



# LECTURE 2:

# REGRESSION & CLASSIFICATION

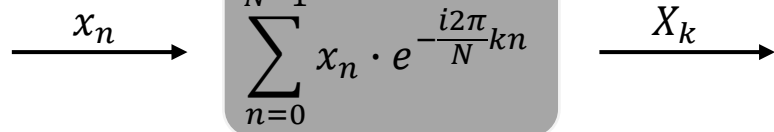
University of Washington, Seattle

Fall 2024

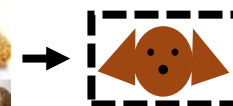
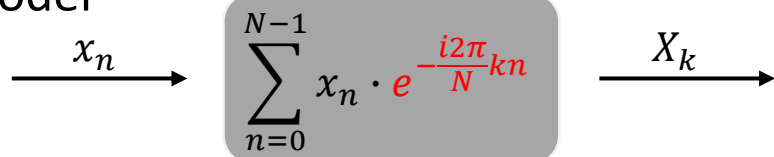


# Previously in EEP 596...

Fixed model



Learned model



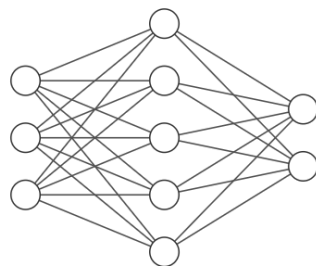
Model



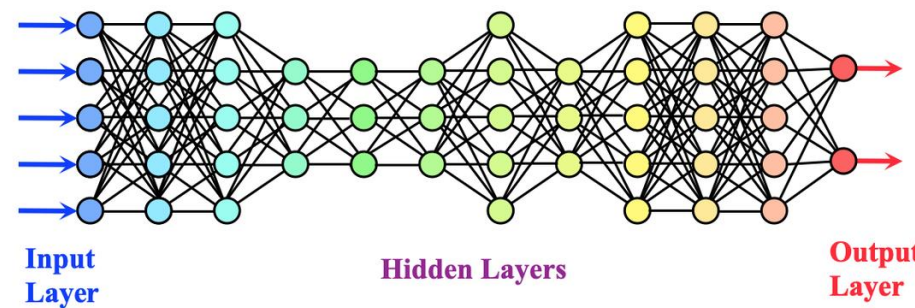
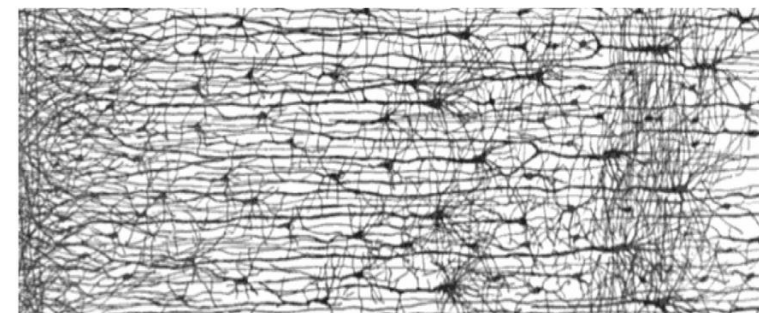
1: Café Poodle  
0: Not Café Poodle



# Previously in EEP 596...



1: Café Poodle  
0: Not Café Poodle





# OUTLINE

## **Part 1: Regression**

- Regression problem
- Linear regression
- Least square formula

## **Part 2: Binary Classification**

- Binary Logistic Regression
- Linear vs Logistic regression
- Logistic regression and Neuron

## **Part 3: Training and Optimization**

- Binary Cross Entropy Loss function
- Training Logistic Regression
- Gradient descent
- Composite Loss function and Chain rule
- Back-propagation algorithm
- Training terminologies



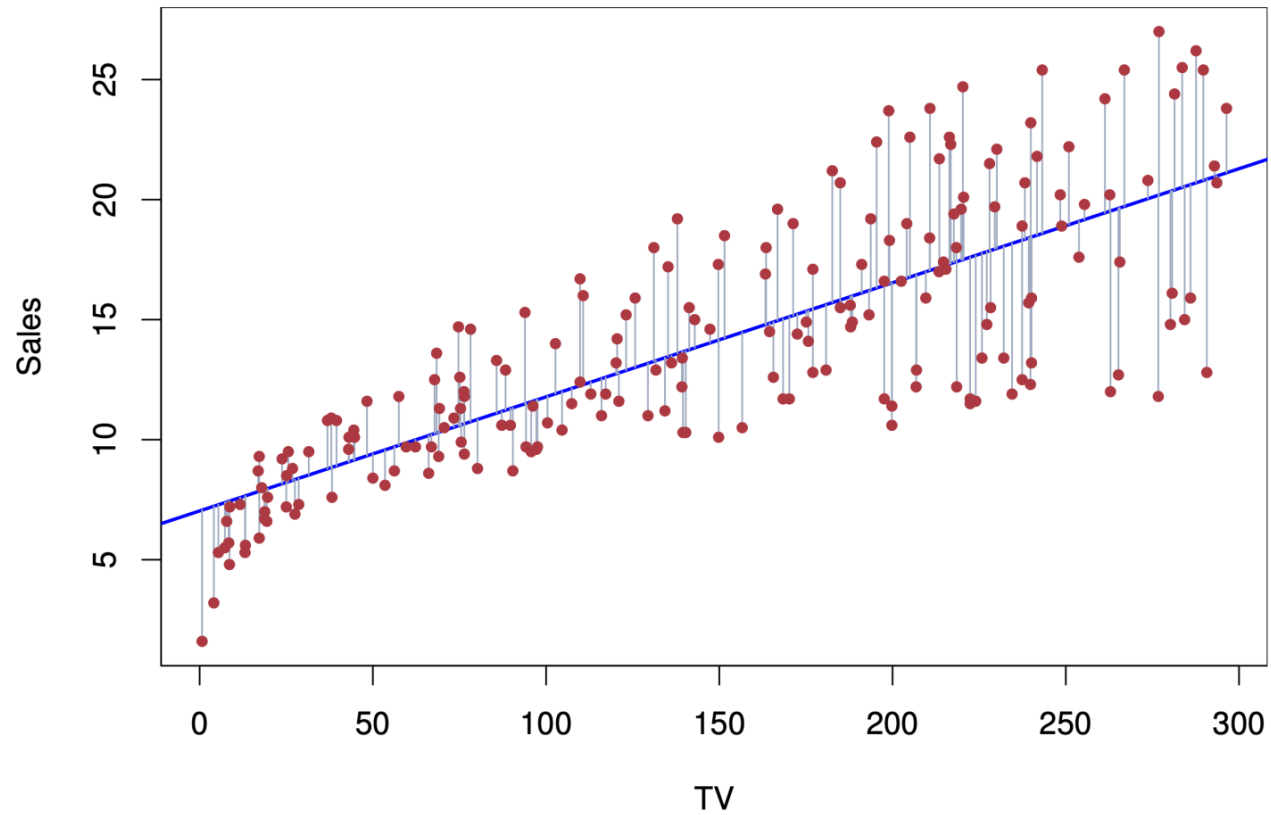
# PART 1:

