

* Loss Function in deep learning.

• Start with loss function!

• Loss function is a model of evaluating how your algorithm is modelling your dataset.

• High value \Rightarrow Algorithm performance is poor.
Low value \Rightarrow Algorithm performing great.

• Loss function is just mathematical function.

• Loss function is function of parameters in machine learning.

e.g. For linear regression $\Rightarrow L(m, b)$

• How we change the loss function value?

• By changing parameter value.

* Why is loss function important?

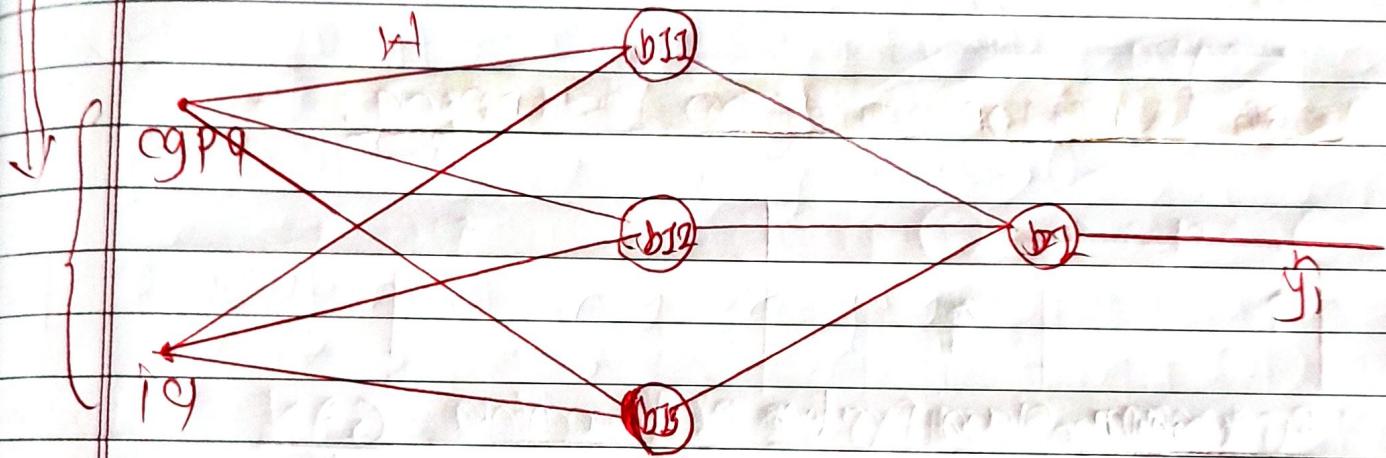
• You can't improve what you can't measure

- Peter Drucker

• This line by Peter Drucker perfectly explains the importance of loss function in machine learning/ deep learning.

What is role of loss function in deep learning?

cgp	iq	package (IPG)
7.1	83	3.2
8.5	91	4.5
6.3	102	6.1
5.1	87	2.7



known,

- To role we will discuss back propagation for now,
- We put first student 'cgp' and 'iq' here and we will put this in neural network by using forward propagation it will predict for first student.
- Let's say for first student if coming out to be 3.7.
- Since this is regression problem, we will

Use mean squared error.

$$(3.2 - 3.7)^2 = (0.5)^2 = 0.25$$

- Apply gradient descent on it and then update weight and biases
- We will do this for every point, after doing this is back propagation.

Idea of updating weight and biases is to reduce loss.

Loss function in deep learning.

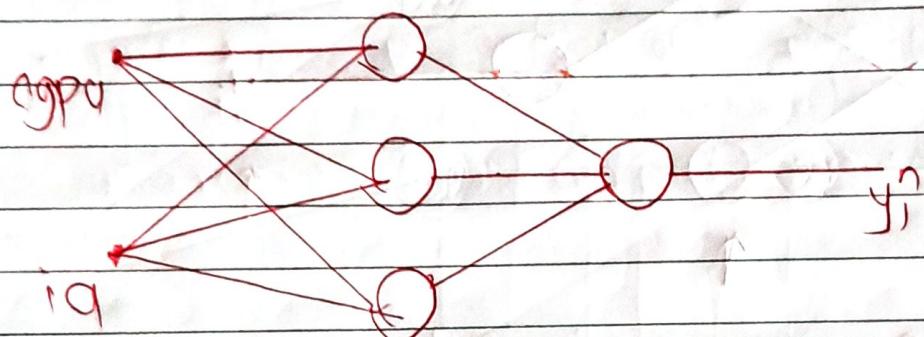
Object detection	Regression	Classification	Autoencoder	GAN	Image editing
• Focal loss	• msp • map • number loss	• Binary- cross-entropy • categorical- cross-entropy • hinge loss	• KLDivergence	• digit- minitor loss	• Triplet loss • max gain loss

- We can create our custom loss function too, keras allow to do that.

