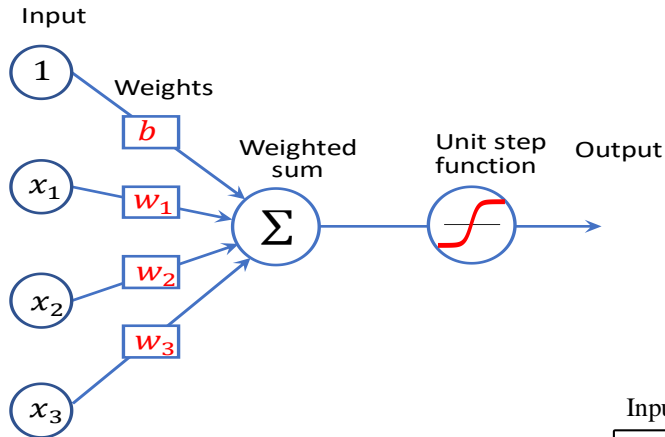
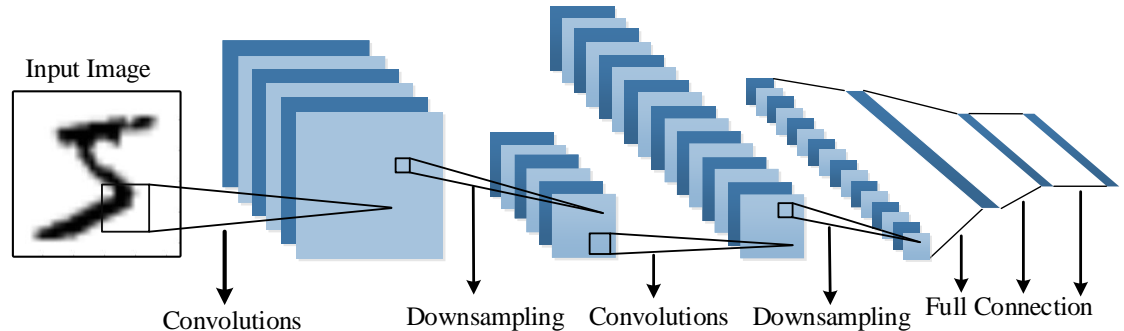


In this Course

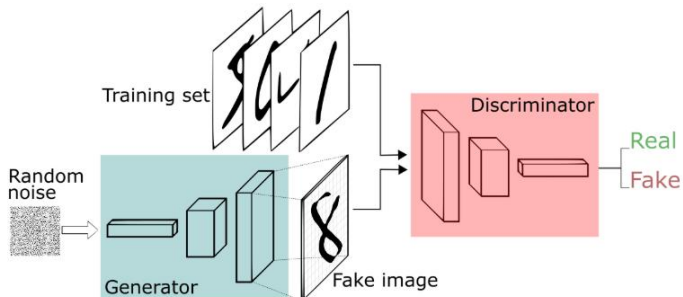
1. DL basics, linear regression, logistic regression etc.



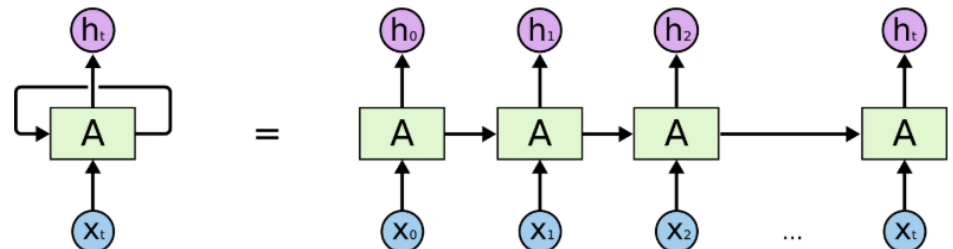
3. Convolutional Neural Networks and Applications



4. Generative Adversarial Networks



5. Recurrent networks and applications



Previous Lecture

- Course Information
- Introduction to Deep Learning
- Deep Learning Basics
 - Linear Regression
 - Loss Function
 - Gradient Descent
 - Regularization

Lecture 2

Gradient Descent

Stochastic Gradient Descent

Logistic Regression

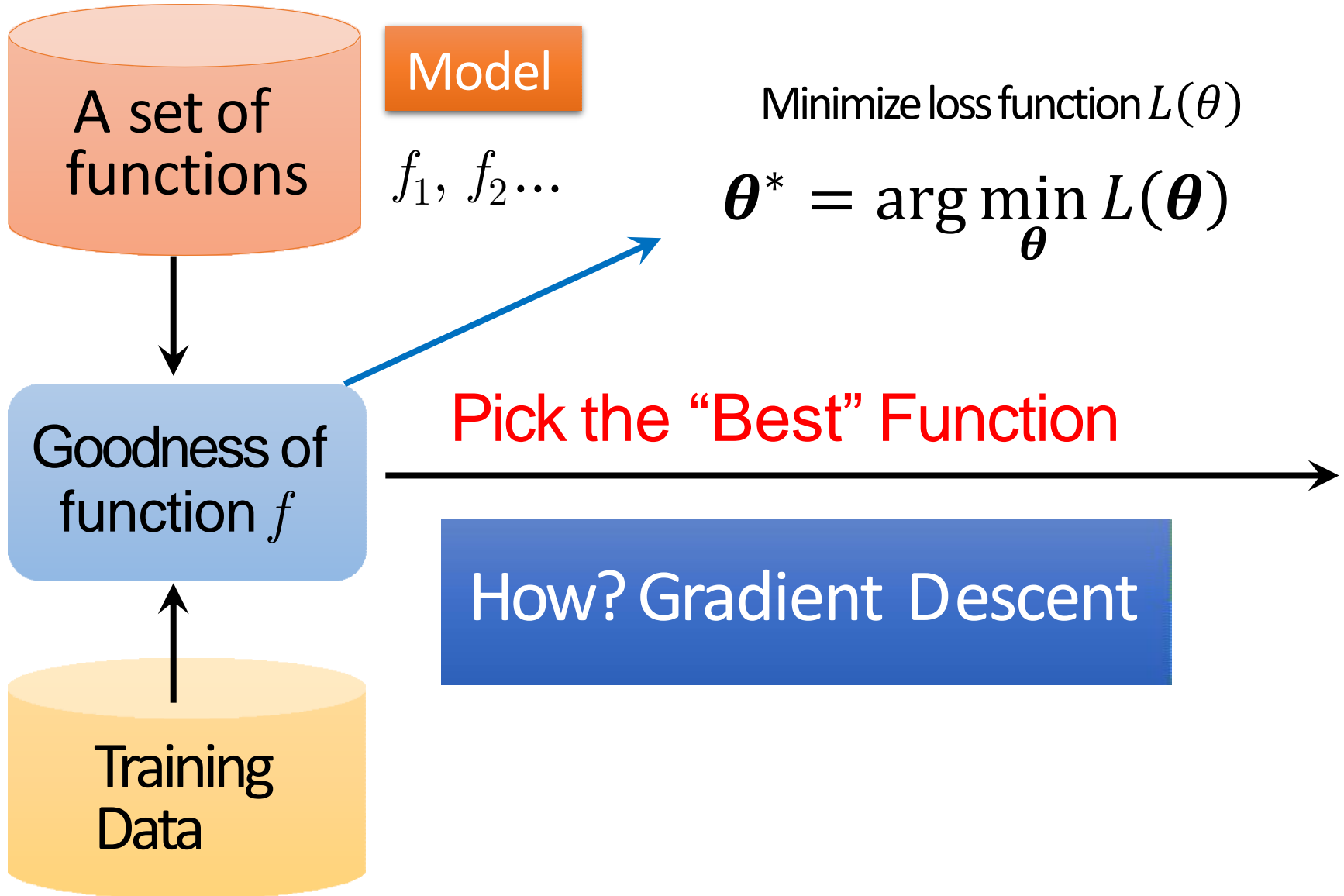
Gradient Descent



Many slides adapted from Andrew Ng and Hungyi Lee

Review: Step 3 in Machine Learning

$$y = f(\theta)$$



Review: Gradient Descent

- In step 3, we have to solve the following optimization problem:

$$\theta^* = \arg \min_{\theta} L(\theta) \quad L: \text{loss function} \quad \theta: \text{parameters}$$

Suppose that θ has two variables $\{\theta_1, \theta_2\}$

Randomly start at $\theta^{(0)} = \begin{bmatrix} \theta_1^{(0)} \\ \theta_2^{(0)} \end{bmatrix}$

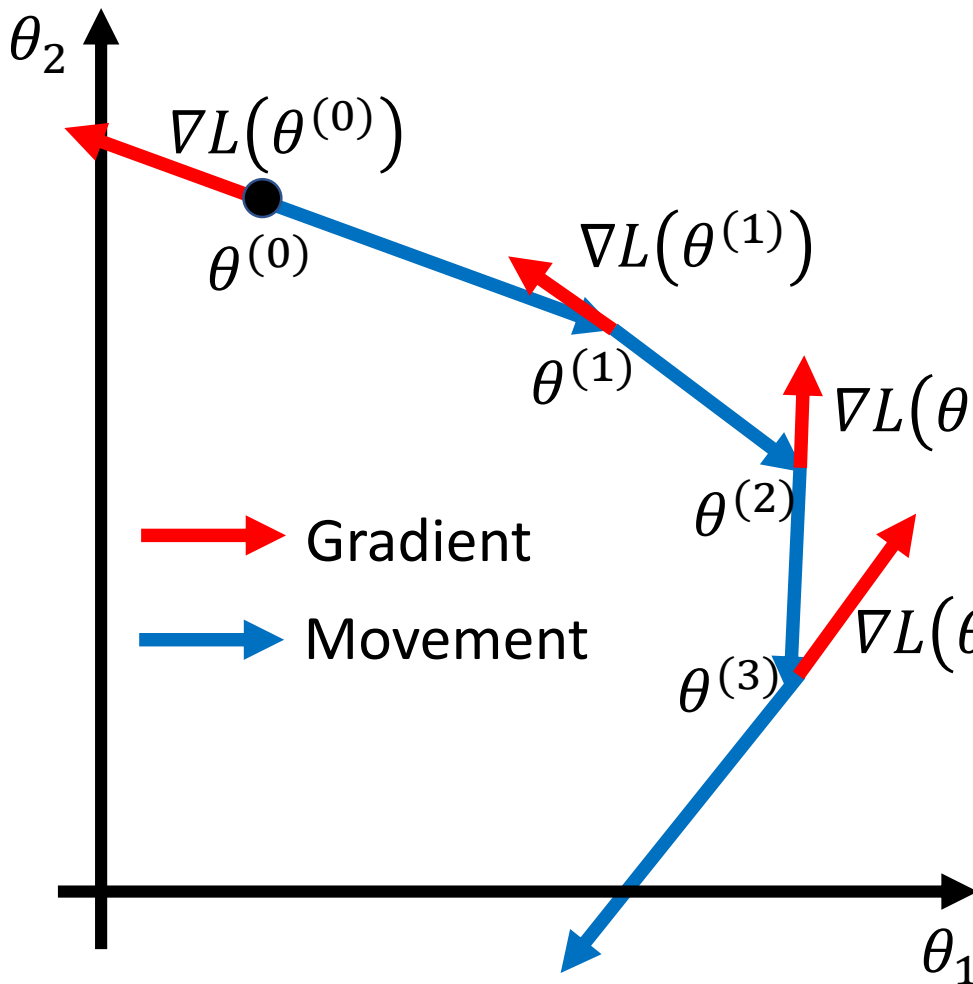
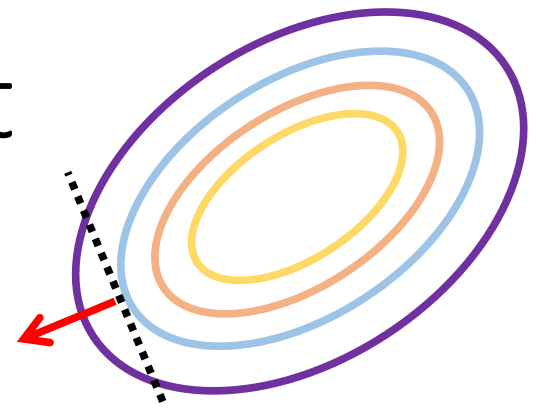
$$\nabla L(\theta) = \begin{bmatrix} \partial L(\theta_1) / \partial \theta_1 \\ \partial L(\theta_2) / \partial \theta_2 \end{bmatrix}$$

$$\begin{bmatrix} \theta_1^{(1)} \\ \theta_2^{(1)} \end{bmatrix} = \begin{bmatrix} \theta_1^{(0)} \\ \theta_2^{(0)} \end{bmatrix} - \eta \begin{bmatrix} \partial L(\theta_1^{(0)}) / \partial \theta_1 \\ \partial L(\theta_2^{(0)}) / \partial \theta_2 \end{bmatrix} \Rightarrow \theta^{(1)} = \theta^{(0)} - \eta \nabla L(\theta^{(0)})$$

$$\begin{bmatrix} \theta_1^{(2)} \\ \theta_2^{(2)} \end{bmatrix} = \begin{bmatrix} \theta_1^{(1)} \\ \theta_2^{(1)} \end{bmatrix} - \eta \begin{bmatrix} \partial L(\theta_1^{(1)}) / \partial \theta_1 \\ \partial L(\theta_2^{(1)}) / \partial \theta_2 \end{bmatrix} \Rightarrow \theta^{(2)} = \theta^{(1)} - \eta \nabla L(\theta^{(1)})$$

Review: Gradient Descent

Gradient: the normal direction of the contour of loss function



Start at position $\theta^{(0)}$

Compute gradient at $\theta^{(0)}$

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

Compute gradient at $\theta^{(1)}$

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

⋮

Gradient Descent

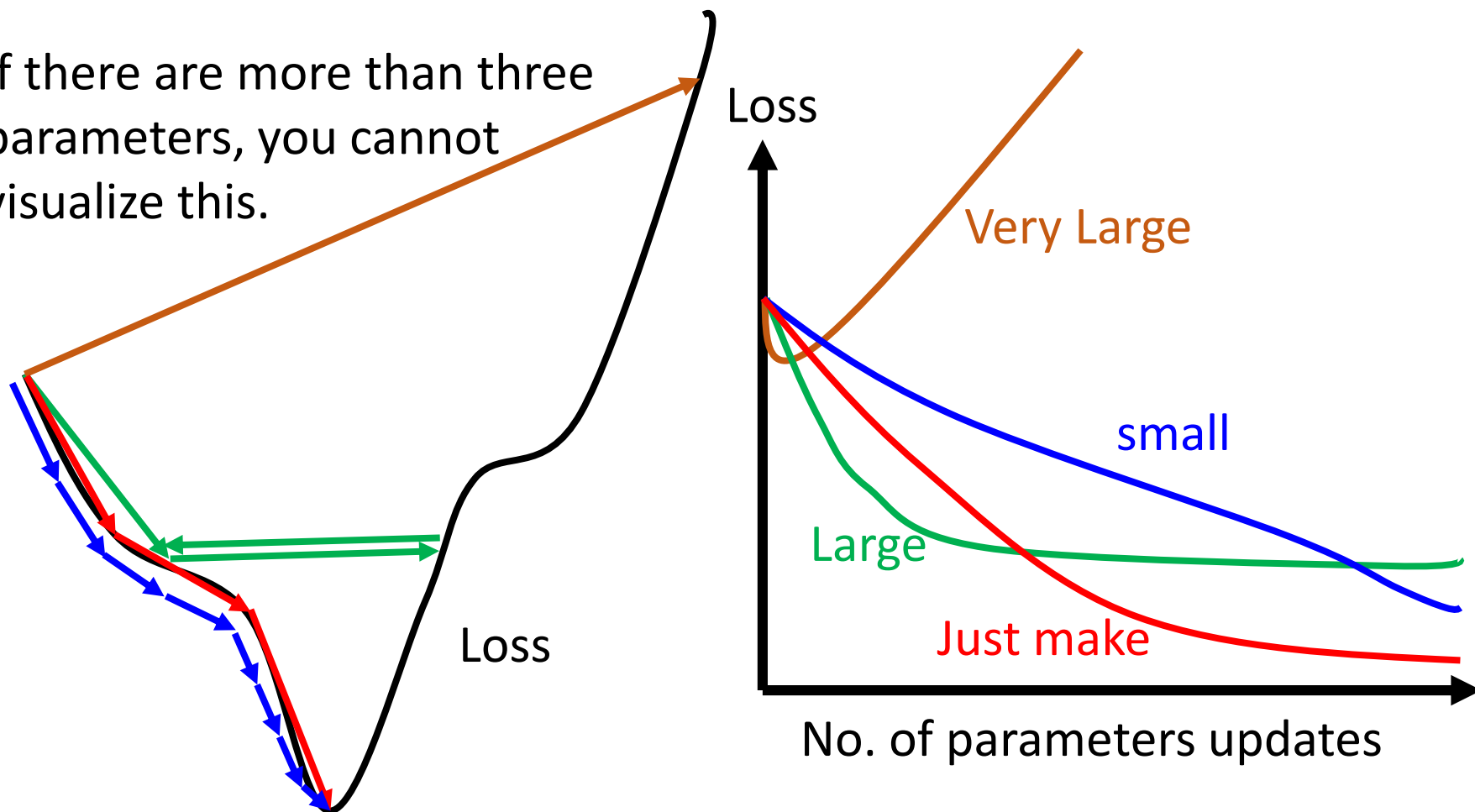
Tip 1: Tuning your
learning rates

Learning Rate

$$\theta^{(i)} = \theta^{(i-1)} - \eta \nabla L(\theta^{(i-1)})$$

Set the learning rate η carefully

If there are more than three parameters, you cannot visualize this.



