



Machine Learning

A. M. Sadeghzadeh, Ph.D.

Sharif University of Technology
Computer Engineering Department (CE)
Data and Network Security Lab (DNSL)

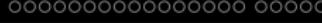
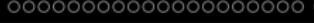
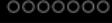
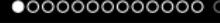


September 24, 2024

Most slides have been adapted from Bhiksha Raj, 11-785, CMU 2020, Fei Fei Li, cs231n, Stanford 2017, and Soleymani, CE40717, Sharif 2020.

Schedule

- 1 Machine Learning
- 2 Supervised Learning
- 3 Neural Networks
- 4 Logistic Regression
- 5 Loss Function
- 6 Optimization
- 7 Regularization



Machine Learning

Definition of machine learning

Tom Mitchell (1998)

A computer program is said to **learn** from **experience E** with respect to some class of **tasks T** and **performance measure P**, if its performance at tasks in **T**, as measured by **P**, improves with experience **E**.

Example

Task: Classifying emails as spam or not spam

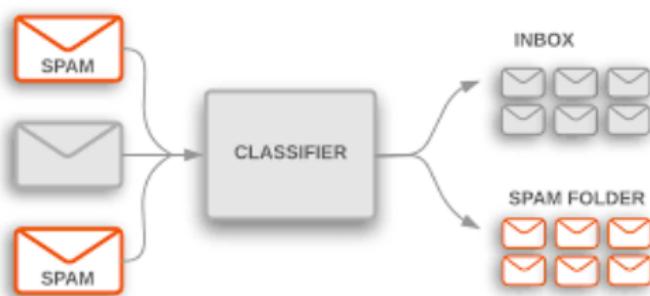
Experience: Watching label of emails as spam or not spam

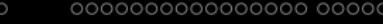
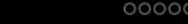
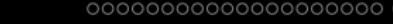
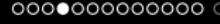
Performance metric: The number (or fraction) of emails correctly classified as spam



The essence of machine learning

- A pattern exist
 - We do not know it mathematically
 - We have data on it





Paradigms of ML

■ Supervised learning

Predicting a target variable for which we get to see examples

Paradigms of ML

- **Supervised learning**

Predicting a target variable for which we get to see examples

- **Unsupervised learning**

Revealing structure in the observed data

Paradigms of ML

- ### ■ Supervised learning

Predicting a target variable for which we get to see examples

- ## ■ Unsupervised learning

Revealing structure in the observed data

- ## ■ Reinforcement learning

Partial (indirect) feedback, no explicit guidance

Given rewards for a sequence of moves to learn a policy and utility functions

Supervised learning

Regression

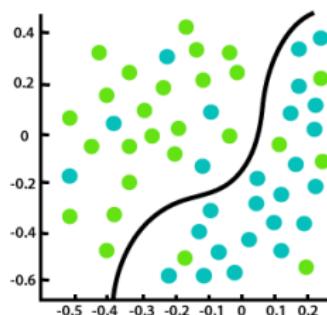
predict a continuous target variable

$$\text{E.g., } y \in [0, 1]$$

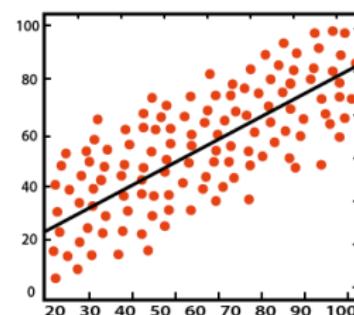
Classification

predict a discrete (unordered) target variable

$$\text{E.g., } y \in \{0, 1, 2, \dots, 9\}$$



Classification



Regression

Supervised learning: housing price prediction

Given: a dataset that contains N samples

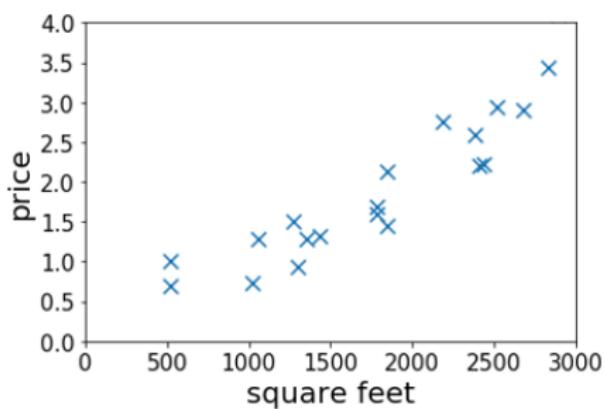
$$\{(x^{(1)}, y^{(1)}), \dots, (x^{(N)}, y^{(N)})\} \quad (1)$$

Supervised learning: housing price prediction

Given: a dataset that contains N samples

$$\{(x^{(1)}, y^{(1)}), \dots, (x^{(N)}, y^{(N)})\} \quad (1)$$

Task: if a residence has x square feet, predict its price?

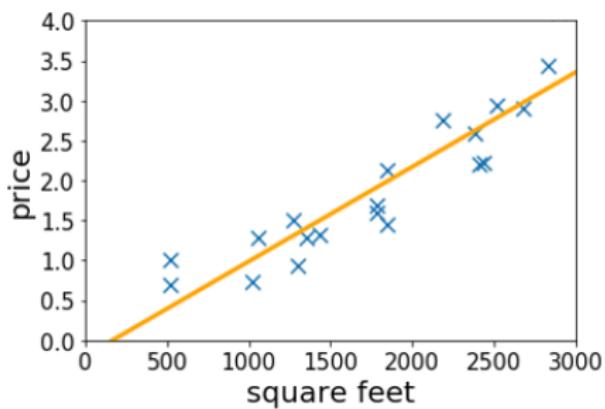


Supervised learning: housing price prediction

Given: a dataset that contains N samples

$$\{(x^{(1)}, y^{(1)}), \dots, (x^{(N)}, y^{(N)})\} \quad (1)$$

Task: if a residence has x square feet, predict its price?

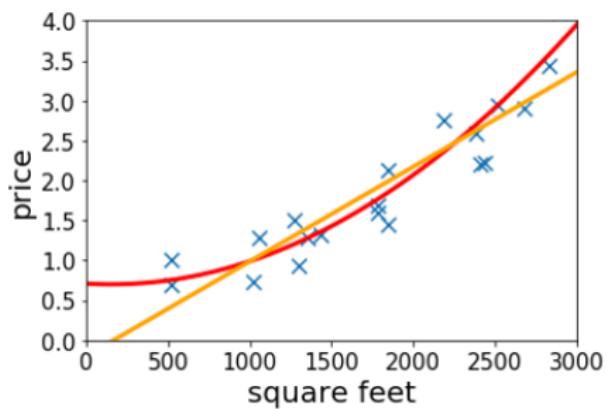


Supervised learning: housing price prediction

Given: a dataset that contains N samples

$$\{(x^{(1)}, y^{(1)}), \dots, (x^{(N)}, y^{(N)})\} \quad (1)$$

Task: if a residence has x square feet, predict its price?



High-dimensional features

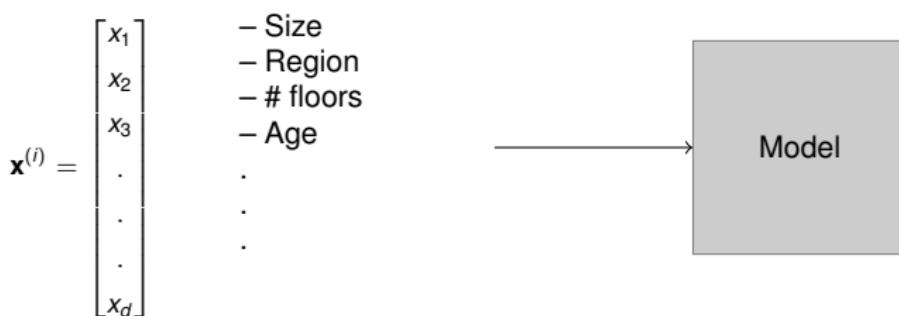
$$\mathbf{x}^{(i)} \in \mathbb{R}^d$$

$$\mathbf{x}^{(i)} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ \vdots \\ \vdots \\ x_d \end{bmatrix}$$

- Size
- Region
- # floors
- Age

High-dimensional features

$$\mathbf{x}^{(i)} \in \mathbb{R}^d$$



High-dimensional features

$$\mathbf{x}^{(i)} \in \mathbb{R}^d$$

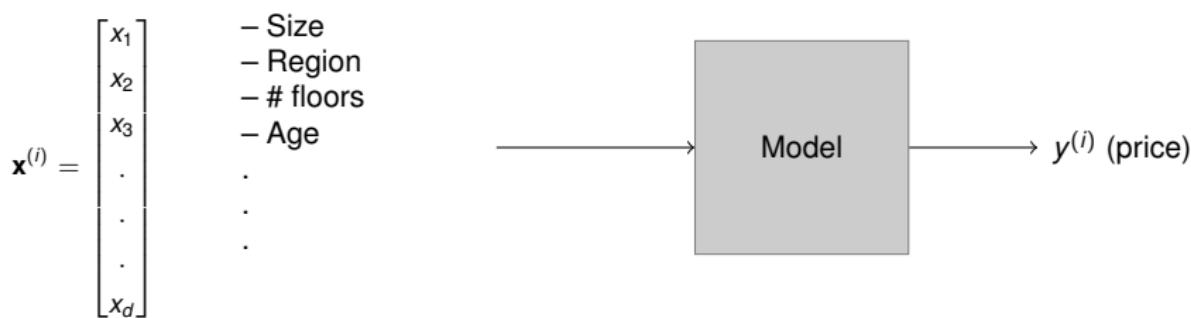
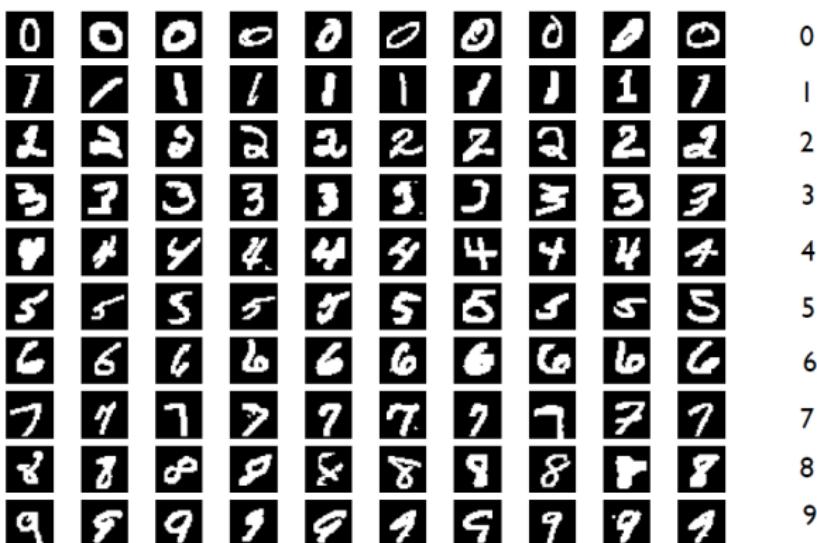
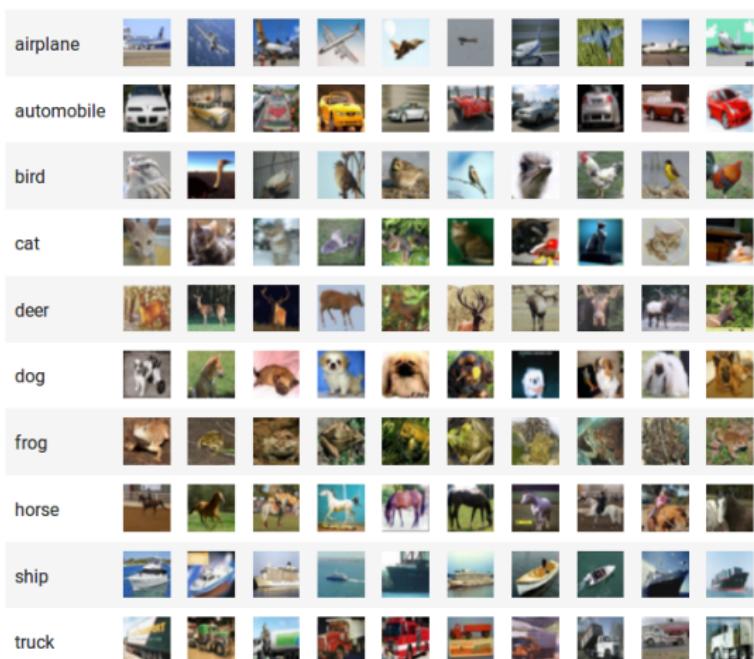


Image classification



Example images from the **MNIST** dataset.

Image classification



Example images from the **CIFAR-10** dataset.

Image classification



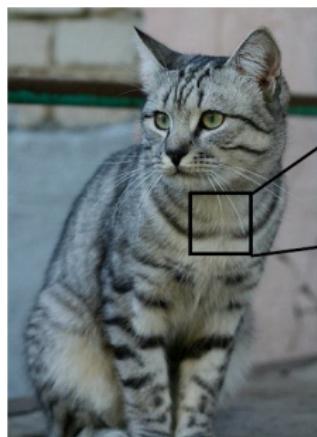
This image by [Nikita](#) is
licensed under [CC-BY 2.0](#).

(assume given a set of labels)
{dog, cat, truck, plane, ...}



cat

Image classification



This image by [Nikita](#) is
licensed under [CC-BY 2.0](#).

[1] 195 112 100	111 184	99 106 99	96 103 112	119 104	97 93 00]
[2] 76 85 98	185 128 185	87 95 95	95 99 115 123	106 103	99 00 00]
[3] 09 81 01	03 128 131	127 109	95 08 102	99 06	93 101 04]
[4] 106 91 01	04 69	91 88 05	181 187 189	98 75 84	96 95 05]
[5] 114 186	140 52	55 59 64	54 67 101	129 98	74 84 93
[6] 133 137	147 143	81 88 09	02 05 74	63 60 69	00 00 00]
[7] 129 137 144	140 109	95 86 79	70 62 65	63 69	73 86 00]
[8] 125 133 148	137 159	121 117	94 95 79	08 65	54 64 72 98]
[9] 127 125 131	147 133	127 126	131 111	98 89	75 61 64 72 84]
[10] 115 114 189	123 150	148 131	118 113	106 109	82 73 65 72 78]
[11] 09 93 99	97 100	147 131	118 113	114 113	109 95 77 00]
[12] 07 97 99	98 100	148 131	118 113	117 113	109 95 77 00]
[13] 62 65 82	89 78	71 88 101	124 126	119 181	107 114 131 119]
[14] 63 65 75	88 89	75 62	81 128	130 135	185 81 98 96 118]
[15] 07 65 71	07 106	95 69	45 76	130 126	187 92 94 185 112]
[16] 118 97 02	08 117	123 116	69 41	51 95	93 89 95 182 107]
[17] 164 146 112	80 82	128 121	104 76	48 45	46 88 181 182 107]
[18] 150 125 112	80 82	128 121	104 76	47 45	46 88 181 182 107]
[19] 129 128 134	161 159	108 149	118 123	134 114	87 65 53 69 06]
[20] 128 112 95	117 158	144 128	113 184	187 182	93 87 81 72 79]
[21] 123 187 96	95 86	83 112	153 149	122 169	184 75 80 187 112 99]
[22] 123 121 182	89 82	86 94	117 115	148 153	182 58 76 92 187]
[23] 122 164 149	103 75	86 70	03 93	183 119	139 102 51 69 04]

What the computer sees

An image is just a tensor of integers between [0, 255]:

e.g. 800 x 600 x 3
(3 channels RGB)

Unsupervised learning

Goal

Discover the intrinsic structure in the data

Data

Dataset contains no labels: $\{x^{(1)}, x^{(2)}, \dots, x^{(N)}\}$

Unsupervised learning

Goal

Discover the intrinsic structure in the data

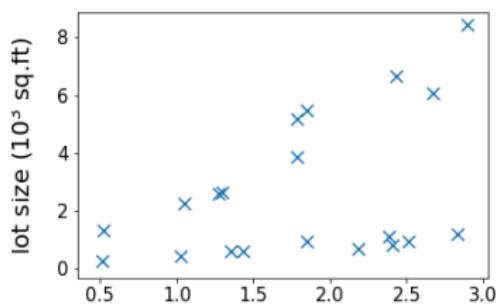
Data

Dataset contains no labels: $\{x^{(1)}, x^{(2)}, \dots, x^{(N)}\}$

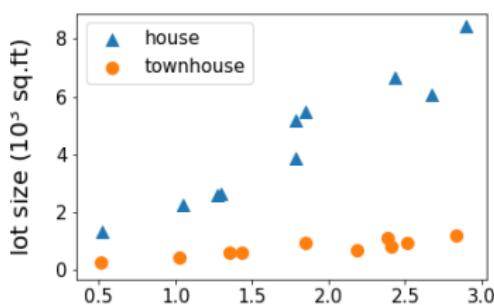
Clustering

Based on a similarity metric

Unsupervised



Supervised



Unsupervised learning

Goal

Discover the intrinsic structure in the data

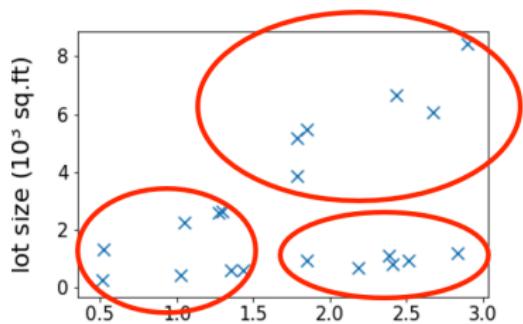
Data

Dataset contains no labels: $\{x^{(1)}, x^{(2)}, \dots, x^{(N)}\}$

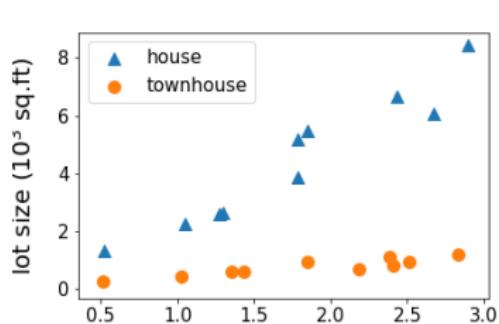
Clustering

Based on a similarity metric

Unsupervised



Supervised



Unsupervised learning

Goal

Discover the intrinsic structure in the data

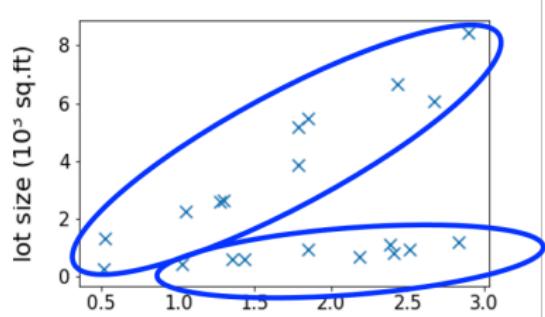
Data

Dataset contains no labels: $\{x^{(1)}, x^{(2)}, \dots, x^{(N)}\}$

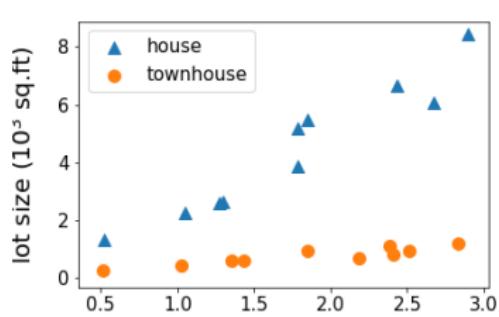
Clustering

Based on a similarity metric

Unsupervised



Supervised



Unsupervised learning

Clustering

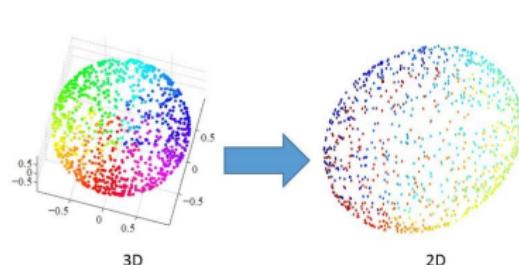
- Partitioning of data into groups of similar data points

Dimensionality reduction

- Data representation using a smaller number of dimensions while preserving (perhaps approximately) some properties of the data

Synthetic data generation

- e.g., Variational AutoEncoders (VAE) and Generative Adversarial Networks (GANs),

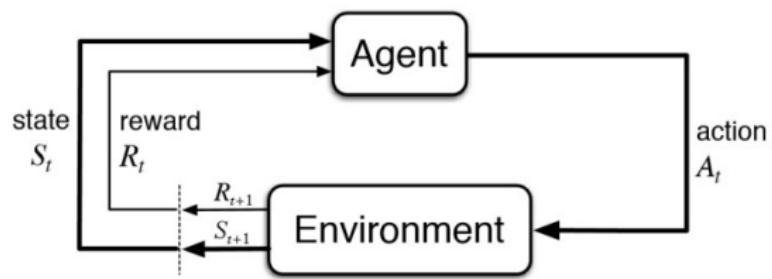


Reinforcement learning

Provides only an indication as to whether an action is correct or not

Data in supervised learning: (input, correct output)

Data in Reinforcement Learning: (input, some output, a reward for this output)



Machine Learning
ooooooooooooooo

Supervised Learning
●ooooooo

Neural Networks
oooooooooooooooooooo

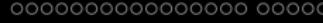
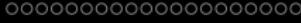
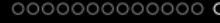
Logistic Regression
ooooo

Loss Function
oooo

Optimization
ooooooooooooooo

Regularization
oooooo

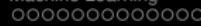
Supervised Learning



Components of supervised learning (function approximation)

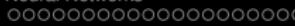
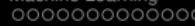
Concept

- Unknown target function $f : \mathcal{X} \rightarrow \mathcal{Y}$
- \mathcal{X} : Input space
- \mathcal{Y} : Output space



Components of supervised learning (function approximation)