# Regression



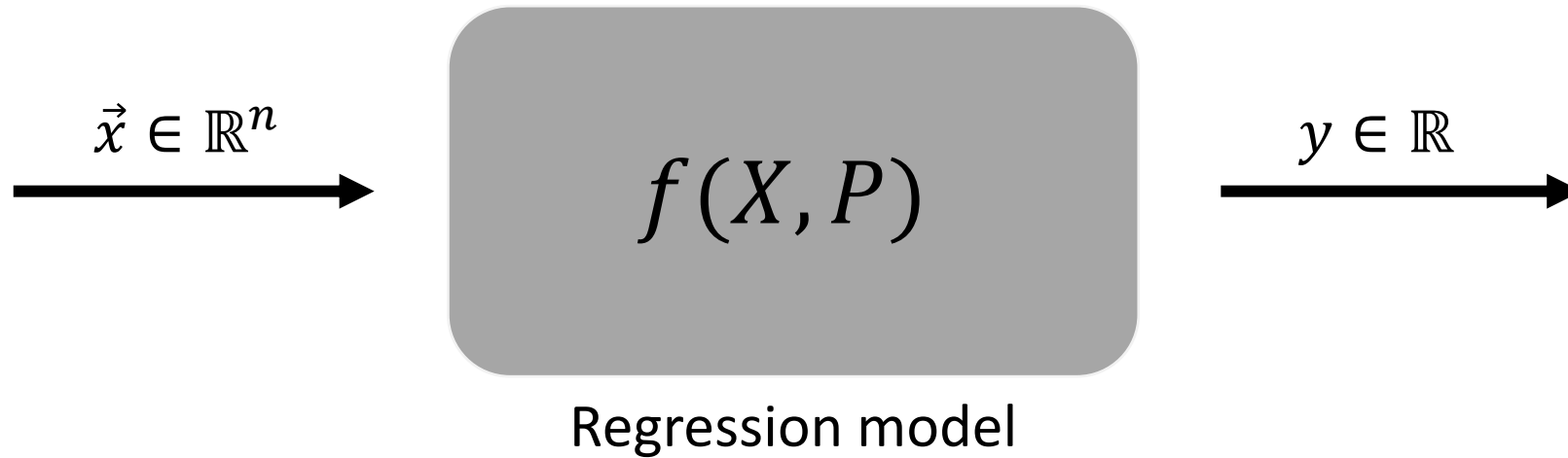
# Regression problem

$$Y = f(X, W)$$





# Regression problem



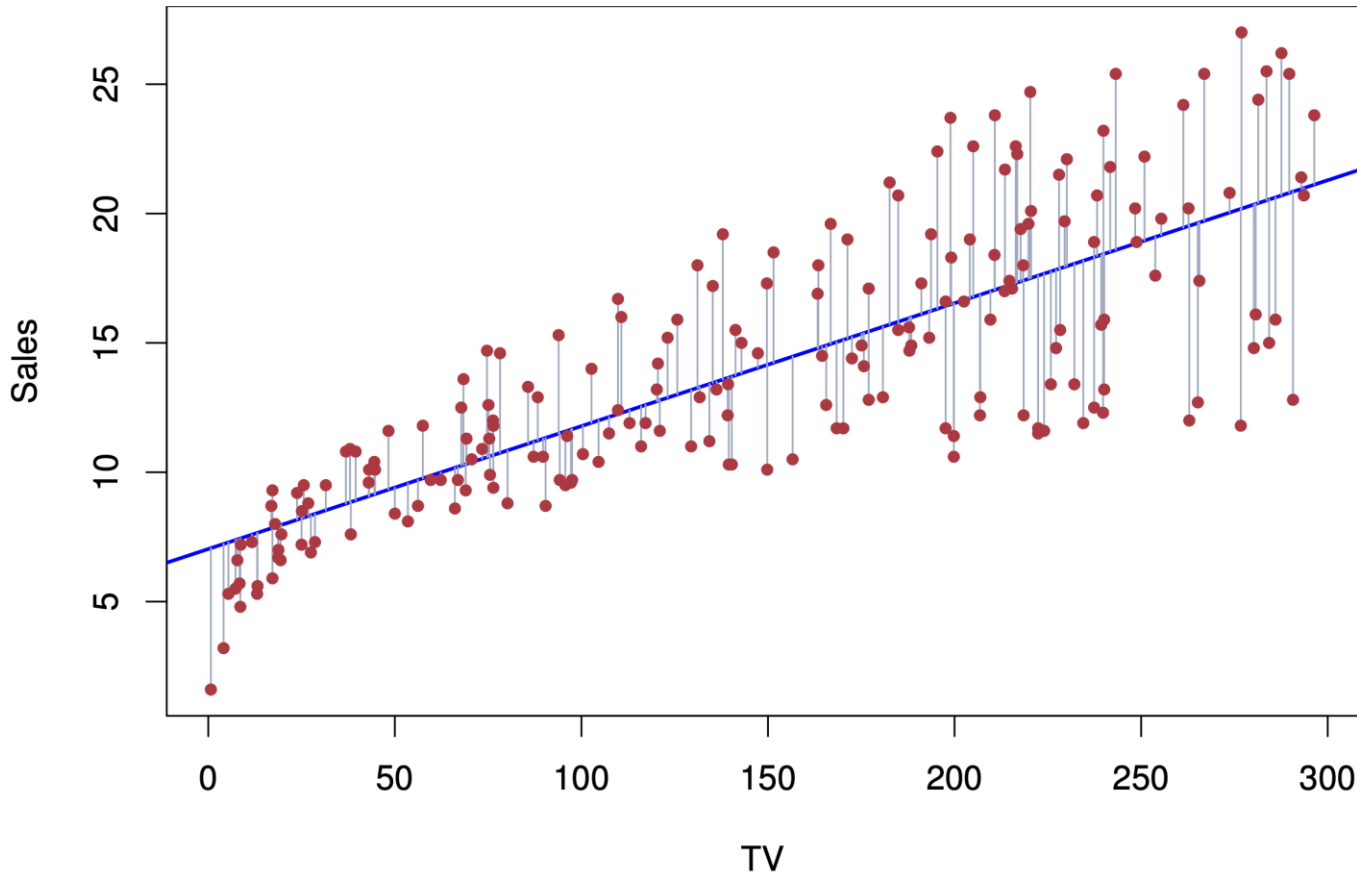
Linear regression (polynomial fit)

$$y(x) = p_0 + p_1x + \cdots + p_mx^m$$

$$y = \vec{p} \cdot \vec{x}$$



# Polynomial fit – Linear regression



$x = \text{input data}$

$y = \text{output data}$

$$\hat{y}(x) = p_0 + p_1x + \cdots + p_mx^m$$

Goal: Minimize  $e = y - \hat{y}$





# Polynomial fit – Linear regression

$$\hat{y}(x) = p_0 + p_1x + \cdots + p_mx^m$$

Goal: Minimize  $e = y - \hat{y}$



# Polynomial fit – Linear regression

$$\hat{y}(x) = p_0 + p_1x + \cdots + p_mx^m$$

Goal: Minimize  $e = y - \hat{y}$

$$(x_0, y_0), (x_1, y_1), \dots, (x_n, y_n)$$

n-given points

$$X = \begin{pmatrix} 1 & x_0 & x_0^2 & \cdots & x_0^m \\ 1 & x_1 & x_1^2 & \cdots & x_1^m \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & x_n & x_n^2 & \cdots & x_n^m \end{pmatrix} \begin{matrix} \vec{\tilde{x}}_0 \\ \\ \\ \vec{\tilde{x}}_n \end{matrix}$$

$$P = \begin{pmatrix} p_0 \\ p_1 \\ \vdots \\ p_m \end{pmatrix}$$



# Polynomial fit – Linear regression

$$\hat{y}(x) = p_0 + p_1x + \cdots + p_mx^m$$

Goal: Minimize  $e = y - \hat{y}$

$(x_0, y_0), (x_1, y_1), \dots, (x_n, y_n)$   
n-given points

$$X = \begin{pmatrix} 1 & x_0 & x_0^2 & \cdots & x_0^m \\ 1 & x_1 & x_1^2 & \cdots & x_1^m \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_n & x_n^2 & \cdots & x_n^m \end{pmatrix} \begin{matrix} \vec{\tilde{x}}_0 \\ \\ \\ \vec{\tilde{x}}_n \end{matrix} \quad P = \begin{pmatrix} p_0 \\ p_1 \\ \vdots \\ p_n \end{pmatrix}$$

Fit (“Learn”) the model:

Least squares error

**Cost**  $J = E_2 = \sum_{i=1}^n (\overbrace{\vec{p} \cdot \vec{\tilde{x}}_i}^{\hat{y}_i} - y_i)^2$   
**Total error**



# Polynomial fit – Linear regression

Find the parameters minimizing the error

$$J = E_2 = \sum_{i=1}^n (\vec{p} \cdot \vec{\tilde{x}}_i - y_i)^2$$

$$\vec{p}^* = \arg \min_{\vec{p}} J(\vec{p})$$



# Polynomial fit – Linear regression

Find the parameters minimizing the error

$$J = E_2 = \sum_{i=1}^n (\vec{p} \cdot \vec{\tilde{x}}_i - y_i)^2$$

$$\vec{p}^* = \arg \min_{\vec{p}} J(\vec{p})$$

$$\forall i : \quad \frac{\partial J}{\partial p_j} = 0 \quad \frac{\partial J}{\partial p_j} = \sum_{i=1}^n 2(\vec{p} \cdot \vec{\tilde{x}}_i - y_i) \vec{\tilde{x}}_i$$

$$\vec{p}^* = (X^T X)^{-1} X^T \vec{y}$$



# Polynomial fit – Linear regression

Find the parameters minimizing the error

$$J = E_2 = \sum_{i=1}^n (\vec{p} \cdot \vec{\tilde{x}}_i - y_i)^2$$

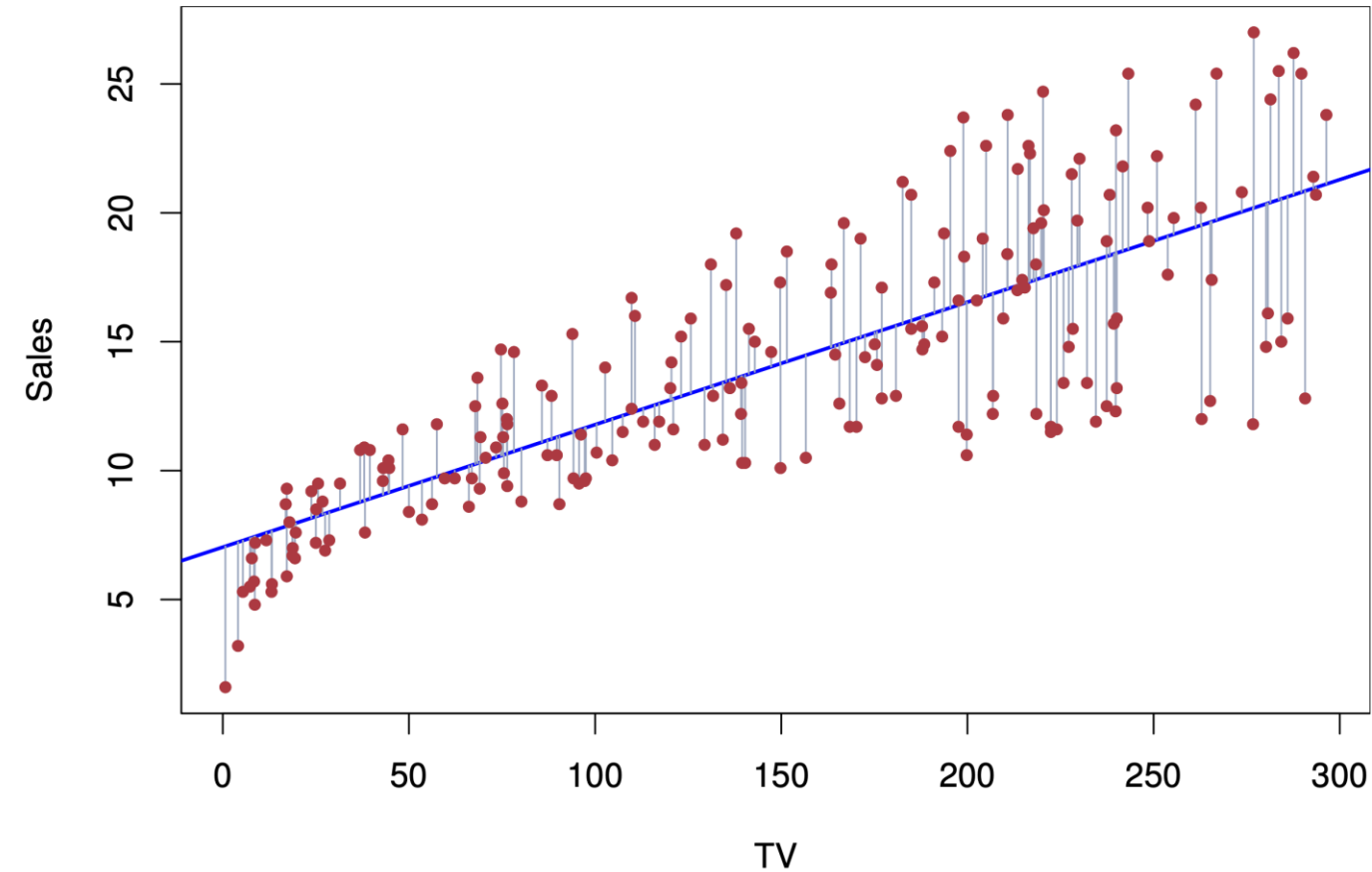
$$\vec{p}^* = \arg \min_{\vec{p}} J(\vec{p})$$

$$\forall i : \quad \frac{\partial J}{\partial p_j} = 0 \qquad \frac{\partial J}{\partial p_j} = \sum_{i=1}^n 2(\vec{p} \cdot \vec{\tilde{x}}_i - y_i) \vec{\tilde{x}}_i$$

$$\vec{p}^* = (X^T X)^{-1} X^T \vec{y} \quad \text{Linear least square formula}$$



