

Seminar 1: Fun with Word Embeddings (3 points)

Today we gonna play with word embeddings; train our own little embedding, load one from gensim model zoo and use it to visualize text corpora.

This whole thing is gonna happen on top of embedding dataset.

Requirements: `pip install --upgrade nltk gensim bokeh`, but only if you're running locally.

```
In [ ]: # download the data:
!wget https://www.dropbox.com/s/obaitrix9jyu84r/quora.txt?dl=1 -O ./quora.txt
# alternative download link: https://yadi.sk/i/BP0du1NaTdUw

--2023-09-24 02:07:54-- https://www.dropbox.com/s/obaitrix9jyu84r/quora.txt?dl=1
Resolving www.dropbox.com (www.dropbox.com)... 162.125.70.18, 2020:100:6026:18::a27d:4612
Connecting to www.dropbox.com (www.dropbox.com)|162.125.70.18|:443... connected.
HTTP request sent, awaiting response... 302 Found
Location: /s/dl/obaitrix9jyu84r/quora.txt [following]
--2023-09-24 02:07:55-- https://www.dropbox.com/s/dl/obaitrix9jyu84r/quora.txt
Reusing existing connection to www.dropbox.com:443.
HTTP request sent, awaiting response... 302 Found
Location: /s/dl/obaitrix9jyu84r/quora.txt [following]
Location: https://ucb74a57ee744538b66124a6b9ac.dl.dropboxusercontent.com/cd/0/get/CET8E5mQLk67mR6x3r1cBaFU91LaZe9gKPxFzFDHDTbTn6Pqr2N9ZXMUTePgWkN3o3j7592A-g18F16qfMbmhVLCm-6fj2AU3jH9K2y-T09u3j2-H9WQ0aU5FwmxvY7f1e7d1e1 [following]
--2023-09-24 02:07:55-- https://ucb74a57ee744538b66124a6b9ac.dl.dropboxusercontent.com/cd/0/get/CET8E5mQLk67mR6x3r1cBaFU91LaZe9gKPxFzFDHDTbTn6Pqr2N9ZXMUTePgWkN3o3j7592A-g18F16qfMbmhVLCm-6fj2AU3jH9K2y-T09u3j2-H9WQ0aU5FwmxvY7f1e7d1e1
hVLCm-6fj2AU3jH9K2y-T09u3j2-H9WQ0aU5FwmxvY7f1e7d1e1
Connecting to ucB74a57ee744538b66124a6b9ac.dl.dropboxusercontent.com (ucB74a57ee744538b66124a6b9ac.dl.dropboxusercontent.com)|162.125.70.15|:443... connected.
HTTP request sent, awaiting response... 200 OK
Length: 33813983 (32M) application/binary
Saving to: './quora.txt'

./quora.txt          100%[=====] 32.25M  6.54MB/s   in 4.8s

2023-09-24 02:08:00 (8.15 MB/s) = './quora.txt' saved [33813983/33813983]

In [ ]: import numpy as np

with open("./quora.txt", encoding="utf-8") as file:
    data = list(file)

data[50]

Out[ ]: "What TV shows or books help you read people's body language?"

Tokenization: a typical first step for a nlp task is to split raw data into words. The text we're working with is in raw format: with all the punctuation and smiles attached to some words, so a simple str.split won't do.

Let's use nltk - a library that handles many nlp tasks like tokenization, stemming or part-of-speech tagging.

In [ ]: from nltk.tokenize import WordPunctTokenizer
tokenizer = WordPunctTokenizer()

print(tokenizer.tokenize(data[50]))

['what', 'TV', 'shows', 'or', 'books', 'help', 'you', 'read', 'people', "'", 's', 'body', ',', 'language', ',', '?']

In [ ]: # TASK: lowercase everything and extract tokens with tokenizer.
# data_tok should be a list of lists of tokens for each line in data.

data_tok = [tokenizer.tokenize(x.lower()) for x in data]

In [ ]: assert all(isinstance(row, (list, tuple)) for row in data_tok), "please convert each line into a list of tokens (strings)"
assert all(all(isinstance(tok, str) for tok in row) for row in data_tok), "please convert each line into a list of tokens (strings)"
is_list_in_tuple = lambda tok: isinstance(tok, list) or isinstance(tok, tuple)
assert all(map(lambda l: not is_list_in(l) or l.islower(), map(lambda tok, l: (tok, l), data_tok)))

In [ ]: print(' '.join(row) for row in data_tok[:2]))

['can i get back with my ex even though she is pregnant with another guy 's baby?', 'what are some ways to overcome a fast food addiction?']

Word vectors: as the saying goes, there's more than one way to train word embeddings. There's Word2Vec and GloVe with different objective functions. Then there's fasttext that uses character-level models to train word embeddings.

The choice is huge, so let's start someplace small: gensim is another nlp library that features many vector-based models including word2vec.
```

```
In [ ]: from gensim.models import Word2Vec
model = Word2Vec(data_tok,
                 vector_size=32,          # embedding vector size
                 min_count=5,            # consider words that occurred at least 5 times
                 window=5)              # define context as a 5-word window around the target word

In [ ]: # now you can get word vectors!
model.get_vector('anything')

Out[ ]: array([-3.5468231,  1.7865689,  0.3841988,  2.7576689,  2.0146688,
        1.9754922,  0.3358841, -0.8200357,  0.48848522,  2.4279432,
        -0.7554728,  2.0798082,  3.2052414,  -0.16385323,  2.6689763,
        -1.8514861,  -0.0295971,  -1.6357832,  0.5121673,  -1.3257158,
        -2.1526686,  -0.4305907,  -1.0003162,  -3.09036,   2.290466,
        -2.3926432,  0.89643824,  1.9492018,  0.8349378,  0.50875283,
        0.8442198,  0.9906684 ], dtype=float32)

In [ ]: # or query similar words directly. Go play with it!
model.most_similar('broad')

Out[ ]: [('price', 0.956332266330719),
 ('sauc', 0.9481489192085266),
 ('vodka', 0.9390178574188232),
 ('cheese', 0.9327408326728823),
 ('fruit', 0.9191921234406851),
 ('fruit', 0.924953798106553),
 ('beans', 0.9199380480363342),
 ('banana', 0.916680897579956),
 ('potato', 0.915650486946106),
 ('wine', 0.9089671674728394)]

Using pre-trained model
```

