

# Big Data HW 2

Parham Gilani - 400101859

سریع ۱)  $(A \rightarrow B)$  میزان ارتباط دو آیتم A و B است. به مقدارهای آن داده dataset پیدا میکند و میزان تکرار آن شده A و B را به میزان شده A و B جدا پیدا میکند. اگر  $(A \rightarrow B)$  باشد یعنی به احتمال زیاد A و B به هم رخ می دهند و اگر کمتر از ۱ باشد پس مقدار A امکان مقدار B را کمتر می کند به عبارتی هر دو با هم ظاهر نمی شوند.

۲)  $conv(A \rightarrow B)$  امکان مقدار A بودن B را نشان می دهد در سایه بازمان که هر دو رخ بدهند. وقتی که  $conv(A \rightarrow B)$  بزرگ است به این معنی امکان مقدار A بودن B کم است و اگر  $conv(A \rightarrow B)$  کم باشد به این معنی که امکان مقدار A بودن B زیاد است.

۳)  $confidence$  امکان وقوع B را با توجه به A پیدا میکند و به این دو نظر گرفتن تعداد دفعات که B رخ داده است و شانس رخ دادن آن به آید که B رخ دهد به نسبت زیاد رخ داده باشد. برای مثال:

که در این حالت به احتمال زیاد A رخ دهد چون به ندرت رخ دهد A رخ دهد مقدار B

به درک  $confidence$  میزان را در نظر نمی گیرد در حالت که  $lift$  و  $conv$  اینها را در نظر می گیرند.

total  $\rightarrow 100$   
 $B \rightarrow 80$   
 $A, B \rightarrow 10$

$conf(A \rightarrow B) = 1$

سریع ۲)  $number - freq = \frac{100}{5} \rightarrow number - freq = 20 \rightarrow 1 \times 20$

۲) به میزان  $support\ threshold$  برابر با تمام حیف اعدادی که کوچکترین مغرب مشترکشان کمتر از ۲۰ باشد  $freq$  هستند.

۳)  $total = 482$

$1 \ 2 \ 3 \ \dots \ 100$   
 $1 \ 12 \ 13 \ \dots \ 2^2, 5^2$

$(6,9) \rightarrow 3: confidence = \frac{sup(6,9,3)}{sup(6,9)} = \frac{[\frac{100}{18}]}{[\frac{100}{18}]} = 1$   $(2,4,5) \rightarrow 3: confidence = \frac{sup(2,3,4,5)}{sup(2,3,4,5)} = \frac{[\frac{100}{60}]}{[\frac{100}{24}]} = \frac{1}{5} \rightarrow$

$interest = confidence - sup(3) = 1 - \frac{1}{3} = \frac{2}{3}$   $interest = confidence - sup(3) = \frac{1}{5} - \frac{1}{3} = -\frac{2}{15}$

سریع ۳)  $انف = \frac{I(I-1)}{2} \rightarrow \text{میزان حافظه} = \frac{I(I-1)}{2} \times 4\text{ byte}$

۲)  $maximum\ pairs = \binom{P}{2} \times B$

۳)  $triangle - matrix - size = \frac{I(I-1)}{2} \times 4\text{ byte}$   $(دو ۱) \rightarrow (دو ۱)$

$triple - size = 3 \times 4\text{ byte} \times num - nonZero - pairs \leq \frac{I(I-1)}{2} \times 4\text{ byte} \rightarrow$

$num - nonZero - pairs \leq \frac{I(I-1)}{8}$

$x$	$h_1$	$h_2$	$h_3$	$S_4: row(2,4,3)$	$S_3: row(4,3,5)$	$S_2: (1,2,4)$	$S_1: (4,2,4)$	سریع ۴) $انف$
0	1	3	2	$h_1 \rightarrow 1$	$h_1 \rightarrow 1$	$h_1 \rightarrow 1$	$h_1 \rightarrow 1$	
1	4	1	1	$h_2 \rightarrow 1$	$h_2 \rightarrow 1$	$h_2 \rightarrow 1$	$h_2 \rightarrow 1$	
2	1	5	0	$h_3 \rightarrow 0$	$h_3 \rightarrow 1$	$h_3 \rightarrow 0$	$h_3 \rightarrow 0$	
3	4	3	5					
4	1	1	4					
5	4	5	3					

  

pair	n	U	lacc	minhash n	minhash sim
۴,۳	3	۵ $\rightarrow$ 5	$\frac{1}{6}$	2 num	$\frac{2}{3}$
۴,۲	2	۱ $\rightarrow$ 5	$\frac{1}{6}$	3 num	1
۴,۱	۲,۴	۵ $\rightarrow$ 2,4	$\frac{1}{2}$	3 num	1
3,2	۱,5	۵ $\rightarrow$ 3,5	$\frac{2}{5}$	2 num	$\frac{2}{3}$
3,1	0	۵ $\rightarrow$ 5	$\frac{1}{6}$	2 num	$\frac{2}{3}$
2,1	2	۵ $\rightarrow$ 1,2	$\frac{1}{6}$	3 num	1

۲) چون  $h_3$  مقادیر یکسان برای  $x$  از ۵ تا ۵ دارد پس مقبورات.

۳) چون  $minhash$  آنها تقریباً یکسان است پس  $hash$  ها نتوانستند آنها را از هم جدا کنند و اگر  $hash$  های  $h_1$  و  $h_2$  مقبورت بودند می توانستیم به حالت بهینه ترسیم.