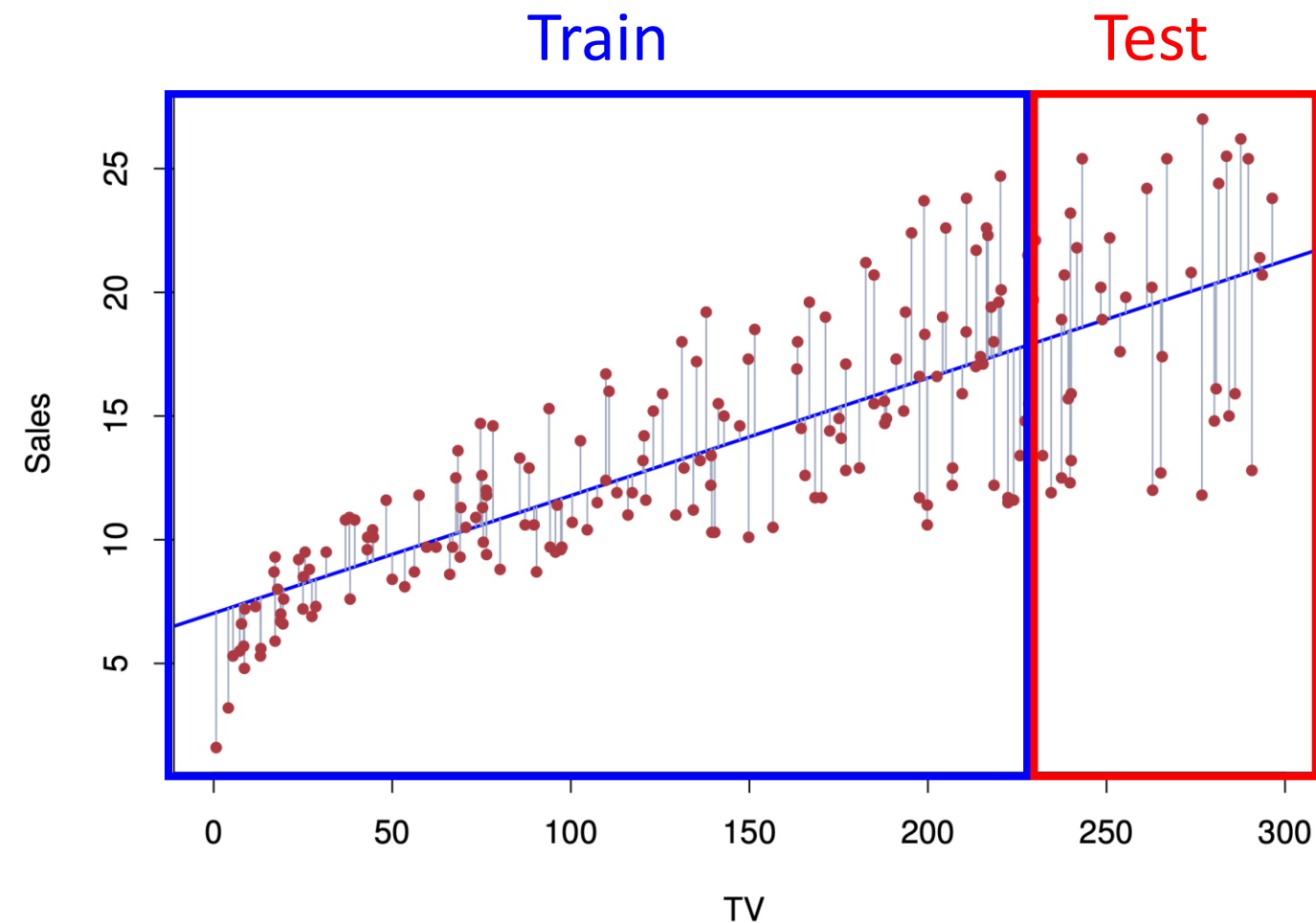
# Linear regression is not perfect



$$\hat{y}(x) = p_0 + p_1x + \cdots + p_mx^m$$



# Linear regression is not perfect



$$\hat{y}(x) = p_0 + p_1x + \cdots + p_mx^m$$

Model might perform worse  
in testing in the presence of  
outliers





# Improving Linear Regression

$$W^* = \arg \min_W L(X, W)$$

$$J_2 = \frac{1}{n} \sum_{i=1}^n (f(\tilde{x}_i, W) - y_i)^2 \quad \text{Mean squared error}$$

$$J_1 = \frac{1}{n} \sum_{i=1}^n |f(\tilde{x}_i, W) - y_i| \quad \text{Mean absolute error}$$

$$J_\infty = \max_{1 \leq i \leq n} |f(\tilde{x}_i, W) - y_i| \quad \text{Max error}$$



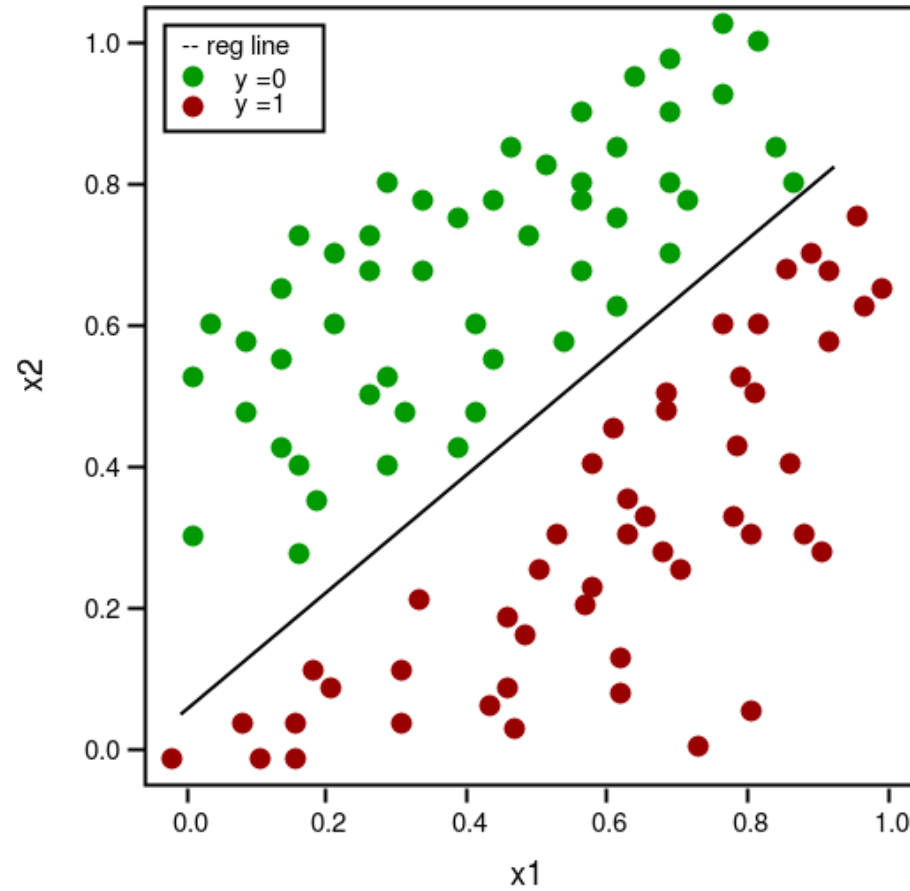
# PART 2:

# Classification



# Binary classification problem

$$y = \sigma(\vec{w}^T \vec{x} + b)$$





# Iris Flower Data

**iris setosa**



petal sepal

**iris versicolor**



petal sepal

**iris virginica**



petal sepal

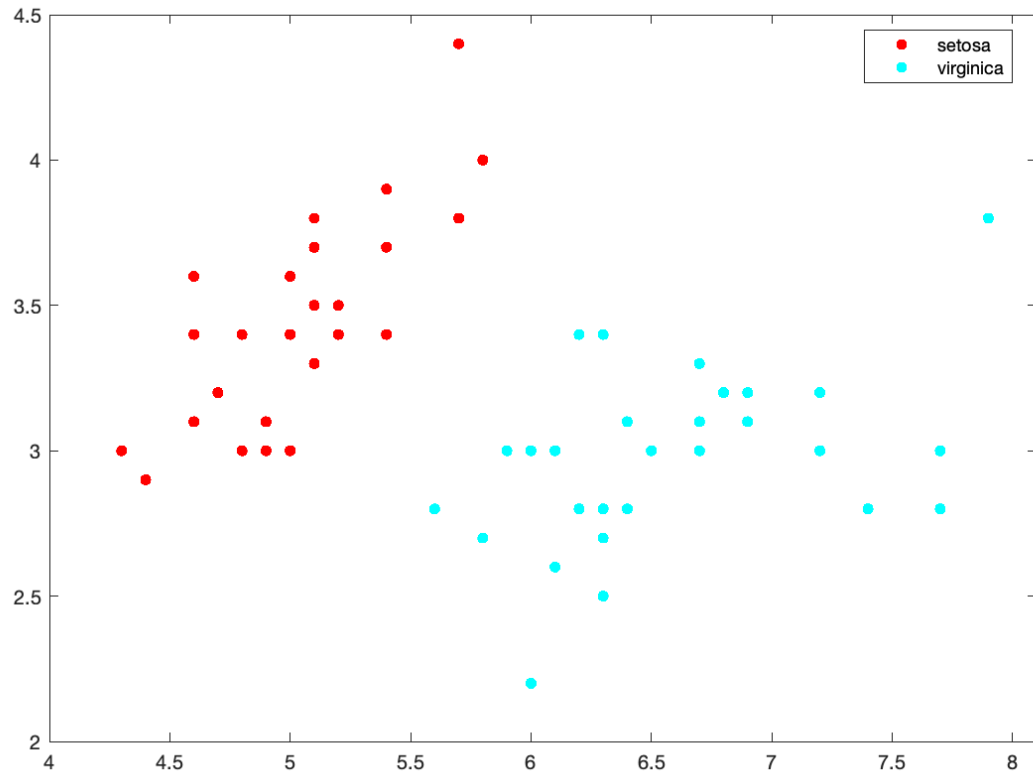
	<code>Id</code>	<code>SepalLengthCm</code>	<code>SepalWidthCm</code>	<code>PetalLengthCm</code>	<code>PetalWidthCm</code>	<code>Species</code>
0	1	5.1	3.5	1.4	0.2	Iris-setosa
1	2	4.9	3.0	1.4	0.2	Iris-setosa
2	3	4.7	3.2	1.3	0.2	Iris-setosa
3	4	4.6	3.1	1.5	0.2	Iris-setosa
4	5	5.0	3.6	1.4	0.2	Iris-setosa
5	6	5.4	3.9	1.7	0.4	Iris-setosa
6	7	4.6	3.4	1.4	0.3	Iris-setosa
7	8	5.0	3.4	1.5	0.2	Iris-setosa
8	9	4.4	2.9	1.4	0.2	Iris-setosa
9	10	4.9	3.1	1.5	0.1	Iris-setosa