Loss function and cost function difference

- This is our input

cgpa	iq	package	Prediction
6.5	100	6.3	6.1
7.1	92	4.1	4
8.5	88	3.5	3.7
9.2	802	7.2	7



- Loss ~~(student)~~ calculate on single training example

For first row loss function would be,

$$(y_i - \hat{y}_i)^2 = (6.3 - 6.1)^2 \Rightarrow \text{loss function}$$

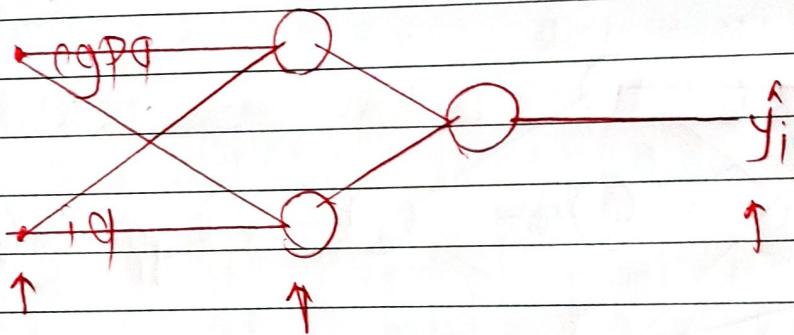
- When we calculate for single row, this is known as loss function
- When we calculate on batch it will become cost function

Cost Function

$$\frac{1}{n} \sum (y_i - \hat{y}_i)^2$$

Mean Squared Error

Also known as, squared loss) (2 logg, Let's study in context of deep learning.



How this model trains?

- We insert rows one by one and we get \hat{y}_i (forward propagation)
- After that we check how off is this prediction, By using loss function $(y_i - \hat{y}_i)^2$
- Based on this loss function we update weight and bias
- We will do this for every point, idea is to minimize value of loss function

- We will discuss advantages of MSE also disadvantages of it to know which domain it could apply on.

Why just $(y_i - \hat{y}_i)$ instead of $(y_i - \hat{y}_i)^2$?

- Because of squaring negative values goes down
- Total error reduced because of minus value, we want all error to be added not subtract.
- After squaring, $(y_i - \hat{y}_i)^2$ becomes quadratic,

$$(true(y_i) - predicted(\hat{y}_i))^2 \Rightarrow 1 \text{ unit}$$

After squaring $\xrightarrow{\quad}$ 1 unit

2 unit After squaring 2 unit

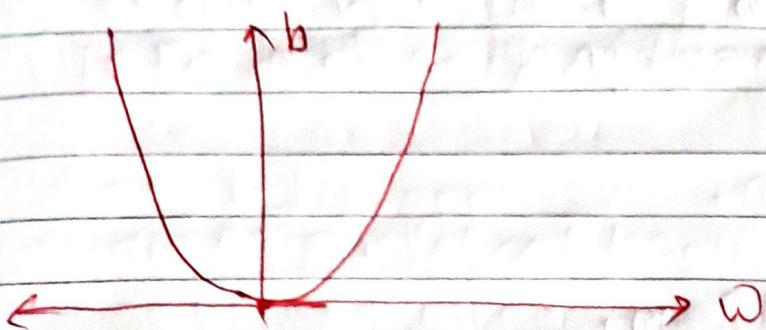
3 unit After squaring 3 unit

- Those points making more error/which are very far from actual line leads to change more weight and bias distinctly.

Wrongly classified $(y_i - \hat{y}_i)^2$ too much \Rightarrow makes more changes in W, b

Advantages

- Easy to interpret
Because of quadratic nature this function is always differentiable.



- There is Only '1' local minima.

Disadvantages

- Prior unit is square,
- Not Robust to outliers.

