$$F(W) = \frac{1}{N} \sum_{i=1}^{N} L\left(y_i, a^L(x_i, W)\right) \longrightarrow \min_{W}$$
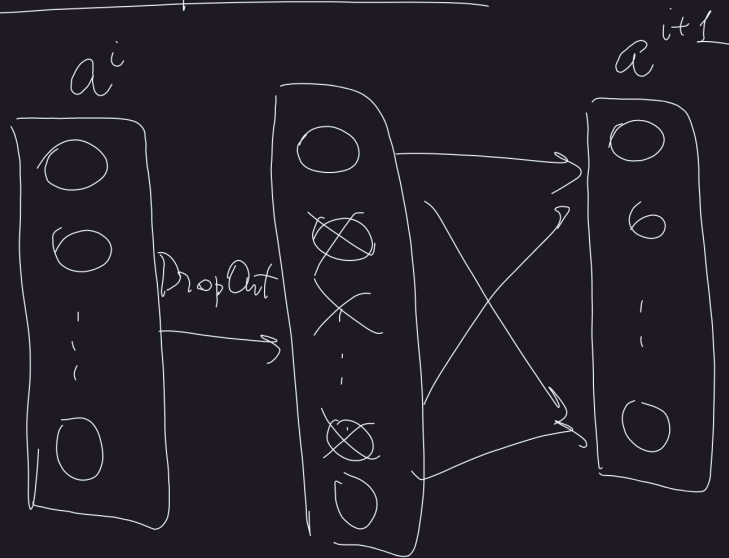
GD:  Init $W_0$
until convergence: $W_{k+1} = W_k - \alpha_t \nabla F(W_k)$

# Regularization

$$\to \quad F(W) = \underbrace{\frac{1}{N} \sum_{i=1}^{N} L(y_i, n^L(x_i, W))}_{L_{data}(W)} + \lambda \sum_{\ell=1}^{L} \| W^\ell \|_F^2 \to \min_W$$

$$f(x) + \sum_i \lambda_i g_i(x) \to \min_x \quad \Longleftrightarrow$$

$$\begin{cases} f(x) \to \min_x \\ g_i(x) \leq \gamma_i \quad \forall i \end{cases}$$
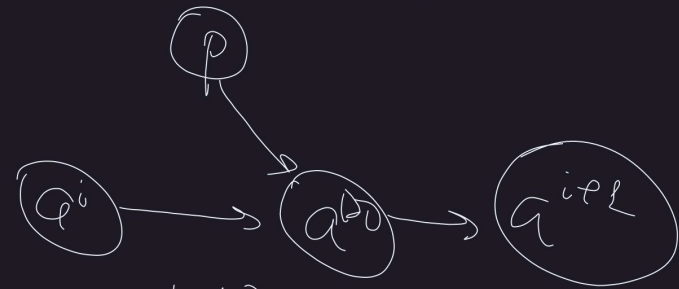
# Drop Out

$a^i$ $a^{i+1}$



DropOut

$y_j$ is diff. for diff. mini-batches

$$\mathbb{E}_{y_j \sim \text{Bern}} a_j^{DO} = \mathbb{E} y_j \left( a_j^i \right) = a_j^i (1-p)$$

$$a_j^{DO} = y_j \cdot a_j^i$$

$$y_j \sim \text{Bern}(y \mid 1-p) : \begin{array}{cc} 0 & 1 \\ p & 1-p \end{array}$$



inverted DO

$$a_j^{invDO} = \frac{1}{1-p} y_j \cdot a_j^i$$

$$\mathbb{E} a_j^{invDO} = a_j^i$$

Gradient Penalty reg. or $R_1$ reg.