Took it a while, huh? Now imagine training life-sized (100-300D) word embeddings on gigabytes of text: wikipedia articles or twitter posts.

Thankfully, nowadays you can get a pre-trained word embedding model in 2 lines of code (no sms required, promise).

```
In [ ]: import gensim.downloader as api
model = api.load('glove840ts3100')

In [ ]: model.most_similar(positive=["codeur", "money"], negative=["brain"])

Out[ ]: [('broker', 0.5820155739784241),
 ('bonuses', 0.5424473285675949),
 ('banker', 0.538512762453172),
 ('designer', 0.5197198390969693),
 ('merchandising', 0.49642333388382855),
 ('street', 0.4822018835286984),
 ('shopper', 0.4928562267698822),
 ('part-time', 0.451282367815993),
 ('freelance', 0.4843311985690901),
 ('upair', 0.4796452522277832)]

Visualizing word vectors
```

One way to see if our vectors are any good is to plot them. Thing is, those vectors are in 300+ space and we humans are more used to 2-3D.

Luckily, we machine learners know about **dimensionality reduction** methods.

Let's use that to plot 1000 most frequent words

```
In [ ]: words = model.index_to_key[:1000]

print(words[:100])

['users', '_', 'please', 'apa', 'justin', 'text', 'hari', 'playing', 'once', 'sei']

In [ ]: # for each word, compute it's vector with model
word_vectors = np.array([model.get_vector(x) for x in words])

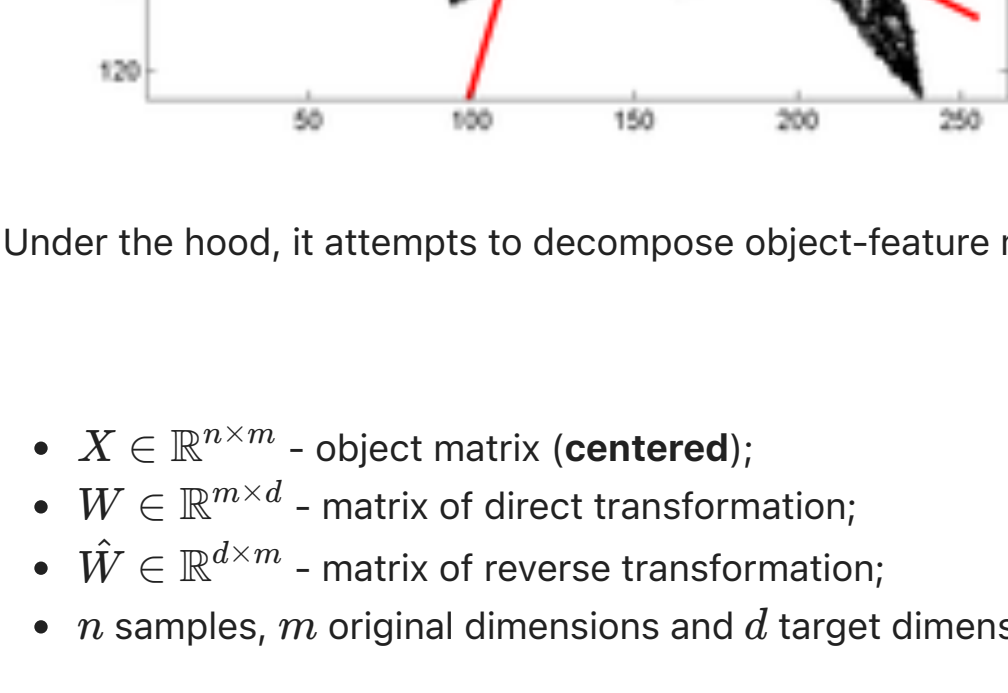
In [ ]: # word_vectors

In [ ]: assert isinstance(word_vectors, np.ndarray)
assert word_vectors.shape == (len(words), 300)
assert np.isfinite(word_vectors).all()

Linear projection: PCA
```

The simplest linear dimensionality reduction method is `__p__rincipal__C__omponent__A__nalysis`.

In geometric terms, PCA tries to find axes along which most of the variance occurs. The "natural" axes, if you wish.



Under the hood, it attempts to decompose object-feature matrix X into two smaller matrices: \hat{W} and \hat{V} minimizing **mean squared error**:

$$\|(XW)\hat{V} - X\|_F^2 \rightarrow_{W, \hat{V}} \min$$

