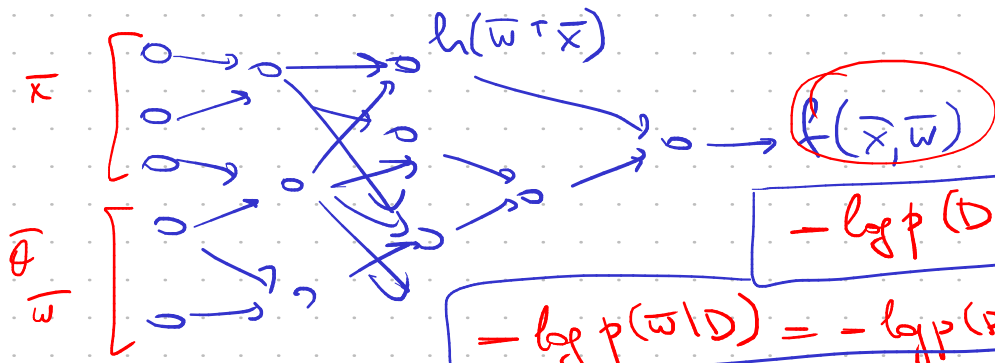


① Gradient descent



$$-\log p(D|\bar{w}) = -\sum_n \log p(x_n|\bar{w})$$

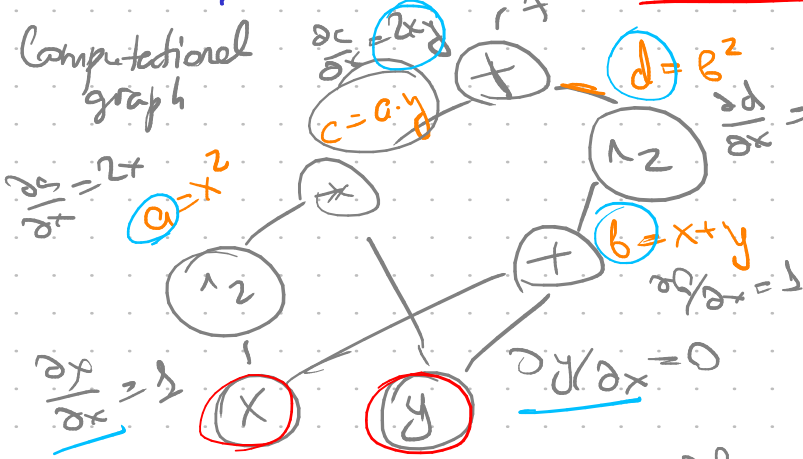
$$-\log p(\bar{w}|D) = -\log p(\bar{w}) - \log p(D|\bar{w})$$

$$\frac{\partial f}{\partial x} = \frac{\partial c}{\partial x} + \frac{\partial d}{\partial x} = 2xy + 2b = 2xy + 2(x+y)$$

$$\bar{w} := \bar{w} - \eta \cdot \nabla_{\bar{w}} F$$

$$F(\bar{w}) \xrightarrow{\bar{w}} \min$$

Computational graph



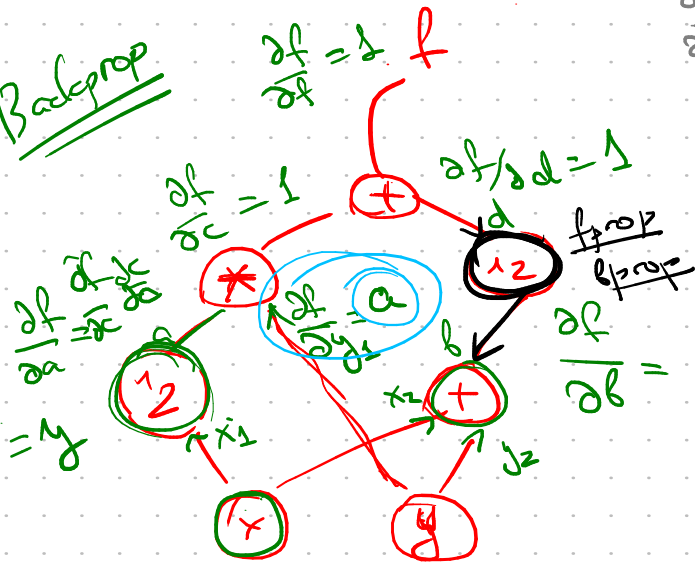
$$f(x, y) = x^2 y + (x+y)^2$$

$$\nabla f = \left(\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y} \right)$$

$$(x^{(k)}, y^{(k)})$$

$$\frac{\partial f}{\partial x} = \frac{\partial c}{\partial x} + \frac{\partial d}{\partial x} = 1 + 0 = 1$$

