

Problem Set 1

Due 11:59pm Thursday, January 26, 2017

Only one late period is allowed for this homework (11:59pm Tuesday 1/31).

General Instructions

Submission instructions: These questions require thought but do not require long answers. Please be as concise as possible. You should submit your answers as a writeup in PDF format via GradeScope and code via the Snap submission site.

Submitting writeup: Prepare answers to the homework questions into a single PDF file and submit it via <http://gradescope.com>. Make sure that the answer to each question is on a *separate page*. This means you should submit a 15-page PDF (1 page for the cover sheet, 1 page for the answers to question 1, 5 pages for answers to question 2, 3 pages for question 3, and 5 pages for question 4). On top of each page write the number of the question you are answering. Please find the cover sheet and the recommended templates located here:

http://web.stanford.edu/class/cs246/homeworks/hw1/hw1_template.tex

http://web.stanford.edu/class/cs246/homeworks/hw1/hw1_template.pdf

Not including the cover sheet in your submission will result in a 2-point penalty. It is also important to tag your answers correctly on Gradescope. We will deduct $5/N$ points for each incorrectly tagged subproblem (where N is the number of subproblems). This means you can lose up to 5 points for incorrect tagging.

Submitting code: Upload your code at <http://snap.stanford.edu/submit>. Put all the code for a single question into a single file and upload it.

Questions

1 MapReduce (25 pts) [Jessica, Michael, Vinaya]

Write a MapReduce program in Hadoop that implements a simple “People You Might Know” social network friendship recommendation algorithm. The key idea is that if two people have a lot of mutual friends, then the system should recommend that they connect with each other.

Input:

Download the input file from the link: <http://snap.stanford.edu/class/cs246-data/hw1q1.zip>.

The input file contains the adjacency list and has multiple lines in the following format:

<User><TAB><Friends>

Here, <User> is a unique integer ID corresponding to a unique user and <Friends> is a comma-separated list of unique IDs corresponding to the friends of the user with the unique ID <User>. Note that the friendships are mutual (i.e., edges are undirected): if A is friend with B then B is also friend with A . The data provided is consistent with that rule as there is an explicit entry for each side of each edge.

Algorithm: Let us use a simple algorithm such that, for each user U , the algorithm recommends $N = 10$ users who are not already friends with U , but have the largest number of mutual friends in common with U .

Output: The output should contain one line per user in the following format:

<User><TAB><Recommendations>

where <User> is a unique ID corresponding to a user and <Recommendations> is a comma-separated list of unique IDs corresponding to the algorithm's recommendation of people that <User> might know, ordered by decreasing number of mutual friends. Even if a user has fewer than 10 second-degree friends, output all of them in decreasing order of the number of mutual friends. If a user has no friends, you can provide an empty list of recommendations. If there are multiple users with the same number of mutual friends, ties are broken by ordering them in a numerically ascending order of their user IDs.

Also, please provide a description of how you are going to use MapReduce jobs to solve this problem. Don't write more than 3 to 4 sentences for this: we only want a very high-level description of your strategy to tackle this problem.

Note: It is possible to solve this question with a single MapReduce job. But if your solution requires multiple MapReduce jobs, then that's fine too.

What to submit

- (i) Submit the source code via the snap electronic submission website.
- (ii) Include in your writeup a short paragraph describing your algorithm to tackle this problem.
- (iii) Include in your writeup the recommendations for the users with following user IDs: 924, 8941, 8942, 9019, 9020, 9021, 9022, 9990, 9992, 9993.

2 Association Rules (30 pts) [Yixin W, Leon, Junwei]

Association Rules are frequently used for Market Basket Analysis (MBA) by retailers to understand the purchase behavior of their customers. This information can be then used for

many different purposes such as cross-selling and up-selling of products, sales promotions, loyalty programs, store design, discount plans and many others.

Evaluation of item sets: Once you have found the frequent itemsets of a dataset, you need to choose a subset of them as your recommendations. Commonly used metrics for measuring significance and interest for selecting rules for recommendations are:

1. **Confidence** (denoted as $\text{conf}(A \rightarrow B)$): *Confidence* is defined as the probability of occurrence of B in the basket if the basket already contains A :

$$\text{conf}(A \rightarrow B) = \Pr(B|A),$$

where $\Pr(B|A)$ is the conditional probability of finding item set B given that item set A is present.

