

# CS246: Mining Massive Datasets

Assignment number: 3\_\_\_\_\_

Submission time: \_\_\_\_\_ and date: \_\_\_\_\_

Fill in and include this cover sheet with each of your assignments. It is an honor code violation to write down the wrong time. Assignments are due at 9:30 am, either handed in at the beginning of class or left in the submission box on the 1<sup>st</sup> floor of the Gates building, near the east entrance. Failure to include the coversheet with your assignment will be penalized by 2 points.

Each student will have a total of *two* free late periods. *One late period expires at the start of each class.* (Assignments are due on Thursdays, which means the first late period expires on the following Tuesday at 9:30am.) Once these late periods are exhausted, any assignments turned in late will be penalized 50% per late period. However, no assignment will be accepted more than *one* late period after its due date. (If an assignment is due to Thursday then we will not accept it after the following Thursday.)

**Your name:** Blaž Sovdat\_\_\_\_\_

**Email:** blaz.sovdat@gmail.com\_\_\_\_\_ **SUNet**

**ID:** \_\_\_\_\_

Collaborators:\_\_\_\_\_

I acknowledge and accept the Honor Code.

(Signed)\_\_\_\_\_

(For CS246 staff only)

Late days: 0 1

Section	Score
1	
2	
3	
4	
Total	

Comments: /

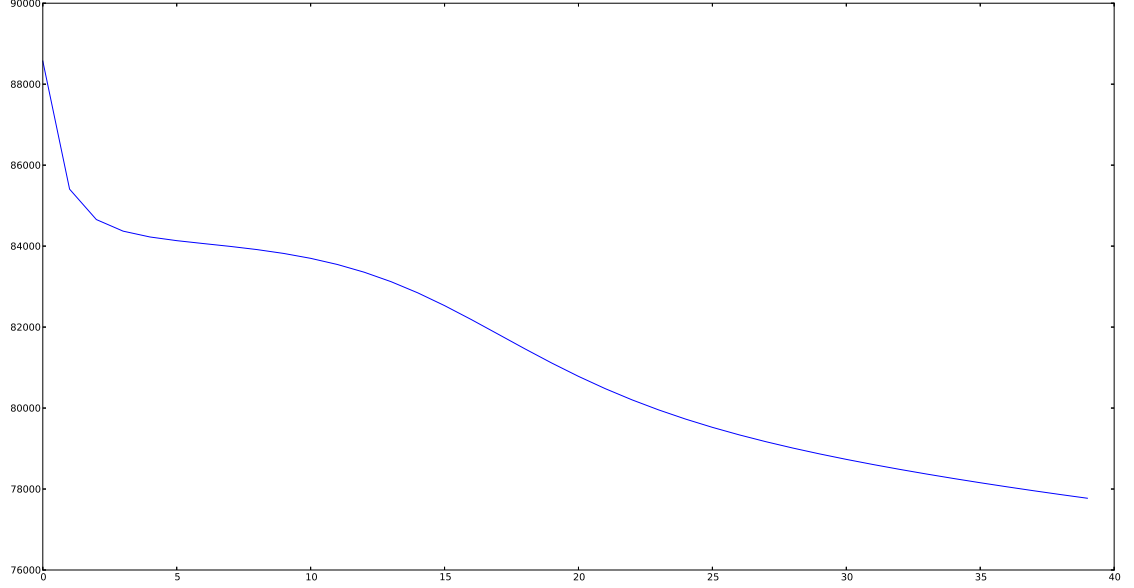


Figure 1: Error in each of the 40 iterations of our algorithm.

**Note.** You can find all source code in appendix A at the end of the document.

## 1 Latent Features for Recommendations

- (a) We have  $\epsilon_{iu} = \frac{\partial}{\partial R_{iu}} E = 2(R_{iu} - q_i p_u^T)$ . This means  $\epsilon_{iu}$  is the prediction error (times two) we make for prediction of rating user  $u$  made for item  $i$ . (Note that in our code we use  $\epsilon_{iu} = R_{iu} - q_i p_u^T$ .) The update equations are as follows:

$$\begin{aligned} q_i &:= q_i + \mu_1(\epsilon_{iu} p_u - \lambda_1 q_i), \\ p_u &:= p_u + \mu_2(\epsilon_{iu} q_i - \lambda_2 p_u). \end{aligned}$$

In our case we have  $\mu_1 = \mu_2 = \eta$  and  $\lambda_1 = \lambda_2 = \lambda$ .

- (b) See listing 1 for our implementation of the stochastic gradient descent algorithm. See figure 1 for plot of the value of the objective function on the training set as a function of the number of iterations. We initialized  $P$  and  $Q$  with values  $[0, \sqrt{5/k}]$ , picked uniformly at random. Furthermore, we found that setting learning rate  $\eta := 0.01565$  works well.
- (c) Define the training error as  $E_{\text{tr}} := \sum_{(i,u) \in \text{train}} (R_{iu} - q_i p_u^T)^2$  and the test error as  $E_{\text{te}} := \sum_{(i,u) \in \text{test}} (R_{iu} - q_i p_u^T)^2$ . See listing 2 for the modified code. See figure 2 for errors with and without regularization on the test set and figure 3 for errors with and without regularization on the training set. True statements are the following ones:

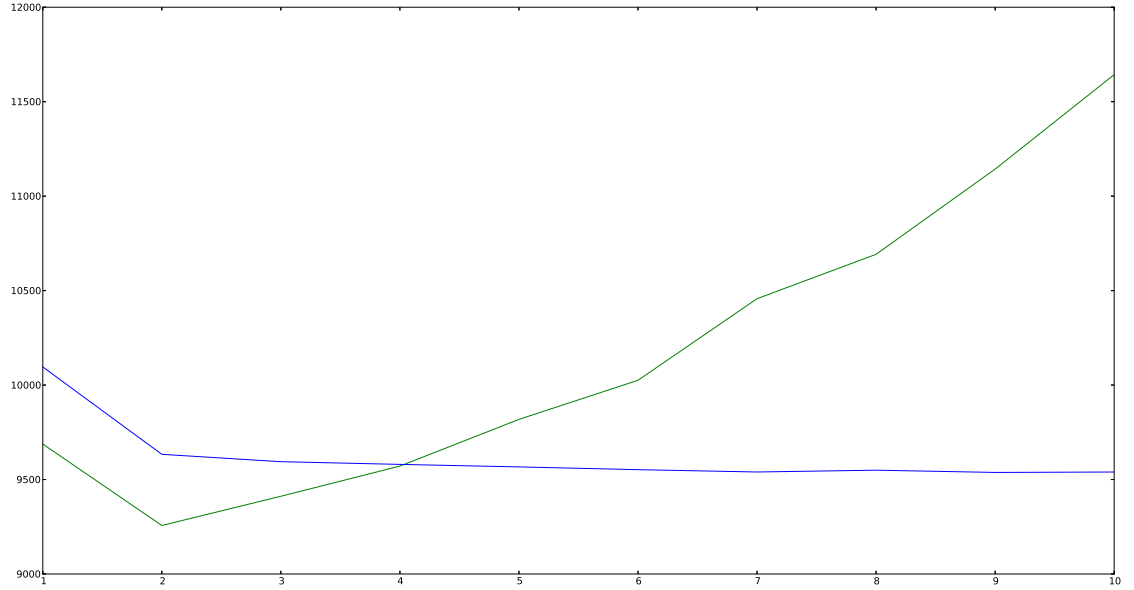


Figure 2: Plot for  $E_{te}$  with no regularization (green) and  $E_{te}$  with  $\lambda = 0.2$  (blue) as a function of  $k \in \{1, 2, \dots, 10\}$ .

- B: Regularization decreases the test error for  $k \geq 5$ .
- D: Regularization increases the training error for all (or almost all)  $k$ .
- H: Regularization decreases overfitting.

