foundations course (Course 1). Each element of the dataset consists of a link to the wikipedia article, the name of the person, and the text of the article (in lowercase).

```
1 wiki = sframe.SFrame('people_wiki.gl/')
2 wiki = wiki.add_row_number() # add row number, starting at 0
```

## Extract word count vectors

For your convenience, we extracted the word count vectors from the dataset. The vectors are packaged in a sparse matrix, where the i-th row gives the word count vectors for the i-th document. Each column corresponds to a unique word appearing in the dataset. The mapping between words and integer indices are given in people\_wiki\_map\_index\_to\_word.gl.

To load in the word count vectors, define the function

```
1 def load_sparse_csr(filename):
2     loader = np.load(filename)
3     data = loader['data']
4     indices = loader['indices']
5     indptr = loader['indptr']
6     shape = loader['shape']
7
8     return csr_matrix((data, indices, indptr), shape)
```

and then run

```
1 word_count = load_sparse_csr('people_wiki_word_count.npz')
```

The word-to-index mapping is given by

```
1 map_index_to_word = sframe.SFrame('people_wiki_map_index_to_word.gl/')
```

(Optional) *Extracting word count vectors yourself.* We provide the pre-computed word count vectors to minimize potential compatibility issues. You are free to experiment with other tools to compute the word count vectors yourself. A good place to start is [sklearn.CountVectorizer](#). Note. Due to variations in [tokenization](#) and other factors, your word count vectors may differ from the ones we provide. For the purpose of the assessment, we ask you to use the vectors from people\_wiki\_word\_count.npz.

## Find nearest neighbors using word count vectors

Let's start by finding the nearest neighbors of the Barack Obama page using the word count vectors to represent the articles and Euclidean distance to measure distance. For this, we will use scikit-learn's implementation of k-nearest neighbors. We first create an instance of the NearestNeighbors class, specifying the model parameters. Then we call the fit() method to attach the training set.

```
1 from sklearn.neighbors import NearestNeighbors  
2  
3 model = NearestNeighbors(metric='euclidean', algorithm='brute')  
4 model.fit(word_count)
```

Run the following cell to obtain the row number for Obama's article:

```
1 print wiki[wiki['name'] == 'Barack Obama']
```

The output should resemble the following:

```
1 +-----+-----+-----+  
2 | id | URI | name |  
3 +---+---+---+  
4 | 35817 | <http://dbpedia.org/resource/Barack_Obama> | Barack Obama |  
5 +---+---+---+  
6 +-----+-----+  
7 | text | word_count |  
8 +-----+-----+  
9 | barack_hussein_obama ii br... | {'operations': 1, 'represe... |  
10 +-----+-----+
```

which locates Obama's article at index 35817.

Let us run the k-nearest neighbor algorithm with Obama's article. Since the NearestNeighbor class expects a vector, we pass the 35817th row of word\_count vector.

```
1 distances, indices = model.kneighbors(word_count[35817], n_neighbors=10) # 1st arg: word count  
vector
```

The query returns the indices of and distances to the 10 nearest neighbors. To display the indices and distances together with the article name, run

```
1 neighbors = sframe.SFrame({'distance':distances.flatten(), 'id':indices.flatten()})  
2 print wiki.join(neighbors, on='id').sort('distance')[['id','name','distance']]
```

which should give us an output in the form of the following:

1	+-----+	+-----+	+-----+
2	id	name	distance
3	+-----+	+-----+	+-----+
4	35817	Barack Obama	0.0
5	24478	Joe Biden	33.0756708171
6	28447	George W. Bush	34.3947670438
7	35357	Lawrence Summers	36.1524549651
8	14754	Mitt Romney	36.1662826401
9	13229	Francisco Barrio	36.3318042492
10	31423	Walter Mondale	36.4005494464
11	22745	Wynn Normington Hugh-Jones	36.4965751818
12	36364	Don Bonker	36.633318168
13	9210	Andy Anstett	36.9594372252
14	+-----+	+-----+	+-----+

For now, treat this piece of code as a black box; we will revisit the **join** operation in the following sections.