- $X \in \mathbb{R}^{n \times m}$ - object matrix (**centered**);
- $\hat{W} \in \mathbb{R}^{m \times d}$ - matrix of direct transformation;
- $\hat{V} \in \mathbb{R}^{d \times m}$ - matrix of reverse transformation;
- n samples, m original dimensions and d target dimensions;

```
In [ ]: from sklearn.decomposition import PCA

# map word vectors onto 2d plane with PCA. Use good old sklearn api (fit, transform)
# after that, normalize vectors to make sure they have zero mean and unit variance
word_vectors_pca = PCA(n_components=2).fit_transform(word_vectors)
# and maybe PCA OF YOUR CODE here :)

In [ ]: word_vectors_pca = (word_vectors_pca - word_vectors_pca.mean(axis=0)) / word_vectors_pca.std(axis=0)

In [ ]: assert word_vectors_pca.shape == (len(word_vectors), 2), "there must be a 2d vector for each word"
assert max(abs(word_vectors_pca.mean(0))) < 1e-5, "points must be zero-centered"
assert max(abs(1.0 - word_vectors_pca.std(0))) < 1e-2, "points must have unit variance"

Let's draw it!
```

```
In [ ]: import bokeh.models as bm, bokeh.plotting as pl
from bokeh.io import output_notebook
output_notebook()

def draw_vectors(x, y, radius=10, alpha=0.25, color='blue',
               width=600, height=400, show=True, **kwargs):
    """ draws an interactive plot for data points with auxiliary info on hover """
    if isinstance(color, str): color = (color,) * len(x)
    data_source = bm.ColumnDataSource({'x': x, 'y': y, 'color': color, **kwargs})

    fig = pl.figure(active_scroll='wheel_zoom', width=width, height=height)
    fig.scatter('x', 'y', size=radius, color='color', alpha=alpha, source=data_source)

    fig.add_tools(bm.HoverTool(tooltips=[(key, "g" + key) for key in kwargs.keys()]))
    if show: pl.show(fig)
    return fig

BokehJS 3.2.2 successfully loaded.
```

```
In [ ]: draw_vectors(word_vectors_pca[:, 0], word_vectors_pca[:, 1], token=words)

# hover a mouse over there and see if you can identify the clusters

Out[ ]: figure(id='p1049', ...)
```