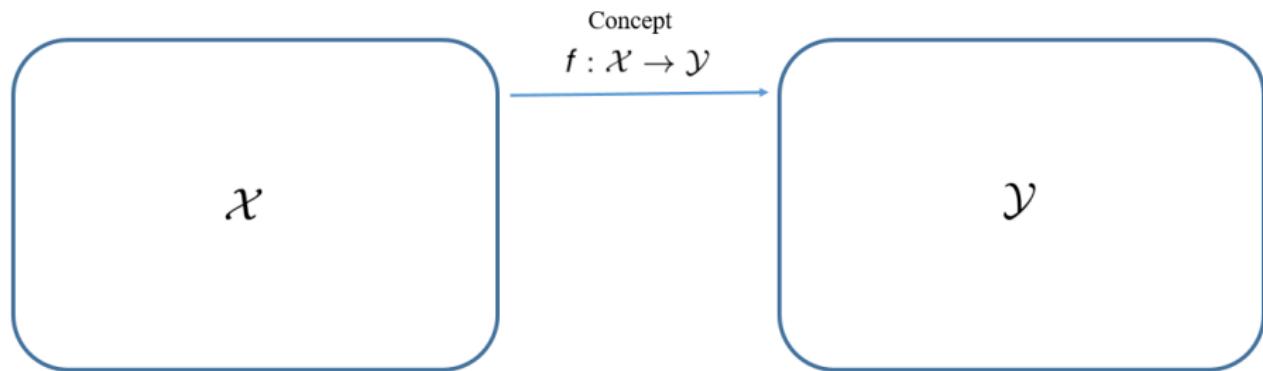
 \mathcal{X}



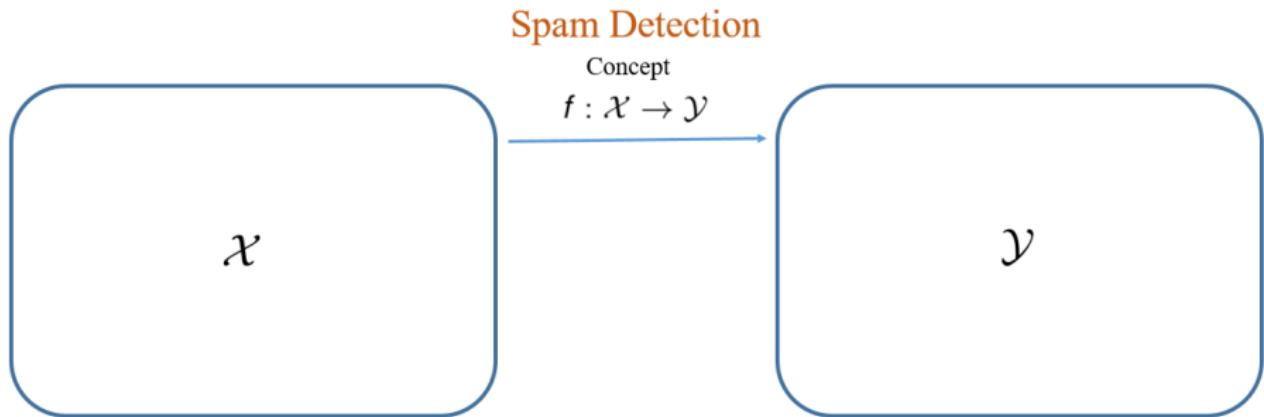
Components of supervised learning (function approximation)

 \mathcal{X} \mathcal{Y}

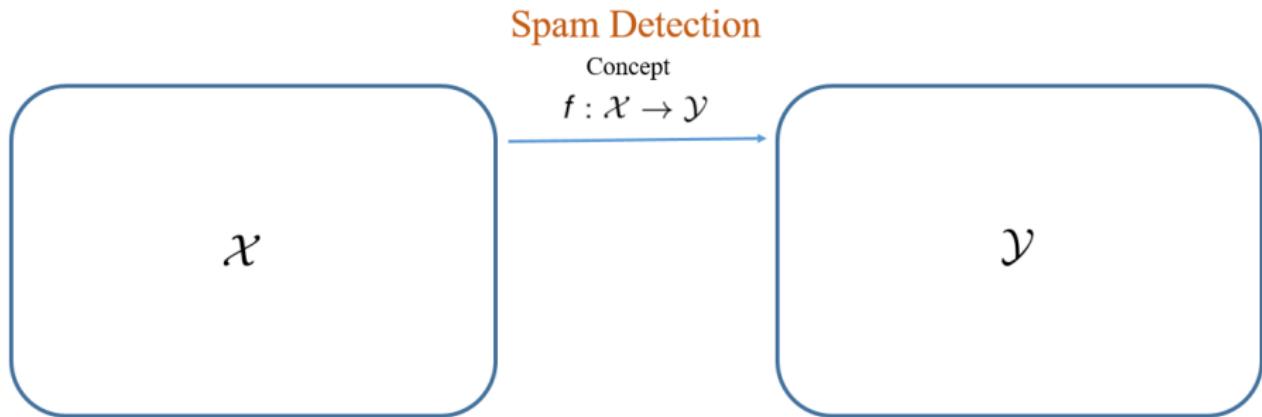
Components of supervised learning (function approximation)



Components of supervised learning (function approximation)

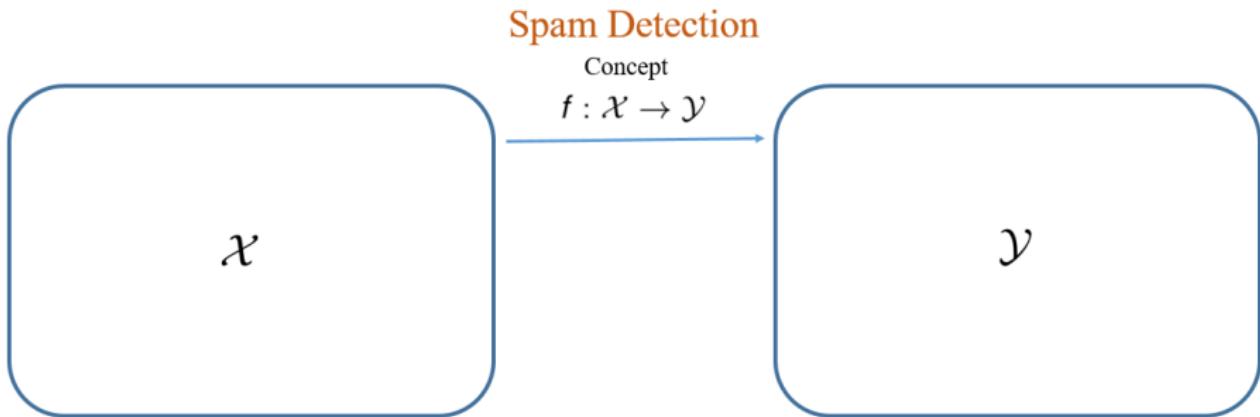


Components of supervised learning (function approximation)



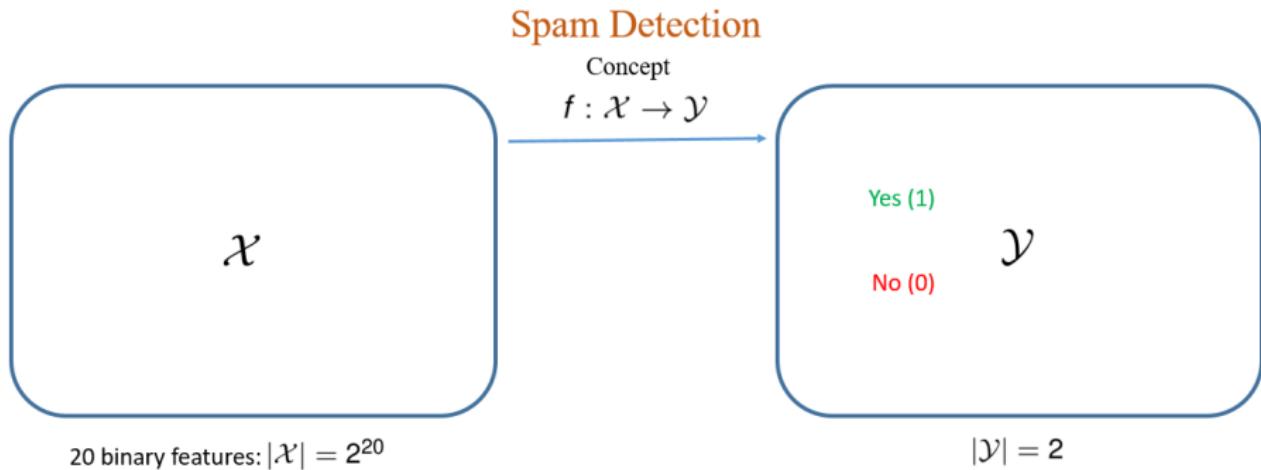
20 binary features:

Components of supervised learning (function approximation)

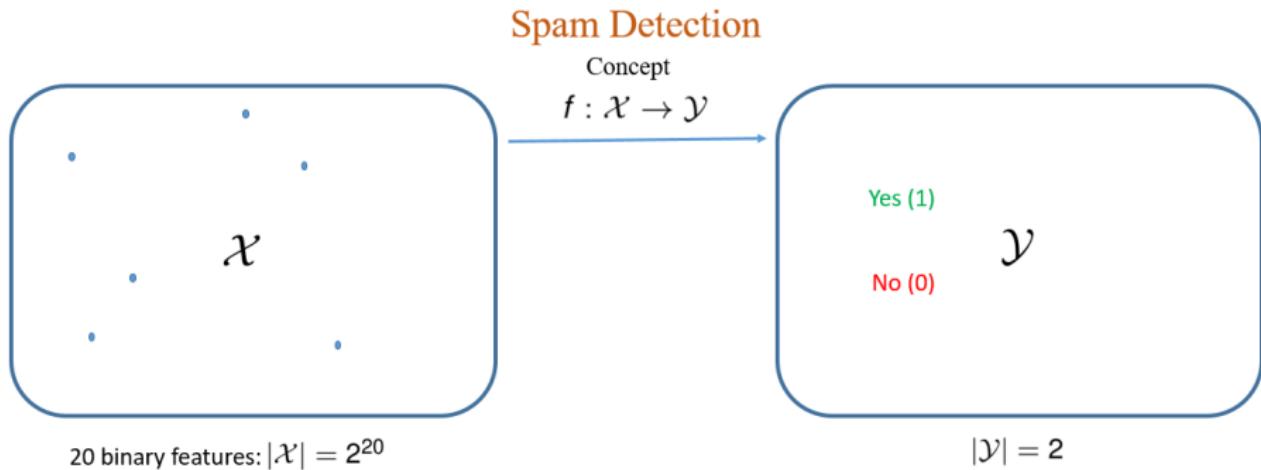


20 binary features: $|\mathcal{X}| = 2^{20}$

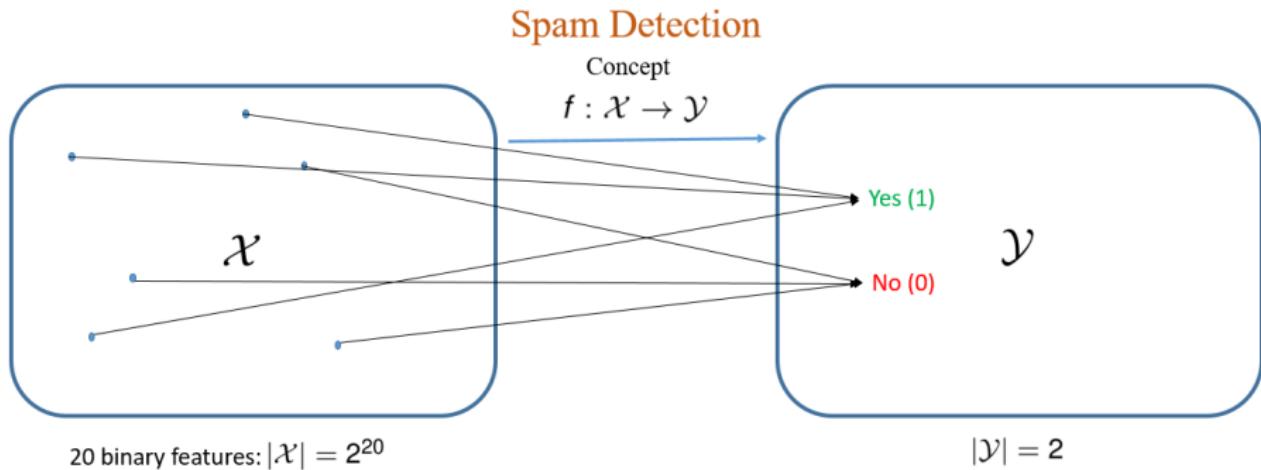
Components of supervised learning (function approximation)



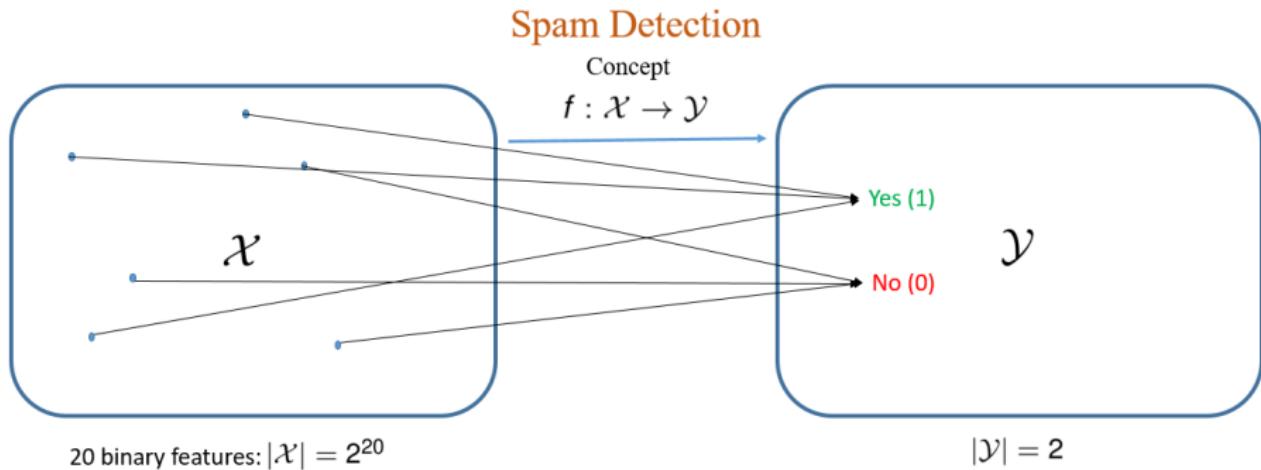
Components of supervised learning (function approximation)



Components of supervised learning (function approximation)

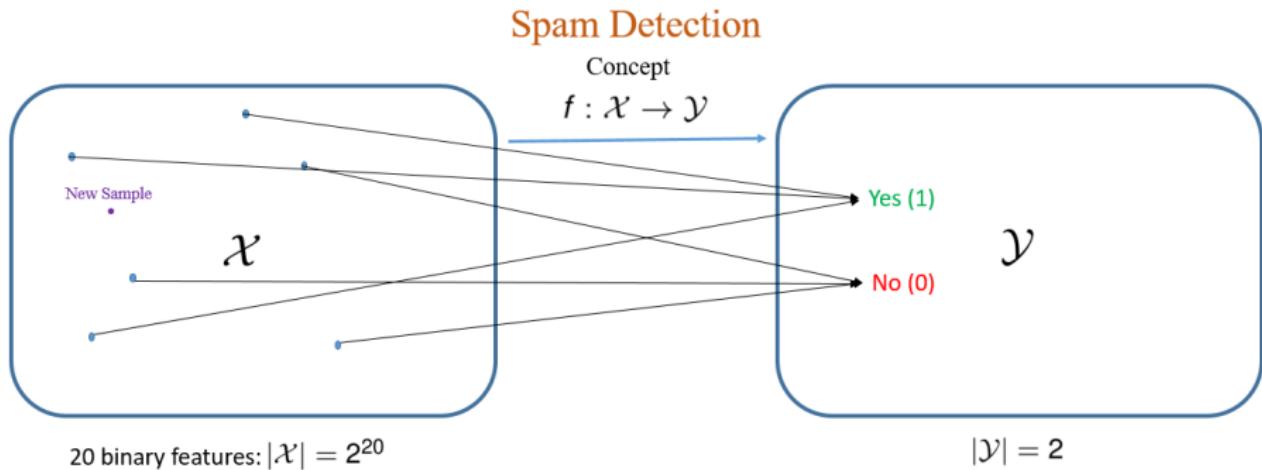


Components of supervised learning (function approximation)



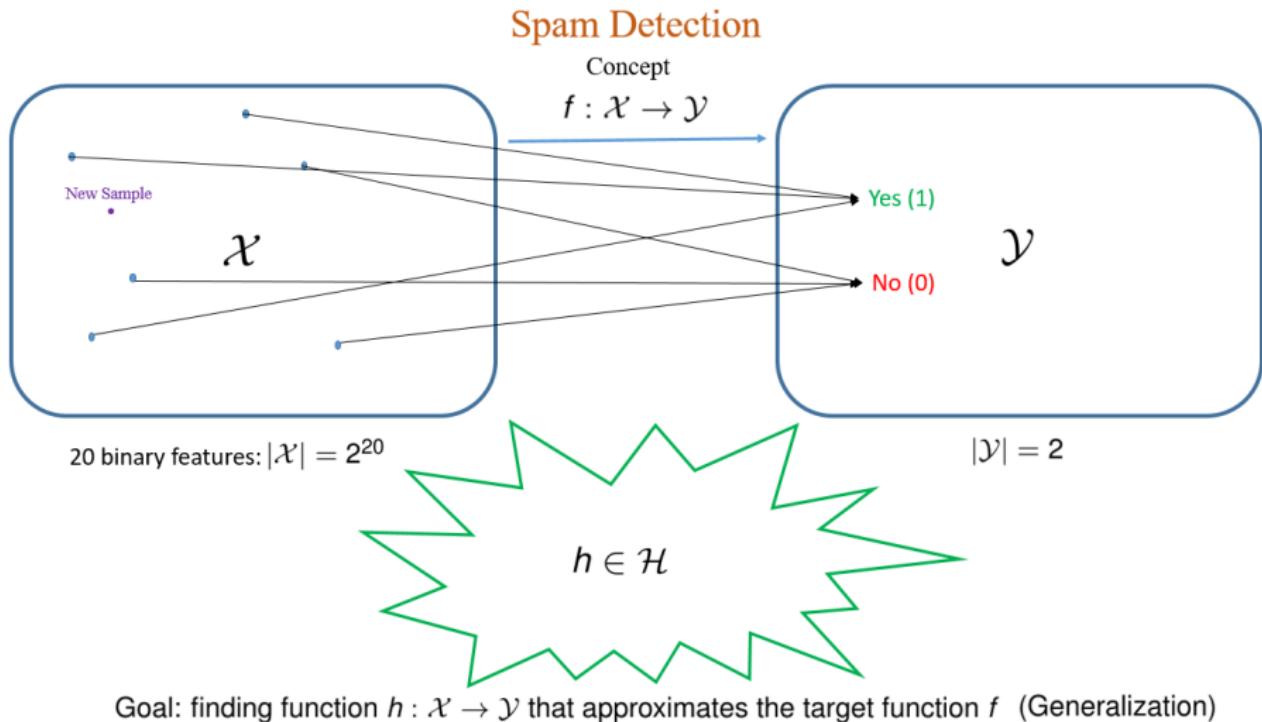
Goal: finding function $h : \mathcal{X} \rightarrow \mathcal{Y}$ that approximates the target function f

Components of supervised learning (function approximation)

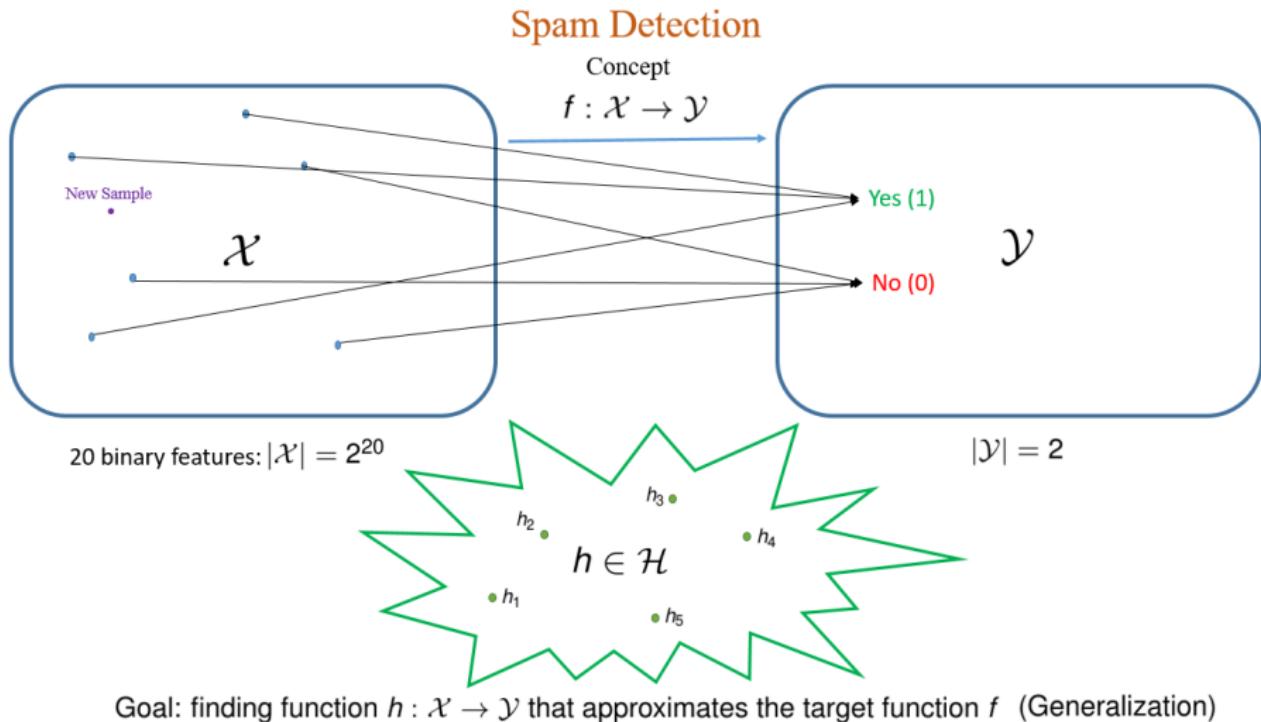


Goal: finding function $h : \mathcal{X} \rightarrow \mathcal{Y}$ that approximates the target function f (Generalization)

