# Neural Network

Ratchainant Thammasudjarit, Ph.D.

# Linear Algebra

*Function*

- A function of univariate is defined as

$$x \mapsto f(x)$$

INPUT x

FUNCTION f:

OUTPUT f(x)

Each input must have only one output

$f(x)$

This is a function

$f(x)$

This is NOT a function

*Wisdom of the Land*

# Linear Algebra

*sigmoid*



- Properties
  - f(x) ∈ [0, 1]
  - x ∈ [-∞, ∞]
  - x → 0, f(x) becomes linear
  - abs(x) > 4, f(x) changes slowly

# Linear Algebra

*tanh*
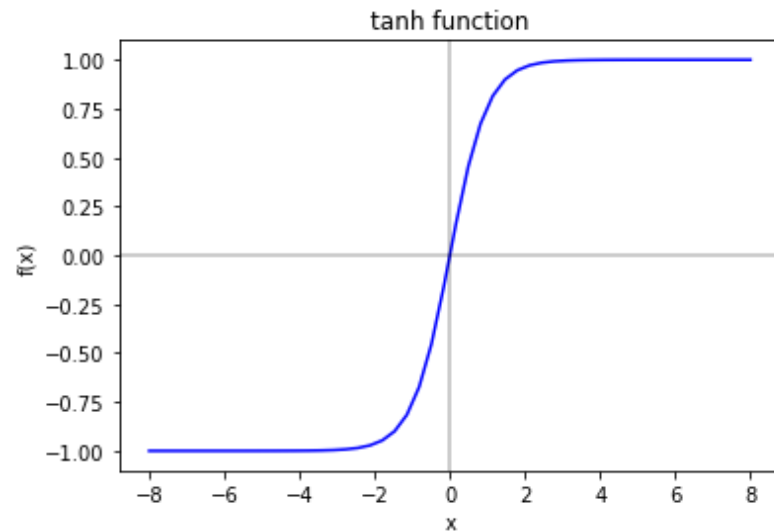


- Properties
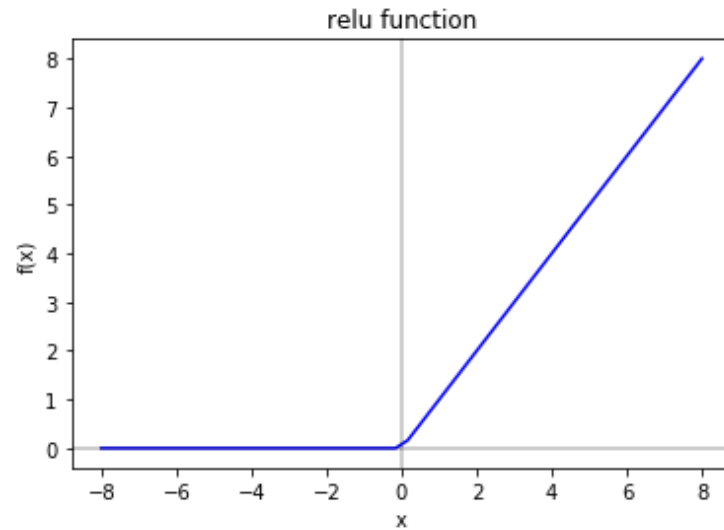  - $f(x) \in [-1, 1]$
  - $x \in [-\infty, \infty]$
  - $x \to 0$, $f(x)$ becomes linear
  - $abs(x) > 2$, $f(x)$ changes slowly

Wisdom of the Land

# Linear Algebra

*ReLu (Rectifier Linear Unit)*



relu function

```
def relu(x):
    y = np.maximum(0, x)
    return y
```

- Properties
  - $f(x) \in [0, \infty]$
  - $x \in [-\infty, \infty]$
  - $x \leq 0, f(x) = 0$
  - $x > 0, f(x) = x$

# Linear Algebra

*Vector*

$$\vec{u} = \langle 3, 4 \rangle$$
$$= 3 \cdot \vec{\mathbf{i}} + 4 \cdot \vec{\mathbf{j}}$$

$$\|\vec{u}\| = \sqrt{3^2 + 4^2} = 5$$

- Dot product

$$\vec{u} = \langle 3, 4 \rangle$$
$$\vec{v} = \langle 4, 3 \rangle$$

$$\vec{u} \cdot \vec{v} = (3 \times 4) + (4 \times 3) = 24$$

- Transpose

if
$$\mathbf{u} = \begin{bmatrix} 3 & 4 \end{bmatrix}$$

then
$$\mathbf{u}^T = \begin{bmatrix} 3 \\ 4 \end{bmatrix}$$

# Linear Algebra

*Matrix*

    *A stack of vectors*

- Pairwise multiplication

$$A = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} \qquad B = \begin{bmatrix} 2 & 2 & 2 \\ 2 & 2 & 2 \end{bmatrix}$$

$$A \odot B = \begin{bmatrix} 1 \times 2 & 2 \times 2 & 3 \times 2 \\ 4 \times 2 & 5 \times 2 & 6 \times 2 \end{bmatrix} = \begin{bmatrix} 2 & 4 & 6 \\ 8 & 10 & 12 \end{bmatrix}$$
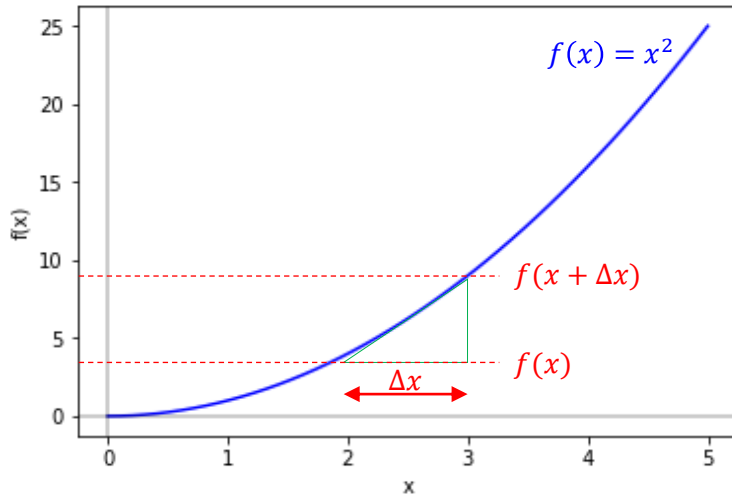
- Transpose

$$A^T = \begin{bmatrix} 1 & 4 \\ 2 & 5 \\ 3 & 6 \end{bmatrix}$$

# Linear Algebra

*Calculus (Derivative)*

- Derivative

$$f'(x) = \lim_{\Delta x \to 0} \frac{f(x + \Delta x) - f(x)}{\Delta x}$$



$f(x) = x^2$

$f(x + \Delta x)$

$f(x)$

$\Delta x$

$\Delta x = 1$

if $x = 2$     then $f(x) = 4$
if $x = 3$     then $f(x) = 9$

$\Delta x = 0.001$

if $x = 2$     then $f(x) = 4$
if $x = 2.001$   then $f(x) \approx 4.004$

if $x = 2$   then slope $= \dfrac{\Delta f(x)}{\Delta x} = \dfrac{0.004}{0.001} = 4$

Mathematically,   $\dfrac{d}{dx} f(x) = 2x = 2(2) = 4$

# Linear Algebra

$$g(z) = \frac{1}{1 + e^{-z}}$$

$$\frac{d}{dz}g(z) = \frac{d}{dz}\left[\frac{1}{1 + e^{-z}}\right]$$

$$= \frac{d}{dz}(1 + e^{-z})^{-1}$$

$$= -(1 + e^{-z})^{-2}\frac{d}{dz}(1 + e^{-z})$$

$$= -(1 + e^{-z})^{-2}\frac{d}{dz}(e^{-z})$$

$$= -(1 + e^{-z})^{-2}(-e^{-z})$$

$$= \frac{e^{-z}}{(1 + e^{-z})^2}$$

$$= \frac{1}{1 + e^{-z}} \cdot \frac{e^{-z}}{1 + e^{-z}}$$

$$= \frac{1}{1 + e^{-z}} \cdot \frac{(1 + e^{-z}) - 1}{1 + e^{-z}}$$

$$= \frac{1}{1 + e^{-z}}\left[\frac{1 + e^{-z}}{1 + e^{-z}} - \frac{1}{1 + e^{-z}}\right]$$

$$= \frac{1}{1 + e^{-z}}\left[1 - \frac{1}{1 + e^{-z}}\right]$$

$$\therefore \frac{d}{dz}g(z) = g(z)(1 - g(z))$$

# Linear Algebra

*Derivative of tanh*

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

$$= \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

$$\frac{d}{dz}g(z) = \frac{(e^z + e^{-z})d(e^z - e^{-z}) - (e^z - e^{-z})d(e^z + e^{-z})}{(e^z + e^{-z})^2}$$

$$= \frac{(e^z + e^{-z})(e^z + e^{-z}) - (e^z - e^{-z})(e^z - e^{-z})}{(e^z + e^{-z})^2}$$

$$= 1 - \left(\frac{e^z - e^{-z}}{e^z + e^{-z}}\right)^2$$

$$\therefore \frac{d}{dz}g(z) = 1 - \tanh^2(z)$$

# Linear Algebra

*Derivative of relu*

$$g(z) = \max(0, z)$$

$$\frac{d}{dz} g(z) = \begin{cases} 0 & if \ z < 0 \\ 1 & if \ z > 0 \\ undefined & if \ z = 0 \end{cases}$$

In neural network practice, z can get close to zero but never be zero, e.g. 0.00000 …

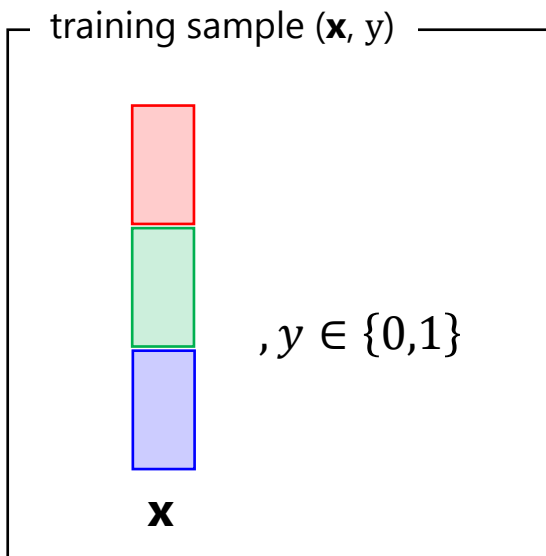$$\therefore \frac{d}{dz} g(z) = \begin{cases} 0 & if \ z < 0 \\ 1 & if \ z > 0 \end{cases}$$
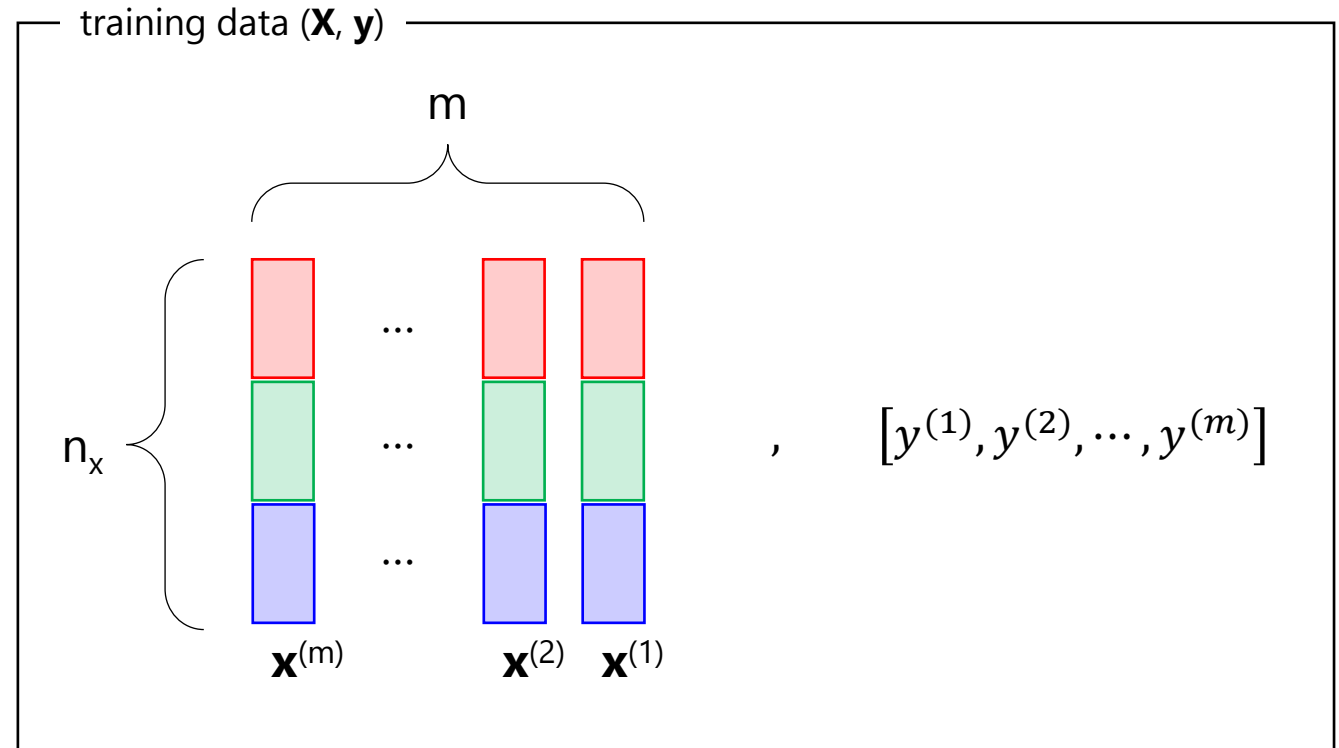
# Perceptron

A smallest unit in neural network

# Training Sample VS Training Data

*Training sample is a pair of feature vector **x** and its class label*

$$(\mathbf{x}, y), \mathbf{x} \in \mathbb{R}^{n_x}, y \in \{0,1\}$$

training sample (**x**, y)



