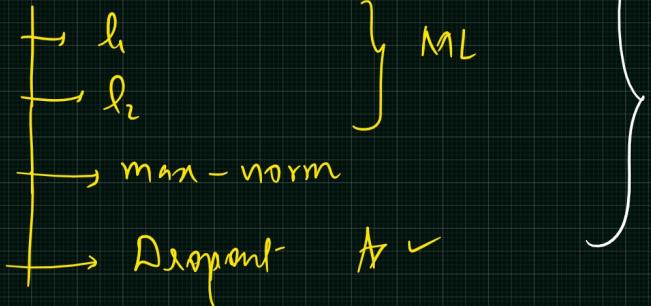


AGENDA

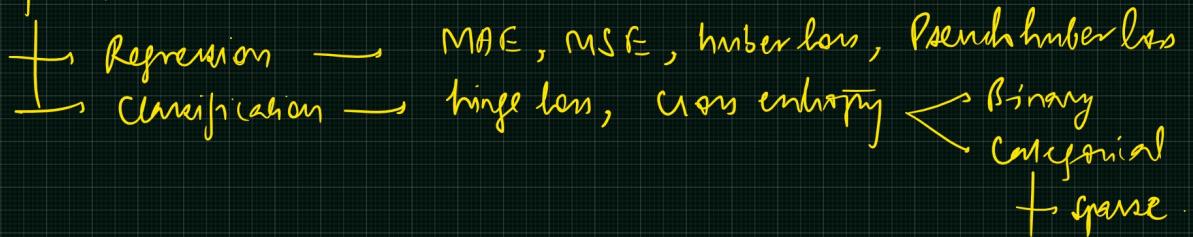
1) Regularization



Why?

To avoid overfitting

2) Loss functions :-



Linear Regression

$$y = \theta_0 \underline{1} + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n = \Theta^T \underline{x} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix}^T \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

↓
parameters
gim

total parameters

$$\frac{n}{n+1} \rightarrow \theta_1 \rightarrow \theta_n$$

$$1 \rightarrow \theta_0$$

$n+1 \uparrow \Rightarrow$ overfitting
 so when \rightarrow large amount which should contain almost all kind of pattern.

1) L_1 regularization | LASSO → Least Absolute Shrinkage & Selection operator

$$J_n(\theta) = J_0(\theta) + \alpha \sum_{i=1}^m |\theta_i|$$

$\Rightarrow m \rightarrow 1 \rightarrow m$
 ⇒ Not taking bias

$\underbrace{\qquad\qquad\qquad}_{L_1 \text{ term}}$ regularization factor $\in [0, 1]$

at some time during the training

$$J_0(\theta) = 0 + \text{noise} \Rightarrow \text{error} \neq 0$$

↓
weight update
 $w = w - \gamma \frac{\partial C}{\partial w}$

while training,

$$J_n(\theta) = J_0(\theta) + L_1 \Rightarrow w^+$$

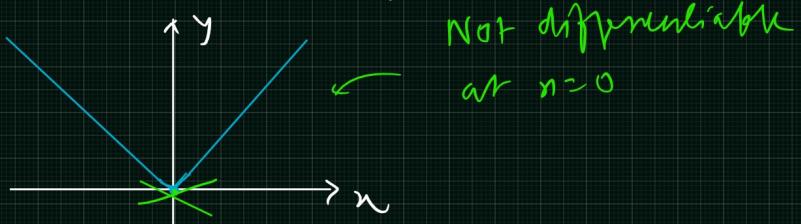
while testing,

$$J_n(\theta) = \tilde{J}_0(\theta)$$

Drawback:-

- i> Latency due to extra computation ✓
- ii> $J_n(\theta)$ is not differentiable everywhere

