

# Theoretical Questions.

## LOSS

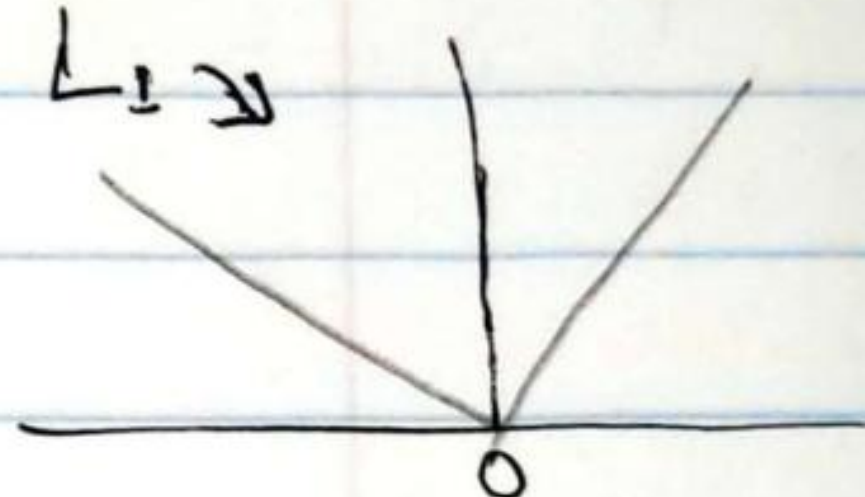
Name → Arinjay Jain

A.no. → A20447307

Email → ajain80@hawk.iit.edu.

1. Mean Absolute Error,  $L_1$  loss → sum of absolute differences between our target and predicted variables.

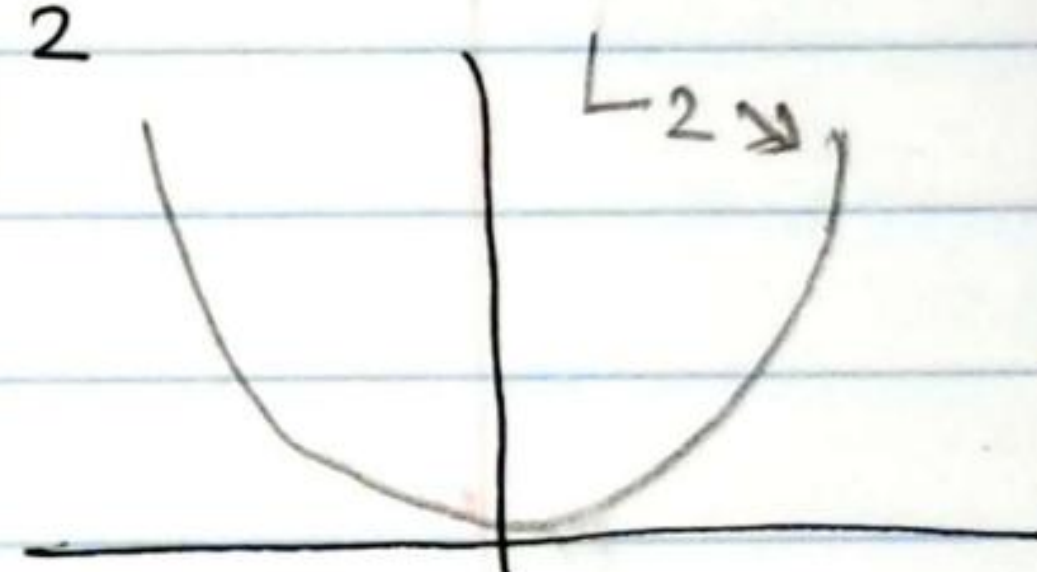
$$L_1: L(\theta) = \frac{\sum_{j=1}^K |\hat{y}_j^{(i)} - y_j^{(i)}|}{n}$$



2 Mean Square Error, Quadratic Loss,  $L_2$  loss.

Commonly used regression loss function.

$$L_2 \equiv L(\theta) = \frac{\sum_{j=1}^K (y_j^{(i)} - \hat{y}_j^{(i)})^2}{n}$$



$L_2$  loss vs  $L_1$  loss.