$, y \in \{0,1\}$

**x**

*Training data is <u>m pair</u> of feature vector **x** and its class label*

training data (**X**, **y**)



m

$n_x$

$\mathbf{x}^{(m)}$   $\mathbf{x}^{(2)}$   $\mathbf{x}^{(1)}$
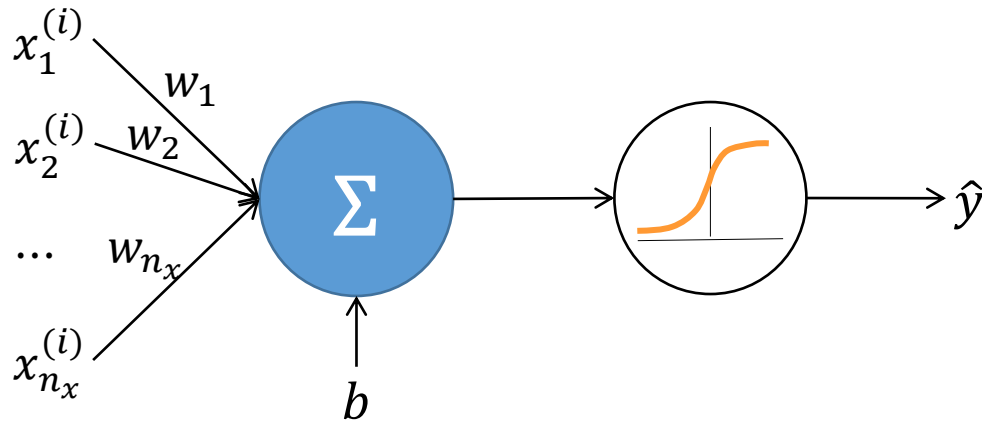
$\left[y^{(1)}, y^{(2)}, \cdots, y^{(m)}\right]$

$\mathbf{x} \in \mathbb{R}^{n_x \times 1}$
$\mathbf{X} \in \mathbb{R}^{n_x \times m}$
$\mathbf{y} \in \mathbb{R}^{1 \times m}$
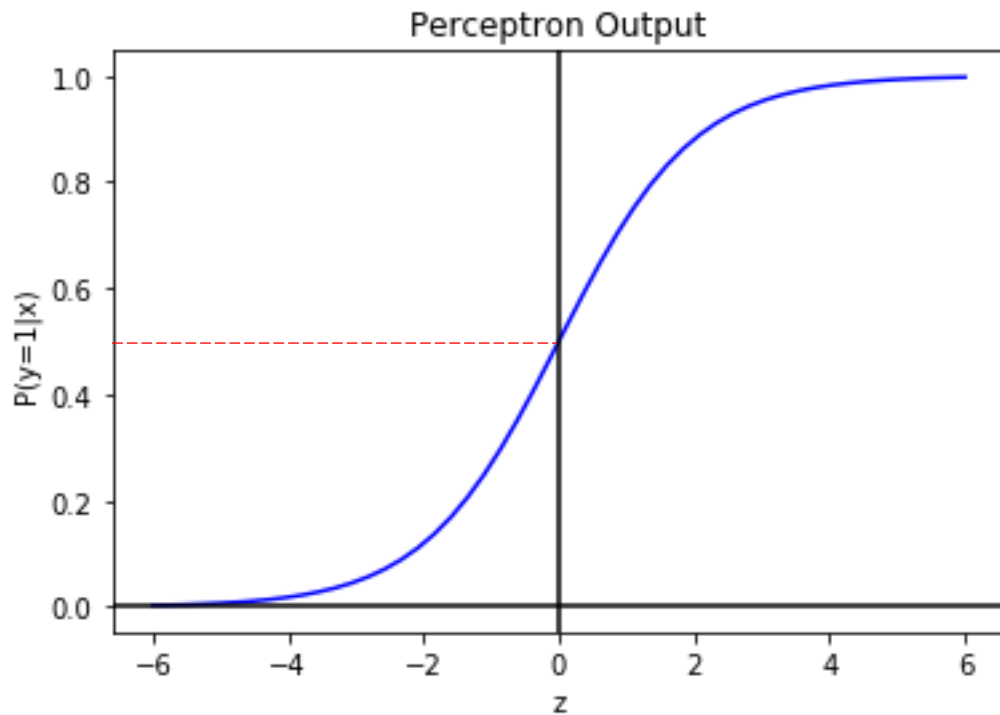
# A Perceptron

*The smallest unit of neural network*



$$\hat{y} = P(y = 1|\mathbf{x}^{(i)}) = \frac{1}{1 + e^{-(\mathbf{w}^T \cdot \mathbf{x}^{(i)} + b)}}$$

- Given $\mathbf{x}^{(i)} \in \mathbb{R}^{n_x \times 1}$ , determine
  - $\hat{y} = P(y = 1|\mathbf{x}^{(i)})$

- Parameters
  - $\mathbf{w} \in \mathbb{R}^{1 \times n_x}$
  - $b \in \mathbb{R}$

- Output
  - $\hat{y} = g(z) = g(\mathbf{w}^T \cdot \mathbf{x}^{(i)} + b)$
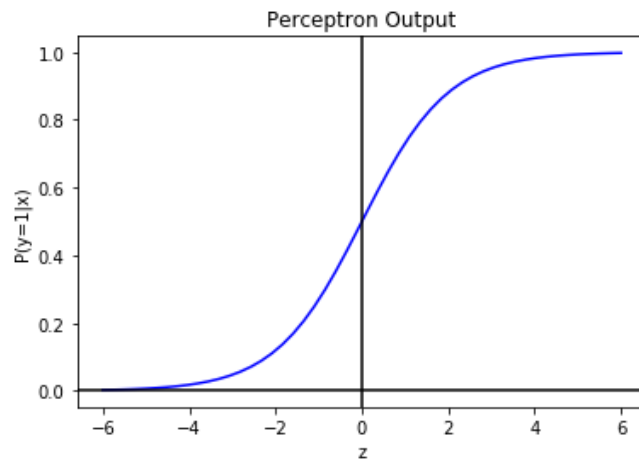
# A Perceptron

*The smallest unit of neural network*

## Perceptron Output



$$g(z) = \frac{1}{1 + e^{-z}}$$

- If $z \gg 0, g(z) \rightarrow 1$

- if $z \ll 0, g(z) \rightarrow 0$

# Loss Function

*Error of a training sample*

The loss function $\mathcal{L}(\hat{y}, y)$ determines how close of $\hat{y}$ to the ground-truth $y$



$$\mathcal{L}(\hat{y}, y) = -(y \log \hat{y} - (1-y) \log(1-\hat{y}))$$

Goal: $\mathcal{L}(\hat{y}, y) \to 0$

if $y = 1$,

$$\therefore z = w^T \cdot x + b, \qquad \to \infty$$

$$\hat{y} = \frac{1}{1 + e^{-z}}, \qquad \to 1$$

$$\mathcal{L}(\hat{y}, y) = -\log \hat{y}, \qquad \to 0$$

if $y = 0$,

$$\therefore z = w^T \cdot x + b, \qquad \to -\infty$$

$$\hat{y} = \frac{1}{1 + e^{-z}}, \qquad \to 0$$

$$\mathcal{L}(\hat{y}, y) = -\log(1 - \hat{y}), \qquad \to 0$$

# Cost Function

*Error of training data*

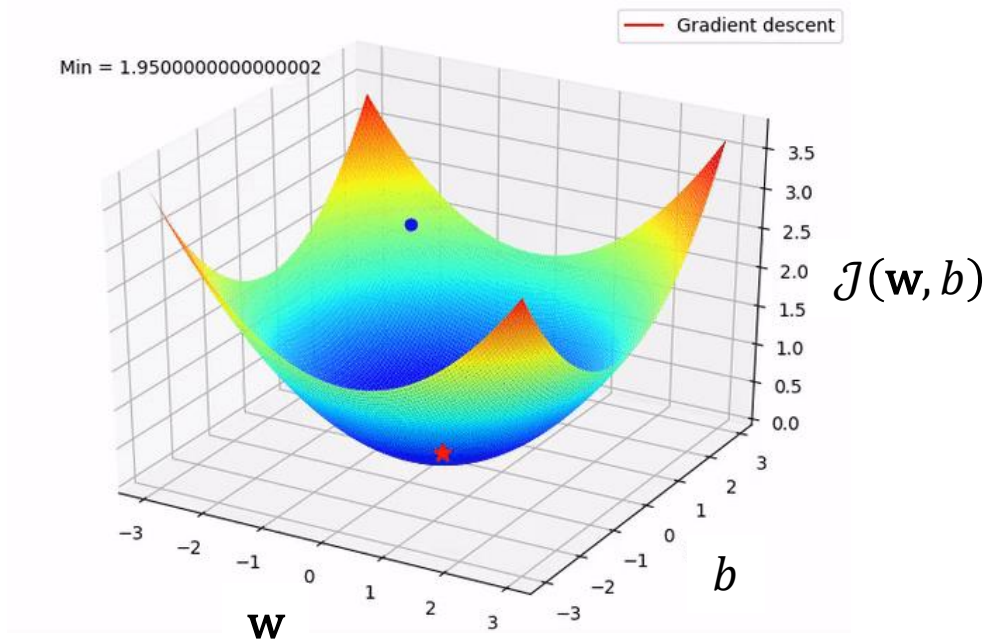The cost function $\mathcal{J}(\mathbf{w}, b)$ indicates how well the model does in entire training samples

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

$$= -\frac{1}{m}\sum_{i=1}^{m}\left[y^{(i)}\log\hat{y}^{(i)} + (1 - y^{(i)})\log(1 - \hat{y}^{(i)})\right]$$

- Training Goal
  - Find $\mathbf{w}$, $b$ that minimize $\mathcal{J}(\mathbf{w}, b)$
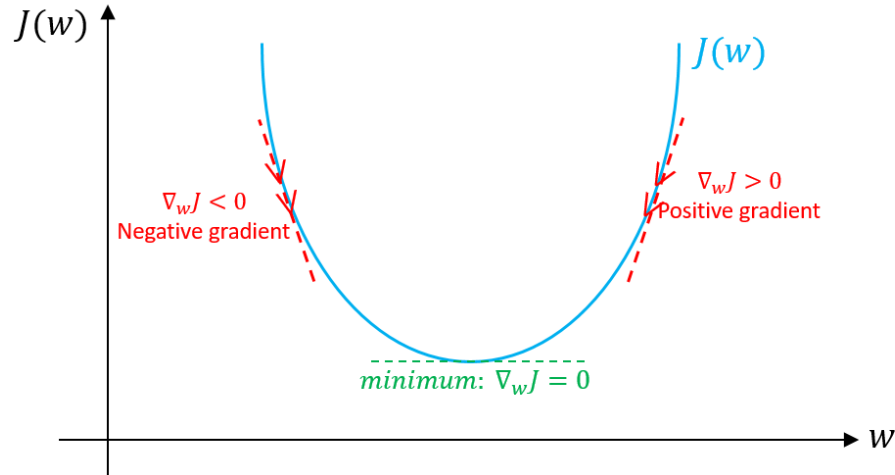
# Gradient Descent

- Goal
  - Determine $\mathbf{w}$, $b$ that minimize $\mathcal{J}(\mathbf{w}, b)$

Wisdom of the Land

# Gradient Descent

Initialize weight and adjust until the cost function approach to minimum



- Procedure

  Repeat {

  $$w := w - \alpha \frac{d}{dw} \mathcal{J}(w)$$

  Until $\mathcal{J}(w) \rightarrow \min(\mathcal{J}(w))$

  }

  where $\alpha$ represents the learning rate (a small positive value)

  On the right

  $$\frac{d}{dw} \mathcal{J}(w) > 0, \qquad \rightarrow w \text{ is adjusted by decreasing } \frac{d}{dw} \mathcal{J}(w)$$

  On the left

  $$\frac{d}{dw} \mathcal{J}(w) < 0, \qquad \rightarrow w \text{ is adjusted by increasing } \frac{d}{dw} \mathcal{J}(w)$$

# Gradient Descent

*Summary*

- Minimizing  $\mathcal{J}(\mathbf{w}, b)$

  - $w := w - \alpha \dfrac{d}{dw} \mathcal{J}(w, b),$    where $w \in \mathbf{w}$

  - $b := b - \alpha \dfrac{d}{db} \mathcal{J}(w, b)$

- In each training iteration, we need to determine

  $$\frac{d}{dw} \mathcal{J}(w, b) \quad \text{and} \quad \frac{d}{db} \mathcal{J}(w, b)$$

Wisdom of the Land

# Computational Graph

*Forward path*

Given a function  $\mathcal{J}(a, b, c) = 3(a + bc)$

$a = 5$

$b = 3$

$c = 2$

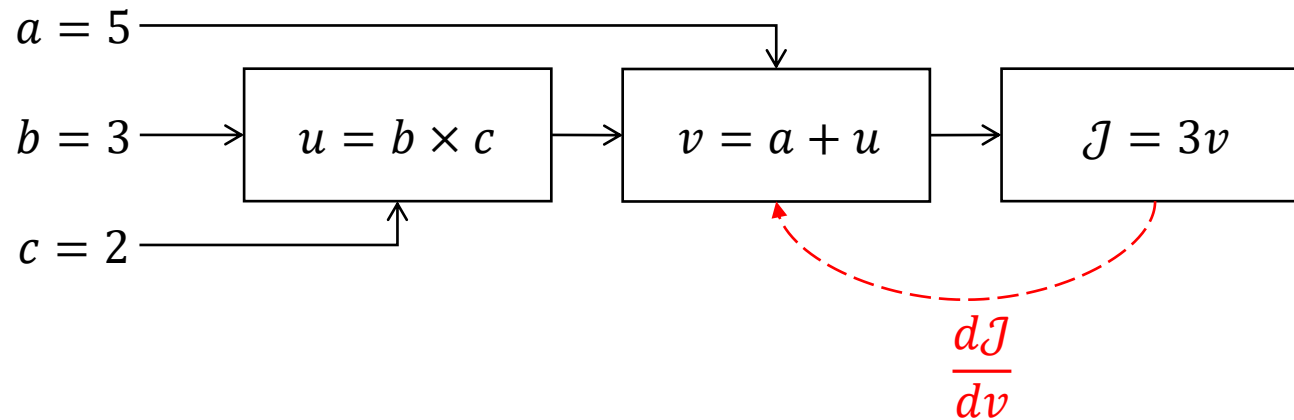| $u = b \times c$ | $v = a + u$ | $\mathcal{J} = 3v$ |
|:---:|:---:|:---:|
| **6** | **11** | **33** |

