

CS246: Mining Massive Datasets

Assignment number: 4_____

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 1st 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: /

1 Support Vector Machine

- (a) We give a sample training set of 5 points $\mathbf{x}_i \in \mathbb{R}^2$ with their respective classes $y_i \in \{-1, 1\}$ such that the set is infeasible under hard constraint SVM, but feasible under soft margin SVM:

- $(0, 1), -$,
- $(4, 0), -$,
- $(2, 2), +$,
- $(5, 1), +$,
- $(3, 3), -$.

This data set is clearly infeasible under hard constraint SVM as it is not linearly separable. However, it is feasible under soft margin SVM. A concrete but trivial solution, assuming $C = 0$, is $\mathbf{w} := \mathbf{0}$ and $b := \xi_1 := \dots \xi_n := 0$ (we can set b and ξ_i 's to anything). (More generally, for any \mathbf{w} , b , and C the problem is feasible under soft margin SVM, because setting $\xi_i := \max(0, 1 - y_i(\mathbf{w}\mathbf{x}_i + b))$ satisfies all inequalities.)

- (b) Let $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$ be an arbitrary dataset and suppose we are doing soft margin SVM. Note that setting i th slack to $\xi_i := \max(0, 1 - y_i(\mathbf{w}\mathbf{x}_i + b))$ for all $1 \leq i \leq n$ makes all the inequalities hold, where \mathbf{w} and b are arbitrary. Thus the problem is feasible.
- (c) Let $\{(\mathbf{x}_i, y_i) \mid 1 \leq i \leq n\}$ be our data set and let $(\mathbf{w}, b, \xi_1, \dots, \xi_n)$ be feasible parameter set. Clearly, if $y_i(\mathbf{w}\mathbf{x}_i + b) \leq 0$ for arbitrary example, i.e., i th example is misclassified, then the corresponding slack will be $\xi_i \geq 1$, so each misclassification contributes at least 1 to the sum $\xi_1 + \xi_2 + \dots + \xi_n$. This means $\sum_{i=1}^n \xi_i$ is an upper bound on the number of misclassifications.
- (d) We compute $\nabla_b f(\mathbf{w}, b)$ for the batch gradient descent:

$$\nabla_b f(\mathbf{w}, b) = \frac{\partial f}{\partial b}(\mathbf{w}, b) = C \sum_{i=1}^n \frac{\partial L}{\partial b} L(\mathbf{x}_i, y_i),$$

where L is the hinge function with the partial derivative of

$$\frac{\partial L}{\partial b}(\mathbf{x}_i, y_i) = \begin{cases} 0 & y_i(\mathbf{w}\mathbf{x}_i + b) \geq 1, \\ -y_i & \text{otherwise.} \end{cases} \quad (1)$$

For stochastic gradient descent we use

$$\nabla_b f_i(\mathbf{w}, b) = C \frac{\partial L}{\partial b}(\mathbf{x}_i, y_i).$$

For mini batch gradient descent we have

$$\nabla_b f_\ell(\mathbf{w}, b) = C \sum_{i=bs \cdot \ell + 1}^{\min(n, (\ell+1)bs)} \frac{\partial L}{\partial b}(\mathbf{x}_i, y_i).$$

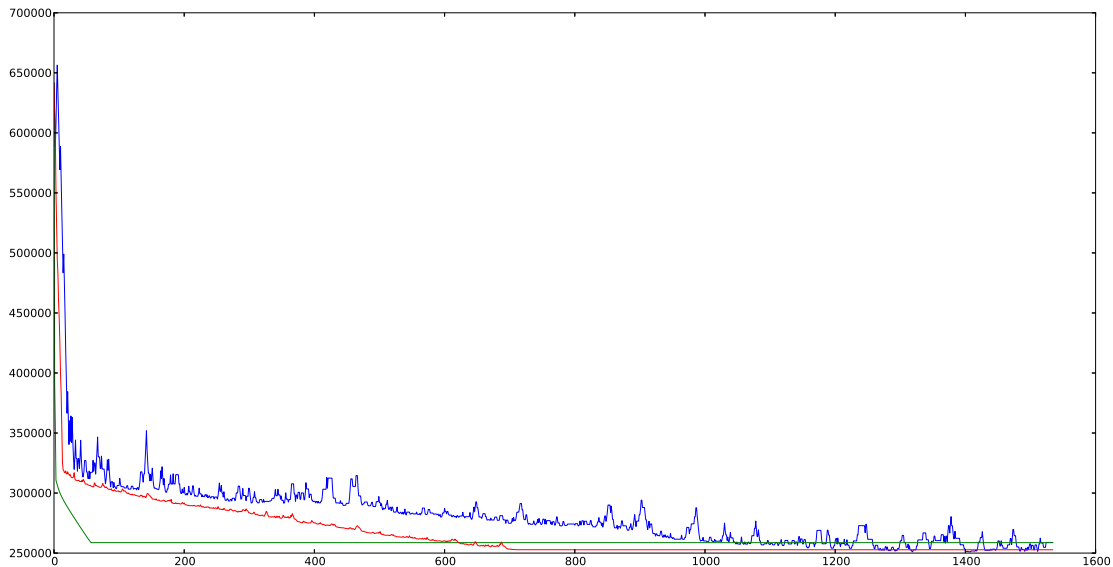


Figure 1: Value of $f_k(\mathbf{w}, b)$ at each iteration k for batch GD (green), stochastic GD (blue), and mini batch GD (red), where $k = 0, 1, \dots, 1534$. See text for detailed explanation.

