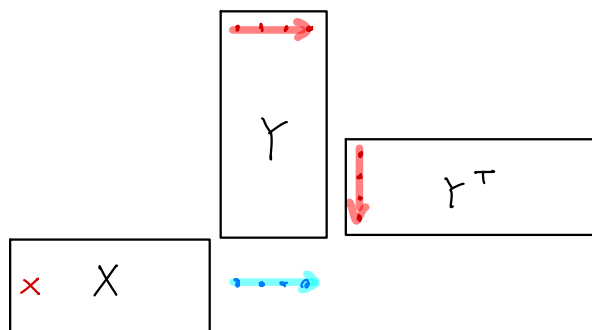we know $\dfrac{\partial L}{\partial w}$ , $\quad w = XY$

then ① $\dfrac{\partial L}{\partial X} = \dfrac{\partial L}{\partial w} \cdot Y^T$

② $\dfrac{\partial L}{\partial Y} = X^T \dfrac{\partial L}{\partial w}$

---

GD (Gradient Descent)

$$\min_{\theta} \mathcal{L}(\theta) \implies \boxed{\theta^{k+1} = \theta^k - \alpha \nabla \mathcal{L}(\theta^k)}$$

---

$\begin{cases} \text{Full Batch GD} & \tilde{\nabla} f(x) = \dfrac{1}{N} \sum\limits_{i=1}^{N} \nabla f_i(x) & \{1, 2, \cdots \cdots N\} \\[2em] \text{Stochastic} \quad `` & \tilde{\nabla} f(x) = \nabla f_i(x) & i \\[2em] \text{Mini Batch} \quad `` & \tilde{\nabla} f(x) = \dfrac{1}{|K|} \sum\limits_{k \in K} \nabla f_k(x) & K \end{cases}$

---

Sigmoid $\implies$ tanh $\implies$ ReLU

$\sigma(z) = \dfrac{1}{1 + e^{-z}}$

$\sigma'(z) = \sigma(z)(1 - \sigma(z))$

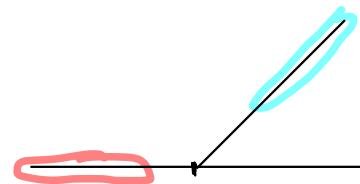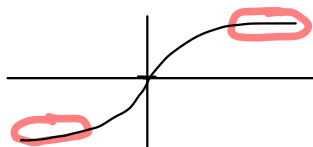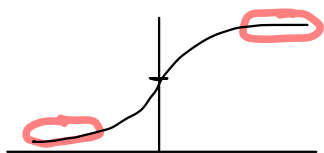$\tanh(z) = \dfrac{e^z - e^{-z}}{e^z + e^{-z}}$

$(\tanh)' = 1 - (\tanh)^2$

$\text{ReLU}(z) = \max(0, z)$

$(\text{ReLU})' = \begin{cases} 1 & z > 0 \\ 0 \end{cases}$

not zero centered.    zero centered.    not zero centered.

Xavier Init   $\boxed{\sigma = \dfrac{1}{\sqrt{D_{in}}}}$

pf.  let  $y = w^T x$

$\qquad = \displaystyle\sum_{i=1}^{D_{in}} w_i x_i$

- $w_i$ : i.i.d
- $x, w$ independent.

- $E[w_i] = 0$ , $Var(w_i) = \sigma^2$
- $E[x_i] = 0$ , $Var[x_i] = V$

---

Our Goal :  $\boxed{Var[y] = Var[x_i]}$

- $Var[y] = E[y^2] - E[\cancel{y}]^2$

- $E[y] = \displaystyle\sum_i E[x_i w_i] = \sum_i \underline{E[x_i] E[w_i]}$

$\qquad\qquad\qquad = 0$

independent.

$\displaystyle\iint x \cdot w \, p(x, w) \, dx \, dw$

$\displaystyle = \iint x \cdot w \, p(x) \cdot p(w) \, dx \cdot dw$

- $Var[y] = E[y^2]$

$\qquad = \displaystyle\sum_i E[w_i^2 x_i^2] + \sum_{i \neq k} E[\cancel{w_i}] E[\cancel{w_k}] E[x_i x_k]$

$\qquad = \displaystyle\sum_i \underline{E[w_i^2]}\ \underline{E[x_i^2]} \quad = \quad \cancel{\sigma^2 V \cdot D_{in}}$

$\qquad\qquad\quad Var(w_i)\quad Var(x_i)$

$\boxed{\begin{array}{c} Note \\[4pt] \sigma^2 \Uparrow \qquad\qquad \sigma^2 \Downarrow \\[4pt] Var(y) \Uparrow \qquad\quad Var(y) \Downarrow \end{array}}$

