Midterm Review Prof. J.C. Kao TAs: T. Monsoor, Y. Liu, S. Rajesh, L. Julakanti, K. Pang

1. Multiple choice (Shreyas)

Please pick the correct answers for each questions, note that each question can have one or more than one correct.

(a) Consider Figure 1 plotting loss values as a function of the number of epochs, select the option that best describe the shaded regions in the plot, and the point where you would stop training to achieve the best generalization.



Figure 1

- i. R1: Overfitting, R2: Underfitting, stop at a.
- ii. R1: Overfitting, R2: Underfitting, stop at b.
- iii. R1: Underfitting, R2: Overfitting, stop at b.
- iv. R1: Underfitting, R2: Overfitting, stop at c.
- (b) When we minimize the negative log likelihood for a classification problem with c classes, which of the following are we inherently performing?
 - i. Maximizing the likelihood of observing the training data.
 - ii. Minimizing the Mean Squared Error.
 - iii. Minimizing the Cross Entropy loss.
- (c) Mark all the correct choices regarding cross validation.
 - i. A 5-fold cross-validation approach results in 5-different model instances being fitted.

- ii. A 5-fold cross-validation approach results in 1 model instance being fitted over and over again 5 times.
- iii. A 5-fold cross-validation approach results in 5-different model instances being fitted over and over again 5 times.
- iv. None of the above.
- (d) Which of the following are considered as hyperparameter choices while training a neural network.
 - i. Loss Function.
 - ii. Learning Rate.
 - iii. Number of Layers.
 - iv. Batch Size.
 - v. All of the above.
- (e) Assuming Stochastic Gradient Descent (SGD) computes gradient using a single sample from the training data, which of the following statements are true.
 - i. Gradient computed using SGD will be noisier than gradient computed using Batch Gradient Descent.
 - ii. Empirically, SGD takes longer (in terms of clock time) to converge than Batch Gradient Descent.
 - iii. SGD usually avoids the trap of poor local minima.
 - iv. SGD is computationally more expensive than Batch Gradient Descent.

2. Short answer (Kaifeng)

- (a) Please explain the difference between batchnormalization during training and testing.
- (b) Your friend designed a novel activation function:

$$f(x) = x^3 \tag{1}$$

Please discuss if this is a good idea to use this activation in a neural network.

(c) Your friend is utilizing a Multi-layer Perceptron (MLP) for a deep learning task and is trying to increase the number of units within each layer to enhance the model's complexity. Please explain potential effect of this action on the model performance.

- (d) Please explain the role of ℓ_1 regularization.
- (e) Please explain the role of the bias correction step in the Adam optimizer.

3. Backpropagation in parallel neural network (Tonmoy)

A parallel neural network consists of twin networks which accept distinct inputs but share the same weights. The outputs of the twin networks are later processed by more hidden layers. Let's assume we have a parallel neural network with the following architecture:

$$h_{p} = W_{1}x_{p}^{(i)} + b_{1}$$

$$z_{1} = ReLU(h_{p})$$

$$h_{q} = W_{1}x_{q}^{(i)} + b_{1}$$

$$z_{2} = ReLU(h_{q})$$

$$z = z_{1} - z_{2}$$

$$z_{3} = W_{2}z + b_{2}$$

$$\hat{y}^{(i)} = \sigma(z_{3})$$

$$L^{(i)} = L_{CE}(y^{(i)}, \hat{y}^{(i)})$$

$$L = -\frac{1}{m}\sum_{i=1}^{m} L^{(i)}$$

In the above architecture, $(x_p^{(i)}, x_q^{(i)})$ represent the pair of i^{th} input example and are each of shape D_x . $y^{(i)}$ represent the label of the i^{th} input example and is a scalar. We also assume z_1 and z_2 have shape of D_z .