## 2 PageRank Computation

- (a) Define  $r^{(0)} = \frac{1}{n}e$  and  $r^{(k)} = \frac{1-\beta}{n}e + \beta M r^{(k-1)}$  for  $k > 0$  and  $\beta \in (0, 1)$ , where  $e$  is the unit vector of dimension  $n$ . Furthermore let  $r$  be the PageRank vector, so we have  $r = \frac{1-\beta}{n}e + \beta M r$ . We now prove  $\|r - r^{(k)}\|_1 \leq 2\beta^k$  for all  $k \geq 0$  by induction. First note  $\|r - r^{(0)}\|_1 \leq 2$ , so we have the base case. (Norm is largest when  $p_i = 1$  for some  $i$  and  $p_j = 0$  for  $i \neq j$ . So  $\|r - r^{(0)}\|_1 \leq 1 - 1/n + n \cdot 1/n = 2(n-1)/n \leq 2$ .) Now suppose

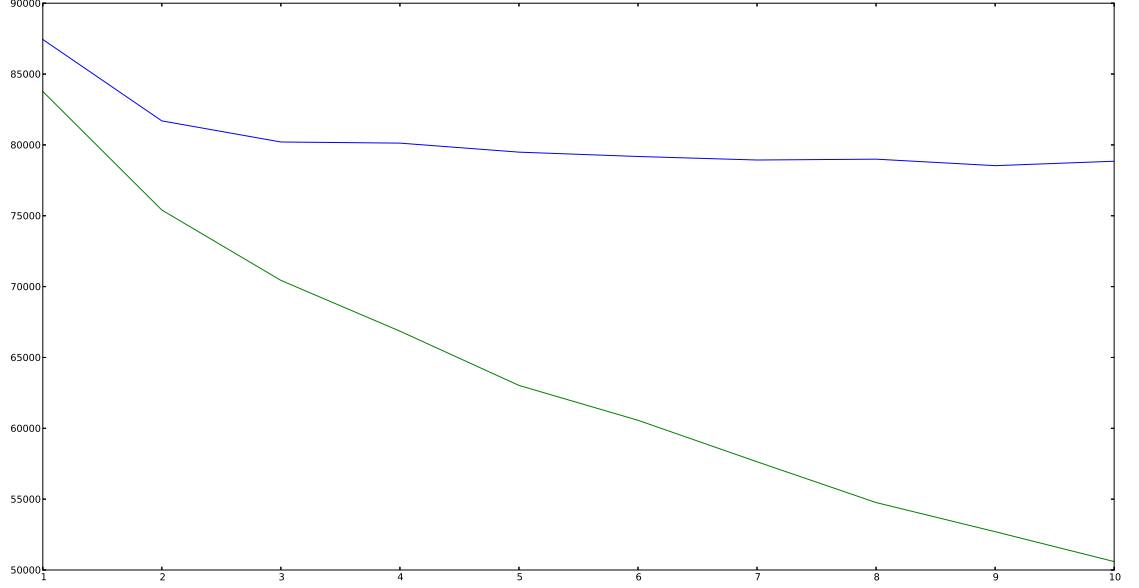


Figure 3: Plot for  $E_{\text{tr}}$  with no regularization (green) and  $E_{\text{tr}}$  with  $\lambda = 0.2$  (blue) as a function of  $k \in \{1, 2, \dots, 10\}$ .

$\|r - r^{(\ell)}\|_1 \leq 2\beta^\ell$  for  $\ell = 0, 1, \dots, k$ . We now do the inductive step:

$$\begin{aligned}
\|r - r^{(k+1)}\|_1 &= \left\| \frac{1-\beta}{n}e + \beta Mr - \frac{1-\beta}{n}e - \beta Mr^{(k)} \right\|_1 \\
&= \|\beta M(r - r^{(k)})\|_1 \\
&= \beta \|M(r - r^{(k)})\|_1 \\
&\leq \beta \|M\|_1 \|r - r^{(k)}\|_1 \\
&\leq 2\beta^{k+1},
\end{aligned}$$

establishing the inequality. (Note that for the matrix  $L_1$  norm we have  $\|M\|_1 = \max_{1 \leq j \leq n} \sum_{i=1}^n |M_{ij}| \leq 1$ , because  $M$  is column stochastic. We also used  $\|Mv\|_1 \leq \|M\|_1 \|v\|_1$  for matrix  $M$  and vector  $v$ .)

(b) Suppose  $\|r - r^{(k)}\|_1 \leq 2\beta^k$  and let  $\delta \in (0, 1)$  be constant. We then solve

$$\|r - r^{(k)}\|_1 \leq 2\beta^k < \delta$$

for  $k$  and get  $k > \log(\delta/2)/\log(1/\beta)$ . Since  $\delta$  is constant this means we need  $k = O(1/\log(1/\beta))$  iterations to get error at most  $\delta$ . Because the cost of each iteration is  $O(m)$ , we have that the total running time required to get the desired error is  $O(m/\log(1/\beta))$ . (We can compute  $(1-\beta)e/n + \beta Mr^{(k)}$  in  $O(m)$ .) This means that after  $O(m/\log(1/\beta))$  steps we have error less than  $\delta$ .

- (c) Show  $\mathbb{E}[\tilde{r}_j] = \frac{1-\beta}{nR} \mathbb{E}[\text{visits}(j)] = r_j$ , where  $r_j$  is  $j$ -th component of the true PageRank vector and  $\mathbb{E}[\tilde{r}_j]$  is the expected number of times the node  $j$  is visited, taken over all  $nR$  random walks. Recall that for PageRank we have

$$r_j = \frac{1-\beta}{n} + \beta \sum_{(i \rightarrow j) \in E} \frac{r_i}{\deg(i)},$$

which is equivalent to  $r = \frac{1-\beta}{n}e + \beta M r$ . Could we somehow get  $\mathbb{E}[\tilde{r}] = \frac{1-\beta}{n}e + \beta M \mathbb{E}[\tilde{r}]$ ?

- (d) We now show that the running time of the MC algorithm is  $O(\frac{nR}{1-\beta})$  if we run it  $R$  times from each node on a graph with  $n$  vertices. First, let  $T(v)$  be running time of the algorithm when run from the vertex  $v$ . After updating the count for  $v$ , it will continue the walk with probability  $\beta$  and terminate with probability  $1 - \beta$ , so we clearly have  $T(v) = O(1) + \beta T(v)$ , which gives us the geometric series

$$T(v) = \sum_{k \geq 0} \beta^k O(1) = O\left(\frac{1}{1-\beta}\right).$$

Since we run the algorithm  $R$  times for each node, we have the total running time of  $T(v)nR = O(\frac{nR}{1-\beta})$ .

- (e) We implemented both PageRank with Power iteration and the MC algorithm in python; for the code see listing 3.
- The running time is 0.008119 seconds. See listing 3.
  - For code for this question see listing 3. To compute the running times we took the mean running time of 50 runs for each  $R = 1, 3, 5$ . For  $R = 1$  the running time is 0.0031599 seconds. For  $R = 3$  the running time is 0.005479 seconds. For  $R = 5$  the running time is 0.00712 seconds. The absolute errors — see listing 3 for implementation — were as in table 1, where error is defined as

$$E_K := \frac{1}{K} \sum_{i=1}^K |\tilde{r}_i - r_i|,$$

for approximate PageRank vector  $\tilde{r}$  and “true” PageRank  $r$ . (See HW text for detailed description of the error metric.)

### 3 Similarity Ranking

- (a) Let  $C_1 := C_2 := 0.8$ . We compute 3 steps of similarity for  $s_A(\text{camera}, \text{phone})$  and  $s_A(\text{camera}, \text{printer})$ .

$R$	$K$	Error
1	10	0.0106990249484
1	30	0.00947132554345
1	50	0.00766502952898
1	All	0.0073999233725
3	10	0.00989902494841
3	30	0.00898243665457
3	50	0.00694048754126
3	All	0.00695841712359
5	10	0.00983703989898
5	30	0.00843815848515
5	50	0.00682750626194
5	All	0.00641664075857

Table 1: Average absolute errors.