$\dfrac{7}{2} = \dfrac{2}{2}$   $\sigma = \dfrac{1}{\sqrt{D_{in}}}$

Xavier Init ♡ ReLU $\Rightarrow$ $\sigma = \sqrt{\dfrac{2}{Din}}$

pf.
- let $z = \sum_{i}^{Din} w_i x_i$

- then $Var(z) = \sigma^2 V \cdot Din. \Rightarrow z \sim N(0, \underline{\sigma^2 V \cdot Din})$

  "$\delta$"

- let $h = \phi(z)$

- Our Goal : $Var(h) = Var(x_i)$

  - $E[h] = \int_0^\infty z \, \rho(z) \, dz$

    $= \int_0^\infty z \cdot \dfrac{1}{\sqrt{2\pi\delta}} \cdot \exp\left(-\dfrac{z^2}{2\delta}\right) dz.$

    $= \int_0^\infty \dfrac{\not z}{\sqrt{2\pi\delta}} \cdot e^{-u} \cdot \dfrac{\delta}{\not z} \, du = \sqrt{\dfrac{\delta}{2\pi}}$

    let $u = \dfrac{z^2}{2\delta}$

    $dz = \dfrac{\delta}{z} \, du$

  - $E[h^2] = \int_0^\infty z^2 \cdot \rho(z) \, dz$

    $= \dfrac{1}{2} \int_{-\infty}^{\infty} z^2 \rho(z) \, dz = \dfrac{1}{2} E[z^2]$

    $= \dfrac{1}{2} Var[z^2] = \dfrac{1}{2}\delta$

  - $Var[h] = \dfrac{1}{2}\delta - \dfrac{\delta}{2\pi} \approx \boxed{\dfrac{1}{2}} Var[z]$

    "균산의 가능성이
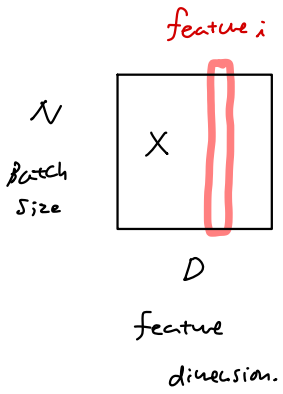
    1/2 배가 된다"

  - $\dfrac{1}{2} \cdot \underbrace{\sigma^2 V \cdot Din}_{\delta} = V$

    $\sigma = \sqrt{\dfrac{2}{Din}}$

# Batch Normalization

inference에서 $N=1$ 일 수도 있다!

feature $i$

$N$
Batch Size

$X$

$D$
feature dimension.

$$\bullet \quad \hat{x}_{ij} = \frac{x_{ij} - \mu_j}{\sqrt{\delta_j^2 + \varepsilon}}$$

$$\mu^{run} \leftarrow m \cdot \mu^{run} + (1-m) \cdot \mu^{batch}$$

$$(\delta^2)^{run} \leftarrow m \cdot (\delta^2)^{run} + (1-m)(\delta^{-2})^{batch}$$

$$\bullet \quad \hat{y}_{ij} = r_j \, \hat{x}_{ij} + \beta_j$$

학습한다.

---

## Why BN works?

① BN $\Rightarrow$ ICS $\Downarrow$

Internal
Convergence Shift

: shift in the
mean/var of hidden activation
during training

② BN $\not\Rightarrow$ ICS $\Downarrow$

but BN $\Rightarrow$ training $\Uparrow$

③ BN $\Rightarrow$ Smooth
loss landscape.

"gradient 의
변화가 안정적"

---

FC 간층에 Normalization 하와면

cost $= 0$ (linear)

Gradient Descent!  $\begin{cases} \text{minimize} \ f(x) \\ \text{update} \quad x' = x - \eta \nabla f(x) \end{cases}$

---

- Suppose $x_{t+1} = x_t + \eta v$

- then by Taylor Approximation

$$f(x_{t+1}) \approx f(x_t) + \langle \nabla f(x_t), \eta v \rangle$$

- minimizing $f(x_{t+1})$, $\quad v = - \dfrac{\nabla f(x_t)}{\| \nabla f(x_t) \|}$

- $x_{t+1} = x_t - \eta \dfrac{\nabla f(x_t)}{\| f(x_t) \|}$

---

Lemma 3.1

- let $f : \mathbb{R}^d \to \mathbb{R}$ continuously differentiable func.