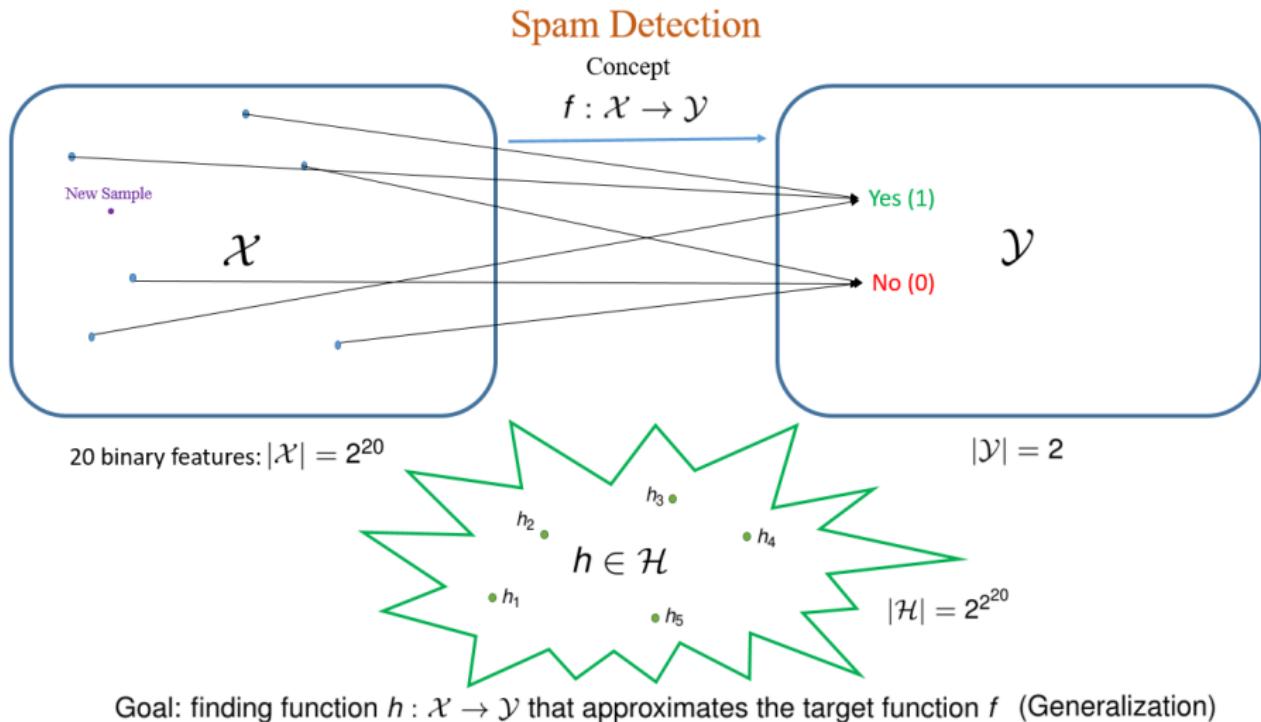
Components of supervised learning (function approximation)



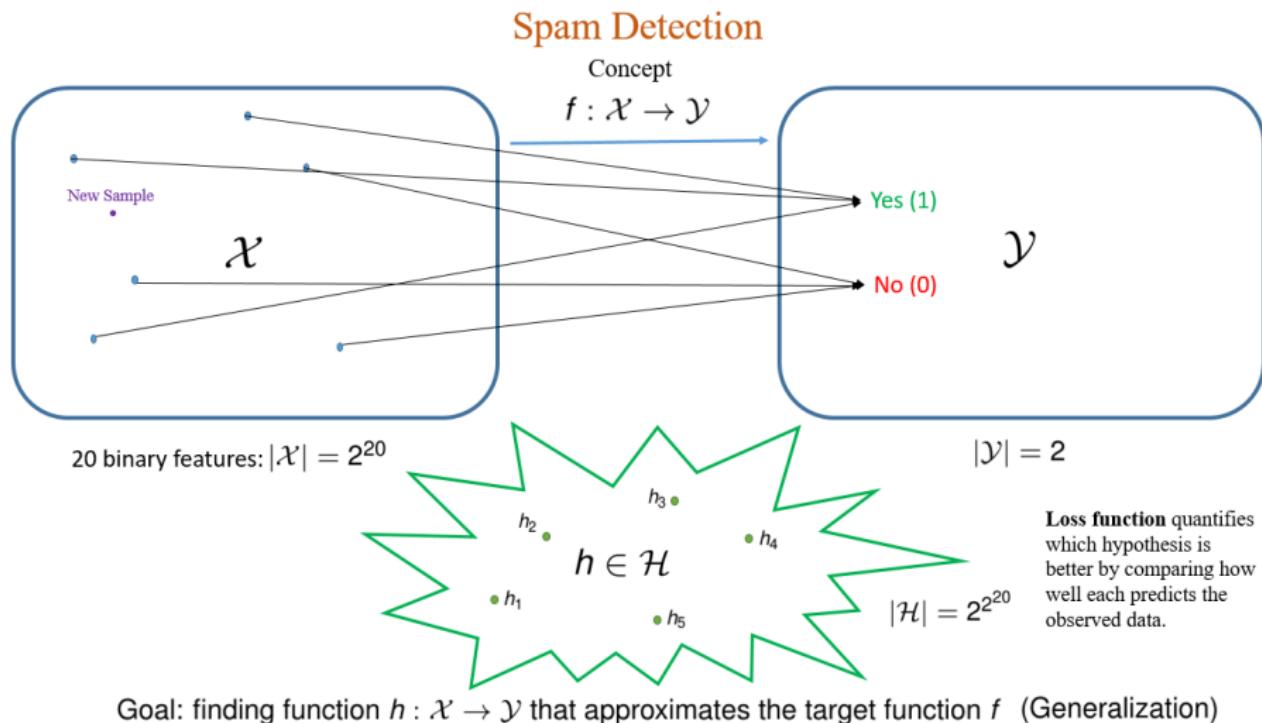
Components of supervised learning (function approximation)



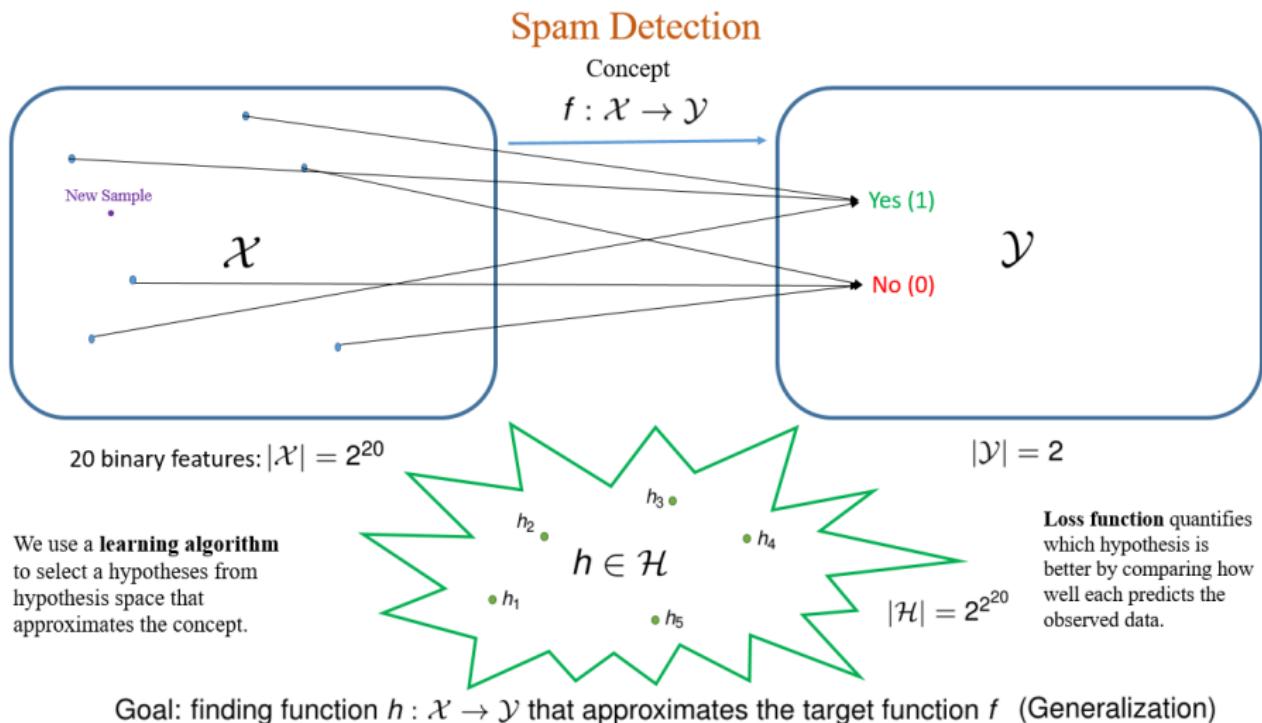
Components of supervised learning (function approximation)

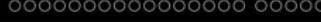
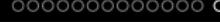


Components of supervised learning (function approximation)



Components of supervised learning (function approximation)





Components of supervised learning (function approximation)

Concept

- Unknown target function $f : \mathcal{X} \rightarrow \mathcal{Y}$
- \mathcal{X} : Input space
- \mathcal{Y} : Output space

Training data

- Samples from the concept $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(N)}, y^{(N)})\}$

Components of supervised learning (function approximation)

Concept

- Unknown target function $f : \mathcal{X} \rightarrow \mathcal{Y}$
- \mathcal{X} : Input space
- \mathcal{Y} : Output space

Training data

- Samples from the concept $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(N)}, y^{(N)})\}$

Hypotheses space

- Pick a hypotheses $h : \mathcal{X} \rightarrow \mathcal{Y}$ from hypotheses space $h \in \mathcal{H}$ (e.g., linear functions) that approximates the target function f

Components of supervised learning (function approximation)

Concept

- Unknown target function $f : \mathcal{X} \rightarrow \mathcal{Y}$
- \mathcal{X} : Input space
- \mathcal{Y} : Output space

Training data

- Samples from the concept $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(N)}, y^{(N)})\}$

Hypotheses space

- Pick a hypotheses $h : \mathcal{X} \rightarrow \mathcal{Y}$ from hypotheses space $h \in \mathcal{H}$ (e.g., linear functions) that approximates the target function f

Learning algorithm

- We use a learning algorithm to select a hypotheses from hypothesis space that approximates the concept

Components of supervised learning

Learning model composed of

- A hypothesis set
- A learning algorithm

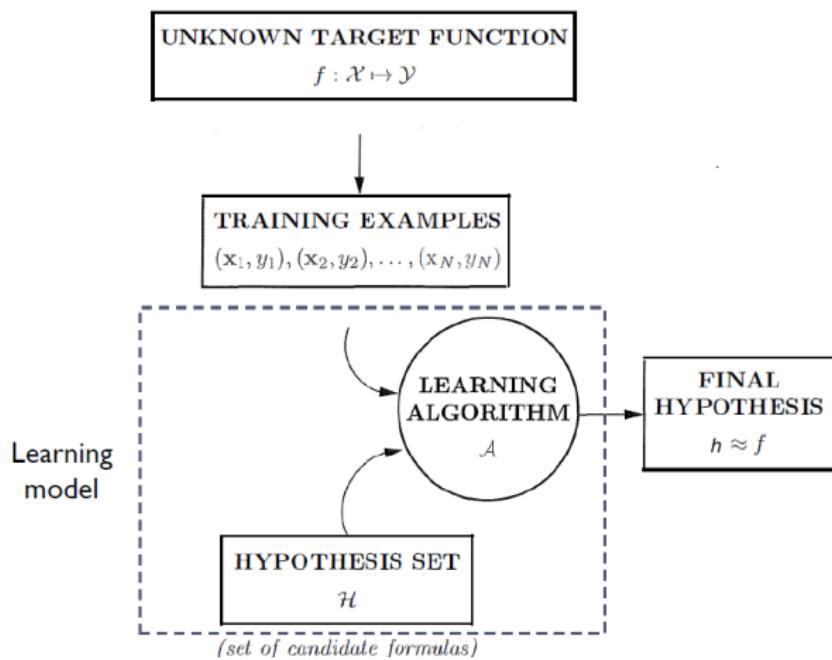


Figure: <https://work.caltech.edu/telecourse.html>

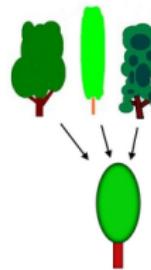
Generalization

Hypotheses evaluation

How well h generalizes to unseen examples.

Test data

Samples from the concept $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(N)}, y^{(N)})\}$
Do not exist in the training data



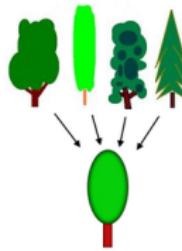
Generalization

Hypotheses evaluation

How well h generalizes to unseen examples.

Test data

Samples from the concept $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(N)}, y^{(N)})\}$
Do not exist in the training data

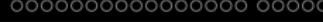
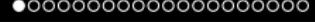
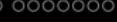
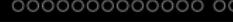


Independent and identically distributed (i.i.d.)

A collection of random variables is independent and identically distributed if each random variable has the same probability distribution as the others and all are mutually independent.

- $\forall i, x^{(i)} \sim \mathcal{D}$ (Identically Distributed)
- $\forall i \neq j, \mathcal{P}(x^{(i)}, x^{(j)}) = \mathcal{P}(x^{(i)})\mathcal{P}(x^{(j)})$ (Independently Distributed)

Machine learning algorithms have focused primarily on sets of data points that were assumed to be independent and identically distributed (i.i.d.).

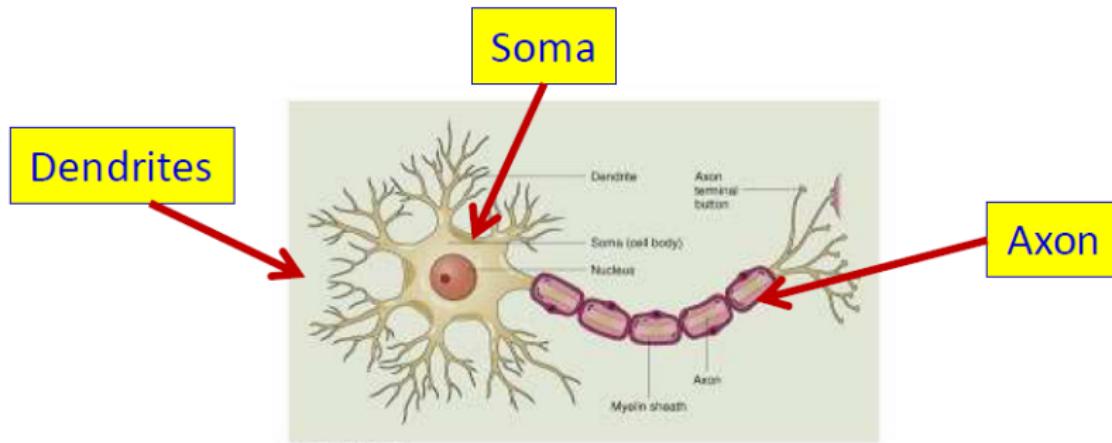


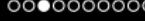
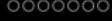
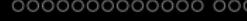
Neural Networks

Modeling the brain

A neuron

- Signals come in through the dendrites into the Soma
- A signal goes out via the axon to other neurons
 - Only one axon per neuron

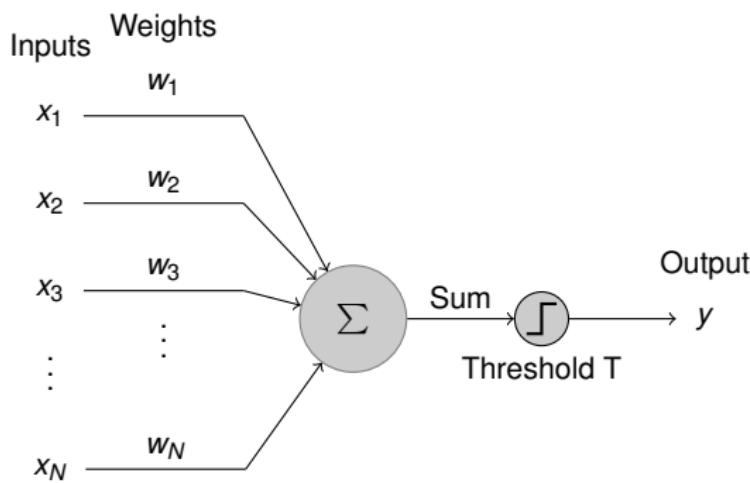




Perceptron

Number of inputs combine linearly

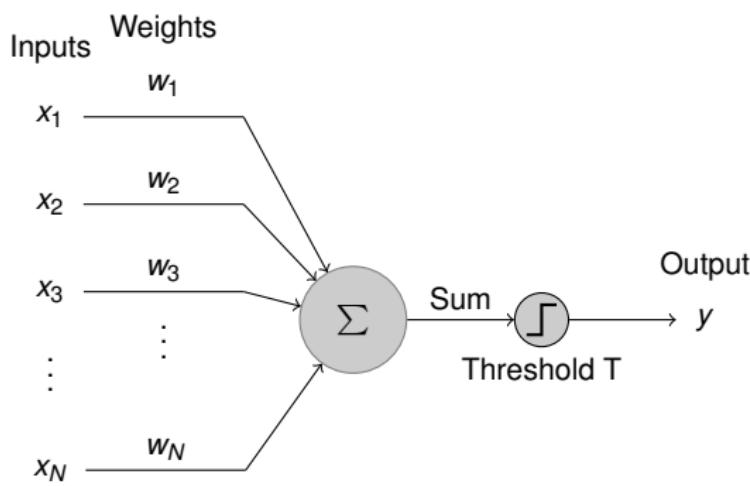
- Threshold logic: Fire if combined input exceeds threshold



Perceptron

Number of inputs combine linearly

- Threshold logic: Fire if combined input exceeds threshold

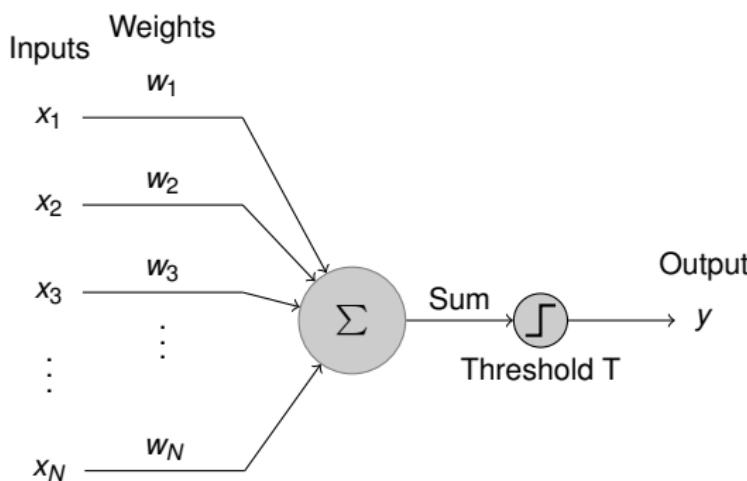


$$y = \begin{cases} 1 & \text{if } \sum_i w_i x_i - T \geq 0 \\ 0 & \text{else} \end{cases}$$

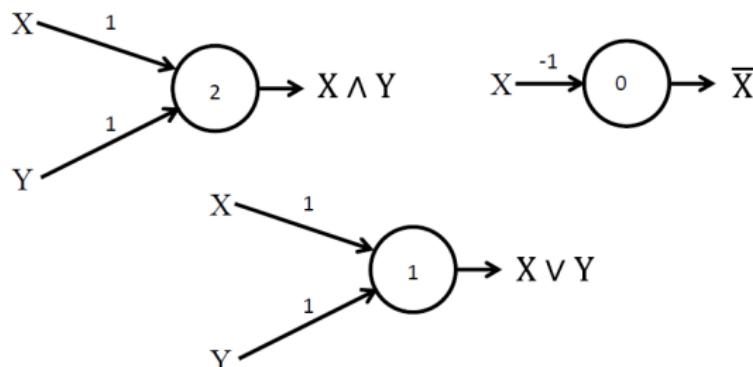
The universal model

Originally assumed could represent any Boolean circuit and perform any logic

- The embryo of an electronic computer that [the Navy] expects will be able to walk, talk, see, write, reproduce itself and be conscious of its existence, New York Times (8 July) 1958
 - Frankenstein Monster Designed by Navy That Thinks, Tulsa, Oklahoma Times 1958



Perceptron

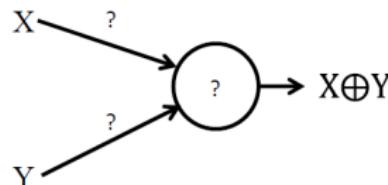


Values shown on edges are weights, numbers in the circles are thresholds ($X, Y \in \{0, 1\}$)

- Easily shown to mimic any Boolean gate
 - But ...

Perceptron

No solution for XOR!
Not universal!

