

Data Science

# Documentation

**Nearest Neighbors - Advanced (LSH)**

January 20, 2018

Eicker Niklas, Halastra Szymon

# Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
1.1	The data set . . . . .	3
<b>2</b>	<b>Data preparation</b>	<b>4</b>
<b>3</b>	<b>Locality Sensitive Hashing</b>	<b>5</b>
3.1	Sorting data into bins . . . . .	5
3.2	First evaluation of a LSH model . . . . .	6
<b>4</b>	<b>Querying the LSH model</b>	<b>7</b>
4.1	A query . . . . .	7
4.2	Realization in a method . . . . .	7

# 1 Introduction

Using GraphLab Create we will load and analyze data on Wikipedia articles about persons. In this advanced project we will use Locality Sensitive Hashing (LSH) to achieve fast and efficient approximate nearest neighbor searches.

## 1.1 The data set

Each element of the original data set consists of a link to a Wikipedia article, the name of the person it is about and the text of the article (in lowercase and without punctuation). There are 59071 entries in total. An excerpt of the data "as is" can be seen in figure 1.1.

Figure 1.1: Excerpt of the data set.

URI	name	text
<http://dbpedia.org/resource/Digby_Morrell> ...	Digby Morrell	digby morrell born 10 october 1979 is a former ...
<http://dbpedia.org/resource/Alfred_J._Lewy> ...	Alfred J. Lewy	alfred j lewy aka sandy lewy graduated from ...
<http://dbpedia.org/resource/Harpdog_Brown> ...	Harpdog Brown	harpdog brown is a singer and harmonica player who ...
<http://dbpedia.org/resource/Franz_Rottensteiner> ...	Franz Rottensteiner	franz rottensteiner born in waldmannsfeld lower ...
<http://dbpedia.org/resource/G-Enka> ...	G-Enka	henry krivits born 30 december 1974 in tallinn ...
<http://dbpedia.org/resource/Sam_Henderson> ...	Sam Henderson	sam henderson born october 18 1969 is an ...
<http://dbpedia.org/resource/Aaron_LaCrate> ...	Aaron LaCrate	aaron lacrate is an american music producer ...
<http://dbpedia.org/resource/Trevor_Ferguson> ...	Trevor Ferguson	trevor ferguson aka john farrow born 11 november ...

## 2 Data preparation

As preparation we add a new column to the data set containing the TF-IDF values as well as a row count acting as unique article ID.

The GraphLab Create method 'text\_analytics.tf\_idf' calculates TF-IDF values for all our articles. Adding the ID is done via an in-build method of the SFrame 'add\_row\_number'. A short preview of the modified data set can be seen in figure 2.1.

Figure 2.1: Preview of the data set with an additional column for TF-IDF and unique IDs.

id	URI	name	text	tf_idf
0	<http://dbpedia.org/resource/Digby_Morrell> ...	Digby Morrell	digby morrell born 10 october 1979 is a former ...	{'selection': 3.836578553093086, ...
1	<http://dbpedia.org/resource/Alfred_J._Lewy> ...	Alfred J. Lewy	alfred j lewy aka sandy lewy graduated from ...	{'precise': 6.44320060695519, ...
2	<http://dbpedia.org/resource/Harpdog_Brown> ...	Harpdog Brown	harpdog brown is a singer and harmonica player who ...	{'just': 2.7007299687108643, ...
3	<http://dbpedia.org/resource/Franz_Rottensteiner> ...	Franz Rottensteiner	franz rottensteiner born in waidmannsfeld lower ...	{'all': 1.6431112434912472, ...
4	<http://dbpedia.org/resource/G-Enka> ...	G-Enka	henry krivits born 30 december 1974 in tallinn ...	{'they': 1.8993401178193898, ...
5	<http://dbpedia.org/resource/Sam_Henderson> ...	Sam Henderson	sam henderson born october 18 1969 is an ...	{'currently': 1.637088969126014, ...

For the rest of the assignment we will use sparse matrices instead of the SFrame. A sparse matrix is a matrix with only a few non zero elements. SciPy supports this format and allows us to easily handle our TF-IDF values in matrix form. First we have to convert the original dictionary structure to the new SciPy sparse matrix format.

We transform the TF-IDF dictionaries in one 59071 x 547979 sparse matrix, where each row is a document and each column a word. Now every document has a TF-IDF value for every word in any document. If the word does not appear in the document, which is the case for most words, the value is zero.

## 3 Locality Sensitive Hashing

LSH is a hashing method that aims to maximize the probability of a "collision" for similar items. Via randomly generated vectors we will assign each document to a specific bin. The number of these bins will be chosen in a way, that the bin of a queried document and it's surrounding bins include the "real" nearest neighbors. Due to the, compared to the complete sample, smaller amount of all of these documents in close bins, it will be possible to calculate the closest neighbors much faster then by brute force.

### 3.1 Sorting data into bins

For the beginning we have to decide on a number of random vectors. This number will define the number of bins. With each vector we will generate one bit for every document by multiplying it with the document's TF-IDF vector and checking the sign of the solution. Each unique combination of bits represents one data bin.

We create 16 Gaussian random vectors, which will yield  $2^{16}$  (65536) bins to sort our data into. In this case the number of bins is comparable to the number of documents in the data set.

Each scalar product with a document and all random vectors gives us an array of Boolean with the length of 16 (= number of vectors). We can transform this array into an ID for the corresponding bin, by taking it as a 16-bit binary number and calculating it's single integer value. At this point we have sorted all our documents into bins and can start to compare single documents to other documents inside the same bin or surrounding bins.

## 3.2 First evaluation of a LSH model

With every document in it's corresponding bin, we can take a look at the actual "closeness" of documents in the same and surrounding bins. To measure the closeness we will be using the cosine distance, which was closer investigated in the primer task.

As an example we will (again) be using the Wikipedia page of Barack Obama. In Obama's bin there are five other pages. Table 3.1 lists the others names and cosine distances to Barack Obama's article. As we can see, most of them are far away from Barack Obama. It is obvious, that LSH is just an approximation and one has to consider documents from surrounding bins to find the nearest neighbors.

Table 3.1: Comparison of the documents that share the same bin as Barack Obama's article.

Rank	Name	Distance to Obama
1	Joe Biden	0.703139
2	Mark Boulware	0.950867
3	John Wells (politician)	0.975966
4	Francis Longstaff	0.978256
5	Madurai T. Srinivasan	0.993092

## 4 Querying the LSH model

### 4.1 Structure of a query

In the previous chapter we found that a query of our model will have to take surrounding bins into account. One can define a range for such a query as the number of different bits in the bins bit representations. Such a query with the range 3 will look like this:

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

### 4.2 Query implementation

For an easy implementation of the query, we use "itertools.combinations". This method takes the number of our vectors (number of bits in a bin ID) and the query range. It returns a collection of lists, which show all possible bit flips that lead to surrounding bins with the given range as maximum number of flips. Our query method takes these lists of bit flips, applies each one of them and saves all documents found in any of the reached bins in a set. This set then includes all possible candidates for being nearest neighbors and can be further examined.