Visualizing neighbors with t-SNE

PCA is nice but it's strictly linear and thus only able to capture coarse high-level structure of the data.

If we instead want to focus on keeping neighboring points near, we could use TSNE, which is itself an embedding method. Here you can read [more on TSNE](#).

```
In [ ]: from sklearn.manifold import TSNE

# map word vectors onto 2d plane with TSNE. hint: don't panic it may take a minute or two to fit.
# normalize them as just like with pca

word_tsne = TSNE(n_components=2).fit_transform(word_vectors)

In [ ]: word_tsne = (word_tsne - word_tsne.mean(axis=0)) / word_tsne.std(axis=0)

In [ ]: draw_vectors(word_tsne[:, 0], word_tsne[:, 1], color='green', token=words)

Out[ ]: figure(id='p1094', ...)
```

Visualizing phrases

Word embeddings can also be used to represent short phrases. The simplest way is to take an **average** of vectors for all tokens in the phrase with some weights.

This trick is useful to identify what data are you working with: find if there are any outliers, clusters or other artefacts.

Let's try this new hammer on our data!

```
In [ ]: def get_phrase_embedding(phrase):
    """
    Convert phrase to a vector by aggregating it's word embeddings. See description above.

    # 1. lowercase phrase
    # 2. tokenize phrase
    # 3. average word vectors for all words in tokenized phrase
    # skip words that are not in model's vocabulary
    # if all words are missing from vocabulary, return zeros

    word = np.zeros([model.vector_size], dtype='float32')

    vectors = np.array([model.get_vector(x, None) for x in tokenizer.tokenize(phrase.lower()) if x in model])
    if len(word_vectors)!=0:
        return vector
    vector = word_vectors.mean(axis=0)

    return vector

In [ ]: vector = get_phrase_embedding('I'm very sure. This never happened to me before...')

assert np.allclose(vector[:10],
                    np.array([-0.31887372, -0.02558171,  0.0933293, -0.1002182, -1.0278669,
                              -0.16621883,  0.05803408,  0.17989802,  1.3701859,  0.08655966],
                              dtype=np.float32))

In [ ]: # let's only consider ~5k phrases for a first run.
chosen_phrases = data[:len(data) // 1000]

# compute vectors for chosen phrases
phrase_vectors = np.array([get_phrase_embedding(x) for x in chosen_phrases])

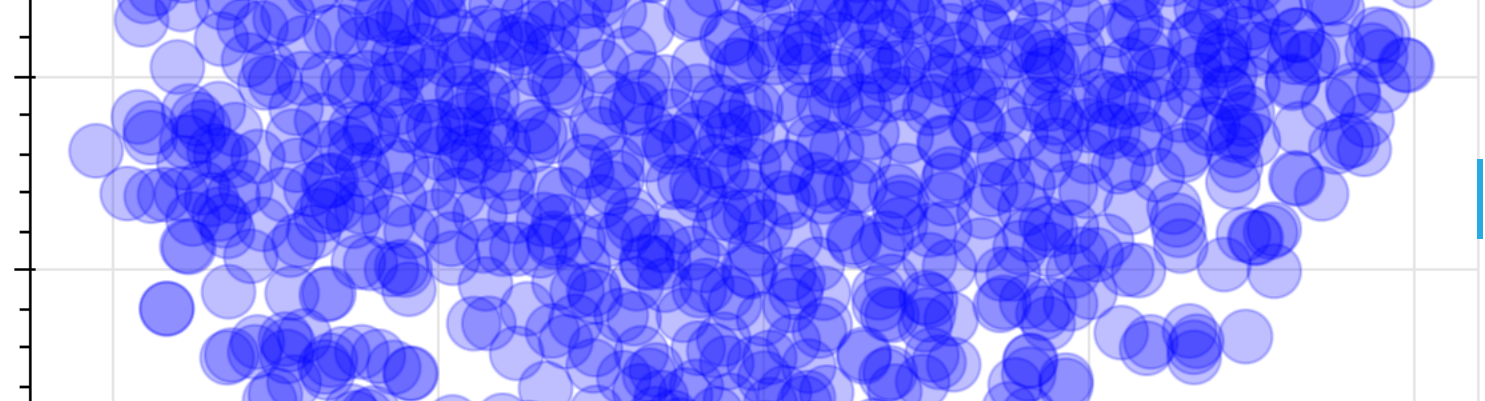
In [ ]: assert isinstance(phrase_vectors, np.ndarray) and np.isfinite(phrase_vectors).all()
assert phrase_vectors.shape == (len(chosen_phrases), model.vector_size)

In [ ]: # map vectors into 2d space with pca, tsne or your other method of choice
# don't forget to normalize

phrase_vectors_2d = TSNE().fit_transform(phrase_vectors)

phrase_vectors_2d = (phrase_vectors_2d - phrase_vectors_2d.mean(axis=0)) / phrase_vectors_2d.std(axis=0)

In [ ]: draw_vectors(phrase_vectors_2d[:, 0], phrase_vectors_2d[:, 1],
                    phrase=[phrase[:150] for phrase in chosen_phrases],
                    radius=20)
```



