

LOSS Functions in MLP (ANN)

* What is Loss Function?

In deep learning, the Loss Function is the feedback mechanism.

It is method for evaluating that how bad our model is performing.

High Loss \rightarrow Model is performing bad.

Low Loss \rightarrow Model is performing better.

Intuition - Imagine 'A person giving an interview for Data science job role, so the feedback of interviewer in this case is the loss function of the model named as Data science Interview'.

* Loss vs Cost :-

① Loss Function - This calculates the error on single training example (one row of data).
For example calculating the difference between the predicted and actual salary for a student.

② Cost function - This is the average of loss functions for entire training dataset (or batch).
It represent total error of model.

* How Loss Functions work in training? :- (16)

① Forward Propagation - Input data (eg: cgpa, iq) enters in the Neural Network.

The network uses the random initial weights to make a prediction.

② Calculate Loss - We compare the prediction against the true value using loss function formula (eg. MSE)

③ Backward Propagation - We use loss value to adjust the weights and biases using an optimizers (eg. gradient descent)

④ Repeat - This happens for every student (data-point) and repeats multiple times (epochs) until the loss is minimized.

* Types of loss Functions :-

In MLP generally we have two types of Loss functions

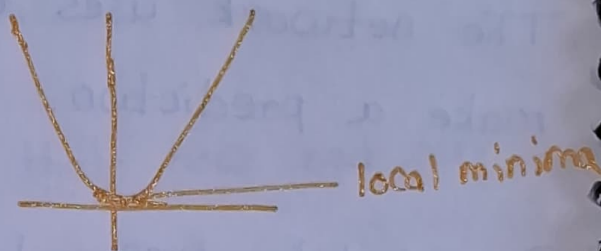
① For Regression

② For classification.

① Regression Loss Functions:-

(i) Mean Squared Error - It is the squared distance of actual values from the predicted value.

$$L = (y_i - \hat{y}_i)^2$$



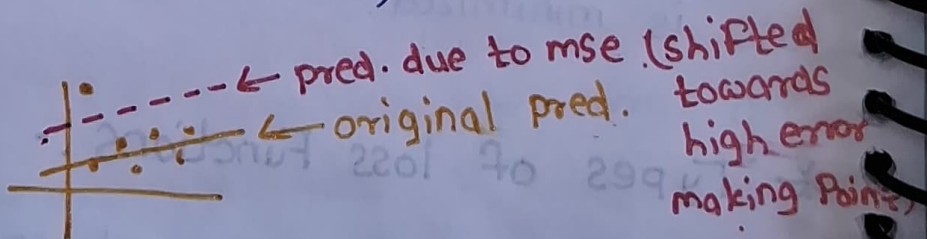
We take squares, because it ensures that during backpropagation the opposite signs does not cancel out each other.

Advantages -

It is easy to interpret

It is differentiable (The curve is smooth so it become easy to use during gradient descent)

It has only one local minima, so has only one Perfect solⁿ.



Disadvantages -

It is not robust to outliers (because due to squares it penalizes the errors more)

Unit mismatch due to squares. The errors are always in squared units.

(ii) Mean Absolute Error (MAE) - It is the positive difference between the actual and predicted value.

$$L = |y_i - \hat{y}_i|$$

$$C = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$



It is used when there are minor outliers in data, because it does not penalizes high errors highly.

Advantages -

Unit remain same as original data

Perform better when some outliers are present.

Disadvantages -

Not differentiable at ~~one~~ like MSE due to it is not a curve. (V shaped)

Time consuming.

(iii) Huber loss - This function attempts to fix the flaws of MSE & MAE by combining them.

It behaves like MSE when no outliers & behaves like MAE when too many outliers (Natural outliers)

Note -

MSE - Use for general regression with No or minimum outliers. ~~(True outliers)~~

MAE - Use when data have outliers (It is robust)

Hubber - The hybrid choice. Best of both world.

Use 'linear or No A.F' in output Layer while using these loss functions.

② Classification Loss Functions:-

① Binary cross Entropy - (Log Loss) -

$$L = -[y \cdot \log(\hat{y}) + (1-y) \cdot \log(1-\hat{y})]$$

It is used for binary classification problem along with 'Sigmoid' A.F. in output layer.

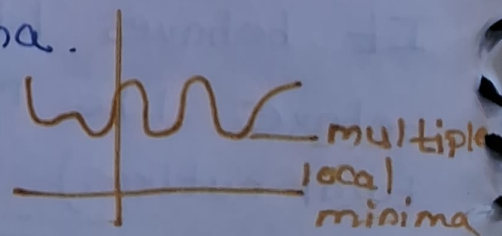
It penalizes the model heavily for making wrong class prediction.

Advantages -

It is differentiable.

Disadvantage -

It can have multiple local minima.



② Categorical Cross Entropy :- It is used for Multiclass classification problem. (more than 2 classes)

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

where, k - No. of classes.

The output layer must have No. of Neurons equal to classes.

i.e. if No. classes = 10 then No. of Neurons in output layer = 10

The A.F must be 'Softmax' in output layer.

Advantages -

Can be useful in Multiclass classification problem.

Disadvantages -

It requires one-hot encoding for target variable.

③ Sparse categorical cross entropy - Same as categorical cross entropy, but does not require One-hot encoding for target variable.

Note -

Binary classification \rightarrow Binary cross entropy \rightarrow
sigmoid A.F. in O/P Layer

Multiclass Classification \rightarrow Categorical cross entropy



Softmax A.F in O/P layer.

(Requires One-Hot(Y))

\uparrow
 \downarrow skip

sparse categorical cross Entropy