- Forward path
  - $\mathcal{J}(a, b, c)$ can be determined

- Backward path
  - Derivative of $\mathcal{J}(a, b, c)$ w.r.t. a or b or c can be determined

# Computational Graph

Given a function $\mathcal{J}(a, b, c) = 3(a + bc)$

$a = 5$

$b = 3 \longrightarrow$ | $u = b \times c$ | $\longrightarrow$ | $v = a + u$ | $\longrightarrow$ | $\mathcal{J} = 3v$ |

$c = 2$

$\dfrac{d\mathcal{J}}{dv}$

- Finding $\dfrac{d\mathcal{J}}{dv}$

  - $v = 11 \xrightarrow{\text{Nudge its value}} v = 11.001$

  - $\mathcal{J} = 33 \xrightarrow{\text{Changed by } v} \mathcal{J} = 33.003$
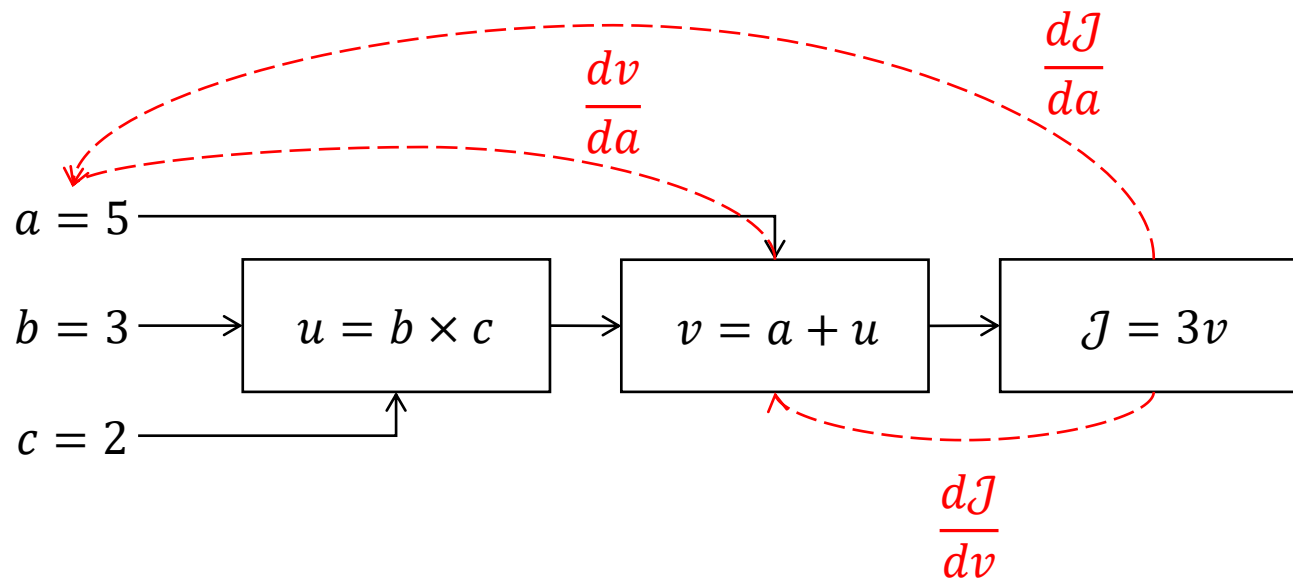
Therefore

$$\frac{d\mathcal{J}}{dv} = \frac{0.003}{0.001} = 3$$

# Computational Graph

*Backward path*

Given a function $J(a, b, c) = 3(a + bc)$



- Finding $\dfrac{dJ}{da}$

  - $a = 5 \xrightarrow{\text{Nudge its value}} a = 5.001$

  - $v = 11 \xrightarrow{\text{Changed by } a} v = 11.001$

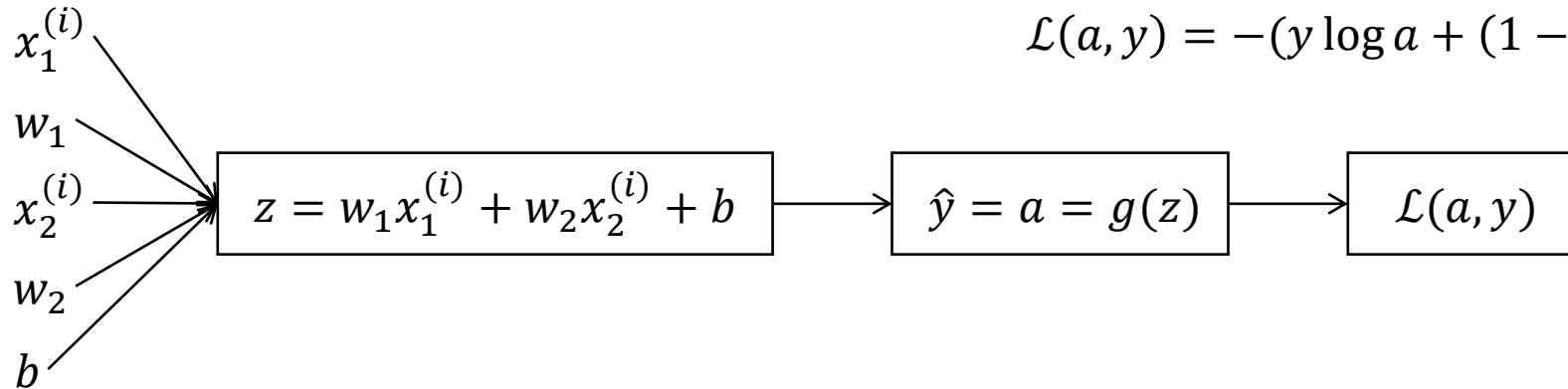  - $J = 33 \xrightarrow{\text{Changed by } a, v} J = 33.003$

Therefore

$$\frac{dJ}{da} = \frac{dJ}{dv} \cdot \frac{dv}{da} = \frac{0.003}{0.001} \cdot \frac{0.001}{0.001} = 3$$

Chain rule

# Perceptron with Gradient Descent

*For a training sample*

Given $\mathbf{w} = [w_1 \quad w_2]$ and $b$



- Recap

$$z = \mathbf{w}^T \cdot \mathbf{x}^{(i)} + b$$
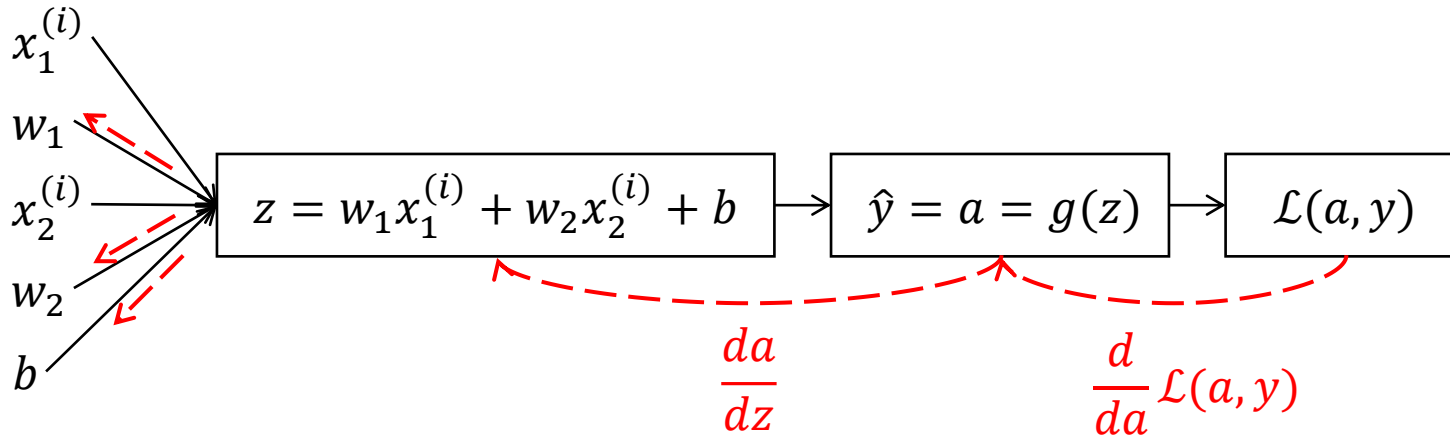
$$\hat{y} = a = g(z)$$

$$g(z) = \frac{1}{1 + e^{-z}}$$

$$\mathcal{L}(a, y) = -(y \log a + (1 - y) \log(1 - a))$$

- Goal
  - To adjust $\mathbf{w} = [w_1 \quad w_2]$ and $b$ to minimize $\mathcal{L}(a, y)$

# Perceptron with Gradient Descent

$$\Delta z = \frac{d}{dz}\mathcal{L}(a, y)$$

$$= \frac{d}{da}\mathcal{L}(a, y) \cdot \frac{da}{dz}$$

since $\dfrac{d}{da}\mathcal{L}(a, y) = \dfrac{d}{da}[-y \log a - (1 - y) \log(1 - a)]$

$$= -\frac{y}{a} + \frac{1 - y}{1 - a}$$

and $\dfrac{da}{dz} = a(1 - a)$

$$\therefore \Delta z = \left(-\frac{y}{a} + \frac{1 - y}{1 - a}\right) \cdot a(1 - a)$$

$$= a - y$$

$$z = w_1 x_1^{(i)} + w_2 x_2^{(i)} + b \rightarrow \hat{y} = a = g(z) \rightarrow \mathcal{L}(a, y)$$

$x_1^{(i)}$

$w_1$

$x_2^{(i)}$

$w_2$

$b$

$\dfrac{da}{dz}$

$\dfrac{d}{da}\mathcal{L}(a, y)$

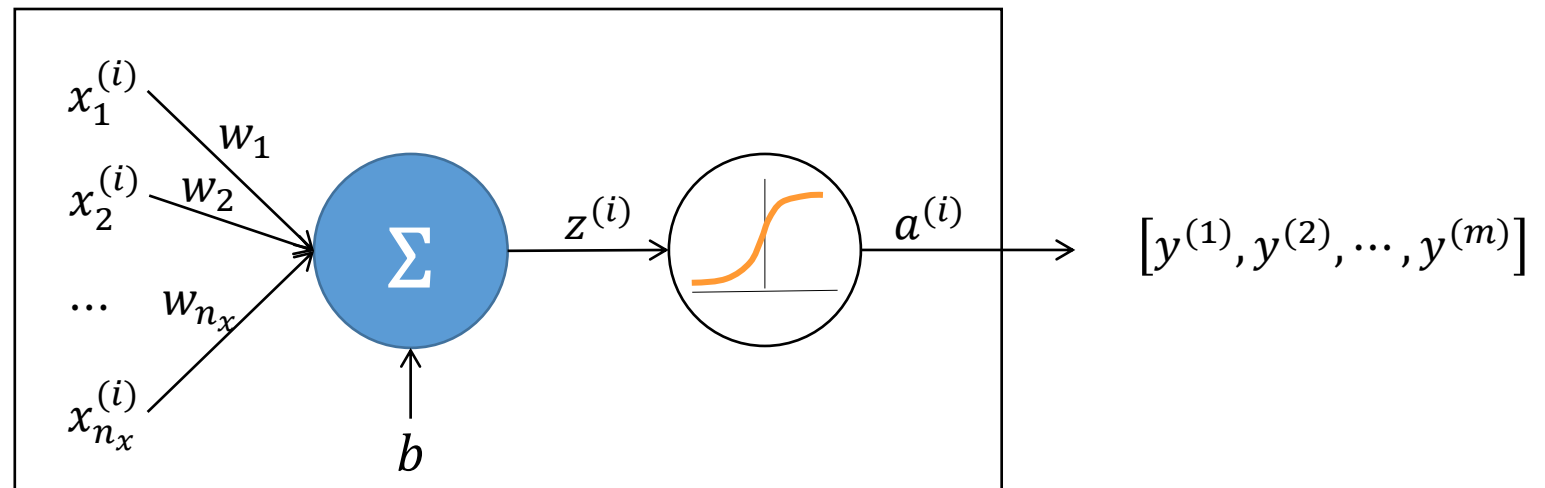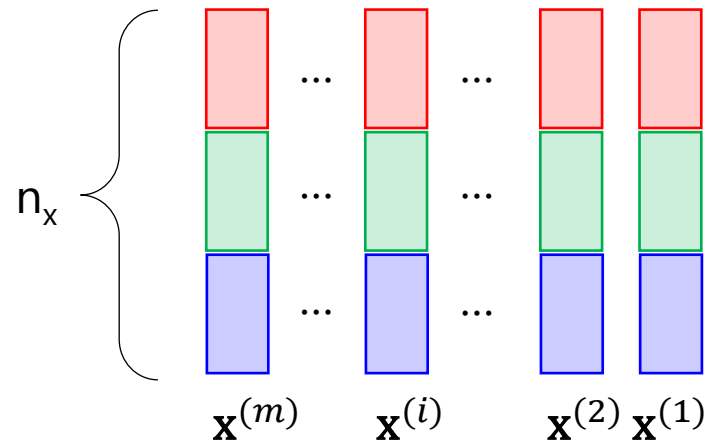In other words, $\Delta z$ is the difference between obtained output and desired output

