

Implementing Locality Sensitive Hashing from scratch

Locality Sensitive Hashing (LSH) provides for a fast, efficient approximate nearest neighbor search. The algorithm scales well with respect to the number of data points as well as dimensions.

In this assignment, you will

- Implement the LSH algorithm for approximate nearest neighbor search
- Examine the accuracy for different documents by comparing against brute force search, and also contrast runtimes
- Explore the role of the algorithm's tuning parameters in the accuracy of the method

NOTICE: This assignment requires SciPy 0.16.0 or later. To upgrade, run

```
1 !conda update -y scipy
```

in a blank notebook.

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: [1_nearest-neighbors-lsh-implementation_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 TF-IDF scores: [people_wiki_tf_idf.npz](#)

Import packages

```
1 import numpy as np                                # dense matrices
2 import sframe                                     # see below for install
3     instruction
4 from scipy.sparse import csr_matrix               # sparse matrices
5 from scipy.sparse.linalg import norm              # norms of sparse matrices
6 from sklearn.metrics.pairwise import pairwise_distances # pairwise distances
7 from copy import copy                            # deep copies
8 import matplotlib.pyplot as plt                  # plotting
9 %matplotlib inline
```

This assignment requires SciPy 0.16.0 or later. To upgrade, run

```
1 !conda update -y scipy
```

in a blank cell.

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

```
1 wiki = sframe.SFrame('people_wiki.g1/')
2 wiki = wiki.add_row_number() # add row number, starting at 0
```

Extract TF-IDF vectors

As in the previous assignment, we extract the TF-IDF vector of each document.

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.

To load in the TF-IDF vectors, run

```
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)
9
10 tf_idf = load_sparse_csr('people_wiki_tf_idf.npz')
```

The word-to-index mapping is given by

```
1 map_index_to_word = sframe.SFrame('people_wiki_map_index_to_word.g1/')
```

(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 of the assessment, we ask you to use the vectors from people_wiki_tf_idf.npz.

Train an LSH model

LSH performs an efficient neighbor search by randomly partitioning all reference data points into different bins. Today we will build a popular variant of LSH known as random binary projection, which approximates cosine distance. There are other variants we could use for other choices of distance metrics.

The first step is to generate a collection of random vectors from the standard Gaussian distribution.

```
1 def generate_random_vectors(num_vector, dim):  
2     return np.random.randn(dim, num_vector)
```

To visualize these Gaussian random vectors, let's look at an example in low-dimensions. Below, we generate 3 random vectors each of dimension 5.

```
1 # Generate 3 random vectors of dimension 5, arranged into a single 5 x 3 matrix.  
2 np.random.seed(0) # set seed=0 for consistent results  
3 print generate_random_vectors(num_vector=3, dim=5)
```

We now generate random vectors of the same dimensionality as our vocabulary size (547979). Each vector can be used to compute one bit in the bin encoding. We generate 16 vectors, leading to a 16-bit encoding of the bin index for each document.

```
1 # Generate 16 random vectors of dimension 547979  
2 np.random.seed(0)  
3 random_vectors = generate_random_vectors(num_vector=16, dim=547979)  
4 print random_vectors.shape
```

Next, we partition data points into bins. Instead of using explicit loops, we'd like to utilize matrix operations for greater efficiency. Let's walk through the construction step by step.

We'd like to decide which bin document 0 should go. Since 16 random vectors were generated in the previous cell, we have 16 bits to represent the bin index. The first bit is given by the sign of the dot product between the first random vector and the document's TF-IDF vector.

```
1 doc = corpus[0, :] # vector of tf-idf values for document 0  
2 print doc.dot(random_vectors[:, 0]) >= 0 # True if positive sign; False if negative sign
```

Similarly, the second bit is computed as the sign of the dot product between the second random vector and the document vector.

```
1 print doc.dot(random_vectors[:, 1]) >= 0 # True if positive sign; False if negative sign
```

We can compute all of the bin index bits at once as follows. Note the absence of the explicit for loop over the 16 vectors. Matrix operations let us batch dot-product computation in a highly efficient manner, unlike the for loop construction. Given the relative inefficiency of loops in Python, the advantage of matrix operations is even greater.

```
1 print doc.dot(random_vectors) >= 0 # should return an array of 16 True/False bits
2 print np.array(doc.dot(random_vectors) >= 0, dtype=int) # display index bits in 0/1's
```

All documents that obtain exactly this vector will be assigned to the same bin. We'd like to repeat the identical operation on all documents in the Wikipedia dataset and compute the corresponding bin indices. Again, we use matrix operations so that no explicit loop is needed.

```
1 print corpus[0:2].dot(random_vectors) >= 0 # compute bit indices of first two documents
2 print corpus.dot(random_vectors) >= 0 # compute bit indices of ALL documents
```

We're almost done! To make it convenient to refer to individual bins, we convert each binary bin index into a single integer:

1	Bin index	integer
2	[0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0]	=> 0
3	[0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,1]	=> 1
4	[0,0,0,0,0,0,0,0,0,0,0,0,0,0,1,0]	=> 2
5	[0,0,0,0,0,0,0,0,0,0,0,0,0,1,0,1]	=> 3
6	...	
7	[1,1,1,1,1,1,1,1,1,1,1,1,1,1,1,0]	=> 65532
8	[1,1,1,1,1,1,1,1,1,1,1,1,1,1,0,1]	=> 65533
9	[1,1,1,1,1,1,1,1,1,1,1,1,1,1,1,0]	=> 65534
10	[1,1,1,1,1,1,1,1,1,1,1,1,1,1,1,1]	=> 65535 (= 2^16-1)

By the [rules of binary number representation](#), we just need to compute the dot product between the document vector and the vector consisting of powers of 2:

```
1 doc = corpus[0, :] # first document
2 index_bits = (doc.dot(random_vectors) >= 0)
3 powers_of_two = (1 << np.arange(15, -1, -1))
4 print index_bits
5 print powers_of_two          # [32768, 16384, 8192, 4096, 2048, 1024, 512, 256, 128, 64, 32, 16,
     8, 4, 2, 1]
6 print index_bits.dot(powers_of_two)
```

Since it's the dot product again, we batch it with a matrix operation:

```
1 index_bits = corpus.dot(random_vectors) >= 0
2 print index_bits.dot(powers_of_two)
```

This array gives us the integer index of the bins for all documents.

