

Artificial Intelligence Foundations and Applications

Machine Learning

Sudeshna Sarkar

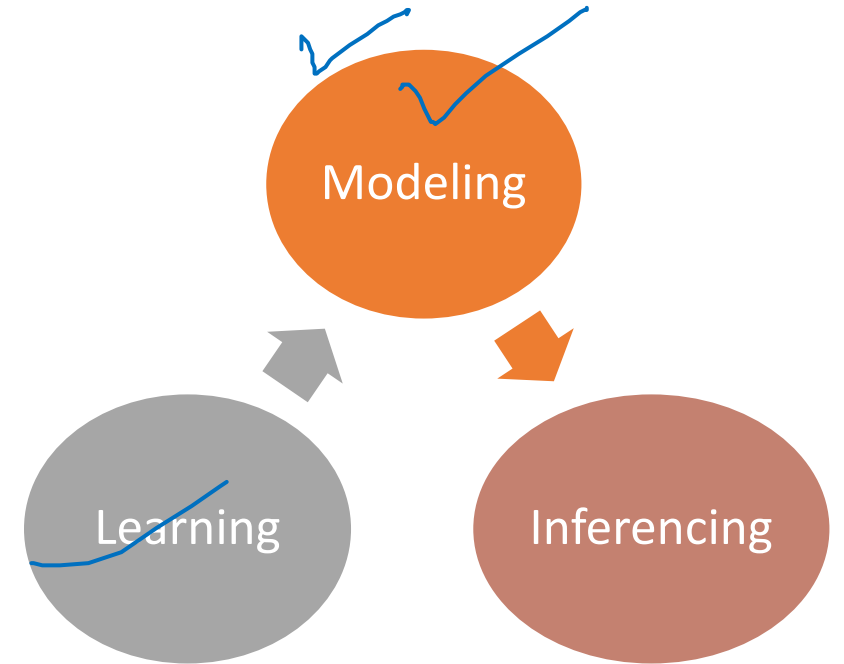
Nov 8 2021

Centre of Excellence in Artificial Intelligence, IIT Kharagpur

Machine Learning : Definition

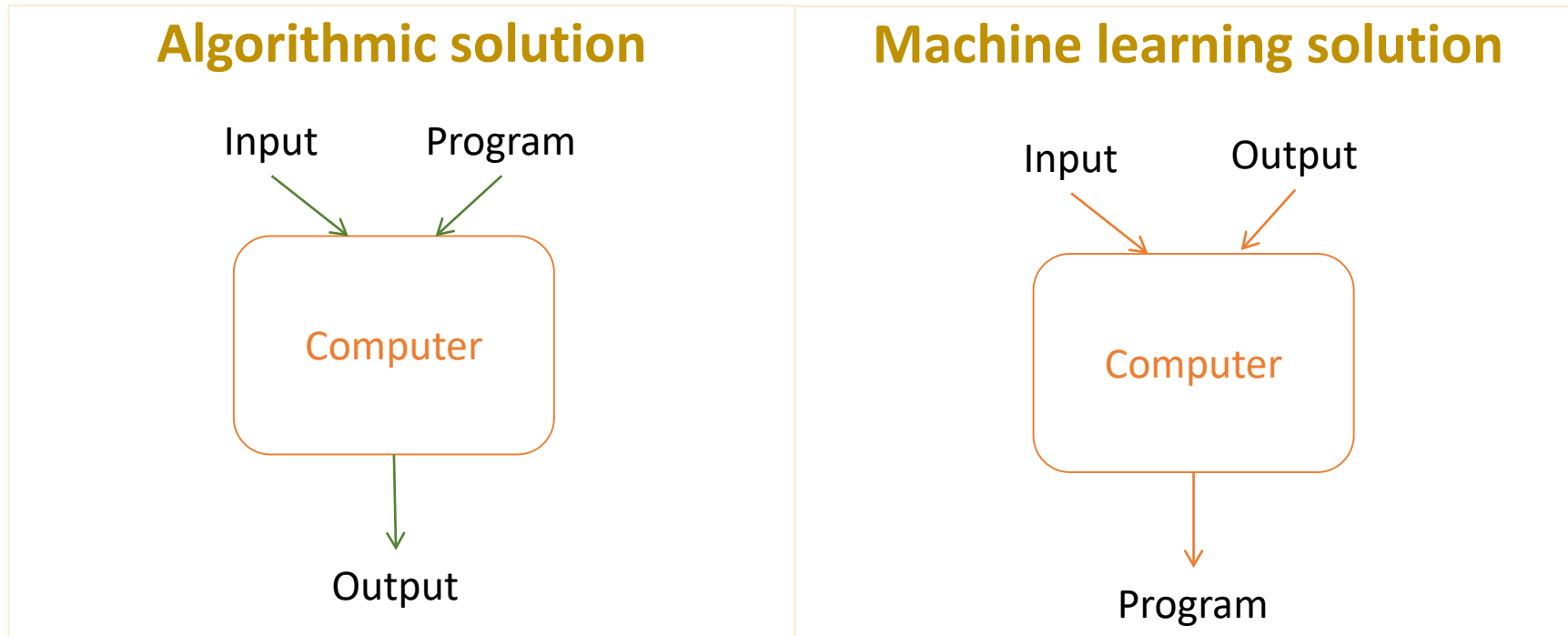
Learning is the ability to evolve behaviours based on data (experience).

- Learn from data such as build a model from data
- Use the model for prediction, decision making or solving some tasks



The Machine Learning Solution

- Collect many examples that specify the correct output for a given input
- ML to get the mapping from input to output



Components of a learning problem

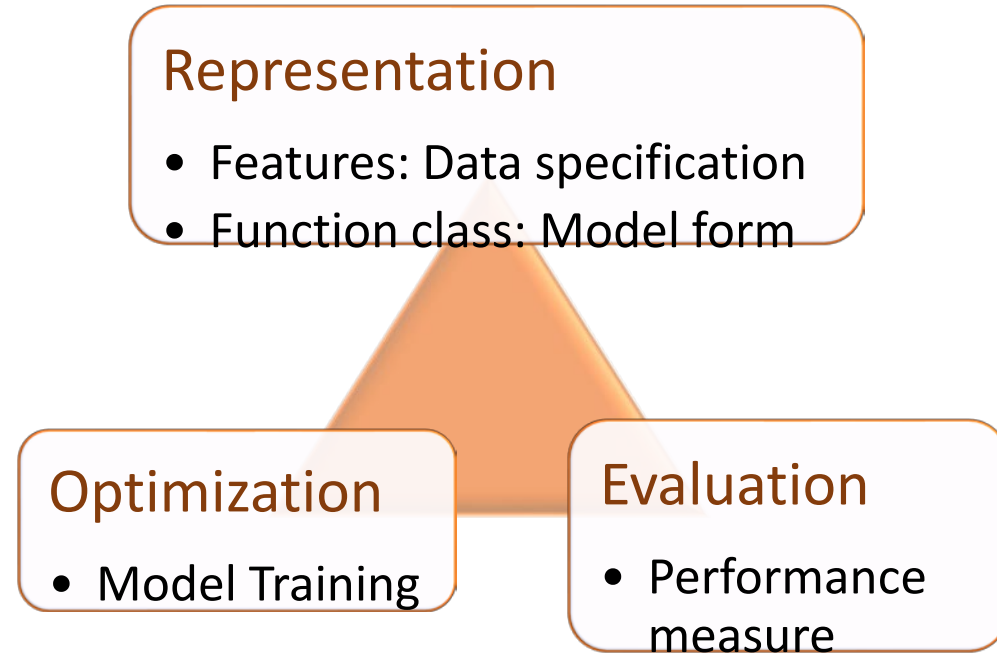
- **Task:** The behaviour or task being improved.
 - For example: classification, acting in an environment
- **Data:** The experiences that are being used to improve performance in the task.
- **Measure of improvement :**
 - For example: increasing accuracy in prediction, improved speed and efficiency

Models for Prediction



Design a Learner

1. Choose the training experience
 - Features
2. Choose the target function f
 - Choose how to represent the target function
3. Choose a learning algorithm to infer the target function



Representation of Data

1. How is the data specified?

A. Features

- Feature vector of n features

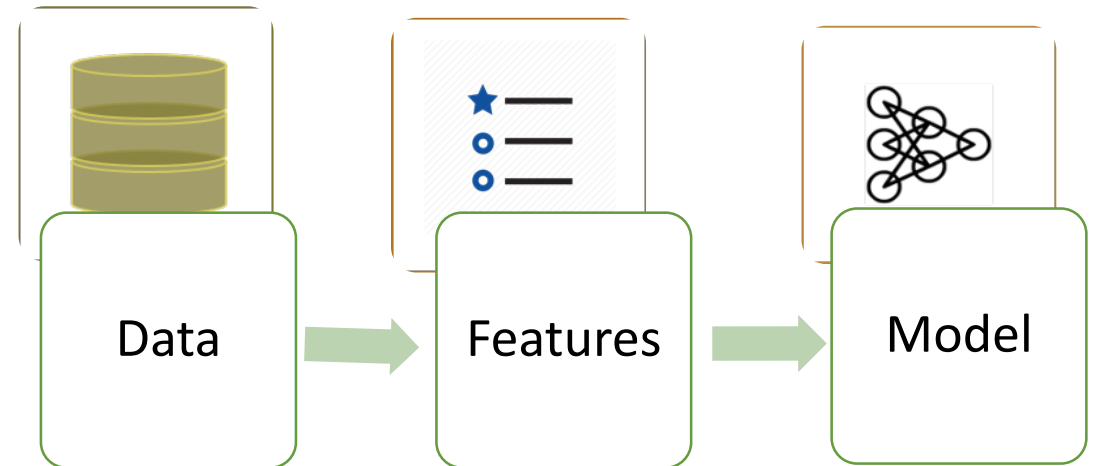
$$\bar{x} = (x_1, x_2, \dots, x_n)$$

B. Convert input to a vector of basis functions

$$(\phi_0(\bar{x}), \phi_1(\bar{x}), \dots, \phi_p(\bar{x}))$$

Feature Choice

- Input Data comprise features
 - Structured features (numerical or categorical values)
 - Unstructured (text, speech, image, video, etc)
- Use only relevant features
- Too many features?
 - Select feature subset (reduction)
 - Extract features: Transform features



1. Model Representation

- The richer the representation, the more useful it is for subsequent problem solving.
 - The richer the representation, the more difficult it is to learn.
- Linear function
 - Decision Tree
 - Graphical Model
 - Neural Network

$$y = f(\bar{x})$$

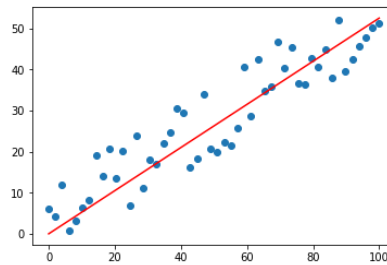
$$y = g(\bar{\phi}(\bar{x}))$$

Model Representation

Hypothesis space

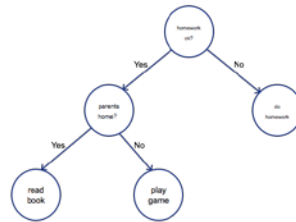
$$y = f(\bar{x})$$

Linear
Function

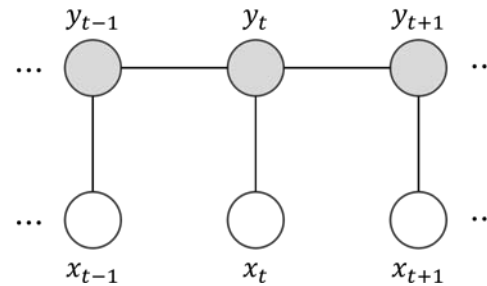


$$Y = w_0 + w_1x_1 + \dots + w_px_p + \epsilon$$

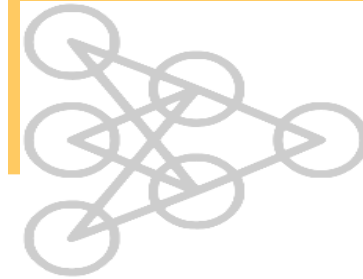
Decision
Tree



Graphical
Model



Neural
Net



2. Evaluation

1. Accuracy = $\frac{\text{\# correctly classified}}{\text{\# all test examples}}$

2. Logarithmic Loss:

$$L_i = -\log(P(Y = y_i | X = x_i))$$

$$L = \sum_{c=1}^M y_{oc} \log(p_{oc})$$

3. Mean Squared error

$$MSE = \frac{1}{m} \sum (y_{pred} - y_{true})^2$$

3. Optimization

1. Define loss function
2. Optimize loss function
 - Stochastic Gradient Descent (Convex functions)
 - Combinatorial optimization
 - E.g.: Greedy search
 - Constrained optimization
 - E.g.: Linear programming

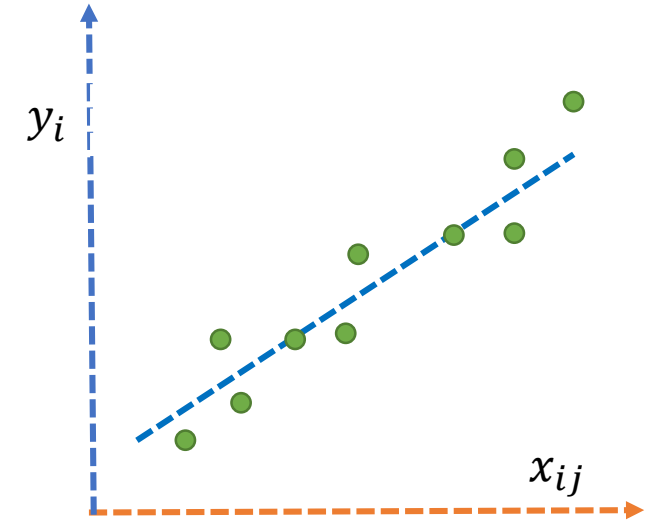
Elements of Optimization

1. Variables
2. Constraints
3. Objective Function

Simple Linear Regression

1. Variables: w_0, w_1, \dots, w_n
2. Constraints: none
3. Objective Function: Minimize

$$\sum_{i=1}^m \left(y_i \left(w_0 + \sum_{j=1}^n w_j x_{ij} \right) \right)^2$$



- m data points, n features
 - x_{ij} : j th attribute of i th instance
 - y_i : output of i th instance

Find coefficients
 w_0, w_1, \dots, w_n to best fit
data

Broad types of machine learning

- **Supervised Learning**

- Training Data with labels: X, y (pre-classified)
- Given an observation x , what is the best label for y ?

- **Unsupervised learning**

- Training Data without labels: X
- Given a set of x 's, find hidden structure

- **Semi-supervised Learning**

- Training Data + some Labels

- **Reinforcement Learning**

- Given: observations and periodic rewards as the agent takes sequential action in an environment
- Determine optimum policy

Supervised Learning

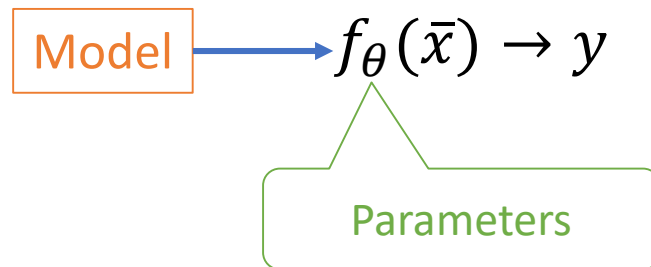
Given data containing the inputs and outputs:

Training Data:

$$\{(\bar{x}_1, y_1), (\bar{x}_2, y_2), \dots, (\bar{x}_m, y_m)\}$$

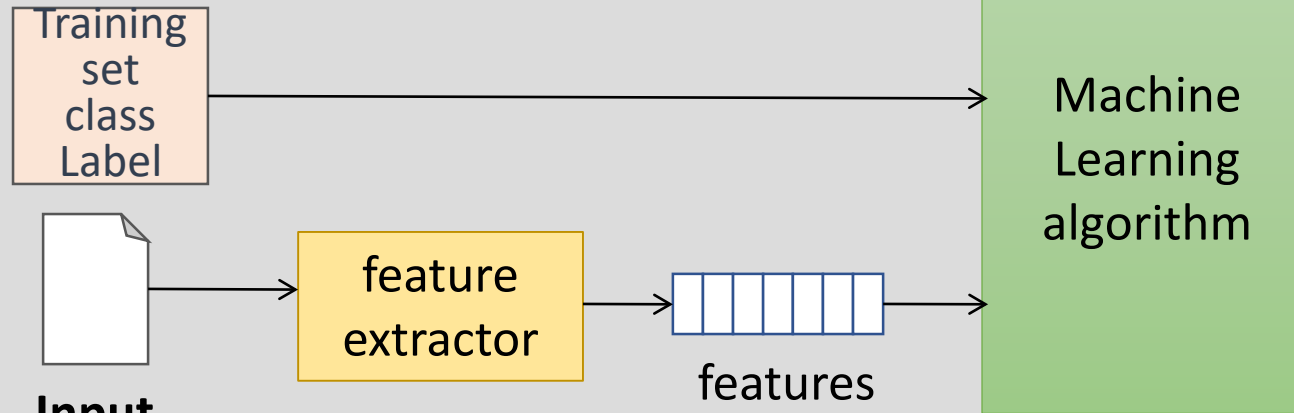
- Learn a function $f(x)$ to predict y given x

| \bar{X} | Y |
|-------------|-------|
| \bar{x}_1 | y_1 |
| \bar{x}_2 | y_2 |
| ... | .. |
| \bar{x}_m | y_m |