# Implementation

# Vectorization

Broadcasting property enable implementation without looping over training data

$$\mathbf{y} = g(\mathbf{z}) = g(\mathbf{w}^T \cdot \boldsymbol{X} + b)$$

# Naïve Implementation

Loop and Loop and Loop

This is NOT an efficient method

Note: Python programming does not have the $\Delta$ character. For simplicity, let's change

$$\Delta w_1 \rightarrow dw_1$$
$$\Delta w_2 \rightarrow dw_2$$
$$\Delta b \rightarrow db$$

Initialize

$$J = 0, dw_1 = 0, dw_2 = 0, db = 0$$

for epoch = 1 to max_epoch:

    for i = 1 to m:

        compute $z^{(i)}, a^{(i)}$

$$J \mathrel{+}= -\left[y^{(i)} \log a^{(i)} + \left(1 - y^{(i)}\right) \log\left(1 - a^{(i)}\right)\right]$$

        compute $dz^{(i)}$

        update $db$

        for j = 1 to $n_x$:

            update $dw_j$

$$J /= \frac{J}{m}$$

# Implement with Vectorization

*The efficient method*

Initialize

$$\mathcal{J} = 0, dw_1 = 0, dw_2 = 0, db = 0$$

for epoch = 1 to max_epoch:

$$\mathbf{Z} = \mathbf{w}^T \cdot \mathbf{X} + b$$

$$\mathbf{A} = g(\mathbf{Z})$$

$$\mathcal{J} = -\frac{1}{m}\sum(\mathbf{y}\log\mathbf{A} - (1-\mathbf{y})\log(1-\mathbf{A}))$$

$$\mathbf{dz} = \mathbf{A} - \mathbf{Y}$$

$$\mathbf{dw} = \frac{1}{m}\mathbf{X} \cdot \mathbf{dz}^T$$

$$db = \frac{1}{m} \cdot np.\,sum(\mathbf{dz})$$

$$\mathbf{w} -= \alpha \cdot \mathbf{dw}$$

$$b -= \alpha \cdot db$$

# Implement with Vectorization

*The efficient method*

Only loop over training iteration

for epoch = 1 to max_epoch:

$$\mathbf{Z}_{1 \times m} = \underbrace{\left(\mathbf{w}_{1 \times n_x}\right)^T \cdot \mathbf{X}_{n_x \times m}}_{1 \times m} + b$$

$$\mathbf{A}_{1 \times m} = g(\mathbf{Z}_{1 \times m})$$

$$\mathcal{J} = -\frac{1}{m} \sum_{i=1}^{m} (\underbrace{\mathbf{y}_{1 \times m} \log \mathbf{A}_{1 \times m}}_{1 \times m} - \underbrace{(1 - \mathbf{y}_{1 \times m}) \log(1 - \mathbf{A}_{1 \times m})}_{1 \times m})$$

(underbraced to $1 \times 1$)

$$\mathbf{dz}_{1 \times m} = \mathbf{A}_{1 \times m} - \mathbf{y}_{1 \times m}$$

$$\mathbf{dw}_{1 \times n_x} = \frac{1}{m} \mathbf{X}_{n_x \times m} \cdot (\mathbf{dz}_{1 \times m})^T$$

$$db = \frac{1}{m} \sum_{i=1}^{m} \mathbf{dz}_{1 \times m}$$

$$\mathbf{w}_{1 \times n_x} \mathrel{-}= \alpha \cdot \mathbf{dw}_{1 \times n_x}$$

$$b \mathrel{-}= \alpha \cdot db$$

*Wisdom of the Land*

# Shallow Neural Network

*None or a few hidden layers*



Hidden layer

Input layer

Output layer

$x_1$

$x_2$

$x_3$

$\hat{y} = a$

1-Layer Network
(0 Hidden Layer)
(perceptron or logistic regression)

$x_1$

$x_2$

$x_3$

$\hat{y} = a$

2-Layer Network
(1 Hidden Layer)

$x_1$

$x_2$

$x_3$

$\hat{y} = a$

3-Layer Network
(2 Hidden Layer)

# Deep Neural Network

*More hidden layers*



4-Layer Network
(3 Hidden Layer)

- Notations

$L$ represents the number of layers (exclude the input layer)

$n^{[l]}$ represents the number of units in the $l^{th}$ layer

Examples:

$$n^{[1]} = 5, \qquad n^{[2]} = 4, \qquad n^{[3]} = 3, \qquad n^{[4]} = 1$$

$g^{[l]}(\cdot)$ represents the activation function in the $l^{th}$ layer

$\mathbf{z}^{[l]}$ represents the linear combination in the $l^{th}$ layer

$\mathbf{a}^{[l]}$ represents the activated values in the $l^{th}$ layer

$\mathbf{W}^{[l]}$ represents the weights matrix in the $l^{th}$ layer

$\mathbf{b}^{[l]}$ represents the bias in the $l^{th}$ layer

$\mathbb{w}$ represents the collection of weights matrix of all layers

$w_{i,j}$ represents the weights connecting between the $i^{th}$ unit in the $l^{th}$ layer to the $j^{th}$ unit in the $(l-1)^{th}$ layer

$\mathbb{b}$ represents the collection of bias vectors of all layers

$\hat{y} = a$

Wisdom of the Land

# Deep Neural Network

*Weights*



$$\hat{y} = a$$

- Weights $\quad \mathbb{W} = \begin{bmatrix} \mathbf{W}^{[1]} \\ \mathbf{W}^{[2]} \\ \mathbf{W}^{[3]} \\ \mathbf{W}^{[4]} \end{bmatrix}$

$$\mathbf{W}^{[l]} = \begin{bmatrix} \mathbf{w}_1^{[l]} \\ \cdots \\ \mathbf{w}_{n^{[l]}}^{[l]} \end{bmatrix}_{n^{[l]}, n^{[l-1]}}$$
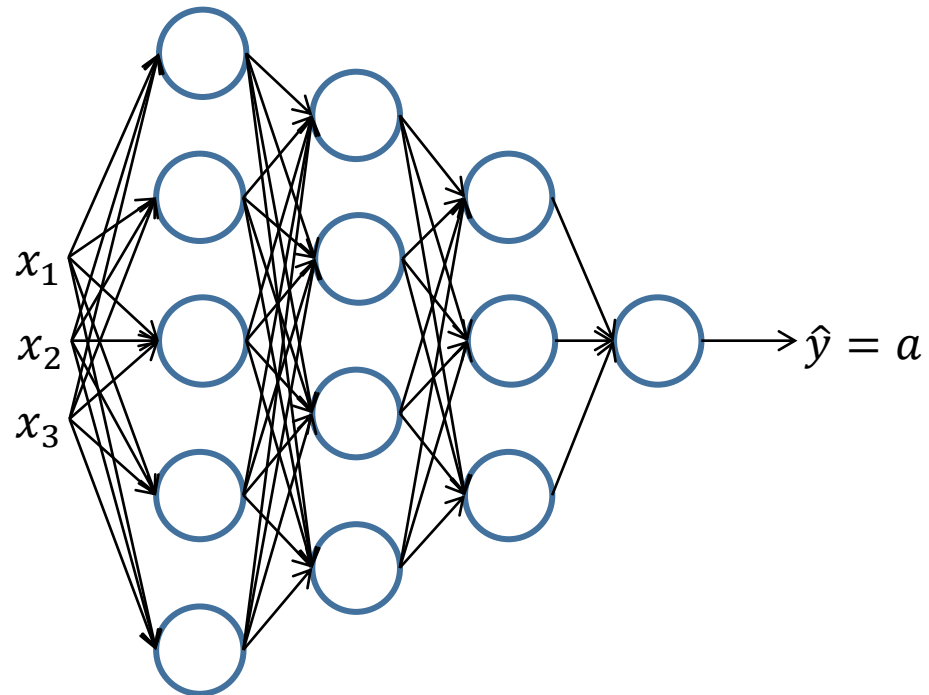
- Example

$$\mathbf{W}^{[3]} = \begin{bmatrix} w_{1,1} & \cdots & w_{1,4} \\ \vdots & \ddots & \vdots \\ w_{3,1} & \cdots & w_{3,4} \end{bmatrix}$$

Dimension of $\mathbf{W}^{[l]}$ is $\left(n^{[l]}, n^{[l-1]}\right)$

# Deep Neural Network

*Biases*



$x_1$

$x_2$

$x_3$

$\hat{y} = a$

- Biases

$$\mathbb{b} = [\mathbf{b}^{[1]} \quad \dots \quad \mathbf{b}^{[4]}] = \begin{bmatrix} b_1^{[1]} & \dots & b_1^{[4]} \\ \vdots & \dots & \vdots \\ b_5^{[1]} & \dots & b_3^{[4]} \end{bmatrix}$$

$$\therefore \mathbf{b}^{[l]} = \begin{bmatrix} b_1^{[l]} \\ \vdots \\ b_{n^{[l]}}^{[l]} \end{bmatrix}$$
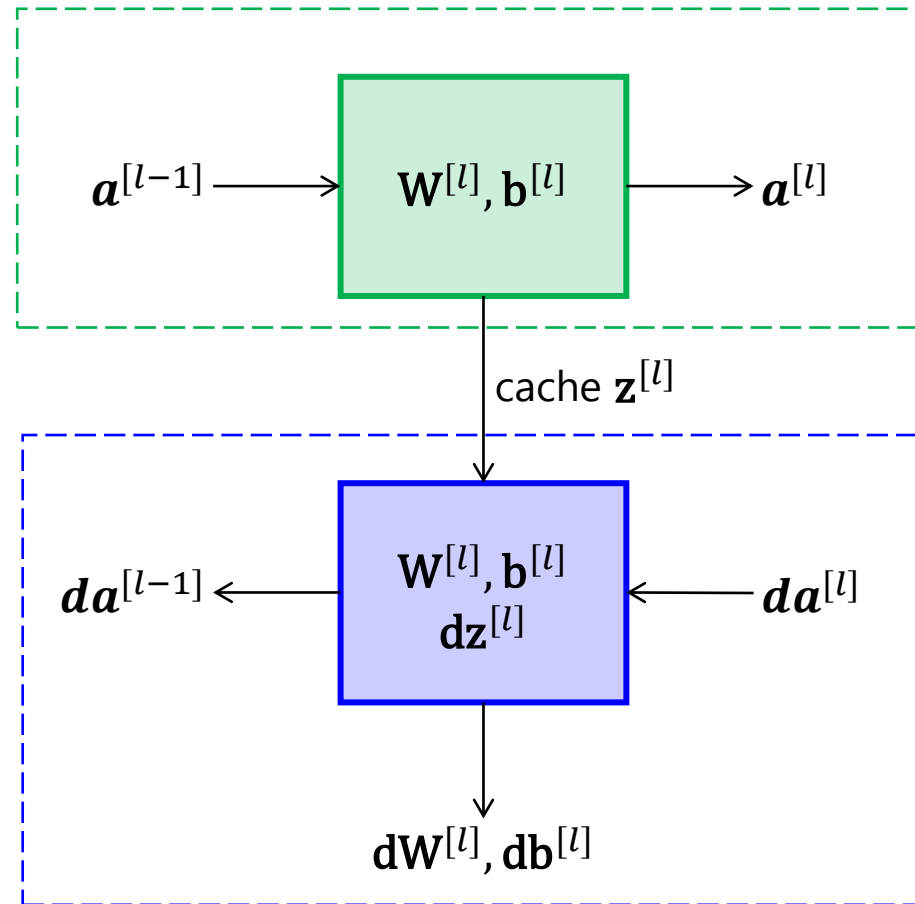
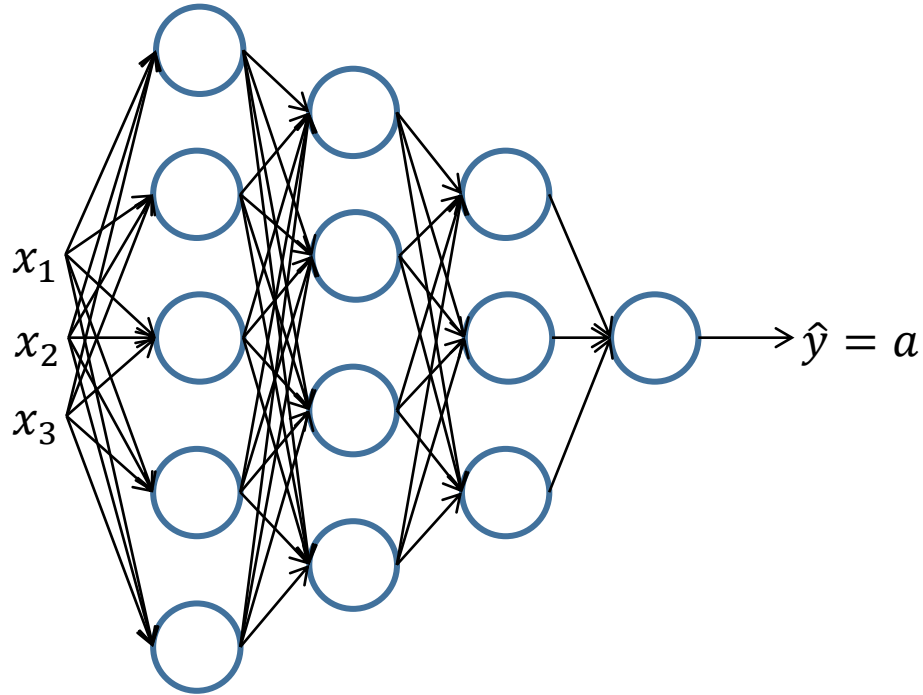Dimension of $\mathbf{b}^{[l]}$ is $(n^{[l]}, 1)$

# Implementation with Vectorization

*The L layers neural network*



- At the layer $l$: $\mathbf{W}^{[l]}, \mathbf{b}^{[l]}$

# Implementation with Vectorization

*The L layers neural network*



$x_1$
$x_2$
$x_3$
$\hat{y} = a$

- At the layer $l : \mathbf{W}^{[l]}, \mathbf{b}^{[l]}$
  - Forward Propagation

    $$\mathbf{z}^{[l]} = \mathbf{W}^{[l]} \cdot \boldsymbol{a}^{[l-1]} + \mathbf{b}^{[l]}$$