Finally, let's build a simple "similar question" engine with phrase embeddings we've built.

```
In [ ]: # compute vector embedding for all lines in data
data_vectors = np.array([get_phrase_embedding(l) for l in data])

In [ ]: model.similar_by_vector(get_phrase_embedding('Hello where are you?'))

Out[ ]: [['you', 0.94973647594519),
 ('there', 0.9299634899806653),
 ('re', 0.924767813682556),
 ('know', 0.928958438444891),
 ('where', 0.9189091324086213),
 ('what', 0.918116731725587),
 ('how', 0.907548246787828),
 ('are', 0.9031845992126405),
 ('n', 0.9030184745788514),
 ('why', 0.8966574668884277)]

In [ ]: from sklearn.metrics.pairwise import cosine_similarity

def find_nearest(query, k=10):
    """
    given text line (query), return k most similar lines from data, sorted from most to least similar
    similarity should be measured as cosine between query and line embedding vectors
    hint: it's okay to use global variables: data and data_vectors. see also: np.argsort, np.argsortort

    """
    # YOUR CODE
    query_emb = get_phrase_embedding(query)
    cosines = cosine_similarity(query_emb.reshape(1, -1), data_vectors)[0]
    top_k_indexes = np.argsort(cosines)[::-1][:k]
    top_k_phrases = np.array(data[top_k_indexes])
    return top_k_phrases

In [ ]: results = find_nearest(query="How do I enter the matrix?", k=10)

print(' '.join(results))

assert len(results) == 10 and isinstance(results[0], str)
assert results[0] == "How do I get to the dark web?"
assert results[1] == "What can I do to save the world?"

How do I get to the dark web?
What should I do to enter hollywood?
How do I use the greenify app?
What can I do to save the world?
How do I win this?
How do I think out of the box? How do I learn to think out of the box?
How do I increase the 5th dimension?
How do I use the pad in PMA?
How do I estimate the competition?
What do I do to enter the line of event management?

In [ ]: find_nearest(query="How does Trump?", k=10)

Out[ ]: array(['What does Donald Trump think about Israel?',
 'What books does Donald Trump think like?',
 'What does Donald Trump think of India?',
 'What does India think of Donald Trump?',
 'What does Donald Trump think of China?',
 'What does Donald Trump think about Pakistan?',
 'What companies does Donald Trump own?',
 'What does Dushka Zapata think about Donald Trump?',
 'How does it feel to date Ivanka Trump?',
 'What does salesforce mean?'], dtype='<U1178')

In [ ]: find_nearest(query="Why don't I ask a question myself?", k=10)

Out[ ]: array(['Why don't I get a date?',
 'Why do you always answer a question with a question? I don't, or do I?',
 'Why can't I ask a question anonymously?',
 'Why don't I get a girlfriend?',
 'Why don't I have a boyfriend?', 'I don't have no question?',
 'Why can't I take a joke?', 'Why don't I ever get a girl?',
 'Can I ask a girl out that I don't know?',
 'Why don't I have a girlfriend?'], dtype='<U1178')

Now what?

• Try running TSNE on all data, not just 1000 phrases
• See what other embeddings are there in the model zoo: gensim.downloader.info()
• Take a look at fastText embeddings
• Optimize find_nearest with locality-sensitive hashing: use nearypy or sklearn.neighbors.
```