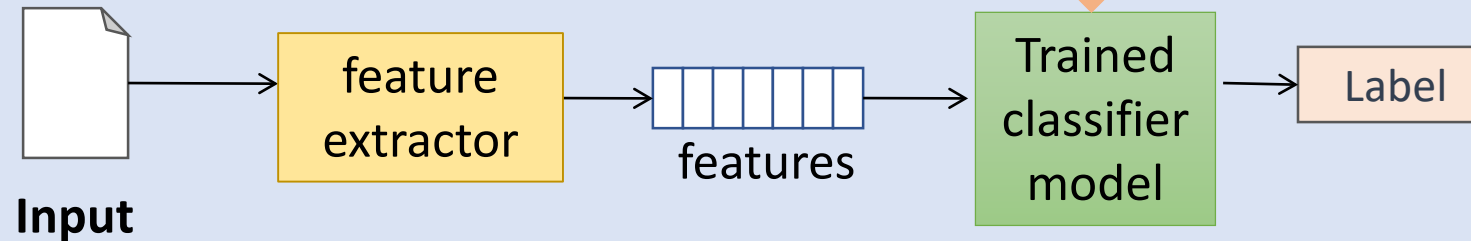


- Training: Learn the model from the Training Data
- Given Test instance \bar{x}' , predict $y' = f_{\theta}(\bar{x}')$

Training phase



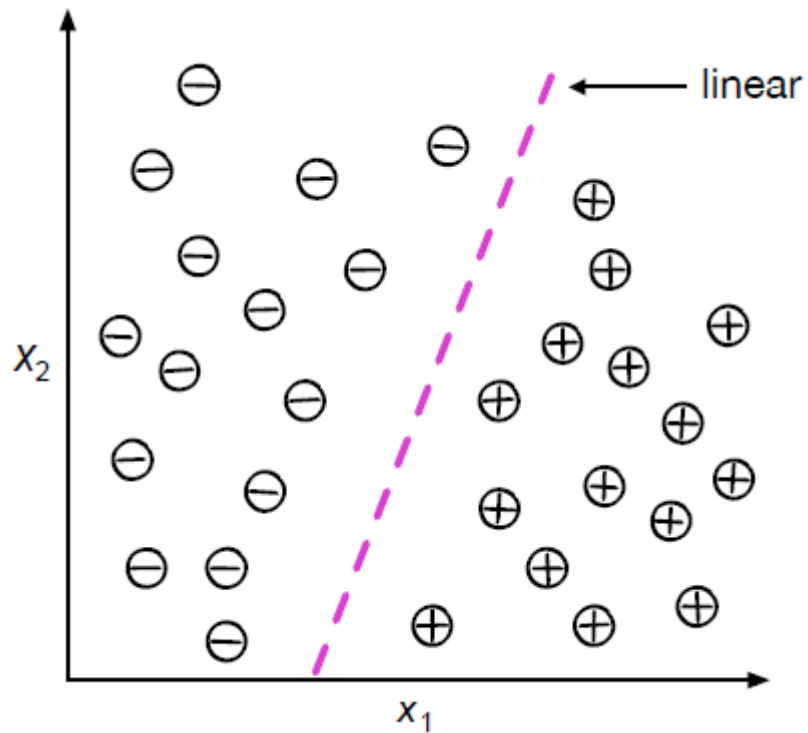
Testing Phase



Supervised Learning

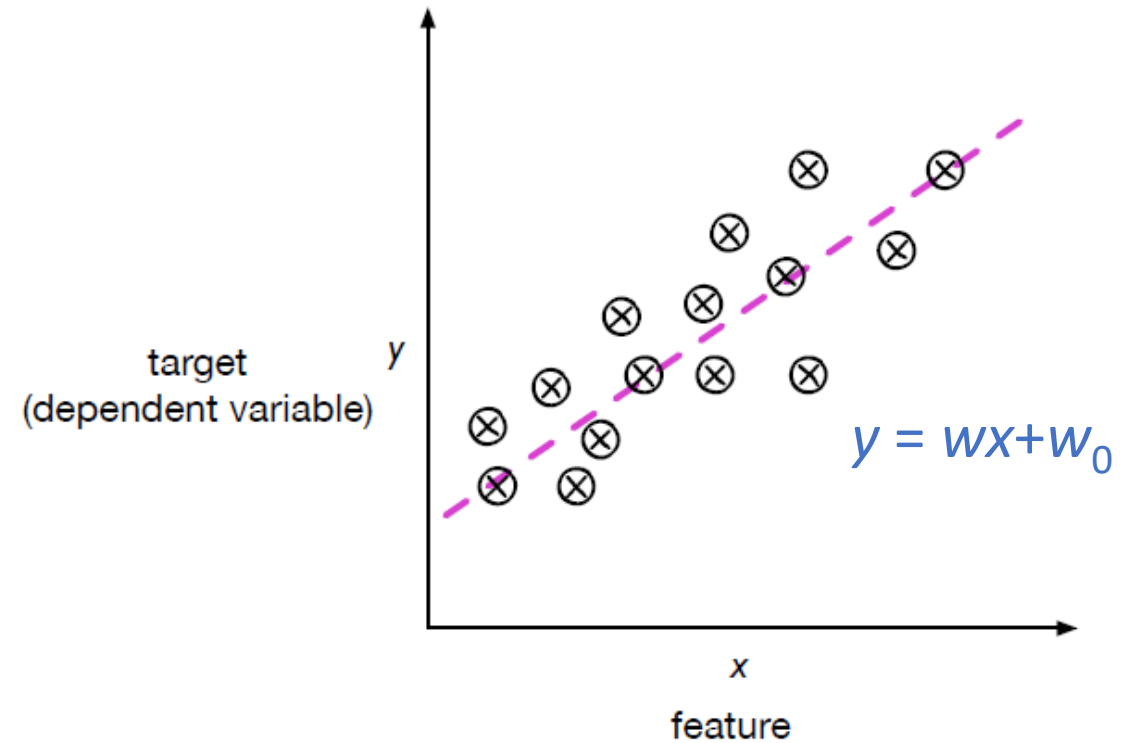
Classification

Y is categorical/ discrete



Regression

Y is numeric / continuous



Example Tasks

Classification

- Object identification from images
- Defect classification
- Credit card transaction fraud or not

Regression

- House price prediction
- Remaining Useful Life Prediction
- Probability of developing cracks
- Demand forecasting

Structured Prediction

- Machine translation: English sentence → Japanese sentence
- Dialogue: conversational history → next utterance
- Image captioning: image → sentence describing image
- Image segmentation: image → segmentation

Supervised Learning

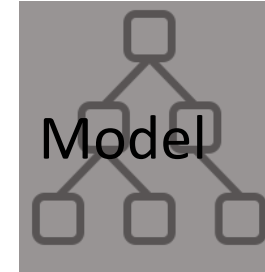
Classification Example

Training
Samples

| | | | | |
|--|--|--|--|--|
| | | | | |
| | | | | |
| | | | | |
| | | | | |
| | | | | |

Test
Instances

| | | | | |
|--|--|--|--|--|
| | | | | |
| | | | | |
| | | | | |
| | | | | |
| | | | | |



Train a model to minimize loss

Probabilistic Classification

| X1 | X2 | X3 | X4 | Category |
|----|----|----|----|----------|
| | | | | Type 1 |
| | | | | Type 3 |
| | | | | Type 1 |
| | | | | Type 2 |

Predict a probability distribution
over the set of classes $\Pr(Y|X)$

| X1 | X2 | X3 | X4 | Type 1 | Type 2 | Type 3 |
|----|----|----|----|--------|--------|--------|
| | | | | 0.4 | 0.15 | 0.45 |
| | | | | 0.2 | 0.7 | 0.1 |
| | | | | 0.1 | 0.2 | 0.7 |
| | | | | 0.5 | 0.1 | 0.4 |

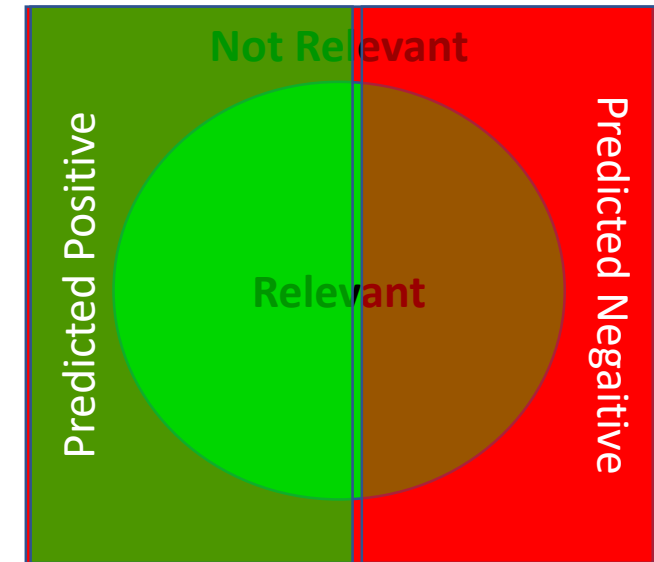
Evaluation for Classification problems

$$\text{Accuracy} = \frac{\# \text{ correctly classified}}{\# \text{ all test examples}}$$

$$= \frac{\# \text{ predicted true } pos + \# \text{ predicted true } neg}{\# \text{ all test examples}}$$

$$\text{Precision} = \frac{\# \text{ predicted true } pos}{\# \text{ predicted pos}}$$

$$\text{Recall} = \frac{\# \text{ predicted true } pos}{\# \text{ True } pos}$$



| | | True Class | |
|-----------------|-----|------------|-----|
| | | Pos | Neg |
| Predicted Class | Pos | TP | FP |
| | Neg | FN | TN |

Loss Function Classification problems

Loss indicates how bad the model's prediction is.

1. Fraction of Misclassifications

$$Error = \sum_{i=1}^m \frac{I(y_i \neq \hat{y}_i)}{m}$$

2. Logarithmic Loss: Maximize the log likelihood. For a loss function, minimize the negative log likelihood of the correct class:

$$L_i = -\log(P(Y = y_i | X = x_i))$$

Logarithmic Loss Function

Logarithmic Loss:

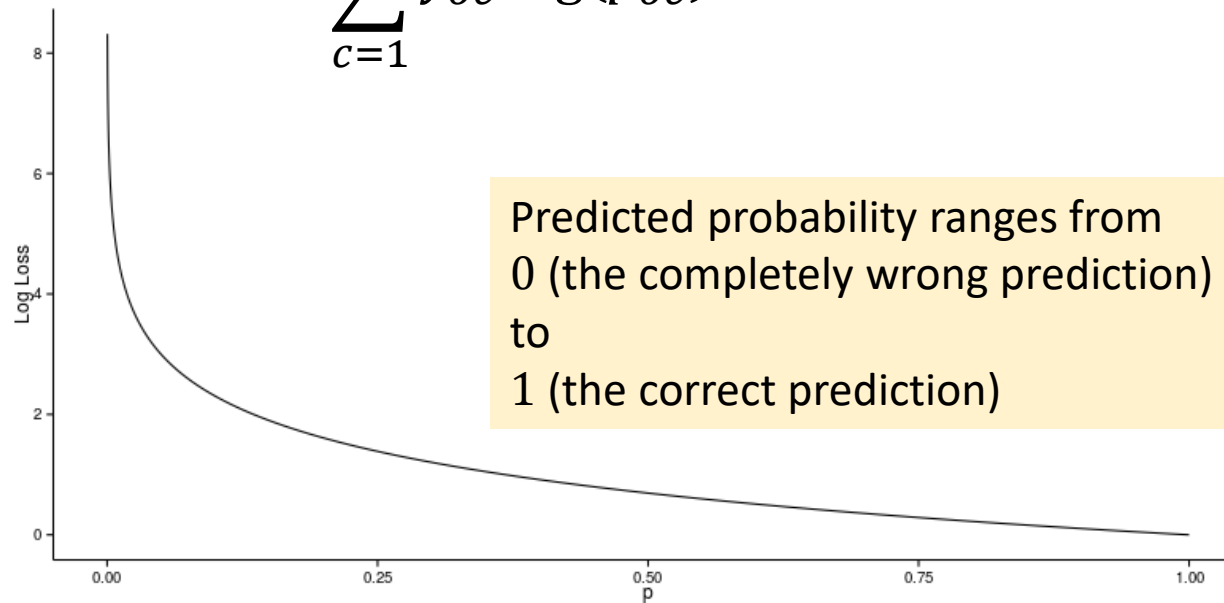
$$L_i = -\log(P(Y = y_i | X = x_i))$$

$$L = \sum_{c=1}^M y_{oc} \log(p_{oc})$$

M - number of classes

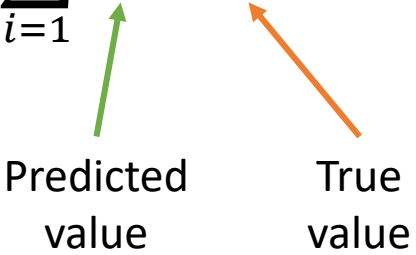
y - binary indicator (0 or 1) if class label **c** is the correct classification for observation **o**

p - predicted probability observation **o** is of class **c**



2. Evaluation for regression problem

Mean Squared error

$$MSE = \frac{1}{m} \sum_{i=1}^m (\hat{y}_i - y_i)^2$$


Predicted value True value

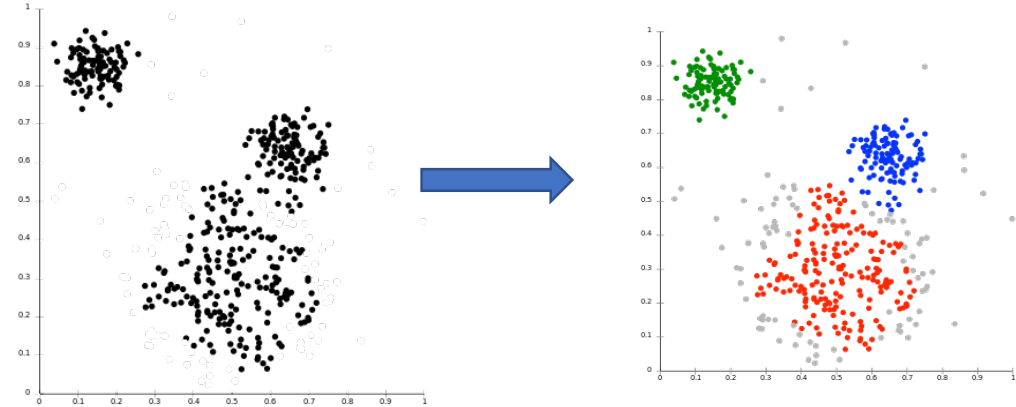
Unsupervised Learning (Clustering)

Given $\{\overline{x_1}, \overline{x_2}, \dots \overline{x_m}, \}$ without labels

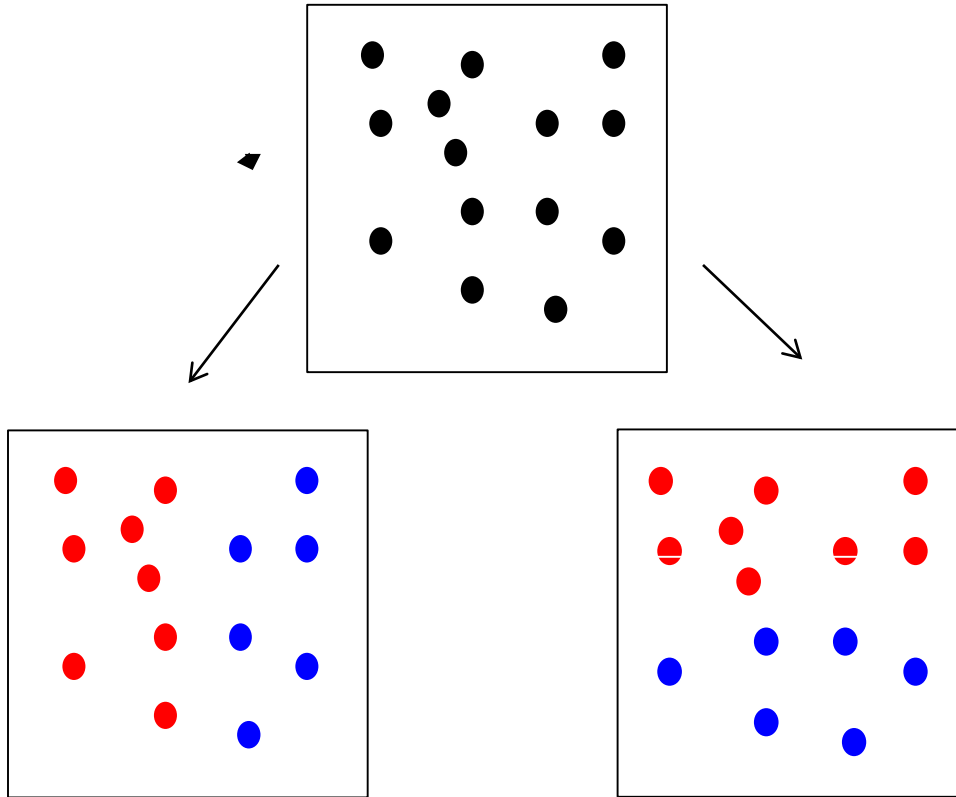
Find hidden structure in the data

- Clustering
- Dimensionality Reduction

Clustering: Grouping similar objects



Clustering Problems



How to evaluate clustering?

