# Announcements

## Assignments:

- HW4
  - Due date Mon, 2/24, 11:59 pm

## Midterm

- Monday the 2nd of March from 5:00pm-6:30pm

## Midterm Conflicts

- See Piazza post
- Due 11:59pm on Wednesday the 19th of February

# Plan

## Last time

- Naïve Bayes Assumptions
- Naïve Bayes MLE and MAP
- MLE vs MAP
- Generative vs Discriminative Models

## Today

- Decision Boundaries
- Gaussian Generative Models
- Neural Networks

# Introduction to Machine Learning

## Generative Models
## then
## Intro to Neural Networks

Instructor: Pat Virtue

# Decision Boundaries

## Decision boundary

- The set of points in the domain of the input ($x$) where the predicted classification changes

## Two class decision boundary

- So far, we have decided to let the decision boundary be all $x$ such that:

$$p(y = 0 \mid x) = p(y = 1 \mid x)$$

- What assumptions are we making here?
  - This assumes that the cost of predicting it wrong is the same for both classes

# Piazza Poll 1

Which of the following also define the decision boundary for two classes when we just want $p(Y = 0 \mid x) = p(Y = 1 \mid x)$?

A. All $x$, s.t. $p(x \mid Y = 0) = p(x \mid Y = 1)$

B. All $x$, s.t. $p(x, Y = 0) = p(x, Y = 1)$

C. All $x$, s.t. $p(Y = 0) = p(Y = 1)$

D. All $x$, s.t. $p(Y = 1 \mid x) = 0.5$

E. All $x$, s.t. $p(x \mid Y = 1) = 0.5$

F. All $x$, s.t. $p(x, Y = 1) = 0.5$

G. All $x$, s.t. $\log p(x, Y = 1) - \log p(x, Y = 0) = 0$

H. None of the above

# Piazza Poll 1

Which of the following also define the decision boundary for two classes when we just want $p(Y = 0 \mid x) = p(Y = 1 \mid x)$?

A. All $x$, s.t. $p(x \mid Y = 0) = p(x \mid Y = 1)$

B. All $x$, s.t. $p(x, Y = 0) = p(x, Y = 1)$

C. All $x$, s.t. $p(Y = 0) = p(Y = 1)$

D. All $x$, s.t. $p(Y = 1 \mid x) = 0.5$

E. All $x$, s.t. $p(x \mid Y = 1) = 0.5$

F. All $x$, s.t. $p(x, Y = 1) = 0.5$

G. All $x$, s.t. $\log p(x, Y = 1) - \log p(x, Y = 0) = 0$

H. None of the above

# Piazza Poll 2

True/False: Logistic regression always produces a linear decision boundary.
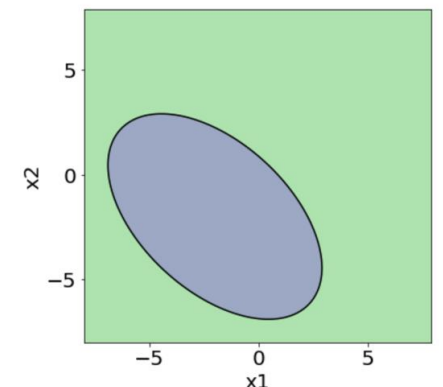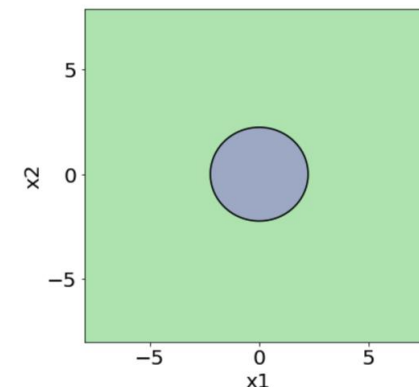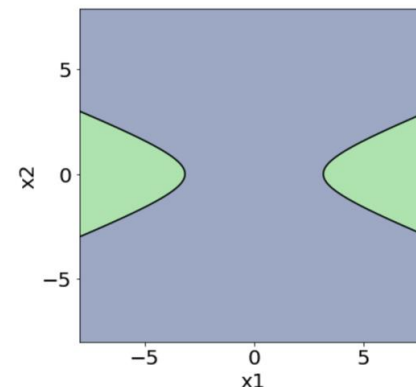
A. I don't know

