

Module X: Probabilistic Blocking

Rebecca C. Steorts

Agenda

- ▶ Data Cleaning Pipeline
- ▶ Blocking
- ▶ Probabilistic Blocking
- ▶ LSH

Load R packages

```
## Loading required package: DBI
## Loading required package: RSQLite
## Loading required package: ff
## Loading required package: bit

##
## Attaching package: 'bit'

## The following object is masked from 'package:base':
##
##      xor

## Attaching package ff

## - getOption("fftempdir")=="/var/folders/bv/xhclmwh90zg08
## - getOption("ffextension")== "ff"
## - getOption("ffdrops")==TRUE
```

Data Cleaning Pipeline



Figure 1: Data cleaning pipeline.

Blocking

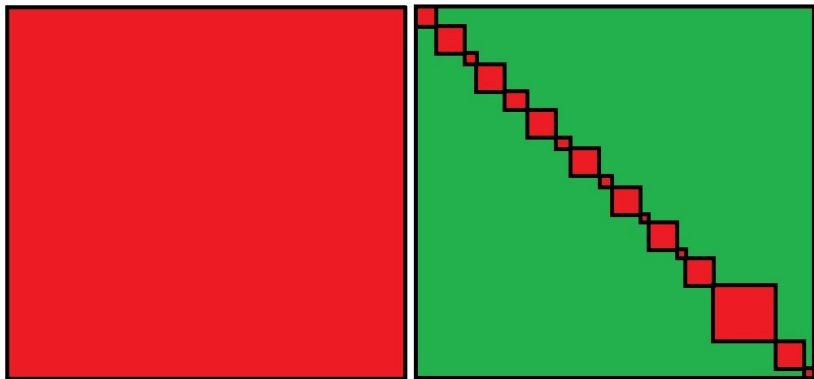


Figure 2: Left: All to all record comparison. Right: Example of resulting blocking partitions.

LSH

Locality sensitive hashing (LSH) is a fast method of blocking for record linkage that originates from the computer science literature.

Finding similar items

- ▶ We want to find similar items
 - ▶ Maybe we are looking for near duplicate documents (plagiarism)
 - ▶ More likely, we are trying to block our data which we can later pass to a record linkage process
- ▶ How do we define *similar*?

Jaccard similarity

As already mentioned there are many ways to define similarity.

In this lecture, we will need the *Jaccard similarity*:

$$Jac(S, T) = \frac{|S \cap T|}{|S \cup T|}.$$

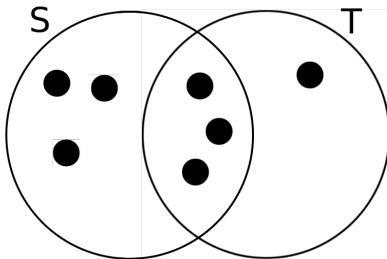


Figure 3: Two sets S and T with Jaccard similarity $3/7$. The two sets share 3 elements in common, and there are 7 elements in total.

How to represent data as sets?

We want to talk about the similarity of our data (records) \Rightarrow we need to compare sets of records!

- ▶ We can construct a set of **short strings** from the data
- ▶ This is useful because similar datasets will have many common elements (common short strings)
- ▶ We can do construct these short strings using *shingling*

k -shingling (how-to)

1. Think of our data set as a string of characters
2. A k -shingle (k -gram) is any sub-string (word) of length k found within the document or record
3. Associate with each data set the set of k -shingles that appear one or more times

Let's try

Suppose our document is the string “Hello world”, then

- ▶ the set of 2-shingles is {he, el, ll, lo, ow, wo, or, rl, ld}
- ▶ the set of 3-shingles is {hel, ell, llo, low, owo, wor, orl, rld}

Your turn

We have the following two records:

```
# load RL data  
data("RLdata500")  
  
# select only 2 records  
records <- RLdata500[129:130, c(1,3)]  
names(records) <- c("First name", "Last name")  
  
# inspect records  
kable(records)
```

	First name	Last name
129	MICHAEL	VOGEL
130	MICHAEL	MEYER

Your turn (continued)

1. Compute the 2-shingles for each record
2. Using Jaccard similarity, how similar are they?
3. What do you learn from this exercise?

