

Tips for training

Tips for Training

1. Batch size

Optimization with Batch

$$\theta^* = \operatorname{argmin}_{\theta} L$$

➤ (Randomly) Pick initial values θ^0

➤ Compute gradient $g^0 = \nabla L^1(\theta^0)$

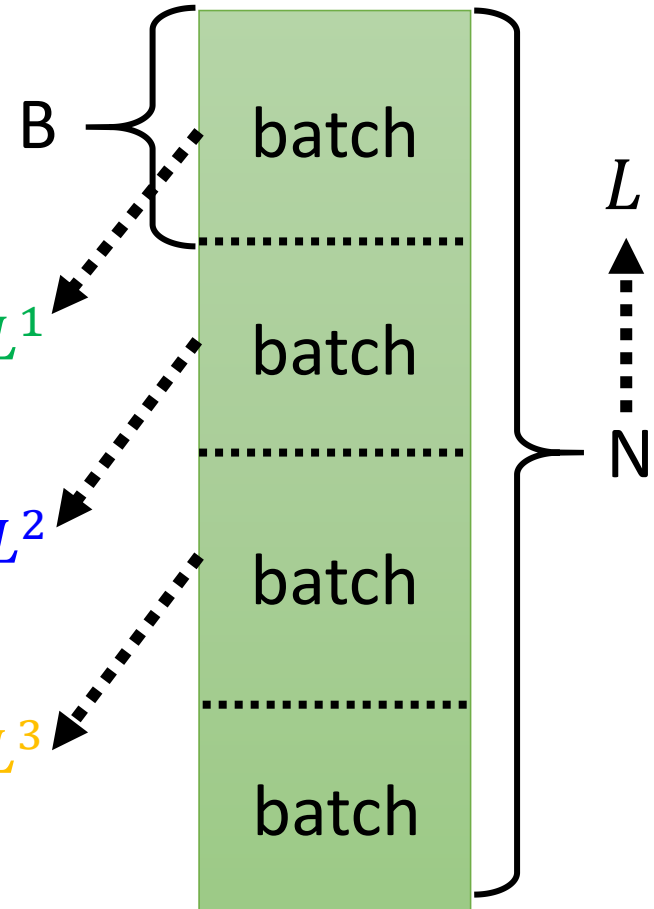
$$\text{update } \theta^1 \leftarrow \theta^0 - \eta g^0$$

➤ Compute gradient $g^1 = \nabla L^2(\theta^1)$

$$\text{update } \theta^2 \leftarrow \theta^1 - \eta g^1$$

➤ Compute gradient $g^3 = \nabla L^3(\theta^2)$

$$\text{update } \theta^3 \leftarrow \theta^2 - \eta g^3$$



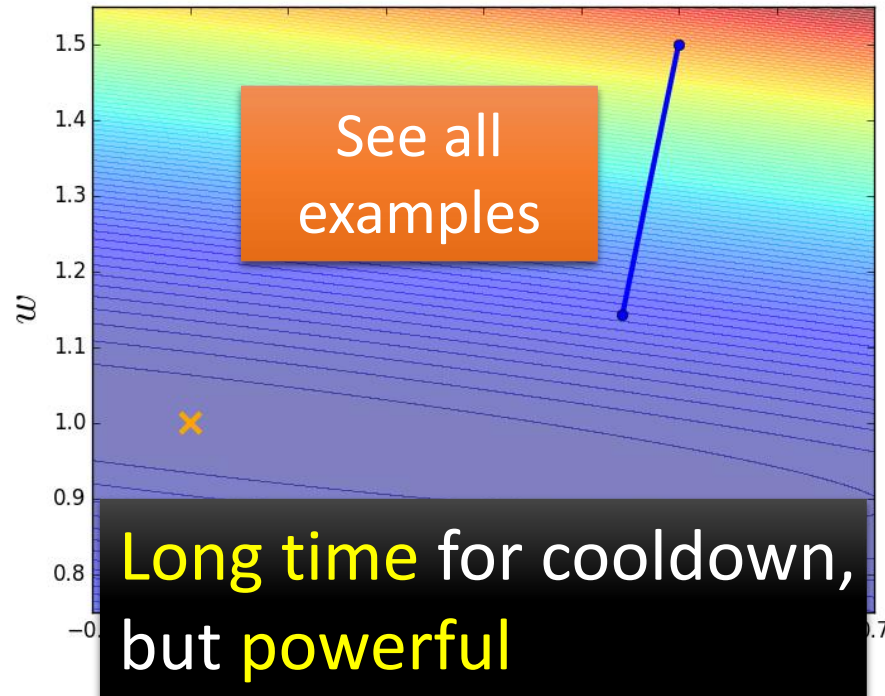
1 **epoch** = see all the batches once → **Shuffle** after each epoch

Small Batch v.s. Large Batch

Consider 20 examples ($N=20$)

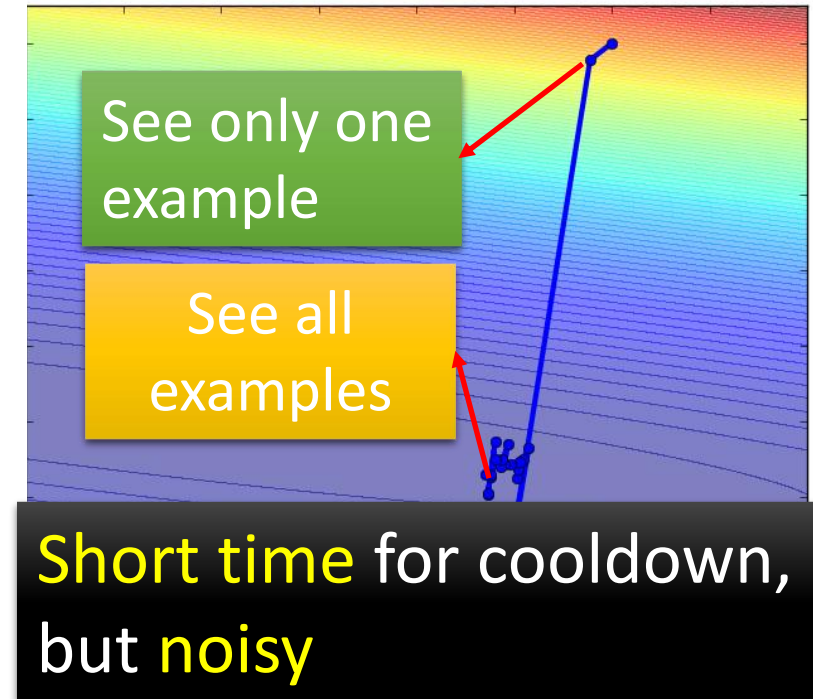
Batch size = N (Full batch)

Update after seeing all
the 20 examples



Batch size = 1

Update for each example
Update 20 times in an epoch

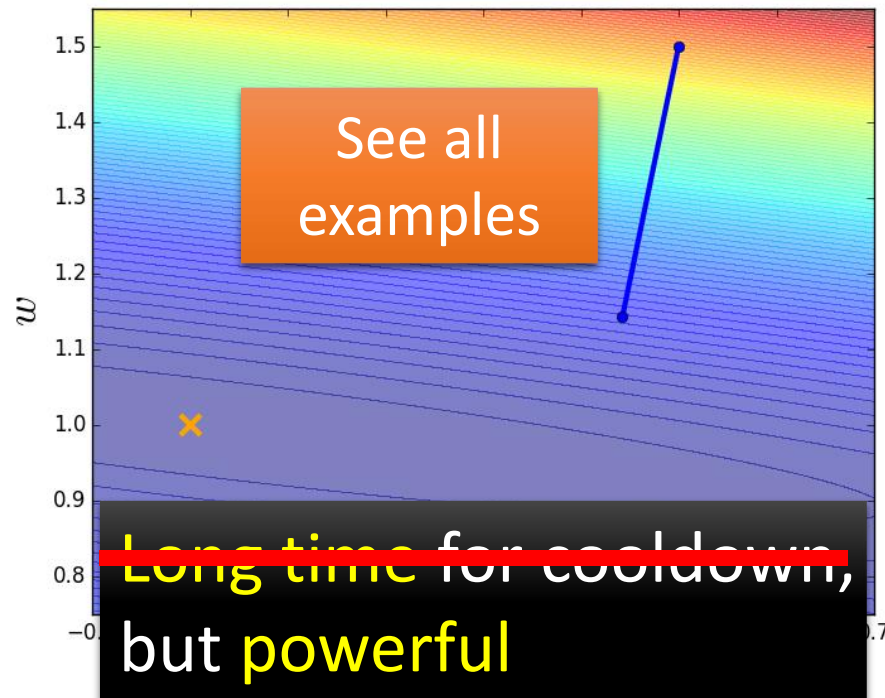


Small Batch v.s. Large Batch

Consider 20 examples ($N=20$)

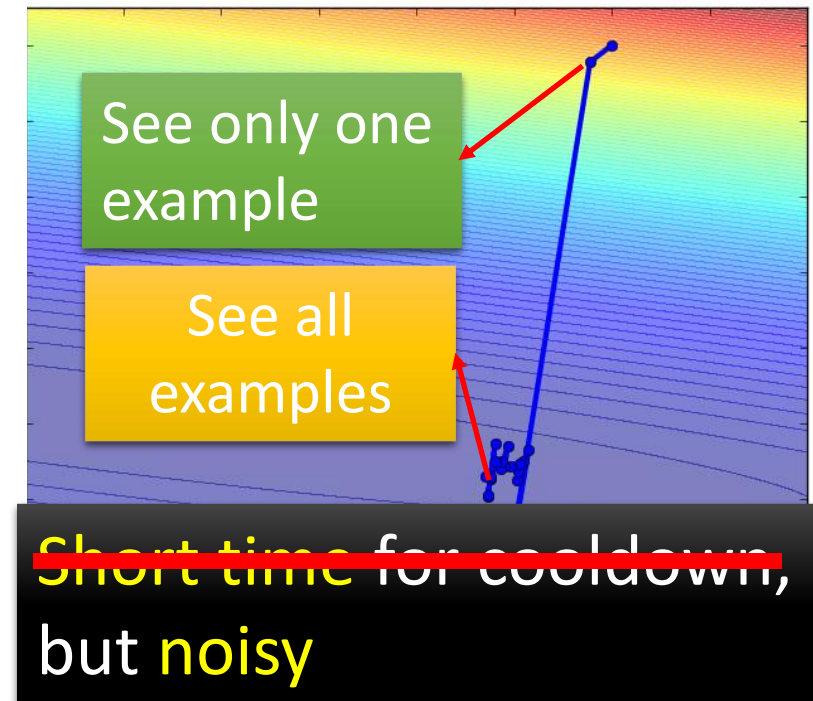
Batch size = N (Full Batch)

Update after seeing all
the 20 examples



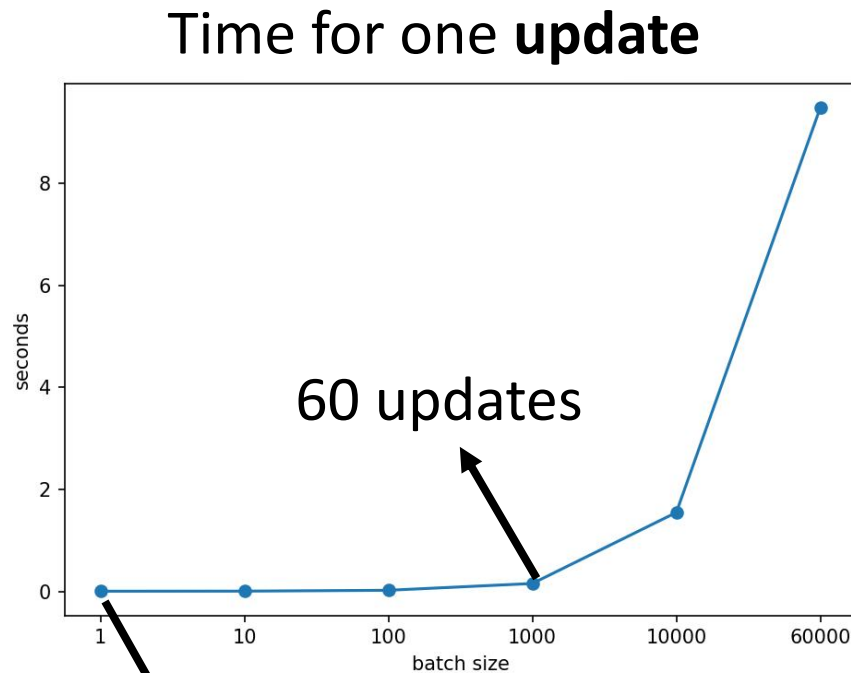
Batch size = 1

Update for each example
Update 20 times in an epoch

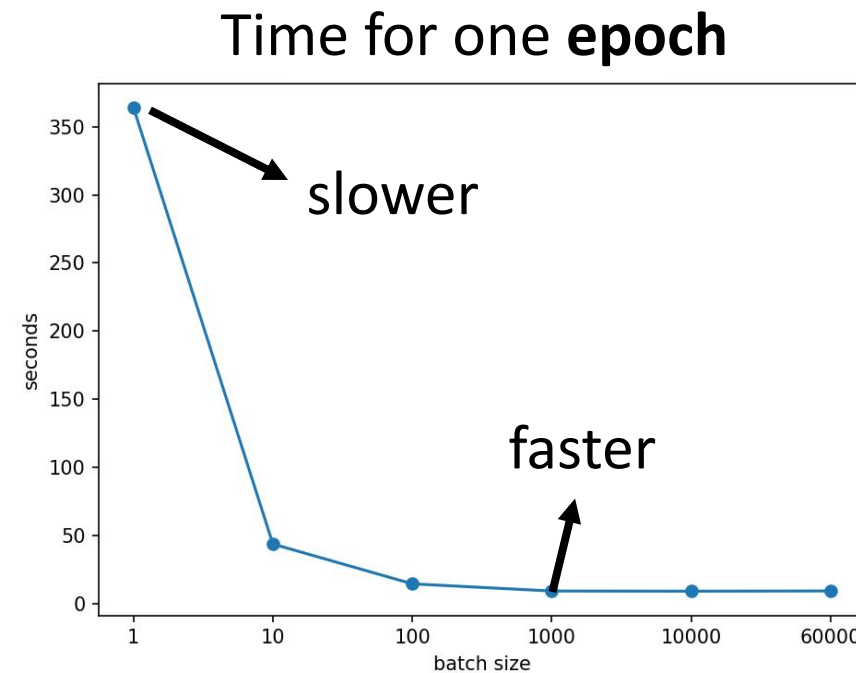


Small Batch v.s. Large Batch

- Smaller batch requires longer time for one epoch (longer time for seeing all data once)



60000 updates in one epoch



MNIST: digit classification
Tesla V100 GPU

Small Batch v.s. Large Batch

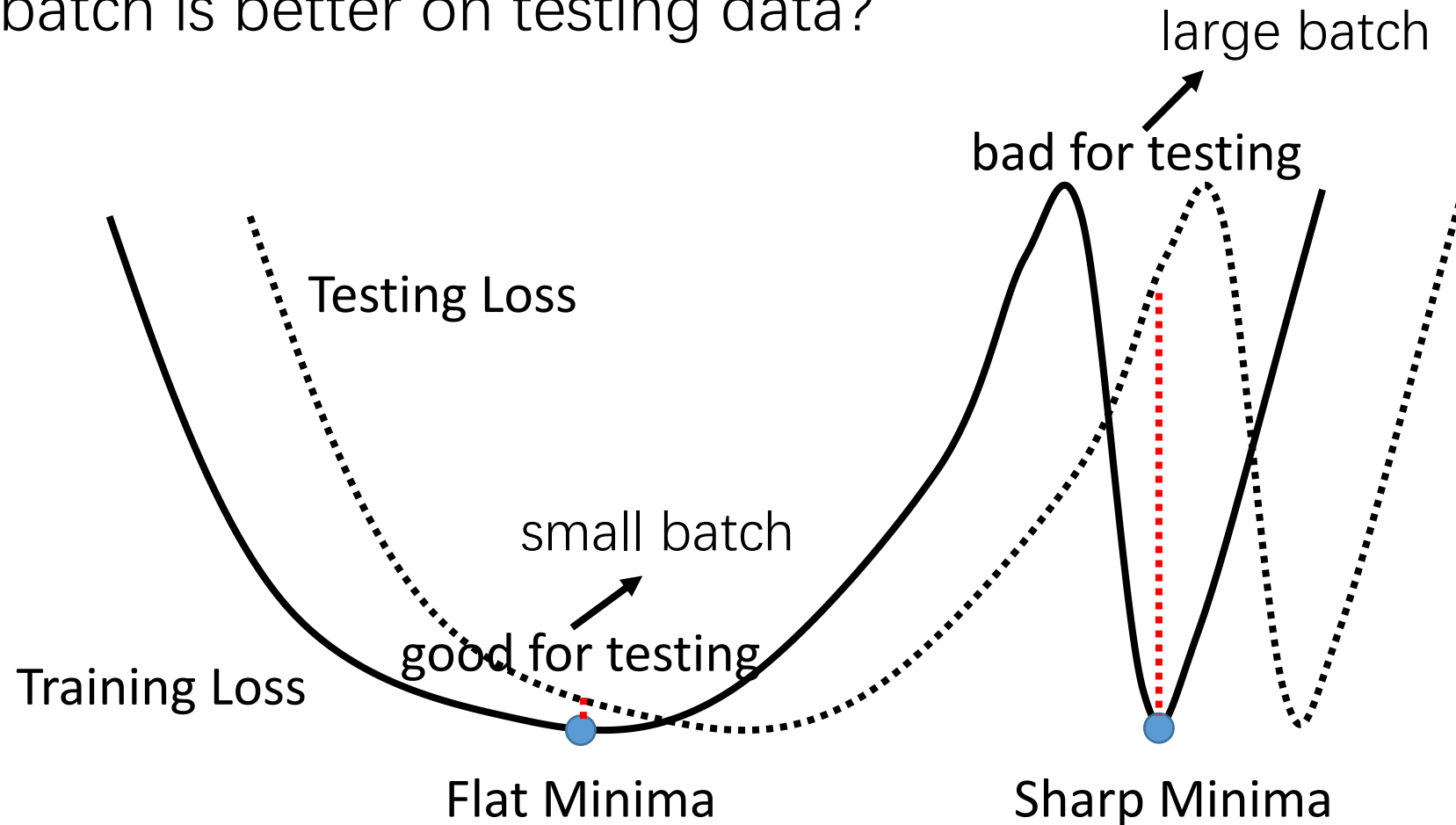
- Small batch is better on testing data?

