

Assignment Number: 2

Student Name: Gurpreet Singh

Roll Number: 150259

Date: October 10, 2017

Part 1

No, name is not a useful attribute in learning a decision tree as splitting on this attribute will have no gain. Non-statistically speaking, there is a lot of variance in the attribute itself, and it will be very inefficient to decide on the basis of name as it can often be a newly observed value.

Part 2

It is not possible to classify the data given perfectly. Looking at data #4 and #6, they have the same values for all fields except for 'name', however have different target values. Since we shall not use 'name' as a decision classifier (as discussed in Part A), therefore we can say that it is not possible to classify the data perfectly.

Part 3

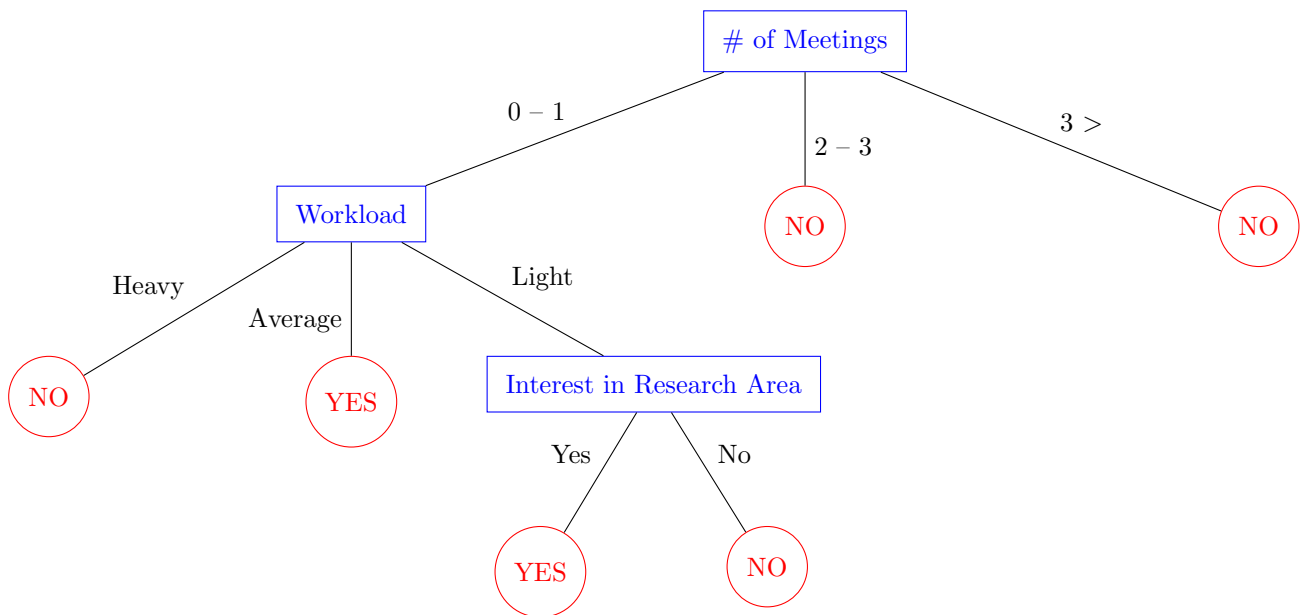


Figure 1: ID3 Tree for the given data

ID3 is a greedy strategy to generate decision trees using Information Gain. The information gain at each level/depth is mentioned below

Level 1

Entropy $S = 0.9183$

IG: Size of Research Group = $S - 0.2600 - 0.2667 - 0.3237 = 0.0679$

IG: Interest in Research Area = $S - 0.2667 - 0.6199 = 0.0317$

IG: Workload = $S - 0.4000 - 0.2163 - 0.2406 = 0.0614$

IG: Number of Meetings = $S - 0.6667 = 0.2516$

Clearly the attribute *Number of Meetings* provides the most IG, and hence we use this attribute to split at the first level.

Level 2 — Number of Meetings

Possible distinctions — ‘0 – 1’, ‘2 – 3’, ‘> 3’

Since for ‘2 – 3’ and ‘> 3’, all labels are ‘No’ (**i.e.** entropy is 0), hence we do need to split here, and we add leaf nodes.

For the samples remaining with value ‘0 – 1’ of the attribute *Number of Meetings*, we will further extend the decision tree.

Entropy $S = 1$

IG: Size of Research Group $= S - 0.0000 - 0.4000 - 0.4855 = 0.1145$

IG: Interest in Research Area $= S - 0.2755 - 0.6897 = 0.0348$

IG: Workload $= S - 0.0000 - 0.0000 - 0.3610 = 0.6390$

Clearly the attribute *Workload* gives the highest gain, and hence we use this attribute to further construct the decision tree.

Level 3 — Workload

Possible distinctions — ‘Light’, ‘Average’, ‘Heavy’

Since there is purity in the samples with ‘Average’ and ‘Heavy’ value in the *Workload* attribute (*i.e.* entropy is 0), we do not split further in this case. Therefore we are only left with two samples, which have the value ‘Light’.

Entropy $S = 1$

IG: Size of Research Group $= S - 1.0000 = 0.0000$

IG: Interest in Research Area $= S - 0.0000 = 1.0000$

Hence we further split using *Interest in Research Area* attribute. Since there is only one sample per distinction, we cannot split further. Hence we have the decision tree showed in Figure 1

Assignment Number: 2

Student Name: Gurpreet Singh

Roll Number: 150259

Date: October 10, 2017

Notations

L: Represents total number of items

N: Represents total number of users

D: Represents length of feature vector for every user *i.e.* length of \mathbf{x}^n

K: 2^L

Part 1

Any user can select any number of items. The probability of selecting an item for a user is represented by a logistic expression.

$$P(\mu_l | \mathbf{x}^n, \mathbf{w}_l) = \sigma(\langle \mathbf{w}_l, \mathbf{x}^n \rangle)^{\mu_l} (1 - \sigma(\langle \mathbf{w}_l, \mathbf{x}^n \rangle))^{(1-\mu_l)} \quad (1)$$

Note: $\sigma(x)$ represents the sigmoid function

The equation 1 represents the probability of user n selecting item l *i.e.* $\mu_l = 1$ if item l has been selected, 0 otherwise. We assume that the probability is represented by the same expression for all users and is independent, since all users are independent and are represented by their feature vector \mathbf{x}^n .