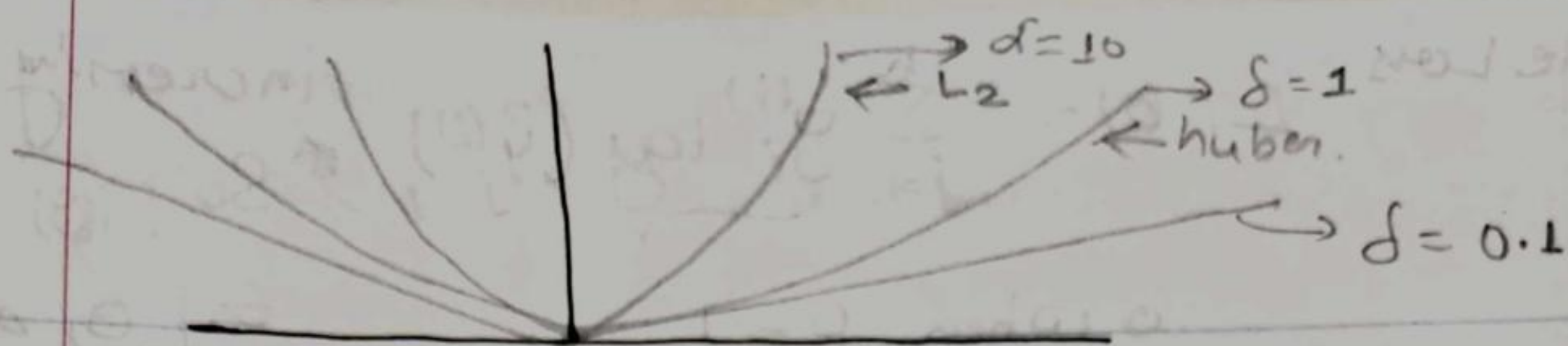
Using the squared error is easier to solve when we have small difference in predictions value ( $\hat{y}$ ) and observed value ( $y$ ). But using ( $L_1$ ) the absolute error is more robust to outliers.

3. Huber Loss, Smooth mean absolute error.

Huber loss is less sensitive to outliers in data than the squared error loss. It's also differentiable at 0. It's basically absolute error which becomes quadratic when error is small. How small that error has to be to make it quadratic depends on a hyperparameter,  $\delta$  (delta), which can be tuned. Huber loss approaches  $L_1$  when  $\delta \rightarrow 0$  and  $L_2$  when  $\delta \rightarrow \infty$  (large  $\delta$ ).

$$L_\delta(y, \hat{y}) = \begin{cases} \frac{1}{2} (y - \hat{y})^2 & \text{for } |y - \hat{y}| \leq \delta \\ \delta |y - \hat{y}| - \frac{1}{2} \delta^2 & \text{otherwise.} \end{cases}$$