$$L_1 \text{ term} = \alpha \sum_{i=1}^m |\theta_i|$$



L_2 regularization | Ridge

$$J_n(\theta) = J_0(\theta) + \frac{\alpha}{2} \sum_{i=1}^m (\theta_i)^2$$

θ_i
↓
 $\frac{\partial \theta_i}{\partial x}$

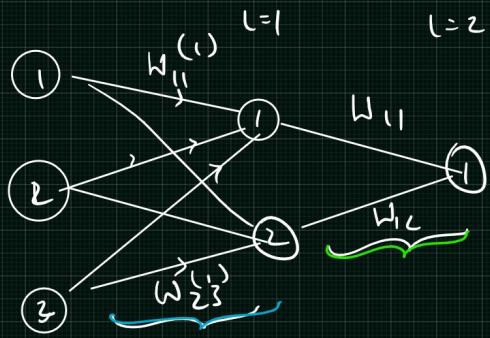
$\underbrace{\hspace{1cm}}$ L_2 term

$$\text{Training} \rightarrow J_n(\theta) = J_0(\theta) + L_2$$

$$\text{Testing} \rightarrow \tilde{J}_n(\theta) = J_0(\theta)$$

Drawback:-

- 1> Latency due to extra computation.
- 2> α as an extra parameter has to be tuned.



$$L_1 = \alpha \left\{ |w_{11}| + \dots + |w_{13}| \right\} + \left\{ |w_{21}| + |w_{22}| \right\}$$

$$J_n(\theta) = J_0(\theta) + L_1$$

$$L_2 = \frac{\alpha}{2} \left\{ (w_{11})^2 + \dots + (w_{23})^2 + (w_{11}^{(2)})^2 + (w_{12}^{(2)})^2 \right\}$$

$\Rightarrow L_1, L_2$ regularization

$$J_n(\theta) = J_0(\theta) + r \underbrace{\alpha \sum_{i=1}^m |\theta_i|}_{L_1} + (1-r) \underbrace{\frac{\alpha}{2} \sum_{i=1}^m \theta_i^2}_{L_2}$$

$r \rightarrow$ controlling factor $r \in [0, 1]$

$r=1 \Rightarrow L_1$ is dominating $J_n(\theta) \rightarrow$ Lasso

$r=0 \Rightarrow L_2$ is dominating $J_n(\theta) \rightarrow$ Ridge

$$r = (0, 1)$$

presence of outliers

rare
or

outliers \rightarrow exponentially less \ll less

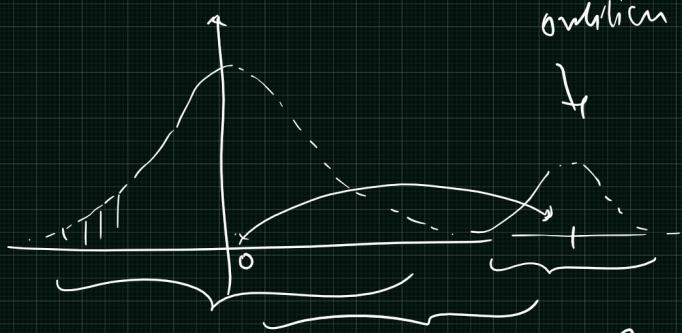
$\rightarrow L_2$

outliers \rightarrow high $\rightarrow L_1$

$$\text{RMSE} = \|V\|_k = \left(|V_0|^k + |V_1|^k + \dots + |V_n|^k \right)^{1/k} \Rightarrow k \text{ norm}$$

$$\|V\|_{k=1} = \left(|V_0| + |V_1| + \dots + |V_n| \right) = L_1$$

$$\|V\|_{k=2} = \left(V_0^2 + V_1^2 + \dots + V_n^2 \right)^{1/2} = L_2$$



$$\Delta S.E. = (y - \hat{y})^2$$

$$A.R.E. = |y - \hat{y}|$$

$$\begin{matrix} 2 \\ y \end{matrix} \quad \begin{matrix} 8 \\ \hat{y} \end{matrix}$$

$$S.E. = 36$$

$$A.R.E. = 6$$

$$\begin{matrix} 2 \\ y \end{matrix} \quad \begin{matrix} 8 \\ \hat{y} \end{matrix}$$

$$-0.5$$

$$S.E. = 0.25$$

$$A.R.E. = 0.5$$

k is large \Rightarrow more sensitive to outliers
 k is less \Rightarrow less sensitive

