

Artificial Intelligence

Lecture 10: Classification

Xiaojin Gong

2021-05-24

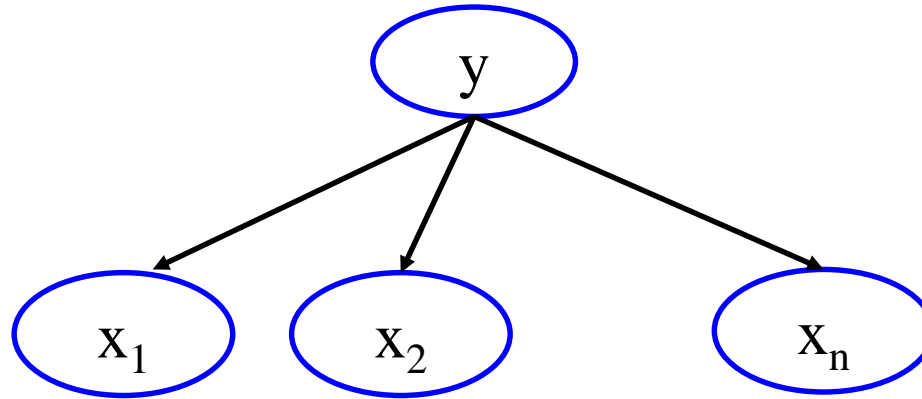
Credits: AI course in Berkeley & MIT

Review

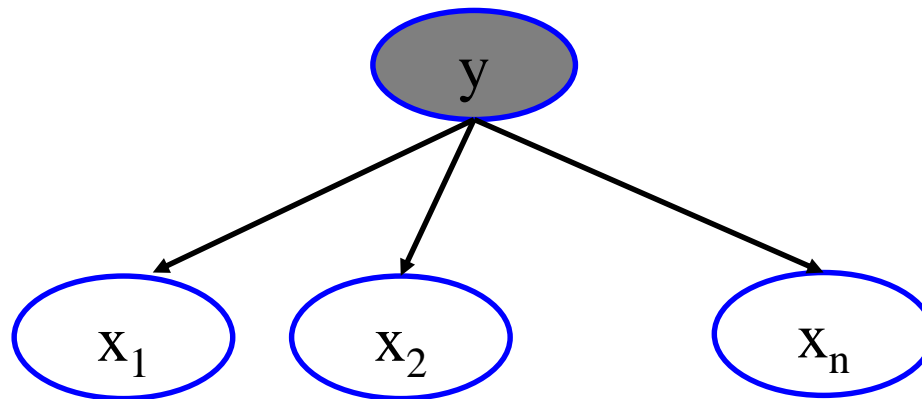
- Supervised learning:
 - Observes example input-output pairs
 - Learns a function that maps from input to output
- Unsupervised learning:
 - Learns patterns in the input even though no explicit feedback is given
- Reinforcement learning:
 - Learns from a series of reinforcements – rewards or punishments

Review

- Supervised classification: naïve Bayes model



- Unsupervised clustering: learning mixtures of Gaussians



Review

- Training & Testing:
 - Training set
 - Validation set
 - Test set
- Generalization & Overfitting:
 - Want a classifier which does well on test data
 - Fitting the training data very closely, but not generalizing well

Outline

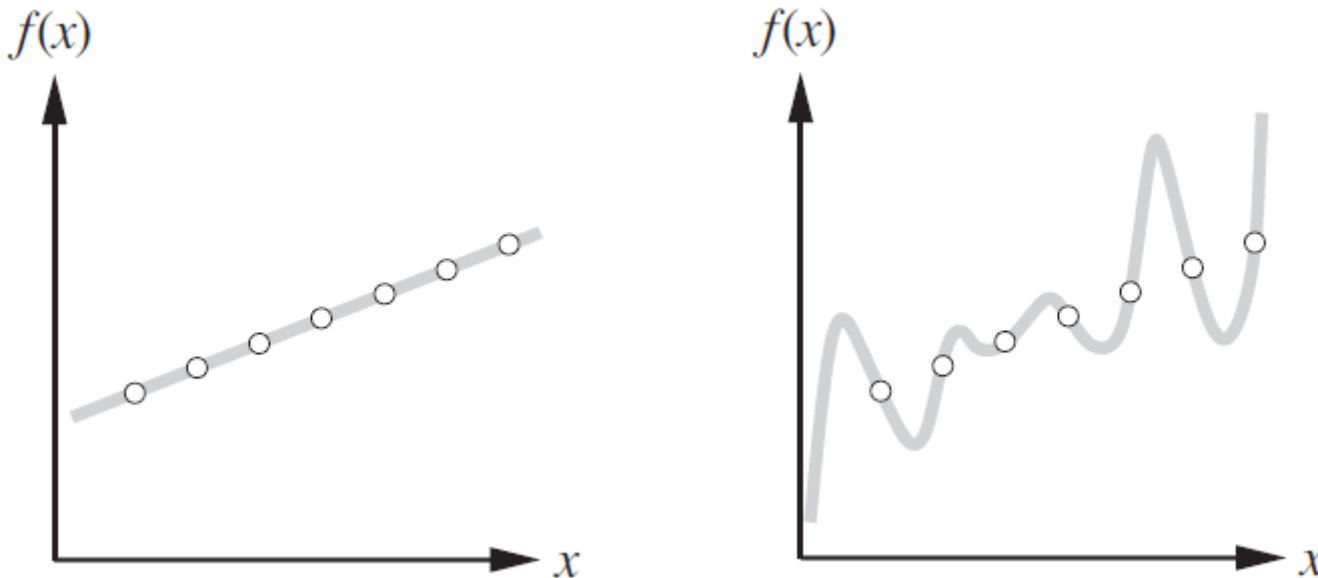
- Classification
 - Decision Trees
 - Support Vector Machines
 - Perceptron
 - Artificial Neural Networks

Supervised Learning

- Given a training set of N example input-output pairs
$$(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$$
 - where each y_j was generated by an unknown function $y = f(x)$,
 - discover a function h that approximates the true function f .
- The function h is a hypothesis
- When y is discrete \Rightarrow classification problem
- When y is continuous \Rightarrow regression problem

Hypothesis Space

- Ockham's razor:
 - Simpler hypotheses tend to generalize to future data better
 - Prefer the simplest hypothesis consistent with data



Generalization & Overfitting

Hypothesis Space

- Choosing the hypothesis h^* that is most probable given the data

$$h^* = \operatorname{argmax}_{h \in \mathcal{H}} P(h|data)$$

$$= \operatorname{argmax}_{h \in \mathcal{H}} P(data|h) P(h)$$

后验 = 先验 * 似然

- A tradeoff between the **expressiveness** of a hypothesis space and the **complexity** of finding a good hypothesis within that space.

Classification

■ Binary

Spam (17 messages, 17 Unread, Mark all as Read)

- Messages that have been in [Spam] folder more than 15 days will be automatically deleted.

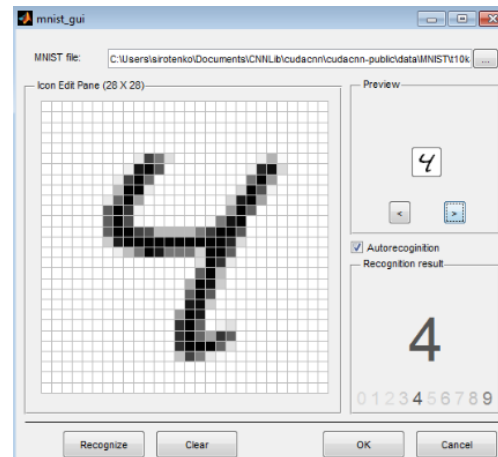
Buttons: Delete, Not spam, Move to, Mark as, More, View

From: Subject

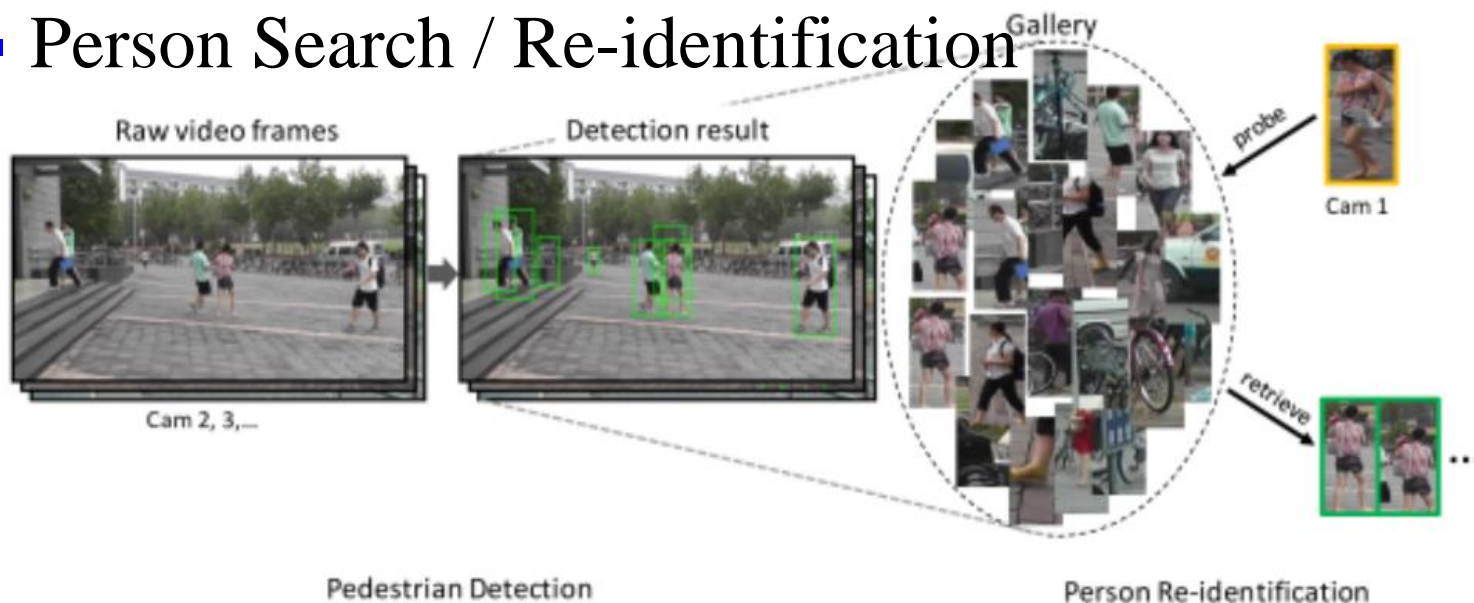
Today (7 messages)

<input type="checkbox"/>	llb2yt	约稿通知:
<input type="checkbox"/>	zhangxuejingla...	计算机AI学术交流会议【EI检索】世界科技出版公司H
<input type="checkbox"/>	ICEIT2016	Togongxj@zju.edu.cn『TPC』电子信会议【2016A
<input type="checkbox"/>	SSME2016	To:gongxj@zju.edu.cn【提交EI+CPCI】【服务科学、管
<input type="checkbox"/>	SSHD2016	To:gongxj@zju.edu.cn 诚邀您加入人文社科学术会议,更
<input type="checkbox"/>	研数	统计软件S.P.SS官方培训与认证(总27/28期)
<input type="checkbox"/>	MEME2016	[MEME2016] 武汉2016/7/5 截稿时间 EI & CPCI收录 0:

■ Multi-class

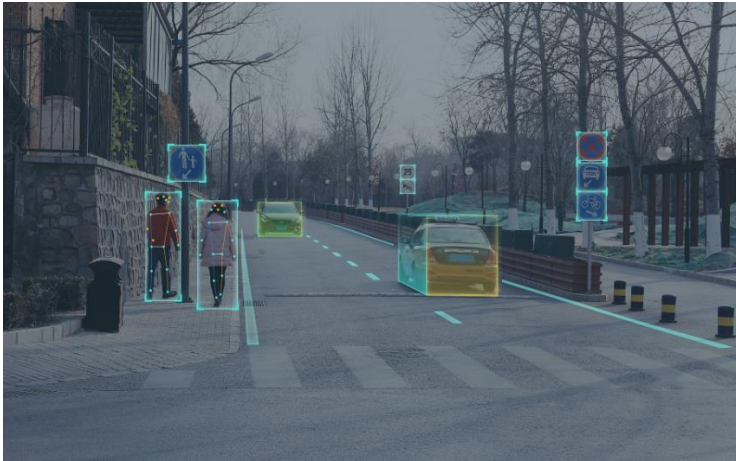


■ Person Search / Re-identification



Classification

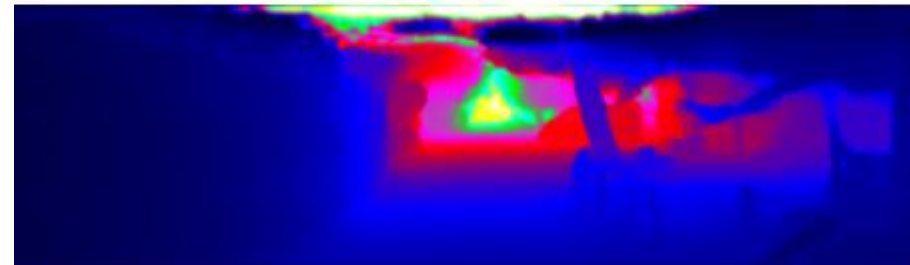
- Object Detection



- Depth Prediction



Input Image

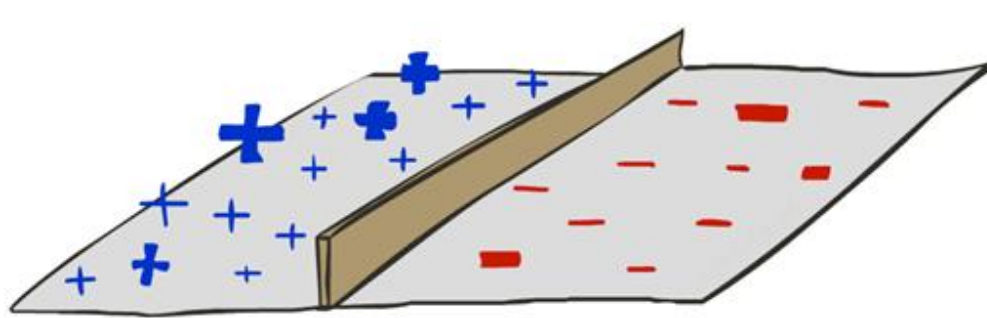


- Semantic Segmentation

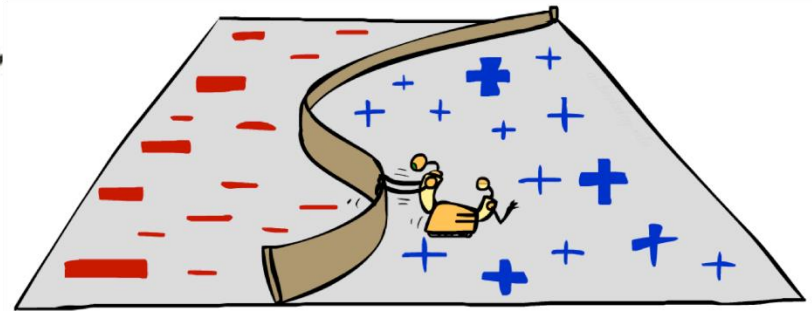


Classification

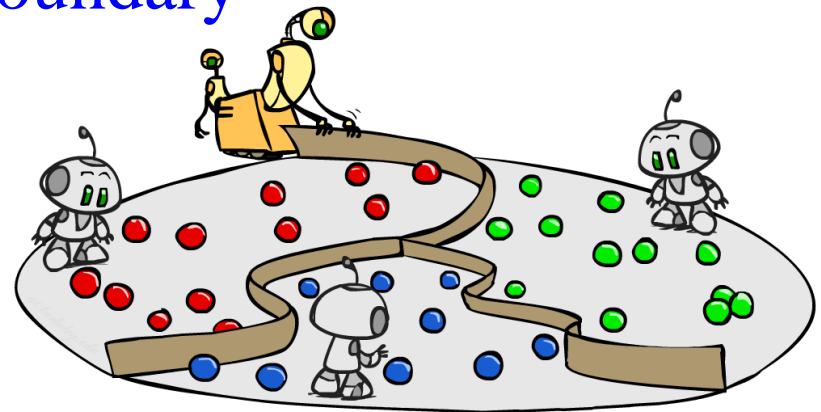
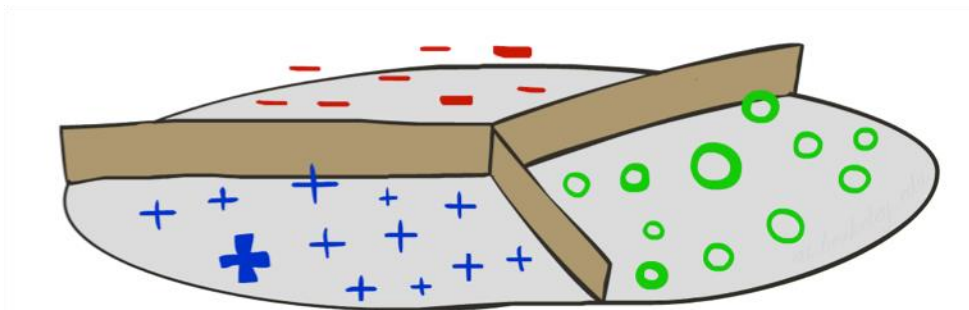
- Linear



- Nonlinear



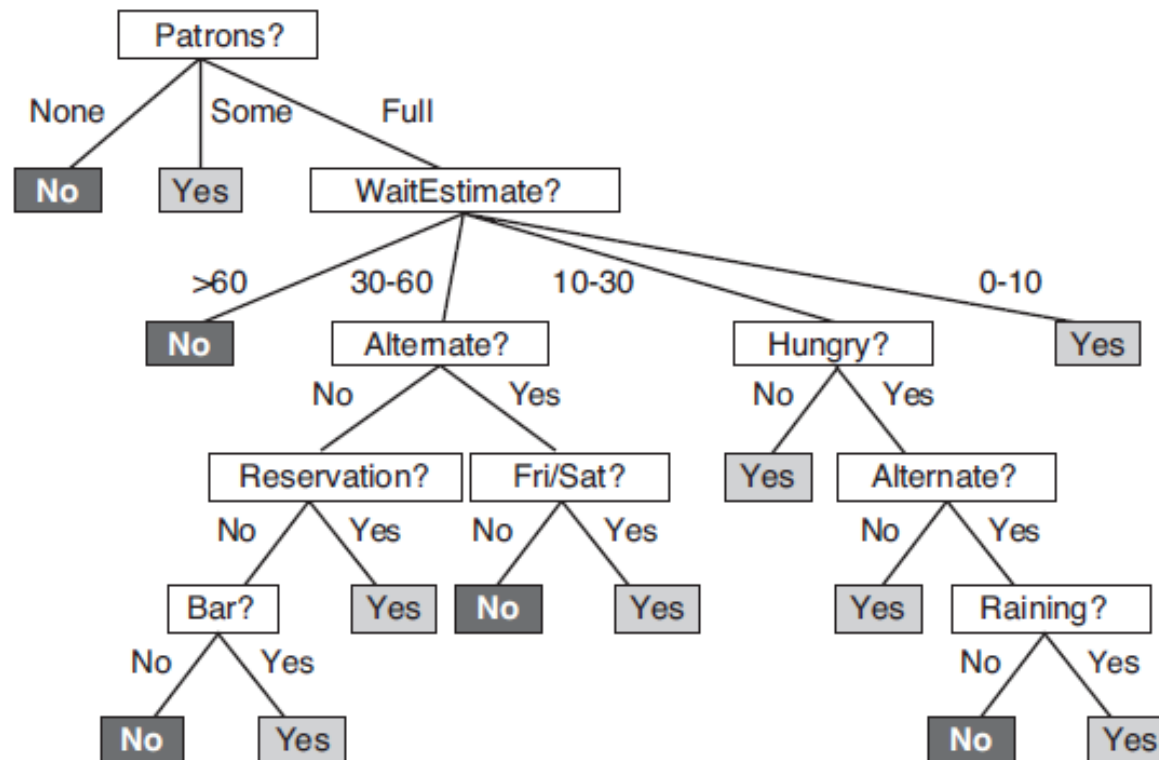
Decision boundary



Decision Trees

■ Representation

- **Input:** a vector of attribute values
- **Output:** a single output value
- Reaches its decision by performing a sequence of tests
- Each node corresponds to a test of the value of one of the input attributes



Decision Trees

- **Learning:** how to learn a good decision tree?

```
function DECISION-TREE-LEARNING(examples, attributes, parent_examples) returns  
a tree
```

```
  if examples is empty then return PLURALITY-VALUE(parent_examples)  
  else if all examples have the same classification then return the classification  
  else if attributes is empty then return PLURALITY-VALUE(examples)
```

```
  else
```

```
     $A \leftarrow \operatorname{argmax}_{a \in \text{attributes}} \text{IMPORTANCE}(a, \text{examples})$ 
```

根结点

```
    tree  $\leftarrow$  a new decision tree with root test A
```

```
    for each value  $v_k$  of A do
```

```
      exs  $\leftarrow \{e : e \in \text{examples} \text{ and } e.A = v_k\}$ 
```

```
      subtree  $\leftarrow$  DECISION-TREE-LEARNING(exs, attributes - A, examples)
```

```
      add a branch to tree with label ( $A = v_k$ ) and subtree subtree
```

```
  return tree
```

递归

Decision Trees

- **Learning:** how to learn a good decision tree?

$$A \leftarrow \operatorname{argmax}_{a \in \text{attributes}} \text{IMPORTANCE}(a, \text{examples})$$

- Maximize the **information gain**: the expected reduction in entropy

$$\text{Gain}(A) = B\left(\frac{p}{p+n}\right) - \text{Remainder}(A)$$

- The entropy of a Boolean random variable:

$$B(q) = -(q \log_2 q + (1 - q) \log_2 (1 - q))$$

- The expected entropy remaining after testing A

$$\text{Remainder}(A) = \sum_{k=1}^d \frac{p_k + n_k}{p + n} B\left(\frac{p_k}{p_k + n_k}\right)$$

