

Deep Learning Foundations and Algorithms

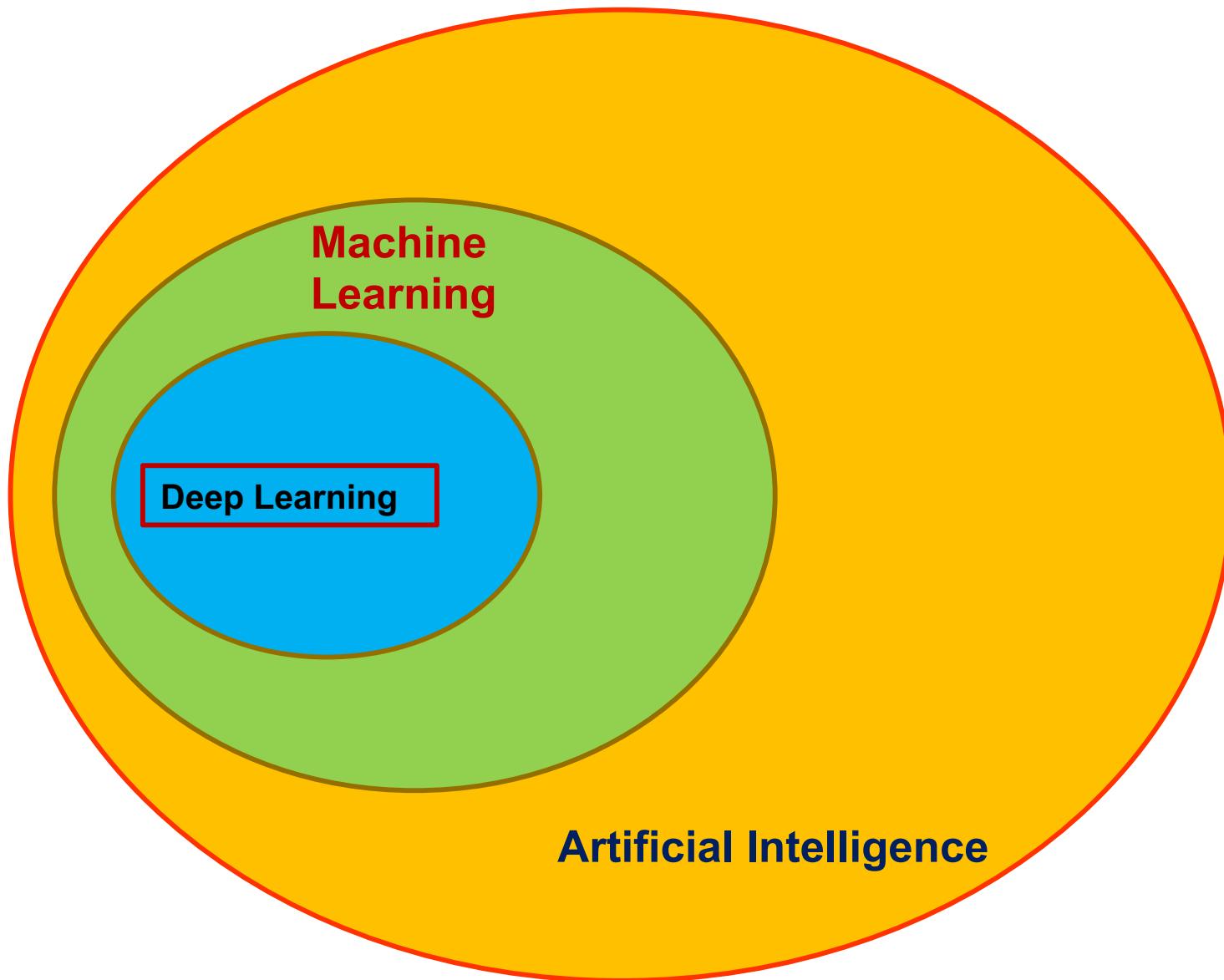
ADMSI -- Talent Sprint – Module 4.1A

Balaji Srinivasan
Mechanical Engineering
IIT-Madras

Topics for Today

- Recap + History of Neural Networks (Biological Inspiration)
- {
 - Gradient Descent → Linear Regression
 - Logistic Regression
- Multiclass Classification
- Neural Networks (FCN, MLP)
 - Not Work - Training (Back prop)
 - Non-linearities
 - Variants of GD
 - Regularization

Recall – The AI Venn diagram

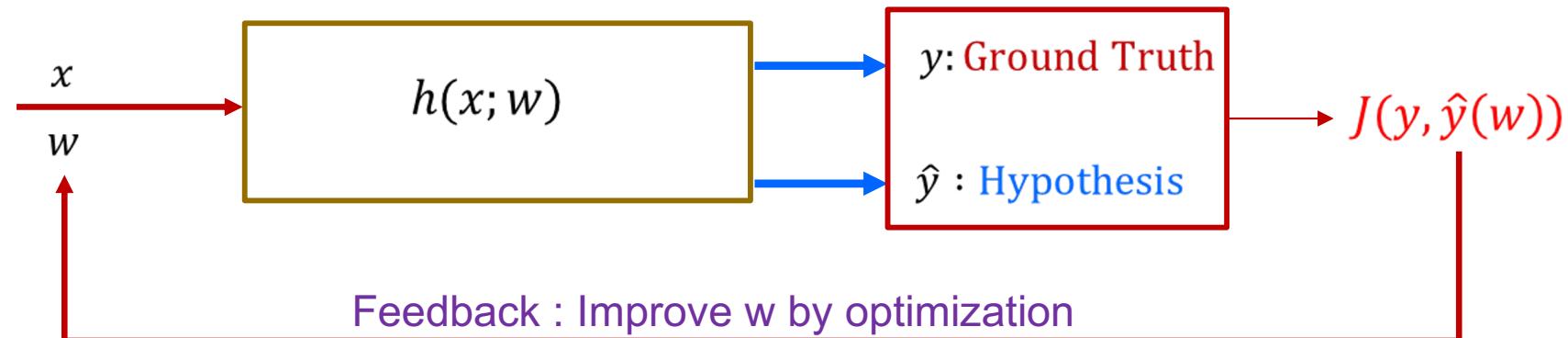


Let us now look at Deep Learning through some simple examples...

A fundamental “trick” in DL

- All problems are data, all solutions are functions/maps
- Cognitive tasks -- Humans get sensory inputs as qualia
 - We must convert these qualitative inputs into **numbers** – Input Vectors
 - Similarly, outputs that humans give must also be converted into numbers – Output/Target vectors
- Determining appropriate inputs and outputs for a machine learning task is an essential part of the process
- Often the “Learning Task” is learning the mapping from input to output.

Learning the parameters via feedback



To learn the parameters (given x, y find w), we follow this paradigm

- Collect lots of data pairs (Input Vector, Output Vector) = (x, y)
- Guess for the form of the **hypothesis function** $h(x; w)$
 - Example : $h(x; w) = w_0 + w_1 x_1 + w_2 x_2$
- For an arbitrary guess for w
 - We will get some $\hat{y} = h(x; w)$ which will not match the ground truth y
- Define a cost function $J(y, \hat{y}(w))$ depending on the difference
- Find optimal w by minimizing $J(w)$ → **Feedback process**
 - By using some optimization procedure such as **Gradient Descent**

Summary

Two main ideas in this session

1. What enables breadth of application?

By converting even qualitative problems into data problems and treating every rule which we wish to learn as a mathematical function.

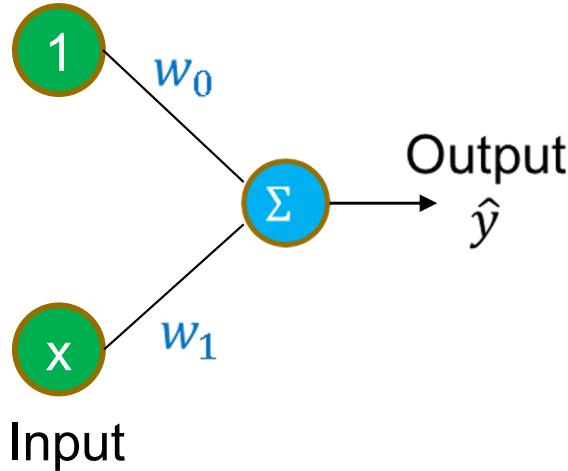


2. How does (supervised) learning happen?

1. We collect data as (Input, Output) pairs.
2. We provide a hypothesis as a form of the function connecting the inputs to the outputs
3. Learn (improve) the parameters based on mathematical feedback from data via loss function.