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

  - Backward Propagation

    $$\mathbf{dz}^{[l]} = \mathbf{da}^{[l]} * g'^{[l]}(\mathbf{z}^{[l]})$$

    $$\mathbf{dW}^{[l]} = \frac{1}{m} \mathbf{dz}^{[l]} \cdot \boldsymbol{a}^{[l-1]^T}$$

    $$\mathbf{db}^{[l]} = \frac{1}{m} \sum_{i=1}^{m} \mathbf{dz}^{[l]}$$

    $$np.\,sum(\mathbf{dz}^{[l]}, axis = 1, keepdims = True)$$

    $$\mathbf{da}^{[l-1]} = \mathbf{W}^{[l]^T} \cdot \mathbf{dz}^{[l]}$$

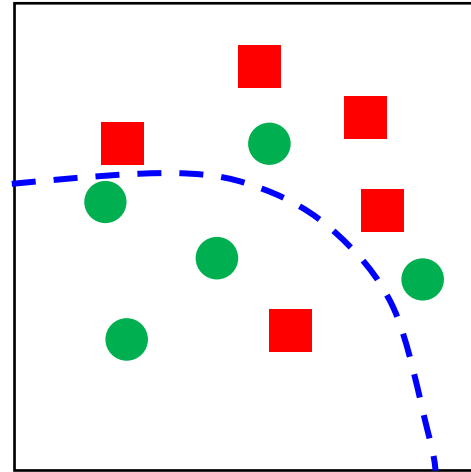*Wisdom of the Land*
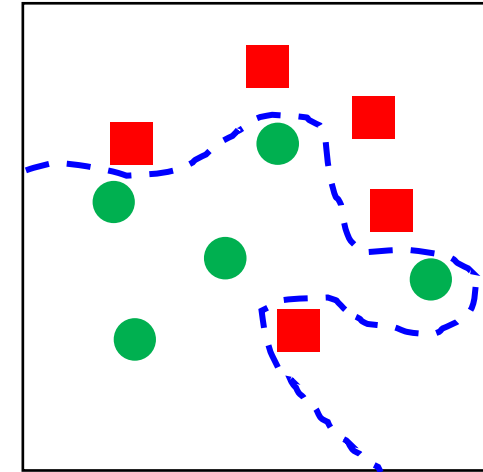
Model Tuning

# Bias and Variance Problem

*Easy to learn but difficult to master*

- Unable to visualize in high dimensional data



High Bias
(Underfitting)

Just right

High Variance
(Overfitting)

# Regularization

L1 and L2

# L1 and L2 Regularization

Recall the cost function of perceptron

Goal:

$$\min_{\mathbf{w},b} \mathcal{J}(\mathbf{w}, b), \qquad \mathbf{w} \in \mathbb{R}^{n_x}, b \in \mathbb{R}$$

- L1 regularization (Lasso)

$$\mathcal{J}(\mathbf{w}, b) = \frac{1}{m} \sum_{i=1}^{m} \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) + \boxed{\frac{\lambda}{2m} \|\mathbf{w}\|_1} \quad \textit{Regularization term}$$

where $\quad \|\mathbf{w}\|_1 = \sum_{j=1}^{n_x} |w_j|$

$\lambda$ represents the regularization parameter

- L1 makes many weights become zeros
- Good for compacting the model
- L1 is not often used

# L1 and L2 Regularization

Recall the cost function of perceptron

Goal:

$$\min_{\mathbf{w},b} \mathcal{J}(\mathbf{w},b), \qquad \mathbf{w} \in \mathbb{R}^{n_x}, b \in \mathbb{R}$$

- L2 regularization (Ridge)

$$\mathcal{J}(\mathbf{w},b) = \frac{1}{m}\sum_{i=1}^{m} \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m}\|\mathbf{w}\|_2^2 \qquad \textit{Regularization term}$$

where $\quad \|\mathbf{w}\|_2^2 = \sum_{j=1}^{n_x} w_j^2 = \mathbf{w}^T \cdot \mathbf{w}$

$\lambda$ represents the regularization parameter

- L2 is much more often used compared to L1

# L2 Regularization

*L2 of Deep Neural Network*

The cost function

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

where

$$\text{Frobenious norm} \quad \|\mathbf{w}\|_F^2 = \sum_{i=1}^{n^{[l-1]}}\sum_{j=1}^{n^{[l]}}\left(w_{ij}^{[l]}\right)^2$$

$$\mathbf{w}^{[l]}: (n^{[l]}, n^{[l-1]})$$

• Modifying back propagation

$$\mathbf{dw}^{[l]} = (\cdot) + \frac{\lambda}{m}\mathbf{w}^{[l]}$$

$$\mathbf{w}^{[l]} = \mathbf{w}^{[l]} - \alpha\mathbf{dw}^{[l]}$$

$$= \mathbf{w}^{[l]} - \alpha\left[(\cdot) + \frac{\lambda}{m}\mathbf{w}^{[l]}\right]$$

$$= \mathbf{w}^{[l]} - \frac{\alpha\lambda}{m}\mathbf{w}^{[l]} - \alpha(\cdot)$$

$$= (1 - \frac{\alpha\lambda}{m})\mathbf{w}^{[l]} - \alpha(\cdot)$$

*Weight Decay*

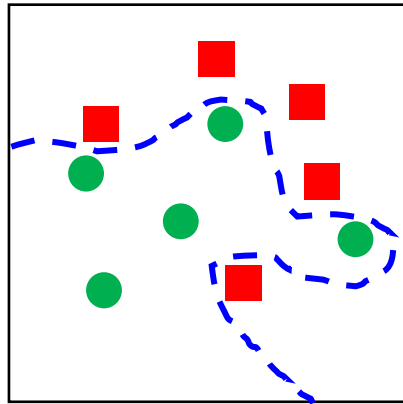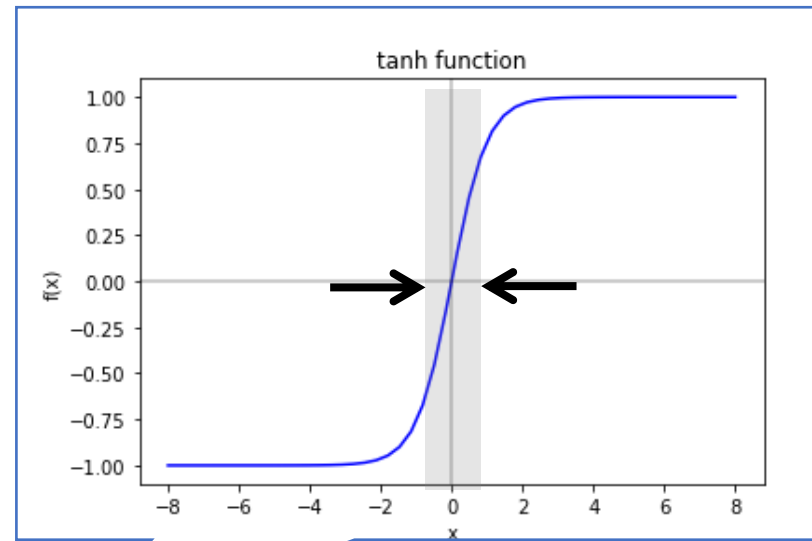where $(\cdot)$ represents the term obtained from original back propagation

*Wisdom of the Land*

# L2 Regularization

*Why does L2 prevent overfitting?*



High Bias
(Underfitting)

Just right
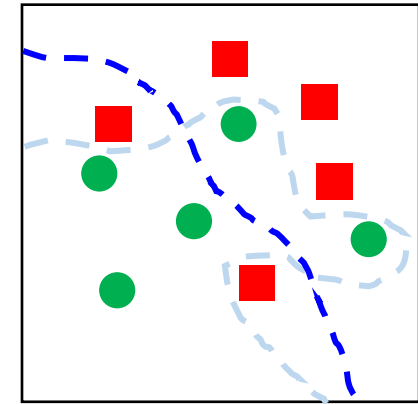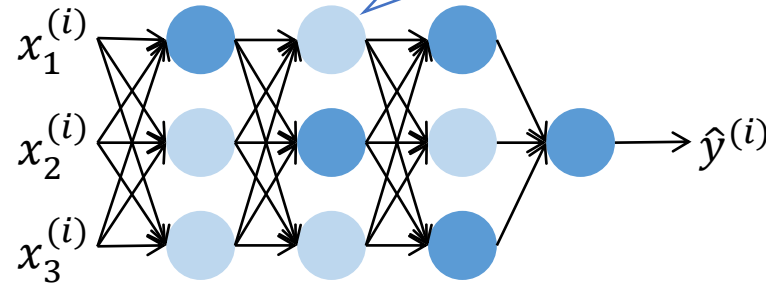
High Variance
(Overfitting)

# L2 Regularization

*Why does L2 prevent overfitting?*



$$\lambda \uparrow$$

$$\mathbf{w}^{[l]} \downarrow$$

$$\mathbf{z}^{[l]} \downarrow$$

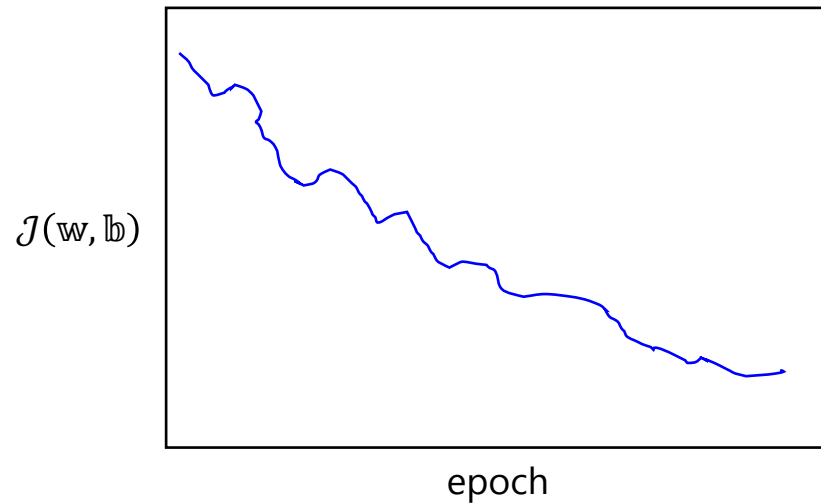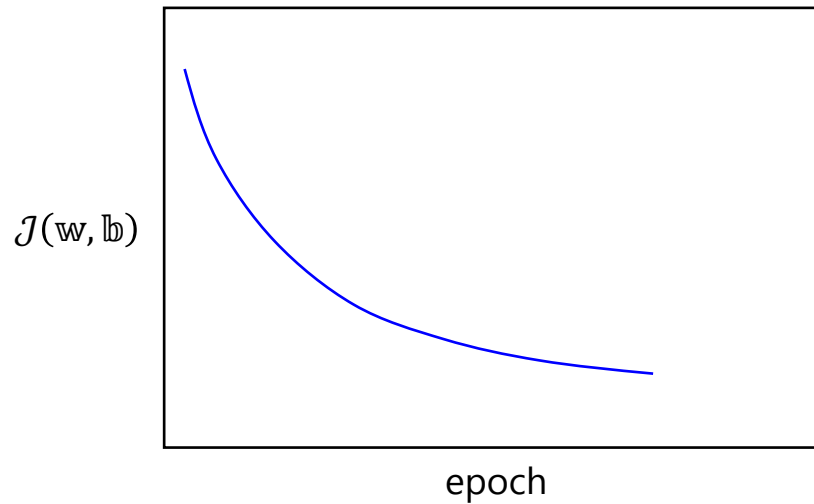$$\because \mathbf{z}^{[l]} = \mathbf{w}^{[l]} \boldsymbol{a}^{[l-1]} + \mathbf{b}^{[l]}$$