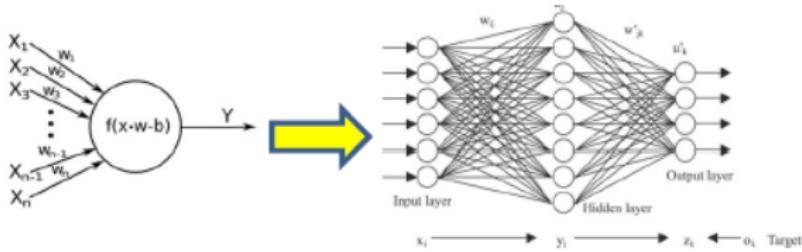
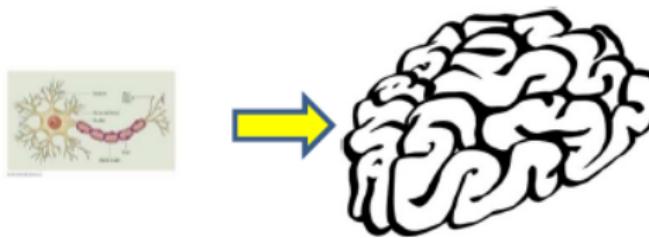


Values shown on edges are weights, numbers in the circles are thresholds ($X, Y \in \{0, 1\}$)

- Easily shown to mimic any Boolean gate
- But ...

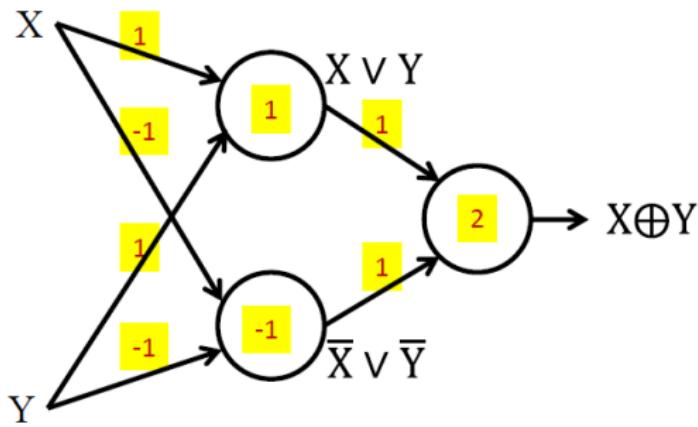
A single neuron is not enough

- Individual elements are weak computational elements
 - Marvin Minsky and Seymour Papert, 1969, Perceptrons: An Introduction to Computational Geometry
- Networked elements are required



Multi-layer perceptron!

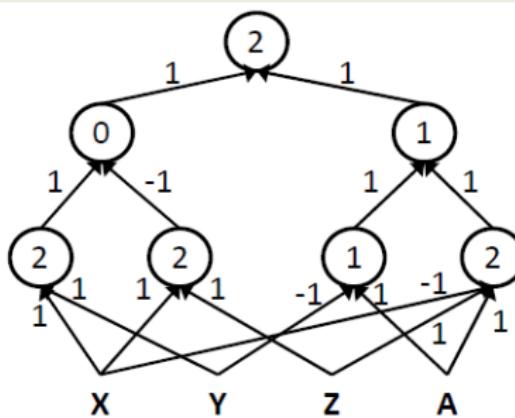
XOR



Multi-layer perceptrons (MLPs)

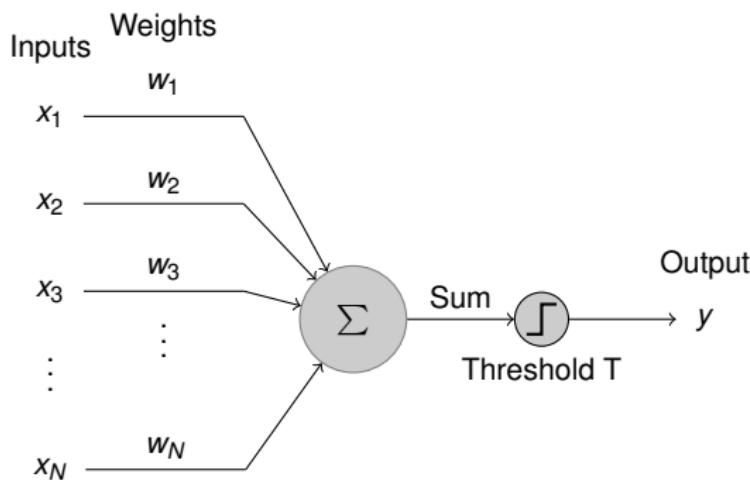
- MLPs can compute *any* Boolean function
 - Any function over any number of inputs and any number of outputs

$$((A \& \bar{X} \& Z) | (A \& \bar{Y})) \& ((X \& Y) | \overline{(X \& Z)})$$



The perceptron with real inputs

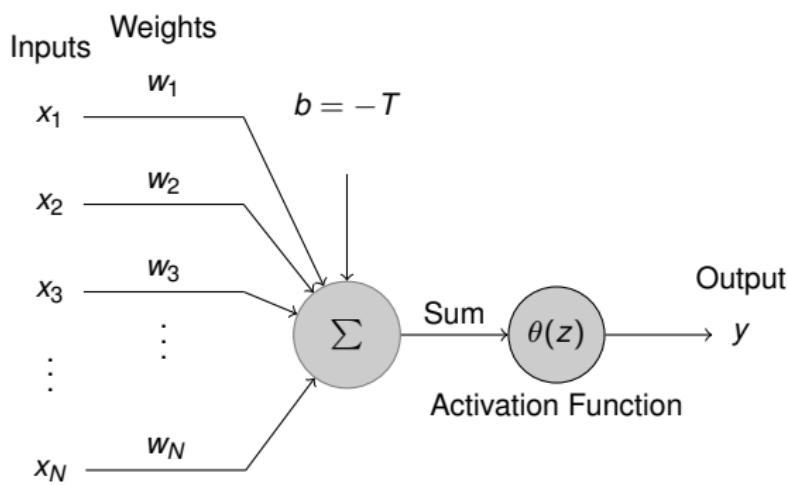
- x_1, x_2, \dots, x_n are real valued
- w_1, w_2, \dots, w_n are real valued
- Unit fires if weighted input exceeds a threshold



$$y = \begin{cases} 1 & \text{if } \sum_i w_i x_i - T \geq 0 \\ 0 & \text{else} \end{cases}$$

The perceptron with real inputs

- An alternate view:
 - A threshold activation $\theta(z)$ operates on the weighted sum of inputs plus a bias
 - $\theta(z)$ outputs 1 if z is non-negative, 0 otherwise
- Unit fires if weighted input exceeds a threshold



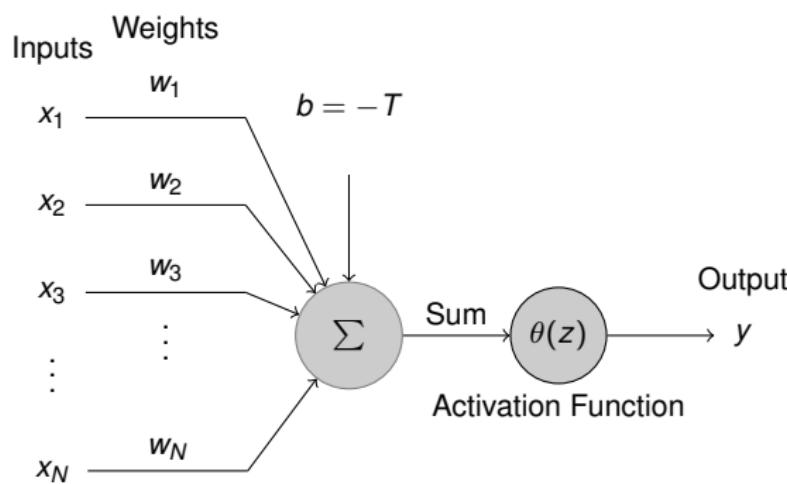
$$z = \sum_i w_i x_i + b$$

$$y = \theta(z) = \begin{cases} 1 & \text{if } z \geq 0 \\ 0 & \text{else} \end{cases}$$

Perceptron

The bias can also be viewed as the weight of another input component that is always set to 1

- If the bias is not explicitly mentioned, we will implicitly be assuming that every perceptron has an additional input that is always fixed at 1



$$z = \sum_i w_i x_i + b$$

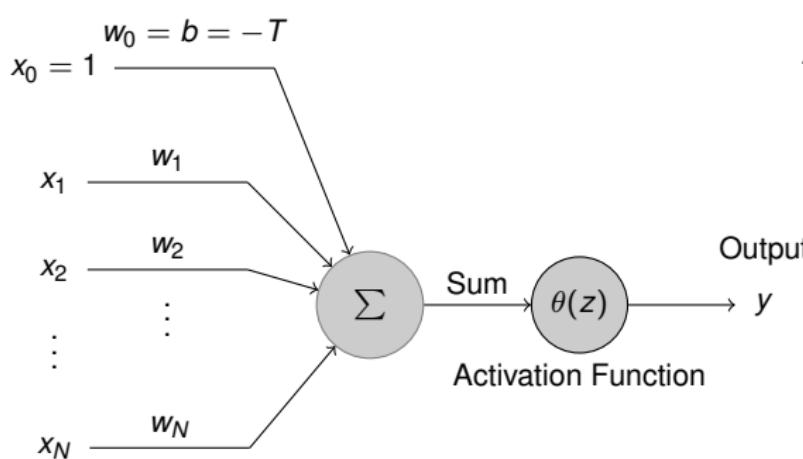
$$y = \begin{cases} 1 & \text{if } z \geq 0 \\ 0 & \text{else} \end{cases}$$

Perceptron

The bias can also be viewed as the weight of another input component that is always set to 1

- If the bias is not explicitly mentioned, we will implicitly be assuming that every perceptron has an additional input that is always fixed at 1

Inputs Weights

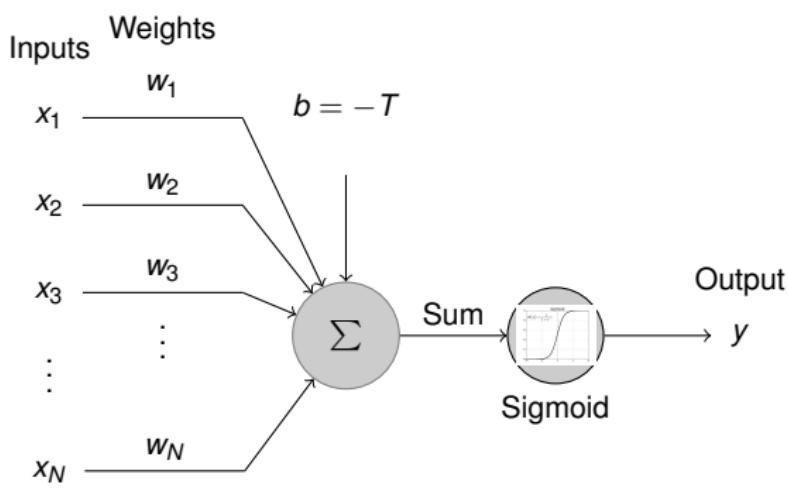


$$z = \sum_{i=1}^N w_i x_i + b \text{ or } z = \sum_{i=0}^N w_i x_i$$

$$y = \begin{cases} 1 & \text{if } z \geq 0 \\ 0 & \text{else} \end{cases}$$

A real output

- An alternate view:
 - A threshold activation $\theta(z)$ operates on the weighted sum of inputs plus a bias
 - $\theta(z)$ outputs 1 if z is non-negative, 0 otherwise
- The output y can also be real valued
 - Sometimes viewed as the probability of firing
- Logistic Regression



$$y = \sigma\left(\sum_i w_i x_i + b\right)$$

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

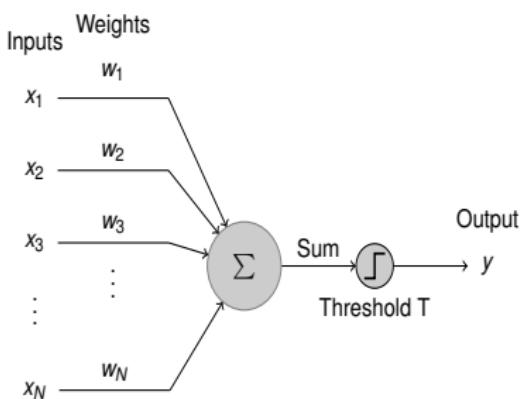
Vectorization

$$z = \sum_i w_i x_i = \mathbf{w}^T \mathbf{x} = \begin{bmatrix} w_0 = b & w_1 & w_2 & \dots & w_{N-1} & w_N \end{bmatrix} \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_{N-1} \\ x_N \end{bmatrix}$$

$$\theta(z) = \begin{cases} 1 & \text{if } z \geq 0 \\ 0 & \text{else} \end{cases}$$

Perceptron

A perceptron is a linear classifier

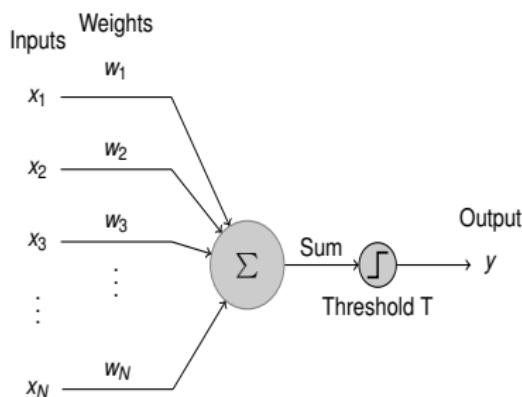


$$y = \theta(\sum_i w_i x_i - T)$$

$$\theta(z) = \begin{cases} 1 & \text{if } z \geq 0 \\ 0 & \text{else} \end{cases}$$

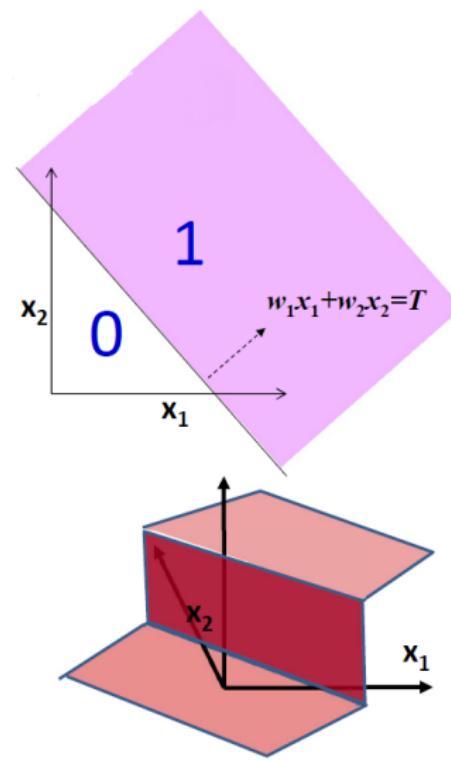
Perceptron

A perceptron is a linear classifier



$$y = \theta\left(\sum_i w_i x_i - T\right)$$

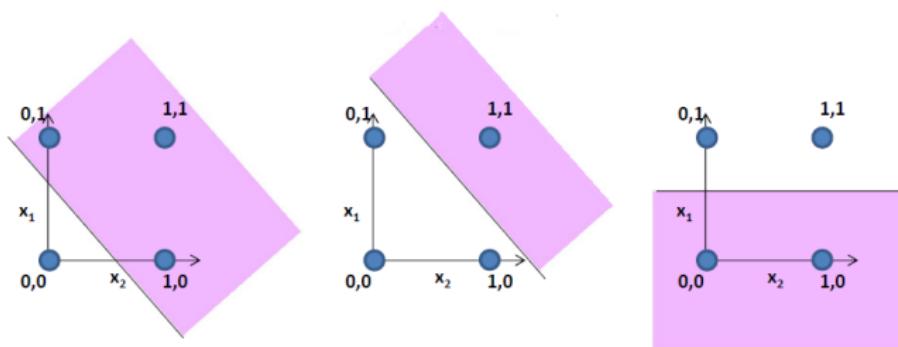
$$\theta(z) = \begin{cases} 1 & \text{if } z \geq 0 \\ 0 & \text{else} \end{cases}$$



Boolean functions with a real perceptron

Boolean perceptrons are also linear classifiers

- Purple regions have output 1 in the figures
- What are these functions?
- Why can we not compose an XOR?

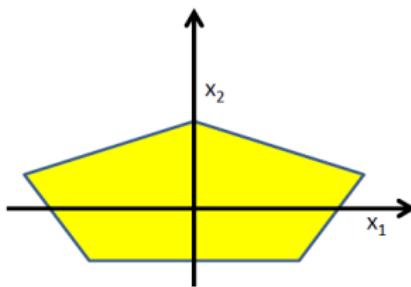


A convex region

Build a network of units with a single output that fires if the input is in the coloured area

- Neurons can now be composed into networks to compute arbitrary classification boundaries

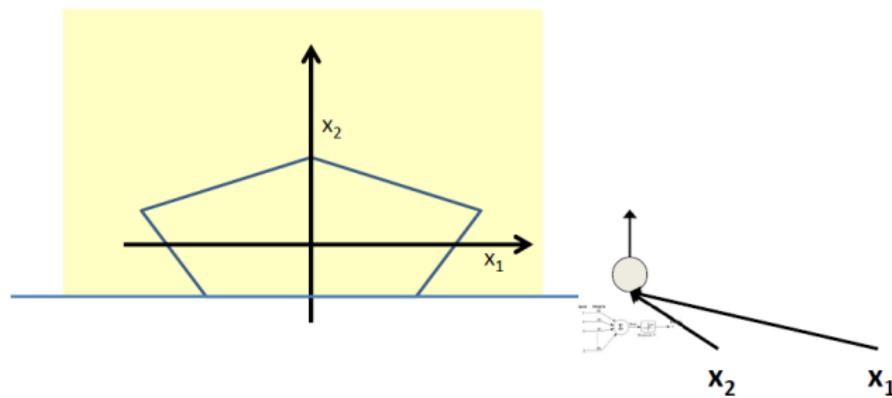
boundaries



A convex region

Build a network of units with a single output that fires if the input is in the coloured area

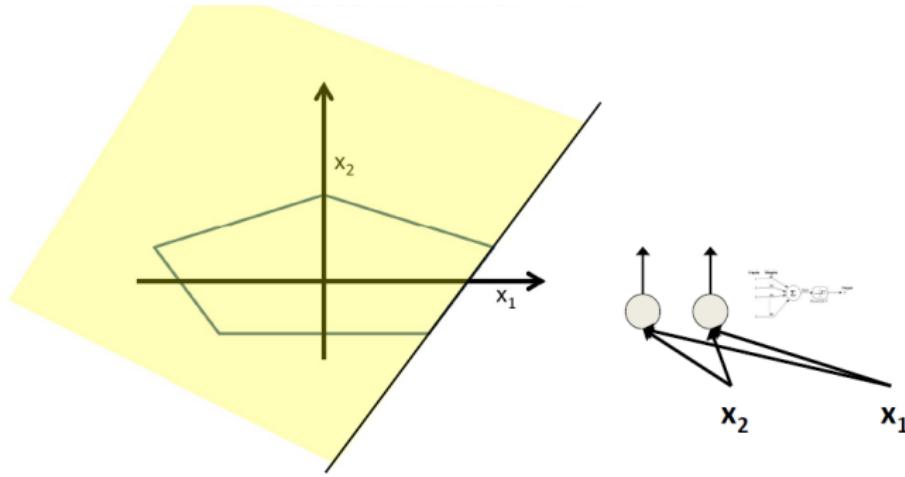
- Neurons can now be composed into networks to compute arbitrary classification boundaries



A convex region

Build a network of units with a single output that fires if the input is in the coloured area

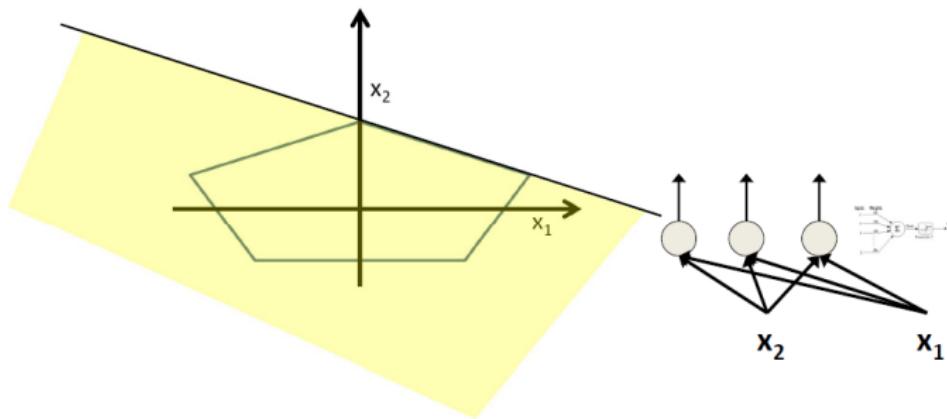
- Neurons can now be composed into networks to compute arbitrary classification boundaries



A convex region

Build a network of units with a single output that fires if the input is in the coloured area

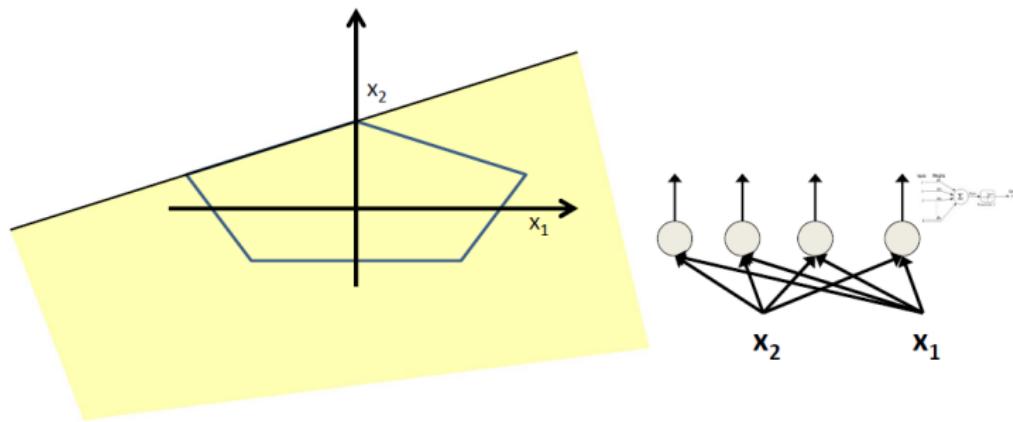
- Neurons can now be composed into networks to compute arbitrary classification boundaries



A convex region

Build a network of units with a single output that fires if the input is in the coloured area