Decision Trees

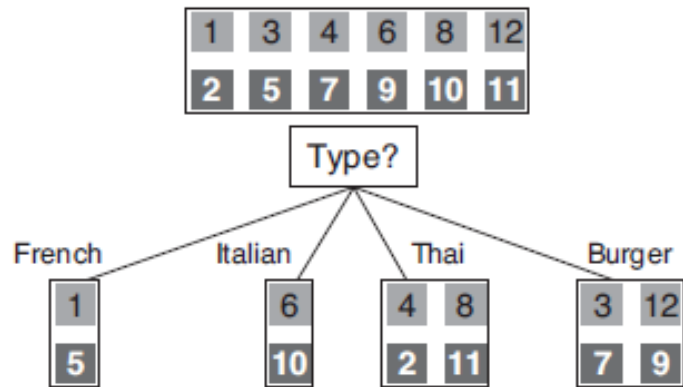
▪ Example

Example	Input Attributes										Goal
	<i>Alt</i>	<i>Bar</i>	<i>Fri</i>	<i>Hun</i>	<i>Pat</i>	<i>Price</i>	<i>Rain</i>	<i>Res</i>	<i>Type</i>	<i>Est</i>	<i>WillWait</i>
x_1	<i>Yes</i>	<i>No</i>	<i>No</i>	<i>Yes</i>	<i>Some</i>	<i>\$\$\$</i>	<i>No</i>	<i>Yes</i>	<i>French</i>	<i>0–10</i>	$y_1 = \text{Yes}$
x_2	<i>Yes</i>	<i>No</i>	<i>No</i>	<i>Yes</i>	<i>Full</i>	<i>\$</i>	<i>No</i>	<i>No</i>	<i>Thai</i>	<i>30–60</i>	$y_2 = \text{No}$
x_3	<i>No</i>	<i>Yes</i>	<i>No</i>	<i>No</i>	<i>Some</i>	<i>\$</i>	<i>No</i>	<i>No</i>	<i>Burger</i>	<i>0–10</i>	$y_3 = \text{Yes}$
x_4	<i>Yes</i>	<i>No</i>	<i>Yes</i>	<i>Yes</i>	<i>Full</i>	<i>\$</i>	<i>Yes</i>	<i>No</i>	<i>Thai</i>	<i>10–30</i>	$y_4 = \text{Yes}$
x_5	<i>Yes</i>	<i>No</i>	<i>Yes</i>	<i>No</i>	<i>Full</i>	<i>\$\$\$</i>	<i>No</i>	<i>Yes</i>	<i>French</i>	<i>>60</i>	$y_5 = \text{No}$
x_6	<i>No</i>	<i>Yes</i>	<i>No</i>	<i>Yes</i>	<i>Some</i>	<i>\$\$</i>	<i>Yes</i>	<i>Yes</i>	<i>Italian</i>	<i>0–10</i>	$y_6 = \text{Yes}$
x_7	<i>No</i>	<i>Yes</i>	<i>No</i>	<i>No</i>	<i>None</i>	<i>\$</i>	<i>Yes</i>	<i>No</i>	<i>Burger</i>	<i>0–10</i>	$y_7 = \text{No}$
x_8	<i>No</i>	<i>No</i>	<i>No</i>	<i>Yes</i>	<i>Some</i>	<i>\$\$</i>	<i>Yes</i>	<i>Yes</i>	<i>Thai</i>	<i>0–10</i>	$y_8 = \text{Yes}$
x_9	<i>No</i>	<i>Yes</i>	<i>Yes</i>	<i>No</i>	<i>Full</i>	<i>\$</i>	<i>Yes</i>	<i>No</i>	<i>Burger</i>	<i>>60</i>	$y_9 = \text{No}$
x_{10}	<i>Yes</i>	<i>Yes</i>	<i>Yes</i>	<i>Yes</i>	<i>Full</i>	<i>\$\$\$</i>	<i>No</i>	<i>Yes</i>	<i>Italian</i>	<i>10–30</i>	$y_{10} = \text{No}$
x_{11}	<i>No</i>	<i>No</i>	<i>No</i>	<i>No</i>	<i>None</i>	<i>\$</i>	<i>No</i>	<i>No</i>	<i>Thai</i>	<i>0–10</i>	$y_{11} = \text{No}$
x_{12}	<i>Yes</i>	<i>Yes</i>	<i>Yes</i>	<i>Yes</i>	<i>Full</i>	<i>\$</i>	<i>No</i>	<i>No</i>	<i>Burger</i>	<i>30–60</i>	$y_{12} = \text{Yes}$

Figure 18.3 Examples for the restaurant domain.

Decision Trees

Example

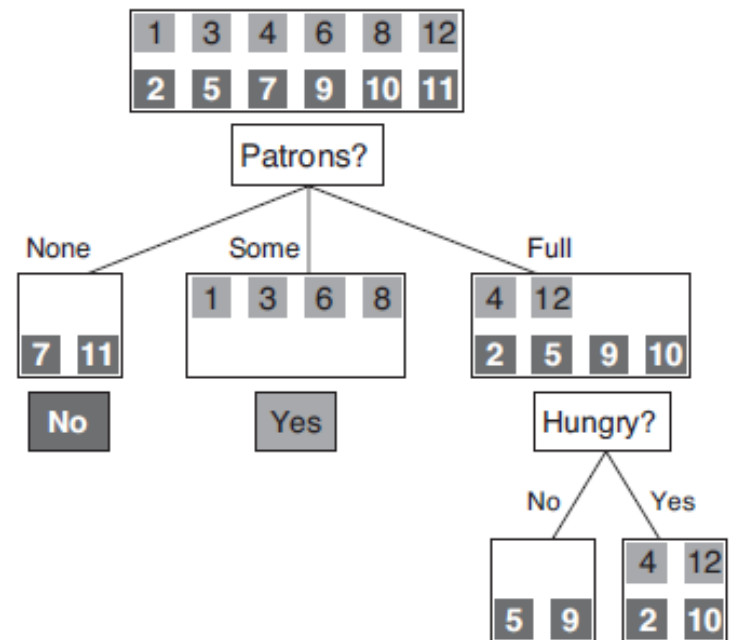


$$\begin{aligned}
 &Gain(Type) \\
 &= 1 - \left[\frac{2}{12} B\left(\frac{1}{2}\right) + \frac{2}{12} B\left(\frac{1}{2}\right) + \frac{4}{12} B\left(\frac{2}{4}\right) + \frac{4}{12} B\left(\frac{2}{4}\right) \right] \\
 &= 0 \text{ bits}
 \end{aligned}$$

$$Gain(A) = B\left(\frac{p}{p+n}\right) - Remainder(A)$$

$$B(q) = -(q \log_2 q + (1 - q) \log_2 (1 - q))$$

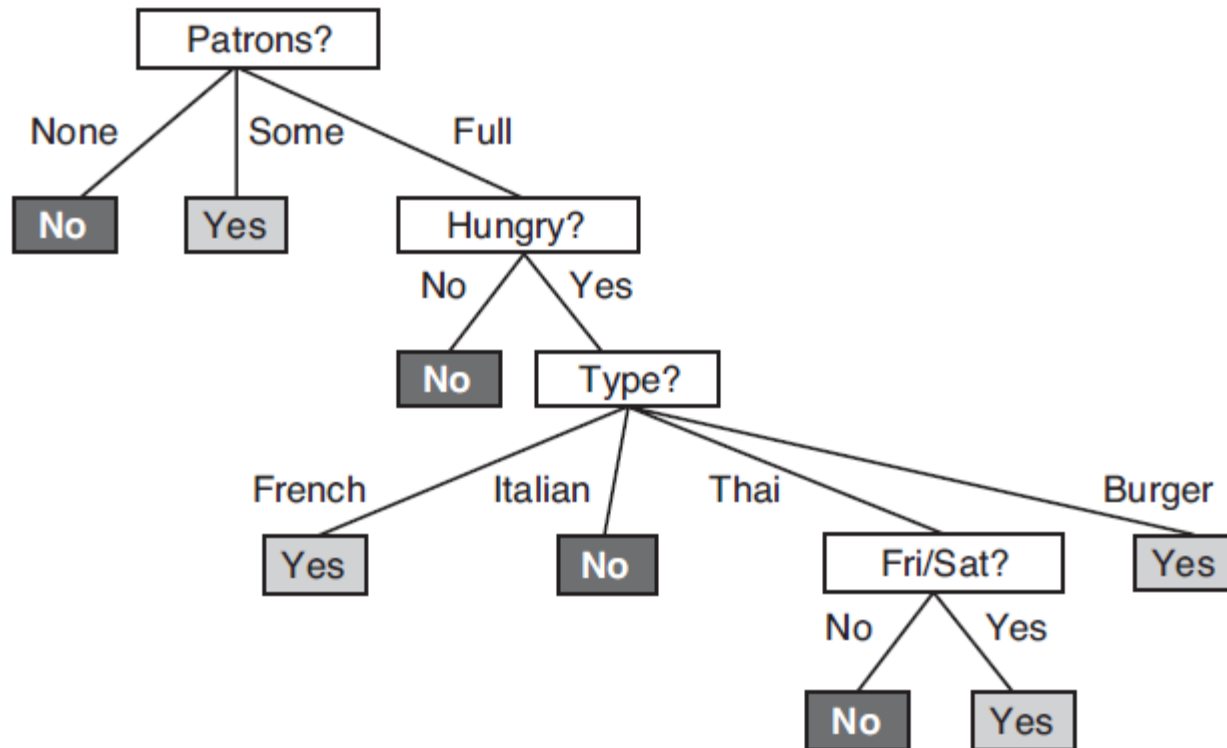
$$Remainder(A) = \sum_{k=1}^d \frac{p_k + n_k}{p + n} B\left(\frac{p_k}{p_k + n_k}\right)$$



$$\begin{aligned}
 &Gain(Patrons) \\
 &= 1 - \left[\frac{2}{12} B\left(\frac{0}{2}\right) + \frac{4}{12} B\left(\frac{4}{4}\right) + \frac{6}{12} B\left(\frac{2}{6}\right) \right] \\
 &\approx 0.541 \text{ bits}
 \end{aligned}$$