But you can always visualize this.

Adaptive Learning Rates

- Popular and Simple Idea: Reduce the learning rate by some factor every few epochs.
 - At the beginning, we are far from the destination, so we use larger learning rate
 - After several epochs, we are close to the destination, so we reduce the learning rate to make sure it converges to the minimum.
 - *e.g.* $1/t$ decay: $\eta^{(t)} = \eta / \sqrt{t + 1}$
- Learning rate **cannot be “one size fits all”**
 - Giving different parameters different learning rates

Adagrad

- Divide the learning rate of each parameter by the *root mean square of all its previous derivatives*

Vanilla Gradient descent

$$w^{(t+1)} \leftarrow w^{(t)} - \eta^{(t)} g^{(t)}$$

w is one parameters

$$\eta^{(t)} = \frac{\eta}{\sqrt{t+1}}$$

$$g^{(t)} = \frac{\partial L(w)}{\partial w} \Big|_{w=w^{(t)}}$$

Adagrad

$$w^{(t+1)} \leftarrow w^{(t)} - \frac{\eta^{(t)}}{\sigma^{(t)}} g^{(t)}$$

$\sigma^{(t)}$: *root mean square* of **all the previous derivatives** of parameter w up to iteration t .

Parameter dependent

Adagrad

$$\eta^{(t)} = \frac{\eta}{\sqrt{t+1}}$$

$$g^{(t)} = \frac{\partial L(w)}{\partial w} \Big|_{w=w^{(t)}}$$

$\sigma^{(t)}$: *root mean square* of **all the previous derivatives** of parameter w up to iteration t .

$$w^{(1)} \leftarrow w^{(0)} - \frac{\eta^{(0)}}{\sigma^{(0)}} g^{(0)}$$

$$\sigma^{(0)} = \sqrt{(g^{(0)})^2}$$

$$w^{(2)} \leftarrow w^{(1)} - \frac{\eta^{(1)}}{\sigma^{(1)}} g^{(1)}$$

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

$$w^{(3)} \leftarrow w^{(2)} - \frac{\eta^{(2)}}{\sigma^{(2)}} g^{(2)}$$

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

\vdots

$$w^{(t+1)} \leftarrow w^{(t)} - \frac{\eta^{(t)}}{\sigma^{(t)}} g^{(t)}$$

$$\sigma^{(t)} = \sqrt{\frac{1}{t+1} \sum_{i=0}^t (g^{(i)})^2}$$

Adagrad

$$\eta^{(t)} = \frac{\eta}{\sqrt{t+1}}$$

$$g^{(t)} = \frac{\partial L(w)}{\partial w} \Big|_{w=w^{(t)}}$$

- Divide the learning rate of each parameter by the ***root mean square of all its previous derivatives***

$$w^{(t+1)} \leftarrow w^{(t)} - \frac{\eta^{(t)}}{\sigma^{(t)}} g^{(t)}$$
$$\eta^{(t)} = \frac{\eta}{\sqrt{t+1}}$$
$$\sigma^{(t)} = \sqrt{\frac{1}{t+1} \sum_{i=0}^t (g^{(i)})^2}$$
$$w^{(t+1)} \leftarrow w^{(t)} - \frac{\eta}{\sqrt{\sum_{i=0}^t (g^{(i)})^2}} g^{(t)}$$

- Other adaptive learning rate methods: Adadelata, Adam,...

Contradiction?

$$\eta^{(t)} = \frac{\eta}{\sqrt{t+1}}$$

$$g^{(t)} = \frac{\partial L(w)}{\partial w} \Big|_{w=w^{(t)}}$$

Vanilla Gradient descent

$$w^{(t+1)} \leftarrow w^{(t)} - \eta^{(t)} g^{(t)}$$

Larger gradient,
larger step

Adagrad

$$w^{(t+1)} \leftarrow w^{(t)} - \frac{\eta}{\sqrt{\sum_{i=0}^t (g^{(i)})^2}} g^{(t)}$$

Larger gradient,
larger step

Larger gradient,
smaller step

Gradient Descent

Tip 2: Stochastic Gradient Descent

Make the training faster

Stochastic Gradient Descent

$$L = \sum_n \left(y^{(n)} - \left(b + \sum_i w_i x_i^{(n)} \right) \right)^2$$

Loss is the summation over all training examples

- **Gradient Descent**

$$\theta^{(i+1)} \leftarrow \theta^{(i)} - \eta \nabla L(\theta^{(i)})$$

- **Stochastic Gradient Descent (SGD)**

Faster!

Randomly pick one example $x^{(n)}$ to update parameters

$$L^{(n)} = \left(y^{(n)} - \left(b + \sum_i w_i x_i^{(n)} \right) \right)^2$$

$$\theta^{(i+1)} \leftarrow \theta^{(i)} - \eta \nabla L(\theta^{(i)})$$

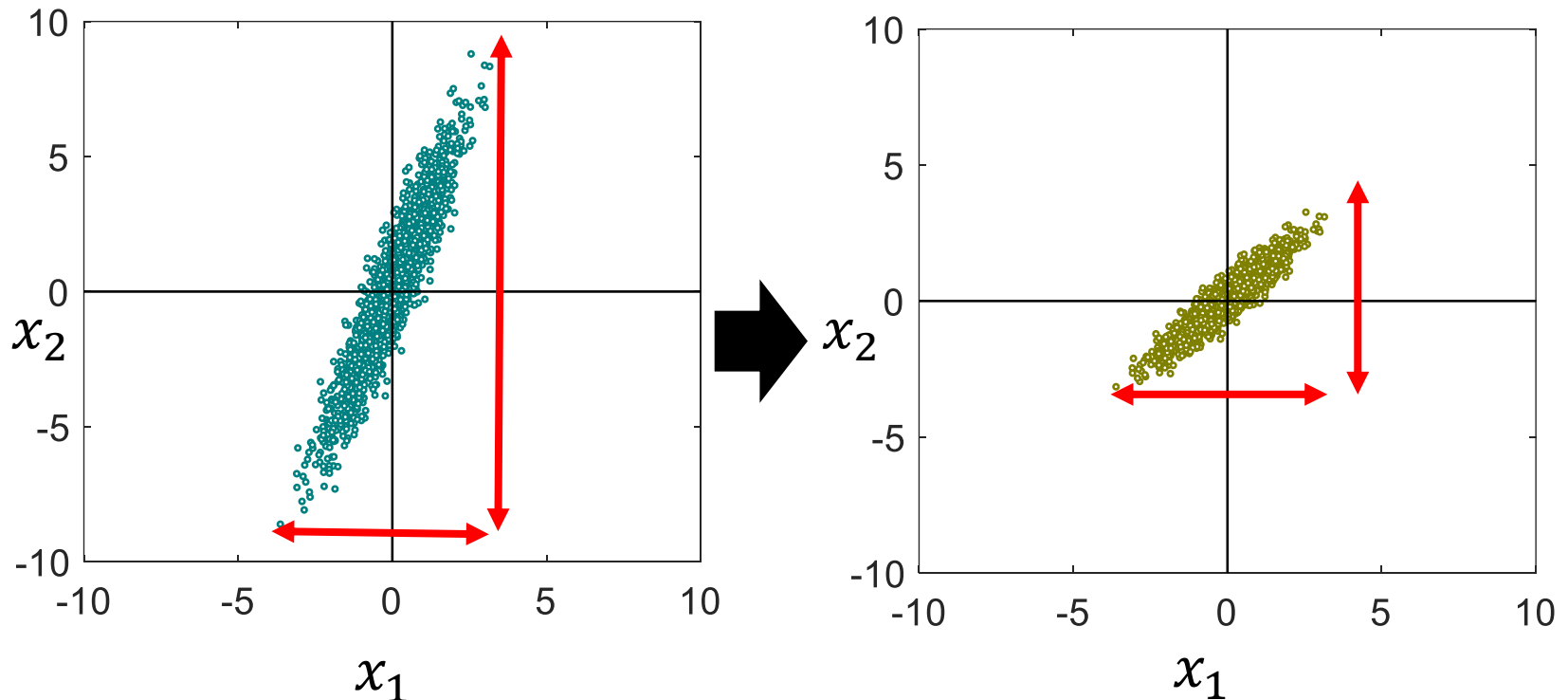
Loss for only one example, i.e. n^{th} sample

Gradient Descent

Tip 3: Feature Scaling

Feature Scaling/Feature Normalization

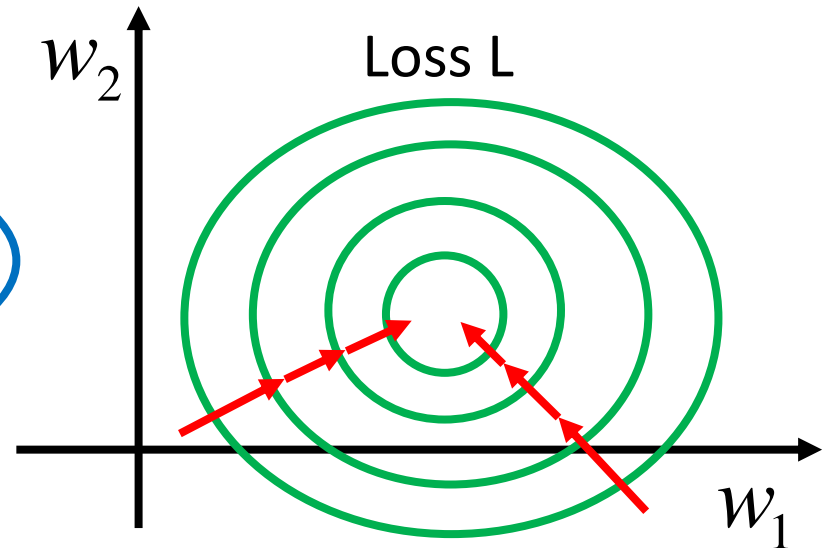
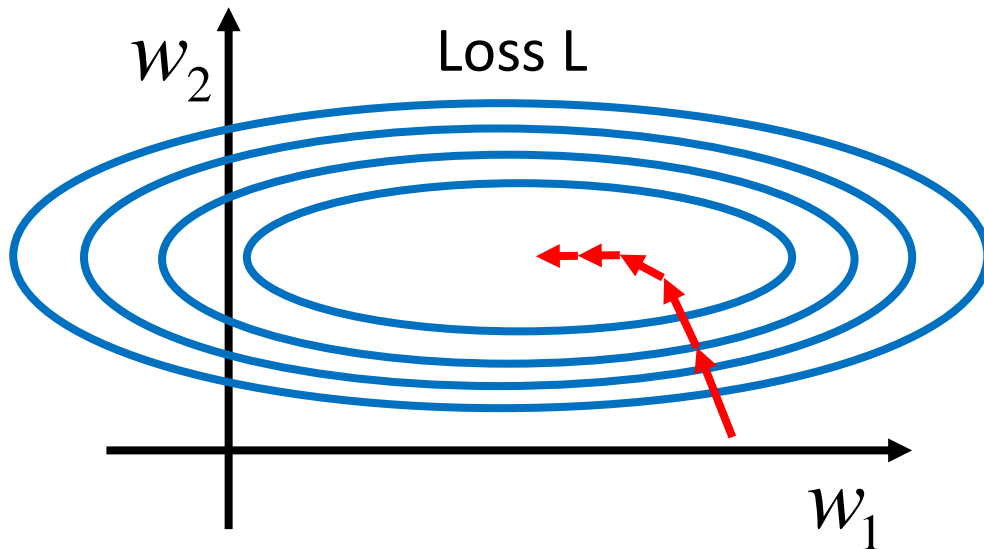
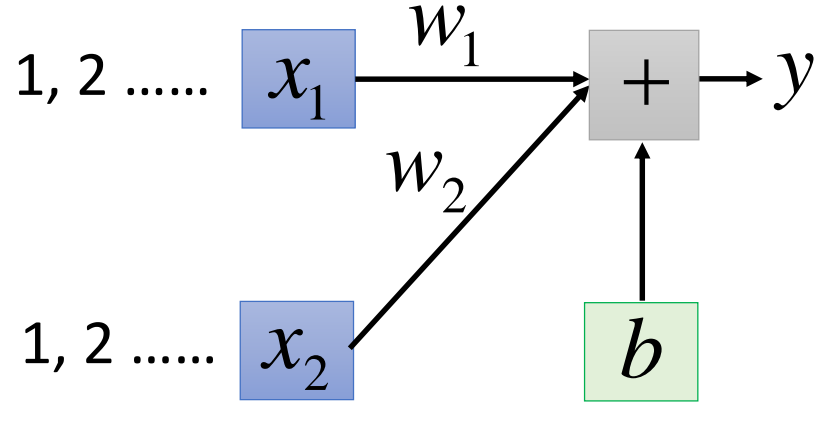
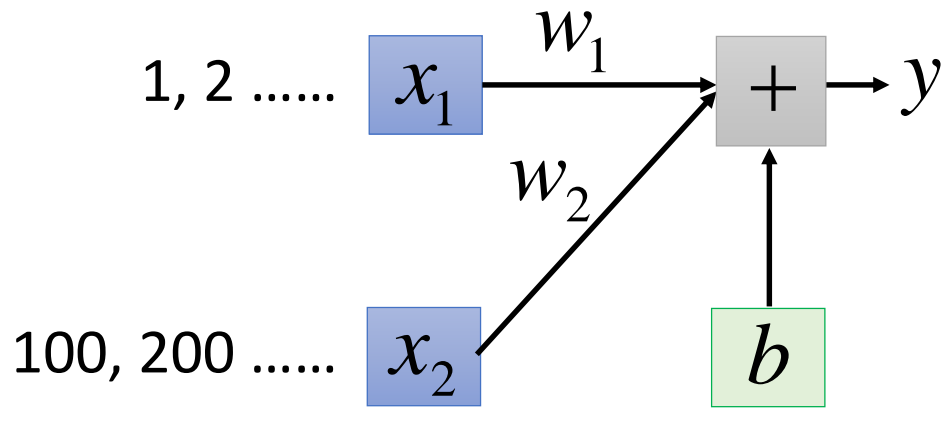
$$y = b + w_1x_1 + w_2x_2$$