Features

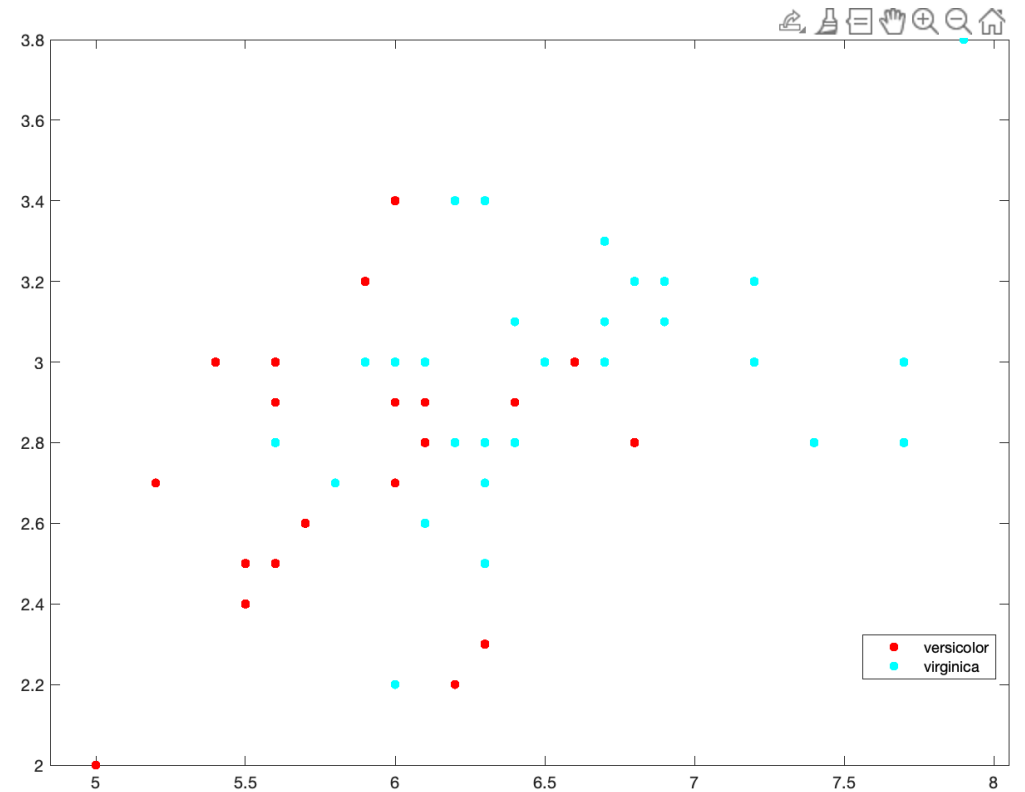
Labels (Targets)



# Linear vs Logistic regression



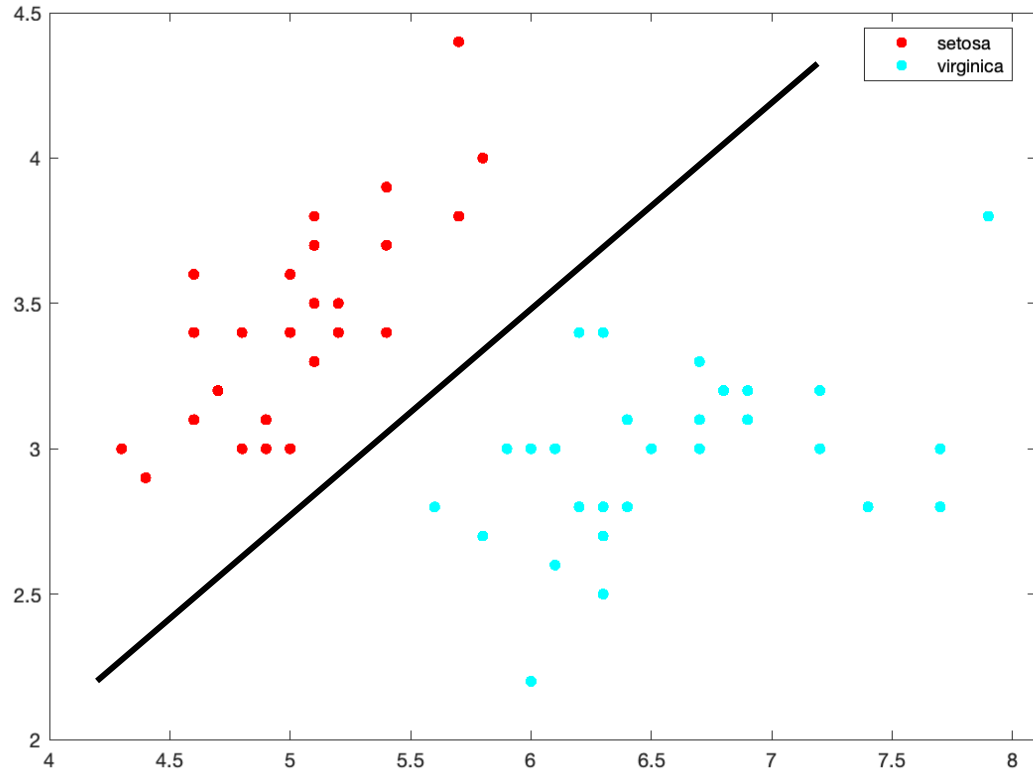
Setosa vs Virginica



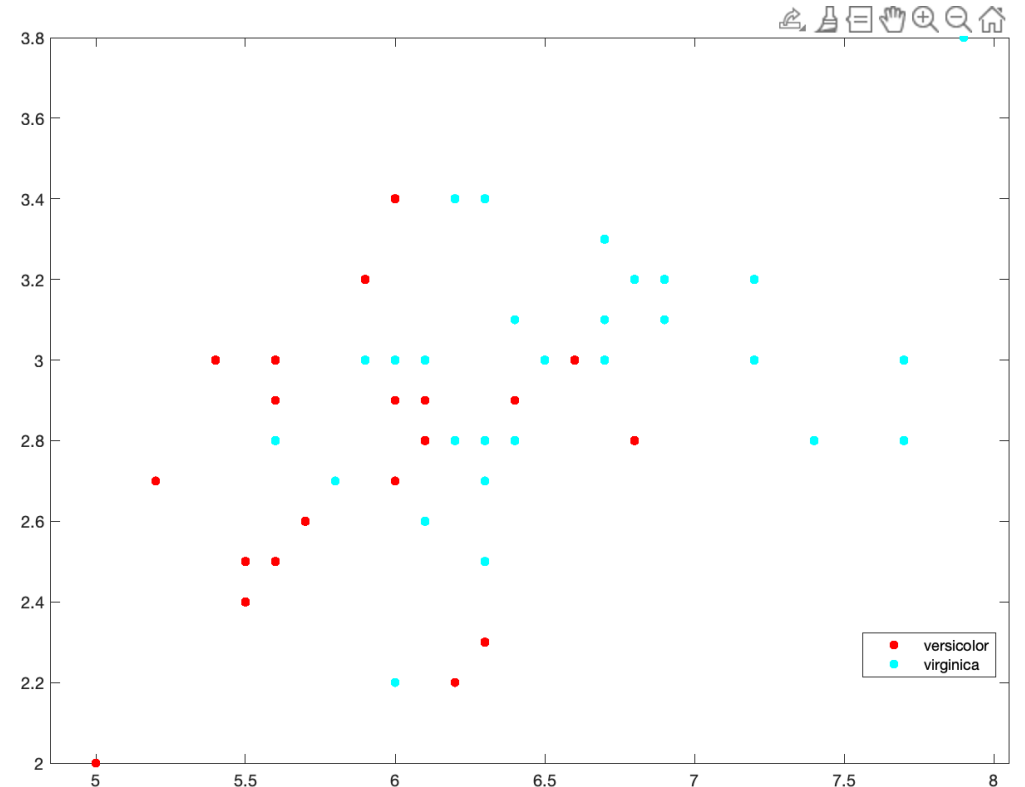
Versicolor vs Virginica



# Linear vs Logistic regression



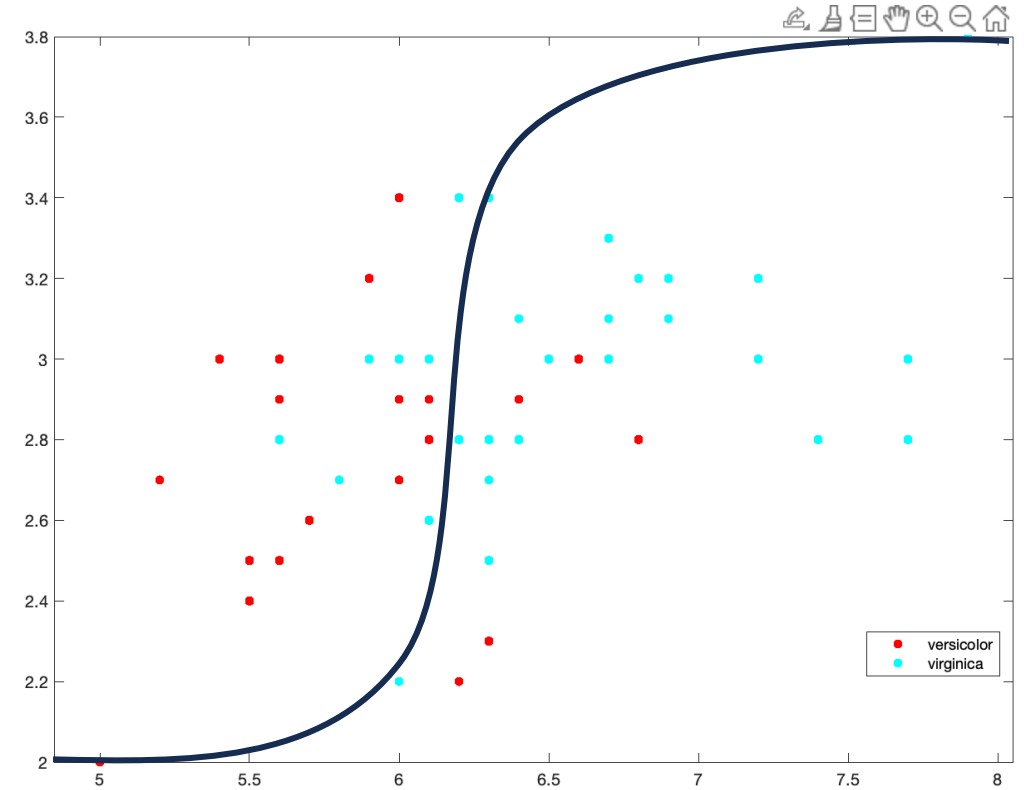
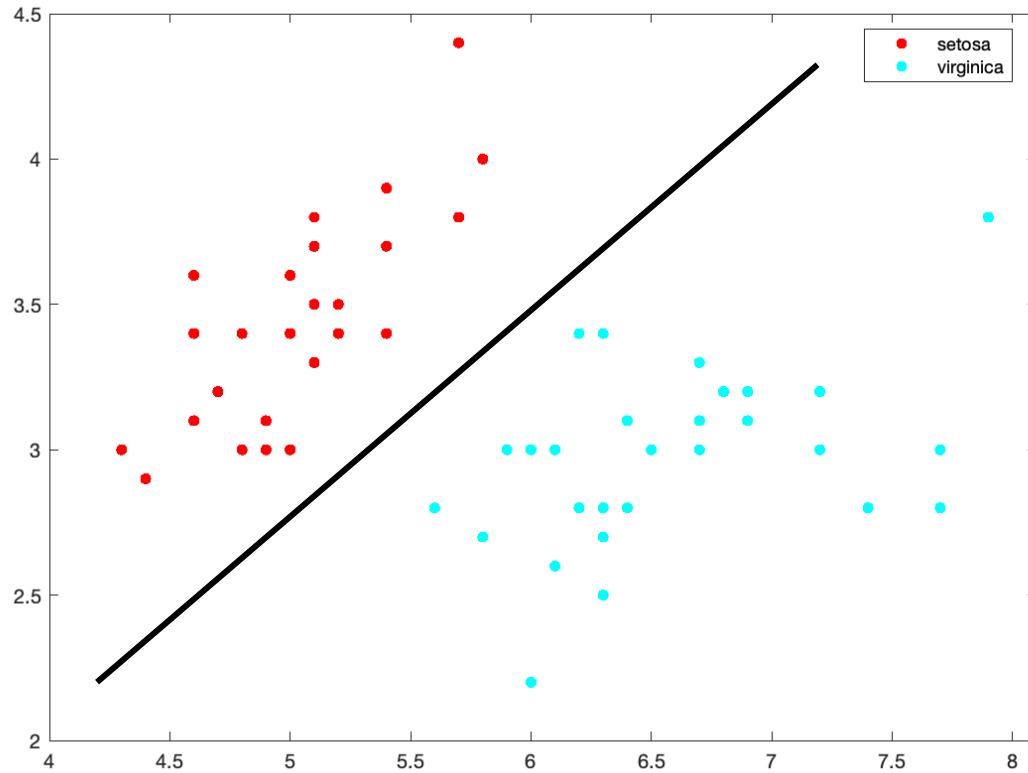
Setosa vs Virginica



