

# ML EXAM - 2

Name: VYOM PATHAK

I.D.: 96703101

Q1

a) It is important to reserve some images for the test dataset in order to evaluate the performance of the model on unseen data. As a whole, we do this to check the generalization capability of the model.

b) One possible explanation that the model performs well on the training and poorly on the test set is that the model is overfitting on the training dataset i.e. model shows high variance.

Q2

a) False b) False c) True

d) False e) False

Q5

→ Considering  $n$  vectors:

$$\phi_1, \phi_2, \dots, \phi_n$$

→ Some of the vectors are pre-assigned to a specific group.

→ The first step of k-means algorithm is partitioning of the vectors in  $K$  groups.

→ We assign each  $\phi_i$  to a group with the nearest representation value.

→ Now, we overwrite this assignment of the vectors which are pre-assignment to a specific group.

→ This helps in faster convergence.

→ Eg: Let there be patients with  $m$  medical features.

→ For some of these patients, the diagnosis is confirmed irrespective of their medical features, hence we already know the corresponding group is already known based on the diagnosis.

→ Now, the first step of the k-means algorithm will be to partition the patients in k-groups with  $f_i$  as the features having close representation of a group classified by a disease.

→ After this, the patients with definitive diagnosis are over-written and are assigned to their pre-set diagnosed disease irrespective of their features.

- Q6 @ SVD - We use SVD instead of SVDS because the size of the data matrix is small ( $100 \times 500$ ) as well as the dataset is described as dense.
- It returns the singular values of  $\Phi$  in descending order.
- Firstly we will center the given data matrix :
- ⇒  $n = \text{size}(\Phi, 1);$
- ⇒  $\text{center} = \text{mean}(\Phi, 1);$
- ⇒  $\tilde{\Phi} = X - \text{repmat}(\text{center}, n, 1);$
- Then, we will perform SVD to obtain  $V$  as the projection

$$\Rightarrow [U, S, \Theta] = \text{svd}(\tilde{\Phi});$$

$\rightarrow$  Then we use the projections to find the embedding matrix  $Y$ .

$$\Rightarrow Y = \underbrace{\tilde{\Phi} \Theta(:, 1:10)}_{\text{SVDS}};$$

⑥ SVDS - We use SVDS instead of SVD because the size of the data matrix is large ( $10^5 \times 10^6$ ) and the dataset is described as sparse ( $10^7$  non-zeros). The function SVDS performs faster calculations under these conditions.

- The steps are similar to the above we just use the svds() function instead of svd().
- Firstly we will center the given data matrix :
- ⇒  $n = \text{size}(\Phi, 1);$
- ⇒  $\text{center} = \text{mean}(\Phi, 1);$
- ⇒  $\tilde{\Phi} = X - \text{repmat}(\text{center}, n, 1);$
- Then, we will perform SVD to obtain V as the projections
- ⇒  $[U, S, \Theta] = \text{svds}(\tilde{\Phi}, 100);$
- Then we use the projections to find the embedding matrix Y.

$$\Rightarrow \gamma \in \overline{\bigcup_{i=1}^m \Theta_i};$$

(P) We need to prove the given inequality using the given assumption

PGD is given as :

$$\theta^{(t+1)} = \text{Prox}_{\gamma L}(\theta^{(t)} - \gamma^{(t)} \nabla L(\theta^{(t)}))$$

→ We define  $G_r(\theta)$  which is a gradient at iteration  $t$ ,

$$G_r(\theta) \triangleq \frac{1}{r} \left( \theta - \text{Prox}_{\gamma L}(\theta - \gamma \nabla L(\theta)) \right)$$

→ At every iteration,

$$\theta^{t+1} = \theta - \gamma^{(t)} G_{\gamma^{(t)}}(\theta^t)$$

→ At optimum,  $\nabla_{\theta} L(\theta^*) = 0$ .

→ We require that the gradient  $L(\cdot)$  is Lipschitz continuous (which is also given).

→ Now we use the Central Lemma & the Descent Lemma (which is also given).