Make different features have the same scaling

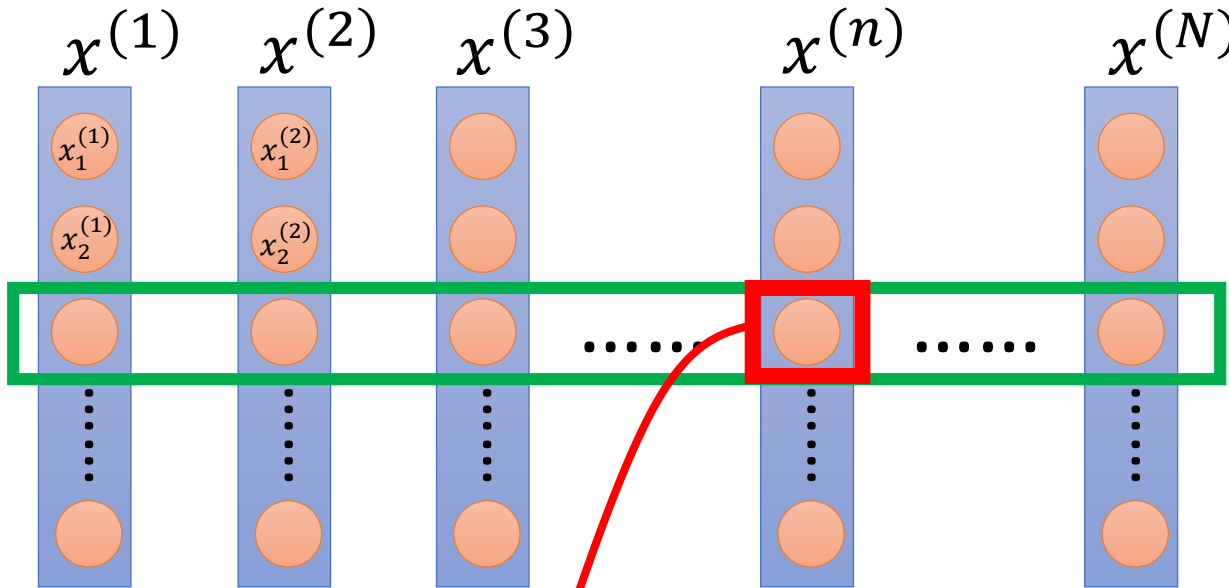
Feature Scaling

$$y = b + w_1x_1 + w_2x_2$$



Feature Scaling

$$m_i = \frac{1}{N} \sum_{i=1}^N x_i^{(n)}$$



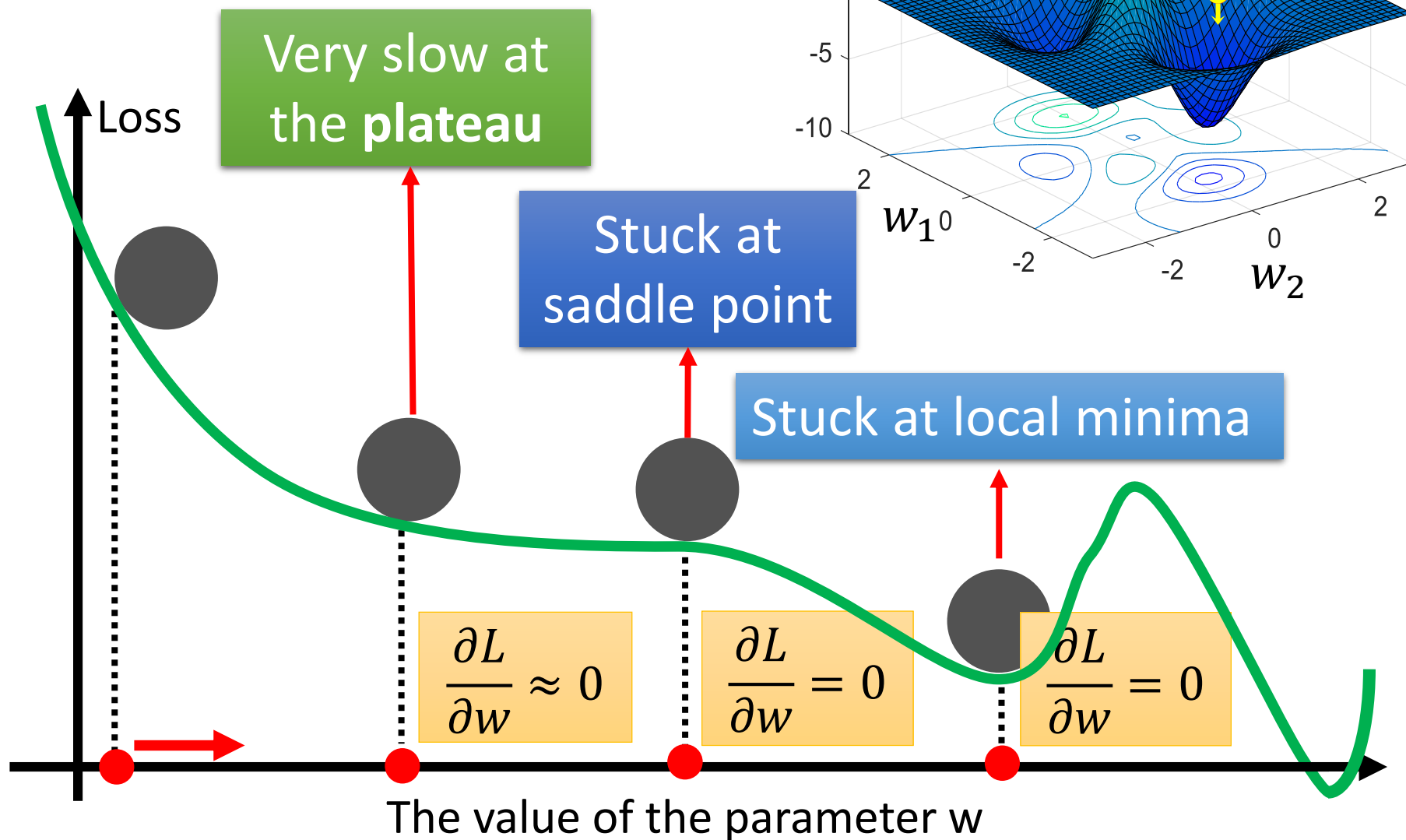
For each dimension i :

- mean: m_i
- standard deviation: σ_i

$$x_i^{(n)} \leftarrow \frac{x_i^{(n)} - m_i}{\sigma_i}$$

Make each feature component zero mean and unit standard deviation.

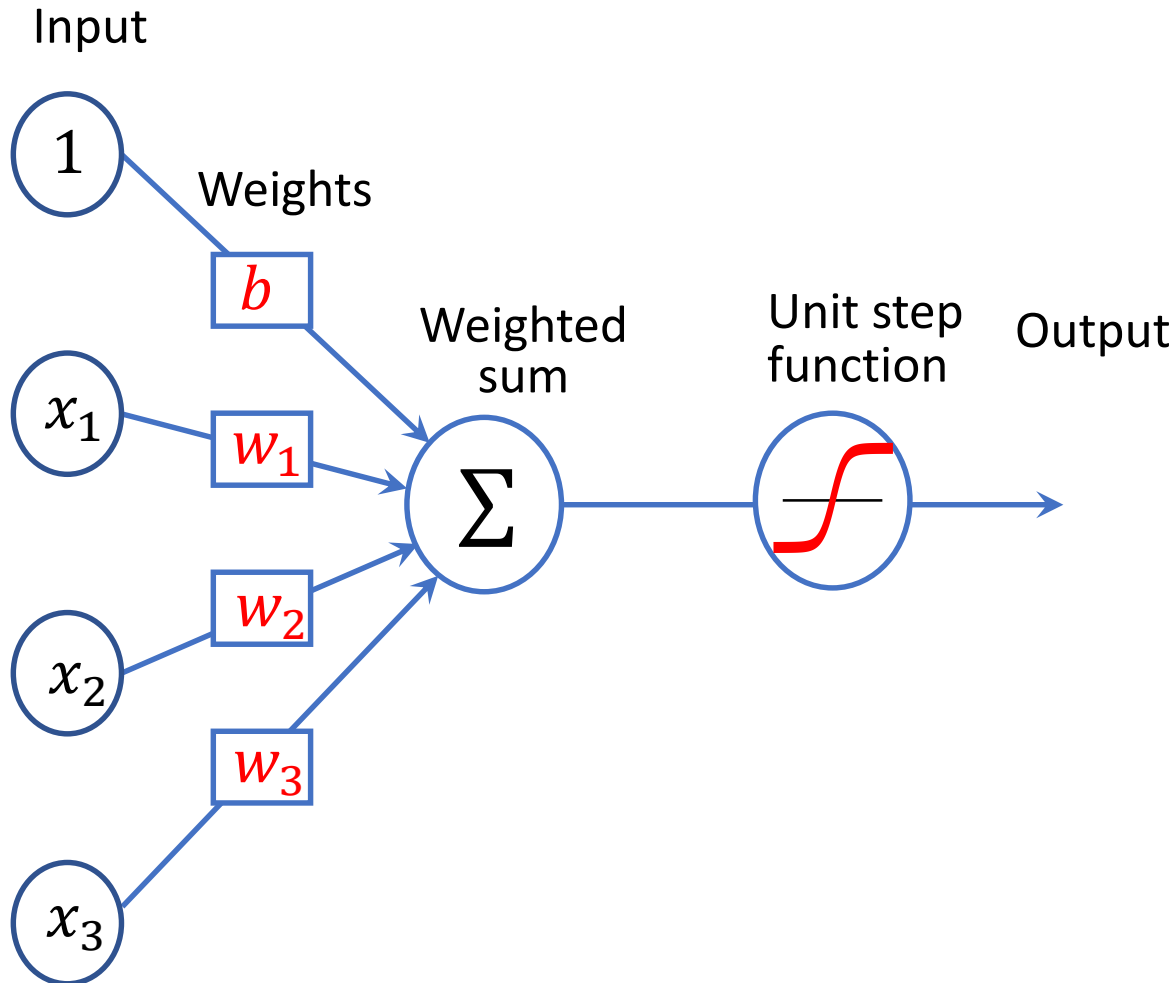
More Limitation of Gradient Descent



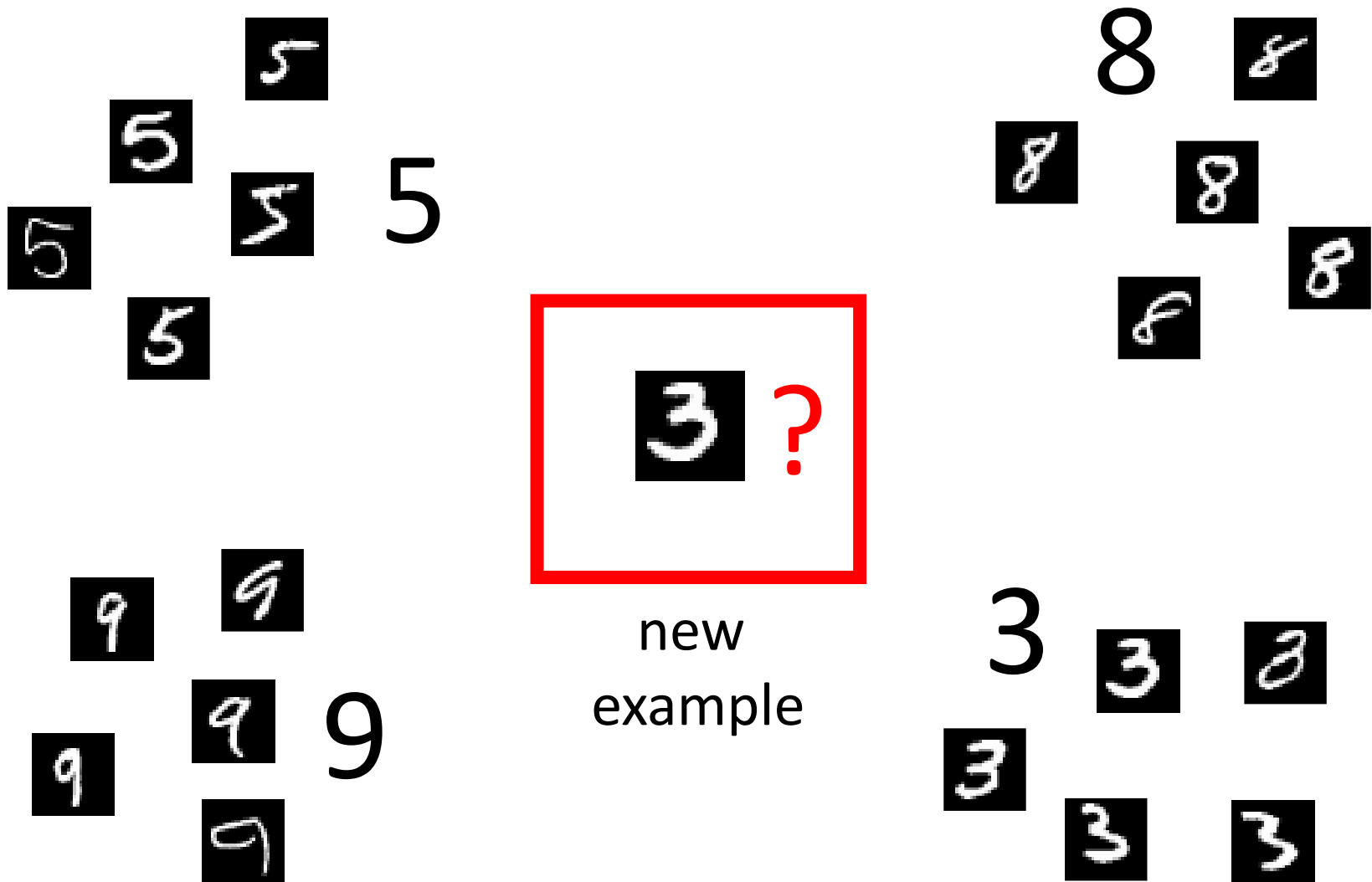
An overview of gradient descent optimization algorithms