Backprop



automatic differentiation libraries

PyTorch
TensorFlow
theano

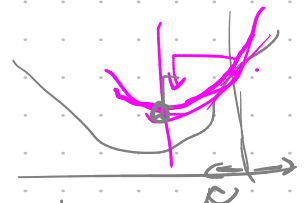
$$\frac{\partial f}{\partial y} = \frac{\partial f}{\partial d} \frac{\partial d}{\partial y} = 1 \cdot 2b = 2b$$

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial a} \frac{\partial a}{\partial x} + \frac{\partial f}{\partial b} \frac{\partial b}{\partial x} = y \cdot 2x + 2b \cdot 1 = 2xy + 2b$$

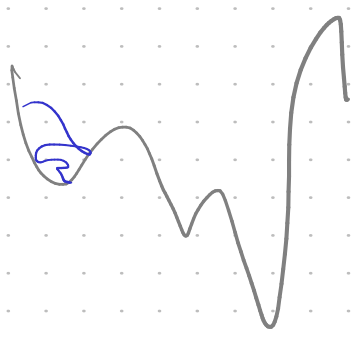
$$\frac{\partial f}{\partial y} = \frac{\partial f}{\partial b} \frac{\partial b}{\partial y} + \frac{\partial f}{\partial c} \frac{\partial c}{\partial y} = 2b \cdot 1 + 1 \cdot a = a + 2b$$

GD:

$$\bar{w} := \bar{w}^{(k)} - \eta \cdot \nabla_{\bar{w}} F |_{\bar{w}^{(k)}}$$



$$F(\bar{w}) \approx F(\bar{w}^{(k)}) + (\bar{w} - \bar{w}^{(k)})^T \nabla_{\bar{w}} F |_{\bar{w}^{(k)}} + \frac{1}{2} (\bar{w} - \bar{w}^{(k)})^T \text{Hess} F |_{\bar{w}^{(k)}} (\bar{w} - \bar{w}^{(k)})$$



$$+ \frac{1}{2} (\bar{w} - \bar{w}^{(k)})^T H_k (\bar{w} - \bar{w}^{(k)}) \quad \bar{w} \rightarrow \min$$

$$\bar{g}_k + H_k (\bar{w}_* - \bar{w}^{(k)}) = 0 \quad \text{second order method}$$

$$\bar{w}_* = \bar{w}^{(k)} - H_k^{-1} \bar{g}_k$$

quasi-Newton algorithms l-BFGS $H^{-1} \approx$ low-rank approx (\bar{w}_k, \bar{g}_k)

$$F(\bar{w}) = \sum_{n=1}^N f(\bar{x}_n, \bar{w}) \quad \mathbb{E}_{\text{unif}(\bar{x}_n)} [f(\bar{x}_n, \bar{w})]$$

SGD - stochastic GD

$$F(\bar{w}) = \mathbb{E}_{q(\bar{x})} [f(\bar{x}, \bar{w})] \quad \bar{w} \rightarrow \min$$

$$G(\bar{w}) = \mathbb{E}_{q(\bar{x})} [\nabla_{\bar{w}} f(\bar{x}, \bar{w})]$$

$$\mathbb{E}_{p(\bar{x})} [f(\bar{x})] \approx \frac{1}{R} \sum f(\bar{x}_i) \quad \bar{x}_i \sim p(\bar{x})$$

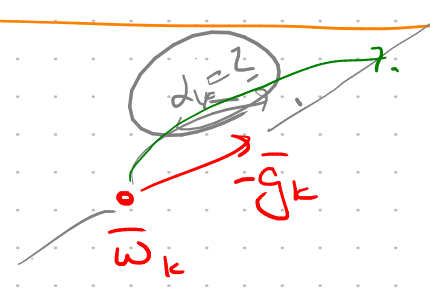
$$\hat{F}(\bar{w}) = \frac{1}{m} \sum_{i=1}^m f(\bar{x}_i, \bar{w}) \quad \text{mini-batch}$$

$$\hat{g}(\bar{w}) = \frac{1}{m} \sum_{i=1}^m \nabla_{\bar{w}} f(\bar{x}_i, \bar{w})$$

$$\text{SGD: } \bar{w}_{k+1} = \bar{w}_k - \gamma \hat{g}_k$$

$$\text{GD: } \bar{w}_{k+1} = \bar{w}_k - \alpha \bar{g}_k$$

$$\varphi_k(\alpha) = F(\bar{w}_k - \alpha \bar{g}_k) \quad \alpha \rightarrow \min$$