Condition to apply MSE in deep learning.

- Last neuron activation function always should be 'linear'.

Mean Absolute error (MAE)

- Also known as 'L1 loss'

$$L = |y_i - \hat{y}_i| \implies \text{Loss Function}$$

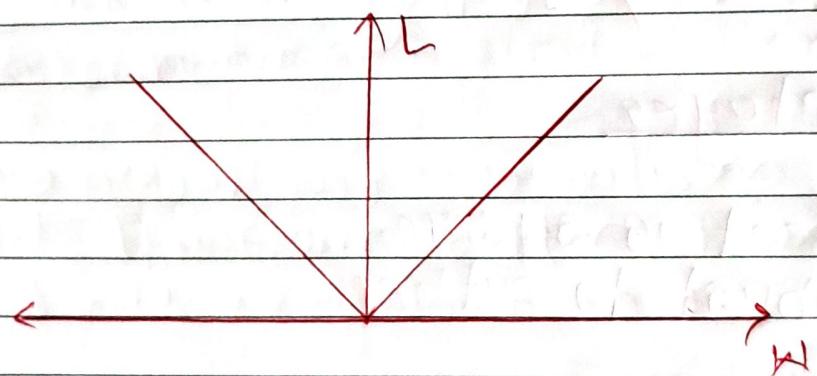
$$L = \frac{1}{n} \sum |y_i - \hat{y}_i| \implies \text{not function.}$$

Advantages

- Easy to understand,
- Unit slope as unit of y ,
- Robust to outlier (so error cannot increase)
(I punish a, b drastically as do in quadratic loss function i.e. MSE).

Disadvantages

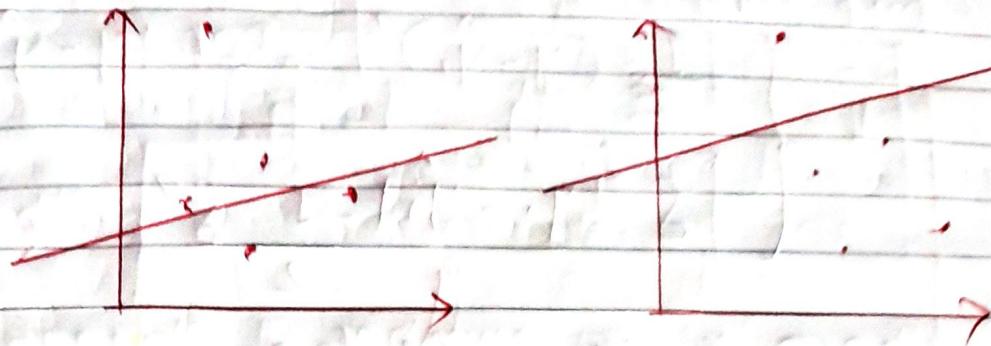
Graph of MAP is,



- This graph is not differentiable, so we cannot apply gradient descent directly we have to calculate sub-gradient descent first.
- Calculating sub-gradient descent computationally expensive.

Huber Loss

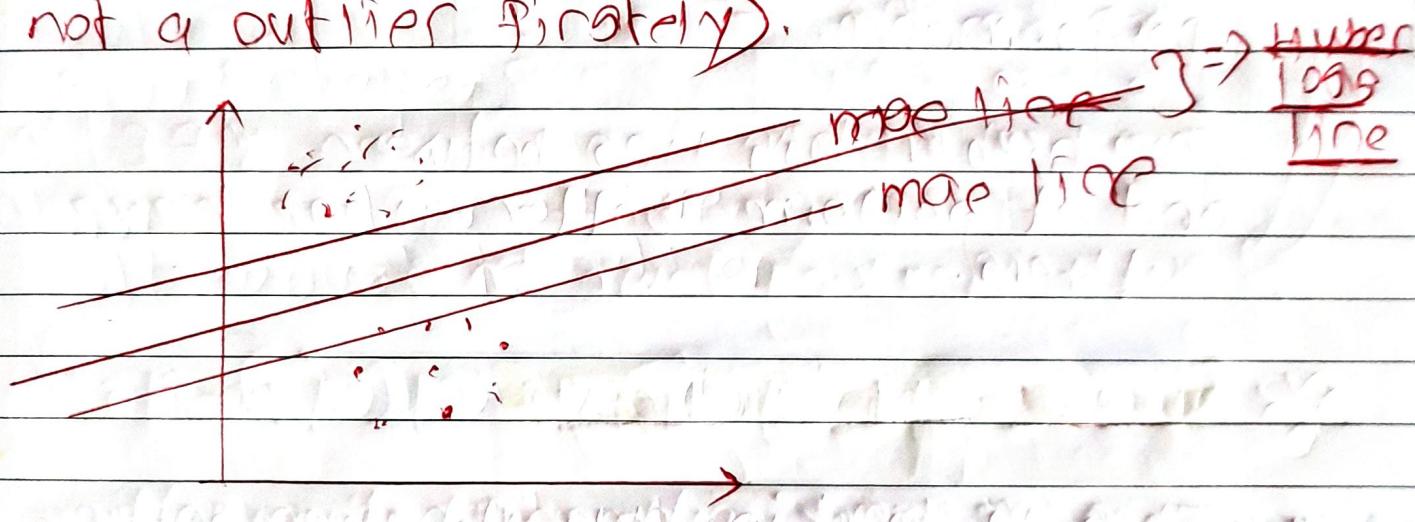
- MSE cause drastic change in weights and map cause less changes in weight & biases