- Internal Evaluation:
 - Intra-cluster distances are minimized
 - Inter-cluster distances are maximized
- External Evaluation

Linear Models

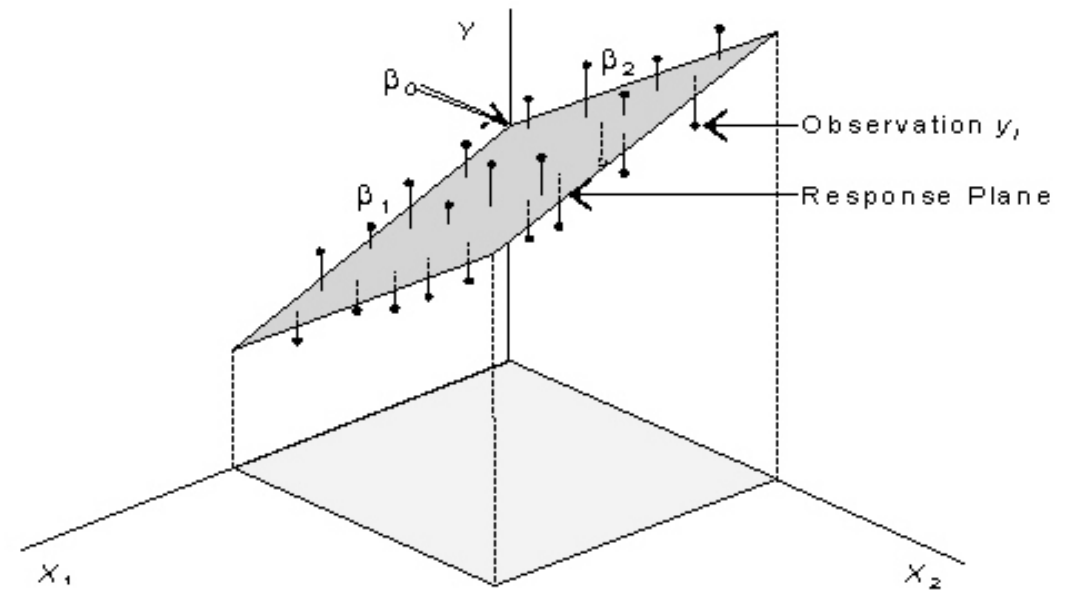
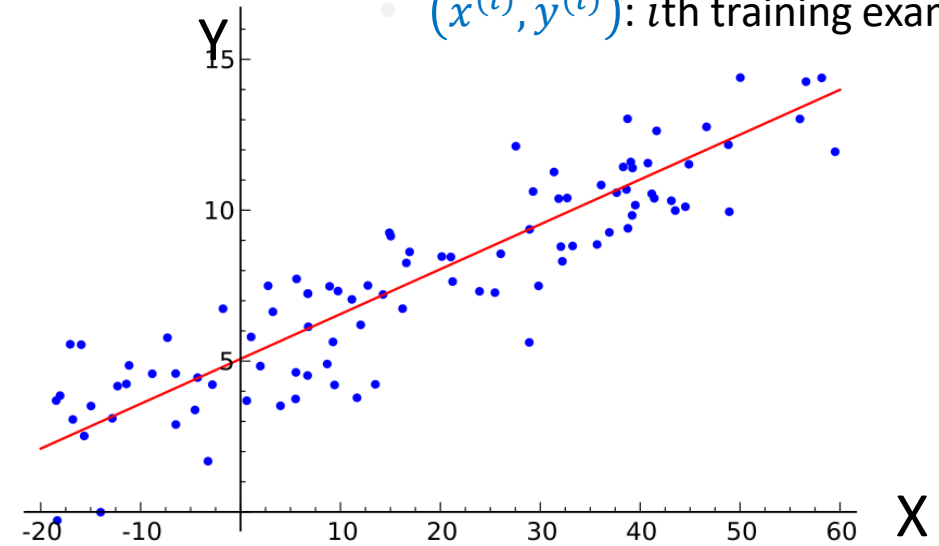
Linear Regression

Understand the relationship between dependent variable x and explanatory variable y

Objective: predict y from x

Linear Model: $Y = \theta X + \theta_0$

- n features
- m training examples
- $(x^{(i)}, y^{(i)})$: i th training example

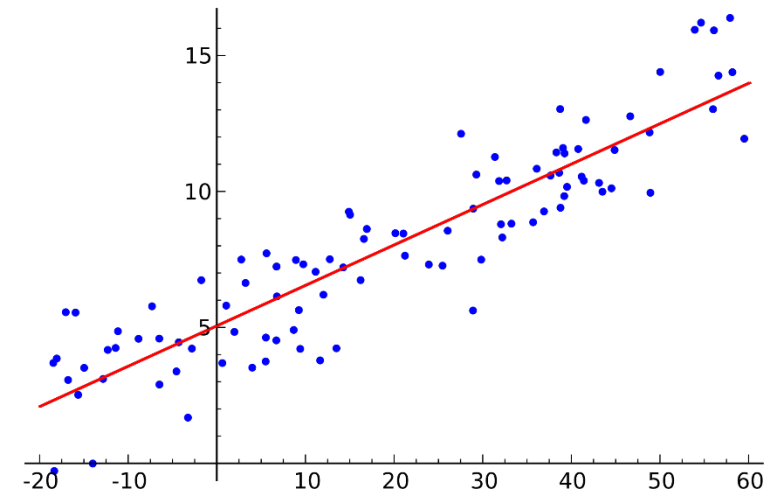


Linear hypothesis function: Intuition

$$h_{\theta}(\mathbf{x}) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_n x_n$$

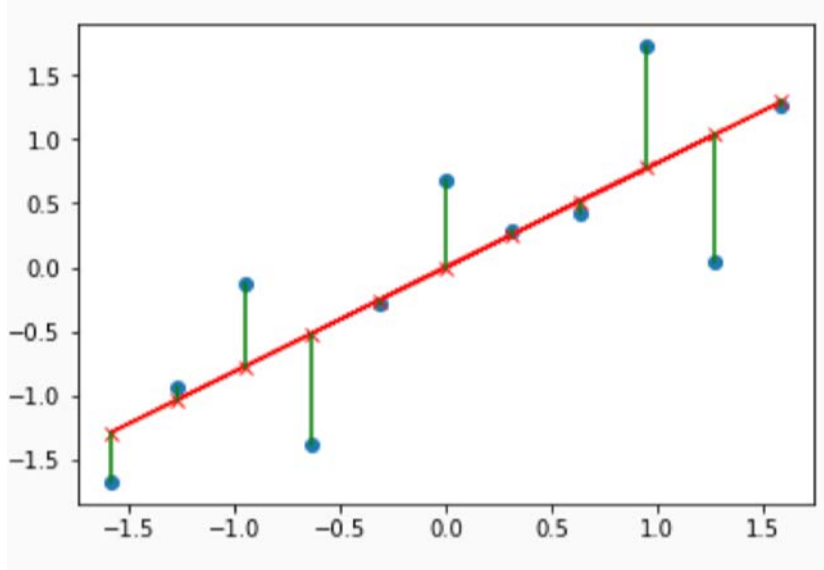
Two equivalent questions:

1. Which is the best straight line to fit the data?
2. How to learn the values of the parameters θ_i ?



Cost function

$e^{(i)} = \widehat{y^{(i)}} - y^{(i)} = h_{\theta}(x^{(i)}) - y^{(i)}$:
prediction error for i th training example



$$e^{(i)} = h_{\theta}(x^{(i)}) - \widehat{y^{(i)}}$$

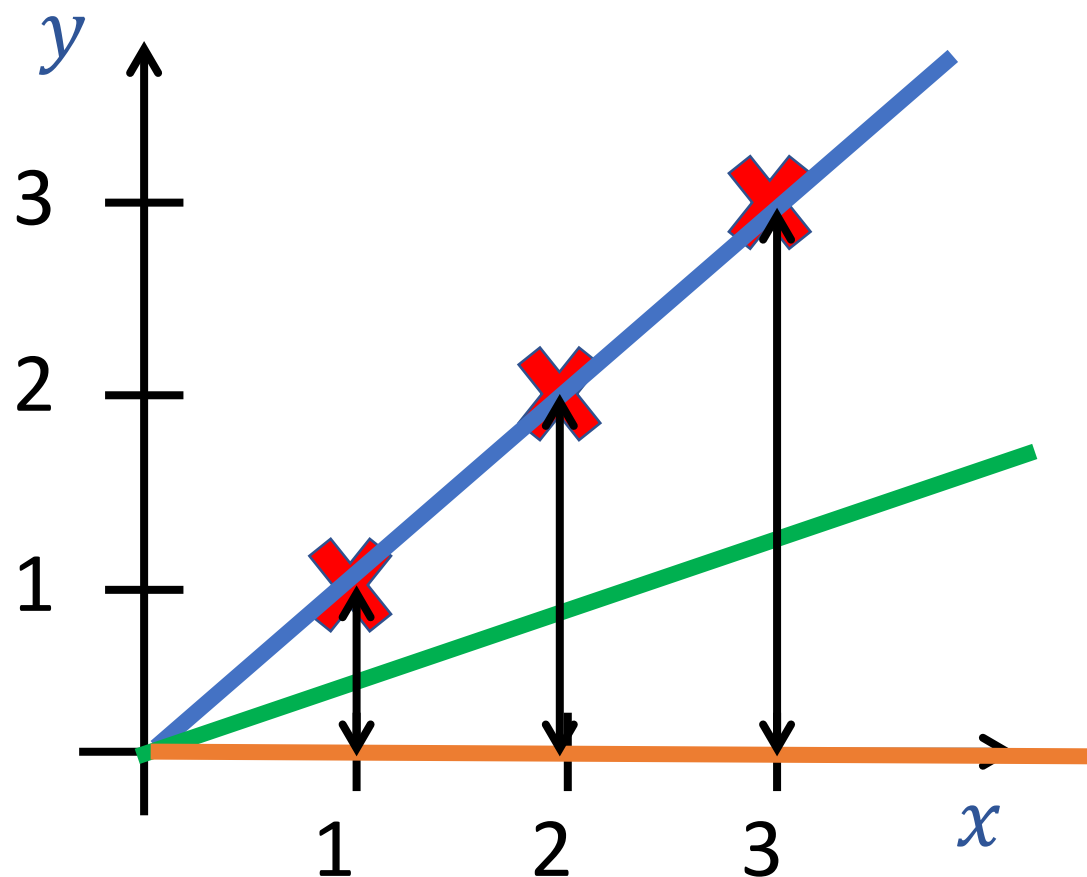
$$J(\bar{\theta}) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

Choose parameters $\bar{\theta}$ so that

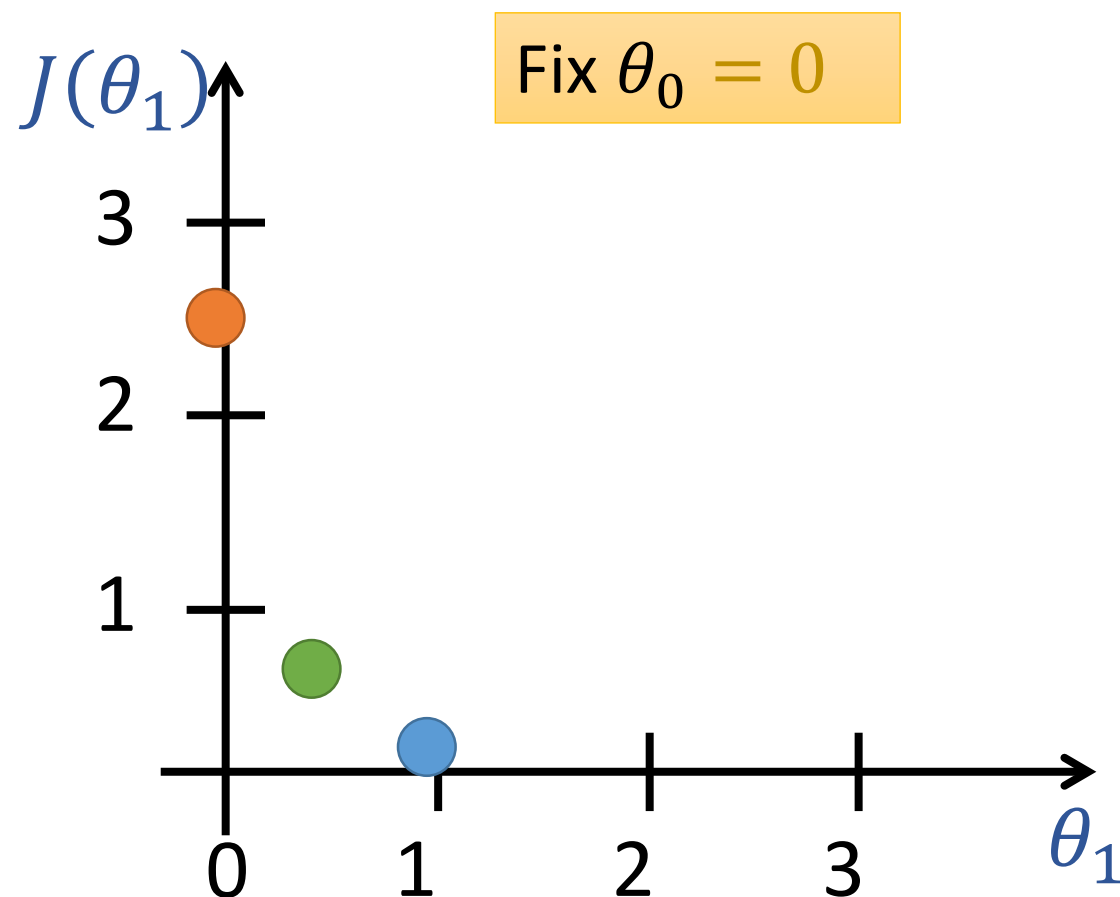
$J(\bar{\theta})$ is minimized

$$\underset{\bar{\theta}}{\text{minimize}} J(\bar{\theta})$$

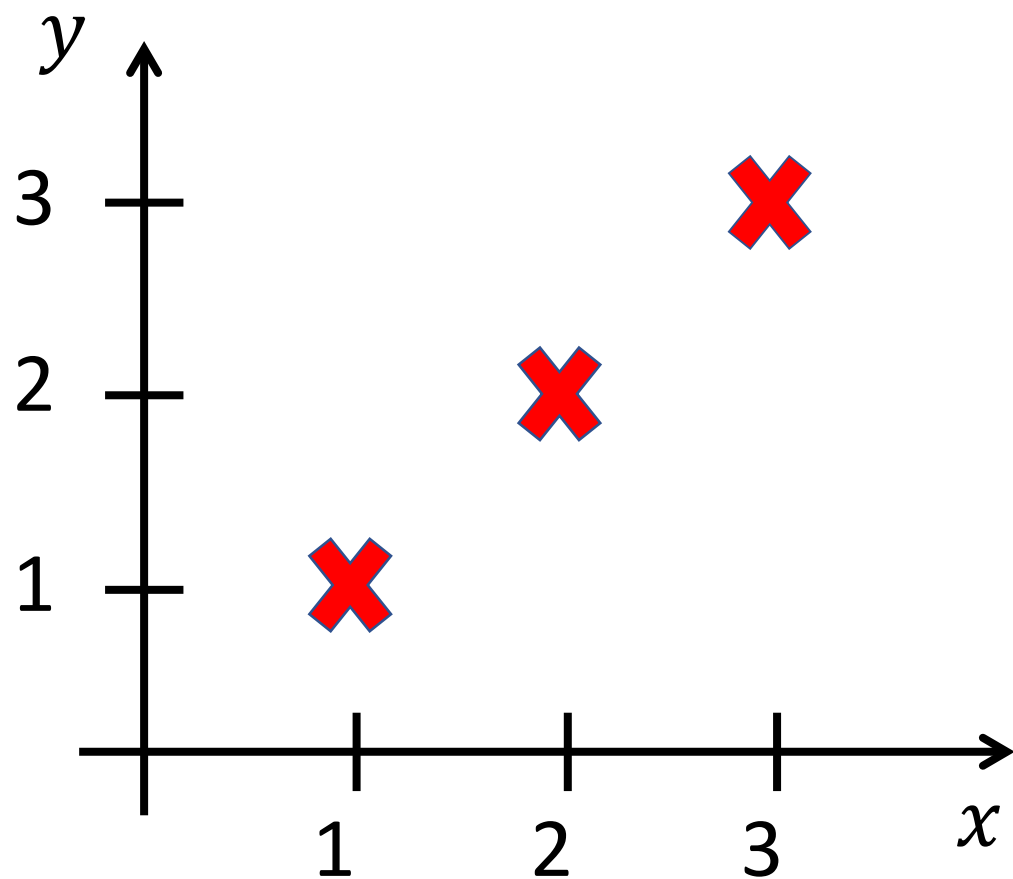
$h_{\theta}(x)$: function of x for fixed θ