Versicolor vs Virginica



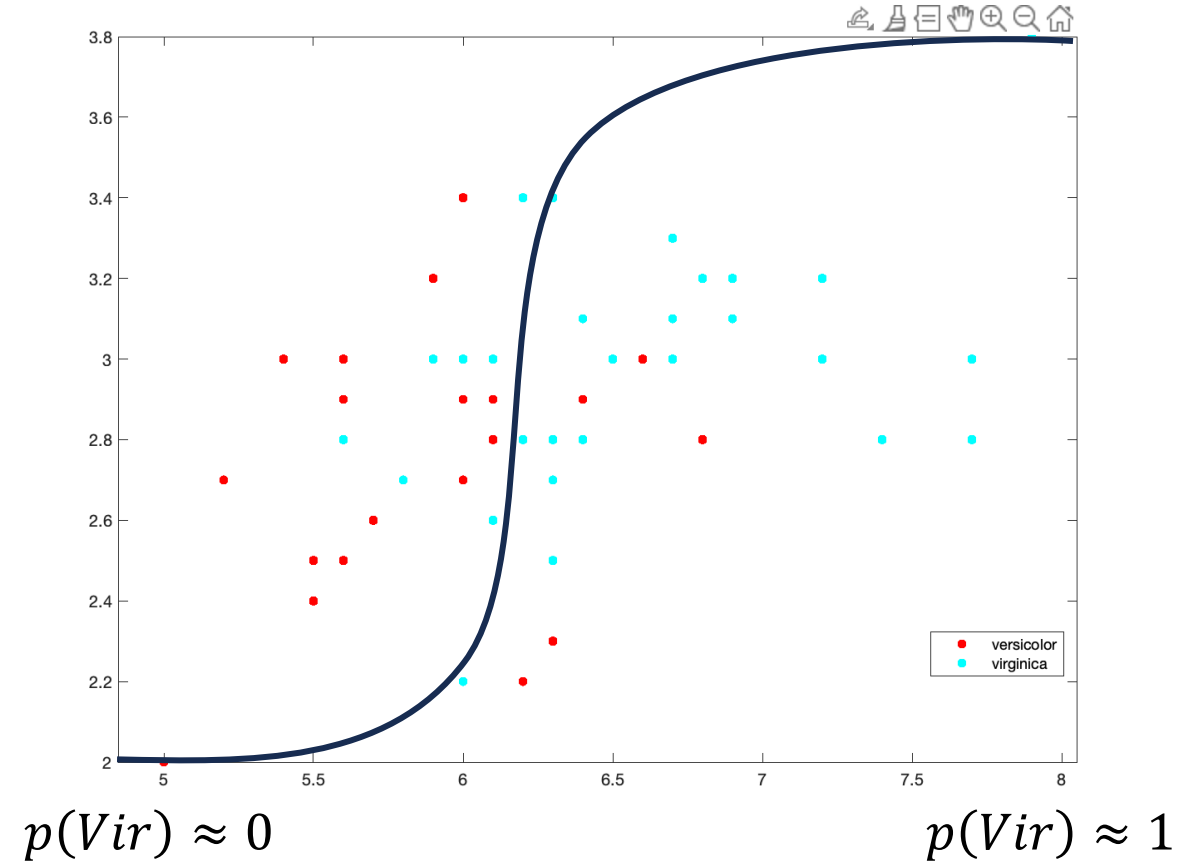
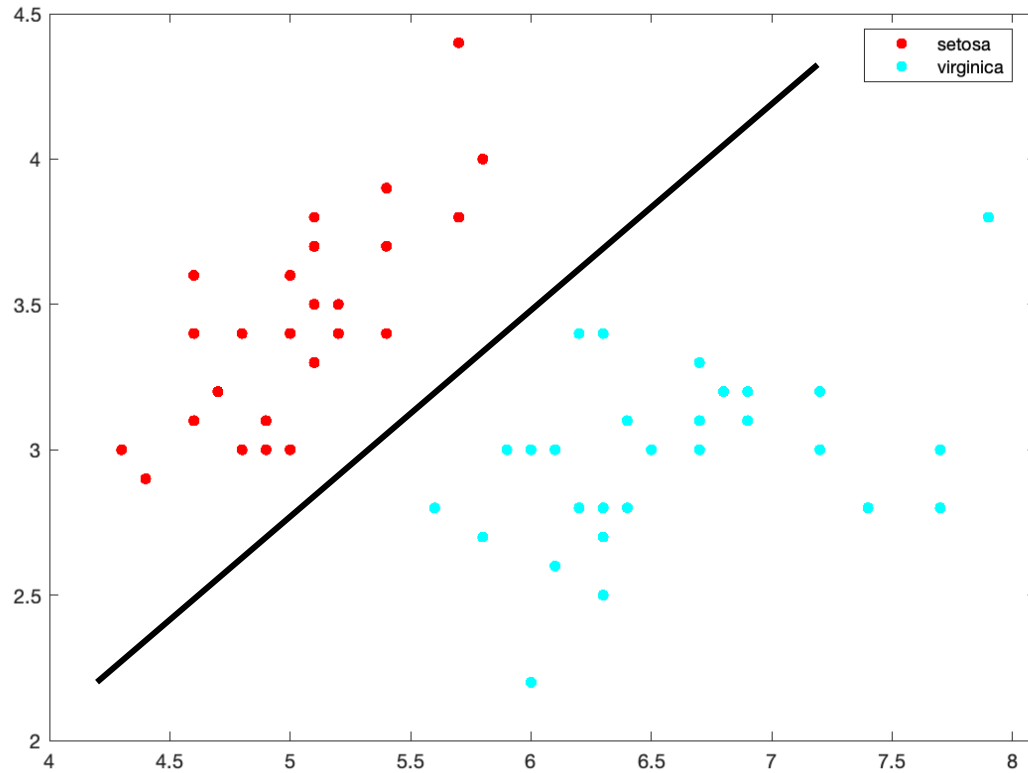
# Linear vs Logistic regression