Your turn solution

1. The 2-shingles for the first record are {mi, ic, ch, ha, ae, el, lv, vo, og, ge, el} and for the second are {mi, ic, ch, ha, ae, el, lm, me, ey, ye, er}
2. There are 6 items in common {mi, ic, ch, ha, ae, el} and 15 items total {mi, ic, ch, ha, ae, el, lv, vo, og, ge, lm, me, ey, ye, er}, so the Jaccard similarity is $\frac{6}{15} = \frac{2}{5} = 0.4$
3. You should have learned that this is very tedious to do by hand!

Useful packages/functions in R

(Obviously) We don't want to do this by hand most times.

Here are some useful packages in R that can help us!

```
library(textreuse) # text reuse/document similarity  
library(tokenizers) # shingles
```

```
##
```

```
## Attaching package: 'tokenizers'
```

```
## The following objects are masked from 'package:textreuse':
```

```
##
```

```
##      tokenize_ngrams, tokenize_sentences, tokenize_skip_n
```

```
##      tokenize_words
```

We can use the following functions to create k -shingles and calculate Jaccard similarity for our data

```
# get k-shingles  
tokenize_character_shingles(x, n)
```

Citation Data Set

Research paper headers and citations, with information on authors, title, institutions, venue, date, page numbers and several other fields

```
library(devtools)
```

```
## Loading required package: usethis
```

```
install_github("resteorts/cora")
```

```
## Skipping install of 'cora' from a github remote, the SHA1 (70e32d5d) has not changed since last install.
```

```
## Use `force = TRUE` to force installation
```

```
library(cora)
```

```
library(ggplot2)
```

```
data(cora) # load the cora data set
```

```
str(cora) # structure of cora
```

```
## 'data.frame': 1879 obs. of 16 variables:
```

```
## $ id : int 1 2 3 4 5 6 7 8 9 10 ...
```

```
## $ title : 'noquote' chr "Inganas and M.R" NA NA NA ...
```

```
## $ book_title : 'noquote' chr NA NA NA NA ...
```

```
## $ authors : 'noquote' chr "M. Ahlskog, J. Paloheimo, H. Stubbs, P. Dyreklev, M. Fahlman, O" "M. Ahlskog, J. Paloheimo, H. Stubbs, P. Dyreklev, M. Fahlman, O" ...
```

```
## $ address : 'noquote' chr NA NA NA NA ...
```

```
## $ date : 'noquote' chr "1994" "1994" "1994" "1994" ...
```

```
## $ year : 'noquote' chr NA NA NA NA ...
```

```
## $ editor : 'noquote' chr NA NA NA NA ...
```

```
## $ journal : 'noquote' chr "Andersson, J Appl. Phys." "JAppl. Phys." "J Appl. Phys." "J Appl. Phys." ...
```

```
## $ volume : 'noquote' chr "76" "76" "76" "76" ...
```

```
## $ pages : 'noquote' chr "893" "893" "893" "893" ...
```

```
## $ publisher : 'noquote' chr NA NA NA NA ...
```

```
## $ institution: 'noquote' chr NA NA NA NA ...
```

```
## $ type : 'noquote' chr NA NA NA NA ...
```

```
## $ tech : 'noquote' chr NA NA NA NA ...
```

```
## $ note : 'noquote' chr NA NA NA NA ...
```


Your turn

Using the title, authors, and journal fields in the cora dataset,

1. Get the 3-shingles for each record (**hint:** use `tokenize_character_shingles`)
2. Obtain the Jaccard similarity between each pair of records (**hint:** use `jaccard_similarity`)

Your turn (solution)

```
# get only the columns we want
n <- nrow(cora) # number of records
dat <- data.frame(id = seq_len(n)) # create id column
dat <- cbind(dat, cora[, c("title", "authors", "journal")]) # get columns we want

# 1. paste the columns together and tokenize for each record
shingles <- apply(dat, 1, function(x) {
  # tokenize strings
  tokenize_character_shingles(paste(x[-1], collapse=" "), n = 3)[[1]]
})

# 2. Jaccard similarity between pairs
jaccard <- expand.grid(record1 = seq_len(n), # empty holder for similarities
                      record2 = seq_len(n))

# don't need to compare the same things twice
jaccard[jaccard$record1 < jaccard$record2,]

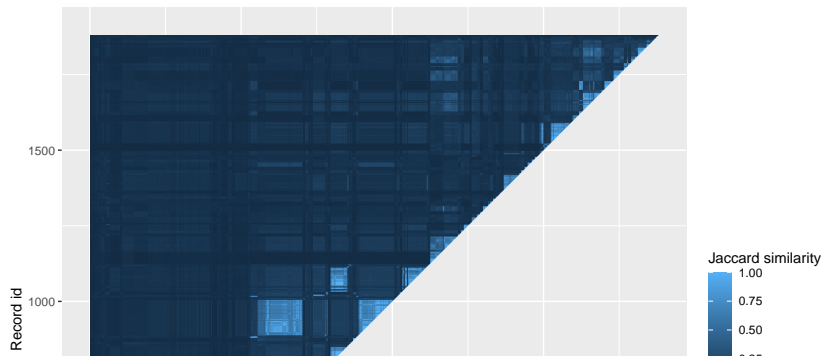
time <- Sys.time() # for timing comparison
jaccard$similarity <- apply(jaccard, 1, function(pair) {
  jaccard_similarity(shingles[[pair[1]]], shingles[[pair[2]]]) # get jaccard for each pair
})
time <- difftime(Sys.time(), time, units = "secs") # timing
```