- (e) See figure 1 for $f_k(\mathbf{w}, b)$ plots. (Parameters of descents are as in HW4Q1 text.) Batch GD converged (green line) in 57th iteration, mini batch GD (red line) in 713th, while stochastic GD converged (blue line) in 1534th iteration. In figure 1 we plotted $f_{57}(\mathbf{w}, n)$ instead of $f_i(\mathbf{w}, b)$ for all $i > 57$ for batch GD; we did similar thing for mini batch GD for $i > 713$ (this is why green and red line are “straight” from some point on). Code for batch GD is in listing 1, code for stochastic GD in listing 2, and code for mini batch GD in listing 3. (Source code is in the appendix at the end of the report.) The running times are in table 1; we computed running times using `timeit` Python module.

We see that stochastic GD takes least time per iteration but takes many iterations to converge; still, its running time is much lower than batch GD’s. Batch GD takes lots of time per iteration as it needs to go through the whole dataset to evaluate the gradient; although it converges very quickly, it is still much slower than stochastic GD. Mini batch GD falls somewhere in between: it converges faster than stochastic GD and takes slightly more time per iteration; in our case it converged faster than both batch and stochastic GD. All variants converge to roughly the same objective function value.

- (f) See listing 5 for code for this task. Figure 2 shows plot of percentage error — the fraction of misclassified examples on the test set — of stochastic GD (with parameters as in HW4Q1 text) as a function of regularization parameter $C = 1, 10, 50, 100, 200, 300, 400, 500$. We see that increasing C (this means we increase punishment for misclassifications) reduces percentage error on the test set. (If we increased C too much, we would overfit:

Method	Time
Batch GD	2125.11 seconds (≈ 35.4 minutes)
Stochastic GD	508.17 seconds (≈ 8.75 minutes)
Mini batch GD	326.584 seconds (≈ 5.4 minutes)

Table 1: Running times for batch, stochastic, and mini batch GD implementations of SVM.

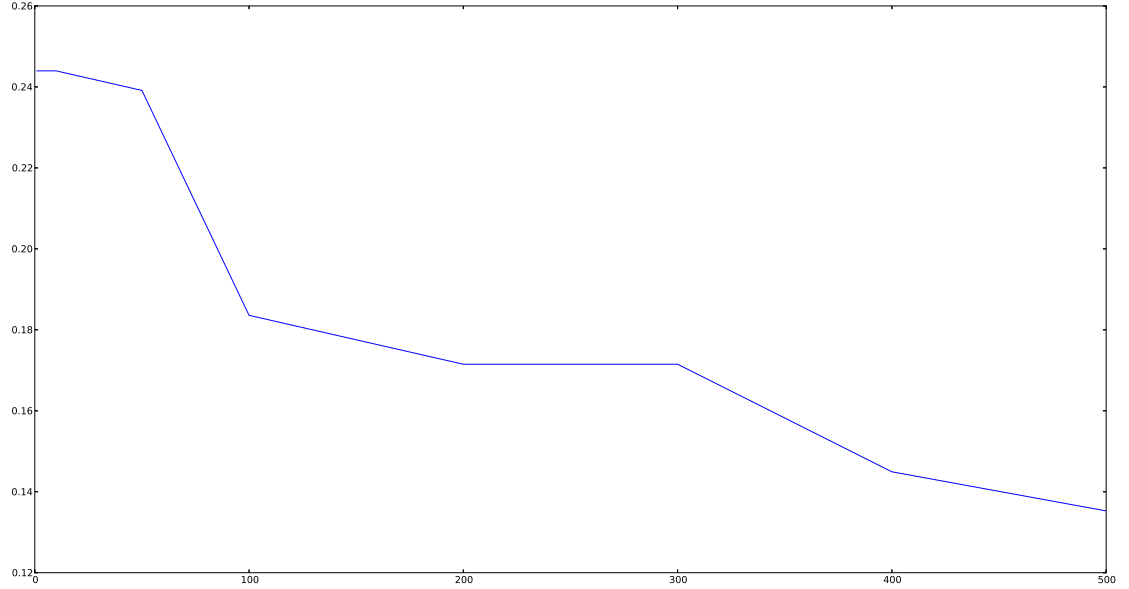


Figure 2: Plot of the percentage error of stochastic GD (parameters are in HW4Q1 text) as a function of regularization parameter C for $C = 1, 10, 50, 100, 200, 300, 400, 500$.

for very large C we practically ignore $\|\mathbf{w}\|^2/2$ factor. Similarly, sending $C \rightarrow 0$ ignores the data.) Percentage error decreases, then is “stable” around $C = 200, 300$ and then decreases some more.

2 Decision Tree Learning

(a) We compute $I(D)$ for each of the binary attributes. First note that we have $I(D) = 100(1 - (60/100)^2 - (40/100)^2) = 48$.

- For “likes wine” attribute $|D_L| = |D_R| = 50$ and $I(D_L) = I(D_R) = 50(1 - (30/50)^2 - (20/50)^2) = 24$, which gives us $I(D) - I(D_L) - I(D_R) = 0$.
- For “likes running” attribute we have $I(D_L) = 30(1 - (20/30)^2 - (10/30)^2) = 13.3$ and $I(D_R) = 70(1 - (40/70)^2 - (30/70)^2) = 34.29$, giving us $I(D) - I(D_L) - I(D_R) =$

0.38.

- For “likes pizza” attribute we have $I(D_L) = 80(1 - (50/80)^2 - (30/80)^2) = 37.5$ and $I(D_R) = 20(1 - (10/20)^2 - (10/20)^2) = 10$, giving us $I(D) - I(D_L) - I(D_R) = 0.5$.

We would use “likes pizza” attribute to split the root, because it has highest value of the Gini index based metric G . Roughly speaking, this means it alone (in the sense that we do not take into account attribute interactions; we just measure how well each single attribute classifies the data set) best classifies the data set.