- <https://arxiv.org/pdf/1609.04747.pdf>

Classification: Logistic Regression



Supervised Classification



Classification: object classification



ImageNet Challenge: classification of
1000 object category

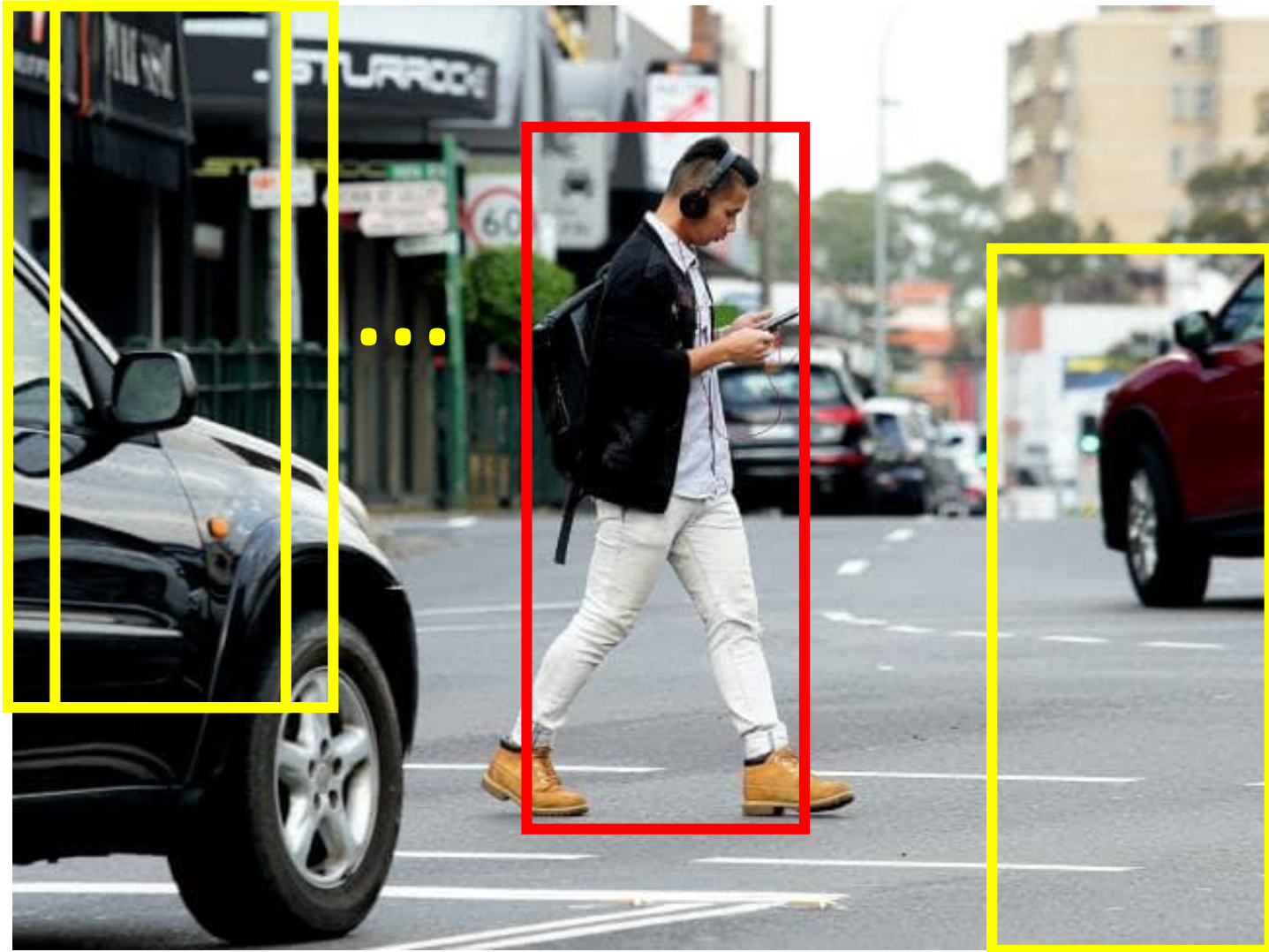
Classification: face recognition



Pedestrian Detection



Pedestrian Detection



Classification

- Email: Spam / Not Spam?
- Pedestrian Detection: Pedestrian / Not Pedestrian?
- Tumor: Malignant / Benign ?

$y \in \{0,1\}$ 0: “Negative Class” (*e.g.* Not Pedestrian)
1: “Positive Class” (*e.g.* Pedestrian)

Values 0 and 1 are somewhat arbitrary.

$y \in \{-1,1\}$ 0: “Negative Class” (*e.g.* Not Pedestrian)
1: “Positive Class” (*e.g.* Pedestrian)

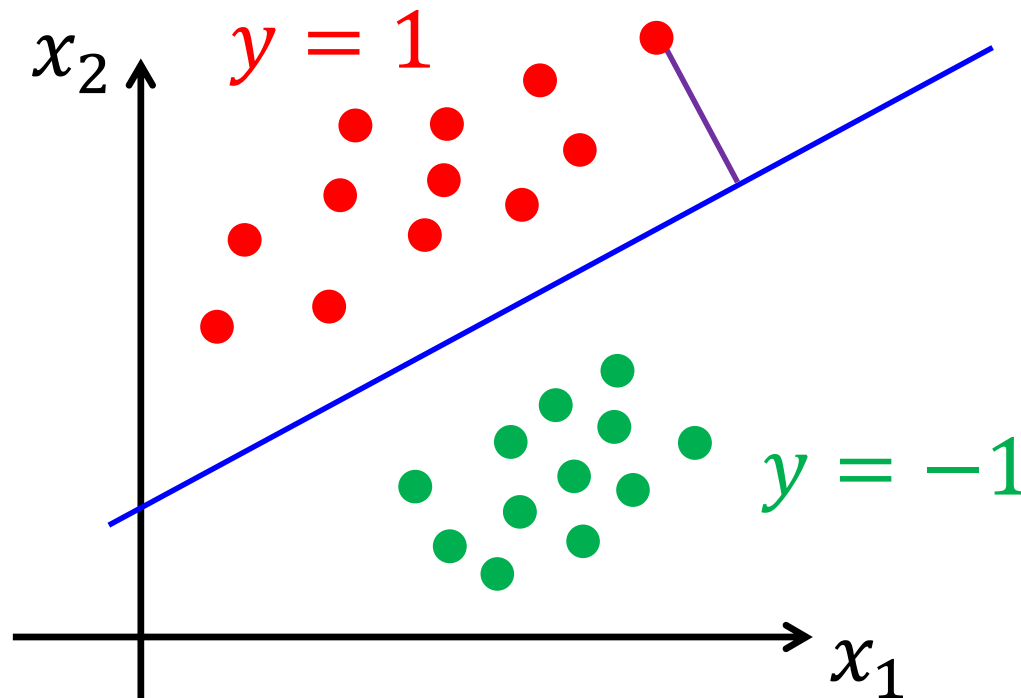
Linear Classifier

- Training data for Classification

$$\{(\mathbf{x}^{(k)}, y^{(k)})\}_{k=1}^N \quad y^{(k)} \in \{-1, 1\}$$

Prediction function:

$$f_{w,b}(\mathbf{x}) = \text{sgn}(\mathbf{w}^T \mathbf{x} + b)$$



$$\frac{|\mathbf{w}^T \mathbf{x} + b|}{\sqrt{\sum_i w_i^2}}$$

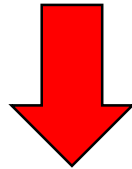
Logistic Regression Model

Want $0 \leq f_{w,b}(x) \leq 1$

$$f_{w,b}(x) = \mathbf{w}^T \mathbf{x} + b?$$

$\sigma(z) \geq 0.5$, class 1

$\sigma(z) < 0.5$, class 2

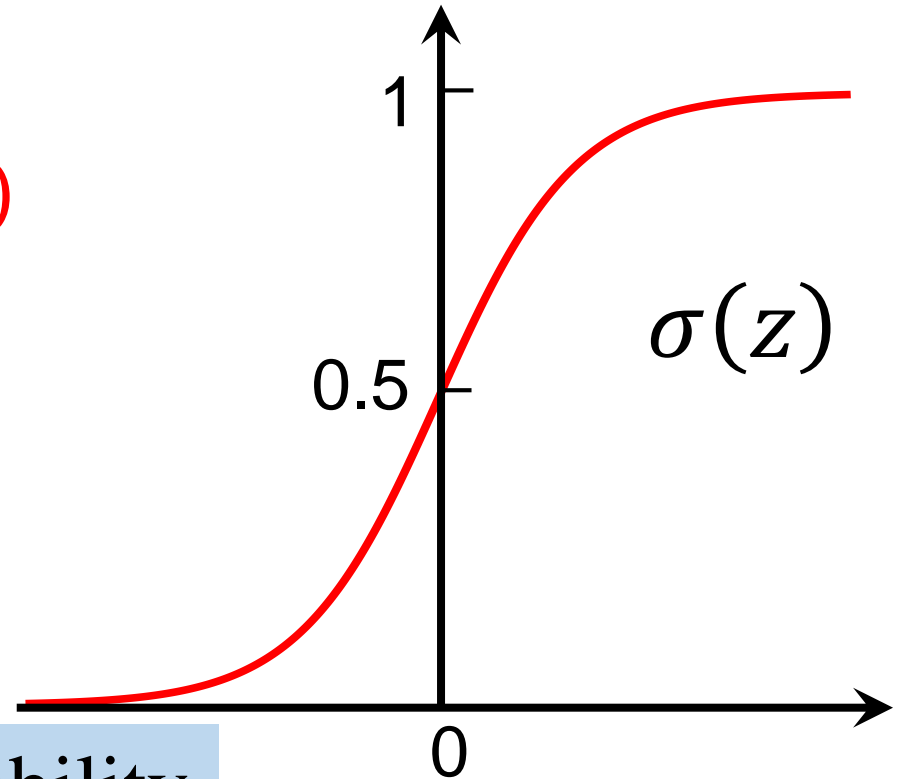


$$f_{w,b}(x) = \sigma(\mathbf{w}^T \mathbf{x} + b)$$

$$\sigma(z) = \frac{1}{1 + \exp(-z)}$$

Sigmoid function

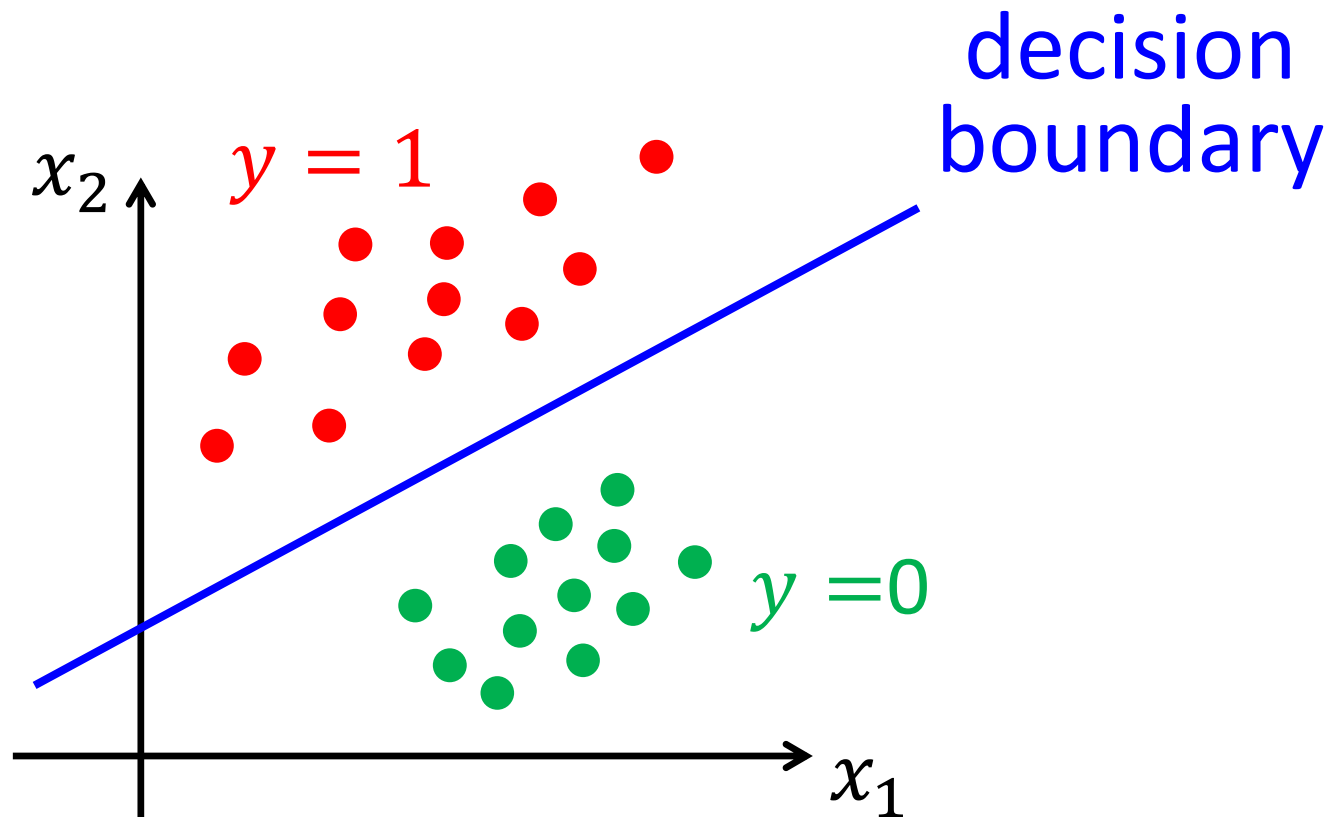
Logistic function



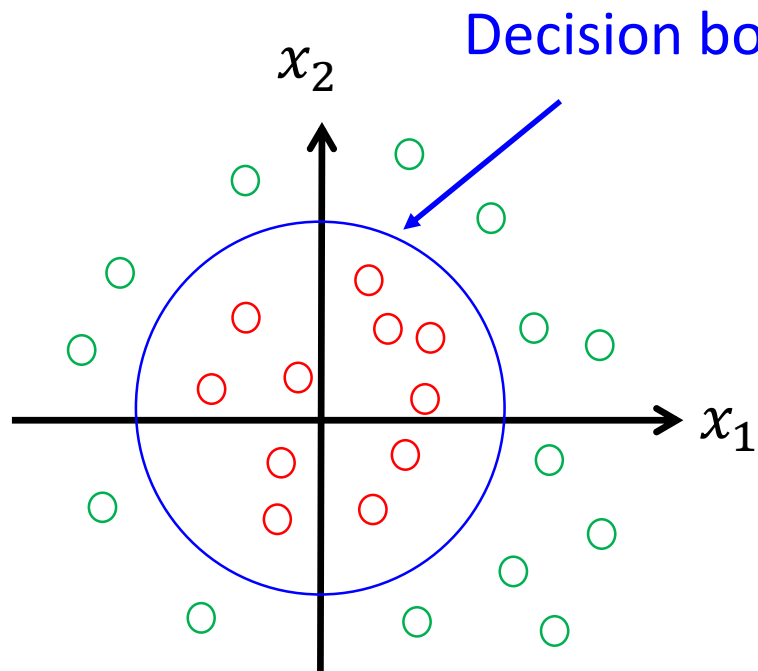
$\sigma(z)$ means posterior Probability

Decision Boundary

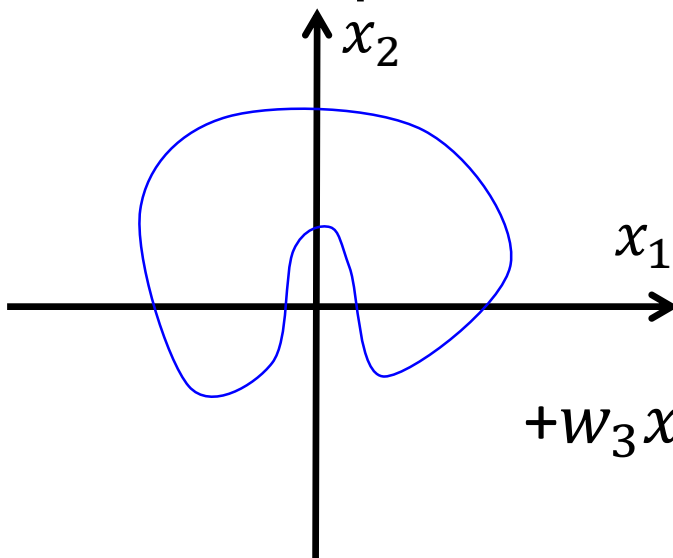
$$f_{w,b}(x) = \sigma(b + w_1x_1 + w_2x_2)$$



Nonlinear decision boundaries



$$f_{w,b}(x) = \sigma(w_0 + w_1x_1 + w_2x_2 + w_3x_1^2 + w_4x_2^2)$$



$$f_w(x) = \sigma(w_0 + w_1x_1 + w_2x_2 + w_3x_1^2 + w_4x_1^2x_2 + w_5x_1^2x_2^2 + w_6x_1^3x_2 + \dots)$$

Step 1: Function Set

- Function set: Including all different w and b

$$\begin{cases} P_{w,b}(C_1|x) \geq 0.5, & \text{class 1} \\ P_{w,b}(C_1|x) < 0.5, & \text{class 2} \end{cases}$$

