**2.1** Change in the learning rate

$$W_t = W_{t-1} - \eta' \cdot \frac{\partial L}{\partial w}$$

$$\eta' = \frac{\eta}{\sqrt{v_t + \varepsilon}} \quad \longrightarrow \quad \text{very small values}$$

$$V_t = V_{t-1} + \left(\frac{\partial L}{\partial w}\right)^2$$

$$V_t = \sum_{i=1}^{n} (\partial L / \partial w)^2$$

$$V_t = \sum_{i=1}^{n} \left(\frac{\partial L}{\partial w}\right)^2 \qquad \left| \frac{\partial L}{\partial w_1}^2 + \frac{\partial L}{\partial w_2}^2 \cdots + \frac{\partial L}{\partial w_n} \right|$$

Adaptive adagt

**3.1** Changing in both GC and GR.

R

⑦ root ... ...ty  (RMS property).

$$W_t = W_{t-1} - \eta' \cdot \frac{\partial L}{\partial w_{t-1}}$$

$$\eta' = \eta \cdot \frac{g}{\sqrt{V_t + \epsilon}}$$

$$V_t = \beta \cdot V_{t-1} + (1-\beta) \cdot \left(\frac{\partial L}{\partial w}\right)^2$$

Adam → Adaptive Momentum
RMSprop ↑

Adv. of RMSprop: Updating gradient values simultaneous.

$$W_t = W_{t-1} + \eta' \cdot \frac{\partial L}{\partial w}$$

$$W_t = W_{t-1} + \eta' \cdot m_t$$

$$m_t = \beta \cdot m_{t-1} + (1-\beta) \cdot \frac{\partial L}{\partial w_{t-1}}$$

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

$$V_t = \beta \cdot v_t + (1-\beta)\left(\partial L / \partial w\right)^2$$

$$\hat{m_t} = \frac{m_t}{1 - \beta_1^t} \quad ; \quad \hat{v_t} = \frac{v_t}{1 - \beta_2^t}$$

Date: 25/02/202

# Training tips in Deep Neural Network

## Normalization

**1) Min-Max Normalization**

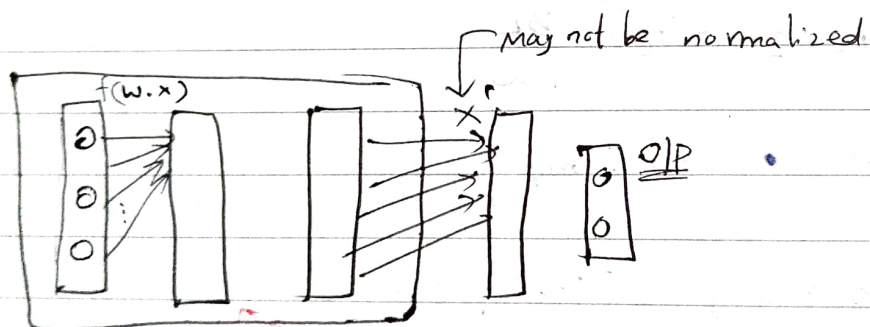$$\boxed{\frac{X - X_{min}}{X_{max} - X_{min}}} \qquad (0-1 \text{ range})$$

**2) Z-score Normalization**

$\downarrow$

$N(0,1)$ $\qquad \dfrac{X - \mu}{\sigma}$

$N(\mu, \sigma)$

**3) Batch Normalization**
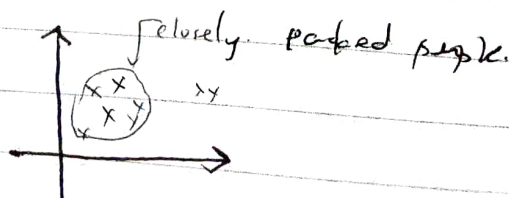


Black box — Normalization done at this layer also

'X' is normalized, but when $w.x$ and non-linear function is used, it (Say, X'here) may no longer be normalized.

"Internal covariate shift" occurs when samples of one distribution has difference within the class.

Eg:- Actual class: cats

When Plotted: Black and orange cats are separate.

It can be avoided if normalized data is used.



closely packed peak.

For mini-batch $B = \{X_1, X_2, ..., X_B\}$

$O/P: y = BN_{\gamma, \beta}(X_B)$

$\overset{*}{x} = BN_{\gamma, \beta}(X_B)$

$x'$ is considered as I/P to next hidden layer.

$$\mu_B = \frac{1}{B} \sum_{i=1}^{B} x_i$$

$$\sigma_B = \frac{1}{B} \sum_{i=1}^{B} (x_i - \mu_B)^2$$

$$\boxed{x_i' = \frac{x_i - \mu_B}{\sqrt{\sigma_B + \varepsilon}}}$$

$$\boxed{x_i' = \gamma x_i' + \beta}$$
<span style="color:red">Scaling ←      → Shifting</span>

## OVERFITTING AND UNDERFITTING :-

Mutually exclusive datasets for training, validation and testing.

Solution : regularization

## Regularization :-

→ used for better model generalization.
→ General cost function with regularization for training is defined as

<span style="color:red">cost function = Loss + Regularization Term.</span>