

# Improving Deep Neural Networks

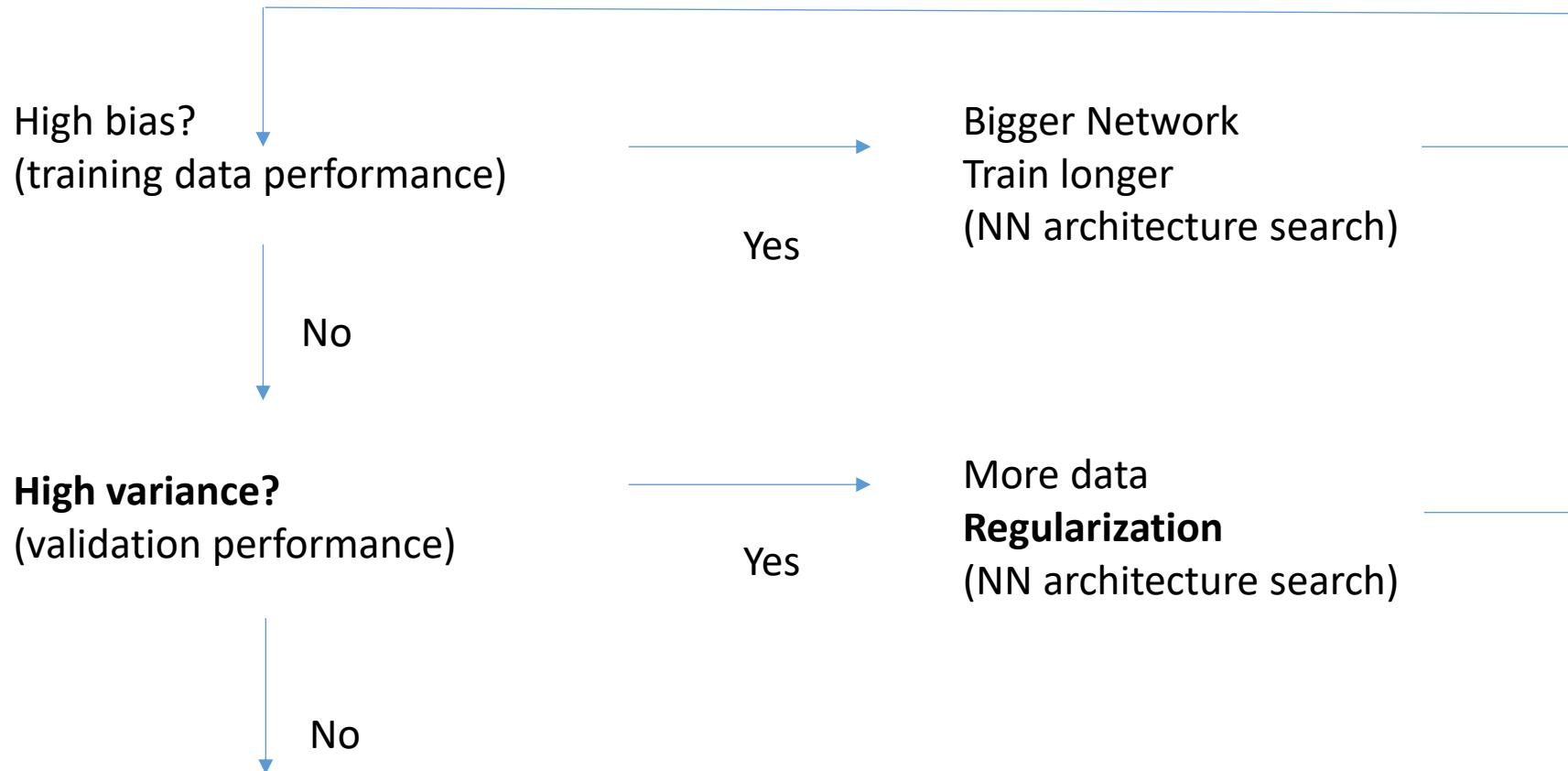
Text and Web Mining (IS6751)

School of Communication and Information

# Overview

- Introduce various approaches for model optimization.
  - Regularization and Early stopping
  - Input normalization and Weights initialization
  - Mini-batch gradient descent
  - Gradient descent with momentum, RMSprop, Adam
  - Learning rate decay
  - Hypermeter tuning process
  - Batch normalization

# Basic recipe for machine learning



# Regularization: Logistic regression

这个slide主要介绍 Logistic Regression 的正则化。

L2 正则化（岭回归）通过惩罚权重的平方和来防止过拟合。

L1 正则化（Lasso 回归）通过惩罚权重的绝对值和来实现特征选择，使模型更稀疏。

正则化参数  $\lambda$  控制正则化强度，影响模型的复杂度和预测能力

- $\min_{\mathbf{w}, b} J(\mathbf{w}, b)$
- Parameters:  $\mathbf{w} \in \mathbb{R}^{n_x}$ ,  $b \in \mathbb{R}$ ,  $\lambda = \text{regularization parameter}$  (positive value)
- $J(\mathbf{w}, b) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \|\mathbf{w}\|_2^2 + \frac{\lambda}{2m} b^2$ 

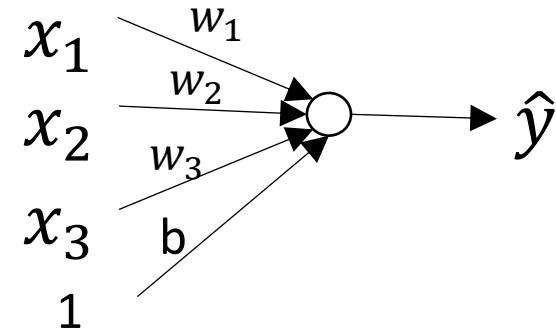
第一项是逻辑回归的损失函数，用于衡量预测值与真实值之间的误差。

omit
- L<sub>2</sub> regularization:  $\|\mathbf{w}\|_2^2 = \sum_{j=1}^{n_x} w_j^2 = \mathbf{w}^\top \mathbf{w}$
- L<sub>1</sub> regularization:  $\|\mathbf{w}\|_1 = \sum_{j=1}^{n_x} |w_j|$  (note:  $\mathbf{w}$  becomes sparse)

$$\begin{aligned} \sum_{j=1}^{n_x} w_j^2 &= \mathbf{w}^\top \mathbf{w} = [w_1 \ w_2 \ w_3] \begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix} \\ &= w_1^2 + w_2^2 + w_3^2 \end{aligned}$$

$$\sum_{j=1}^{n_x} |w_j| = |w_1| + |w_2| + |w_3|$$

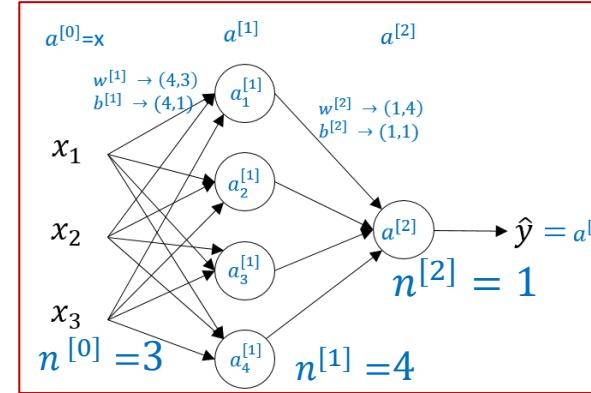
Rationale: When  $\lambda$  approaches infinity, shrink the estimates of  $w_i$  towards zero.



# Regularization: Neural network

$$J(\mathbf{W}^{[1]}, \mathbf{b}^{[1]}, \dots, \mathbf{W}^{[L]}, \mathbf{b}^{[L]}) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \sum_{l=1}^L \|\mathbf{W}^{[l]}\|_2^2$$

$$\|\mathbf{W}^{[l]}\|_2^2 = \sum_{i=1}^{n^{[l]}} \sum_{j=1}^{n^{[l-1]}} (w_{ij}^{[l]})^2 \quad \text{where } \mathbf{W}^{[l]} \text{ is } (n^{[l]}, n^{[l-1]})$$



When  $\mathbf{W}^{[l]} = \begin{bmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \end{bmatrix}$ ,

$$\|\mathbf{W}^{[l]}\|_2^2 = w_{11}^2 + w_{12}^2 + w_{21}^2 + w_{22}^2$$

$$\frac{\partial J}{\partial \mathbf{W}^{[l]}} = d\mathbf{W}^{[l]} = (\text{from backprop}) + \frac{\lambda}{m} \mathbf{W}^{[l]}$$

$$\mathbf{W}^{[l]} = \mathbf{W}^{[l]} - \alpha d\mathbf{W}^{[l]} = \mathbf{W}^{[l]} - \alpha [(\text{from backprop}) + \frac{\lambda}{m} \mathbf{W}^{[l]}]$$

$$= \mathbf{W}^{[l]} - \alpha \frac{\lambda}{m} \mathbf{W}^{[l]} - \alpha (\text{from backprop})$$

$$= (1 - \alpha \frac{\lambda}{m}) \mathbf{W}^{[l]} - \alpha (\text{from backprop})$$

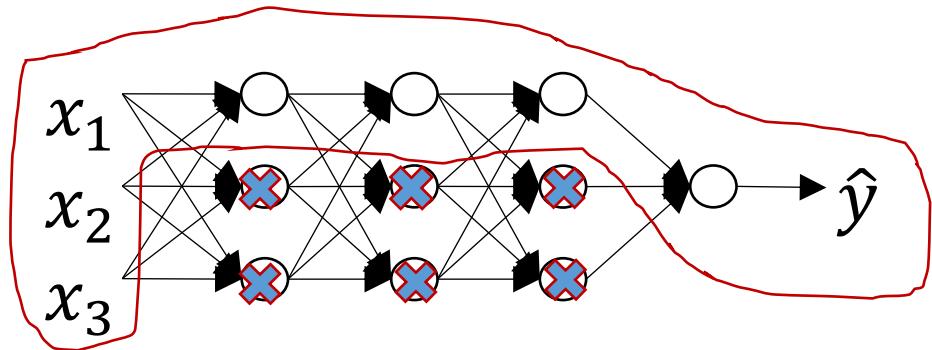
Without Regularization:  
 $\mathbf{W}^{[l]} = \mathbf{W}^{[l]} - \alpha (\text{from backprop})$

深度学习训练时，通常会使用 L2 正则化 (Weight Decay) 来防止神经网络模型过拟合。

在 Adam, SGD, RMSProp 等优化算法中，通常可以设置 weight\_decay 参数来控制 L2 正则化强度。

$0 \leq \alpha \frac{\lambda}{m} < 1$ , so called “Weight decay”

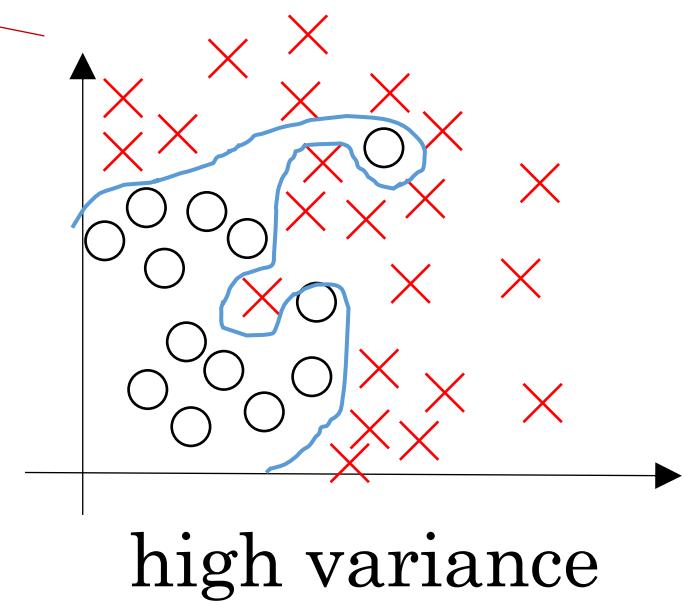
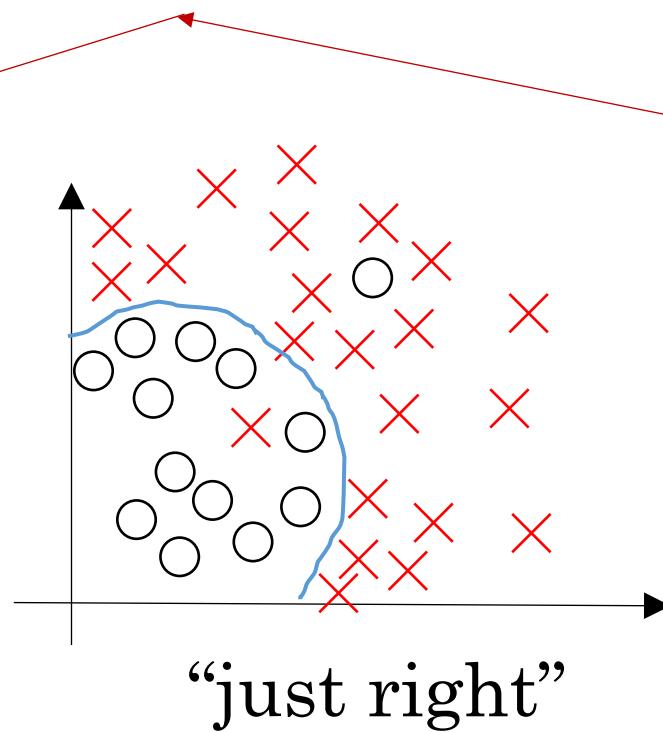
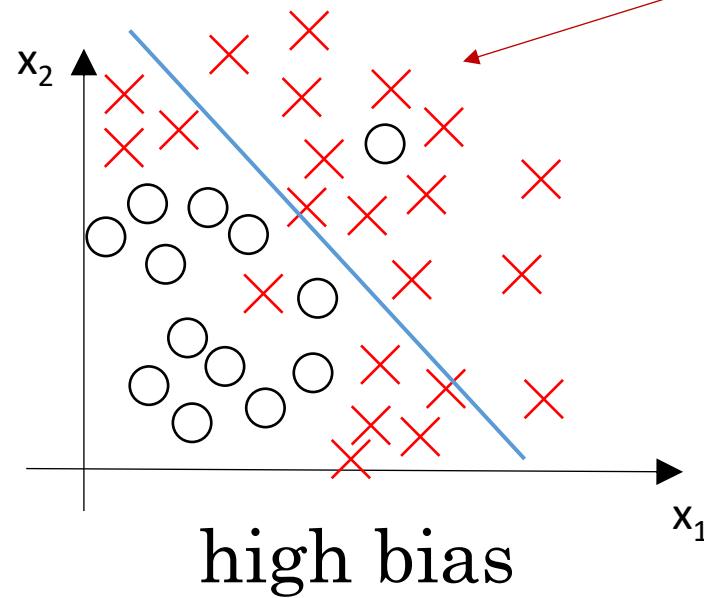
# How does regularization prevent overfitting?



$$J(\mathbf{W}^{[L]}, \mathbf{b}^{[L]}) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \sum_{l=1}^L \|\mathbf{W}^{[l]}\|_2^2$$

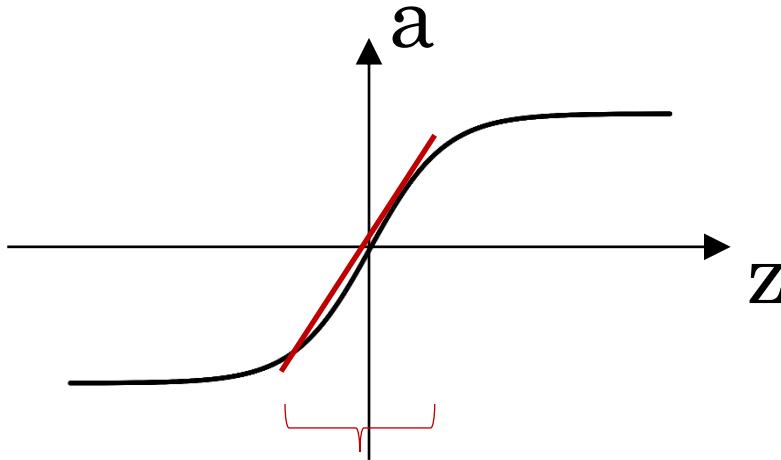
$$\mathbf{W}^{[l]} \approx 0$$

$\lambda$  is a regularization parameter.