map to
outlier

map to
outlier

- Let's we have this kind of data that it have 25% of outlier (so that 25% is not a outlier specially).



- In this case we want middle line which is we get after updating weight and biases by using Huber loss as loss function.

- If outliers not present \Rightarrow behave like MLE.
- If outliers present \Rightarrow behave like MAD.
- Mathematically answer is,

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

↓ mae.

' δ ' parameter used to determine which point is outlier or not.

' δ ' by adjusting this parameter we can check performance.

- This huber loss lies between MAE / MSE and it perform better if we have more outliers.

2) Binary Cross Entropy

- This is used in classification setting
- It will be only used when you have only two classes.

$$\text{Loss Function} \Rightarrow -y \log(\hat{y}) - (1-y) \log(1-\hat{y})$$

$\hat{y} \Rightarrow$ neural net architecture prediction
 $y \Rightarrow$ Actual value / target.

- Activation function for last neuron is sigmoid

Cost Function $\Rightarrow \text{log}$

$$-\frac{1}{n} \left[\sum_{i=1}^n y_i \log y_i + (1-y_i) \log(1-y_i) \right]$$

- After calculating log we update weight and biases by using gradient descent.

Advantages

- This log function is differentiable (we can apply gradient descent).

Disadvantages

- Multiple minima could present.
- This is not intuitive.
- Some people called it log-log too.

Categorical cross entropy

- Used in classification scenarios, only when you have multi-class classification problem.
- ~~This is not up, as this same log function used in softmax regression too.~~
- To calculate above for single student, this is

Our formula,

$$L = - \sum_{j=1}^k y_j \log(\hat{y}_j)$$

- Where 'k' is number of classes in data

- This is our opt up, (Close after one-hot encoding)

<u>Copy</u>	<u>1g</u>	<u>Planned</u>	<u>yes</u>	<u>no</u>	<u>maybe</u>
8	80	yes	1	0	0
6	60	no	0	1	0
7	70	maybe	0	0	1

- After breaking down this,

$$L = -y_1 \log(\hat{y}_1) - y_2 \log(\hat{y}_2) - y_3 \log(\hat{y}_3)$$

- When we work with multi-class classification some changes occurs, in output layers we have no. of neurons = no. of class.

- Activation Function for output layers is softmax, This is formula,

$$f(z_i) = \frac{e^{z_i}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

$$\begin{aligned} z_1 &\Rightarrow \text{yes} \\ z_2 &\Rightarrow \text{no} \\ z_3 &\Rightarrow \text{maybe} \end{aligned}$$