It is clear that there are K possible subsets of the items. Hence the user can choose any one of the subsets. Let these subsets be denoted by S_k . Hence the powerset of the items set is $\{S_k\}_{k=1}^K$. Now, we can write the expression for a user selecting a subset of items.

$$P(S_k | \mathbf{x}^n, \{\mathbf{w}_l\}_{l=1}^L) = \prod_{l \in S_k} P(\mu_l = 1 | \mathbf{x}^n) \prod_{l \notin S_k} P(\mu_l = 0 | \mathbf{x}^n) \quad (2)$$

Here, $\{\mathbf{w}_l\}_{l=1}^L$ is a parameter for the model.

Now we define our latent variable $z^n \in [K]$. This is the index of the subset of items, in the powerset, that the user n has chosen *i.e.* $z^n = k$ represents the subset S_k . Hence, we can define the conditional probability of z^n as follows

$$P(z^n | \mathbf{x}^n, \{\mathbf{w}_l\}_{l=1}^L) = \prod_{k=1}^K P(S_k | \mathbf{x}^n, \{\mathbf{w}_l\}_{l=1}^L)^{\mathbb{I}[z^n=k]} \quad (3)$$

We also define a function $h : [K] \rightarrow \mathbb{R}$ such that

$$h(z^n) = \sum_{l \in S_{z^n}} c_l \quad (4)$$

We can now write the expression for the conditional probability of our output variable b^n as follows

$$b^n | z^n, \sigma \sim \mathcal{N}(b^n | h(z^n), \sigma^2) \quad (5)$$

Note: $b^n | z^n, \sigma \sim b^n | z^n, \mathbf{x}^n, \sigma$. This is because once we assume a value of z^n , \mathbf{x}^n no longer plays a role in the conditional probability

Note: We can collectively define $\sigma, \{\mathbf{w}_l\}_{l=1}^L$ as Θ

Part 2

We can represent the Complete Likelihood (CLE) as follows

$$CLE \sim \mathbf{b}, \mathbf{z} | \mathbf{X}, \Theta$$

Since all users are independent, we can write the above as follows

$$P(\mathbf{b}, \mathbf{z} | \mathbf{X}, \Theta) = \prod_{n=1}^N P(b^n, z^n | \mathbf{x}^n, \Theta) \quad (6)$$

Using bayes rule, we can further expand this probability term

$$\begin{aligned} P(b^n, z^n | \mathbf{x}^n, \Theta) &= P(b^n | z^n, \Theta) P(z^n | \mathbf{x}^n, \Theta) \\ \implies P(\mathbf{b}, \mathbf{z} | \mathbf{X}, \Theta) &= \prod_{n=1}^N P(b^n | z^n, \Theta) P(z^n | \mathbf{x}^n, \Theta) \end{aligned}$$

We can also write the Complete Log Likelihood (CLL) from the above expression

$$CLL = \sum_{n=1}^N \log(P(b^n | z^n, \sigma)) + \sum_{n=1}^N \log(P(z^n | \mathbf{x}^n, \{\mathbf{w}_l\}_{l=1}^L)) \quad (7)$$

Part 3

Since directly finding Θ_{MLE} is a NP hard problem, therefore, we use alternating optimization to find out an approximate MLE estimate of Θ . The algorithm for alternating optimization is as follows

Time Complexity Analysis

We need to compute $\langle \mathbf{w}_l, \mathbf{x}^n \rangle$ for all items and users. It is optimal to save this and store it before starting the updation steps. This can be done in $\mathcal{O}(NLD)$ steps. We can also compute the values of the isomorphism h for all subsets *i.e.* compute $h(k)$ for all k . This needs to be done only once, and can be done in the preprocessing part. Hence we need not add the time for this in the iteration time.

Therefore, for all z^n updations, we need to check for all subsets, and we need to compute probability of the user choosing that subset. Since we have already computed some of the terms, we will only require $\mathcal{O}(NKL)$ time for this updation.

The MLE update of σ requires only $\mathcal{O}(N)$.

For the MLE estimate of \mathbf{w}_l , we need to find the MLE estimate of a logistic expression, using function approximation. As mentioned, we can assume that each FA takes $\mathcal{O}(ND)$ time, therefore, we can do this step in $\mathcal{O}(NLD)$.

Algorithm 1: *Alternating Optimization — Hard Assignment*

Initialize Θ_{MLE} to Θ_{MLE}^0
while **no convergence**:

Update \mathbf{z}

$$\begin{aligned}
\forall \quad n \in [N] : \quad z^n &= \arg \max_k P(z^n = k \mid b^n, \mathbf{x}^n, \Theta) \\
&= \arg \max_k \log (P(b^n, z^n = k \mid \mathbf{x}^n, \Theta)) \\
&= \arg \min_k (b^n - h(k))^2 - \sum_{l \in S_k} \log (\sigma(\langle \mathbf{w}_l, \mathbf{x}^n \rangle)) - \sum_{l \notin S_k} \log (1 - \sigma(\langle \mathbf{w}_l, \mathbf{x}^n \rangle))
\end{aligned}$$