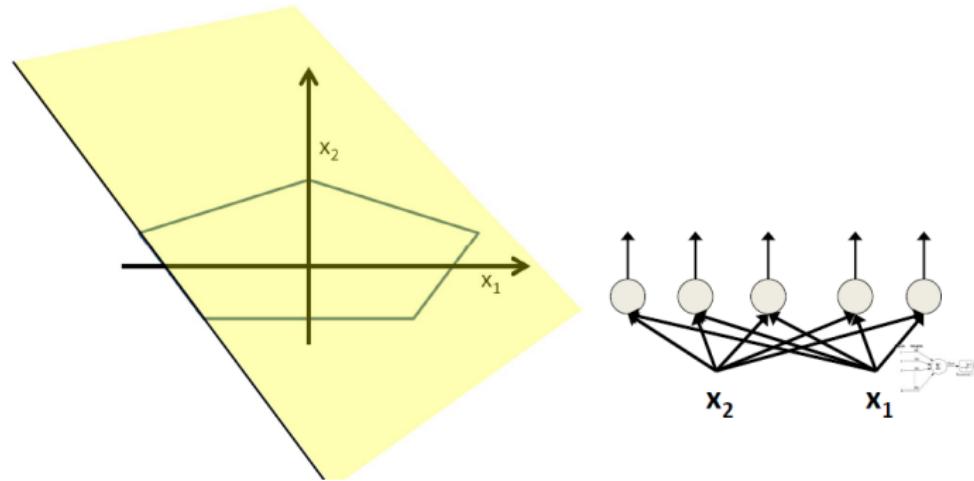
- Neurons can now be composed into networks to compute arbitrary classification boundaries



A convex region

Build a network of units with a single output that fires if the input is in the coloured area

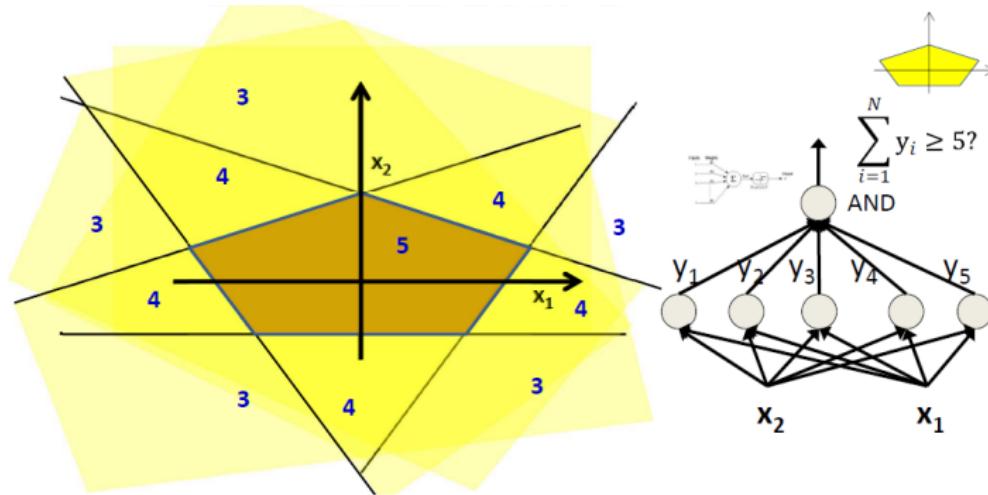
- Neurons can now be composed into networks to compute arbitrary classification boundaries



A convex region

Build a network of units with a single output that fires if the input is in the coloured area

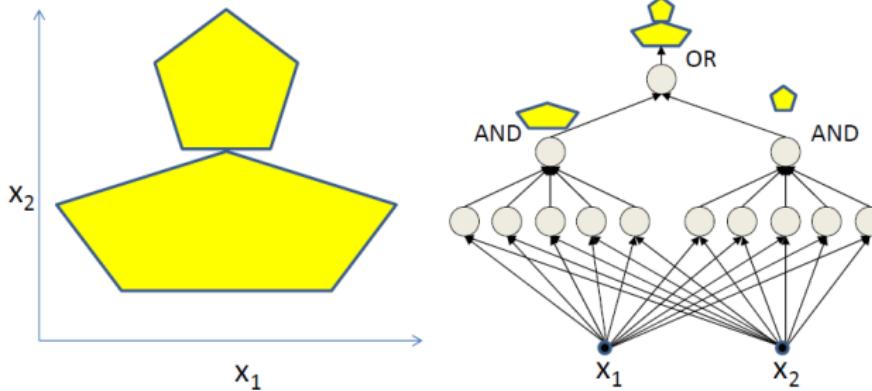
- Neurons can now be composed into networks to compute arbitrary classification boundaries

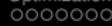
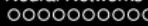
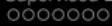
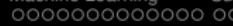


A convex region

Build a network of units with a single output that fires if the input is in the coloured area

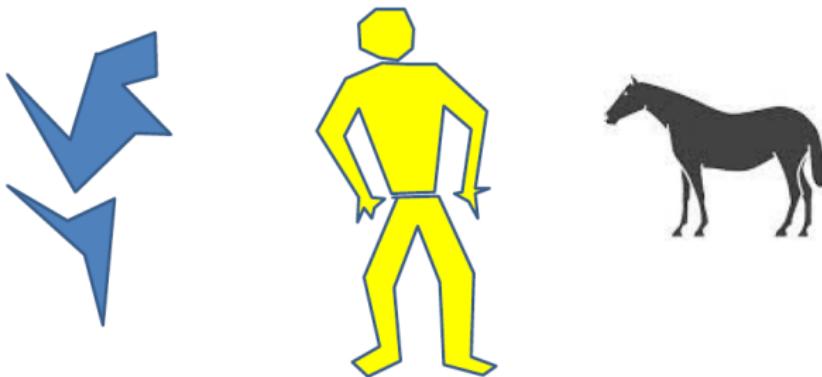
- Neurons can now be composed into networks to compute arbitrary classification boundaries





Complex decision boundaries

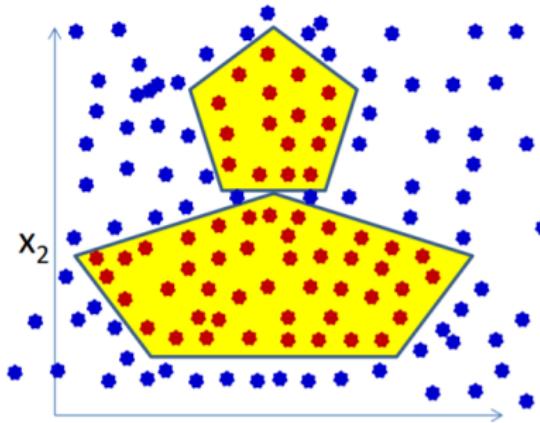
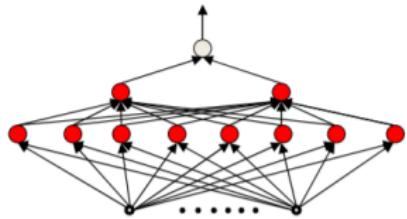
Can compose very complex decision boundaries



A visual proof that neural nets can compute any function

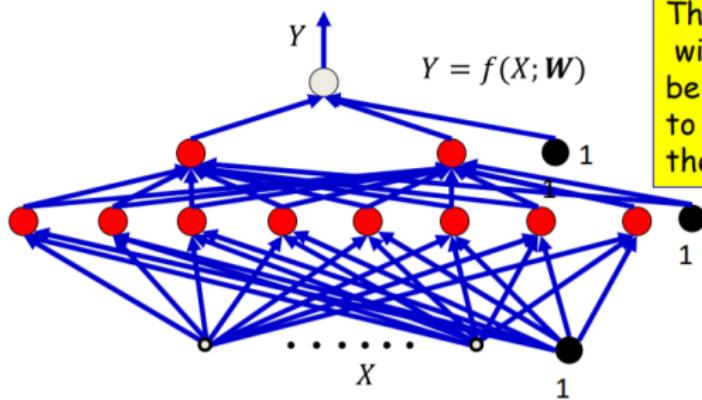
A real problem

- Learn an MLP for this function
 - 1 in the yellow regions, 0 outside
- Using just the samples
- We know this can be perfectly represented using an MLP



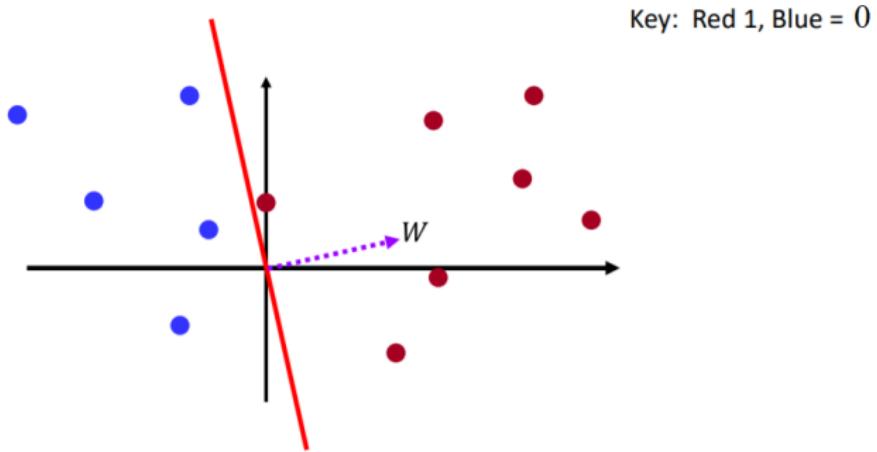
What we learn: The parameters of the network

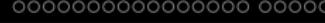
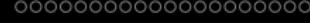
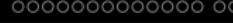
- Given: the architecture of the network
- The parameters of the network: The weights and biases
 - The weights associated with the blue arrows in the picture
- Learning the network: Determining the values of these parameters such that the network computes the desired function



The Perceptron Problem

- Find the hyperplane $w^T x = 0$ that perfectly separates the two groups of points
 - $w^T x = 0$ is the hyperplane comprising all x 's orthogonal to vector w
 - Learning the perceptron = finding the weight vector w for the separating hyperplane
 - w points in the direction of the positive class





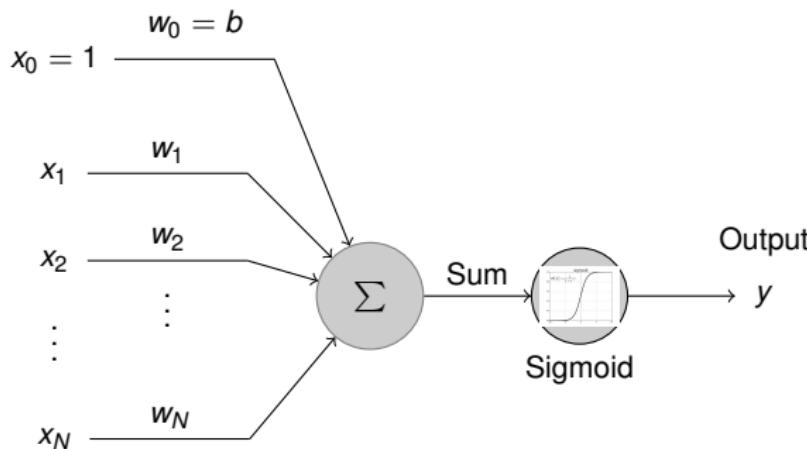
Logistic Regression

Logistic regression

It is the perceptron with a sigmoid activation

- It actually computes the probability that the input belongs to class 1

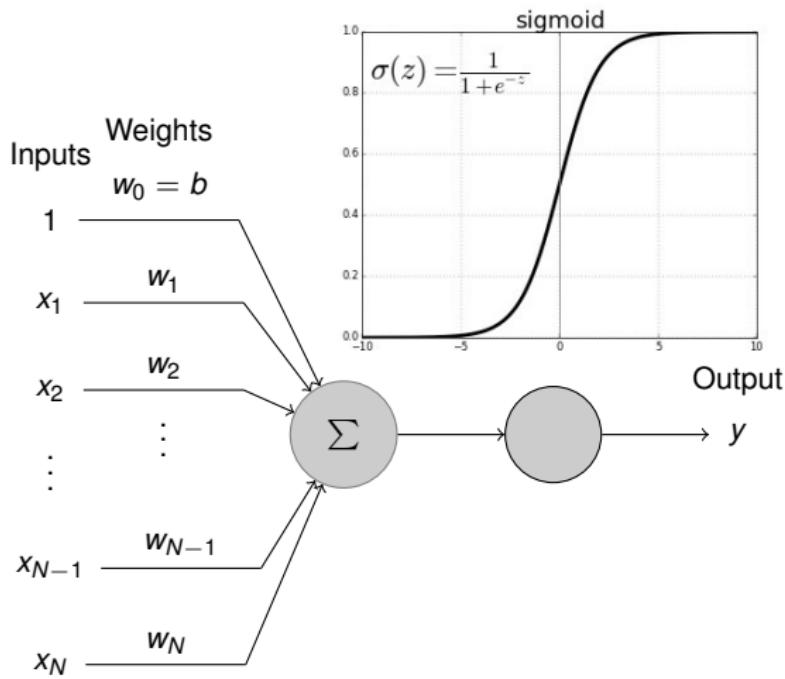
Inputs Weights



Logistic regression

It is the perceptron with a sigmoid activation

- It actually computes the probability that the input belongs to class 1



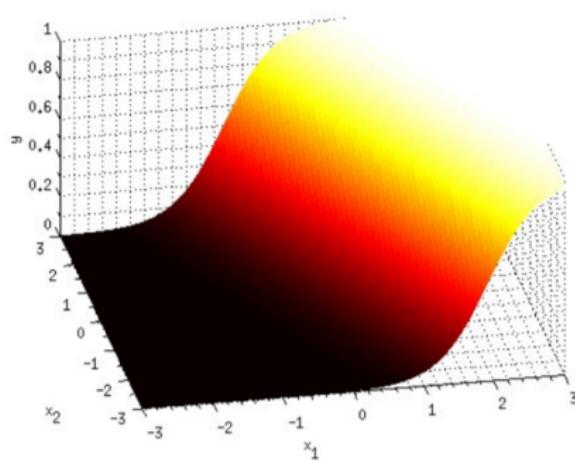
$$p(y = 1|x) = \sigma(z) = \frac{1}{1 + \exp(-z)}$$

$$y = \begin{cases} 1 & \text{if } \sigma(z) \geq \frac{1}{2} \\ 0 & \text{if } \sigma(z) < \frac{1}{2} \end{cases}$$

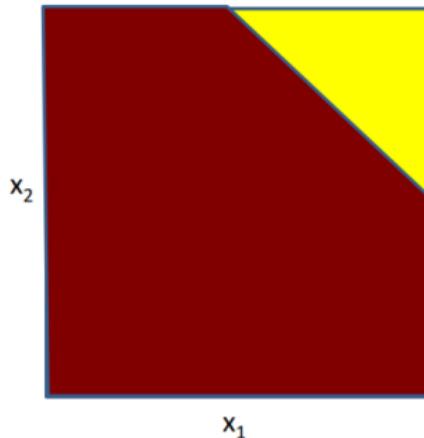
Logistic regression

This is the perceptron with a sigmoid activation

- It actually computes the probability that the input belongs to class 1



Decision: $y > 0.5?$



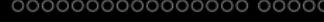
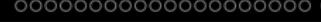
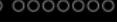
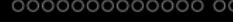
When X is a 2-D variable

$$P(Y = 1|X) = \frac{1}{1 + \exp(-\sum_i w_i x_i - b)}$$

Vectorization

$$z = \sum_i w_i x_i = \mathbf{w}^T \mathbf{x} = \begin{bmatrix} w_0 = b & w_1 & w_2 & \dots & w_{N-1} & w_N \end{bmatrix} \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_{N-1} \\ x_N \end{bmatrix}$$

$$f(\mathbf{x}, \mathbf{w}) = \sigma(\mathbf{w}^T \mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^T \mathbf{x})}$$



Loss Function



Loss function

- Loss function quantifies our unhappiness with the scores across the training data.

Loss function

- Loss function quantifies our unhappiness with the scores across the training data.
- A loss function tells how good our current classifier is
 - Given a dataset of examples $\{(x^{(i)}, y^{(i)})\}_{i=1}^N$
 - $x^{(i)}$ is data and $y^{(i)}$ is binary label.

Loss function

- Loss function quantifies our unhappiness with the scores across the training data.
- A loss function tells how good our current classifier is
 - Given a dataset of examples $\{(x^{(i)}, y^{(i)})\}_{i=1}^N$
 - $x^{(i)}$ is data and $y^{(i)}$ is binary label.
- Loss over the dataset is a average of loss over examples:

$$L(f(X, \mathbf{w}), \mathbf{y}) = \frac{1}{N} \sum_i L_i(f(\mathbf{x}^{(i)}, \mathbf{w}), y^{(i)})$$

Loss function

- Loss function quantifies our unhappiness with the scores across the training data.
- A loss function tells how good our current classifier is
 - Given a dataset of examples $\{(x^{(i)}, y^{(i)})\}_{i=1}^N$
 - $x^{(i)}$ is data and $y^{(i)}$ is binary label.
- Loss over the dataset is a average of loss over examples:

$$L(f(X, \mathbf{w}), \mathbf{y}) = \frac{1}{N} \sum_i L_i(f(\mathbf{x}^{(i)}, \mathbf{w}), y^{(i)})$$

- Empirical loss L is only an empirical approximation of the true loss

Loss function

- Loss function quantifies our unhappiness with the scores across the training data.
- A loss function tells how good our current classifier is
 - Given a dataset of examples $\{(x^{(i)}, y^{(i)})\}_{i=1}^N$
 - $x^{(i)}$ is data and $y^{(i)}$ is binary label.
- Loss over the dataset is a average of loss over examples:

$$L(f(X, \mathbf{w}), \mathbf{y}) = \frac{1}{N} \sum_i L_i(f(\mathbf{x}^{(i)}, \mathbf{w}), y^{(i)})$$

- Empirical loss L is only an empirical approximation of the true loss
- **Learning**
 - Estimate the parameters to minimize the empirical loss L

$$\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w}} L(f(X, \mathbf{w}), \mathbf{y})$$

Cross entropy

If we have two separate probability distributions $P(x)$ and $Q(x)$ over the same random variable x , we can measure how different these two distributions are using the Kullback-Leibler (KL) divergence:

$$D_{KL}(P||Q) = \mathbb{E}_{x \sim P} [\log \frac{P(x)}{Q(x)}] = \mathbb{E}_{x \sim P} [\log P(x) - \log Q(x)]$$

$$D_{KL}(P||Q) = \sum_i p_i(x) \log p_i(x) - \sum_i p_i(x) \log q_i(x)$$

Cross entropy

If we have two separate probability distributions $P(x)$ and $Q(x)$ over the same random variable x , we can measure how different these two distributions are using the Kullback-Leibler (KL) divergence:

$$D_{KL}(P||Q) = \mathbb{E}_{x \sim P} [\log \frac{P(x)}{Q(x)}] = \mathbb{E}_{x \sim P} [\log P(x) - \log Q(x)]$$

$$D_{KL}(P||Q) = \sum_i p_i(x) \log p_i(x) - \sum_i p_i(x) \log q_i(x)$$

- It is the extra amount of information needed to send a message containing symbols drawn from probability distribution P , when we use a code that was designed to minimize the length of messages drawn from probability distribution Q .

Cross entropy

If we have two separate probability distributions $P(x)$ and $Q(x)$ over the same random variable x , we can measure how different these two distributions are using the Kullback-Leibler (KL) divergence:

$$D_{KL}(P||Q) = \mathbb{E}_{x \sim P} [\log \frac{P(x)}{Q(x)}] = \mathbb{E}_{x \sim P} [\log P(x) - \log Q(x)]$$

$$D_{KL}(P||Q) = \sum_i p_i(x) \log p_i(x) - \sum_i p_i(x) \log q_i(x)$$

- It is the extra amount of information needed to send a message containing symbols drawn from probability distribution P , when we use a code that was designed to minimize the length of messages drawn from probability distribution Q .
- KL divergence is non-negative and measures the difference between two distributions, it is often conceptualized as measuring some sort of distance between these distributions. However, it is not a true distance measure because it is not symmetric

Cross entropy

- Minimizing loss function

$$\operatorname{argmin}_{Q(x)} D_{KL}(P||Q) = \sum_i p_i(x) \log p_i(x) - \sum_i p_i(x) \log q_i(x)$$

Cross entropy

- Minimizing loss function

$$\operatorname{argmin}_{Q(x)} D_{KL}(P||Q) = \sum_i p_i(x) \log p_i(x) - \sum_i p_i(x) \log q_i(x)$$

- $P(x)$ is the true distribution of the random variable x and it is fixed.