- (b) Under reasonable attribute estimation measures — like information gain and Gini index — the learner will identify a_1 as the most important attribute and use it as the split in the root node; all other attributes will appear in “lower layers” of the tree; the complete binary decision tree overfits the data.

It is obvious that the desired tree is the one with a single split on a_1 , with leaf in the left branch ($a_1 = 1$) predicting 1 and leaf in the right branch ($a_1 = 0$) predicting 0. The reason for this is that simple models tend to generalize well; they are also easier to understand. (Very roughly, Occam’s razor tells us we should favor simple hypotheses over complicated ones.)

- (c) We use the following criterion for pruning:

$$\alpha = \frac{\text{error}(\text{pruned}(T, t), S) - \text{error}(T, S)}{|\text{nodes}(T)| - |\text{nodes}(\text{pruned}(T, t))|},$$

where $\text{error}(T, S)$ is error rate of T over the sample S . See figure 3 for original tree; figure 4 for T_1 ; figure 5 for T_2 ; figure 6 for T_3 ; and figure 7 for the final tree. Error rates (computed on the data from HW4Q2 figure 1) are included in the figure captions.

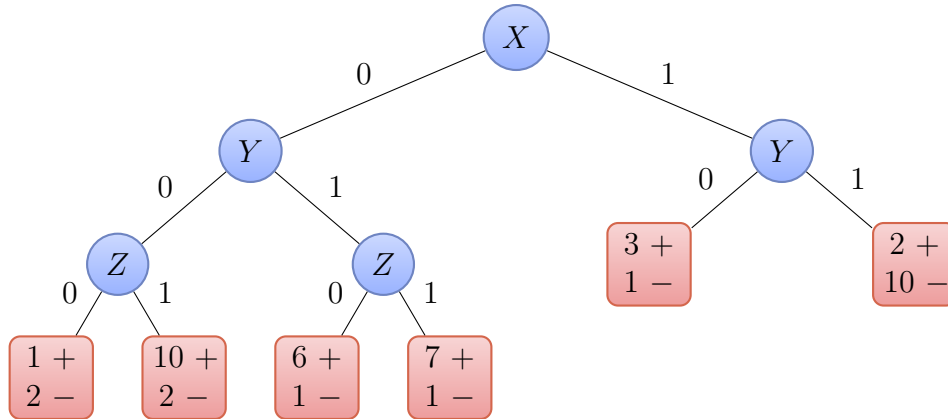


Figure 3: Decision tree T_0 with $err = (1 + 2 + 1 + 1 + 1 + 2)/46$.

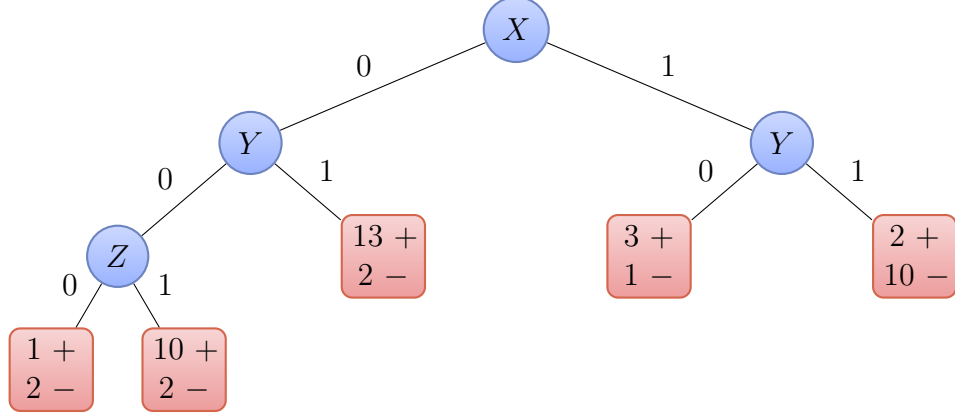


Figure 4: Decision tree T_1 with $err = (1 + 2 + 2 + 1 + 2)/46$.

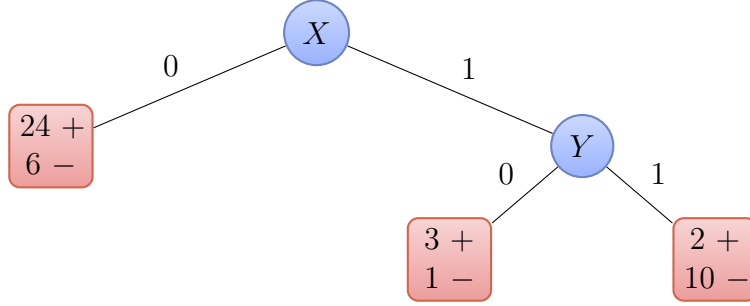


Figure 5: Decision tree T_2 with $err = (6 + 1 + 2)/46$.

- (d) We now compute generalization errors for trees from figures 4–7 on the test data. We found that T_1 and T_3 have generalization error of $1/2$; tree T_2 has smaller generalization error of $1/4$; tree T_3 has zero generalization error; trivial tree T_4 has $1/2$ generalization error. Thus tree T_3 from figure 6 has the smallest generalization error.

3 Clustering Data Streams

We first remember the notation. Let $d : \mathbb{R}^p \times \mathbb{R}^p \rightarrow \mathbb{R}^+$ be given by $d(x, y) := \|x - y\|_2$ and for $x \in \mathbb{R}^p$ and $T \subset \mathbb{R}^p$ we let

$$d(x, T) := \min_{z \in T} d(x, z).$$

For subsets $S, T \subset \mathbb{R}^p$ and a weight function $w : S \rightarrow \mathbb{R}^+$, define

$$\text{cost}_w(S, T) := \sum_{x \in S} w(x) d(x, T)^2.$$

Also, let $T^* := \arg \min_{|T|=k} \text{cost}_w(S, T)$. We now turn to problems.