② Baruanan SGD $F(\bar{x}) \quad f(x, y) = x^2 + y - y^2$

$$\bar{x}_t \xrightarrow{-\bar{g}_t} \bar{x}_{t+1} = \bar{x}_t - \alpha \cdot \bar{g}_t$$

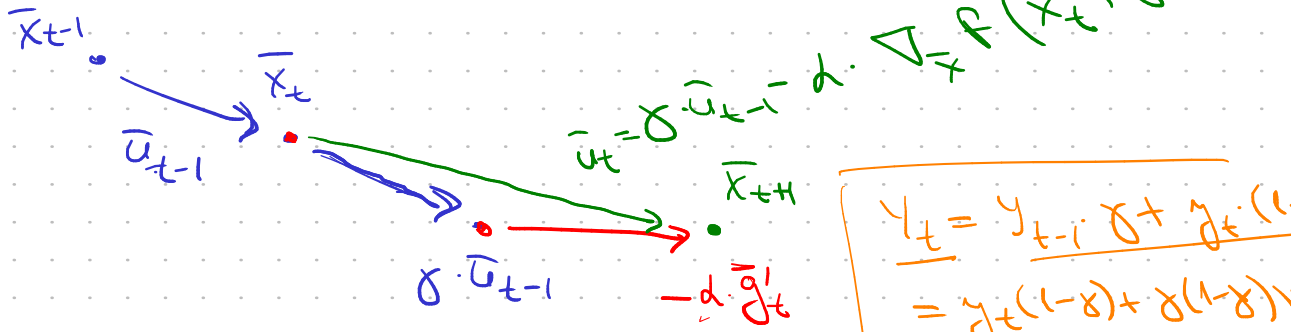
$$\bar{u}_t = \bar{x}_{t+1} - \bar{x}_t = -\alpha \cdot \bar{g}_t$$

1) SGD w/ momentum

$$\gamma = 0.99, 0.995, 0.95$$

$$\bar{u}_t = \gamma \bar{u}_{t-1} - \alpha \bar{g}_t = \gamma \bar{u}_{t-1} - \alpha \nabla_{\bar{x}} f(\bar{x}_t)$$

2) Nesterov accelerated gradients (NAG)



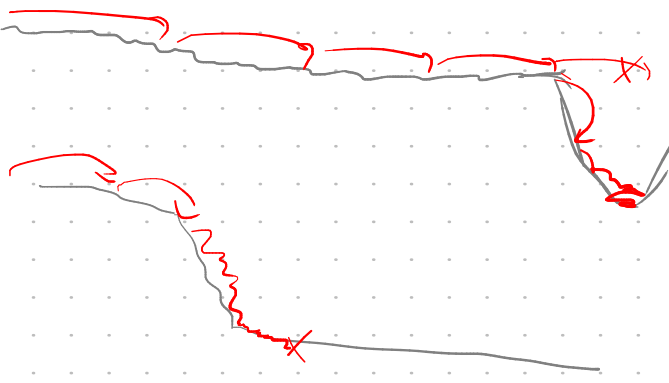
$$\begin{aligned}
 y_t &= y_{t-1} \delta + g_t (1-\delta) \\
 &= y_t (1-\delta) + \delta (1-\delta) y_{t+1} \\
 &+ \delta^2 (1-\delta) y_{t+2} + \dots
 \end{aligned}$$

3) Adaptive SGD

- Adagrad



$$\begin{aligned}
 G_{0,i} &= 0 \\
 G_{t,i} &= G_{t-1,i} + g_{t,i}^2 \\
 x_{t+1,i} &= x_{t,i} - \frac{\alpha}{\sqrt{G_{t,i} + \epsilon}} g_{t,i}
 \end{aligned}$$



