classification — sigmoid, Tanh, ReLU
softmax

# Regression — Linear Activation function
we use it on the o/p layer.

$$loss = (Y - \hat{Y}) \quad - \text{single datapoint}$$

$$\text{cost funct} = \frac{1}{n} \sum_{i=1}^{n} (Y_i - \hat{Y})^2 \quad - \text{Batch datapoint}$$

Loss function —

ANN

classification

① Binary Cross Entropy

② Categorical cross entropy
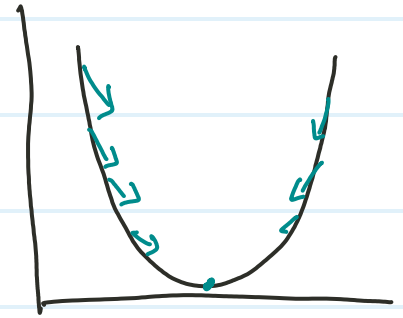
Regression

① MSE

② MAE

③ Huber loss

# Regression

## ① MSE

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

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

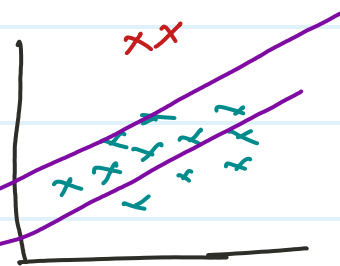$$(a-b)^2 = a^2 - 2ab + b^2$$

$$\boxed{ax^2 + bx + c}$$

quadratic eqn

### Adv.

① MSE is differentiable
② It has one local and one global minima
③ It converg faster

### Disadv

① Not Robust to outlier.
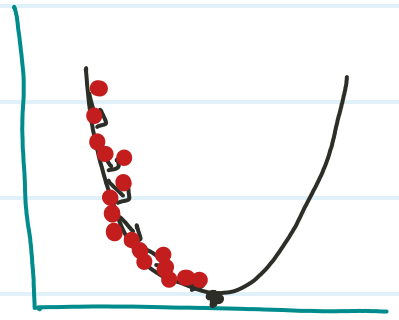
xx

we are just penalizing error

## ② MAE

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

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

**Adv.**

① Robust to outlier condition

**Disadv.**

① Convergen is slow

## ③ Huber Loss —

Combination of MSE and MAE

$$cost fun. = \begin{cases} \frac{1}{2} \frac{1}{n} \sum_{i=1}^{n} (y_i - \hat{y})^2 & \text{if } |y - \hat{y}| \leq \delta \\ \delta |y - \hat{y}| - \frac{1}{2} \delta^2, \end{cases}$$
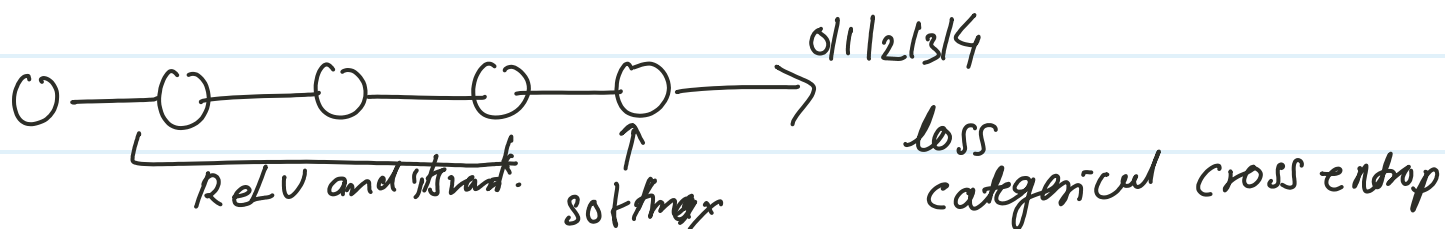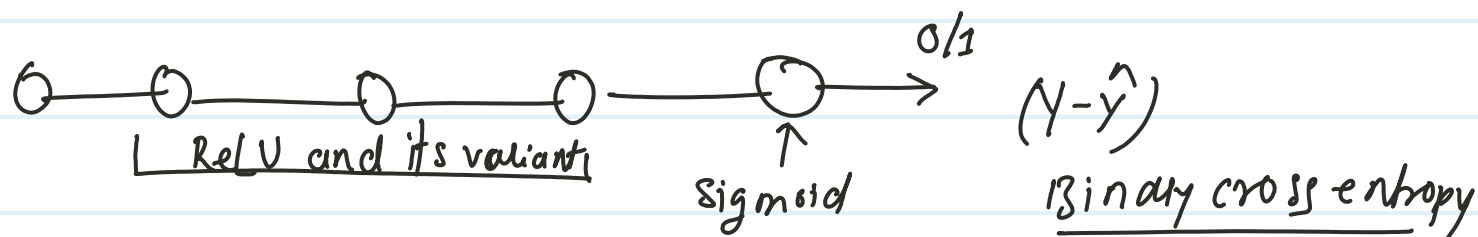
→ when outlier is present