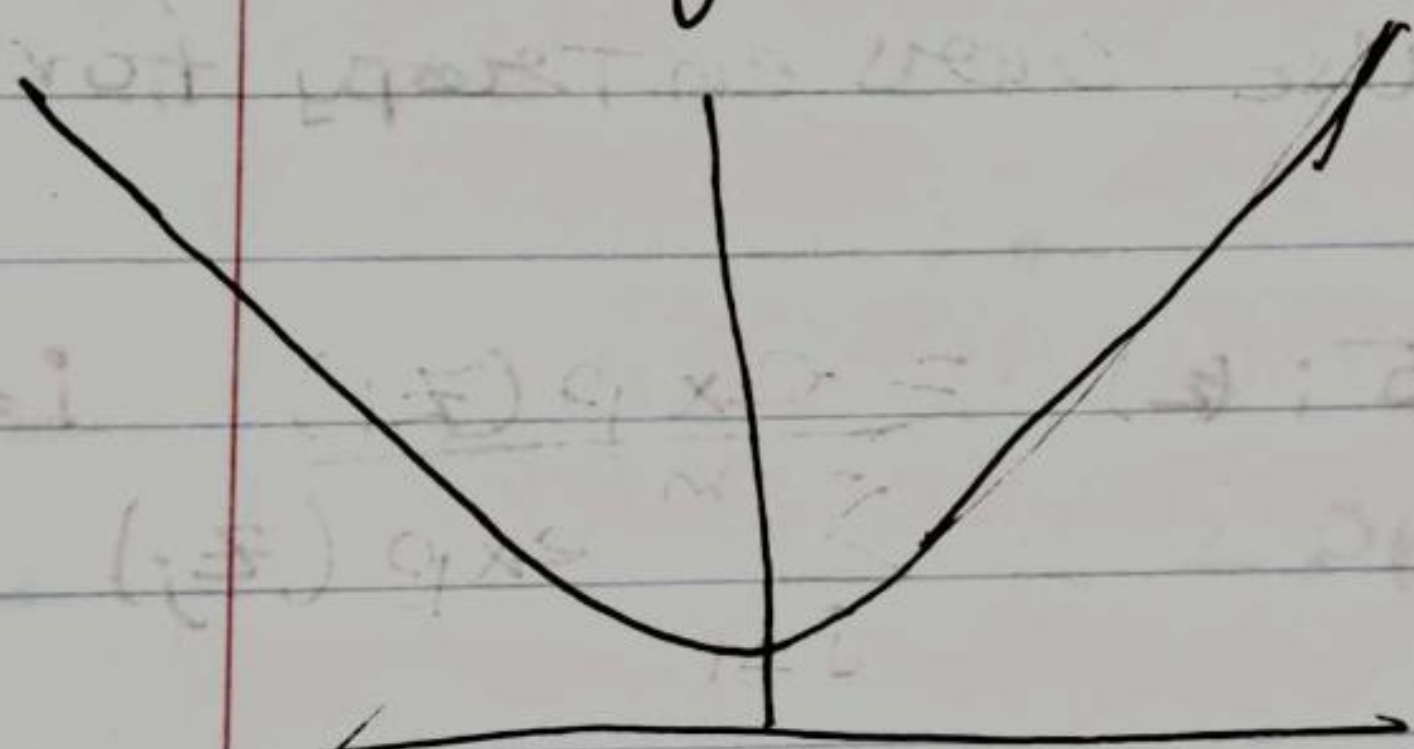


④ Log-Cosh Loss  $\Rightarrow \kappa$

$$L(\theta) = \sum_{i=1}^n \log(\cosh(\hat{y}_i - y_i))$$

Advantage:-  $\log(\cosh(x))$  is approximately equal to  $\frac{x^2}{2}$  for small  $x$  and to  $|x| - \log(2)$  for large  $x$ . It has all the advantages of huber loss and is twice differentiable everywhere unlike Huber loss.

$$\log(\cosh(x)) \approx \begin{cases} x^2/2 & \text{if } x \text{ is small} \\ |x| - \log(2) & \text{otherwise,} \end{cases}$$



Q2

Cross entropy loss  $\Rightarrow$

$$y_j^{(i)} = P(y=j | x^{(i)}) \quad j \in [1, \kappa]$$

We know the likelihood:  $L(\theta) = \prod_{i=1}^m \prod_{j=1}^{\kappa} (P(y=j | x^{(i)}))^{y_j^{(i)}}$   
to derive cross entropy loss taking the "log" of  $L(\theta) \rightarrow \ell(\theta) \leftarrow \log\text{-likelihood}$

$$\begin{aligned} \ell(\theta) &= -\log L(\theta) = -\sum_{i=1}^m \sum_{j=1}^{\kappa} y_j^{(i)} \log(P(y=j | x^{(i)})) \\ &= -\sum_{i=1}^m \sum_{j=1}^{\kappa} y_j^{(i)} \log(\hat{y}_j^{(i)}) \end{aligned}$$



Simple Loss  $L_i(\theta) = - \sum_{j=1}^k y_j^{(i)} \log(\hat{y}_j^{(i)})$  increasing function  
So  $L_i(\theta) \in [0, \infty]$   
 $0$  when  $p \rightarrow 1$   
 $\infty$  when  $p \rightarrow 0$ .

# for the worst case  $\Rightarrow$  Random class assignment  
 Probability  $P(y=j|x^{(i)}) = 1/k \rightarrow$  for equal classes.

$$L(\theta) = - \log(P(y=j|x^{(i)}))$$

$$= - \log(1/k) = \log(k)$$

$$\boxed{L(\theta) = \log(k)} \leftarrow \text{bad logs value.}$$

Q.3

Softmax loss.

is a combination of Softmax as activation in the output layer and use cross entropy for loss.

equation  $\Rightarrow \sigma_i(z) = \frac{\exp(z_i)}{\sum_{j=1}^m \exp(z_j)} \quad i=1, \dots, m$

$$Z_i = W_i \cdot x + B_i$$

Softmax can make  $Z_i$  nonnegative by letting them become exponential, then the sum of all items is normalized, now each  $O_i = \sigma_i(z)$  can be interpreted as the probability of data  $x$  belong to the Category  $i$  or the likelihood.

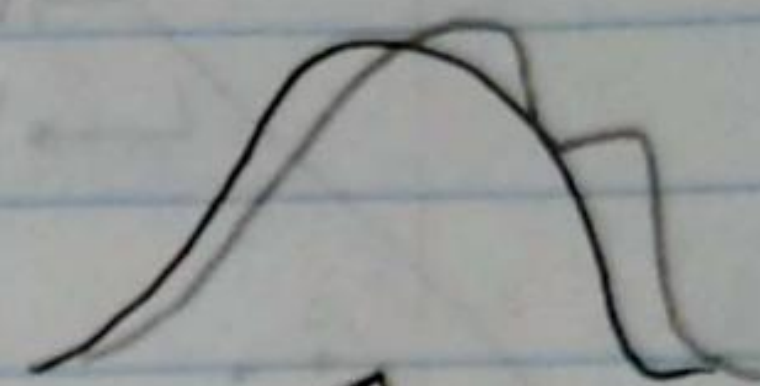
\* In Sigmoid we have probabilities of every output like every  $k$ -class outputs. but so in this case will get greater the 1 sum of all probabilities, but if we use softmax will get the probability in  $(0 \leftrightarrow 1)$ .