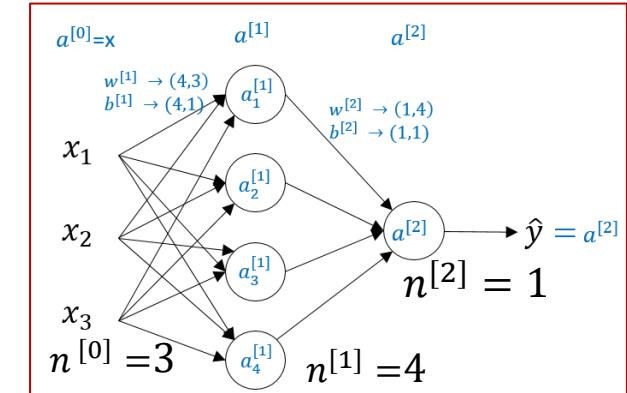


# How does regularization prevent overfitting?

tanh:  $g(z) = \tanh(z)$



$$J(\mathbf{W}^{[L]}, \mathbf{b}^{[L]}) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \sum_{l=1}^L \|\mathbf{W}^{[l]}\|_2^2$$



When  $\lambda$  increases,  $\mathbf{W}^{[l]}$  decreases. So,  $z^{[l]}$  decreases.  $z^{[l]} = \mathbf{W}^{[l]} \mathbf{a}^{[l-1]} + \mathbf{b}^{[l]}$

Every layer  $\approx$  linear

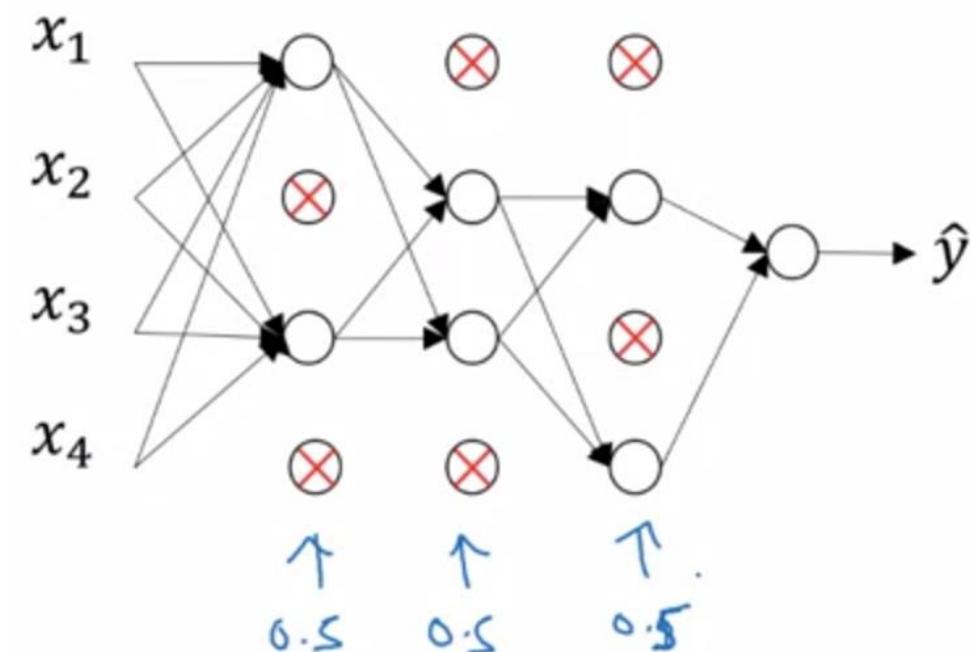
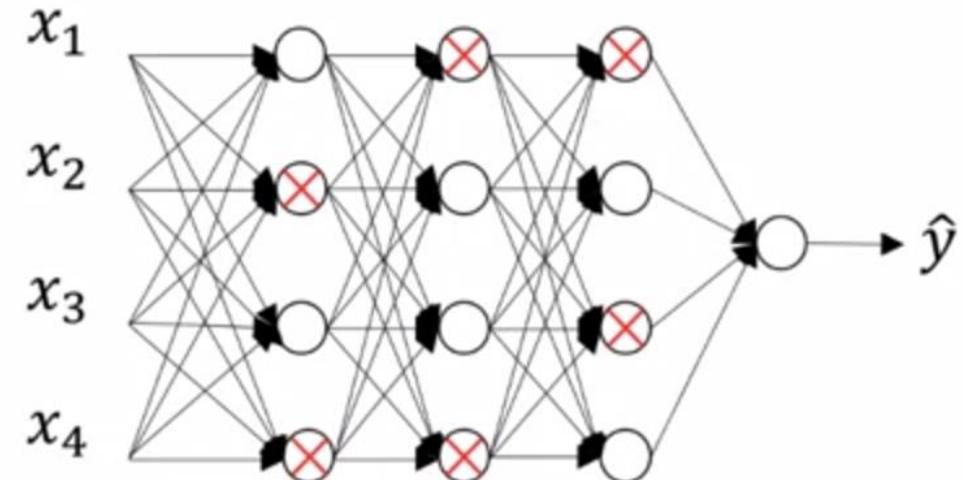
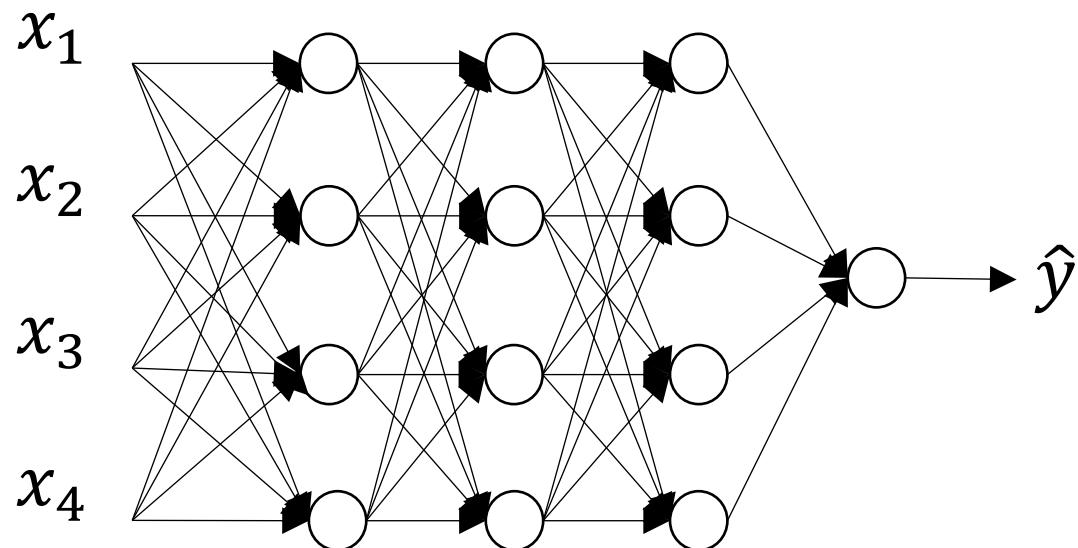
# Dropout regularization

Dropout 是一种正则化方法，用于防止神经网络过拟合。

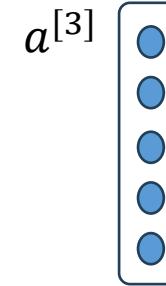
在训练时随机丢弃部分神经元，使得网络不会过分依赖某些特定神经元，提高模型的泛化能力。

在测试时不进行 Dropout，但要对神经元的激活值进行缩放，以匹配训练时的期望值。

Dropout rate（丢弃率）通常设为 0.2 - 0.5，根据不同的网络结构进行调整。



# Implementing dropout



$$\text{E.g., } d3 = \begin{bmatrix} 0.9 \\ 0.6 \\ 0.4 \\ 0.2 \\ 0.3 \end{bmatrix} < 0.8 = \begin{bmatrix} 0 \\ 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}$$

- Illustrate with layer  $l = 3$  and **Keep-prob** = 0.8

```
d3 = torch.rand(a3.shape[0],a3.shape[1]) < keep-prob
```

```
a3 = torch.mul(a3,d3) # a3 *= d3
```

```
a3 /= keep-prob
```

- When 5 units, 1 unit shuts off.

$$z^{[4]} = W^{[4]}a^{[3]} + b^{[4]}$$

Reduced by 20%  
So apply “/= 0.8”

E.g.,  $\begin{bmatrix} 0.5 \\ 0.6 \\ 0.9 \\ 0.3 \\ 0.5 \end{bmatrix} * \begin{bmatrix} 1 \\ 1 \\ 0 \\ 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 0.0 \\ 0.6 \\ 0.0 \\ 0.3 \\ 0.5 \end{bmatrix}$

$$\begin{bmatrix} 0.5 \\ 0.6 \\ 0.0 \\ 0.3 \\ 0.5 \end{bmatrix} / 0.8 = \begin{bmatrix} 0.625 \\ 0.750 \\ 0.000 \\ 0.375 \\ 0.625 \end{bmatrix}$$

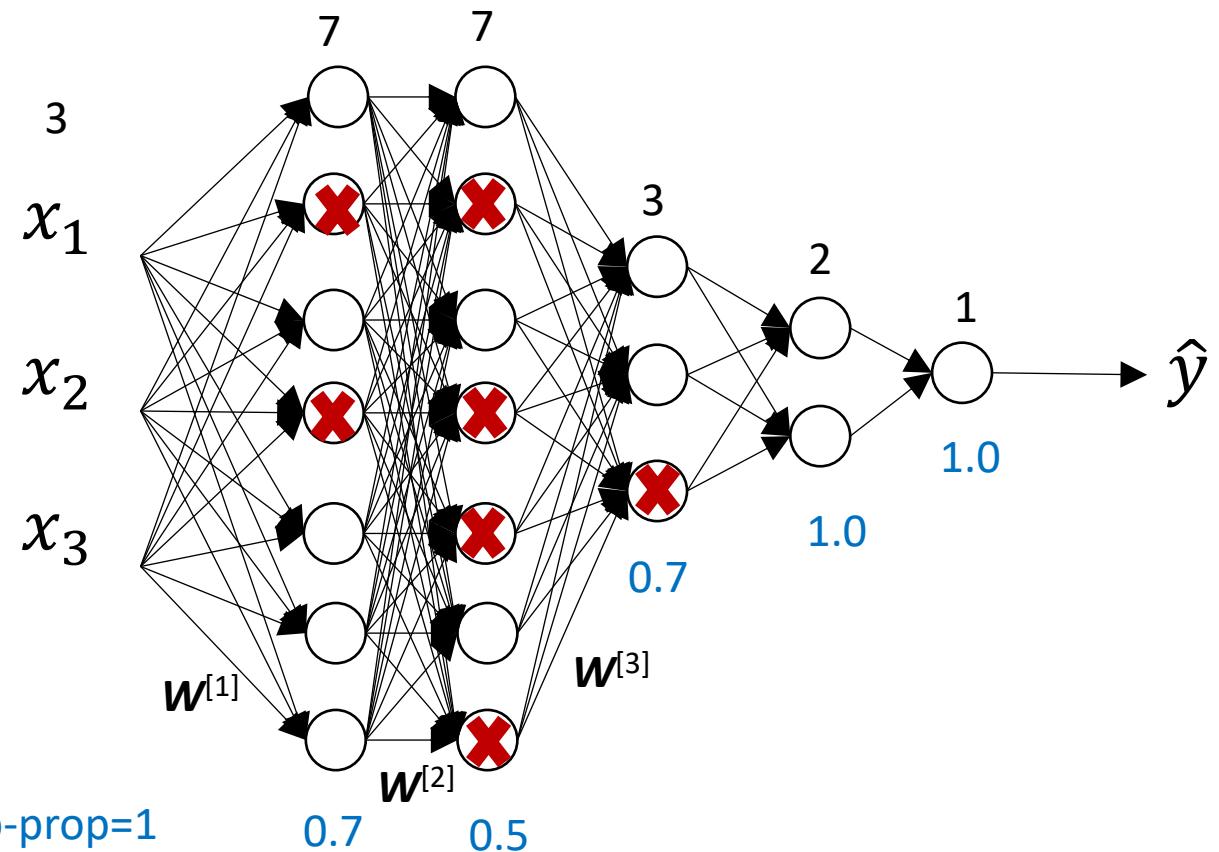
No drop out during test time

E.g.,  $d3 = \begin{bmatrix} 0.7 \\ 0.5 \\ 0.9 \\ 0.5 \\ 0.6 \end{bmatrix} < 0.8 = \begin{bmatrix} 1 \\ 1 \\ 0 \\ 1 \\ 1 \end{bmatrix}$