Now we are ready to complete the following function. Given the integer bin indices for the documents, you should compile a list of document IDs that belong to each bin. Since a list is to be maintained for each unique bin index, a dictionary of lists is used.

1. Compute the integer bin indices. This step is already completed.

2. For each document in the dataset, do the following:

- Get the integer bin index for the document.
- Fetch the list of document ids associated with the bin; if no list yet exists for this bin, assign the bin an empty list.
- Add the document id to the end of the list.

```
1 def train_lsh(data, num_vector=16, seed=None):
2
3     dim = data.shape[1]
4     if seed is not None:
5         np.random.seed(seed)
6     random_vectors = generate_random_vectors(num_vector, dim)
7
8     powers_of_two = 1 << np.arange(num_vector-1, -1, -1)
9
10    table = []
11
12    # Partition data points into bins
13    bin_index_bits = (data.dot(random_vectors) >= 0)
14
15    # Encode bin index bits into integers
16    bin_indices = bin_index_bits.dot(powers_of_two)
17
18    # Update `table` so that `table[i]` is the list of document ids with bin index equal to i.
19    for data_index, bin_index in enumerate(bin_indices):
20        if bin_index not in table:
21            # If no list yet exists for this bin, assign the bin an empty list.
22            table[bin_index] = ... # YOUR CODE HERE
23            # Fetch the list of document ids associated with the bin and add the document id to the
24            # end.
25            ... # YOUR CODE HERE
26
27    model = {'data':
28             'bin_index_bits': bin_index_bits,
29             'bin_indices': bin_indices,
30             'table': table,
31             'random_vectors': random_vectors,
32             'num_vector': num_vector}
33
34 return model
```

Checkpoint.

```
1 model = train_lsh(corpus, num_vector=16, seed=143)
2 table = model['table']
3 if 0 in table and table[0] == [39583] and \
4     143 in table and table[143] == [19693, 28277, 29776, 30399]:
5     print 'Passed!'
6 else:
7     print 'Check your code.'
```

Note. We will be using the model trained here in the following sections, unless otherwise indicated.

Inspect bins

Let us look at some documents and see which bins they fall into.

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

Quiz Question. What is the document id of Barack Obama's article?

Quiz Question. Which bin contains Barack Obama's article? Enter its integer index.

Recall from the previous assignment that Joe Biden was a close neighbor of Barack Obama.

Quiz Question. Examine the bit representations of the bins containing Barack Obama and Joe Biden. In how many places do they agree?

Compare the result with a former British diplomat, whose bin representation agrees with Obama's in only 8 out of 16 places.

```
1 print wiki[wiki['name']=='Wynn Normington Hugh-Jones']
2
3 print np.array(model['bin_index_bits'][22745], dtype=int) # list of 0/1's
4 print model['bin_index_bits'][35817] == model['bin_index_bits'][22745]
```

How about the documents in the same bin as Barack Obama? Are they necessarily more similar to Obama than Biden? Let's look at which documents are in the same bin as the Barack Obama article.

```
1 print model['table'][model['bin_indices'][35817]]
```

There are four other documents that belong to the same bin. Which documents are they?

```
1 doc_ids = list(model['table'][model['bin_indices'][35817]])
2 doc_ids.remove(35817) # display documents other than Obama
3
4 docs = wiki.filter_by(values=doc_ids, column_name='id') # filter by id column
5 print docsIt turns out that Joe Biden is much closer to Barack Obama than any of the four
      documents, even though Biden's bin representation differs from Obama's by 2 bits.
```

It turns out that Joe Biden is much closer to Barack Obama than any of the four documents, even though Biden's bin representation differs from Obama's by 2 bits.

```
1 def cosine_distance(x, y):
2     xy = x.dot(y.T)
3     dist = xy/(norm(x)*norm(y))
4     return 1-dist[0,0]
5
6 obama_tf_idf = corpus[35817,:]
7 biden_tf_idf = corpus[24478,:]
8
9 print '===== Cosine distance from Barack Obama'
10 print 'Barack Obama - {0:24s}: {1:f}'.format('Joe Biden',
11                                              cosine_distance(obama_tf_idf, biden_tf_idf))
12 for doc_id in doc_ids:
13     doc_tf_idf = corpus[doc_id,:]
14     print 'Barack Obama - {0:24s}: {1:f}'.format(wiki[doc_id]['name'],
15                                              cosine_distance(obama_tf_idf, doc_tf_idf))
```

Moral of the story. Similar data points will in general *tend* to fall into *nearby* bins, but that's all we can say about LSH. In a high-dimensional space such as text features, we often get unlucky with our selection of only a few random vectors such that dissimilar data points go into the same bin while similar data points fall into different bins. **Given a query document, we must consider all documents in the nearby bins and sort them according to their actual distances from the query.**

Query the LSH model

Let us first implement the logic for searching nearby neighbors, which goes like this:

```
1 1. Let L be the bit representation of the bin that contains the query documents.
2 2. Consider all documents in bin L.
3 3. Consider documents in the bins whose bit representation differs from L by 1 bit.
4 4. Consider documents in the bins whose bit representation differs from L by 2 bits.
5 ...
```

To obtain candidate bins that differ from the query bin by some number of bits, we use `itertools.combinations`, which produces all possible subsets of a given list. See [this documentation](#) for details.