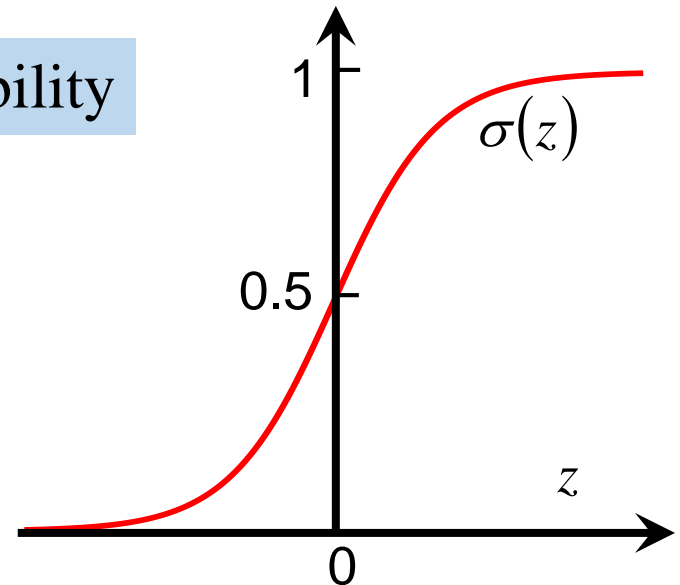
$$P_{w,b}(C_1|x) = \sigma(z) \quad \text{Posterior Probability}$$

- Hypothesis

sigmoid function (logistic function)

$$z = \mathbf{w}^T \mathbf{x} + b = \sum_i w_i x_i + b$$

$$\sigma(z) = \frac{1}{1 + \exp(-z)}$$



Step 1: Function Set

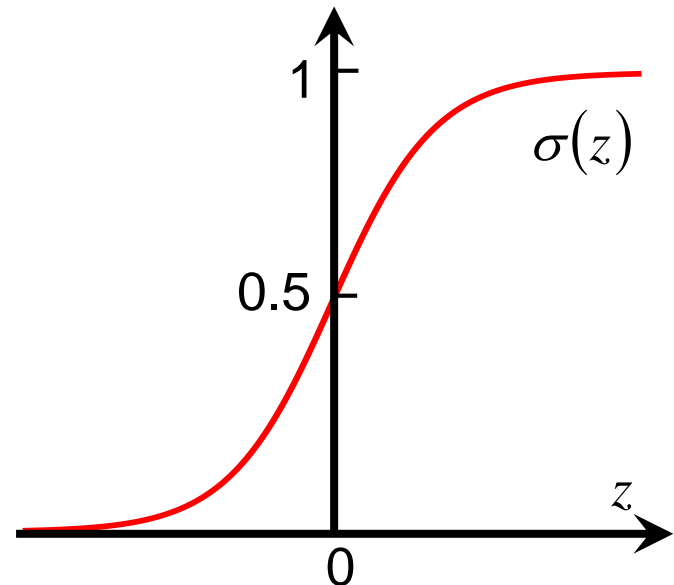
Function set: Including all different w and b

$$\left\{ \begin{array}{ll} z \geq 0 & \text{class 1} \\ z < 0 & \text{class 2} \end{array} \right.$$

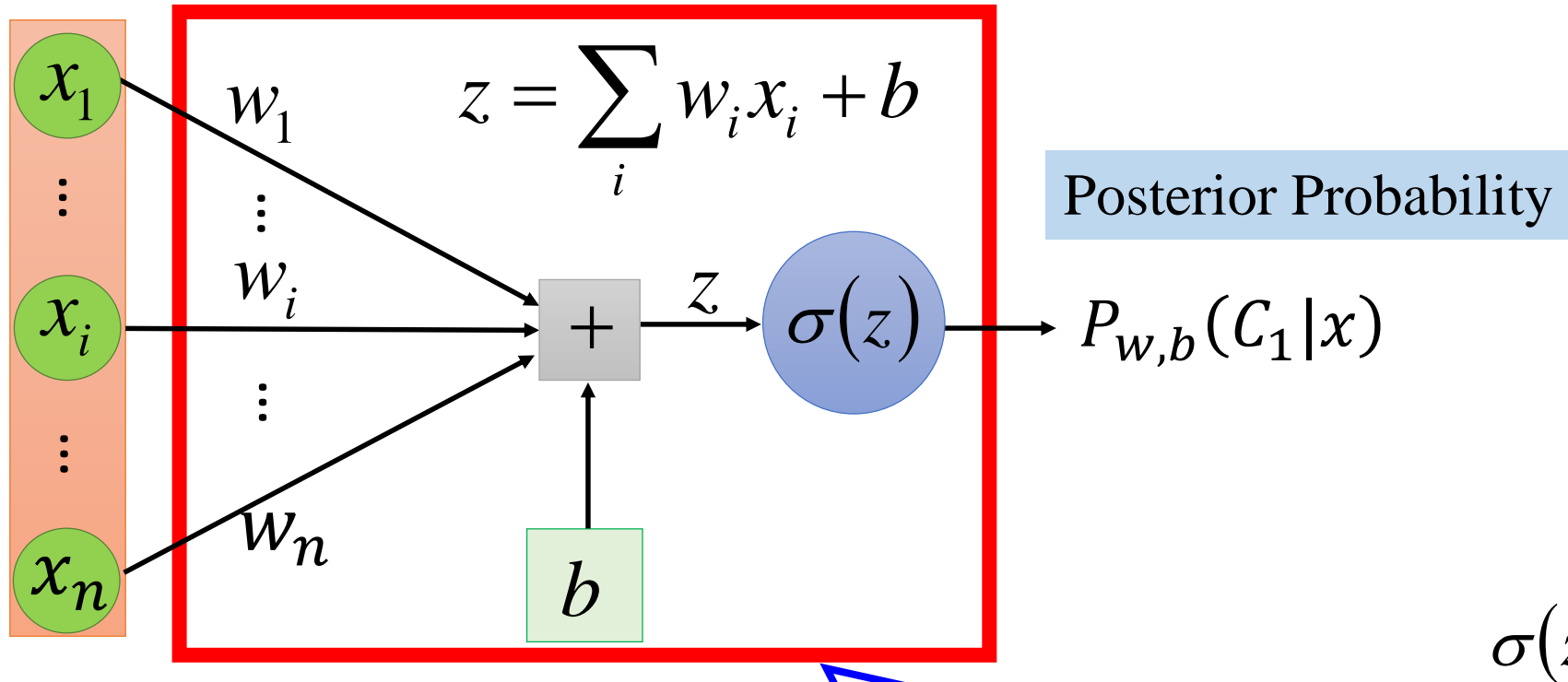
$$P_{w,b}(C_1|x) = \sigma(z)$$

$$z = \mathbf{w}^T \mathbf{x} + b = \sum_i w_i x_i + b$$

$$\sigma(z) = \frac{1}{1 + \exp(-z)}$$



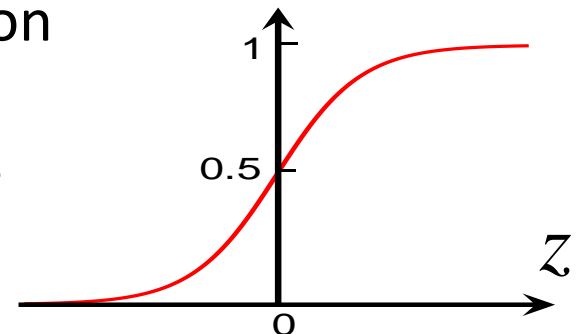
Step 1: Function Set



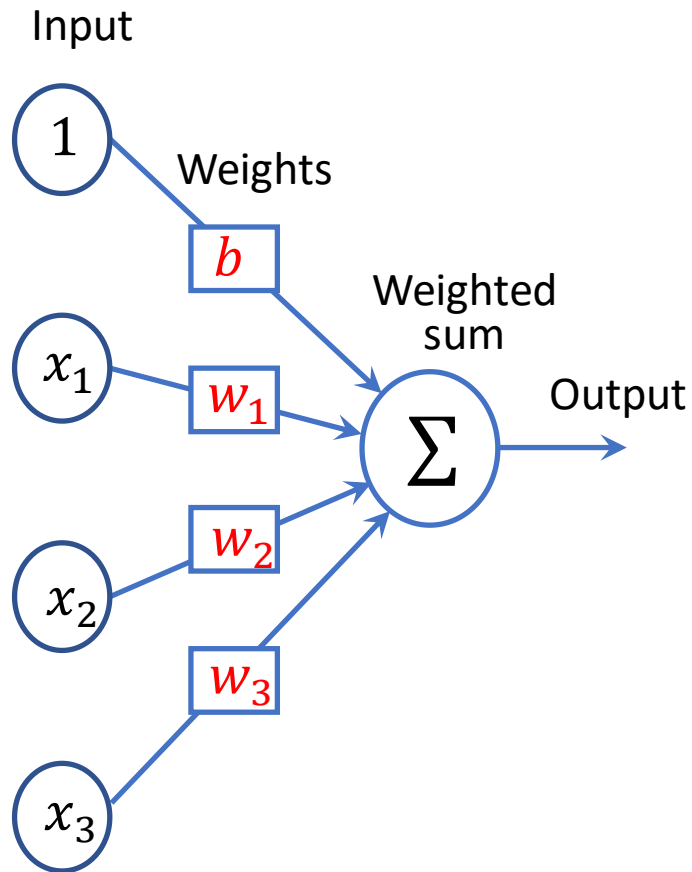
Logistic
Regression

Sigmoid Function

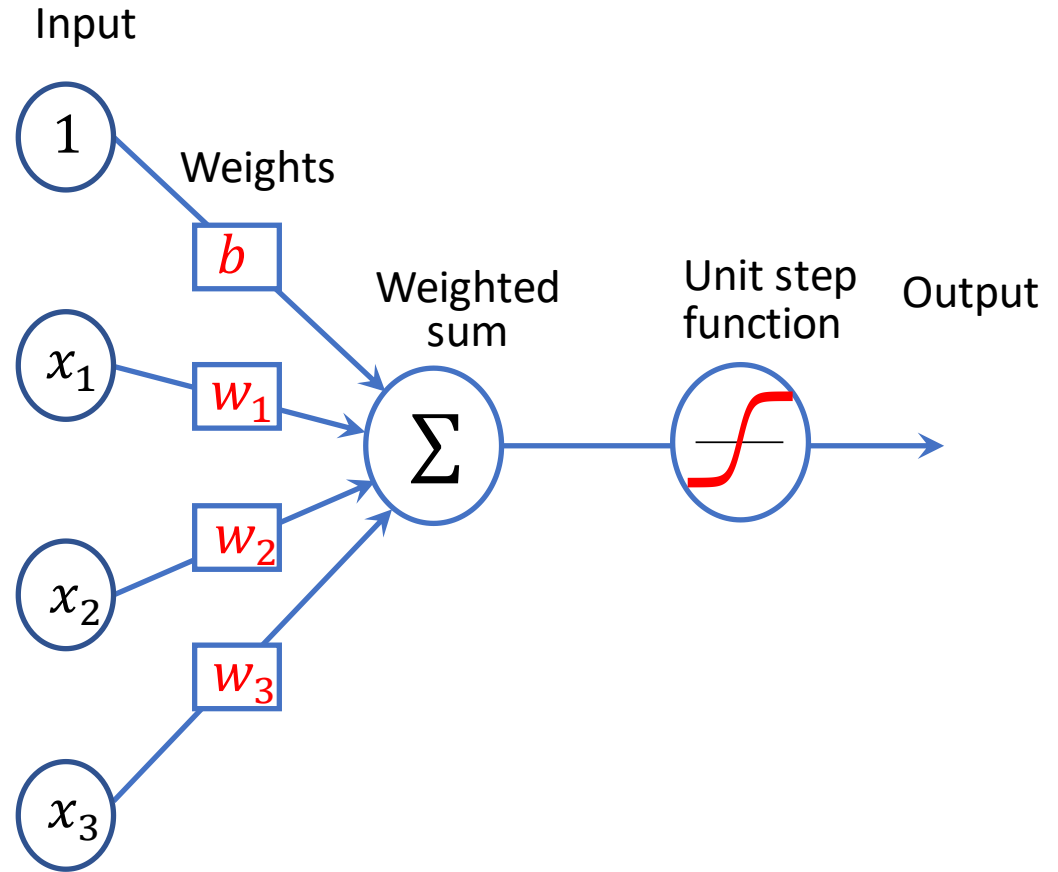
$$\sigma(z) = \frac{1}{1 + e^{-z}}$$



Linear Regression (Prediction)



vs. Logistic Regression (Classification)



Step 2: Goodness of a Function

Training
Data

$x^{(1)}$	$x^{(2)}$	$x^{(3)}$	$\dots \dots$	$x^{(N)}$
C_1	C_1	C_2		C_1

$$P_{w,b}(C_1|x) = \frac{1}{1+\exp(-z)}$$

$$z = \mathbf{w}^T \mathbf{x} + b$$

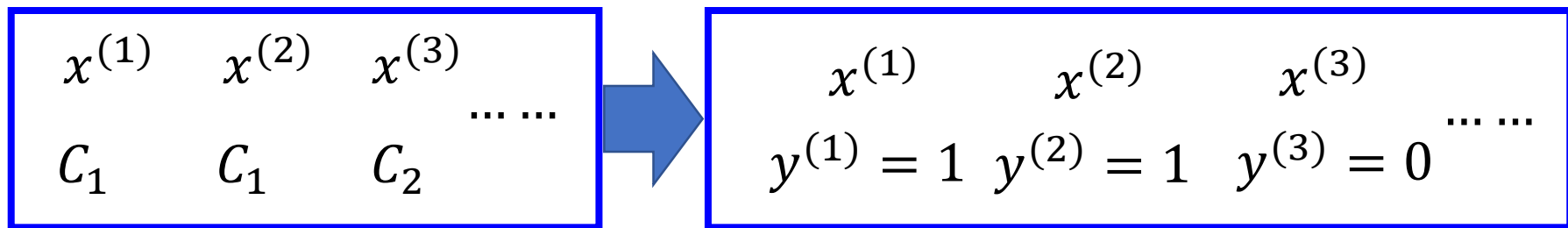
Assume the data is generated based on $f_{w,b}(x) = P_{w,b}(C_1|x)$

Given a set of \mathbf{w} and b , what is its probability of generating the data?

$$L(w, b) = f_{w,b}(x^{(1)})f_{w,b}(x^{(2)})\left(1 - f_{w,b}(x^{(3)})\right)\cdots f_{w,b}(x^{(N)})$$

The most likely w^* and b^* is the one with the largest $L(w, b)$.

$$w^*, b^* = \arg \max_{w,b} L(w, b)$$



$$L(w, b) = f_{w,b}(x^{(1)}) f_{w,b}(x^{(2)}) (1 - f_{w,b}(x^{(3)})) \dots$$

$$w^*, b^* = \arg \max_{w, b} L(w, b) = w^*, b^* = \arg \min_{w, b} -\ln L(w, b)$$