## Interpreting the nearest neighbors

All of the 10 people are politicians, but about half of them have rather tenuous connections with Obama, other than the fact that they are politicians.

- Francisco Barrio is a Mexican politician, and a former governor of Chihuahua.
- Walter Mondale and Don Bonker are Democrats who made their career in late 1970s.
- Wynn Normington Hugh-Jones is a former British diplomat and Liberal Party official.
- Andy Anstett is a former politician in Manitoba, Canada.

Nearest neighbors with raw word counts got some things right, showing all politicians in the query result, but missed finer and important details.

For instance, let's find out why Francisco Barrio was considered a close neighbor of Obama. To do this, let's look at the most frequently used words in each of Barack Obama and Francisco Barrio's pages.

First, run the following cell to obtain the word\_count column, which represents the word count vectors in the dictionary form. This way, we can quickly recognize words of great importance.

**Note.** Understanding this code is not required for completing this assignment. Feel free to treat it as a black box.

```
1 def unpack_dict(matrix, map_index_to_word):
2     table = list(map_index_to_word.sort('index')['category'])
3     data = matrix.data
4     indices = matrix.indices
5     indptr = matrix.indptr
6
7     num_doc = matrix.shape[0]
8
9     return [{k:v for k,v in zip([table[word_id] for word_id in indices[indptr[i]:indptr[i+1]]], 
10                                data[indptr[i]:indptr[i+1]].tolist())} \
11            for i in xrange(num_doc) ]
12
13 wiki['word_count'] = unpack_dict(word_count, map_index_to_word)
```

Now printing the SFrame produces

```
1 +-----+-----+
2 | id | URI | name |
3 +-----+-----+
4 | 0 | <http://dbpedia.org/resource/Digby_Morrell> | Digby Morrell |
5 | 1 | <http://dbpedia.org/resource/Alfred_J._Lewy> | Alfred J. Lewy |
6 | 2 | <http://dbpedia.org/resource/Harpdog_Brown> | Harpdog Brown |
7 | 3 | <http://dbpedia.org/resource/Franz_Rottensteiner> | Franz Rottensteiner |
8 | 4 | <http://dbpedia.org/resource/G-Enka> | G-Enka |
9 | 5 | <http://dbpedia.org/resource/Sam_Henderson> | Sam Henderson |
10 | 6 | <http://dbpedia.org/resource/Aaron_LaCrate> | Aaron LaCrate |
11 | 7 | <http://dbpedia.org/resource/Trevor_Ferguson> | Trevor Ferguson |
12 | 8 | <http://dbpedia.org/resource/Grant_Nelson> | Grant Nelson |
13 | 9 | <http://dbpedia.org/resource/Cathy_Caruth> | Cathy Caruth |
14 +-----+-----+
15 +-----+-----+
16 | text | word_count |
17 +-----+-----+
18 | digby morrell born 10 octo... | {'selection': 1, 'carltons... |
19 | alfred j lewy aka sandy le... | {'precise': 1, 'thomas': 1... |
20 | harpdog brown is a singer ... | {'just': 1, 'issued': 1, '... |
21 | franz rottensteiner born i... | {'englishreading': 1, 'all... |
22 | henry krvits born 30 decem... | {'they': 1, 'gangstergenka... |
23 | sam henderson born october... | {'now': 1, 'currently': 1,... |
24 | aaron lacrate is an americ... | {'exclusive': 2, 'producer... |
25 | trevor ferguson aka john f... | {'taxi': 1, 'salon': 1, 'g... |
26 | grant nelson born 27 april... | {'houston': 1, 'frankie': ... |
27 | cathy caruth born 1955 is ... | {'phenomenon': 1, 'deboras... |
28 +-----+-----+
```

To make things even easier, we provide a utility function that displays a dictionary in tabular form:

```
1 def top_words(name):
2     """
3         Get a table of the most frequent words in the given person's wikipedia page.
4     """
5     row = wiki[wiki['name'] == name]
6     word_count_table = row[['word_count']].stack('word_count', new_column_name=['word', 'count'])
7     return word_count_table.sort('count', ascending=False)
8
9 obama_words = top_words('Barack Obama')
10 print obama_words
11
12 barrio_words = top_words('Francisco Barrio')
13 print barrio_words
```

Let's extract the list of most frequent words that appear in both Obama's and Barrio's documents. We've so far sorted all words from Obama and Barrio's articles by their word frequencies. We will now use a dataframe operation known as **join**. The **join** operation is very useful when it comes to playing around with data: it lets you combine the content of two tables using a shared column (in this case, the word column). See [the documentation](#) for more details.