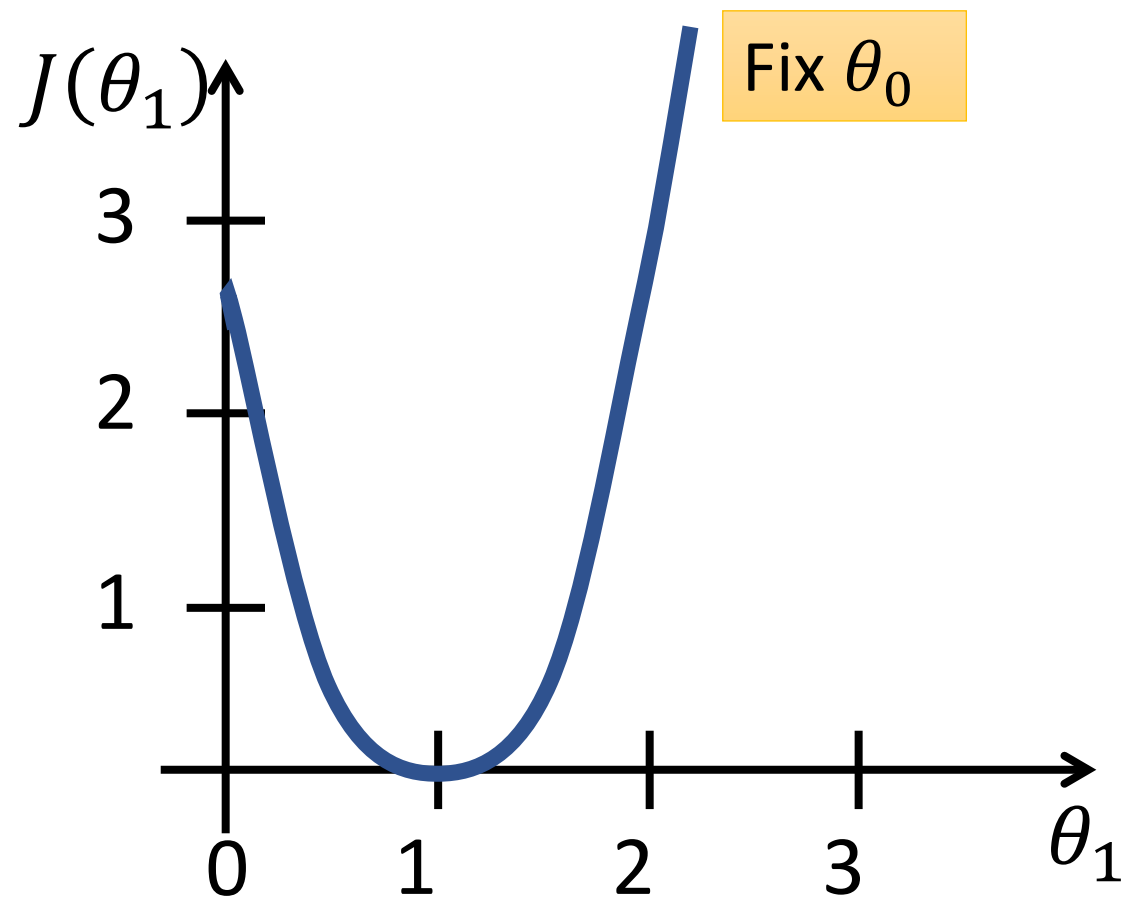
$J(\theta_1)$: function of θ_1



$h_{\theta}(x)$, function of x for fixed θ



$J(\theta_1)$, function of θ_0, θ_1

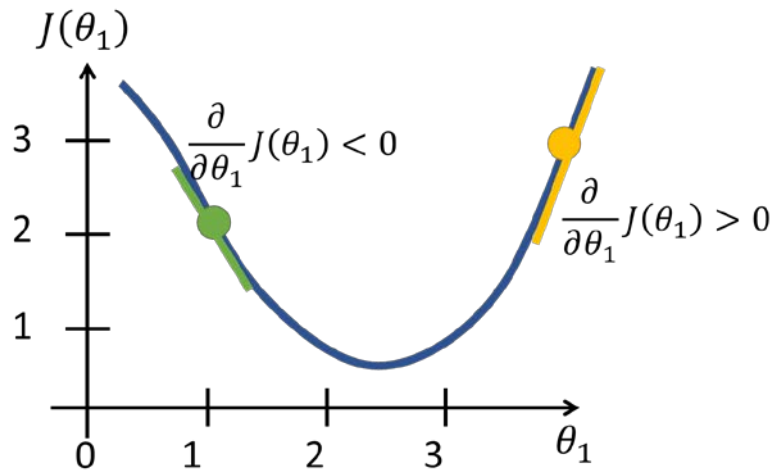
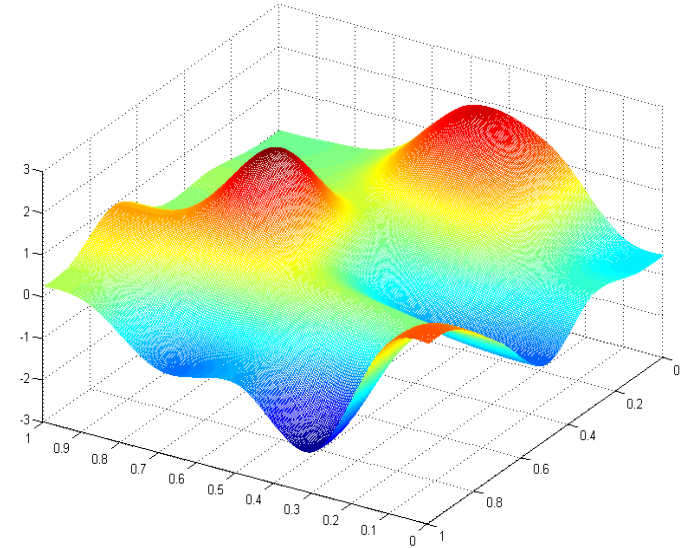


Minimizing cost function & Gradient Descent

Minimizing function $J(\theta_0, \theta_1)$

Outline:

- Start with some θ_0, θ_1
- Keep changing θ_0, θ_1 to reduce $J(\theta_0, \theta_1)$
- until we end up at a minimum



$$\theta_1 := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_1)$$

$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

Computing partial derivatives

Repeat until convergence{

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\bar{\theta})$$

Equivalently

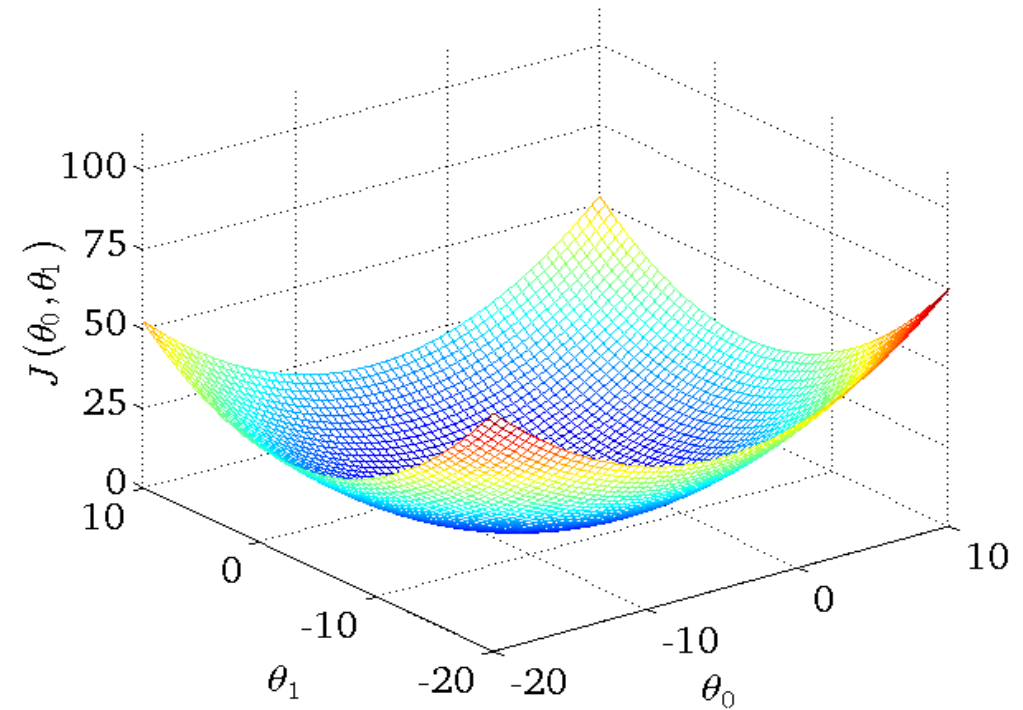
$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

}

$$\frac{\partial}{\partial \theta_j} J(\bar{\theta}) = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

Convergence

- The cost function in linear regression is always a convex function – always has a single global minimum
- So, gradient descent will always converge



Batch gradient descent

“Batch”: Each step of gradient descent uses all the training examples

Repeat until convergence{

m : Number of training examples

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})$$

$$\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x^{(i)}$$

}

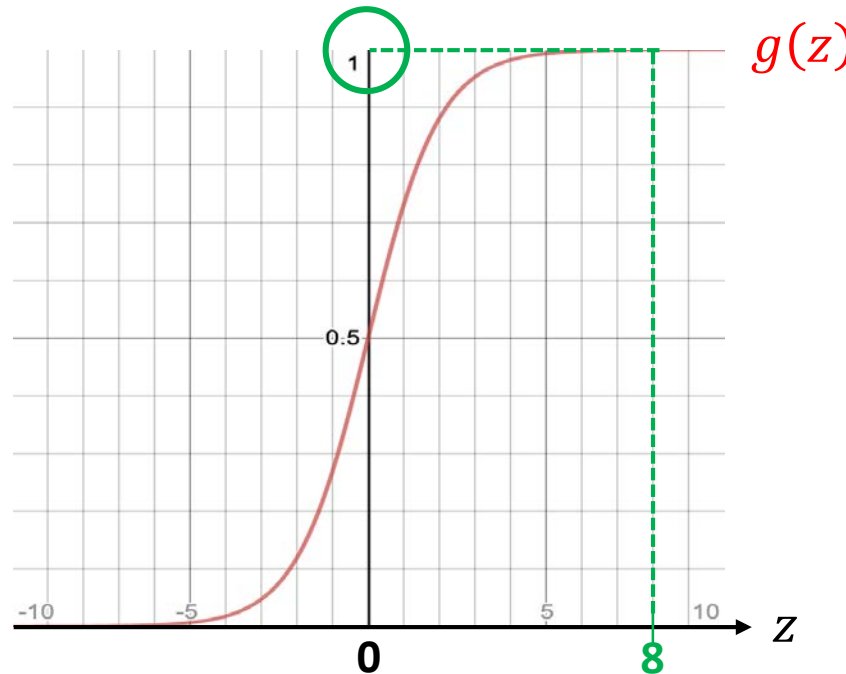
Logistic Regression for Classification

Regression vs. Classification

- How can we make the possible outputs of $\mathbf{h}_\theta(\mathbf{x}) = \boldsymbol{\theta}^T \mathbf{x}$ discrete-valued (as opposed to real-valued)?
 - By using an **activation function** (e.g., **sigmoid or logistic function**)

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

$z \in \mathbb{R}$, but
 $g(z) \in [0,1]$



Assume a labeled example (\mathbf{x}, y) :

If $y = \mathbf{1}$, we want $g(z) \approx 1$ (i.e., we want a correct prediction)

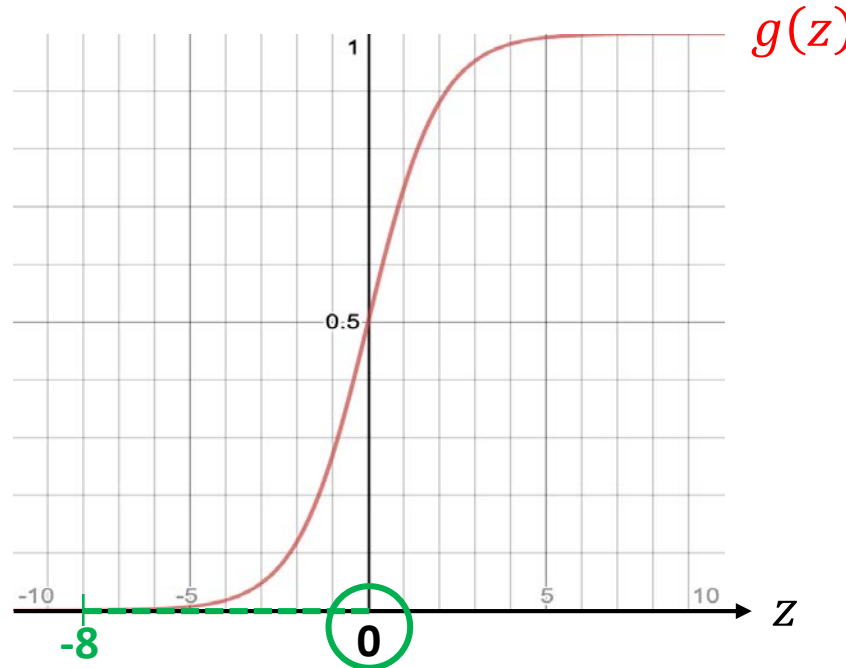
For this to happen, $\mathbf{z} \gg \mathbf{0}$

Regression vs. Classification

- How can we make the possible outputs of $\mathbf{h}_\theta(\mathbf{x}) = \boldsymbol{\theta}^T \mathbf{x}$ discrete-valued (as opposed to real-valued)?
 - By using an **activation function** (e.g., **sigmoid or logistic function**)

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

$z \in \mathbb{R}$, but
 $g(z) \in [0,1]$



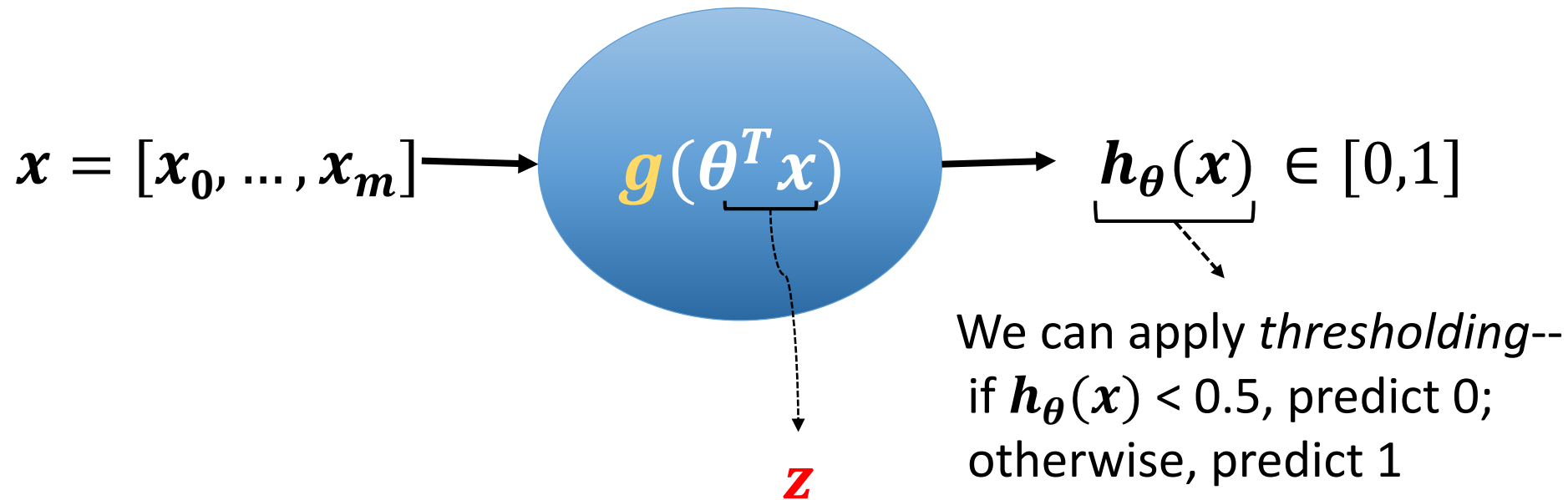
Assume a labeled example (\mathbf{x}, y) :

If $y = \mathbf{0}$, we want $g(z) \approx 0$ (i.e., we want a correct prediction)

For this to happen, $z \ll \mathbf{0}$

Regression vs. Classification

- How can we make the possible outputs of $\mathbf{h}_\theta(\mathbf{x}) = \boldsymbol{\theta}^T \mathbf{x}$ discrete-valued (as opposed to real-valued)?
 - By using an **activation function** (e.g., **sigmoid or logistic function**)



Alternative Interpretation: $\mathbf{h}_\theta(\mathbf{x})$ = estimated probability that $y = 1$ on input \mathbf{x}

Hypothesis representation

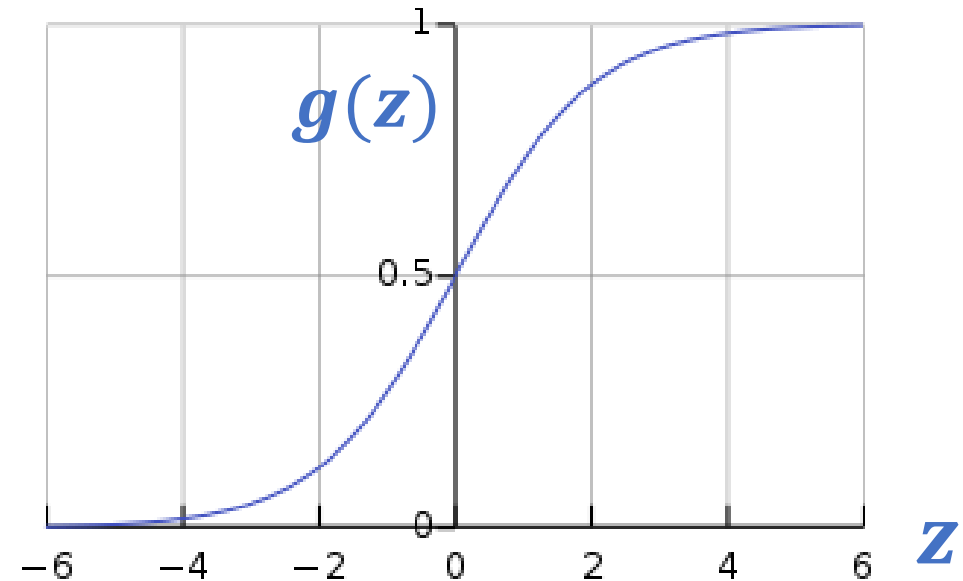
Want $0 \leq h_{\theta}(x) \leq 1$

$$h_{\theta}(x) = g(\theta^{\top} x),$$

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

- Sigmoid function
- Logistic function

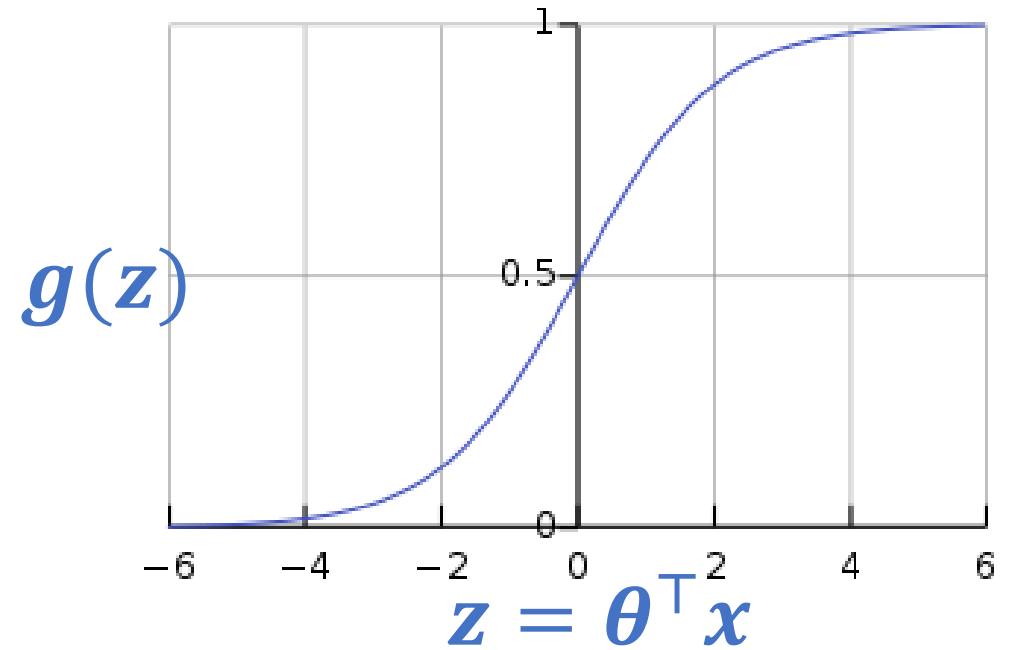
$$h_{\theta}(x) = \frac{1}{1 + e^{-\theta^{\top} x}}$$