- 1 1. Decide on the search radius r . This will determine the number of different bits between the two vectors.
- 2 2. For each subset (n_1, n_2, \dots, n_r) of the list $[0, 1, 2, \dots, \text{num_vector}-1]$, do the following:
 - 3 * Flip the bits (n_1, n_2, \dots, n_r) of the query bin to produce a new bit vector.
 - 4 * Fetch the list of documents belonging to the bin indexed by the new bit vector.
 - 5 * Add those documents to the candidate set.

Each line of output from the following cell is a 3-tuple indicating where the candidate bin would differ from the query bin. For instance,

```
1 (0, 1, 3)
```

indicates that the candidate bin differs from the query bin in first, second, and fourth bits. For illustrations, inspect the output of the following piece of code:

```
1 from itertools import combinations
2
3 num_vector = 16
4 search_radius = 3
5
6 for diff in combinations(range(num_vector), search_radius):
7     print diff
```

With this output in mind, implement the logic for nearby bin search:

```
1 def search_nearby_bins(query_bin_bits, table, search_radius=2, initial_candidates=set()):
2     """
3         For a given query vector and trained LSH model, return all candidate neighbors for
4         the query among all bins within the given search radius.
5
6     Example usage
7     -----
8     >>> model = train_lsh(corpus, num_vector=16, seed=143)
9     >>> q = model['bin_index_bits'][0] # vector for the first document
10
11    >>> candidates = search_nearby_bins(q, model['table'])
12    """
13    num_vector = len(query_bin_bits)
14    powers_of_two = 1 << np.arange(num_vector-1, -1, -1)
15
16    # Allow the user to provide an initial set of candidates.
17    candidate_set = copy(initial_candidates)
18
19    for different_bits in combinations(range(num_vector), search_radius):
20        # Flip the bits (n_1,n_2,...,n_r) of the query bin to produce a new bit vector.
21        ## Hint: you can iterate over a tuple like a list
22        alternate_bits = copy(query_bin_bits)
23        for i in different_bits:
24            alternate_bits[i] = ... # YOUR CODE HERE
25
26        # Convert the new bit vector to an integer index
27        nearby_bin = alternate_bits.dot(powers_of_two)
28
29        # Fetch the list of documents belonging to the bin indexed by the new bit vector.
30        # Then add those documents to candidate_set
31        # Make sure that the bin exists in the table!
32        # Hint: update() method for sets lets you add an entire list to the set
33        if nearby_bin in table:
34            ... # YOUR CODE HERE: Update candidate_set with the documents in this bin.
35
36    return candidate_set
```

Checkpoint. Running the function with `search_radius=0` should yield the list of documents belonging to the same bin as the query.

```
1 obama_bin_index = model['bin_index_bits'][35817] # bin index of Barack Obama
2 candidate_set = search_nearby_bins(obama_bin_index, model['table'], search_radius=0)
3 if candidate_set == set([35817, 21426, 53937, 39426, 50261]):
4     print 'Passed test'
5 else:
6     print 'Check your code'
7 print 'List of documents in the same bin as Obama: 35817, 21426, 53937, 39426, 50261'
```

Checkpoint. Running the function with search_radius=1 adds more documents to the fore.

```
1 candidate_set = search_nearby_bins(obama_bin_index, model['table'], search_radius=1,
2     initial_candidates=candidate_set)
3 if candidate_set == set([39426, 38155, 38412, 28444, 9757, 41631, 39207, 59050, 47773, 53937,
4     21426, 34547,
5     23229, 55615, 39877, 27404, 33996, 21715, 50261, 21975, 33243, 58723,
6     35817, 45676,
7     19699, 2804, 20347]):
8     print 'Passed test'
9 else:
10    print 'Check your code'
```

Note. Don't be surprised if few of the candidates look similar to Obama. This is why we add as many candidates as our computational budget allows and sort them by their distance to the query.

Now we have a function that can return all the candidates from neighboring bins. Next we write a function to collect all candidates and compute their true distance to the query.

```
1 def query(vec, model, k, max_search_radius):
2
3     data = model['data']
4     table = model['table']
5     random_vectors = model['random_vectors']
6     num_vector = random_vectors.shape[1]
7
8
9     # Compute bin index for the query vector, in bit representation.
10    bin_index_bits = (vec.dot(random_vectors) >= 0).flatten()
11
12    # Search nearby bins and collect candidates
13    candidate_set = set()
14    for search_radius in xrange(max_search_radius+1):
15        candidate_set = search_nearby_bins(bin_index_bits, table, search_radius,
16                                           initial_candidates=candidate_set)
17
18    # Sort candidates by their true distances from the query
19    nearest_neighbors = sframe.SFrame({'id':candidate_set})
20    candidates = data[np.array(list(candidate_set)),:]
21    nearest_neighbors['distance'] = pairwise_distances(candidates, vec, metric='cosine').flatten()
22
23    return nearest_neighbors.topk('distance', k, reverse=True), len(candidate_set)
```

Let's try it out with Obama:

```
1 print query(corpus[35817,:], model, k=10, max_search_radius=3)
```