## Q4 Kullback-Liebler loss $\Rightarrow$

— Kullback-Liebler loss divergen measures similarity between distribution.

$$L(\theta) = - \sum \sum y^{(i)} \log \left( \frac{y^{(i)}}{\hat{y}^{(i)}} \right)$$



↑ find the similarity.

for single loss

$$L_i(\theta) = - \sum_{j=1}^K y_j^{(i)} \log \left( \frac{y_j^{(i)}}{\hat{y}_j^{(i)}} \right)$$

$\Rightarrow$  Determin the likelihood ratio for entire data set.

$$\frac{P(x_1, \dots, x_m)}{Q(x_1, \dots, x_m)} = \prod_{i=1}^m \frac{P(x_i)}{Q(x_i)}$$

$> 1 \rightarrow P(x)$  better model

$< 1 \rightarrow Q(x)$  better model.

taking log to change  $\pi \rightarrow \Sigma$

$$\Sigma \log \left( \frac{P(x_i)}{Q(x_i)} \right)$$

$> 0 \rightarrow P(x)$  better.

$< 0 \rightarrow Q(x)$  better.

Expected value to compute

$$E_{x \sim p} \left[ \log \frac{P(x_i)}{Q(x_i)} \right] = \int_{-\infty}^{\infty} P(x) \log \left( \frac{P(x)}{Q(x)} \right) dx$$

$$\approx \Sigma P(x_i) \log \left( \frac{P(x_i)}{Q(x_i)} \right) = KL(P||Q)$$

$$= \Sigma P(x_i) \log P(x_i) - \Sigma P(x_i) \log Q(x_i)$$

— entropy

cross entropy.