TSNE on all data

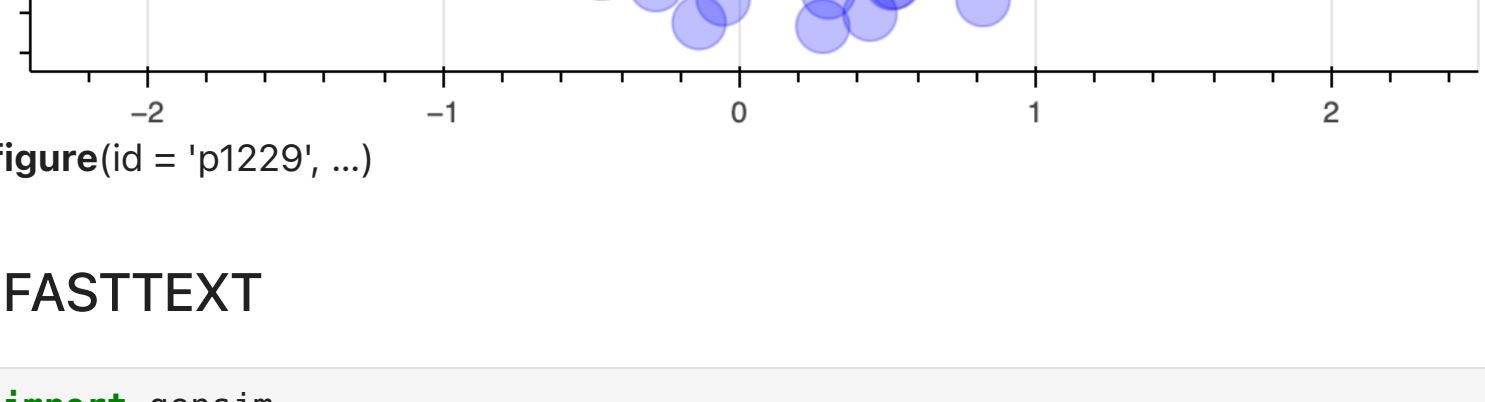
```
In [ ]: # TSNE
phrase_vectors = np.array([get_phrase_embedding(x) for x in data])

In [ ]: phrase_vectors_2d = TSNE().fit_transform(phrase_vectors)

phrase_vectors_2d = (phrase_vectors_2d - phrase_vectors_2d.mean(axis=0)) / phrase_vectors_2d.std(axis=0)

In [ ]: draw_vectors(phrase_vectors_2d[:, 0], phrase_vectors_2d[:, 1],
                    phrase=[phrase[:150] for phrase in chosen_phrases],
                    radius=20)

BokehJSWarning: ColumnDataSource's columns must be of the same length. Current lengths: ('color', 537272), ('phrase', 1001), ('x', 537272), ('y', 537272)
```