For instance, running

```
1 combined_words = obama_words.join(barrio_words, on='word')
```

will extract the rows from both tables that correspond to the common words.

Since both tables contained the column named count, SFrame automatically renamed one of them to prevent confusion. Let's rename the columns to tell which one is for which. By inspection, we see that the first column (count) is for Obama and the second (count.1) for Barrio.

```
1 combined_words = combined_words.rename({'count':'Obama', 'count.1':'Barrio'})
```

**Note.** The **join** operation does not enforce any particular ordering on the shared column. So to obtain, say, the five common words that appear most often in Obama's article, sort the combined table by the Obama column. Don't forget ascending=False to display largest counts first.

```
1 combined_words.sort('Obama', ascending=False)
```

**Quiz Question.** Among the words that appear in both Barack Obama and Francisco Barrio, take the 5 that appear most frequently in Obama. How many of the articles in the Wikipedia dataset contain all of those 5 words?

Hint:

- Refer to the previous paragraph for finding the words that appear in both articles. Sort the common words by their frequencies in Obama's article and take the largest five.
- Each word count vector is a Python dictionary. For each word count vector in SFrame, you'd have to check if the set of the 5 common words is a subset of the keys of the word count vector. Complete the function has\_top\_words to accomplish the task.
- Convert the list of top 5 words into set using the syntax "set(common\_words)", where common\_words is a Python list. See [this link](#) if you're curious about Python sets.
- Extract the list of keys of the word count dictionary by calling the [keys\(\)](#) method.
- Convert the list of keys into a set as well.

- Use `issubset()` method to check if all 5 words are among the keys.
- Now apply the `has_top_words` function on every row of the SFrame.
- Compute the sum of the result column to obtain the number of articles containing all the 5 top words.

```

1 common_words = ... # YOUR CODE HERE
2
3 def has_top_words(word_count_vector):
4     # extract the keys of word_count_vector and convert it to a set
5     unique_words = ... # YOUR CODE HERE
6     # return True if common_words is a subset of unique_words
7     # return False otherwise
8     return ... # YOUR CODE HERE
9
10 wiki['has_top_words'] = wiki['word_count'].apply(has_top_words)
11
12 # use has_top_words column to answer the quiz question
13 ... # YOUR CODE HERE

```

**Checkpoint.** Check your `has_top_words` function on two random articles:

```

1 print 'Output from your function:', has_top_words(wiki[32]['word_count'])
2 print 'Correct output: True'
3 print 'Also check the length of unique_words. It should be 167'
4
5 print 'Output from your function:', has_top_words(wiki[33]['word_count'])
6 print 'Correct output: False'
7 print 'Also check the length of unique_words. It should be 188'

```

**Quiz Question.** Measure the pairwise distance between the Wikipedia pages of Barack Obama, George W. Bush, and Joe Biden. Which of the three pairs has the smallest distance?

Hint: For this question, take the row vectors from the TF-IDF matrix that correspond to Obama, Bush, and Biden. To compute the Euclidean distance between any two sparse vectors, use `sklearn.metrics.pairwise.euclidean_distances`.

**Quiz Question.** Collect all words that appear both in Barack Obama and George W. Bush pages. Out of those words, find the 10 words that show up most often in Obama's page.

**Note.** Even though common words are swamping out important subtle differences, commonalities in rarer political words still matter on the margin. This is why politicians are being listed in the query result instead of musicians, for example. In the next subsection, we will introduce a different metric that will place greater emphasis on those rarer words.

