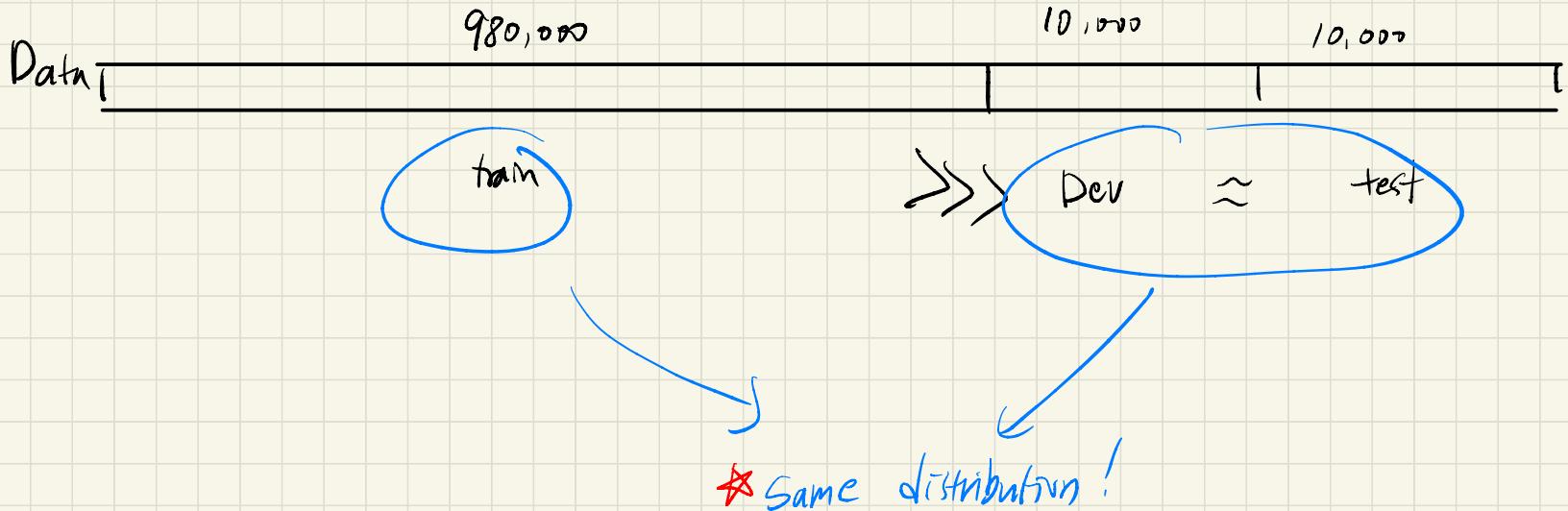
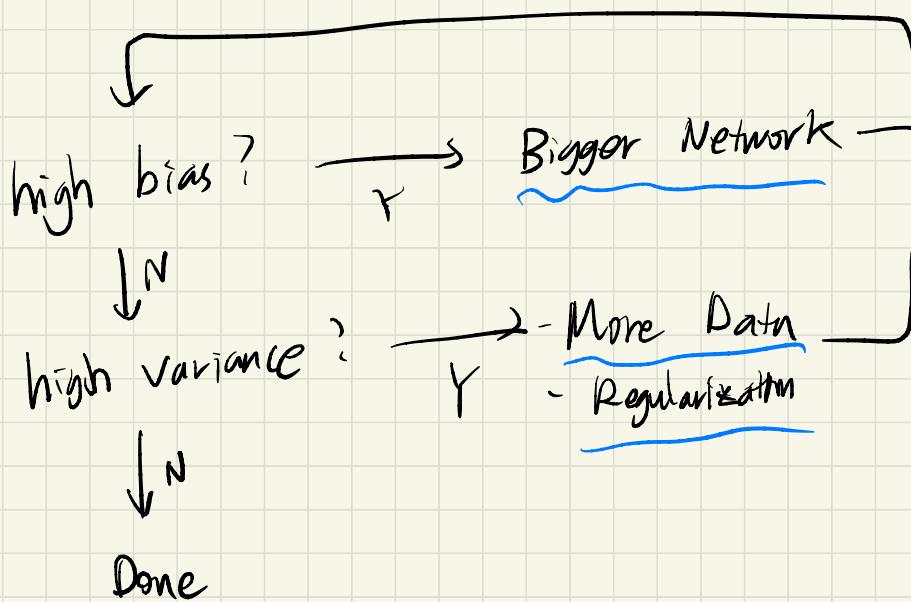