# Why does drop-out work?



- Intuition: Can't rely on any one feature, so must spread out weights.
  - shrink weights like L<sub>2</sub> regularization



Keep-prop=1

0.7

0.5

[since  $W^{[2]}$  has more parameters, a small keep-prop value (0.5) is used]

# Early stopping

## 图表解析

横轴 (# iterations) : 表示训练迭代次数 (或模型复杂度)，从  $w \approx 0$  (初始权重接近 0) 到较大的  $W$  (训练后权重)。

纵轴 (Error or Loss) : 表示损失函数值，包括：

Training error (训练误差) : 随着训练的进行，训练误差不断下降。

Validation set error (验证集误差) : 起初下降，但在某个点后开始上升。

Early Stopping 关键点

在训练早期，训练误差和验证误差都在下降，表示模型正在学习数据模式。

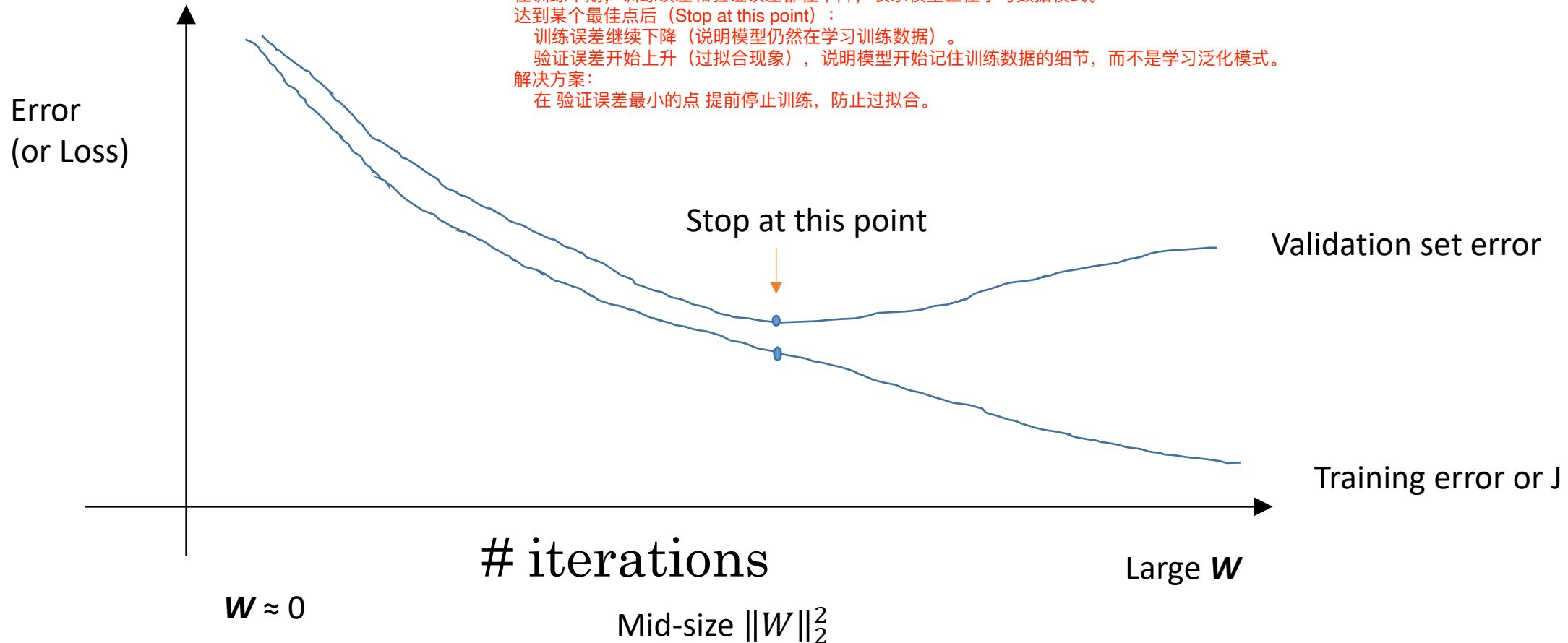
达到某个最佳点后 (Stop at this point) :

训练误差继续下降 (说明模型仍然在学习训练数据)。

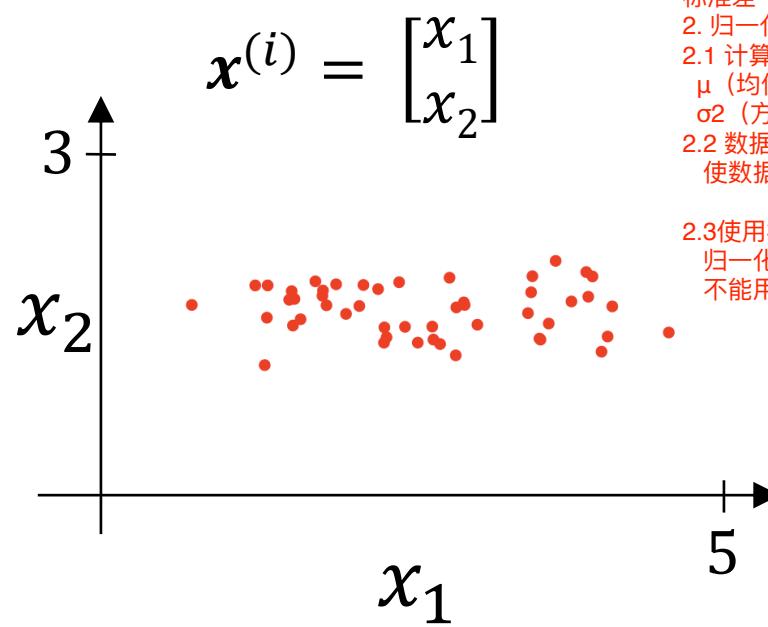
验证误差开始上升 (过拟合现象)，说明模型开始记住训练数据的细节，而不是学习泛化模式。

解决方案：

在验证误差最小的点 提前停止训练，防止过拟合。



# Normalizing inputs: Normalizing training sets



均值 0: 所有特征值围绕 0 分布

标准差 1: 数据的尺度被统一, 避免某些特征值过大而主导模型训练

2. 归一化的步骤

2.1 计算均值和方差:

$\mu$  (均值) : 计算所有样本的平均值

$\sigma^2$  (方差) : 计算数据的方差, 衡量数据的分布范围

2.2 数据标准化 (Normalization)

使数据均值变成 0, 标准差变成 1

2.3 使用相同的均值和方差归一化测试集

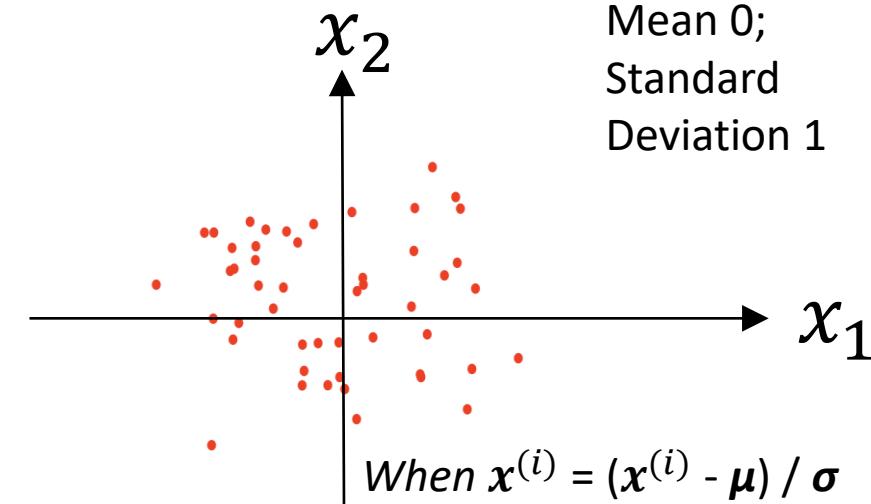
归一化不仅应用于训练集, 测试集也必须使用相同的  $\mu$  和  $\sigma$  进行归一化

不能用测试集的均值和标准差, 否则会导致数据分布不匹配。

$$\begin{aligned} x^{(i)} - \mu &= \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} - \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix} \end{aligned}$$

Mean 0

When  $x^{(i)} = x^{(i)} - \mu$



E.g., use  $m$  number of  $x^{(1)}$ ,  $x^{(2)}, \dots, x^{(m)}$  to calculate  $\mu_1$  and  $\mu_2$

$$\begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix} \leftarrow \mu = \frac{1}{m} \sum_{i=1}^m x^{(i)}$$

$$\begin{bmatrix} \sigma_1^2 \\ \sigma_2^2 \end{bmatrix} \leftarrow \sigma^2 = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu)^2$$

Calculate mean and variance:

Normalize inputs:

$$x^{(i)} = (x^{(i)} - \mu) / \sigma$$

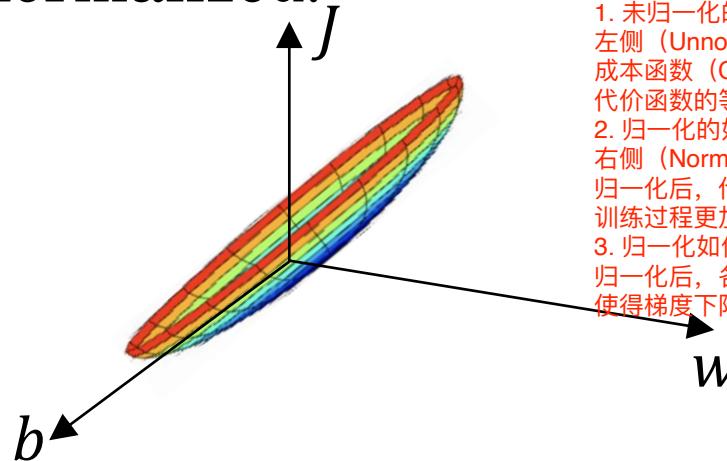
$$\begin{aligned} x^{(i)} &= (x^{(i)} - \mu) / \sigma \\ &= \left( \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} - \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix} \right) / \begin{bmatrix} \sigma_1 \\ \sigma_2 \end{bmatrix} \end{aligned}$$

- Use the same  $\mu$  and  $\sigma$  to normalize test set:  $x^{(i)} = (x^{(i)} - \mu) / \sigma$

# Why normalize inputs?

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

Unnormalized:



## 1. 未归一化的数据如何影响优化?

左侧 (Unnormalized) :

成本函数 (Cost function) 变成拉长的椭圆形 (elongated shape)，梯度下降时可能会在狭窄区域内震荡，使收敛速度变慢。代价函数的等高线呈现出长条形，导致梯度下降路径很曲折，优化器需要花费更多时间找到最优解。

## 2. 归一化的好处

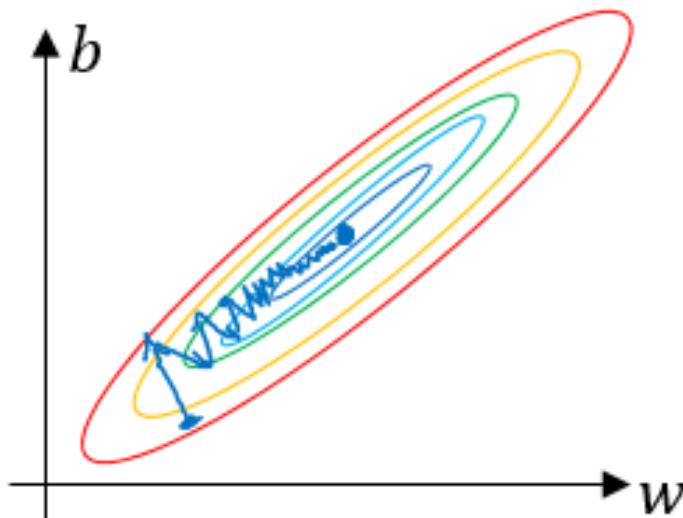
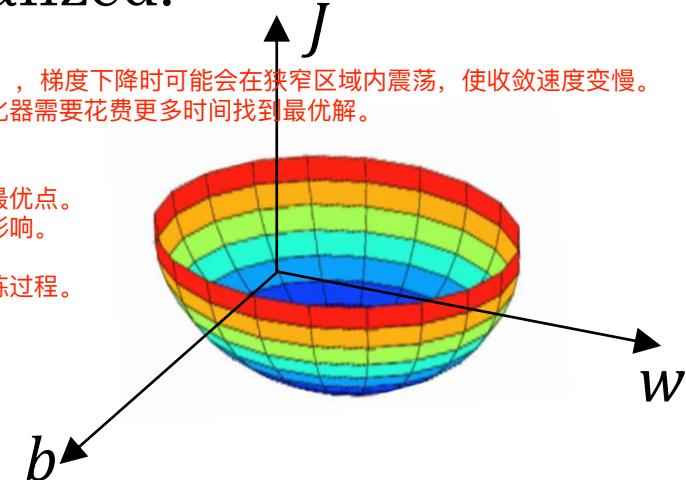
右侧 (Normalized) :

归一化后，代价函数变成对称的圆形，梯度下降能够快速收敛到最优点。训练过程更加平稳，参数更新不会受到某些特征值过大或过小的影响。

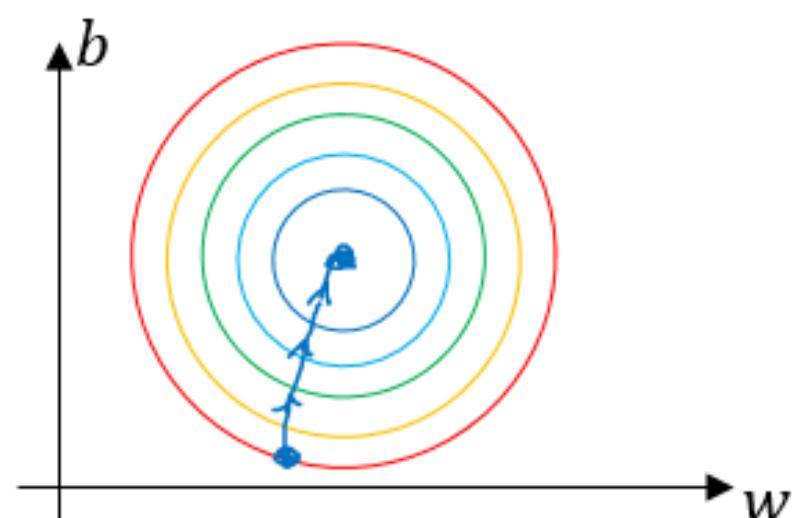
## 3. 归一化如何帮助神经网络?