## Extract the TF-IDF vectors

Much of the perceived commonalities between Obama and Barrio were due to occurrences of extremely frequent words, such as "the", "and", and "his". So the nearest neighbors algorithm is recommending plausible results sometimes for the wrong reasons.

To retrieve articles that are more relevant, we should focus more on rare words that don't happen in every article. **TF-IDF** (term frequency-inverse document frequency) is a feature representation that penalizes words that are too common. Let us load in the TF-IDF vectors and repeat the nearest neighbor search.

For your convenience, we extracted the TF-IDF vectors from the dataset. The vectors are packaged in a sparse matrix, where the i-th row gives the TF-IDF vectors for the i-th document. Each column corresponds to a unique word appearing in the dataset. The mapping between words and integer indices are given in `people_wiki_map_index_to_word.gl`.

To load in the TF-IDF vectors, run

```
1 tf_idf = load_sparse_csr('people_wiki_tf_idf.npz')
```

In addition to the sparse matrix, we also store the TF-IDF vectors in dictionary form as well, to allow for easy interpretation.

```
1 wiki['tf_idf'] = unpack_dict(tf_idf, map_index_to_word)
```

(Optional) Extracting TF-IDF vectors yourself. We provide the pre-computed TF-IDF vectors to minimize potential compatibility issues. You are free to experiment with other tools to compute the TF-IDF vectors yourself. A good place to start is `sklearn.TfidfVectorizer`. Note. Due to variations in `tokenization` and other factors, your TF-IDF vectors may differ from the ones we provide. For the purpose the assessment, we ask you to use the vectors from `people_wiki_tf_idf.npz`.

## Find nearest neighbors using TF-IDF vectors

Since we are now using a different set of features, we should create a new nearest neighbor model. Create another instance of the `NearestNeighbor` class as follows. Then call the `fit()` method to associate it with the TF-IDF vectors.

```
1 model_tf_idf = NearestNeighbors(metric='euclidean', algorithm='brute')
2 model_tf_idf.fit(tf_idf)
```

Perform the nearest neighbor search by running

```
1 distances, indices = model_tf_idf.kneighbors(tf_idf[35817], n_neighbors=10)
```

To print the names of the articles, we perform join using the indices:

```
1 neighbors = sframe.SFrame({'distance':distances.flatten(), 'id':indices.flatten()})
2 print wiki.join(neighbors, on='id').sort('distance')[['id', 'name', 'distance']]
```

which produces

1	id	name	distance
2	35817	Barack Obama	0.0
3	7914	Phil Schiliro	106.861013691
4	46811	Jeff Sessions	108.871674216
5	44681	Jesse Lee (politician)	109.045697909
6	38376	Samantha Power	109.108106165
7	6507	Bob Menendez	109.781867105
8	38714	Eric Stern (politician)	109.95778808
9	44825	James A. Guest	110.413888718
10	44368	Roland Grossenbacher	110.4706087
11	33417	Tulsi Gabbard	110.696997999
12			
13			
14			

Let's determine whether this list makes sense.

- With a notable exception of Roland Grossenbacher, the other 8 are all American politicians who are contemporaries of Barack Obama.
- Phil Schiliro, Jesse Lee, Samantha Power, and Eric Stern worked for Obama.

Clearly, the results are more plausible with the use of TF-IDF. Let's take a look at the word vector for Obama and Schiliro's pages. Notice that TF-IDF representation assigns a weight to each word. This weight captures relative importance of that word in the document. Let us sort the words in Obama's article by their TF-IDF weights; we do the same for Schiliro's article as well.

```
1 def top_words_tf_idf(name):
2     row = wiki[wiki['name'] == name]
3     word_count_table = row[['tf_idf']].stack('tf_idf', new_column_name=['word', 'weight'])
4     return word_count_table.sort('weight', ascending=False)
5
6 obama_tf_idf = top_words_tf_idf('Barack Obama')
7 print obama_tf_idf
8
9 schiliro_tf_idf = top_words_tf_idf('Phil Schiliro')
10 print schiliro_tf_idf
```

Using the **join** operation we learned earlier, try your hands at computing the common words shared by Obama's and Schiliro's articles. Sort the common words by their TF-IDF weights in Obama's document. The first 10 words should say: Obama, law, democratic, Senate, presidential, president, policy, states, office, 2011.

