

Binary class classification

Output layer = sigmoid

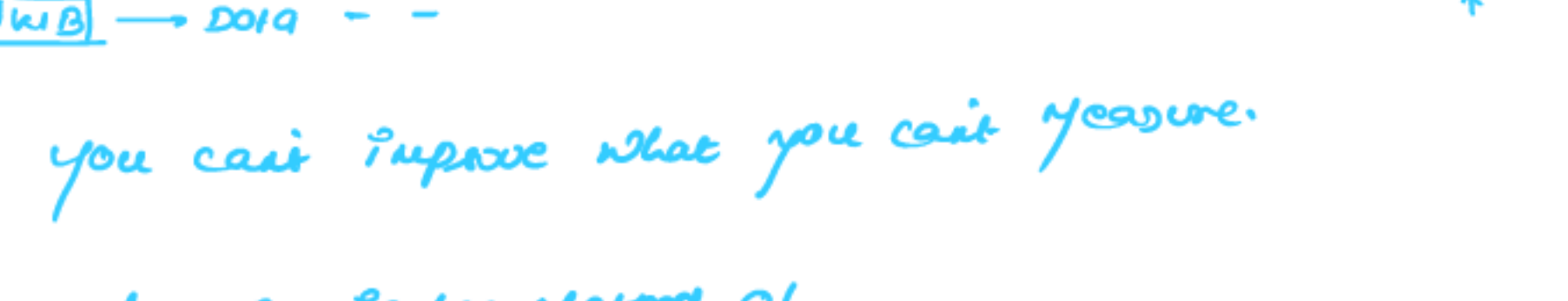
Multiclass classification

Output layer = softmax

Dead neuron  $\Rightarrow$  associated with one neuron.

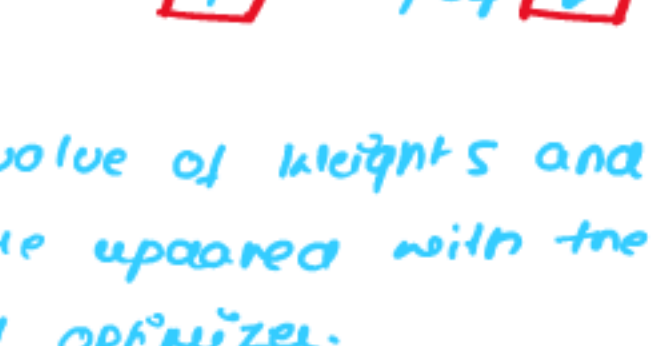
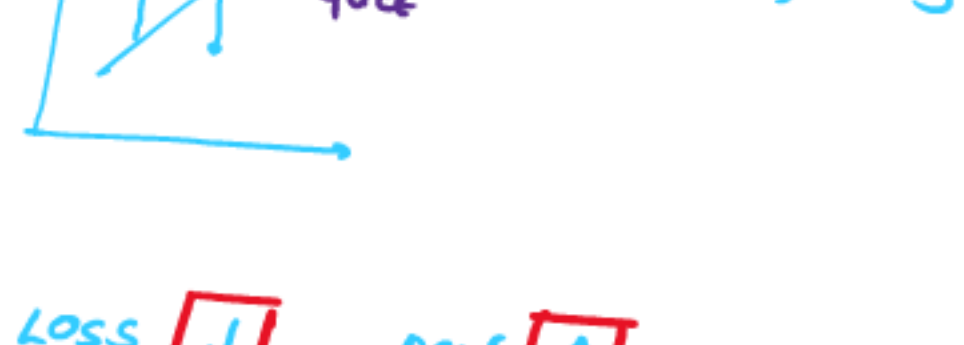
vanishing gradient  $\rightarrow$  layer specific problem

## Loss function.

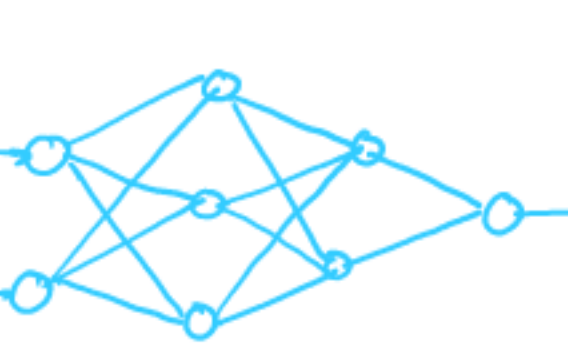


You can't improve what you can't measure.

Loss function is the method of evaluating how well our algorithm performing on given dataset.



value of weights and biases are updated with the help of optimizer.

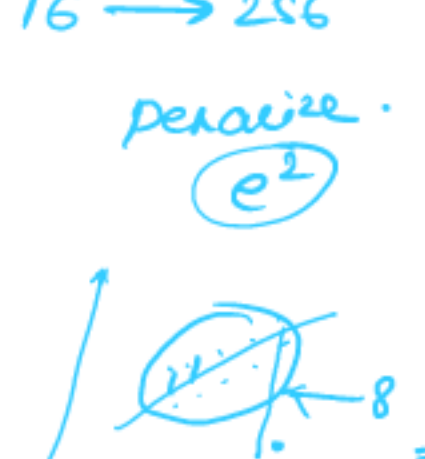


## Mean squared error

loss fn.  $\Rightarrow$  single record  
cost fn.  $\Rightarrow$  for whole batch.