Step 1. We have

$$\begin{aligned}
s_A(\text{camera, phone}) &= \frac{0.8}{6} (s_B(\text{nokia, nokia}) + s_B(\text{nokia, apple}) + s_B(\text{kodak, nokia}) + s_B(\text{kodak, apple}) + s_B(\text{canon, nokia}) + s_B(\text{canon, apple})) \\
&= \frac{0.8}{6} \cdot 1 = \frac{4}{5 \cdot 6} = 2/15 \\
s_A(\text{camera, printer}) &= \frac{0.8}{3} (s_B(\text{nokia, hp}) + s_B(\text{kodak, hp}) + s_B(\text{canon, hp})) = 0 \\
s_A(\text{phone, printer}) &= \frac{0.8}{2} (s_B(\text{nokia, hp}) + s_B(\text{apple, hp})) = 0 \\
s_B(\text{nokia, apple}) &= \frac{0.8}{2} (s_A(\text{phone, phone}) + s_A(\text{phone, camera})) = 2/5 \\
s_B(\text{nokia, hp}) &= \frac{0.8}{2} (s_A(\text{phone, printer}) + s_A(\text{camera, printer})) = 0 \\
s_B(\text{nokia, kodak}) &= \frac{0.8}{2} (s_A(\text{phone, camera}) + s_A(\text{camera, camera})) = 2/5 \\
s_B(\text{nokia, canon}) &= \frac{0.8}{2} (s_A(\text{phone, camera}) + s_A(\text{camera, camera})) = 2/5 \\
s_B(\text{kodak, apple}) &= \frac{0.8}{1} (s_A(\text{phone, camera})) = 0 \\
s_B(\text{kodak, hp}) &= \frac{0.8}{1} (s_A(\text{printer, camera})) = 0 \\
s_B(\text{kodak, canon}) &= \frac{0.8}{1} (s_A(\text{camera, camera})) = 4/5 \\
s_B(\text{canon, apple}) &= \frac{0.8}{1} (s_A(\text{phone, camera})) = 0 \\
s_B(\text{canon, hp}) &= \frac{0.8}{1} (s_A(\text{printer, camera})) = 0 \\
s_B(\text{apple, hp}) &= \frac{0.8}{1} (s_A(\text{phone, printer})) = 0
\end{aligned}$$

Step 2. We now use similarities from step 1.

$$\begin{aligned}
s_A(\text{camera}, \text{phone}) &= \frac{0.8}{6} \left(1 + 3 \frac{2}{5}\right) = 22/75 & s_A(\text{camera}, \text{printer}) &= 0 \\
s_A(\text{phone}, \text{printer}) &= 0 & s_B(\text{nokia}, \text{apple}) &= 34/75 \\
s_B(\text{nokia}, \text{hp}) &= 0 & s_B(\text{nokia}, \text{kodak}) &= 34/75 \\
s_B(\text{nokia}, \text{canon}) &= 34/75 & s_B(\text{kodak}, \text{apple}) &= 8/75 \\
s_B(\text{kodak}, \text{hp}) &= 0 & s_B(\text{kodak}, \text{canon}) &= 4/5 \\
s_B(\text{canon}, \text{apple}) &= 8/75 & s_B(\text{canon}, \text{hp}) &= 0 \\
s_B(\text{apple}, \text{hp}) &= 0 & &
\end{aligned}$$

Step 3. Finally we compute 3 step using similarities from the previous one:

$$s_A(\text{camera}, \text{phone}) = \frac{4}{5 \cdot 6} \left(1 + \frac{2 \cdot 8}{75} + \frac{3 \cdot 34}{75}\right) = 0.343 \quad s_A(\text{camera}, \text{printer}) = 0$$

- (b) The similarity scores above assume that all edges are equally relevant. Let's say we also have information about how many times a URL is clicked for a given query (or equivalently, many users bought how many items). For a pair  $(X, y)$ , we denote this information by the weight  $W_{(X, y)}$ . We propose the following:

$$\begin{aligned}
s_A(X, Y) &= \frac{C_1}{\sum_{i=1}^{|O(X)|} \sum_{j=1}^{|O(Y)|} W_{X, O_i(X)} W_{Y, O_j(Y)}} \sum_{i=1}^{|O(X)|} \sum_{j=1}^{|O(Y)|} W_{X, O_i(X)} W_{Y, O_j(Y)} s_B(O_i(X), O_j(Y)), \\
s_B(x, y) &= \frac{C_2}{\sum_{i=1}^{|I(x)|} \sum_{j=1}^{|I(y)|} W_{x, I_i(x)} W_{y, I_j(y)}} \sum_{i=1}^{|I(x)|} \sum_{j=1}^{|I(y)|} W_{x, I_i(x)} W_{y, I_j(y)} s_B(I_i(x), I_j(y)).
\end{aligned}$$

- (c) Let  $C_1 := C_2 := 0.8$ . We compute 3 iterations of  $s_A(a_1, a_2)$  for complete bipartite graphs  $K_{2,1}$  and  $K_{2,2}$ . For  $K_{2,1}$  with  $A = \{a_1, a_2\}$  and  $B = \{b_1\}$  we have:

Step 1) Note  $s_B(b_1, b_1) = 1$  by “definition”. (We would not need to even write these down as they never change; but the author felt including them would make no harm.)

$$s_A(a_1, a_2) = 0.8 s_B(b_1, b_1) = 0.8 = 4/5 \quad s_B(b_1, b_1) = 1.$$

Step 2) We now use previous scores:

$$s_A(a_1, a_2) = 0.8 s_B(b_1, b_1) = 4/5 \quad s_B(b_1, b_1) = 1.$$

Step 3) Similarly for the final step

$$s_A(a_1, a_2) = 0.8 s_B(b_1, b_1) = 4/5 \quad s_B(b_1, b_1) = 1.$$

We now repeat the procedure for  $K_{2,2}$  with  $A = \{a_1, a_2\}$  and  $B = \{b_1, b_2\}$ :

Step 1) Using initialization values we have

$$\begin{aligned} s_A(a_1, a_2) &= \frac{0.8}{4}(s_B(b_1, b_1) + s_B(b_1, b_2) + s_B(b_2, b_2)) = \frac{0.8}{4}2 = 4/10, \\ s_B(b_1, b_2) &= \frac{0.8}{4}(s_A(a_1, a_1) + s_A(a_1, a_2) + s_A(a_2, a_2)) = \frac{0.8}{4}2 = 4/10. \end{aligned}$$

Step 2) We now use scores from step 1:

$$\begin{aligned} s_A(a_1, a_2) &= \frac{0.8}{4}(s_B(b_1, b_1) + s_B(b_1, b_2) + s_B(b_2, b_2)) = \frac{0.8}{4}(2 + 4/10) \\ &= 24/50, \\ s_B(b_1, b_2) &= \frac{0.8}{4}(s_A(a_1, a_1) + s_A(a_1, a_2) + s_A(a_2, a_2)) = \frac{0.8}{4}(2 + 4/10) \\ &= 24/50. \end{aligned}$$

Step 3) Finally, compute third step using scores from step 2:

$$\begin{aligned} s_A(a_1, a_2) &= \frac{0.8}{4}(s_B(b_1, b_1) + s_B(b_1, b_2) + s_B(b_2, b_1) + s_B(b_2, b_2)) = \frac{0.8}{4}(2 + \frac{24}{50}) \\ &= 124/250, \\ s_B(b_1, b_2) &= \frac{0.8}{4}(s_A(a_1, a_1) + s_A(a_1, a_2) + s_A(a_2, a_1) + s_A(a_2, a_2)) = \frac{0.8}{4}(2 + \frac{24}{50}) \\ &= 124/250. \end{aligned}$$

For  $K_{2,1}$  the scores  $s_A(a_1, a_2)$  do not change through iterations, while for  $K_{2,2}$  they converge to  $1/2$ .

## 4 Dense Communities in Networks

We here use the well-known handshaking lemma, stating  $2|E(G)| = \sum_{v \in V(G)} \deg_G(v)$ . (For example,  $\rho(S) = \frac{1}{2|S|} \sum_{v \in S} \deg_S(v)$ .)

(a) We first prove lower bounds on the size of  $A(S)$  and on the number of iterations of the algorithm.

(i) We show  $|A(S)| \geq \frac{\epsilon}{1+\epsilon}|S|$ . Note that we have

$$2|E[S]| = \sum_{v \in S} \deg_S(v) \geq \sum_{v \in S \setminus A(S)} \deg_S(v) \geq 2(1 + \epsilon) \frac{|E[S]|}{|S|} (|S| - |A(S)|),$$

implying  $|S| \geq (|S| - |A(S)|)(1 + \epsilon)$ , which immediately yields  $|A(S)| \geq \frac{\epsilon}{1+\epsilon}|S|$ . (We have  $\sum_{v \in S \setminus A(S)} \deg_S(v) \geq 2(1 + \epsilon) \frac{|E[S]|}{|S|} (|S| - |A(S)|)$  because by definition each node in  $S \setminus A(S)$  has degree at least  $2(1 + \epsilon)\rho(S)$  and because  $S$  and  $A(S)$  are disjoint there are  $|S| - |A(s)|$  nodes.)