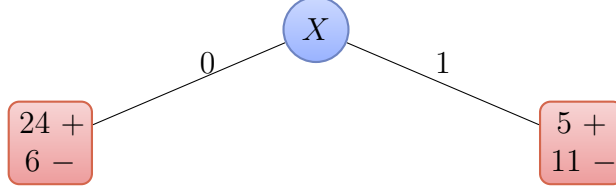


Figure 6: Decision tree T_3 with $err = (6 + 5)/46$.



Figure 7: Decision tree T_4 with $err = 17/46$.

(a) Let $(a, b) \in \mathbb{R}_+$. Clearly we have

$$2a^2 + 2b^2 - (a + b)^2 = a^2 + b^2 - 2ab = (a - b)^2 \geq 0,$$

giving us $(a + b)^2 \leq 2a^2 + 2b^2$.

(b) We now show

$$\text{cost}(S, T) \leq 2 \text{cost}_w(\hat{S}, T) + 2 \sum_{i=1}^{\ell} \text{cost}(S_i, T_i).$$

First note that $\sum_{t_{ij} \in \hat{S}} |S_{ij}| = |S|$ and that each $x \in S$ belongs to exactly one S_i and exactly one S_{ij} . Pick arbitrary $x \in S$ and suppose $x \in S_i$ with $z := \arg \min_{z \in T} d(x, z)$. Now apply triangle inequality to get $d(x, z) \leq d(x, t_{ix}) + d(t_{ix}, z)$ for $t_{ix} := \arg \min_{t_{ix} \in T_i} d(x, t_{ix})$.

Taking squares we get

$$\begin{aligned} d(x, z)^2 &\leq (d(x, t_{ix}) + d(t_{ix}, z))^2 \\ &\leq 2d(x, t_{ix})^2 + 2d(t_{ix}, z)^2, \end{aligned}$$

with the last inequality following by previous bullet. If $x \in S_i$ “belongs to” center $t_{ix} \in T_i$, then $\arg \min_{z \in T} d(x, z) = \arg \min_{z \in T} d(t_{ix}, z)$, so the proof is practically finished.

We now sum over all elements and get what we want:

$$\begin{aligned} \sum_{x \in S} \text{cost}(x, T) &= \sum_{i=1}^{\ell} \sum_{x \in S_i} \min_{z \in T} d(x, z)^2 \\ &\leq \sum_{i=1}^{\ell} \sum_{x \in S_i} \left(2 \min_{t_{ix} \in T_i} d(x, t_{ix})^2 + 2 \min_{z \in T} d(t_{ix}, z)^2 \right) \\ &= 2 \text{cost}_w(\hat{S}, T) + 2 \sum_{i=1}^{\ell} \text{cost}(S_i, T_i). \end{aligned}$$

- (c) We now prove $\sum_{i=1}^{\ell} \text{cost}(S_i, T_i) \leq \alpha \text{cost}(S, T^*)$. Note that $S_1, S_2, \dots, S_{\ell}$ form partition of S and by ALG's design $\text{cost}(S_i, T_i) \leq \alpha \text{cost}(S_i, T')$ for all $|T'| = k$. We thus have

$$\begin{aligned}
\sum_{i=1}^{\ell} \text{cost}(S_i, T_i) &\leq \alpha \sum_{i=1}^{\ell} \text{cost}(S_i, T^*) \\
&= \alpha \sum_{i=1}^{\ell} \sum_{x \in S_i} d(x, T^*)^2 \\
&= \alpha \sum_{x \in S} d(x, T^*)^2 \\
&= \alpha \text{cost}(S, T^*).
\end{aligned}$$

- (d) We now show $\text{cost}_w(\hat{S}, T) \leq \alpha \text{cost}_w(\hat{S}, T^*)$. Let $T_{\star} := \arg \min_{|T'|=k} \text{cost}_w(\hat{S}, T')$, so $\text{cost}_w(\hat{S}, T) \leq \alpha \text{cost}_w(\hat{S}, T_{\star})$. Now since $\text{cost}_w(\hat{S}, T^*) \geq \text{cost}_w(\hat{S}, T_{\star})$, we put all together and get

$$\alpha \text{cost}_w(\hat{S}, T^*) \geq \alpha \text{cost}_w(\hat{S}, T_{\star}) \geq \text{cost}_w(\hat{S}, T),$$

so we “proved” $\text{cost}_w(\hat{S}, T) \leq \alpha \text{cost}_w(\hat{S}, T^*)$.

- (e) In this bullet we prove $\text{cost}_w(\hat{S}, T^*) \leq 2 \sum_{i=1}^{\ell} \text{cost}(S_i, T_i) + 2 \text{cost}(S, T^*)$, using similar arguments as in (b). Fix $t_{ij} \in \hat{S}$. (Recall that each center t_{ij} has associated $S_{ij} \subseteq S_i$.) Then for each $x \in S_{ij}$ we have $d(t_{ij}, z) \leq d(t_{ij}, x) + d(x, z)$ for $z = \arg \min_{z \in T^*} d(t_{ij}, z) = d(x, T^*)$. Using (a) we get $d(t_{ij}, z)^2 \leq 2d(t_{ij}, x)^2 + 2d(x, z)^2$. We now sum over all elements to get the desired inequality:

$$\begin{aligned}
\text{cost}_w(\hat{S}, T^*) &= \sum_{t_{ij} \in \hat{S}} |S_{ij}| d(t_{ij}, T^*)^2 \\
&\leq \sum_{t_{ij} \in \hat{S}} \sum_{x \in S_{ij}} (2d(t_{ij}, x)^2 + 2d(x, T^*)^2) \\
&= 2 \sum_{i=1}^{\ell} \sum_{x \in S_i} d(t_{ij}, x)^2 + 2 \sum_{x \in S} d(x, T^*)^2 \\
&= 2 \sum_{i=1}^{\ell} \text{cost}(S_i, T_i) + 2 \text{cost}(S, T^*)
\end{aligned}$$

(Above we are just summing over different partitions of the same set and rearranging the summation terms, so everything is fine.)

