

# Homework 1

36-350: Data Mining

Due at start of class, Friday, 4 September 2009

1.
  - (a) What is the bag-of-words representation of the sentence “To be or not to be”?
  - (b) Suppose we search for the above sentence via the keyword “be”. What is the bag-of-words representation for this query, and what is its Euclidean distance from the sentence?
  - (c) How would inverse-document-frequency (IDF) weighting help when making a Web query for “The Principles of Data Mining”?
  - (d) Describe a simple text search that can’t be done using a bag-of-words representation (no matter what distance measure is used). “Simple” means no actual understanding of English is required.
2.
  - (a) What is the Euclidean distance between each of the vectors  $(1, 0, 0)$ ,  $(1, 4, 5)$ , and  $(10, 0, 0)$ ?
  - (b) Divide each vector by its sum. How do the relative distances change?
  - (c) Divide each vector by its Euclidean length. How do the relative distances change?
  - (d) Suppose we’re using the bag-of-words representation for similarity searching with a Euclidean metric. Describe how the previous parts of the question illustrate a potential problem if we do not normalize for document length.
  - (e) Consider the conventional searching scheme where the user picks a set of keywords and the system returns all documents containing those keywords. Describe how the previous parts of the question illustrate a potential problem with this type of search.

The remaining questions are this week's computer exercise. They use some a some pre-written R functions in a file called `01.R`, available from Blackboard or from <http://www.stat.cmu.edu/~cshalizi/350>, and a fragment of the *Times* corpus, `nyt_corpus.zip`. See the end of this document for some notes about the functions. Also, please read the “Minimal Advice on Programming” handout (linked to in both places), if you haven't already — especially the part on commenting your code, and the part on how to complain about this assignment.

3. (a) Create document vectors for each of the stories in the `music` and `art` folders. Give the commands you used.  
*Note:* Feel free to use the `read.directory` function in `01.R` if you comment it.
  - (b) What command would you use to extract the 37<sup>th</sup> word of story number 1595645 in `art`? (That word is “experiencing”.) Give a command to count the number of times the word “the” appears in that story. (There are at least two ways to do this. The correct answer is 103.)
  - (c) Give the commands you would use to construct a bag-of-words data-frame from the document vectors for the `art` and `music` stories. *Hint:* consider `lapply`.
  - (d) Create distance matrices from this data frame for (a) the straight Euclidean distance, (b) the distance with word-count normalization and (c) the distance with vector-length scaling, and then for all three again with inverse-document-frequency weighting. Give the commands you use.
  - (e) For each of the six different difference measures, what is the average distance between stories in the same category and between stories in different categories? (Include the R command you use to compute this — don't do it by hand!)
  - (f) Create multidimensional scaling plots for the different distances, and describe what you see. Include the code you used, the plots, and explanations for the code.
4. Comment the `sq.Euc.dist` function — that is, go over it and explain, in English, what it each does, and how the lines work together to calculate the function.
  5. (a) Explain what the “cosine distance” has to do with cosines.
  - (b) Calculate, by hand, the cosine distances between the three vectors in question 2.
  - (c) Write a function to calculate the matrix of cosine distances (really, similarities) between all the vectors in a data-frame. *Hint:* you may want to use the `distances` function in `01.R`. Check that your function agrees with your answer to the previous part.

6. Write a function to find the document which best matches a given query string. The function should take two arguments:

- The query, as a single character string
- The bag-of-words matrix

and return the row number corresponding to the best-matching document. You can pick the distance measurement, but you should include inverse document-frequency weighting.

*Hint 1:* The first step should be to turn the query into a bag-of-words vector. Make sure it's comparable to the bag-of-words matrix!

*Hint 2:* Look at the `nearest.points` function in `O1.R`.

*Hint 3:* Test that if the “query” one of the original documents, that's the document which is returned. Do this for at least two of the original documents.

## Some Notes on the Functions

**XML** The stories from the *Times* corpus are in a language called XML. You don't have to know it, but you do need to install and load the XML package from CRAN.

**Reading documents into R** The function `read.doc` reads documents into R, as follows:

```
music.1 = read.doc("music//0023931.xml")
```

This loads the file `0023931.xml` in the directory `music`, and then turns the text of the news story into a vector of word instances in the order in which they appeared in the story. `read.doc` removes the annotations, removes all punctuation, shifts all letters to lower case, and turns all numbers into the pound sign `#`. You can access the  $n^{\text{th}}$  word as `music.1[n]`.

`table(music.1)` creates the bag-of-words representation from the vector. The `table` function is an extremely useful part of the basic R package, and you should familiarize yourself with it.

**Creating a Bag-of-Words Data-Frame** The function `make.Bow.frame` converts a list of bag-of-word vectors into a data frame, with one row for each document and one column for each word. By default, words which appear in only a single document are removed from the unified list of columns; this can be suppressed by running it with the argument `remove.singletons=FALSE`.

**Normalization** The functions `div.by.sum` and `div.by.euc.length` normalize an array by the sum of each row and by the Euclidean length of each row, respectively. For instance,

```
x = div.by.sum(docs)
```

would create a new data-frame, `x`, in which each row of `docs` was normalized by the sum of entries in that row.

**Computing distances** The function `distances` computes a matrix of distances between the different bag-of-word vectors in a data frame.

```
d = distances(x)
```

creates a new matrix, `d`, where `d[i, j]` is the distance between `x[i, ]` and `x[j, ]`.

**Multidimensional scaling** There are three standard multi-dimensional scaling functions in R, `cmdscale`, which is part of the default package `stats`, and `isoMDS` and `sammon`, which are part of the package `MASS`. See their help files for details.