Cross entropy

- Minimizing loss function

$$\operatorname{argmin}_{Q(x)} D_{KL}(P||Q) = \sum_i p_i(x) \log p_i(x) - \sum_i p_i(x) \log q_i(x)$$

- $P(x)$ is the true distribution of the random variable x and it is fixed.

$$\operatorname{argmin}_{Q(x)} - \sum_i p_i(x) \log q_i(x)$$

Cross entropy

- Minimizing loss function

$$\operatorname{argmin}_{Q(x)} D_{KL}(P||Q) = \sum_i p_i(x) \log p_i(x) - \sum_i p_i(x) \log q_i(x)$$

- $P(x)$ is the true distribution of the random variable x and it is fixed.

$$\operatorname{argmin}_{Q(x)} - \sum_i p_i(x) \log q_i(x)$$

$$CrossEntropy(P(x), Q(x)) = - \sum_i p_i(x) \log q_i(x)$$

Cross entropy

- Minimizing loss function

$$\operatorname{argmin}_{Q(x)} D_{KL}(P||Q) = \sum_i p_i(x) \log p_i(x) - \sum_i p_i(x) \log q_i(x)$$

- $P(x)$ is the true distribution of the random variable x and it is fixed.

$$\operatorname{argmin}_{Q(x)} - \sum_i p_i(x) \log q_i(x)$$

$$CrossEntropy(P(x), Q(x)) = - \sum_i p_i(x) \log q_i(x)$$

Example

The label of data x is 0. The output of Logistic Regression for x is 0.7

Cross entropy

- Minimizing loss function

$$\operatorname{argmin}_{Q(x)} D_{KL}(P||Q) = \sum_i p_i(x) \log p_i(x) - \sum_i p_i(x) \log q_i(x)$$

- $P(x)$ is the true distribution of the random variable x and it is fixed.

$$\operatorname{argmin}_{Q(x)} - \sum_i p_i(x) \log q_i(x)$$

$$\text{CrossEntropy}(P(x), Q(x)) = - \sum_i p_i(x) \log q_i(x)$$

Example

The label of data x is 0. The output of Logistic Regression for x is 0.7

$$\text{True Dist. } P(x) = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \quad -- \quad \text{Predicted Dist. } Q(x) = \begin{bmatrix} 1 - 0.7 \\ 0.7 \end{bmatrix}$$

Cross entropy

- Minimizing loss function

$$\operatorname{argmin}_{Q(x)} D_{KL}(P||Q) = \sum_i p_i(x) \log p_i(x) - \sum_i p_i(x) \log q_i(x)$$

- $P(x)$ is the true distribution of the random variable x and it is fixed.

$$\operatorname{argmin}_{Q(x)} - \sum_i p_i(x) \log q_i(x)$$

$$\text{CrossEntropy}(P(x), Q(x)) = - \sum_i p_i(x) \log q_i(x)$$

Example

The label of data x is 0. The output of Logistic Regression for x is 0.7

$$\text{True Dist. } P(x) = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \quad -- \quad \text{Predicted Dist. } Q(x) = \begin{bmatrix} 1 - 0.7 \\ 0.7 \end{bmatrix}$$

$$\text{CrossEntropy}(P(x), Q(x)) = -1 \times \log(1 - 0.7) - 0 \times \log 0.7 = 1.2$$

Logistic regression loss function

■ Logistic Regression

$$P(y = 1 | \mathbf{x}^{(i)}; \mathbf{w}) = f(\mathbf{x}^{(i)}, \mathbf{w})$$

$$P(y = 0 | \mathbf{x}^{(i)}; \mathbf{w}) = 1 - f(\mathbf{x}^{(i)}, \mathbf{w})$$

$$y^{(i)} \in \{0, 1\}$$

Logistic regression loss function

■ Logistic Regression

$$P(y=1|\mathbf{x}^{(i)}; \mathbf{w}) = f(\mathbf{x}^{(i)}, \mathbf{w})$$

$$P(y=0|\mathbf{x}^{(i)}; \mathbf{w}) = 1 - f(\mathbf{x}^{(i)}, \mathbf{w})$$

$$y^{(i)} \in \{0, 1\}$$

■ Cross entropy as the loss function

$$L(f(\mathbf{X}, \mathbf{w}), \mathbf{y}) = \frac{1}{N} \sum_i L_i(f(\mathbf{x}^{(i)}, \mathbf{w}), y^{(i)})$$

$$L(f(\mathbf{X}, \mathbf{w}), \mathbf{y}) = \frac{1}{N} \sum_i -y^{(i)} \log f(\mathbf{x}^{(i)}, \mathbf{w}) - (1 - y^{(i)}) \log(1 - f(\mathbf{x}^{(i)}, \mathbf{w}))$$

Logistic regression loss function

- Logistic Regression

$$P(y = 1 | \mathbf{x}^{(i)}; \mathbf{w}) = f(\mathbf{x}^{(i)}, \mathbf{w})$$

$$P(y = 0 | \mathbf{x}^{(i)}; \mathbf{w}) = 1 - f(\mathbf{x}^{(i)}, \mathbf{w})$$

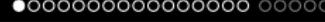
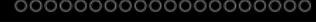
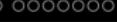
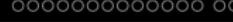
$$y^{(i)} \in \{0, 1\}$$

- Cross entropy as the loss function

$$L(f(\mathbf{X}, \mathbf{w}), \mathbf{y}) = \frac{1}{N} \sum_i L_i(f(\mathbf{x}^{(i)}, \mathbf{w}), y^{(i)})$$

$$L(f(\mathbf{X}, \mathbf{w}), \mathbf{y}) = \frac{1}{N} \sum_i -y^{(i)} \log f(\mathbf{x}^{(i)}, \mathbf{w}) - (1 - y^{(i)}) \log(1 - f(\mathbf{x}^{(i)}, \mathbf{w}))$$

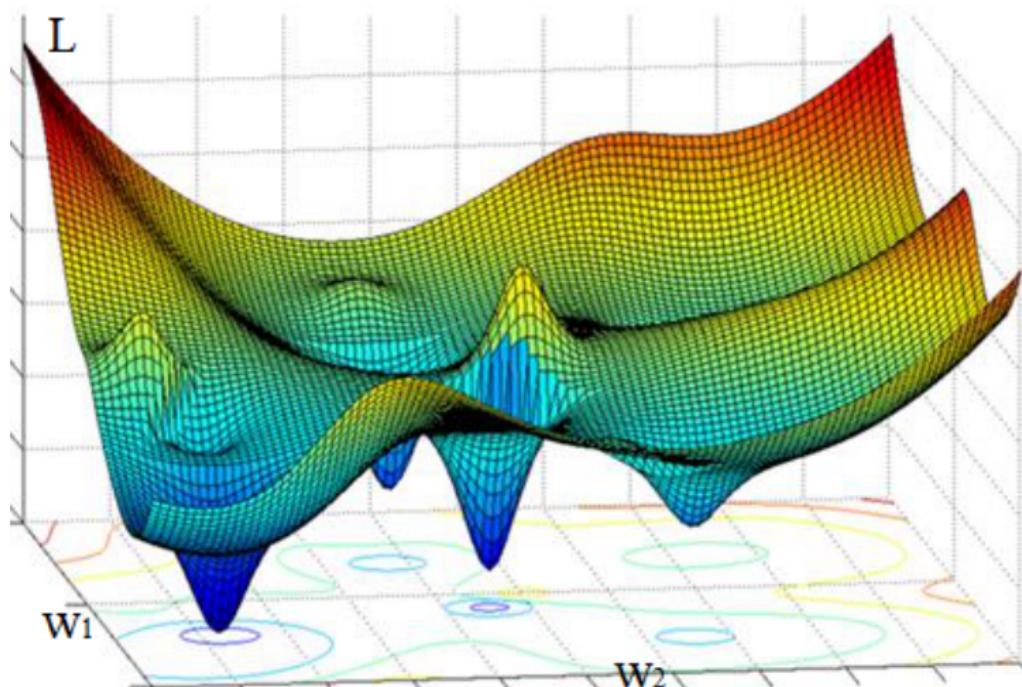
How do we find the best w ?



Optimization

Optimization

$$\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w}} L(f(\mathbf{X}, \mathbf{w}), \mathbf{y})$$



Finding the minimum of a function

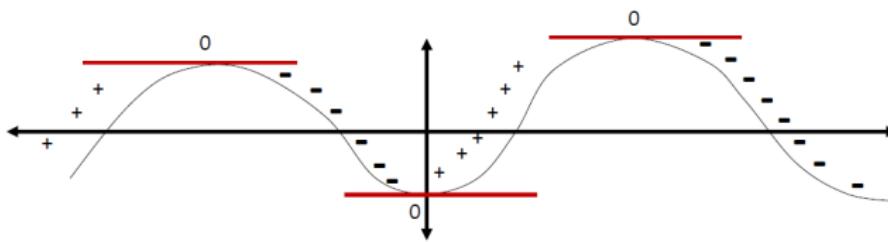
- Solve:

$$\frac{df(x)}{dx} = 0$$

- The solution x^* is a turning point
- Check the double derivative at x^*

$$\frac{d^2 f(x^*)}{dx^{*2}}$$

- If it is positive x^* is a minimum, otherwise it is a maximum



Finding the minimum of a function

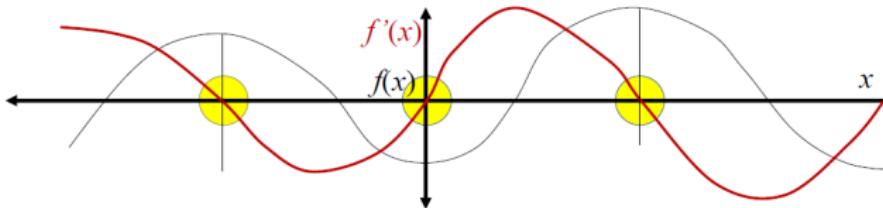
- Solve:

$$\frac{df(x)}{dx} = 0$$

- The solution x^* is a turning point
- Check the double derivative at x^*

$$\frac{d^2 f(x^*)}{dx^{*2}}$$

- If it is positive x^* is a minimum, otherwise it is a maximum



Finding the minimum of a function

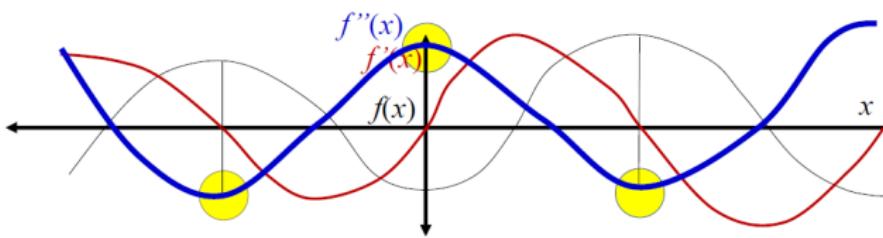
- Solve:

$$\frac{df(x)}{dx} = 0$$

- The solution x^* is a turning point
- Check the double derivative at x^*

$$\frac{d^2 f(x^*)}{dx^{*2}}$$

- If it is positive x^* is a minimum, otherwise it is a maximum



Finding the minimum of a function

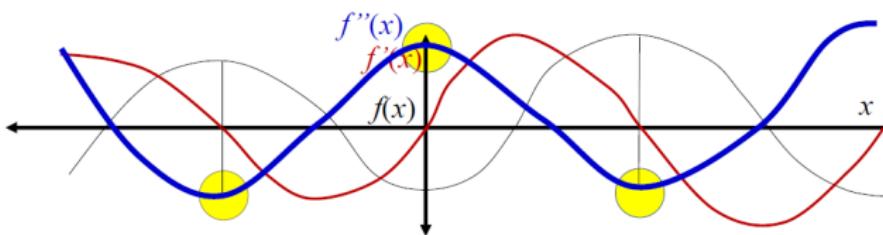
- Solve:

$$\frac{df(x)}{dx} = 0$$

- The solution x^* is a turning point
- Check the double derivative at x^*

$$\frac{d^2 f(x^*)}{dx^{*2}}$$

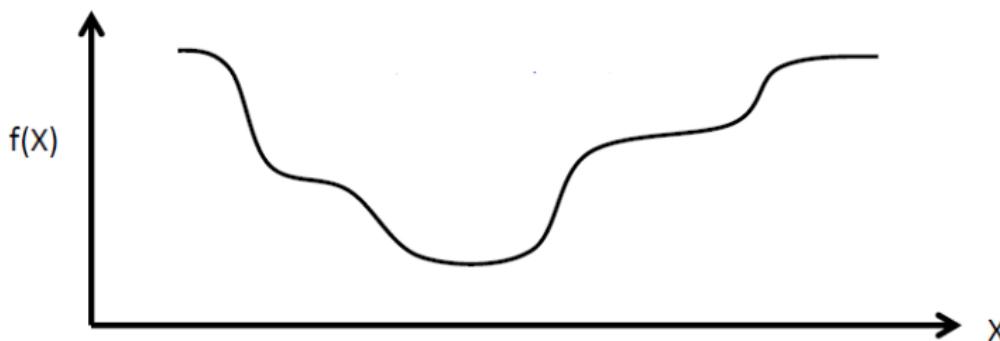
- If it is positive x^* is a minimum, otherwise it is a maximum



It does not work!!!

Closed form solutions are not always available

- Often it is not possible to simply solve $\nabla_X f(X) = 0$
 - The function to minimize/maximize may have an intractable form
- In these situations, iterative solutions are used
 - Begin with a guess for the optimal X and refine it iteratively until the correct value is obtained



Follow the slope



Follow the slope



Follow the slope



Derivative

- In 1-dimension, the derivative of a function gives the slope:

$$\frac{df(x)}{dx} = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$

- In multiple dimensions, the gradient is the vector of (partial derivatives) along each dimension
- The direction of steepest descent is the negative gradient

- To decrease f at x :

if $\frac{df(x)}{dx} < 0 \Rightarrow f(x+h) < f(x) \Rightarrow f$ is decreasing in $x \Rightarrow$ take a step in **positive** direction

if $\frac{df(x)}{dx} > 0 \Rightarrow f(x+h) > f(x) \Rightarrow f$ is increasing in $x \Rightarrow$ take a step in **negative** direction

- Take a step in the reverse direction of derivative at x .

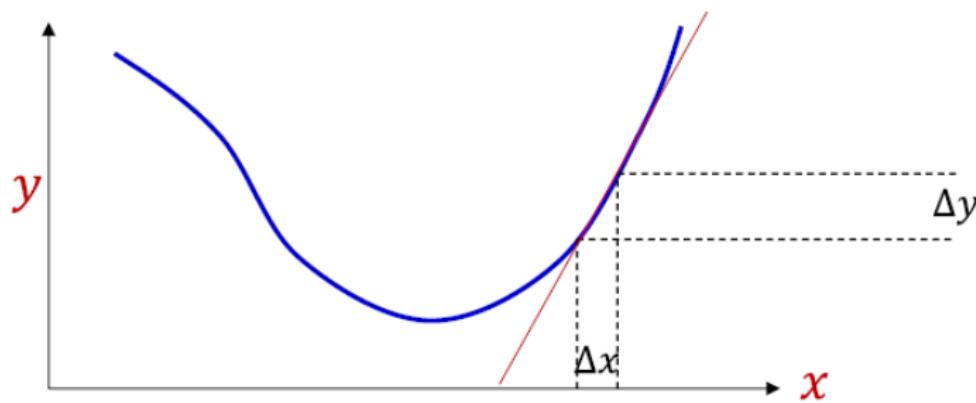
A brief note on derivatives...

- A derivative of a function at any point tells us how much a minute increment to the argument of the function will increment the value of the function

- For any $y = f(x)$, expressed as a multiplier α to a tiny increment Δx to obtain the increments Δy to the output

$$\Delta y = \alpha \Delta x$$

- Based on the fact that at a fine enough resolution, any smooth, continuous function is locally linear at any point



Multivariate scalar function

- Scalar function of vector argument

$$y = f(\mathbf{x}) \quad (f : \mathbb{R}^d \rightarrow \mathbb{R})$$

$$\Delta y = \alpha \Delta \mathbf{x}$$

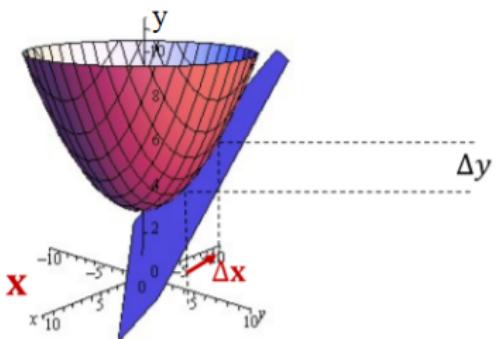
Multivariate scalar function

- Scalar function of vector argument
 $y = f(\mathbf{x})$ ($f : \mathbb{R}^d \rightarrow \mathbb{R}$)

$$\Delta y = \alpha \Delta \mathbf{x}$$

- $\Delta \mathbf{x}$ is now a vector

$$\Delta \mathbf{x} = \begin{bmatrix} \Delta x_1 \\ \Delta x_2 \\ \vdots \\ \Delta x_d \end{bmatrix}$$



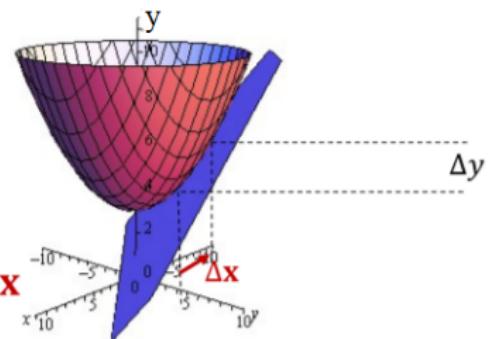
Multivariate scalar function

- Scalar function of vector argument
 $y = f(\mathbf{x})$ ($f : \mathbb{R}^d \rightarrow \mathbb{R}$)

$$\Delta y = \alpha \Delta \mathbf{x}$$

- $\Delta \mathbf{x}$ is now a vector

$$\Delta \mathbf{x} = \begin{bmatrix} \Delta x_1 \\ \Delta x_2 \\ \vdots \\ \Delta x_d \end{bmatrix}$$



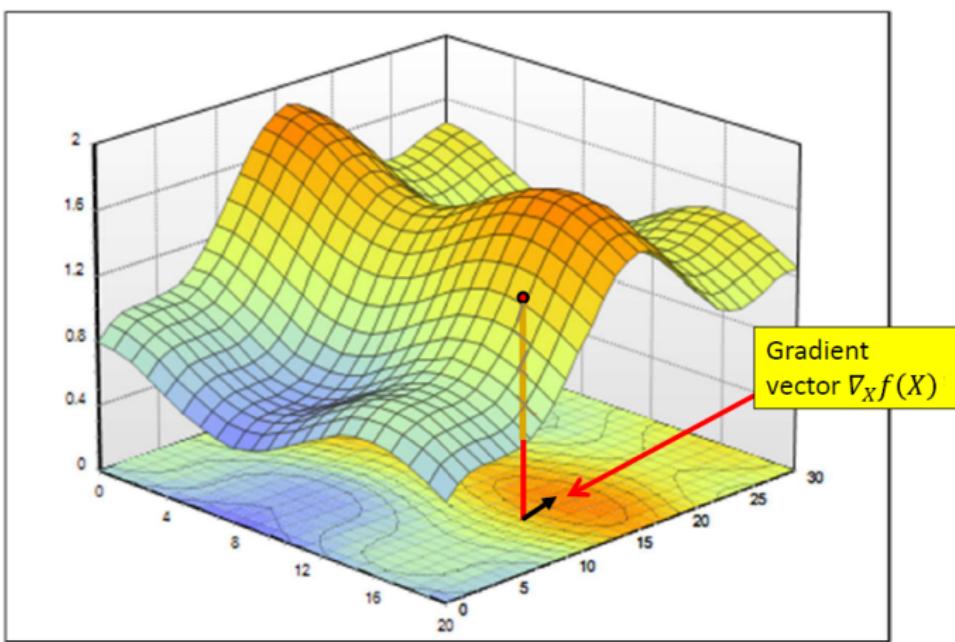
- Giving us that α is a row vector: $\alpha = [\alpha_1 \quad \alpha_2 \quad \dots \quad \alpha_d]$

$$\Delta y = \alpha_1 \Delta x_1 + \alpha_2 \Delta x_2 + \dots + \alpha_d \Delta x_d$$

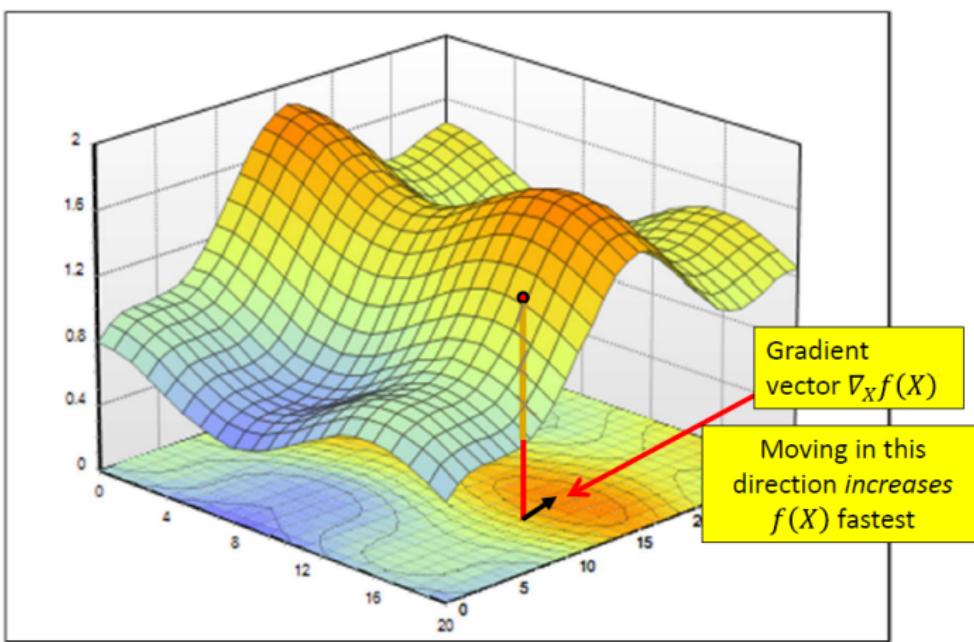
- The partial derivative α_i gives us how y increments when only x_i is incremented.
- Often represented as $\frac{\partial y}{\partial x_i}$

$$\Delta y = \frac{\partial y}{\partial x_1} \Delta x_1 + \frac{\partial y}{\partial x_2} \Delta x_2 + \dots + \frac{\partial y}{\partial x_d} \Delta x_d$$