(f) Using previous inequalities we have:

$$\begin{aligned}
\text{cost}(S, T) &\leq 2 \text{cost}_w(\hat{S}, T) + 2 \sum_{i=1}^{\ell} \text{cost}(S_i, T_i) \\
&\leq 2\alpha \text{cost}_w(\hat{S}, T^*) + 2\alpha \text{cost}(S, T^*) \\
&\leq 2\alpha \left(2 \sum_{i=1}^{\ell} \text{cost}(S, T_i) + 2 \text{cost}(S, T^*) \right) + 2\alpha \text{cost}(S, T^*) \\
&\leq 2\alpha (2\alpha \text{cost}(S, T^*) + 2 \text{cost}(S, T^*)) + 2\alpha \text{cost}(S, T^*) \\
&= (4\alpha^2 + 6\alpha) \text{cost}(S, T^*).
\end{aligned}$$

(g) We show that we with a good choice of partitioning, ALGSTR works with $O(\sqrt{nk})$ memory. If we pick $\ell = O(\sqrt{n/k})$ then we use $n/\ell = O(\sqrt{nk})$ memory in each of ℓ steps and $k\ell = O(\sqrt{nk})$ memory for the final clustering. We thus use $O(\sqrt{nk})$ memory in total.

4 Data Streams

Let $S = \langle a_1, a_2, \dots, a_t \rangle$ be a data stream of items from $\{1, 2, \dots, n\}$. Assume for any $1 \leq i \leq n$, $F[i]$ is the number of times item i has appeared in S . For given parameters $\epsilon > 0$ and $\delta > 0$ the algorithm picks $\lceil \log(1/\delta) \rceil$ independent hash functions $h_j : \{1, \dots, n\} \rightarrow \{1, 2, \dots, \lceil e/\epsilon \rceil\}$. Let $\tilde{F}[i] := \min_j c_{j, h_j(i)}$. (Note: Our notation here is somewhat cumbersome; sorry for inconvenience.)

(a) We claim $\tilde{F}[i] \geq F[i]$ for all $1 \leq i \leq n$. To see this note that if item i appeared $F[i]$ times, then $c_{j, h_j(i)} \geq F[i]$ for all $1 \leq j \leq n$, because each hash function h_j will increment the count $c_{j, h_j(i)}$ at each of the $F[i]$ appearances of i . (Count may be greater, depending on the number of collisions.)

(b) For any $1 \leq i \leq n$ and $1 \leq j \leq \lceil \log(1/\delta) \rceil$ we will prove

$$\mathbb{E}[c_{j, h_j(i)}] \leq F[i] + \frac{\epsilon}{e}(t - F[i]).$$

First fix item $i \in \{1, 2, \dots, n\}$ and define random variable $X_j := c_{j, h_j(i)}$. We already know $X_j \geq F[i]$; we just need to bound “surplus”¹ in our count. To do this, we define indicator random variable $Y_{j,k} := 1$ if $h_j(i) = h_k(i)$ for $k \neq i$ and $k \in \{1, 2, \dots, n\}$ and $Y_{j,k} := 0$ otherwise. (Intuitively, sum of $Y_{j,k}$ ’s counts collisions, and thus surplus in our counts.) We can now write X_j in terms of $F[i]$ and indicator random variables:

$$X_j = F[i] + \sum_{1 \leq k \neq i \leq n} Y_{j,k} F[k].$$

¹This is translation of “presežek”.

Note $\mathbb{E}[Y_{j,k}] = \Pr[Y_{j,k} = 1] = \Pr[h_j(k) = h_j(i)] \leq \epsilon/e$ because hash functions are uniform and independent. By linearity of expectation we have

$$\begin{aligned}\mathbb{E}[X_j] &= F[i] + \sum_{1 \leq k \neq i \leq n} \mathbb{E}[Y_{j,k}] F[k] \\ &\leq F[i] + \frac{\epsilon}{e} \sum_{1 \leq k \neq i \leq n} F[k] \\ &= F[i] + \frac{\epsilon}{e} (t - F[i]),\end{aligned}$$

establishing the inequality, because we fixed arbitrary item i . (We used $t = \sum_i F[i]$ and $\sum_{k \neq i} = t - F[i]$.)

(c) We now show

$$\Pr[\tilde{F}[i] \leq F[i] + \epsilon t] \geq 1 - \delta.$$

Recall that the Markov's inequality says $\Pr[X \geq a] \leq \mathbb{E}[X]/a$ for real constant $a > 0$ and nonnegative random variable X . For X_j from previous bullet we have $\Pr[X_j - F[i] \geq \epsilon t] \leq \mathbb{E}[X_j - F[i]]/(t\epsilon) \leq 1/e$. First note that $\Pr[\tilde{F}[i] \leq F[i] + \epsilon t] = \Pr[\tilde{F}[i] - F[i] \leq \epsilon t]$ and (because $\tilde{F}[i] = \min_j c_{j,h_j(i)}$) that $\tilde{F}[i] - F[i] \leq \epsilon t$ implies $c_{j,h_j(i)} - F[i] \leq \epsilon t$ for all $1 \leq j \leq n$. We will “work with” event ‘ $\tilde{F}[i] \geq F[i] + \epsilon t$ ’. Thus