B. True

C. False

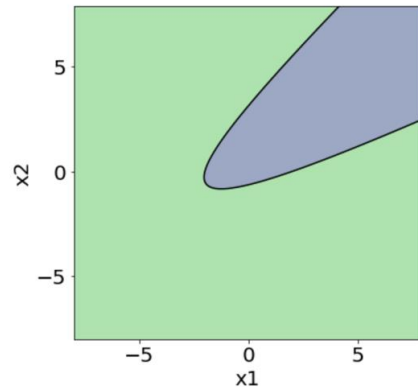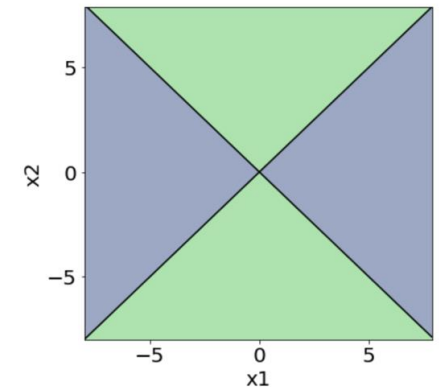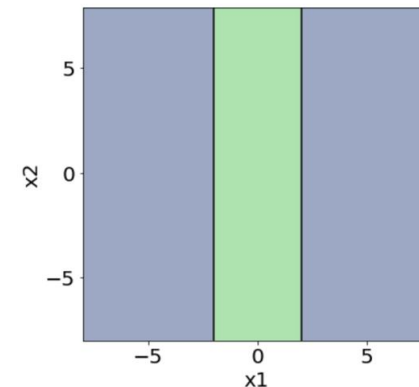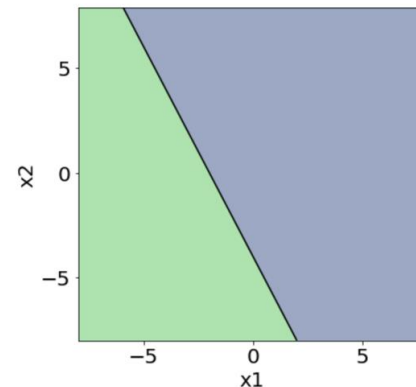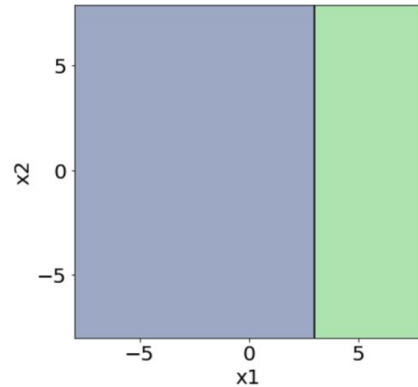# Piazza Poll 2

True/False: Logistic regression always produces a linear decision boundary.

A. I don't know

B. True

C. False
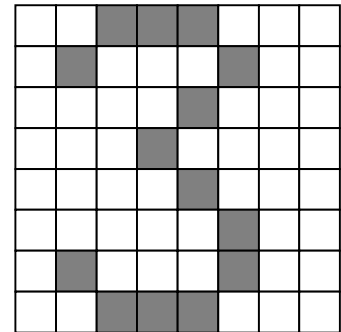
# Generative Models

- Class distribution: $Y \sim Bern(\phi)$
- Class conditional distribution: $X_m \sim Bern(\theta_{m,y})$
- Naïve Bayes $X_i$ conditionally independent $X_j$ given $Y$ for all $i \neq j$
$$p(X_i, X_j \mid Y) = p(X_i \mid Y) \mid p(X_j \mid Y)$$

Digits:

- Class distribution: $Y \sim Multinomial(\phi, 1)$
- Class conditional distribution: $X_m \sim Bern(\theta_{m,y})$
- Naïve Bayes $X_i$ conditionally independent $X_j$ given $Y$ for all $i \neq j$
$$p(X_i, X_j \mid Y) = p(X_i \mid Y) \mid p(X_j \mid Y)$$



Recitation?

# Fisher Iris Dataset

https://en.wikipedia.org/wiki/Iris_flower_data_set

# Fisher Iris Dataset

https://en.wikipedia.org/wiki/Iris_flower_data_set

Fisher (1936) used 150 measurements of flowers from 3 different species: Iris setosa (0), Iris virginica (1), Iris versicolor (2) collected by Anderson (1936)

| Species | Sepal Length | Sepal Width | Petal Length | Petal Width |
|---------|--------------|-------------|--------------|-------------|
| 0 | 4.3 | 3.0 | 1.1 | 0.1 |
| 0 | 4.9 | 3.6 | 1.4 | 0.1 |
| 0 | 5.3 | 3.7 | 1.5 | 0.2 |
| 1 | 4.9 | 2.4 | 3.3 | 1.0 |
| 1 | 5.7 | 2.8 | 4.1 | 1.3 |
| 1 | 6.3 | 3.3 | 4.7 | 1.6 |
| 1 | 6.7 | 3.0 | 5.0 | 1.7 |

# Fisher Iris Dataset



Iris data

# Fisher Iris Dataset

# Generative Models with Continuous Features

Iris dataset:

- Class distribution: $Y \sim Bern(\phi)$
- Class conditional distribution: $\boldsymbol{X} \sim \mathcal{N}(\boldsymbol{\mu}_y, \boldsymbol{\Sigma}_y)$
- Naïve Bayes assumption?