Man-norm Regularization :-

\hookrightarrow We scale weights.

w for incoming connection is constrained based on below condition

$$\frac{\|w\|_2}{l_2 \text{ norm}} \leq r \quad \hookrightarrow \text{man-norm hyperparameter}$$

After each training step

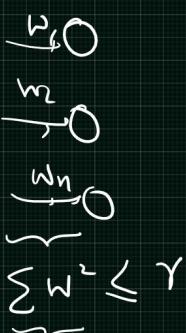
$$w' \leftarrow w - \eta \frac{\partial f}{\partial w} \quad \hookrightarrow \quad w'' \leftarrow \frac{w' r}{\|w'\|_2}$$

conditions :-

1) if $\|w\|_2 = r \Rightarrow$

$$w'' \leftarrow w' \frac{r}{\|w'\|_2}$$

$w'' \leftarrow w'$ No regularization or scaling



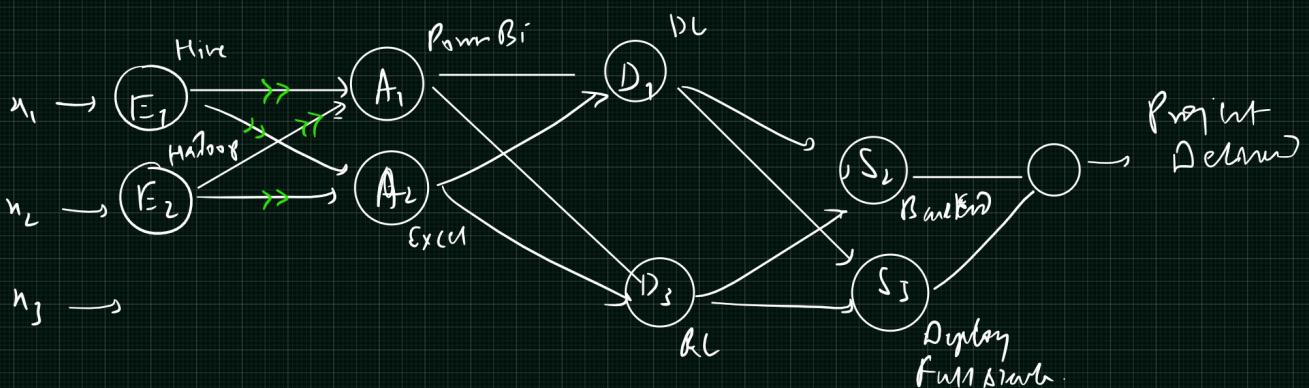
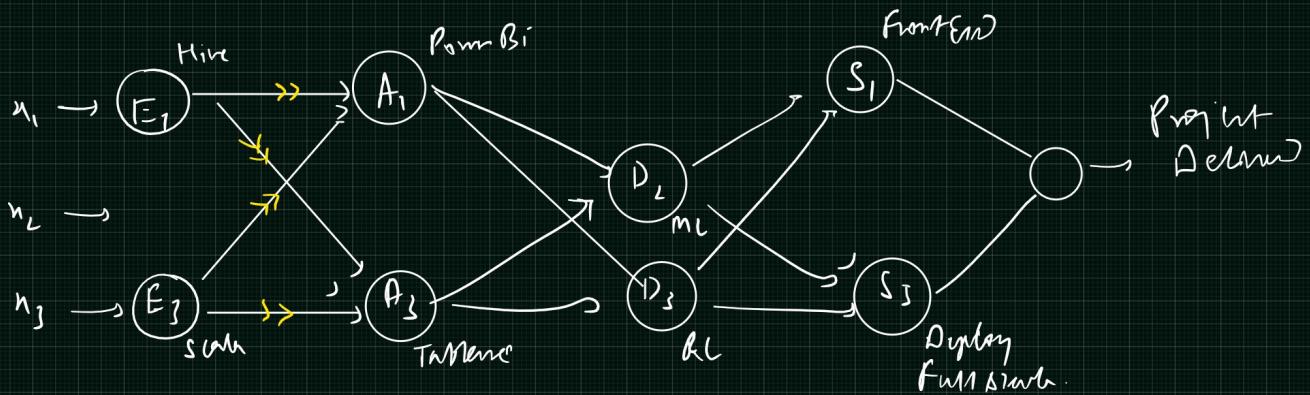
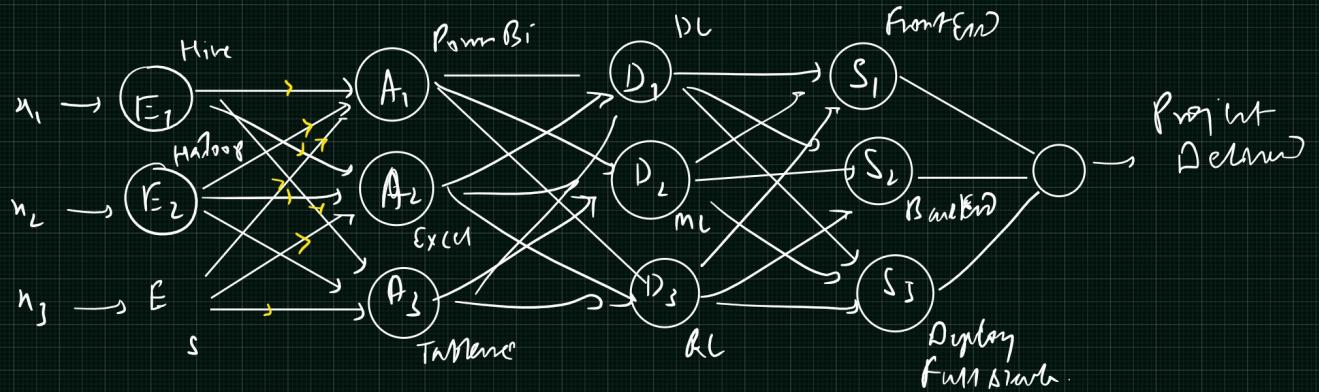
2) if $\|w\|_2 < r$