# Loss and cost function for classification

## ① Binary cross Entropy -

$$\text{Loss fn} = -Y * \log(\hat{y}) - (1-y) * \log(1-\hat{y})$$

$$\text{Loss} = \begin{cases} -\log(1-\hat{y}), & \text{if } y = 0 \\ -\log(\hat{y}), & \text{if } y = 1 \end{cases}$$

It only used to find loss in Binary classification



$\lfloor$ ReLU and its variants  Sigmoid  $(Y-\hat{y})$  Binary cross entropy

o/1

o/1/2/3/4  loss  categorical cross entropy

$\lfloor$ ReLU and its variant.  softmax

## ② Categorical cross Entropy

| fi | f2 | f3 | O/p J | j=1 | J=2 | J=3 |
|----|----|----|-------|-----|-----|-----|
| i | 2 | 3 | 4  Good | ! | 0 | 0 |
| | 5 | 6 | 7  Bad | 0 | 1 | 0 |
| | 8 | 9 | 10  Neutral | 0 | 0 | 1 |

No. of category   $c = 3$

$$\text{Loss}(x_i, y_i) = \sum_{j=1}^{C} y_{ij} * \log(\hat{y}_{ij})$$

Actual value $= y_{ij} = \left[ y_1, y_2, y_3 \text{-------} y_c \right]$

$$\left[ y_{21}, y_{22}, y_{23} \text{-----} y_{2c} \right]$$

$$y_{ii} = \begin{cases} 1, & \text{if the element is in the class} \\ 0, & \text{otherwise} \end{cases}$$

it used for multiclass classification
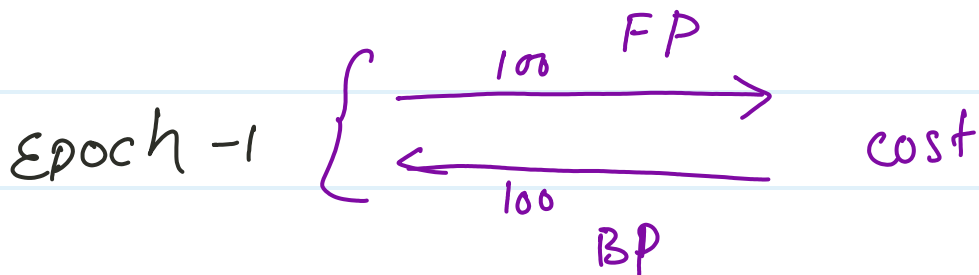
# Optimizer

## ① Gradient Decent Optimizer



weight update formula

$$w_{new} = w_{old} - \eta \boxed{\frac{\partial E}{\partial w_{old}}} - \text{slop}$$

## Epoch —

$\underline{10000}$ datapoint

traning batch $\frac{10000}{100} =) \underline{100}$

$$\text{Epoch -1} \begin{cases} \xrightarrow[]{100 \quad FP} \\ \xleftarrow[100]{} \end{cases} \quad cost$$

$$BP$$

Epoch $=$ I reacheel on global minima

| Adv. | Dis adv. |
|---|---|
| ① Conversion will happen | ① Require huge ram and GPU |

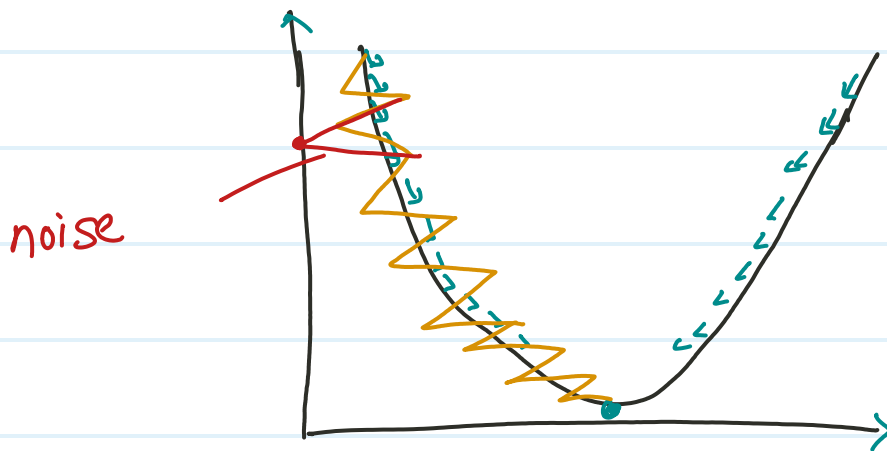| | |
|---|---|
| 100 Data — 20% | 10000 Data over capacity |

