$$\frac{1}{N} \sum_{i=1}^{N} L\left(y_i, f(x_i, \Theta)\right) + \frac{\lambda}{2} \|\Theta\|_2^2 \longrightarrow \min_{\Theta}$$

$$\Big\uparrow$$

$$\begin{cases} \frac{1}{N} \sum_{i=1}^{N} L\left(y_i, f(x_i, \Theta)\right) \longrightarrow \min_{\Theta} \\ \|\Theta\|_2^2 \leq \eta \end{cases}$$

# Drop Out

$z = Wx + b$

$a = g(z)$

$y_i \sim Bern$

$a_i^{DO} = a_i \cdot y_i$

|   | 0 | 1 |
|---|---|---|
|   | $p$ | $1-p$ |

$\mathbb{E} \, a_i^{DO} = \mathbb{E}_{y_i \sim Bern} \, a_i \cdot y_i = p(0 \cdot a_i) + (1-p) 1 \cdot a_i = a_i(1-p)$

## Classic DO

Train: $y_i \sim Bern, \; a_i^{DO} = a_i \cdot y_i$

Test: $a_i^{DO} = \dfrac{a_i}{1-p}$

## Inverted DO:

Train: $y_i \sim Bern, \; a_i^{DO} = \dfrac{1}{1-p} a_i \cdot y_i$

Test: $a_i^{DO} = a_i$

# Batch Normalization

$W_2$ (small scale feature)

Loss

GD

Norm →

$W_2$

Loss

GD

$W_1$

$W_1$ (large scale feature)

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

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

$$m_j = \frac{1}{N_{batch}} \sum_{i=1}^{N_{batch}} x_{ij} \qquad \forall j$$

$$\sigma_j^2 = \frac{1}{N_{batch}} \sum_i (x_{ij} - m_j)^2 \qquad \forall j$$

$$z_{ij} = \frac{x_{ij} - m_j}{\sqrt{\sigma_j^2 + \varepsilon}}$$

$$\hat{x}_{ij} = \gamma_j z_{ij} + \delta_j$$

Train: compute $m_j, \sigma_j$
on mini-batch and
learn $\gamma_j, \delta_j$

Test: use saved $m_j, \sigma_j$
for trainset

$$\text{Image } x \in \mathbb{R}^{H \times W \times C}$$

$$\text{Batch } \in \mathbb{R}^{B \times H \times W \times C}$$

BatchNorm: for each $C$ average $B \times H \times W$

LayerNorm: for each $B$ average $H \times W \times C$

InstanceNorm: for each $B, C$ average $H \times W$

# Weight Initialization

$$z = Wx$$

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

$$x_j \sim N(0, 1), \quad w_{ij} \sim N(0, Var(W))$$

$$z_i = \sum_{j=1}^{n_{input}} w_{ij} x_j \; ; \quad E z_i = E \sum w_{ij} x_j = \sum E w_{ij} x_j =$$

$$= \sum_j E w_{ij} E x_j = 0$$

$$i = 1, \ldots, n_{output}$$
$$j = 1, \ldots, n_{input}$$

$$Var(z_i) = \sum_{j=1}^{n_{input}} \overset{Var(W)}{Var(w_{ij})} \overset{1}{Var(x_j)} =$$

$$\frac{\partial f}{\partial x_j} = \sum_{i=1}^{n_{outputs}} (W^T)_{ji} \frac{\partial f}{\partial z_i} = n_{input} \cdot Var(W) = 1$$

$$Var\left(\frac{\partial f}{\partial x_j}\right) = n_{outputs} \cdot Var(W)$$

$$Var(W) = \frac{2}{n_{input} + n_{output}}$$

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

$$W_{ij} \sim R\left(-\sqrt{\frac{6}{n_{input} + n_{output}}}, \sqrt{\frac{6}{n_{input} + n_{output}}}\right)$$

Xavier (Glorot) init.

Activations: Sigmoid, tanh

$p(x)$, uniform cont. dists.