$$w'' \leftarrow w' \frac{r}{\|w'\|_2} \quad r > \text{denominator}$$

=> Scaling up

Dropout layers! -

Data science Team



$n \rightarrow$ no. of neurons

$2^m \rightarrow$ number is trained.

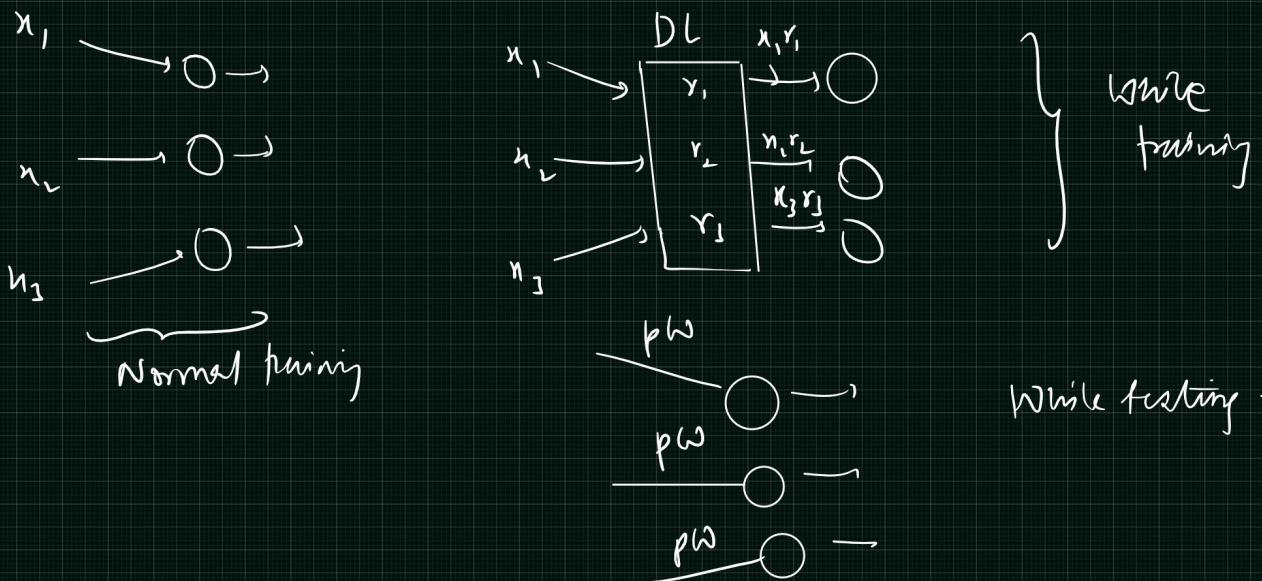
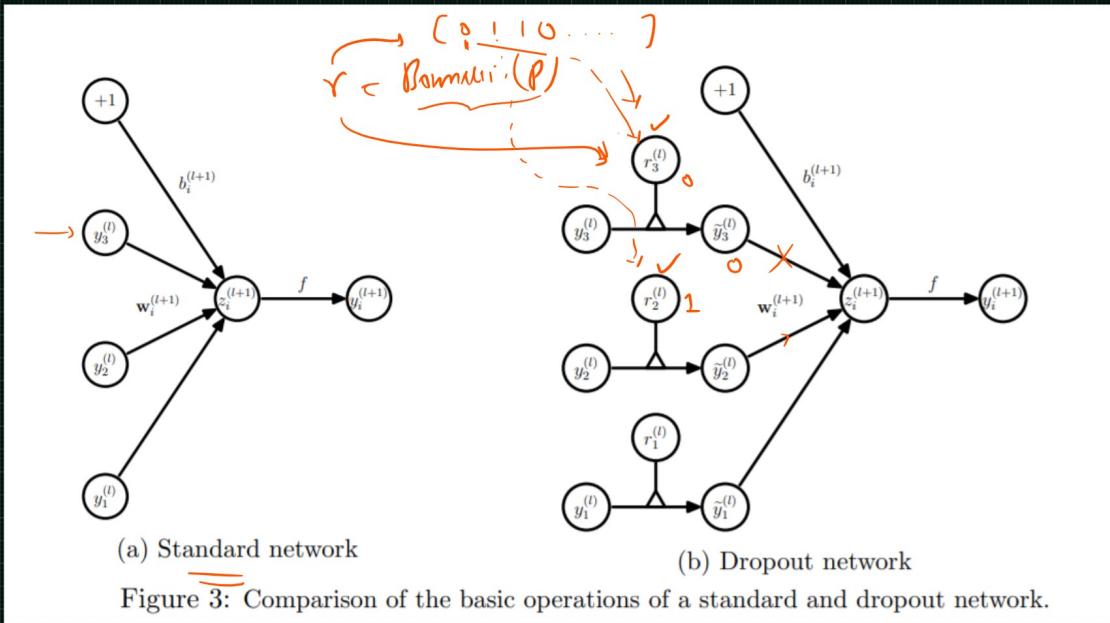
This section describes the dropout neural network model. Consider a neural network with L hidden layers. Let $l \in \{1, \dots, L\}$ index the hidden layers of the network. Let $\mathbf{z}^{(l)}$ denote the vector of inputs into layer l , $\mathbf{y}^{(l)}$ denote the vector of outputs from layer l ($\mathbf{y}^{(0)} = \mathbf{x}$ is the input). $W^{(l)}$ and $\mathbf{b}^{(l)}$ are the weights and biases at layer l . The feed-forward operation of a standard neural network (Figure 3a) can be described as (for $l \in \{0, \dots, L-1\}$ and any hidden unit i)

$$\begin{aligned} z_i^{(l+1)} &= \mathbf{w}_i^{(l+1)} \mathbf{y}^l + b_i^{(l+1)}, \\ y_i^{(l+1)} &= f(z_i^{(l+1)}), \end{aligned} \quad z = w \mathbf{x} + b$$