- (ii) Let  $|S|$  be the current size. In next iteration we are left with at most  $\frac{\epsilon}{1+\epsilon}|S|$  nodes. This means  $|S| - \frac{\epsilon}{1+\epsilon}|S| = \frac{1}{1+\epsilon}|S|$  for each iteration, i.e., in each iteration the size decreases by factor of at least  $(1 + \epsilon)$ . Since  $|S| = |V| = n$  initially, this means we will obviously make at most  $1 + \log_{1+\epsilon} n$  steps before  $S = \emptyset$ . (If  $n = 1$  the algorithm clearly needs one step, while  $\log_{1+\epsilon} n = 0$ .)
- (b) We now show that the density of the set the algorithm returns is at most  $2(1 + \epsilon)$  times smaller than  $\rho^*(G)$ .

- (i) Let  $S^* := \arg \max_{S \subseteq V(G)} \frac{|E[S]|}{|S|}$  and for the sake of contradiction assume there exists  $v \in S^*$  such that  $\deg_{S^*}(v) \leq \rho^*(G) - 1$ . Now note that

$$\begin{aligned} \rho(S^* \setminus \{v\}) &= \frac{|E[S^*]| - \deg_{S^*}(v)}{|S^*| - 1} \\ &\geq \frac{|E[S^*]| - \rho^*(G) + 1}{|S^*| - 1} \\ &= \frac{\rho^*(G)(|S^*| - 1) + 1}{|S^*| - 1} \\ &= \rho^*(G) + \frac{1}{|S^*| - 1}, \end{aligned}$$

meaning that  $S^* \setminus \{v\}$  is denser than  $S^*$ . This contradicts the maximality assumption, establishing  $\deg_{S^*}(v) \geq \rho^*(G)$  for all  $v \in S^*$ . We used  $|E[S^*]| = |S^*|\rho^*(G)$ .

- (ii) Consider the first iteration of the while loop so that there is a node  $v \in S^* \cap A(S)$ . We show that then  $2(1 + \epsilon)\rho(S) \geq \rho^*(G)$ . Since  $v \in A(S)$  we have  $\deg_S(v) \geq 2(1 + \epsilon)\rho(S)$ . Because this is the first iteration such that  $S^* \cap A(S) \neq \emptyset$ , we have  $S^* \subseteq S$  and thus  $\deg_S(v) \geq \deg_{S^*}(v)$ . (Because  $S$  is superset of  $S^*$ , the vertex  $v$  can only “touch” more edges in the graph spanned by  $S$  than in the graph spanned by  $S^*$ , thus  $\deg_S(v) \geq \deg_{S^*}(v)$ .) We now trivially have

$$\rho^*(G) \leq \deg_{S^*}(v) \leq \deg_S(v) \leq 2(1 + \epsilon)\rho(S),$$

establishing the desired inequality.

- (iii) We now conclude  $\rho(\tilde{S}) \geq \frac{1}{2(1+\epsilon)}\rho^*(G)$ . This must be the case, because  $\tilde{S}$  is the set  $S$  that maximizes  $\rho(S)$  and we know from previous “bullet” that there exists (we compute it in one of the iterations of the algorithm) at least one  $S$  such that  $2(1 + \epsilon)\rho(S) \geq \rho^*(G)$ , because  $S^* \subseteq V$  and we start with  $S = V$  and keep removing vertices until we are left with  $S = \emptyset$ .

- (c) See listing 4 for the implementation of the algorithm.

- (i) See figure 4 for iterations needed (blue) and the theoretical bounds (green). We see that in practice the algorithm performs — in terms of number of iterations as

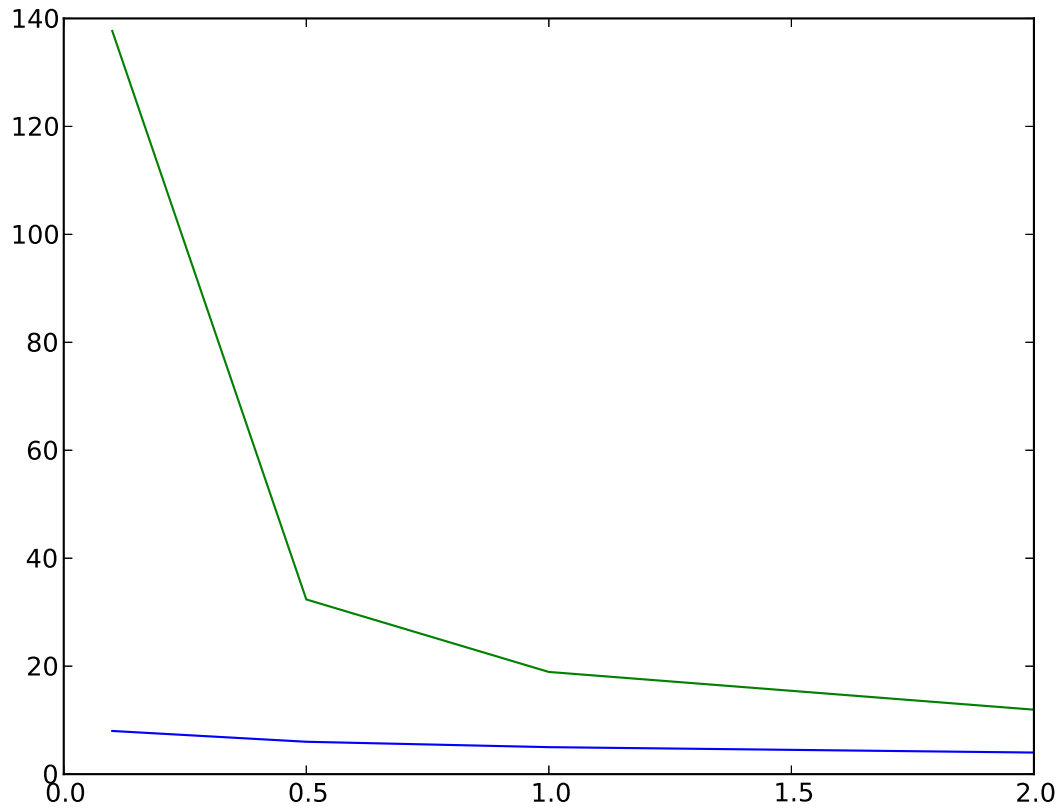


Figure 4: Number of steps (blue) and theoretical bounds (green) for various epsilons (epsilons are on  $x$ -axis).

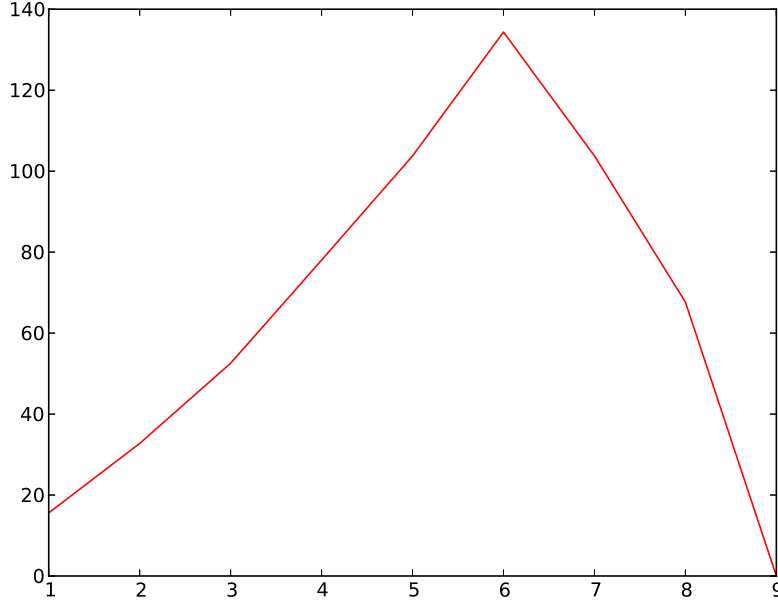


Figure 5: Density  $\rho(S_i)$  as a function of the number of steps.

a function of  $\epsilon$  — much better than what the theoretical upper bound guarantees us. (Note: We plotted  $1 + \log_{1+\epsilon}(n)$ , not  $\log_{1+\epsilon}(n)$ .)