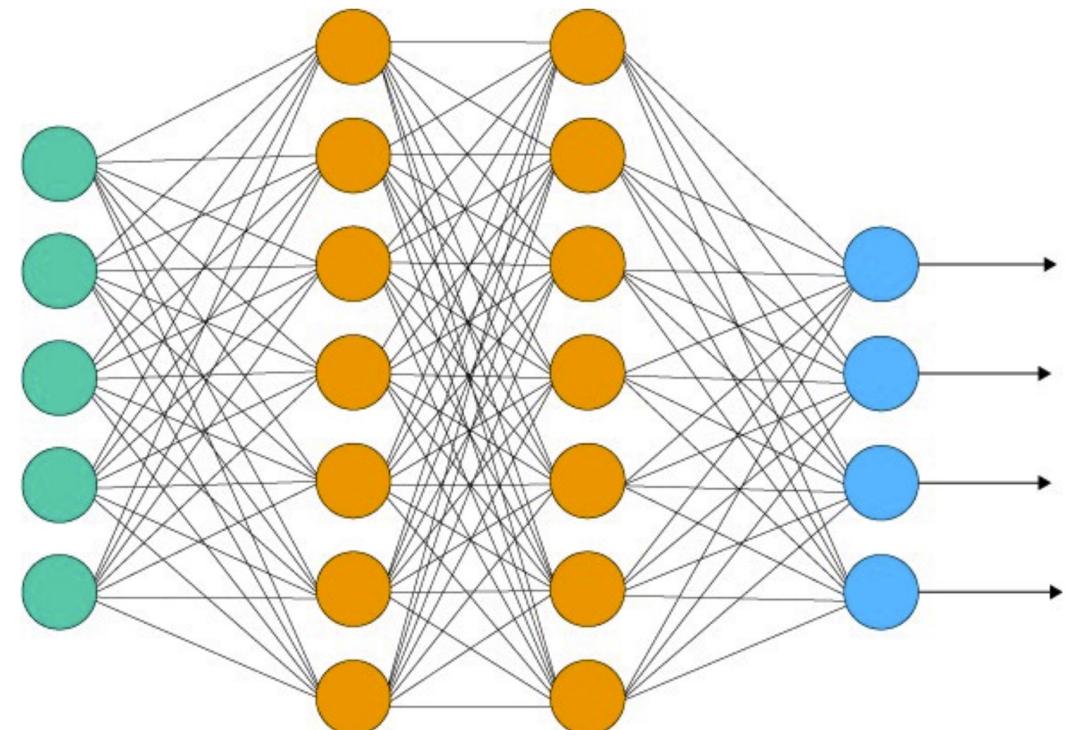
Neural Networks

Linear Model



$$\hat{y} = w_0 + w_1 x = \sum w_i x_i$$

Neural Network Model



Input Layer

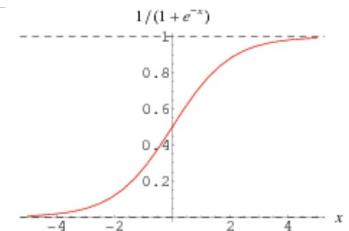
Hidden Layer

Output Layer

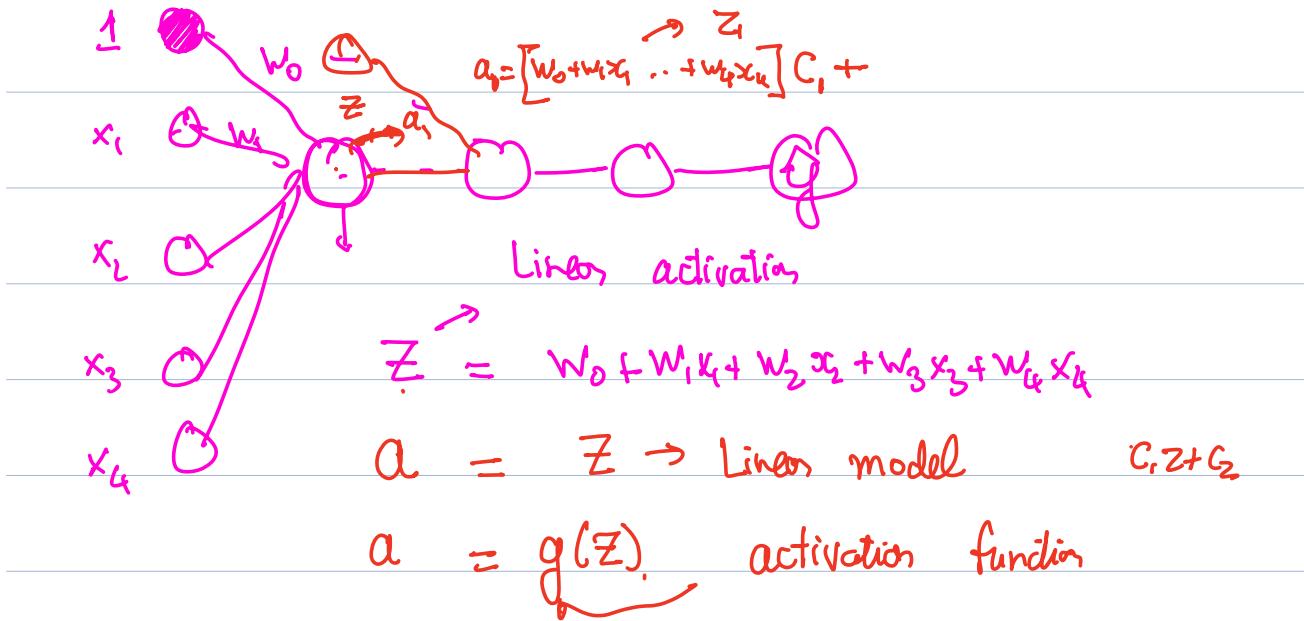
$$a = g(\sum w_i x_i)$$

$$a = \Sigma \rightarrow g$$

A Neuron has a linear and a nonlinear operation



Example nonlinear function is **Sigmoid(x)** = $\frac{1}{1+\exp(-x)}$



Examples $g(Z) = Z^2$

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

$$= \sin(Z)$$

Why?