- (a) Draw the computational graph for the parallel neural network described above. You can start from $L^{(i)}$ as your output variable and then backtrack to the input variables $x_p^{(i)}$ and $x_q^{(i)}$.
- (b) Compute $\nabla_{\hat{y}^{(i)}} L^{(i)}$ and denote it as $\delta_{\hat{y}^{(i)}}$. For all the following parts, you can refer to this computed gradient as $\delta_{\hat{y}^{(i)}}$.
- (c) Compute $\nabla_{z_3} L^{(i)}$ and denote it as δ_{z_3} . For all the following parts, you can refer to this computed gradient as δ_{z_3} .
- (d) Compute $\nabla_{b_2} L^{(i)}$ and denote it as δ_{b_2} . For all the following parts, you can refer to this computed gradient as δ_{b_2} .
- (e) Compute $\nabla_{W_2} L^{(i)}$ and denote it as δ_{W_2} . For all the following parts, you can refer to this computed gradient as δ_{W_2} .
- (f) Compute $\nabla_z L^{(i)}$ and denote it as δ_z . For all the following parts, you can refer to this computed gradient as δ_z .
- (g) Compute $\nabla_{z_1} L^{(i)}$ and denote it as δ_{z_1} . For all the following parts, you can refer to this computed gradient as δ_{z_1} .
- (h) Compute $\nabla_{z_2} L^{(i)}$ and denote it as δ_{z_2} . For all the following parts, you can refer to this computed gradient as δ_{z_2} .
- (i) Compute $\nabla_{h_q} L^{(i)}$ and denote it as δ_{h_q} . For all the following parts, you can refer to this computed gradient as δ_{h_q} .
- (j) Compute $\nabla_{h_p} L^{(i)}$ and denote it as δ_{h_p} . For all the following parts, you can refer to this computed gradient as δ_{h_p} .
- (k) Compute $\nabla_{b_1} L^{(i)}$.

(l) Compute $\nabla_{W_1} L^{(i)}$.

4. Regularization techniques (Yang)

- (a) True or False: Regularization is intended to reduce training error but not validation error.
- (b) Consider a model $\hat{\mathcal{L}}(\theta) = \mathcal{L}(\theta; \mathbf{X}, \mathbf{y}) + \alpha \Omega(\theta)$ where $\mathcal{L}(\theta; \mathbf{X}, \mathbf{y})$ is some loss function and $\Omega(\theta)$ is some norm penalty. What are the effects on the model when $\alpha = 0$ and $\alpha \to \infty$?
- (c) Mathematically show that ℓ_2 regularization shrinks the weight in gradient descent. Hint: start with $\tilde{\mathcal{L}}(\theta; \mathbf{X}, \mathbf{y}) = \mathcal{L}(\theta; \mathbf{X}, \mathbf{y}) + \frac{\alpha}{2} ||\theta||_2^2$ and derive the gradient descent step for θ .
- (d) List two dataset augmentation techniques for image classification.
- (e) How did you implement dropout in homework 4? Please comment on both training and testing.

5. Optimization techniques (Lahari)

(a) In lecture, we have learnt about Nesterov Momentum and it's update rule for parameters. The update rule for parameter θ is given by:

$$v_t = \alpha v_{t-1} - \epsilon \nabla_{\theta} L(\theta_{t-1} + \alpha v_{t-1})$$

$$\theta_t = \theta_{t-1} + v_t$$
(2)

Prove that the update rule in (3) is equivalent to the update rule in (2)

$$v_{new} = \alpha v_{old} - \epsilon \nabla_{\tilde{\theta}_{old}} L(\theta_{old})$$

$$\tilde{\theta}_{new} = \tilde{\theta}_{old} + v_{new} + \alpha (v_{new} - v_{old})$$
(3)

Explain one advantage of using the update update rule in (3) over the update rule in (2).

- (b) Consider the two loss curves $L_1(x)$ and $L_2(x)$ shown in Figure 2. Which loss curve has a saddle point? Which loss curve has a poor local minima? In which of the loss curves, is an optimizer more likely to escape the trap of a saddle point or a poor local minima? And what property does the optimizer require for it to escape these traps in this case?
- (c) Consider the contour plot shown in Figure 3, where the loss surface is plotted with respect to just 2 weights w1 and w2, where $w1, w2 \in R$ (scalars). Assume you are given a hypothetical scenario, where you start from point A and use vanilla gradient descent in many iterations to get to point B. During this process, we have started accumulating momentum based on the following equation,

$$g_t = \nabla_{\theta_t} L(\theta_t)$$
$$v_t = v_{t-1} - \epsilon g_t$$

Comment on which direction does the weight update occur if we use the following optimizers : vanilla gradient descent, gradient descent with momentum, gradient descent with Nesterov momentum, Adagrad.