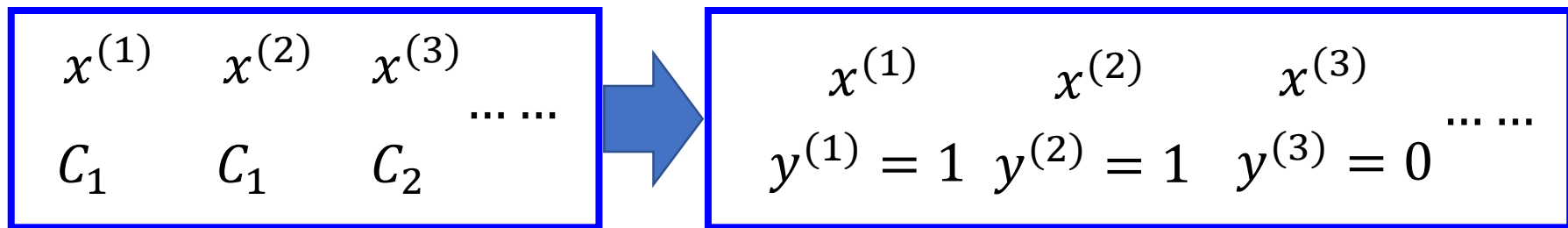
$$-\ln L(w, b)$$

$$= -\ln f_{w,b}(x^{(1)}) \Rightarrow -\left[\underbrace{1}_{\text{orange}} \ln f(x^{(1)}) + \underbrace{0}_{\text{orange}} \ln (1 - f(x^{(1)})) \right]$$

$$-\ln f_{w,b}(x^{(2)}) \Rightarrow -\left[\underbrace{1}_{\text{orange}} \ln f(x^{(2)}) + \underbrace{0}_{\text{orange}} \ln (1 - f(x^{(2)})) \right]$$

$$-\ln (1 - f_{w,b}(x^{(3)})) \Rightarrow -\left[\underbrace{0}_{\text{orange}} \ln f(x^{(3)}) + \underbrace{1}_{\text{orange}} \ln (1 - f(x^{(3)})) \right]$$

\vdots



$$L(w, b) = f_{w,b}(x^{(1)})f_{w,b}(x^{(2)}) \left(1 - f_{w,b}(x^{(3)})\right) \dots$$

$$w^*, b^* = \arg \max_{w,b} L(w, b) = w^*, b^* = \arg \min_{w,b} -\ln L(w, b)$$

$$\begin{aligned}
 & -\ln L(w, b) \\
 & = -\ln f_{w,b}(x^{(1)}) \Rightarrow -\left[y^{(1)} \ln f(x^{(1)}) + (1 - y^{(1)}) \ln (1 - f(x^{(1)}))\right] \\
 & \quad -\ln f_{w,b}(x^{(2)}) \Rightarrow -\left[y^{(2)} \ln f(x^{(2)}) + (1 - y^{(2)}) \ln (1 - f(x^{(2)}))\right] \\
 & \quad -\ln (1 - f_{w,b}(x^{(3)})) \Rightarrow -\left[y^{(3)} \ln f(x^{(3)}) + (1 - y^{(3)}) \ln (1 - f(x^{(3)}))\right] \\
 & \quad \vdots \\
 & \quad \quad \quad \parallel \\
 & \quad \quad \quad \sum_{k=1}^N -\left[y^{(k)} \ln f_{w,b}(x^{(k)}) + (1 - y^{(k)}) \ln (1 - f_{w,b}(x^{(k)}))\right]
 \end{aligned}$$

Step 2: Goodness of a Function

$$L(w, b) = f_{w,b}(x^{(1)})f_{w,b}(x^{(2)})\left(1 - f_{w,b}(x^{(3)})\right)\cdots f_{w,b}(x^{(N)})$$

$$-\ln L(w, b) = \ln f_{w,b}(x^{(1)}) + \ln f_{w,b}(x^{(2)}) + \ln\left(1 - f_{w,b}(x^{(3)})\right) \cdots$$

$y^{(k)}$: 1 for class 1, 0 for class 2

$$= \sum_{k=1}^N \left[y^{(k)} \ln f_{w,b}(x^{(k)}) + (1 - y^{(k)}) \ln (1 - f_{w,b}(x^{(k)})) \right]$$

Cross entropy between two Bernoulli distribution

Distribution p :

$$p(x = 1) = y^{(n)}$$

$$p(x = 0) = 1 - y^{(n)}$$



cross

entropy

Distribution q :

$$q(x = 1) = f(x^{(n)})$$

$$q(x = 0) = 1 - f(x^{(n)})$$

$$H(p, q) = - \sum_x p(x) \ln(q(x))$$

Step 2: Goodness of a Function

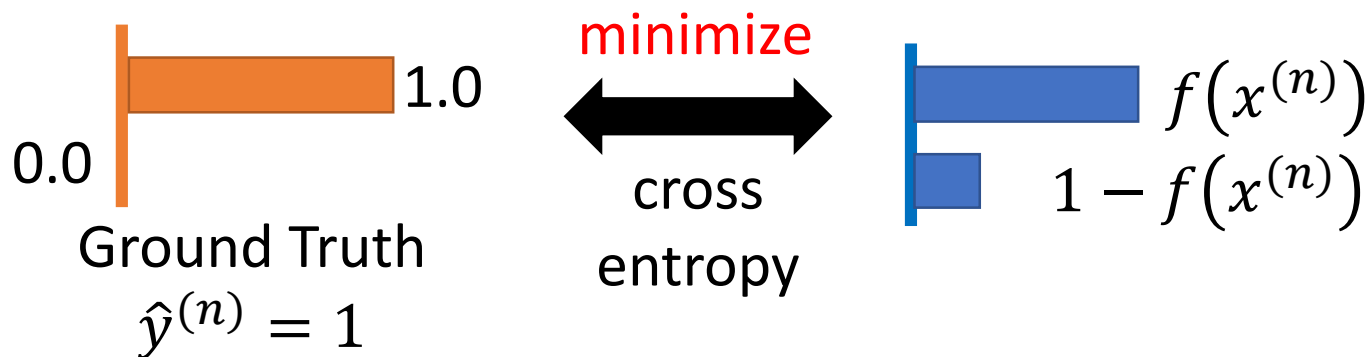
$$L(w, b) = f_{w,b}(x^{(1)})f_{w,b}(x^{(2)})\left(1 - f_{w,b}(x^{(3)})\right)\cdots f_{w,b}(x^{(N)})$$

$$-\ln L(w, b) = \ln f_{w,b}(x^{(1)}) + \ln f_{w,b}(x^{(2)}) + \ln\left(1 - f_{w,b}(x^{(3)})\right) \cdots$$

$y^{(k)}$: 1 for class 1, 0 for class 2

$$= \sum_{k=1}^N - \left[y^{(k)} \ln f_{w,b}(x^{(k)}) + (1 - y^{(k)}) \ln (1 - f_{w,b}(x^{(k)})) \right]$$

Minimize cross entropy between two Bernoulli distribution



Step 3: Find the best function

chain rule

$$\frac{d \ln x}{dx} = \frac{1}{x}$$

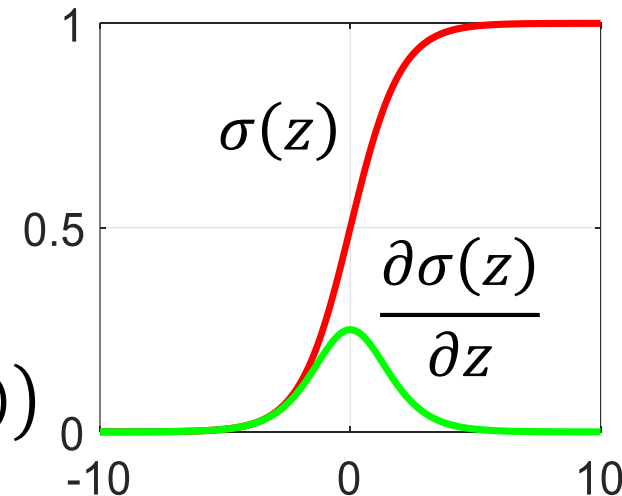
$$\left(1 - f_{w,b}(x^{(k)})\right) x_i^{(k)}$$

$$\frac{-\ln L(w, b)}{\partial w_i} = \sum_n - \left[y^{(k)} \frac{\ln f_{w,b}(x^{(n)})}{\partial w_i} + (1 - y^{(k)}) \ln \left(1 - f_{w,b}(x^{(k)})\right) \right] \frac{\partial}{\partial w_i}$$

$$\frac{\partial \ln f_{w,b}(x)}{\partial w_i} = \frac{\partial \ln f_{w,b}(x)}{\partial z} \frac{\partial z}{\partial w_i}$$

$$\frac{\partial z}{\partial w_i} = x_i$$

$$\frac{\partial \ln \sigma(z)}{\partial z} = \frac{1}{\sigma(z)} \frac{\partial \sigma(z)}{\partial z} = \frac{1}{\cancel{\sigma(z)}} \cancel{\sigma(z)} (1 - \sigma(z))$$



$$f_{w,b}(x) = \sigma(z) = \frac{1}{1 + \exp(-z)}$$

$$z = \mathbf{w}^T \mathbf{x} + b = \sum_i w_i x_i + b$$

Step 3: Find the best function

$$\left(1 - f_{w,b}(x^{(k)})\right) x_i^{(k)}$$

$$-f_{w,b}(x^{(k)}) x_i^{(k)}$$

$$\frac{-\ln L(w, b)}{\partial w_i} = \sum_k - \left[y^{(k)} \frac{\ln f_{w,b}(x^{(k)})}{\partial w_i} + (1 - y^{(k)}) \frac{\ln (1 - f_{w,b}(x^{(k)}))}{\partial w_i} \right]$$

$$\frac{\partial \ln (1 - f_{w,b}(x))}{\partial w_i} = \frac{\partial \ln (1 - f_{w,b}(x))}{\partial z} \frac{\partial z}{\partial w_i} \quad \frac{\partial z}{\partial w_i} = x_i$$

$$\frac{\partial \ln (1 - \sigma(z))}{\partial z} = - \frac{1}{1 - \sigma(z)} \frac{\partial \sigma(z)}{\partial z} = - \frac{1}{1 - \sigma(z)} \sigma(z) (1 - \sigma(z))$$

$$f_{w,b}(x) = \sigma(z) = \frac{1}{1 + \exp(-z)}$$

$$z = \mathbf{w}^T \mathbf{x} + b = \sum_i w_i x_i + b$$

Step 3: Find the best function

$$\frac{-\ln L(w, b)}{\partial w_i} = \sum_k - \left[y^{(k)} \frac{\left(1 - f_{w,b}(x^{(k)})\right) x_i^{(k)}}{\partial w_i} + (1 - y^{(k)}) \frac{-f_{w,b}(x^{(k)}) x_i^{(k)}}{\partial w_i} \right]$$

$$= \sum_n - \left[y^{(k)} \left(1 - f_{w,b}(x^{(k)})\right) x_i^{(k)} - (1 - y^{(k)}) f_{w,b}(x^{(k)}) x_i^{(k)} \right]$$

$$= \sum_k - \left[y^{(k)} - \cancel{y^{(k)} f_{w,b}(x^{(k)})} - f_{w,b}(x^{(k)}) + \cancel{y^{(k)} f_{w,b}(x^{(k)})} \right] x_i^{(k)}$$

$$= \sum_k - \left(y^{(k)} - f_{w,b}(x^{(k)}) \right) x_i^{(k)}$$

Larger difference, larger update

$$w_i \leftarrow w_i - \eta \sum_k - \left(y^{(k)} - f_{w,b}(x^{(k)}) \right) x_i^{(k)}$$

$$P_{w,b}(C_1|x) = f_{w,b}(x) = \sigma(z)$$

Multiclass Classification (3 classes as example)

$$C_1: \mathbf{w}^{(1)}, b_1; \quad z_1 = \mathbf{w}^{(1)} \cdot \mathbf{x} + b_1$$

$$C_2: \mathbf{w}^{(2)}, b_2; \quad z_2 = \mathbf{w}^{(2)} \cdot \mathbf{x} + b_2$$

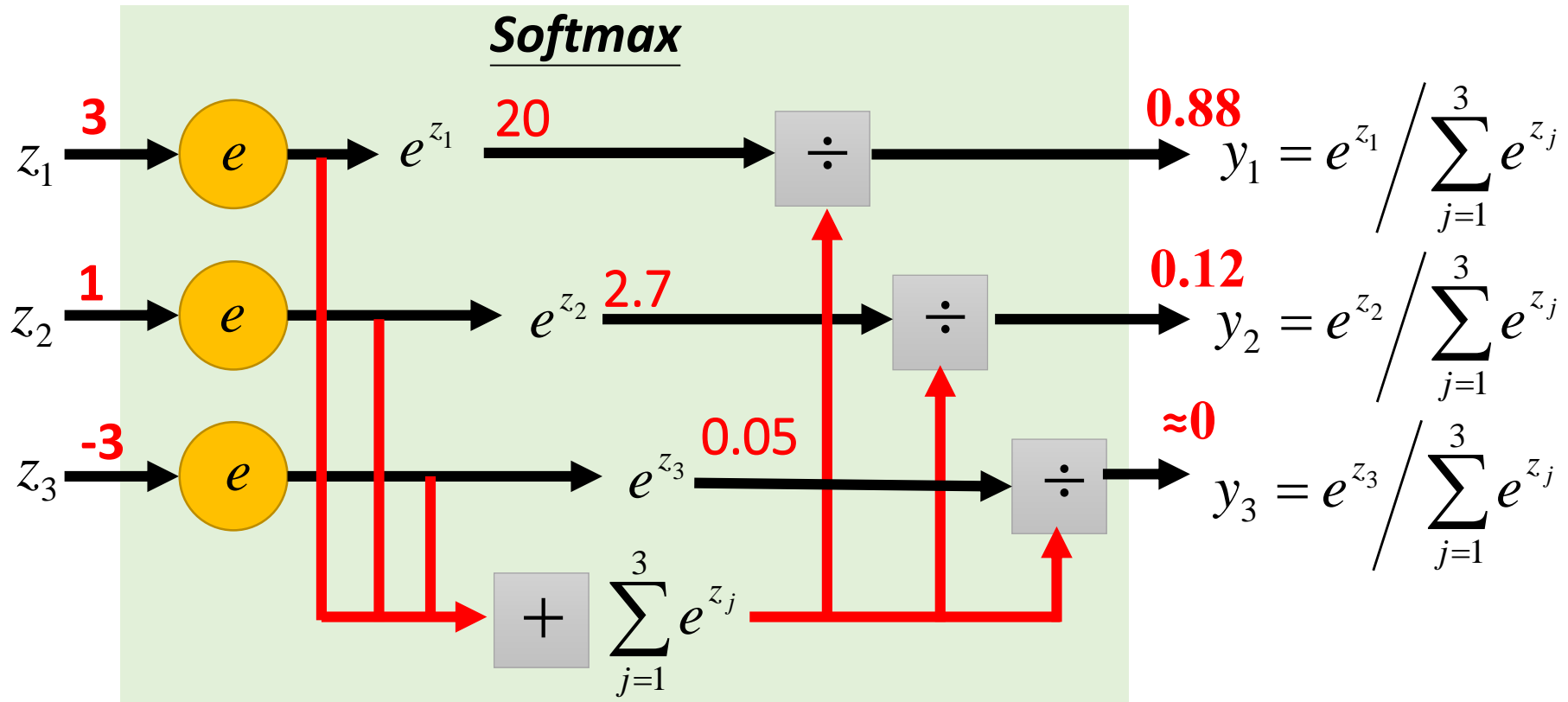
$$C_3: \mathbf{w}^{(3)}, b_3; \quad z_3 = \mathbf{w}^{(3)} \cdot \mathbf{x} + b_3$$

Probability:

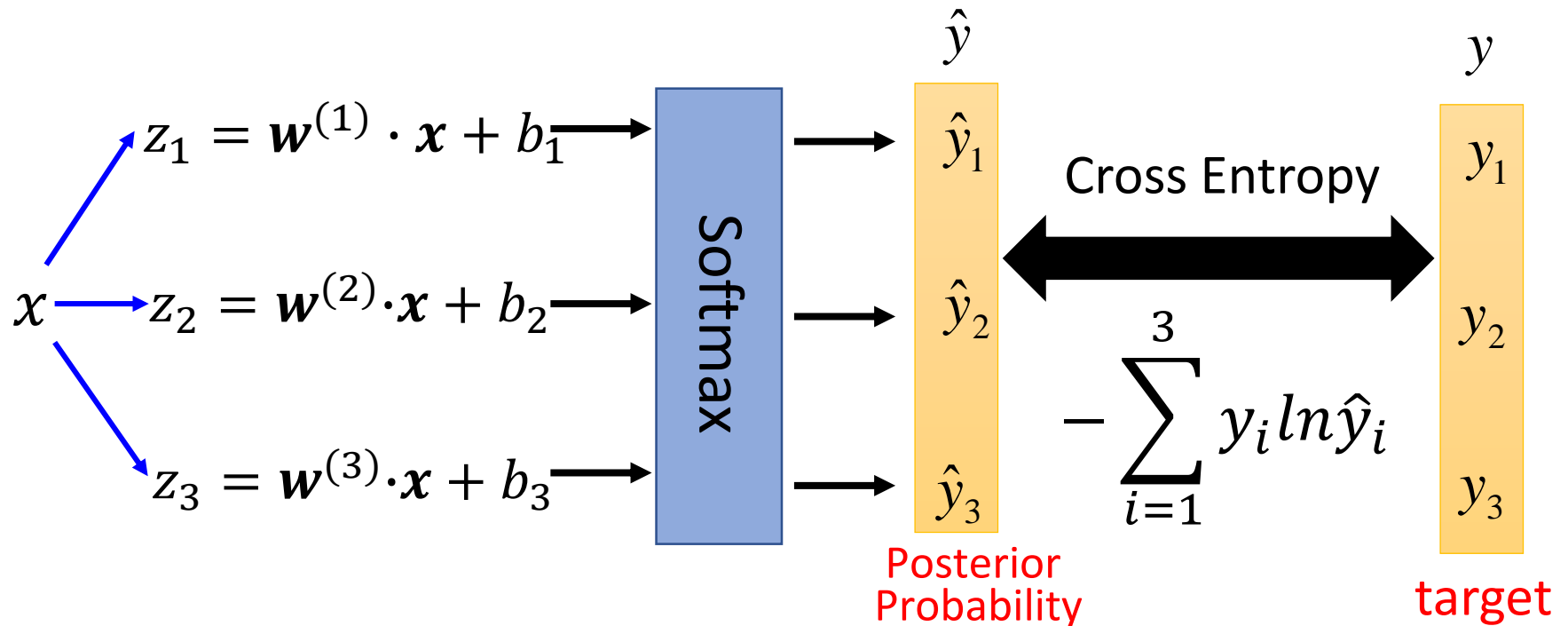
$$\blacksquare 1 > y_i > 0$$

$$\blacksquare \sum_i y_i = 1$$

$$y_i = P(C_i | x)$$



Multiclass Classification (3 classes as example)



If $\mathbf{x} \in \text{class 1}$

$$y = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$$

$$-\ln \hat{y}_1$$

If $\mathbf{x} \in \text{class 2}$

$$y = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$$

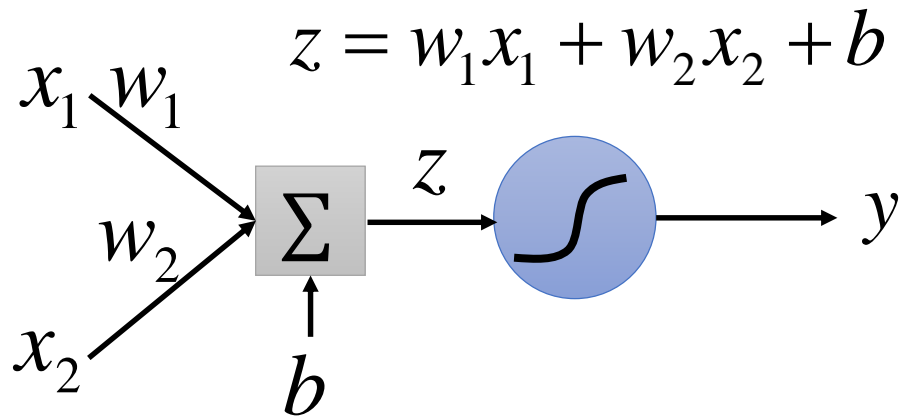
$$-\ln \hat{y}_2$$

If $\mathbf{x} \in \text{class 3}$

$$y = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$

$$-\ln \hat{y}_3$$

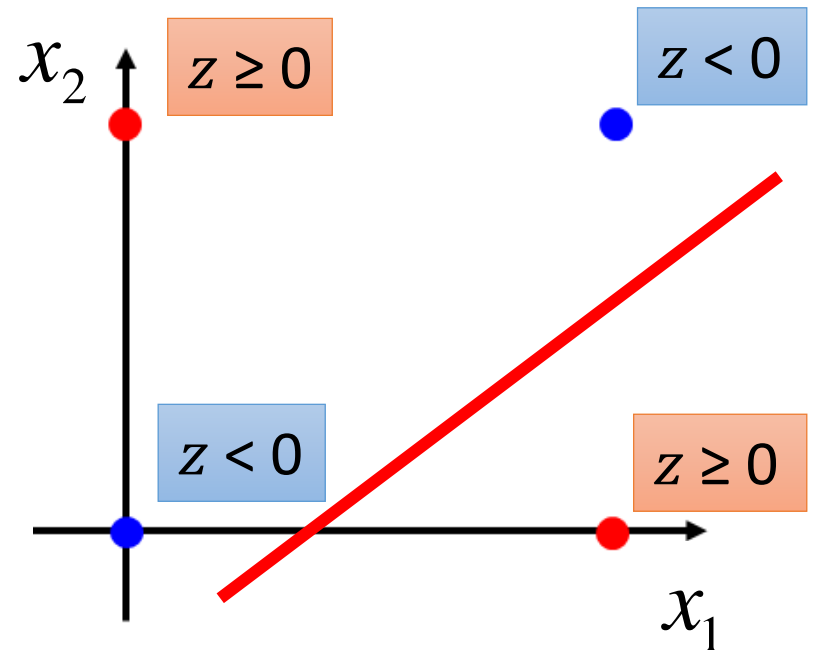
Limitation of Logistic Regression



$$\begin{cases} \text{Class1} & y \geq 0.5 \quad (z \geq 0) \\ \text{Class2} & y < 0.5 \quad (z < 0) \end{cases}$$

Can we?

Input Feature		Label
x_1	x_2	
0	0	Class 2
0	1	Class 1
1	0	Class 1
1	1	Class 2

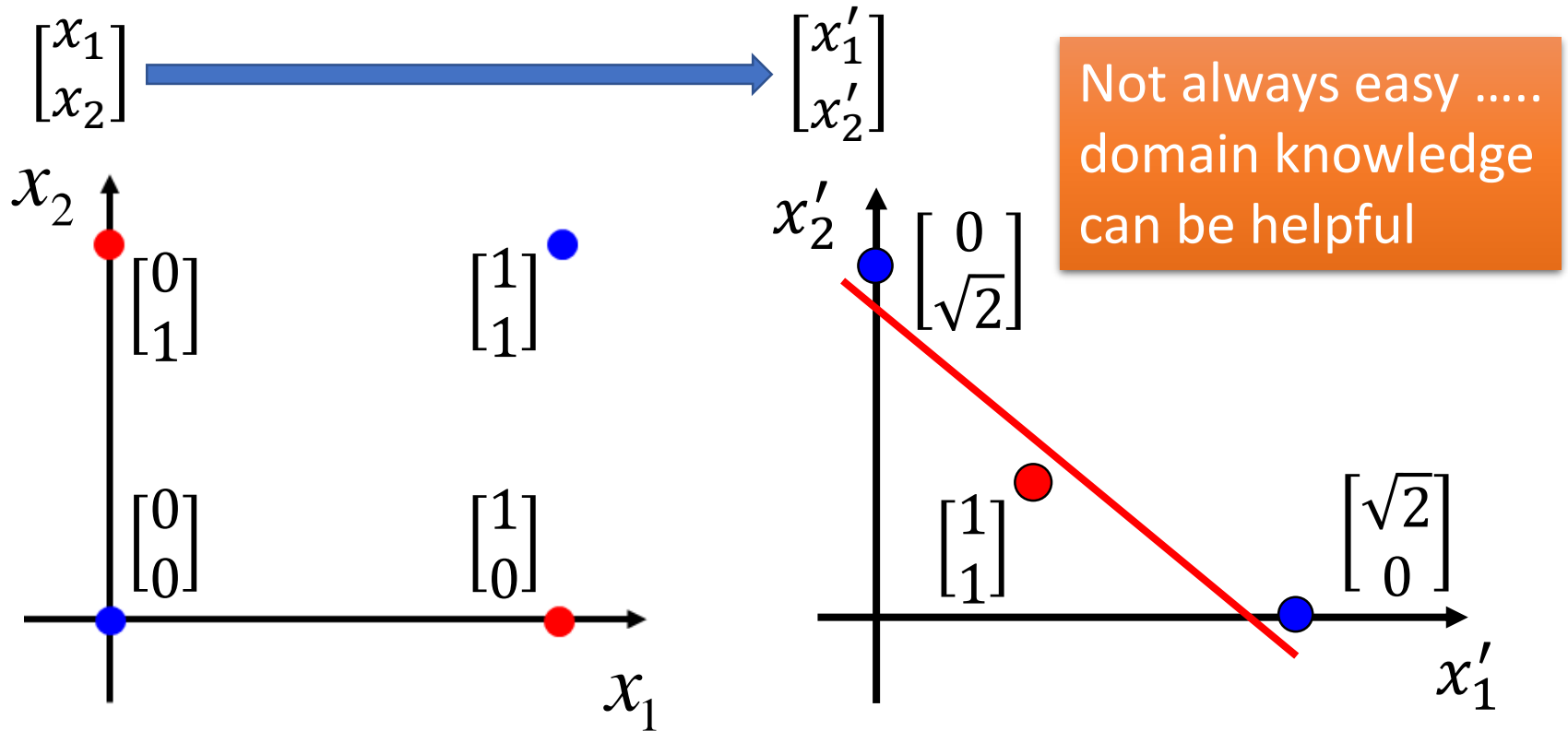


Limitation of Logistic Regression

- Feature Representation

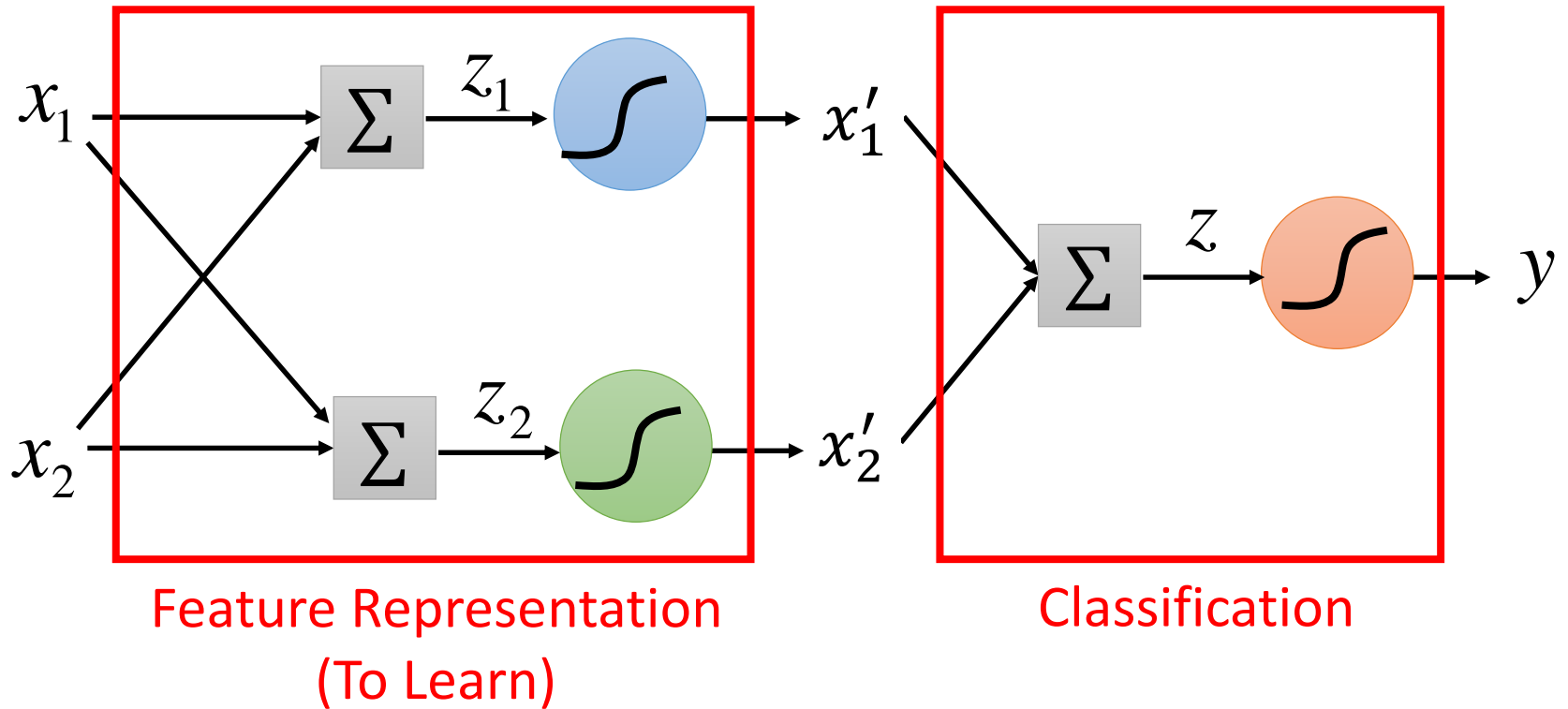
x'_1 : distance to $\begin{bmatrix} 0 \\ 0 \end{bmatrix}$