Big data era : $1,000,000 \uparrow$ data



[Basic Recipe]



[Regularization]

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m f(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \|w\|_2^2$$

$$L_2 \text{ Regularization: } \|w\|_2^2 = \sum_{j=1}^{n_x} w_j^2 = w^T w$$

$$L_1 \quad " \quad : \frac{\lambda}{2m} \sum_{i=1}^{n_o} |w_i| = \frac{\lambda}{2m} \|w\|_1$$

[Normalize inputs]

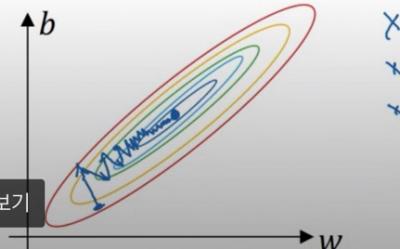
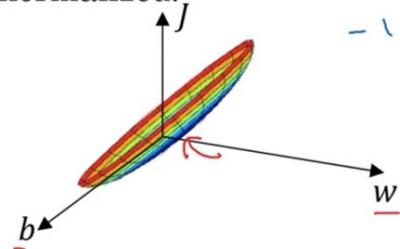
Normalizing Inputs (C2W1L09)

Why normalize inputs?

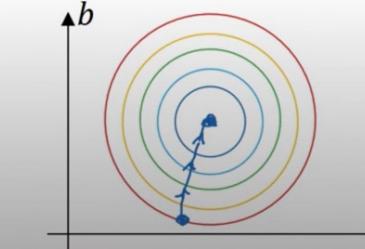
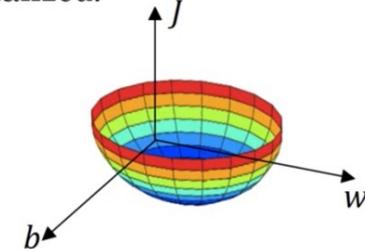
$J(w, b) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)})$

Unnormalized:

$w_1 \quad x_1: \underline{1 \dots 1000} \leftarrow$
 $w_2 \quad x_2: \underline{0 \dots 1} \leftarrow$
 $b \quad -1 \dots 1$



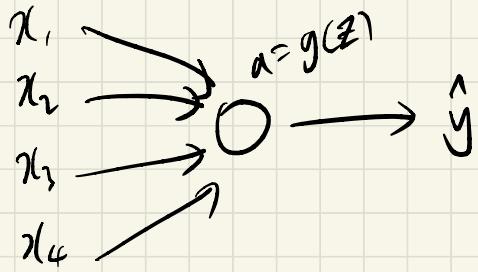
Normalized:



동영상 더보기

▶ 🔍 4:54 / 5:30

Andrew Ng



[weights initialization]

$$z = w_1x_1 + w_2x_2 + \dots + w_nx_n + b$$

$$w^{[l]} = \text{np.random.randn(shape)} \times \text{np.sqrt}\left(\frac{2}{n^{[l-1]}}\right)$$

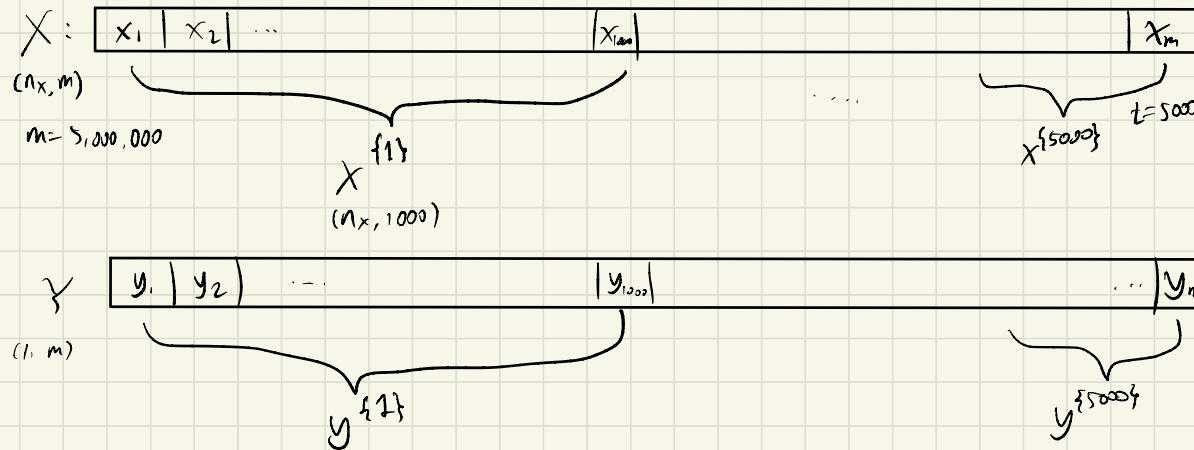
when $g^{[l]}(z) = \text{ReLU}(z)$

initialize
randomly

$$\tanh, \quad \sqrt{\frac{1}{n^{[l-1]}}} \quad \} \quad \text{Xavier initialization}$$