Logistic regression

$$h_{\theta}(x) = g(\theta^{\top}x)$$

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



Suppose predict “y = 1” if $h_{\theta}(x) \geq 0.5$

$$z = \theta^{\top}x \geq 0$$

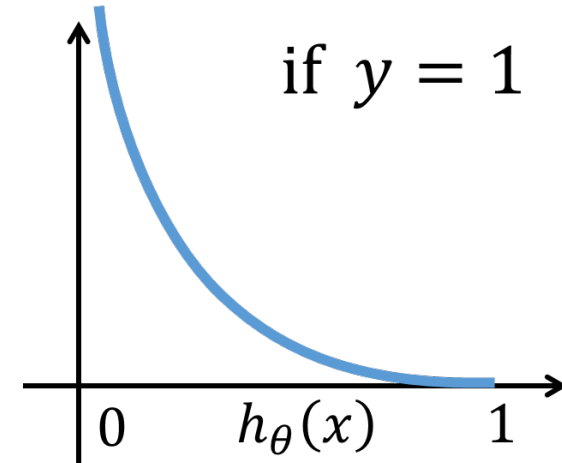
predict “y = 0” if $h_{\theta}(x) < 0.5$

$$z = \theta^{\top}x < 0$$

Cost function for Logistic Regression

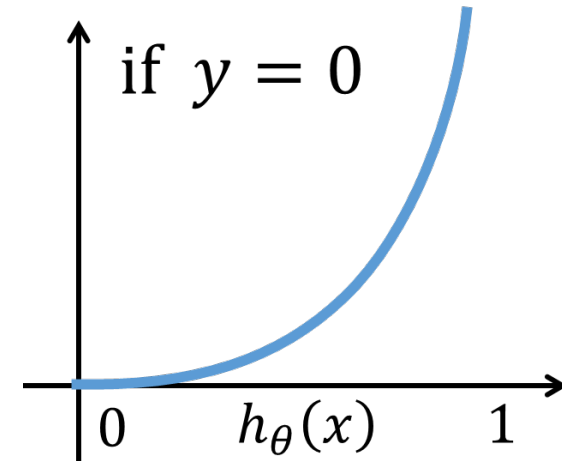
Linear Regression

$$\text{Cost}(h_{\theta}(x), y) = \frac{1}{2} (h_{\theta}(x) - y)^2$$



Logistic Regression

$$\begin{aligned} \text{Cost}(h_{\theta}(x), y) &= \begin{cases} -\log(h_{\theta}(x)) & \text{if } y = 1 \\ -\log(1 - h_{\theta}(x)) & \text{if } y = 0 \end{cases} \\ &= -y \log(h_{\theta}(x)) - (1 - y) \log(1 - h_{\theta}(x)) \end{aligned}$$



Logistic regression

$$\begin{aligned} J(\theta) &= \frac{1}{m} \sum_{i=1}^m \text{Cost}(h_{\theta}(x^{(i)}), y^{(i)}) \\ &= -\frac{1}{m} \left[\sum_{i=1}^m y^{(i)} \log(h_{\theta}(x^{(i)})) + (1 - y^{(i)}) \log(1 - h_{\theta}(x^{(i)})) \right] \end{aligned}$$

Learning: fit parameter θ

$$\min_{\theta} J(\theta)$$

Prediction: given new x

$$\text{Output } h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}}$$

Gradient descent

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

Goal: $\min_{\theta} J(\theta)$

Good news: Convex function!

Bad news: No analytical solution

Repeat {

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

}

(Simultaneously update all θ_j)

$$\frac{\partial}{\partial \theta_j} J(\theta) = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

Gradient descent

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

Goal: $\min_{\theta} J(\theta)$

Repeat { (Simultaneously update all θ_j)

$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

}

Gradient descent for **Linear Regression**

Repeat {

$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

}

$$h_{\theta}(x) = \theta^{\top} x$$

Gradient descent for **Logistic Regression**

Repeat {

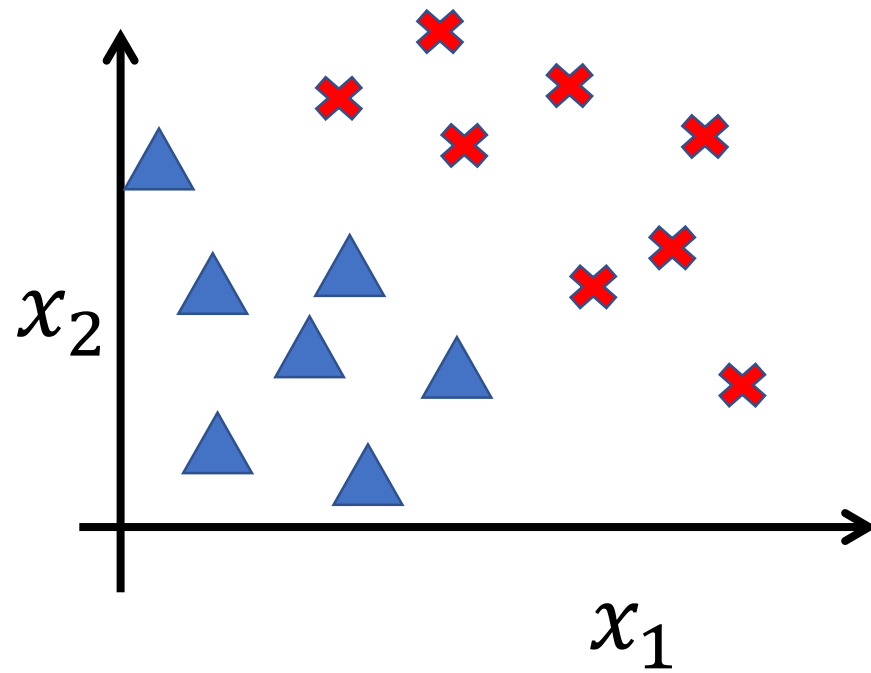
$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

}

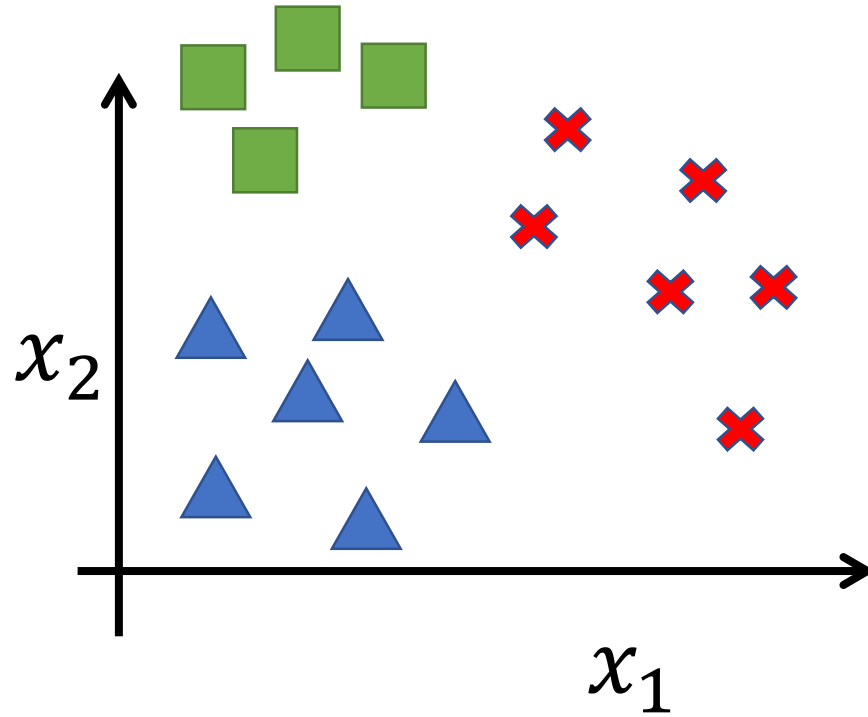
$$h_{\theta}(x) = \frac{1}{1 + e^{-\theta^{\top} x}}$$

Multiclass classification

Binary classification



Multiclass classification



Multi-class Classification

- Binary classification: $h_{\theta}(x) = \theta^T x$
- Multi-class Classification: y can take on K different values $y^{(i)} \in \{1, 2, \dots, K\}$
- $h_{\theta}(x)$ estimates the probability of belonging to each class

$$P(y = k|x, \theta) \propto \exp(\theta_k^T x)$$

$$\theta = \begin{bmatrix} \vdots & \vdots & \vdots \\ \theta_1 & \theta_2 & \theta_k \\ \vdots & \vdots & \vdots \end{bmatrix}$$

$$P(y = k|x, \theta) = \frac{\exp(\theta_k^T x)}{\sum_{j=1}^K \exp(\theta_j^T x)}$$

Multi-class Classification Cost Function

$$P(y = k|x, \theta) = \frac{\exp(\theta_k^T x)}{\sum_{j=1}^K \exp(\theta_j^T x)}$$

$$J(\theta) = - \left[\sum_{i=1}^m \sum_{j=1}^K 1\{y^{(i)} = k\} \log \frac{\exp(\theta_k^T x^{(i)})}{\sum_{j=1}^K \exp(\theta_j^T x^{(i)})} \right]$$

Machine Learning– Part 2

Neural Networks

Biological Neural Network

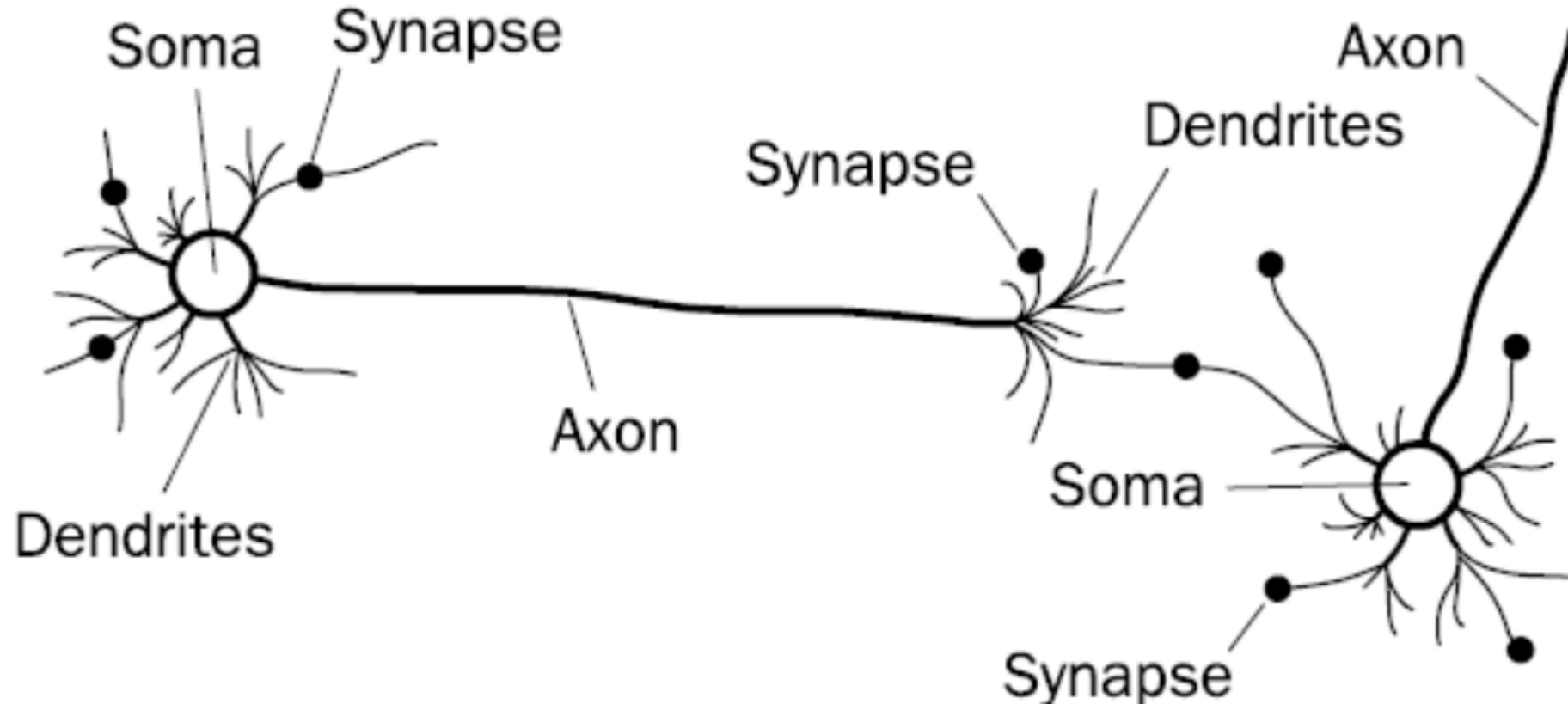
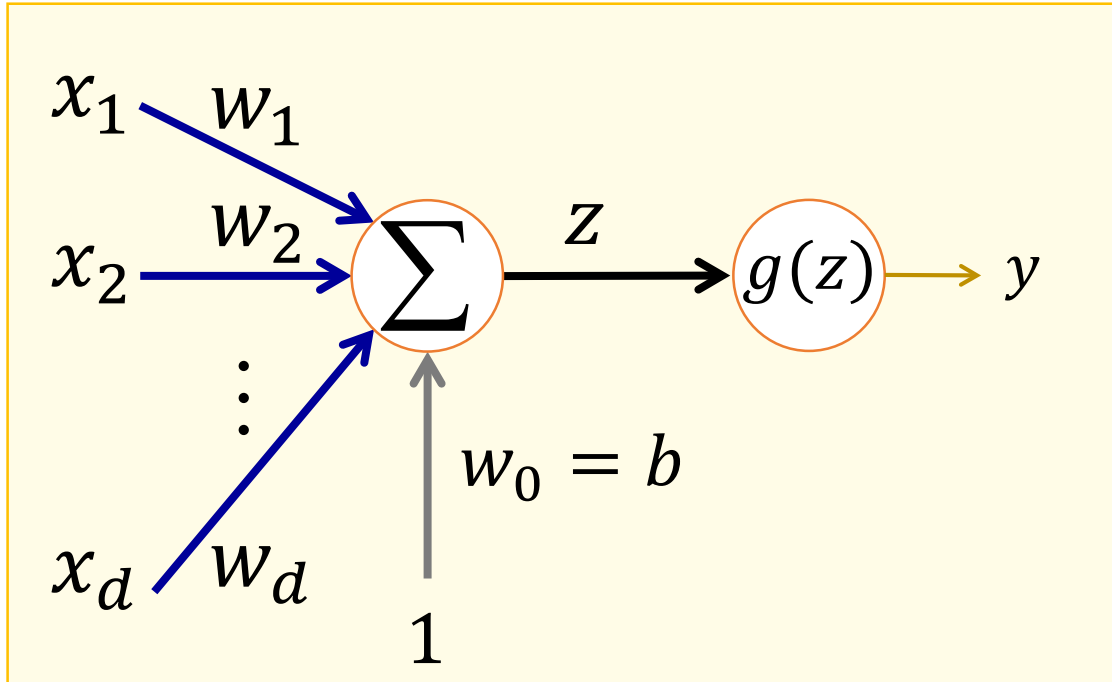


Image courtesy: F. A. Makinde et. al., "Prediction of crude oil viscosity using feed-forward back-propagation neural network (FFBPNN)". Petroleum & Coal 2012

Artificial Neuron

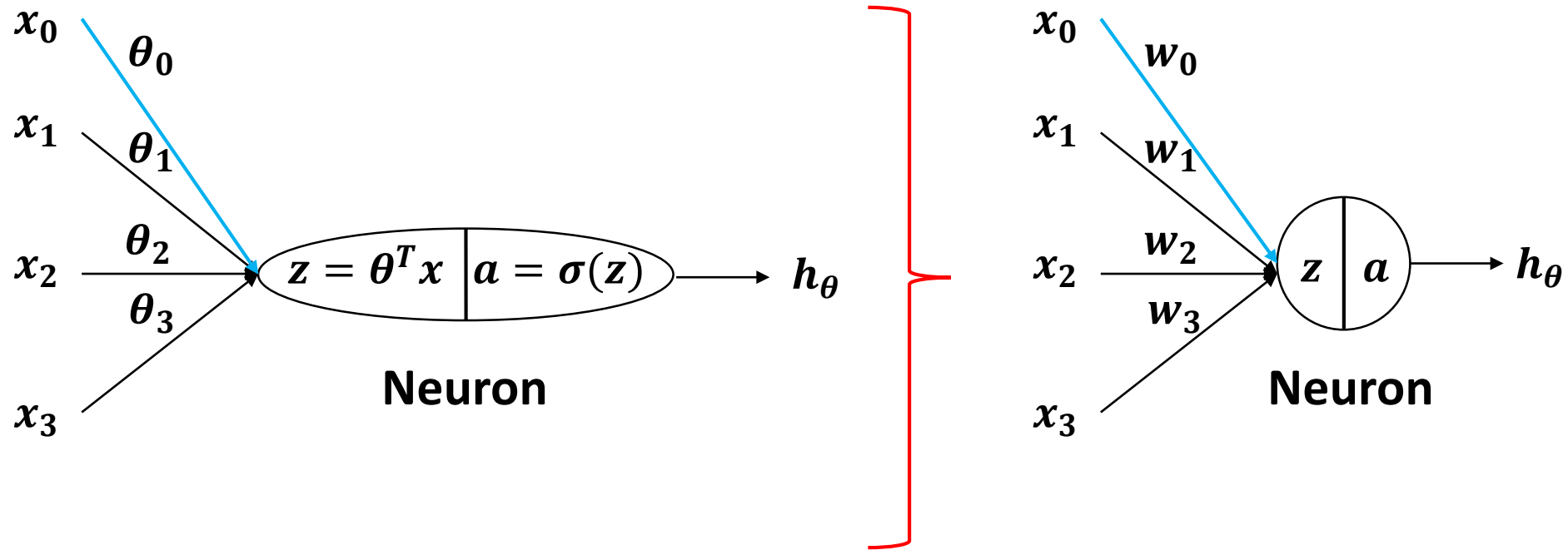


Terminologies

x : input, w : weights, b : bias
 z : pre-activation (input activation)
 g : activation function
 y : activation for output units