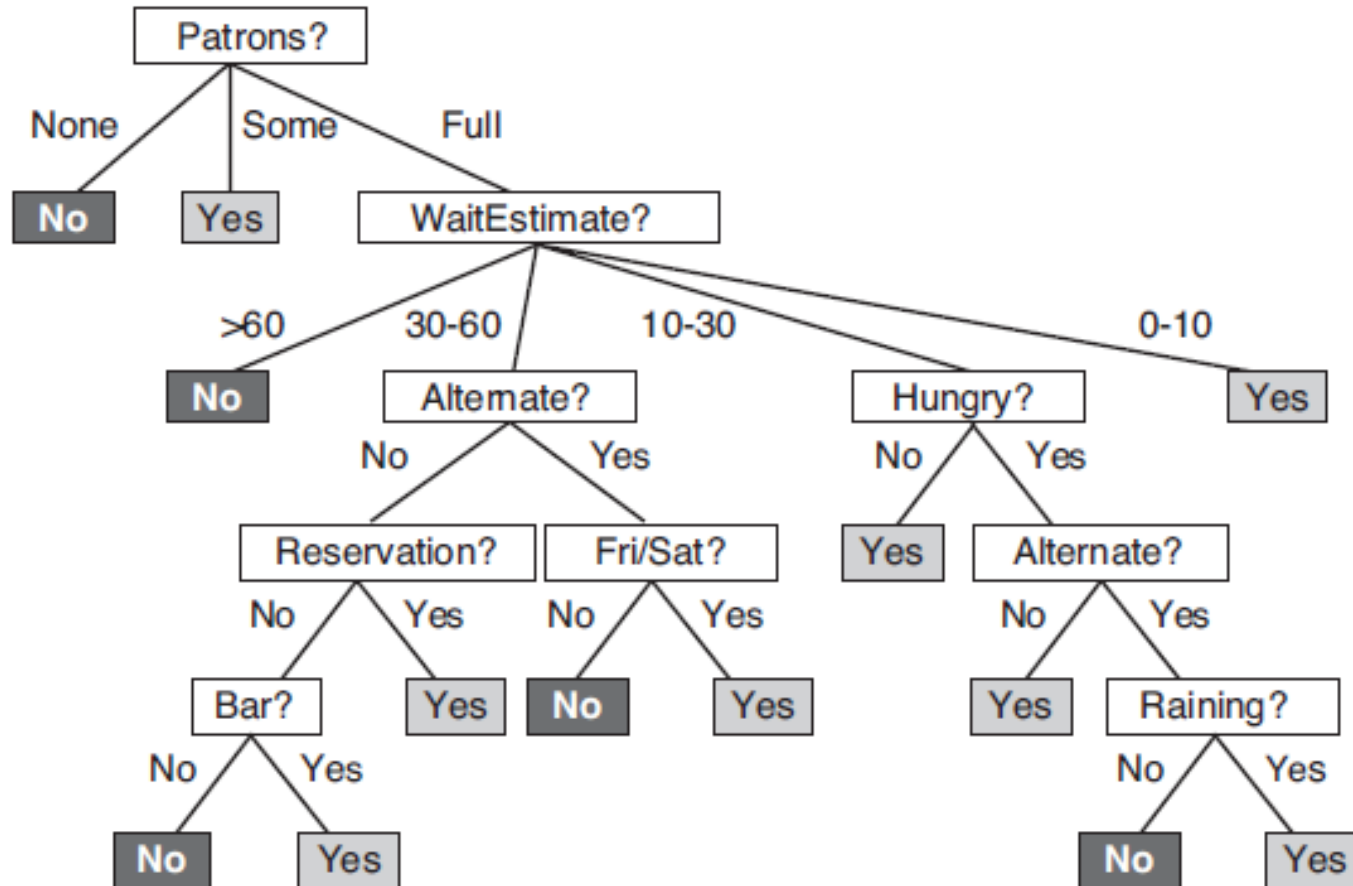
Decision Trees

- Example



Decision Trees

- Example

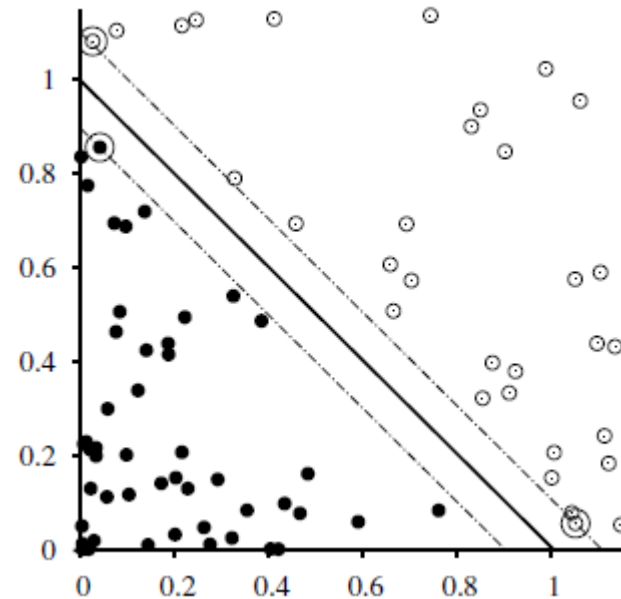
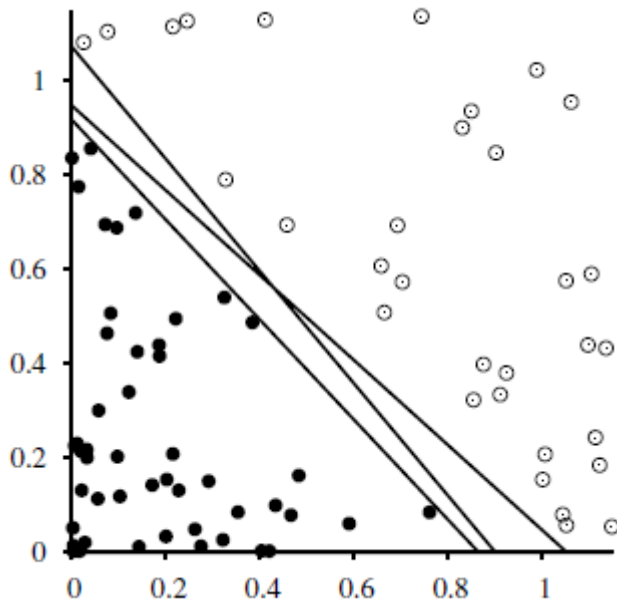


Decision Trees

- The decision-tree-learning algorithm adopts a **greedy divide-and-conquer strategy**.
- The learning algorithm looks at the **examples**, the set of examples is crucial for constructing the tree.
- For decision trees, **decision tree pruning** combats overfitting.
- DT is possible for a human to understand the reason for the output of the learning algorithm.

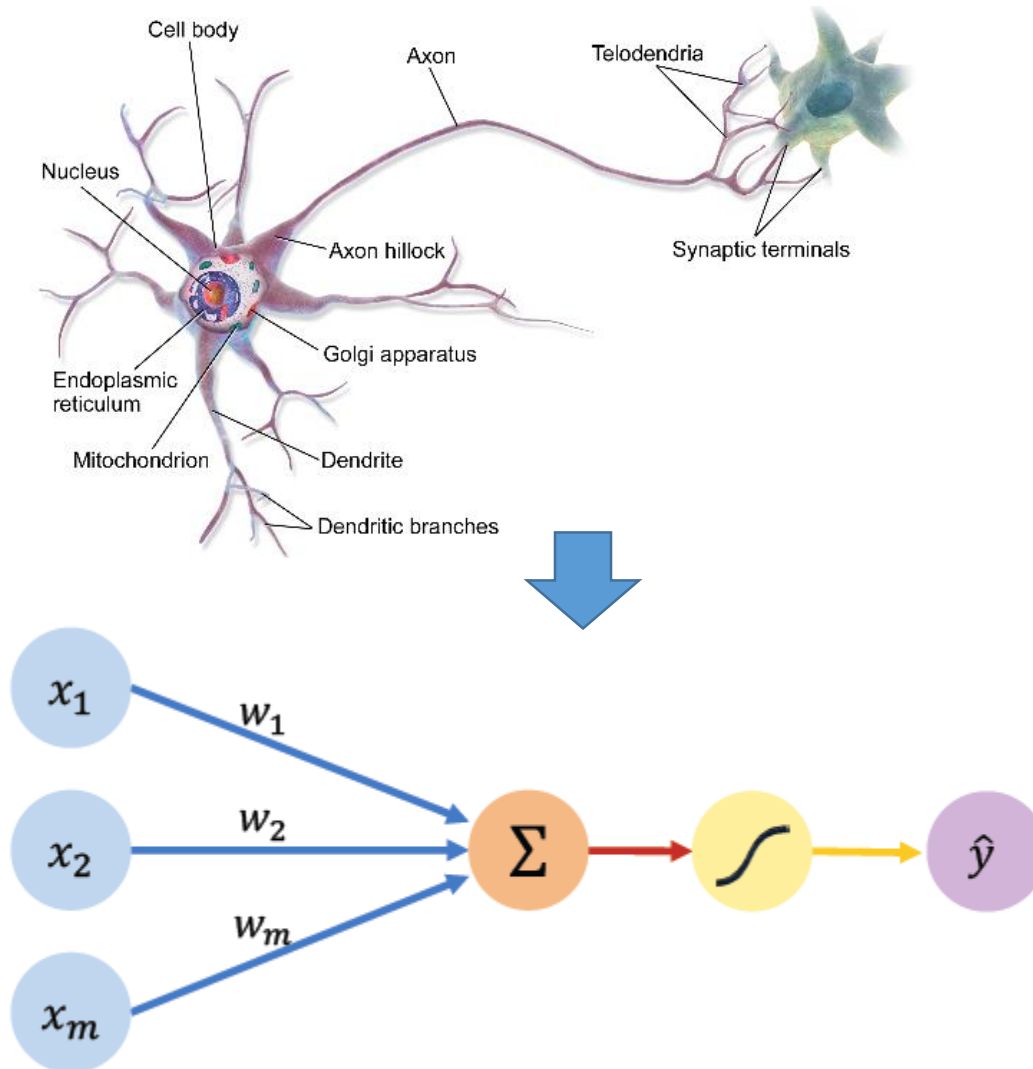
Support Vector Machines

- Construct a **maximum margin separator** —find the separator with max margin
- **Only support vectors matter**; other training examples are ignorable



Perceptron

- Inspired by human neuron



Perceptron

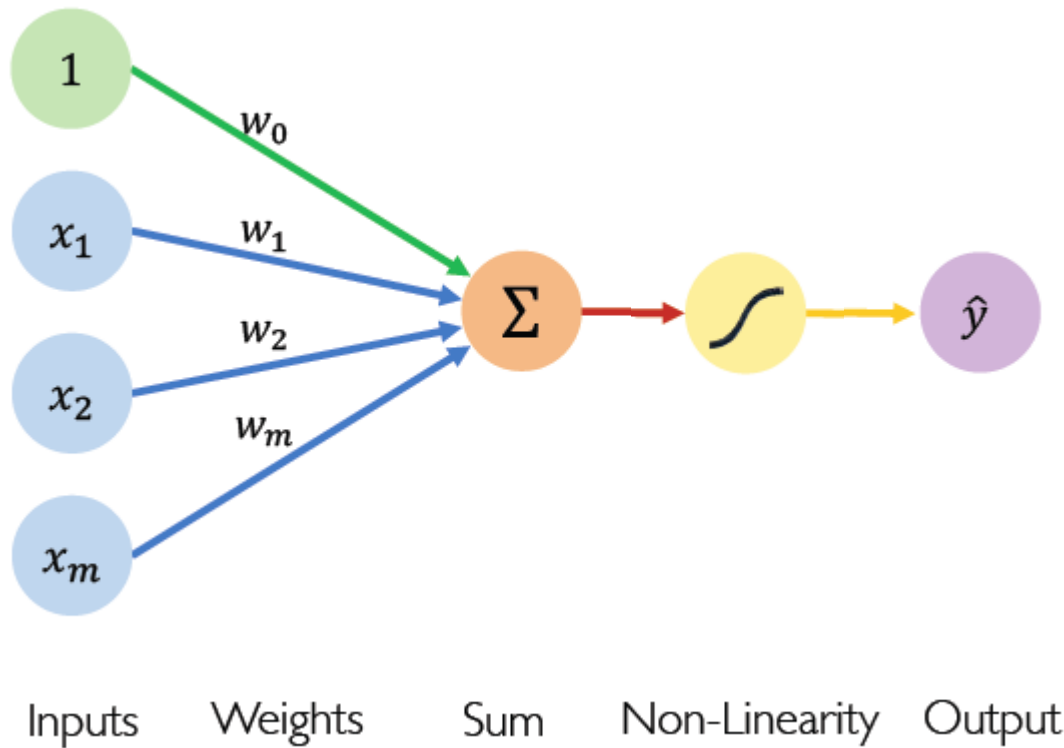


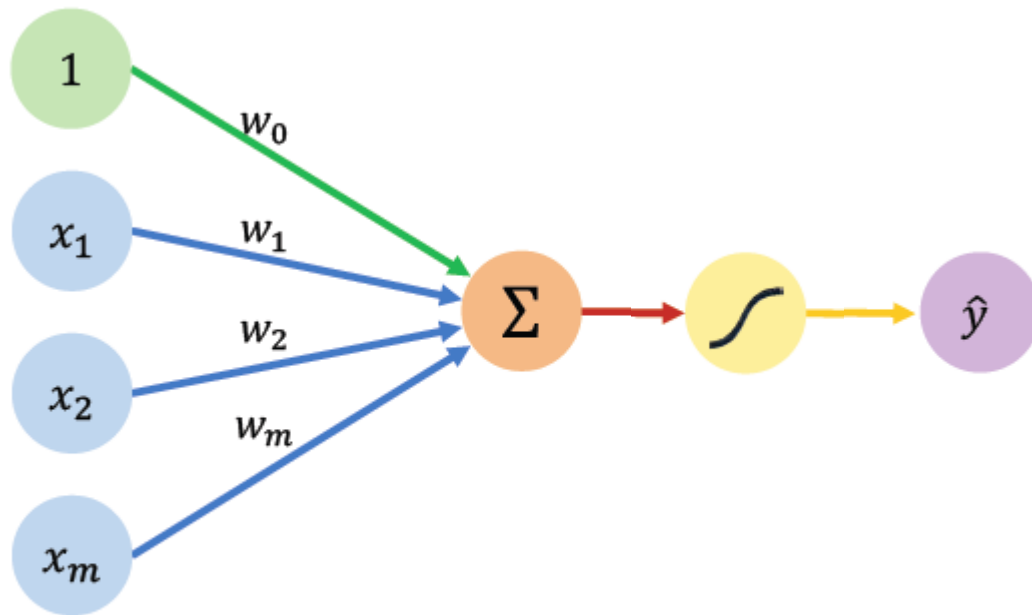
Diagram illustrating the mathematical formula for the perceptron output:

$$\hat{y} = g \left(w_0 + \sum_{i=1}^m x_i w_i \right)$$

Labels and arrows pointing to the formula components:

- Output:** Points to \hat{y} (purple arrow).
- Non-linear activation function:** Points to g (orange arrow).
- Bias:** Points to w_0 (green arrow).
- Linear combination of inputs:** Points to the summation term $\sum_{i=1}^m x_i w_i$ (red arrow).