	Name	Network Type	Data set
SB = 256	F_1	Fully Connected	MNIST (LeCun et al., 1998a)
	F_2	Fully Connected	TIMIT (Garofolo et al., 1993)
LB =	C_1	(Shallow) Convolutional	CIFAR-10 (Krizhevsky & Hinton, 2009)
0.1 x data set	C_2	(Deep) Convolutional	CIFAR-10
	C_3	(Shallow) Convolutional	CIFAR-100 (Krizhevsky & Hinton, 2009)
	C_4	(Deep) Convolutional	CIFAR-100




Name	Training Accuracy		Testing Accuracy	
	SB	LB	SB	LB
F_1	99.66% \pm 0.05%	99.92% \pm 0.01%	98.03% \pm 0.07%	97.81% \pm 0.07%
F_2	99.99% \pm 0.03%	98.35% \pm 2.08%	64.02% \pm 0.2%	59.45% \pm 1.05%
C_1	99.89% \pm 0.02%	99.66% \pm 0.2%	80.04% \pm 0.12%	77.26% \pm 0.42%
C_2	99.99% \pm 0.04%	99.99% \pm 0.01%	89.24% \pm 0.12%	87.26% \pm 0.07%
C_3	99.56% \pm 0.44%	99.88% \pm 0.30%	49.58% \pm 0.39%	46.45% \pm 0.43%
C_4	99.10% \pm 1.23%	99.57% \pm 1.84%	63.08% \pm 0.5%	57.81% \pm 0.17%

Small Batch v.s. Large Batch

- Small batch is better on testing data?



Small Batch v.s. Large Batch

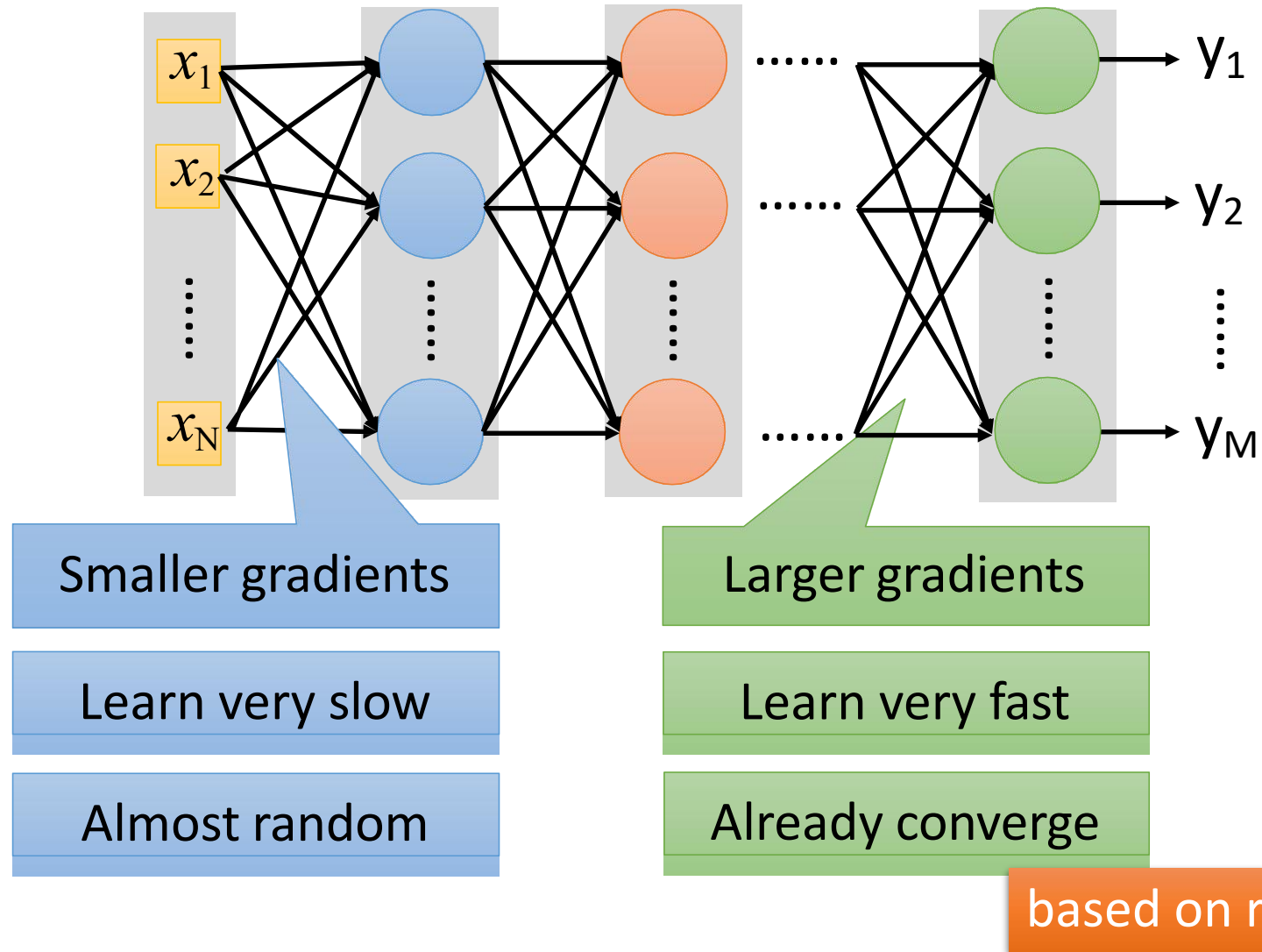
	Small	Large
Speed for one update (no parallel)	Faster	Slower
Speed for one update (with parallel)	Same	Same (not too large)
Time for one epoch	Slower	Faster 
Gradient	Noisy	Stable
Optimization	Better 	Worse
Generalization	Better 	Worse

Batch size is a hyperparameter you have to decide.

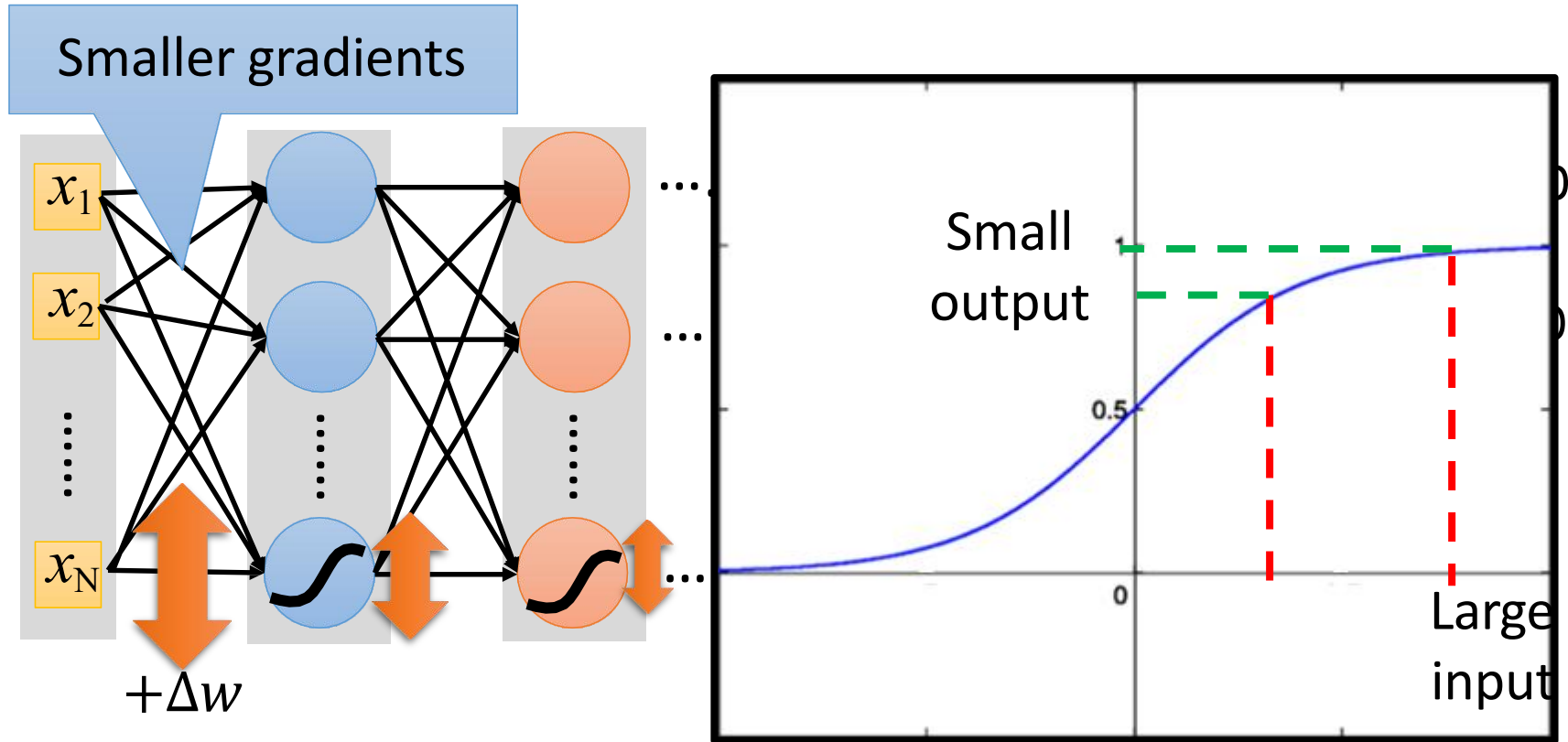
Tips for Training

2. Activation Function

Vanishing Gradient Problem

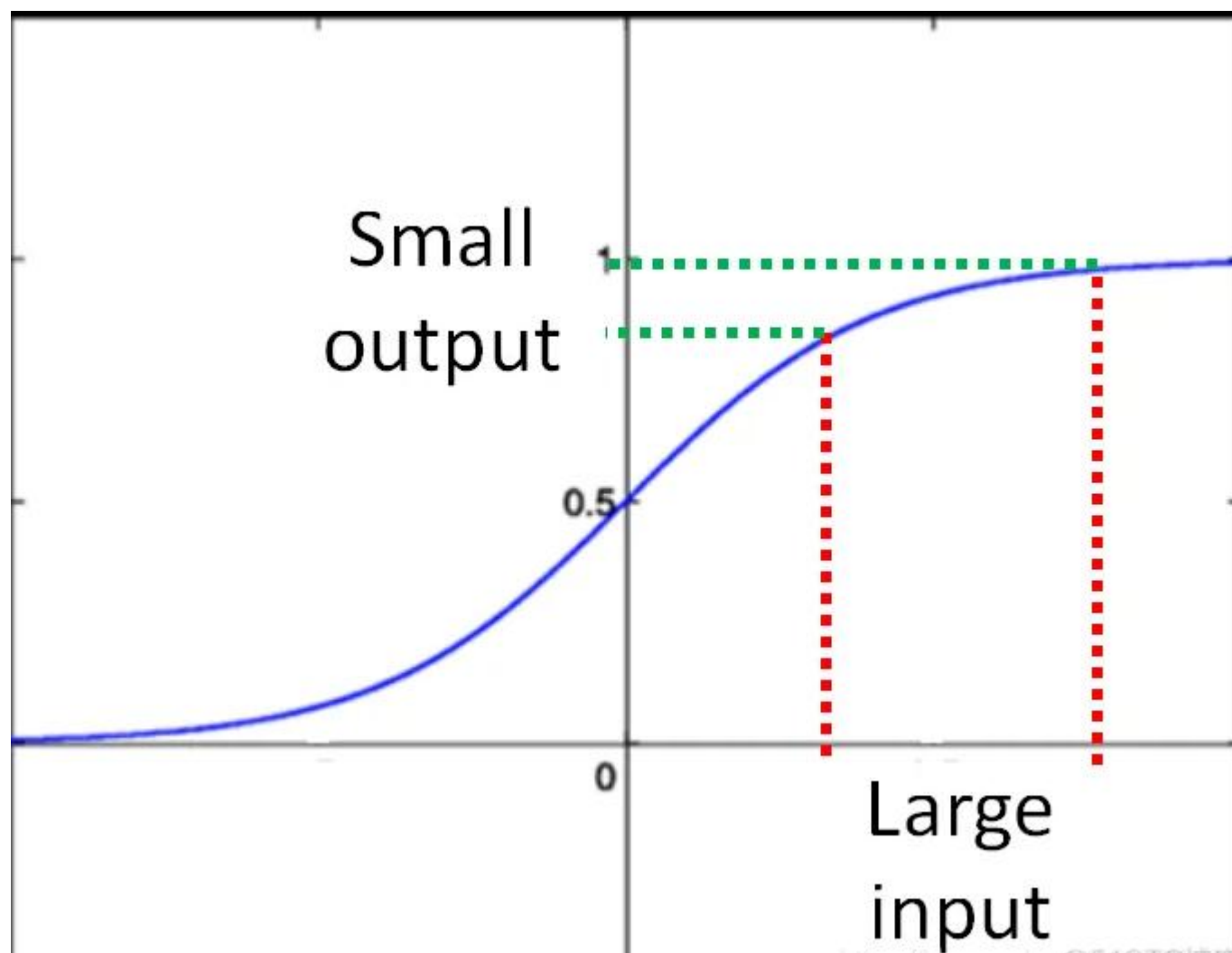


Vanishing Gradient Problem



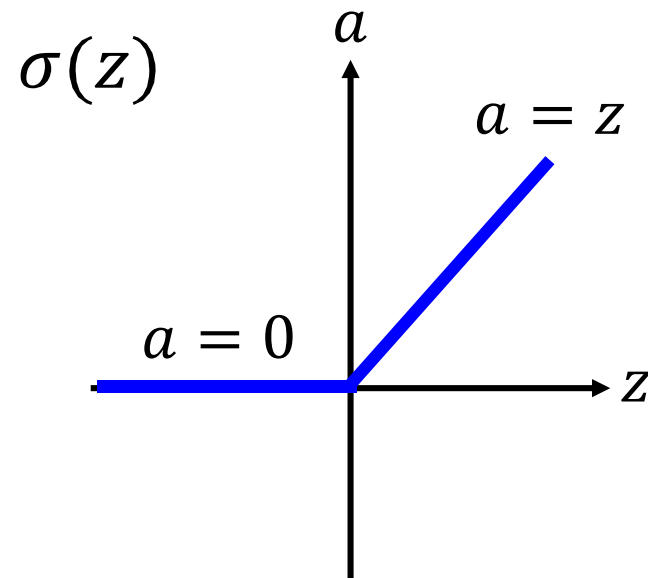
Intuitive way to compute the derivatives ...

$$\frac{\partial l}{\partial w} = ? \quad \frac{\Delta l}{\Delta w}$$



ReLU

- Rectified Linear Unit (ReLU)

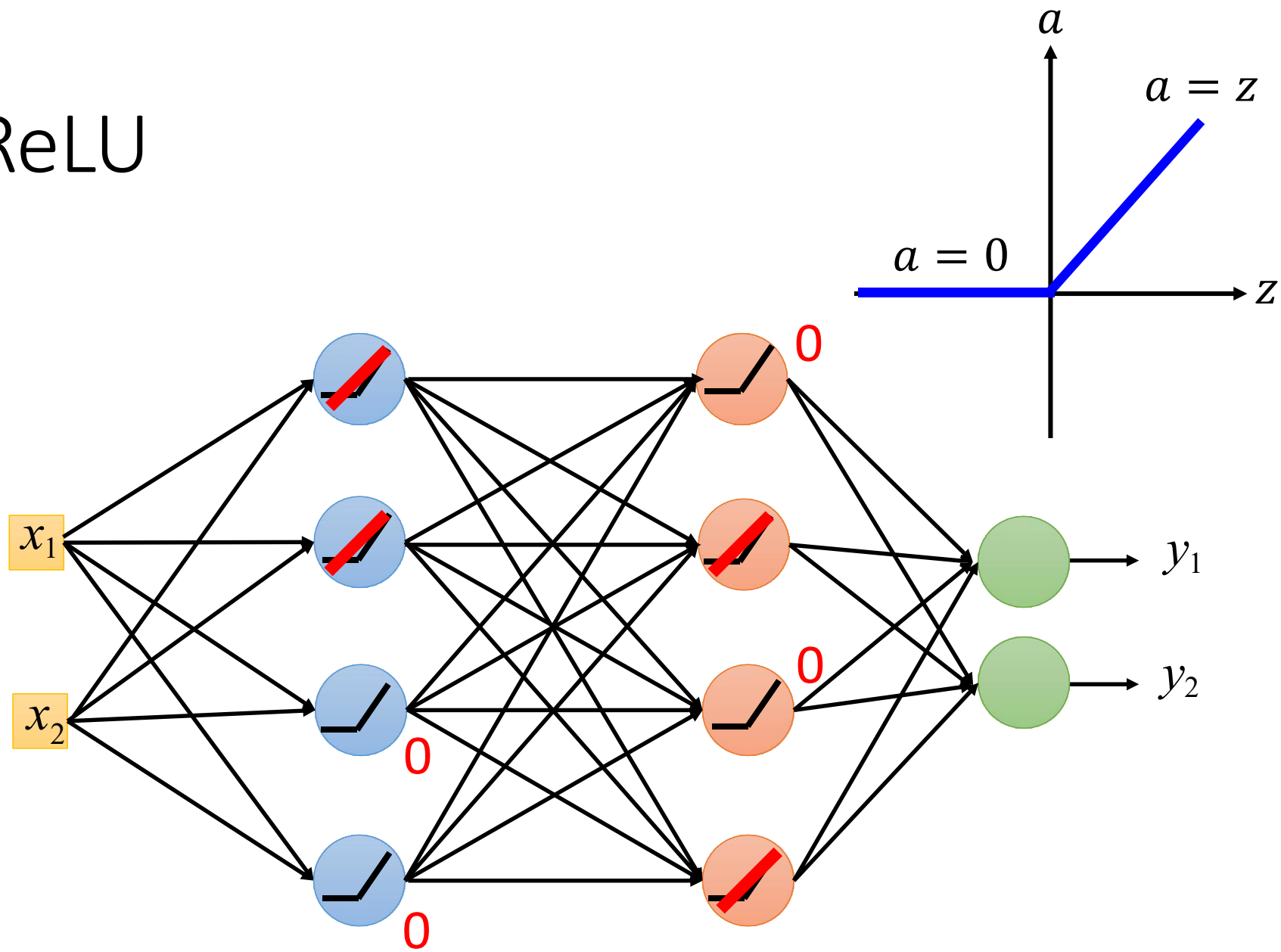


[Xavier Glorot, AISTATS'11]
[Andrew L. Maas, ICML'13]
[Kaiming He, arXiv'15]