# Piazza Poll 3

Iris dataset:
- Class distribution: $Y \sim Bern(\phi)$
- Class conditional distribution: $\boldsymbol{X} \sim \mathcal{N}(\boldsymbol{\mu}_y, \boldsymbol{\Sigma}_y)$
- Naïve Bayes assumption?

Which of the following pairs of Gaussian class conditional distributions satisfy the Naïve Bayes assumptions? Select ALL that apply.

A. $\boldsymbol{\mu}_{y=0} = \begin{bmatrix} -1 \\ 0 \end{bmatrix}, \boldsymbol{\Sigma}_{y=0} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, \qquad \boldsymbol{\mu}_{y=1} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \boldsymbol{\Sigma}_{y=1} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$

B. $\boldsymbol{\mu}_{y=0} = \begin{bmatrix} -1 \\ 0 \end{bmatrix}, \boldsymbol{\Sigma}_{y=0} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, \qquad \boldsymbol{\mu}_{y=1} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \boldsymbol{\Sigma}_{y=1} = \begin{bmatrix} 3 & 0 \\ 0 & 3 \end{bmatrix}$

C. $\boldsymbol{\mu}_{y=0} = \begin{bmatrix} -1 \\ 0 \end{bmatrix}, \boldsymbol{\Sigma}_{y=0} = \begin{bmatrix} 1 & -1 \\ 1 & 1 \end{bmatrix}, \qquad \boldsymbol{\mu}_{y=1} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \boldsymbol{\Sigma}_{y=1} = \begin{bmatrix} 1 & -1 \\ 1 & 1 \end{bmatrix}$

D. $\boldsymbol{\mu}_{y=0} = \begin{bmatrix} -1 \\ 0 \end{bmatrix}, \boldsymbol{\Sigma}_{y=0} = \begin{bmatrix} 1 & -1 \\ 1 & 1 \end{bmatrix}, \qquad \boldsymbol{\mu}_{y=1} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \boldsymbol{\Sigma}_{y=1} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$

# Piazza Poll 3

A. $\boldsymbol{\mu}_{y=0} = \begin{bmatrix} -1 \\ 0 \end{bmatrix}, \boldsymbol{\Sigma}_{y=0} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, \quad \boldsymbol{\mu}_{y=1} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \boldsymbol{\Sigma}_{y=1} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$

B. $\boldsymbol{\mu}_{y=0} = \begin{bmatrix} -1 \\ 0 \end{bmatrix}, \boldsymbol{\Sigma}_{y=0} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, \quad \boldsymbol{\mu}_{y=1} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \boldsymbol{\Sigma}_{y=1} = \begin{bmatrix} 3 & 0 \\ 0 & 3 \end{bmatrix}$

C. $\boldsymbol{\mu}_{y=0} = \begin{bmatrix} -1 \\ 0 \end{bmatrix}, \boldsymbol{\Sigma}_{y=0} = \begin{bmatrix} 1 & -1 \\ 1 & 1 \end{bmatrix}, \quad \boldsymbol{\mu}_{y=1} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \boldsymbol{\Sigma}_{y=1} = \begin{bmatrix} 1 & -1 \\ 1 & 1 \end{bmatrix}$

D. $\boldsymbol{\mu}_{y=0} = \begin{bmatrix} -1 \\ 0 \end{bmatrix}, \boldsymbol{\Sigma}_{y=0} = \begin{bmatrix} 1 & -1 \\ 1 & 1 \end{bmatrix}, \quad \boldsymbol{\mu}_{y=1} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \boldsymbol{\Sigma}_{y=1} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$

# Piazza Poll 3

Iris dataset:
- Class distribution: $Y \sim Bern(\phi)$
- Class conditional distribution: $X \sim \mathcal{N}(\boldsymbol{\mu}_y, \boldsymbol{\Sigma}_y)$
- Naïve Bayes assumption?

Which of the following pairs of Gaussian class conditional distributions satisfy the Naïve Bayes assumptions? Select ALL that apply.

A. $\boldsymbol{\mu}_{y=0} = \begin{bmatrix} -1 \\ 0 \end{bmatrix}, \boldsymbol{\Sigma}_{y=0} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, \quad \boldsymbol{\mu}_{y=1} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \boldsymbol{\Sigma}_{y=1} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$

B. $\boldsymbol{\mu}_{y=0} = \begin{bmatrix} -1 \\ 0 \end{bmatrix}, \boldsymbol{\Sigma}_{y=0} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, \quad \boldsymbol{\mu}_{y=1} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \boldsymbol{\Sigma}_{y=1} = \begin{bmatrix} 3 & 0 \\ 0 & 3 \end{bmatrix}$

C. $\boldsymbol{\mu}_{y=0} = \begin{bmatrix} -1 \\ 0 \end{bmatrix}, \boldsymbol{\Sigma}_{y=0} = \begin{bmatrix} 1 & -1 \\ 1 & 1 \end{bmatrix}, \quad \boldsymbol{\mu}_{y=1} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \boldsymbol{\Sigma}_{y=1} = \begin{bmatrix} 1 & -1 \\ 1 & 1 \end{bmatrix}$