To identify the documents, it's helpful to join this table with the Wikipedia table:

```
1 result, num_candidates_considered = query(corpus[35817,:], model, k=10, max_search_radius=3)
2 print result.join(wiki[['id', 'name']], on='id').sort('distance')
```

which produces a table of the form

1	id	distance	name
2			
3			
4	35817	-6.66133814775e-16	Barack Obama
5	24478	0.703138676734	Joe Biden
6	56008	0.856848127628	Nathan Cullen
7	37199	0.874668698194	Barry Sullivan (lawyer)
8	40353	0.890034225981	Neil MacBride
9	9267	0.898377208819	Vikramaditya Khanna
10	55909	0.899340396322	Herman Cain
11	9165	0.900921029925	Raymond F. Clevenger
12	57958	0.903003263483	Michael J. Malbin
13	49872	0.909532800353	Lowell Barron
14			
15	[10 rows x 3 columns]		

We have shown that we have a working LSH implementation!

Experimenting with your LSH implementation

In the following sections we have implemented a few experiments so that you can gain intuition for how your LSH implementation behaves in different situations. This will help you understand the effect of searching nearby bins and the performance of LSH versus computing nearest neighbors using a brute force search.

Effect of nearby bin search

How does nearby bin search affect the outcome of LSH? There are three variables that are affected by the search radius:

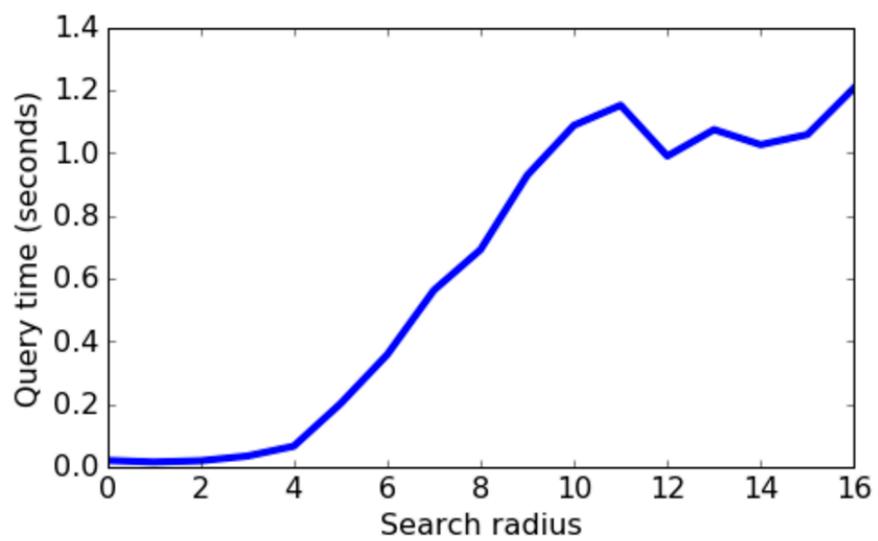
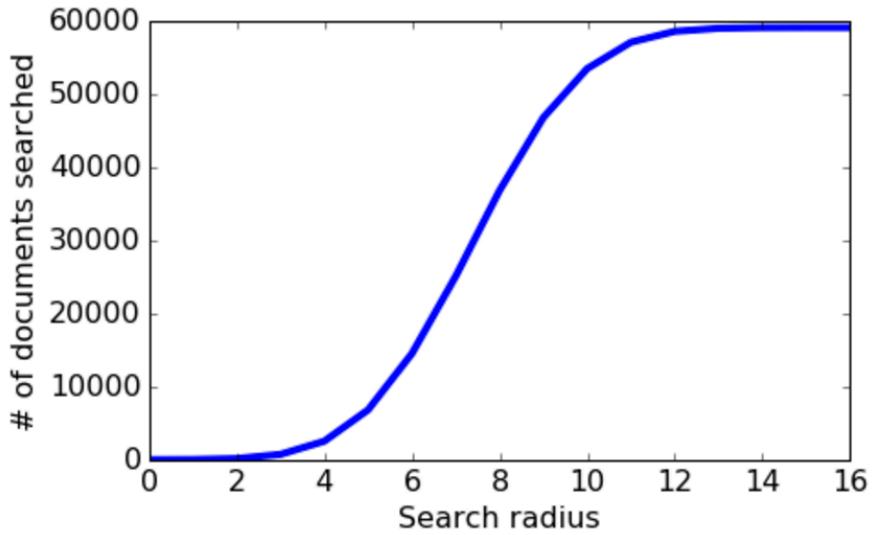
- Number of candidate documents considered
- Query time
- Distance of approximate neighbors from the query