**Quiz Question.** Among the words that appear in both Barack Obama and Phil Schiliro, take the 5 that have largest weights in Obama. How many of the articles in the Wikipedia dataset contain all of those 5 words?

```
1 common_words = ... # YOUR CODE HERE
2
3 def has_top_words(word_count_vector):
4     # extract the keys of word_count_vector and convert it to a set
5     unique_words = ... # YOUR CODE HERE
6     # return True if common_words is a subset of unique_words
7     # return False otherwise
8     return ... # YOUR CODE HERE
9
10 wiki['has_top_words'] = wiki['word_count'].apply(has_top_words)
11
12 # use has_top_words column to answer the quiz question
13 ... # YOUR CODE HERE
```

Notice the huge difference in this calculation using TF-IDF scores instead of raw word counts. We've eliminated noise arising from extremely common words.

### Choosing metrics

You may wonder why Joe Biden, Obama's running mate in two presidential elections, is missing from the query results of model\_tf\_idf. Let's find out why. First, compute the distance between TF-IDF features of Obama and Biden.

**Quiz Question.** Compute the Euclidean distance between TF-IDF features of Obama and Biden.

The distance is larger than the distances we found for the 10 nearest neighbors. But one may wonder, is Biden's article that different from Obama's, more so than, say, Schiliro's? It turns out that, when we compute nearest neighbors using the Euclidean distances, we unwittingly favor short articles over long ones. Let us compute the length of each Wikipedia document, and examine the document lengths for the 100 nearest neighbors to Obama's page.

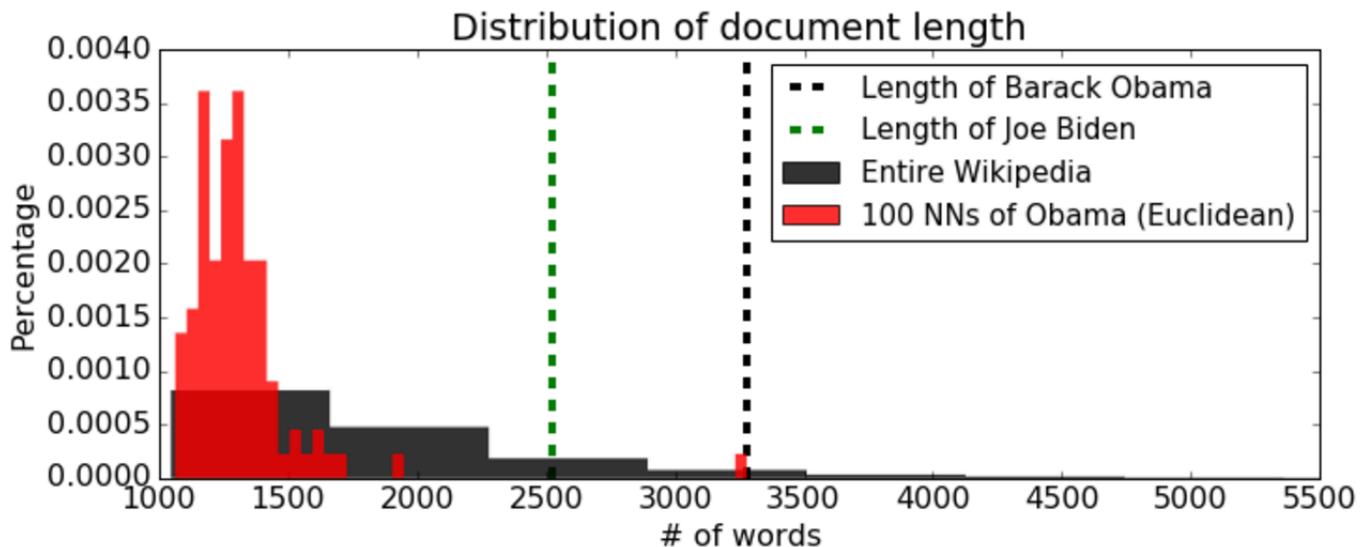
```
1 # Comptue length of all documents
2 def compute_length(row):
3     return len(row['text'])
4 wiki['length'] = wiki.apply(compute_length)
5
6 # Compute 100 nearest neighbors and display their lengths
7 distances, indices = model_tf_idf.kneighbors(tf_idf[35817], n_neighbors=100)
8 neighbors = sframe.SFrame({'distance':distances.flatten(), 'id':indices.flatten()})
9 nearest_neighbors_euclidean = wiki.join(neighbors, on='id')[['id', 'name', 'length',
10     'distance']].sort('distance')
```