Q1)

A. In this section we have found the most frequent words in first layer and making triplets in second layer. Here is the details:

### **First Pass: Identifying Frequent Words**

- **Tokenization and Counting:**
  - Each abstract is split into individual words.
  - The occurrences of each word across all abstracts are counted using PySpark transformations (``flatMap``, ``map``, ``reduceByKey``).
- **Filtering Frequent Words:**
  - Words that appear fewer times than the specified ``first_pass_support_threshold`` (e.g., 300 occurrences) are discarded.
  - A set of frequent words is returned, which is used in the second pass to restrict the scope of analysis.

### **Second Pass: Identifying Frequent Trigrams**

- **Trigram Generation:**
  - For each abstract, all possible trigrams (combinations of three unique words) are generated using Python's ``itertools.combinations``.
  - Only trigrams where all three words are in the frequent word set (from the first pass) are retained, ensuring efficiency.
- **Counting Trigram Occurrences:**
  - The occurrences of each trigram across all abstracts are counted using similar transformations.
- **Filtering Frequent Trigrams:**
  - Trigrams with occurrences below the ``second_pass_support_threshold`` (e.g., 1000) are discarded.
  - The remaining frequent trigrams are collected as the final result.

B. In this section the PCY algorithm is a little bit different than A-priori algorithm. Which means in the first layer after filtering the words to find the frequent words we have to use a bitmap on them and in the second layer we should make triplets and find if they are frequent or not. I must indicate that we have a parameter named `hash_table_size` which means the hash table isn't unlimited and if it exceeded it will be set by the remaining of that number into the `hash_table_size`. Here is further details:

### Steps in the Algorithm

- **Parameters**
  - **support\_threshold\_first\_pass**: A lower threshold used in the first pass to identify potential candidates.
  - **support\_threshold\_second\_pass**: A higher threshold used in the second pass to filter the final frequent trigrams.
  - **hash\_table\_size**: Size of the hash table to reduce collisions when hashing trigrams.
- **Utility Functions**
  - **hash\_trigram**: Hashes a trigram using SHA-256 and maps it to a bucket in the hash table.
  - **normalize\_trigram**: Sorts the words in a trigram to ensure consistency, treating `(A, B, C)` and `(C, B, A)` as the same.

### First Pass: Generating Frequent Words and Populating the Hash Table

- **Word Counting:**
  - Tokenizes abstracts into words and counts their occurrences using PySpark transformations.
  - Filters words with counts below `'support_threshold_first_pass'` and stores the frequent words in a set.

- **Trigram Generation and Hashing:**
  - Generates all trigrams from abstracts.
  - Filters trigrams to include only those composed of frequent words from the first pass.
  - Hashes each trigram into a bucket using `'hash_trigram'` and populates the hash table with bucket counts.
- **Bitmap Creation:**
  - Converts the hash table into a bitmap where each bucket is marked as frequent (`'True'`) if its count meets the `'support_threshold_first_pass'`. This bitmap reduces memory usage and helps prune infrequent trigrams early.

## **Second Pass: Filtering and Counting Candidate Trigrams**

- **Candidate Filtering Using Bitmap:**
  - Generates trigrams from abstracts, ensuring:
  - All words are frequent (from the first pass).
  - The trigram maps to a bucket marked as frequent in the bitmap.
- **Trigram Counting:**
  - Counts occurrences of the filtered trigrams across all abstracts.
- **Final Filtering**
  - Retains trigrams that meet the `support_threshold_second_pass`.

C. In this section I have simply used jaccard similarity to find If both algorithms are doing a single task and it returns 1 which means the answers are the same.

Q2) In this question I have implemented an LSH that finds similar papers according to their abstracts. At first, we must generate MinHash signatures in order to make Shingles, then generate multiple hash functions and use it to make some Signatures and then for each item in the Signatures we must use LSH to get the answers and then find the similar articles using them and the Jaccard similarity. Here is the further code detail:

## 1. Preprocessing

- **Tokenization:**
  - Splits the text of each paper (e.g., abstract) into tokens (words).
  - Converts all text to lowercase for case-insensitive comparisons.
- **Shingles Generation:**
  - Creates overlapping substrings (shingles) of length (k) from the tokens.
  - Example: For tokens ["`this`", "`is`", "`a`", "`test`"] and (k=3), shingles are: {"`thisisatest`"}

## 2. MinHashing

- **Generate Random Hash Functions:**
  - $h(x) = (a.x + b) \% p$
- **Create MinHash Signatures:**
  - For each paper's shingles, compute the minimum hash value under each of the hash functions.
  - A signature is an array of these minimum hash values.

MinHashing ensures that the similarity between two sets of shingles is preserved.

### 3. Locality-Sensitive Hashing (LSH)

- **Purpose:**
  - Efficiently group and retrieve similar items by hashing similar MinHash signatures into the same buckets.
- **Steps:**
  - **Split MinHash Signatures into Bands:**
    - Divide the signature into (b) bands, each containing (r) rows.
  - **Hash Each Band:**
    - For each band (subarray of the signature), compute a hash.
    - Papers with the same hash for a band are likely to be similar.

### 4. Query for Similar Papers

- **Jaccard Similarity:**
  - Measures the similarity between two sets of shingles:  $J(A, B) = \frac{\text{Intersection}(A,B)}{\text{Union}(A,B)}$
  - If the similarity exceeds a threshold (e.g., 0.5 ), the papers are considered similar.
- **Steps:**
  - Identify "candidate" papers using LSH:
    - For the query paper, look up its bands in the LSH index to find potential matches.
  - Compute Jaccard Similarity for these candidates.

### 5. Main Workflow

- **Data Preparation:**
  - Tokenize and shingle all papers in the dataset.
  - Compute MinHash signatures and build the LSH index.
- **Query Similarity:**
  - For each paper, retrieve candidate matches using LSH.
  - Calculate and rank similar papers based on Jaccard Similarity.