$$\begin{aligned}\Pr[\tilde{F}[i] \geq F[i] + \epsilon t] &= \Pr[\tilde{F}[i] - F[i] \geq \epsilon t] \\ &= \Pr \left[\bigwedge_{j=1}^{\lceil \log(1/\delta) \rceil} (c_{j,h_j(i)} - F[i] \geq \epsilon t) \right] \\ &= \prod_{j=1}^{\lceil \log(1/\delta) \rceil} \Pr[c_{j,h_j(i)} - F[i] \geq \epsilon t] \\ &\leq (1/e)^{\lceil \log(1/\delta) \rceil} \leq (1/e)^{\log(1/\delta)} = (1/e)^{-\log \delta} = e^{\log \delta} = \delta,\end{aligned}$$

because we used natural logarithms for this question, i.e., $\log x := \log_e x$. Finally, note that $\Pr[\tilde{F}[i] \leq F[i] + \epsilon t] = 1 - \Pr[\tilde{F}[i] \geq F[i] + \epsilon t] \geq 1 - \delta$. This finishes the proof.

(d) In HW3Q4 we had a streaming implementation of algorithm for finding dense communities in networks, where we processed the input graph edge-by-edge and counted degrees of each of the nodes; it uses at least $4|V| = O(|V|)$ bytes of memory for counting vertex degrees, assuming counters are 4B integers. We could use algorithm from this question to estimate (two times) the degree of each vertex with $\epsilon > 0$ parameter of HW3Q4 algorithm, because $2 \deg_G(v)$ equals the number of edges that v “touches”. Setting, for instance, $\delta := 1/|V|^2$ we would this way use only $O(\log |V|)$ space (because ϵ is just some constant) for vertex degree counts. The total asymptotic space complexity of the HW3Q4 algorithm is unchanged.

A Source code for SVM

This appendix contains source code we wrote for the first question. (There is lots of debugging output; please ignore it.)

A.1 SVM via batch GD, stochastic GD, and mini batch GD

Listing 1: SVM via batch gradient descent.

```
1  #!/usr/bin/python
2
3  import os, sys, time
4  import numpy as np
5  from numpy import linalg
6  import timeit
7
8  # derivative of hinge by wj
9  def hinge(X, Y, w, b, j):
10     return 100*sum([0 if Y[i]*(np.dot(X[i], w)+b) >= 1 else -X[i][j]*Y[i] for i
11                   in range(len(X))]) # C=100
12
13  # compute cost function — iterates over the whole dataset
14  def f(X, Y, w, b):
15     return linalg.norm(w)**2/2 + 100*sum([max(0, 1.0-Y[i]*(np.dot(w, X[i]))+b))
16     for i in range(len(X))]) # C=100
17
18  # derivative of f by b
19  def grad_b(X, Y, w, b):
20     return 100*sum([0 if Y[i]*(np.dot(X[i], w)+b) >= 1 else -Y[i] for i in range
21                   (len(X))]) # C=100
22
23  # batch gradient descent
24  def bgd(X, Y, w, b, ni = 0.0000003, eps = 0.25):
25     assert len(X) == len(Y), "Counts_do_not_match"
26     k = 0
27     crr_f = prev_f = f(X, Y, w, b)
28     delta_cost = 1.0
29     print "crr_cost: ", crr_f
30     while True:
31         start = time.time()
32         # do the update
33         tw = list(w)
34         for j in range(len(X[0])):
35             tw[j] = w[j] - ni*(w[j] + hinge(X, Y, w, b, j))
36         b = b - ni * grad_b(X, Y, w, b)
37         w = list(tw)
38         end = time.time()
39         print "secs: ", end-start
40         print "iter: ", k
41         print "—" * 80
```

```

39     k = k+1
40     crr_f = f(X, Y, w, b)
41     print "crr_cost:", crr_f
42     delta_cost = (abs(prev_f - crr_f))*100/(prev_f)
43     print "delta_cost:", delta_cost
44     if delta_cost < eps:
45         print "Converged after", k, "steps!"
46         break
47     prev_f = crr_f
48     print "w=", w
49     print "b=", b
50     return (w, b)
51
52 def load_train(fname):
53     return [[int(el) for el in line.split(',')]] for line in open(fname).read().
        split()]
54
55 def load_test(fname):
56     return [int(el) for el in open(fname).read().split()]
57
58 # entry point
59 if __name__ == '__main__':
60     #X = load_train('features.txt')
61     #Y = load_test('target.txt')
62     #w = [0 for x in range(len(X[0]))]
63     #b = 0
64     #(w, b) = bgd(X, Y, w, b)
65     s = "\X=\load_train('features.txt')\n\
66 \Y=\load_test('target.txt')\n\
67 \w=[0 for x in range(len(X[0]))]\n\
68 \b=0\n\
69 \ (w, b) = bgd(X, Y, w, b)"
70     print(timeit.timeit(s, setup="from __main__ import *", number=1))

```

Listing 2: SVM via stochastic gradient descent.

```

1  #!/usr/bin/python
2
3  import os, sys, time
4  import numpy as np
5  from numpy import linalg
6  import random as rnd
7  import timeit
8
9  # derivative of hinge by wj; here, i is index of the training example
10 def hinge(X, Y, w, b, j, i):
11     return 100*(0 if Y[i]*(np.dot(X[i], w)+b) >= 1 else -X[i][j]*Y[i]) # C=100
12
13 # compute cost function — iterates over the whole dataset
14 def f(X, Y, w, b):
15     return linalg.norm(w)**2/2 + 100*sum([max(0, 1.0-Y[i]*(np.dot(w, X[i])+b))
        for i in range(len(X))]) # C=100