Update MLE of Θ

$$\begin{aligned}
\sigma_{MLE} &= \arg \max_{\sigma} CLL \\
&= \arg \max_{\sigma} \log (P(\mathbf{b} \mid \mathbf{z}, \sigma)) \\
&= \sqrt{\frac{1}{N} \sum_{n=1}^N (b^n - h(z^n))^2}
\end{aligned}$$

$$\begin{aligned}
\forall \quad l \in [L] : \quad \mathbf{w}_{lMLE} &= \arg \max_{\mathbf{w}_l} CLL \\
&= \arg \max_{\mathbf{w}_l} \log (P(\mathbf{z} \mid \mathbf{X}, \{\mathbf{w}_l\}_{l=1}^L)) \\
&= \arg \max_{\mathbf{w}_l} \sum_{n=1}^N \left(\mathbb{I}[l \in S_{z^n}] \log (\sigma(\langle \mathbf{w}_l, \mathbf{x}^n \rangle)) + \mathbb{I}[l \notin S_{z^n}] \log (1 - \sigma(\langle \mathbf{w}_l, \mathbf{x}^n \rangle)) \right)
\end{aligned}$$

Hence the total time for each updation is $\mathcal{O}(NLD + N + NKL)$. Since $L \leq D$, we can clearly say that this is bounded by $\mathcal{O}(NKD)$ or $\mathcal{O}(2^L ND)$.

Part 4

In case of soft assignment, we cannot directly use the mode of the posterior (point estimate) as an estimate of \mathbf{z} . Therefore, we take the expected value of \mathbf{z} and use EM algorithm to compute the MLE estimate of Θ .

It will be easier to use the one-hot representation of the z^n , as we will simply assign the expected probability in the k^{th} index of the user buying the k^{th} subset of items. This is similar to the soft k-means algorithm.

Note: \mathbf{z}^n represents the one-hot representation, whereas z^n represents the numerical value

Algorithm 2: *Alternating Optimization — Soft Assignment*

Initialize Θ_{MLE} to Θ_{MLE}^0
while **no convergence**:

Update \mathbf{z}

$\forall n \in [N] :$

$$\forall k \in [K] : \quad \mathbb{E}[\mathbf{z}_k^n] = \frac{P(S^k | \mathbf{x}^n, \{\mathbf{w}_l\}_{l=1}^L) \mathcal{N}(b^n | h(k), \sigma^2)}{\sum_{k'=1}^K P(S^{k'} | \mathbf{x}^n, \{\mathbf{w}_l\}_{l=1}^L) \mathcal{N}(b^n | h(k'), \sigma^2)}$$

Update MLE of Θ

$$\begin{aligned} \sigma_{MLE} &= \arg \max_{\sigma} \mathbb{E}[CLL] \\ &= \arg \max_{\sigma} \mathbb{E}[\log(P(\mathbf{b} | \mathbf{Z}, \sigma))] \\ &= \sqrt{\frac{1}{N} \sum_{n=1}^N \mathbb{E}[(b^n - h(z^n))^2]} \\ &= \sqrt{\frac{1}{N} \sum_{n=1}^N \left((b^n)^2 - 2 \sum_{k=1}^K \mathbf{z}_k^n h(k) + \sum_{k=1}^K \mathbf{z}_k^n h(k)^2 \right)} \end{aligned}$$

$$\forall l \in [L] : \quad \mathbf{w}_{lMLE} = \arg \max_{\mathbf{w}_l} \mathbb{E}[CLL]$$

$$\begin{aligned} &= \arg \max_{\mathbf{w}_l} \mathbb{E}[\log(P(\mathbf{Z} | \mathbf{X}, \{\mathbf{w}_l\}_{l=1}^L))] \\ &= \arg \max_{\mathbf{w}_l} \mathbb{E} \left[\log \left(\prod_{n=1}^N \prod_{k=1}^K \left[\prod_{l' \in S_k} \sigma(\langle \mathbf{w}_{l'}, \mathbf{x}^n \rangle) \prod_{l' \notin S_k} (1 - \sigma(\langle \mathbf{w}_{l'}, \mathbf{x}^n \rangle)) \right]^{\mathbf{z}_k^n} \right) \right] \\ &= \arg \max_{\mathbf{w}_l} \sum_{k:l \in S_k} \sum_{n=1}^N \mathbf{z}_k^n \sigma(\langle \mathbf{w}_l, \mathbf{x}^n \rangle) + \sum_{k:l \notin S_k} \sum_{n=1}^N \mathbf{z}_k^n (1 - \sigma(\langle \mathbf{w}_l, \mathbf{x}^n \rangle)) \end{aligned}$$

Assignment Number: 2

Student Name: Gurpreet Singh

Roll Number: 150259

Date: October 10, 2017

Part 1

We have the optimization problem (P1) as follows

$$\begin{aligned} \min_{\mathbf{w}, \{\xi_i\}} \quad & \|\mathbf{w}\|_2^2 + \sum_{i=1}^n \xi_i^2 \\ \text{s.t. } \forall i \in [n] \quad & 1 - y^i \langle \mathbf{w}, \mathbf{x}^i \rangle \leq \xi_i \\ & \xi_i \geq 0 \end{aligned} \tag{P1}$$

Now consider an optimization problem (P2), which is similar to the optimization problem given in equation (P1), except that we do not have the condition $\xi_i \geq 0$ for any $i \in [n]$. Clearly, the convexity still holds, and hence we will have a solution for (P2).