$$p = \frac{1}{1 + e^{-x}} \quad 0 \leq p \leq 1$$



# Linear vs Logistic regression

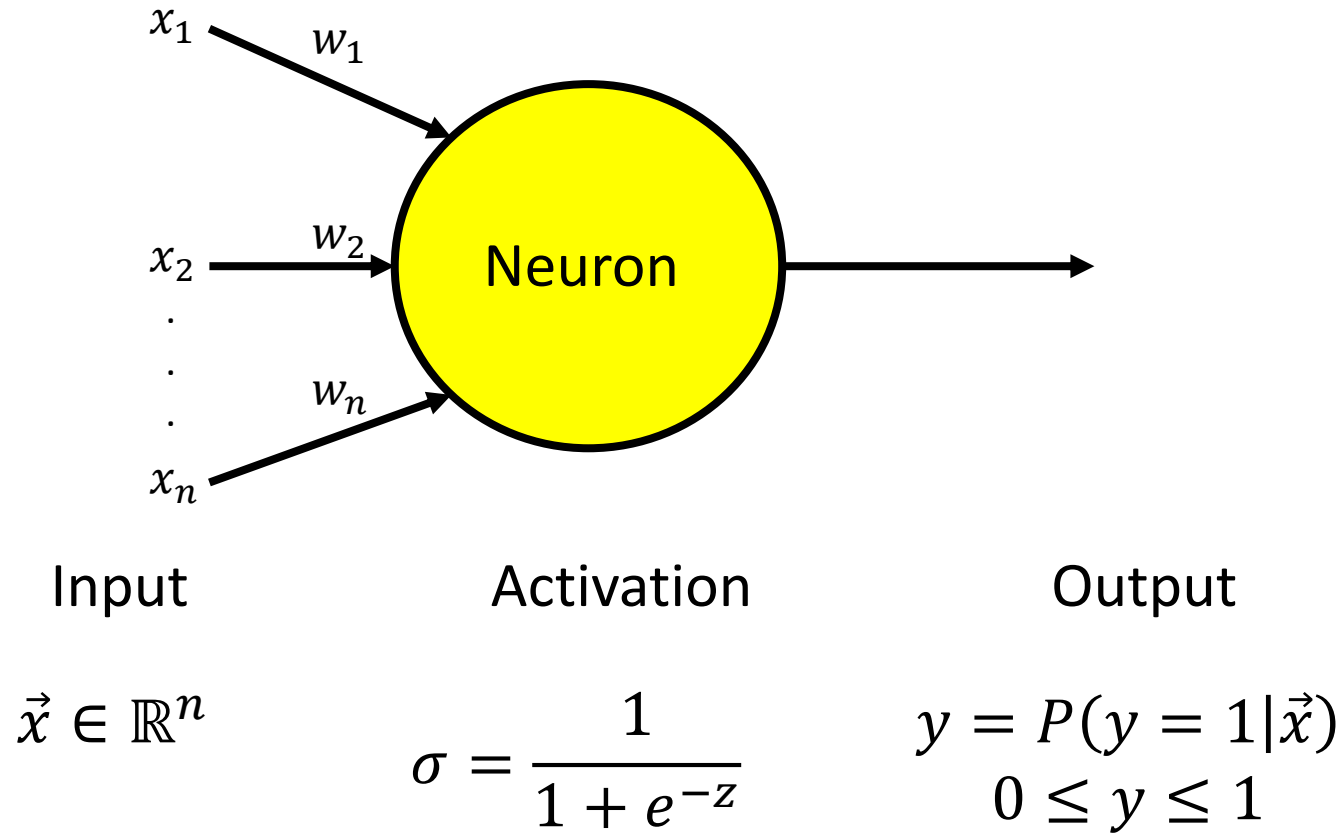


$$p = \frac{1}{1 + e^{-x}} \quad 0 \leq p \leq 1$$



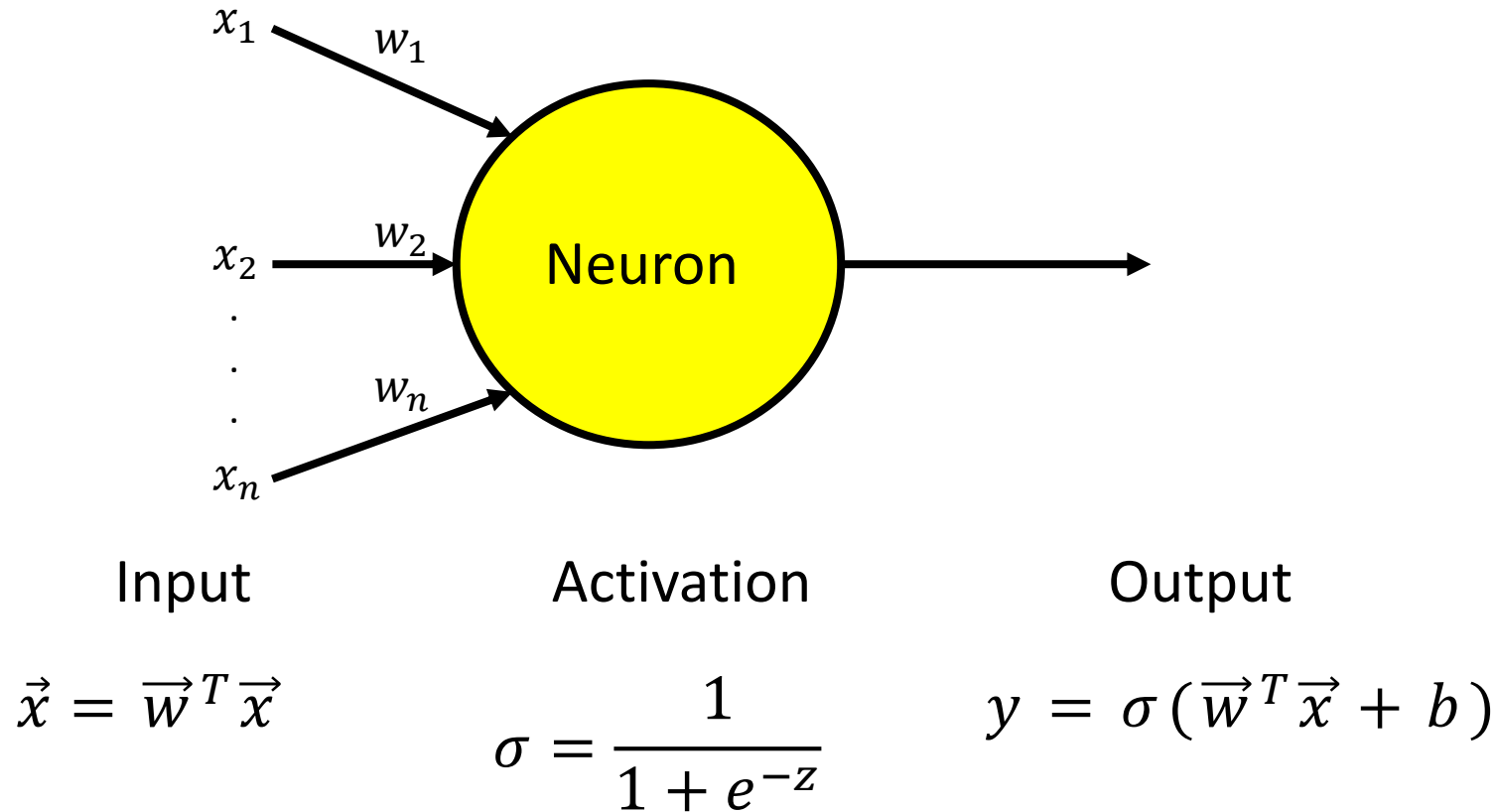


# Artificial Neuron as Logistic Regression Model





# Artificial Neuron as Logistic Regression Model



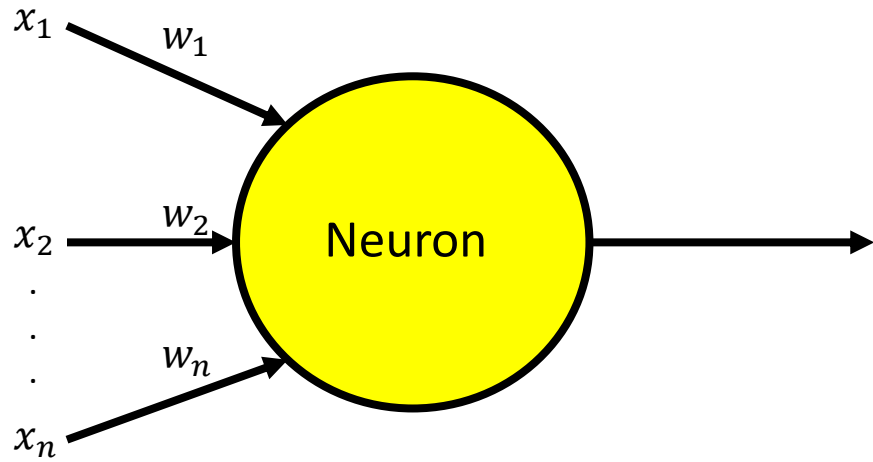


# PART 3:

## Training and Optimization



# Cross Entropy Loss function (Binary)



Input

Activation

Output

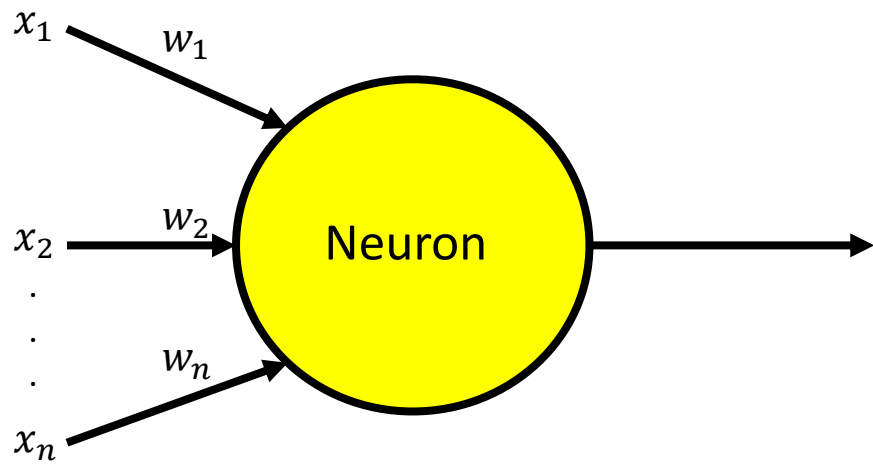
$$\vec{x} = \vec{w}^T \vec{x}$$

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

$$y = \sigma(\vec{w}^T \vec{x} + b)$$



# Cross Entropy Loss function (Binary)



Input

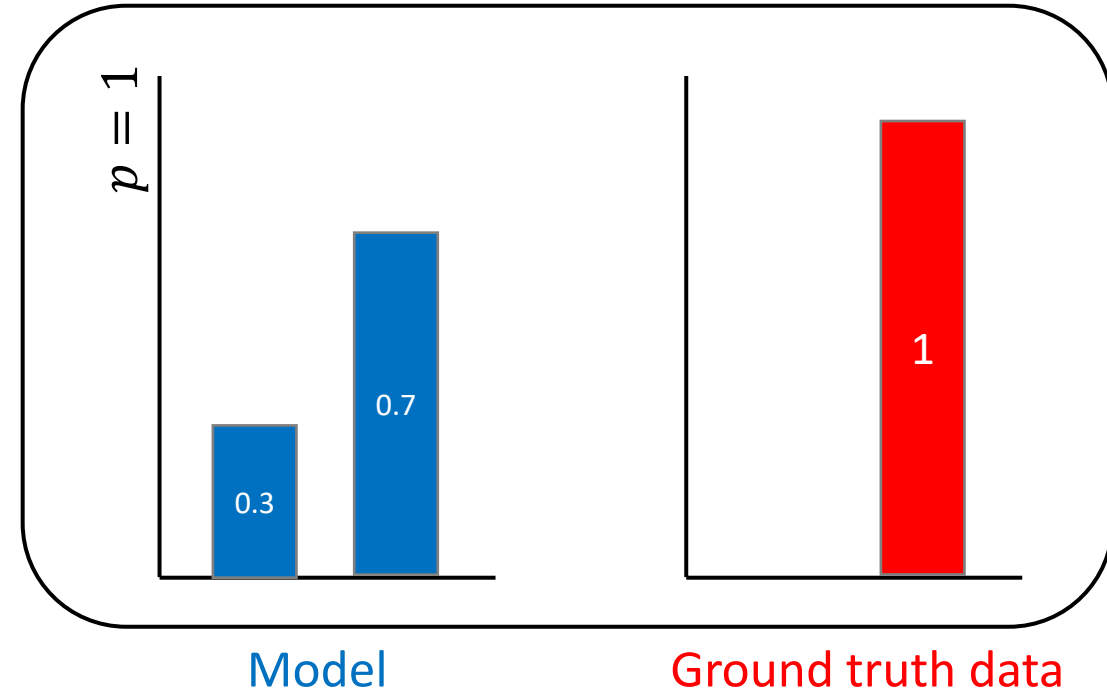
Activation

Output

$$\vec{x} = \vec{w}^T \vec{x}$$

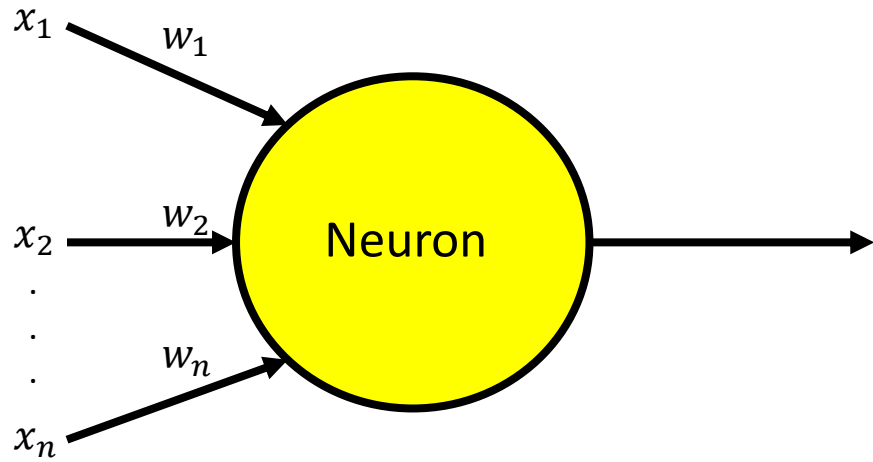
$$\sigma = \frac{1}{1 + e^{-z}}$$

$$y = \sigma(\vec{w}^T \vec{x} + b)$$





# Cross Entropy Loss function (Binary)



Input

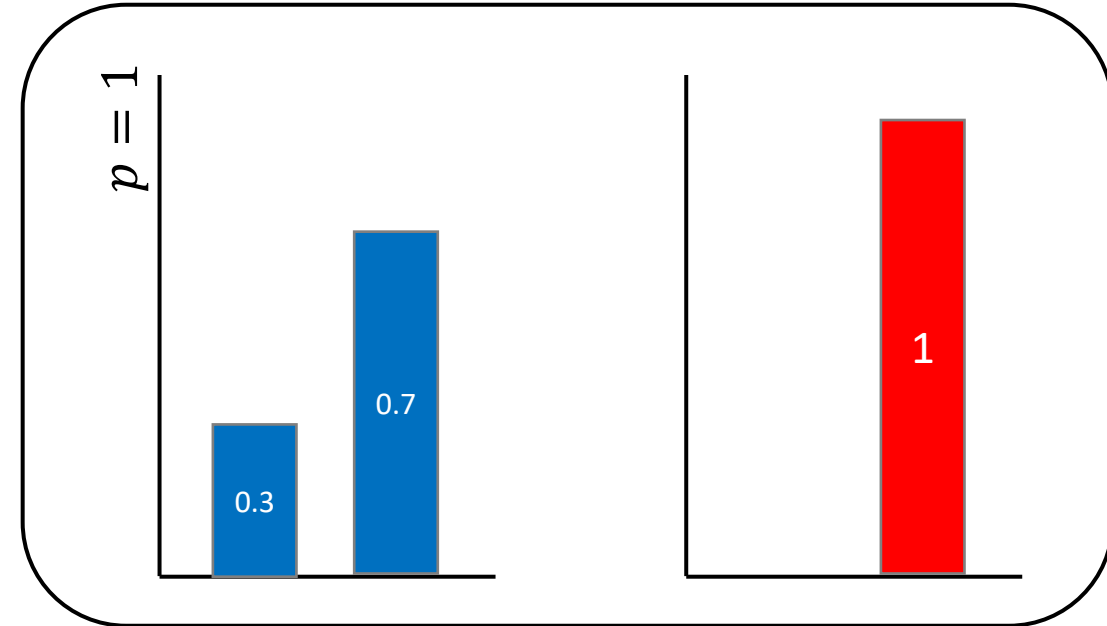
Activation

Output

$$\vec{x} = \vec{w}^T \vec{x}$$

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

$$y = \sigma(\vec{w}^T \vec{x} + b)$$



Model

Ground truth data

$J \approx$  Measure of difference in distributions