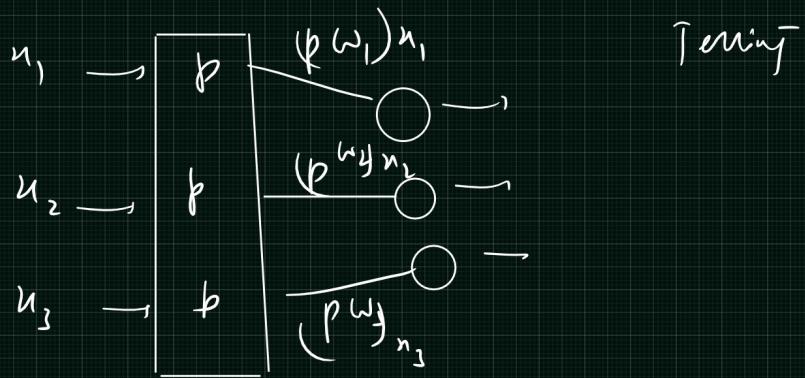
where f is any activation function, for example, $f(x) = 1 / (1 + \exp(-x))$.

With dropout, the feed-forward operation becomes (Figure 3b)

$$\begin{aligned} r_j^{(l)} &\sim \text{Bernoulli}(p), \Rightarrow \text{prob } 1 \text{ or not } 1 \\ \tilde{\mathbf{y}}^{(l)} &= \mathbf{r}^{(l)} * \mathbf{y}^{(l)}, \\ z_i^{(l+1)} &= \mathbf{w}_i^{(l+1)} \tilde{\mathbf{y}}^l + b_i^{(l+1)}, \\ y_i^{(l+1)} &= f(z_i^{(l+1)}). \end{aligned} \quad z = w \mathbf{x}' + b$$