- (ii) See figure 5 for density, figure 6 for the number of edges, and figure 7 for the number of vertices.
- (iii) See figure 8 for density, figure 9 for the number of edges, and figure 10 for the number of vertices.

## A Source code

This appendix includes all the source code that we reference throughout the document.

### A.1 Latent Features for Recommendations

Listing 1: SGD code for HW3Q1 item b.

```

1 //
2 // Compile with 'g++ -std=c++0x -O3 sgd.cpp -o sgd '
3 // — Blaz Sovdat, 18 feb 2014
4 //
5
6 #include <algorithm>
```

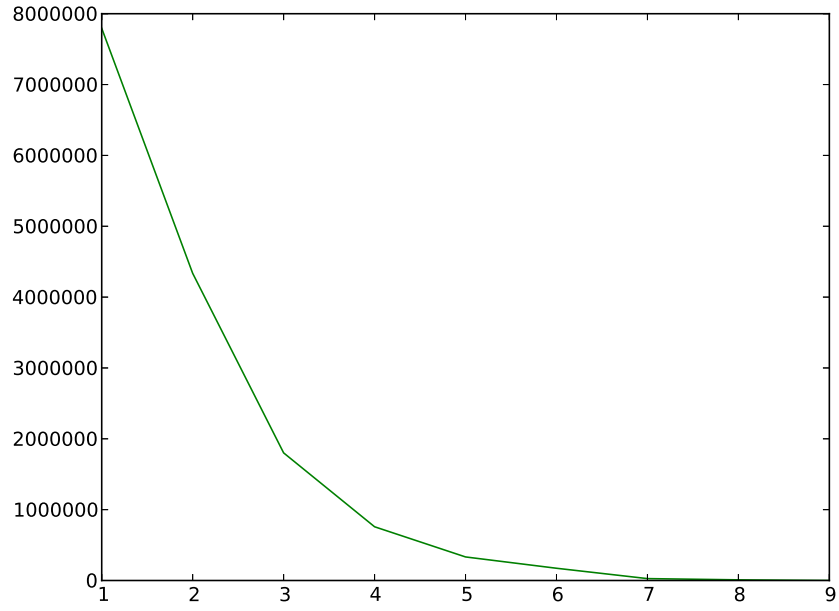


Figure 6: Number of edges  $|E[S_i]|$  as a function of the number of steps.

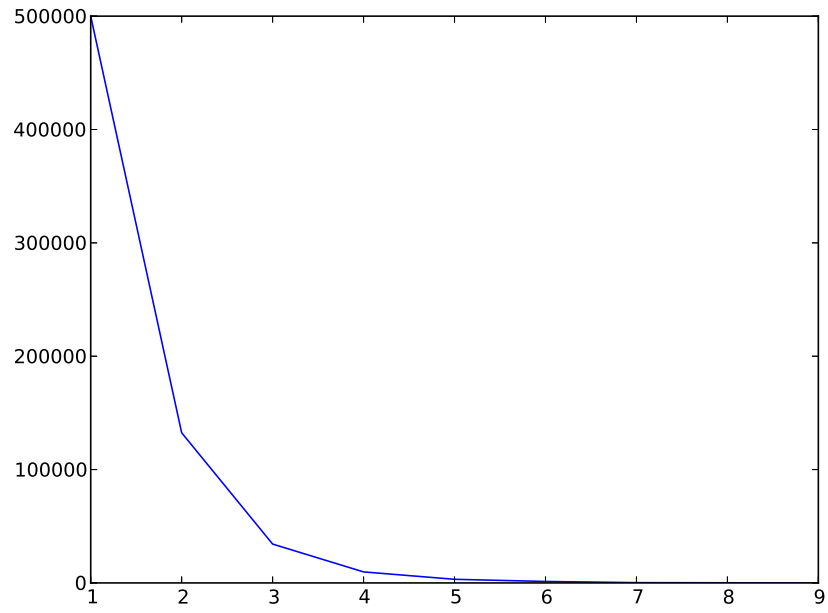


Figure 7: Number of vertices  $|S_i|$  as a function of the number of steps.

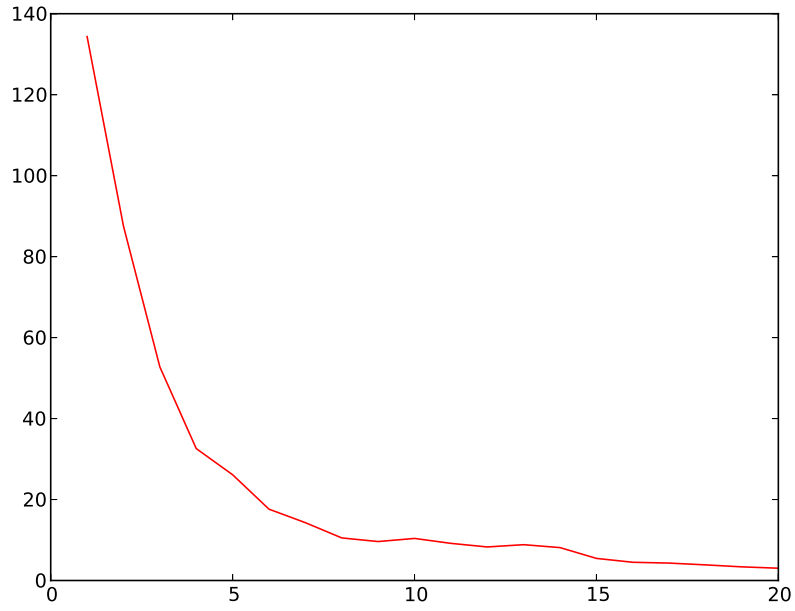


Figure 8: Densities  $\rho(S_i)$  for each of 20 communities.

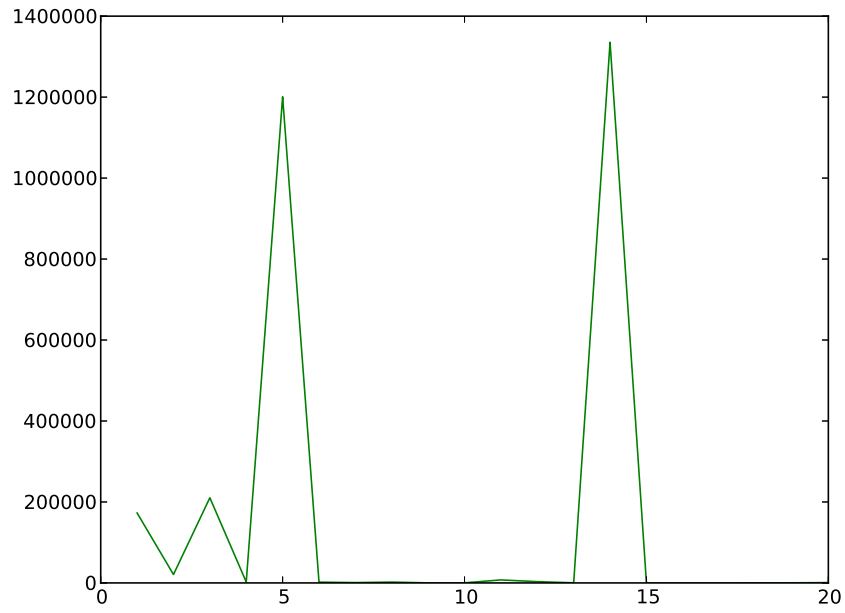


Figure 9: Edge sizes  $|E[S_i]|$  for each of 20 communities.

