

C247 Midterm

ASHISH KUMAR SINGH

UID : 105479019

(i) (a)

O = Answer i-1 1 = & ref

1a)

(ii) (a)

4 outer points will give wrong class

Wicn tso of error = $\frac{4}{14} = 0.2857$

02. skob mewm of misclass. tan

mewm of blots mewm report pribut
skob mewm report no

(i) L

C, 8

denied mewm avst Wicn & mewm

skob tan, 2013 29612 from Atico

(i) (b)

skob with diff. case (8)

(i) (b)

, C, A

Answer :- 2

2(a) i)

$$L = -y \cdot s$$

$$\frac{\partial L}{\partial s} = -y$$

$$s = \frac{z^T \hat{z}}{\|z\|_2 \|\hat{z}\|_2}$$

Using derivative quotient rule,

$$\frac{\partial s}{\partial z} = \frac{(\|z\|_2 \|\hat{z}\|_2) \hat{z} - (z^T \hat{z}) \left(\frac{\|z\|_2}{\|\hat{z}\|_2} \right) z}{\|z\|_2^2 \|\hat{z}\|_2^2}$$

$$\text{and } \frac{\partial s}{\partial \hat{z}} = \frac{(\|z\|_2 \|\hat{z}\|_2) z - (z^T \hat{z}) \left(\frac{\|z\|_2}{\|\hat{z}\|_2} \right) \hat{z}}{\|z\|_2^2 \|\hat{z}\|_2^2}$$

Using chain rule,

$$\nabla_z L = \frac{\partial L}{\partial z} = \frac{\partial s}{\partial z} \cdot \frac{\partial L}{\partial s} = -\left(\frac{\partial s}{\partial z}\right) \cdot y$$

$$\nabla_{\hat{z}} L = \frac{\partial L}{\partial \hat{z}} = \frac{\partial s}{\partial \hat{z}} \cdot \frac{\partial L}{\partial s} = -\left(\frac{\partial s}{\partial \hat{z}}\right) \cdot y$$

②

2 a) ii)

Using chain rule for matrix derivatives,

$$\nabla_{w_2} L = \frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial z} h_i^T + \frac{\partial L}{\partial \hat{z}} \hat{h}_i^T$$

$$\nabla_{w_2} L = \delta_z h_i^T + \delta_{\hat{z}} \hat{h}_i^T$$

2 a) iii), then writing indicator prob

$$\nabla_{h_i} L = \frac{\partial L}{\partial h_i} = w_2^T \frac{\partial L}{\partial z}$$

$$\delta_{h_i} = w_2^T \delta_z$$

$$\nabla_{\hat{h}_i} L = \frac{\partial L}{\partial \hat{h}_i} = w_2^T \frac{\partial L}{\partial \hat{z}}$$

$$\delta_{\hat{h}_i} = w_2^T \delta_{\hat{z}}$$

2 a) iv)

$$\nabla_m L = \frac{\partial L}{\partial m} = \mathbb{I}(w, x \geq 0) \odot \delta_{h_i}$$

$$\nabla_n L = \frac{\partial L}{\partial n} = \mathbb{I}(w, \hat{x} \geq 0) \odot \delta_{\hat{h}_i}$$

where indicator func $\mathbb{I}(w, x \geq 0) = 1$ when
 $(w, x \geq 0)$ else 0

2a) v)

$$\nabla_{W_1} L = \frac{\partial L}{\partial W_1} = \frac{\partial L}{\partial m} x^T + \frac{\partial L}{\partial n} \hat{x}^T$$

$$\nabla_{W_1} L = \delta_m x^T + \delta_n \hat{x}^T$$

2 b) i)

Given $y = +1$, $z^{(g)} = \hat{z}^{(g)}$

$$S = \frac{z^T \hat{z}}{\|z\|_2 \|\hat{z}\|_2} \Rightarrow \cancel{z^{(g)}}$$

Putting $z = S$ $\frac{(z^{(g)})^T z^{(g)}}{\|z^{(g)}\|_2 \|z^{(g)}\|_2} = \frac{\|z^{(g)}\|_2^2}{\|z^{(g)}\|_2^2} = 1$

$$L = -y \cdot S = -1$$

$$\boxed{L = -1}$$

2 b) ii)

When $z^{(g)}$ & $\hat{z}^{(g)}$ are orthogonal

$$(z^{(g)})^T \cdot \hat{z}^{(g)} = 0$$

$$S = 0$$

$$L = -y \cdot S = 0$$

$$\boxed{L = 0}$$

$$2 b) iii)$$

$$z^{(g)} = -\hat{z}^{(g)}$$

$$s = \frac{(z^{(g)})^T (-z^{(g)})}{\|z^{(g)}\|_2 \|-\hat{z}^{(g)}\|_2}$$

$$s = -\frac{\|z^{(g)}\|_2^2}{\|z^{(g)}\|_2^2} = -1$$

$$L = -y \cdot s = 1$$

$$\boxed{L = 1}$$

2 c)

Yes,

Since $y=+1$, both signatures are genuine, so, when embedding are same loss should be minimum.

And when embedding are $\hat{z}^{(g)} = -z^{(g)}$ loss is maximum.

For orthogonal embeddings, ideally loss should be positive but it doesn't matter much as small perturbations will lead to non-orthogonal embeddings in next iterations.