$$\text{or } \sqrt{\frac{2}{n^{[l-1]} + n^{[l]}}}$$

[mini batch GD]



for $t = 1, \dots, 5000$

1. Forward Prop. on $X^{[t:500]}$

$$\begin{aligned} Z^{[0]} &= W^{[0]} X^{[t:500]} + b^{[0]} \\ A^{[0]} &= g^{[0]}(Z^{[0]}) \\ &\vdots \\ A^{[L-1]} &= g^{[L-1]}(Z^{[L-1]}) \end{aligned}$$

} Vectorize
(1,000 examples)

2. Compute cost $J^{[t:500]}$

$$= \frac{1}{1000} \sum_i^{} J(g^{[i]}(x), y^{[i]}) + \frac{\lambda}{2 \cdot 1000} \sum_j^{} \|w^{[0j]}\|^2$$

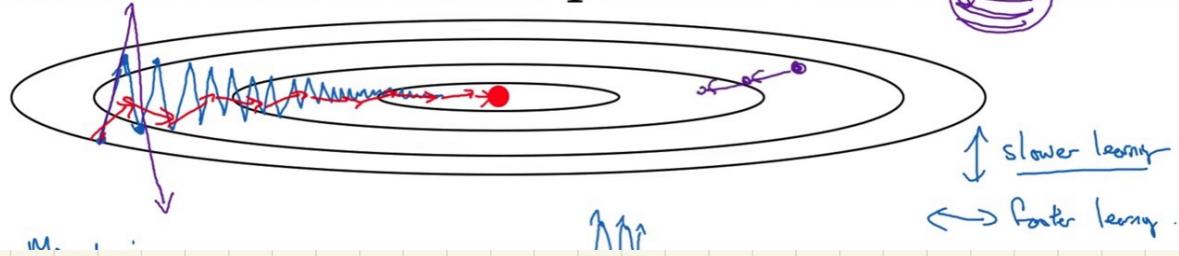
1 epoch

3. Back Prop.

$$4. w^{[0]} := w^{[0]} - \alpha \Delta w^{[0]}, \quad b^{[0]} := b^{[0]} - \alpha \Delta b^{[0]}$$

[GD with Momentum]

Gradient descent example



On iteration t :

Compute dW, db on current minibatch.

$$Vdw = \beta Vdw + (1-\beta)dw \quad \rightarrow \text{Moving Average of } dw$$
$$Vdb = \beta Vdb + (1-\beta)db \quad \begin{matrix} \downarrow \text{velocity} \\ \downarrow \text{acceleration} \end{matrix}$$

$$w := w - \alpha Vdw \quad b := b - \alpha Vdb$$

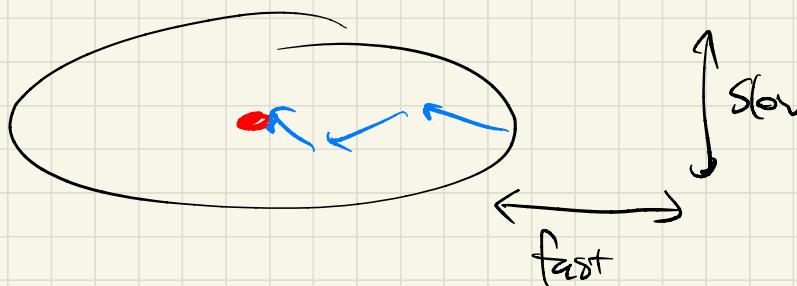
[RMS Prop]

$$S_{dw} = \beta S_{dw} + (1-\beta) dw^2$$

$$S_{db} = \beta S_{db} + (1-\beta) db^2$$

$$w := w - \alpha \frac{dw}{\sqrt{S_{dw}}}$$

$$b := b - \alpha \frac{db}{\sqrt{S_{db}}}$$



[Adam]

momentum

V_{dw}^{corr} = $V_{dw} / (1 - \beta_1^t)$, $V_{db}^{corr} = V_{db} / (1 - \beta_1^t)$