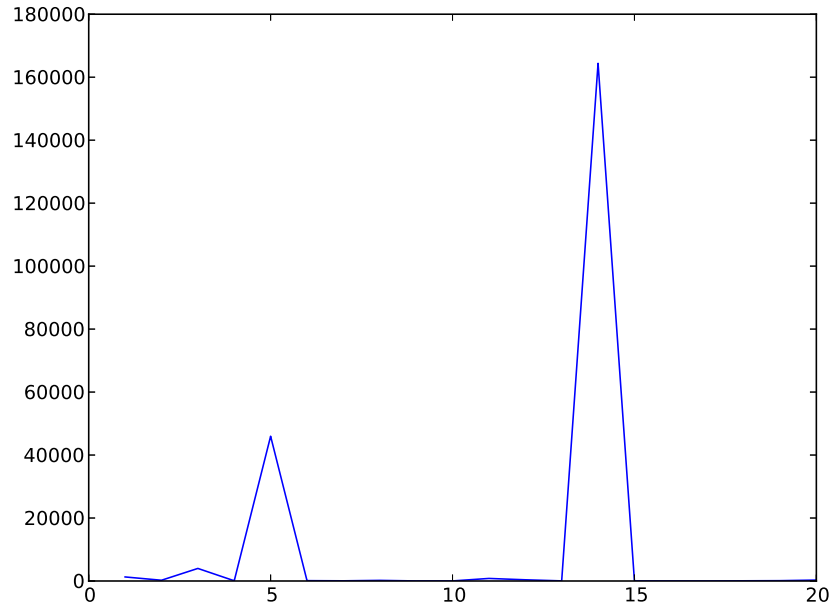


Figure 10: Number of vertices  $|S_i|$  for each of 20 communities.

```

7 #include <iostream>
8 #include <fstream>
9 #include <string>
10 #include <vector>
11 #include <cmath>
12 #include <map>
13 #include <random>
14 #include <ctime>
15
16 class Params {
17 public:
18     Params(const int& k_ = 20) : k(k_) { }
19
20     void init(const int& m, const int& n); // pass the number of movies
21     double dot(const int& qi, const int& pu) const;
22     double all_norms() const {
23         double norms = 0.0;
24         for (auto it = p_idx.begin(); it != p_idx.end(); ++it) { norms += norm_p
25             (it->second); }
26         for (auto it = q_idx.begin(); it != q_idx.end(); ++it) { norms += norm_q
27             (it->second); }
28         return norms;
29     }
30 }
31
32 double norm_q(int qi) const { // square of L2 norm
33     double norm = 0.0;

```

```

30     for (int i = 0; i < k; ++i) { norm += Q[qi+i]*Q[qi+i]; }
31     return norm;
32 }
33 double norm_p(int pi) const { // square of L2 norm
34     double norm = 0.0;
35     for (int i = 0; i < k; ++i) { norm += P[pi+i]*P[pi+i]; }
36     return norm;
37 }
38 void update(const double* q, const int& qi, const double* p, const int& pi);
39 public:
40     const int k;
41     double* P;
42     double* Q;
43     std::map<int, int> p_idx;
44     std::map<int, int> q_idx;
45 };
46
47 void Params::init(const int& m_, const int& n_) {
48     std::default_random_engine generator((int)time(0));
49     std::uniform_real_distribution<double> dist(0.0, sqrt(5.0/k));
50     const int m = m_+1, n = n_+1;
51     Q = new double[m*k];
52     for (int i = 0; i < m*k; ++i) { Q[i] = dist(generator); }
53     P = new double[n*k];
54     for (int i = 0; i < n*k; ++i) { P[i] = dist(generator); }
55 }
56
57 double Params::dot(const int& qi, const int& pu) const {
58     double s = 0.0;
59     for (int j = 0; j < k; ++j) { s += Q[qi+j]*P[pu+j]; }
60     return s;
61 }
62
63 void Params::update(const double* q, const int& qi, const double* p, const int
    & pi) {
64     for (int j = 0; j < k; ++j) { Q[qi+j] = q[j]; P[pi+j] = p[j]; }
65 }
66
67 void sgd(const std::string& fname, Params& params, const int& k, const double&
    lambda, const int& iterations, const double& eta);
68 double estimate_error(const std::string& fname, Params& params, const int& k,
    const double& lambda);
69
70 // entry point
71 int main(int argc, char** argv) {
72     Params params(20);
73     sgd("ratings.train.txt", params, 20, 0.2, 40, 0.01565);
74     return 0;
75 }
76

```

```

77 void sgd(const std::string& fname, Params& params, const int& k = 20, const
    double& lambda = 0.2, const int& iterations = 40, const double& eta = 0.1)
    {
78     double crr_eps = 0.0, tmp_eps = 0.0;
79     int p_idx = 0, q_idx = 0;
80     int crr_user = 0, crr_movie = 0, crr_rating = 0;
81     // preliminary pass to initialize params
82     std::string line;
83     std::ifstream infile(fname.c_str());
84     getline(infile, line);
85     do {
86         int t = line.find_first_of('\t', 0);
87         crr_user = atoi(line.substr(0, t).c_str());
88         crr_movie = atoi(line.substr(t+1, line.find_last_of('\t')).c_str());
89         if (params.p_idx.find(crr_user) == params.p_idx.end()) { params.p_idx[
            crr_user] = p_idx++; }
90         if (params.q_idx.find(crr_movie) == params.q_idx.end()) { params.q_idx[
            crr_movie] = q_idx++; }
91         getline(infile, line);
92     } while (!infile.eof());
93     infile.close();
94
95     params.init(q_idx-1, p_idx-1);
96
97     double* p = new double[k]();
98     double* q = new double[k]();
99
100    // do k iterations of stochastic gradient descent
101    for (int i = 0; i < iterations; ++i) {
102        // for each r_{iu} in the data
103        std::ifstream infile(fname.c_str());
104        getline(infile, line);
105        do {
106            int t = line.find_first_of('\t', 0);
107            crr_user = atoi(line.substr(0, t).c_str());
108            crr_movie = atoi(line.substr(t+1, line.find_last_of('\t')).c_str());
109            crr_rating = atoi(line.substr(line.find_last_of('\t')+1).c_str());
110            // compute epsilon
111            p_idx = k*params.p_idx[crr_user];
112            q_idx = k*params.q_idx[crr_movie];
113            // update epsilon for next iteration
114            crr_eps = crr_rating - params.dot(q_idx, p_idx);
115            // now update the vectors q and p
116            for (int j = 0; j < k; ++j) {
117                q[j] = params.Q[q_idx+j] + eta*(crr_eps*params.P[p_idx+j] - lambda*
                    params.Q[q_idx+j]);
118                p[j] = params.P[p_idx+j] + eta*(crr_eps*params.Q[q_idx+j] - lambda*
                    params.P[p_idx+j]);
119                if (std::isnan(p[j]) || std::isnan(q[j])) { std::cout << "ISNAN:_" <<
                    i << "-" << j << std::endl; }
120            }

```



```

121     params.update(q, q_idx, p, p_idx);
122
123     getline(infile, line);
124 } while (!infile.eof() && line.length() > 2);
125 infile.close();
126 std::cout << i << ":" << estimate_error(fname, params, k, lambda) << std::
    endl;
127 }
128 std::cout << "Using_eta=" << eta << std::endl;
129 }
130
131 double estimate_error(const std::string& fname, Params& params, const int& k,
    const double& lambda) {
132     double err = 0.0, crr_err = 0.0;
133     std::ifstream infile(fname.c_str());
134     std::string line;
135     int crr_movie = 0, crr_user = 0, crr_rating = 0;
136     int p_idx = 0, q_idx = 0;
137     getline(infile, line);
138     do {
139         int t = line.find_first_of('\t', 0);
140         crr_user = atoi(line.substr(0, t).c_str());
141         crr_movie = atoi(line.substr(t+1, line.find_last_of('\t')).c_str());
142         crr_rating = atoi(line.substr(line.find_last_of('\t')+1).c_str());
143         p_idx = k*params.p_idxs[crr_user];
144         q_idx = k*params.q_idxs[crr_movie];
145         crr_err = crr_rating - params.dot(q_idx, p_idx);
146         err += crr_err*crr_err;
147         getline(infile, line);
148     } while (!infile.eof() && line.length() > 2);
149     infile.close();
150     err += lambda*(params.all_norms());
151     return err;
152 }