归一化后，各个特征的值处于相同的尺度，防止某些特征主导训练过程。使得梯度下降的路径更加直接，提高模型收敛速度。

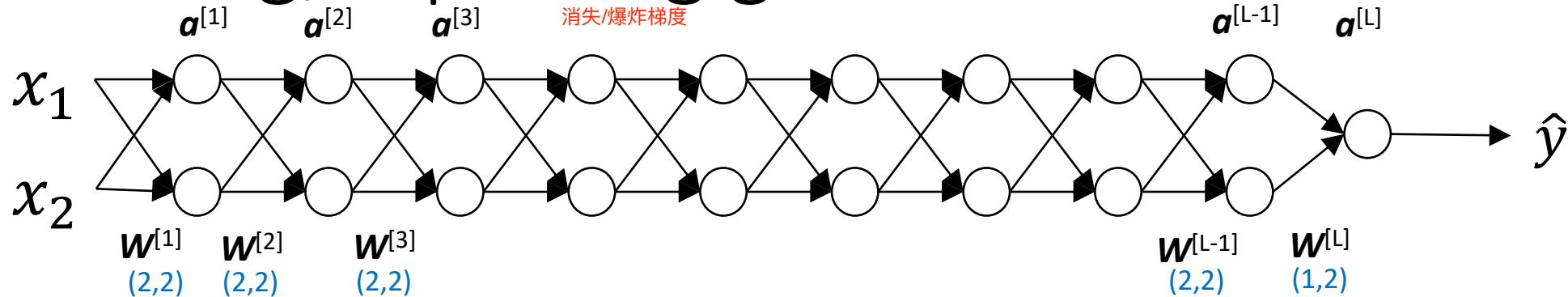
Normalized:



- Cost function becomes an **elongated shape** (long and thin)
- Take longer to learn weights or reach to an optimal point.



# Vanishing/exploding gradients



When  $g(z) = z$ ,  $b^{[l]} = 0$

$$\hat{y} = \mathbf{W}^{[L]} \mathbf{W}^{[L-1]} \mathbf{W}^{[L-2]} \dots \mathbf{W}^{[3]} \mathbf{W}^{[2]} \mathbf{W}^{[1]} \mathbf{x}$$

1. 什么是梯度消失和梯度爆炸?

梯度消失 (Vanishing Gradient) :

发生在深层网络的前层权重很小 ( $<1$ ) 时, 使得梯度在反向传播过程中指数级缩小, 导致前层权重更新几乎为零, 网络无法学习。

梯度爆炸 (Exploding Gradient) :

发生在前层权重很大 ( $>1$ ) 时, 使得梯度在反向传播过程中指数级增长, 导致梯度更新过大, 模型发散 (不收敛)。

$$z^{[1]} = \mathbf{W}^{[1]} \mathbf{x}$$

$$\mathbf{a}^{[1]} = g(z^{[1]}) = z^{[1]}$$

$$z^{[2]} = \mathbf{W}^{[2]} \mathbf{a}^{[1]}$$

$$\mathbf{a}^{[2]} = g(z^{[2]}) = z^{[2]}$$

ReLU / Leaky ReLU / GELU 代替 Sigmoid/Tanh 作为激活函数。

使用 He/Xavier 初始值进行权重初始化, 防止初始梯度过大或过小。

Batch Normalization 归一化数据, 使梯度更稳定。

梯度裁剪 (Gradient Clipping) 控制梯度大小, 避免梯度爆炸。

ResNet 的残差连接 让梯度可以更容易传播, 防止梯度消失。

When  $\mathbf{W}^{[l]} = \begin{vmatrix} 1.5 & 0 \\ 0 & 1.5 \end{vmatrix}$  (i.e.,  $\mathbf{W}^{[l]} > I$ ),  $\hat{y} = \mathbf{W}^{[L]} \begin{vmatrix} 1.5 & 0 \\ 0 & 1.5 \end{vmatrix}^{L-1} \mathbf{x} = \mathbf{W}^{[L]} 1.5^{L-1} \mathbf{x}$  ( $\hat{y}$  increases exponentially)

When  $\mathbf{W}^{[l]} = \begin{vmatrix} 0.5 & 0 \\ 0 & 0.5 \end{vmatrix}$  (i.e.,  $\mathbf{W}^{[l]} < I$ ),  $\hat{y} = \mathbf{W}^{[L]} 0.5^{L-1} \mathbf{x}$  ( $\hat{y}$  decreases exponentially)

Computing gradients:

$$dz^{[2]} = a^{[2]} - y$$

$$d\mathbf{W}^{[2]} = dz^{[2]} a^{[1]T}$$

$$db^{[2]} = dz^{[2]}$$

$$dz^{[1]} = \boxed{\mathbf{W}^{[2]T}} dz^{[2]} * g^{[1]'}(\mathbf{z}^{[1]})$$

$$d\mathbf{W}^{[1]} = dz^{[1]} \mathbf{x}^T$$

$$db^{[1]} = dz^{[1]}$$

Through back propagation,  $\mathbf{W}$  values are multiplied. So, derivatives can increase or decrease exponentially.

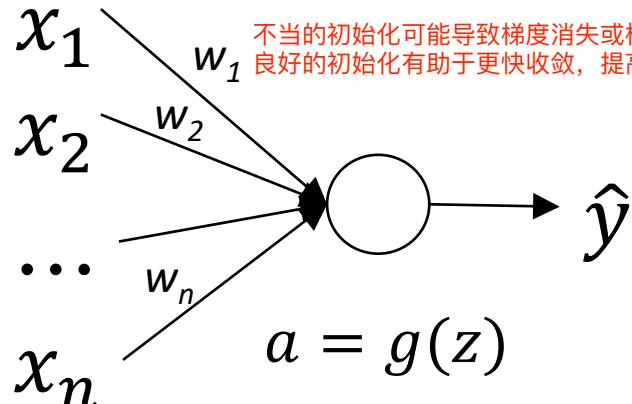
$$\begin{vmatrix} 1.5 & 0 \\ 0 & 1.5 \end{vmatrix} \begin{vmatrix} 1.5 & 0 \\ 0 & 1.5 \end{vmatrix} = \begin{vmatrix} 1.5^2 & 0 \\ 0 & 1.5^2 \end{vmatrix}$$

$$\begin{vmatrix} 1.5 & 0 \\ 0 & 1.5^{l-1} \end{vmatrix} \begin{vmatrix} x_1 \\ x_2 \end{vmatrix} = 1.5^{l-1} \begin{vmatrix} x_1 \\ x_2 \end{vmatrix} = 1.5^{l-1} \mathbf{x}$$

# Weights Initialization

$$n^{[l-1]} = n \quad n^{[l]} = 1$$

在神经网络中，权重  $w$  的初始化对训练效果有很大的影响：



Single  
Neuron  
Example

$$z = w_1x_1 + w_2x_2 + \dots + w_nx_n + b \text{ (ignore } b\text{)}$$

When  $n^{[l-1]}$  is large,  $w_i$  becomes smaller.

$$\text{Variance}(w_i) = \frac{2}{n^{[l-1]}}$$

shape of  $W$ :  $(n^{[l]}, n^{[l-1]})$

$$W^{[l]} = \text{np.random.randn(shape of } W\text{)} * \text{np.sqrt} \left( \frac{2}{n^{[l-1]}} \right) \text{ (with Relu)}$$

$$W^{[l]} = \text{np.random.randn(shape of } W\text{)} * \text{np.sqrt} \left( \frac{1}{n^{[l-1]}} \right) \text{ (with tanh, called Xavier initialization)}$$

$$W^{[l]} = \text{np.random.randn(shape of } W\text{)} * \text{np.sqrt} \left( \frac{2}{n^{[l-1]} + n^{[l]}} \right) \text{ (another approach)}$$

`w1=np.random.randn(3,5) * np.sqrt (2/5)`  
w1

When  $n^{[l-1]} = 5$ ,  $n^{[l]} = 3$

array([[-0.22929147, 0.6993235 , -0.10430622, -0.2149747 , 1.06357577],  
[ 0.33875831, 0.62626452, -0.94417563, -0.72660319, -0.79194391],  
[ 0.15885362, 0.46539225, -0.76860632, 0.08878817, 0.2415148 ]])

`np.var(w1)`

$$2/5 = 0.4$$

0.3454423104406922

$$\text{np.sqrt}(2/5) \approx 0.63$$

`w2=np.random.randn(3,100) * np.sqrt (2/100)`  
np.var(w2)

When  $n^{[l-1]} = 100$ ,  $n^{[l]} = 3$

0.018709112797490716

$$2/100 = 0.02$$

$$\text{np.sqrt}(2/100) \approx 0.14$$

✓ 好的权重初始化方法可以：

防止梯度消失或爆炸，提高训练稳定性。

加速网络收敛，减少训练时间。

适配不同激活函数，确保信息能够有效传递。

在深度学习中，合理选择 权重初始化方式 是训练成功的重要一步！

# Mini-batch gradient descent: Batch vs. mini-batch gradient descent

Vectorization allows you to efficiently compute on  $m$  examples.

What if  $m = 5,000,000$  ?

-> 5,000 mini-batches of 1,000 each

-> Mini-batch  $t$ :  $\mathbf{X}^{\{t\}}, \mathbf{Y}^{\{t\}}$

📌 Mini-batch Gradient Descent 是深度学习中最常用的优化方法：

数据集分成多个 mini-batch，每个 batch 计算一次梯度更新，比 Batch 计算更快，比 SGD 更稳定。  
利用向量化计算提升 GPU 计算效率。  
适用于大规模数据集训练（如 500 万样本），可以有效加速收敛。

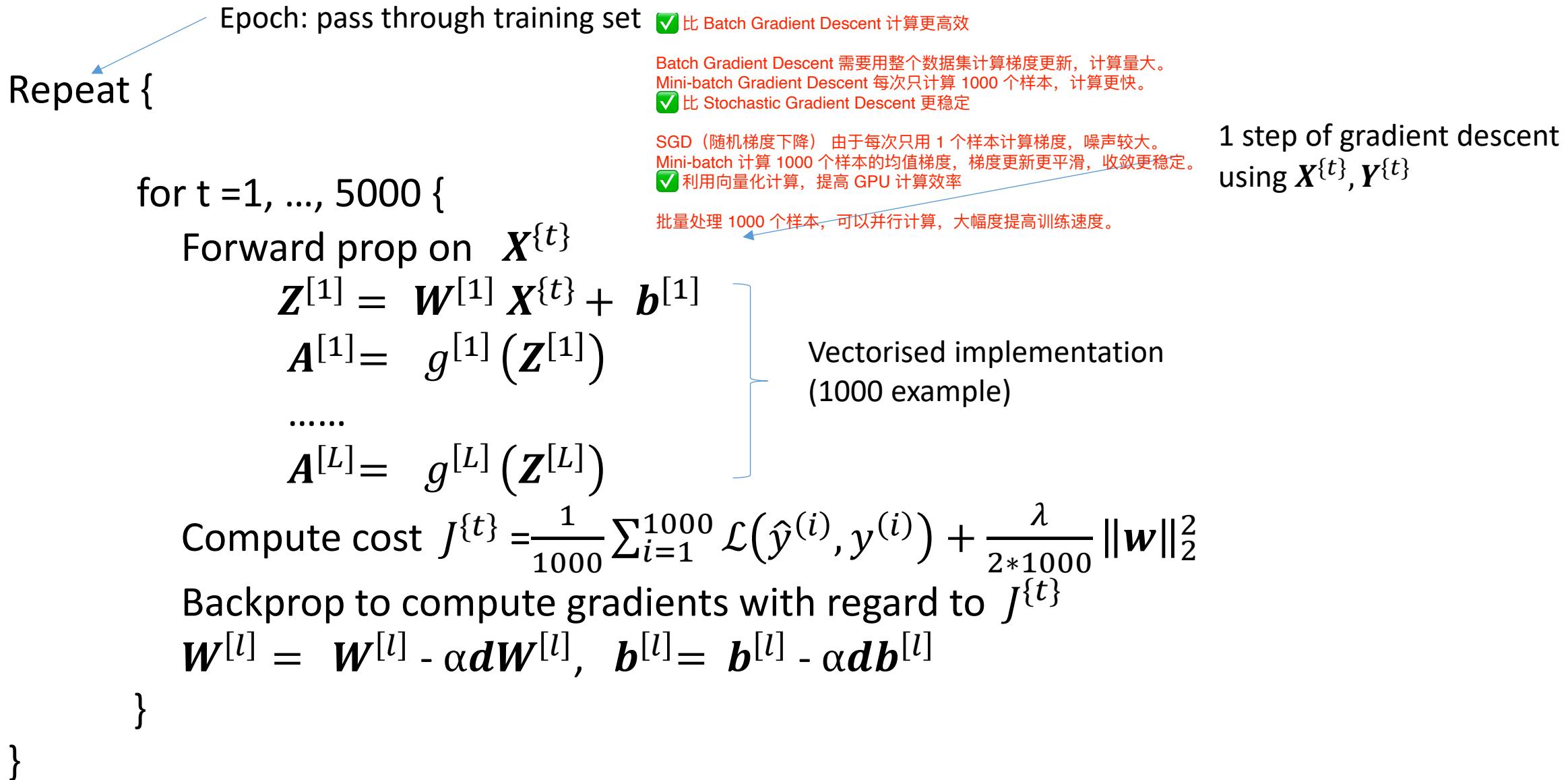
$$\mathbf{X} = [\underbrace{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(1000)}}_{(\mathbf{n}_x, m)} \mid \underbrace{\mathbf{x}^{(1001)}, \dots, \mathbf{x}^{(2000)}}_{\mathbf{X}^{\{2\}} (\mathbf{n}_x, 1000)} \mid \dots \mid \dots \mathbf{x}^{(m)}]$$

$\mathbf{X}^{\{1\}} (\mathbf{n}_x, 1000)$        $\mathbf{X}^{\{2\}} (\mathbf{n}_x, 1000)$        $\dots$        $\mathbf{X}^{\{5000\}} (\mathbf{n}_x, 1000)$

$$\mathbf{Y} = [\underbrace{y^{(1)}, y^{(2)}, \dots, y^{(1000)}}_{(1, m)} \mid \underbrace{y^{(1001)}, \dots, y^{(2000)}}_{\mathbf{Y}^{\{2\}} (1, 1000)} \mid \dots \mid \dots y^{(m)}]$$