- RMSprop

$$G_{t,i} := G_{t-1,i} \cdot \delta + (1-\delta) \cdot g_{t,i}^2$$

- Adadelta

$$\begin{aligned}
 \bar{g} &= \mu_c, \quad \bar{x} = \text{cek}, \quad P(x) = \mu \\
 \bar{x}_c &:= \bar{x}_c - \frac{\alpha}{\sqrt{G_{t,i} + \epsilon}} g_{t,i} \\
 \bar{x}_c &:= \bar{x}_c - \frac{H^{-1} \bar{g}}{(\mu_c)^2} \mu_c = \frac{c^2}{\mu} - \mu_c
 \end{aligned}$$

$$x_{t+1,i} = x_{t,i} - \alpha \frac{\sqrt{R_{t,i} + \epsilon}}{\sqrt{G_{t,i} + \epsilon}} g_{t,i}, \quad R_{t,i} = \beta R_{t-1,i} + (1-\beta) \cdot u_{t,i}^2$$

- Adam (2014)

$$\begin{aligned}
 G_{t,i} &= \beta_2 G_{t-1,i} + (1-\beta_2) g_{t,i}^2 \\
 m_{t,i} &= \beta_1 m_{t-1,i} + (1-\beta_1) g_{t,i}
 \end{aligned}$$

$$x_{t+1,i} = x_{t,i} - \frac{\alpha}{\sqrt{G_{t,i} + \epsilon}} \cdot m_{t,i}$$

$$\beta_1 = 0.9, \quad \beta_2 = 0.999, \quad \epsilon = 10^{-7}$$

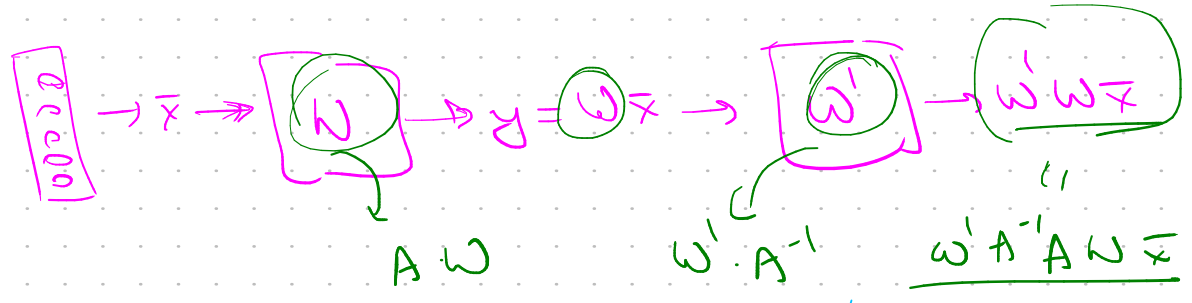
- AdamW

weight decay

$$\begin{aligned}
 \bar{w} &= (1-\eta) \bar{w} - \alpha \cdot \nabla_{\bar{w}} L \\
 &= \bar{w} - \eta \cdot \bar{w} - \alpha \cdot \nabla_{\bar{w}} L(\bar{w}) \\
 &= \bar{w} - \alpha \cdot \nabla_{\bar{w}} \left[L(\bar{w}) + \frac{\eta}{2\alpha} \|\bar{w}\|^2 \right]
 \end{aligned}$$