Consider this solution to be $\mathbf{w}_0, \{\xi_i^0\}$. If we do not have any $k \in [n]$ such that $\xi_k^0 < 0$. Then we can safely claim that the solution for both the optimization problems (P1) and (P2) is the same. Otherwise, we will have at least one $k \in [n]$ such that the condition holds, i.e. $\xi_k^0 < 0$. Since this is a solution of the (P2), we can say that this satisfies the condition $1 - y^k \langle \mathbf{w}, \mathbf{x}^k \rangle \leq \xi_k^0$. Now consider the pair $\mathbf{w}, \{\xi_i^1\}$ where

$$\xi_i^1 = \begin{cases} \xi_i^0 & \text{if } i \neq k \\ 0 & \text{if } i = k \end{cases}$$

Since $0 > \xi_k^0$, we can say that it follows all constraints of (P2) and $0 < (\xi_k^0)^2$. For all other ξ_i^1 , they are the same and hence do not change anything (since all ξ_i are independent). Hence, we can say that $w, \{\xi_i^1\}$ is a solution of (P2), and $val(P2, (w, \{\xi_i^1\})) < val(P2, (w, \{\xi_i^0\}))$.

However, since $w, \{\xi_i^0\}$ was claimed to be a valid solution, this is not possible. Hence the claim that can exist a $k \in [n]$ such that $\xi_k^0 < 0$ is false. Hence, the solution of (P2) is always the same as that of (P1), which suggests that the second constraint is vacuous.

Part 2

$$\begin{aligned} \min_{\mathbf{w}, \{\xi_i\}} \quad & \|\mathbf{w}\|_2^2 + \sum_{i=1}^n \xi_i^2 \\ \text{s.t. } \forall i \in [n] \quad & 1 - y^i \langle \mathbf{w}, \mathbf{x}^i \rangle \leq \xi_i \quad \text{and} \quad -\xi_i \leq 0 \end{aligned} \tag{P1}$$

We can convert this problem to an unconstrained optimization problem using Lagrange Multipliers.

$$\min_{\mathbf{w}, \{\xi_i\}} \max_{\lambda, \gamma \geq 0} \|\mathbf{w}\|_2^2 + \sum_{i=1}^n \xi_i^2 + \sum_{i=1}^n \lambda_i (1 - y^i \langle \mathbf{w}, \mathbf{x}^i \rangle - \xi_i) - \sum_{i=1}^n \gamma_i \xi_i \quad (P1)$$

Part 3

We can convert the optimization problem (P1) given in equation (P1) to a dual problem. Assuming strong duality, we can switch the optimization variables.

$$\max_{\lambda, \gamma \geq 0} \min_{\mathbf{w}, \{\xi_i\}} \|\mathbf{w}\|_2^2 + \sum_{i=1}^n \xi_i^2 + \sum_{i=1}^n \lambda_i (1 - y^i \langle \mathbf{w}, \mathbf{x}^i \rangle - \xi_i) - \sum_{i=1}^n \gamma_i \xi_i \quad (P1)$$

Now using this, we can first optimize based on the inner optimization (using differentiation). That is, we can partially differentiate w.r.t. to \mathbf{w} and $\{\xi_i\}$ separately, in order to obtain a dual problem.

$$\begin{aligned} \frac{\partial P1}{\partial \mathbf{w}} &= 2\mathbf{w} - \sum_{i=1}^n \lambda_i y^i \mathbf{x}^i = 0 \\ \implies \mathbf{w}' &= \frac{1}{2} \sum_{i=1}^n \lambda_i y^i \mathbf{x}^i \\ \\ \frac{\partial P1}{\partial \xi_k} &= 2\xi_k - \lambda_k - \gamma_i = 0 \\ \implies \xi'_k &= \frac{\lambda_k + \gamma_i}{2} \end{aligned}$$

Hence, we can replace these values in the optimization problem, while solving the inner optimization problem.

$$\begin{aligned} & \max_{\lambda, \gamma \geq 0} \|\mathbf{w}'\|_2^2 + \sum_{i=1}^n \xi_i'^2 + \sum_{i=1}^n \lambda_i (1 - y^i \langle \mathbf{w}', \mathbf{x}^i \rangle - \xi'_i) - \sum_{i=1}^n \gamma_i \xi'_i \\ &= \max_{\lambda, \gamma \geq 0} \left\| \frac{1}{2} \sum_{i=1}^n \lambda_i y^i \mathbf{x}^i \right\|_2^2 + \frac{1}{4} \sum_{i=1}^n (\lambda_i + \gamma_i)^2 - \sum_{i=1}^n \lambda_i - \sum_{i=1}^n \lambda_i y^i \left\langle \frac{1}{2} \sum_{j=1}^n \lambda_j y^j \mathbf{x}^j, \mathbf{x}^i \right\rangle - \frac{1}{2} \sum_{i=1}^n (\lambda_i + \gamma_i)^2 \\ &= \max_{\lambda, \gamma \geq 0} \left\| \frac{1}{2} \sum_{i=1}^n \lambda_i y^i \mathbf{x}^i \right\|_2^2 - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \lambda_i \lambda_j y^i y^j \langle \mathbf{x}^j, \mathbf{x}^i \rangle - \frac{1}{4} \sum_{i=1}^n (\lambda_i + \gamma_i)^2 - \sum_{i=1}^n \lambda_i \end{aligned}$$

We can expand the first term as

$$\begin{aligned} & \left\| \frac{1}{2} \sum_{i=1}^n \lambda_i y^i \mathbf{x}^i \right\|_2^2 = \frac{1}{4} \left\| \sum_{i=1}^n \lambda_i y^i \mathbf{x}^i \right\|_2^2 \\ \implies & \left\| \frac{1}{2} \sum_{i=1}^n \lambda_i y^i \mathbf{x}^i \right\|_2^2 = \frac{1}{4} \sum_{i=1}^n \sum_{j=1}^n \lambda_i \lambda_j y^i y^j \langle \mathbf{x}^i, \mathbf{x}^j \rangle \end{aligned}$$