→ By convexity of  $L(\cdot)$  we have,

$$L(z) \geq L(\theta) + \langle \nabla L(\theta), z - \theta \rangle.$$

$\theta = \theta^t$  &  $\theta' = \theta^{t+1}$ . When we combine this with the Descent Lemma we get the following:

$$L(\theta') \leq L(z) - \langle \nabla L(\theta), z - \theta \rangle$$

$$+ \langle \nabla L(\theta), \theta' - \theta \rangle$$

$$+ \frac{M}{2} \|\theta' - \theta\|^2$$

where  $L(z) = L(\theta) + \gamma \varphi(\theta^t)$

$\Rightarrow$  This simplifies to the following:

$$L(\hat{\theta}^{(t+1)}) + \lambda R(\hat{\theta}^{(t+1)}) \leq L(\hat{\theta}) + \lambda R(\hat{\theta})$$

$$+ \frac{1}{\sqrt{t}} (\hat{\theta}^{(t)} - \hat{\theta}^{(t+1)})^T (\hat{\theta}^{(t+1)} - \hat{\theta}) \\ + \frac{M}{2} \| \hat{\theta}^{(t)} - \hat{\theta}^{(t+1)} \|^2$$

Q3

a) Yes, Huber loss is differentiable

$$l(\varepsilon) = \begin{cases} \varepsilon^2/2 & |\varepsilon| \leq a \\ a(|\varepsilon| - a) & |\varepsilon| > a \end{cases}$$

$$= \begin{cases} \varepsilon^2/2 & |\varepsilon| \leq a \\ a(\varepsilon - a) & |\varepsilon| > a \text{ & } \varepsilon > 0 \\ a(-\varepsilon - a) & |\varepsilon| > a \text{ & } \varepsilon < 0 \end{cases}$$

$$\nabla l(\varepsilon) = \begin{cases} \varepsilon & |\varepsilon| \leq a \\ a & \varepsilon > a \\ -a & \varepsilon < -a \end{cases}$$

$$g_i = \nabla (\lambda(\phi_i^T \theta - y_i)) = \begin{cases} \phi_i^T (\phi_i^T \theta - y_i) & |\epsilon| \leq a \\ a & \epsilon > a \\ -a & \epsilon < -a \end{cases}$$

⑥ Algo. for PGD,

Initialize  $\theta^{(0)}$ ,  $\gamma$  (step size)

for  $t=1, \dots, d\theta$

for  $i=1, \dots, d\theta$

compute  $\theta_c = \arg \min_j (\lambda(\phi_{c,i}^T \theta_j - y_i))$

end for

$$\theta_c^{(t+1)} \leftarrow \theta^{(t)} + \gamma(a)$$

$$\gamma|\theta_c| > a \Rightarrow \theta_c > 0$$

$$\theta_i^{(t+1)} \leftarrow \theta^{(t)} - \gamma(a)$$

$$\gamma|\theta_c| > a \Rightarrow \theta_c < 0$$

for  $i=1 \text{ to } N$

$$v = \| \theta_j \|$$

$$\text{if } v \leq \sqrt{\lambda}, \quad \theta_j^{(t+1)} = 0$$

use

$$\theta_j^{(t+1)} \leftarrow \left(1 - \frac{\gamma\lambda}{V}\right) \theta_j^{(t+1)}$$

endif

end for

(c) Stochastic PGD ,

$\theta^{(0)}$ ,  $V$  (step size)

$g_i \in \nabla l(\theta)$

$$\theta_i \leftarrow \arg \min_{\theta_i} \left( f(\theta_i) + \frac{1}{2} \| \theta_i - \theta_i^{(t)} - V g_i \| \right)$$

for  $j=1$  to  $M$

$$V \in \| \theta_j \|$$

if  $V \leq \gamma\lambda$  then  $\theta_j^{(t+1)} = 0$

use

$$\theta_j^{(t+1)} \leftarrow \left(1 - \frac{\gamma\lambda}{V}\right) \theta_j^{(t+1)}$$

endif

end for