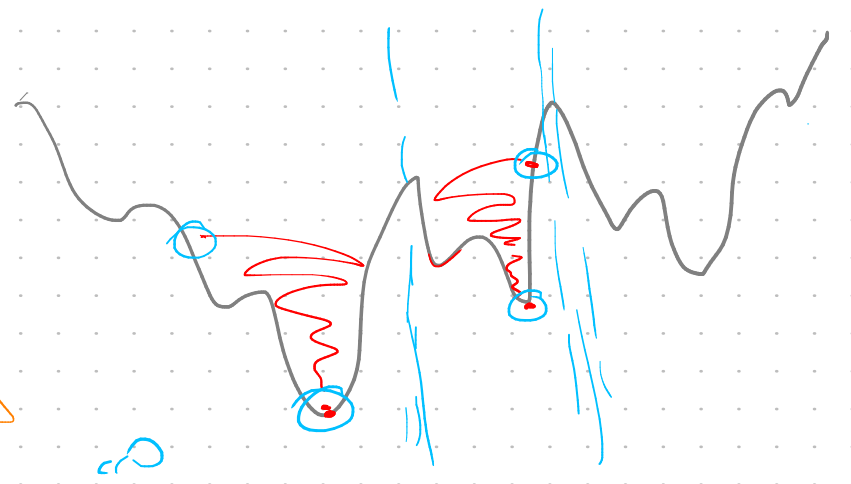
SGD

④ Dropout



⑤ Weight Initialization

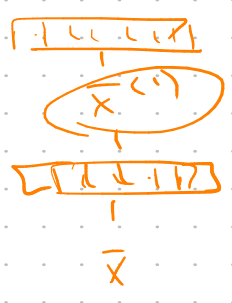
pre-training
unsupervised pre-training



$$y = \bar{w}^T x = \sum_i w_i x_i$$

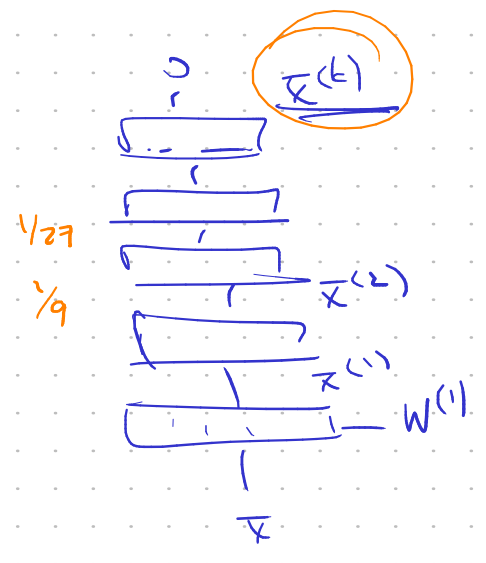
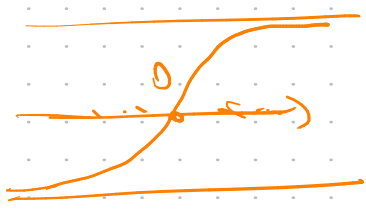
$$\text{Var}[y] = \sum_i \text{Var}[y_i]$$

$$\text{Var}[y_i] = \cancel{E[x_i]^2 \text{Var}[w_i]} + \cancel{E[w_i]^2 \text{Var}[x_i]} + \text{Var}[w_i] \text{Var}[x_i]$$



$$E[y] = 0$$

$$\bar{x}^{(1)} = h(y)$$



$$\text{Var}[y_i] = \text{Var}[w_i] - \text{Var}[x_i]$$

$$\text{Var}[y] = n \cdot \text{Var}[w_i] \cdot \text{Var}[x_i]$$

$$n \cdot \text{Var}[w_i] \approx 1$$

$$\text{Var}[w_i] \approx 1/n$$

$$\text{Var}[w_i] = O(1/n)$$

$$w_i \sim \text{Unit}\left(\left[-\frac{1}{\sqrt{n}}, \frac{1}{\sqrt{n}}\right]\right)$$

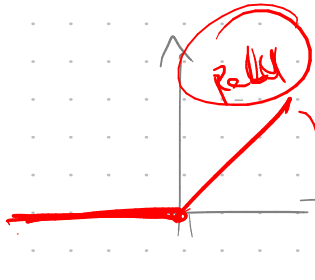
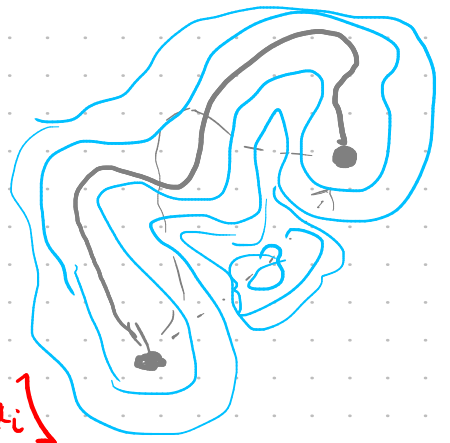
Xavier Glorot (2010)
Bengio

$$\text{Var}(\text{Unit}(a,b)) = \frac{(b-a)^2}{12}$$

$$\text{Var}\left(\text{Unit}\left(\left[-\frac{1}{\sqrt{n}}, \frac{1}{\sqrt{n}}\right]\right)\right) = \frac{(2/\sqrt{n})^2}{12} = \frac{1}{3n}, \quad n \cdot \text{Var}[w_i] \approx \frac{1}{3}$$