Decision boundary is stretched out

Wisdom of the Land

# L2 Regularization

*Smoother cost function*

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

**Without regularization**

$\mathcal{J}(\mathbb{w}, \mathbb{b})$

epoch

**With regularization**

$\mathcal{J}(\mathbb{w}, \mathbb{b})$

epoch

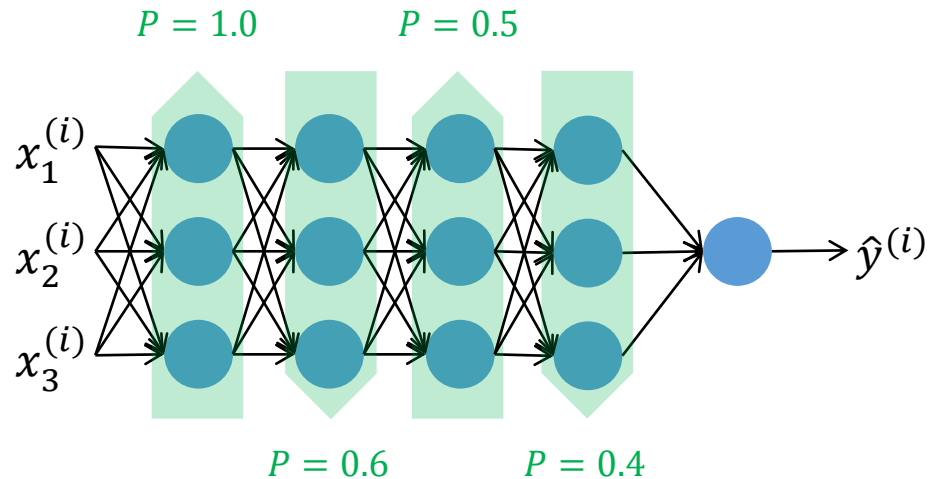*Wisdom of the Land*

# Regularization

## DropOut & DropConnect
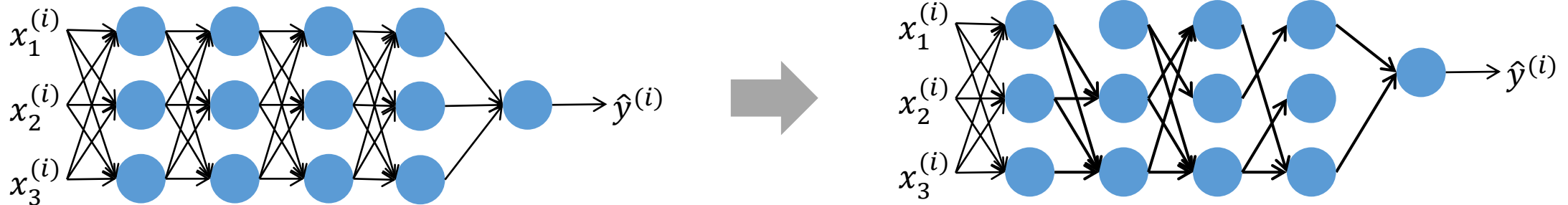
# DropOut

*Ensemble Neural Network*

- Set the probability of enable hidden units in each hidden layer
- The edges connected to any disable hidden unit will be removed
- DropOut is applied ONLY training but NOT testing
- DropOut produces many possible combination of neuron network
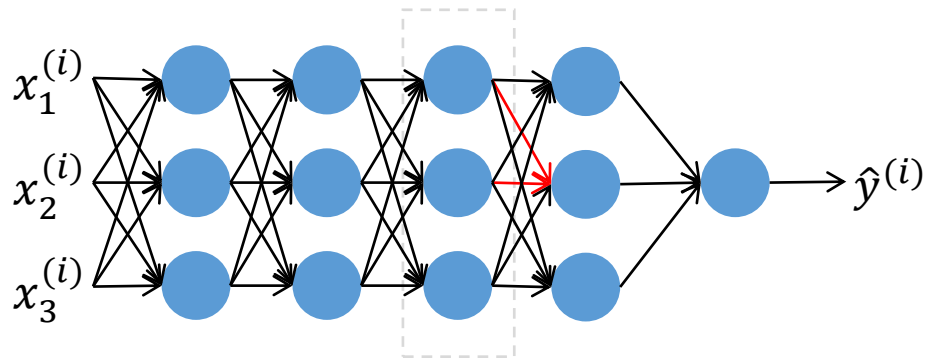- DropOut is very popular in computer vision

# DropConnect

*Generalization of DropOut*

- Set the Boolean mask to randomly disable connections
- Rescale output on active connections
- DropConnect is applied ONLY for training but NOT testing
- DropConnect makes even more possible combination than dropout
- DropConnect is very popular in computer vision



Wisdom of the Land

# **Implementation**
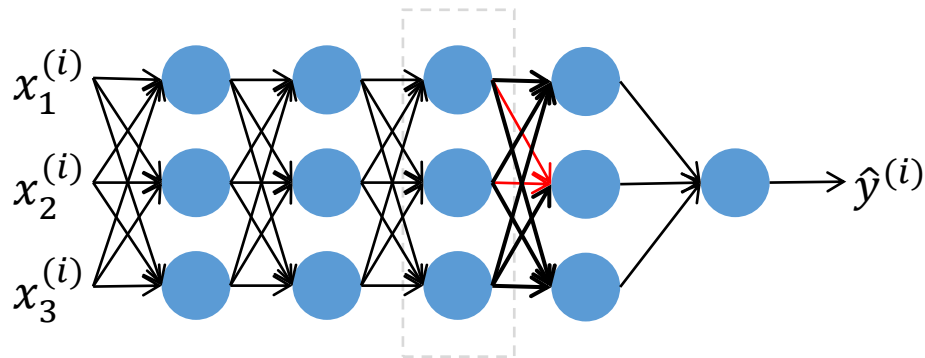


*Illustration of DropConnect at the layer 3*

- Define theoretical variable to Python variable
  - P as keepprob

- Create a Boolean mask to randomly disable weights

```
keepprob = 0.8
d3 = np.random.rand(a3.shape[0], a3.shape[1]) < keepprob


In [43]: d3
Out[43]:
array([[ True,  True,  True],
       [False,  True,  True],
       [ True, False,  True]])
```

# Implementation

*Illustration of DropConnect at the layer 3*



- Rescale the $a^{[3]}$

```
a3 /= keepprob
```

- Example: If $a^{[3]}$ is

```
array([[1., 2., 3.],
       [4., 5., 6.],
       [7., 8., 9.]])
```
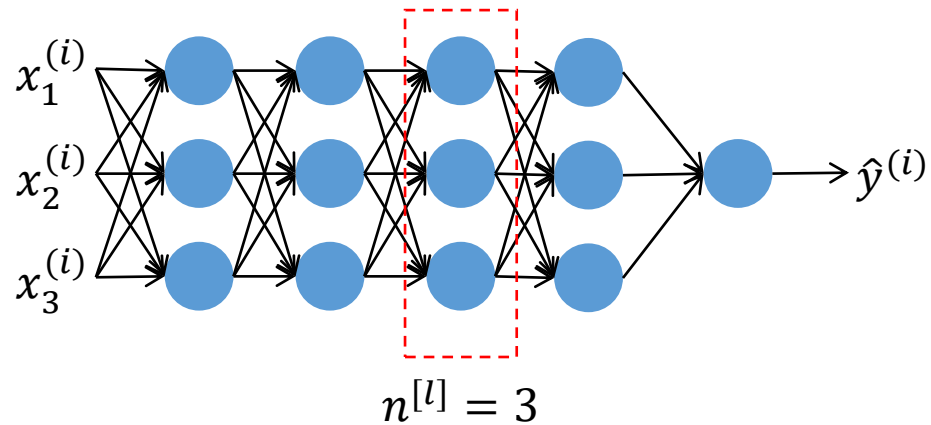
- Rescale $a^{[3]}$ after drop connect will be

```
array([[ 1.25,  2.5 ,  3.75],
       [ 0.  ,  6.25,  7.5 ],
       [ 8.75,  0.  , 11.25]])
```

- In practice a3 is obtained from

$$a^{[3]} = g^{[3]}(z^{[3]}) = g^{[3]}(W^{[3]} \cdot a^{[2]} + b^{[3]})$$

*Wisdom of the Land*

# Recommended Settings

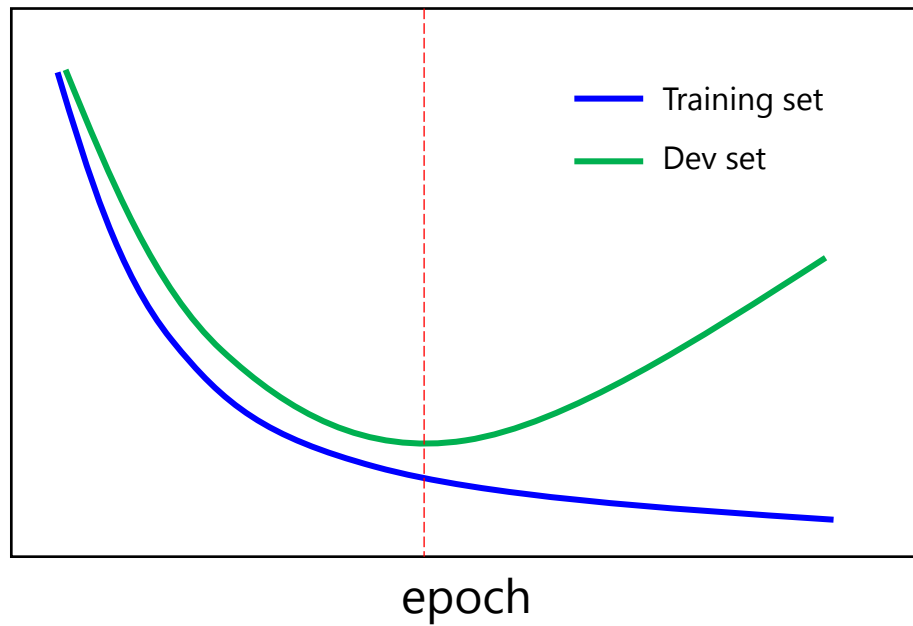*Some idea to set your keepprob*



$$n^{[l]} = 3$$

- The keepprob can be different value in different layer
- The lower $n^{[l]}$ (number of hidden units in the layer $l$ ) the higher keepprob
- Input layer has keepprob equal to 1

# Other Regularizations

Early stopping

$\mathcal{J}(\mathbb{W}, \mathbb{b})$



epoch

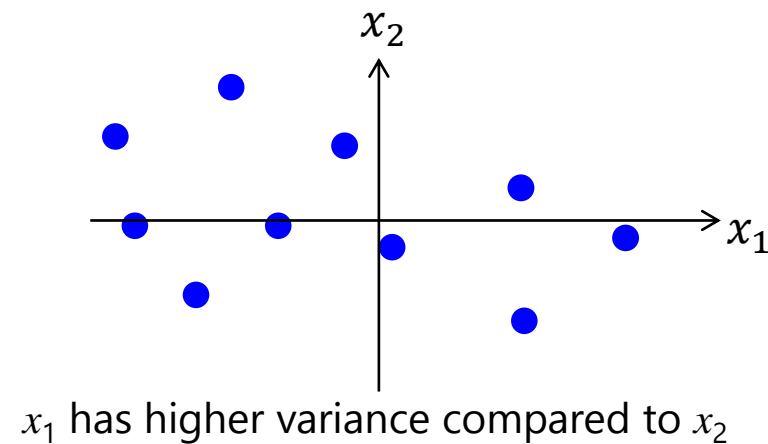- Stop training based on the cost convergence
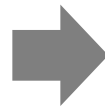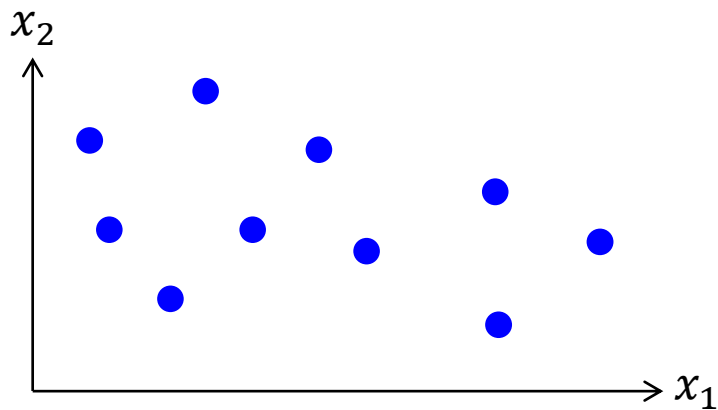
# Normalized Input

*Speed up your training*

- Subtract mean
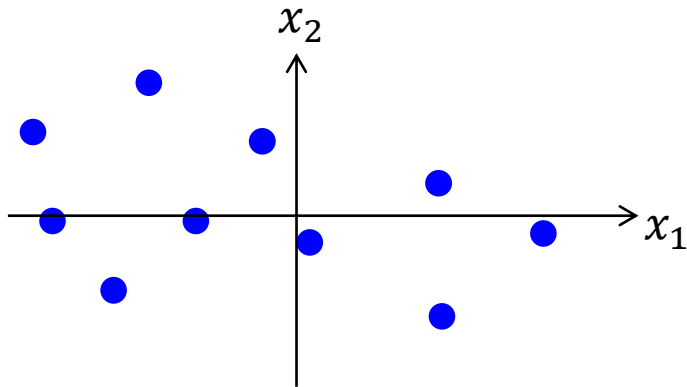
$$\mu = \frac{1}{m}\sum_{i=1}^{m} \mathbf{x}^{(i)}$$