# Cross Entropy Loss function (Binary)

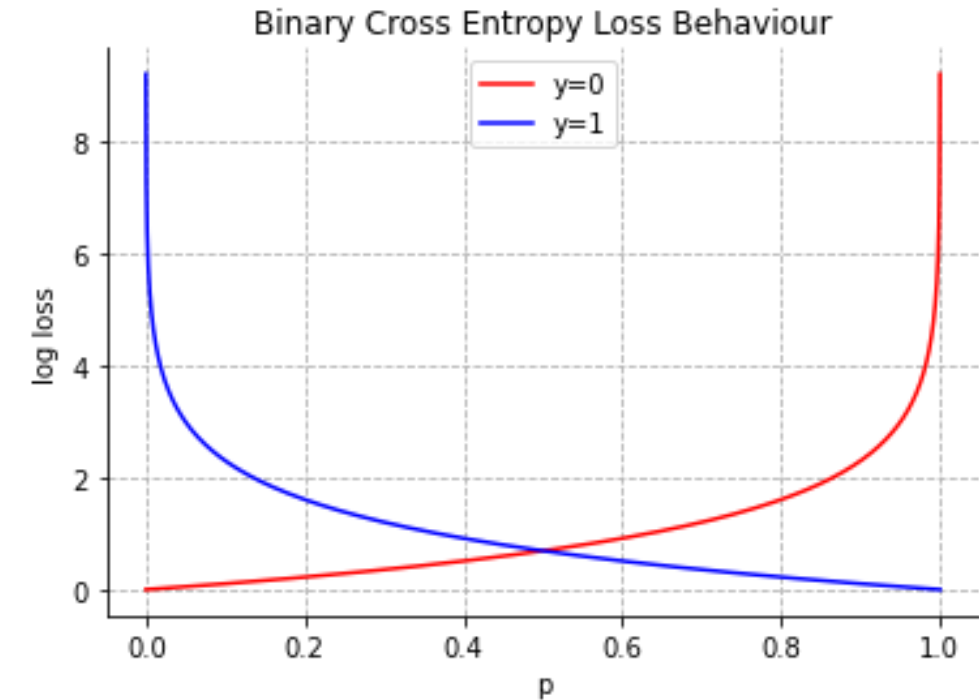
$$L(\hat{y}, y) = f(x) = \begin{cases} -\log(\hat{y}), & y = 1 \\ -\log(1 - \hat{y}), & y = 0 \end{cases}$$



# Cross Entropy Loss function (Binary)

$$L(\hat{y}, y) = f(x) = \begin{cases} -\log(\hat{y}), & y = 1 \\ -\log(1 - \hat{y}), & y = 0 \end{cases}$$

$$L(\hat{y}, y) = -(y \log \hat{y} + (1 - y) \log (1 - \hat{y}))$$







# Cross Entropy Loss function (Binary)

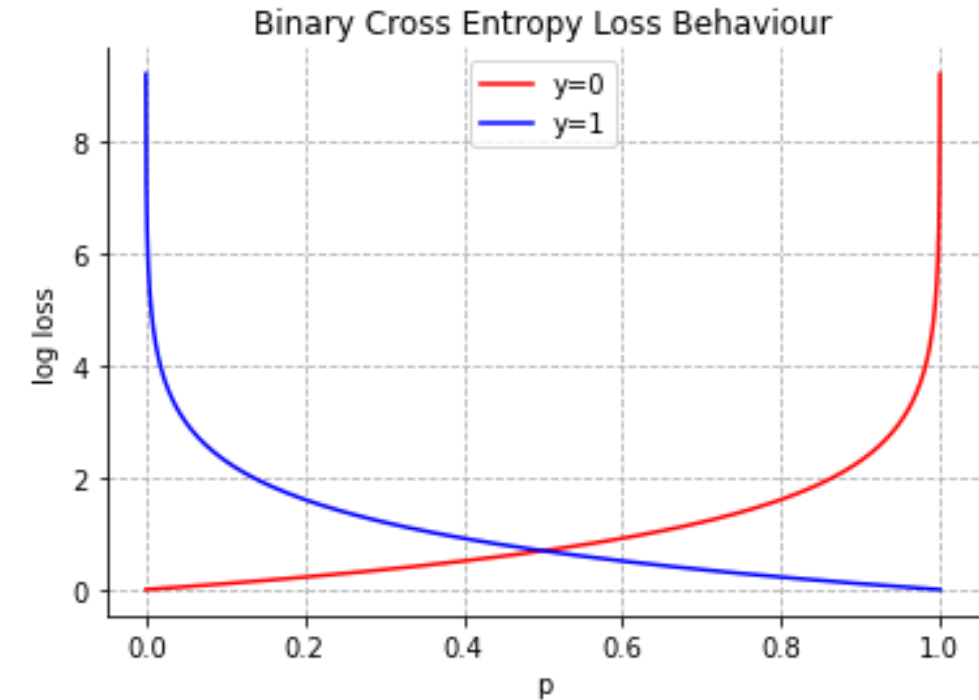
$$L(\hat{y}, y) = f(x) = \begin{cases} -\log(\hat{y}), & y = 1 \\ -\log(1 - \hat{y}), & y = 0 \end{cases}$$

$$L(\hat{y}, y) = -(y \log \hat{y} + (1 - y) \log (1 - \hat{y}))$$

Cross Entropy Loss is **Convex**

$$\Leftrightarrow L(\hat{y}, y)'' \geq 0$$

The line segment between any two points does not lie below the graph





# Logistic Regression Training

Training Set  $D$

$$D : \{(\vec{x}^{(1)}, y^{(1)}), \dots, (\vec{x}^{(i)}, y^{(i)}), \dots, (\vec{x}^{(m)}, y^{(m)})\}$$

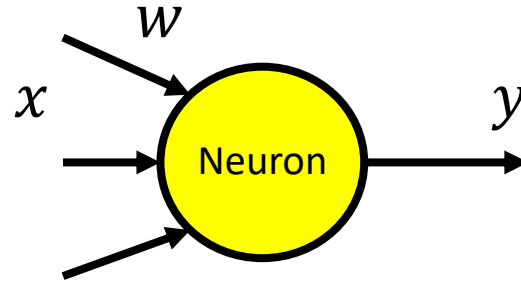
Cost  $J$

$$\begin{aligned} J(\{\hat{y}\}^m, \{y\}^m; \{\vec{x}\}^m) &= \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) \\ &= -\frac{1}{m} \sum_{i=1}^m (y^{(i)} \log \hat{y}^{(i)} + (1 - y^{(i)}) \log (1 - \hat{y}^{(i)})) \end{aligned}$$