D. $\boldsymbol{\mu}_{y=0} = \begin{bmatrix} -1 \\ 0 \end{bmatrix}, \boldsymbol{\Sigma}_{y=0} = \begin{bmatrix} 1 & -1 \\ 1 & 1 \end{bmatrix}, \quad \boldsymbol{\mu}_{y=1} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \boldsymbol{\Sigma}_{y=1} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$

# Decision Boundaries

Iris dataset:

- Class distribution: $Y \sim Bern(\phi)$
- Class conditional distribution: $X \sim \mathcal{N}(\boldsymbol{\mu}_y, \boldsymbol{\Sigma}_y)$

- Naïve Bayes assumption:

- Linear Decision Boundary:

- Quadradic Decision Boundary:

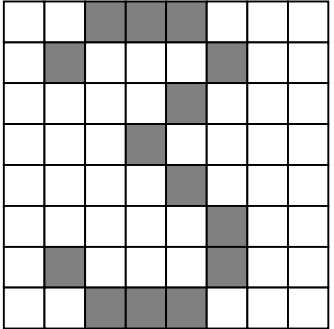# Introduction to Machine Learning

## Intro to Neural Networks

Instructor: Pat Virtue

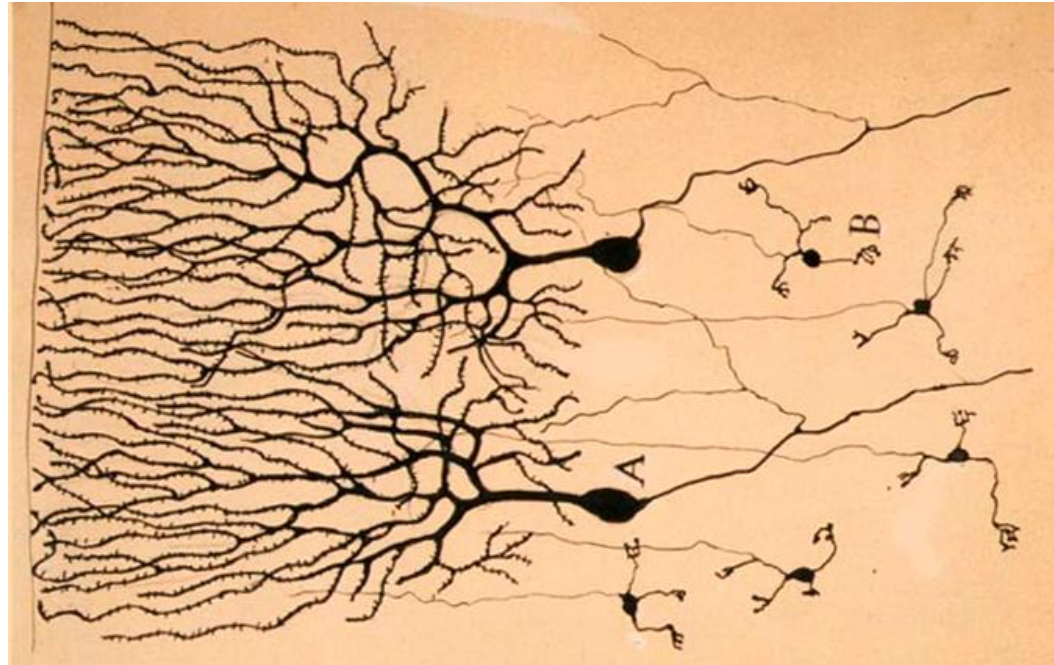# Neural Networks from HW2

1-D Regression

# Neural Networks from HW2

## Digit Classification

# Neural Networks
## Inspired by actual human brain
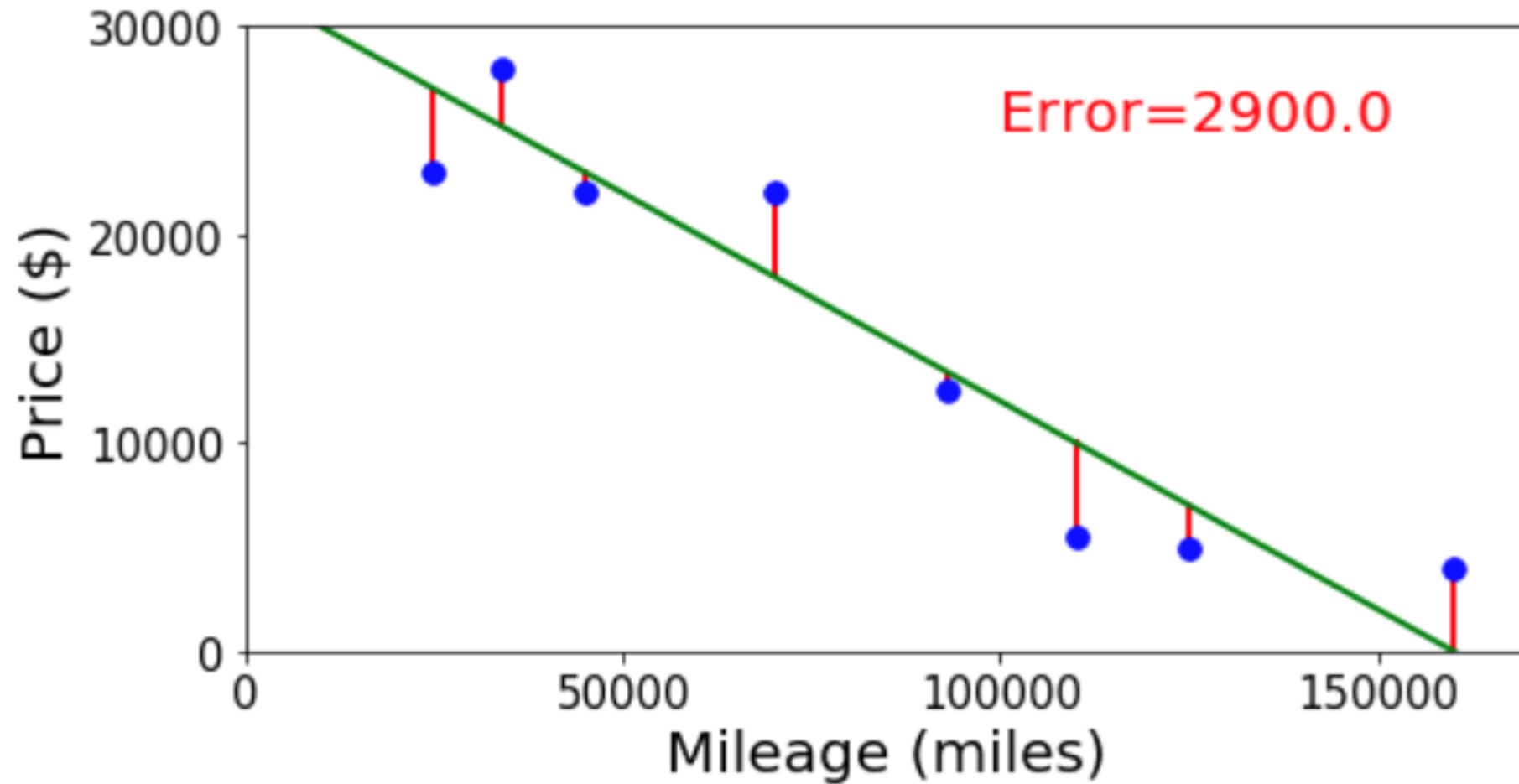
**Input Signal**

**Output Signal**



DOG

**CAT**

TREE

CAR

SKY

Image: https://en.wikipedia.org/wiki/Neuron

# Neural Networks

Simple single neuron example:

- Selling my car

# Neural Networks
## Many layers of neurons, millions of parameters
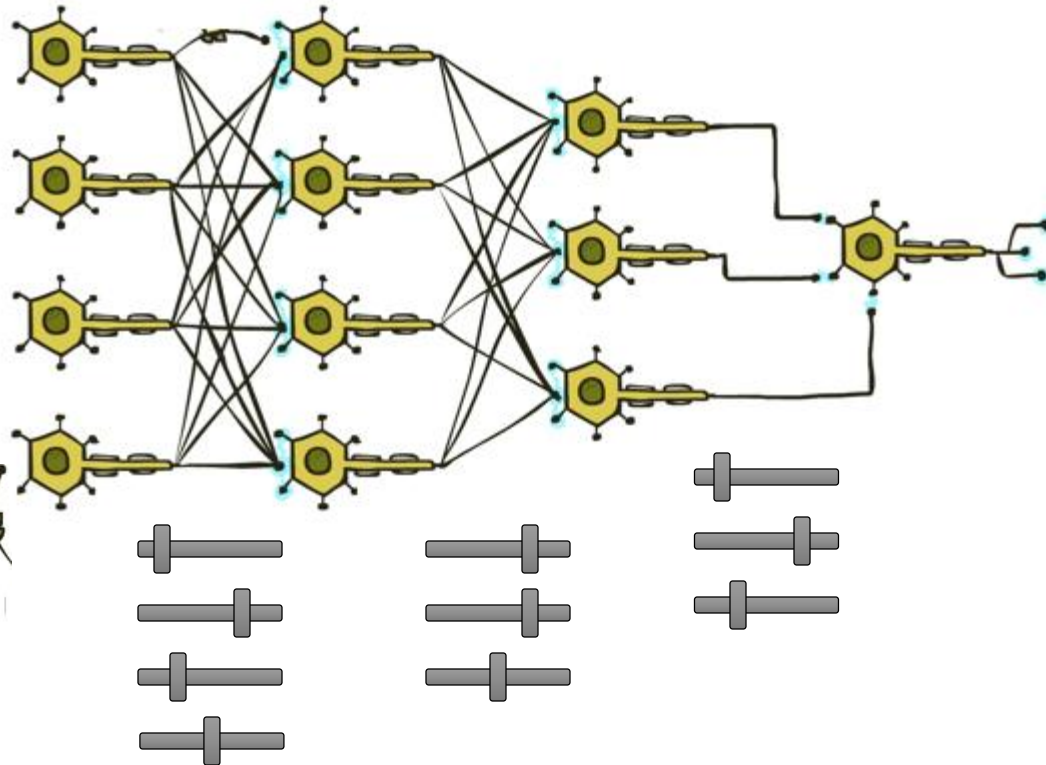
Input
Signal