$$\mathbf{X} := \mathbf{X} - \mu$$

- Data will have zero mean



$x_1$ has higher variance compared to $x_2$

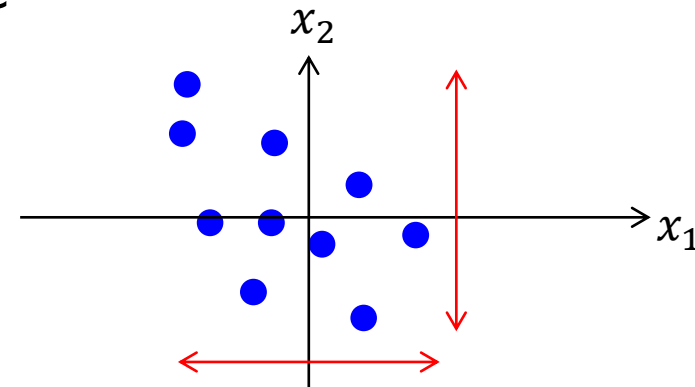# Normalized Input

*Speed up your training*



- Normalized variance

$$\sigma^2 = \frac{1}{m}\sum_{i=1}^{m}\left(\mathbf{x}^{(i)}\right)^2$$

$$\mathbf{X} = \frac{\mathbf{X}}{\sigma^2}$$
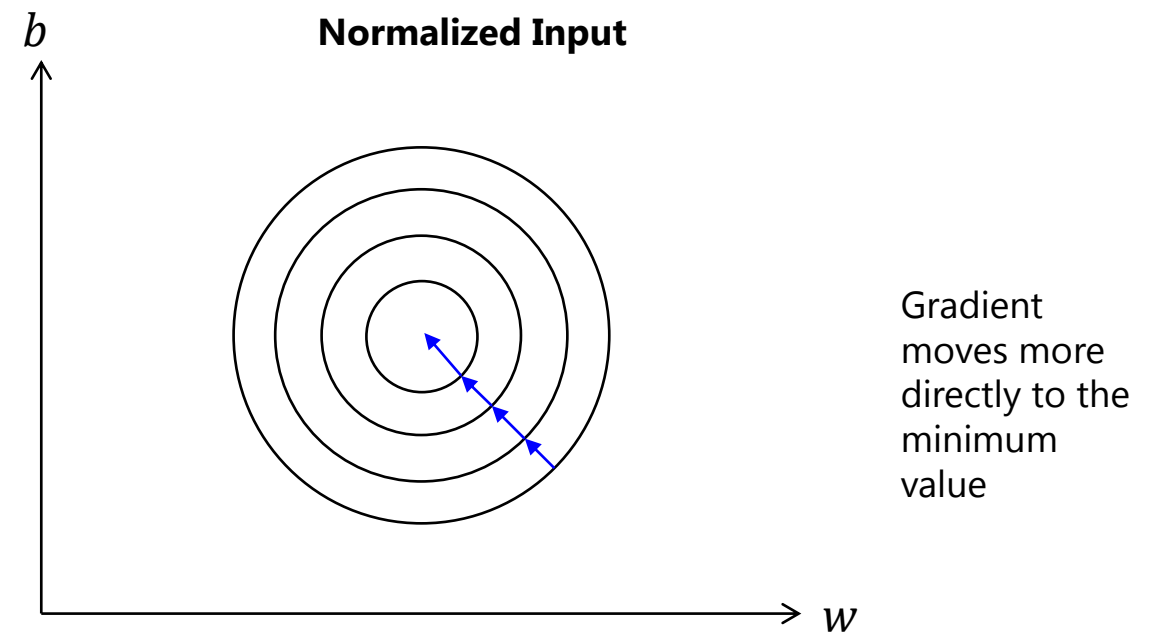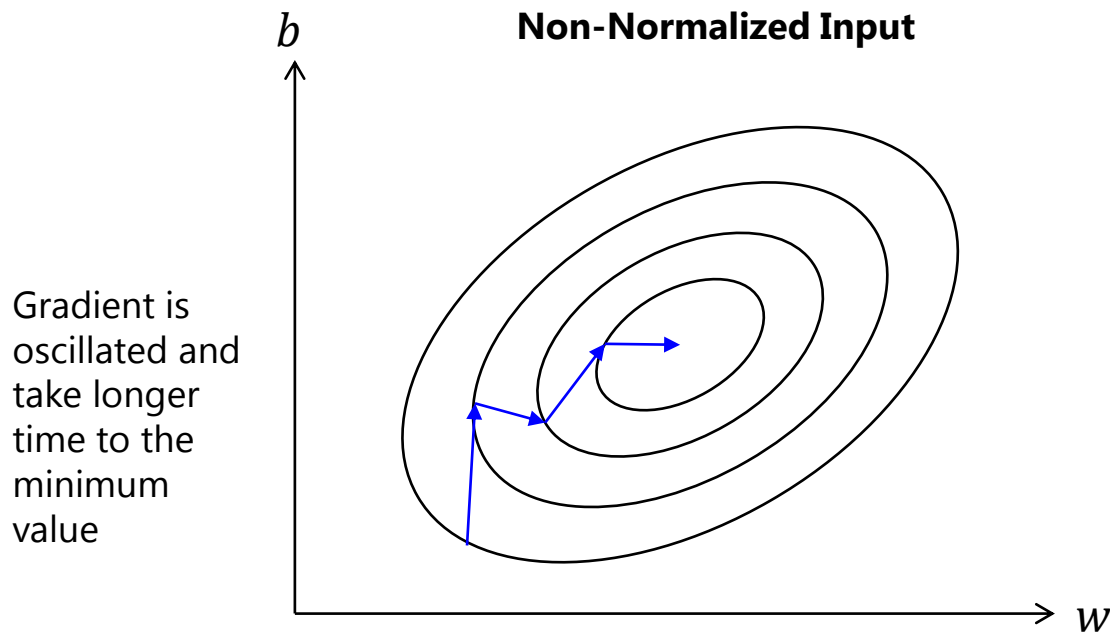
- Data will have zero mean and normalized variance



- $\mu$ and $\sigma^2$ from training set will be used in test set

# Normalized Input

- Given a contour of the cost function

*Speed up your training*



**Non-Normalized Input**

$b$

$w$

Gradient is oscillated and take longer time to the minimum value

**Normalized Input**

$b$

$w$

Gradient moves more directly to the minimum value

# Batch VS Mini-batch

*Dividing training data into minibatch*

Notations

$\mathbf{X}^{\{t\}}$ and $\mathbf{y}^{\{t\}}$ for training samples of mini-batch t

- Vectorization enables effective computation over $m$ samples

$$\mathbf{X}_{n_x \times m} = \left[\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(m)}\right]$$

$$\mathbf{y}_{1 \times m} = \left[y^{(1)}, y^{(2)}, \dots, y^{(m)}\right]$$

- What if m is extremely large, e.g. 5,000,000
  - Use mini-batch

$$\mathbf{X}_{n_x \times m} = \left[\underbrace{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(1000)}}_{\mathbf{X}^{\{1\}}} | \underbrace{\mathbf{x}^{(1001)}, \mathbf{x}^{(1002)}, \dots, \mathbf{x}^{(2000)}}_{\mathbf{X}^{\{2\}}} | \cdots\cdots | \underbrace{\dots\dots \mathbf{x}^{(5,000,000)}}_{\mathbf{X}^{\{5000\}}}\right]$$

$$\mathbf{y}_{1 \times m} = \left[\underbrace{y^{(1)}, y^{(2)}, \dots, y^{(1000)}}_{\mathbf{y}^{\{1\}}} | \underbrace{y^{(1001)}, y^{(1002)}, \dots, y^{(2000)}}_{\mathbf{y}^{\{2\}}} | \cdots\cdots | \underbrace{\dots\dots y^{(5,000,000)}}_{\mathbf{y}^{\{5000\}}}\right]$$

# Mini-Batch Implementation

*From the previous example*

while {
  for t in range(0, 5000):

*# Forward propagation on X^{t}*

$$\mathbf{z}^{[1]} = \mathbf{w}^{[1]} \cdot \mathbf{X}^{\{t\}} + \mathbf{b}^{[1]}$$

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

$$\vdots$$

$$\mathbf{z}^{[L]} = \mathbf{w}^{[L]} \cdot \boldsymbol{a}^{[L-1]} + \mathbf{b}^{[L]}$$

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

*Use vectorization implementation*

*from* $\mathbf{X}^{\{t\}}, \mathbf{y}^{\{t\}}$

$$\mathcal{J}^{\{t\}} = \frac{1}{1000} \sum_{i=1}^{1000} \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2 \times 1000} \sum_{l} \left\| \mathbf{w}^{[l]} \right\|_F^2$$

Backward propagation w.r.t. $\mathcal{J}^{\{t\}}$ using $\mathbf{X}^{\{t\}}$ and $\mathbf{y}^{\{t\}}$

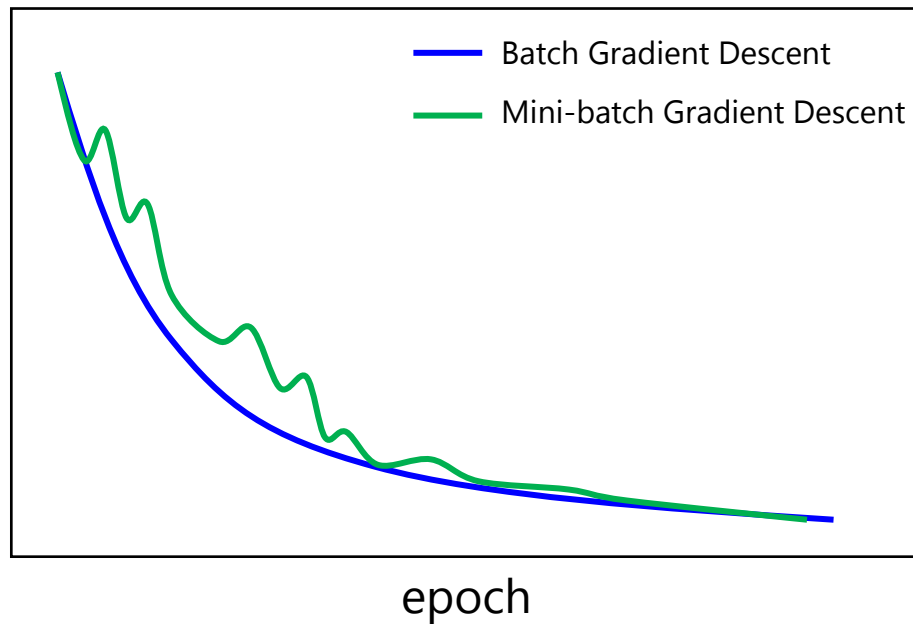$$\mathbf{w}^{[l]} = \mathbf{w}^{[l]} - \alpha \mathbf{dw}^{[l]}$$
$$\mathbf{b}^{[l]} = \mathbf{b}^{[l]} - \alpha \mathbf{db}^{[l]}$$

} until converge

# Mini-batch gradient descent

*Cost function comparison*

$\mathcal{J}(\mathbb{W}, \mathbb{b})$
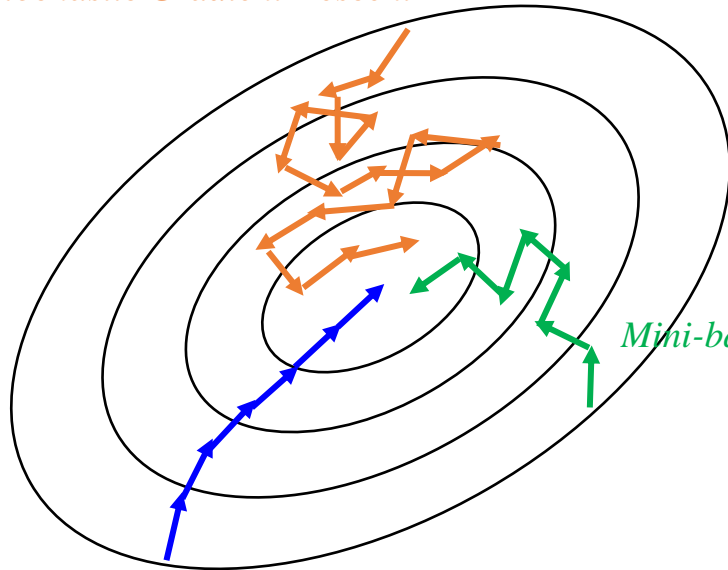


epoch

- The cost function of mini-batch gradient descent could be oscillated
  - $\mathbf{X}^{\{1\}}$, $\mathbf{y}^{\{1\}}$ might be an easy gradient
  - $\mathbf{X}^{\{2\}}$, $\mathbf{y}^{\{2\}}$ might be a harder gradient
  - Overall should be downward

Legend:
— Batch Gradient Descent
— Mini-batch Gradient Descent

# Choosing Mini-batch size

*The extreme cases*

*Stochastic Gradient Descent*

*Mini-batch Gradient Descent*

*Batch Gradient Descent*

- Mini-batch size = $m$
  - Batch gradient descent

$$\left(\mathbf{X}^{\{1\}}, \mathbf{y}^{\{1\}}\right) = (\mathbf{X}, \mathbf{y})$$

- Mini-batch size = 1
  - Stochastic gradient descent

$$\left(\mathbf{X}^{\{1\}}, \mathbf{y}^{\{1\}}\right) = \left(\mathbf{x}^{(1)}, y^{(1)}\right)$$

| Characteristics | Batch | Mini-batch | SGD | |
|---|---|---|---|---|
| Speed | Very Slow | Fast | Slow | *Does not leverage vectorization* |
| Convergence to minimum | Guarantee | Not Guarantee | Not Guarantee | |

# Choosing Mini-batch size

*Conclusion*

- Small Training Set ($m \leq 2000$)
  - Use Batch gradient descent

- Otherwise
  - Use Mini-batch

- Suggested mini-batch size
  - $2^k$ to perfectly fit in CPU/GPU memory, e.g. 64, 128, 256, 512, etc

# Gradient Optimization

Momentum Gradient Descent, RMSprop, and ADAM

Wisdom of the Land

# Momentum Gradient Descent

*Idea*



→ Gradient Descent

→ Momentum Gradient Descent

→ Overshoot (too much learning rate)

- Desired movement directions
  - Fast to the minimum point

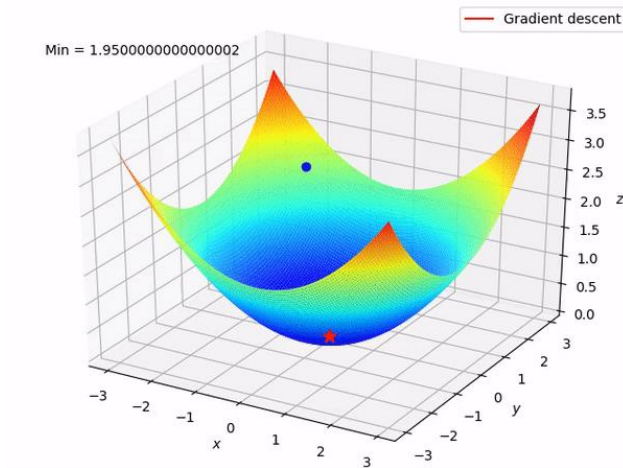    From our example, movement on this direction (——→) should be fast

  - Slow to others

    From our example, movement on this direction ( ↕ ) should be slow

# Momentum Gradient Descent

*Velocity and Acceleration*

Similar to a ball rolling in a bowl



- On iteration t:
  Compute **dW**, **db** on a set of mini-batch

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

*friction* — *velocity*      *acceleration*

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

$$W = W - \alpha v_{dW}$$

$$b = b - \alpha v_{db}$$

# Momentum Gradient Descent

*Implementation*

Initialize $\mathbf{v_{dW}} = 0, \mathbf{v_{db}} = 0$

Note: dimension of $\mathbf{v_{dW}} = \mathbf{dW}$

dimension of $\mathbf{v_{db}} = \mathbf{db}$

Hyper-parameters are $\alpha, \beta$

- On iteration t:

  Compute **dW**, **db** on a set of mini-batch

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

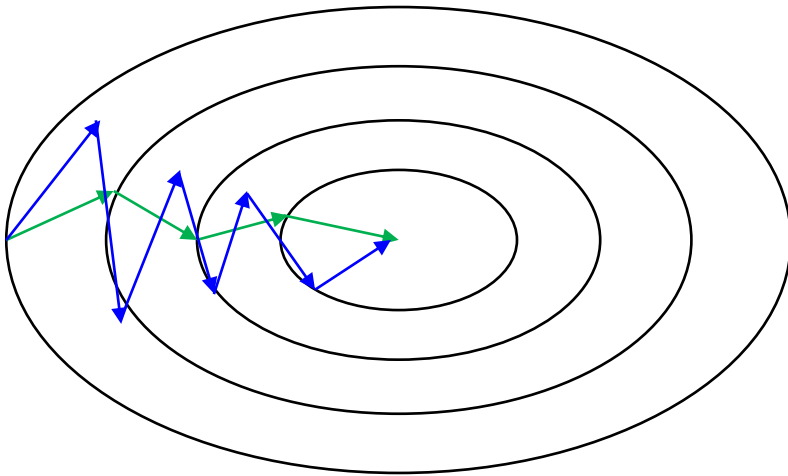  $$\mathbf{v_{db}} = \beta \mathbf{v_{db}} + (1 - \beta)\mathbf{db}$$

  $$\mathbf{W} = \mathbf{W} - \alpha \mathbf{v_{dW}}$$

  $$\mathbf{b} = \mathbf{b} - \alpha \mathbf{v_{db}}$$

- In practice, $\beta = 0.9$
- Some literature omitted $1 - \beta$

# RMSprop

*Speed up gradient descent*



→ Gradient Descent

→ RMSprop

- On iteration t:
  Compute **dW**, **db** on a set of mini-batch

  $$s_{\mathbf{dW}} = \beta s_{\mathbf{dW}} + (1 - \beta)\mathbf{dW}^2$$

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

  $$\mathbf{W} = \mathbf{W} - \frac{\alpha \mathbf{dW}}{\sqrt{s_{\mathbf{dW}}} + \varepsilon}$$
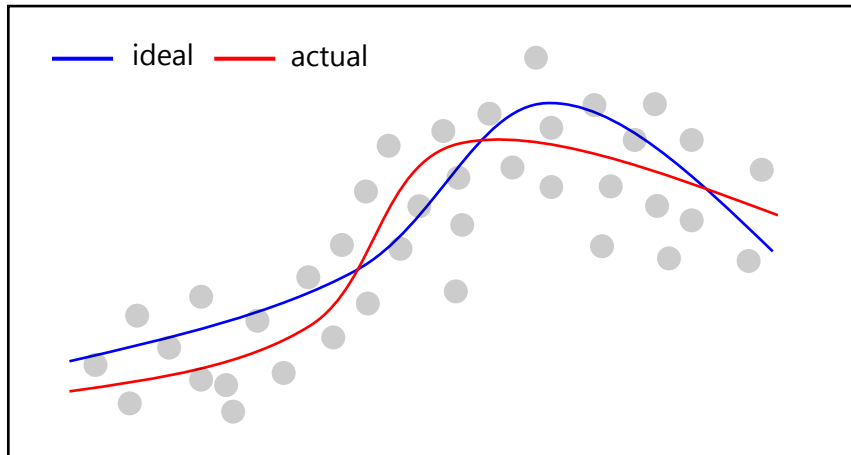
  *a small positive number for numerical stabilization (preventing divided by zero)*

  $$\mathbf{b} = \mathbf{b} - \frac{\alpha \mathbf{db}}{\sqrt{s_{\mathbf{db}}} + \varepsilon}$$

- RMSprop enables to use the larger learning rate without overshoot movement

*Wisdom of the Land*

# Bias Correction

_Idea_



ideal — actual

- Bias correction of data point t is defined as

$$v_t = \beta v_{t-1} + (1-\beta)\theta_t \quad \xrightarrow{\text{data point t}}$$

set $\beta = 0.98,\qquad v_0 = 0$

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

$$v_2 = 0.98 v_1 + 0.02\theta_2 = 0.0196\theta_1 + 0.02\theta_2$$

- Use $\dfrac{v_t}{1-\beta^t}$ to estimate the ideal data point t

- Not accurate at the beginning of t but later will be much more accurate

# Adaptive Moment Estimation (Adam)

*Combining momentum and RMSprop*

Initialize $\mathbf{v_{dW}} = 0, s_{dW} = 0, \mathbf{v_{db}} = 0, s_{db} = 0$

Note: Adam works well in wide range across deep learning applications

- On iteration t:

  Compute **dw**, **db** on a set of mini-batch

$$\mathbf{v_{dW}} = \beta_1 \mathbf{v_{dW}} + (1 - \beta_1)\mathbf{dW}$$

$$\mathbf{v_{db}} = \beta_1 \mathbf{v_{db}} + (1 - \beta_1)\mathbf{db}$$

*momentum*

$$s_{dW} = \beta_2 s_{dW} + (1 - \beta_2)\mathbf{dW}^2$$

$$s_{db} = \beta_2 s_{db} + (1 - \beta_2)\mathbf{db}^2$$

*RMSprop*

$$\mathbf{v_{dW}^{correct}} = \frac{\mathbf{v_{dW}}}{1 - \beta_1^t}, \qquad \mathbf{v_{db}^{correct}} = \frac{\mathbf{v_{db}^{correct}}}{1 - \beta_1^t}$$

$$s_{dW}^{correct} = \frac{s_{dW}}{1 - \beta_2^t}, \qquad s_{db}^{correct} = \frac{s_{db}^{correct}}{1 - \beta_2^t}$$

*bias correction*

$$\mathbf{W}^{[l]} = \mathbf{W}^{[l]} - \alpha \frac{\mathbf{v_{dW}^{correct}}}{\sqrt{s_{dW}^{correct}} + \varepsilon}$$
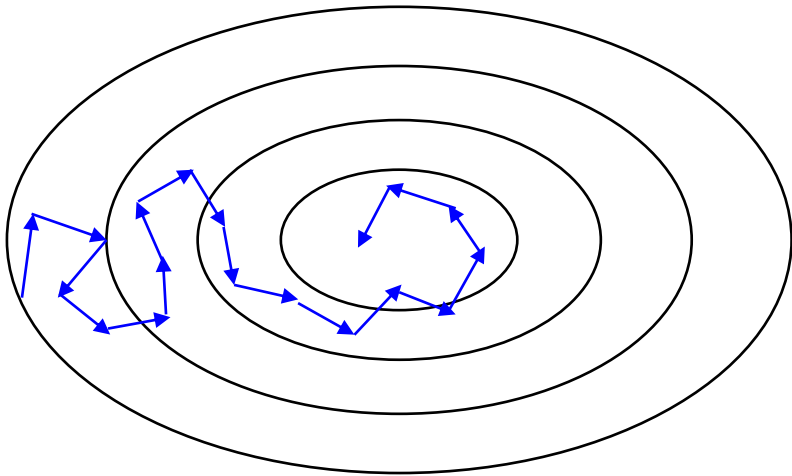
$$\mathbf{b}^{[l]} = \mathbf{b}^{[l]} - \alpha \frac{\mathbf{v_{db}^{correct}}}{\sqrt{s_{db}^{correct}} + \varepsilon}$$

*updating*

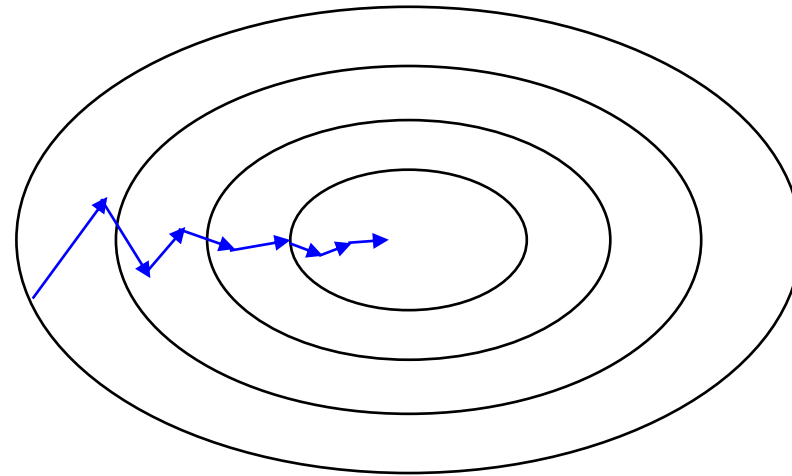# Suggested Hyper-parameters

*Fixed learning rate VS Decayed Learning rate*

- $\alpha$: Need to be tuned
- $\beta_1$: 0.9
- $\beta_2$: 0.999
- $\varepsilon$: $10^{-8}$



**Fixed learning rate**

**Decayed Learning rate**

# Deep learning hyperparameters

*Ranked by importance*

- 1$^{st}$ priority
  - $\alpha$: Learning rate

- 2$^{nd}$ priority
  - $\beta$: Adam's parameters only ($\beta_1, \beta_2, \varepsilon$) (Use default values)
  - Mini-batch size
  - $n^{[l]}$: Number of hidden units in the layer $l$

- 3$^{rd}$ priority
  - $L$: Number of layers
  - $\beta$: Momentum and RMSprop

# Tuning Strategies

*Grid approach*

Hyper-parameter 1

Hyper-parameter 2
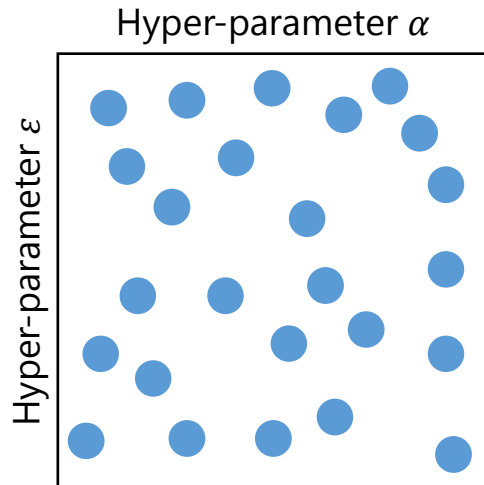


- Grid approach is acceptable for traditional ML

- Limitations
  - Important of some hyperparameters might not be fully address

  - Ex. Grid allows trial on 5 values of $\alpha$ and $\varepsilon$ from 25 experiments

*Wisdom of the Land*

# Tuning Strategies

*For deep learning random grid is better*



Hyper-parameter $\alpha$
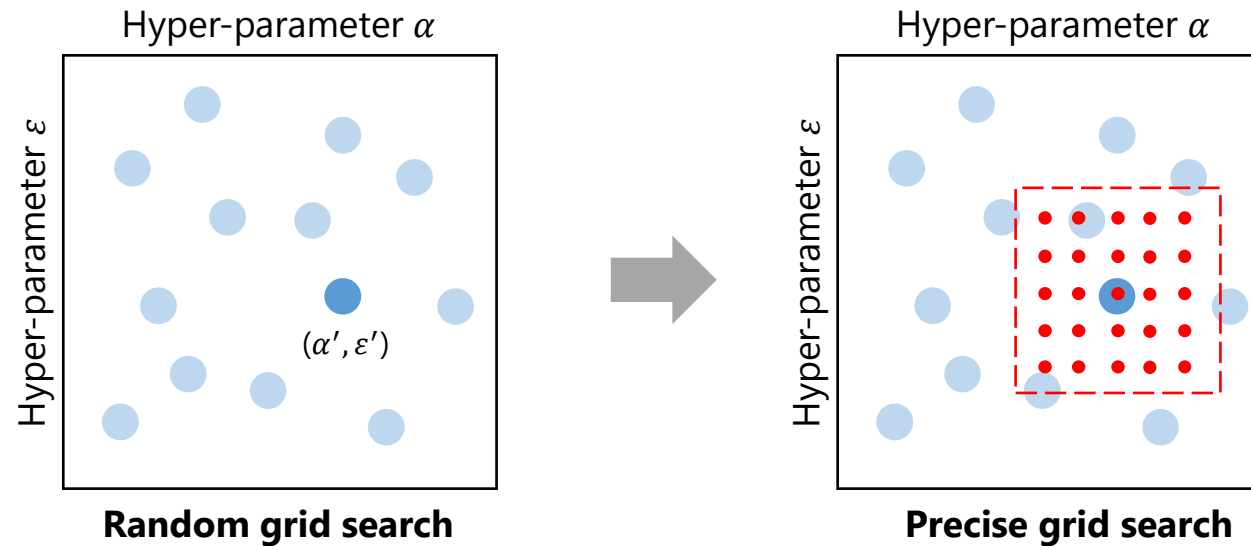
Hyper-parameter $\varepsilon$

- Random grid allows trial on 25 values of $\alpha$ from 25 experiments

- Random grid gives us more richly to explore sets of possible hyperparameters

*Wisdom of the Land*

# Tuning Strategies

*Coarse to fine*

- Start from randomly coarse
- Randomly fine later



Hyper-parameter $\alpha$

Hyper-parameter $\varepsilon$

$(\alpha', \varepsilon')$

**Random grid search**

Hyper-parameter $\alpha$

Hyper-parameter $\varepsilon$

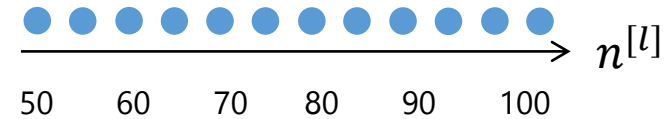**Precise grid search**

# Tuning Strategies

*Use the appropriate scale to pitch hyperparameters*
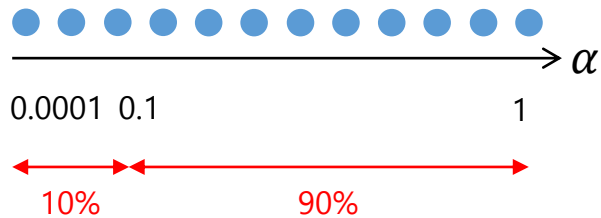
- Normal scale works for some hyperparameters
  - $n^{[l]}$
  - $L$



- Log-scale is good for learning rate





Spending 10% of resource to search between 0.0001 to 0.1 whereas 90% of resources to search between 0.1 to 1