

Today's Agenda

- 1) Loss Functions
- 2) Activation Functions

Backprop / Optimizers

What are loss functions?

Method of evaluating how well the algorithm performs on top of the dataset. $\{f(x) = x^2 + 2\}$

high \rightarrow poor

low \rightarrow good

$L(\text{parameters})$

$$y = mx + c$$

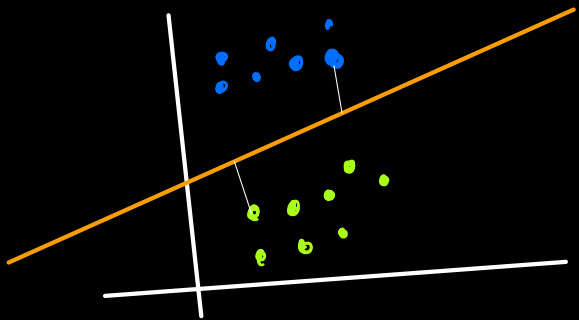
(m, c) are the parameters

Parameters = Weights + Bias

| | |
|--------------------------|---|
| <u>LOOP</u> | $(y - \hat{y})^2 = ?$ |
| For each training sample | updating the line update the parameters Optimizers (Gradient Descent) |

Condition to stop :- least loss

Eye of the Algorithm



$$L(m, c) = \text{minimum}$$

Errors

Far

Close

MORE

LESS

$$(y - \hat{y})^2$$

= not always positive

$$= 0.7 + 0.1 + 0.2 - 0.7$$

SAT | GRE | 12 | LPA

$$y = 6.5$$

$$\hat{y} = 7.2$$

$$(6.5 - 7.2) = -0.7$$

Reduce the overall loss

Quadratic Term

$$x^2 + x + 2$$

1) High Penalization

$$(7)^2 = \underline{49} = \sqrt{49}$$

Mean Square Error

Error / Loss / Cost

- 1) Error / Loss \rightarrow single data point
- 2) Cost \rightarrow for the entire dataset / batch

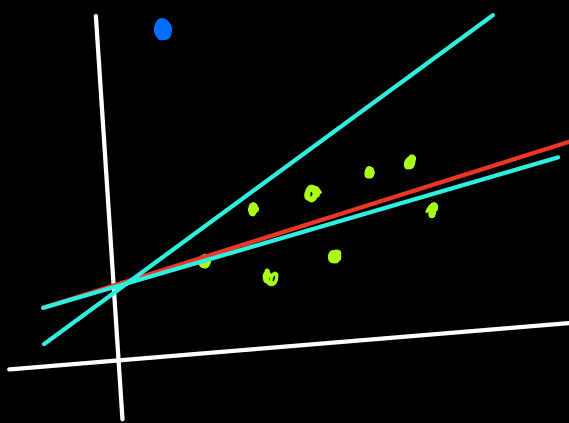
Mean Square Error

$$LF = (y - \hat{y})^2$$

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

Cons

- 1) Cannot handle outliers



Type of Loss functions

Regression

- ① MSE
- ② MAE

Classification

- B

 - ① Log loss
BCE

M

 - ① C.C.E
 - ② S.C.C.E

BCE

③ Huber Loss

③ Hinge loss

Auto encoder

① KL Divergence

GAN

① Discriminator Loss

Object Detection

Focal loss

Embeddings

Triplet loss

Object Detection

1) Co-ordinate of the object

→ Regression

2) Class of the object

→ Classification

Cost Function = Regression + Classification

Mean Square Error

| 10 th | 12 th | GRE | LPA (y) | (\hat{y}) |
|------------------|------------------|-----|-------------|---------------|
| 63 | 73 | 320 | 7.2 | 6.5 |
| 51 | 61 | 260 | 3.1 | 2.4 |

$$(\text{True} - \text{predicted})^2$$

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

Example :- $(7.2 - 6.5)^2 = (.7)^2 = \underline{.49} \text{ (loss)}$

$$C.F = \frac{1}{2} (.49) + (3.1 - 2.4)^2$$

$$= \frac{1}{2} .49 + .49$$

$$= \frac{.98}{2}$$

$$C.F = \underline{\underline{.49}}$$

Pros

- 1) Easy to interpret
- 2) Differentiable (SGD)

Cons

- 1) Not Robust to outliers

3) 1 local minima

Mean Absolute Error (L1 Loss)

$$L.F = |y - \hat{y}|$$

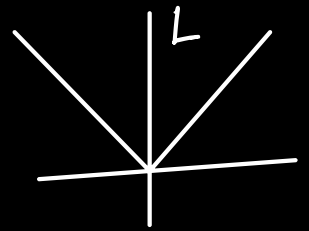
$$C.F = \frac{1}{n} \sum_{i=0}^n |y - \hat{y}|$$