Perceptron



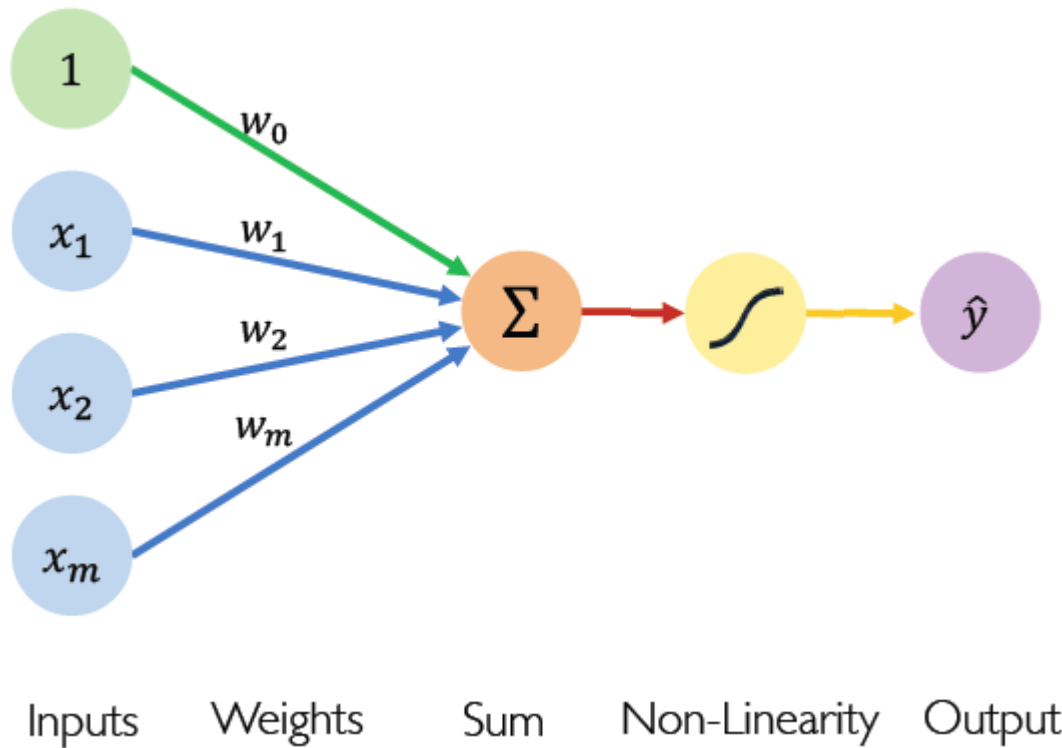
Inputs Weights Sum Non-Linearity Output

$$\hat{y} = g \left(w_0 + \sum_{i=1}^m x_i w_i \right)$$

$$\hat{y} = g (w_0 + \mathbf{X}^T \mathbf{W})$$

$$\text{where: } \mathbf{X} = \begin{bmatrix} x_1 \\ \vdots \\ x_m \end{bmatrix} \text{ and } \mathbf{W} = \begin{bmatrix} w_1 \\ \vdots \\ w_m \end{bmatrix}$$

Perceptron

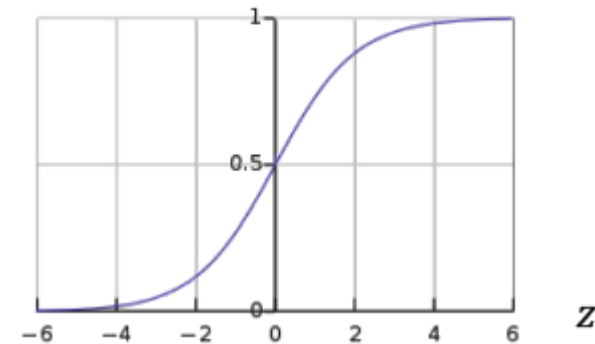


Activation Functions

$$\hat{y} = g(w_0 + \mathbf{X}^T \mathbf{W})$$

- Example: sigmoid function

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

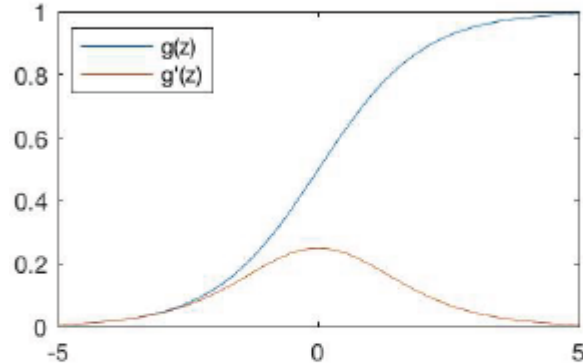


Perceptron

- Activation function:

- The purpose of activation functions is to introduce nonlinearities into the network
- All activation functions are non-linear

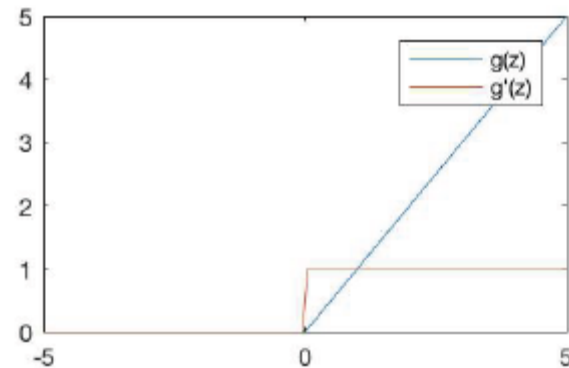
Sigmoid Function



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

$$g'(z) = g(z)(1 - g(z))$$

Rectified Linear Unit (ReLU)

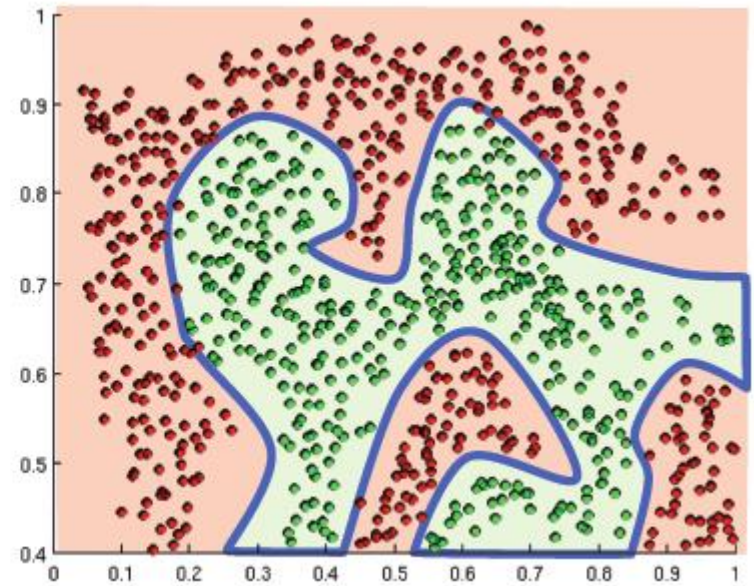
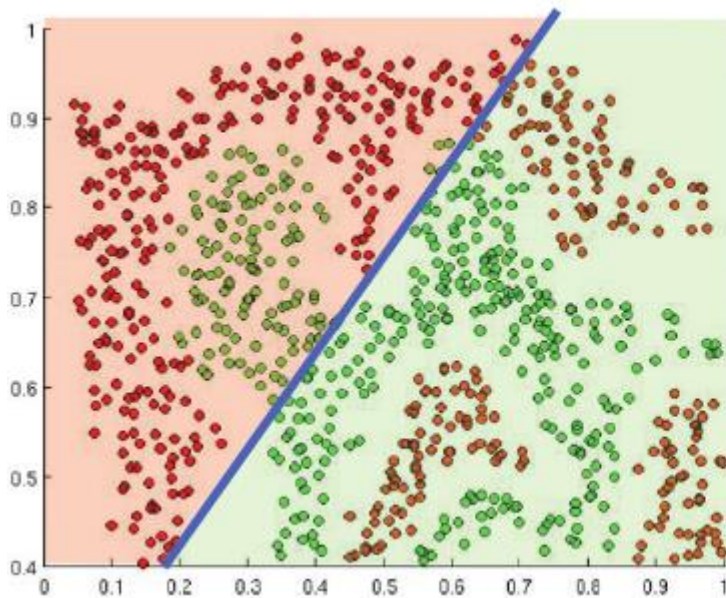


$$g(z) = \max(0, z)$$

$$g'(z) = \begin{cases} 1, & z > 0 \\ 0, & \text{otherwise} \end{cases}$$

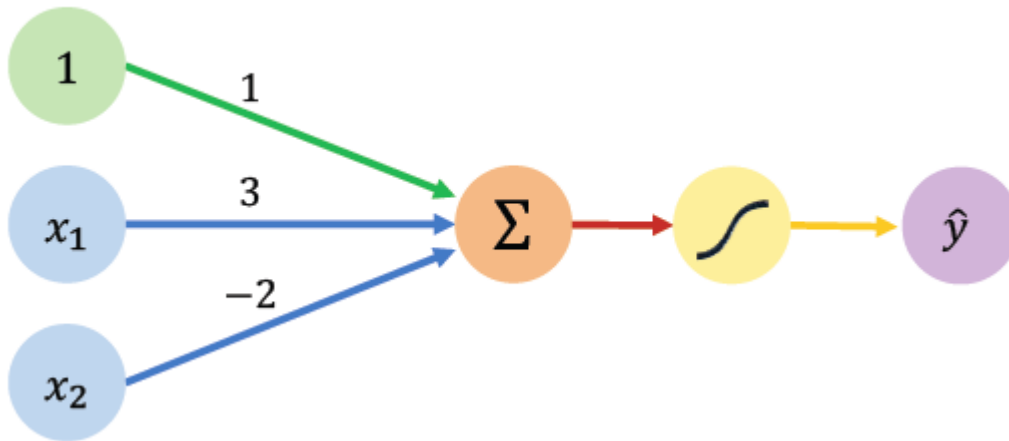
Perceptron

- Activation function:
 - Linear activation functions produce linear decisions no matter the network size
 - Nonlinearities allow us to approximate arbitrarily complex functions



Perceptron

- Example



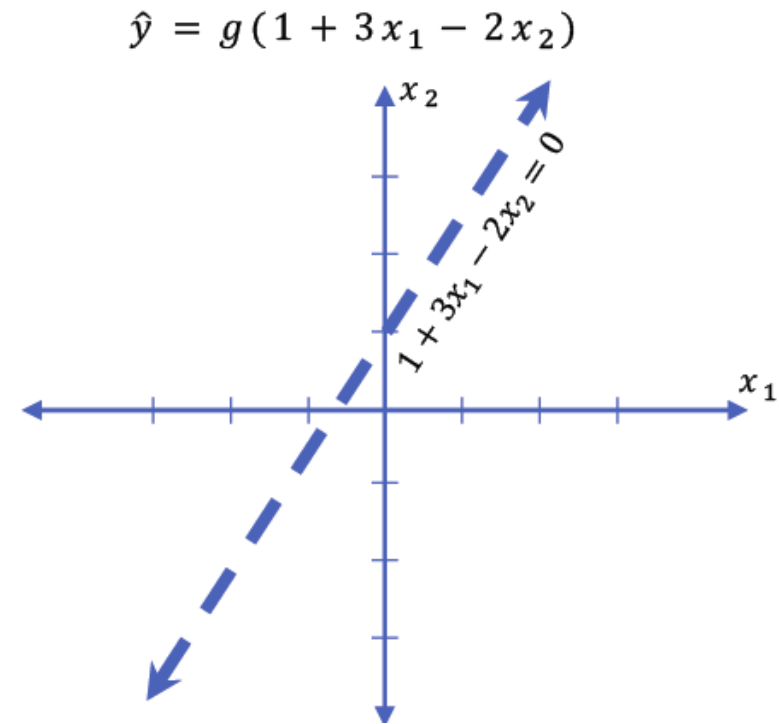
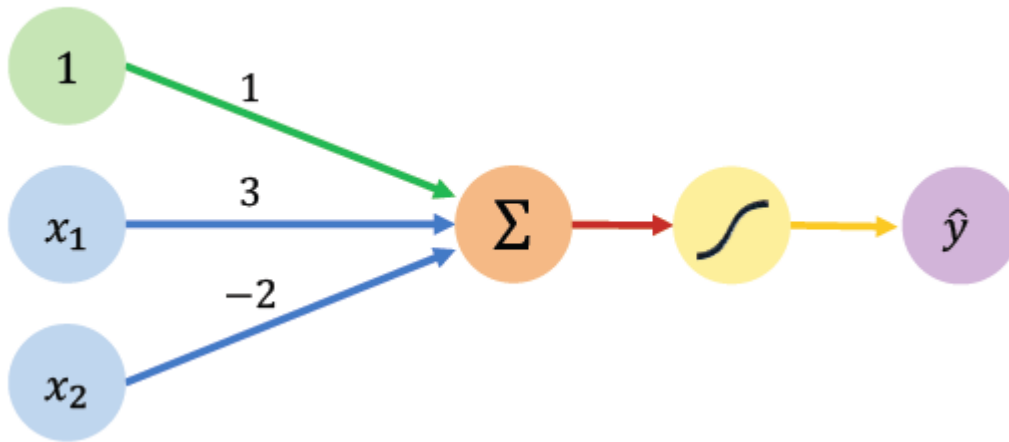
We have: $w_0 = 1$ and $\mathbf{W} = \begin{bmatrix} 3 \\ -2 \end{bmatrix}$