2. **Lift** (denoted as $\text{lift}(A \rightarrow B)$): *Lift* measures how much more “ A and B occur together” than “what would be expected if A and B were statistically independent”:

$$\text{lift}(A \rightarrow B) = \frac{\text{conf}(A \rightarrow B)}{S(B)},$$

where $S(B) = \frac{\text{Support}(B)}{N}$ and N = total number of transactions (baskets).

3. **Conviction** (denoted as $\text{conv}(A \rightarrow B)$): *Conviction* compares the “probability that A appears without B if they were independent” with the “actual frequency of the appearance of A without B ”:

$$\text{conv}(A \rightarrow B) = \frac{1 - S(B)}{1 - \text{conf}(A \rightarrow B)}.$$

(a) [3pts]

A drawback of using *confidence* is that it ignores $\Pr(B)$. Why is this a drawback? Explain why *lift* and *conviction* do not suffer from this drawback?

(b) [3pts]

A measure is *symmetrical* if $\text{measure}(A \rightarrow B) = \text{measure}(B \rightarrow A)$. Which of the measures presented here are symmetrical? For each measure, please provide either a proof that the measure is symmetrical, or a counterexample that shows the measure is not symmetrical.

(c) [4pts]

A measure is *desirable* if its value is maximal for rules that hold 100% of the time (such rules are called *perfect implications*). This makes it easy to identify the best rules. Which of the above measures have this property? Explain why.

Product Recommendations: The action or practice of selling additional products or services to existing customers is called cross-selling. Giving product recommendation is one of the examples of cross-selling that are frequently used by online retailers. One simple method to give product recommendations is to recommend products that are frequently browsed together by the customers.

Suppose we want to recommend new products to the customer based on the products they have already browsed on the online website. Write a program using the A-priori algorithm to find products which are frequently browsed together. Fix the support to $s = 100$ (i.e. product pairs need to occur together at least 100 times to be considered frequent) and find itemsets of size 2 and 3.

Use the online browsing behavior dataset at: <http://snap.stanford.edu/class/cs246-data/browsing.txt>. Each line represents a browsing session of a customer. On each line, each string of 8 characters represents the id of an item browsed during that session. The items are separated by spaces.

Note: for parts (d) and (e), the writeup will require a specific rule ordering but the program need not sort the output.

(d) [10pts]

Identify pairs of items (X, Y) such that the support of $\{X, Y\}$ is at least 100. For all such pairs, compute the confidence scores of the corresponding association rules: $X \Rightarrow Y, Y \Rightarrow X$. Sort the rules in decreasing order of *confidence* scores and list the top 5 rules in the writeup. Break ties, if any, by lexicographically increasing order on the left hand side of the rule.

(e) [10pts]

Identify item triples (X, Y, Z) such that the support of $\{X, Y, Z\}$ is at least 100. For all such triples, compute the *confidence* scores of the corresponding association rules: $(X, Y) \Rightarrow Z, (X, Z) \Rightarrow Y, (Y, Z) \Rightarrow X$. Sort the rules in decreasing order of *confidence* scores and list the top 5 rules in the writeup. Order the left-hand-side pair lexicographically and break ties, if any, by lexicographical order of the first then the second item in the pair.

What to submit

Upload all the code on snap and include the following in your writeup:

- (i) Explanation for 2(a).
- (ii) Proofs and/or counterexamples for 2(b).
- (iii) Explanation for 2(c).

(iv) Top 5 rules with confidence scores [2(d)].

(v) Top 5 rules with confidence scores [2(e)].

3 Locality-Sensitive Hashing (15 pts) [Luda, Yixin C, Sachin]

When simulating a random permutation of rows, as described in Sect. 3.3.5 of MMDS, we could save a lot of time if we restricted our attention to a randomly chosen k of the n rows, rather than hashing all the row numbers. The downside of doing so is that if none of the k rows contains a 1 in a certain column, then the result of the minhashing is “don’t know,” i.e., we get no row number as a minhash value. It would be a mistake to assume that two columns that both minhash to “don’t know” are likely to be similar. However, if the probability of getting “don’t know” as a minhash value is small, we can tolerate the situation, and simply ignore such minhash values when computing the fraction of minhashes in which two columns agree.