Figure 2: Loss curves $L_1(x)$ (Left), $L_2(x)$ (Right)



Figure 3: Contour plot of a Loss function L(w1, w2)

- (d) In the Gradient descent + momentum scheme, find a general expression of v_t in terms of gradients $g_1, g_2, ..., g_t$ and ϵ (learning rate), considering an initial value of momentum $v_0 = 0$.
- (e) Consider that the gradients $g_1, g_2, ..., g_t$ in part (d) are i.i.d. random variables with mean μ and variance σ . Find the expected value of weights θ_t at t = 3.

6. ℓ_{∞} regularization (Tonmoy)

Let $x \in \mathbb{R}^n$, then we define the ℓ_{∞} norm and the Log-Sum-Exponent (LSE) of it as follows:

$$\|x\|_{\infty} = \max_{i} |x_{i}|$$
$$LSE(x) = \ln\left(\sum_{i=1}^{n} e^{|x_{i}|}\right)$$

(a) Show that the following inequality holds for $n \ge 1$,

$$\|x\|_{\infty} \le LSE(x) \le \|x\|_{\infty} + \ln(n) \tag{4}$$

- (b) Is the lower bound in (4) strict for n > 1?
- (c) Under what condition on x, will the upper bound in (4) be satisfied with equality.
- (d) Use the result from (4) to show that the following inequality holds,

$$||x||_{\infty} \le \frac{1}{t} LSE(tx) \le ||x||_{\infty} + \frac{\ln(n)}{t}$$
 (5)

for some scaling constant t > 0.

b)
$$\mathcal{L}(\theta) = (\eta \overline{1}) p(\eta = 1 | \theta)^{3} p(\eta = 0 | \theta)^{1-3}$$

 $\mathcal{L}(\theta) = (\eta \overline{1}) p(\eta = 1 | \theta)^{3} p(\eta = 0 | \theta)^{1-3}$
labels = $\eta = 0$
labels = 0
labels = $\eta = 0$
labels = $\eta = 0$
labels = (\eta = 0)
labels = (\eta = 0)

BCE :
$$\lambda(0) = -\frac{1}{n} \frac{2}{3} y_{1} \log p(y_{2}, 10) + (1-y_{1}) \log p(y_{2}, 0)$$

5 fold cross volidation fits S diffuent models. Answer : (1)

(ک

d) Anything that is not learnt through grodient durcent is a hyperparameter. Answer: (V) Ally the above.

C) a single sample i) True, b'cos its estimate will take more II) False, SUD steps but converges foster f cloch time. in terms (ii) True, every sample is unlikely to have the same minima 04 the batch. False, Botch 6D requires (v)loading the entire data. Answers: c, cii//

Short Answers

(a) Please explain the difference between batchnormalization during training and testing.

Solution: During training, it first computes the mean and variance of the mini-batch data in the unit-wise. Then BatchNorm normalizes the layer's input based on the mean and variance. After the normalization, BatchNorm applies two learnable parameters for each unit: γ for scaling and β for shifting, which are learned during training and allow the network to adaptively adjust the output distribution. Meanwhile, it also keeps tracking the running averages for mean and variance through all batches.

During testing, BatchNorm uses the fixed accumulated running averages from training for normalizing the testing data. The learned scaling and shifting parameters γ and β are then applied to the normalized data.

(b) Your friend designed a novel activation function:

$$f(x) = x^3 \tag{1}$$

Please discuss if this is a good idea to use this activation in a neural network. **Solution:** This activation function is nonlinear and differentiable everywhere, which satisfy some requirements of a good activation function. However, it is likely to cause exploding gradients as the gradient can be very large for inputs with large absolute values. Also, small inputs can lead to vanishing gradients, e.g. inputs that are close to zero.