- let $f$ $\beta$-smooth. $\Rightarrow$ $\forall x, y$ $\boxed{| \nabla f(y) - \nabla f(x) |} \le \boxed{\beta \| x - y \|}$

then $f(y) \le f(x) + \underbrace{\langle \nabla f(x), y - x \rangle}_{\text{접선}} + \underbrace{\dfrac{\beta}{2} \| y - x \|^2}_{\text{2통면.}}$

---

let $x_{t+1} = x_t - \eta \nabla f(x_t)$

then $f(x_{t+1}) \le f(x_t) + \underbrace{\langle \nabla f(x_t), -\eta \nabla f(x_t) \rangle}_{- \eta \| \nabla f(x_t) \|^2} + \underbrace{\dfrac{\beta}{2} \| -\eta \nabla f(x_t) \|^2}_{\frac{\beta}{2} \eta^2 \| \nabla f(x_t) \|^2}$

$= f(x_t) - \boxed{\left[ \eta - \dfrac{\beta}{2} \eta^2 \right]} \| \nabla f(x_t) \|^2$

$\eta \ll 1$ then positive

Proof for ③.1

let $g(t) = f(x + t(y-x))$

then $g'(t) = \langle \nabla f(x + t(y-x)), y-x \rangle$

$g''(t) = (y-x)^T \nabla^2 f(x + t(y-x))(y-x) \leq \beta \| y-x \|^2$

---

$$\int_0^1 (1-s) g''(s) \, ds = \left[ (1-s) g'(s) \right]_0^1 + \int_0^1 g'(s) \, ds$$

$$g(1) = g(0) + g'(0) + \int_0^1 (1-s) g''(s) \, ds$$

---

$$f(y) = f(x) + \langle \nabla f(x), y-x \rangle + \int_0^1 (1-s)(y-x)^T \nabla^2 f(x + s(y-x))(y-x) \, ds.$$

$$\leq \int_0^1 (1-s) \beta \| y-x \|^2 \, ds$$

$$= \frac{\beta}{2} \| y-x \|^2 .$$

let $\|\nabla f(y) - \nabla f(x)\| \leq \beta \|x - y\|$

then $\quad v^T \nabla^2 f(x) v \leq \beta \|v\|^2$

let $\quad \phi(t) = \nabla f(x + t(y-x))$

$\quad \phi'(t) = \nabla^2 f(x + t(y-x))\,(y-x)$

then $\quad \nabla f(y) - \nabla f(x) = \int_0^1 \phi'(t)\,dt$

$$= \int_0^1 \nabla^2 f(x + t(y-x))\,dt \cdot (y-x)$$

then $\quad \langle \nabla f(y) - \nabla f(x),\ y-x \rangle \leq \|\nabla f(y) - \nabla f(x)\| \, \|y-x\|$

$$\leq \beta \|y-x\|^2$$

too small

$$(y-x)^T \int_0^1 \nabla^2 f(x + t(y-x))\,dt\ (y-x) \leq \beta \|y-x\|^2.$$

$$v^T \nabla^2 f(x) v \leq \beta \|v\|^2.$$

SGD

$$f(y) \leq f(x) + \langle \nabla f(x), y - x \rangle + \frac{\beta}{2} \| y - x \|^2$$

let $\quad x_{t+1} = x_t - \eta \, \tilde{\nabla} f(x_t)$ $\qquad$ $E_t \left[ \tilde{\nabla} f(x_t) \right] = \nabla f(x_t)$

then $\quad f(x_{t+1}) \leq f(x_t) - \eta \langle \nabla f(x_t), \tilde{\nabla} f(x_t) \rangle + \frac{\beta}{2} \cdot \eta^2 \| \tilde{\nabla} f(x_t) \|^2$

---

$$E_t [f(x_{t+1})] \leq f(x_t) - \eta \| \nabla f(x_t) \|^2 + \frac{\beta}{2} \eta^2 \underbrace{E_t \left[ \| \tilde{\nabla} f(x_t) \|^2 \right]}_{\leq G}$$

$$E \left[ \| \nabla f(x_t) \|^2 \right] \leq \frac{1}{\eta} \left( E [f(x_t)] - E [f(x_{t+1})] \right) + \frac{\beta}{2} \eta G$$

$$\sum_{t=0}^{T-1} E \left[ \| \nabla f(x_t) \|^2 \right] \leq \frac{1}{\eta} \left( f_0 - \underset{\text{lower bound}}{\boxed{f^*}} \right) + \frac{\beta}{2} \eta G T$$

$$\min E \left[ \| \nabla f(x_t) \|^2 \right] \leq \frac{1}{\eta T} [f_0 - f^*] + \frac{\beta}{2} \eta G$$

let $\quad \eta = \frac{1}{\sqrt{T}}$ $\quad$ then $\quad \min E \left[ \| \nabla f(x_t) \|^2 \right] = \Theta \left( \frac{1}{\sqrt{T}} \right)$