The output should be in the form

1	id	name	length	distance
2	35817	Barack Obama	3278	0.0
3	7914	Phil Schiliro	1288	106.861013691
4	46811	Jeff Sessions	1398	108.871674216
5	44681	Jesse Lee (politician)	1374	109.045697909
6	38376	Samantha Power	1911	109.108106165
7	6507	Bob Menendez	1222	109.781867105
8	38714	Eric Stern (politician)	1589	109.95778808
9	44825	James A. Guest	1251	110.413888718
10	44368	Roland Grossenbacher	1099	110.4706087
11	33417	Tulsi Gabbard	1352	110.696997999
12				
13				
14				

To see how these document lengths compare to the lengths of other documents in the corpus, let's make a histogram of the document lengths of Obama's 100 nearest neighbors and compare to a histogram of document lengths for all documents.

```
1 plt.figure(figsize=(10.5,4.5))
2 plt.hist(wiki['length'], 50, color='k', edgecolor='None', histtype='stepfilled', normed=True,
3          label='Entire Wikipedia', zorder=3, alpha=0.8)
4 plt.hist(nearest_neighbors_euclidean['length'], 50, color='r', edgecolor='None', histtype
5          ='stepfilled', normed=True,
6          label='100 NNs of Obama (Euclidean)', zorder=10, alpha=0.8)
7 plt.axvline(x=wiki['length'][wiki['name']] == 'Barack Obama'][0], color='k', linestyle='--',
8             linewidth=4,
9             label='Length of Barack Obama', zorder=2)
10 plt.axvline(x=wiki['length'][wiki['name']] == 'Joe Biden'][0], color='g', linestyle='--',
11              linewidth=4,
12              label='Length of Joe Biden', zorder=1)
13 plt.axis([1000, 5500, 0, 0.004])
14
15 plt.legend(loc='best', prop={'size':15})
16 plt.title('Distribution of document length')
17 plt.xlabel('# of words')
18 plt.ylabel('Percentage')
19 plt.rcParams.update({'font.size':16})
20 plt.tight_layout()
```



Relative to the rest of Wikipedia, nearest neighbors of Obama are overwhelmingly short, most of them being shorter than 2000 words. The bias towards short articles is not appropriate in this application as there is really no reason to favor short articles over long articles (they are all Wikipedia articles, after all). Many Wikipedia articles are 2500 words or more, and both Obama and Biden are over 2500 words long.

Note: Both word-count features and TF-IDF are proportional to word frequencies. While TF-IDF penalizes very common words, longer articles tend to have longer TF-IDF vectors simply because they have more words in them.

To remove this bias, we turn to **cosine distances**:

$$d(\mathbf{x}, \mathbf{y}) = 1 - \frac{\mathbf{x}^T \mathbf{y}}{\|\mathbf{x}\| \|\mathbf{y}\|}$$

Cosine distances let us compare word distributions of two articles of varying lengths.