(c) Your friend is utilizing a Multi-layer Perceptron (MLP) for a deep learning task and is trying to increase the number of units within each layer to enhance the model's complexity. Please discuss potential effect of this action on the model performance.

Solution: (i) If the model is originally underfitting on the training data, adding more units in layers allows the MLP to capture more complex patterns in the data, which can improve the model performance, e.g. decrease both training and testing error. (ii) On the other hand, it may also cause overfitting as increasing the model's capacity is likely to make the model sensitive to the training data. As a result, the training loss may still keep decreasing while the testing error increase.

(d) Please explain the role of ℓ_1 regularization.

Solution: ℓ_1 regularization introduces a penalty which is equal to the sum of absolute values of the model weights. It encourages the model to optimize some of the feature weights to zero. This property allows the model to learn a simpler and sparser patterns by pushing less important feature weights to zero, which can help prevent overfitting. Further, by observing the trained weights, we can implement feature selection to simplify the model and save computation resource.

(e) Please explain the role of the bias correction step in the Adam optimizer. **Solutions:** Since Adam optimizer uses running averages to estimate the gradient and its square, these estimations are biased towards zero at the start of training because we initialize them to zero. Therefore, the optimizer is likely to take larger steps in the initial several updates of the model, leading to unstable training and slower convergence. The bias correction step adjusts these estimations to be more accurate.

After the early phase of training, the estimations tend to be accurate, so the bias correction factor will gradually approach 1, reducing the impact of bias correction.

```
4. Regularization techniques
   - lecture 9
  (a) Folde

f^{0} loss fr. [For norm penalty

(b) \tilde{\chi}(\theta) = \chi(\theta; \chi, y) + \kappa \Omega(\theta)
           ~= 0? L(0) = L(0; 3, 4)
           a→w? Ž(0) → a Ω(0)
                                 he learning of L(+; X, y)
  (c) \tilde{L}(\underline{0}; \underline{x}, \underline{y}) = L(\underline{0}; \underline{x}, \underline{y}) + \underline{x} \|\underline{0}\|_{2}^{2}
          e is a vector
         Ve L(e; x, y) - 赤 ん(e; x, y)
  ₹ 11012 > $ 0TO
                           2 0 ( L( 1 ; ×, y) + × 11 0 1 ,
2 = 0.0 = 40
                           = Vo 1(2) 2.3), a e
        € ← E - CVe ž (E; Ž, J)
          € 0 - ( [ Ve 1(2; ¥, y) + « 2)
          ←(1- for) e - s Ve ((e) ¥,y)
            weight decay
(d) fispping / mirrorsug - roundous crops
rointron, brightness
        (ons correction
                                        adjustments
        pooling (CAN)
      label smoothing
        Subsampling
        color space trans.
(c) inverted dropont
        training : generate binary mask
                      * values after affine transform
                     1 p e the prob. of keeping
       test. nothing
```

(85) (a) Given equation of nesteror momentum $V_{t} = KV_{t-1} - E \nabla_{0}L(\theta_{t-1} + KV_{t-1})$ $\Theta_{t} = \Theta_{t-1} + V_{t}$ (\mathcal{D}) To prove it is equivalent to: Vnew = XVord - EV Lload) Vnew = Ord + Vnew + X(Vnew - Vord) -(2) => Let us assume that new parameter S $\tilde{O} = O + \alpha V$ This is : Odd = Odd + X Vold | Odd = Odd - X Vold Ones = Ones + X Vhers | Ones = Ones - X Vnu => From (D, we have: Vnew = XVold - EV_0 K(Odd + XVold) $\nabla_{\Theta} \mathcal{L}(\tilde{\Theta}_{dd}) = \frac{\partial \mathcal{L}(\tilde{\Theta}_{dd})}{\partial \tilde{\Theta}} = \frac{\partial \tilde{\Theta}}{\partial \tilde{\Theta}} (\frac{\partial \mathcal{L}(\tilde{\Theta}_{dd})}{\partial \tilde{\Theta}})$ $= \frac{\partial \theta^{1}}{\partial \theta} = \frac{\partial (\theta^{1} + \theta^{2})}{\partial \theta} = \frac{\partial (\theta^{2} + \theta^{2})}{\partial \theta} = \frac{\partial (\theta^{2} + \theta^{2})}{\partial \theta}$ $\begin{bmatrix} 2 & -2 \\ -2$