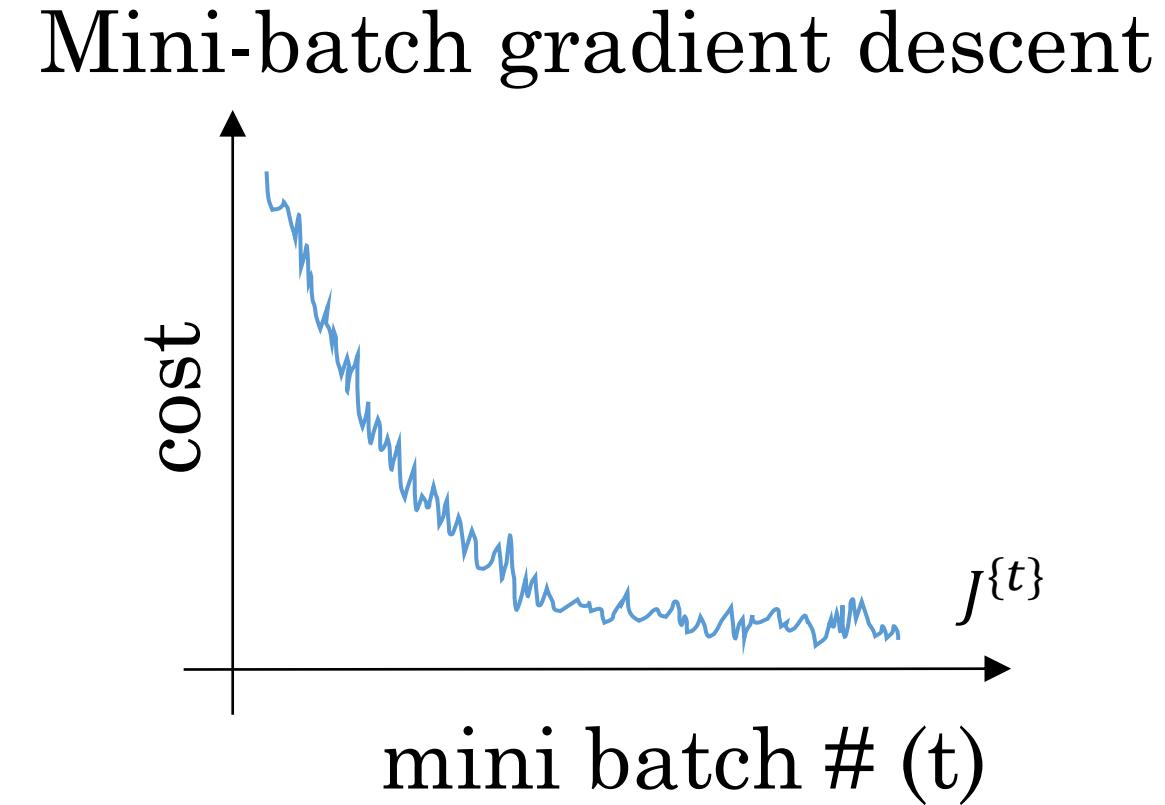
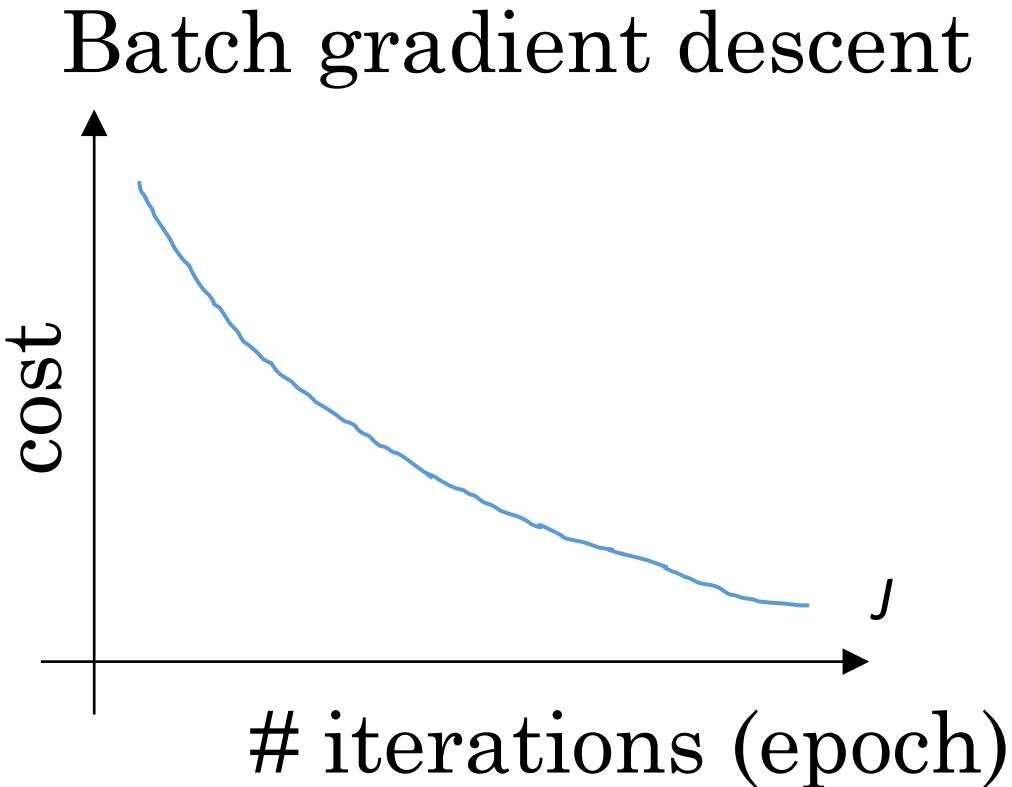
$\mathbf{Y}^{\{1\}} (1, 1000)$        $\mathbf{Y}^{\{2\}} (1, 1000)$        $\dots$        $\mathbf{Y}^{\{5000\}} (1, 1000)$

# Mini-batch gradient descent



# Understanding mini-batch gradient descent:

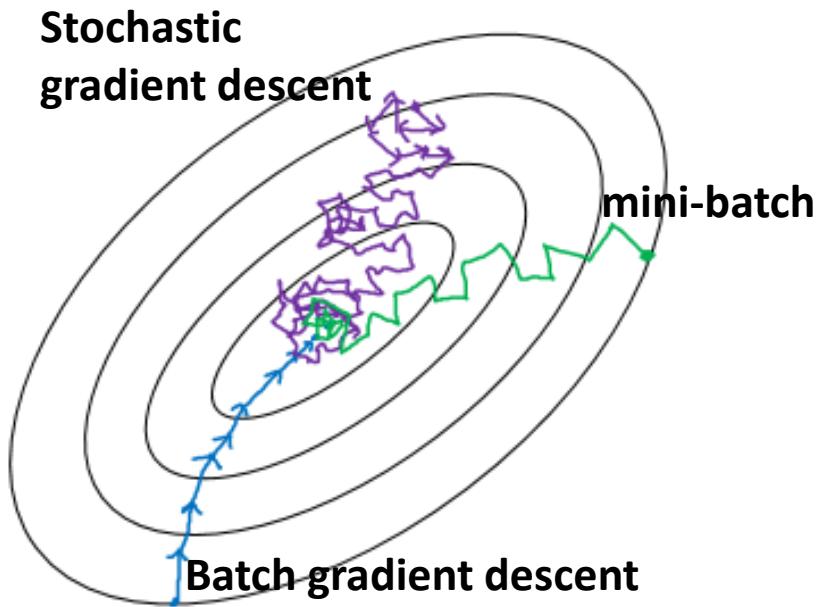
## Training with mini batch gradient descent



*Plot  $J^{(t)}$  computed using  $\mathbf{X}^{(t)}, \mathbf{Y}^{(t)}$*

# Choosing your mini-batch size

- If mini-batch size =  $m$ , **batch** gradient descent  $(\mathbf{X}^{\{t\}}, \mathbf{Y}^{\{t\}}) = (\mathbf{X}, \mathbf{Y})$
- If mini-batch size = 1, **stochastic** gradient descent  $(\mathbf{X}^{\{t\}}, \mathbf{Y}^{\{t\}}) = (\mathbf{x}^{(i)}, \mathbf{y}^{(i)})$
- In practice, somewhere in-between 1 and  $m$



Stochastic  
gradient descent

↓  
Lose speedup  
from  
vectorization

In-between (mini-batch  
size not too big/small)

↓  
Fastest learning

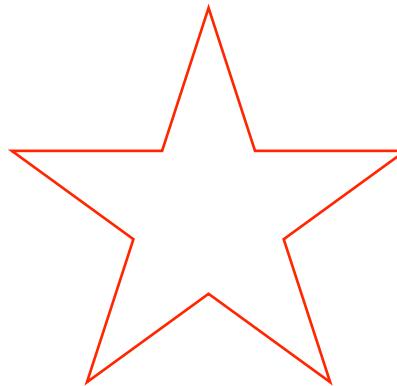
- vectorization (e.g., 1000)
- Make process without processing entire training set

Batch gradient  
descent

↓  
Too long per  
iteration

# Choosing your mini-batch size

- If small training set, use **batch** gradient descent ( $m \leq 2000$ )
- Typical mini-batch sizes:
  - 64 ( $2^6$ ), 128 ( $2^7$ ), 256 ( $2^8$ ), and 512 ( $2^9$ )
- Make sure mini-batch  $(\mathbf{X}^{\{t\}}, \mathbf{Y}^{\{t\}})$  fit in CPU/GPU memory.



# Exponentially weighted (running) averages: Temperature in a city

$$\theta_1 = 40^{\circ}\text{F}$$

$$\theta_2 = 49^{\circ}\text{F}$$

$$\theta_3 = 45^{\circ}\text{F}$$

:

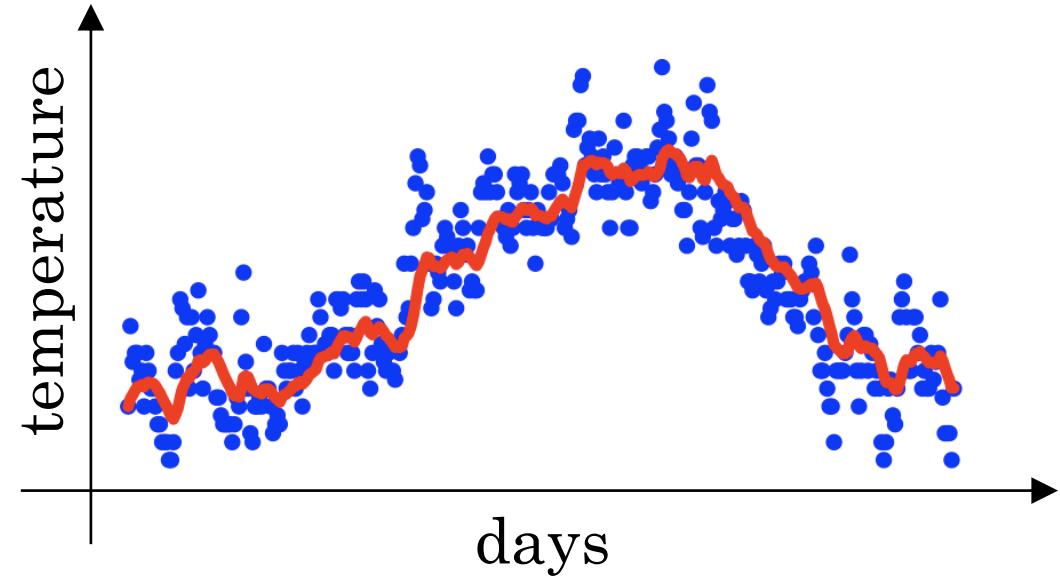
:

$$\theta_{180} = 60^{\circ}\text{F}$$

$$\theta_{181} = 56^{\circ}\text{F}$$

:

:



$$V_0 = 0$$

$$V_1 = 0.9 V_0 + 0.1 \theta_1$$

$$V_2 = 0.9 V_1 + 0.1 \theta_2$$

...

$$V_t = 0.9 V_{t-1} + 0.1 \theta_t$$

# Exponentially weighted averages

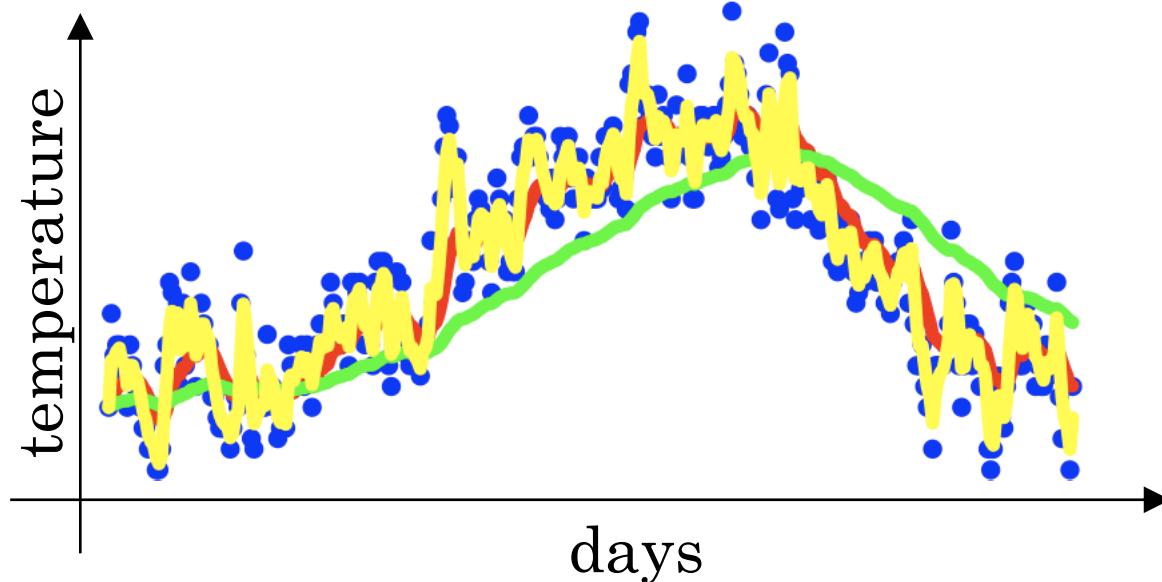
$$V_t = \beta V_{t-1} + (1 - \beta) \theta_t$$

$V_t$  is approximately averaged over  
 $\approx 1/(1-\beta)$  days' temperature