x'_2 : distance to $\begin{bmatrix} 1 \\ 1 \end{bmatrix}$

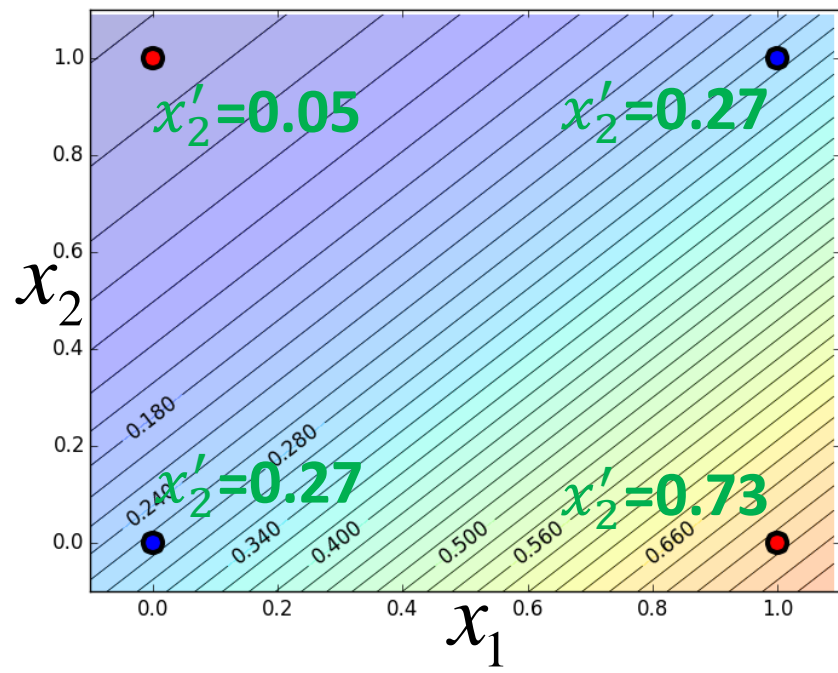
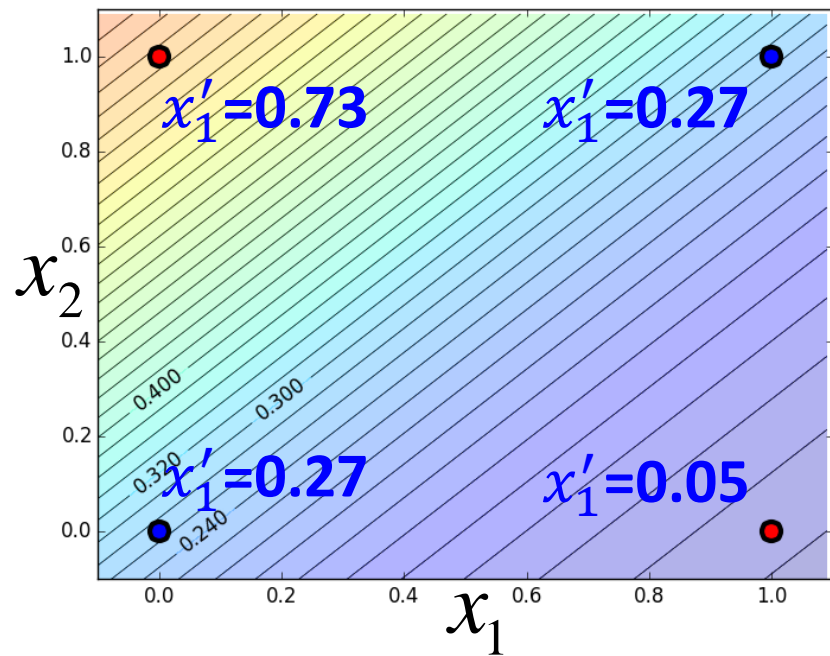
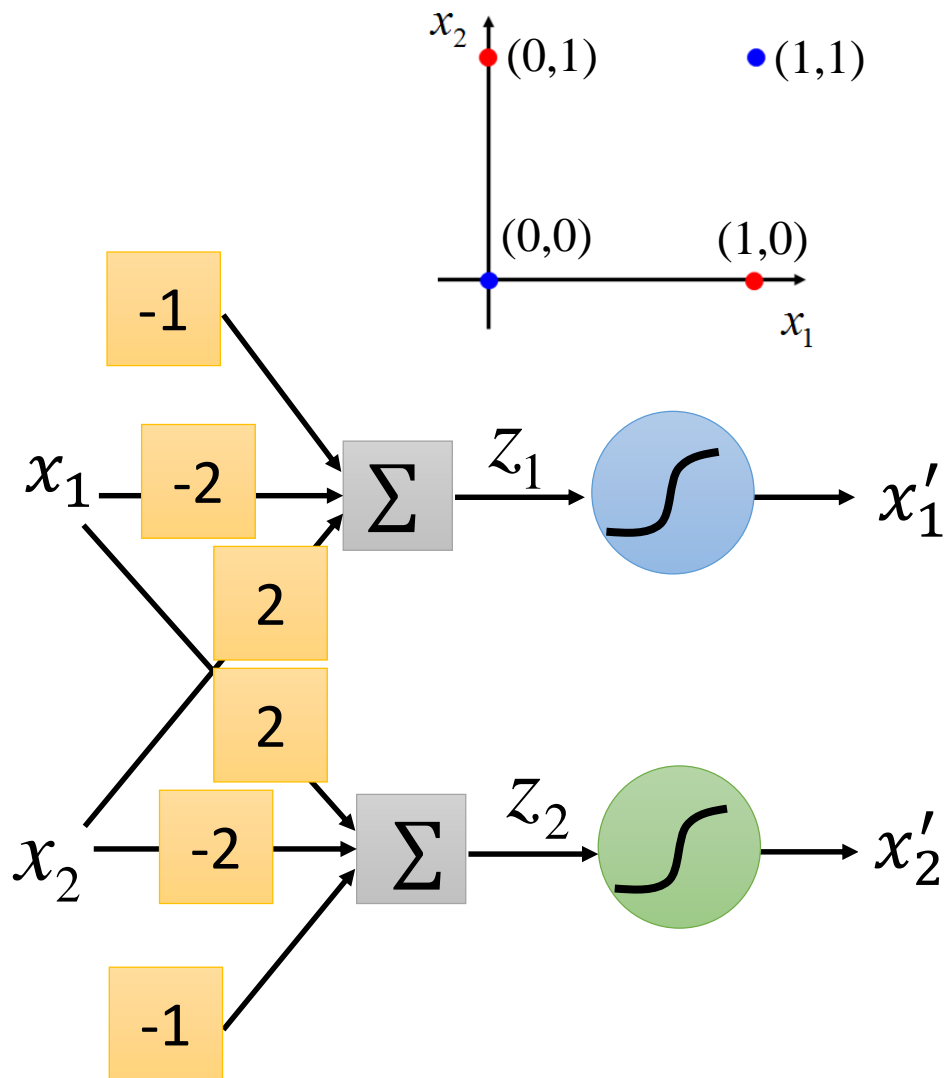


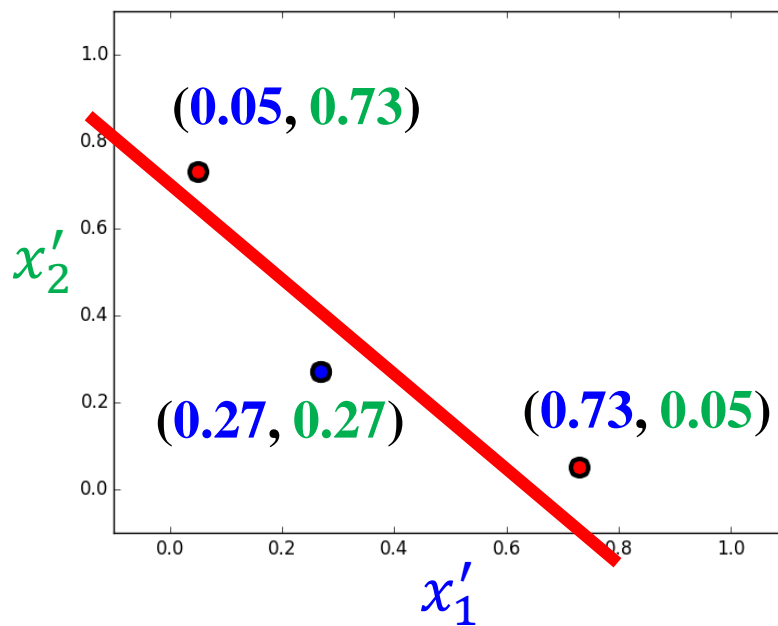
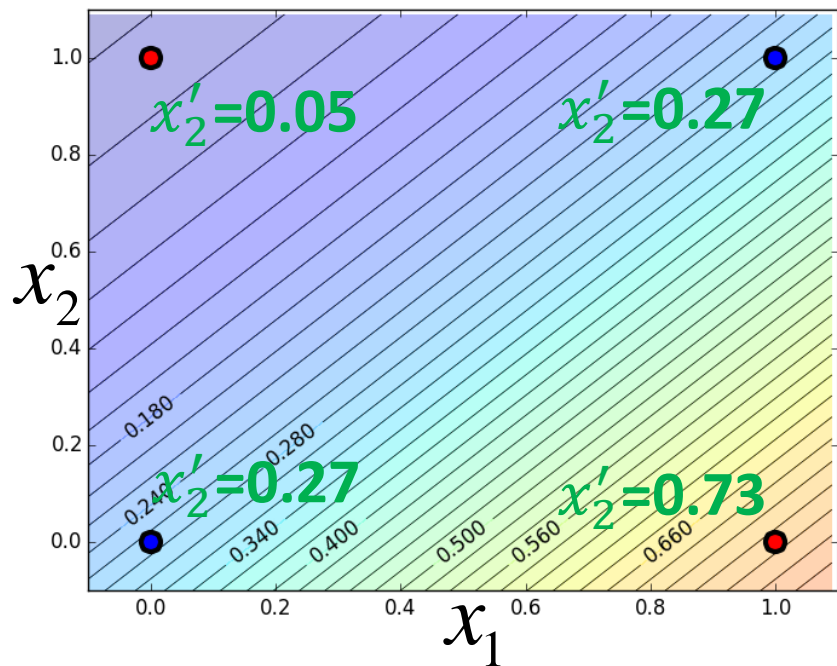
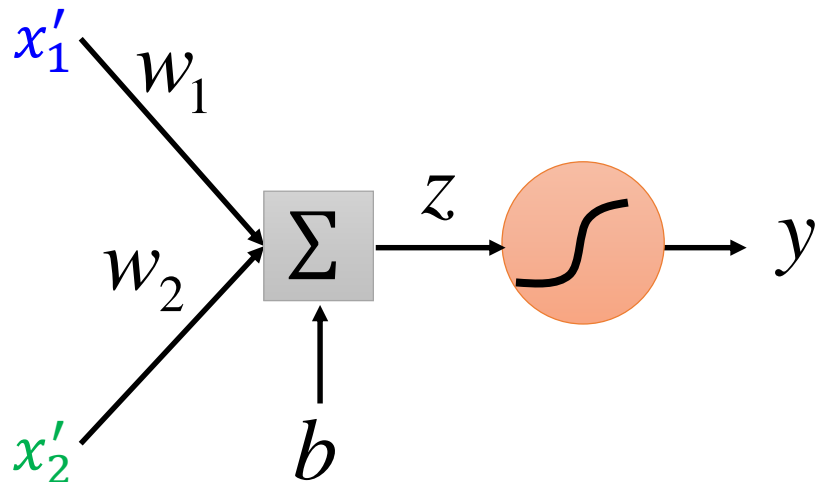
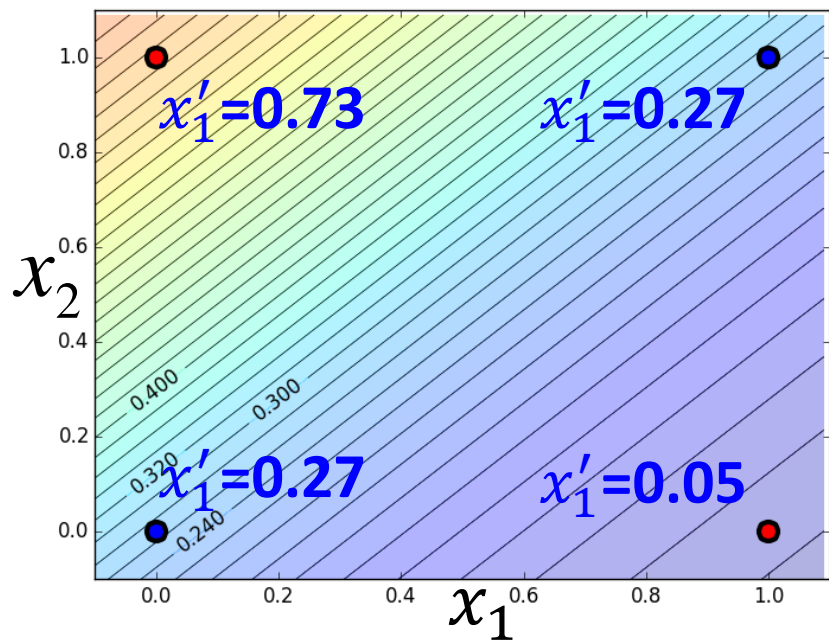
Limitation of Logistic Regression

- Cascading logistic regression models



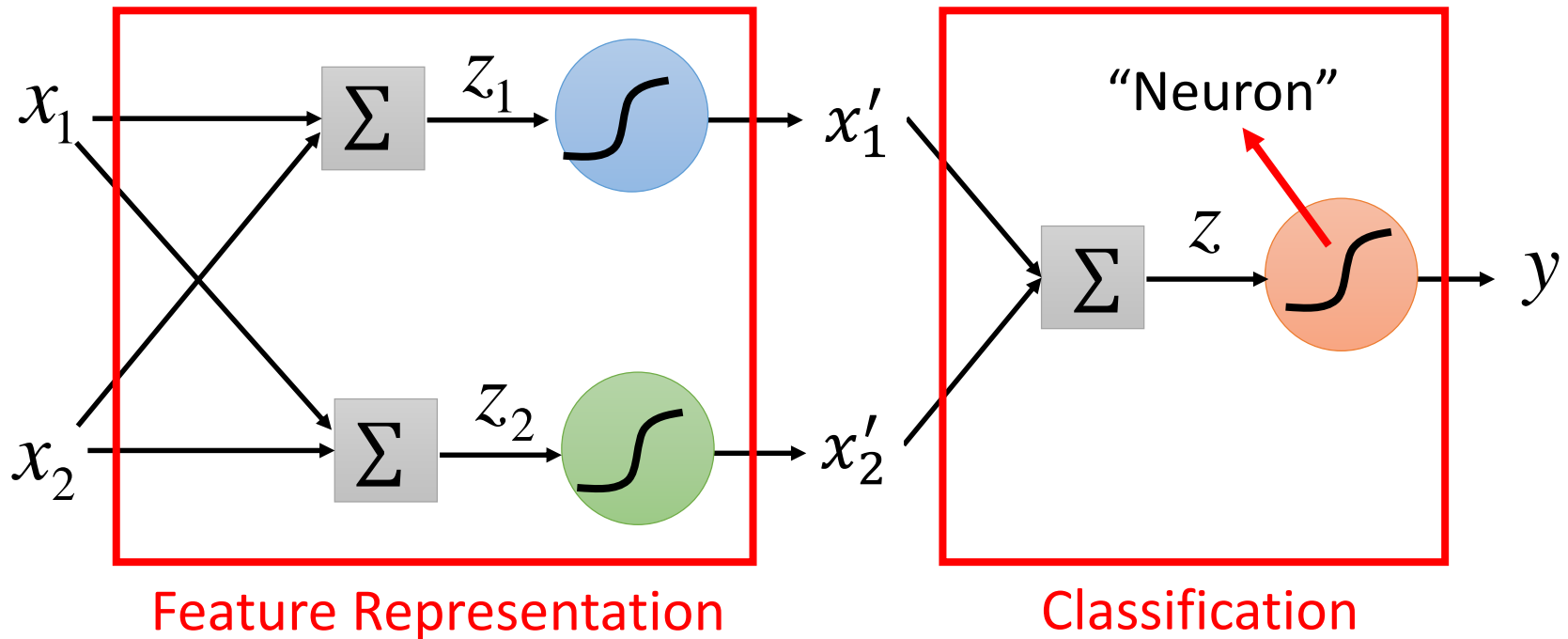
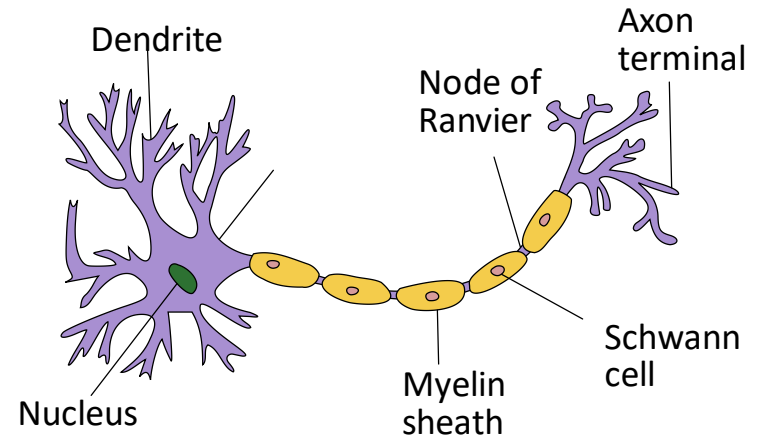
(ignore bias in this figure)





Deep Learning!

All the parameters of the logistic regressions are jointly learned.



Neural Network

Summary and Next Lecture

(Check WebOodi for time)

- Neural Networks
- Multilayer Neural Networks
- Backpropagation