```

```

16
17 # derivative of f by b
18 def grad_b(X, Y, w, b, i):
19     return 100*(0 if Y[i]*(np.dot(X[i], w)+b) >= 1 else -Y[i]) # C=100
20
21 # stochastic gradient descent
22 def sgd(X, Y, w, b, ni = 0.0001, eps = 0.001):
23     I = range(len(X))
24     rnd.shuffle(I)
25     X = [X[i] for i in I]
26     Y = [Y[i] for i in I]
27     i, k = 1, 0
28     n = len(X)
29     crr_f = prev_f = f(X, Y, w, b)
30     print "crr_cost: ", crr_f
31     delta_cost = 0.0
32     while True:
33         tw = list(w)
34         # update
35         for j in range(len(X[0])):
36             tw[j] = w[j] - ni*(w[j] + hinge(X, Y, w, b, j, i)) # derivative of hinge
37                 by w[j]
38             b = b - ni*grad_b(X, Y, w, b, i)
39             w = list(tw)
40             i = (i % n) + 1
41             k = k+1
42             crr_f = f(X, Y, w, b)
43             print "crr_cost: ", crr_f
44             delta_cost = 0.5*delta_cost + 0.5*(abs(prev_f - crr_f))*100/(prev_f)
45             print "k=", k
46             print "delta_cost: ", delta_cost
47             print "-"*80
48             if delta_cost < eps:
49                 print "[DONE] Converged after", k, "steps"
50                 break
51             prev_f = crr_f
52     return (w, b)
53
54 def load_train(fname):
55     return [[int(el) for el in line.split(',')] for line in open(fname).read().
56         split()]
57
58 def load_test(fname):
59     return [int(el) for el in open(fname).read().split()]
60
61 # entry point
62 if __name__ == '__main__':
63     #X = load_train('features.txt')
64     #Y = load_test('target.txt')
65     #w = [0 for x in range(len(X[0]))]
66     #b = 0

```

```

65     #(w, b) = sgd(X, Y, w, b)
66     s = "\X=\load_train('features.txt')\n\
67     \Y=\load_test('target.txt')\n\
68     \w=\[0\for x\in\range(len(X[0]))]\n\
69     \b=\0\n\
70     \w, \b) =sgd(X, \Y, \w, \b)"
71     print(timeit.timeit(s, setup="from __main__ import \", number=1))

```

Listing 3: SVM via mini batch gradient descent.

```

1  #!/usr/bin/python
2
3  import os, sys, time
4  import numpy as np
5  from numpy import linalg
6  import random as rnd
7  import timeit
8
9  # derivative of hinge by wj
10 def hinge(X, Y, w, b, j, l, bs):
11     return 100*sum([0 if Y[i]*(np.dot(X[i], w)+b) >= 1 else -X[i][j]*Y[i] for i
12                     in range(l*bs+1, min(len(X), (l+1)*bs))]) # C=100
13
14 # compute cost function — iterates over the whole dataset
15 def f(X, Y, w, b):
16     return linalg.norm(w)**2/2 + 100*sum([max(0, 1.0-Y[i]*(np.dot(w, X[i])+b))
17     for i in range(len(X))]) # C=100
18
19 # derivative of f by b
20 def grad_b(X, Y, w, b, l, bs):
21     return 100*sum([0 if Y[i]*(np.dot(X[i], w)+b) >= 1 else -Y[i] for i in range
22                     (l*bs+1, min(len(X), (l+1)*bs))]) # C=100
23
24 # mini batch gradient descent
25 def mini_bgd(X, Y, w, b, ni = 0.00001, eps = 0.01, bs = 20):
26     I = range(len(X))
27     rnd.shuffle(I)
28     # reorder
29     X = [X[i] for i in I]
30     Y = [Y[i] for i in I]
31     n = len(X)
32     l, k = 0, 0
33     crr_f = prev_f = f(X, Y, w, b)
34     delta_cost = 0.0
35     print "crr_cost:\n", crr_f
36     while True:
37         start = time.time()
38         # do the update
39         tw = list(w)
40         for j in range(len(X[0])):
41             tw[j] = w[j] - ni*(w[j] + hinge(X, Y, w, b, j, l, bs))
42         b = b - ni * grad_b(X, Y, w, b, l, bs)

```

```

40     w = list(tw)
41     end = time.time()
42     print "secs:", end-start
43     print "iter:", k
44     print "-"*80
45     k = k+1
46     crr_f = f(X, Y, w, b)
47     print "crr_cost:", crr_f
48     delta_cost = 0.5*delta_cost + 0.5*(abs(prev_f - crr_f))*100/(prev_f)
49     print "delta_cost:", delta_cost
50     if delta_cost < eps:
51         print "[DONE] Converged", k, "steps!"
52         break
53     l = (l + 1) % ((n + bs - 1)/bs)
54     prev_f = crr_f
55     print "w=", w
56     print "b=", b
57     return (w, b)
58
59 def load_train(fname):
60     return [[int(el) for el in line.split(',')]] for line in open(fname).read().
        split()]
61
62 def load_test(fname):
63     return [int(el) for el in open(fname).read().split()]
64
65 # entry point
66 if __name__ == '__main__':
67     #X = load_train('features.txt')
68     #Y = load_test('target.txt')
69     #w = [0 for x in range(len(X[0]))]
70     #b = 0
71     #(w, b) = mini_bgd(X, Y, w, b)
72     s = "X=\nload_train('features.txt')\n\
73     Y=\nload_test('target.txt')\n\
74     w=\n[0 for x in range(len(X[0]))]\n\
75     b=\n0\n\
76     (w, b) =\nmini_bgd(X, Y, w, b)"
77     print(timeit.timeit(s, setup="from __main__ import *", number=1))

```

Plots were generated by running `python hw4q1-bgd.py > hw4q1-bgd-a.txt` (similarly run mini batch GD and stochastic GD) and then running `plot.py` from listing 4.