When  $\beta = 0.9$ ,  $\approx 10$  days temperature (red)

$\beta = 0.98$ ,  $\approx 50$  days temperature (green)

$\beta = 0.5$ ,  $\approx 2$  days temperature (yellow)



# Bias correction in exponentially weighted average

when  $\beta = 0.98$

$$v_t = \beta v_{t-1} + (1 - \beta) \theta_t$$

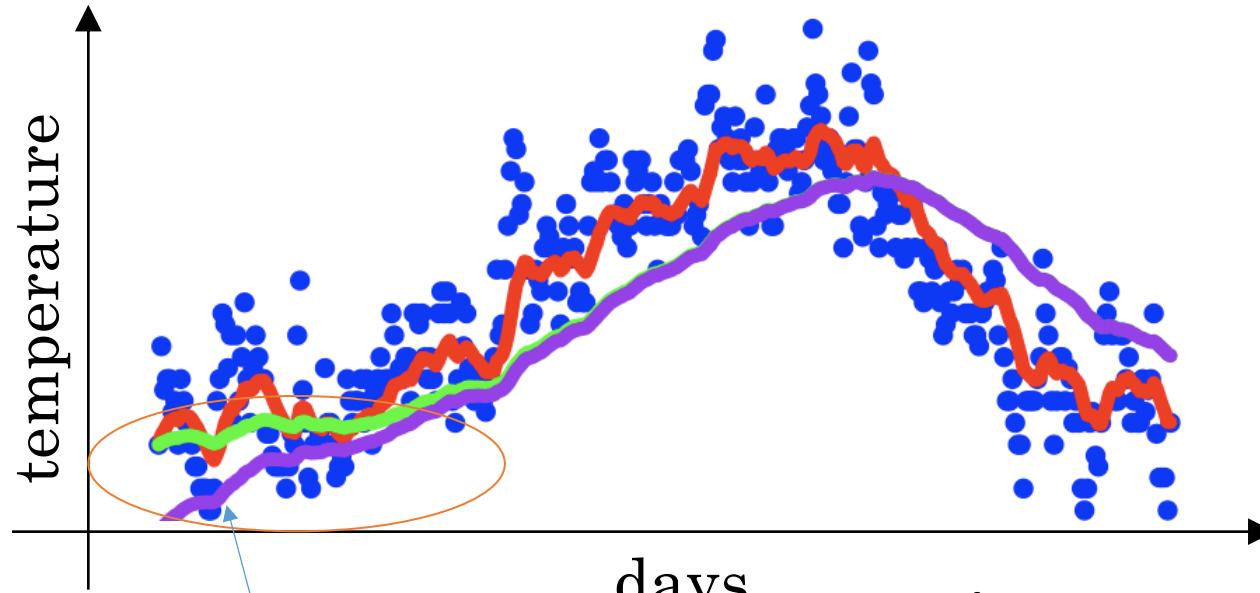
$$v_0 = 0$$

$$v_1 = 0.98v_0 + 0.02\theta_1 = 0.02\theta_1$$

$$v_2 = 0.98v_1 + 0.02\theta_2$$

$$= 0.0196\theta_1 + 0.02\theta_2$$

(lower values in early stage, purple line)



$$\text{Bias correction: } v_t = \frac{v_t}{1-\beta^t}$$

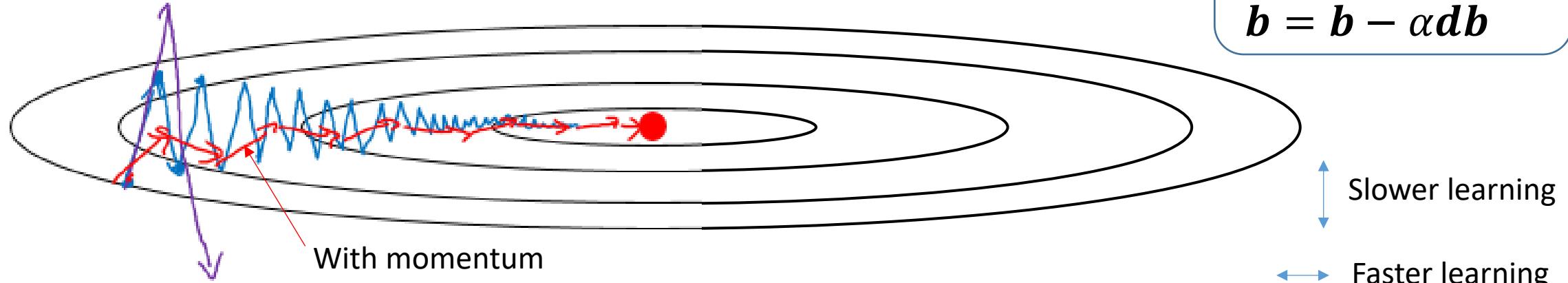
$$\text{When } t=1, 1 - \beta^1 = 0.02$$

$$\text{corrected } v_1 = \frac{v_1}{1-\beta^1} = \frac{0.02\theta_1}{0.02} = \theta_1$$

$$\text{When } t=2, 1 - \beta^2 = 0.0396$$

$$\text{corrected } v_2 = \frac{v_2}{1-\beta^2} = \frac{0.0196\theta_1 + 0.02\theta_2}{0.0396}$$

# Gradient descent with **momentum** uses exponentially weighted average



$V_{dW} = 0, v_{db} = 0$  # initialize running averages

On iteration  $t$ :

Compute  $dW, db$  on the current mini-batch

$$V_{dW} = \beta V_{dW} + (1 - \beta) dW$$

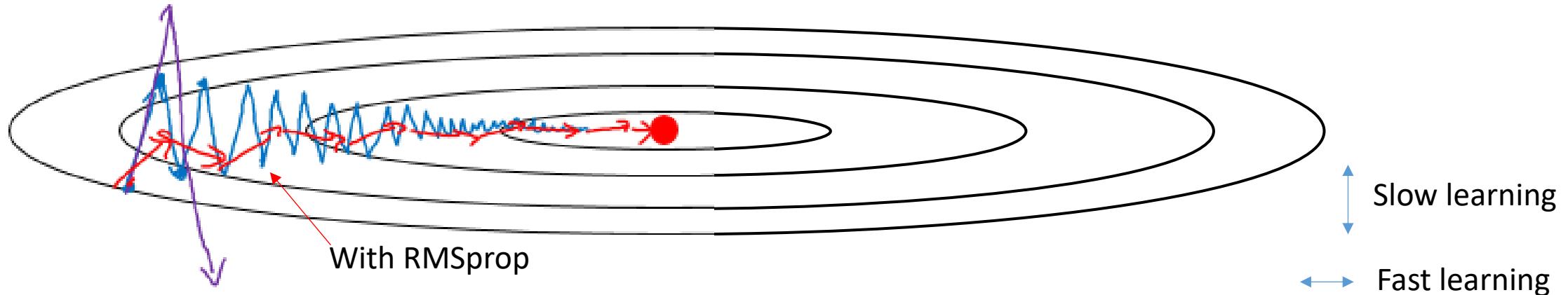
$$v_{db} = \beta v_{db} + (1 - \beta) db$$

$$W = W - \alpha V_{dW}, \quad b = b - \alpha v_{db}$$

Hyperparameters:  $\alpha, \beta$   
 $\beta = 0.9$

Average over last  $\approx 10$  gradients

# RMSprop (Root Mean Square Propagation)



$S_{dw} = 0, S_{db} = 0$  # initialize running averages

On iteration  $t$ :

Compute  $dW, db$  on the current mini-batch

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

Element-wise square of  $dW$  (or  $db$ )  
small

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

large

$$W = W - \alpha \frac{dW}{\sqrt{S_{dw} + \varepsilon}}, \quad b = b - \alpha \frac{db}{\sqrt{S_{db} + \varepsilon}}, \quad \varepsilon = 10^{-8}$$

## Adam (Adaptive moment estimation) optimization algorithm

$$V_{dw} = 0, S_{dw} = 0, v_{db} = 0, s_{db} = 0$$

On iteration  $t$ :

Compute  $dW, db$  on the current mini-batch

$$\begin{aligned} V_{dw} &= \beta_1 V_{dw} + (1 - \beta_1) dW, \quad v_{db} = \beta_1 v_{db} + (1 - \beta_1) db && \text{(from momentum } \beta_1\text{)} \\ S_{dw} &= \beta_2 S_{dw} + (1 - \beta_2) dW^2, \quad s_{db} = \beta_2 s_{db} + (1 - \beta_2) db^2 && \text{(from RMSprop } \beta_2\text{)} \end{aligned}$$

$$V_{dw}^{corrected} = \frac{V_{dw}}{1 - \beta_1^t}, \quad v_{db}^{corrected} = \frac{v_{db}}{1 - \beta_1^t}$$

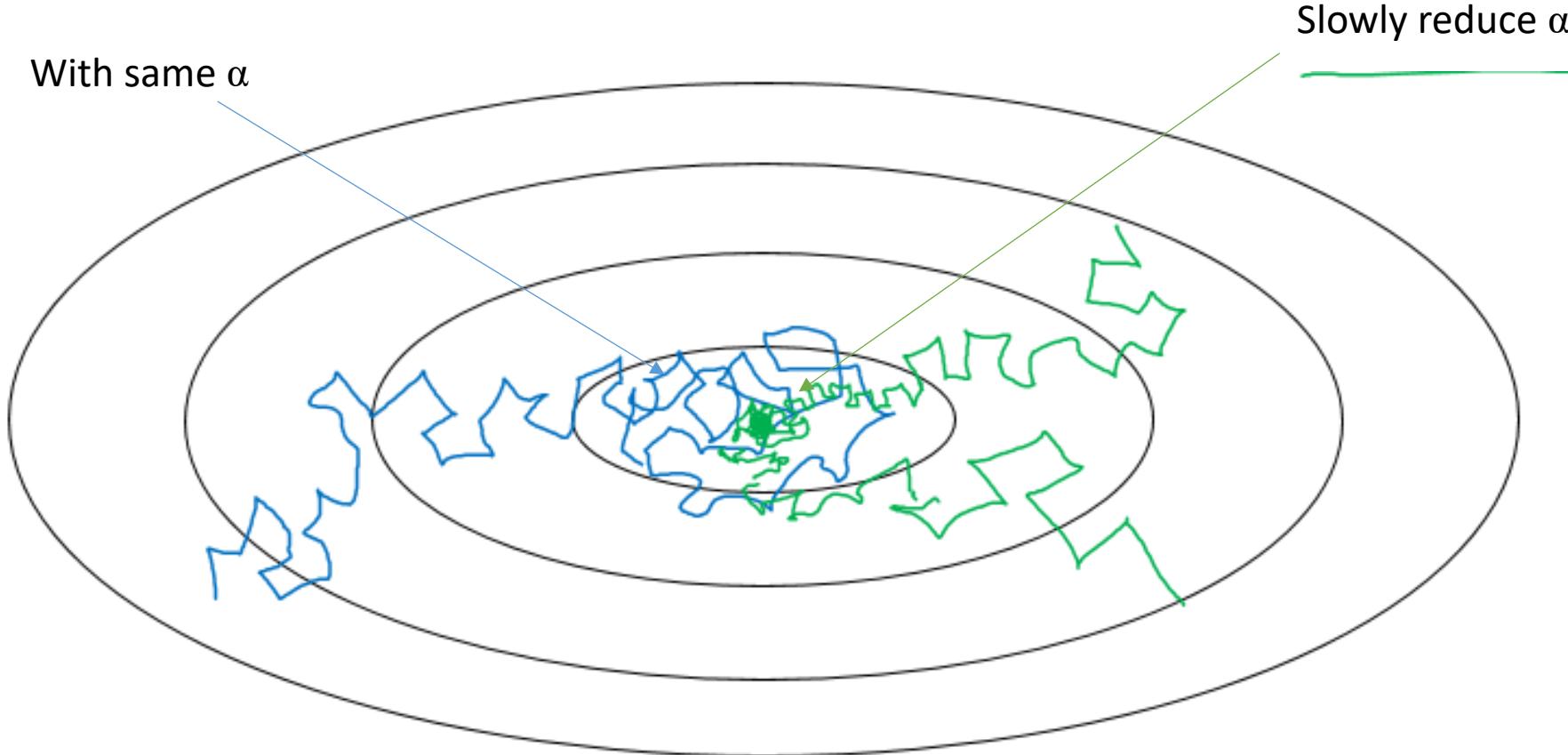
$$S_{dw}^{corrected} = \frac{S_{dw}}{1 - \beta_2^t}, \quad s_{db}^{corrected} = \frac{s_{db}}{1 - \beta_2^t}$$

Hyperparameters choice:  
 $\alpha$  : needs to be tune  
 $\beta_1$ : 0.9  
 $\beta_2$ : 0.999  
 $\epsilon$  :  $10^{-8}$

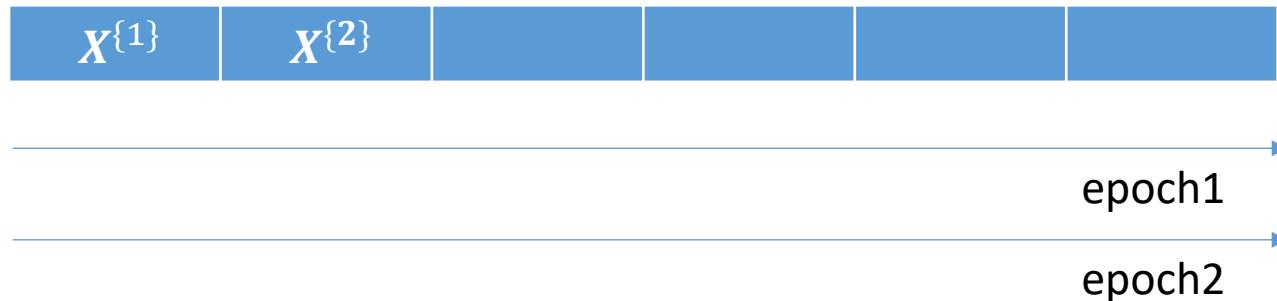
$$W = W - \alpha \frac{V_{dw}^{corrected}}{\sqrt{S_{dw}^{corrected} + \epsilon}}, \quad b = b - \alpha \frac{v_{db}^{corrected}}{\sqrt{s_{db}^{corrected} + \epsilon}}$$

# Learning rate decay

**Basic approach:**  
 $W = W - \alpha dW$   
 $b = b - \alpha db$



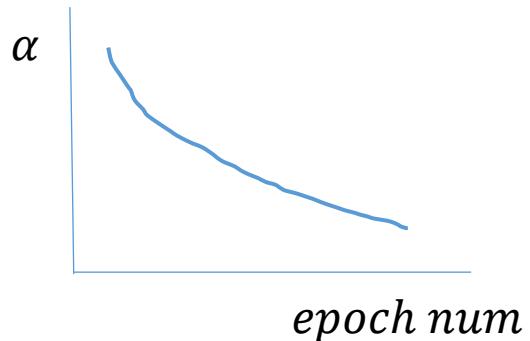
# Learning rate decay



1 epoch = 1 pass through data

$$\alpha = \frac{1}{1 + \text{decay-rate} * \text{epoch-num}} \alpha_0, \text{ where } \alpha_0=0.2 ; \text{decay-rate}=1$$

Epoch	$\alpha$
1	0.1
2	0.067
3	0.05
4	0.04
...	
...	