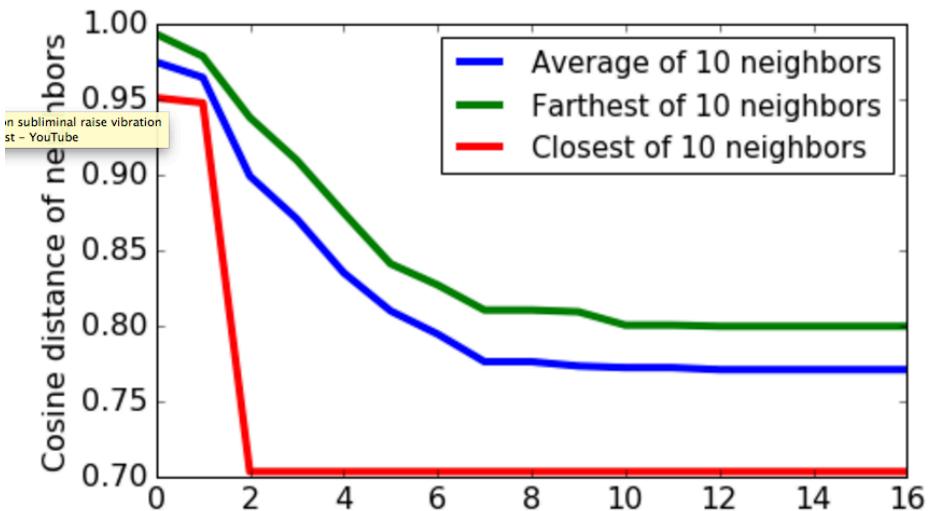
Let us run LSH multiple times, each with different radii for nearby bin search. We will measure the three variables as discussed above.

```
1 num_candidates_history = []
2 query_time_history = []
3 max_distance_from_query_history = []
4 min_distance_from_query_history = []
5 average_distance_from_query_history = []
6
7 for max_search_radius in xrange(17):
8     start=time.time()
9     # Perform LSH query using Barack Obama, with max_search_radius
10    result, num_candidates = query(corpus[35817,:], model, k=10,
11                                     max_search_radius=max_search_radius)
12    end=time.time()
13    query_time = end-start # Measure time
14
15    print 'Radius:', max_search_radius
16    # Display 10 nearest neighbors, along with document ID and name
17    print result.join(wiki[['id', 'name']], on='id').sort('distance')
18
19    # Collect statistics on 10 nearest neighbors
20    average_distance_from_query = result['distance'][1:].mean()
21    max_distance_from_query = result['distance'][1:].max()
22    min_distance_from_query = result['distance'][1:].min()
23
24    num_candidates_history.append(num_candidates)
25    query_time_history.append(query_time)
26    average_distance_from_query_history.append(average_distance_from_query)
27    max_distance_from_query_history.append(max_distance_from_query)
28    min_distance_from_query_history.append(min_distance_from_query)
```

Notice that the top 10 query results become more relevant as the search radius grows. Let's plot the three variables:

```
1 plt.figure(figsize=(7,4.5))
2 plt.plot(num_candidates_history, linewidth=4)
3 plt.xlabel('Search radius')
4 plt.ylabel('# of documents searched')
5 plt.rcParams.update({'font.size':16})
6 plt.tight_layout()
7
8 plt.figure(figsize=(7,4.5))
9 plt.plot(query_time_history, linewidth=4)
10 plt.xlabel('Search radius')
11 plt.ylabel('Query time (seconds)')
12 plt.rcParams.update({'font.size':16})
13 plt.tight_layout()
14
15 plt.figure(figsize=(7,4.5))
16 plt.plot(average_distance_from_query_history, linewidth=4, label='Average of 10 neighbors')
17 plt.plot(max_distance_from_query_history, linewidth=4, label='Farthest of 10 neighbors')
18 plt.plot(min_distance_from_query_history, linewidth=4, label='Closest of 10 neighbors')
19 plt.xlabel('Search radius')
20 plt.ylabel('Cosine distance of neighbors')
21 plt.legend(loc='best', prop={'size':15})
22 plt.rcParams.update({'font.size':16})
23 plt.tight_layout()
```





Some observations:

- As we increase the search radius, we find more neighbors that are a smaller distance away.
- With increased search radius comes a greater number documents that have to be searched. Query time is higher as a consequence.
- With sufficiently high search radius, the results of LSH begin to resemble the results of brute-force search.

Quiz Question. What was the smallest search radius that yielded the correct nearest neighbor, namely Joe Biden?

Quiz Question. Suppose our goal was to produce 10 approximate nearest neighbors whose average distance from the query document is within 0.01 of the average for the true 10 nearest neighbors. For Barack Obama, the true 10 nearest neighbors are on average about 0.77. What was the smallest search radius for Barack Obama that produced an average distance of 0.78 or better?

Quality metrics for neighbors

The above analysis is limited by the fact that it was run with a single query, namely Barack Obama. We should repeat the analysis for the entirety of data. Iterating over all documents would take a long time, so let us randomly choose 10 documents for our analysis.

For each document, we first compute the true 25 nearest neighbors, and then run LSH multiple times. We look at two metrics:

- Precision@10: How many of the 10 neighbors given by LSH are among the true 25 nearest neighbors?
- Average cosine distance of the neighbors from the query

Then we run LSH multiple times with different search radii.