# Minimizing Loss using Gradient Descent

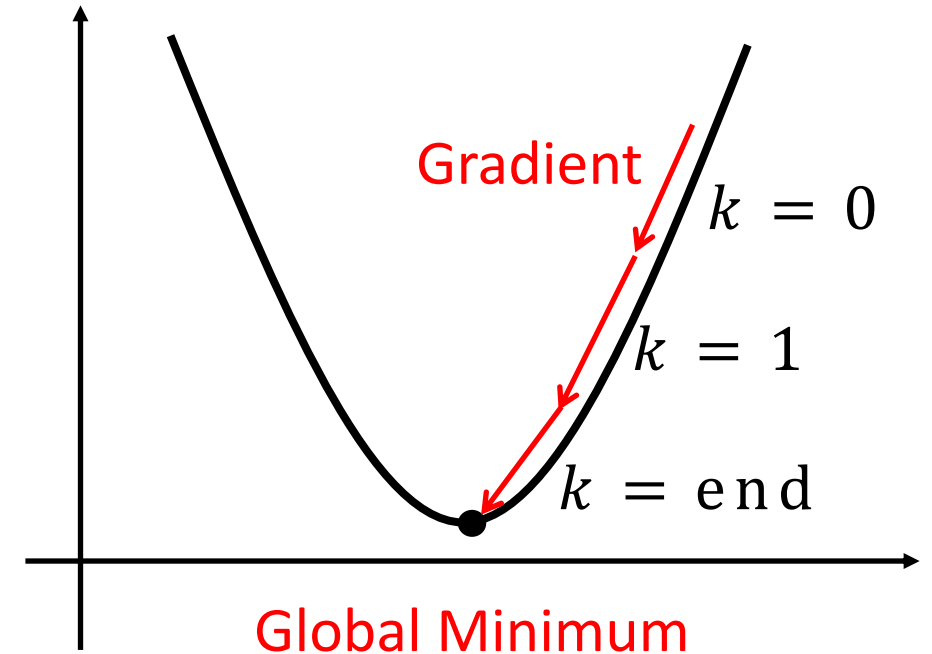
$$y = \sigma(\vec{w}^T \vec{x} + b)$$



$$J = L(\vec{w}, b, y)$$

$$\vec{w}_{k+1} = \vec{w}_k - \alpha \nabla_{\vec{w}} J(\vec{w}_k; b)$$

$$b_{k+1} = b_k - \alpha \frac{\partial}{\partial b} J(\vec{w}; b_k)$$





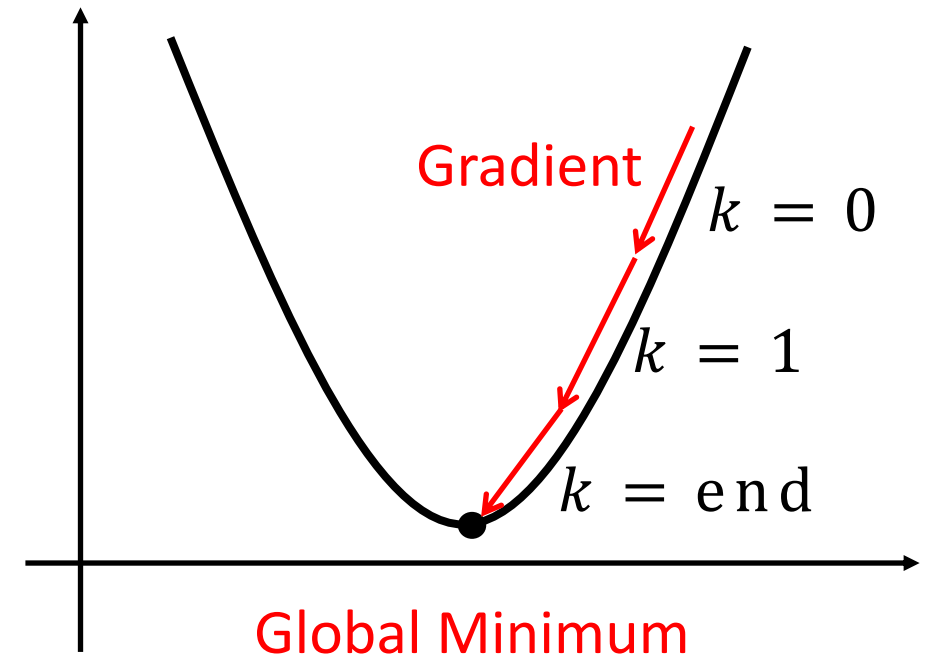
# Minimizing Loss using Gradient Descent

$$y = \sigma(\vec{w}^T \vec{x} + b)$$

$$\vec{w}_{k+1} = \vec{w}_k - \underset{\substack{\text{Learning rate} \\ \downarrow}}{\alpha} \nabla_{\vec{w}} J(\vec{w}_k; b)$$

$$b_{k+1} = b_k - \alpha \frac{\partial}{\partial b} J(\vec{w}; b_k)$$

$$J = L((\vec{w}, b), y)$$



Iteratively adjust  $(\vec{w}, b)$  until we reach the global minimum



# Finding Gradient w.r.t. weights and biases

$$J = L(\hat{y}, y)$$



# Finding Gradient w.r.t. weights and biases

$$J = L(\hat{y}, y)$$

$$\underbrace{\hat{y}}_{\hat{y} = \sigma(z)} \text{ Activation function}$$



# Finding Gradient w.r.t. weights and biases

$$J = L(\hat{y}, y)$$

$$\hat{y} = \sigma(z) \quad \text{Activation function}$$

$$z = \vec{w}^T \vec{x} + b \quad \text{Integrate Inputs}$$



# Finding Gradient w.r.t. weights and biases

$$J = L(\hat{y}, y)$$

$$\hat{y} = \sigma(z)$$

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

$$J = L(\sigma(\vec{w}^T \vec{x} + b), y)$$





# Finding Gradient w.r.t. weights and biases

$$J = L(\hat{y}, y)$$

$$\hat{y} = \sigma(z)$$

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

$$J = L(\sigma(\boxed{\vec{w}}^T \vec{x} + b), y)$$



# Using Chain Rule to Compute Gradients

$$J = L(\hat{y}, y)$$

$\underbrace{\hspace{1.5cm}}$

$$\hat{y} = \sigma(z)$$

$\underbrace{\hspace{1.5cm}}$

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

$$J = L(\sigma(\boxed{\vec{w}}^T \vec{x} + b), y)$$

$$P = f(g(h(x)))$$

$$\frac{dP}{dx} = \frac{df}{dg} * \frac{dg}{dh} * \frac{dh}{dx}$$



# Using Chain Rule to Compute Gradients

$$J = L(\overbrace{\sigma(\underbrace{\vec{w}^T \vec{x} + b}_Z)}^{\hat{y}}, y)$$

$$P = f(g(h(x)))$$

$$\frac{dP}{dx} = \frac{df}{dg} * \frac{dg}{dh} * \frac{dh}{dx}$$



# Using Chain Rule to Compute Gradients

$$J = L(\underbrace{\sigma(\underbrace{\vec{w}^T \vec{x} + b}_{z})}_{\hat{y}}, y)$$

$$P = f(g(h(x)))$$

$$\frac{\partial J}{\partial \vec{w}} = \frac{\partial L}{\partial \hat{y}} * \frac{\partial \hat{y}}{\partial z} * \frac{\partial z}{\partial \vec{w}}$$
$$\frac{\partial L}{\partial \hat{y}} * \frac{\partial \hat{y}}{\partial z} * \nabla_{\vec{w}} z$$

$$\frac{dP}{dx} = \frac{df}{dg} * \frac{dg}{dh} * \frac{dh}{dx}$$



# Using Chain Rule to Compute Gradients

$$\nabla_{\vec{w}} L(\hat{y}, y) = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z} \nabla_{\vec{w}} z$$

FWD

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

→

$$\hat{y} = \sigma(z)$$

→

$$L(\hat{y}, y)$$

BWD

$$\nabla_{\vec{w}} z$$

←

×

$$\frac{\partial \hat{y}}{\partial z}(z)$$

←

×

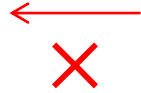
$$\frac{\partial L}{\partial \hat{y}}(\hat{y}, y)$$



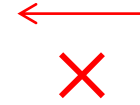
# Using Chain Rule to Compute Gradients

BWD

$$\nabla_{\vec{w}} z$$



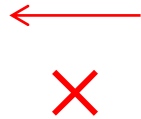
$$\frac{\partial \hat{y}}{\partial z}(z)$$



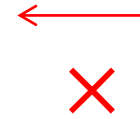
$$\frac{\partial L}{\partial \hat{y}}(\hat{y}, y)$$

$$\nabla_{\vec{w}} L$$

$$\vec{x}$$



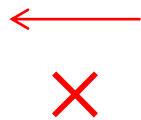
$$\sigma(z)(1 - \sigma(z))$$



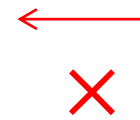
$$-\frac{y}{\hat{y}} + \frac{1 - y}{1 - \hat{y}}$$

$$\frac{\partial L}{\partial b}$$

$$1$$



$$\sigma(z)(1 - \sigma(z))$$



$$-\frac{y}{\hat{y}} + \frac{1 - y}{1 - \hat{y}}$$



# Training Terminologies

- **Forward Propagation:**  
Computing the **loss through forward pass** for a single training example
- **Backward Propagation:**  
Computing **gradients of parameters** through backward pass for a single training example
- **Batch:**  
Training set could be divided into **smaller sets** called batches
- **Iteration:**  
When an **entire batch** is passed both **forward and backward**
- **Epoch:**  
When an **entire dataset** is passed both **forward and backward** through the NN once



# Batch Gradient Descent

Training Set  $D$

$$D : \{(\vec{x}^{(1)}, y^{(1)}), \dots, (\vec{x}^{(i)}, y^{(i)}), \dots, (\vec{x}^{(m)}, y^{(m)})\}$$

Cost  $J$

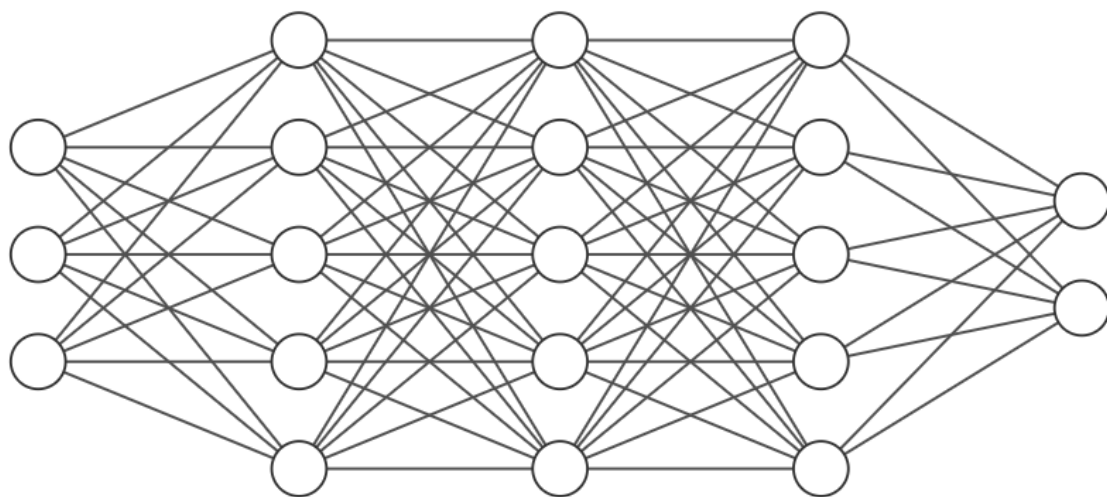
$$\begin{aligned} J(\{\hat{y}\}^m, \{y\}^m; \{\vec{x}\}^m) &= \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) \\ &= -\frac{1}{m} \sum_{i=1}^m (y^{(i)} \log \hat{y}^{(i)} + (1 - y^{(i)}) \log (1 - \hat{y}^{(i)})) \end{aligned}$$

Sum the gradients for all  $m$ -examples (Epoch)





# Next episode in EE P 596...



Optimizations in Deep Learning