Xavier init

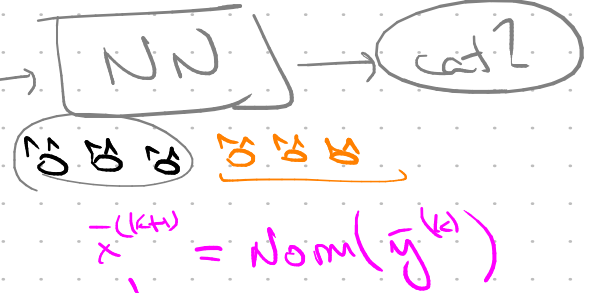
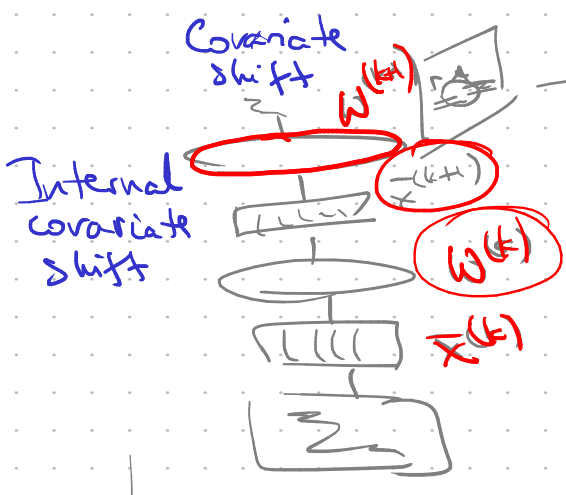
win. unit $([-\sqrt{\frac{3}{n}}, \sqrt{\frac{3}{n}}])$
 $N(0, \frac{1}{n})$



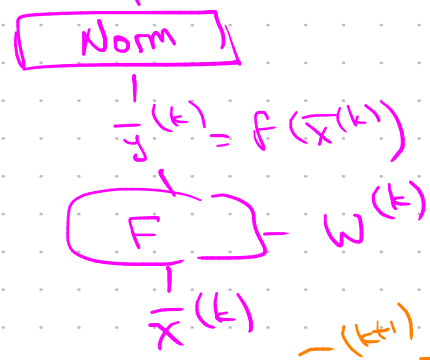
$$\text{Var}[y_i] = E[x_i]^2 \text{Var}[w_i]$$

$$\text{Var}[w_i] = \frac{2}{n} + \text{Var}[w_i] \text{Var}[x_i]$$

⑥ Batch normalization

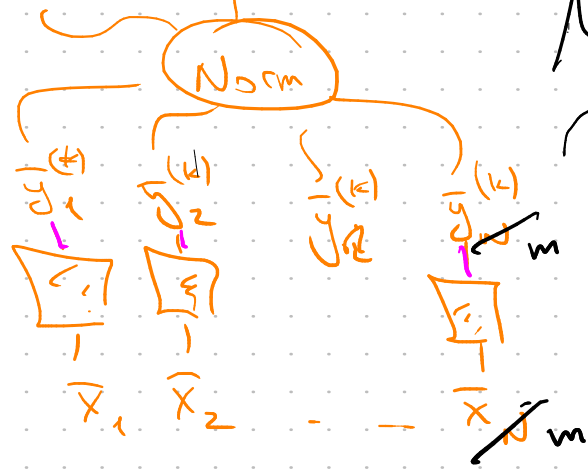
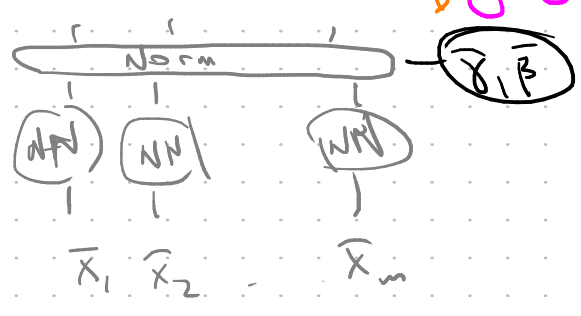


$$\hat{x}^{(k+1)} = \text{Norm}(\hat{y}^{(k)})$$



$$\hat{x}_n^{(k+1)} = \frac{\hat{y}_n^{(k)} - \frac{1}{m} \sum_i \hat{y}_i^{(k)}}{\sqrt{\text{Var}(\dots)}}$$

$$\hat{x}^{(k+1)} = \frac{\hat{y}^{(k)} - E[\hat{y}^{(k)}]}{\sqrt{\text{Var}(\hat{y}^{(k)})}}$$



$$\text{BN}(\hat{y}^{(k)}) = \frac{\hat{y}^{(k)} - E_m[\hat{y}^{(k)}]}{\sqrt{\text{Var}_m[\hat{y}^{(k)}]}} \odot \bar{\gamma} + \bar{\beta}$$