Pros

- 1) Easy same unit
- 2) Robust to outliers

Cons

- 1) Not differentiable
↳ subgradients



Huber Loss

$$L = \begin{cases} \frac{1}{2} (y - \hat{y})^2 \\ \delta |y - \hat{y}| - \frac{1}{2} \delta^2 \end{cases}$$

for $|y - \hat{y}| \leq \delta$

otherwise

hyperparameter

Combination of MSE and MAE

Pros

1) More Robust to Outliers

Classification

1) Binary Classification

Log loss / Binary Cross Entropy

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

$$\text{Cost Function} = -\frac{1}{n} \left[\sum_{i=1}^n y \log(\hat{y}) + (1-y) \log(1-\hat{y}) \right]$$

| CRPA | IO | Placement (y) | \hat{y} |
|------|----|-------------------|-----------|
| | | | |
| | 10 | | |
| | 68 | 1 | 0.7 |
| 8 | | | |
| | 61 | 0 | 0.3 |
| 7 | | | |

Case 1

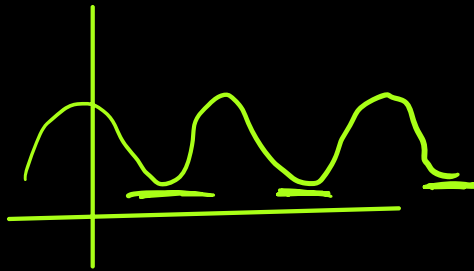
Case 2

$$\begin{aligned}
 &= -\gamma \log(\hat{\gamma}) \\
 &= -1 \log(0.7) \\
 &= -1.0
 \end{aligned}$$

$$\begin{aligned}
 &= -(1-\gamma) \log(1-\hat{\gamma}) \\
 &= -1 \cdot \log(0.7) \\
 &= -1.0
 \end{aligned}$$

Advantages

1) Differentiable



Cons

1) M.C.C Problem

2) Multiple local minima

Multi-Class Classification

1) Categorical Cross Entropy

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

$K = \text{no of classes in data}$

OHE

| 12 | CUA | Placed |
|----|-----|--------|
| 80 | 8.2 | Yes |
| 60 | 6.1 | No |
| 70 | 7.1 | Maybe |

| Yes | No | Maybe |
|-----|----|-------|
| 1 | 0 | 0 |
| 0 | 1 | 0 |
| 0 | 0 | 1 |

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

For Case 1 (Yes)

$$= -y \log(\hat{y}_1)$$

For Case 3 (Maybe)

$$= -y \log(\hat{y}_3)$$

For Case 2 (No)

$$= -y \log(\hat{y}_2)$$

2) Sparse Categorical Cross Entropy

| | H_1 | H_2 | 0 |
|-------|-----------------------|-----------------------|--|
| x_1 | <input type="radio"/> | <input type="radio"/> | <input checked="" type="radio"/> $\rightarrow 0.4$ |
| x_2 | <input type="radio"/> | <input type="radio"/> | <input checked="" type="radio"/> $\rightarrow 0.3$ |
| | <input type="radio"/> | <input type="radio"/> | <input checked="" type="radio"/> $\rightarrow 0.3$ |

Activation : Softmax

Class

$$\text{Softmax} = \frac{e^{z_1}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

Class 1

$$\frac{e^{z_2}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

Case 1

$$\begin{bmatrix} 0.5 & 0.3 & 0.2 \\ 1 & 0 & 0 \end{bmatrix}$$

$$= -y \log(\hat{y})$$
$$= -1 \log(0.5)$$

Case 2

$$\begin{bmatrix} 0.3 & 0.6 & 0.1 \\ 0 & 1 & 0 \end{bmatrix}$$

$$= -1 \log(0.6)$$

| IQ | CGPA | Placed |
|----|------|--------|
| 80 | 8.2 | Yes |
| 60 | 6.1 | No |
| 70 | 7.1 | Maybe |

Integer Encoding

1

2

3

$$= -\gamma \log(\hat{y})$$

Very Fast

$$L.F = - \sum_{j=1}^K \gamma_j \log(\hat{\gamma}_j)$$

$$C.F = -\frac{1}{n} \sum_{i=1}^n \sum_{j=1}^K \gamma_{ij} \log(\hat{\gamma}_{ij})$$

Activation Functions

Relu

$$f(x) = \max(0, x)$$

if

$$\begin{cases} x, & x > 0 \\ 0, & x \leq 0 \end{cases}$$

<https://keras.io/api/layers/activations/#creating-custom-activations>

<https://www.v7labs.com/blog/neural-networks-activation-functions>