Towards Neural Networks

- Logistic regression is a **neural network** with only 1 neuron



Output Units: Linear

$$\hat{y} = w^T a + b$$

Used to produce the mean of a conditional Gaussian distribution:

$$p(\mathbf{y} | \mathbf{x}) = N(\mathbf{y}; \hat{\mathbf{y}}, \sigma)$$

Maximizing log-likelihood \Rightarrow Minimizing squared error

Output Units: Sigmoid

$$\hat{y} = \sigma(w^T a + b)$$

$$\begin{aligned} J(\theta) &= -\log p(y|x) \\ &= -\log \sigma((2y - 1)(\mathbf{w}^T \mathbf{a} + b)) \end{aligned}$$

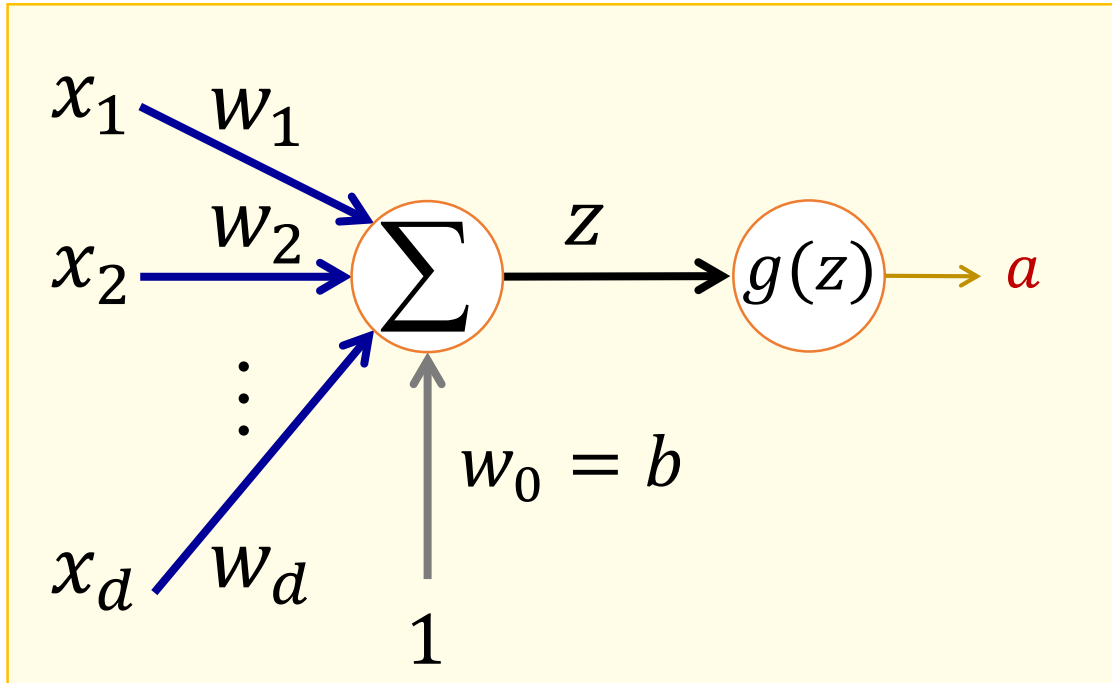
Output Softmax Units

Need to produce a vector \hat{y} with $\hat{y}_i = p(y = i|x)$

$$\text{softmax}(z)_i = \frac{\exp(z_i)}{\sum_j \exp(z_j)}$$

$$\log \text{softmax}(z)_i = z_i - \log \sum_j \exp(z_j)$$

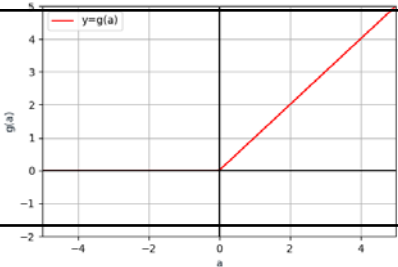
Artificial Neuron – hidden unit



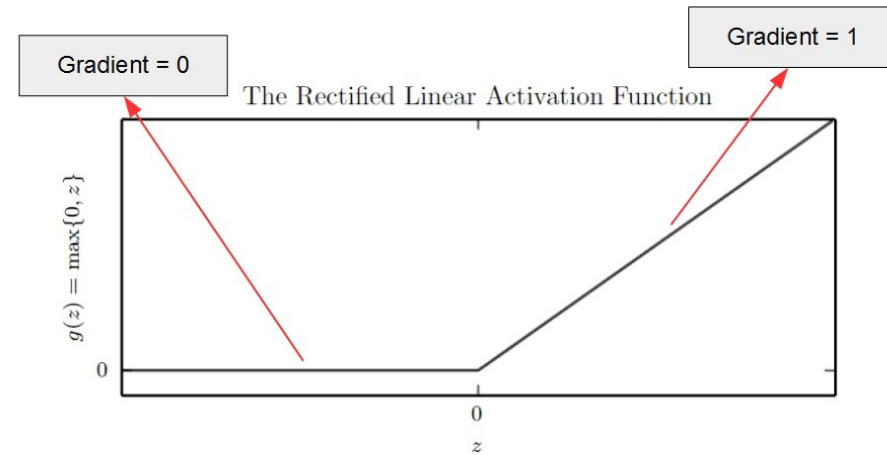
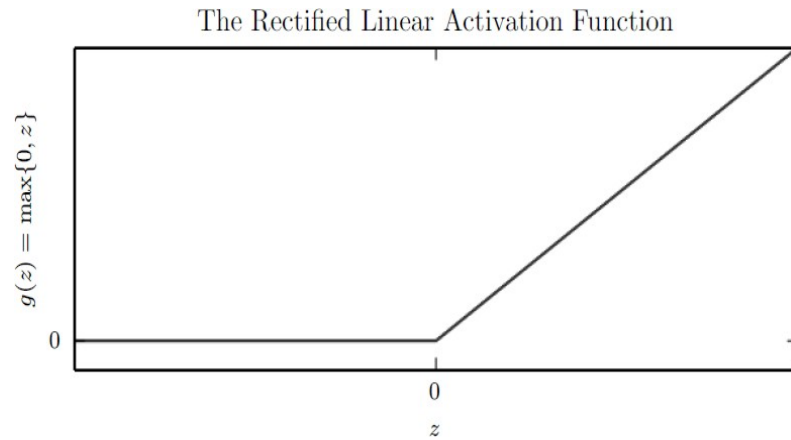
$$\mathbf{w} = [w_1 \ w_2 \ \dots \ w_d]^T \text{ and } \mathbf{x} = [x_1 \ x_2 \ \dots \ x_d]^T$$

$$\mathbf{z} = b + \sum_{i=1}^d w_i x_i = [\mathbf{w}^T \ b] \begin{bmatrix} \mathbf{x} \\ 1 \end{bmatrix}$$
$$\mathbf{a} = g(z)$$

Activation Functions for Hidden Nodes

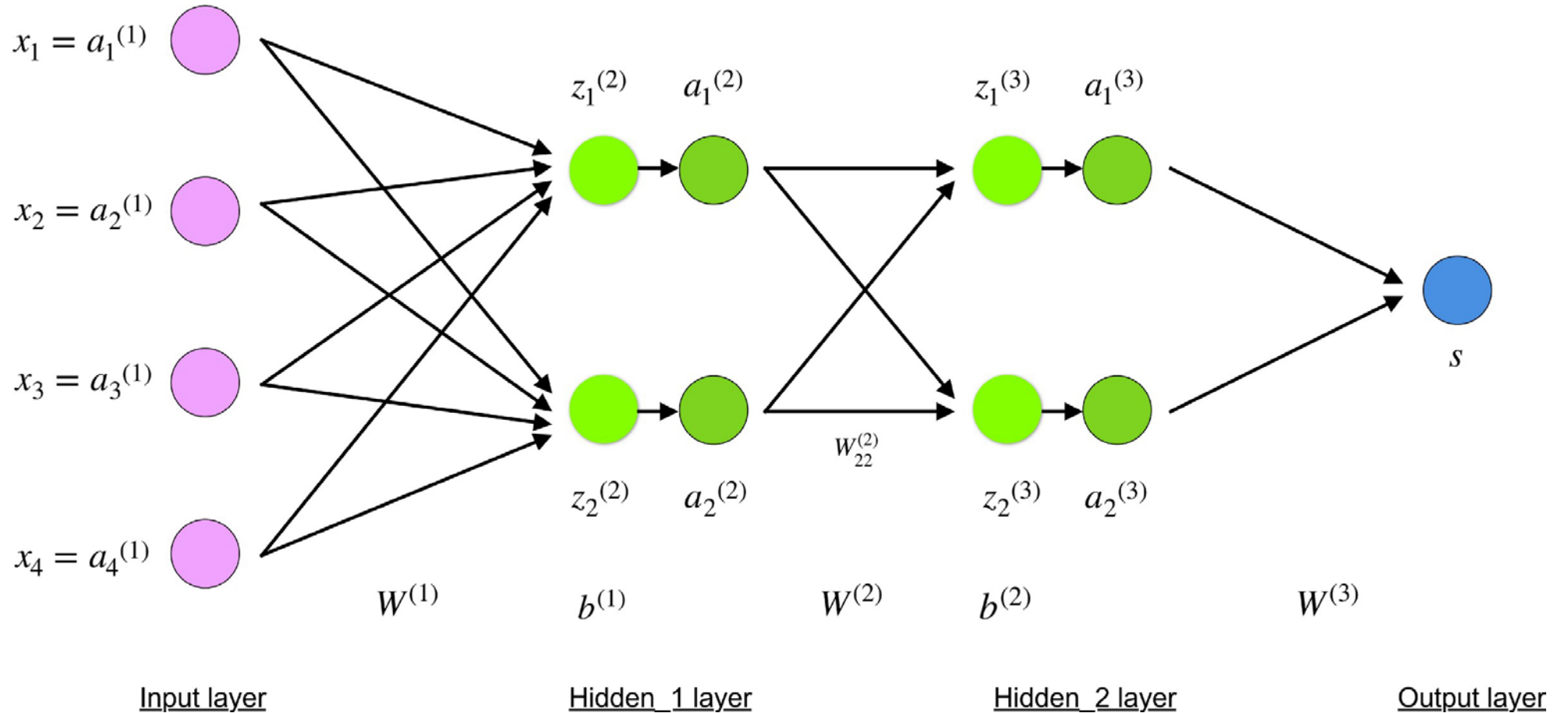
| Name | Function | Gradient | Graph |
|----------|--|---|---|
| Sigmoid | $\sigma(z) = \frac{1}{1 + \exp(-z)}$ | $g'(z) = g(z)(1 - g(z))$ | |
| Tanh | $\tanh(z) = \frac{\exp(z) - \exp(-z)}{\exp(z) + \exp(-z)}$ | $g'(z) = 1 - g^2(z)$ | |
| ReLU | $g(z) = \max(0, z)$ | $g'(z) = \begin{cases} 1, & z \geq 0 \\ 0, & z < 0 \end{cases}$ |  |
| softplus | $g(z) = \ln(1 + e^z)$ | | |

Rectified Linear Units

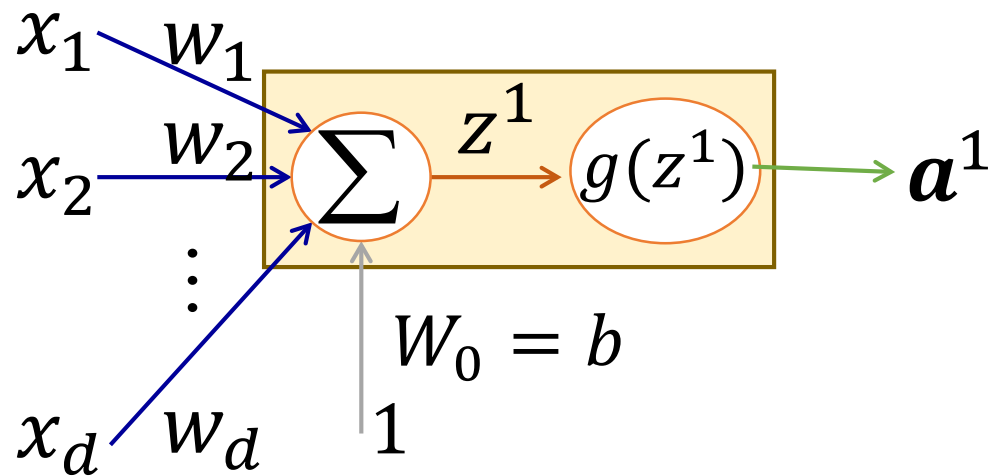


- Activation function: $g(z) = \max\{0, z\}$
 - Gives large and *consistent* gradients (does not saturate) when active
 - Efficient to optimize, converges much faster than sigmoid or tanh

Multilayer Neural Network

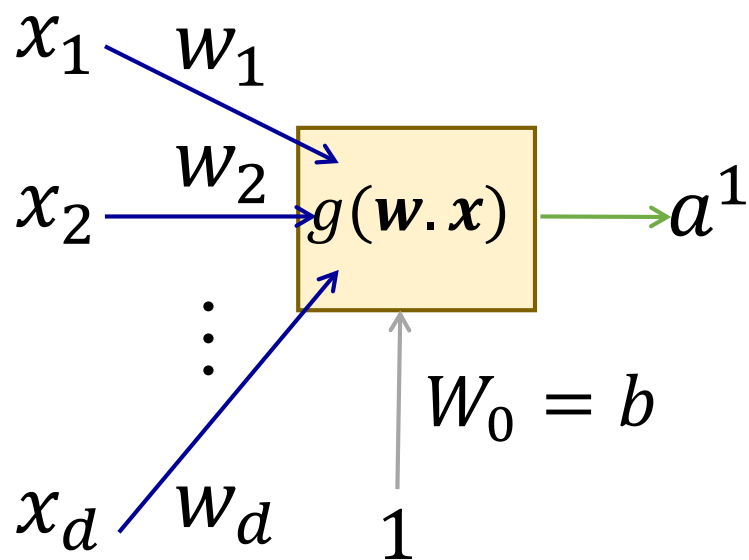


Basic Neural Units



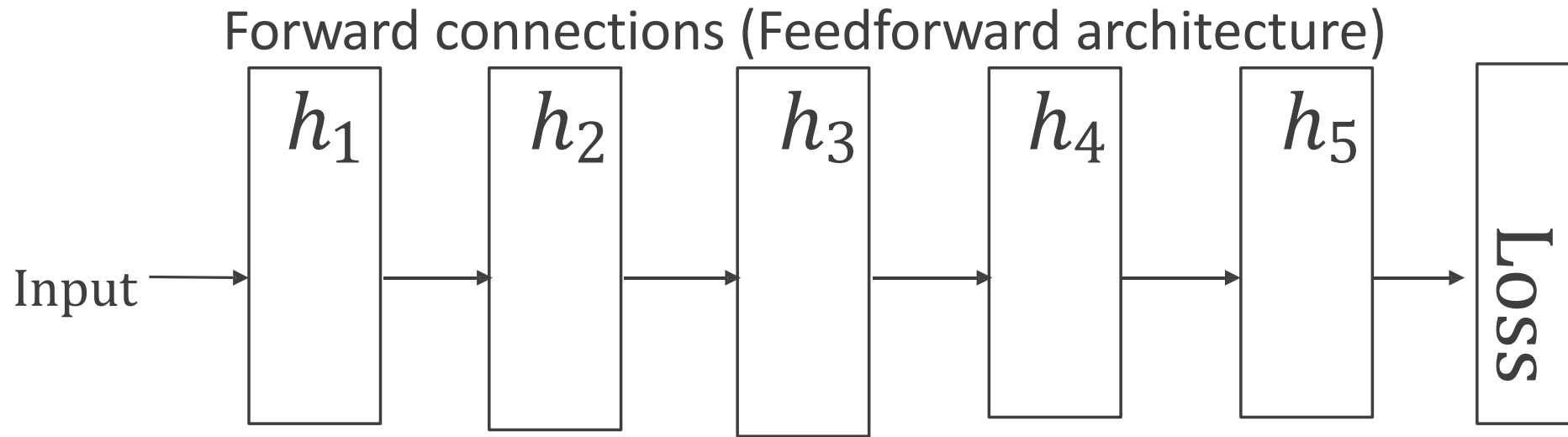
$$z_1^1 = b_1^1 + \sum_{i=1}^d w_{1,i}^1 x_i$$

$$a_1^1 = g(z_1^1)$$



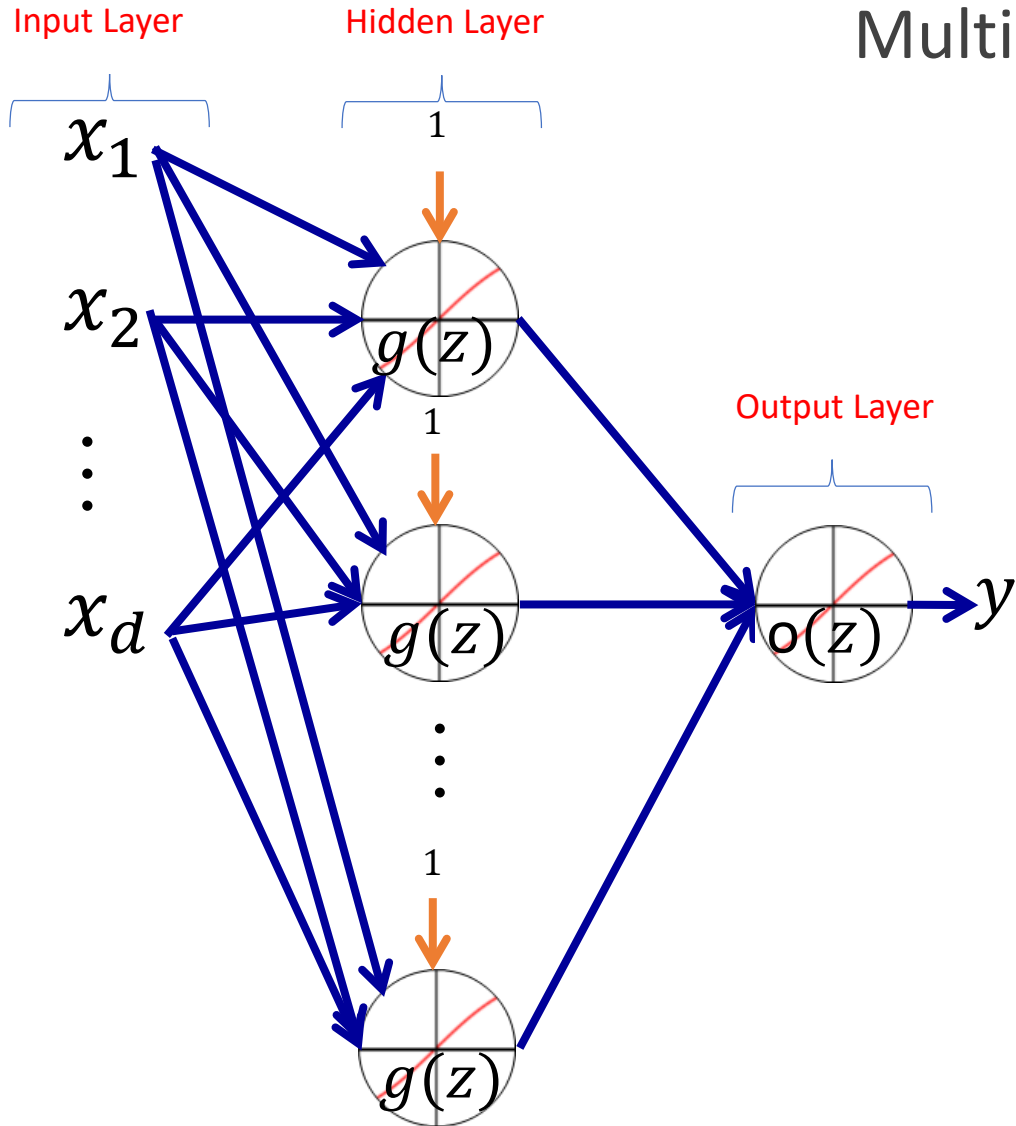
Neural networks in blocks

We can visualize $a_L = h_L \circ h_{L-1} \circ \dots \circ h_1(x)$ as a cascade of blocks.