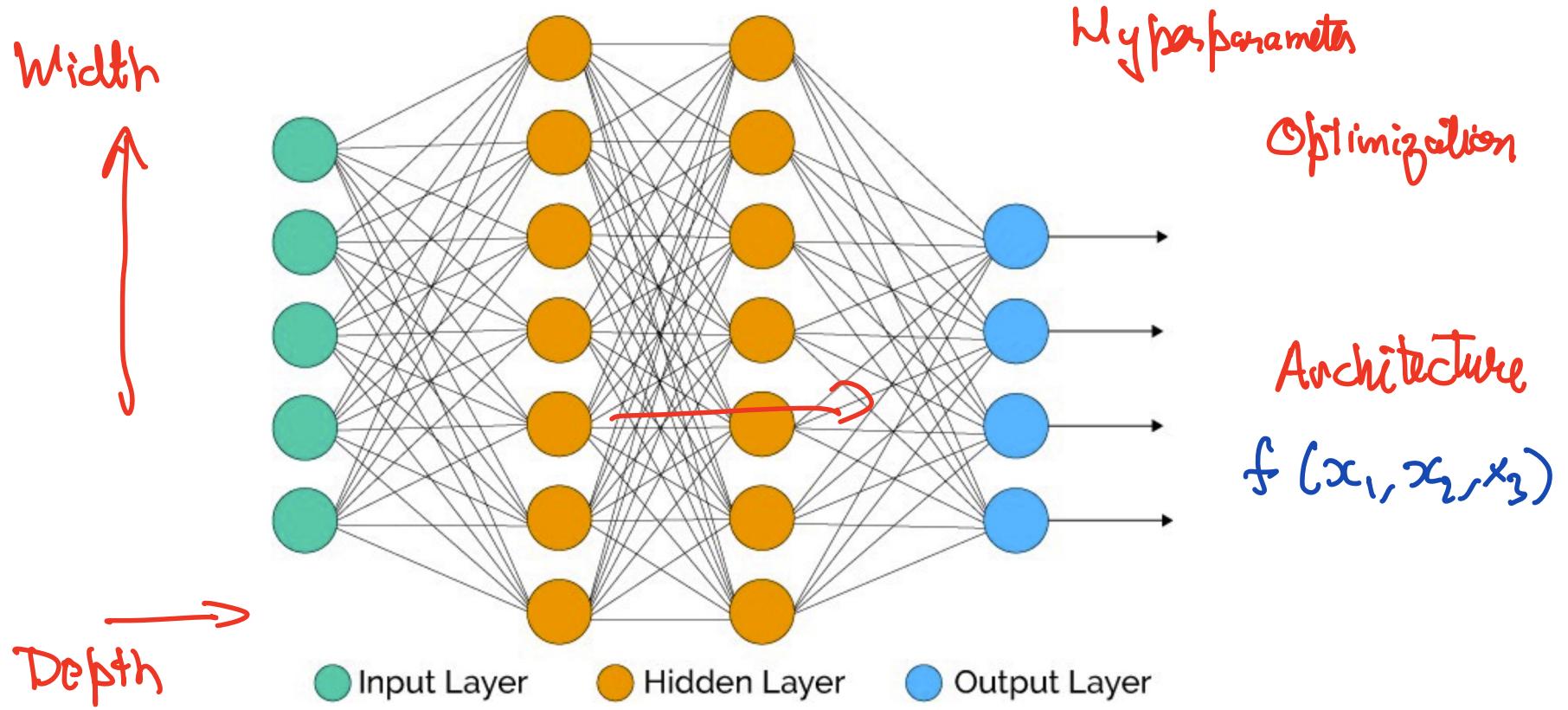


No nonlinearity $a_1 = (w_1 x_1 + w_2 x_2 + w_3 x_3 + w_4 x_4)$

$$a_2 = w_5 a_1 = (w_5 w_1 x_1 + w_5 w_2 x_2 + w_5 w_3 x_3 + w_5 w_4 x_4)$$

$w_5 w_1 x_1 \rightarrow$ Still linear in x

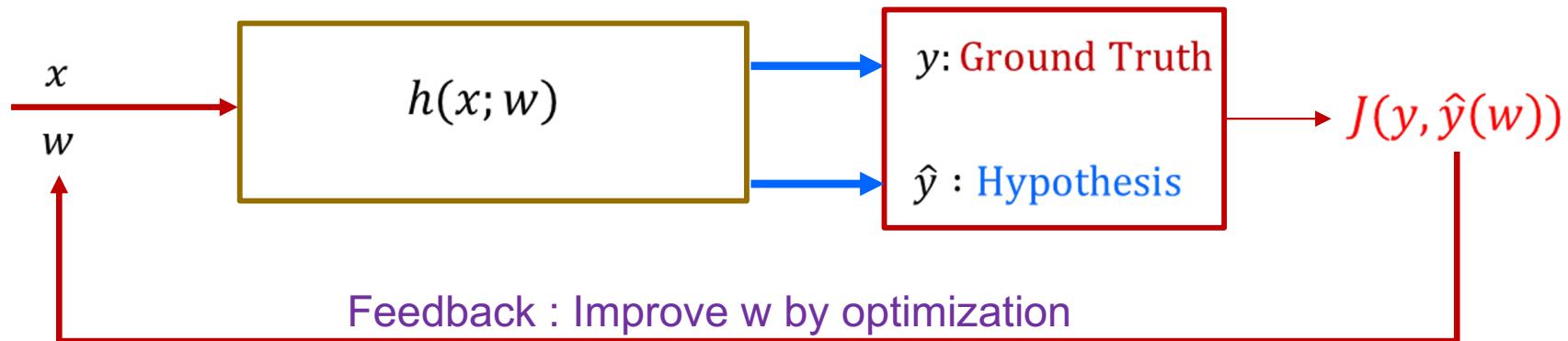
Universal Approximation Theorem



Theorem -- Given sufficient data and neurons, a Neural Network can approximate any function to any desired accuracy

Heart of modern Deep Learning – Lots of data, lots of layers => Lots of learning

Recall : Learning the parameters via feedback



- To learn the parameters, we follow this paradigm
 - Guess for the form of the **hypothesis function** $h(x; w)$
 - For an arbitrary guess for w
 - We will get some $\hat{y} = h(x; w)$ which will not match the ground truth y
 - Define a cost function $J(y, \hat{y}(w))$ depending on the difference
 - Find optimal w by **minimizing** $J(w)$ → Feedback process
 - Requires **Backpropagation** for $\frac{\partial J}{\partial w}$

Even complicated Neural Networks work in exactly the same way!

What the DL engineer must provide

- Appropriate decisions for input and output vectors (x, y)
 - Recall : All problems are data, all solutions are functions/maps

- Choosing appropriate datasets



- Some appropriate form of the forward model $\hat{y} = h(x; w)$
 - Example : Linear Model $\hat{y} = w_0 + w_1x_1 + w_2x_2 \dots + w_nx_n$

- Form of the Loss function

- Example : Least Squares $J(y, \hat{y}) = (y - \hat{y})^2$

- Optimization algorithm – Example: Gradient Descent

- And associated hyperparameters such as α

Machine Learning is not magic! It requires a lot of input from the engineer

Summary of Deep Learning

- Learning algorithms involve various types of hypothesis functions
 $\hat{y} = h(x; w)$.
- Each of these have their own purpose and domain where they work well
 - Linear Regression – For simple polynomial regression problems
 - Logistic Regression – For two-class (binary) classification problems
 - **Deep Neural Networks – For any general, non-linear problem**
 - There is a theorem that assures us that sufficiently large neural network will approximate *any* function
 - Convolutional Neural Networks – For vision/image based problems
 - Recurrent Neural Networks – For sequential/time-series problems
- There are also appropriate loss functions for each.
- It is possible that you might have a better model than these for your problem. The general procedure outlined here remains the same.

Calculation in a single neuron

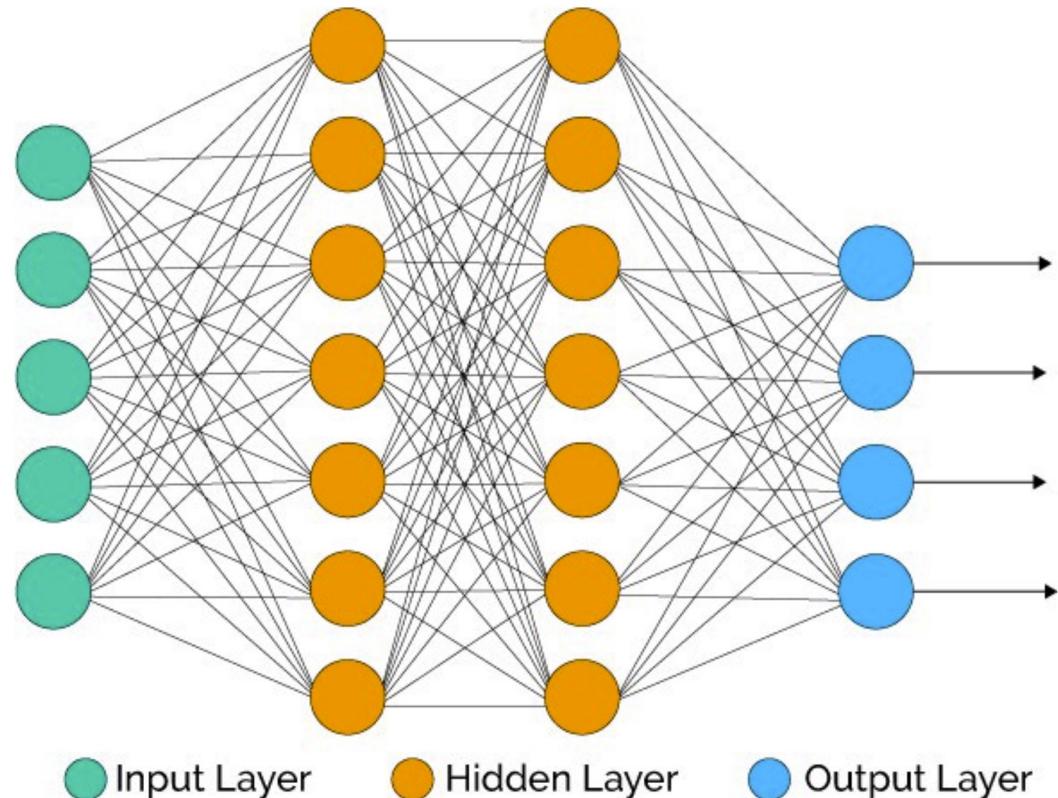
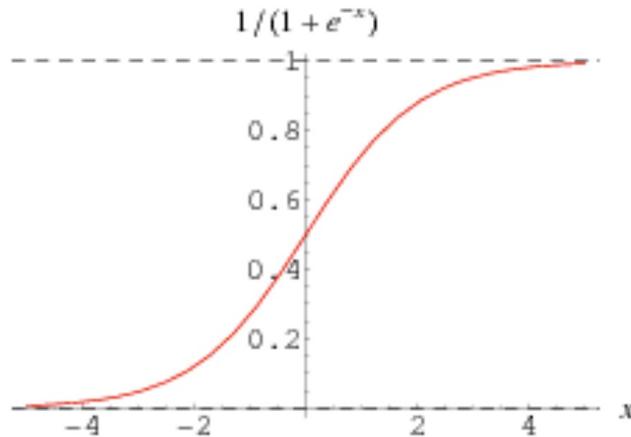
Recall : Neural Networks