$$F(W) = L_{data}(W) + \lambda \| \nabla_W L_{data}(W) \|_2^2 \to \min_W$$

$L_{data}(W)$

narrow loc. min.    wide loc. min. $W$

# Batch Normalization
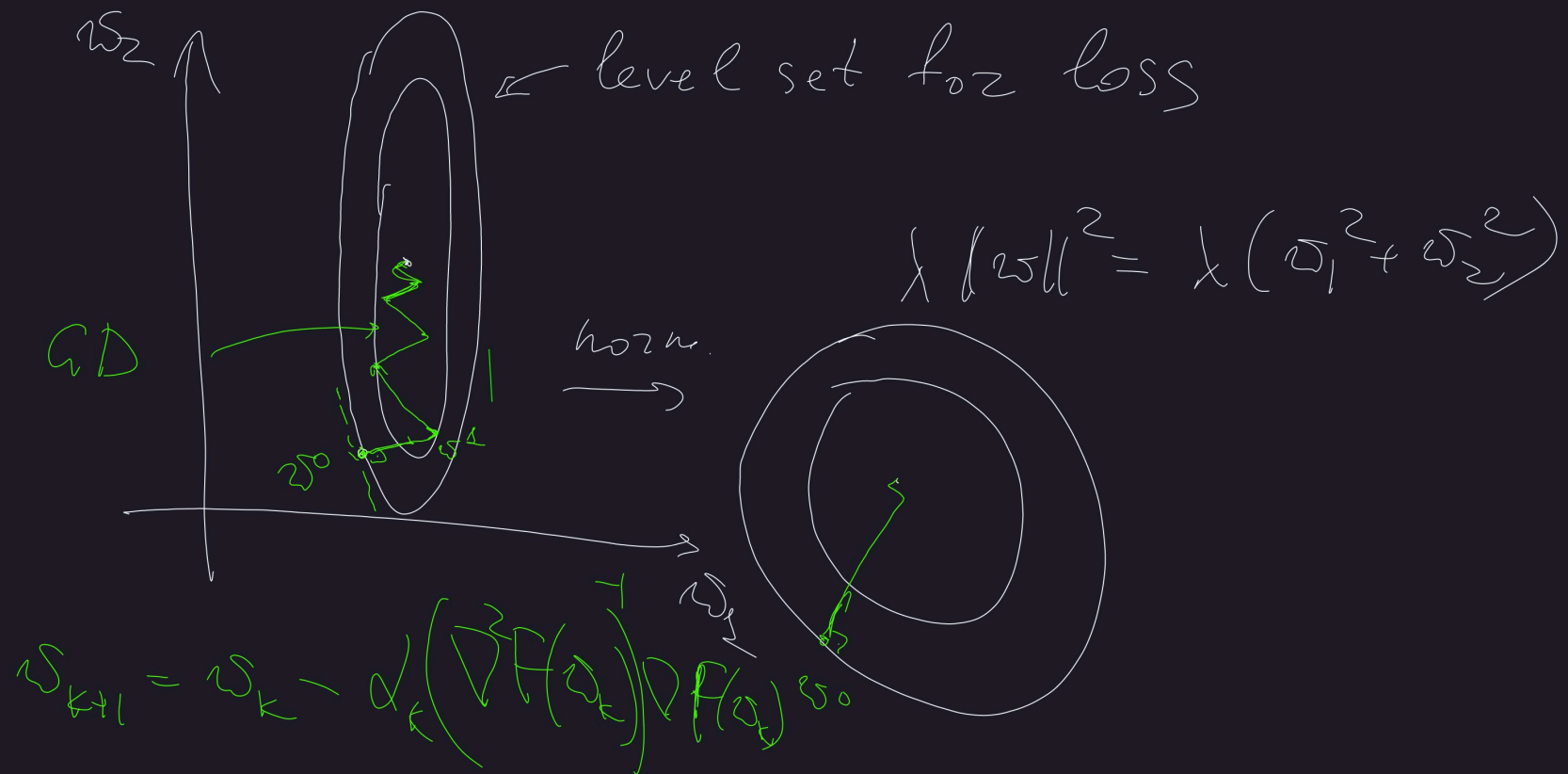
$$fl(x) = S \cdot M \cdot 2^E$$

$$\{x_1, \ldots, x_N\} \xrightarrow{\text{norm.}} \{\hat{x}_1, \ldots, \hat{x}_N\}$$

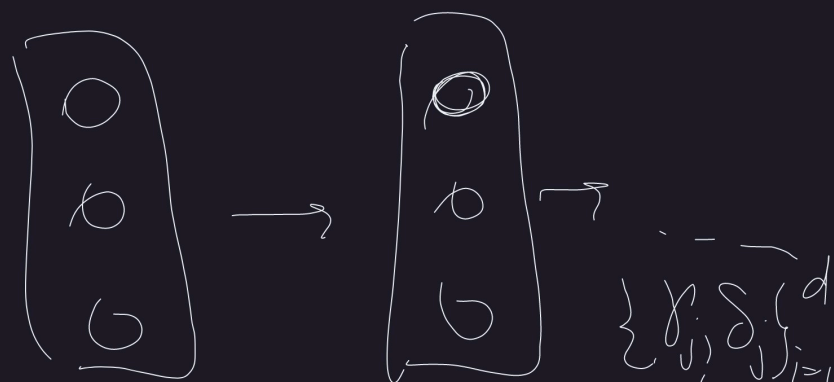$$\hat{x}_{ij} = \frac{x_{ij} - m_j}{\sqrt{\sigma_j^2 + \varepsilon}}, \qquad m_j = \frac{1}{N} \sum_{i=1}^{N} x_{ij}$$

$$\sigma_j^2 = \frac{1}{N} \sum_{i=1}^{N} \left(x_{ij} - m_j\right)^2$$

$$w^T x = \sum_{j=1}^{d} w_j x_j = \sum_j \hat{w}_j \hat{x}_j$$

$w_2$

← level set for loss

GD

$\lambda \|w\|^2 = \lambda(w_1^2 + w_2^2)$

norm.