The activation functions must be **1st-order differentiable (almost) everywhere**

Multilayer Neural Network

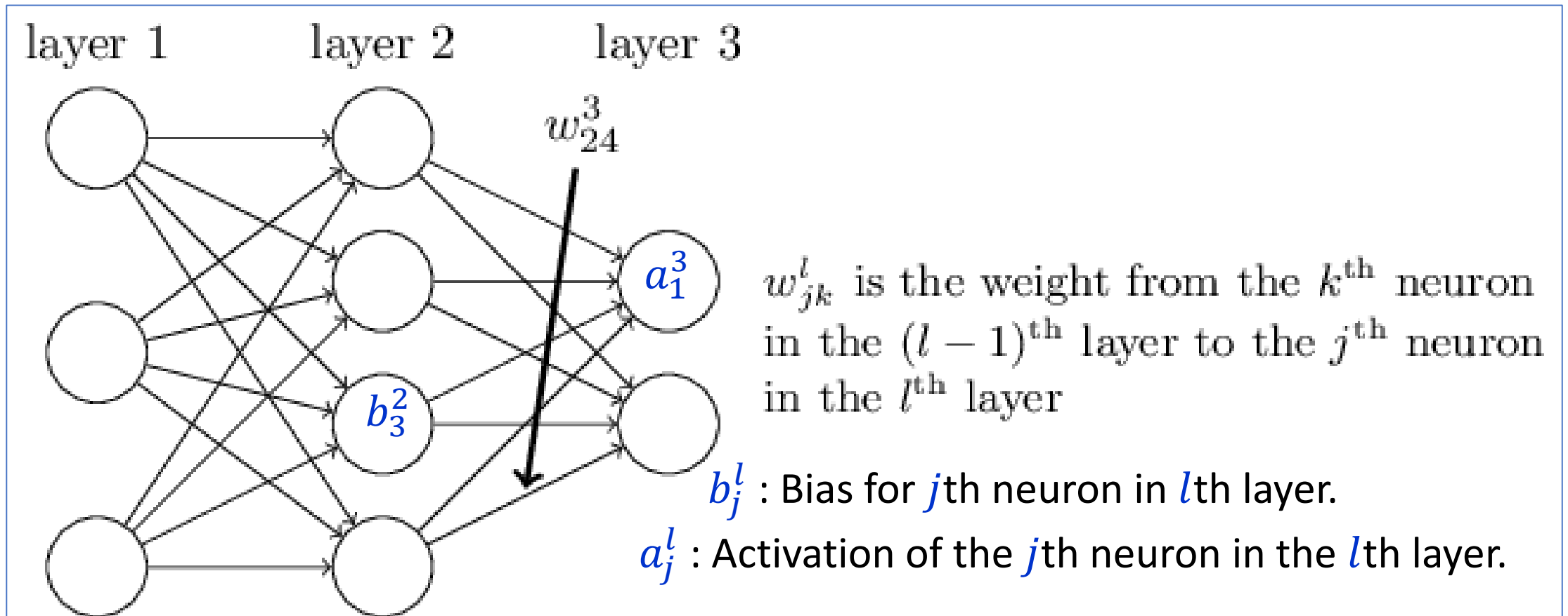


Intuition:

Hidden Layer: Extracts better representation of the input data

Output layer: Does the classification

Notations



$$a_j^l = g(\sum_k w_{jk}^l a_k^{l-1} + b_j^l)$$

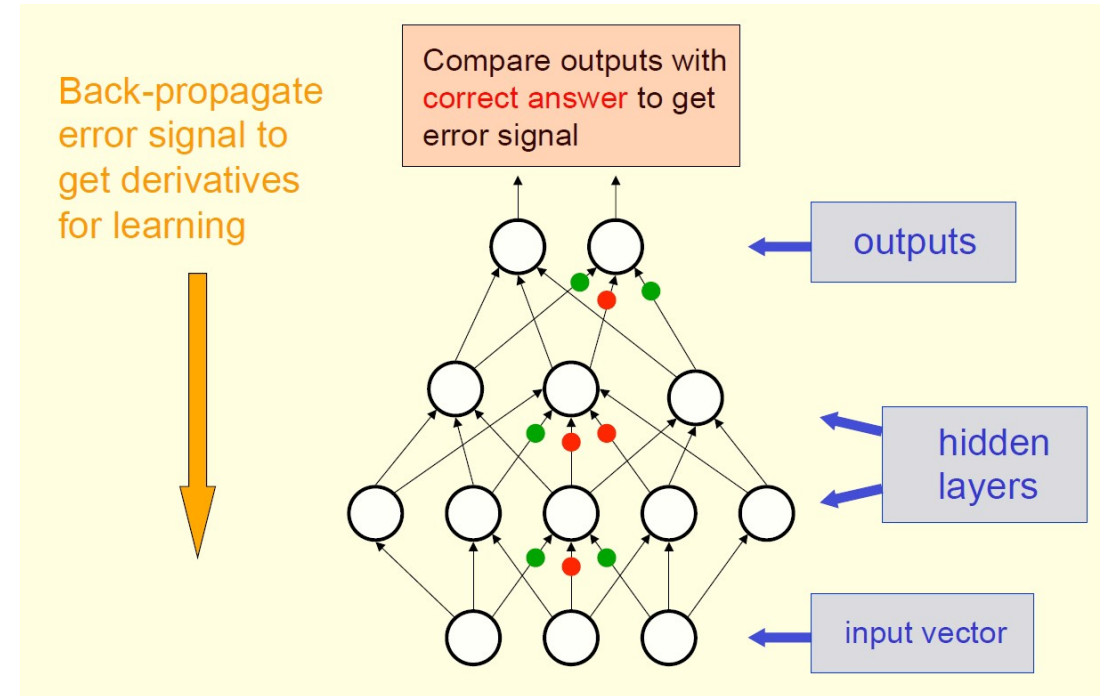
Vectorized form: $a^l = g(w^l a^{l-1} + b^l)$

$$z^l = w^l a^{l-1} + b^l$$

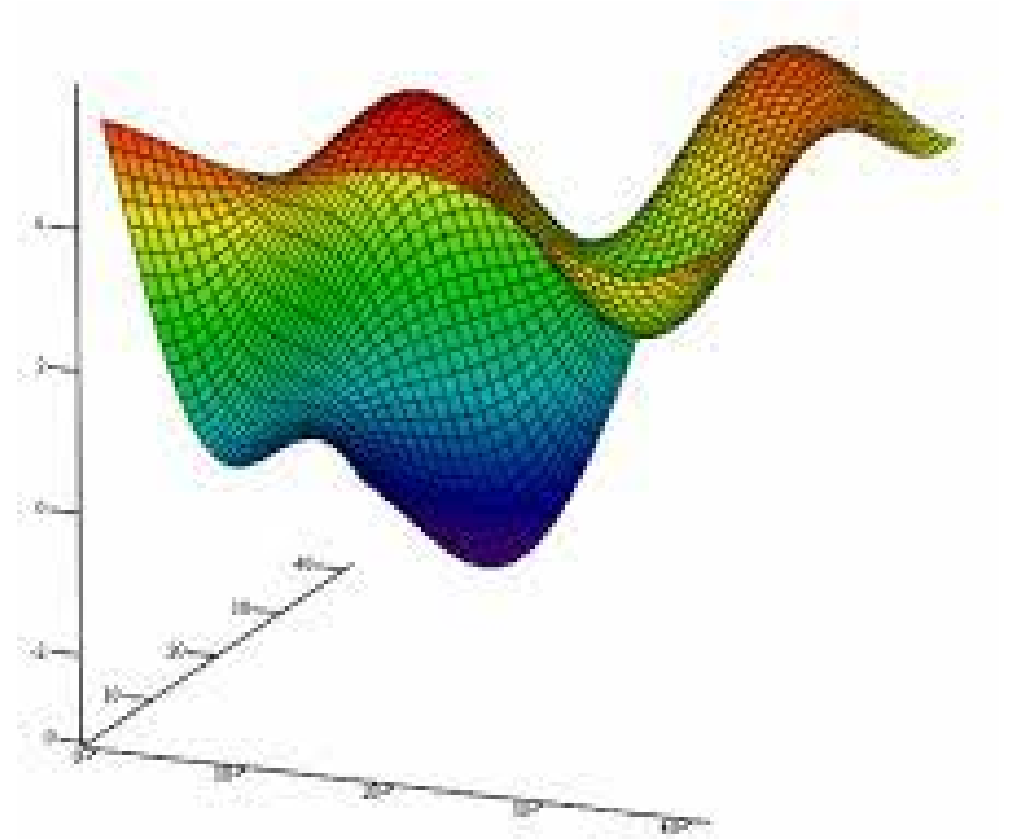
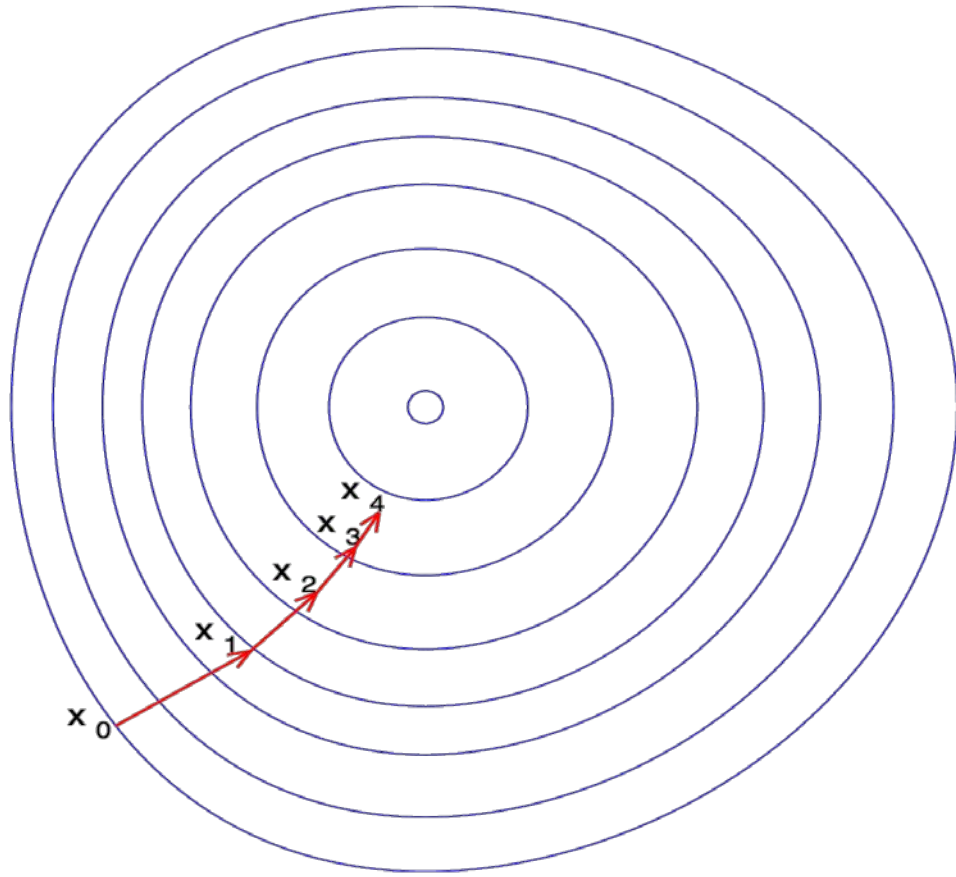
$$a^l = g(z^l)$$

Backpropagation

- Feedforward Propagation: Accept input $x^{(i)}$, pass through intermediate stages and obtain output $\hat{y}^{(i)}$
- During Training: Compute scalar cost $J(\theta)$
$$J(\theta) = \sum_i L(NN(x^{(i)}; \theta), y^{(i)})$$
- Backpropagation allows information to flow backwards from cost to compute the gradient



Gradient Descent



Stochastic gradient descent

Sample rather than computing the full sum

While not converged:

Select a single datapoint in order from the data

Compute gradient with just one point

Update parameters

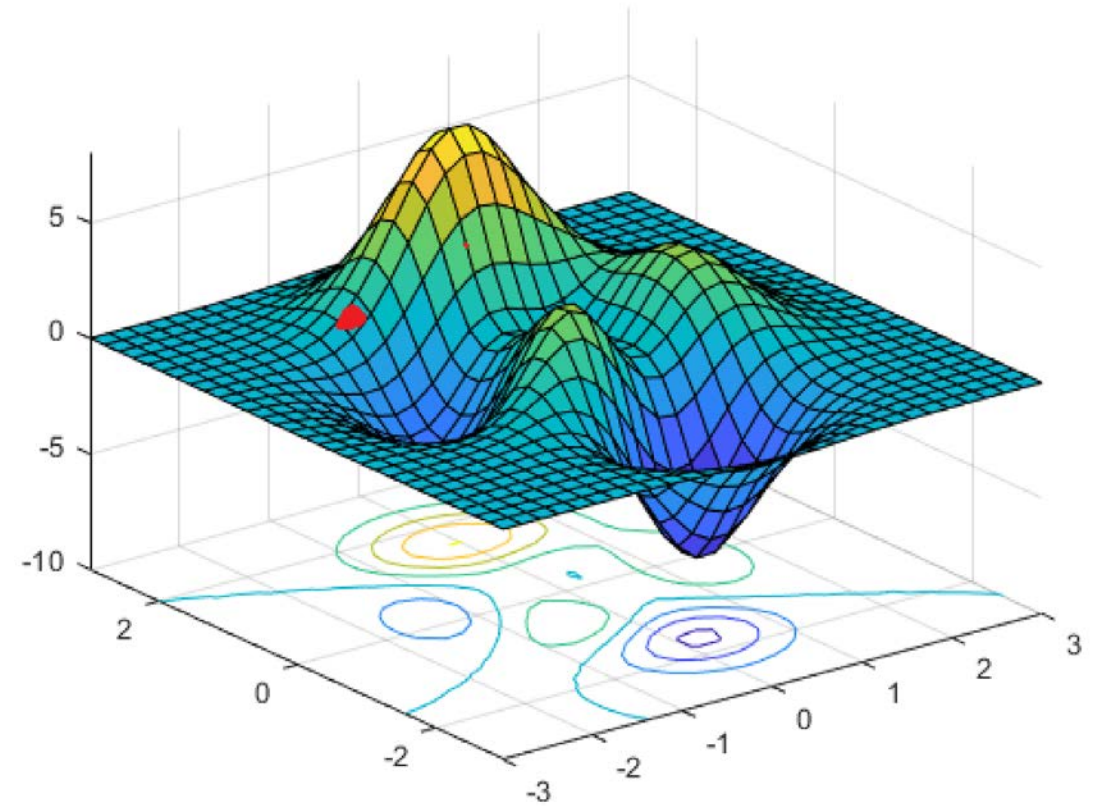
Mini-batch SGD:

While not converged:

Select a random batch of datapoints

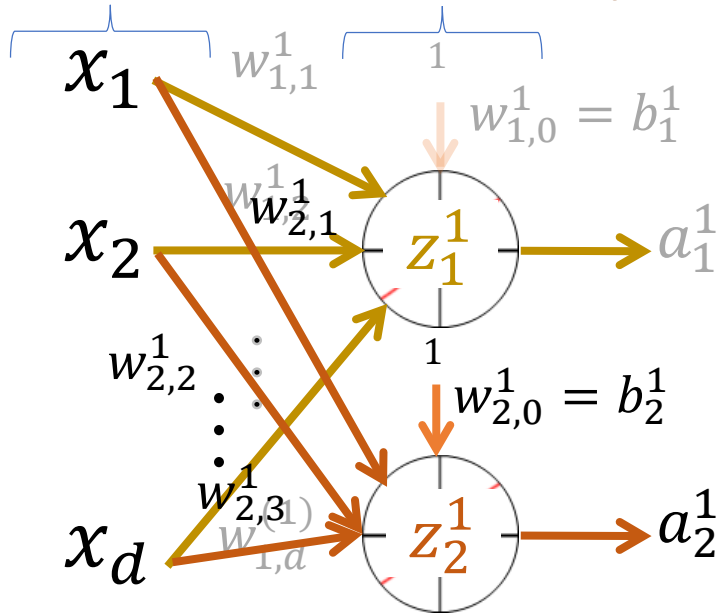
Compute gradient with the batch

Update parameters



Multilayer Neural Network

Input Layer Hidden Layer



$$z_1^1 = b_1^1 + \sum_{i=1}^d w_{1,i}^1 x_i \quad a_1^1 = g(z_1^1)$$

$$z_2^1 = b_2^1 + \sum_{i=1}^d w_{2,i}^1 x_i \quad a_2^1 = g(z_2^1)$$

... ..