Output
Signal

DOG

**CAT**

TREE

CAR

SKY

# Neural Networks
Many layers of neurons, millions of parameters
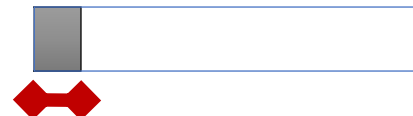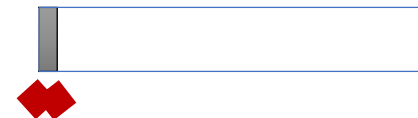
Input
Signal

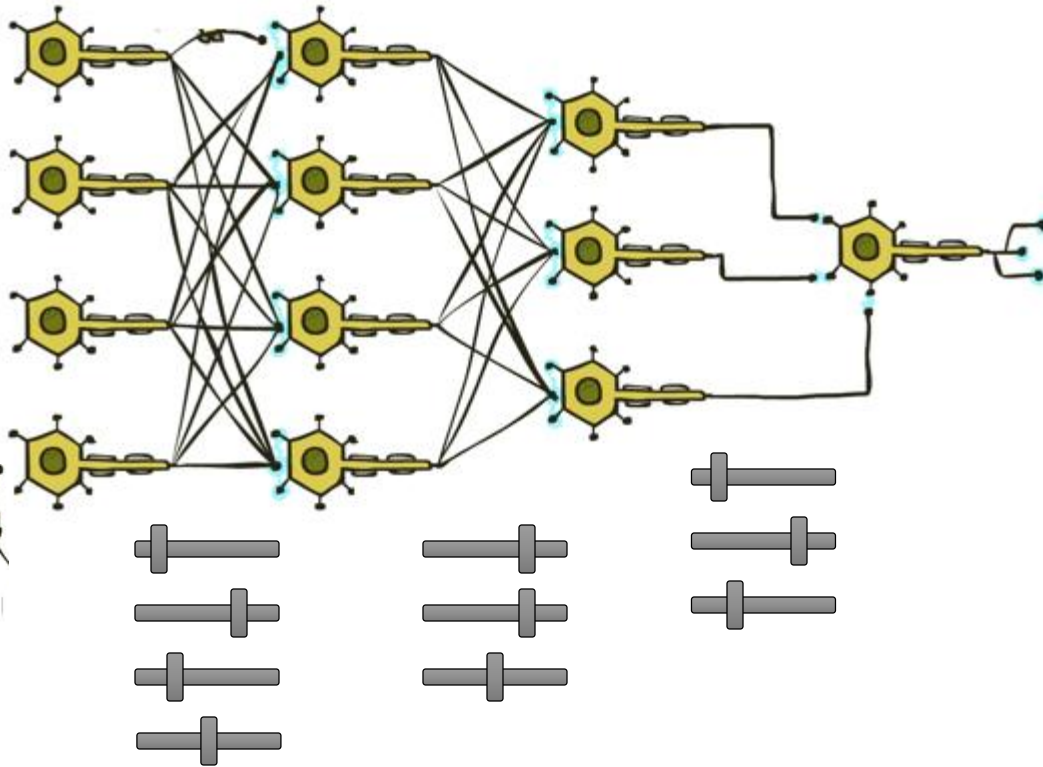Output
Signal



DOG
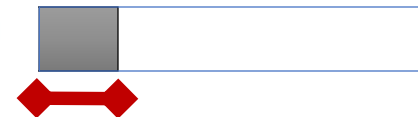
**CAT**

TREE

CAR

SKY

# Neural Networks
## Many layers of neurons, millions of parameters

Input
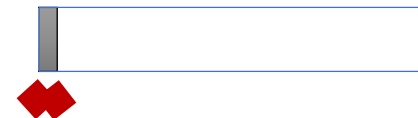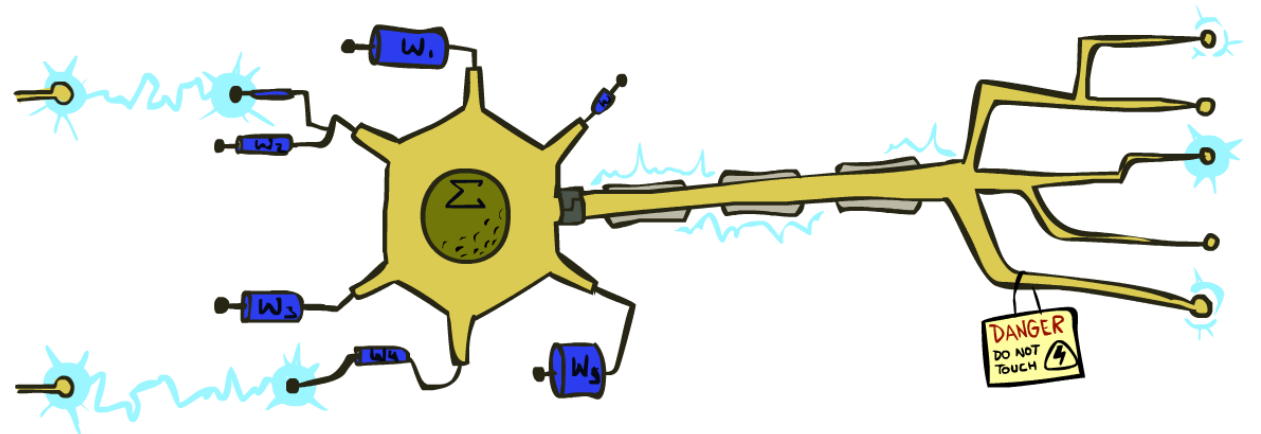Signal

Output
Signal

LEFT

**RIGHT**

UP

DOWN

BUTTON

# Very Loose Inspiration: Human Neurons

# Simple Model of a Neuron (McCulloch & Pitts, 1943)



Inputs $a_i$ come from the output of node i to this node j (or from "outside")

Each input link has a ***weight*** $w_{i,j}$

There is an additional fixed input $a_0$ with ***bias*** weight $w_{0,j}$

The total input is $in_j = \Sigma_i\, w_{i,j}\, a_i$

The output is $a_j = g(in_j) = g(\Sigma_i\, w_{i,j}\, a_i) = g(\mathbf{w.a})$

# Single Neuron

## Single neuron system

- Perceptron (if $g$ is step function)
- Logistic regression (if $g$ is sigmoid)

$x_1$

$w_1$

Computed Value

True Label

$\Sigma$ $g$

$z_1$

$y$

$x_2$

$w_2$

$h_{\boldsymbol{w}}(\boldsymbol{x}) = z_1$

$$h_{\boldsymbol{w}}(\boldsymbol{x}) = g\left(\sum_i w_i x_i\right)$$

# Optimizing

How do we find the "best" set of weights?