Replacing back in the original equation

$$= \max_{\lambda, \gamma \geq 0} -\frac{1}{4} \sum_{i=1}^n \sum_{j=1}^n \lambda_i \lambda_j y^i y^j \langle \mathbf{x}^j, \mathbf{x}^i \rangle - \frac{1}{4} \sum_{i=1}^n (\lambda_i + \gamma_i)^2 - \sum_{i=1}^n \lambda_i$$

We can alter this optimization problem by multiplying with -4 and inverting the optimizing condition. Hence the dual can be given as follows

$$\min_{\lambda, \gamma \geq 0} \lambda^T Q \lambda + \sum_{i=1}^n (\lambda_i + \gamma_i)^2 - 4\lambda_i \quad (D1)$$

where

$$Q_{ij} = y^i y^j \langle \mathbf{x}^i, \mathbf{x}^j \rangle$$

Note that $\lambda_i, \gamma_i \geq 0$ for all $i \in [n]$. Hence, $(\lambda_i + \gamma_i)^2 \geq (\lambda_i)^2$, hence it is very trivial to see that $\gamma_i = 0$ for all $i \in [n]$. Therefore, we can further simplify the dual as follows

$$\min_{\lambda \geq 0} \lambda^T Q \lambda + \sum_{i=1}^n (\lambda_i^2 - 4\lambda_i) \quad (D1)$$

where

$$Q_{ij} = y^i y^j \langle \mathbf{x}^i, \mathbf{x}^j \rangle$$

Part 4

Let us state the original dual problem.

$$\min_{\alpha \in [0, C]^n} \alpha^T Q \alpha - \sum_{i=1}^n \alpha_i$$

The main difference that seems to be there between the two objectives is that the dual variable (each value) in the original problem is upper bounded by a constant C .

It seems that the dual variable in (D1) has no upper bound, however there is still a tradeoff between the high and low values of our dual variable. Consider $\lambda > 2$. In this case, the second term will become positive. Hence, we do not want very high values of λ . Looking at the primal problem of the original SVM, we can say that in our case (by similarity of terms), $C = 2$. Hence we are trying to keep the value of λ within the range $[0, C]^n$, however it is a loose bound.

Therefore, the only difference is that in the original SVM problem, we have a strict upper bound on the value of the dual variable, whereas in our case, the upper bound is loose, but still existant.

Part 5

No. From our approach, we have proved that $\xi_i < 0$ only increases the value of the objective function while $\xi_i = 0$ is a solution, and hence the lower bound is vacuous, however, this is not the case with the original SVM. In the original problem, we are hoping to get negative values of ξ_i which would mean that we are giving weightage to the points which are far from the margin, which is not a correct optimization problem, as it essentially ignores the large margin problem. Hence, in that case, we need to explicitly add the lower bound.

Assignment Number: 2

Student Name: Gurpreet Singh

Roll Number: 150259

Date: October 10, 2017

Part 3

The averaging trick did not give better results in the case of Vanilla Gradient Descent. Determined on the basis of held-out validation, the model with w_bar gave slightly lesser accuracy as well as higher value of the objective function for all iterations. This is evident in the validation data provided in figure 2

Part 4

I took $h(n)$ as $\frac{1}{n}$, as it allows for averaging the sum of the support vectors. As for the constant multiplied to $\frac{h(n)}{\sqrt{t+1}}$, I used held-out validation to find the fastest converging as well as the one with the best accuracy. Amongst the set of values $\{1, 8, 16, 32, 64, 100\}$, the best value came out to be 32.

Hence, the best value of η for me was $\frac{32}{n\sqrt{t+1}}$. This is inspite of the fact that initially the value of the objective function increases *i.e.* the objective function first diverges, then converges. Using a smaller value of the constant (5) avoids this problem, however gives lower final accuracy and slower convergence.