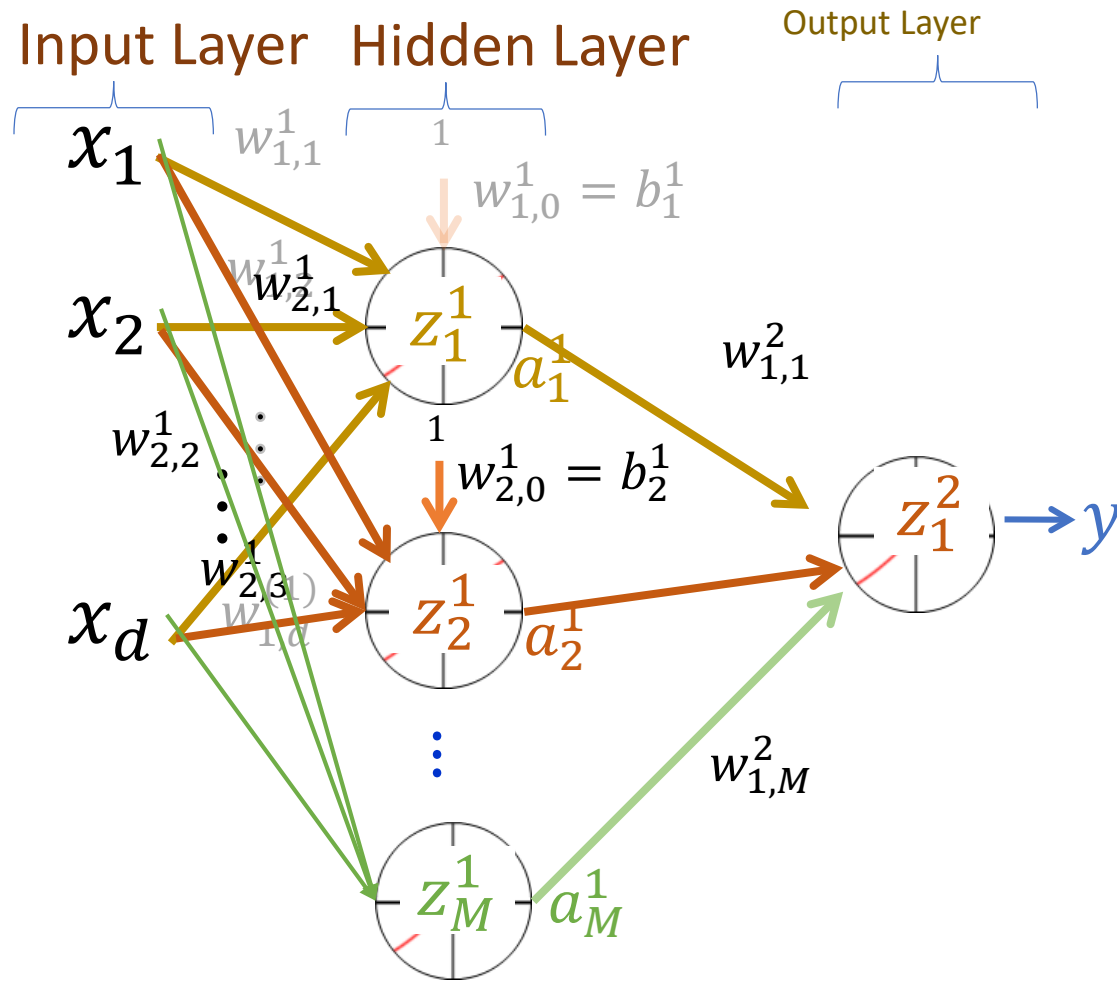
$$z_m^1 = b_m^1 + \sum_{i=1}^d w_{m,i}^1 x_i \quad a_m^1 = g(z_m^1)$$

$$a^{(0)} = x$$

$$z^{(1)} = \mathbf{w}^{(1)} \mathbf{a}^{(0)}$$

$$a^{(1)} = g(z^{(1)})$$

Multilayer Neural Network



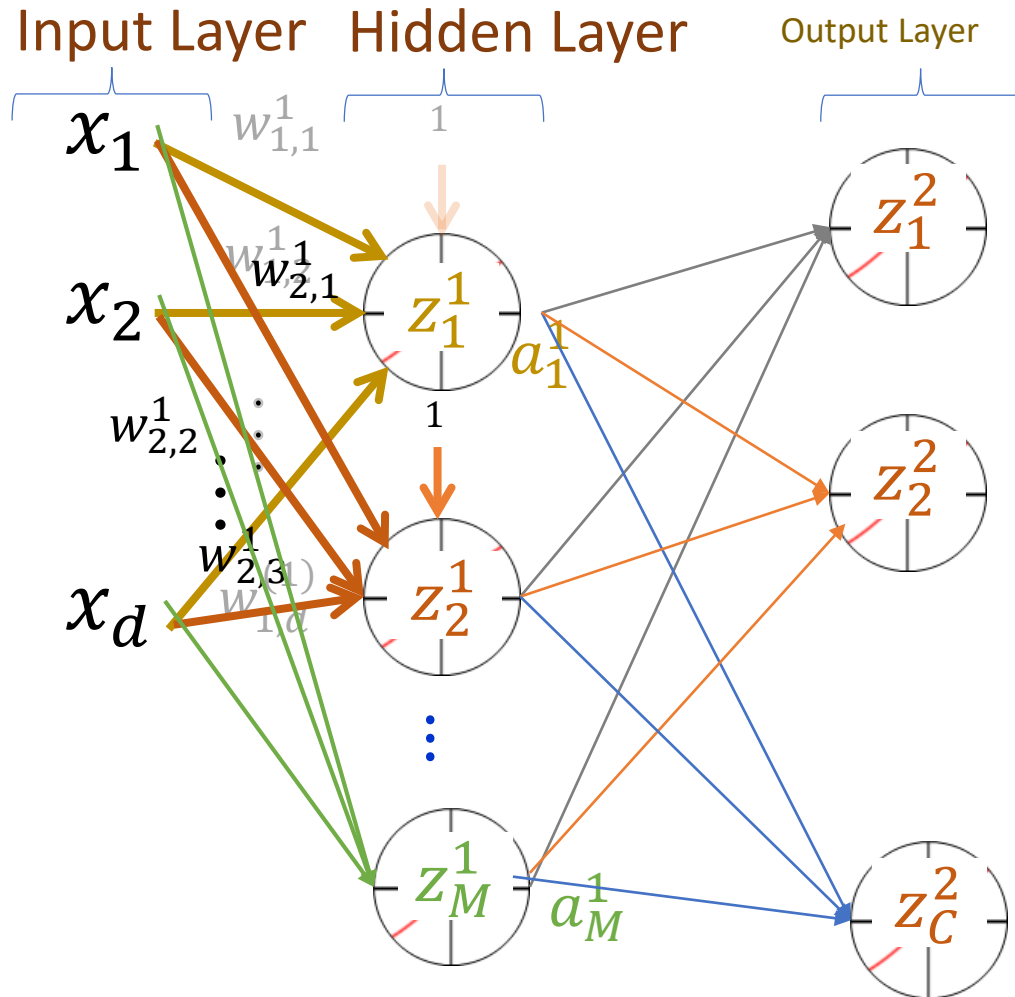
Output Layer Activation

$$y_1 = o(z_1^2)$$

output

- Sigmoid for 2-class classification
- Softmax for multi-class classification
- Linear for regression

Multilayer Neural Network

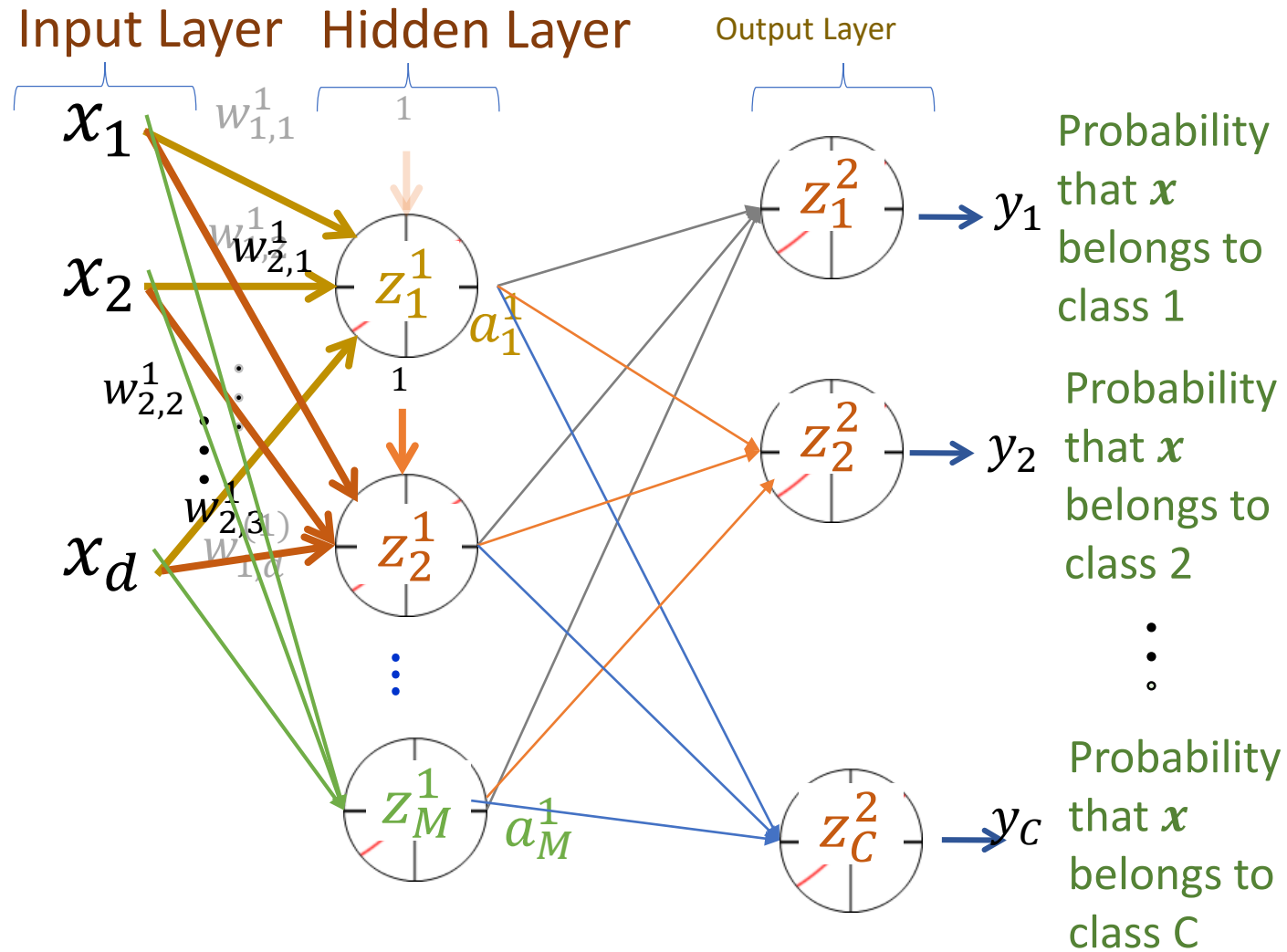


$$\rightarrow y_1 = o_1(z_1^2) = \frac{\exp(z_1^2)}{\sum_c \exp(z_c^2)}$$

$$\rightarrow y_2 = o_2(z_2^2) = \frac{\exp(z_2^2)}{\sum_c \exp(z_c^2)}$$

$$\rightarrow y_c = o_c(z_c^2) = \frac{\exp(z_c^2)}{\sum_c \exp(z_c^2)}$$

Training a Neural Network – Loss Function



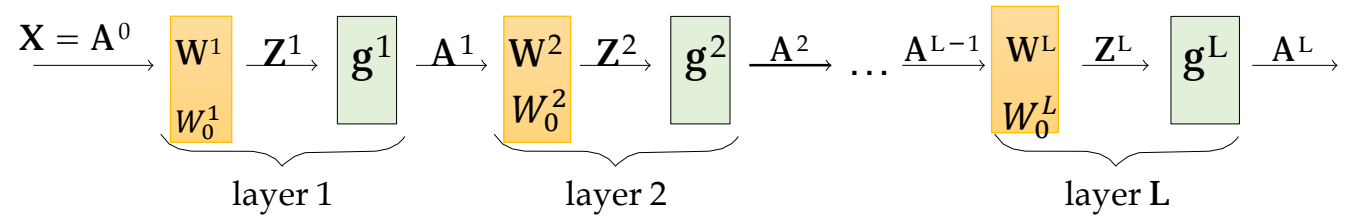
Aim to maximize the probability corresponding to the correct class for any example x

$$\begin{aligned} \max y_c \\ \equiv \max (\log y_c) \\ \equiv \min (-\log y_c) \end{aligned}$$

Can be equivalently expressed as

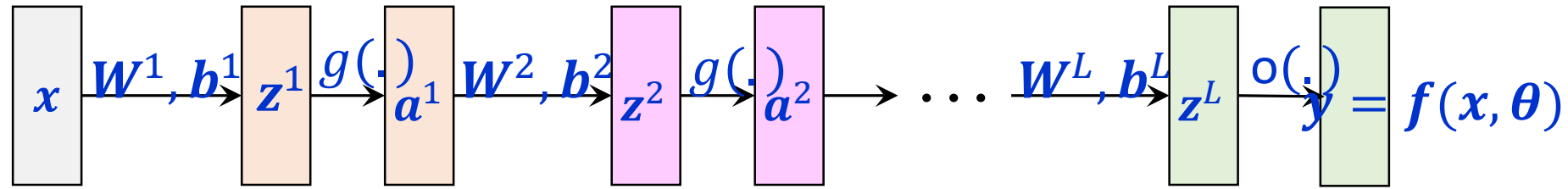
$-\sum_i \prod_{i=c} \log(y_i)$
known as cross-entropy loss

Multi layered network



Forward Pass in a Nutshell

θ is the collection of all learnable parameters i.e., all W and b



Hidden layer pre-activation:

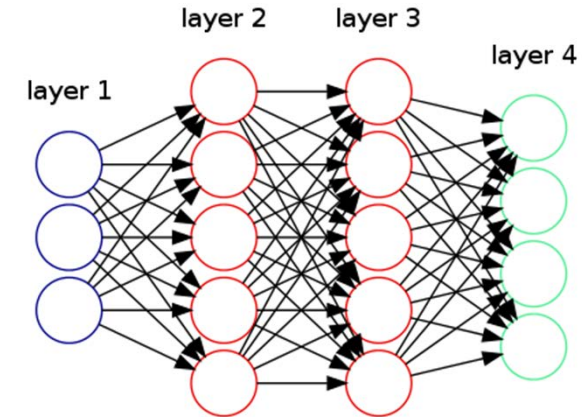
For $l = 1, \dots, L$; $\mathbf{z}^{(l)} = \mathbf{W}^{(l)} \mathbf{a}^{(l-1)} + \mathbf{b}^{(l)}$

Hidden layer activation:

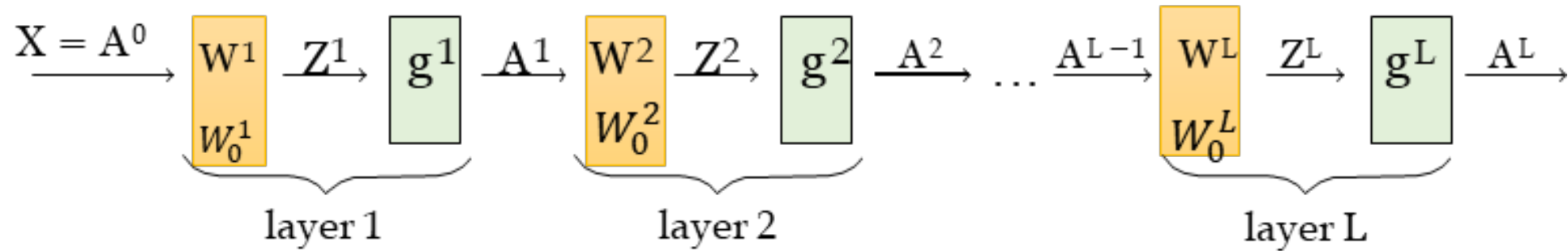
For $l = 1, \dots, L - 1$; $\mathbf{a}^{(l)} = g(\mathbf{z}^{(l)})$

Output layer activation:

For $l = L$; $\mathbf{y} = \mathbf{a}^{(L)} = o(\mathbf{z}^{(L)}) = f(\mathbf{x}, \theta)$



Error back-propagation



- We will train neural networks using gradient descent methods.
- To do SGD for a training example (x, y) , we need to compute

$$\nabla_W \text{Loss}(NN(x; W), y)$$

where W represents all weights W^l, W_0^l in all the layers $l = (1, \dots, L)$.

$$\frac{\partial \text{Loss}}{\partial W^L} = \underbrace{\frac{\partial \text{Loss}}{\partial A^L}}_{\text{Depends on Loss function}} \cdot \underbrace{\frac{\partial A^L}{\partial Z^L}}_{g^{L'}} \cdot \underbrace{\frac{\partial Z^L}{\partial W^L}}_{A^{L-1}}$$