$$h_w(\boldsymbol{x}) = g\left(\sum_i w_i x_i\right)$$

# Multilayer Perceptrons

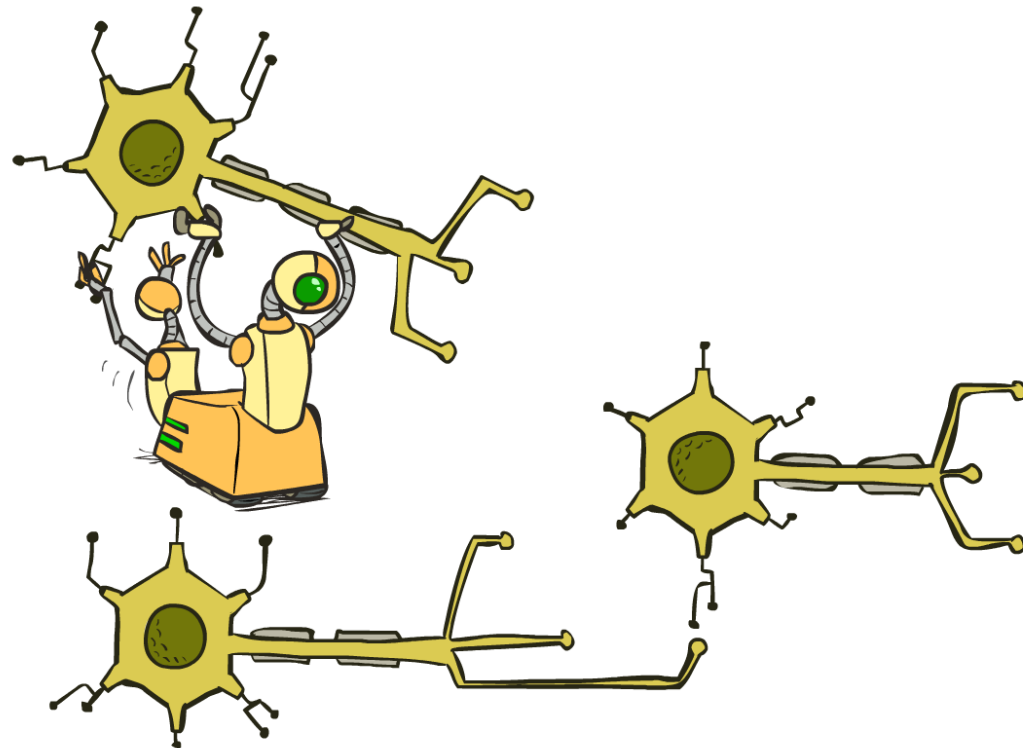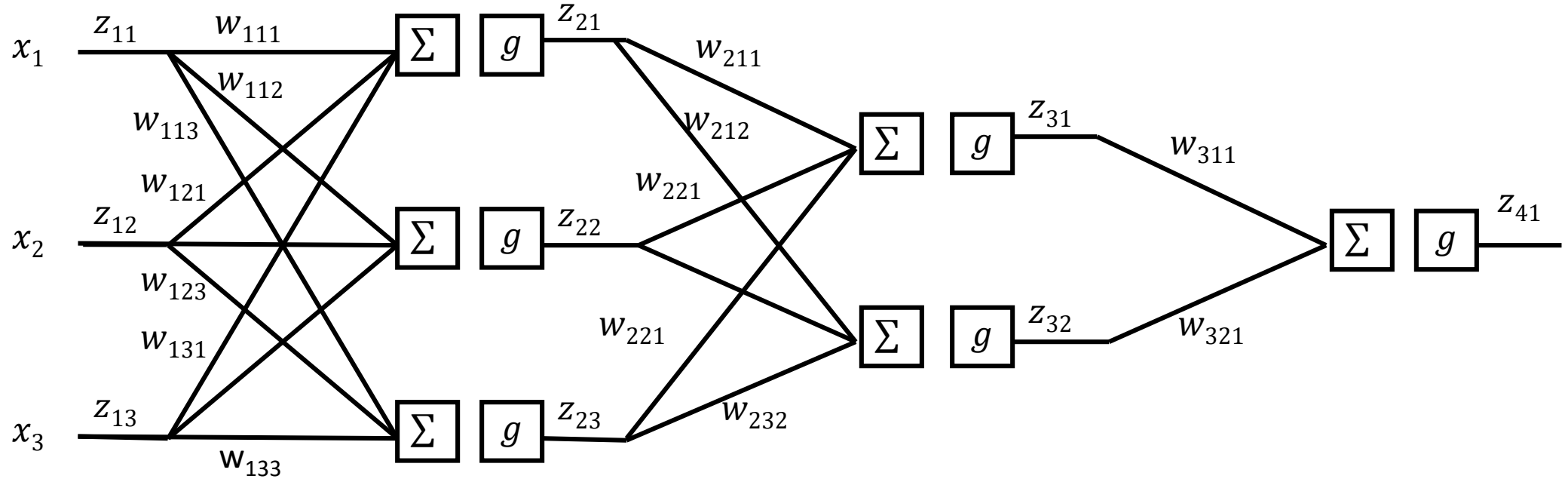A ***multilayer perceptron*** is a feedforward neural network with at least one ***hidden layer*** (nodes that are neither inputs nor outputs)

MLPs with enough hidden nodes can represent any function

# Neural Network Equations

$x_1$ $z_{11}$ $w_{111}$ $\Sigma$ $g$ $z_{21}$ $w_{211}$

$w_{112}$

$w_{113}$ $w_{212}$

$w_{121}$ $z_{31}$ $w_{311}$

$x_2$ $z_{12}$ $\Sigma$ $g$ $z_{22}$ $w_{221}$ $\Sigma$ $g$

$z_{41}$

$w_{123}$ $\Sigma$ $g$

$w_{131}$ $w_{221}$ $z_{32}$ $w_{321}$

$x_3$ $z_{13}$ $\Sigma$ $g$ $z_{23}$ $w_{232}$

$w_{133}$

$$h_w(\boldsymbol{x}) = z_{4,1}$$

$$z_{1,1} = x_1$$

$$z_{4,1} = g\left(\sum_i w_{3,i,1}\, z_{3,i}\right)$$

$$z_{3,1} = g\left(\sum_i w_{2,i,1}\, z_{2,i}\right)$$

$$h_w(x) = g\left(\sum_k w_{3,k,1}\; g\left(\sum_j w_{2,j,k}\; g\left(\sum_i w_{1,i,j}\; x_i\right)\right)\right)$$

$$z_{d,1} = g\left(\sum_i w_{d-1,i,1}\, z_{d-1,i}\right)$$

# Optimizing

How do we find the "best" set of weights?

$$h_w(x) = g\left(\sum_k w_{3,k,1}\ g\left(\sum_j w_{2,j,k}\ g\left(\sum_i w_{1,i,j}\ x_i\right)\right)\right)$$

# Neural Networks Properties

- Large number of neurons
    - Danger for overfitting
- Modelling assumptions vs data assumptions trade-off

- Gradient descent can get stuck in bad local optima

What if there are not non-linear activations?

- A deep neural network with only linear layers can be reduced to an exactly equivalent single linear layer

Universal Approximation Theorem:

- A two-layer neural network with a sufficient number of neurons can approximate any continuous function to any desired accuracy.