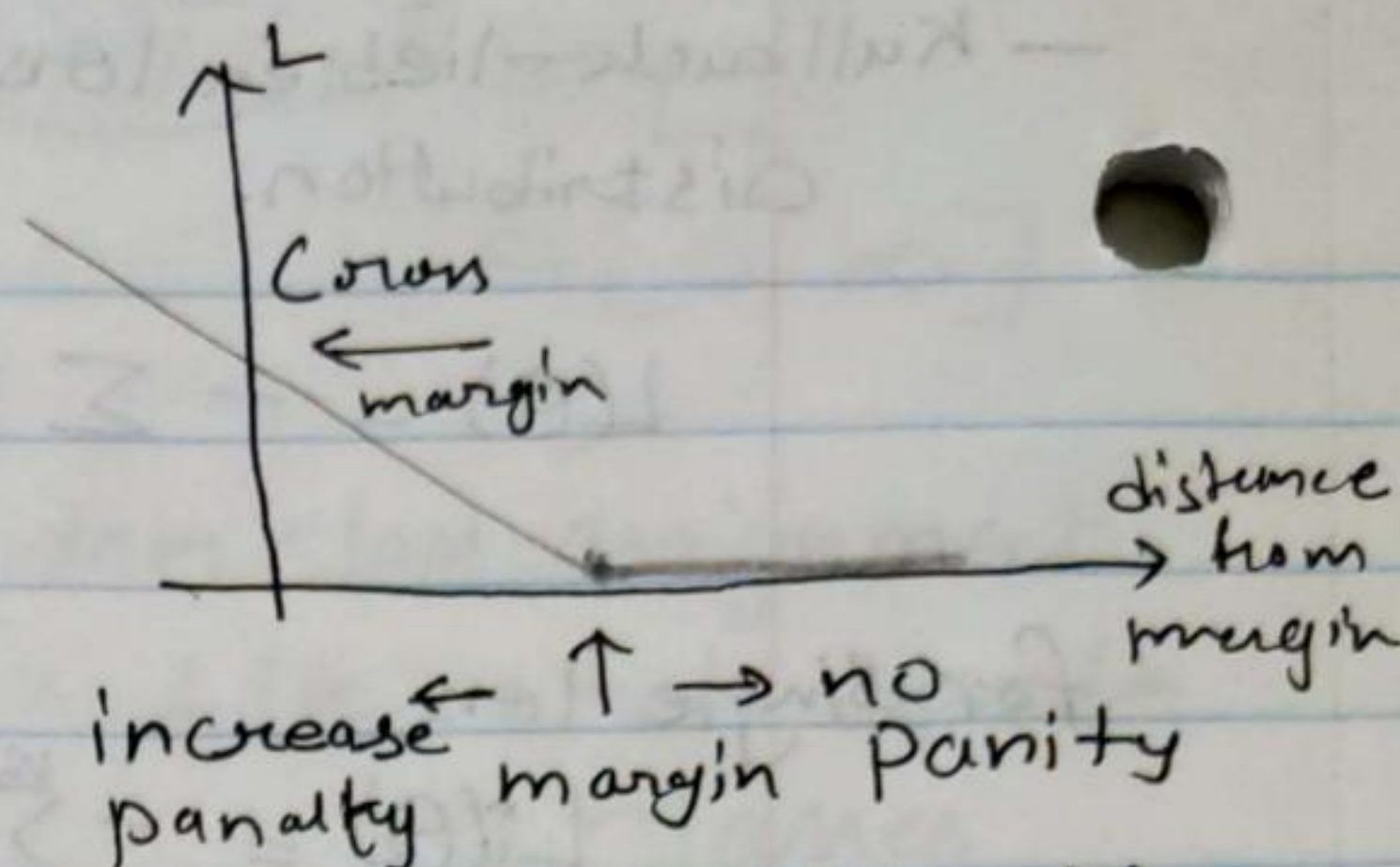
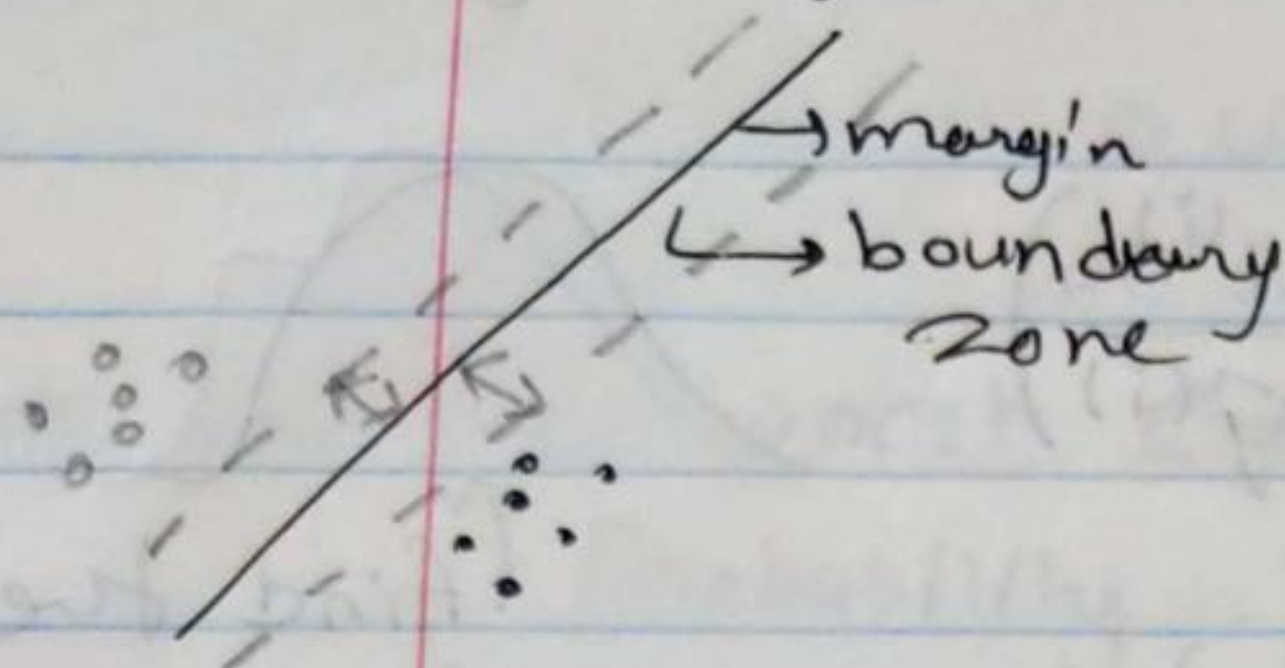
in this case  $P(x_i) = y^{(i)}$  and  $Q(x_i) = \hat{y}^{(i)}$

if  $y^{(i)}$  is fixed then  $\text{mean}(-\text{entropy})$  is not changing

In this case Kullback-Liebler Loss is same as cross entropy loss.



Q5  $\Rightarrow$  Hinge loss  $\Rightarrow$  SVM (Support vector machine) take margin of SVM.