Reason:

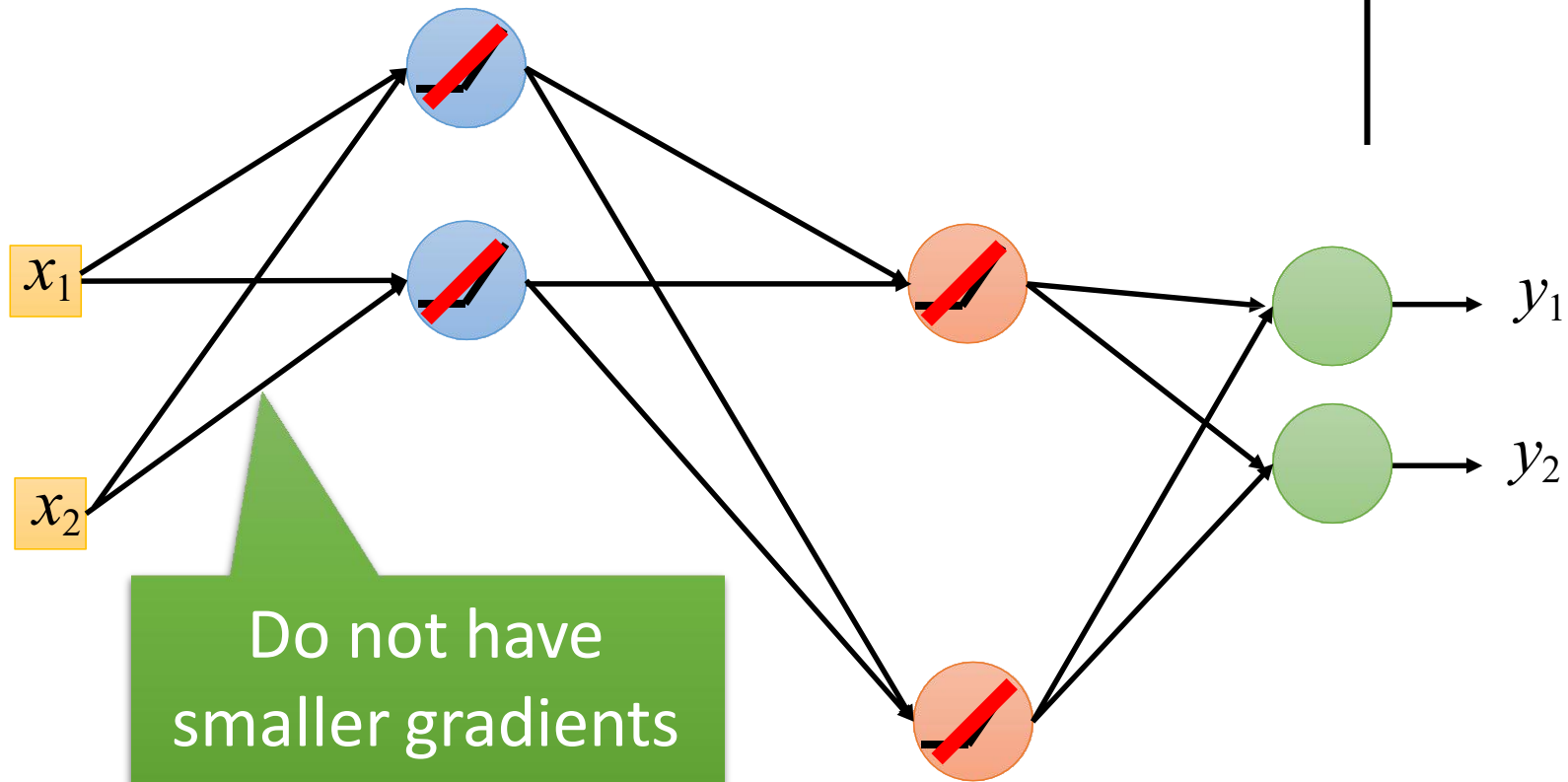
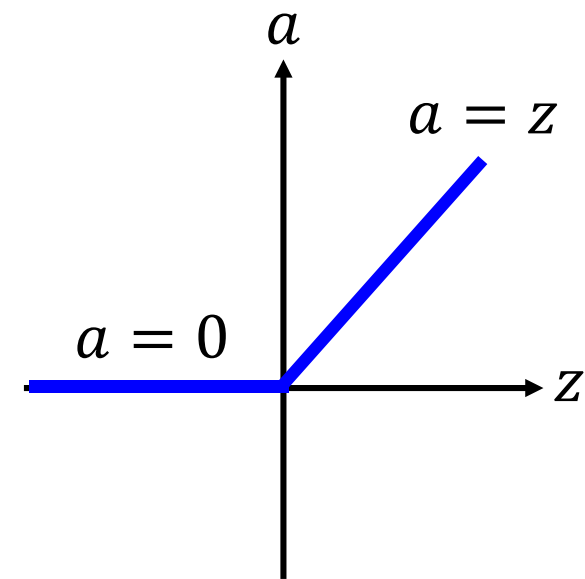
1. Fast to compute
2. Biological reason
3. Infinite sigmoid with different biases
4. Vanishing gradient problem

ReLU



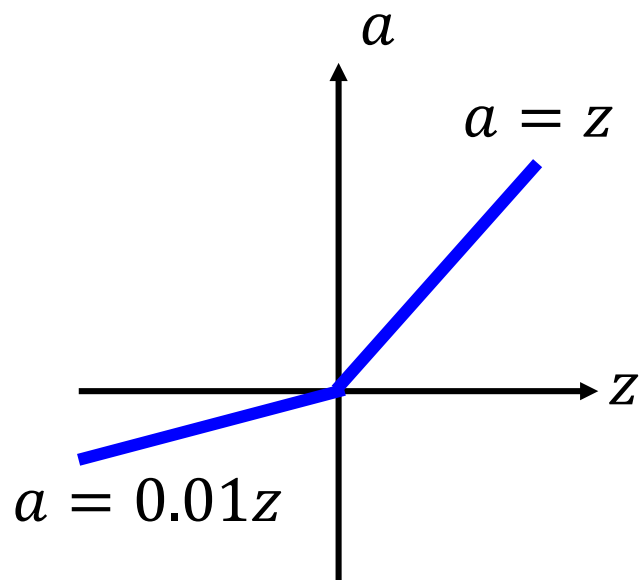
ReLU

A Thinner linear network

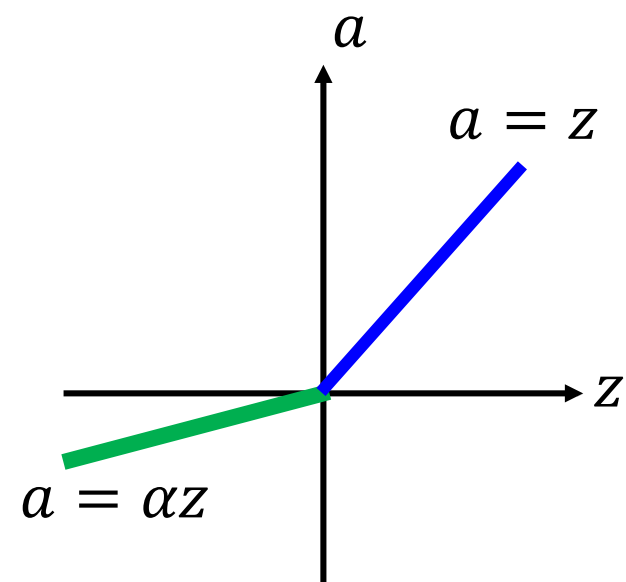


ReLU - variant



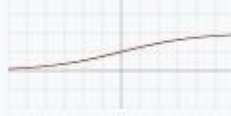

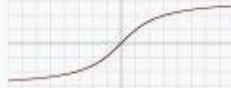

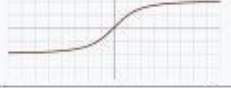




Leaky ReLU





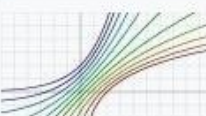




Parametric ReLU



α also learned by
gradient descent

Name	Plot	Equation
Identity		$f(x) = x$
Binary step		$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$
Logistic (a.k.a. Sigmoid or Soft step)		$f(x) = \sigma(x) = \frac{1}{1 + e^{-x}}^{[1]}$
TanH		$f(x) = \tanh(x) = \frac{(e^x - e^{-x})}{(e^x + e^{-x})}$
ArcTan		$f(x) = \tan^{-1}(x)$
Softsign ^{[9][10]}		$f(x) = \frac{x}{1 + x }$
Inverse square root unit (ISRU) ^[11]		$f(x) = \frac{x}{\sqrt{1 + \alpha x^2}}$
Rectified linear unit (ReLU) ^[12]		$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$
Leaky rectified linear unit (Leaky ReLU) ^[13]		$f(x) = \begin{cases} 0.01x & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$
Parametric rectified linear unit (PReLU) ^[14]		$f(\alpha, x) = \begin{cases} \alpha x & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$
Randomized leaky rectified linear unit (RReLU) ^[15]		$f(\alpha, x) = \begin{cases} \alpha x & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}^{[3]}$

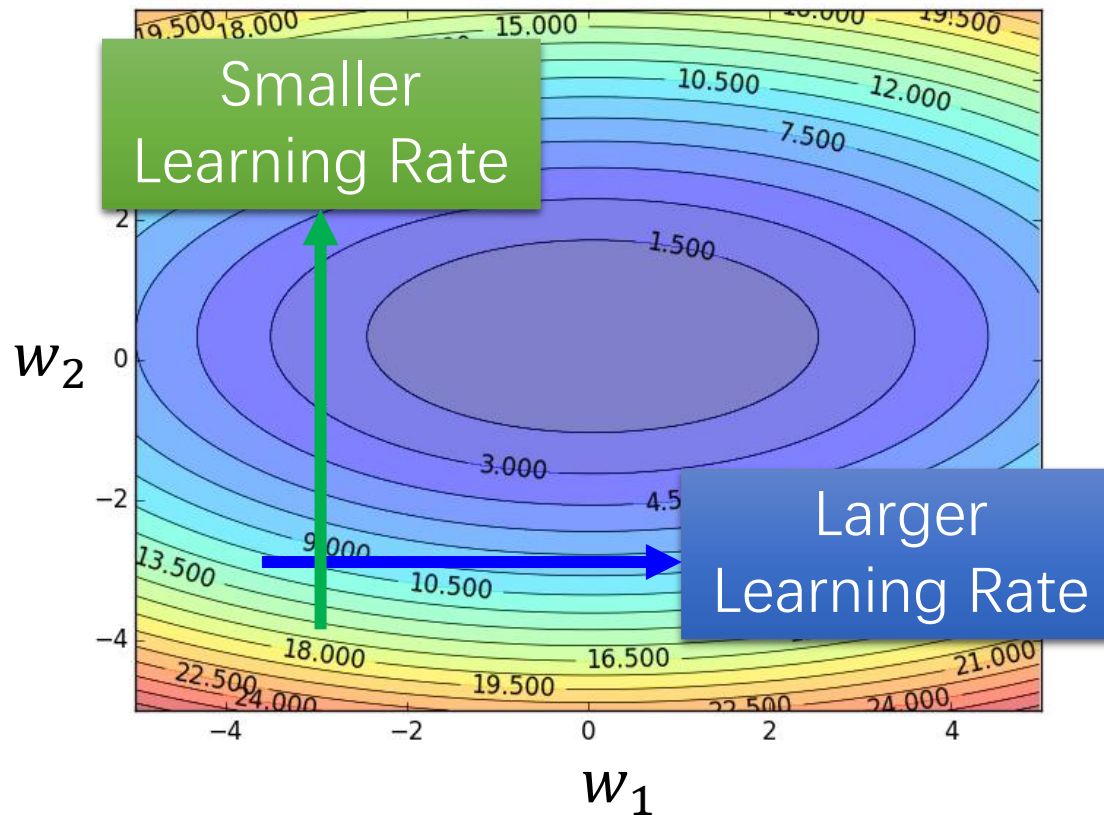
Exponential linear unit (ELU) ^[16]		$f(\alpha, x) = \begin{cases} \alpha(e^x - 1) & \text{for } x \leq 0 \\ x & \text{for } x > 0 \end{cases}$
Scaled exponential linear unit (SELU) ^[17]		$f(\alpha, x) = \lambda \begin{cases} \alpha(e^x - 1) & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$ with $\lambda = 1.0507$ and $\alpha = 1.67326$
S-shaped rectified linear activation unit (SReLU) ^[18]		$f_{t_l, a_l, t_r, a_r}(x) = \begin{cases} t_l + a_l(x - t_l) & \text{for } x \leq t_l \\ x & \text{for } t_l < x < t_r \\ t_r + a_r(x - t_r) & \text{for } x \geq t_r \end{cases}$ t_l, a_l, t_r, a_r are parameters.
Inverse square root linear unit (ISRLU) ^[11]		$f(x) = \begin{cases} \frac{x}{\sqrt{1+\alpha x^2}} & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$
Adaptive piecewise linear (APL) ^[19]		$f(x) = \max(0, x) + \sum_{s=1}^S a_i^s \max(0, -x + b_i^s)$
SoftPlus ^[20]		$f(x) = \ln(1 + e^x)$
Bent identity		$f(x) = \frac{\sqrt{x^2 + 1} - 1}{2} + x$
Sigmoid-weighted linear unit (SiLU) ^[21] (a.k.a. Swish ^[22])		$f(x) = x \cdot \sigma(x)^{[5]}$
SoftExponential ^[23]		$f(\alpha, x) = \begin{cases} -\frac{\ln(1-\alpha(x+\alpha))}{\alpha} & \text{for } \alpha < 0 \\ x & \text{for } \alpha = 0 \\ \frac{e^{\alpha x} - 1}{\alpha} + \alpha & \text{for } \alpha > 0 \end{cases}$
Sinusoid ^[24]		$f(x) = \sin(x)$
Sinc		$f(x) = \begin{cases} 1 & \text{for } x = 0 \\ \frac{\sin(x)}{x} & \text{for } x \neq 0 \end{cases}$

Tips for Training

3. Optimizer

Different parameters need different learning rate

Formulation for **one** parameter:



$$\theta_i^{t+1} \leftarrow \theta_i^t - \boxed{\eta} g_i^t$$

$$g_i^t = \frac{\partial L}{\partial \theta_i} \bigg|_{\theta = \theta^t}$$

$$\theta_i^{t+1} \leftarrow \theta_i^t - \boxed{\frac{\eta}{\sigma_i^t}} g_i^t$$

Parameter
dependent

Root Mean Square

$$\boldsymbol{\theta}_i^{t+1} \leftarrow \boldsymbol{\theta}_i^t - \boxed{\frac{\eta}{\sigma_i^t}} \mathbf{g}_i^t$$

$$\boldsymbol{\theta}_i^1 \leftarrow \boldsymbol{\theta}_i^0 - \frac{\eta}{\sigma_i^0} \mathbf{g}_i^0$$

$$\sigma_i^0 = \sqrt{(\mathbf{g}_i^0)^2} = |\mathbf{g}_i^0|$$

$$\boldsymbol{\theta}_i^2 \leftarrow \boldsymbol{\theta}_i^1 - \frac{\eta}{\sigma_i^1} \mathbf{g}_i^1$$

$$\sigma_i^1 = \sqrt{\frac{1}{2} [(\mathbf{g}_i^0)^2 + (\mathbf{g}_i^1)^2]}$$

$$\boldsymbol{\theta}_i^3 \leftarrow \boldsymbol{\theta}_i^2 - \frac{\eta}{\sigma_i^2} \mathbf{g}_i^2$$

$$\sigma_i^2 = \sqrt{\frac{1}{3} [(\mathbf{g}_i^0)^2 + (\mathbf{g}_i^1)^2 + (\mathbf{g}_i^2)^2]}$$

⋮

$$\boldsymbol{\theta}_i^{t+1} \leftarrow \boldsymbol{\theta}_i^t - \frac{\eta}{\sigma_i^t} \mathbf{g}_i^t$$

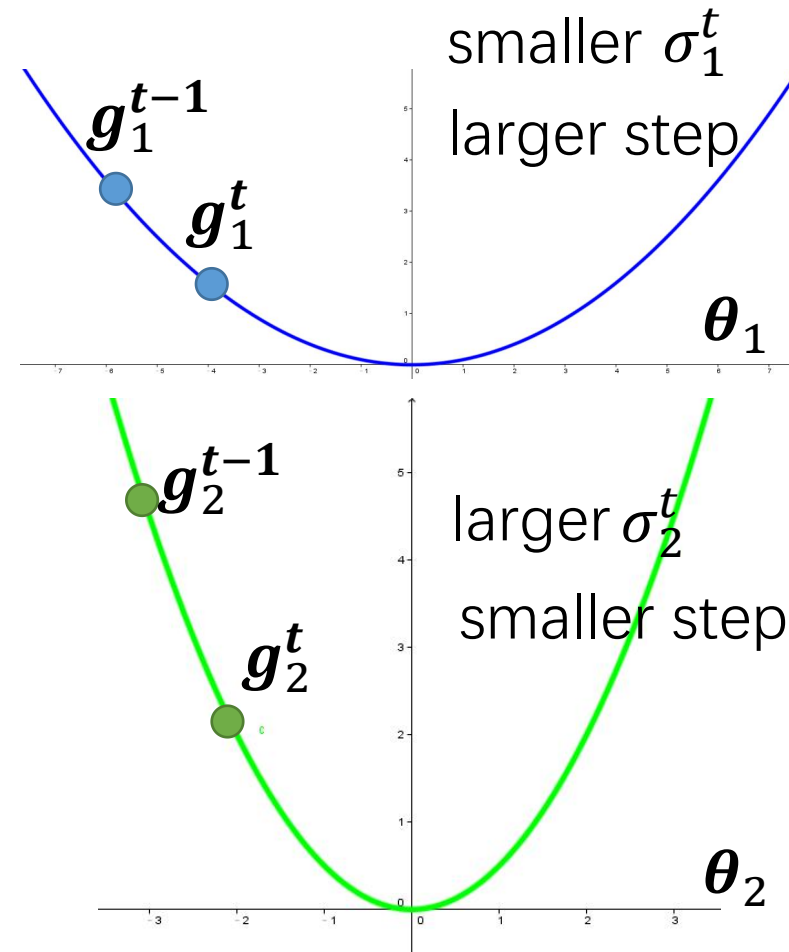
$$\sigma_i^t = \sqrt{\frac{1}{t+1} \sum_{j=0}^t (\mathbf{g}_i^j)^2}$$