C = 1.000000		C = 8.000000		C = 16.000000		C = 32.000000		C = 64.000000		C = 100.000000	
w	\hat{w}	w	\hat{w}	w	\hat{w}	w	\hat{w}	w	\hat{w}	w	\hat{w}
240000.0000	240000.0000	240000.0000	240000.0000	240000.0000	240000.0000	240000.0000	240000.0000	240000.0000	240000.0000	240000.0000	240000.0000
185361.6590	196634.2028	185361.6590	196634.2028	185361.6590	196634.2028	185361.6590	196634.2028	185361.6590	196634.2028	185361.6590	196634.2028
177999.8025	188166.3755	177999.8025	188166.3755	177999.8025	188166.3755	177999.8025	188166.3755	177999.8025	188166.3755	177999.8025	188166.3755
172602.1798	183499.6096	172602.1798	183499.6096	172602.1798	183499.6096	172602.1798	183499.6096	172602.1798	183499.6096	172602.1798	183499.6096
168201.4599	180009.1548	168201.4599	180009.1548	168201.4599	180009.1548	168201.4599	180009.1548	168201.4599	180009.1548	168201.4599	180009.1548
165127.4616	177155.1981	165127.4616	177155.1981	165127.4616	177155.1981	165127.4616	177155.1981	165127.4616	177155.1981	165127.4616	177155.1981
163076.0913	174792.9345	163076.0913	174792.9345	163076.0913	174792.9345	163076.0913	174792.9345	163076.0913	174792.9345	163076.0913	174792.9345
161658.9783	172825.5421	161658.9783	172825.5421	161658.9783	172825.5421	161658.9783	172825.5421	161658.9783	172825.5421	161658.9783	172825.5421
160677.2284	171169.0698	160677.2284	171169.0698	160677.2284	171169.0698	160677.2284	171169.0698	160677.2284	171169.0698	160677.2284	171169.0698
159910.2211	169760.2540	159910.2211	169760.2540	159910.2211	169760.2540	159910.2211	169760.2540	159910.2211	169760.2540	159910.2211	169760.2540
159289.1608	168548.6620	159289.1608	168548.6620	159289.1608	168548.6620	159289.1608	168548.6620	159289.1608	168548.6620	159289.1608	168548.6620
158767.1894	167495.9523	158767.1894	167495.9523	158767.1894	167495.9523	158767.1894	167495.9523	158767.1894	167495.9523	158767.1894	167495.9523
158338.4925	166573.9946	158338.4925	166573.9946	158338.4925	166573.9946	158338.4925	166573.9946	158338.4925	166573.9946	158338.4925	166573.9946
157982.8953	165761.7493	157982.8953	165761.7493	157982.8953	165761.7493	157982.8953	165761.7493	157982.8953	165761.7493	157982.8953	165761.7493
157680.8009	165042.7061	157680.8009	165042.7061	157680.8009	165042.7061	157680.8009	165042.7061	157680.8009	165042.7061	157680.8009	165042.7061
157416.3340	164402.9631	157416.3340	164402.9631	157416.3340	164402.9631	157416.3340	164402.9631	157416.3340	164402.9631	157416.3340	164402.9631
157179.2033	163831.2193	157179.2033	163831.2193	157179.2033	163831.2193	157179.2033	163831.2193	157179.2033	163831.2193	157179.2033	163831.2193
156965.1224	163317.5408	156965.1224	163317.5408	156965.1224	163317.5408	156965.1224	163317.5408	156965.1224	163317.5408	156965.1224	163317.5408
156769.3786	162853.8360	156769.3786	162853.8360	156769.3786	162853.8360	156769.3786	162853.8360	156769.3786	162853.8360	156769.3786	162853.8360
156588.3341	162432.9860	156588.3341	162432.9860	156588.3341	162432.9860	156588.3341	162432.9860	156588.3341	162432.9860	156588.3341	162432.9860
156421.4966	162049.2815	156421.4966	162049.2815	156421.4966	162049.2815	156421.4966	162049.2815	156421.4966	162049.2815	156421.4966	162049.2815
156268.2215	161697.8960	156268.2215	161697.8960	156268.2215	161697.8960	156268.2215	161697.8960	156268.2215	161697.8960	156268.2215	161697.8960
156125.7199	161374.8761	156125.7199	161374.8761	156125.7199	161374.8761	156125.7199	161374.8761	156125.7199	161374.8761	156125.7199	161374.8761
155991.5272	161076.4946	155991.5272	161076.4946	155991.5272	161076.4946	155991.5272	161076.4946	155991.5272	161076.4946	155991.5272	161076.4946
155863.8713	160799.4608	155863.8713	160799.4608	155863.8713	160799.4608	155863.8713	160799.4608	155863.8713	160799.4608	155863.8713	160799.4608
155741.7933	160541.3166	155741.7933	160541.3166	155741.7933	160541.3166	155741.7933	160541.3166	155741.7933	160541.3166	155741.7933	160541.3166
155624.1287	160300.2136	155624.1287	160300.2136	155624.1287	160300.2136	155624.1287	160300.2136	155624.1287	160300.2136	155624.1287	160300.2136
155510.6101	160074.3320	155510.6101	160074.3320	155510.6101	160074.3320	155510.6101	160074.3320	155510.6101	160074.3320	155510.6101	160074.3320
155400.8141	159861.9535	155400.8141	159861.9535	155400.8141	159861.9535	155400.8141	159861.9535	155400.8141	159861.9535	155400.8141	159861.9535
155294.2074	159661.6260	155294.2074	159661.6260	155294.2074	159661.6260	155294.2074	159661.6260	155294.2074	159661.6260	155294.2074	159661.6260
155190.4724	159471.9985	155190.4724	159471.9985	155190.4724	159471.9985	155190.4724	159471.9985	155190.4724	159471.9985	155190.4724	159471.9985
155089.3755	159292.3037	155089.3755	159292.3037	155089.3755	159292.3037	155089.3755	159292.3037	155089.3755	159292.3037	155089.3755	159292.3037
154990.6788	159121.7076	154990.6788	159121.7076	154990.6788	159121.7076	154990.6788	159121.7076	154990.6788	159121.7076	154990.6788	159121.7076
154894.1411	158959.3671	154894.1411	158959.3671	154894.1411	158959.3671	154894.1411	158959.3671	154894.1411	158959.3671	154894.1411	158959.3671
154799.6058	158804.7193	154799.6058	158804.7193	154799.6058	158804.7193	154799.6058	158804.7193	154799.6058	158804.7193	154799.6058	158804.7193
154706.9148	158657.0877	154706.9148	158657.0877	154706.9148	158657.0877	154706.9148	158657.0877	154706.9148	158657.0877	154706.9148	158657.0877
154615.9505	158515.8929	154615.9505	158515.8929	154615.9505	158515.8929	154615.9505	158515.8929	154615.9505	158515.8929	154615.9505	158515.8929
154526.5797	158380.7643	154526.5797	158380.7643	154526.5797	158380.7643	154526.5797	158380.7643	154526.5797	158380.7643	154526.5797	158380.7643
154438.6799	158251.1541	154438.6799	158251.1541	154438.6799	158251.1541	154438.6799	158251.1541	154438.6799	158251.1541	154438.6799	158251.1541
154352.2039	158126.7122	154352.2039	158126.7122	154352.2039	158126.7122	154352.2039	158126.7122	154352.2039	158126.7122	154352.2039	158126.7122
154267.1603	158007.1840	154267.1603	158007.1840	154267.1603	158007.1840	154267.1603	158007.1840	154267.1603	158007.1840	154267.1603	158007.1840
154183.1841	157891.8918	154183.1841	157891.8918	154183.1841	157891.8918	154183.1841	157891.8918	154183.1841	157891.8918	154183.1841	157891.8918
154100.4909	157780.7716	154100.4909	157780.7716	154100.4909	157780.7716	154100.4909	157780.7716	154100.4909	157780.7716	154100.4909	157780.7716
154018.9477	157673.4428	154018.9477	157673.4428	154018.9477	157673.4428	154018.9477	157673.4428	154018.9477	157673.4428	154018.9477	157673.4428
153938.5330	157569.6030	153938.5330	157569.6030	153938.5330	157569.6030	153938.5330	157569.6030	153938.5330	157569.6030	153938.5330	157569.6030
153859.2947	157469.1997	153859.2947	157469.1997	153859.2947	157469.1997	153859.2947	157469.1997	153859.2947	157469.1997	153859.2947	157469.1997
153781.6408	157371.8352	153781.6408	157371.8352	153781.6408	157371.8352	153781.6408	157371.8352	153781.6408	157371.8352	153781.6408	157371.8352
153705.4358	157277.3368	153705.4358	157277.3368	153705.4358	157277.3368	153705.4358	157277.3368	153705.4358	157277.3368	153705.4358	157277.3368
153630.6142	157185.5358	153630.6142	157185.5358	153630.6142	157185.5358	153630.6142	157185.5358	153630.6142	157185.5358	153630.6142	157185.5358
153557.2351	157096.2574	153557.2351	157096.2574	153557.2351	157096.2574	153557.2351	157096.2574	153557.2351	157096.2574	153557.2351	157096.2574