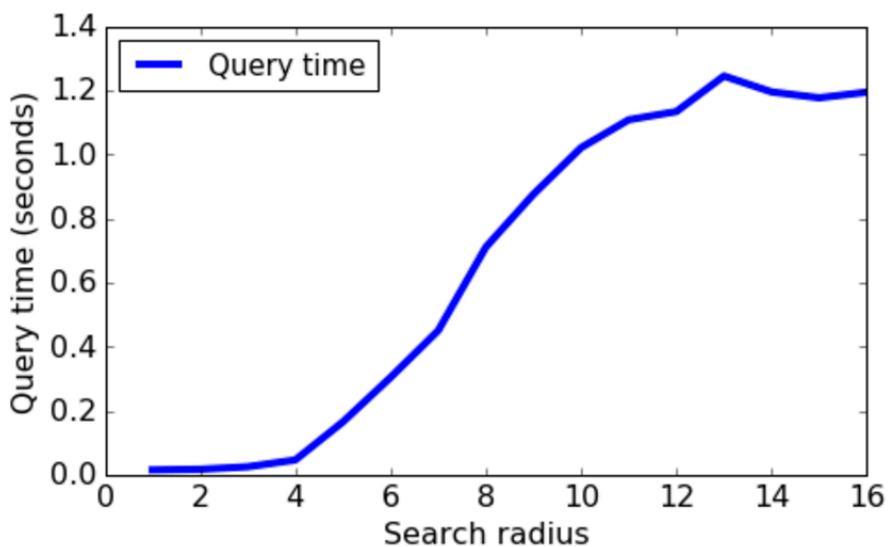
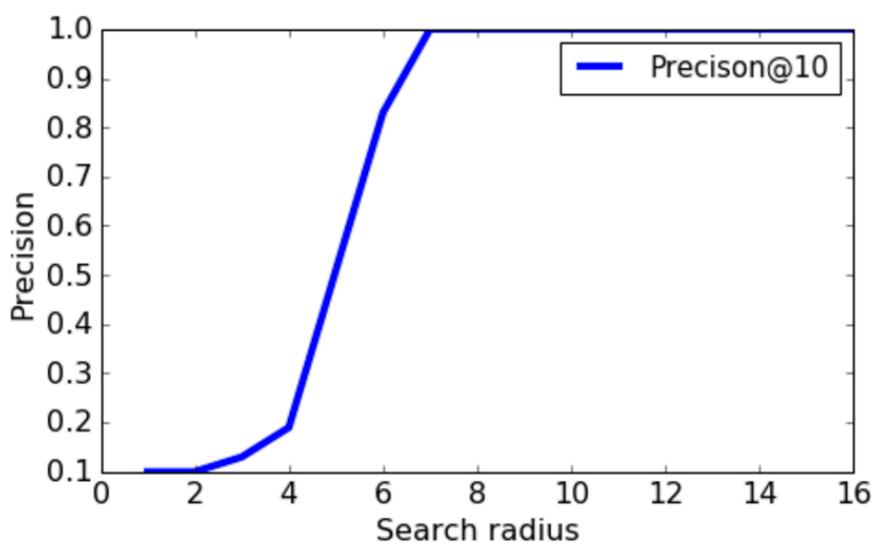
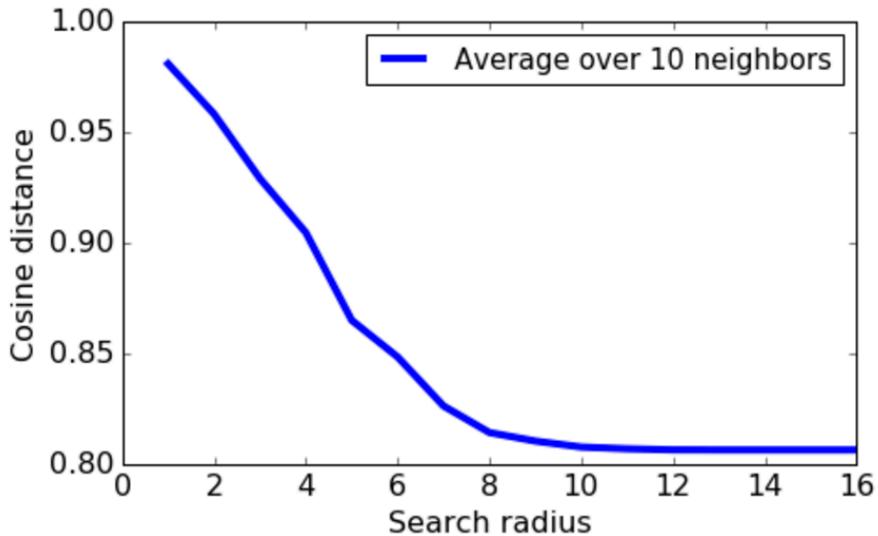
```
1 def brute_force_query(vec, data, k):
2     num_data_points = data.shape[0]
3
4     # Compute distances for ALL data points in training set
5     nearest_neighbors = sframe.SFrame({'id':range(num_data_points)})
6     nearest_neighbors['distance'] = pairwise_distances(data, vec, metric='cosine').flatten()
7
8     return nearest_neighbors.topk('distance', k, reverse=True)
```

The following cell will run LSH with multiple search radii and compute the quality metrics for each run. Allow a few minutes to complete.

```
1 max_radius = 17
2 precision = {i:[] for i in xrange(max_radius)}
3 average_distance = {i:[] for i in xrange(max_radius)}
4 query_time = {i:[] for i in xrange(max_radius)}
5
6 np.random.seed(0)
7 num_queries = 10
8 for i, ix in enumerate(np.random.choice(corpus.shape[0], num_queries, replace=False)):
9     print('%s / %s' % (i, num_queries))
10    ground_truth = set(brute_force_query(corpus[ix,:], corpus, k=25)['id'])
11    # Get the set of 25 true nearest neighbors
12
13    for r in xrange(1,max_radius):
14        start = time.time()
15        result, num_candidates = query(corpus[ix,:], model, k=10, max_search_radius=r)
16        end = time.time()
17
18        query_time[r].append(end-start)
19        # precision = (# of neighbors both in result and ground_truth)/10.0
20        precision[r].append(len(set(result['id']) & ground_truth)/10.0)
21        average_distance[r].append(result['distance'][1:].mean())
```

Let's plot:

```
1 plt.figure(figsize=(7,4.5))
2 plt.plot(range(1,17), [np.mean(average_distance[i]) for i in xrange(1,17)], linewidth=4, label
   ='Average over 10 neighbors')
3 plt.xlabel('Search radius')
4 plt.ylabel('Cosine distance')
5 plt.legend(loc='best', prop={'size':15})
6 plt.rcParams.update({'font.size':16})
7 plt.tight_layout()
8
9 plt.figure(figsize=(7,4.5))
10 plt.plot(range(1,17), [np.mean(precision[i]) for i in xrange(1,17)], linewidth=4, label
   ='Precision@10')
11 plt.xlabel('Search radius')
12 plt.ylabel('Precision')
13 plt.legend(loc='best', prop={'size':15})
14 plt.rcParams.update({'font.size':16})
15 plt.tight_layout()
16
17 plt.figure(figsize=(7,4.5))
18 plt.plot(range(1,17), [np.mean(query_time[i]) for i in xrange(1,17)], linewidth=4, label='Query
   time')
19 plt.xlabel('Search radius')
20 plt.ylabel('Query time (seconds)')
21 plt.legend(loc='best', prop={'size':15})
22 plt.rcParams.update({'font.size':16})
23 plt.tight_layout()
```



The observations for Barack Obama generalize to the entire dataset.