Find $\sigma(z)$

Neural Network Model

Common Activation function

$$\text{Sigmoid}(x) = \frac{1}{1+\exp(-x)}$$

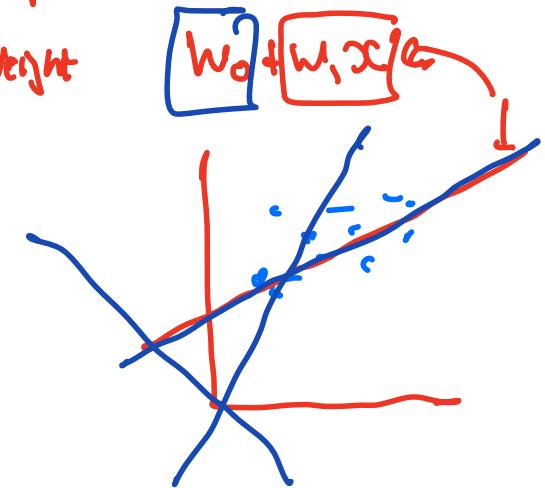
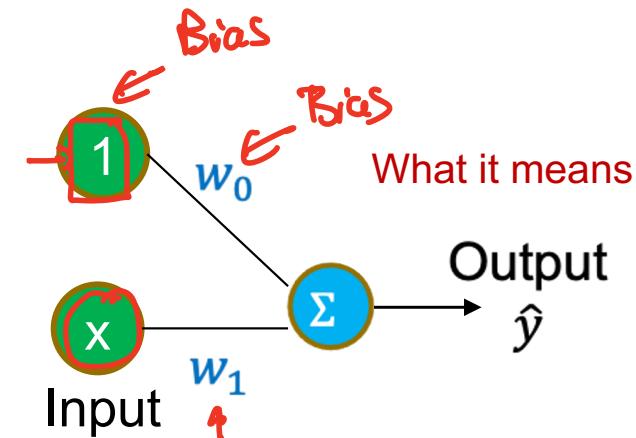
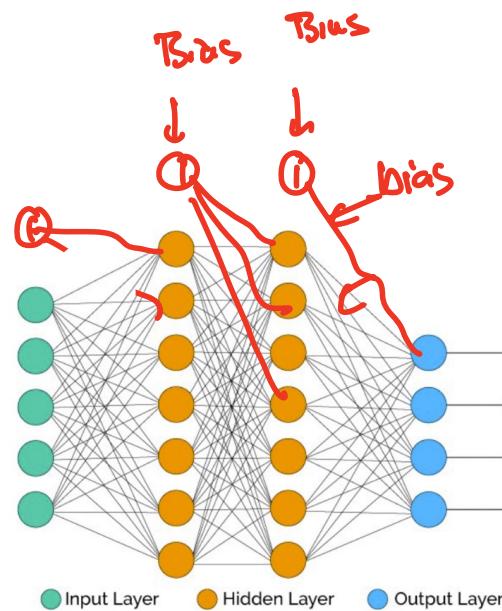
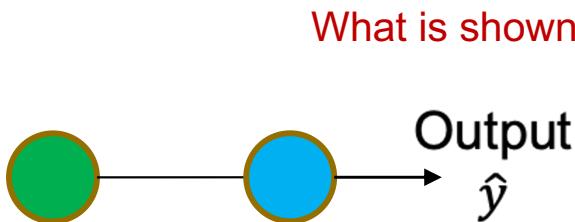


$$a = g(\sum w_i x_i)$$

A diagram showing the computation of a neuron's output. It consists of three circles: a yellow circle on the left labeled a , a blue circle in the middle labeled Σ , and a green circle on the right labeled g . Arrows point from the Σ circle to the g circle, and from the g circle to the a circle.

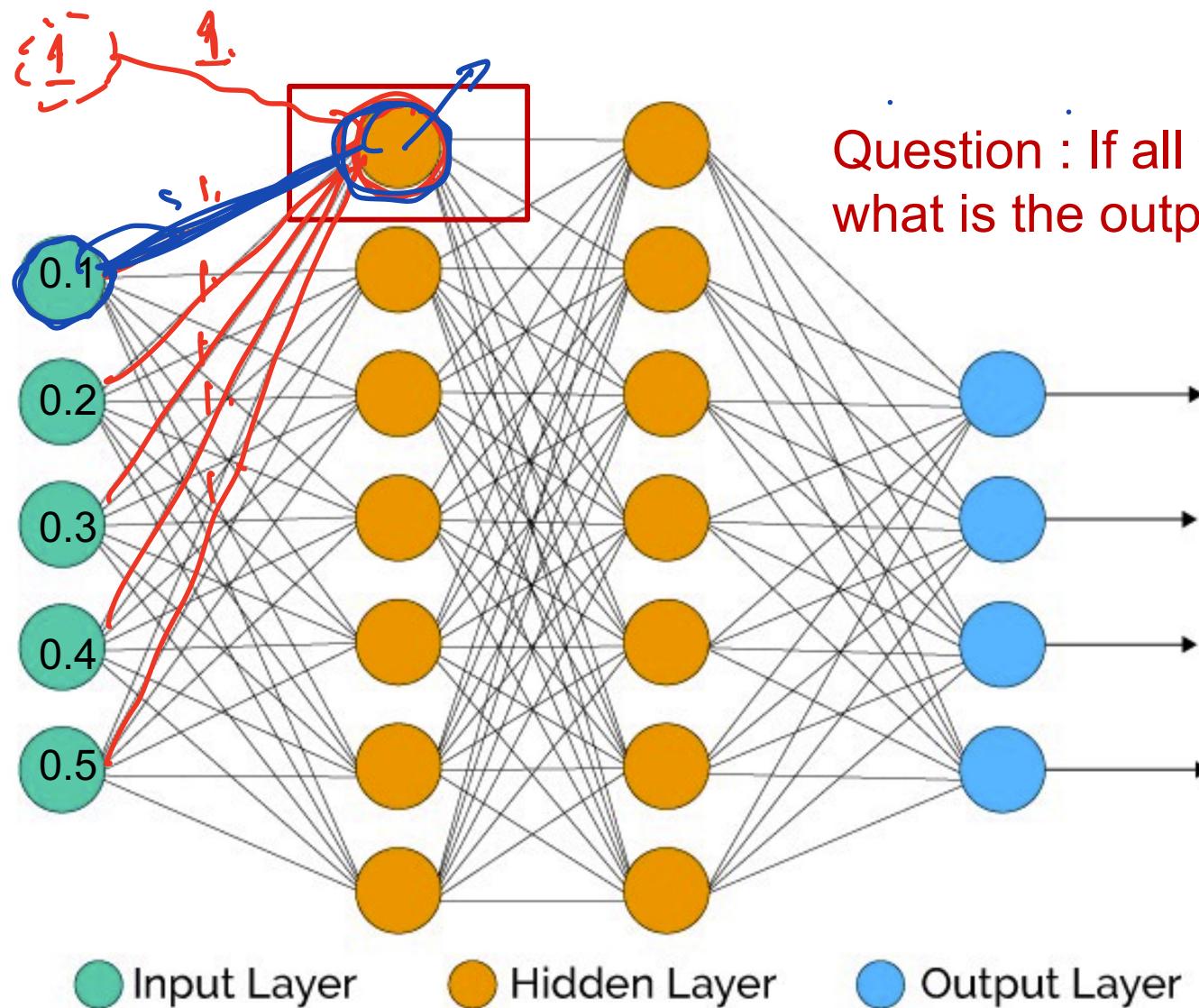
A Neuron has a linear and a nonlinear operation

Weights and biases



- Constant neurons representing a “feature” with 1 are usually not depicted in the pictorial representation.
- The parameters going from the constant neuron is called a bias and the neuron is called a “bias unit”
- The other parameters are called “weights”

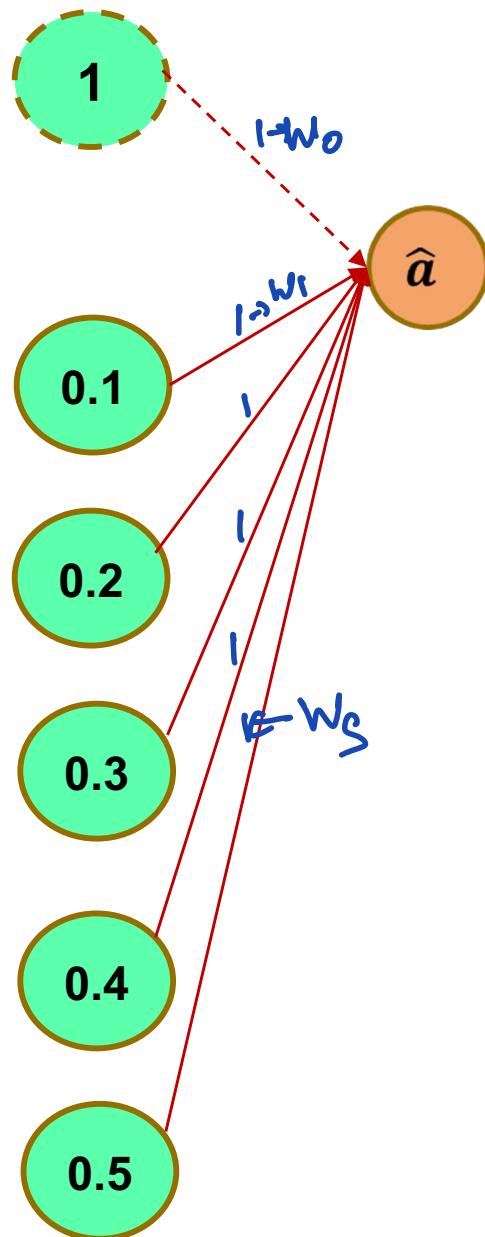
The output of a single neuron



Question : If all weights are 1, then what is the output of the shown neuron?