Let us train a new nearest neighbor model, this time with cosine distances. We then repeat the search for Obama's 100 nearest neighbors.

```

1 model2_tf_idf = NearestNeighbors(algorithm='brute', metric='cosine')
2 model2_tf_idf.fit(tf_idf)
3 distances, indices = model2_tf_idf.kneighbors(tf_idf[35817], n_neighbors=100)
4 neighbors = sframe.SFrame({'distance':distances.flatten(), 'id':indices.flatten()})
5 nearest_neighbors_cosine = wiki.join(neighbors, on='id')[['id', 'name', 'length', 'distance']]
    .sort('distance')
6 print nearest_neighbors_cosine

```

which prints

1	+-----+	id	name	length	distance	+
2		35817	Barack Obama	3278	-1.11022302463e-15	
3	+-----+	24478	Joe Biden	2523	0.703138676734	
4		38376	Samantha Power	1911	0.742981902328	
5		57108	Hillary Rodham Clinton	3472	0.758358397887	
6		38714	Eric Stern (politician)	1589	0.770561227601	
7		46140	Robert Gibbs	1572	0.784677504751	
8		6796	Eric Holder	1430	0.788039072943	
9		44681	Jesse Lee (politician)	1374	0.790926415366	
10		18827	Henry Waxman	1607	0.798322602893	
11		2412	Joe the Plumber	1422	0.799466360042	
12	+-----+-----+-----+-----+-----+-----+-----+					
13						
14						

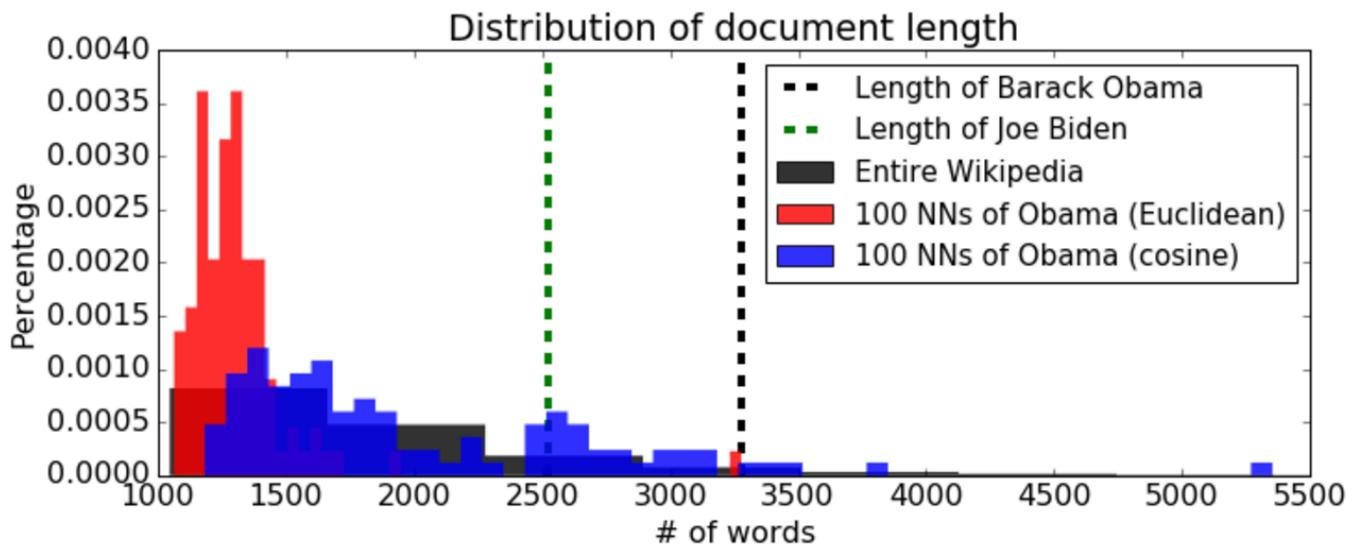
From a glance at the above table, things look better. For example, we now see Joe Biden as Barack Obama's nearest neighbor! We also see Hillary Clinton on the list. This list looks even more plausible as nearest neighbors of Barack Obama.

