## Weight Initialization:
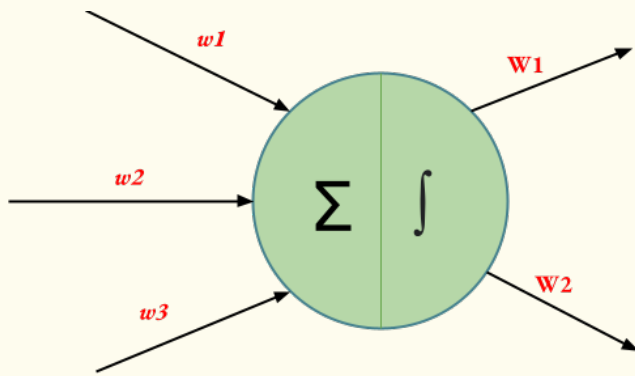
→ to tackle with vanishing and Exploding gradient problem.

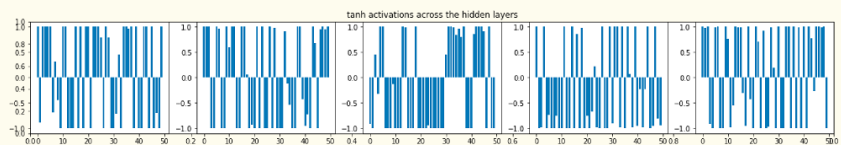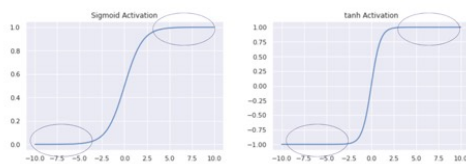→ faster convergence



- fan_in = 3
- fan_out = 2

## 1. zero initialization:

→ highly ineffective as neurons learn the same feature during first few iterations

→ slow convergence

## 2. Random Initialization:

→ better then zero initialization or const initialization.

→ very high or low value may lead to vanishing or Exploding Gradient problem.



For random values with relatively larger magnitude, the tanh and sigmoid activations get saturated!

Sigmoid Activation    tanh Activation

tanh activations across the hidden layers

→ saturating tanh activations for large random weights.

a) random uniform $\Rightarrow w_i \sim U(0,1)$
b) random normal $\Rightarrow w_i \sim N(0,1)$
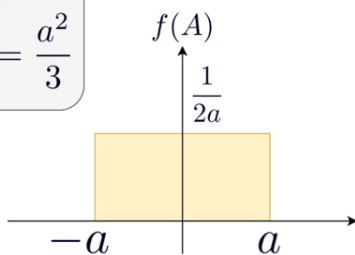
## 3. Xavier/Glorot Initialization:

→ $w_i \sim U\left[-\sqrt{\dfrac{\sigma}{f_{in}+f_{out}}}, \sqrt{\dfrac{\sigma}{f_{in}+f_{out}}}\right]$

→ uniform distribution

→ takes account of input and output counctions

→ mostly used for sigmoid.

→ weights do not saturate or vanish. during forward pass.

### a) Glorot Normal:

$$W \sim N(0,\sigma) \qquad \sigma = \sqrt{\dfrac{6}{f_{in}+f_{out}}}$$

### b) Glorot uniform:

$$\boxed{\begin{aligned}\mathbb{E}(A) &= 0 \\ Var(A) &= \frac{(2a)^2}{12} = \frac{a^2}{3}\end{aligned}}$$

$f(A)$

$\frac{1}{2a}$

$-a \qquad a$

$$Var(A) = \frac{(2a)^2}{12} = \frac{a^2}{3}$$

We know that the variance should be equal to $\frac{2}{fan_{in}+fan_{out}}$; we can work backward to find the endpoints of the interval.

$$Var(w) = \frac{(2a)^2}{12} = \frac{a^2}{3}$$

We have,

$$Var(w) = \frac{2}{fan_{in}+fan_{out}}$$

$$\implies \frac{a^2}{3} = \frac{2}{fan_{in}+fan_{out}}$$

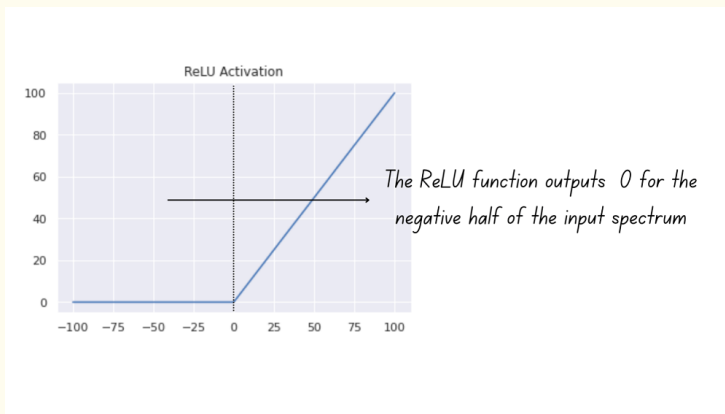$$a^2 = \frac{6}{fan_{in}+fan_{out}}$$

$$\implies a = \sqrt{\frac{6}{fan_{in}+fan_{out}}}$$

$$w \in \mathcal{U}\left[-\sqrt{\frac{6}{fan_{in}+fan_{out}}}, \sqrt{\frac{6}{fan_{in}+fan_{out}}}\right]$$

# 4. He Weight Initialization: (kaiming)

→ mostly used for ReLU and Leaky-ReLU.

→ helps avoiding slow convergence

→ avoid oscillations while reach--ing minima.