**Basic approach:**  
 $W = W - \alpha dW$   
 $b = b - \alpha db$

# Tuning Process: Hyperparameters

$\alpha$  (the most important)

$\beta$  (the second most important)

# of hidden units (the second most important)

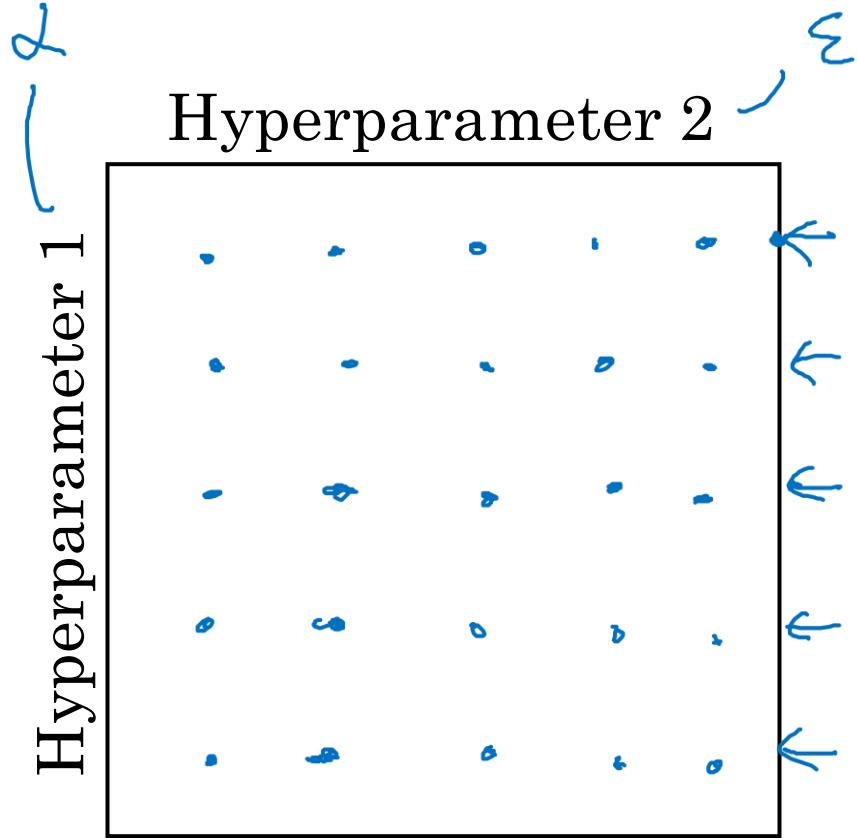
Mini-batch size (the second most important)

# of layers (the third most important)

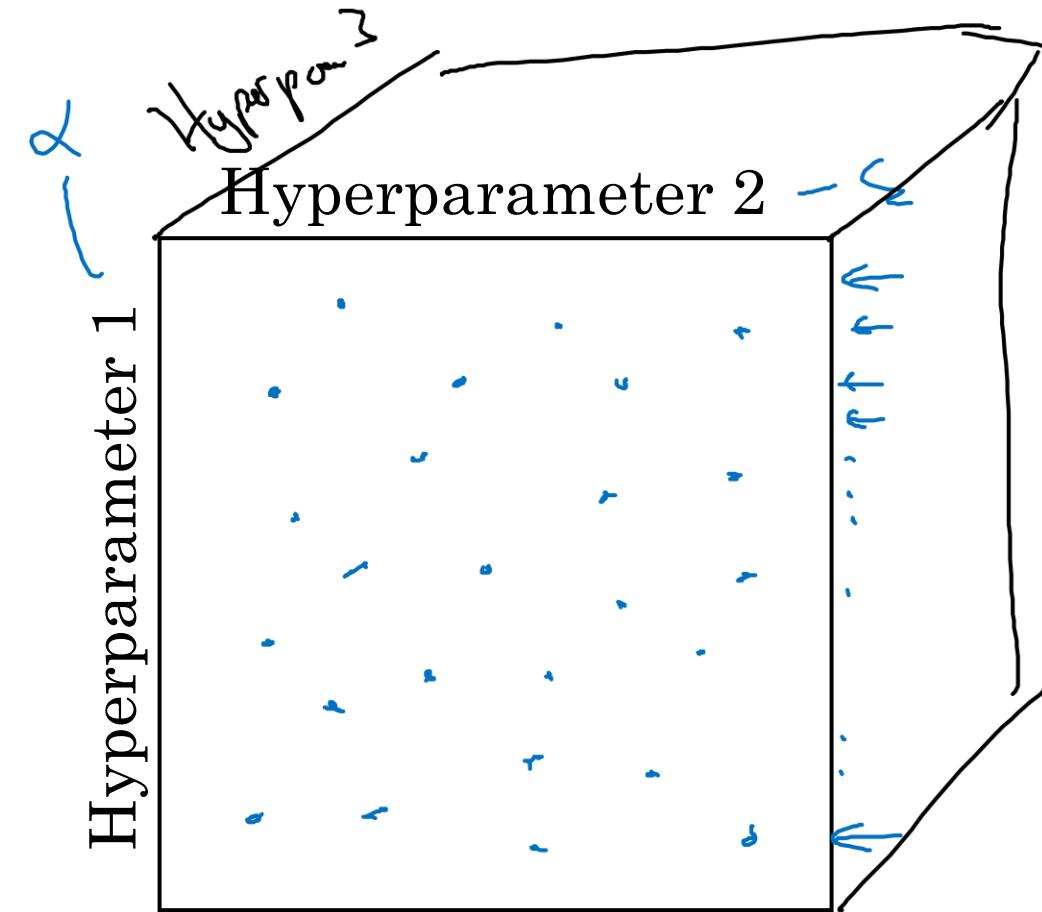
Learning rate decay (the third most important)

$\beta_1, \beta_2, \epsilon$  in adam

# Try random values: Don't use a grid

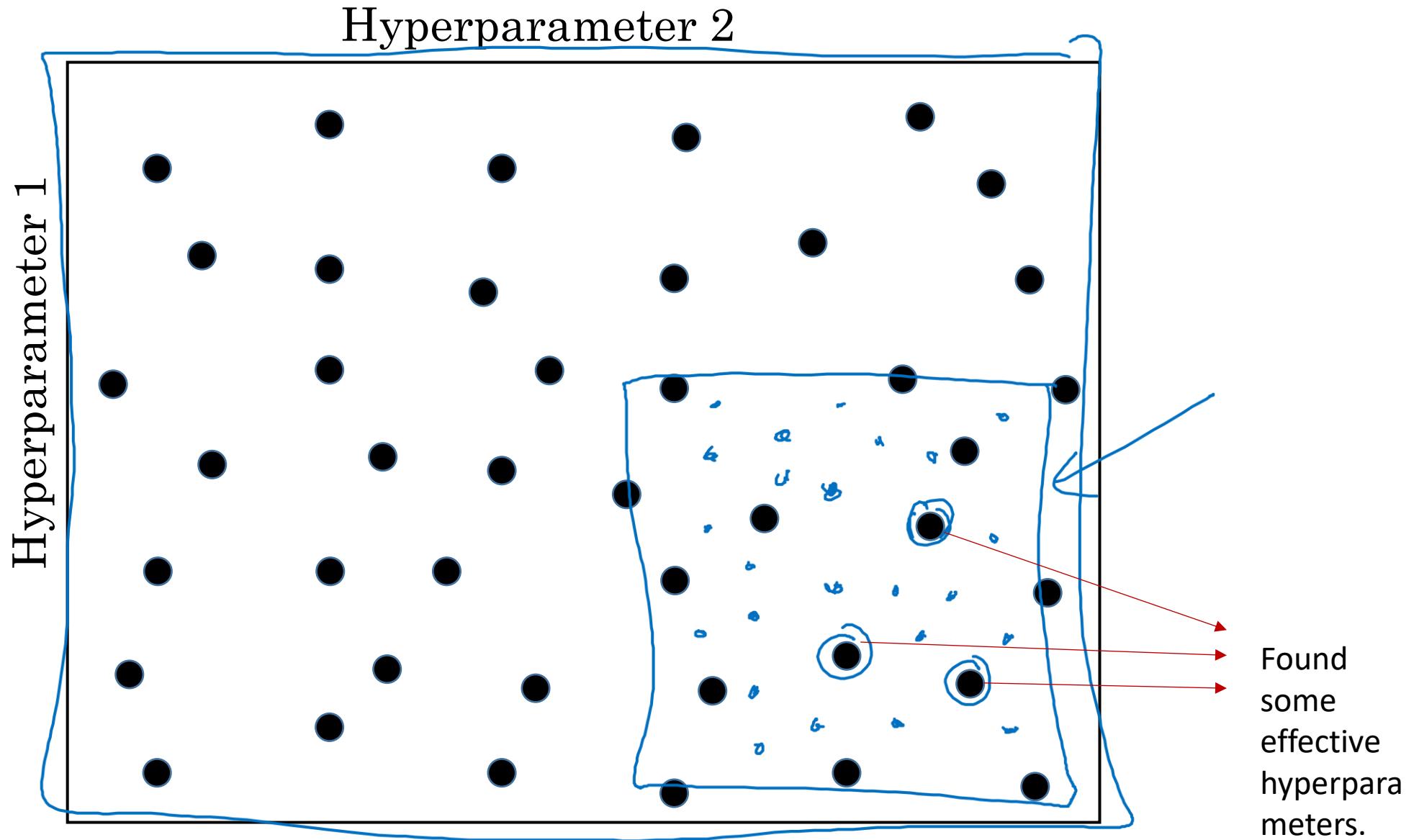


Consider only 5 different  $\alpha$  values



Consider 25 different  $\alpha$  values  
( $\alpha$  is more important than  $\varepsilon$ )

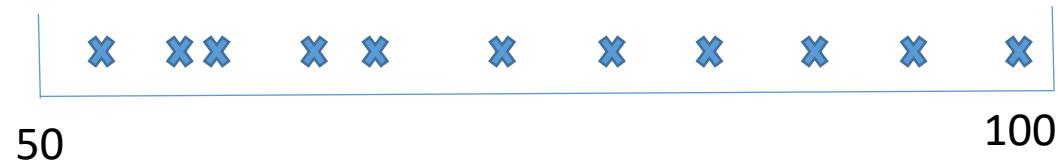
# Coarse to fine



# Using an appropriate scale to pick hyperparameters:

## Picking hyperparameters at random

# of hidden units:  $N^{[l]} = 50, \dots, 100$

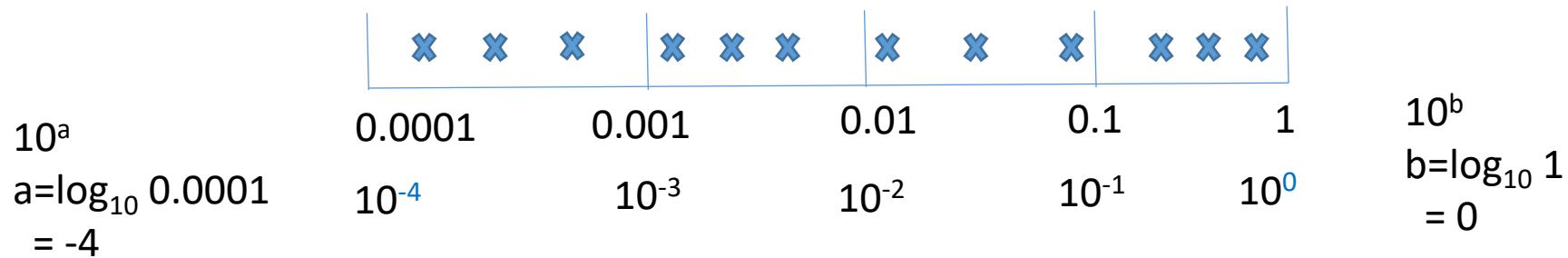


# of layers:  $L = 2 - 4$

2, 3, 4

# Appropriate scale for hyperparameters

Learning Rate:  $\alpha = 0.0001, \dots, 1$

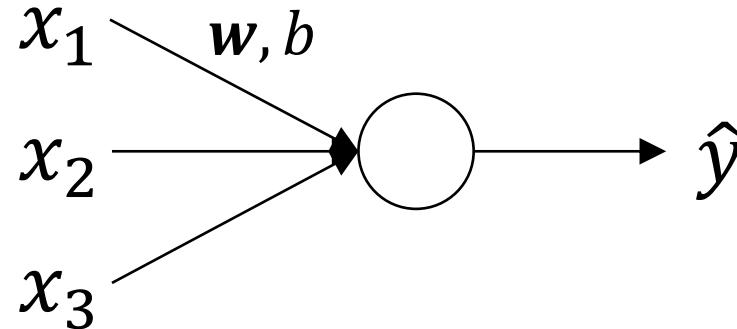


$$r = -4 * \text{np.random.rand}() \quad \leftarrow r \in [-4, 0]$$
$$\alpha = 10^r \quad \leftarrow 10^{-4} \dots 10^0$$

$10^a \dots 10^b ; r \in [a, b] ; \alpha = 10^r$  where  $r = \text{numpy.random.uniform}(low=a, high=b)$

# Batch Normalization: Normalizing activations in a network

Normalizing inputs to speed up learning



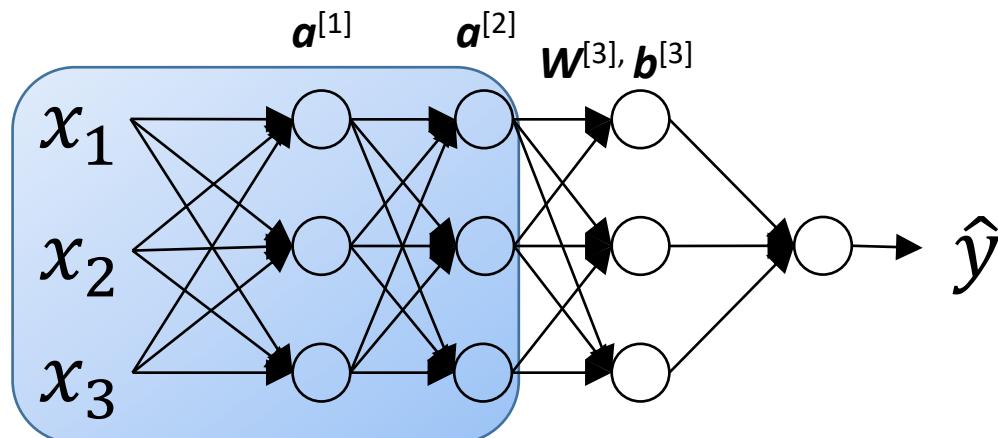
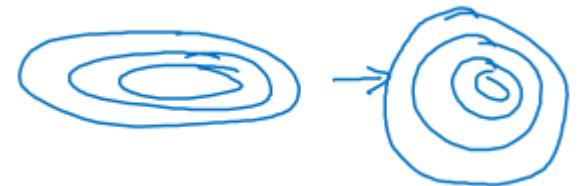
$$\mu = \frac{1}{m} \sum_{i=1}^m x^{(i)}$$

E.g., use  $m$  number of  $x_1, x_2, x_3$  to calculate  $\mu_1, \mu_2, \mu_3$

$$\sigma^2 = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu)^2$$

Element-wise

$$x^{(i)} = (x^{(i)} - \mu) / \sigma$$



Can we normalize  $a^{[2]}$  to train  $w^{[3]}, b^{[3]}$  faster?