$$g(z) = \sigma(z)$$

$$x \xrightarrow{w} z \xrightarrow{g} a$$



Find the output of the neuron \hat{a}

All weights are equal to 1

Ans: The output has two parts / calculations / sum

Linear -- $z = \sum w_i x_i \rightarrow$ Single number

$$z = w_0 x_0 + w_1 x_1 + w_2 x_2 + w_3 x_3 + w_4 x_4 + w_5 x_5$$

$$z = 1 + 0.1 + 0.2 + 0.3 + 0.4 + 0.5 = 2.5$$

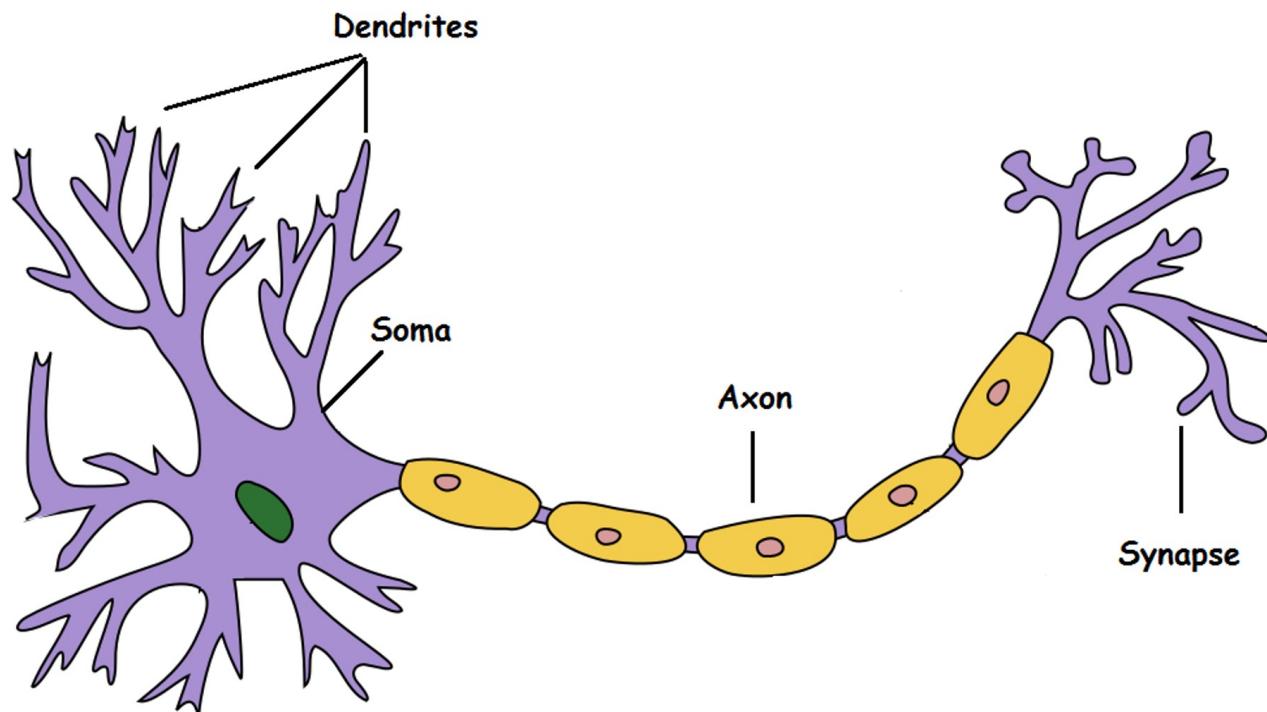
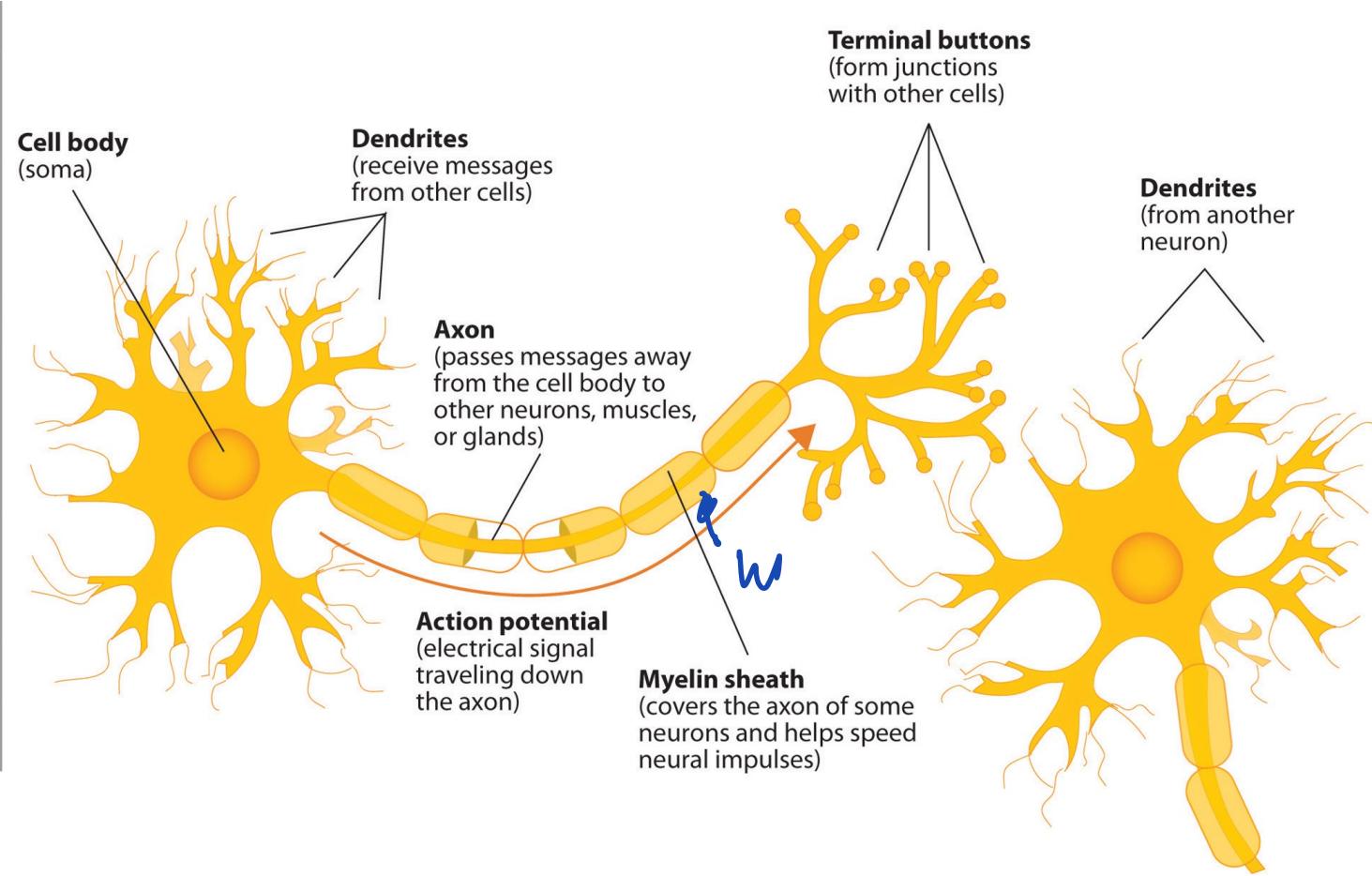
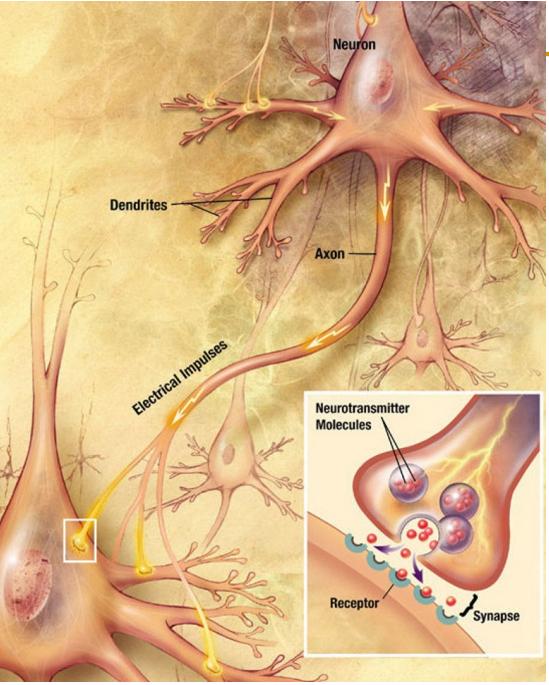
Nonlinear Activation -- $\hat{a} = g(z)$