$$\begin{aligned}\hat{y} &= g(w_0 + \mathbf{X}^T \mathbf{W}) \\ &= g\left(1 + \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}^T \begin{bmatrix} 3 \\ -2 \end{bmatrix}\right) \\ \hat{y} &= g(1 + 3x_1 - 2x_2)\end{aligned}$$

This is just a line in 2D!

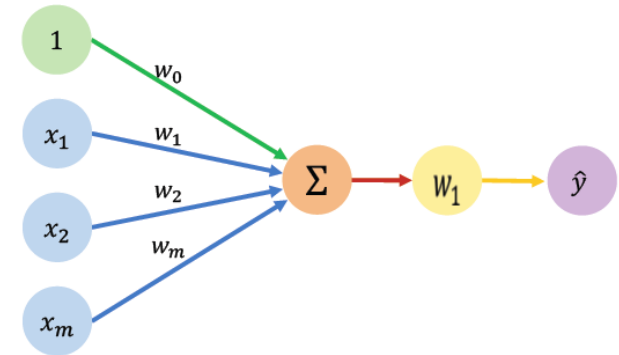
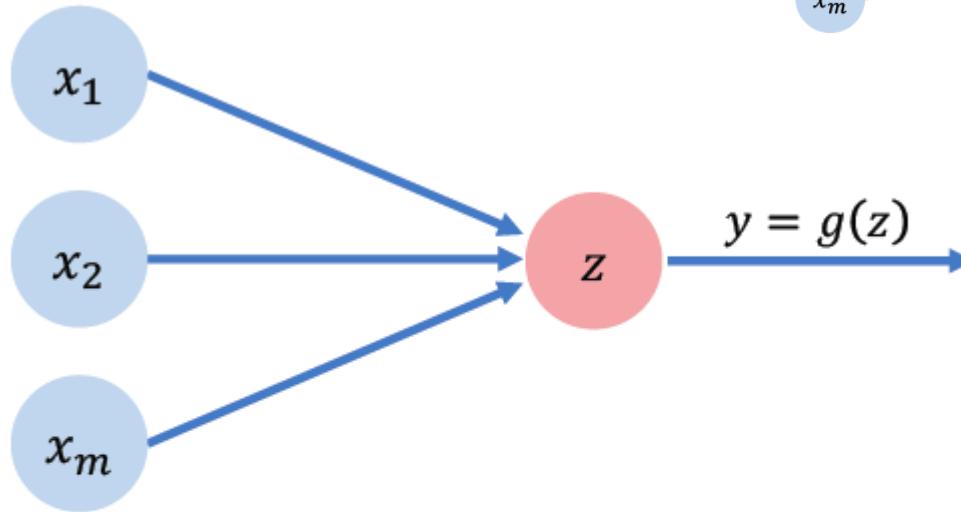
Perceptron

- Example



Perceptron

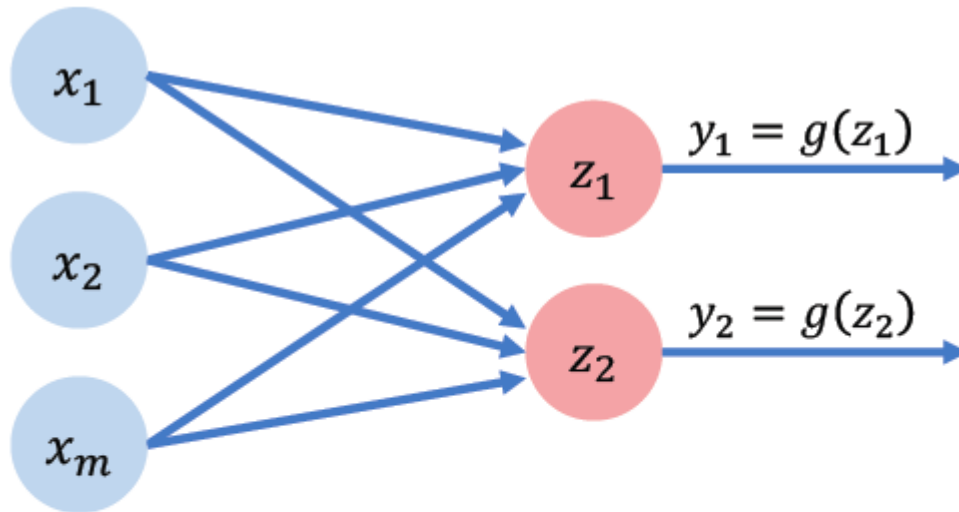
- The Simplified Perceptron



$$z = w_0 + \sum_{j=1}^m x_j w_j$$

Perceptron

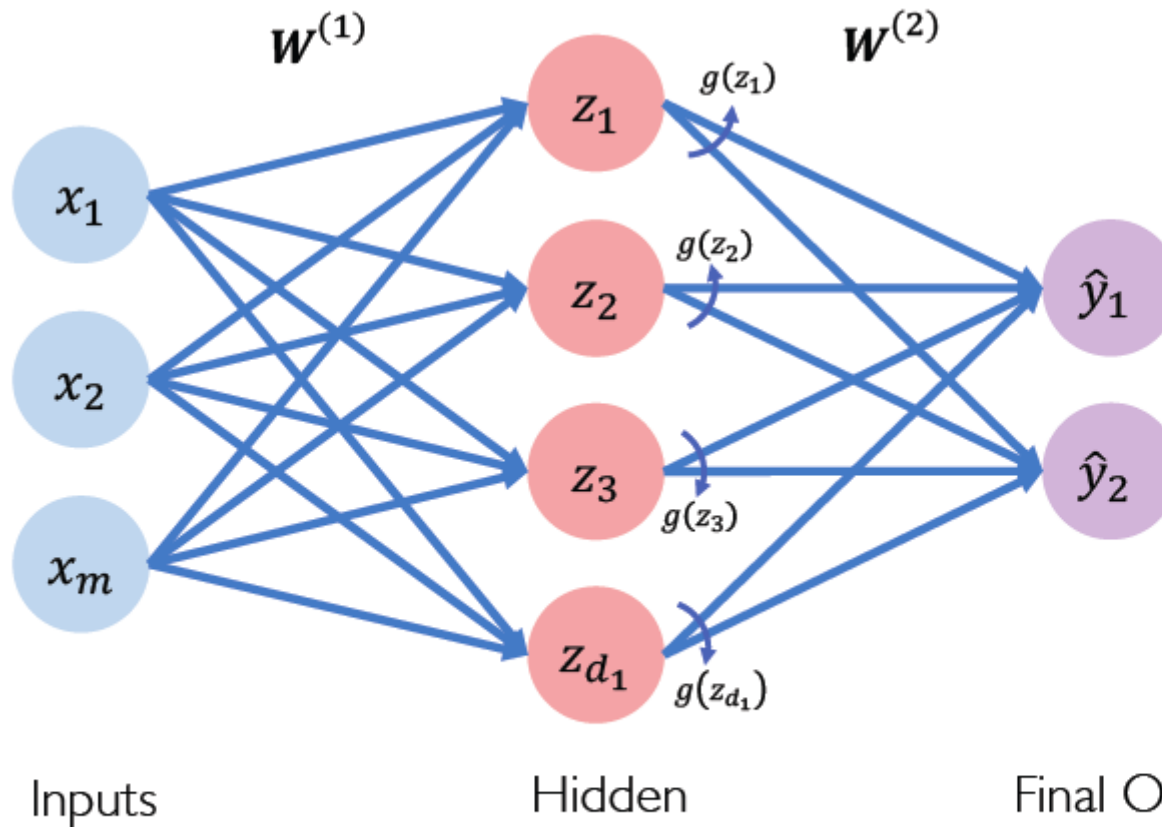
- Multi Output Perceptron



$$z_i = w_{0,i} + \sum_{j=1}^m x_j w_{j,i}$$

Artificial Neural Network

- Single Layer Neural Network



$$z_i = w_{0,i}^{(1)} + \sum_{j=1}^m x_j w_{j,i}^{(1)} \quad \hat{y}_i = g \left(w_{0,i}^{(2)} + \sum_{j=1}^{d_1} z_j w_{j,i}^{(2)} \right)$$

Artificial Neural Network

- Loss Optimization

$$\mathbf{W}^* = \operatorname{argmin}_{\mathbf{W}} \frac{1}{n} \sum_{i=1}^n \mathcal{L}(f(x^{(i)}; \mathbf{W}), y^{(i)})$$

$$\mathbf{W}^* = \operatorname{argmin}_{\mathbf{W}} J(\mathbf{W})$$



Remember:

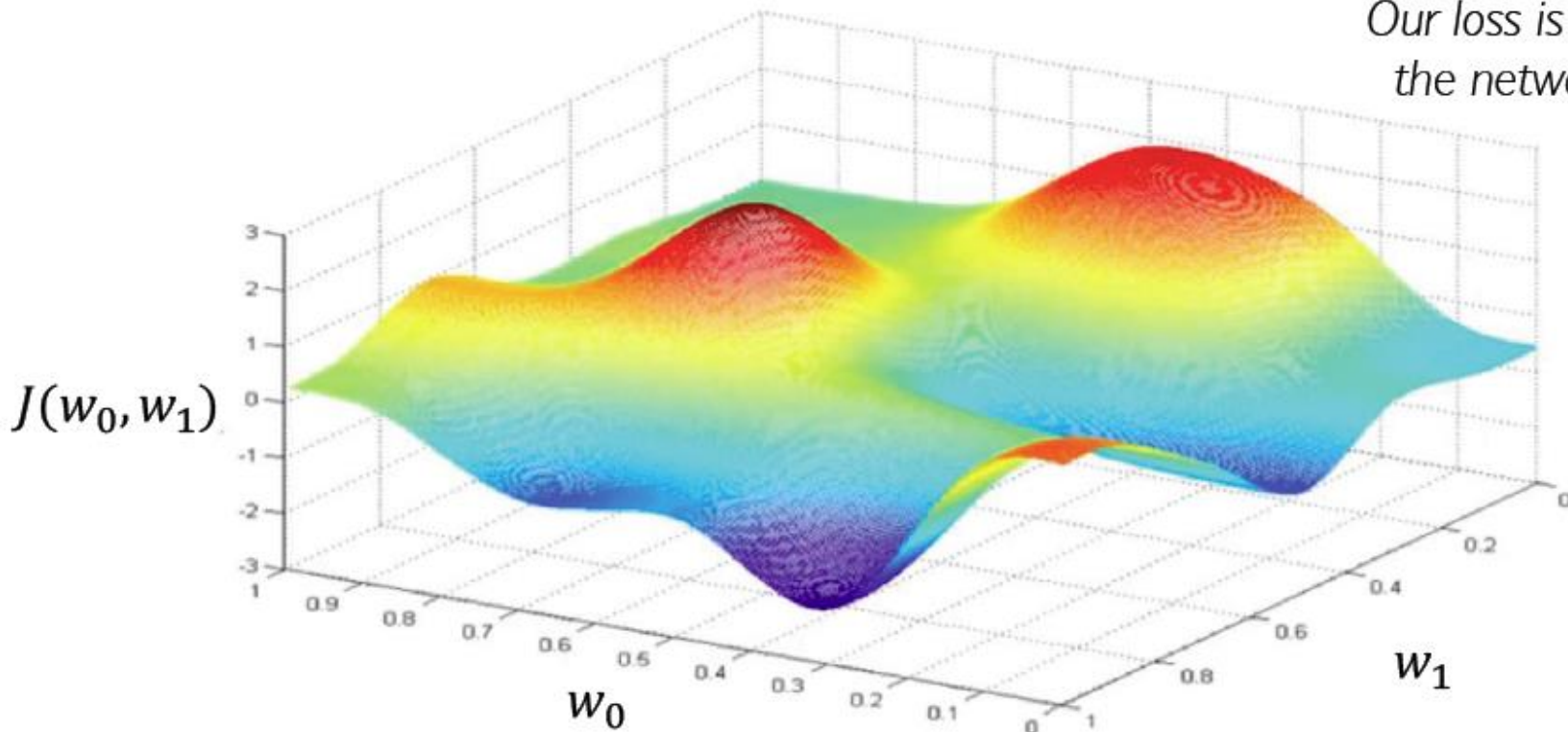
$$\mathbf{W} = \{\mathbf{W}^{(0)}, \mathbf{W}^{(1)}, \dots\}$$