- $f(z_1) + f(z_2) + f(z_3) = 1$

- This is loss function for single row,

$$L = -y_1 \log(\hat{y}_1) - y_2 \log(\hat{y}_2) - y_3 \log(\hat{y}_3)$$

- For first row, all neuron's which basically three neuron output is,

$$\begin{bmatrix} 0.2 & 0.3 & 0.5 \end{bmatrix}$$

↓ ↓ ↓
yes no maybe

$$\begin{bmatrix} 1 & 0 & 0 \end{bmatrix} \Rightarrow \text{true table}$$

↓ ↓ ↓
y₁ y₂ y₃

+ loss

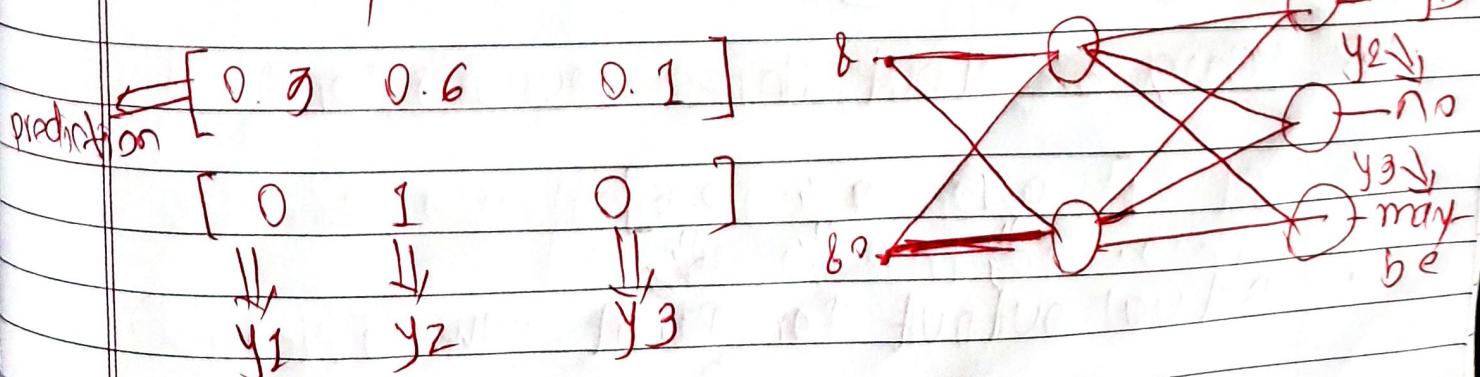
- We apply loss function,

$$L = -y_1 \log(\hat{y}_1) - 0 + 0$$

$$L = -1 \log(0.2) \Rightarrow \text{We update weight & bias accordingly.}$$

- This is loss for first row.

- Similarly for second row,



- Loss becомe 10,

$$\begin{aligned} L &= -y_2 \log(y_2) \\ &= -1 \log(0.6) \end{aligned}$$

- similarly for third point/ row;

$$L = -y_3 \log(y_3)$$

After
that
we update
weight
biases
according

Sparse categorical cross Entropy

cgpa	iq	planned	AFTER integer encoding
7	70	yes	1
8	80	no	2
6	60	maybe	3

- If we use categorical cross entropy we use one-hot encoding

- For sparse categorical cross entropy we use integer encoding

- This is output three neurons for first row,

$$\begin{bmatrix} 0.1 & 0.4 & 0.5 \\ \downarrow & \downarrow & \downarrow \\ y_1 & y_2 & y_3 \end{bmatrix}$$

- Actual output for first row is ($y_2 = 1$)

- so loss becomes,

$$L_1 = -Y_1 \log(Y_1)$$

$$= -1 \log(0.2) \Rightarrow$$

Forgot about them.

$$-Y_2 \log(Y_2) \quad \& \quad Y_3 \log(Y_3)$$

After this weights are updated accordingly.

- Similarly for second and third row,

$$L_2 = -Y_2 \log(Y_2) = -\cancel{Y_2} \log(Y_2)$$

$$L_3 = -Y_3 \log(Y_3) = -\cancel{Y_3} \log(Y_3)$$

- If we have too much classes we should use sparse categorical entropy since other calculation not needed and we don't have to one-hot encode it.

$$-Y_3 \log(Y_3) \quad \& \quad Y_2 \log(Y_2)$$

Forgot about them.

$$\rightarrow -Y_2 \log(Y_2) \quad \& \quad Y_1 \log(Y_1)$$

- Cost Function,

$$C = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^k y_{ij} \log(Y_{ij})$$

Summary

- Regression

- use more
- use map if outliers present. } for combination of outlier and normal distribution
- classification
 - Binary classification \Rightarrow Binary cross entropy.
 - Multiclass classification \Rightarrow categorical cross entropy. (less category).
 - If we have more category/classes then use sparse categorical cross entropy.