Root Mean Square

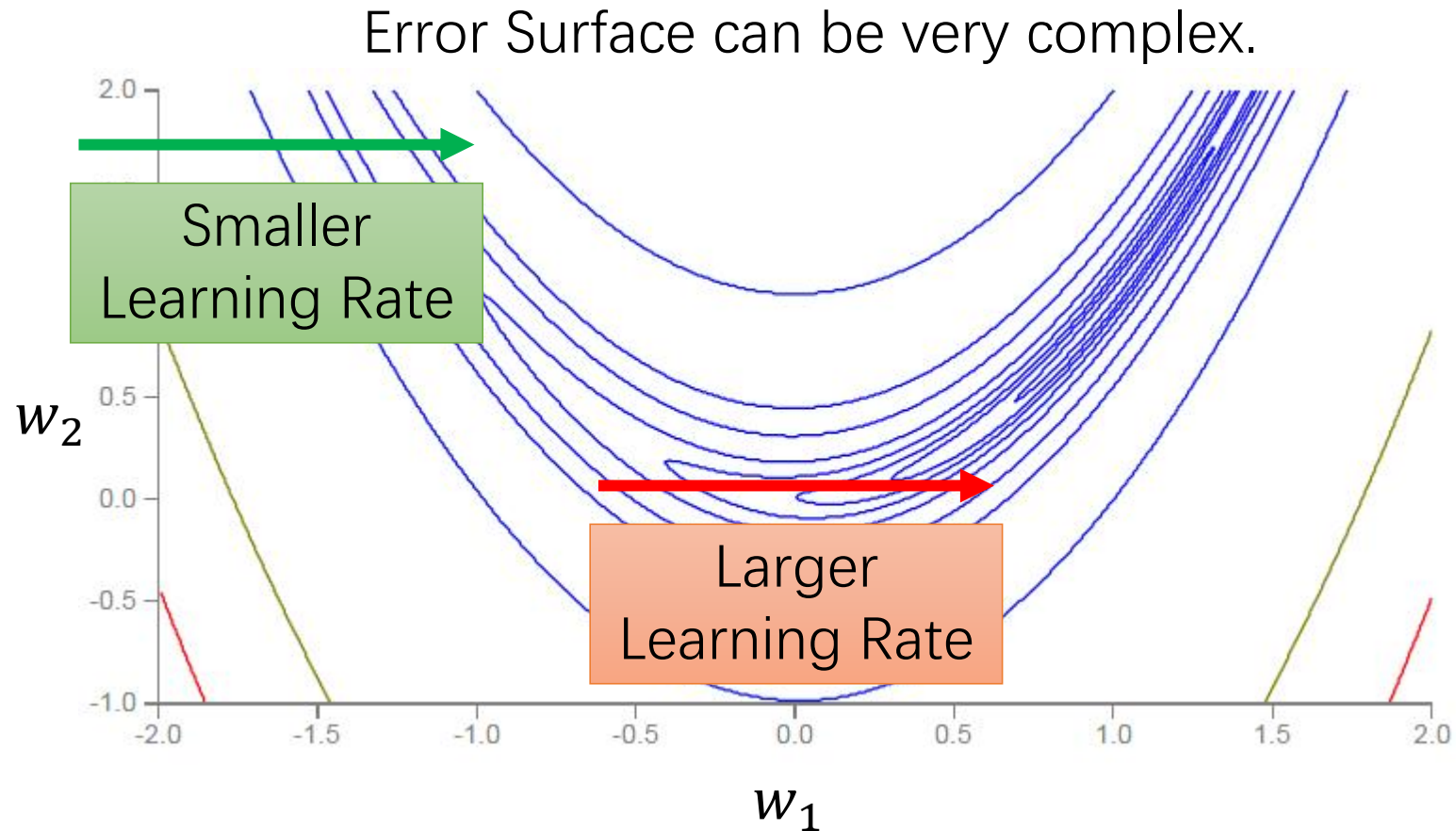
$$\theta_i^{t+1} \leftarrow \theta_i^t - \boxed{\frac{\eta}{\sigma_i^t}} g_i^t$$

$$\sigma_i^t = \sqrt{\frac{1}{t+1} \sum_{j=0}^t (g_i^j)^2}$$

Used in **Adagrad**



Learning rate adapts dynamically



RMSProp

$$\boldsymbol{\theta}_i^{t+1} \leftarrow \boldsymbol{\theta}_i^t - \boxed{\frac{\eta}{\sigma_i^t}} \mathbf{g}_i^t$$

$$\boldsymbol{\theta}_i^1 \leftarrow \boldsymbol{\theta}_i^0 - \frac{\eta}{\sigma_i^0} \mathbf{g}_i^0 \quad \sigma_i^0 = \sqrt{(\mathbf{g}_i^0)^2} \quad 0 < \alpha < 1$$

$$\boldsymbol{\theta}_i^2 \leftarrow \boldsymbol{\theta}_i^1 - \frac{\eta}{\sigma_i^1} \mathbf{g}_i^1 \quad \sigma_i^1 = \sqrt{\alpha(\sigma_i^0)^2 + (1 - \alpha)(\mathbf{g}_i^1)^2}$$

$$\boldsymbol{\theta}_i^3 \leftarrow \boldsymbol{\theta}_i^2 - \frac{\eta}{\sigma_i^2} \mathbf{g}_i^2 \quad \sigma_i^2 = \sqrt{\alpha(\sigma_i^1)^2 + (1 - \alpha)(\mathbf{g}_i^2)^2}$$

⋮

$$\boldsymbol{\theta}_i^{t+1} \leftarrow \boldsymbol{\theta}_i^t - \frac{\eta}{\sigma_i^t} \mathbf{g}_i^t \quad \sigma_i^t = \sqrt{\alpha(\sigma_i^{t-1})^2 + (1 - \alpha)(\mathbf{g}_i^t)^2}$$

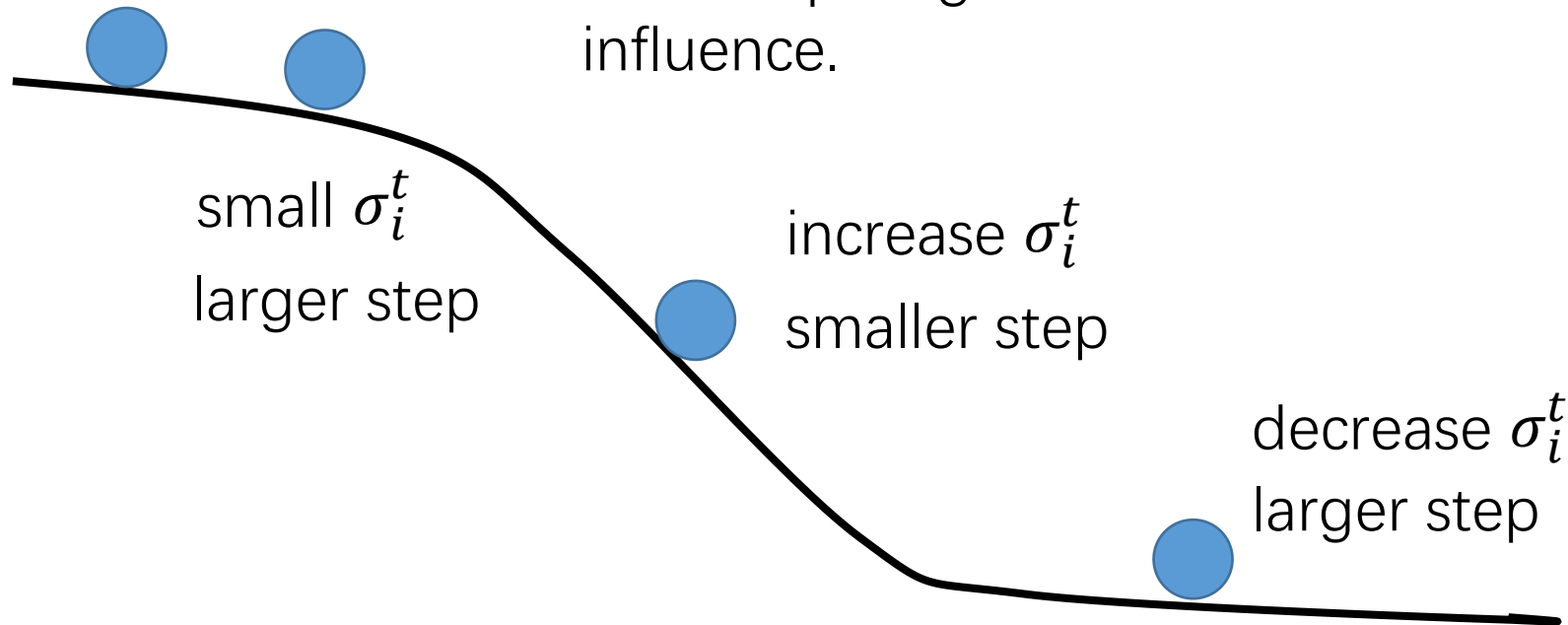
RMSProp

$$\theta_i^{t+1} \leftarrow \theta_i^t - \frac{\eta}{\sigma_i^t} g_i^t \quad \sigma_i^t = \sqrt{\alpha (\sigma_i^{t-1})^2 + (1 - \alpha) (g_i^t)^2} \quad 0 < \alpha < 1$$

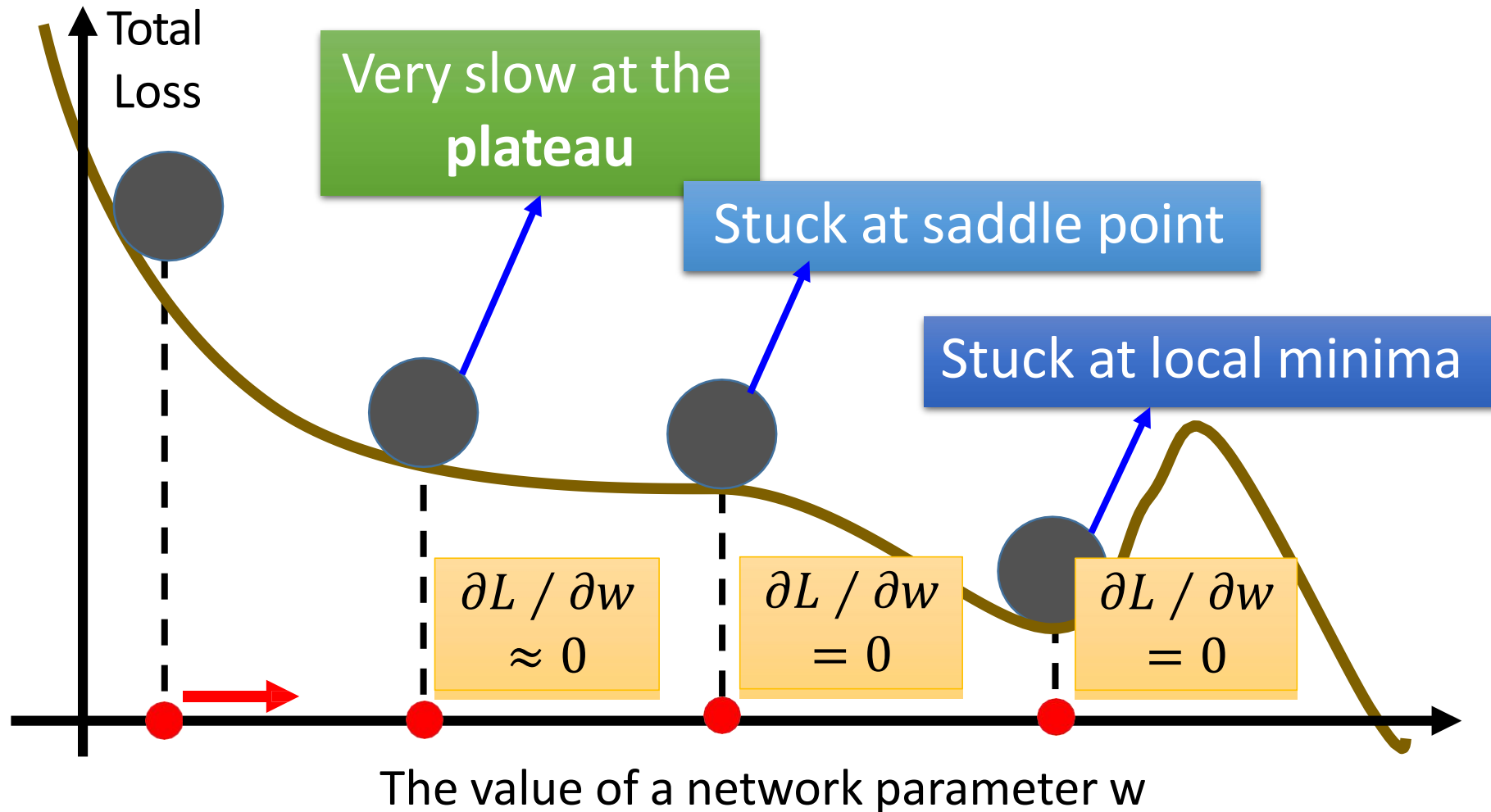
$g_i^1 \quad g_i^2 \quad \dots \quad g_i^{t-1}$

↑

The recent gradient has larger influence, and the past gradients have less influence.