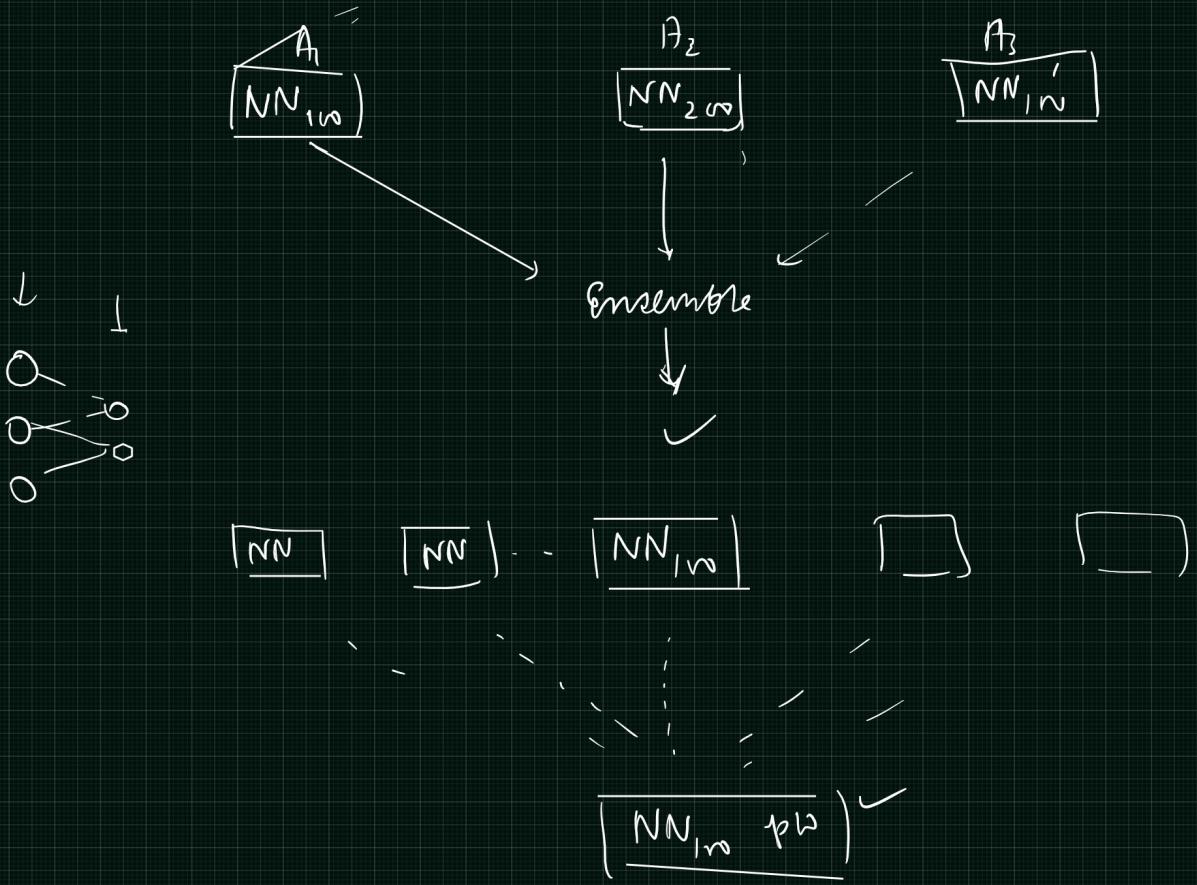


DL



Tuning

$p \rightarrow$ keeping probability \approx
probability of keeping the neurons
 $p \in [0.5, 1]$



Loss functions.

$$\frac{\partial \mathcal{L}}{\partial w} \quad \overbrace{\nabla_{\theta} \mathcal{L}(\theta)}$$

A) Regression :-

(i) MAE mean absolute error | l₁ loss | l₁ norm loss |

$$MAE = \frac{1}{m} \sum_{i=1}^m |y^{(i)} - \hat{y}^{(i)}|$$

Disadvantage:

Not continuously differentiable at every point.

(ii) MSE mean squared error | l₂ loss | l₂ norm loss |

$$MSE = \frac{1}{m} \sum_{i=1}^m (y^{(i)} - \hat{y}^{(i)})^2$$

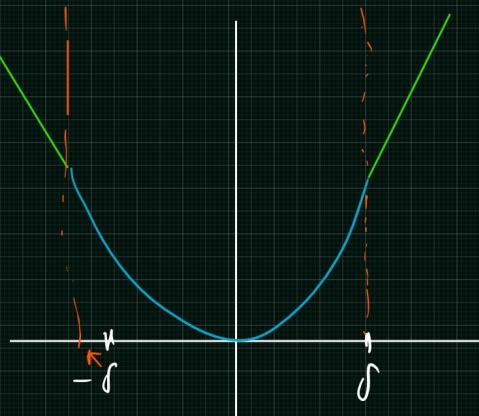
Drawbacks:-

Sensitive to outliers

(iii) Huber loss

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

$$\begin{aligned} |y - \hat{y}| &\leq \delta \\ \Rightarrow -\delta &\leq (y - \hat{y}) \leq \delta \end{aligned}$$



(iv) Pseudo Huber loss

Refer Notebook

B> For Classification

(i) Hinge loss (for Binary classification)