Hinge Loss

$$L_i(\theta) = \max(0, 1 - y^{(i)} \hat{y}_t^{(i)})$$

positive if prediction agree in sign with margin  $< 1$  or disagree in sign  
negative if predictions agree in sign with margin  $> 1$

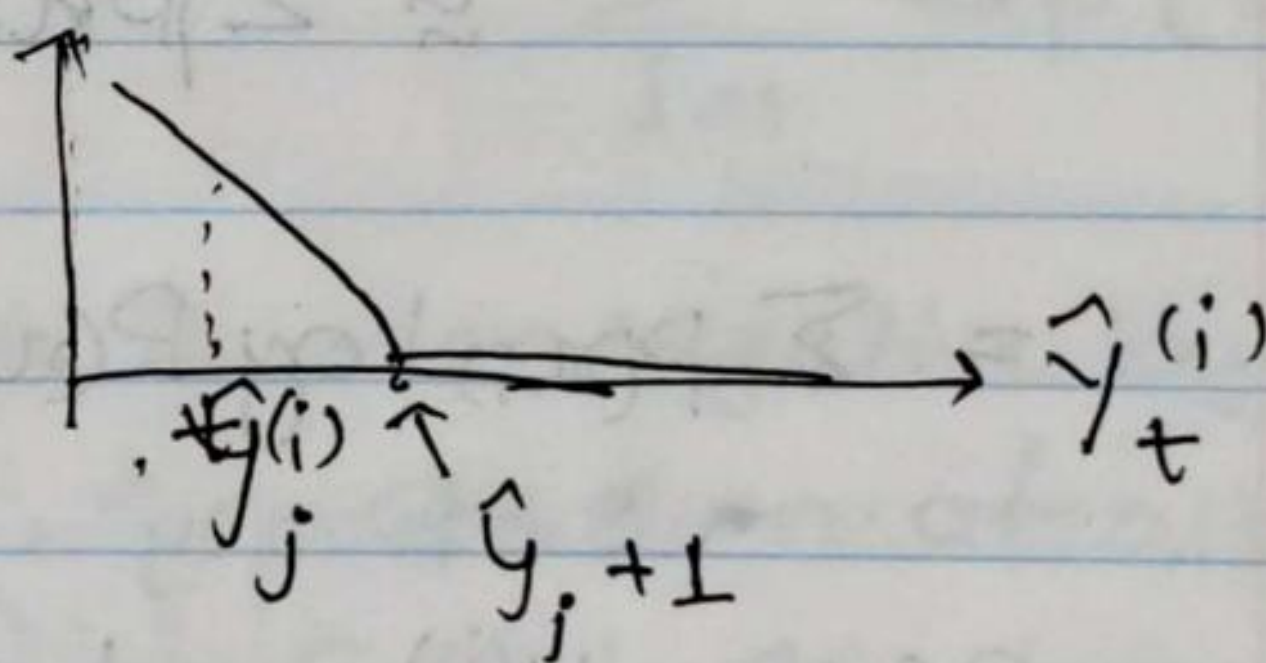
in distance form  $\hat{y}$  - distance from discriminant function

$$L_i = \sum_{j \neq t} \max(0, \hat{y}_j^{(i)} - \hat{y}_t^{(i)} + 1)$$

Sum over incorrect class labels

positive if  $\hat{y}_t^{(i)} < \hat{y}_j^{(i)} + 1$  (incorrect)  
negative if  $\hat{y}_t^{(i)} > \hat{y}_j^{(i)} + 1$  (correct)