## ② SGD (stochastic Gradient Decent)

It use epoch to train batch data as a optimizer
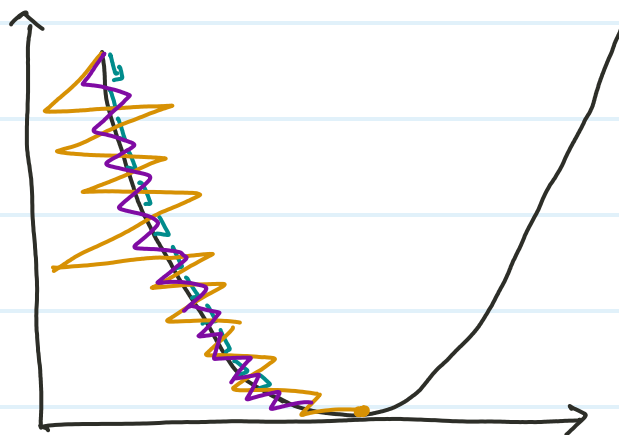
noise

| Adv | Disadv. |
|---|---|
| ① solve Resource issue | ① Convergen will take time |
| | ② noise will get introduce. |

## ③ mini batch SGD



**Adv.**

① conversion speed will increas

② noise will be less

③ Efficient for resouce

**Dis adv.**

① Noise still exist

## ④ SGD with momentum

$$W_{new} = W_{old} - \eta \frac{\partial L}{\partial w_{old}}$$
$$b_{new} = b_{old} - \eta \frac{\partial L}{\partial b_{old}}$$

$\Big\}$ original formal

SGD

$$W_t = W_{t-1} - \eta \frac{\partial L}{\partial w_{t-1}}$$

$$b_t = b_{t-1} - \eta \frac{\partial L}{\partial t-1}$$

# ☆ Exponetial weight Averuge

suppose

| time | $t_1$ | $t_2$ | $t_3$ | $t_4$ | $t_5$ | - - - - - | $t_n$ |
|------|-------|-------|-------|-------|-------|-----------|-------|
| value | $a_1$ | $a_2$ | $a_3$ | $a_4$ | $a_5$ | - - - - - | $a_n$ |

time value $\quad V_{t_1} = a_1$
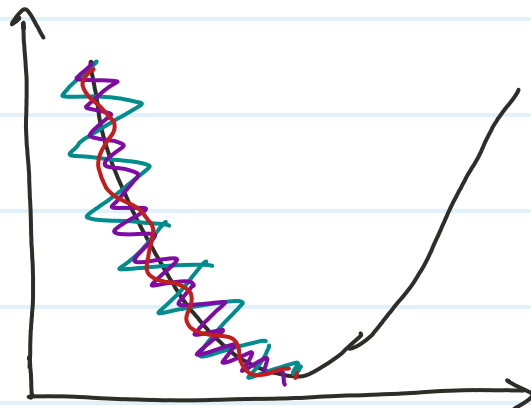
$$V_{t_2} = \beta * V_{t_1} + (1-\beta) a_2$$

$\beta = 0.95$

$$= 0.95 * a_1 + (1 - 0.95) a_2$$

$$V_{t_3} = \beta * V_{t_2} + (1-\beta) a_3$$

Exponetial weight Averuge

$$W_t = W_{t-1} - \eta \, Vd\omega$$

$$\boxed{Vd\omega_t = \beta * Vd\omega_t + (1-\beta) * \frac{dL}{d\omega_{t-1}}}$$

| Adv. | Disadv. |
|------|---------|
| ① Reduce the noise | ① We do not have dynamic learning rate |
| ② smoothen the noise | |
| ③ Quick conversion | |
| ④ working for mini batch | |

⑤ **Adagrad** (Adaptive Gradient Decent)

$$w_t = w_{t-1} - \eta \frac{dL}{dw_{t-1}}$$

$$w_t = w_{t-1} - \eta' \frac{dLoss}{dw_{t-1}}$$

$$\eta' = \frac{\eta}{\sqrt{\alpha_t + \epsilon}}$$

denominator should not become zero

$$\alpha_t = \sum_{i=1}^{t} \left( \frac{dL}{dw_t} \right)^2$$

t is current time stamp

Disadv.

① initialy faster conversion after few time become slow.

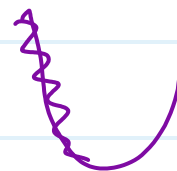② $\eta'$ = possibility to become small value $\approx 0$

⑥ **Adadelta and Rms prop**

We are bringing Exponetial weighted Avg.
in learning.

⑦ **Adam optimizer**

This is combination of SGD with momentum
and Rms prop

$$w_f = w_{f-1} - \eta' \, Vdw$$

$$Vdw_t = \beta * Vdw_{t-1} + (1-\beta) \frac{\partial Loss}{\partial w_{t-1}}$$