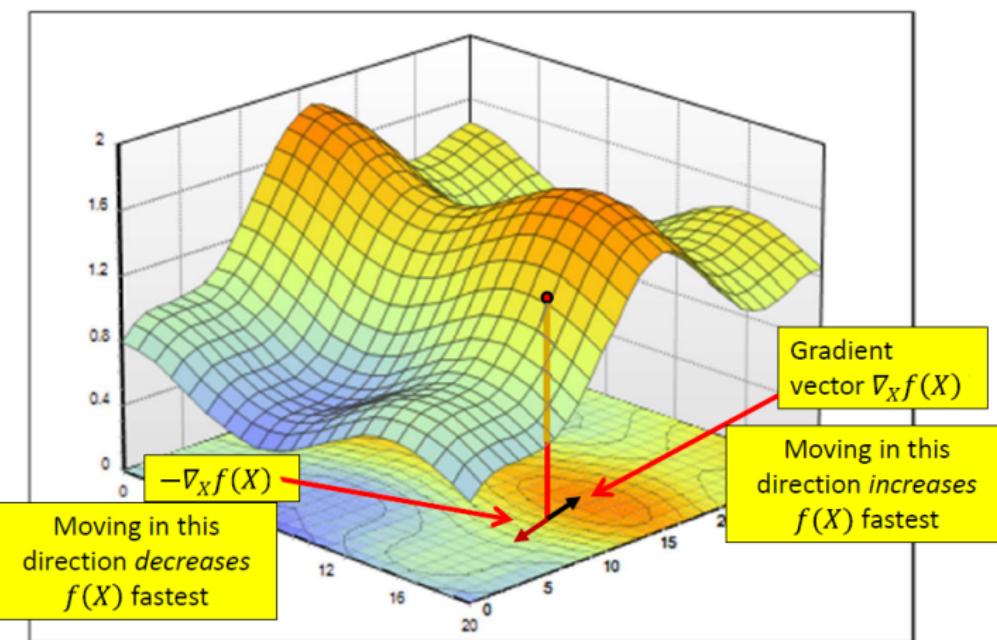
Gradient

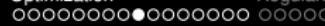
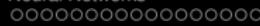
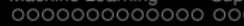


Gradient

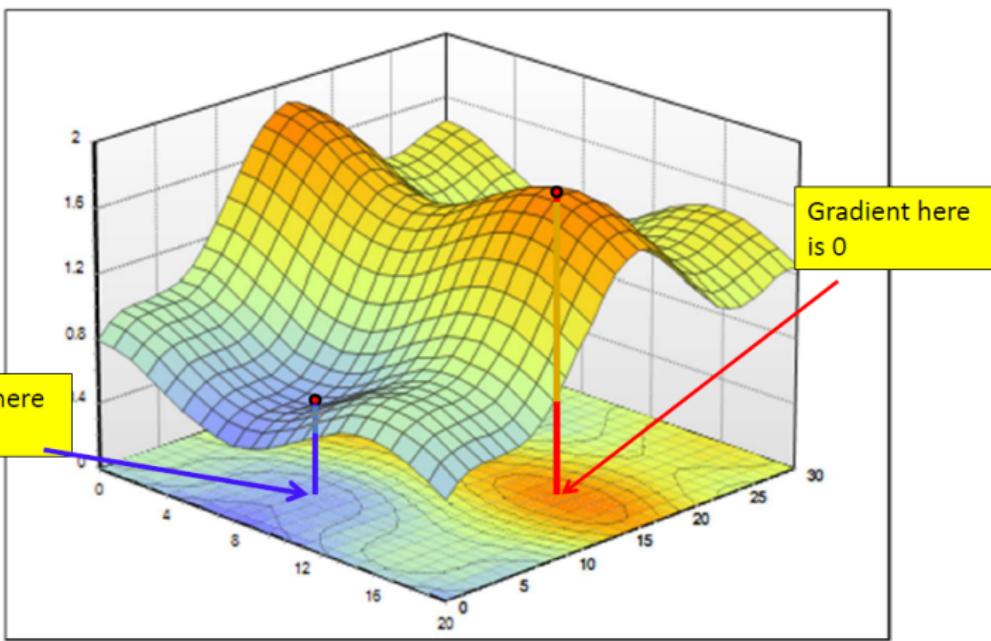


Gradient





Gradient



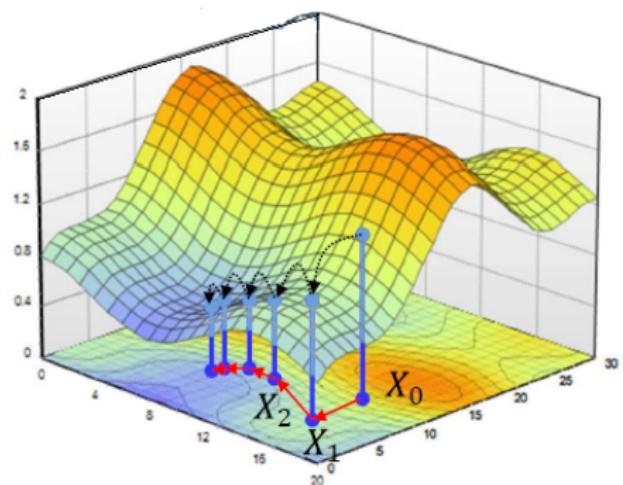
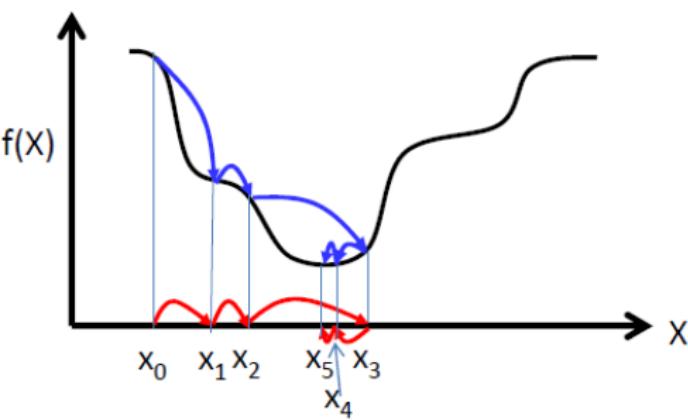
Iterative solutions (Gradient Descent)

- Iterative solutions

- Start from an initial guess X_0 for the optimal X
- Update the guess towards a (hopefully) better value of $f(X)$
- Stop when $f(X)$ no longer decreases

- Problems

- Which direction to step in
- How big must the steps be



Gradient descent

- The gradient descent method to find the minimum of a function iteratively

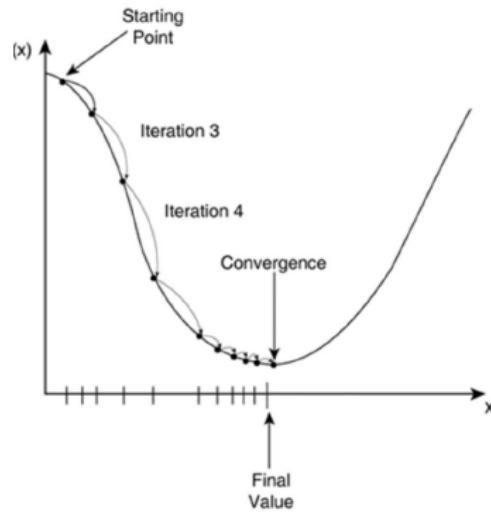
$$\mathbf{x}^{t+1} = \mathbf{x}^t - \eta \nabla_{\mathbf{x}} f(\mathbf{x}^t)$$

- η is the "step size" (Also called "learning rate")

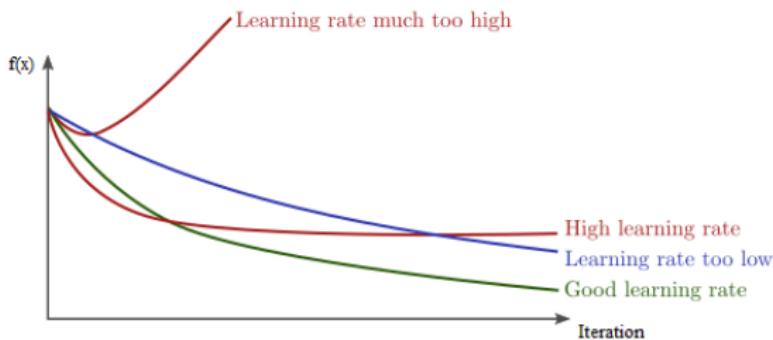
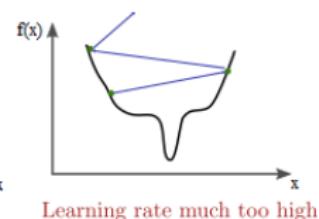
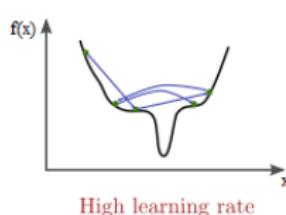
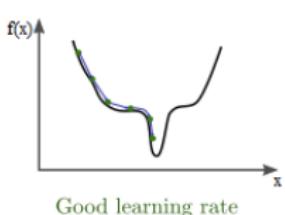
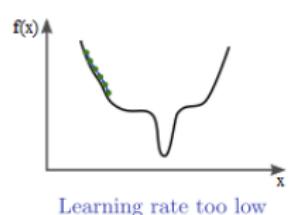
- The gradient descent algorithm converges when the following criterion is satisfied.

$$|f(\mathbf{x}^{t+k}) - f(\mathbf{x}^t)| < \epsilon$$

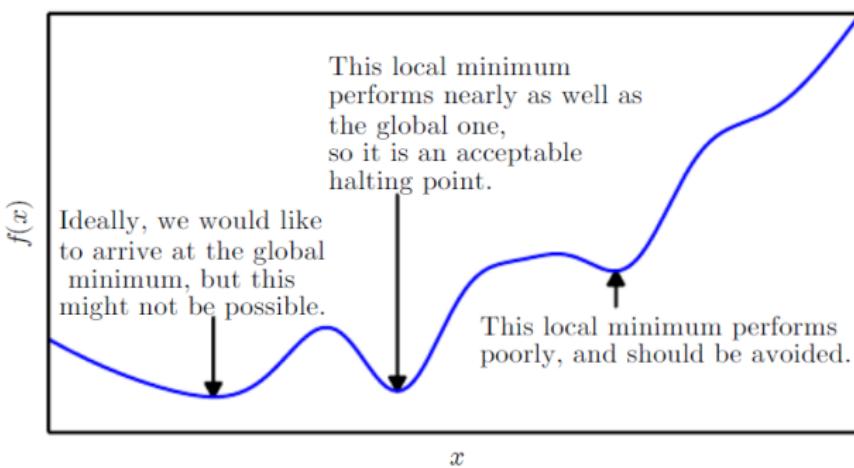
- k is a hyperparameter.



Influence of step size



Local minimum



Logistic regression weights update

■ Logistic regression loss function

$$L(f(X, \mathbf{w}), \mathbf{y}) = \frac{1}{N} \sum_i L_i(f(\mathbf{x}^{(i)}, \mathbf{w}), y^{(i)})$$

$$L(f(X, \mathbf{w}), \mathbf{y}) = \frac{1}{N} \sum_i -y^{(i)} \log f(\mathbf{x}^{(i)}, \mathbf{w}) - (1 - y^{(i)}) \log(1 - f(\mathbf{x}^{(i)}, \mathbf{w}))$$

$$L(f(X, \mathbf{w}), \mathbf{y}) = \frac{1}{N} \sum_i -y^{(i)} \log \sigma(\mathbf{w}^T \mathbf{x}^{(i)}) - (1 - y^{(i)}) \log(1 - \sigma(\mathbf{w}^T \mathbf{x}^{(i)}))$$



Logistic regression weights update

■ Logistic regression loss function

$$L(f(X, \mathbf{w}), \mathbf{y}) = \frac{1}{N} \sum_i L_i(f(\mathbf{x}^{(i)}, \mathbf{w}), y^{(i)})$$

$$L(f(X, \mathbf{w}), \mathbf{y}) = \frac{1}{N} \sum_i -y^{(i)} \log f(\mathbf{x}^{(i)}, \mathbf{w}) - (1 - y^{(i)}) \log(1 - f(\mathbf{x}^{(i)}, \mathbf{w}))$$

$$L(f(X, \mathbf{w}), \mathbf{y}) = \frac{1}{N} \sum_i -y^{(i)} \log \sigma(\mathbf{w}^T \mathbf{x}^{(i)}) - (1 - y^{(i)}) \log(1 - \sigma(\mathbf{w}^T \mathbf{x}^{(i)}))$$

How do we find the best \mathbf{w} ?

Logistic regression weights update

- Logistic regression loss function

$$L(f(X, \mathbf{w}), \mathbf{y}) = \frac{1}{N} \sum_i L_i(f(\mathbf{x}^{(i)}, \mathbf{w}), y^{(i)})$$

$$L(f(X, \mathbf{w}), \mathbf{y}) = \frac{1}{N} \sum_i -y^{(i)} \log f(\mathbf{x}^{(i)}, \mathbf{w}) - (1 - y^{(i)}) \log(1 - f(\mathbf{x}^{(i)}, \mathbf{w}))$$

$$L(f(X, \mathbf{w}), \mathbf{y}) = \frac{1}{N} \sum_i -y^{(i)} \log \sigma(\mathbf{w}^T \mathbf{x}^{(i)}) - (1 - y^{(i)}) \log(1 - \sigma(\mathbf{w}^T \mathbf{x}^{(i)}))$$

How do we find the best \mathbf{w} ? Gradient Descent

$$\frac{\partial L_i}{\partial w_j} = -(y^{(i)} \frac{1}{\sigma(\mathbf{w}^T \mathbf{x}^{(i)})} - (1 - y^{(i)}) \frac{1}{1 - \sigma(\mathbf{w}^T \mathbf{x}^{(i)})}) \frac{\partial}{\partial w_j} \sigma(\mathbf{w}^T \mathbf{x}^{(i)})$$

$$\frac{\partial L_i}{\partial w_j} = -(y^{(i)} \frac{1}{\sigma(\mathbf{w}^T \mathbf{x}^{(i)})} - (1 - y^{(i)}) \frac{1}{1 - \sigma(\mathbf{w}^T \mathbf{x}^{(i)})}) \sigma(\mathbf{w}^T \mathbf{x}^{(i)}) (1 - \sigma(\mathbf{w}^T \mathbf{x}^{(i)})) \frac{\partial}{\partial w_j} \mathbf{w}^T \mathbf{x}^{(i)}$$

$$\frac{\partial L_i}{\partial w_j} = -(y^{(i)} (1 - \sigma(\mathbf{w}^T \mathbf{x}^{(i)})) - (1 - y^{(i)}) \sigma(\mathbf{w}^T \mathbf{x}^{(i)})) x_j^{(i)}$$

$$\frac{\partial L_i}{\partial w_j} = -(y^{(i)} - f(\mathbf{x}^{(i)}, \mathbf{w})) x_j^{(i)}$$

Logistic regression weights update

- Gradient descent to train logistic regression

$$w_j^{t+1} = w_j^t - \eta \frac{\partial L_i}{\partial w_j} \quad (*)$$

$$\frac{\partial L_i}{\partial w_j} = -(y^{(i)} - f(\mathbf{x}^{(i)}, \mathbf{w}))x_j^{(i)} \quad (**)$$

$$\begin{aligned} (*) \& \Rightarrow w_j^{t+1} = w_j^t - \eta \frac{1}{N} \sum_{i=1}^N (-(y^{(i)} - f(\mathbf{x}^{(i)}, \mathbf{w}^t))x_j^{(i)}) \\ & w_j^{t+1} = w_j^t + \eta \frac{1}{N} \sum_{i=1}^N (y^{(i)} - f(\mathbf{x}^{(i)}, \mathbf{w}^t))x_j^{(i)} \\ \Rightarrow \mathbf{w}^{t+1} &= \mathbf{w}^t + \eta \frac{1}{N} \sum_{i=1}^N (y^{(i)} - f(\mathbf{x}^{(i)}, \mathbf{w}^t))\mathbf{x}^{(i)} \end{aligned}$$

- Perceptron learning algorithm

For misclassified sample $\mathbf{x}^{(i)}$:

$$\text{If } y^{(i)} = 0 : \quad \mathbf{w}^{t+1} = \mathbf{w}^t - \mathbf{x}^{(i)}$$

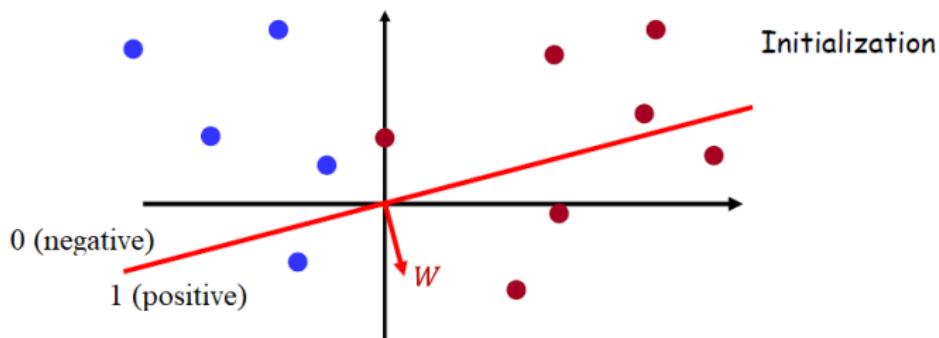
$$\text{If } y^{(i)} = 1 : \quad \mathbf{w}^{t+1} = \mathbf{w}^t + \mathbf{x}^{(i)}$$

Perceptron learning algorithm

For misclassified sample $\mathbf{x}^{(i)}$:

$$\text{If } y^{(i)} = 0 : \quad \mathbf{w}^{t+1} = \mathbf{w}^t - \mathbf{x}^{(i)}$$

$$\text{If } y^{(i)} = 1 : \quad \mathbf{w}^{t+1} = \mathbf{w}^t + \mathbf{x}^{(i)}$$

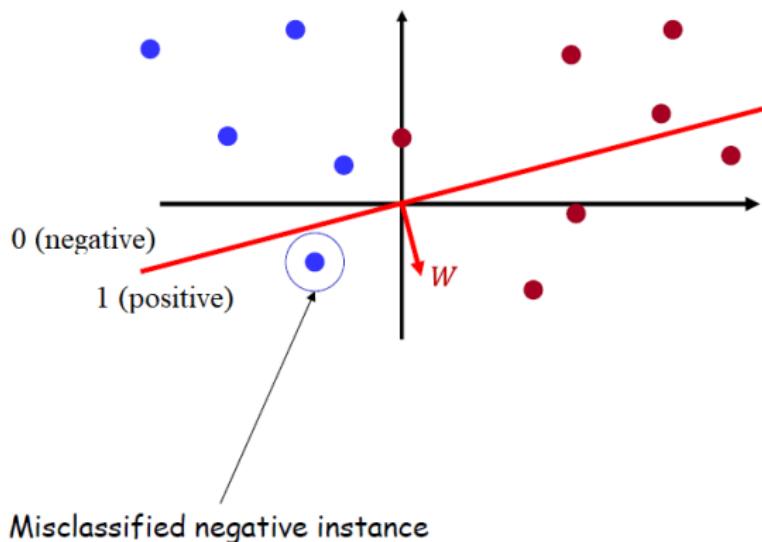


Perceptron learning algorithm

For misclassified sample $\mathbf{x}^{(i)}$:

$$\text{If } y^{(i)} = 0 : \quad \mathbf{w}^{t+1} = \mathbf{w}^t - \mathbf{x}^{(i)}$$

$$\text{If } y^{(i)} = 1 : \quad \mathbf{w}^{t+1} = \mathbf{w}^t + \mathbf{x}^{(i)}$$

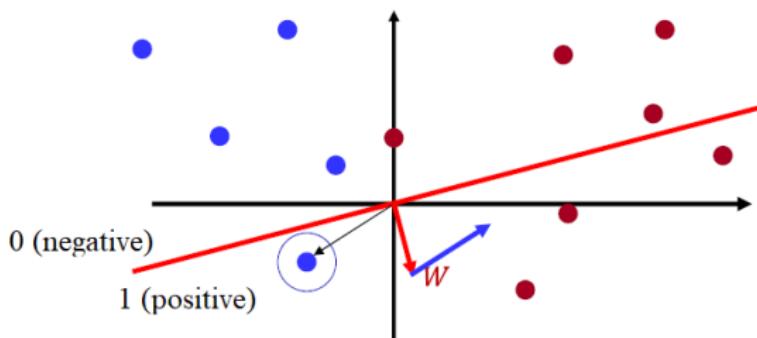


Perceptron learning algorithm

For misclassified sample $\mathbf{x}^{(i)}$:

$$\text{If } y^{(i)} = 0 : \quad \mathbf{w}^{t+1} = \mathbf{w}^t - \mathbf{x}^{(i)}$$

$$\text{If } y^{(i)} = 1 : \quad \mathbf{w}^{t+1} = \mathbf{w}^t + \mathbf{x}^{(i)}$$



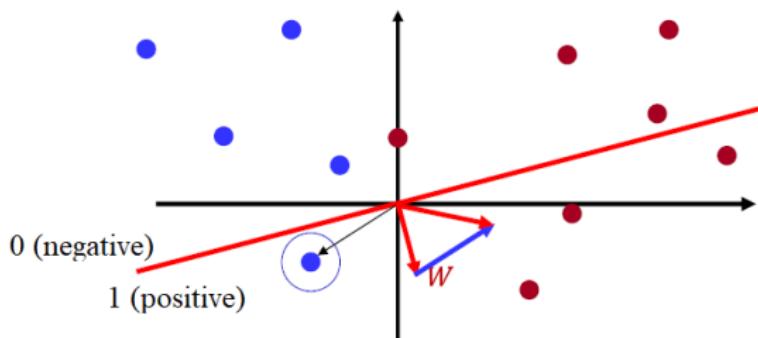
Misclassified *negative* instance, subtract it from W

Perceptron learning algorithm

For misclassified sample $\mathbf{x}^{(i)}$:

$$\text{If } y^{(i)} = 0 : \quad \mathbf{w}^{t+1} = \mathbf{w}^t - \mathbf{x}^{(i)}$$

$$\text{If } y^{(i)} = 1 : \quad \mathbf{w}^{t+1} = \mathbf{w}^t + \mathbf{x}^{(i)}$$