So, we have, Vor L(Bord) = Vor L(Dord + dV ora) $V_{new} = K V_{old} - E \nabla_{\overline{o}} \chi [\overline{O}_{old}] \longrightarrow This is (2)$ (qualiton -> From (), we have Onew = Od + Vnew Oold - When = Oold - XVold + Vhen Ozd = Odd + Vnew + ox (Vnew - Vold) -> This is 2) equation Hence, we have shown that both $O \in Q$ are equivalent as we can get Q from O by a change in raviable. This representation helps us in implementation of nestoror momentum as this doem't require us to calculate gradient at a different value of 0+KV. (Also, both 0 and $\overline{0}$ start from the same value of parameter initialization as $0 = \overline{0}$ initially as V=0 at start _

(b) L(x) has a saddle point. The curve decreases on one side of the saddle and increases on other side of saddle point. The gradient at a saddle point is equal to zero. L2(2) has a poor local minima. This is a local minima as it is a minimum value of the function in it's surrounding until a certain limit in all the directions. The gradient is zero at a local minima. caddle point 10cal We need momentum to get out of saddle point or local minima because the gradient becomes zero at these points. It is also more likely to escape the saddle point than the local minima because after crossing the local minine, the momentum decreases due to gradient being in opposite direction to the momentum accumulated benotes down the slope.

(C) · A ' 22 400 \mathcal{W}_{1} From the above contour plot, we can make some initial observations: (1) The gradient along W2 direction is higher than the gradient along w, direction. This is because the contour lines along We directions are very close to each other which indicates a steep curve in we direction. The contour lines along w, direction are further apart and hence has slower descent. (or low gradient) Made with Goodr

(2) Vanilla gradient descent is in the direction perpendicular to a contour line. À B (3) Sgd+ momentum A veg csgd B momentum KJ 47

(A) sgd+ nesteror momentum. A B Sigd + MP In this, first we take a step along the momentum and then calculate the gradient at that point in graph. (5) Adagrad. The update equation: [a ccumulates] gradient] $\alpha \leftarrow \alpha + 909$ $0 \leftarrow 0 + \underbrace{e}_{\sqrt{\alpha+\nu}} 0 q$ If gradient is accumulated more in a direction, then the step <u>e</u> is len in that direction.

Made with Goodnotes

If gradient accumulated is len in a direction, then the step $\frac{E}{Vatu}$ is more in that direction

From plot, we can see that when it travevsed from A to B, the gradients accumulated along W₂ direction is more and gradient along w₁ direction is lener than it.

Hence, the step along w, direction is more and step along we direction is len.



(d) Sgd + momentum update:

$$V_{t} = KV_{t-1} - \xi g_{t}$$

$$\Theta_{t} = \Theta_{t-1} + V_{t}$$

$$\rightarrow V_{0} = 0 \quad [Given]$$

$$V_{1} = KV_{0} - \xi g_{1} = -\xi g_{1}$$

$$V_{2} = KV_{1} - \xi g_{2} = \kappa(-\xi g_{1}) - \xi g_{2} = -\xi(g_{2} + Ng_{1})$$

$$V_{2} = \kappa V_{2} - \xi g_{2} = \kappa(-\xi(g_{1} + \kappa g_{1})) - \xi g_{3}$$

$$= -\xi(g_{3} + \kappa g_{2} + \kappa^{2}g_{1})$$

$$\vdots$$
From thus, we can observe that

$$\Rightarrow V_{t} = -\xi(g_{1} + \kappa g_{t-1} + \kappa^{2}g_{1} + \xi^{-1}g_{1})$$

$$E(V_{t}) = E[-\xi(g_{t} + \kappa g_{t-1} + \kappa^{2}g_{1} + \xi^{-1}g_{1})]$$

$$= -\xi[E(g_{t}) + \kappa E(g_{t-1} + \kappa^{2}g_{t-2} + \dots + \kappa^{t-1}\xi(g_{1})]$$

$$(e) \quad \Theta_{0} = \Theta_{0}$$

$$\Theta_{1} = \Theta_{0} + V_{1}$$

$$\Theta_{2} = \Theta_{1} + V_{2} = \Theta_{0} + V_{1} + V_{2}$$

$$\vdots$$
mean contents 0, for general $\Theta_{t} = \Theta_{0} + V_{1} + V_{2} + \dots + V_{t}$