Artificial Neural Network

- Loss Optimization – Gradient Descent

$$W^* = \underset{W}{\operatorname{argmin}} J(W)$$

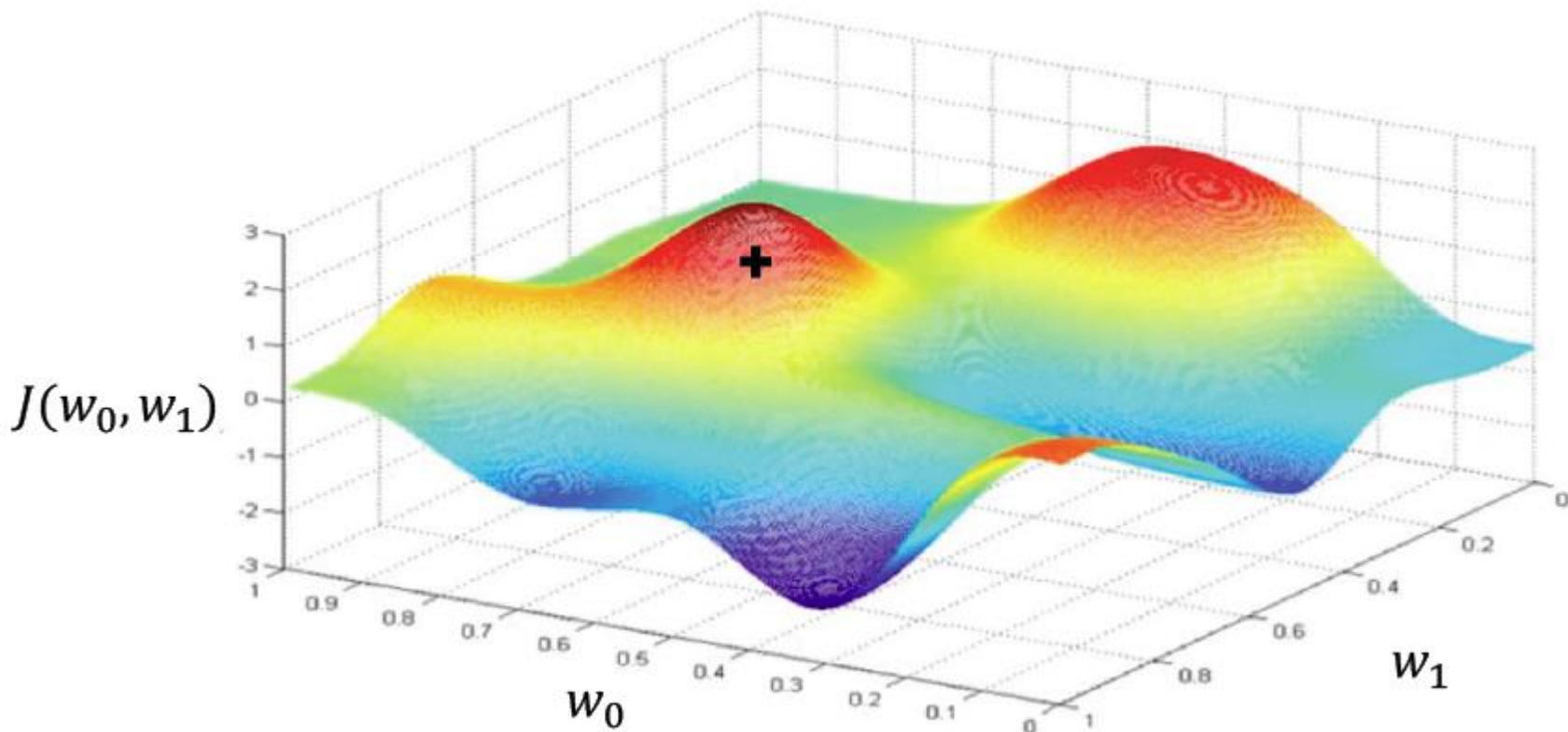
Remember:
*Our loss is a function of
the network weights!*



Artificial Neural Network

- Loss Optimization – Gradient Descent

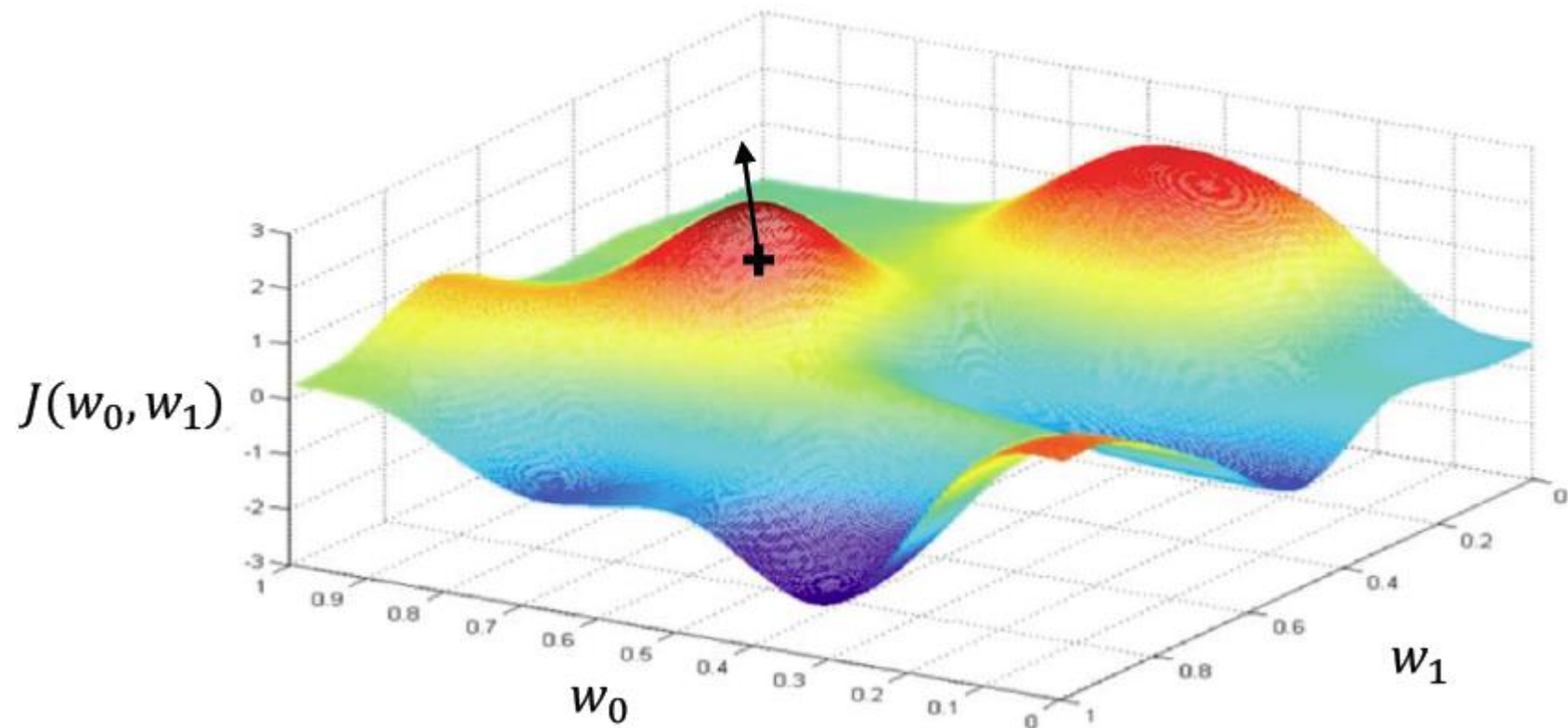
Randomly pick an initial (w_0, w_1)



Artificial Neural Network

- Loss Optimization – Gradient Descent

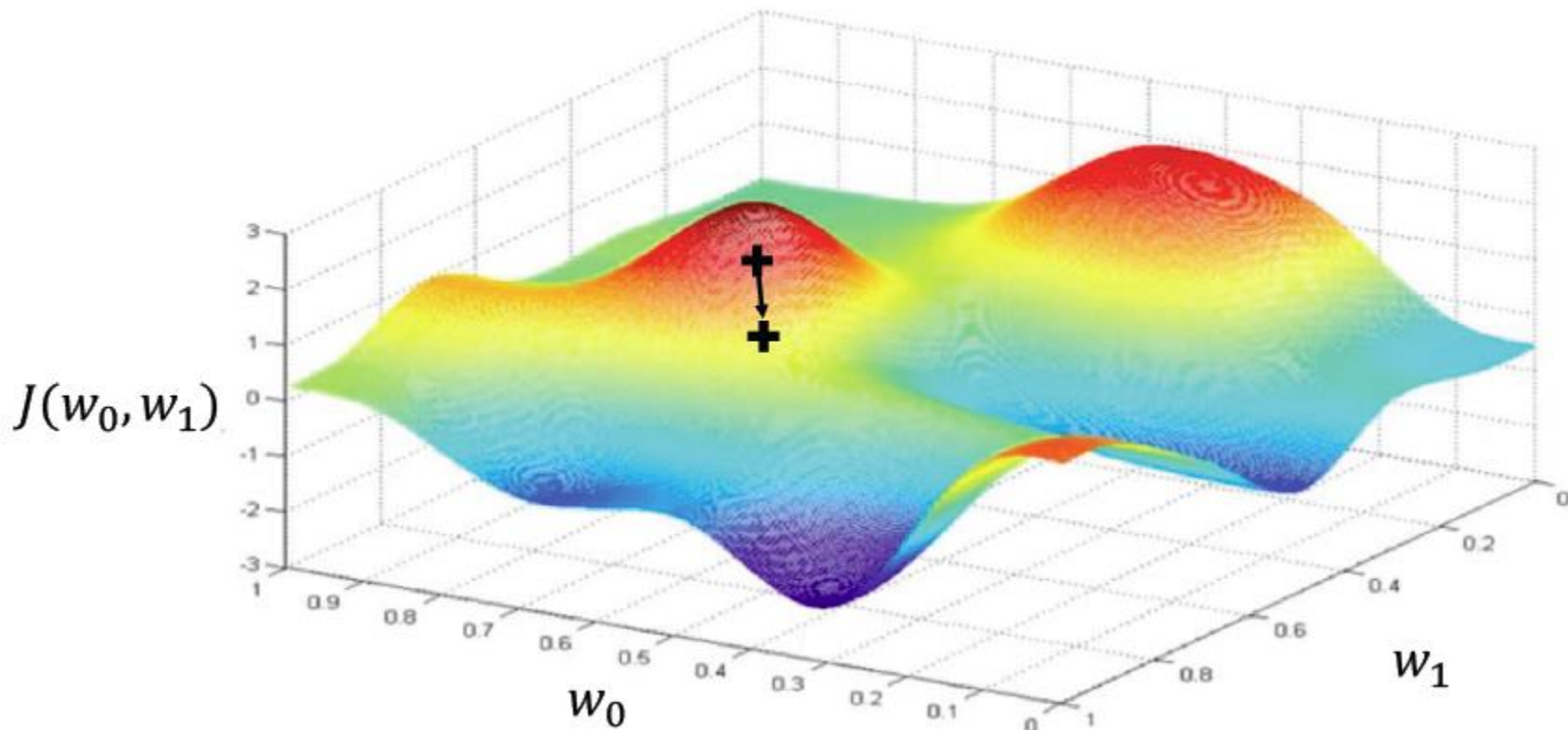
Compute gradient, $\frac{\partial J(W)}{\partial W}$