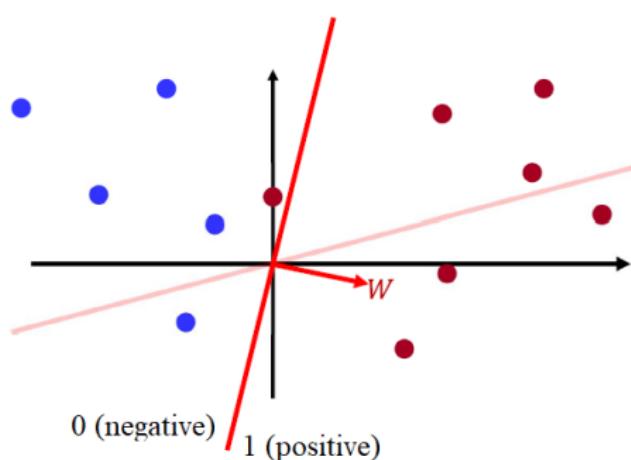
The new weight

Perceptron learning algorithm

For misclassified sample $\mathbf{x}^{(i)}$:

$$\text{If } y^{(i)} = 0 : \quad \mathbf{w}^{t+1} = \mathbf{w}^t - \mathbf{x}^{(i)}$$

$$\text{If } y^{(i)} = 1 : \quad \mathbf{w}^{t+1} = \mathbf{w}^t + \mathbf{x}^{(i)}$$



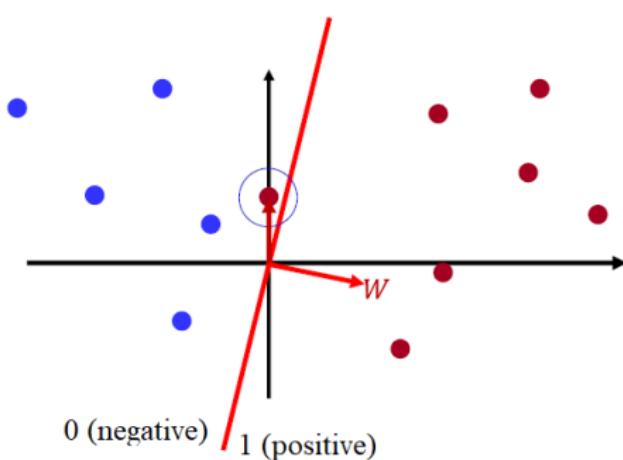
The new weight (and boundary)

Perceptron learning algorithm

For misclassified sample $\mathbf{x}^{(i)}$:

$$\text{If } y^{(i)} = 0 : \quad \mathbf{w}^{t+1} = \mathbf{w}^t - \mathbf{x}^{(i)}$$

$$\text{If } y^{(i)} = 1 : \quad \mathbf{w}^{t+1} = \mathbf{w}^t + \mathbf{x}^{(i)}$$



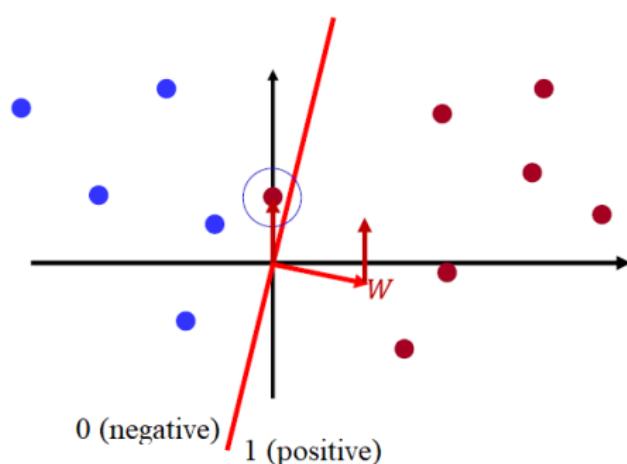
Misclassified *positive* instance

Perceptron learning algorithm

For misclassified sample $\mathbf{x}^{(i)}$:

$$\text{If } y^{(i)} = 0 : \quad \mathbf{w}^{t+1} = \mathbf{w}^t - \mathbf{x}^{(i)}$$

$$\text{If } y^{(i)} = 1 : \quad \mathbf{w}^{t+1} = \mathbf{w}^t + \mathbf{x}^{(i)}$$



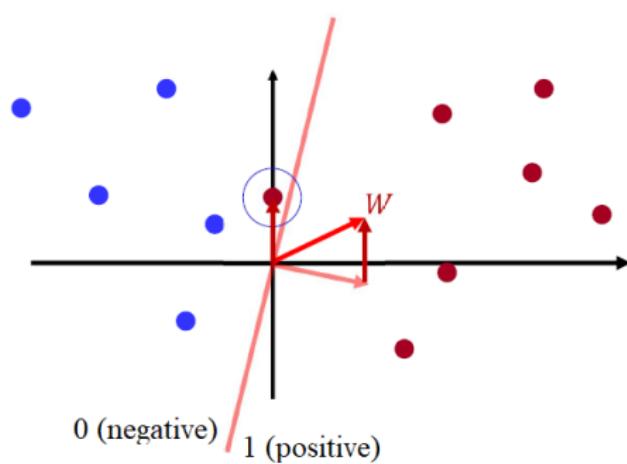
Misclassified *positive* instance, add it to W

Perceptron learning algorithm

For misclassified sample $\mathbf{x}^{(i)}$:

$$\text{If } y^{(i)} = 0 : \quad \mathbf{w}^{t+1} = \mathbf{w}^t - \mathbf{x}^{(i)}$$

$$\text{If } y^{(i)} = 1 : \quad \mathbf{w}^{t+1} = \mathbf{w}^t + \mathbf{x}^{(i)}$$



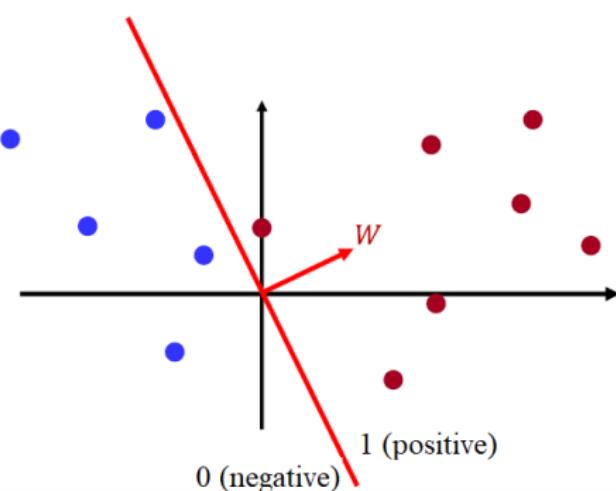
The new weight vector

Perceptron learning algorithm

For misclassified sample $\mathbf{x}^{(i)}$:

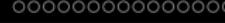
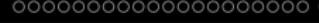
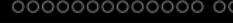
$$\text{If } y^{(i)} = 0 : \quad \mathbf{w}^{t+1} = \mathbf{w}^t - \mathbf{x}^{(i)}$$

$$\text{If } y^{(i)} = 1 : \quad \mathbf{w}^{t+1} = \mathbf{w}^t + \mathbf{x}^{(i)}$$

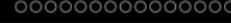
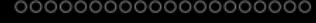
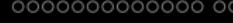


The new decision boundary

Perfect classification, no more updates, we are done



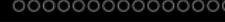
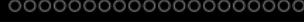
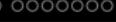
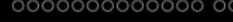
Regularization



Learning VS. Optimization

Generalization VS. Overfitting

- Ability to perform well on previously unobserved inputs is called **generalization**.



Learning VS. Optimization

Generalization VS. Overfitting

- Ability to perform well on previously unobserved inputs is called **generalization**.
- What separates machine learning from optimization is that we want the **generalization error**, also called the test error, to be low as well.

Learning VS. Optimization

Generalization VS. Overfitting

- Ability to perform well on previously unobserved inputs is called **generalization**.
- What separates machine learning from optimization is that we want the **generalization error**, also called the test error, to be low as well.
- We typically estimate the generalization error of a machine learning model by measuring its performance on a **test set** of examples that were collected separately from the training set.
- How can we affect performance on the test set when we get to observe only the training set? The field of **statistical learning theory** provides some answers.

Learning VS. Optimization

Generalization VS. Overfitting

- Ability to perform well on previously unobserved inputs is called **generalization**.
- What separates machine learning from optimization is that we want the **generalization error**, also called the test error, to be low as well.
- We typically estimate the generalization error of a machine learning model by measuring its performance on a **test set** of examples that were collected separately from the training set.
- How can we affect performance on the test set when we get to observe only the training set? The field of **statistical learning theory** provides some answers.
- If the training and the test set are collected arbitrarily, there is indeed little we can do.

Learning VS. Optimization

Generalization VS. Overfitting

- Ability to perform well on previously unobserved inputs is called **generalization**.
- What separates machine learning from optimization is that we want the **generalization error**, also called the test error, to be low as well.
- We typically estimate the generalization error of a machine learning model by measuring its performance on a **test set** of examples that were collected separately from the training set.
- How can we affect performance on the test set when we get to observe only the training set? The field of **statistical learning theory** provides some answers.
- If the training and the test set are collected arbitrarily, there is indeed little we can do.
- The train and test data are generated by a probability distribution over datasets called the data generating process. We typically make a set of assumptions known collectively as the **i.i.d.** assumptions.

Learning VS. Optimization

Generalization VS. Overfitting

- Ability to perform well on previously unobserved inputs is called **generalization**.
- What separates machine learning from optimization is that we want the **generalization error**, also called the test error, to be low as well.
- We typically estimate the generalization error of a machine learning model by measuring its performance on a **test set** of examples that were collected separately from the training set.
- How can we affect performance on the test set when we get to observe only the training set? The field of **statistical learning theory** provides some answers.
- If the training and the test set are collected arbitrarily, there is indeed little we can do.
- The train and test data are generated by a probability distribution over datasets called the data generating process. We typically make a set of assumptions known collectively as the **i.i.d.** assumptions.
- These assumptions are that the examples in each dataset are independent from each other, and that the train set and test set are identically distributed, drawn from the same probability distribution as each other.

(Deep Learning, Ian Goodfellow, MIT press, 2016)

Machine learning

underfitting and overfitting

- The factors determining how well a machine learning algorithm will perform are its ability to:
 - Make the training error small.
 - Make the gap between training and test error small.

Machine learning

underfitting and overfitting

- The factors determining how well a machine learning algorithm will perform are its ability to:
 - Make the training error small.
 - Make the gap between training and test error small.
- These two factors correspond to the two central challenges in machine learning: **underfitting** and **overfitting**. Underfitting occurs when the model is not able to obtain a sufficiently low error value on the training set. Overfitting occurs when the gap between the training error and test error is too large.

Machine learning

underfitting and overfitting

- The factors determining how well a machine learning algorithm will perform are its ability to:
 - Make the training error small.
 - Make the gap between training and test error small.
- These two factors correspond to the two central challenges in machine learning: **underfitting** and **overfitting**. Underfitting occurs when the model is not able to obtain a sufficiently low error value on the training set. Overfitting occurs when the gap between the training error and test error is too large.
- We can control whether a model is more likely to overfit or underfit by altering its **capacity**.

Machine learning

underfitting and overfitting

- The factors determining how well a machine learning algorithm will perform are its ability to:
 - Make the training error small.
 - Make the gap between training and test error small.
- These two factors correspond to the two central challenges in machine learning: **underfitting** and **overfitting**. Underfitting occurs when the model is not able to obtain a sufficiently low error value on the training set. Overfitting occurs when the gap between the training error and test error is too large.
- We can control whether a model is more likely to overfit or underfit by altering its **capacity**.
- Models with low capacity may struggle to fit the training set. Models with high capacity can overfit by memorizing properties of the training set that do not serve them well on the test set.

Machine learning

underfitting and overfitting

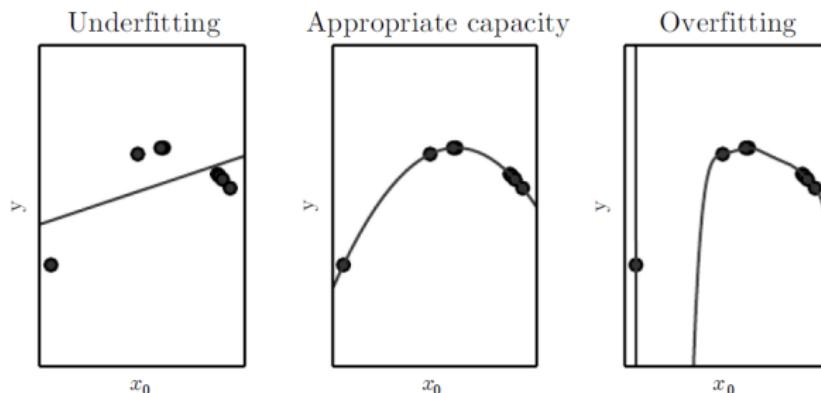
- The factors determining how well a machine learning algorithm will perform are its ability to:
 - Make the training error small.
 - Make the gap between training and test error small.
- These two factors correspond to the two central challenges in machine learning: **underfitting** and **overfitting**. Underfitting occurs when the model is not able to obtain a sufficiently low error value on the training set. Overfitting occurs when the gap between the training error and test error is too large.
- We can control whether a model is more likely to overfit or underfit by altering its **capacity**.
- Models with low capacity may struggle to fit the training set. Models with high capacity can overfit by memorizing properties of the training set that do not serve them well on the test set.
- One way to control the capacity of a learning algorithm is by choosing its **hypothesis space**, the set of functions that the learning algorithm is allowed to select as being the solution.

(Deep Learning, Ian Goodfellow, MIT press, 2016)

Machine learning

underfitting and overfitting

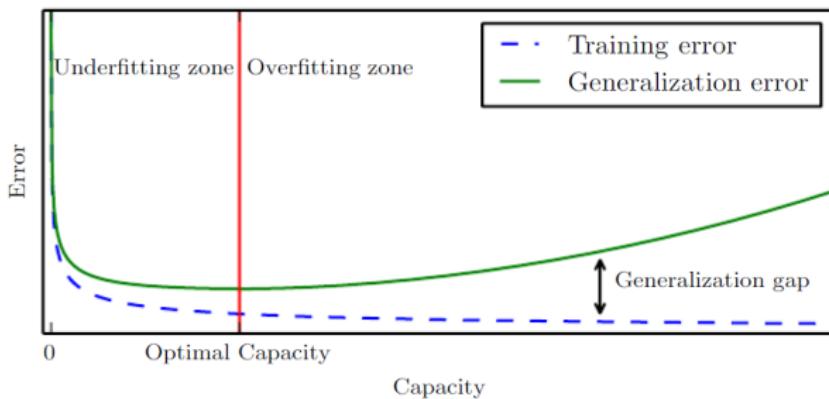
- Machine learning algorithms will generally perform best when their capacity is appropriate for the true complexity of the task they need to perform and the amount of training data they are provided with.



(Deep Learning, Ian Goodfellow, MIT press, 2016)

Capacity

Typically, generalization error has a U-shaped curve as a function of model capacity.

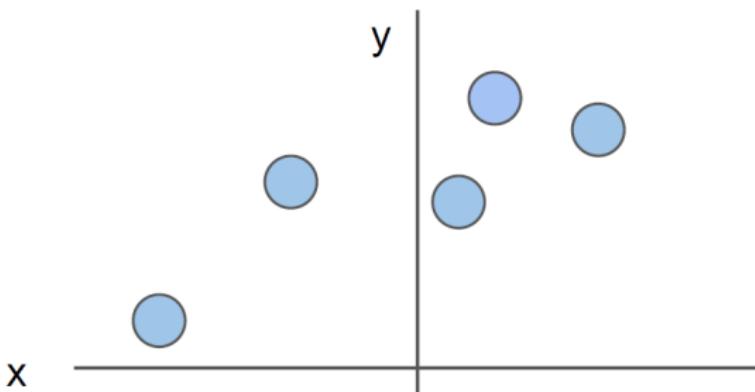


Regularization

- **Regularization** is any modification we make to a learning algorithm that is intended to reduce its generalization error but not its training error.
 - We can give a learning algorithm a **preference** for one solution in its hypothesis space to another. This means that both functions are eligible, but one is preferred.

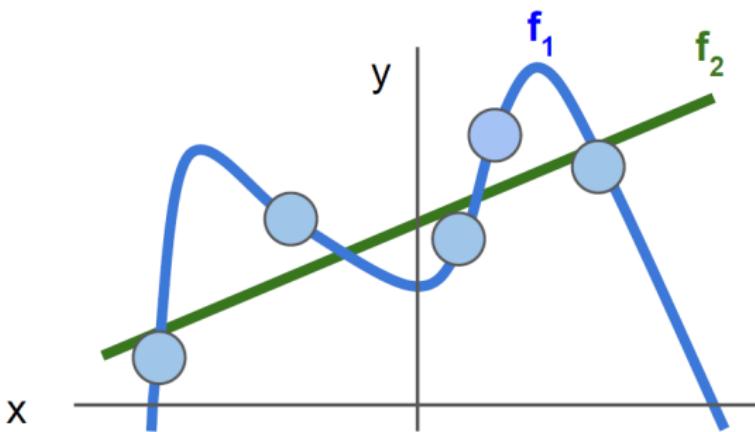
Regularization

- How do we choose between hypotheses?
 - Occams Razor: Among multiple competing hypotheses, the simplest is the best, (William of Ockham 1285-1347)
 - Avoid overfitting: prefer simple models that generalize better
- Regularization intuition: toy example training data



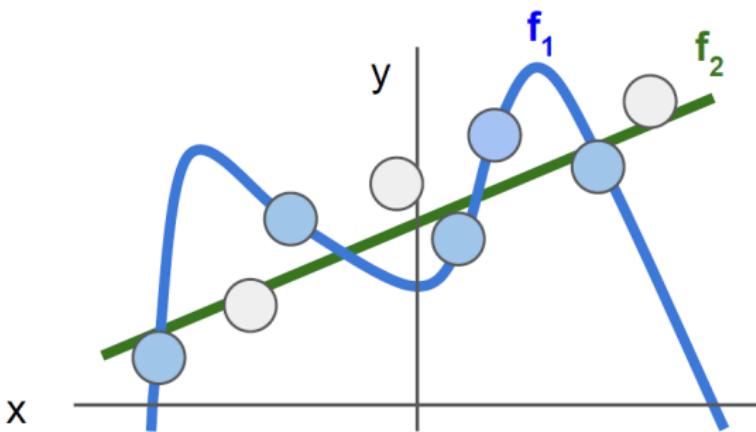
Regularization

- How do we choose between hypotheses?
 - Occams Razor: Among multiple competing hypotheses, the simplest is the best, (William of Ockham 1285-1347)
 - Avoid overfitting: prefer simple models that generalize better
- Regularization intuition: Prefer Simpler Models



Regularization

- How do we choose between hypotheses?
 - Occams Razor: Among multiple competing hypotheses, the simplest is the best, (William of Ockham 1285-1347)
 - Avoid overfitting: prefer simple models that generalize better
- Regularization: Prefer Simpler Models



- Regularization pushes against fitting the data too well so we don't fit noise in the data

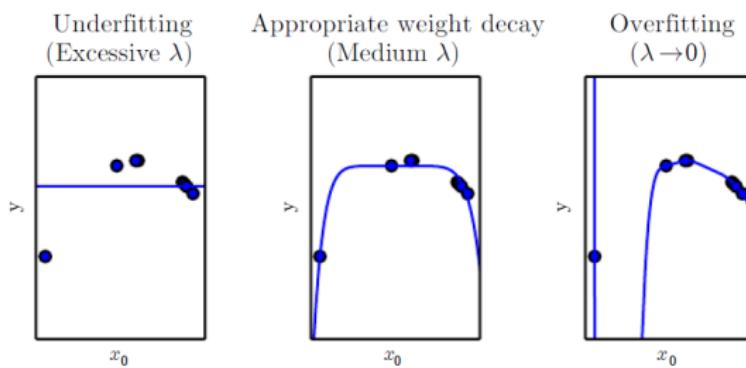
Regularization

- **Training loss:** Model predictions should match training data
- **Regularization:** Prevent the model from doing too well on training data $R(\mathbf{w})$

$$L(f(X, \mathbf{w}), \mathbf{y}) = \frac{1}{N} \sum_i L_i(f(\mathbf{x}^{(i)}, \mathbf{w}), y^{(i)}) + \lambda R(\mathbf{w})$$

- λ : regularization strength (hyperparameter)
- Simple examples

- L_1 regularization: $R(\mathbf{w}) = \sum_i |w_i|$
- L_2 regularization: $R(\mathbf{w}) = \sum_i w_i^2$



Regularization

- Which of w_1 or w_2 will the L2 regularizer prefer? (Training loss is the same)
 - L2 regularization likes to spread out the weights

$$\mathbf{x}^{(i)} = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} \quad \mathbf{w}_1 = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} \quad \mathbf{w}_2 = \begin{bmatrix} 0.25 \\ 0.25 \\ 0.25 \\ 0.25 \end{bmatrix}$$

$$\mathbf{w}_1^T \mathbf{x}^{(i)} = \mathbf{w}_2^T \mathbf{x}^{(i)} = 1$$

Regularization

- Which of w_1 or w_2 will the L2 regularizer prefer? (Training loss is the same)
 - L2 regularization likes to spread out the weights

$$\mathbf{x}^{(i)} = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} \quad \mathbf{w}_1 = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} \quad \mathbf{w}_2 = \begin{bmatrix} 0.25 \\ 0.25 \\ 0.25 \\ 0.25 \end{bmatrix}$$

$$\mathbf{w}_1^T \mathbf{x}^{(i)} = \mathbf{w}_2^T \mathbf{x}^{(i)} = 1$$

$$L2(\mathbf{w}_1) = 1, \quad L2(\mathbf{w}_2) = 0.25$$

Regularization

- Which of w_1 or w_2 will the $L2$ regularizer prefer? (Training loss is the same)
 - $L2$ regularization likes to spread out the weights

$$\mathbf{x}^{(i)} = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} \quad \mathbf{w}_1 = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} \quad \mathbf{w}_2 = \begin{bmatrix} 0.25 \\ 0.25 \\ 0.25 \\ 0.25 \end{bmatrix}$$

$$\mathbf{w}_1^T \mathbf{x}^{(i)} = \mathbf{w}_2^T \mathbf{x}^{(i)} = 1$$

$$L2(\mathbf{w}_1) = 1, \quad L2(\mathbf{w}_2) = 0.25$$

- More complex regularization methods
 - Dropout
 - MixUp
 - Stochastic depth