Answer :- 3

(b)

3 a)

B. & C. are not differentiable at zero

B. & C. are not differentiable at zero

~~Both linear and sigmoid have zero for zero~~

Derivative of tanh & sigmoid is close to 0 for large feature values. which when multiplied in chain rule leads to vanishing gradients

3 b)

C. & D.

no more iterations, efficient

have learnable parameters γ, β

non-linear unit stats

BN changes statistic to ~~var~~ mean which can be seen as adding noise

BN is applicable at test time

3 c)

A, C, D

L1 leads to some weights

add reg might increase loss

3d)

A, B, D

We can augment imbalanced class more
Multitask can be applicable even if we
have small data for a task as we can
augment.

3e)

- (1) Loss
- (2) Epochs / Number of Iterations
- (3) Best validation accuracy / least val loss
- (4) Early stopping iteration
- (5) Validation loss
- (6) Training Loss

Answer :- 4

4a) i)

GD with momentum and Adam will get out of the trap because both of them have momentum component which ~~will~~ will push θ towards right even when gradients are small in plateau region.

4a) ii)

GD momentum > Adam > Adagrad
GD momentum uses vanilla momentum which will give greatest update to θ .
Adam uses exponentially weight momentum so magnitude will be less.
Adagrad will ~~stop~~ decrease magnitude in that direction as previous step were in same dire.

4b)

$$E[g_t] = \mu, \quad t=1, 2, \dots$$

$$E[v_t] = E\left[\left(1-\beta_1\right) \sum_{i=1}^t \beta_1^{t-i} g_i\right]$$

Using Linearity of expectations,

$$E[v_t] = \left(1-\beta_1\right) \sum_{i=1}^t \beta_1^{t-i} E[g_i]$$

$$E[v_t] = \left(1-\beta_1\right) \sum_{i=1}^t \beta_1^{t-i} \mu$$

$$E[v_t] = \mu \left(1-\beta_1\right) \left(1 + \beta_1 + \beta_1^2 + \dots + \beta_1^{t-1}\right)$$

Using GP sum formula

$$E[v_t] = \mu \left(1-\beta_1\right) \frac{\left(1-\beta_1^t\right)}{\left(1-\beta_1\right)}$$

$$E[v_t] = \mu \left(1-\beta_1^t\right)$$

$$\gamma_1 = \frac{1}{\left(1-\beta_1^t\right)}$$

Similar calculation can be done for γ_2

$$\gamma_2 = \frac{1}{\left(1-\beta_2^t\right)}$$

Answer :- 5.

$$\begin{aligned}\hat{\theta}_{\text{new}} &= \hat{\theta}_{\text{old}} + v_{\text{new}} + \alpha(v_{\text{new}} - v_{\text{old}}) \\&= \hat{\theta}_{\text{old}} + v_{\text{new}}(1+\alpha) - \alpha v_{\text{old}} \\&= \hat{\theta}_{\text{old}} + (1+\alpha)(\alpha v_{\text{old}} - \epsilon \nabla_{\theta} L(\hat{\theta}_{\text{old}}) - \alpha v_{\text{old}}) \\&\leftarrow \cancel{\hat{\theta}_{\text{old}} + (1+\alpha)}$$

$$\begin{aligned}\hat{\theta}_{\text{new}} &= \hat{\theta}_{\text{old}} + \alpha v_{\text{old}} - \epsilon \nabla_{\theta} L(\hat{\theta}_{\text{old}}) + \alpha^2 v_{\text{old}} \\&\quad - \epsilon \alpha \nabla_{\theta} L(\hat{\theta}_{\text{old}}) - \alpha v_{\text{old}}$$

$$\hat{\theta}_{\text{new}} = \hat{\theta}_{\text{old}} + \alpha^2 v_{\text{old}} - (1+\alpha) \epsilon \nabla_{\theta} L(\hat{\theta}_{\text{old}})$$

Substituting $\hat{\theta}_{\text{old}} = \theta_{\text{old}} + \alpha v_{\text{old}}$

$$\begin{aligned}&= \theta_{\text{old}} + \alpha v_{\text{old}} + \alpha^2 v_{\text{old}} - (1+\alpha) \epsilon \nabla_{\theta} L(\hat{\theta}_{\text{old}}) \\&= \theta_{\text{old}} + (1+\alpha) \left[\alpha v_{\text{old}} - \epsilon \nabla_{\theta} L(\theta_{\text{old}} + \alpha v_{\text{old}}) \right]\end{aligned}$$

$$\hat{\theta}_{\text{new}} = \theta_{\text{old}} + (1+\alpha) v_{\text{new}}$$

$$\theta_{\text{new}} + \alpha v_{\text{new}} = \theta_{\text{old}} + (1+\alpha) v_{\text{new}}$$

$$\theta_{\text{new}} = \theta_{\text{old}} + v_{\text{new}}$$

which is same as vanilla nesterov momentum.

1 b) i)

False

¶

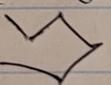
for $k=1$ training error = 0

1 b) ii)

Method 2
because tuning params to test will
not generalize for unseen data. So
tuning hyperparams should be done
on validation data.

1 c)

B, D

Increasing k will have decision boundary
with more sides like , not smooth

1 d) i)

(3) overfit the data

1 d) ii)

A, D,