Let's make a plot to better visualize the effect of having used cosine distance in place of Euclidean on our TF-IDF vectors.

```

1 plt.figure(figsize=(10.5,4.5))
2 plt.figure(figsize=(10.5,4.5))
3 plt.hist(wiki['length'], 50, color='k', edgecolor='None', histtype='stepfilled', normed=True,
4         label='Entire Wikipedia', zorder=3, alpha=0.8)
5 plt.hist(nearest_neighbors_euclidean['length'], 50, color='r', edgecolor='None', histtype
6         ='stepfilled', normed=True,
7         label='100 NNs of Obama (Euclidean)', zorder=10, alpha=0.8)
8 plt.hist(nearest_neighbors_cosine['length'], 50, color='b', edgecolor='None', histtype
9         ='stepfilled', normed=True,
10        label='100 NNs of Obama (cosine)', zorder=11, alpha=0.8)
11 plt.axvline(x=wiki['length'][wiki['name']] == 'Barack Obama'][0], color='k', linestyle='--',
12             linewidth=4,
13             label='Length of Barack Obama', zorder=2)
14 plt.axvline(x=wiki['length'][wiki['name']] == 'Joe Biden'][0], color='g', linestyle='--',
15             linewidth=4,
16             label='Length of Joe Biden', zorder=1)
17 plt.axis([1000, 5500, 0, 0.004])
18 plt.legend(loc='best', prop={'size':15})
19 plt.title('Distribution of document length')
20 plt.xlabel('# of words')
21 plt.ylabel('Percentage')
22 plt.rcParams.update({'font.size': 16})
23 plt.tight_layout()

```



Indeed, the 100 nearest neighbors using cosine distance provide a sampling across the range of document lengths, rather than just short articles like Euclidean distance provided.

**Moral of the story:** In deciding the features and distance measures, check if they produce results that make sense for your particular application.

## Problem with cosine distances: tweets vs. long articles

Happily ever after? Not so fast. Cosine distances ignore all document lengths, which may be great in certain situations but not in others. For instance, consider the following (admittedly contrived) example.

```
1 +-----+
2 | | Follow |
3 | One that shall not be named | |
4 | @username | |
5 | | |
6 | Democratic governments control law in response to | |
7 | popular act. | |
8 | | |
9 | 8:05 AM - 16 May 2016 | |
10 | | |
11 | Reply Retweet (1,332) Like (300) | |
12 | | |
13 +-----+
```

How similar is this tweet to Barack Obama's Wikipedia article? Let's transform the tweet into TF-IDF features, using an encoder fit to the Wikipedia dataset. (That is, let's treat this tweet as an article in our Wikipedia dataset and see what happens.)

```
1 tweet = {'act': 3.4597778278724887,
2   'control': 3.721765211295327,
3   'democratic': 3.1026721743330414,
4   'governments': 4.167571323949673,
5   'in': 0.0009654063501214492,
6   'law': 2.4538226269605703,
7   'popular': 2.764478952022998,
8   'response': 4.261461747058352,
9   'to': 0.04694493768179923}
```

Let's look at the TF-IDF vectors for this tweet and for Barack Obama's Wikipedia entry, just to visually see their differences.

```
1 word_indices = [map_index_to_word[map_index_to_word['category']==word][0]['index'] for word in  
     tweet.keys()]  
2 tweet_tf_idf = csr_matrix( (list(tweet.values()), ([0]*len(word_indices), word_indices)),  
3                             shape=(1, tf_idf.shape[1]) )
```

Now, compute the cosine distance between the Barack Obama article and this tweet:

```
1 from sklearn.metrics.pairwise import cosine_distances  
2  
3 obama_tf_idf = tf_idf[35817]  
4 print cosine_distances(obama_tf_idf, tweet_tf_idf)
```

Let's compare this distance to the distance between the Barack Obama article and all of its Wikipedia 10 nearest neighbors:

```
1 distances, indices = model2_tf_idf.kneighbors(obama_tf_idf, n_neighbors=10)  
2 print distances
```

With cosine distances, the tweet is "nearer" to Barack Obama than everyone else, except for Joe Biden! This probably is not something we want. If someone is reading the Barack Obama Wikipedia page, would you want to recommend they read this tweet? Ignoring article lengths completely resulted in nonsensical results. In practice, it is common to enforce maximum or minimum document lengths. After all, when someone is reading a long article from *The Atlantic*, you wouldn't recommend him/her a tweet.