Listing 4: Code snippet used for plotting.

```

1  #!/usr/bin/python
2
3  import os, sys, time
4  import numpy as np
5  from matplotlib import pyplot
6
7  if __name__ == '__main__':

```

```

8  # Batch gradient descent
9  Y_bgd = []
10 X = []
11 k = 0
12 for ln in open('hw4q1-bgd-a.txt').read().split('\n'):
13     if ln[:8] == 'crr_cost':
14         X.append(k)
15         Y_bgd.append(float(ln.split()[1]))
16         k = k+1
17 pyplot.plot(X, Y_bgd)
18 pyplot.show()
19 print "X=_", X
20 print "Y=_", Y_bgd
21
22 # Stochastic gradient descent
23 k = 0
24 X = []
25 Y_sgd = []
26 for ln in open('hw4q1-sgd-a.txt').read().split('\n'):
27     if ln[:8] == 'crr_cost':
28         X.append(k)
29         Y_sgd.append(float(ln.split()[1]))
30         k = k+1
31 for i in range(len(Y_sgd) - len(Y_bgd)): Y_bgd.append(Y_bgd[-1])
32 pyplot.plot(X, Y_sgd, X, Y_bgd)
33 pyplot.show()
34 print "X=_", X
35 print "Y=_", Y_sgd
36
37 # Mini Batch Gradient Descent
38 k = 0
39 #X = []
40 Y_mbgd = []
41 for ln in open('mini-bgd-a.txt').read().split('\n'):
42     if ln[:8] == 'crr_cost':
43         #X.append(k)
44         Y_mbgd.append(float(ln.split()[1]))
45         k = k+1
46 for i in range(len(Y_sgd) - len(Y_mbgd)): Y_mbgd.append(Y_mbgd[-1])
47 pyplot.plot(X, Y_sgd, X, Y_bgd, X, Y_mbgd)
48 pyplot.show()
49 print "X=_", X
50 print "Y=_", Y_mbgd

```

A.2 Code used to experiment with regularization

Listing 5: Code for HW4Q1 item (f).

```

1  #!/usr/bin/python
2

```



```

3 import os, sys, time
4 import numpy as np
5 from numpy import linalg
6 import random as rnd
7 from matplotlib import pyplot
8
9 # derivative of hinge by wj; here, i is index of the training example
10 def hinge(X, Y, w, b, j, i, C):
11     return C*(0 if Y[i]*(np.dot(X[i], w)+b) >= 1 else -X[i][j]*Y[i])
12
13 # compute cost function — iterates over the whole dataset
14 def f(X, Y, w, b, C):
15     return linalg.norm(w)**2/2 + C*sum([max(0, 1.0-Y[i]*(np.dot(w, X[i])+b)) for
16         i in range(len(X))])
17
18 # derivative of f by b
19 def grad_b(X, Y, w, b, i, C):
20     return C*(0 if Y[i]*(np.dot(X[i], w)+b) >= 1 else -Y[i]) # C=100
21
22 # stochastic gradient descent
23 def sgd(X, Y, w, b, ni = 0.0001, eps = 0.001, C = 100):
24     I = range(len(X))
25     rnd.shuffle(I)
26     X = [X[i] for i in I]
27     Y = [Y[i] for i in I]
28     i, k = 1, 0
29     n = len(X)
30     crr_f = prev_f = f(X, Y, w, b, C)
31     #print "crr_cost: ", crr_f
32     delta_cost = 0.0
33     while True:
34         tw = list(w)
35         # update
36         for j in range(len(X[0])):
37             tw[j] = w[j] - ni*(w[j] + hinge(X, Y, w, b, j, i, C)) # derivative of
38                 hinge by w[j]
39         b = b - ni*grad_b(X, Y, w, b, i, C)
40         w = list(tw)
41         i = ((i+1) % n)
42         k = k+1
43         crr_f = f(X, Y, w, b, C)
44         #print "crr_cost: ", crr_f
45         delta_cost = 0.5*delta_cost + 0.5*(abs(prev_f - crr_f))*100/(prev_f)
46         if k%500 == 0: print "k=", k
47         #print "delta_cost: ", delta_cost
48         #print "-"*80
49         if delta_cost < eps:
50             #print "[DONE] Converged after", k, "steps"
51             break
52     prev_f = crr_f
53     print "Final_f=", crr_f, "after", k, "steps"

```

```

52     return (w, b)
53
54 def load_ftr(fname):
55     return [[int(el) for el in line.split(',')] for line in open(fname).read().
56             split()]
57
58 def load_trg(fname):
59     return [int(el) for el in open(fname).read().split()]
60
61 def est_err(X, Y, w, b):
62     n = len(X)
63     miss = 0
64     for i in range(n):
65         if Y[i]*(np.dot(X[i], w)+b) <= 0: miss = miss+1
66     print "[DEBUG] Misclassified", miss, "out_of", n
67     return 1.0*miss/n
68
69 # entry point
70 if __name__ == '__main__':
71     Cs = [1, 10, 50, 100, 200, 300, 400, 500]
72     # train data
73     X = load_ftr('features.train.txt')
74     Y = load_trg('target.train.txt')
75     # test data
76     X_test = load_ftr('features.test.txt')
77     Y_test = load_trg('target.test.txt')
78     E = []
79     for C in Cs:
80         X = load_ftr('features.train.txt')
81         Y = load_trg('target.train.txt')
82         w = [0 for x in range(len(X[0]))]
83         b = 0
84         (w, b) = sgd(X, Y, w, b, 0.0001, 0.001, C)
85         print "w=", w
86         print "b=", b
87         E.append(est_err(X_test, Y_test, w, b))
88         print "For C=", C, "the_percent_error_equals", E[-1]
89         print "-" * 80
90     print "E=", E
91     pyplot.plot(Cs, E)
92     pyplot.show()

```

References