Effect of number of random vectors

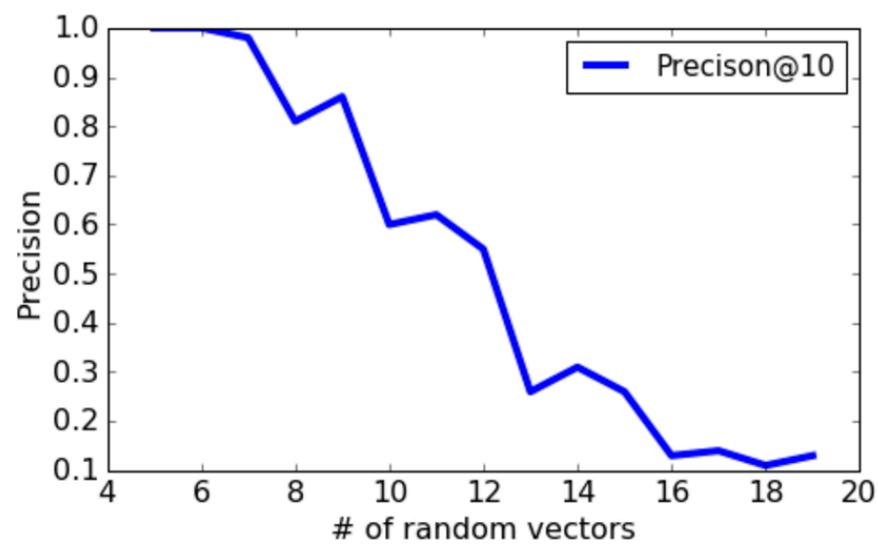
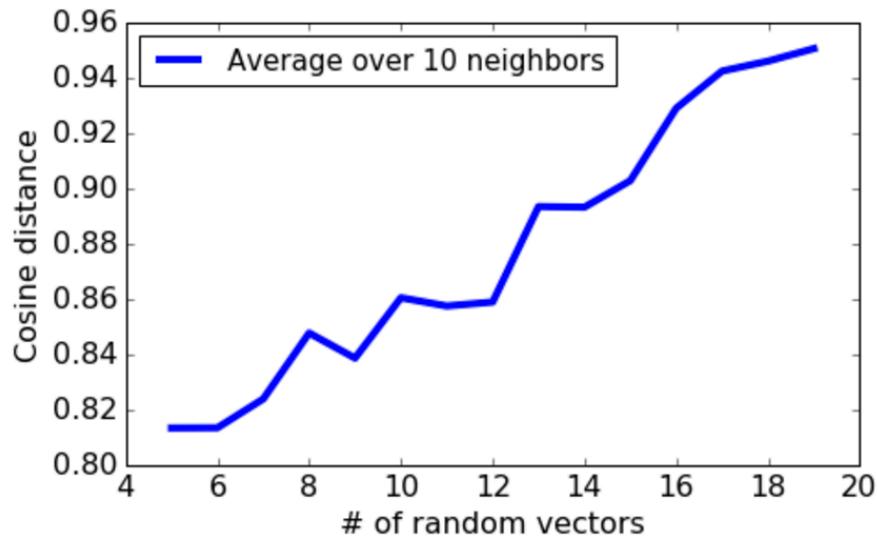
Let us now turn our focus to the remaining parameter: the number of random vectors. We run LSH with different number of random vectors, ranging from 5 to 20. We fix the search radius to 3.