$\leftarrow x \sim R[a, b]$



$a \quad b \quad x$

# Kaiming (He) init.

$$Var(w) = \frac{2}{n_{input}}$$

$$W_{oj} \sim N(0, \; Var(w))$$

$$W_{oj} \sim R\left(-\sqrt{\frac{6}{n_{input}}}, \; \sqrt{\frac{6}{n_{input}}}\right)$$

Activations: ReLu, Leaky ReLu

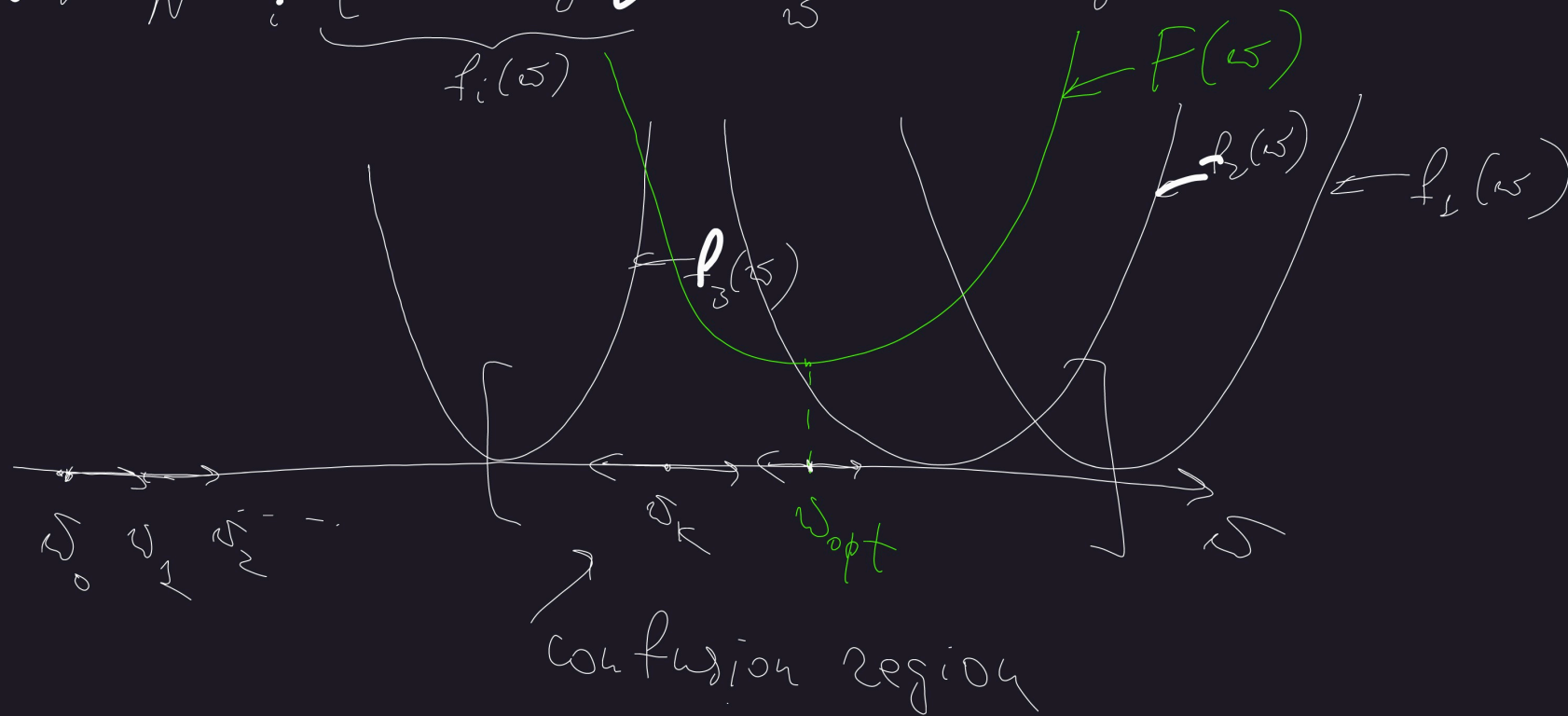$$F(w) = \frac{1}{N} \sum_{i=1}^{N} f_i(w) \longrightarrow \min_{w} \quad ; \quad N \gg 1$$

$$DF(w) = \frac{1}{N} \sum_{i=1}^{N} Df_i(w)$$

$$\underline{SGD} \quad \begin{cases} I_k \subset Unif(1, \ldots, N) \\ g_k = \frac{1}{|I_k|} \sum_{i \in I_k} Df_i(w_k) \\ w_{k+1} = w_k - d_k g_k \end{cases}$$

$$F(w) = \frac{1}{N} \sum_i \underbrace{\left( w x_i - y_i \right)^2}_{f_i(w)} \longrightarrow \min_w \qquad x_i, y_i \in \mathbb{R}$$