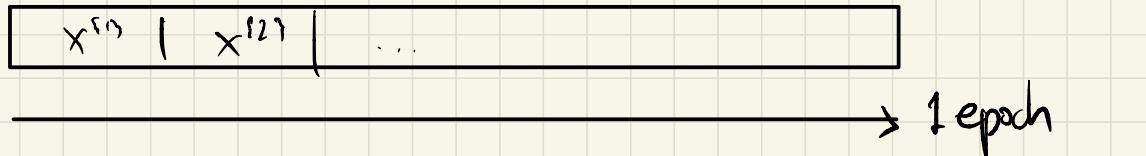
RMS prop

S_{dw}^{corr} = $S_{dw} / (1 - \beta_2^t)$, $S_{db}^{corr} = S_{db} / (1 - \beta_2^t)$

$$w := w - \alpha \frac{V_{dw}^{corr}}{\sqrt{S_{dw}^{corr}} + \epsilon}$$

$$b := b - \alpha \frac{V_{db}^{corr}}{\sqrt{S_{db}^{corr}} + \epsilon}$$

[Leaving Rate Decay]

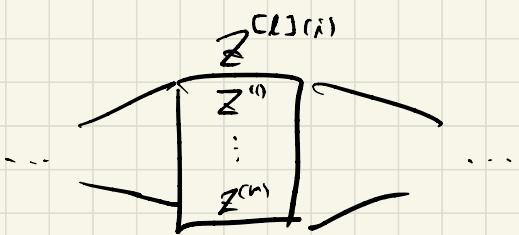


$$\alpha = \frac{1}{1 + \text{decay rate} \times \text{epoch}} \alpha_0$$

$\alpha_0 = 0.2$
 $\text{d.r.} = 1$

Epoch	α
1	$\frac{1}{1 + (1 \times 1)} \cdot 0.2 = 0.1$
2	$\frac{1}{1 + (1 \times 2)} \cdot 0.2 = 0.067$
3	$= 0.05$
4	$= 0.04$
.	
.	

[Batch Normalization]



$$\tilde{Z}^{(l)} = \gamma Z_{\text{norm}}^{(i)} + \beta$$

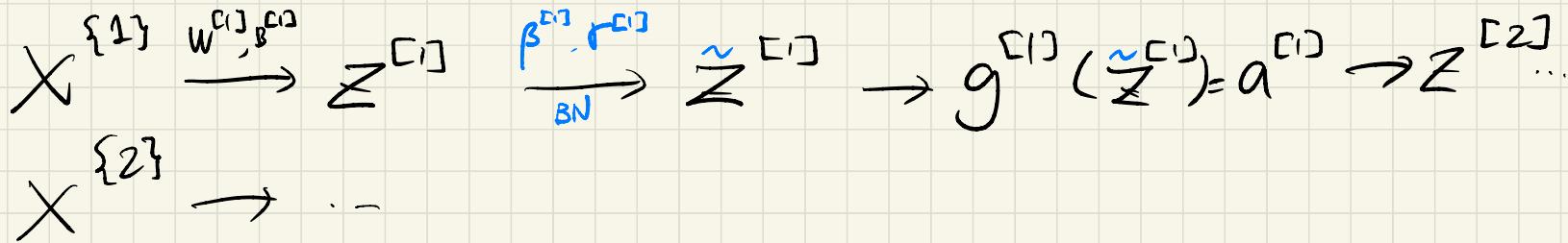
where $Z_{\text{norm}}^{(i)} = \sqrt{\sigma^2 + \epsilon}$

then $\tilde{Z}^{(l)} = Z^{(i)}$

$$\mu = \frac{1}{m} \sum_i Z^{(i)}$$

$$\sigma^2 = \frac{1}{m} \sum_i (Z_i - \mu)^2$$

$$Z_{\text{norm}}^{(i)} = \frac{Z^{(i)} - \mu}{\sqrt{\sigma^2 + \epsilon}}$$



$$\sum_{(n^{[l]}, 1)}^{[l]} = w^{[l]} \cdot a^{[l-1]} + \cancel{b^{[l]}}$$

$$\sum_{\text{norm}}^{[l]}$$

$$\tilde{\sum}^{[l]} = r^{[l]} \cdot \sum_{\text{norm}}^{[l]} + \beta^{[l]}$$