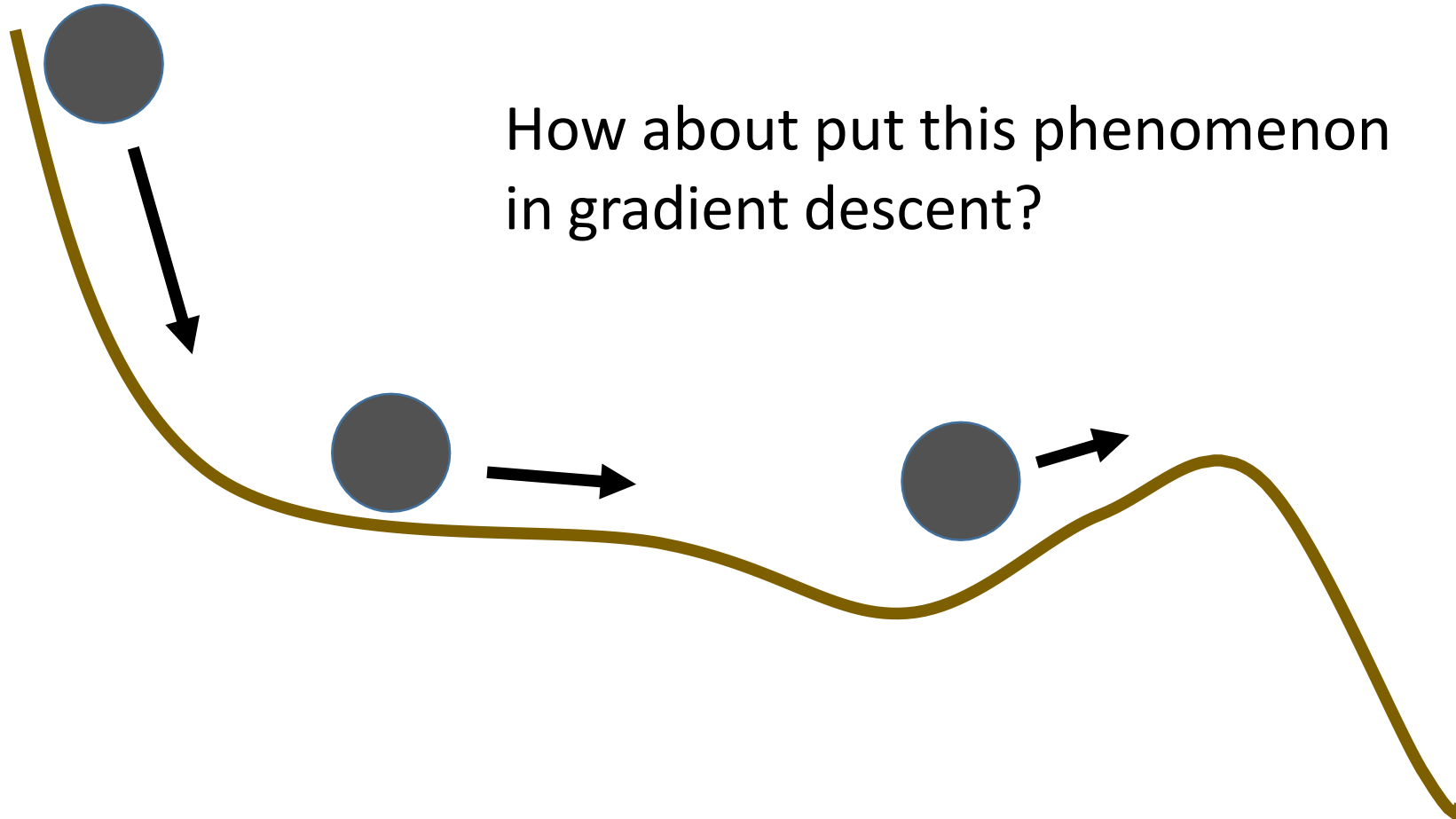


Hard to find optimal network parameters

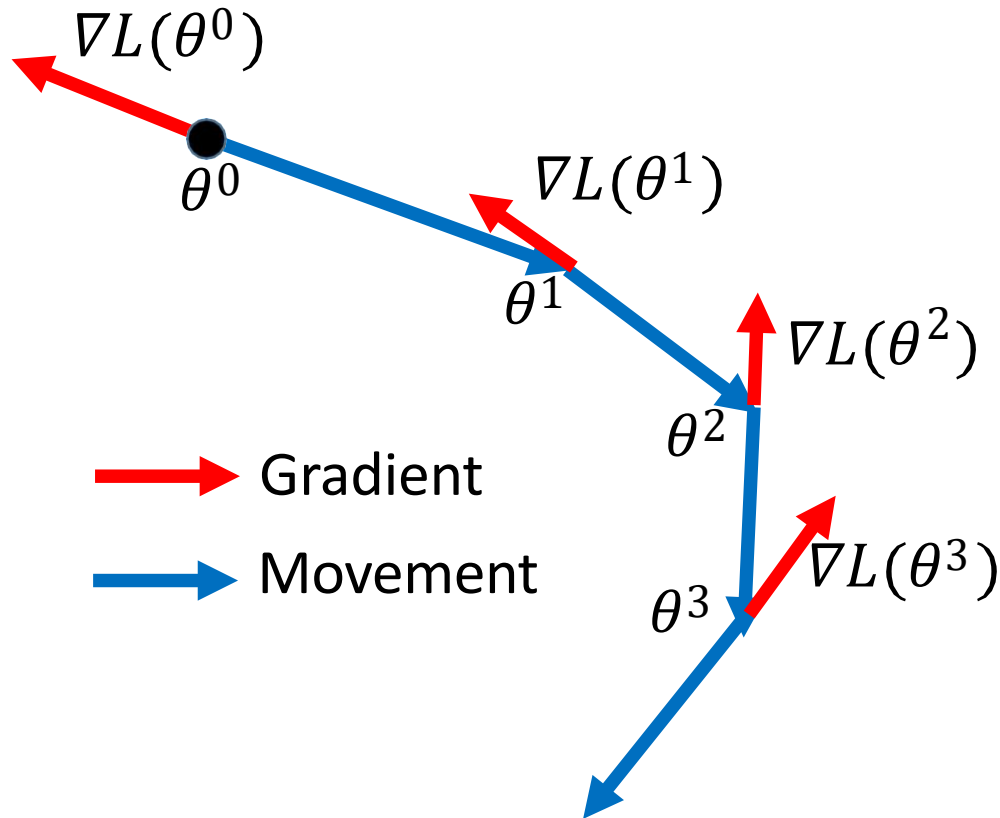


In physical world

- Momentum



Review: Vanilla Gradient Descent



Start at position θ^0

Compute gradient at θ^0

Move to $\theta^1 = \theta^0 - \eta \nabla L(\theta^0)$

Compute gradient at θ^1

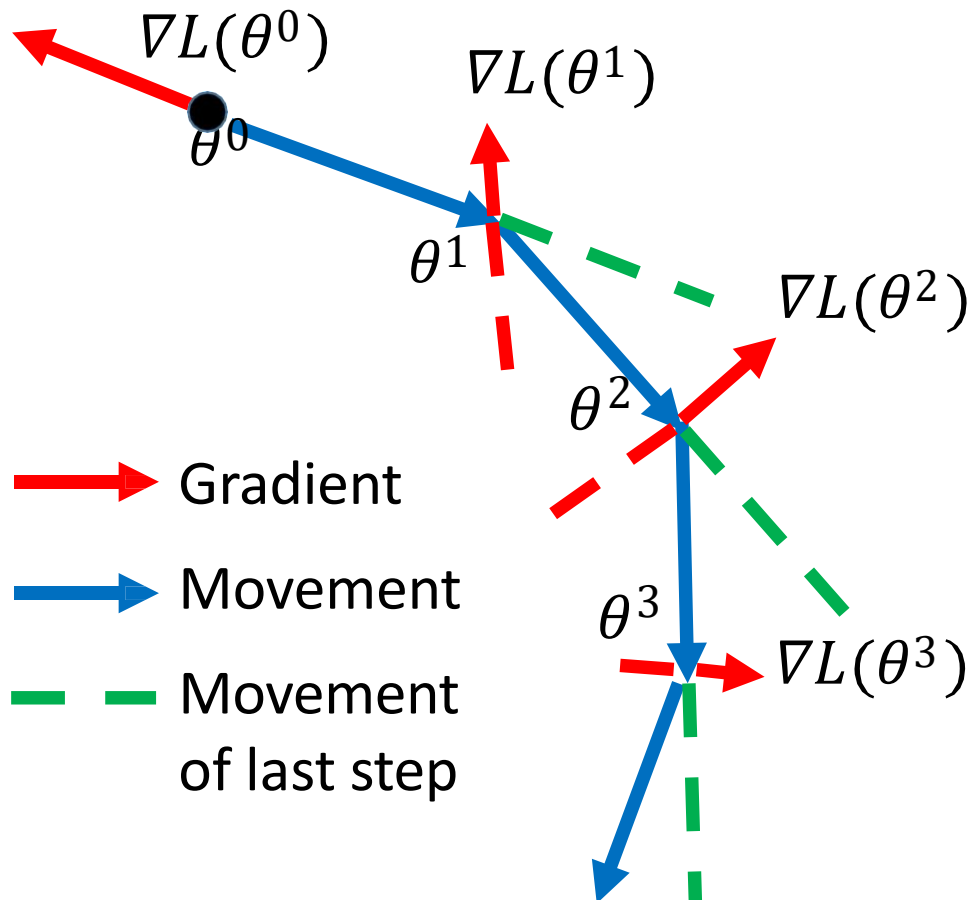
Move to $\theta^2 = \theta^1 - \eta \nabla L(\theta^1)$

⋮

Stop until $\nabla L(\theta^t) \approx 0$

Momentum

Movement: movement of last step minus gradient at present



- Start at point θ^0 Movement $v^0=0$

- Compute gradient at θ^0

Movement $v^1 = \lambda v^0 - \eta \nabla L(\theta^0)$

Move to $\theta^1 = \theta^0 + v^1$

- Compute gradient at θ^1

Movement $v^2 = \lambda v^1 - \eta \nabla L(\theta^1)$

- Move to $\theta^2 = \theta^1 + v^2$

Movement not just based on gradient, but previous movement.

Momentum

Movement: movement of last step minus gradient at present

v^i is actually the weighted sum of all the previous gradient:

$$\nabla L(\theta^0), \nabla L(\theta^1), \dots \nabla L(\theta^{i-1})$$

$$v^0 = 0$$

$$v^1 = -\eta \nabla L(\theta^0)$$

$$v^2 = -\lambda \eta \nabla L(\theta^0) - \eta \nabla L(\theta^1)$$

\vdots

- Start at point θ^0 Movement $v^0=0$

- Compute gradient at θ^0

$$\text{Movement } v^1 = \lambda v^0 - \eta \nabla L(\theta^0)$$

$$\text{Move to } \theta^1 = \theta^0 + v^1$$

- Compute gradient at θ^1

$$\text{Movement } v^2 = \lambda v^1 - \eta \nabla L(\theta^1)$$

- Move to $\theta^2 = \theta^1 + v^2$

Movement not just based on gradient, but previous movement

Adam: RMSProp + Momentum

Algorithm 1: *Adam*, our proposed algorithm for stochastic optimization. See section 2 for details, and for a slightly more efficient (but less clear) order of computation. g_t^2 indicates the elementwise square $g_t \odot g_t$. Good default settings for the tested machine learning problems are $\alpha = 0.001$, $\beta_1 = 0.9$, $\beta_2 = 0.999$ and $\epsilon = 10^{-8}$. All operations on vectors are element-wise. With β_1^t and β_2^t we denote β_1 and β_2 to the power t .

Require: α : Stepsize

Require: $\beta_1, \beta_2 \in [0, 1)$: Exponential decay rates for the moment estimates

Require: $f(\theta)$: Stochastic objective function with parameters θ

Require: θ_0 : Initial parameter vector

$m_0 \leftarrow 0$ (Initialize 1st moment vector) \rightarrow for momentum

$v_0 \leftarrow 0$ (Initialize 2nd moment vector) \rightarrow for RMSprop

$t \leftarrow 0$ (Initialize timestep)

while θ_t not converged **do**

$t \leftarrow t + 1$

$g_t \leftarrow \nabla_{\theta} f_t(\theta_{t-1})$ (Get gradients w.r.t. stochastic objective at timestep t)

$m_t \leftarrow \beta_1 \cdot m_{t-1} + (1 - \beta_1) \cdot g_t$ (Update biased first moment estimate)

$v_t \leftarrow \beta_2 \cdot v_{t-1} + (1 - \beta_2) \cdot g_t^2$ (Update biased second raw moment estimate)

$\hat{m}_t \leftarrow m_t / (1 - \beta_1^t)$ (Compute bias-corrected first moment estimate)

$\hat{v}_t \leftarrow v_t / (1 - \beta_2^t)$ (Compute bias-corrected second raw moment estimate)

$\theta_t \leftarrow \theta_{t-1} - \alpha \cdot \hat{m}_t / (\sqrt{\hat{v}_t} + \epsilon)$ (Update parameters)

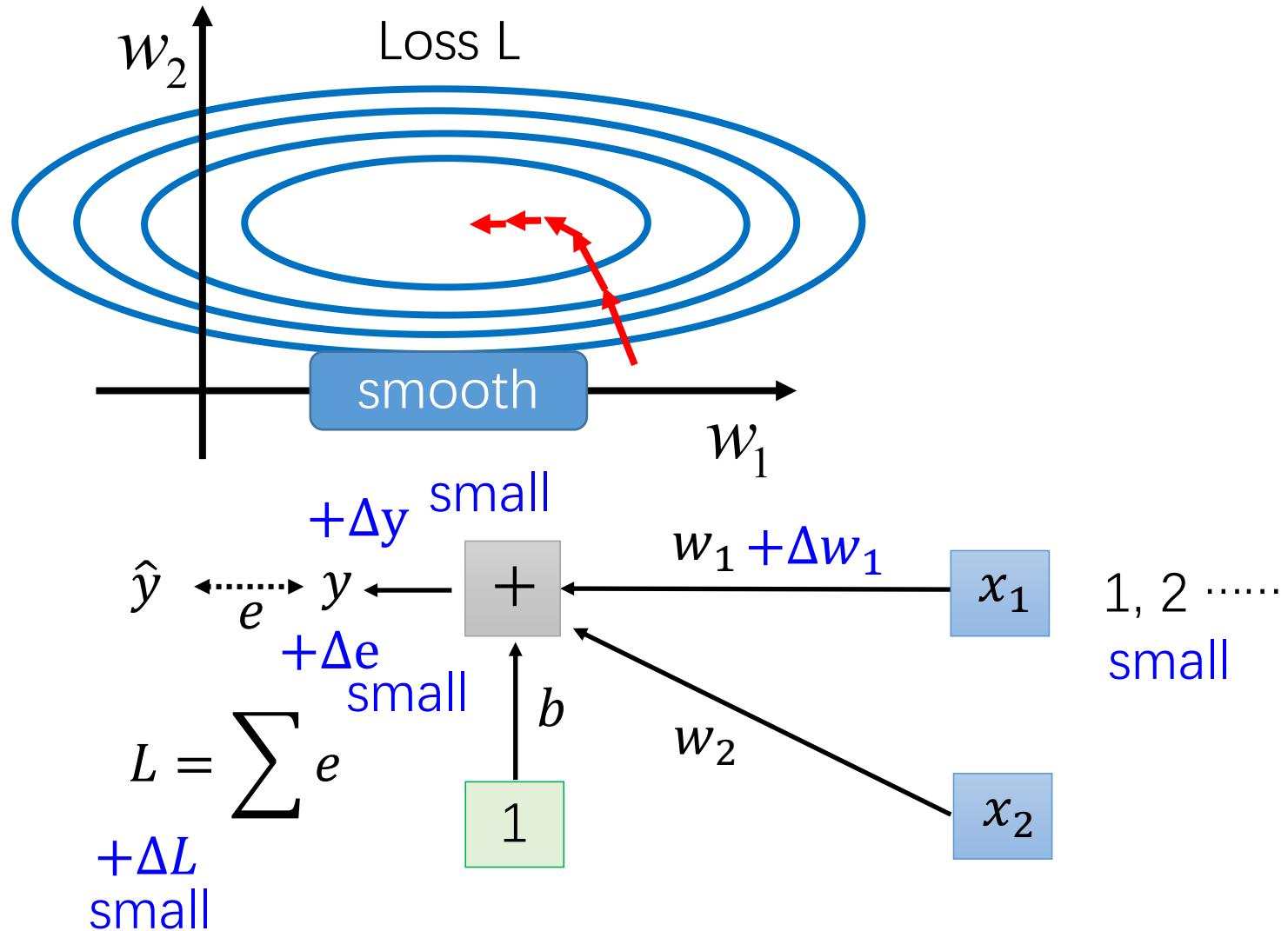
end while

return θ_t (Resulting parameters)