Artificial Neural Network

- Loss Optimization – Gradient Descent

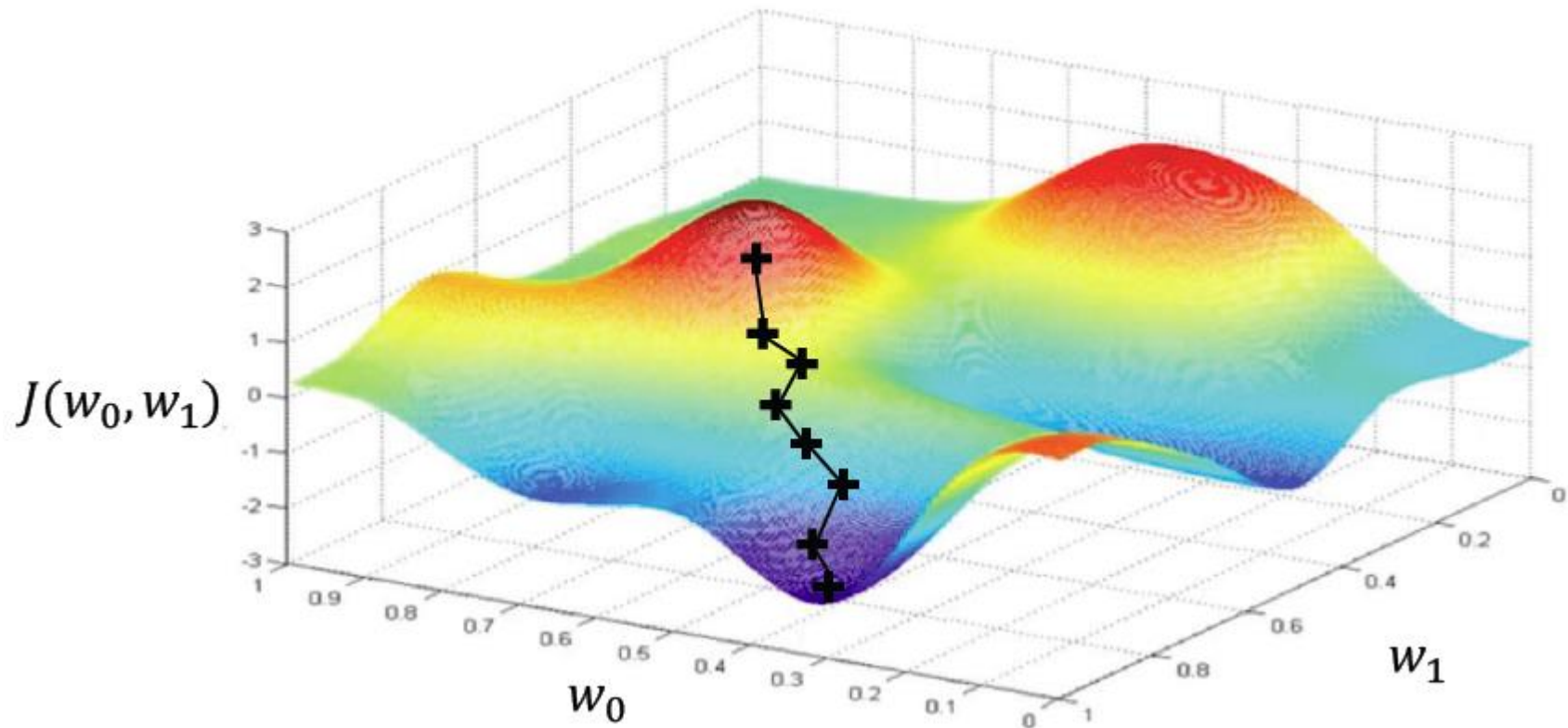
Take small step in opposite direction of gradient



Artificial Neural Network

- Loss Optimization – Gradient Descent

Repeat until convergence



Artificial Neural Network

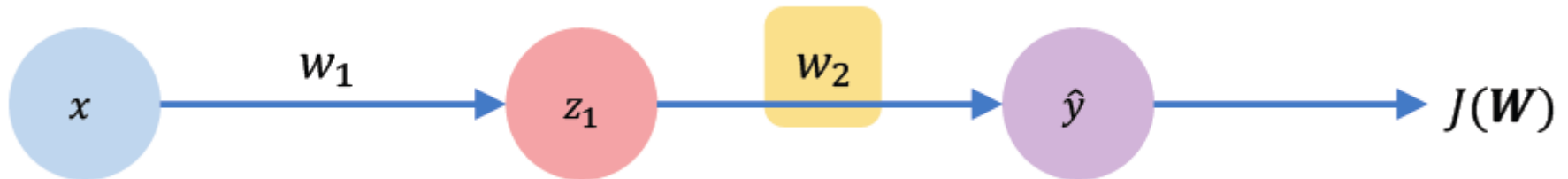
▪ Loss Optimization – Gradient Descent

Algorithm

1. Initialize weights randomly $\sim \mathcal{N}(0, \sigma^2)$
2. Loop until convergence:
3. Compute gradient, $\frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
4. Update weights, $\mathbf{W} \leftarrow \mathbf{W} - \eta \frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
5. Return weights

Artificial Neural Network

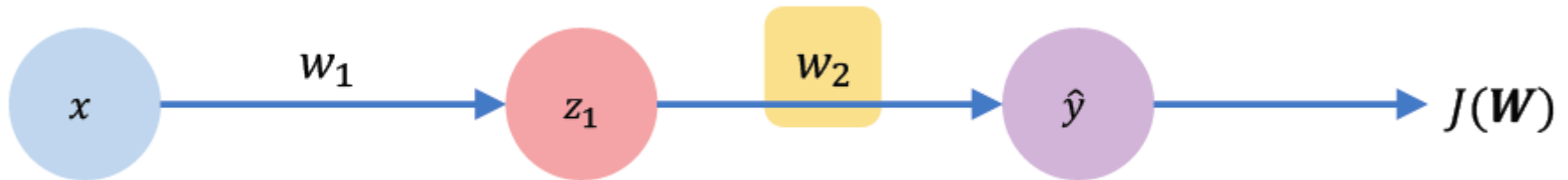
- Computing Gradients: Backpropagation



How does a small change in one weight (ex. w_2) affect the final loss $J(\mathbf{W})$?

Artificial Neural Network

- Computing Gradients: **Backpropagation**

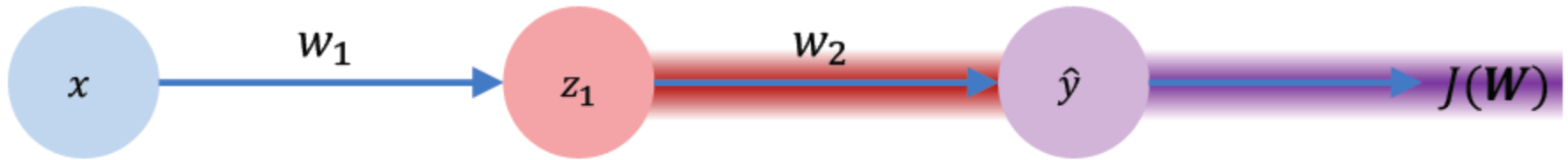


$$\frac{\partial J(\mathbf{W})}{\partial w_2} =$$

Let's use the chain rule!

Artificial Neural Network

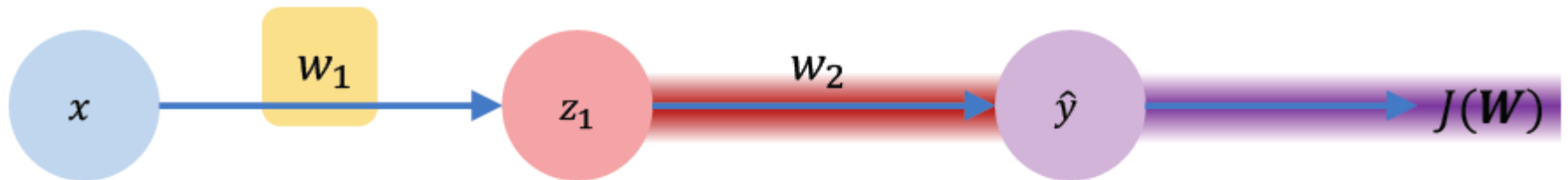
- Computing Gradients: Backpropagation



$$\frac{\partial J(W)}{\partial w_2} = \underbrace{\frac{\partial J(W)}{\partial \hat{y}}}_{\text{purple}} * \underbrace{\frac{\partial \hat{y}}{\partial w_2}}_{\text{orange}}$$

Artificial Neural Network

- Computing Gradients: **Backpropagation**

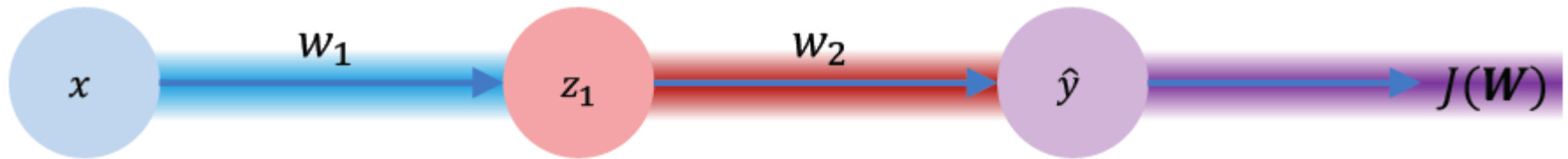


$$\frac{\partial J(W)}{\partial w_1} = \frac{\partial J(W)}{\partial \hat{y}} * \frac{\partial \hat{y}}{\partial w_1}$$

Apply chain rule! Apply chain rule!

Artificial Neural Network

- Computing Gradients: **Backpropagation**



$$\frac{\partial J(W)}{\partial w_1} = \underbrace{\frac{\partial J(W)}{\partial \hat{y}}}_{\text{purple}} * \underbrace{\frac{\partial \hat{y}}{\partial z_1}}_{\text{orange}} * \underbrace{\frac{\partial z_1}{\partial w_1}}_{\text{blue}}$$

*Repeat this for **every weight in the network** using gradients from later layers*

Artificial Neural Network

- Stacking Perceptrons to form neural networks
- Optimization through backpropagation

Reducing the Dimensionality of Data with Neural Networks

G. E. Hinton^{*}, R. R. Salakhutdinov

+ See all authors and affiliations

Science 28 Jul 2006:
Vol. 313, Issue 5786, pp. 504-507
DOI: 10.1126/science.1127647

Readings

- Artificial Intelligence
 - Chapter 18.1 -18.5, 18.7, 18.9
- Homework 6