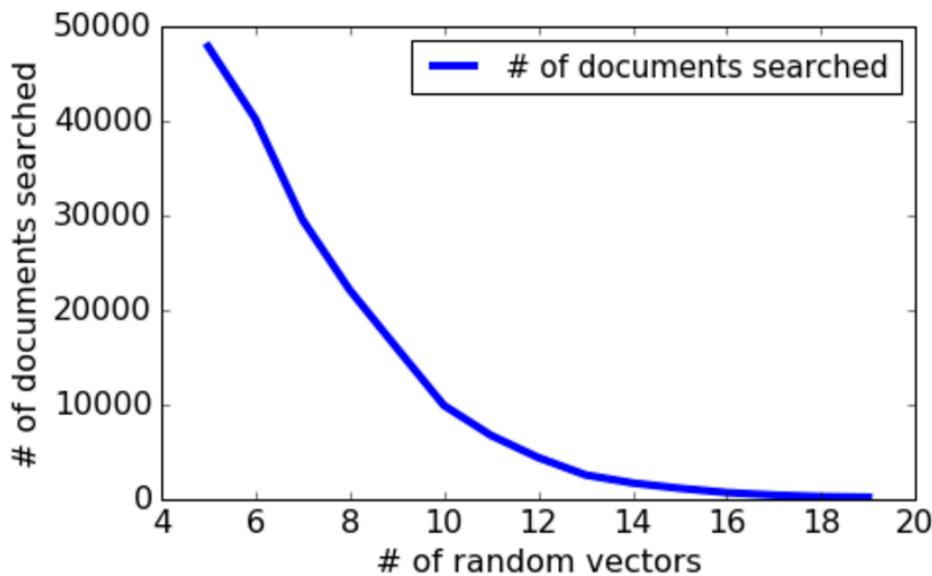
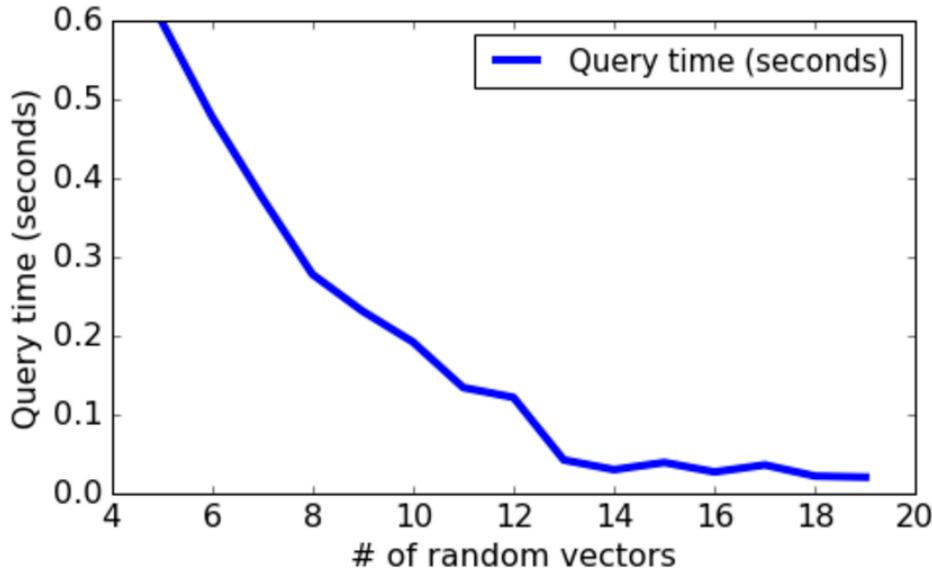
Allow a few minutes for the following cell to complete.

```
1 precision = {i:{} for i in xrange(5,20)}
2 average_distance = {i:{} for i in xrange(5,20)}
3 query_time = {i:{} for i in xrange(5,20)}
4 num_candidates_history = {i:{} for i in xrange(5,20)}
5 ground_truth = []
6
7 np.random.seed(0)
8 num_queries = 10
9 docs = np.random.choice(corpus.shape[0], num_queries, replace=False)
10
11 for i, ix in enumerate(docs):
12     ground_truth[ix] = set(brute_force_query(corpus[ix,:], corpus, k=25)['id'])
13     # Get the set of 25 true nearest neighbors
14
15 for num_vector in xrange(5,20):
16     print('num_vector = %s' % (num_vector))
17     model = train_lsh(corpus, num_vector, seed=143)
18
19 for i, ix in enumerate(docs):
20     start = time.time()
21     result, num_candidates = query(corpus[ix,:], model, k=10, max_search_radius=3)
22     end = time.time()
23
24     query_time[num_vector].append(end-start)
25     precision[num_vector].append(len(set(result['id'])) & ground_truth[ix])/10.0
26     average_distance[num_vector].append(result['distance'][1:].mean())
27     num_candidates_history[num_vector].append(num_candidates)
```

Plot the metrics as a function of the number of random vectors using the following piece of code :

```
1 plt.figure(figsize=(7,4.5))
2 plt.plot(range(5,20), [np.mean(average_distance[i]) for i in xrange(5,20)], linewidth=4, label
   ='Average over 10 neighbors')
3 plt.xlabel('# of random vectors')
4 plt.ylabel('Cosine distance')
5 plt.legend(loc='best', prop={'size':15})
6 plt.rcParams.update({'font.size':16})
7 plt.tight_layout()
8
9 plt.figure(figsize=(7,4.5))
10 plt.plot(range(5,20), [np.mean(precision[i]) for i in xrange(5,20)], linewidth=4, label
   ='Precision@10')
11 plt.xlabel('# of random vectors')
12 plt.ylabel('Precision')
13 plt.legend(loc='best', prop={'size':15})
14 plt.rcParams.update({'font.size':16})
15 plt.tight_layout()
16
17 plt.figure(figsize=(7,4.5))
18 plt.plot(range(5,20), [np.mean(query_time[i]) for i in xrange(5,20)], linewidth=4, label='Query
   time (seconds)')
19 plt.xlabel('# of random vectors')
20 plt.ylabel('Query time (seconds)')
21 plt.legend(loc='best', prop={'size':15})
22 plt.rcParams.update({'font.size':16})
23 plt.tight_layout()
24
25 plt.figure(figsize=(7,4.5))
26 plt.plot(range(5,20), [np.mean(num_candidates_history[i]) for i in xrange(5,20)], linewidth=4,
   label='# of documents searched')
27 plt.xlabel('# of random vectors')
28 plt.ylabel('# of documents searched')
29 plt.legend(loc='best', prop={'size':15})
30 plt.rcParams.update({'font.size':16})
31 plt.tight_layout()
32 plt.tight_layout()
```





We see a similar trade-off between quality and performance: as the number of random vectors increases, the query time goes down as each bin contains fewer documents on average, but on average the neighbors are likewise placed farther from the query. On the other hand, when using a small enough number of random vectors, LSH becomes very similar brute-force search: Many documents appear in a single bin, so searching the query bin alone covers a lot of the corpus; then, including neighboring bins might result in searching all documents, just as in the brute-force approach.