$2^0$

$w_1$

$w_{k+1} = w_k - \alpha_k \left(\nabla F(w_k)\right) \nabla F(w_k) \, \text{so}$

$$M_j = \frac{1}{N_{batch}} \sum_{i=1}^{N_{batch}} x_{ij}$$

$$\sigma_j^2 = \frac{1}{N_{batch}} \sum_i \left( x_{ij} - M_j \right)^2$$

$$\{\gamma_j, \delta_j\}_{j=1}^d$$

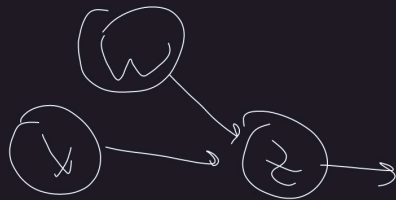$$y_{ij} = \frac{x_{ij} - m_j}{\sqrt{\sigma_j^2 + \varepsilon}}, \quad \hat{x}_{ij} = \gamma_j y_{ij} + \delta_j$$

$$\{x_{ij}\}_{i,j=1}^{N_{batch}, d} \longrightarrow \boxed{BN} \longrightarrow \{\hat{x}_{ij}\}_{ij=1}^{N_{batch}, d}$$

$$M_j^{new} = \alpha M_j + (1 - \alpha) m_j^{prev}$$

$$\gamma_j = \sqrt{\sigma_j^2 + \varepsilon}$$

$$\delta_j = m_j$$

# Weight init.

$x, \quad z = Wx$



$\nabla_z L$

$(\nabla_x L)_j = \sum_{i=1}^{n_{outputs}} W_{ij} (\nabla_z L)_i$

$dL = \nabla_z L^\top dz =$

$= \underbrace{\nabla_z L^\top W}_{\nabla_x L^\top} dx$

$Var(W_{ij}) = \dfrac{1}{n_{outputs}}$

$\nabla_x L = W^\top \nabla_z L$

$Var((\nabla_x L)_j) = n_{outputs} \, Var(W_{ij}) Var((\nabla_z L)_i)$

$$z = Wx$$

$$z_i = \sum_j W_{ij} x_j$$

$$x_j \sim N(x_j | 0, 1)$$

$$W_{ij} \sim N(W_{ij} | 0, Var(W_{ij}))$$

Xavier Glorot

Kaiming He

$$Var(W_{ij}) = \frac{1}{n_{inputs}}$$

$$Var(W_{ij}) = \frac{2}{n_{inputs} + n_{outputs}}$$

$$\mathbb{E} z_i = \mathbb{E} \sum_j W_{ij} x_j = \sum_j \mathbb{E} W_{ij} \cdot \underbrace{\mathbb{E} x_j}_{=0} = 0$$

$$Var(z_i) = Var\left(\sum_j W_{ij} x_j\right) = \sum_{j=1}^{n_{input}} Var(W_{ij}) \cdot \overbrace{Var(x_j)}^{iid} = n_{inputs} \cdot Var(W_{ij})$$

$$W_{ij} \sim R\left[-\frac{1}{\sqrt{n}}, \frac{1}{\sqrt{n}}\right]$$

$$Var(W_{ij}) = \frac{1}{3n} = \frac{2}{n_{inputs} + n_{outputs}}$$

$$\Rightarrow n = \frac{n_{inputs} + n_{outputs}}{6} \quad , \quad W_{ij} \sim R\left[-\sqrt{\frac{6}{n_{inputs} + n_{outputs}}}, \sqrt{\frac{6}{\cdots}}\right]$$

$$F(W) \longrightarrow \min_{W \in Orth.}$$

## Orth. init.

$$W \in \mathbb{R}^{n \times d}$$

$$W^T W = I$$

$$z = Wx$$

$$\|z\|_2 = \|Wx\|_2 = \|x\|_2$$

$$\nabla_x L = W^T \nabla_z L$$

$$\|\nabla_x L\|_2 = \|W^T \nabla_z L\|_2 \stackrel{n=d}{=} \|\nabla_z L\|_2$$

$$f(x) \to \min_{x \in \Omega}$$

$$f(g(y)) \to \min_y$$

$$f((x+y)^2) \to \min_{x,y}$$

$$x \in \Omega \iff x = g(y), \quad y \in \mathbb{R}^n$$

$$[a,b] \ni g(y) = a + (b-a)\sigma(y)$$

$$g(y) = \log(1 + \exp(y))$$

$$g(y) = \exp(y) > 0$$

$$g(y) = y^2$$

$$W = expm(V - V^\top)$$

$$p(y|x) = N(y \mid NN_1(x, \theta), NN_2(x, \theta))$$