rediction Accuracy:

C = 1.000000		C = 8.000000		C = 16.000000		C = 32.000000		C = 64.000000		C = 100.000000	
w	\hat{w}	w	\hat{w}	w	\hat{w}	w	\hat{w}	w	\hat{w}	w	\hat{w}
68.25%	68.42%	71.64%	69.05%	76.09%	74.93%	76.32%	76.18%	76.32%	76.08%	76.32%	75.80%

Figure 2: Cross validation data for some values of η

Part 5 and Part 6

The plots have been generated after running GD and SCD on the complete training dataset *i.e.* 300K examples.

For GD, I have taken $n_iter = 2000$ and $spacing = 10$, whereas for SCD, $n_iter = 100000$ and $spacing = 5000$.

From the figures, it is evident that Gradient Descent gives lower objective function value, and thus a better solution, whereas SCD gives a lot of fluctuations. However, this is not actually the case. In the case of SCD, there is obviously a tradeoff between the objective value and margin, and the speed. The SCD method very quickly gives good accuracy, whereas the GD takes much more time to give the same accuracy.

Also, looking at the theoretical plots, it seems that the work done by SCD is much smaller than that of GD. However, this is not in relevance to the observed time for the same set of training data. The theoretical time and wall-clock time are definitely correlated, however, it seems that there is a lot more overhead on $\mathcal{O}(1)$ operations than expected by the work done per iteration in the SCD method. The difference in the methods might become more relevant in case of much larger values of the size of the dataset used.

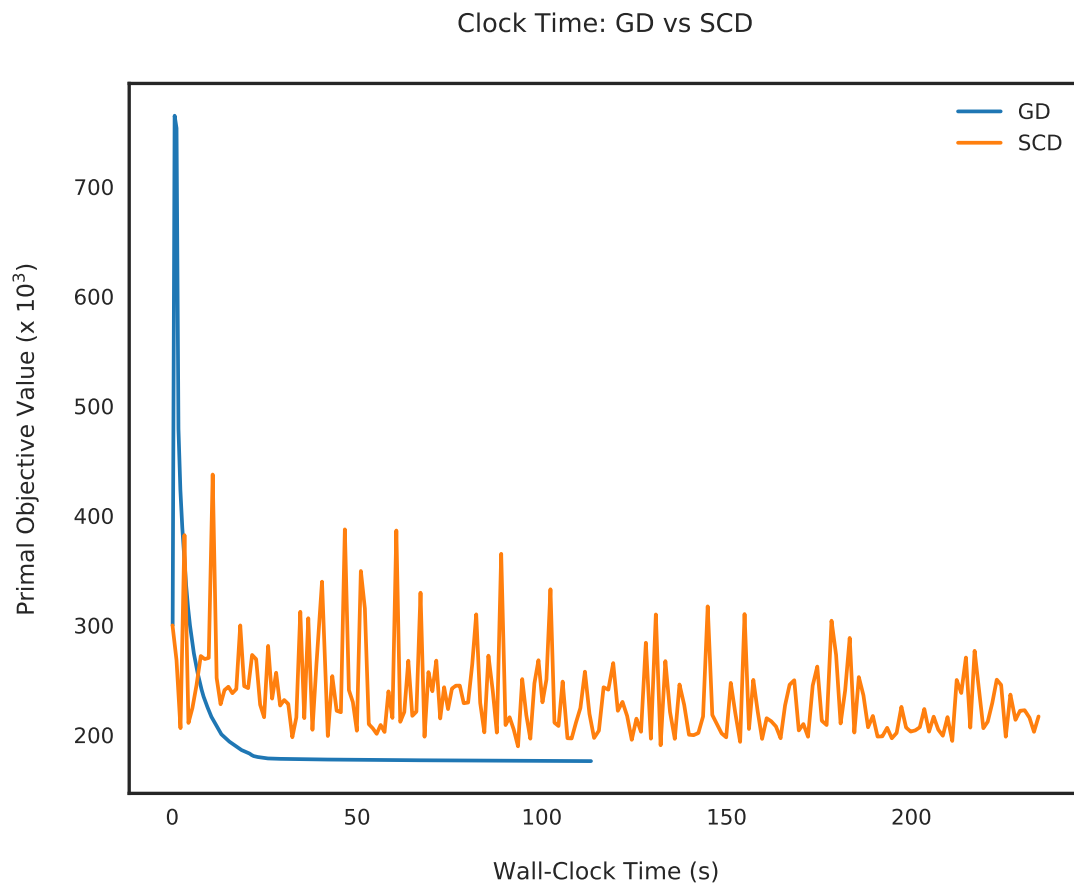


Figure 3: Clock Time Analysis of Gradient Descent vs Stochastic Coordinate Descent

