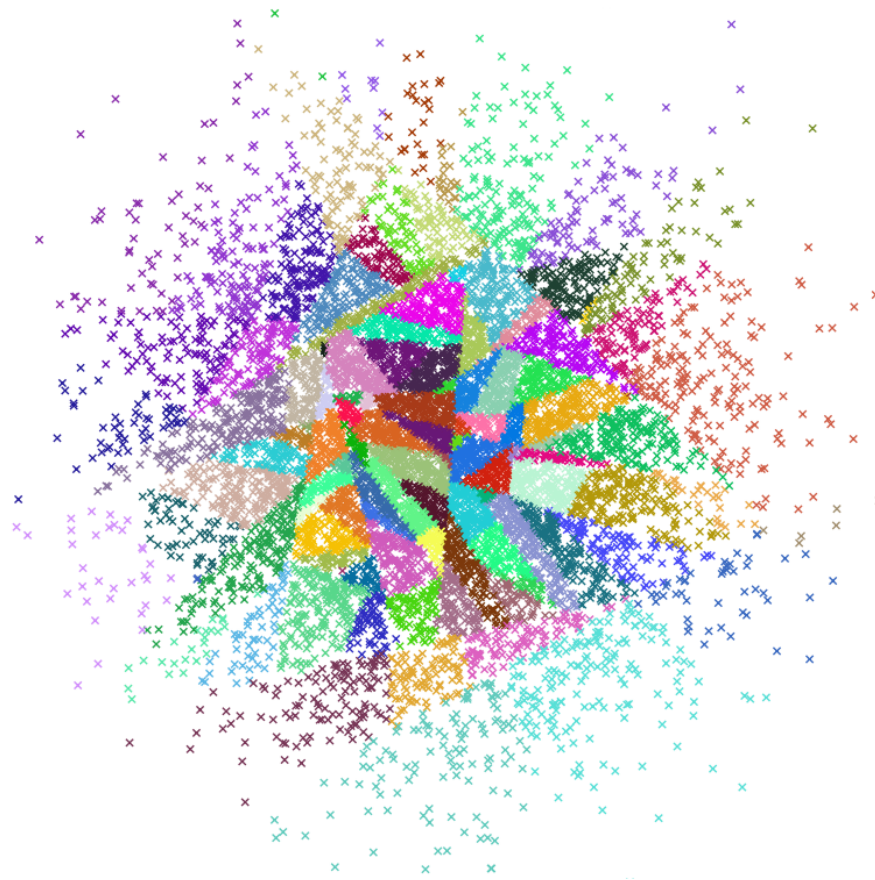


Simple Approximate Nearest Neighbors in Python with Annoy and lmdb



Kevin Yang [Follow](#)

Jan 4, 2018 · 4 min read



Colorful points in space

Recently, I've been playing around with adding and subtracting word embeddings learned from GloVe. For example, we can take the embedding for the word “king” subtract the embedding for “man” and then add the embedding for “woman” and end up with a resulting embedding (vector). Next, if we have a corpus of these word-embedding pairs we can search through it to find the most similar embedding and retrieve the corresponding word. If we did this for the query above we would get:

King + (Woman - Man) = Queen

There are many ways to search through this corpus of word-embedding pairs for the nearest neighbors to a query. One way to guarantee that you find the optimal vector is to iterate through your corpus and compare how similar each pair is to the query. This however is incredibly time consuming and not often used. A better exact technique would be to use a vectorized cosine distance shown below:

```
1 vectors = np.array(embeddingmodel.embeddings)
2 ranks = np.dot(query,vectors.T)/np.sqrt(np.sum(vectors*
3 mostSimilar = []
4 [mostSimilar.append(idx) for idx in ranks.argsort()[::-1]]
```

For more info on cosine distance check out this great resource: [Cosine Similarity](#)

Vectorized cosine distance is notably faster than the iterative method but still may be too slow. This is where approximate nearest neighbors shines: returning approximate results but blazingly quickly. Many times you don't need exact optimal results, for example: what even is the exact similar word for "Queen?" In these cases where you need good enough results quickly, you should use approximate nearest neighbors.