(a) [5pts]

Suppose a column has m 1’s and therefore $n - m$ 0’s. Prove that the probability we get “don’t know” as the minhash value for this column is at most $(\frac{n-k}{n})^m$.

(b) [5pts]

Suppose we want the probability of “don’t know” to be at most e^{-10} . Assuming n and m are both very large (but n is much larger than m or k), give a simple approximation to the smallest value of k that will assure this probability is at most e^{-10} . Hints: (1) You can use $(\frac{n-k}{n})^m$ as the exact value of the probability of “don’t know.” (2) Remember that for large x , $(1 - \frac{1}{x})^x \approx 1/e$.

(c) [5pts]

Note: This question should be considered separate from the previous two parts, in that we are no longer restricting our attention to a randomly chosen subset of the rows.

When minhashing, one might expect that we could estimate the Jaccard similarity without using all possible permutations of rows. For example, we could only allow cyclic permutations i.e., start at a randomly chosen row r , which becomes the first in the order, followed by rows $r + 1$, $r + 2$, and so on, down to the last row, and then continuing with the first row, second row, and so on, down to row $r - 1$. There are only n such permutations if there are n rows. However, these permutations are not sufficient to estimate the Jaccard similarity correctly.

Give an example of two columns such that the probability (over cyclic permutations only) that their minhash values agree is not the same as their Jaccard similarity. In your answer, please provide (a) an example of a matrix with two columns (let the two columns correspond to sets denoted by $S1$ and $S2$) (b) the Jaccard similarity of $S1$ and $S2$ (c) the probability that a random cyclic permutation yields the same minhash value for both $S1$ and $S2$.

What to submit

Include the following in your writeup:

- (i) Proof for 3(a)
- (ii) Derivation and final answer for 3(b)
- (iii) Example for 3(c)

4 LSH for Approximate Near Neighbor Search (30 pts) [Nihit, Rishabh, Anthony, Naveen]

In this problem, we study the application of LSH to the problem of finding approximate near neighbors.

Assume we have a dataset \mathcal{A} of n points in a metric space with distance metric $d(\cdot, \cdot)$. Let c be a constant greater than 1. Then, the (c, λ) -Approximate Near Neighbor (ANN) problem is defined as follows: Given a query point z , assuming that there is a point x in the dataset with $d(x, z) \leq \lambda$, return a point x' from the dataset with $d(x', z) \leq c\lambda$ (this point is called a (c, λ) -ANN). The parameter c therefore represents the maximum approximation factor allowed in the problem.

Let us consider a LSH family \mathcal{H} of hash functions that is $(\lambda, c\lambda, p_1, p_2)$ -sensitive for the distance measure $d(\cdot, \cdot)$. Let¹ $\mathcal{G} = \mathcal{H}^k = \{g = (h_1, \dots, h_k) | h_i \in \mathcal{H}, \forall 1 \leq i \leq k\}$, where $k = \log_{1/p_2}(n)$.

Let us consider the following procedure:

1. Select $L = n^\rho$ random members g_1, \dots, g_L of \mathcal{G} , where $\rho = \frac{\log(1/p_1)}{\log(1/p_2)}$.
2. Hash all the data points as well as the query point using all g_i ($1 \leq i \leq L$).
3. Retrieve at most² $3L$ data points (chosen uniformly at random) from the set of L buckets to which the query point hashes.

¹The equality $\mathcal{G} = \mathcal{H}^k$ is saying that every function of \mathcal{G} is an AND-construction of k functions of \mathcal{H} , so $g(x) = g(y)$ only if $h_i(x) = h_i(y)$ for every h_i underlying g .

²If there are fewer than $3L$ data points hashing to the same buckets as the query point, just take all of them.

4. Among the points selected in phase 3, report the one that is the closest to the query point as a (c, λ) -ANN.

The goal of the first part of this problem is to show that this procedure leads to a correct answer with constant probability.