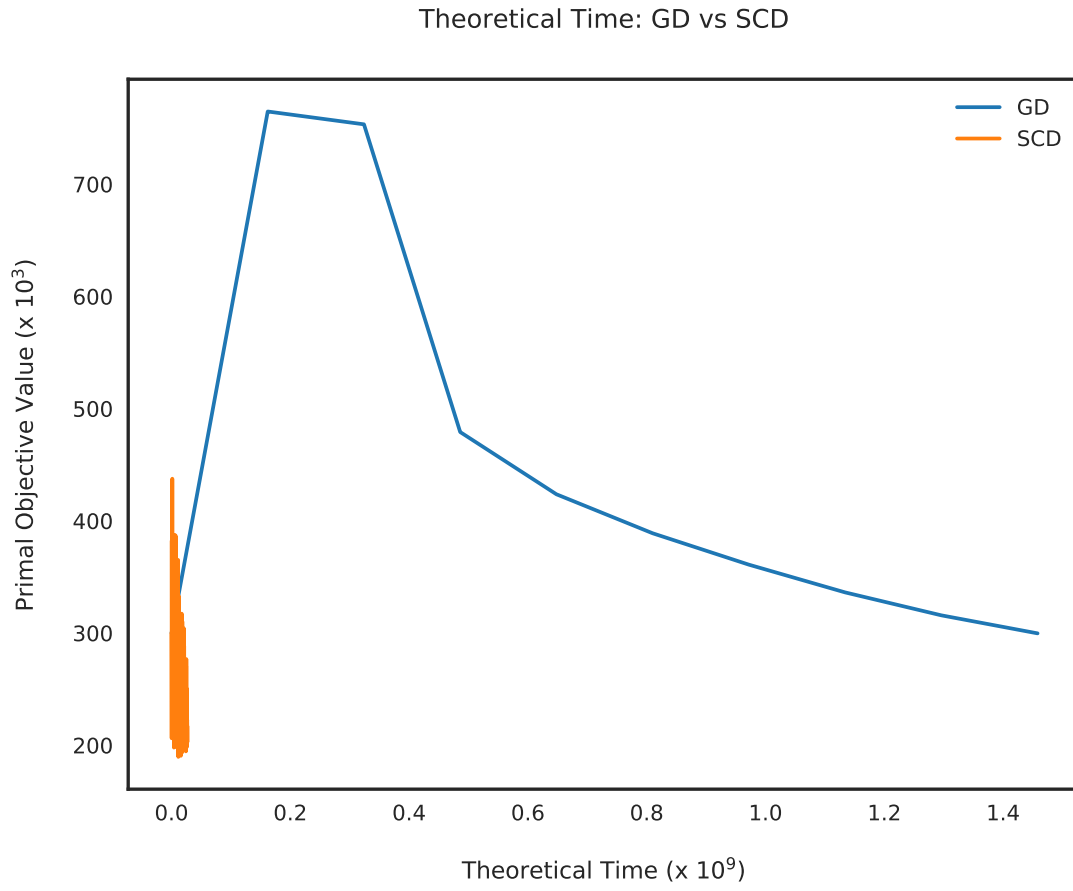


Figure 4: Theoretical Time Analysis of Gradient Descent vs Stochastic Coordinate Descent

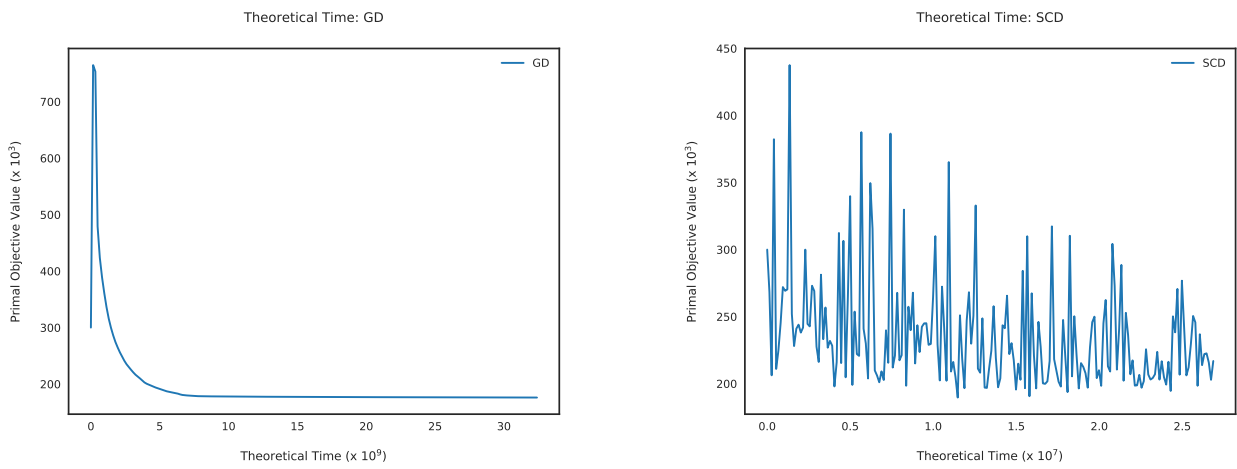


Figure 5: Theoretical Time Analysis of (a) Gradient Descent and (b) Stochastic Coordinate Descent