ReLU Activation

The ReLU function outputs 0 for the negative half of the input spectrum

a) Normal:

$$W_i \sim N(0, \sigma)$$

$$\sigma = \sqrt{\frac{2}{fin}}$$

b) Uniform:

$$W_i \sim U\left[-\sqrt{\frac{6}{fin}}, \sqrt{\frac{6}{fin}}\right]$$

* Another method to prevent Exploding gradient problem is — Gradient Clipping
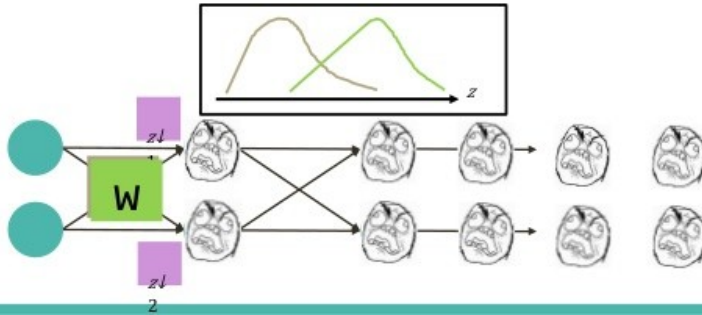
# 💡 Batch Normalization

→ faster convergence
- ↳ increase in LR
- ↳ less oscillation
- ↳ reduce training time.

→ remove necessity dropout

→ prevents from vanishing or exploding gradient problem.

→ saves from dead activation for sigmoid & tanh.

→ reduce dependency on Hyper-parameters.

→ saves NN from Internal covariate shift

→ smoothens the loss function curve.

## Internal Covariate Shift

- During training, layers need to continuously adapt to the new distribution of their inputs

**ICS:** change in the distribution of the network activations output or feature map due to change in w/w parameters during training back propagation.

**solution:**

**Input:** Values of $x$ over a mini-batch: $\mathcal{B} = \{x_{1...m}\}$;
Parameters to be learned: $\gamma, \beta$
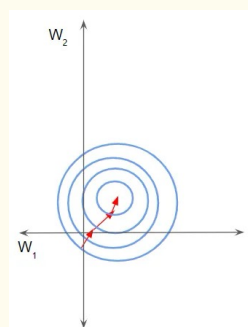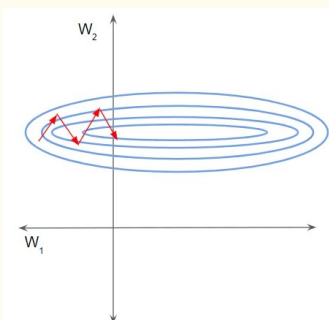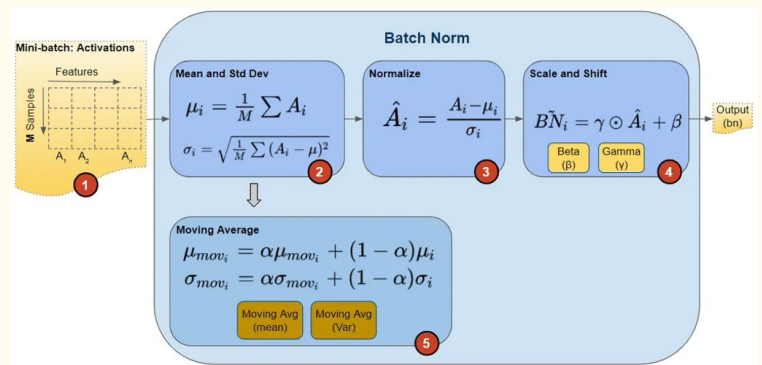**Output:** $\{y_i = BN_{\gamma,\beta}(x_i)\}$

$$\mu_\mathcal{B} \leftarrow \frac{1}{m}\sum_{i=1}^{m} x_i \qquad \text{// mini-batch mean}$$

$$\sigma_\mathcal{B}^2 \leftarrow \frac{1}{m}\sum_{i=1}^{m}(x_i - \mu_\mathcal{B})^2 \qquad \text{// mini-batch variance}$$
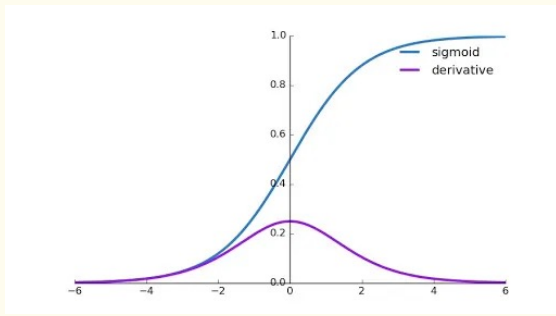
$$\widehat{x}_i \leftarrow \frac{x_i - \mu_\mathcal{B}}{\sqrt{\sigma_\mathcal{B}^2 + \epsilon}} \qquad \text{// normalize}$$

$$y_i \leftarrow \gamma\widehat{x}_i + \beta \equiv BN_{\gamma,\beta}(x_i) \qquad \text{// scale and shift}$$

**Algorithm 1:** Batch Normalizing Transform, applied to activation $x$ over a mini-batch.

$\longrightarrow$ smoothens the loss surface that leads to faster convergence.

sigmoid
derivative

→ prevent from
saturation



**Batch Norm**

Activations (a)

*Learnable Params* — Beta (β), Gamma (γ)

*Saved Params* — Moving Avg (mean), Moving Avg (Var)

Output (bn)



Figure 2. Single crop validation accuracy of Inception and its batch-normalized variants, vs. the number of training steps.

Legend:
- Inception
- BN–Baseline
- BN–x5
- BN–x30
- BN–x5–Sigmoid
- ◆ Steps to match Inception