```

Listing 2: SGD code for HW3Q1 item c.

```

1 //
2 // Compile with 'g++ -std=c++0x -O3 sgd.cpp -o sgd'
3 // — Blaz Sovdat, 18 feb 2014
4 //
5
6 #include <algorithm>
7 #include <iostream>
8 #include <fstream>
9 #include <string>
10 #include <vector>
11 #include <cmath>
12 #include <map>
13 #include <random>
14 #include <ctime>
15

```

```

16 class Params {
17 public:
18     Params(const int& k_ = 20) : k(k_) { }
19
20     void init(const int& m, const int& n); // pass the number of movies
21     double dot(const int& qi, const int& pu) const;
22     double all_norms() const {
23         double norms = 0.0;
24         for (auto it = p_idx.begin(); it != p_idx.end(); ++it) { norms += norm_p
25             (k*it->second); }
26         for (auto it = q_idx.begin(); it != q_idx.end(); ++it) { norms += norm_q
27             (k*it->second); }
28         return norms;
29     }
30     double norm_q(const int& qi) const { // square of L2 norm
31         double norm = 0.0;
32         for (int i = 0; i < k; ++i) { norm += Q[qi+i]*Q[qi+i]; }
33         return norm;
34     }
35     double norm_p(const int& pi) const { // square of L2 norm
36         double norm = 0.0;
37         for (int i = 0; i < k; ++i) { norm += P[pi+i]*P[pi+i]; }
38         return norm;
39     }
40     void update(const double *q, const int& qi, const double* p, const int& pi);
41 public:
42     const int k;
43     double* P;
44     double* Q;
45     std::map<int, int> p_idx;
46     std::map<int, int> q_idx;
47 };
48
49 void Params::init(const int& m_, const int& n_) {
50     std::default_random_engine generator((int)time(0));
51     std::uniform_real_distribution<double> dist(0.0, sqrt(5.0/k));
52     const int m = m_+1, n = n_+1;
53     Q = new double[m*k];
54     for (int i = 0; i < m*k; ++i) { Q[i] = dist(generator); }
55     P = new double[n*k];
56     for (int i = 0; i < n*k; ++i) { P[i] = dist(generator); }
57 }
58
59 double Params::dot(const int& qi, const int& pu) const {
60     double s = 0.0;
61     for (int j = 0; j < k; ++j) { s += Q[qi+j]*P[pu+j]; }
62     return s;
63 }
64
65 void Params::update(const double* q, const int& qi, const double* p, const int
    & pi) {

```

```

64   for (int j = 0; j < k; ++j) { Q[qi+j] = q[j]; P[pi+j] = p[j]; }
65 }
66
67 void sgd(const std::string& fname, Params& params, const int& k, const double&
        lambda, const int& iterations, const double& eta);
68 double estimate_error(const std::string& fname, Params& params, const int& k,
        const double& lambda);
69
70 // entry point
71 int main(int argc, char** argv) {
72     for (int k = 1; k <= 10; ++k) {
73         Params params(k);
74         // sgd("ratings.train.txt", params, k, 0.0, 40, 0.03);
75         sgd("ratings.train.txt", params, k, 0.2, 40, 0.03);
76         // std::cout << k << ":" << estimate_error("ratings.val.txt", params, k,
77             0.0) << std::endl;
78         std::cout << k << ":" << estimate_error("ratings.val.txt", params, k, 0.2)
79             << std::endl;
80     }
81     return 0;
82 }
83
84 void sgd(const std::string& fname, Params& params, const int& k = 20, const
        double& lambda = 0.2, const int& iterations = 40, const double& eta = 0.1)
85 {
86     double crr_eps = 0.0, tmp_eps = 0.0;
87     int p_idx = 0, q_idx = 0;
88     int crr_user = 0, crr_movie = 0, crr_rating = 0;
89     // preliminary pass to initialize params
90     std::string line;
91     std::ifstream infile(fname.c_str());
92     getline(infile, line);
93     do {
94         int t = line.find_first_of('\t', 0);
95         crr_user = atoi(line.substr(0, t).c_str());
96         crr_movie = atoi(line.substr(t+1, line.find_last_of('\t')).c_str());
97         // set indices
98         if (params.p_idx.find(crr_user) == params.p_idx.end()) { params.p_idx[crr_user] = p_idx++; }
99         if (params.q_idx.find(crr_movie) == params.q_idx.end()) { params.q_idx[crr_movie] = q_idx++; }
100        getline(infile, line);
101    } while (!infile.eof());
102    infile.close();
103    // initialize matrices P and Q
104    params.init(q_idx-1, p_idx-1);
105    // temporary k-dimensional vectors
106    double* p = new double[k]();
107    double* q = new double[k]();
108
109    // do k iterations of stochastic gradient descent

```

```

107 for (int i = 0; i < iterations; ++i) {
108     // for each  $r_{iu}$  in the data
109     std::ifstream infile(fname.c_str());
110     getline(infile, line);
111     do {
112         int t = line.find_first_of('\t', 0);
113         // retrieve indices
114         crr_user = atoi(line.substr(0, t).c_str());
115         crr_movie = atoi(line.substr(t+1, line.find_last_of('\t')).c_str());
116         crr_rating = atoi(line.substr(line.find_last_of('\t')+1).c_str());
117         // compute epsilon
118         p_idx = k*params.p_idx[crr_user];
119         q_idx = k*params.q_idx[crr_movie];
120         // update epsilon for next iteration
121         crr_eps = crr_rating - params.dot(q_idx, p_idx);
122         // now update the vectors q and p
123         for (int j = 0; j < k; ++j) {
124             q[j] = params.Q[q_idx+j] + eta*(crr_eps*params.P[p_idx+j] - lambda*
125                 params.Q[q_idx+j]);
126             p[j] = params.P[p_idx+j] + eta*(crr_eps*params.Q[q_idx+j] - lambda*
127                 params.P[p_idx+j]);
128             if (std::isnan(p[j]) || std::isnan(q[j])) { std::cout << "ISNAN:" <<
129                 i << "-" << j << std::endl; }
130         }
131         // assign them to parameters matrices
132         params.update(q, q_idx, p, p_idx);
133     } while (!infile.eof() && line.length() > 2);
134     infile.close();
135     std::cout << "Error_for_k=" << k << ":" << estimate_error(fname, params, k,
136         lambda) << std::endl;
137 }
138 double estimate_error(const std::string& fname, Params& params, const int& k,
139     const double& lambda) {
140     double err = 0.0, crr_err = 0.0;
141     std::ifstream infile(fname.c_str());
142     std::string line;
143     int crr_movie = 0, crr_user = 0, crr_rating = 0;
144     int p_idx = 0, q_idx = 0;
145     getline(infile, line);
146     do {
147         int t = line.find_first_of('\t', 0);
148         crr_user = atoi(line.substr(0, t).c_str());
149         crr_movie = atoi(line.substr(t+1, line.find_last_of('\t')).c_str());
150         crr_rating = atoi(line.substr(line.find_last_of('\t')+1).c_str());
151         if (params.p_idx.find(crr_user) != params.p_idx.end() && params.q_idx.
152             find(crr_movie) != params.q_idx.end()) {
153             p_idx = k*params.p_idx[crr_user];

```

```

152     q_idx = k*params.q_idxes[crr_movie];
153     crr_err = crr_rating - params.dot(q_idx, p_idx);
154     err += crr_err*crr_err;
155 }
156 getline(infile, line);
157 } while (!infile.eof() && line.length() > 2);
158 infile.close();
159 // err += lambda*(params.all_norms()); // error function was redefined for
    this bullet
160 return err;
161 }

```

## A.2 PageRank Computation

Listing 3: PageRank code for HW3Q2.

```

1  #!/usr/bin/python
2
3  import scipy as sp
4  from scipy.sparse import *
5  import numpy as np
6  from scipy import *
7  import random as rnd
8  from time import time
9
10 avg = lambda L: 1.0*sum(L)/len(L)
11
12 def readfile(fname):
13     for edge in open(fname):
14         yield [int(v) for v in edge.split('\t')]
15
16 def error(P, R, n):
17     IP = sorted(range(len(P)), key = lambda k: P[k])
18     IR = sorted(range(len(R)), key = lambda k: R[k])
19     IP.reverse()
20     IR.reverse()
21     return 1.0*sum([abs(P[IP[i]] - R[IR[i]]) for i in range(n)])/n
22
23 def mc(fname, R, beta = 0.8):
24     G = {}
25     for T in readfile(fname):
26         if T[0] not in G: G[T[0]] = [T[1]]
27         else: G[T[0]].append(T[1])
28
29     n = len(G)
30     # frequency count
31     f = [0 for k in range(n)]
32     # simulate random walk
33     # for each vertex
34     for v in G:

```

```

35     # simulate it R times
36     for iter in range(R):
37         # walk terminates at each node with prob. 1-beta
38         node = v
39         while True:
40             f[node-1] = f[node-1]+1
41             if rnd.uniform(0, 1) > 1-beta: # continue the walk
42                 node = G[node][rnd.randint(0, len(G[node])-1)]
43             else: # game over
44                 break
45         # print f
46         r = [(1-beta)*f[i]/(n*R) for i in range(n)]
47     return r
48
49 def pagerank(fname, beta = 0.8):
50     row = []
51     col = []
52     for T in readfile(fname):
53         row.append(T[0]-1)
54         col.append(T[1]-1)
55     n = 100
56     row = array(row)
57     col = array(col)
58     dat = array([1.0 for k in row])
59     # convert it to full matrix to simplify our life
60     M = csc_matrix((dat, (row, col)), shape = (n, n)).todense()
61     # make it column stochastic
62     for i in range(n): M[:, i] = M[:, i].dot(1.0/M[:, i].sum())
63     # initialize page rank vector
64     r = np.transpose(array([1.0/n for k in range(n)]))
65     t = np.transpose(array([(1-beta)/n for k in range(n)]))
66     # now do power iteration
67     for i in range(40):
68         r = array(t+beta*M.dot(r))
69         r = r[0,:]
70     return r
71
72 # entry point
73 if __name__ == '__main__':
74     # compute PageRank with Power iteration
75     pr = pagerank('graph.txt', 0.8)
76     #T = []
77     #for i in range(50):
78     # start = time()
79     # pr = pagerank('graph.txt', 0.8)
80     # T.append(time() - start)
81     #print avg(T)
82
83     # now run the MC algortihm
84     r1 = mc('graph.txt', 1, 0.8)
85     r3 = mc('graph.txt', 3, 0.8)

```

```

86 r5 = mc('graph.txt', 5, 0.8)
87
88 print "*" * 3, "R=_1", "*" * 3
89 print "Top_10:", error(pr, r1, 10)
90 print "Top_30:", error(pr, r1, 30)
91 print "Top_50:", error(pr, r1, 50)
92 print "All_:", error(pr, r1, len(r1))
93
94 print "*" * 3, "R=_3", "*" * 3
95 print "Top_10:", error(pr, r3, 10)
96 print "Top_30:", error(pr, r3, 30)
97 print "Top_50:", error(pr, r3, 50)
98 print "All_:", error(pr, r3, len(r3))
99
100 print "*" * 3, "R=_5", "*" * 3
101 print "Top_10:", error(pr, r5, 10)
102 print "Top_30:", error(pr, r5, 30)
103 print "Top_50:", error(pr, r5, 50)
104 print "All_:", error(pr, r5, len(r5))

```

### A.3 Dense Communities in Networks

Listing 4: Streaming implementatoin of the algorithm from HW3Q4.

```

1  #!/usr/bin/python
2
3  from time import time
4  import matplotlib.pyplot as plt
5  import math
6
7  error = lambda P, R, I, n: 1.0*sum([abs(P[I[i]] - R[I[i]]) for i in range(n)])
      /n
8  avg = lambda L: 1.0*sum(L)/len(L)
9
10 # we count each edge twice, 2|E(G)| = sum_v deg(v), so we need to factor out 2
11 density = lambda S: avg(S.values())/2 if len(S) > 0 else 0
12
13 def readfile(fname):
14     for edge in open(fname):
15         yield [int(v) for v in edge.split()]
16
17 # basic version of the algorithm
18 def algo(fname, eps):
19     S = {}
20     T = {}
21     for e in readfile(fname):
22         # print e
23         S[e[0]] = 1 if e[0] not in S else S[e[0]]+1
24         S[e[1]] = 1 if e[1] not in S else S[e[1]]+1
25     dS = density(S)

```

```

26 print dS
27 i = 1
28 plot_rhos = [dS]
29 plot_edgs = [sum(S.values())/2]
30 plot_sizs = [len(S)]
31 plot_X = [i]
32 while len(S) > 0:
33     tS = {}
34     for v in S:
35         if S[v] > 2*(1+eps)*dS: tS[v] = 0
36     for e in readfile(fname): # compute S \ A(S)
37         if e[0] in tS and e[1] in tS:
38             tS[e[0]] = tS[e[0]]+1
39             tS[e[1]] = tS[e[1]]+1
40     S = dict(tS)
41     dS = density(tS)
42     if dS > density(T): T = dict(tS)
43     i = i+1
44     plot_rhos.append(dS)
45     plot_edgs.append(sum(S.values())/2 if len(S) > 0 else 0)
46     plot_sizs.append(len(S))
47     plot_X.append(i)
48     print "Current_density: ", dS
49     print " |S|= ", len(S), " |T|= ", len(T)
50     print "-"*80
51 # print T
52 print density(T)
53 print "Done"
54
55 # print "rhos=", plot_rhos
56 # print "edgs=", plot_edgs
57 # print "sizs=", plot_sizs
58
59 plt.plot(plot_X, plot_rhos, 'r')
60 plt.draw()
61
62 plt.figure()
63 plt.plot(plot_X, plot_edgs, 'g')
64 plt.draw()
65
66 plt.figure()
67 plt.plot(plot_X, plot_sizs, 'b')
68 plt.draw()
69
70 plt.show()
71
72 return i
73
74 # stops after it finds 20 communities
75 # NOTE: Due to time constraints the code is not efficient
76 def multiple_components(fname, eps = 0.05):

```



```

77 D = set()
78 plot_rhos = []
79 plot_sizs = []
80 plot_edgs = []
81 plot_X = []
82 for i in range(1, 21):
83     start = time()
84     S = {}
85     T = {}
86     for e in readfile(fname):
87         if e[0] not in D and e[1] not in D:
88             S[e[0]] = 1 if e[0] not in S else S[e[0]]+1
89             S[e[1]] = 1 if e[1] not in S else S[e[1]]+1
90     dS = density(S)
91     while len(S) > 0:
92         tS = {}
93         for v in S:
94             if S[v] > 2*(1+eps)*dS: tS[v] = 0
95         for e in readfile(fname): # compute S \ A(S)
96             if e[0] in tS and e[1] in tS:
97                 tS[e[0]] = tS[e[0]]+1
98                 tS[e[1]] = tS[e[1]]+1
99         S = dict(tS)
100        dS = density(tS)
101        if dS > density(T): T = dict(tS)
102        plot_rhos.append(density(T))
103        plot_edgs.append(sum(T.values())/2 if len(T) > 0 else 0)
104        plot_sizs.append(len(T))
105        plot_X.append(i)
106        for v in T: D.add(v)
107        print "Current_density: ", density(T)
108        print "|T|=", len(T)
109        print "|E[T]|=", sum(T.values())/2
110        print "Community_found_in ", time()-start, "seconds"
111        print "-"*80
112
113    plt.plot(plot_X, plot_rhos, 'r')
114    plt.draw()
115
116    plt.figure()
117    plt.plot(plot_X, plot_edgs, 'g')
118    plt.draw()
119
120    plt.figure()
121    plt.plot(plot_X, plot_sizs, 'b')
122    plt.draw()
123
124    plt.show()
125
126 # entry point
127 if __name__ == '__main__':

```

```

128 # algo('livejournal-undirected-small.txt ', 0.01)
129 """
130 n = 499923
131 X = [0.1, 0.5, 1.0, 2.0]
132 for eps in X:
133     start = time()
134     # algo returns the number of steps
135     Y.append(algo('livejournal-undirected.txt ', eps))
136     print "time: ", time()-start, "seconds"
137 print "X=", X
138 print "Y=", Y
139 plt.plot(X, Y)
140 plt.plot(X, [math.log(n, 1+x) for x in X])
141 plt.show()
142 """
143
144 # print algo('livejournal-undirected.txt ', 0.05)
145
146 multiple_components('livejournal-undirected.txt ', 0.05)

```

## References