$$= \frac{1}{(1 + \exp(-2.5))}$$

$$= 0.9241$$

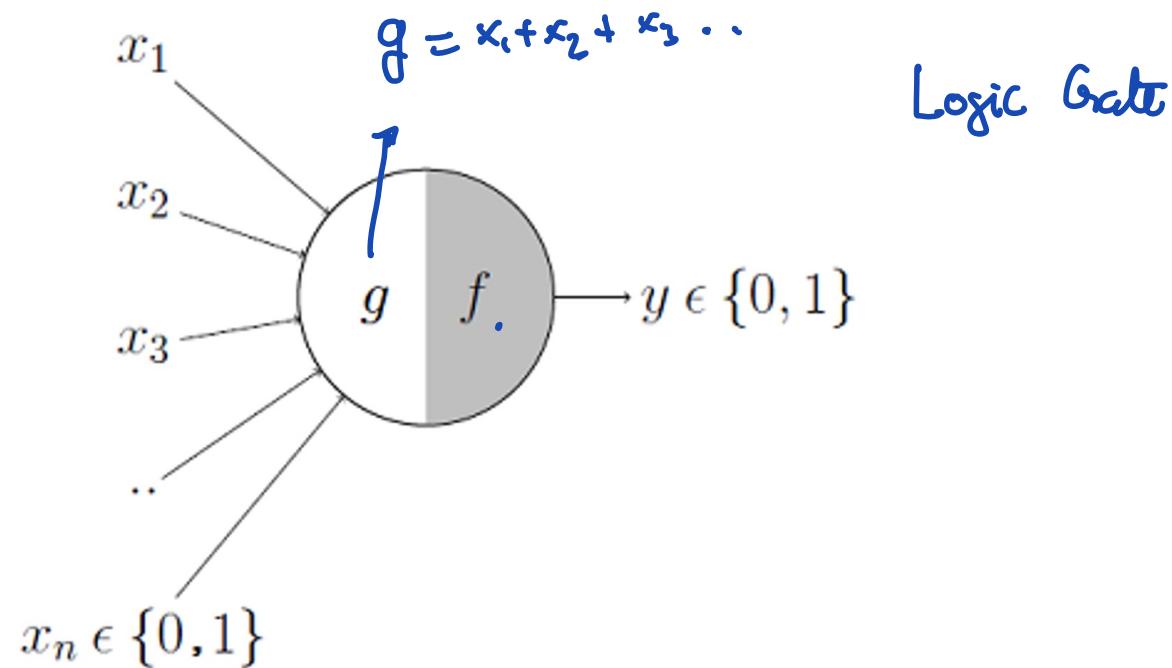
Don't forget the bias unit!

From biology to computation



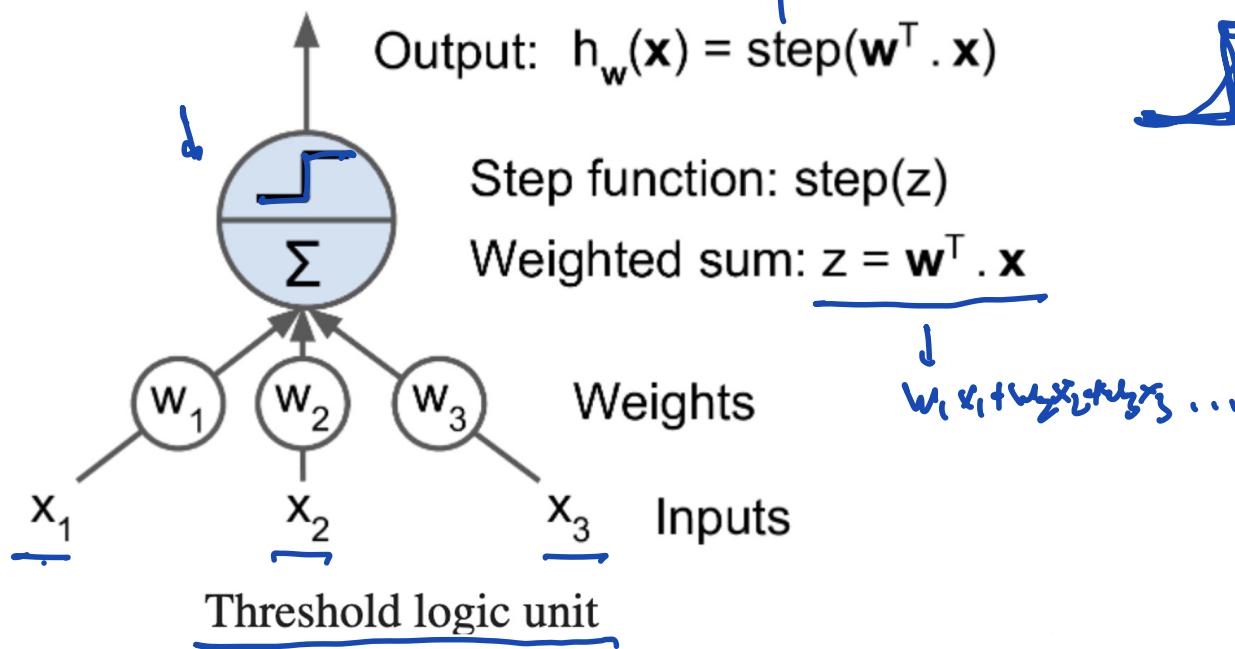
McCulloch-Pitts Neuron

- First Mathematical model of Neuron (1943)

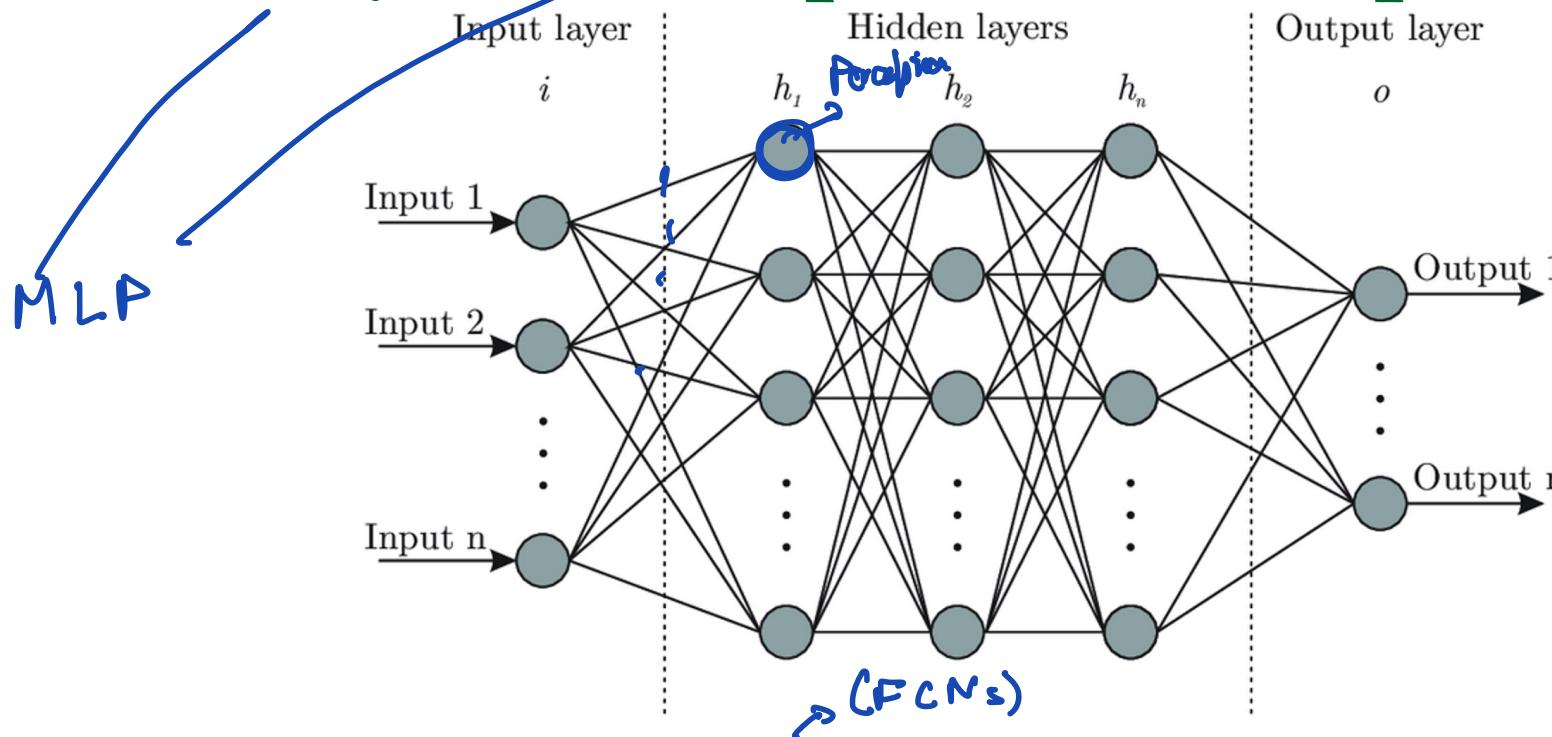


Rosenblatt Perceptron (1957)