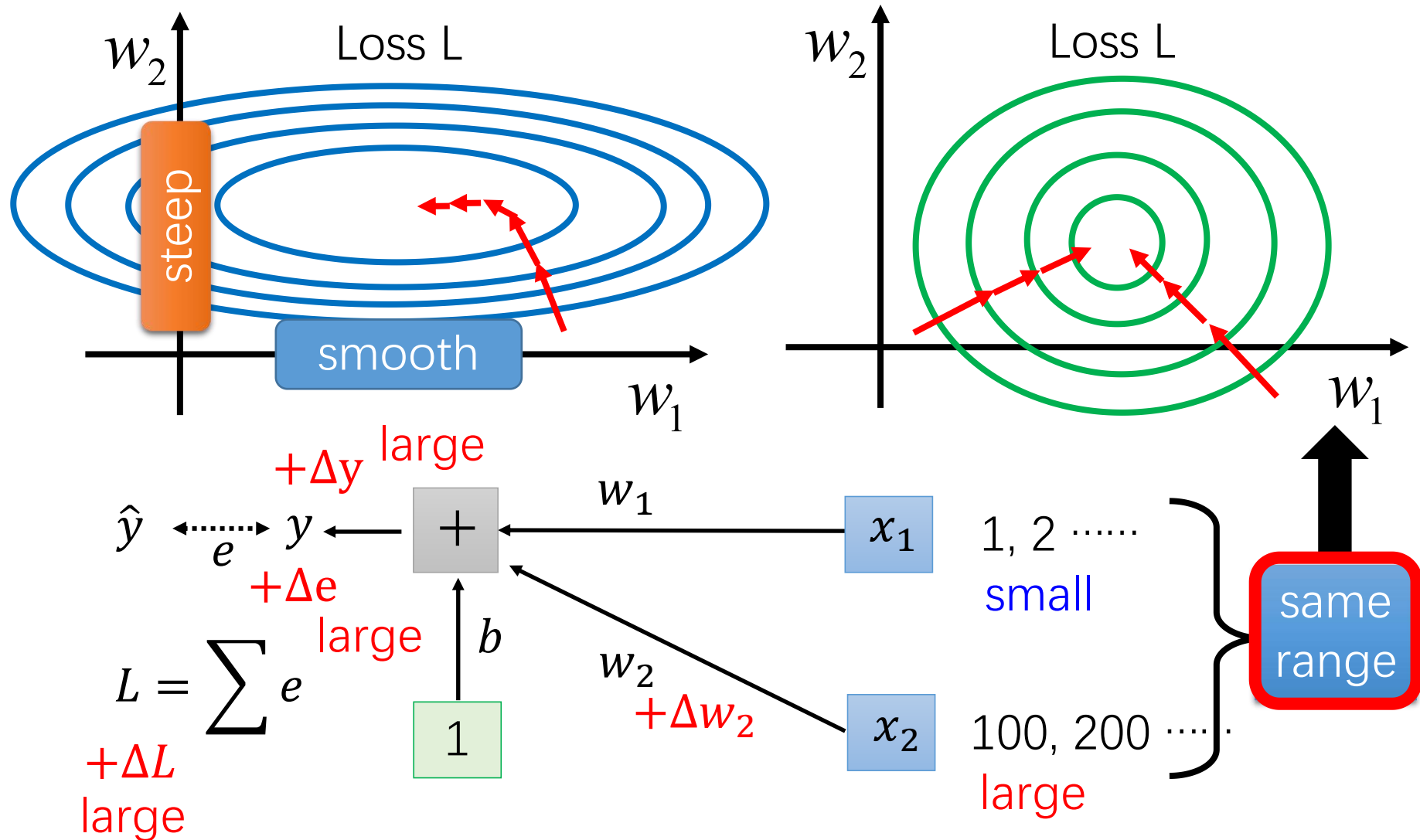
Tips for Training

4. Normalization

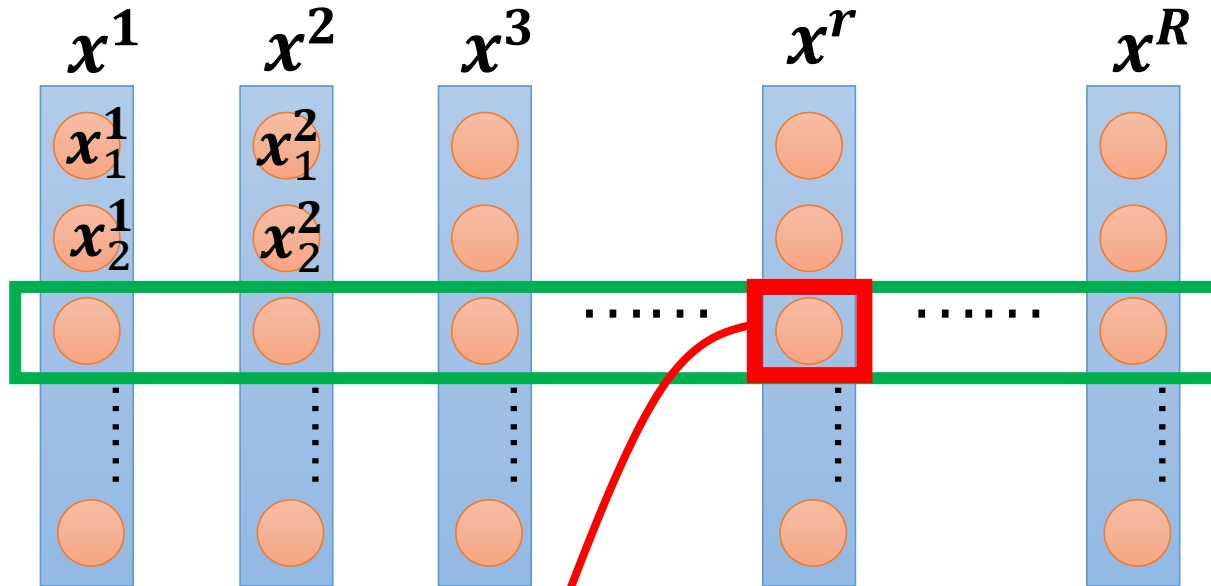
Changing Landscape



Changing Landscape



Feature Normalization



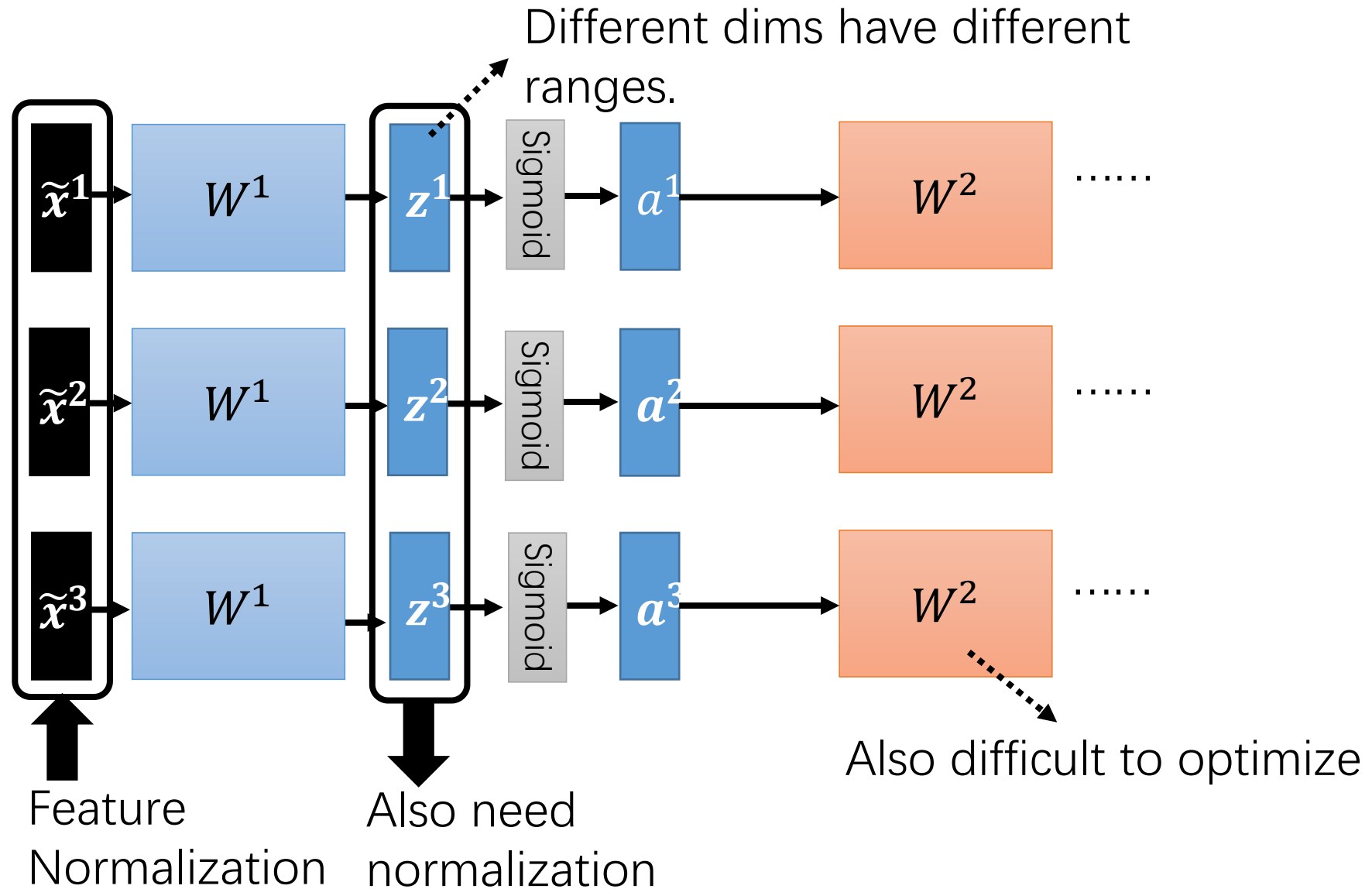
For each dimension i :
mean: m_i
standard deviation: σ_i

$$\tilde{x}_i^r \leftarrow \frac{x_i^r - m_i}{\sigma_i}$$

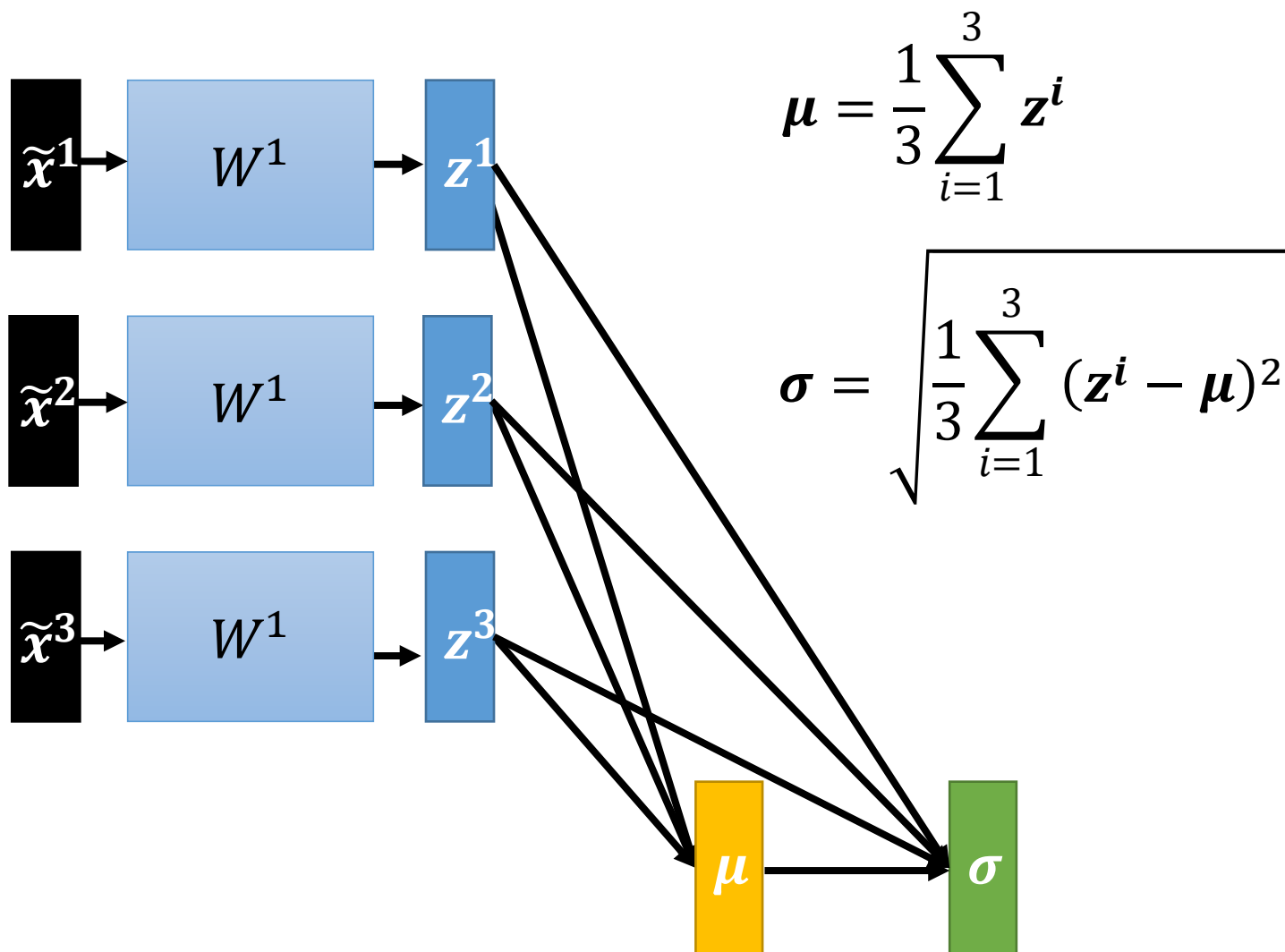
The means of all dims are 0,
and the variances are all 1

In general, feature normalization makes gradient descent converge faster.

Considering Deep Learning



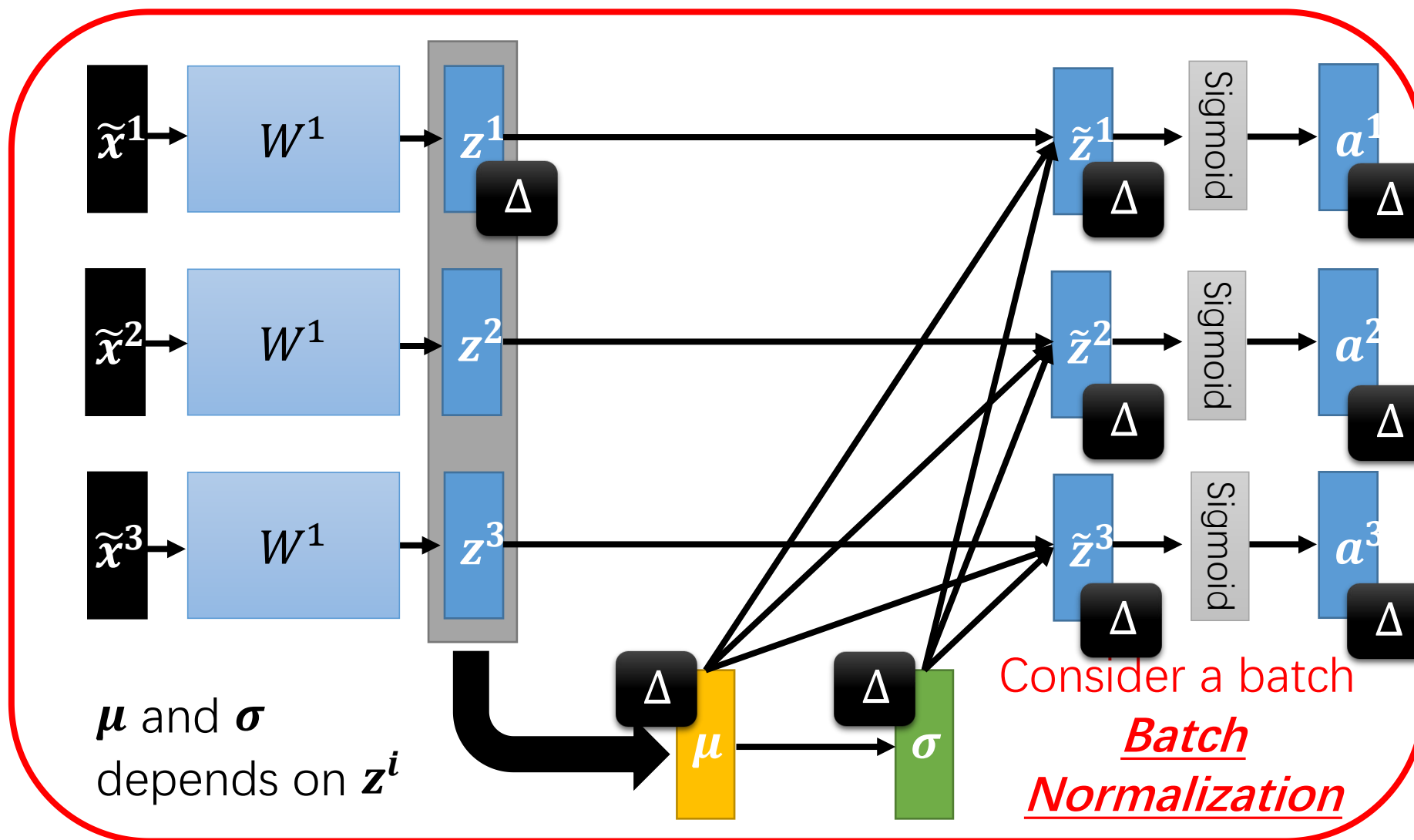
Considering Deep Learning



Considering Deep Learning

$$\tilde{z}^i = \frac{z^i - \mu}{\sigma}$$

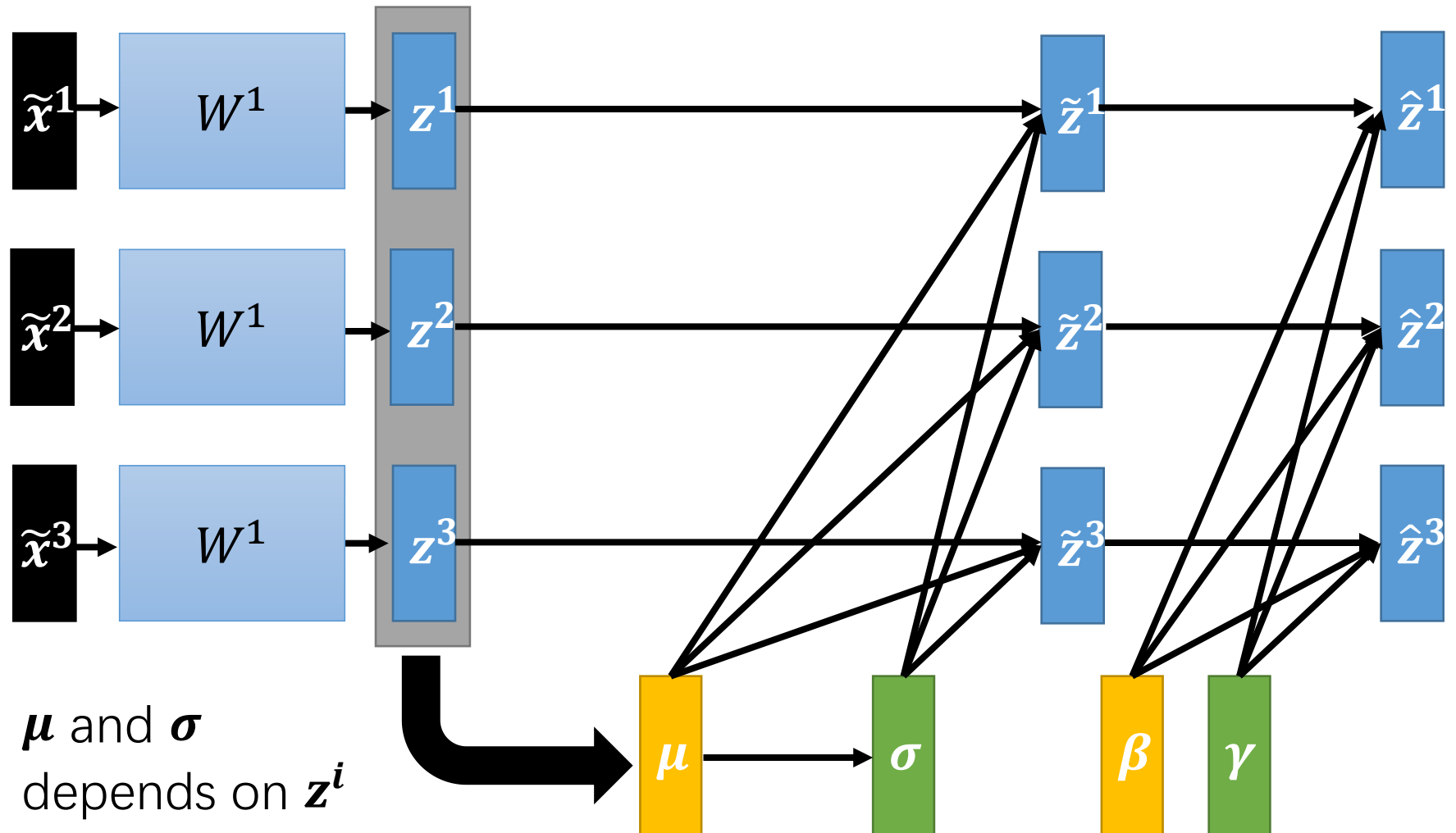
This is a large network!



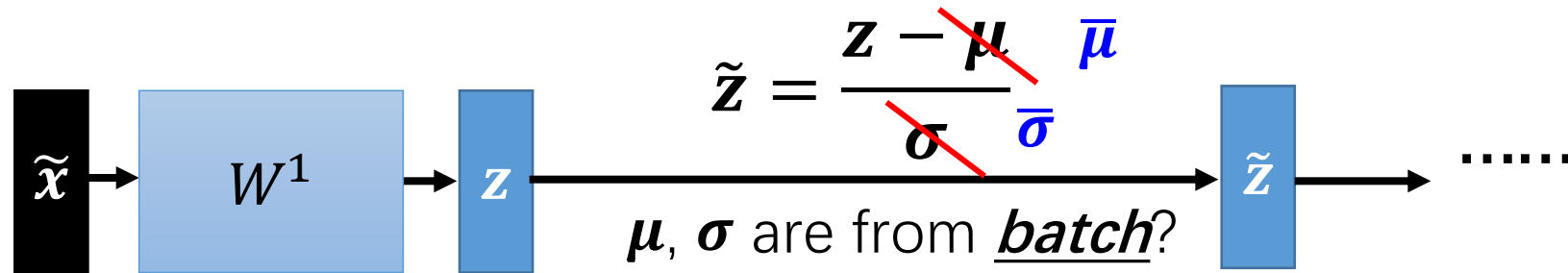
Batch normalization

$$\tilde{\mathbf{z}}^i = \frac{\mathbf{z}^i - \boldsymbol{\mu}}{\boldsymbol{\sigma}}$$

$$\hat{\mathbf{z}}^i = \boldsymbol{\gamma} \odot \tilde{\mathbf{z}}^i + \boldsymbol{\beta}$$



Batch normalization – Testing



We do not always have batch at testing stage.