(a) [5 pts]

Let $W_j = \{x \in \mathcal{A} | g_j(x) = g_j(z)\}$ ($1 \leq j \leq L$) be the set of data points x mapping to the same value as the query point z by the hash function g_j . Define $T = \{x \in \mathcal{A} | d(x, z) > c\lambda\}$. Prove:

$$\Pr \left[\sum_{j=1}^L |T \cap W_j| \geq 3L \right] \leq \frac{1}{3}.$$

(Hint: Markov's Inequality.)

(b) [5 pts]

Let $x^* \in \mathcal{A}$ be a point such that $d(x^*, z) \leq \lambda$. Prove:

$$\Pr [\forall 1 \leq j \leq L, g_j(x^*) \neq g_j(z)] < \frac{1}{e}.$$

(c) [5 pts]

Conclude that with probability greater than some fixed constant the reported point is an actual (c, λ) -ANN.

(d) [15 pts]

A dataset of images,³ `patches.mat`, is provided in: <http://snap.stanford.edu/class/cs246-data/lsh.zip>. For this problem, if you don't have matlab on your computer, you may want to use matlab on corn. To do so execute

```
ssh -X <SUNET ID>@corn.stanford.edu
(Your stanford email password)
module load matlab
matlab
```

³Dataset and code adopted from Brown University's Greg Shakhnarovich

Each column in this dataset is a 20×20 image patch represented as a 400-dimensional vector. We will use the L_1 distance metric on \mathbb{R}^{400} to define similarity of images. We would like to compare the performance of LSH-based approximate near neighbor search with that of linear search.⁴ You should use the code provided with the dataset for this task. The included `ReadMe.txt` file explains how to use the provided code. In particular, you will need to use the functions `lsh` and `lshlookup`. The parameters $L = 10, k = 24$ work for this exercise, but feel free to use other parameter values as long as you explain the reason behind your parameter choice.

- For each of the image patches in columns 100, 200, 300, \dots , 1000, find the top 3 near neighbors⁵ (excluding the original patch itself) using both LSH and linear search. What is the average search time for LSH? What about for linear search?
- Assuming $\{z_j \mid 1 \leq j \leq 10\}$ to be the set of image patches considered (*i.e.*, z_j is the image patch in column $100j$), $\{x_{ij}\}_{i=1}^3$ to be the approximate near neighbors of z_j found using LSH, and $\{x_{ij}^*\}_{i=1}^3$ to be the (true) top 3 near neighbors of z_j found using linear search, compute the following error measure:

$$error = \frac{1}{10} \sum_{j=1}^{10} \frac{\sum_{i=1}^3 d(x_{ij}, z_j)}{\sum_{i=1}^3 d(x_{ij}^*, z_j)}$$

Plot the error value as a function of L (for $L = 10, 12, 14, \dots, 20$, with $k = 24$). Similarly, plot the error value as a function of k (for $k = 16, 18, 20, 22, 24$ with $L = 10$). Briefly comment on the two plots (one sentence per plot would be sufficient).

- Finally, plot the top 10 near neighbors found⁶ using the two methods (using the default $L = 10, k = 24$ or your alternative choice of parameter values for LSH) for the image patch in column 100, together with the image patch itself. You may find the functions `reshape()` and `mat2gray()` useful to convert the matrices to images; you can also use the functions `imshow()` and `subplot()` to display the images.
How do they compare visually?

What to submit

- Include the proof for 4(a) in your writeup.
- Include the proof for 4(b) in your writeup.
- Include the reasoning for why the reported point is an actual (c, λ) -ANN in your writeup [4(c)].

⁴By linear search we mean comparing the query point z directly with every database point x .

⁵Sometimes, the function `nnlsh` may return less than 3 nearest neighbors. You can use a `while` loop to check that `lshlookup` returns enough results, or you can manually run the program multiple times until it returns the correct number of neighbors.

⁶Same remark, you may sometimes have less than 10 nearest neighbors in your results; you can use the same hacks to bypass this problem.

(iv) Include the following in your writeup for 4(d):

- Average search time for LSH and linear search.
- Plots for error value vs. L and error value vs. K , and brief comments for each plot
- Plot of 10 nearest neighbors found by the two methods (also include the original image) and brief visual comparison

(v) Upload the matlab code for 4(d) on snap.