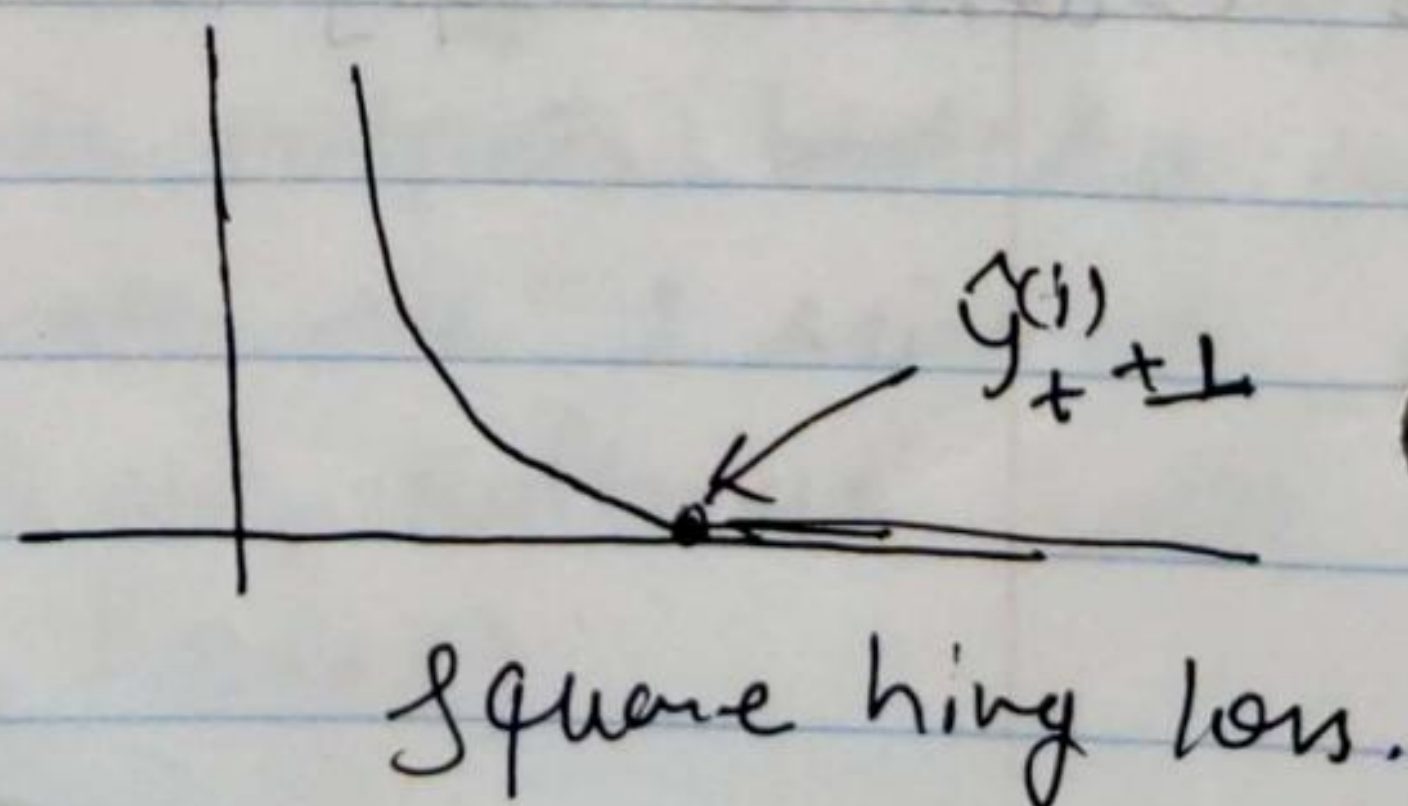
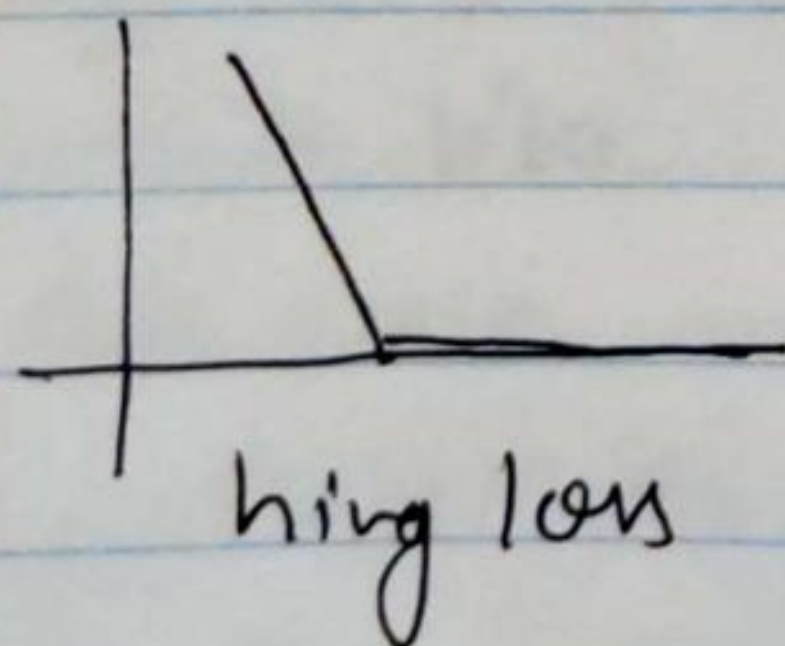
$$L_i(\theta) = \sum_{j \neq y^{(i)}} \max(0, \hat{y}_j^{(i)} - \hat{y}_t^{(i)} + 1)$$



Squared hinge loss

$\hookrightarrow$  Just a square of the elements.

$$L_i(\theta) = \frac{1}{2} \sum_{j=1}^K \max(0, \hat{y}_j^{(i)} - \hat{y}_t^{(i)} + 1)^2$$





the worst value or low random score before learning.

when  $\hat{y}_i \leq 0$  so all  $L_i = \max(0, 0 - 0 + 1) \dots \max(0, 0 - 0 + 1)$

#K  $\rightarrow$  no of classes.  $= \underline{K-1}$   $\leftarrow$  bad value.