 $E(0_3) = E[0_0 + v_1 + v_2 + v_3]$ $= \mathbb{E}(\Theta) + \mathbb{E}(v_1) + \mathbb{E}(v_2) + \mathbb{E}(v_3)$ Griven, g1, g2..., gt are i.i.d over with mean 11 and var of? =) From part (d), we can $\mathbb{E}(V_1) = \mathbb{E}(-2g_1) = -2\mathbb{E}(g_1) = -2\mathbb{E}(g_1) = -2\mathbb{E}(g_1)$ $\mathbb{E}(v_2) = \mathbb{E}\left(-\mathbb{E}(g_2 + \alpha g_1)\right) = -\mathbb{E}\left[\mathcal{M} + \alpha \mathcal{M}\right]$ = - E N [1+ N] $E[v_2] = E(-\epsilon(g_2 + \alpha g_2 + \alpha^2 g_1))$ $= -\varepsilon \left[M + \alpha M + \kappa^2 M \right]$ $= - \varepsilon M (1 + \alpha + \alpha^2)$ $\mathbb{E}(0_3) = \mathbb{E}(0) - \mathbb{E}\left(1 + (1 + \alpha) + (1 + \alpha + \alpha^2)\right)$ * geometric series with factor is and first term a is $a + \alpha r + \alpha r^{n-1} = \frac{\alpha (1-r^n)}{1-r}$

Made with Goodnotes

Alternately, $\mathbb{E}(V_t) = -\mathcal{E}[\mathcal{U} + \mathcal{K}\mathcal{U} + \dots + \mathcal{K}\mathcal{U}]$ = $-\mathcal{E}\mathcal{U}[1 + \mathcal{K} + \dots + \mathcal{K}^{t-1}]$ $= -2M\left[\frac{1-x^{t}}{1-x}\right]$ $E(0_{2}) = E(0_{0}) + \sum_{t=1}^{3} \left[-\varepsilon u \left[1 - \alpha^{t} \right]^{2} \right]$ $= \mathbb{E}(Q_{D}) - \frac{\mathcal{E}\mathcal{M}}{1-\alpha} \left(\sum_{t=1}^{2} (1-\alpha^{t})\right)$ $= \mathbb{E}(\Theta_0) - \mathbb{E}\mathcal{N}\left[3 - \frac{1}{2}\mathbf{x}_{+}\right]$ $= \mathbb{E}(\varphi_{0}) - \underbrace{\mathbb{E}}_{1-\alpha} \left[3 - \left(\frac{1-\alpha}{1-\alpha} \right) \right]$

Xal We have drawn the computational graph abore. Now, we will use backpropagation to compute the lerivatives reavined for learning the parameters through gradient descent. ** Always write dimensions $X_{p}^{(i)} \in \mathbb{R}^{p_{X}}, X_{a}^{(i)} \in \mathbb{R}^{p_{X}}, \hat{\mathcal{Y}}^{(i)} \in \mathbb{R}$ $W_1 \in \mathbb{R}^{D_2 \times D_X}$, $b_1 \in \mathbb{R}^{D_2}$ WZEIR'XDZ, bZER

b) from the computational graph,

$$L^{(i)} = L_{CE} (y^{(i)}, f^{(i)})$$

$$= y^{(i)} \log (f^{(i)}) + (1 - y^{(i)}) \log (1 - f^{(i)})$$

$$= y^{(i)} \log (g^{(i)}) + (1 - y^{(i)}) \log (1 - f^{(i)})$$

$$= y^{(i)} \log (1 - f^{(i)}) + (1 - y^{(i)}) \log (1 - f^{(i)})$$

$$= y^{(i)} - \frac{1}{1 - g^{(i)}} + \frac{1}{1 - g^{(i)}}$$

$$= y^{(i)} - \frac{1}{1 - g^{(i)}} + \frac{y^{(i)}}{1 - g^{(i)}}$$

$$= y^{(i)} - \frac{1 - y^{(i)}}{1 - g^{(i)}}$$





D) from the computational graph, backpropating through D gate $\frac{dL^{(i)}}{dh} = \frac{dL}{dz_3} = \frac{Sz_3}{z_3}$