$(y_{act} - \hat{y})^2$  is loss fn.

$$\frac{1}{n} \sum_{i=1}^n (y_{act} - \hat{y})^2 = CF$$



## Advantages.

- It is in the form of quadratic eq., we will get the gradient descent with only one global minima.
- no local minima.
- easy to interpret
- differentiable



## Disadvantage.

- Error unit is not same. It is in sq form.
- not Robust to outlier.
- true nature of our data is not captured and it got deviated.

## Mean Absolute Error.

$$Loss = |y_{act} - \hat{y}|$$

$$CF = \frac{1}{n} \sum |y_{act} - \hat{y}|$$



## Advantages.

- Error unit is same.
- Easy to interpret
- Outliers are better handling than MSE, it is not penalizing the model by sq. the errors



## Disadvantage

- Graph Differentiable
- there may have multiple local minima.

## Huber Loss

$$Loss = \begin{cases} \frac{1}{2} (y_{act} - \hat{y})^2 & \text{for } |y_{act} - \hat{y}| \leq \delta \\ \delta |y_{act} - \hat{y}| - \frac{1}{2} \delta^2 & \text{for } |y_{act} - \hat{y}| > \delta \end{cases}$$



It is less sensitive to outliers  
Huber Loss is used in regression problem compared with MSE and MAE.

Whereas data point is not an outlier it will behave like MSE.

Whereas data point is an outlier it will behave like MAE.

$$\delta = 1.0$$

## Advantages.

- Outliers are handled properly.
- Local minima situations are also handled fine.

## Disadvantage.

- complex.

## Classification

### Binary Cross Entropy.

$$y = \begin{matrix} 1 & 0 \end{matrix} \quad \hat{y} = \begin{matrix} 1 & 0 \end{matrix}$$

$$y = \begin{matrix} 1 & 0 \end{matrix} \quad \hat{y} = \begin{matrix} 0 & 1 \end{matrix}$$

$$y = \begin{matrix} 1 & 0 \end{matrix} \quad \hat{y} = \begin{matrix} 0.5 & 0.5 \end{matrix}$$

$$Entropy \Rightarrow 1$$

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

$$CF = \frac{1}{n} \sum_{i=1}^n [y \log(\hat{y}) + (1-y) \log(1-\hat{y})]$$

$$class 0 \Rightarrow -\log(1-\hat{y})$$

$$class 1 \Rightarrow -\log(\hat{y})$$

## Categorical Cross Entropy.

used for multiclass classification problem.



$$Loss = - \sum_{i=1}^K y_{act} \log(\hat{y}_i)$$

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

$$K = \text{no. of classes}$$

all the probabilities of class will sum up to 1

$$y = \begin{matrix} 1 & 0 & 0 \end{matrix} \quad \hat{y} = \begin{matrix} 0.3 & 0.2 & 0.5 \end{matrix}$$

$$-1 \times \log(0.3) - 0 \times \log(0.2) - 0 \times \log(0.5)$$

$$-1 \times \log(0.3) - 0 - 0$$

$$\Rightarrow -1 \times \log(0.3) = 1.101$$

$$\Rightarrow \text{Result}$$

produce one hot array.

containing to probability given for each category.

$$Cost = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^K y_{ij} \log(\hat{y}_{ij})$$

## Sparsity categorical Cross Entropy.

It is frustrating when you are dealing with large classification problems. w/ 1000 classes.

And when we do one hot encoding of 1000 classes  $\rightarrow$  1000 vectors

increased significant amount of memory.

$$y = \begin{matrix} 1 & 0 & 0 \end{matrix} \quad \hat{y} = \begin{matrix} 0.3 & 0.2 & 0.5 \end{matrix}$$

$$-1 \times \log(0.3) - 0 \times \log(0.2) - 0 \times \log(0.5)$$

$$-1 \times \log(0.3) - 0 - 0$$

$$\Rightarrow -1 \times \log(0.3) = 1.101$$

$$\Rightarrow \text{Result}$$

## Dropout Ratio.

Dropout ratio is the technique that drops the neuron from the neural net as figure then during training.

Some neurons are involved in training and rest are turned off.

"to prevent overfitting"

these units are not connected during particular forward as backward propagation

$$P = 0.5$$

Dropout is an approach in neural networks which helps in reducing independent learning among neurons.

## Exploding gradient



$$w_{new} = w_{old} - \alpha \frac{\partial L}{\partial w_{old}}$$

$$w_{11}^{new} = w_{11}^{old} - \alpha \left[ \frac{\partial L}{\partial w_{11}} \right]$$

$$= 0.2 \times 10 \times 0.5 \times 1$$

$$= 1$$

## Solution.

proper weight initialization.

non saturating act fn.

Model is exploding with learning rate.

Exponential growth in model parameters.