This took took 99.51 seconds \approx 1.66 minutes

Your turn (solution, cont'd)

plot the jaccard similarities for each pair of records

```
ggplot(jaccard) +  
  geom_raster(aes(x = record1, y = record2,  
                  fill=similarity)) +  
  theme(aspect.ratio = 1) +  
  scale_fill_gradient("Jaccard similarity") +  
  xlab("Record id") + ylab("Record id")
```



Your turn (solution, cont'd)

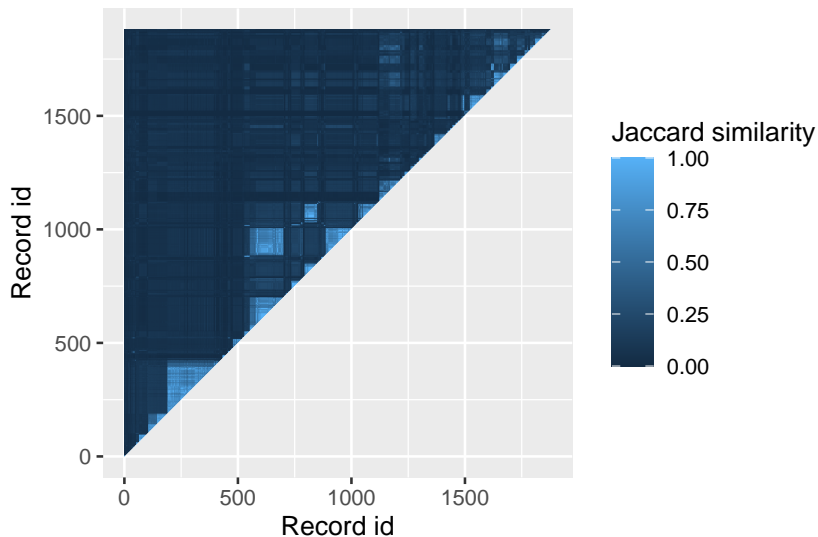


Figure 5: Jaccard similarity for each pair of records. Light blue indicates the two records are more similar and dark blue indicates less similar.

Hashing

For a dataset of size n , the number of comparisons we must compute is

$$\frac{n(n-1)}{2}$$

- ▶ For our set of records, we needed to compute 1,764,381 comparisons
- ▶ A better approach for datasets of any realistic size is to use *hashing*

Hash functions

- ▶ Traditionally, a *hash function* maps objects to integers such that similar objects are far apart
- ▶ Instead, we want special hash functions that do the **opposite** of this, i.e. similar objects are placed closed together!

Definition: Hash function

Hash functions $h()$ are defined such that

*If records A and B have high similarity, then the probability that $h(A) = h(B)$ is **high** and if records A and B have low similarity, then the probability that $h(A) \neq h(B)$ is **low**.*

Hashing shingles

Instead of storing the strings as shingles, we can instead store *hashed values*

These are integers, and they take up less space.