Computing the moving average of μ and σ of the batches during training.

$$\mu^1 \quad \mu^2 \quad \mu^3 \quad \dots \quad \mu^t$$

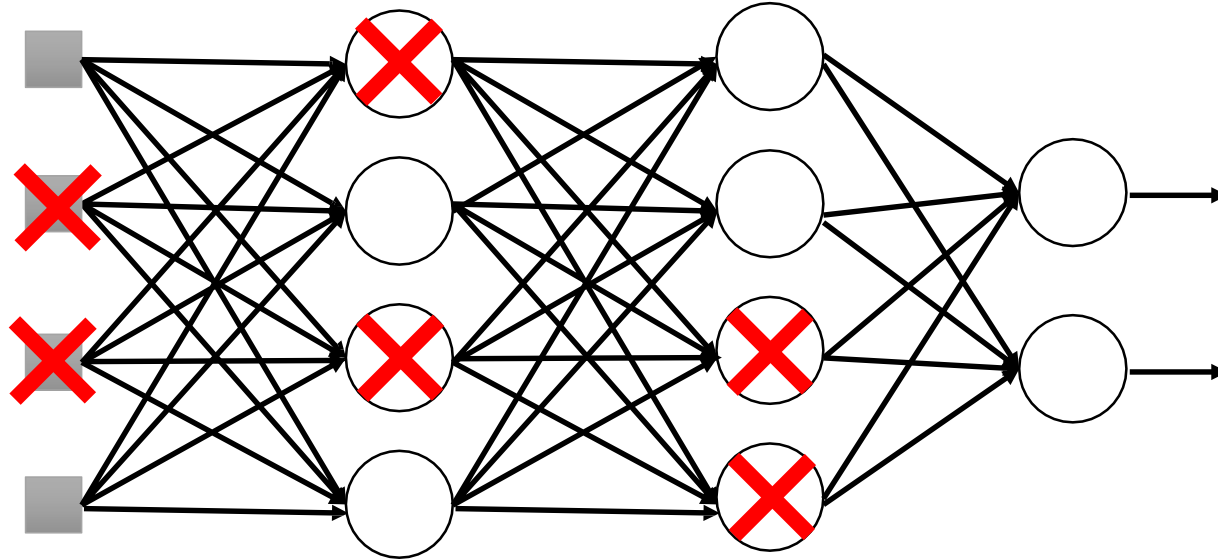
$$\bar{\mu} \leftarrow p\bar{\mu} + (1 - p)\mu^t$$

Tips for Training

5. Regularization

Dropout

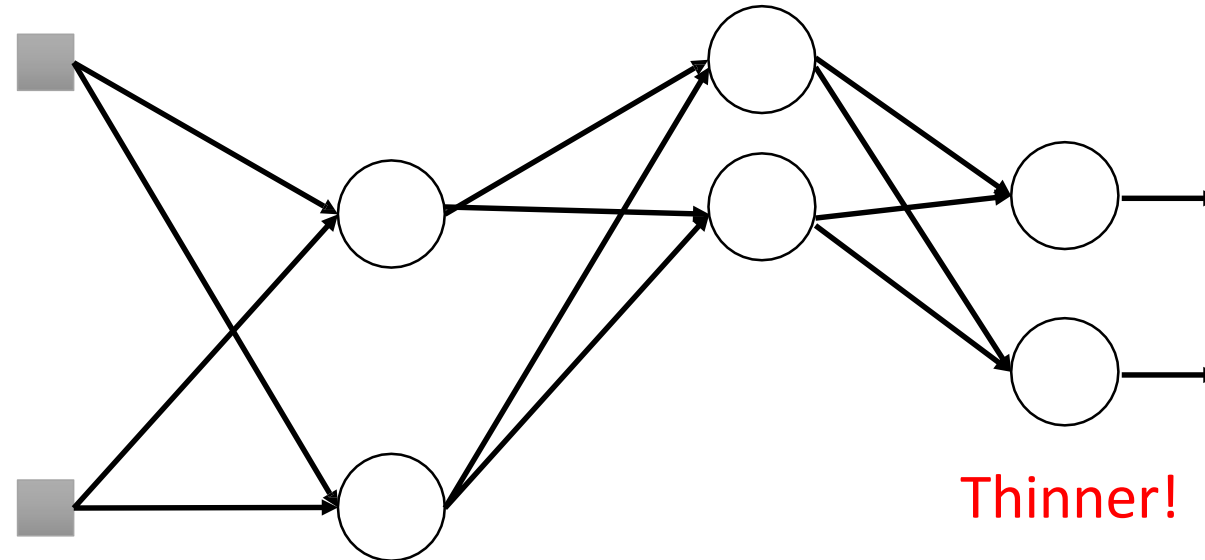
Training:



- **Each time before updating the parameters**
 - Each neuron has $p\%$ to dropout

Dropout

Training:

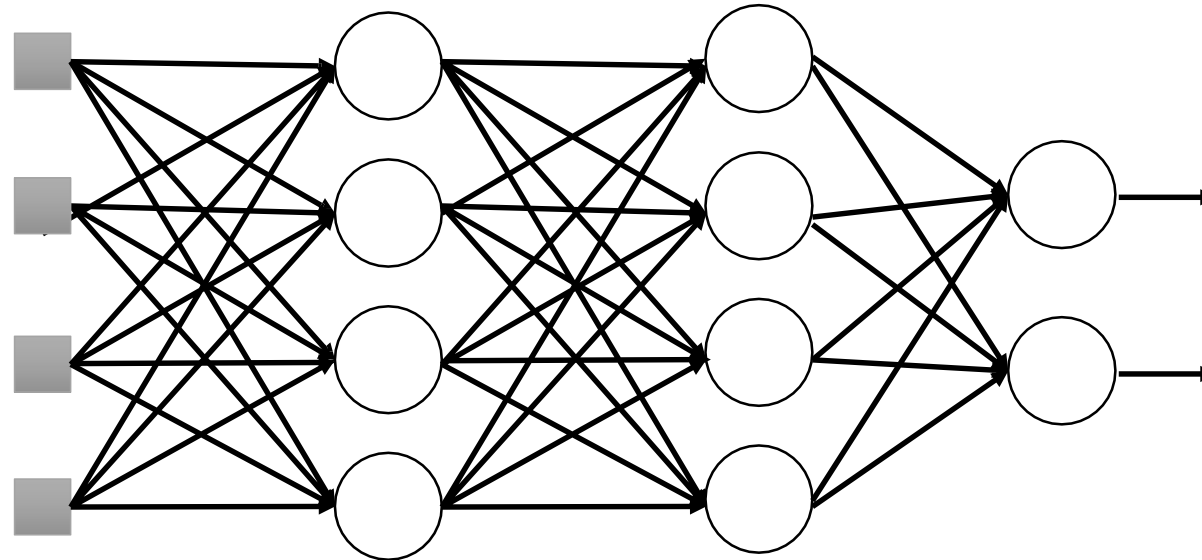


- **Each time before updating the parameters**
 - Each neuron has $p\%$ to dropout
 - ➡ **The structure of the network is changed.**
 - Using the new network for training

For each batch, we resample the dropout neurons

Dropout

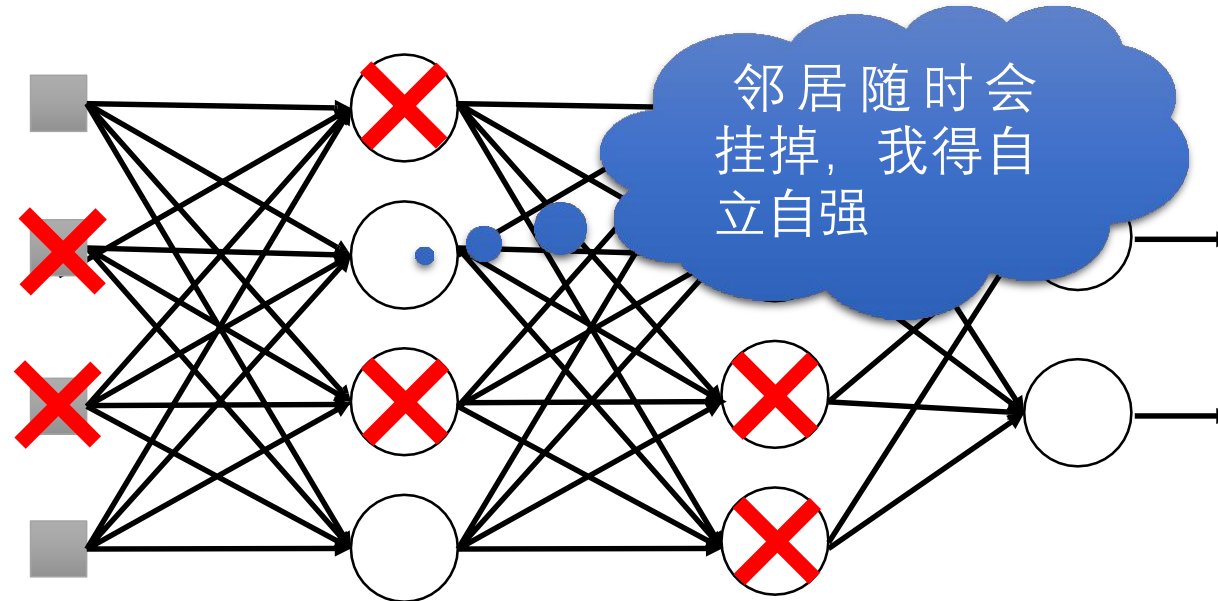
Testing:



➤ No dropout

- If the dropout rate at training is $p\%$, all the weights times $1-p\%$
- Assume that the dropout rate is 50%.
If a weight $w = 1$ by training, set $w = 0.5$ for testing.

Dropout - Intuitive Reason



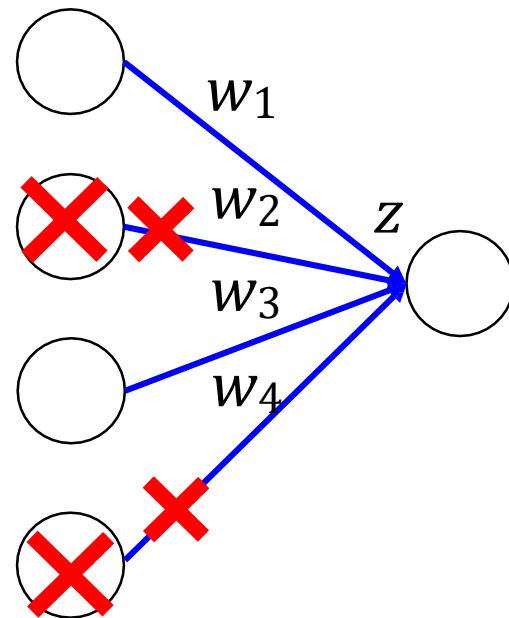
- When teams up, if everyone expect the partner will do the work, nothing will be done finally.
- However, if you know your partner will dropout, you will do better.
- When testing, no one dropout actually, so obtaining good results eventually.

Dropout - Intuitive Reason

- Why the weights should multiply (1-p)% (dropout rate) when testing?

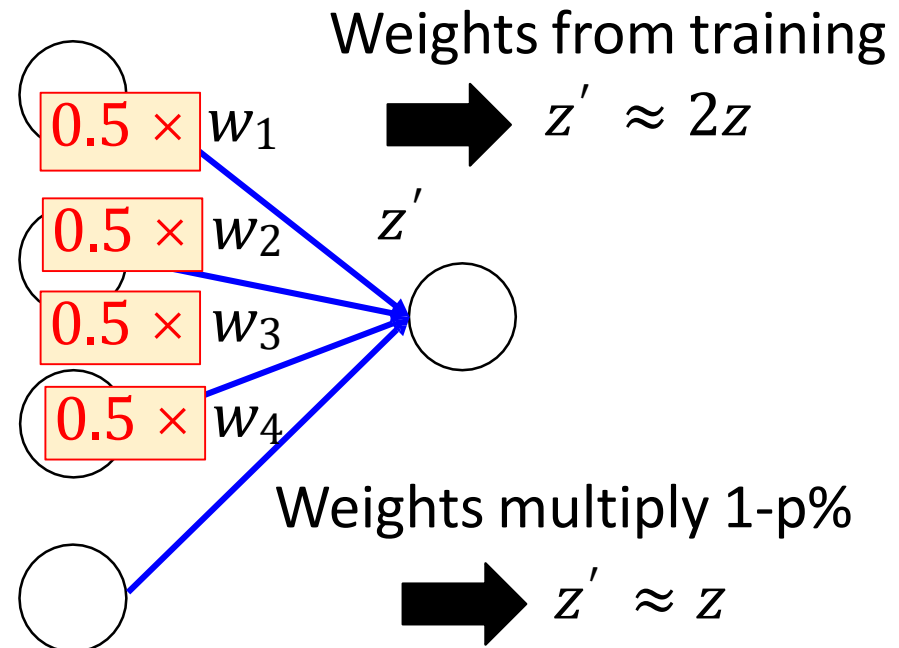
Training of Dropout

Assume dropout rate is 50%

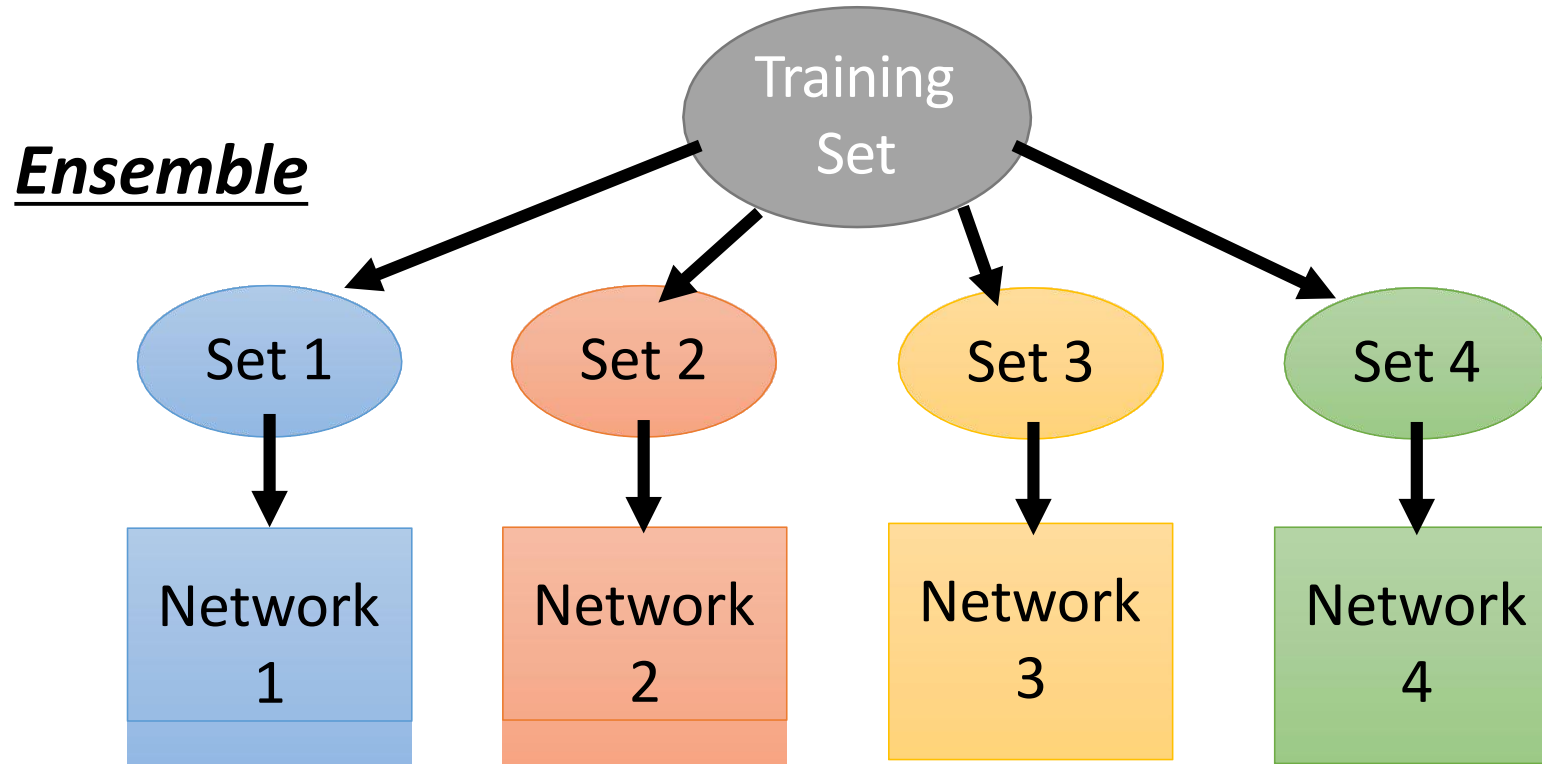


Testing of Dropout

No dropout



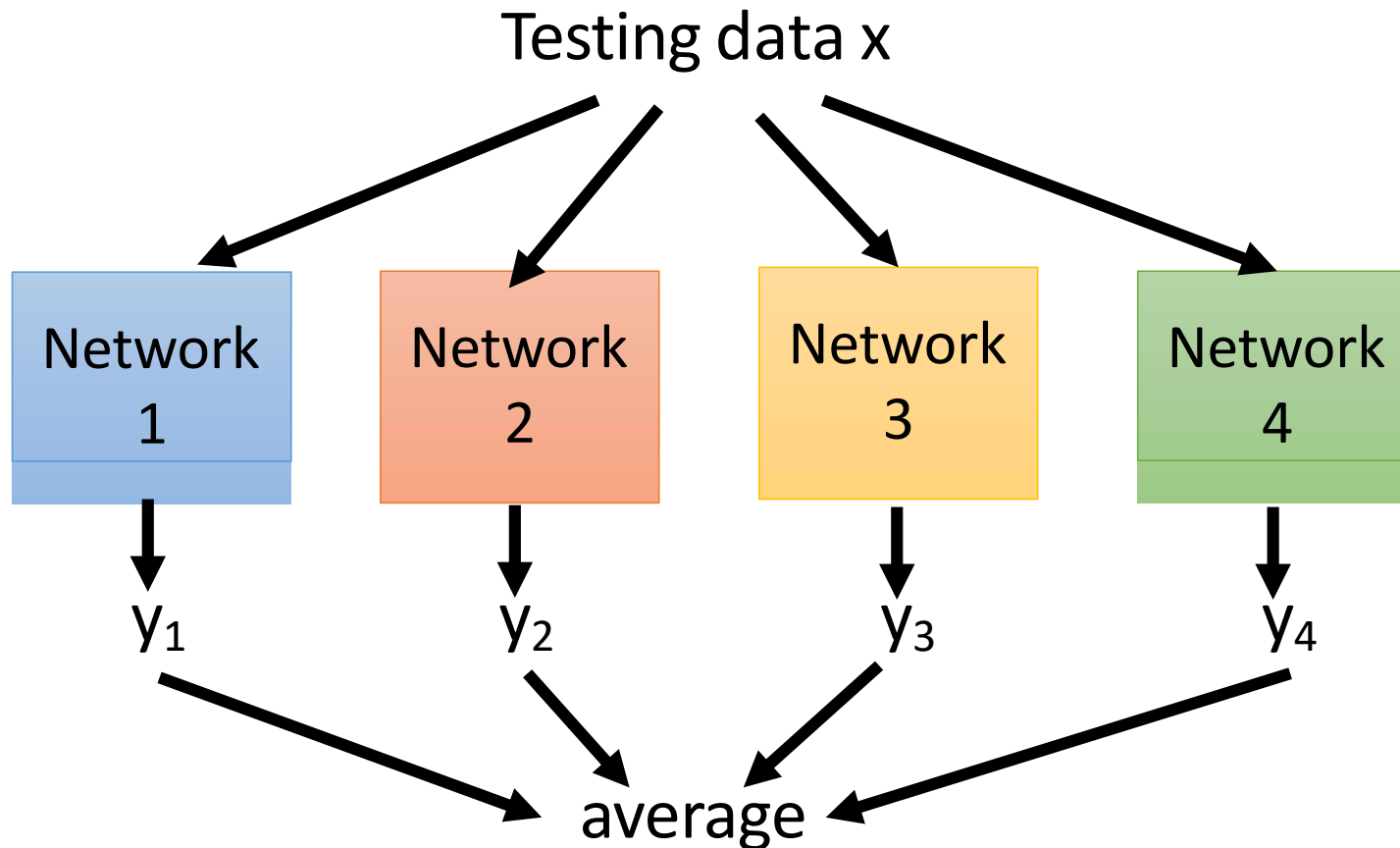
Dropout is a kind of ensemble.