Differentiable
Smooth \Rightarrow
Easy to train
Easier to interpret
as prob



Multilayer Perceptrons/Deep NNs



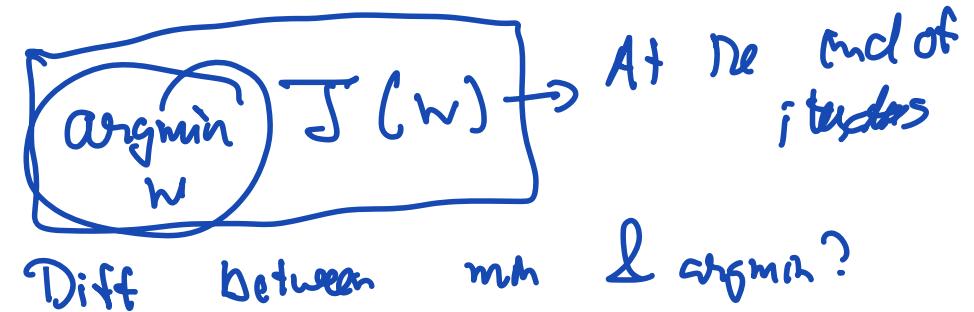
- MLPs or Fully connected Networks are very common
 - At least as parts of other, more sparse networks → less # parameters
- Satisfy the **Universal Approximation Theorem**
 - Can approximate any function to arbitrary accuracy with just one layer.

Types of Neural Networks

- Artificial Neural Networks }
 - For general numerical data
 - Convolutional Neural Networks (CNNs) }
 - For image based data
 - Recurrent Neural Networks (RNNs) }
 - For sequential data
- Transformers

Detailed Gradient Descent Derivation For Linear Regression

Introduction



- Gradient Descent is widely used in machine learning algorithms.
- Used to minimize cost functions and optimize models.
- In this presentation, we delve deep into its derivation for linear regression with one variable.

Value that J takes

$$\min \quad J(w_1, w_2) = (w_1 - 1)^2 + (w_2 - 2)^2$$

\rightarrow What is the w that led to min?

$$(w_1 = 1, w_2 = 2) \rightarrow \text{argmin}$$

Linear Regression Model with One Variable

X

S.No

X
Area

Grand Total
 y
 \hat{y}

1

1000 $\rightarrow x^0$
50 L

2

2000 $\rightarrow x^0$
70 L

3

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

$$\hat{y} = w_0 + w_1 x$$

The hypothesis for our model:

Aiming to minimize the cost function $J(w_0, w_1)$:

MSE avg cost

$$J(w_0, w_1) = \frac{1}{2m} \sum_{i=1}^m (\hat{y}^i - y^i)^2$$

m \rightarrow 100
Original Data

ith data pt.

$x^{(i)}$

ith data pt.

$y^{(i)}$

Grand truth

$m = \#$ of data samples
in the training set

Derivative with respect to w_0

$$\frac{\partial J}{\partial w_0}, \frac{\partial J}{\partial w_1}$$

Differentiating the squared error with respect to w_0 :

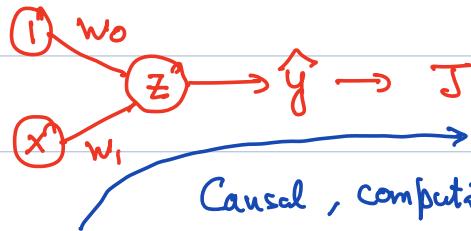
$$\frac{\partial}{\partial w_0} (\hat{y}^i - y^i)^2 = 2(\hat{y}^i - y^i) \frac{\partial \hat{y}^i}{\partial w_0}$$

Given:

$$\frac{\partial \hat{y}^i}{\partial w_0} = 1$$

We have:

$$\frac{\partial}{\partial w_0} J(w_0, w_1) = \frac{1}{m} \sum_{i=1}^m (\hat{y}^i - y^i)$$



Causal, computational chain which led
from $w_0, w_1 \rightarrow J$.

$$\begin{aligned} z &= w_0 + w_1 x & \frac{\partial z}{\partial w_0} &= 1 \\ \hat{y} &= z \rightarrow & \frac{\partial z}{\partial w_1} &= x \\ J &= \frac{1}{2} (\hat{y} - y)^2 & \frac{\partial J}{\partial w_0} &= \frac{\partial J}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z} \frac{\partial z}{\partial w_0} \\ &&& (\hat{y} - y) & 1 & 1 \\ &&& \frac{\partial J}{\partial w_0} &= (\hat{y} - y) * 1 \end{aligned}$$

"DELTA" RULE

$$\frac{\partial J}{\partial w_j} = e * x_j \quad \text{Correspondingly rule}$$

$$\frac{\partial J}{\partial w_1} = \frac{\partial J}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z} \frac{\partial z}{\partial w_1}$$

$$\frac{\partial J}{\partial w_j} = (\hat{y} - y) x_j \quad \text{Activation of the input}$$

Error in the Output

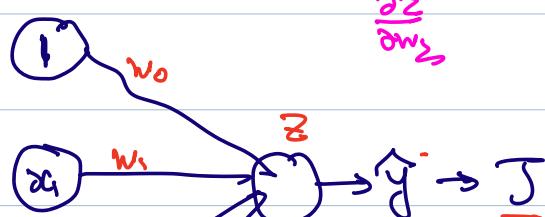
$$z = w_0 + w_1 x_1 + w_2 x_2 + w_3 x_3$$

$$\frac{\partial z}{\partial w_2} = x_2$$

$$\frac{\partial J}{\partial w_0} = (\hat{y} - y)$$

$$z = w_0 + w_1 x_1 + w_2 x_2 + w_3 x_3$$

$$\frac{\partial z}{\partial w_3}$$



$$\frac{\partial J}{\partial w_0} = e * 1$$

$$\frac{\partial J}{\partial w_2} = e * x_2$$

$$\frac{\partial J}{\partial w_1} = e * x_1$$

$$\frac{\partial J}{\partial w_3} = e * x_3$$

$$\frac{\partial J}{\partial w_2} = \frac{\partial J}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z} \frac{\partial z}{\partial w_2}$$

Derivative with respect to w_1

Using the chain rule:

$$\frac{\partial}{\partial w_1}(\hat{y}^i - y^i)^2 = 2(\hat{y}^i - y^i)x^i$$

Given:

$$\frac{\partial}{\partial w_1}\hat{y}^i = x^i$$

We get:

$$\frac{\partial}{\partial w_1}J(w_0, w_1) = \frac{1}{m} \sum_{i=1}^m (\hat{y}^i - y^i)x^i$$

Derivative with respect to w_1

Using the chain rule:

$$\frac{\partial}{\partial w_1}(\hat{y}^i - y^i)^2 = 2(\hat{y}^i - y^i)x^i$$

Given:

$$\frac{\partial}{\partial w_1}\hat{y}^i = x^i$$

We get:

$$\frac{\partial}{\partial w_1}J(w_0, w_1) = \frac{1}{m} \sum_{i=1}^m (\hat{y}^i - y^i)x^i$$

Final Expressions for linear regression

$$\text{Hypothesis: } \hat{y} = w_0 + w_1 x$$

$$\text{Cost Fn: } J = \frac{1}{2m} \sum (\hat{y}^i - y^i)^2$$

$$\text{Grad w.r.t } w_0 : \frac{\partial J}{\partial w_0} = \frac{1}{m} \sum (\hat{y}^i - y^i)$$

$$\text{Grad w.r.t } w_1 : \frac{\partial J}{\partial w_1} = \frac{1}{m} \sum (\hat{y}^i - y^i) x^i$$

Gradient Descent Updates

$$w_0 \leftarrow w_0 - \alpha \frac{\partial J}{\partial w_0}$$

$$w_1 \leftarrow w_1 - \alpha \frac{\partial J}{\partial w_1}$$

Quadratic Regression: Final Expressions

Lin. Model

$$\begin{matrix} x_1 \\ x_2 \\ \vdots \\ x_1^2 \end{matrix}$$

S.No	x_1	x_2	y	\hat{y}
.
.
.
.
.