Q6  $\Rightarrow y_1 = 1, y_2 = 2, y_3 = 3$

$\hat{y}^1 = (0.5, 0.4, 0.3)$   $\hat{y}^{(2)} = (1.3, 0.8, -0.6)$

$\hat{y}^3 = (1.4, -0.4, 2.7)$

$x_1 \ x_2 \ x_3$   
 $\hat{y}_1 \ 0.5 \ 0.4 \ 0.3$

$\hat{y}_2 \ 1.3 \ 0.8 \ -0.6$

$\hat{y}_3 \ 1.4 \ -0.4 \ 2.7$

$1 \ 2 \ 3$

we assuming that  $x_1$  belong to class 1

$x_2$  belong to class 2,  $x_3$  belong to class 3.

$x_3$  will be the best because it has large

distance not only from decision boundary but also from other values.

$$L_i = \sum_{j \neq t} \max(0, \hat{y}_j - \hat{y}_t + 1)$$

$$L_1 = \max(0, 0.5 - 1.4 + 1) + \max(0, 1.3 - 1.4 + 1) = 0.1 + 0.9 = 1$$

$$L_2 = \max(0, 0.4 - 0.8 + 1) + \max(0, -0.4 - 0.8 + 1) = 0.6$$

$$L_3 = \max(0, 0.3 - 2.7 + 1) + \max(0, -0.6 - 2.7 + 1) = 0 + 0 = 0$$

For  $x_3$  we found best hinge loss value.



Q7

By adding regularization term in loss it become low weight means, or its called weight decay.

Smaller coefficients means more stable solution that will generalize better our model.

$$L(\theta) = \underbrace{\frac{1}{m} \sum_{i=1}^m L_i(f(x_i, \omega), y_i)}_{\text{data loss}} + \underbrace{\lambda R(\theta)}_{\substack{\text{weight of regularization} \\ \text{(hyper parameter)}}} \quad \text{Regulation term.}$$

$L_1$  Regularization :  $R(\theta) = \sum_{ij} |\theta_{ij}|$  absolute.

$L_2$  Regularization :  $R(\theta) = \sum_{ij} (\theta_{ij})^2$

\*  $L_1$  makes weights sparse (concentrate weights), more likely to get some units are 1 and some 0 basically it's shutdown some coefficient.

eg  $\Rightarrow (0.5) + (0.5) = 1 + 0$

\*  $L_2$  Can be drive it using try to make them smaller and spreading them equally e.g.  $0.5^2 + 0.5^2 < 1^2 + 0^2$ .

Can be drive it using MAP inference using a Gaussian prior for  $\omega$ .

So will choose the  $L_2$  because it not shut the coefficients, it divide the weights equally over the model/network, better way to minimize the loss value.



Q8 ⇒

$$L_1: R(\theta) = \sum_{i,j} |\theta_{ij}| \quad \text{On gradient descent } \frac{\partial L}{\partial \theta} = \frac{\partial R(\theta)}{\partial \theta} \\ = \text{Sign}(\theta)$$

during gradient descent subtract  $\lambda (\text{Sign}(\theta))$   
Just add or subtract " $\lambda$ "

$$L_2: R(\theta) = \sum \theta^2 = \theta^T \theta \Rightarrow \frac{\partial R(\theta)}{\partial \theta} = 2\theta$$

during gradient descent subtract " $\lambda \theta$ "

$$\frac{\partial L}{\partial \theta} = \frac{\partial}{\partial \theta} \frac{1}{m} \sum_{i=1}^m L_i(\underbrace{f(x^{(i)}, \theta)}_{\hat{y}}, y^{(i)}) + \lambda \underbrace{\frac{\partial R(\theta)}{\partial \theta}}_{\substack{\left\{ \begin{array}{l} L_2 \\ = 2\theta \end{array} \right\} \quad \left\{ \begin{array}{l} L_1 = \lambda \text{Sign}(\theta) \end{array} \right\}}}$$

Q9 Kernel, bias and activity regularization ⇒

Kernel regularization: regularize  $w \rightarrow$  reduce  $w$

Bias regularization: regularize  $b \rightarrow$  reduce  $b$

activity regularization: regularize  $\hat{y} \rightarrow$  reduce  $\hat{y}$   
(including  $w$  and  $b$ )

- ① majorly we use kernel regularization because it's  
we multiply the input, and not depends on input.
- ② activity ( $\hat{y}$ ) is also effective but it's depends on  
of input. (i.e.  $\hat{y} = w\underline{x} + b$ .) Implicitly we cover  
kernel ( $w$ ) and bias ( $b$ ) with this.
- ③ Some time bias regularization if function expected  
to output small values.