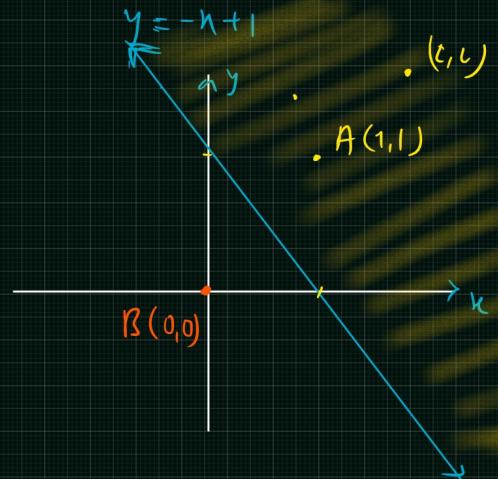
$$L = \max(0, (1 - t \cdot \hat{y}))$$

predicted value
target

here $t = \begin{cases} +1, & \text{Class A} \\ -1, & \text{Class B} \end{cases}$

case 1: $A^{+1} (x_1, y_1)$

$$f(x, y) = \hat{y} = x_1 + x_2 - 1 = 3$$



$$\begin{aligned} L &= \max(0, 1 - (+1)3) \\ &= \max(0, 1 - 3) \\ &= \max(0, -2) \\ L &= 0 \quad \Rightarrow \text{correct prediction} \end{aligned}$$

$$\begin{aligned} y + x - 1 &= 0 \\ \hat{y} &= f(x, y) = 0 \end{aligned}$$

case 2: $\hat{y} = f^{-1}(x_1, y_1) = -3 \quad \text{Assumption}$

$$L = \max(0, 1 - (+1)(-3))$$

$$= \max(0, 1 + 3)$$

$$= \max(0, 4)$$

$$L = 4 \quad \Rightarrow \text{Wrong prediction}$$

Cross Entropy Loss

$$CE = - \sum_{i=1}^C p_i \log(q_i)$$

↑
actual
↓
predicted
↓
probability

Binary Cross Entropy loss



$$\begin{aligned} BCE &= - \sum_{i=1}^2 y_i \log(\sigma(z_i)) \\ &= - \sum_{i=1}^2 y_i \log \hat{y}_i \end{aligned}$$

$$BCE = -y \log \hat{y} - (1-y) \log(1-\hat{y})$$

Categorical Cross Entropy loss (class > 2)



$$\begin{aligned} CCE &= - \sum_{i=1}^n y_i \log(\text{softmax}(z_i)) \\ &= - \sum_{i=1}^n y_i \log(\hat{y}_i) \end{aligned}$$

If $n=3$.

Actual value	y_A	y_B	y_C	\rightarrow	1
predicted value	\hat{y}_A	\hat{y}_B	\hat{y}_C	\rightarrow	1

$$CCE = -y_A \log(\hat{y}_A) - y_B \log(\hat{y}_B) - y_C \log(\hat{y}_C)$$

	y_A	y_B	y_C
actual	0	1	0
pred	0.3	1.0	0
	\hat{y}_A	\hat{y}_B	\hat{y}_C

$$\text{CE} = - \underbrace{0 \cdot \log(0)}_{=0} - \underbrace{1 \cdot \log(1)}_{=0} - \underbrace{0 \cdot \log(0)}_{=0}$$

$$\text{CE} = 0$$

pred 0.7 0.3 0

$$\begin{aligned}\text{CE} &= -1 \cdot \log(\hat{y}_A) \\ &= -1 \cdot \log(0.7) \\ &= -(-0.5) \\ \text{CE} &= +0.5\end{aligned}$$

0 $\xrightarrow{\text{true}}$ 1
-ve true

Categorical cross entropy

$$\begin{array}{ccc}y_A & y_B & y_C \\ 0 & 1 & 0 \\ \ln = 1 & \ln = 2 & \ln = 3\end{array}$$

— Sparse categorical cross entropy

$$\begin{array}{ccc}y_A & y_B & y_C \\ 1 & 2 & 3\end{array}$$

sparse array \rightarrow [list of zeros, followed by non-zero values]

$$\begin{array}{ccc}10^b & \downarrow & \left(\begin{array}{cccccc} 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{array} \right) & \left\{ \begin{array}{l} (0, 3) \rightarrow 1 \\ (1, 2) \rightarrow 1 \end{array} \right\} \end{array}$$