So, by Chain rule $\frac{dL^{(i)}}{dw_2} = \frac{dm}{dw_2} \frac{dL^{(i)}}{dm}$ $S_{W_2} = S_{Z_3} Z^T$



h) Since we backpropagate - terminal of (F) through gate for Zz, So $\frac{dL''}{dz_2} = -\frac{dL^{(i)}}{dz}$ dl(i) $\delta_{2_1} = -\delta_z$ i) from the Computational graph, Z2 = ReLU(har) From Discussion, we know that back propagating through ReLU(.) leads to a hadamard product, 50

 $\frac{dL^{(i)}}{dha} = \Pi(h_{q}>0) \odot \frac{dL^{(i)}}{dz_2}$ $S_{hay} = \prod (hay > a) \odot \delta_{22}$ j) Using the same logic as i) we have Shp= I(hp>0) · SZI K) By law of total Derivatives, $\frac{dL^{(i)}}{db_{i}} = \frac{dL^{(i)}}{dhp} + \frac{dL^{(i)}}{dhar}$ $S_{5} = S_{hp} + S_{har}$







6 a) Suppose, $P = \max_{i} \left\{ \sum_{j=1}^{n} |X_{i}| \right\} \left[\sum_{j=1}^{n} |X_{i}| \right]$ Then, we have the following lover bound $e^{P} \leq \frac{2}{i=1} e^{i \times e^{i}}$ We also have the following upper band ŜelXil≤ nep i=1 Combining the upper and lover bounds we have

$$e^{P} \leq \sum_{i=1}^{n} e^{ix_{i1}} \leq ne^{P}$$

Taking natural logarithm of the above
Ineavaility we have
 $\ln(e^{P}) \leq \ln\left(\frac{2}{1-1}e^{ix_{i1}}\right) \leq \ln(ne^{P})$
 $\ln(e^{P}) \leq \ln\left(\frac{2}{1-1}e^{ix_{i1}}\right) \leq \ln(n)^{+}$
 $p \ln e \leq \ln\left(\frac{2}{1-1}e^{ix_{i1}}\right) \leq \ln(n)^{+}$
 $p \ln(e)$
Since $P = 11 \times 11_{60}$, so we have the
reavired in equality
 $11 \times 11_{60} \leq LSE(X) \leq 11 \times 11_{60} + \ln(n)$

5) Clearly for n>1, eP< zelxil Taking natural logarithm of both $ln(e^{p}) \leq ln\left(\sum_{i=1}^{n} e^{|Xi|}\right)$ sides, $Pln(e) \leq ln\left(\frac{\hat{s}e^{|Xi|}}{\hat{s}=1}\right)$ $||X||_{\infty} \leq LSE(X)$

c) Suppose,

$$|Xi| = |Xi|$$
for all i, i \in n . Then we have
$$\begin{array}{c}
 2 e^{|Xi|} \\
 = & 2 e^{P} \\
 i \in i
 = & ne^{P} \\
 Idence, the upper bound will
 be satisfied with equality
 LSE(X) = ||X||_{o} + ln(n)$$

a) Let $t \neq 0$ be some scaling Constant. Substituting X with tx in (2), we get $||tx||_{\infty} \leq LSE(tx) \leq ||tx||_{\infty} + ln(n)$

Now, $||tx||_{\infty} = t||x||_{\infty}$

So, $t||X||_{\infty} \leq LSE(tX) \leq t||X||_{\infty} + ln(n)$ $t||X||_{\infty} \leq t$, we get the required t_{1} inequality, $1|X||_{\infty} \leq t$, $LSE(tX) \leq ||X||_{\infty} + ln(n)$ t