Hashing shingles

```
# instead store hash values (less memory)
hashed_shingles <- apply(dat, 1, function(x) {
  string <- paste(x[-1], collapse=" ") # get the string
  shingles <- tokenize_character_shingles(string, n = 3)[[1]] # 3-shing
  hash_string(shingles) # return hashed shingles
})
```

```
# Jaccard similarity on hashed shingles
hashed_jaccard <- expand.grid(record1 = seq_len(n), record2 = seq_len(n))

# don't need to compare the same things twice
hashed_jaccard <- hashed_jaccard[hashed_jaccard$record1 < hashed_jaccard$record2, ]

time <- Sys.time() # see how long this takes
hashed_jaccard$similarity <- apply(hashed_jaccard, 1, function(pair) {
  jaccard_similarity(hashed_shingles[[pair[1]]], hashed_shingles[[pair[2]]) # get jaccard for each hashed pair
})
time <- difftime(Sys.time(), time, units = "secs") # timing
```

This took up 6.53296×10^5 bytes, while storing the shingles took 8.411816×10^6 bytes; the whole pairwise comparison still took the same amount of time (≈ 1.69 minutes)

Similarity preserving summaries of sets

- ▶ Sets of shingles are large (larger than the original document)
- ▶ If we have millions of documents, it may not be possible to store all the shingle-sets in memory
- ▶ We can replace large sets by smaller representations, called *signatures*
- ▶ And use these signatures to **approximate** Jaccard similarity

Characteristic matrix

In order to get a signature of our data set, we first build a *characteristic matrix*

Columns correspond to records and the rows correspond to all hashed shingles

```
# return if an item is in a list
item_in_list <- function(item, list) {
  as.integer(item %in% list)
}

# get the characteristic matrix
# items are all the unique hash values
# columns will be each record
# we want to keep track of where each hash is included
char_mat <- data.frame(item = unique(unlist(hashed_shingles))

# for each hashed shingle, see if it is in each row
contained <- lapply(hashed_shingles, function(col) {
```

Minhashing

Want create the signature matrix through minhashing

1. Permute the rows of the characteristic matrix m times
2. Iterate over each column of the permuted matrix
3. Populate the signature matrix, row-wise, with the row index from the first 1 value found in the column

The signature matrix is a hashing of values from the permuted characteristic matrix and has one row for the number of permutations calculated (m), and a column for each record

Minhashing (cont'd)

```
# set seed for reproducibility
set.seed(02082018)

# function to get signature for 1 permutation
get_sig <- function(char_mat) {
  # get permutation order
  permute_order <- sample(seq_len(nrow(char_mat)))

  # get min location of "1" for each column (apply(2, ...))
  t(apply(char_mat[permute_order, ], 2, function(col) min(v

}

# repeat many times
m <- 360
sig_mat <- matrix(NA, nrow=m, ncol=ncol(char_mat)) #empty matrix
for(i in 1:m) {
  sig_mat[i, ] <- get_sig(char_mat) #fill matrix
}
```

Signature matrix and Jaccard similarity

The relationship between the random permutations of the characteristic matrix and the Jaccard Similarity is

$$Pr\{\min[h(A)] = \min[h(B)]\} = \frac{|A \cap B|}{|A \cup B|}$$

We use this relationship to **approximate** the similarity between any two records

We look down each column of the signature matrix, and compare it to any other column

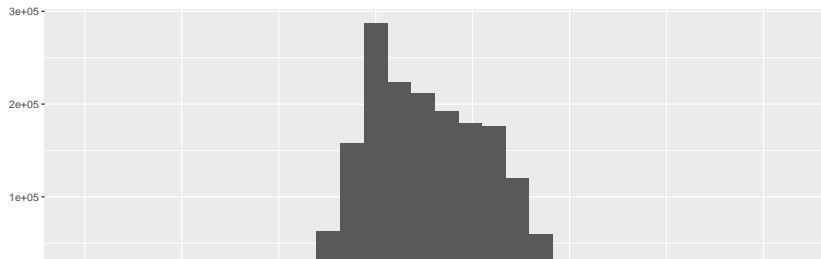
The number of agreements over the total number of combinations is an approximation to Jaccard measure

Jaccard similarity approximation

```
# add jaccard similarity approximated from the minhash to  
# number of agreements over the total number of combination  
hashed_jaccard$similarity_minhash <- apply(hashed_jaccard,  
      sum(sig_mat[, row[["record1"]]] == sig_mat[, row[["record2"]]] /  
    })
```

```
# how far off is this approximation? plot differences  
qplot(hashed_jaccard$similarity_minhash - hashed_jaccard$similarity_minhash,  
      xlab("Difference between Jaccard similarity and minhash approximation"))
```

```
## `stat_bin()` using `bins = 30`. Pick better value with `bins`
```



Locality Sensitive Hashing (LSH)

Locality Sensitive Hashing (LSH)