In this article we will be writing a simple python script to quickly find approximate nearest neighbors. We will use the Python library Annoy and lmbd. For my corpus, I will be using word-embedding pairs, but the instructions below work with any type of embeddings: such as embeddings of songs to create a recommendation engine for similar songs or even photo embeddings to enable reverse image search.

Making an Annoy Index

Lets create a python script called: "make_annoy_index." To start off lets include the dependencies that we will use:

```

1  '''
2  Usage: python2 make_annoy_index.py \
3      --embeddings=<embedding path> \
4      --num_trees=<int> \
5      --verbose
6
7  Generate an Annoy index and lmbd map given an embedding
8
9  Embedding file can be
10     1. A .bin file that is compatible with word2vec binary
11        There are pre-trained vectors to download at http://
12     2. A .gz file with the GloVe format (item then a list of
13     3. A plain text file with the same format as above
14
15  '''
16

```

The last line we have is an import from “vector_utils”. We will be writing “vector_utils” later on so don’t worry now!

Next, on to the meat of our script: the “create_index” function. Here we will be generating our lmbd map and our Annoy index.

1. First we find the length of our embedding which is used to instantiate an Annoy index.
2. Next we instantiate an lmbd map with the line: “env = lmbd.open(fn_lmbd, map_size=int(1e9)).”
3. Make sure we don’t have an Annoy index or lmbd map in the current path
4. Add every key and vector from our embeddings file to both our lmbd map and our Annoy index
5. Build and save our Annoy index

```

1  '''
2  function create_index(fn, num_trees=30, verbose=False)
3  -----
4  Creates an Annoy index and lmbd map given an embedding
5
6  Input:
7      fn            - filename of the embedding file
8      num_trees     - number of trees to build Annoy index
9      verbose       - log status
10
11 Return:
12     Void
13 '''
14 def create_index(fn, num_trees=30, verbose=False):
15     fn_annoy = fn + '.annoy'
16     fn_lmbd = fn + '.lmbd' # stores word <-> id mapping
17
18     word, vec = get_vectors(fn).next()
19     size = len(vec)
20     if verbose:
21         print("Vector size: {}".format(size))
22
23     env = lmbd.open(fn_lmbd, map_size=int(1e9))
24     if not os.path.exists(fn_annoy) or not os.path.exists(fn_lmbd):
25         i = 0
26         a = annoy.AnnoyIndex(size)
27         with env.begin(write=True) as txn:
28             for word, vec in get_vectors(fn):
29                 a.add_item(i, vec)
30                 id = 'i%d' % i
31                 word = 'w' + word

```

I've included "argparse" so we can call our script from the command line:

```

1  '''
2  private function _create_args()
3  -----
4  Creates an argparse object for CLI for create_index()
5
6  Input:
7      Void
8
9  Return:
10     args object with required arguments for threshold_
11
12  '''
13  def _create_args():
14     parser = argparse.ArgumentParser()

```

Add a main function to call our script and we are done with “make_annoy_index.py”:

```

1  if __name__ == '__main__':
2     args = _create_args()
3     create_index(args.embeddings, num_trees=args.num_tr

```

Now we can just call our new script from the command line to generate an Annoy index and a corresponding lmdb map!

```

python2 make_annoy_index.py \
    --embeddings=<embedding path> \
    --num_trees=<int> \
    --verbose

```

Writing Vector Utils

In “make_annoy_index.py” we included a python script “vector_utils.” We will write this script now. “Vector_utils” is used to help read in vectors from .txt, .bin, and .pkl files.