FASTTEXT

```
In [ ]: import gensim
gensim.downloader.info().keys()

Out[ ]: dict_keys(['corpora', 'models'])

In [ ]: gensim.downloader.info()['models'].keys()

Out[ ]: dict_keys(['fasttext-wiki-news-subwords-300', 'conceptnet-numberbatch-17-06-300', 'word2vec-ruscorpora-300', 'word2vec-google-news-300', 'glove-wiki-gigaword-50', 'glove-wiki-gigaword-100', 'glove-wiki-gigaword-200', 'glove-wiki-gigaword-300', 'glove-twitter-100', 'glove-twitter-200', 'testing_word2vec-matrix-synop sis'])

In [ ]: model = api.load('fasttext-wiki-news-subwords-300')

In [ ]: get_phrase_embedding('Hello, how are you?').shape

Out[ ]: (300,)

In [ ]: data_vectors = np.array([get_phrase_embedding(l) for l in data])

In [ ]: results = find_nearest(query="How do I enter the matrix?", k=10)

print(' '.join(results))

How do I evaluate the integral?
How do I estimate the competition?
How do I simplify the expression?
How do I manage the business?
How do I choose the proper profession?
How do I access the Extralorent website?
How do I prepare the resume?
How do I increase the vocabulary?
How do I crack the GMAT?
How do I crack the CLAT?

In [ ]: data_vectors.shape

Out[ ]: (537272, 300)

Local-Sensitivity Hashing (approximate neighbours)
```

```
In [ ]: import numpy
import tqdm

from neapy import Engine
from neapy.hashes import RandomBinaryProjections

# Dimension of our vector space
dimension = 300

# Create a random binary hash with 10 bits
rbp = RandomBinaryProjections('rbp', 10)

# Create engine with pipeline configuration
engine = Engine(dimension, lshashes=[rbp])

# Index 2000000 random vectors (set their data to a unique string)
for index, vec in tqdm.tqdm(enumerate(data_vectors)):
    engine.store_vector(vec, 'data_{}'.format(index))

537272it [00:05, 89959.33it/s]

In [ ]: def find_nearest_neappy(query):
    """
    given text line (query), return k most similar lines from data, sorted from most to least similar
    similarity should be measured as cosine between query and line embedding vectors
    hint: it's okay to use global variables: data and data_vectors. see also: np.argsort, np.argsortort

    """
    # YOUR CODE
    query_emb = get_phrase_embedding(query)
    top_k_indexes = np.array([int(x) for x in engine.neighbours(query_emb)])
    top_k_neighbours = np.array(data[top_k_indexes])
    return top_k_neighbours

In [ ]: find_nearest_neappy(query="How does Trump?", k=10)

Out[ ]: array(['How does GCM generate registration_ids?',
 'How does Donald Trump persuade?',
 'How does Informatica?',
 'How does politics originate?',
 'How does Websockets operate?',
 'How does PPF works?',
 'How does iCefaces works?',
 'How does anything exist?',
 'How does TheTake works?',
 'How does SHAREit work?'], dtype='<U1178')

In [ ]: find_nearest_neappy(query="How does Trump?", k=10)

Out[ ]: array(['How does Donald Trump persuade?',
 'How does Donald Trump treat waitstaff?',
 'Why does everybody hate Trump?',
 'What does nationalism mean?',
 'What does Tiffany Trump do?',
 'What does manipulation mean?',
 'What does SENSEX mean?',
 'What does 0.8026 mean?',
 'What does eljeto mean?'], dtype='<U1178')

In [ ]: find_nearest_neappy(query="Why don't I ask a question myself?", k=10)

Out[ ]: array(['Why don't I ask a question anonymously?',
 'Why don't I have a female friend?',
 'Why don't I love myself?',
 'Why don't I have a girlfriend?',
 'Why don't I take a joke?',
 'Why don't I ever get a girl?',
 'Why don't I have a girlfriend?'], dtype='<U1178')

In [ ]: find_nearest_neappy(query="Why don't I ask a question myself?", k=10)

Out[ ]: array(['Why can't I ask a question anonymously?',
 'Why don't I have a female friend?',
 'Why don't I love myself?',
 'Why don't I want friends?',
 'Why can't I understand myself?',
 'Why can't I feel myself?',
 'Why can't I keep a conversation going?'], dtype='<U1178')
```