$F(w)$

$f_2(w)$

$f_1(w)$

$f_3(w)$

$w_0 \quad w_1 \quad w_2 \cdots$

$w_k$

$w_{opt}$

$w$

confusion region

$$\hat{g}_k = \frac{1}{N} \sum_{i=1}^{N} \nabla f_i(\omega_k)$$

$$g_k = \frac{1}{|I_k|} \sum_{i \in I_k} \nabla f_i(\omega_k)$$

$$g_k \sim \mathcal{N}\left(\hat{g}_k, \hat{\Sigma}_k\right), \quad \hat{\Sigma}_k \to 0 \text{ if } |I_k| \to N$$
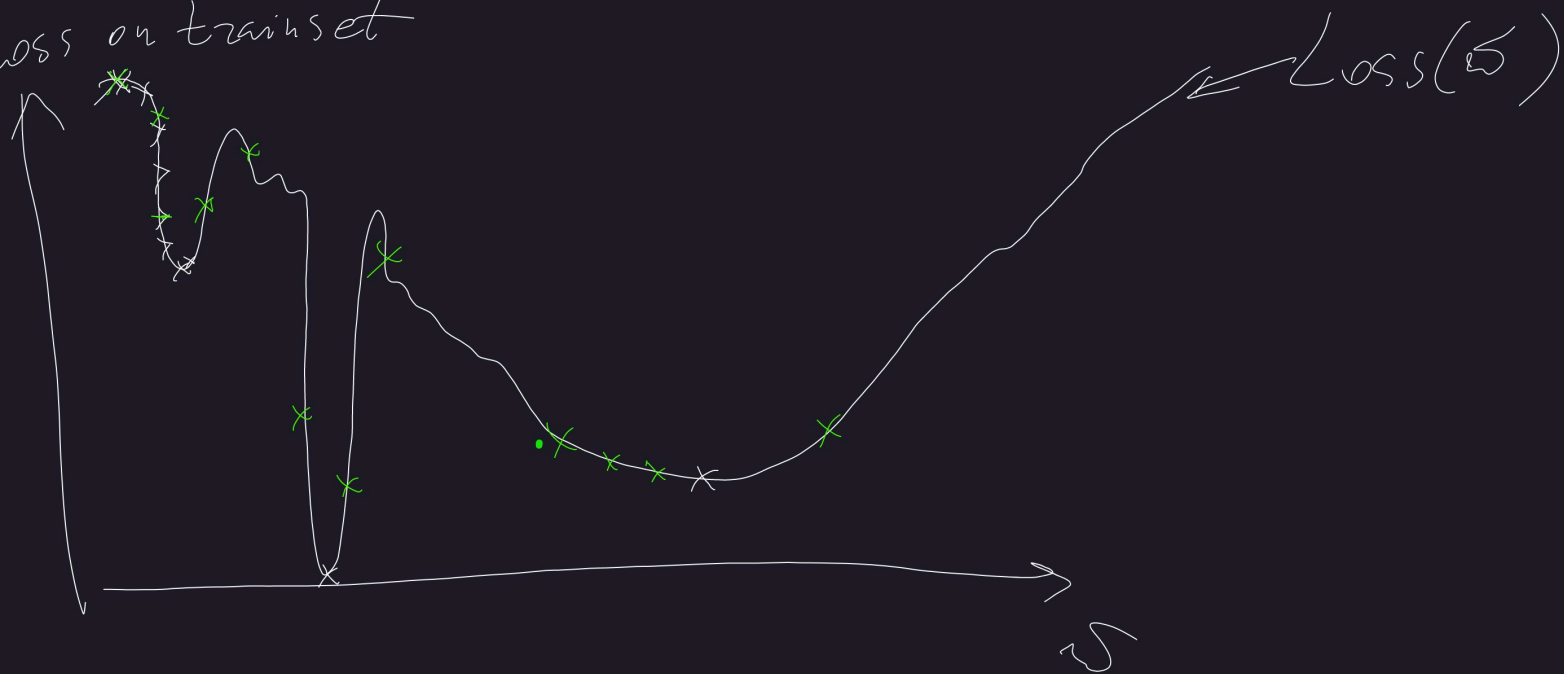
$$X \sim \mathcal{N}(x | \mu, \sigma^2)$$

$$X = \mu + \sigma \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, 1)$$

$$SGD: \omega_{k+1} = \omega_k - \alpha_k g_k = \omega_k - \alpha_k \left( \hat{g}_k + \varepsilon_k \right) \overset{\text{Noisy full}}{GD}$$

$$\varepsilon_k \sim \mathcal{N}(0, \Sigma_k)$$

loss on trainset

Loss($w$)

$S$

# Learning Rates Schedule

$d_k$ exponential decay

cosine decay

epoch

epoch

SWA

$$S = \frac{1}{4} \sum_{i=1}^{4} S_i$$

$S_1 \quad S_2 \quad S_3 \quad S_4$