Writing this script is not that relevant to what we are doing so I’ve just included the whole script below:

```

1  '''
2  Vector Utils
3
4  Utils to read in vectors from txt, .bin, or .pkl.
5
6  Taken from Erik Bernhardsson
7  Source: https://github.com/erikbern/ann-presentation/blob/master/
8  '''
9  import gzip
10 import struct
11 import cPickle
12
13 def _get_vectors(fn):
14     if fn.endswith('.gz'):
15         f = gzip.open(fn)
16         fn = fn[:-3]
17
18     else:
19         f = open(fn)
20
21     if fn.endswith('.bin'): # word2vec format
22         words, size = (int(x) for x in f.readline().split())
23
24         t = 'f' * size
25
26         while True:
27             pos = f.tell()
28             buf = f.read(1024)
29             if buf == '' or buf == '\n': return
30             i = buf.index(' ')
31             word = buf[:i]
32             f.seek(pos + i + 1)
33
34             vec = struct.unpack(t, f.read(4 * size))
35
36             yield word.lower(), vec
37

```

Testing our Index and lmdb Map

Now that we have generated an Annoy index and lmdb map lets write a script to use them to do inference.

Name our file “annoy_inference.py” and include the dependencies:

```

1  '''
2  Usage: python2 annoy_inference.py \
3      --token='hello' \
4      --num_results=<int> \
5      --verbose
6
7  Query an Annoy index to find approximate nearest neighbors
8
9  '''

```

Now we need to load in our Annoy index and our lmdb map; we will load them globally for easy access. Note that for me “VEC_LENGTH” is 50. Make sure that your “VEC_LENGTH” matches the length of your embedding, otherwise Annoy will not be happy.

```

1  VEC_LENGTH = 50
2  FN_ANNNOY = 'glove.6B.50d.txt.annoy'
3  FN_LMDB = 'glove.6B.50d.txt.lmdb'
4
5  a = annoy.AnnoyIndex(VEC_LENGTH)
6  a.load(FN_ANNNOY)

```

The fun part, the “calculate” function.

1. Get the index for our query from our lmdb map
2. Get the corresponding vector from Annoy with “get_item_vector(id)”
3. Get the nearest neighbors from Annoy with “a.get_nns_by_vector(v, num_results)”

```

1  '''
2  private function calculate(query, num_results)
3  -----
4  Queries a given Annoy index and lmdb map for num_results
5
6  Input:
7      query          - query to be searched
8      num_results    - the number of results
9
10 Return:
11     ret_keys        - list of num_results nearest neighbors
12
13 '''
14 def calculate(query, num_results, verbose=False):
15     ret_keys = []
16     with env.begin() as txn:
17         id = int(txn.get('w' + query)[1:])
18         ...

```

Again, here is an “argparse” object to make reading command line arguments easier

```

1  '''
2  private function _create_args()
3  -----
4  Creates an argparse object for CLI for calculate() function
5
6  Input:
7      Void
8
9  Return:
10     args object with required arguments for threshold_
11
12 '''
13 def _create_args():
14     parser = argparse.ArgumentParser()

```

And, a main function to call “annoy_inference.py” from the command line:


```
1  if __name__ == '__main__':  
2      args = _create_args()  
3      print(calculate(args.token, args.num_results, args.
```

Great! Now we can use our Annoy index and lmdb map to get the nearest neighbors to a query!

```
python2 annoy_inference.py --token="test" --num_results=30  
  
['test', 'tests', 'determine', 'for', 'crucial', 'only',  
'preparation', 'needed', 'positive', 'guided', 'time',  
'performance', 'one', 'fitness', 'replacement', 'stages',  
'made', 'both', 'accuracy', 'deliver', 'put',  
'standardized', 'best', 'discovery', '.', 'a', 'diagnostic',  
'delayed', 'while', 'side']
```

Code

You can clone the corresponding repo to get all the code in this tutorial:

https://github.com/kyang6/annoy_tutorial

Conclusion

Boy, that was a lot of code. But, now we have scripts to take any embedding data and generate an Annoy index and lmdb map to find approximate nearest neighbors!

If you have questions or concerns please email me!

kyang6@cs.stanford.edu