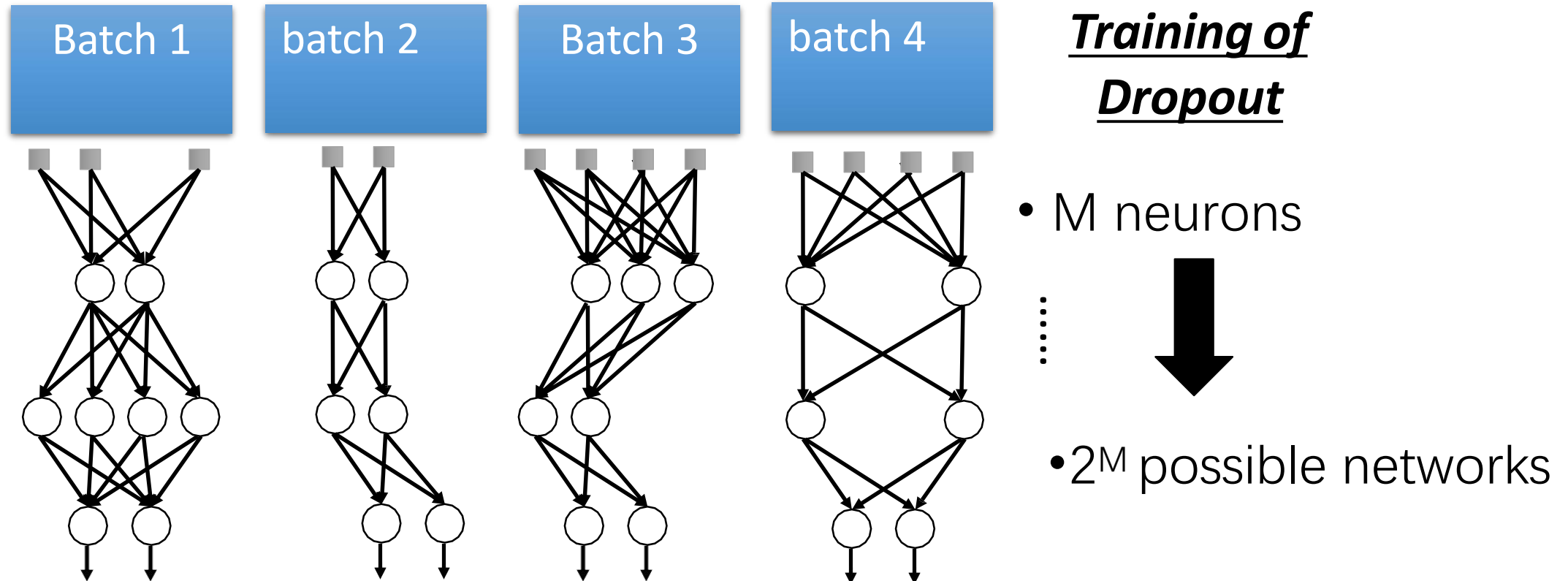
Train a bunch of networks with different structures

Dropout is a kind of ensemble.

Ensemble



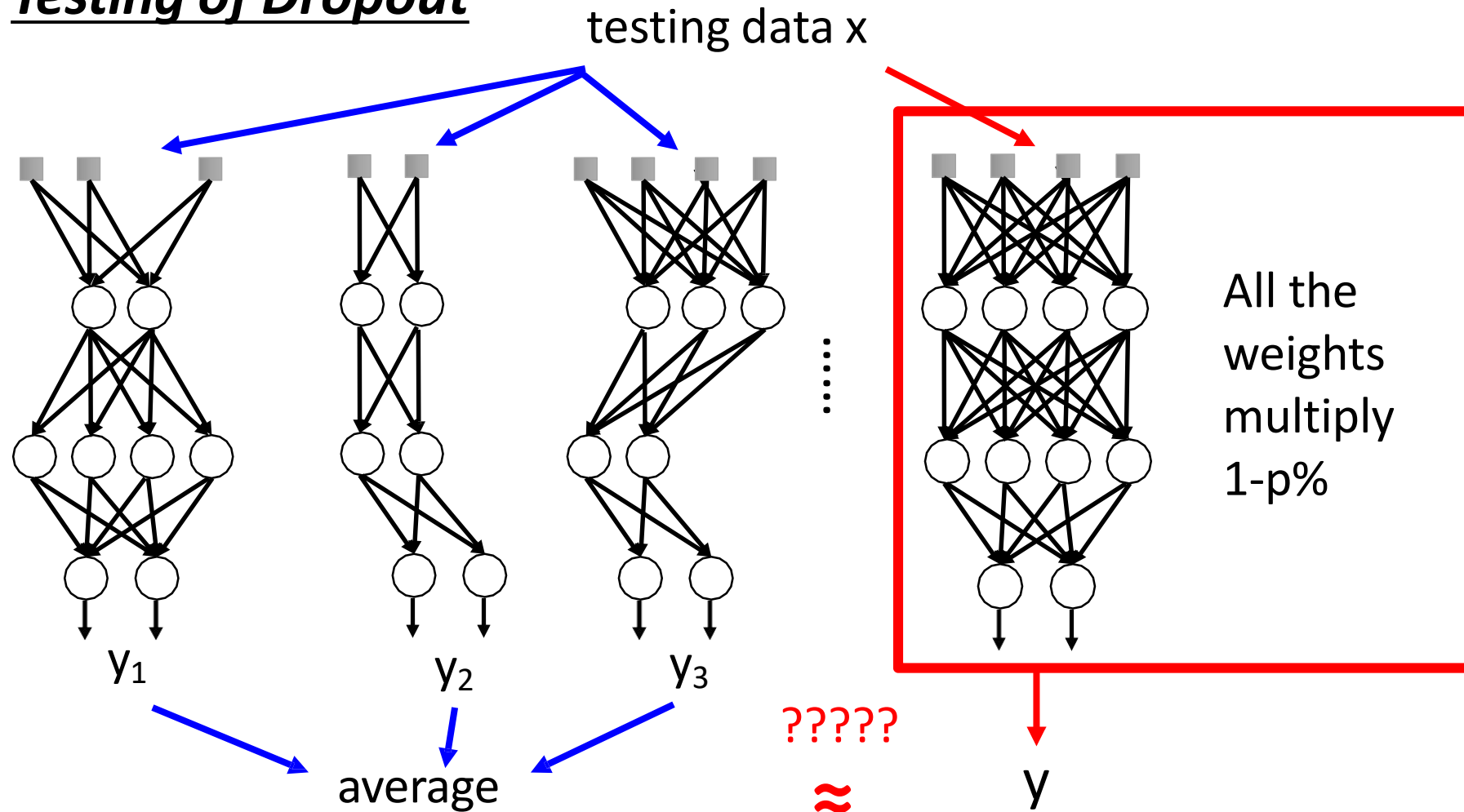
Dropout is a kind of ensemble.



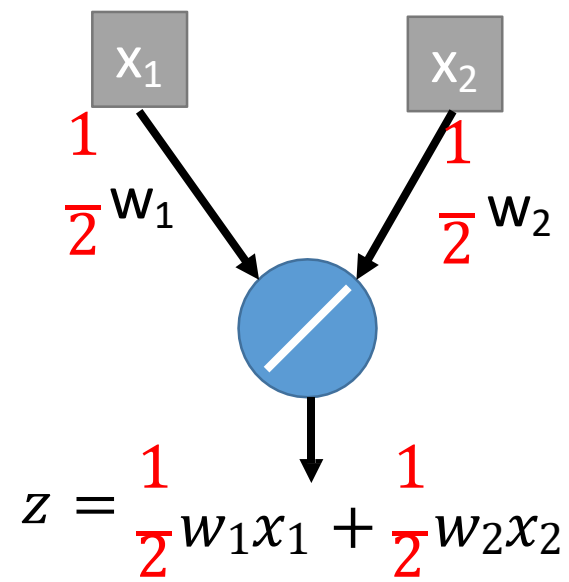
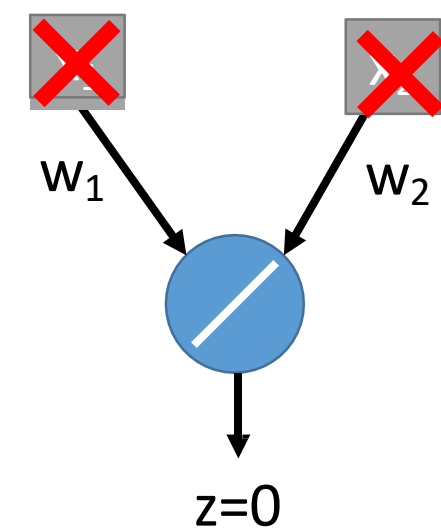
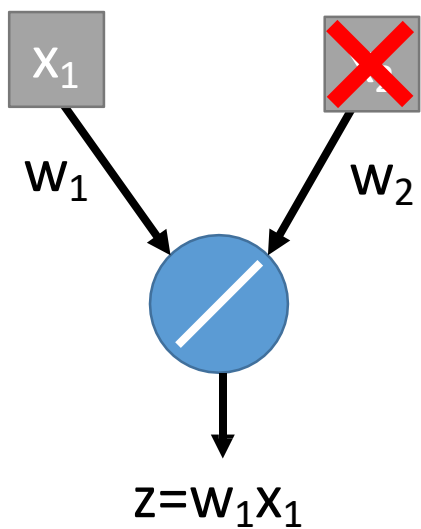
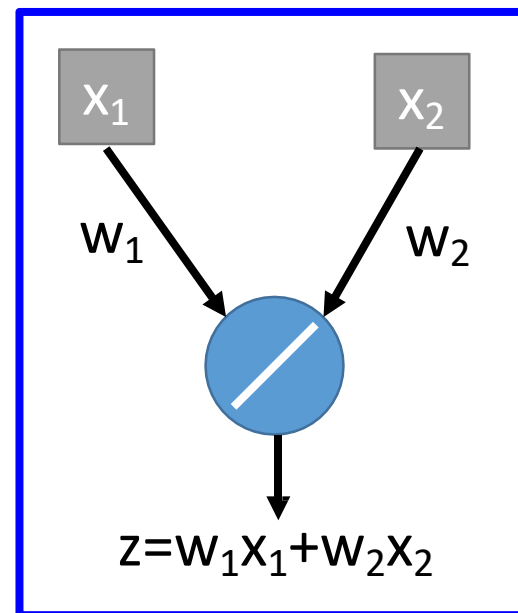
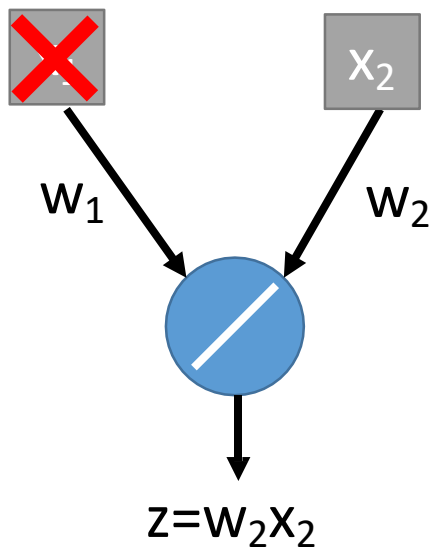
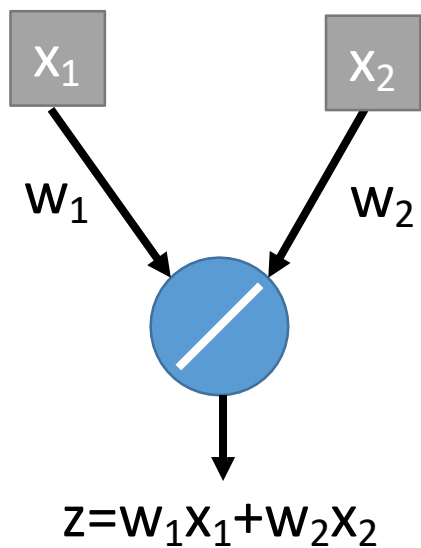
- Using one mini-batch to train one network
- Some parameters in the network are shared

Dropout is a kind of ensemble.

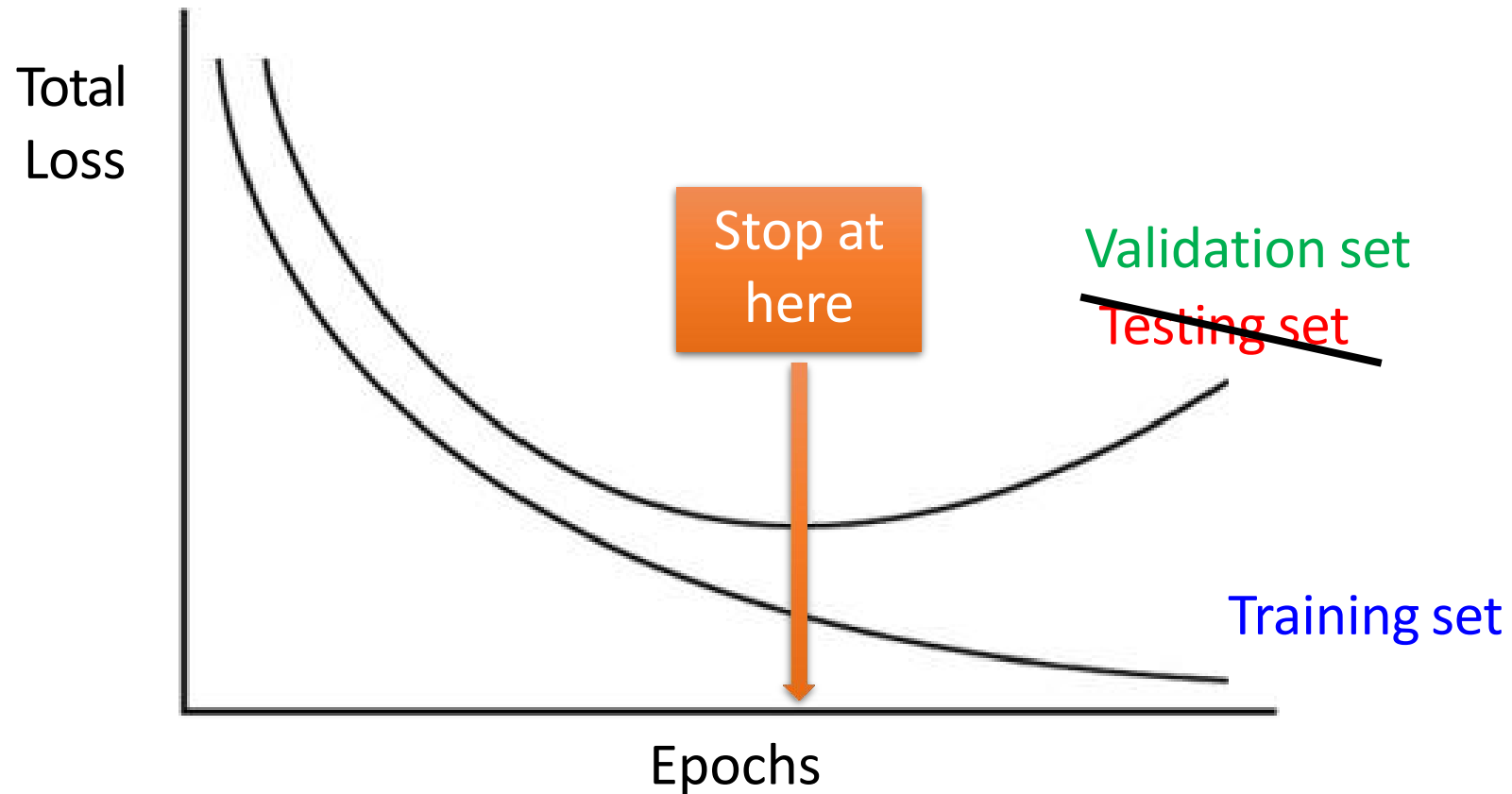
Testing of Dropout



Testing of Dropout



Early Stopping



Keras: <http://keras.io/getting-started/faq/#how-can-i-interrupt-training-when-the-validation-loss-isnt-decreasing-anymore>

Why Deep?

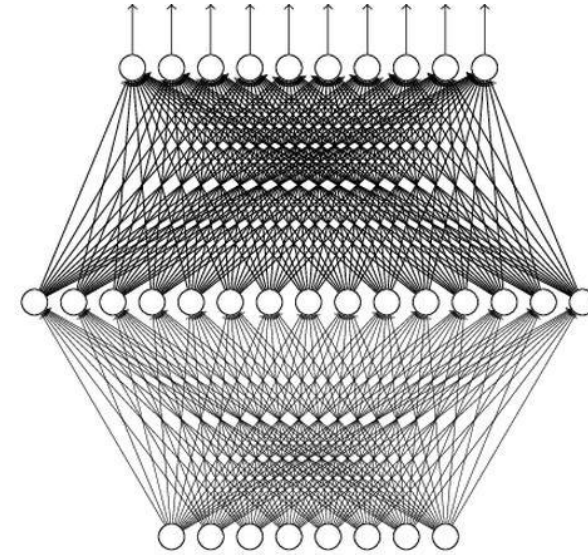
Universality Theorem

Any continuous function f

$$f : R^N \rightarrow R^M$$

Can be realized by a network
with one hidden layer

(given **enough** hidden
neurons)



Reference for the reason:
<http://neuralnetworksanddeeplearning.com/chap4.html>

Why “Deep” neural network not “Fat” neural network?

Deeper is Better?

Layer X Size	Word Error Rate (%)
1 X 2k	24.2
2 X 2k	20.4
3 X 2k	18.4
4 X 2k	17.8
5 X 2k	17.2
7 X 2k	17.1

Not surprised, more parameters, better performance

1 X 3772 22.5

1 X 4634 22.6

1 X 16k 22.1

Seide, Frank, Gang Li, and Dong Yu. "Conversational Speech Transcription Using Context-Dependent Deep Neural Networks." *Interspeech*. 2011.

Fat + Short v.s. Thin + Tall

Layer X Size	Word Error Rate (%)	Layer X Size	Word Error Rate (%)
1 X 2k	24.2		
2 X 2k	20.4		
3 X 2k	18.4		
4 X 2k	17.8		
5 X 2k	17.2	1 X 3772	22.5
7 X 2k	17.1	1 X 4634	22.6
		1 X 16k	22.1

Seide Frank, Gang Li, and Dong Yu. "Conversational Speech Transcription Using Context-Dependent Deep Neural Networks." *Interspeech*. 2011.