-> need to normalize  $z^{[2]}$

# Implementing Batch Norm

Given some intermediate values in NN,  $\mathbf{z}^{(1)}, \mathbf{z}^{(2)}, \dots, \mathbf{z}^{(m)}$ , in a layer L.

$$\mu = \frac{1}{m} \sum_{i=1}^m \mathbf{z}^{(i)}$$

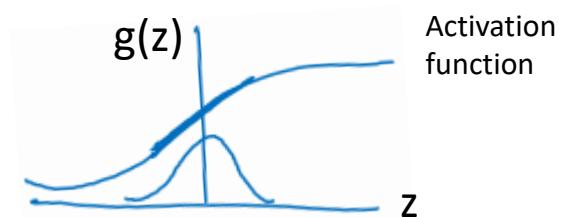
$$\sigma^2 = \frac{1}{m} \sum_{i=1}^m (\mathbf{z}^{(i)} - \mu)^2$$

$$\mathbf{z}_{norm}^{(i)} = \frac{\mathbf{z}^{(i)} - \mu}{\sqrt{\sigma^2 + \epsilon}} \quad (\text{mean 0 and standard deviation 1})$$

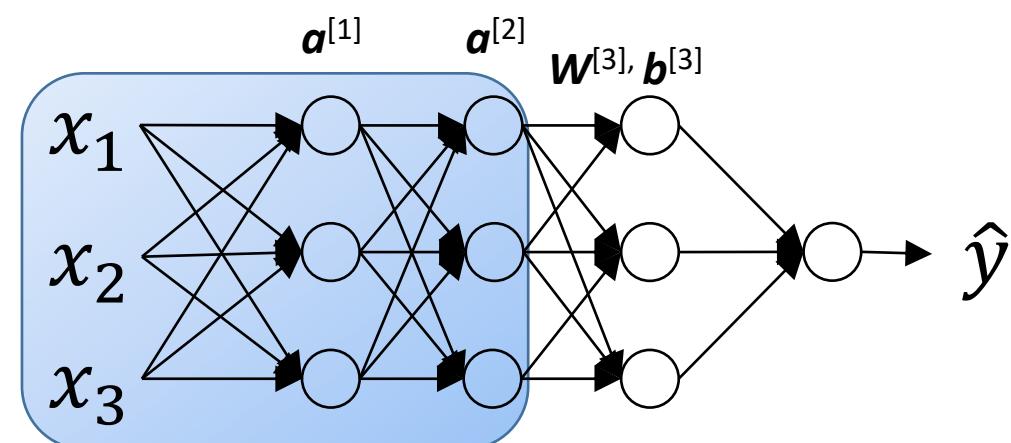
$$\mathbf{Z} = \begin{bmatrix} & & & \\ | & | & | & | \\ \mathbf{z}^{(1)} & \mathbf{z}^{(2)} & \dots & \mathbf{z}^{(m)} \\ | & | & & | \end{bmatrix}$$

$$\tilde{\mathbf{z}}^{(i)} = \gamma \mathbf{z}_{norm}^{(i)} + \beta \quad \text{where } \gamma \text{ and } \beta \text{ are learnable parameters}$$

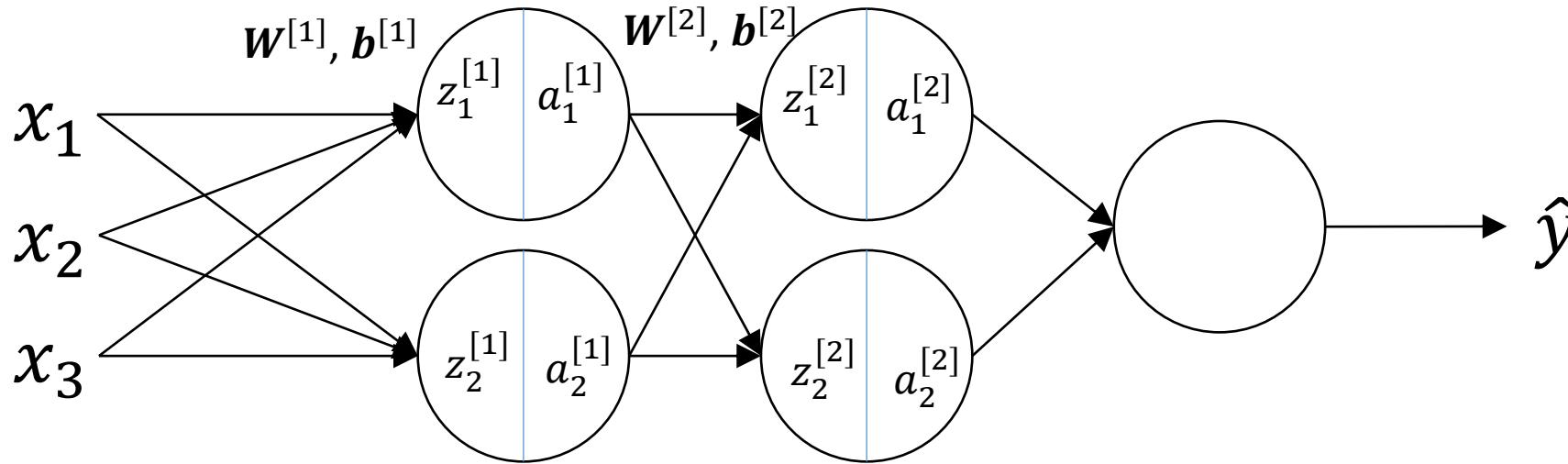
(adjust mean and standard deviation).



If  $\gamma = \sqrt{\sigma^2 + \epsilon}$  and  $\beta = \mu$ , then  $\tilde{\mathbf{z}}^{(i)} = \mathbf{z}^{(i)}$



# Fitting Batch Norm into a neural network: Adding Batch Norm to a network



$$x \rightarrow Z^{[1]} \xrightarrow{\text{Batch Norm (BN)}} \tilde{Z}^{(1)} \rightarrow A^{[1]} = g^{[1]}(\tilde{Z}^{(1)}) \rightarrow Z^{[2]} \xrightarrow{\text{BN}} \tilde{Z}^{(2)} \rightarrow A^{[2]} = g^{[2]}(\tilde{Z}^{(2)})$$

Parameters:  $\mathbf{W}^{[1]}, \mathbf{b}^{[1]}, \mathbf{W}^{[2]}, \mathbf{b}^{[2]}, \dots, \mathbf{W}^{[L]}, \mathbf{b}^{[L]}$   
 $\boldsymbol{\beta}^{[1]}, \boldsymbol{\gamma}^{[1]}, \boldsymbol{\beta}^{[2]}, \boldsymbol{\gamma}^{[2]}, \dots, \boldsymbol{\beta}^{[L]}, \boldsymbol{\gamma}^{[L]}$

$$\boldsymbol{\beta}^{[l]} = \boldsymbol{\beta}^{[l]} - \alpha d\boldsymbol{\beta}^{[l]}$$
$$\boldsymbol{\gamma}^{[l]} = \boldsymbol{\gamma}^{[l]} - \alpha d\boldsymbol{\gamma}^{[l]}$$

# Working with mini-batches

$m$  is the size of a mini batch

$$\boldsymbol{\mu} = \frac{1}{m} \sum_{i=1}^m \mathbf{z}^{(i)}$$

$$\sigma^2 = \frac{1}{m} \sum_{i=1}^m (\mathbf{z}^{(i)} - \boldsymbol{\mu})^2$$

$$\mathbf{z}_{norm}^{(i)} = \frac{\mathbf{z}^{(i)} - \boldsymbol{\mu}}{\sqrt{\sigma^2 + \epsilon}}$$

$$\tilde{\mathbf{z}}^{(i)} = \gamma \mathbf{z}_{norm}^{(i)} + \beta$$

$$W^{[1]}, b^{[1]} \quad \beta^{[1]}, \gamma^{[1]}$$

$$\mathbf{x}^{\{1\}} \rightarrow \mathbf{z}^{[1]} \rightarrow \tilde{\mathbf{z}}^{(1)}$$

Batch Norm (BN)

$$\rightarrow A^{[1]} = g^{[1]}(\tilde{\mathbf{z}}^{(1)}) \rightarrow \mathbf{z}^{[2]} \xrightarrow{\text{BN}} \tilde{\mathbf{z}}^{(2)} \rightarrow A^{[2]}$$

$$W^{[1]}, b^{[1]} \quad \beta^{[1]}, \gamma^{[1]}$$

$$\mathbf{x}^{\{2\}} \rightarrow \mathbf{z}^{[1]} \rightarrow \tilde{\mathbf{z}}^{(1)}$$

Batch Norm (BN)

$$\rightarrow A^{[1]} = g^{[1]}(\tilde{\mathbf{z}}^{(1)}) \rightarrow \mathbf{z}^{[2]} \xrightarrow{\text{BN}} \tilde{\mathbf{z}}^{(2)} \rightarrow A^{[2]}$$

.....

$$W^{[1]}, b^{[1]} \quad \beta^{[1]}, \gamma^{[1]}$$

$$\mathbf{x}^{\{n\}} \rightarrow \mathbf{z}^{[1]} \rightarrow \tilde{\mathbf{z}}^{(1)}$$

Batch Norm (BN)

$$\rightarrow A^{[1]} = g^{[1]}(\tilde{\mathbf{z}}^{(1)}) \rightarrow \mathbf{z}^{[2]} \xrightarrow{\text{BN}} \tilde{\mathbf{z}}^{(2)} \rightarrow A^{[2]}$$

Parameters:  $\mathbf{W}^{[L]}, \mathbf{b}^{[L]}, \dots, \beta^{[L]}, \gamma^{[L]}$        $\mathbf{z}^{[L]} = \mathbf{W}^{[L]} \mathbf{A}^{[L-1]} + \mathbf{b}^{[L]}$

( $n^{[L]}, 1$ )    ( $n^{[L]}, 1$ )    ( $n^{[L]}, 1$ )

# Implementing gradient descent

For  $t=1 \dots$  number of Mini Batches

compute forward prop on  $\mathbf{X}^{\{t\}}$

in each hidden layer, use BN to replace  $\mathbf{Z}^{[L]}$  with  $\tilde{\mathbf{z}}^{(L)}$

Use backprop to compute  $dW^{[L]}, db^{[L]}, \dots, d\beta^{[L]}, d\gamma^{[L]}$

Update parameter

$$\mathbf{W}^{[l]} = \mathbf{W}^{[l]} - \alpha d\mathbf{W}^{[l]}$$

$$\boldsymbol{\beta}^{[l]} = \boldsymbol{\beta}^{[l]} - \alpha d\boldsymbol{\beta}^{[l]}$$

$$\boldsymbol{\gamma}^{[l]} = \boldsymbol{\gamma}^{[l]} - \alpha d\boldsymbol{\gamma}^{[l]}$$

$$\mathbf{b}^{[l]} = \mathbf{b}^{[l]} - \alpha db^{[l]}$$

Also works with momentum, RMSprop, and Adam

# Batch Norm at test time: one test data

In each layer L **at training time**:

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

*m* is the size of a mini batch

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

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

$$\tilde{z}^{(i)} = \gamma z_{\text{norm}}^{(i)} + \beta$$

$\mu, \sigma^2$  : estimate using exponentially weighted average (across mini-batch) during training time

$$\mu^{[l]} = \beta \mu^{[l]} + (1 - \beta) \mu^{\{\text{cur-batch}\}[l]}$$

$X^{[1]}, X^{[2]}, X^{[3]}, \dots$  mini batch

$\mu^{[1][l]}, \mu^{[2][l]}, \mu^{[3][l]}, \dots \rightarrow \mu^{[l]}$  in each layer L

$\sigma^{2[1][l]}, \sigma^{2[2][l]}, \sigma^{2[3][l]}, \dots \rightarrow \sigma^{2[l]}$  in each layer L

---

Test time for each test data: (in each layer L)

$$z_{\text{norm}}^{(i)} = \frac{z^{(i)} - \mu^{[l]}}{\sqrt{\sigma^{2[l]} + \epsilon}} \quad \tilde{z}^{(i)} = \gamma z_{\text{norm}}^{(i)} + \beta$$

# Summary of today's lecture

- Learned various approaches for model optimization.
  - Regularization and Early stopping
  - Input normalization and Weights initialization
  - Mini-batch gradient descent
  - Gradient descent with momentum, RMSprop, Adam
  - Learning rate decay
  - Hyperparameter tuning process
  - Batch normalization



# Reference

- Neural Network and Deep Learning, Andrew Ng,  
<https://www.coursera.org/learn/neural-networks-deep-learning>