Hypothesis: $\hat{y} = w_0 + w_1 x + w_2 x^2$

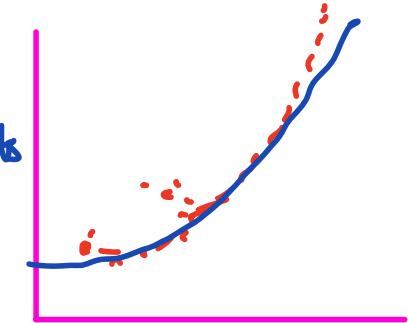
Features are nonlin

in weights

$$\text{Grad w.r.t } w_0 : \frac{\partial J}{\partial w_0} = \frac{1}{m} \sum (\hat{y}^i - y^i)$$

$$\text{Grad w.r.t } w_1 : \frac{\partial J}{\partial w_1} = \frac{1}{m} \sum (\hat{y}^i - y^i) x^i$$

$$\text{Grad w.r.t } w_2 : \frac{\partial J}{\partial w_2} = \frac{1}{m} \sum (\hat{y}^i - y^i) (x^i)^2$$



Nonlin model
 w_1, w_2

$$\hat{y} = w_0 + w_1 x + w_2 \ln(x) + w_3 t^2 \rightarrow \text{Lin as model}$$

Feature Engineering

Gradient Descent Updates

$$w_0 \leftarrow w_0 - \alpha \frac{\partial J}{\partial w_0}$$

$$w_1 \leftarrow w_1 - \alpha \frac{\partial J}{\partial w_1}$$

$$w_2 \leftarrow w_2 - \alpha \frac{\partial J}{\partial w_2}$$

General n Feature Case: Final Expressions

Hypothesis: $\hat{y} = w_0 + \sum_{j=1}^n w_j x_j$

n features

Cost Fn: $J(\mathbf{w}) = \frac{1}{2m} \sum (\hat{y}^i - y^i)^2$

Gradient w.r.t w_j : $\frac{\partial J}{\partial w_j} = \frac{1}{m} \sum (\hat{y}^i - y^i) x_j^i$ for $j = 0, 1, \dots, n$

¶

Data pts

Gradient Descent Updates

$$w_j \leftarrow w_j - \alpha \frac{\partial J}{\partial w_j} \quad \text{for each } j = 0, 1, \dots, n$$

Matrix Notation for Linear Regression

Design Matrix (X)

The design matrix X organizes the feature values of our dataset:

- ▶ Each row corresponds to an observation.
- ▶ Each column corresponds to a feature.

$X_{ij} \rightarrow j^{th}$ feature
 $X_{32} \leftarrow 3^{rd}$ data pt
 $X_{32} \leftarrow 2^{nd}$ feature

1^{st} data pt
Const feature \rightarrow Bias

$$X = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1n} \\ 1 & x_{21} & x_{22} & \dots & x_{2n} \\ \vdots & \vdots & x_{32} & \ddots & \vdots \\ 1 & x_{m1} & x_{m2} & \dots & x_{mn} \end{bmatrix}$$

S.No | x_1 | x_2 | x_3 | y
1 | | | |
..... | | | |
..... | | | |

Design Matrix

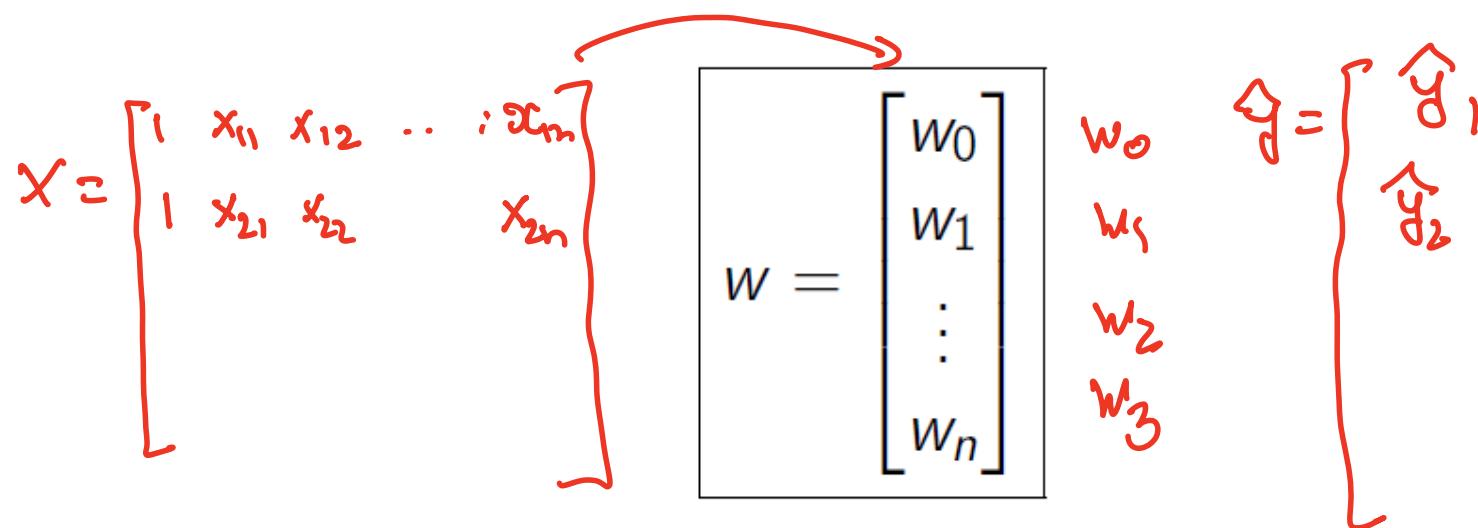
Matrix of inputs
arranged column wise

The column of ones allows for the inclusion of the bias term in matrix notation.

Parameters (w)

$$y = w^T x$$

The parameters or weights, including the bias term, are represented as:



- ▶ w_0 is the bias term.
- ▶ w_1, \dots, w_n are the weights for each feature.

Hypothesis Function (\hat{y})

The model prediction or hypothesis function is denoted by \hat{y} and is given by:

$$\hat{y} = Xw$$

Matrix mult

Matrix multiplication

Entries defined

Dimensions:
 \hat{y} : 100×1
 X : 100×4
 w : 4×1
 $m \times 1$ $m \times (n+1)$ $(n+1) \times 1$

This is a vector where each element is the predicted response for an observation.

$$\hat{y}^{(1)} = w_0 + w_1 x_1^{(1)} + w_2 x_2^{(1)} \dots + w_n x_n^{(1)}$$
$$= w_0 + w_1 x_{11} + w_2 x_{12} + w_3 x_{13} \dots + w_n x_{1n}$$

$$m = 100$$

$$n = 3$$

↳ Does not include bias

Model Representation with Bias Included

With the bias w_0 included in the parameters:

$$\hat{y} = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1n} \\ 1 & x_{21} & x_{22} & \dots & x_{2n} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{m1} & x_{m2} & \dots & x_{mn} \end{bmatrix} \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_n \end{bmatrix}$$

This compact form is efficient for computations and modeling.

Model Representation with Separate Bias

Alternatively, keeping the bias term separate:

$$\hat{y} = Xw + b$$

Here, **b** is a vector containing the bias term repeated for each observation.

- ▶ This form is mathematically equivalent to the previous one.
 - ▶ The choice of representation depends on computational or conceptual preferences.
-

Derivation of the General Gradient Expression for the Linear Model in Linear Regression

Linear Regression Model

- Let's consider a dataset with m training examples and n features.
- Our design matrix X is $m \times (n + 1)$ including a bias term.
- The parameter vector w is $(n + 1) \times 1$.
- The observed value vector y is $m \times 1$.

The linear hypothesis function is given by:

$$\hat{y} = Xw$$

$m \times 1$

The goal is to minimize the cost function, which in the case of linear regression is often the Mean Squared Error (MSE).

Cost Function

The Mean Squared Error (MSE) cost function is defined as:

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

$$J(w) = \frac{1}{2m} (Xw - y)^T e$$

$$\vec{v} = \begin{bmatrix} v_1 \\ v_2 \\ v_3 \end{bmatrix}$$

$$v^T v = [v_1 \ v_2 \ v_3] \begin{bmatrix} v_1 \\ v_2 \\ v_3 \end{bmatrix}$$

$$= v_1^2 + v_2^2 + v_3^2$$

To minimize $J(w)$, we will calculate its gradient $\nabla_w J(w)$.

Gradient Derivation: Step 1

Start with applying the chain rule to the cost function:

$$\nabla_w J(w) = \frac{1}{2m} \nabla_w [(Xw - y)^\top (Xw - y)]$$



Gradient Derivation: Step 2

Exploit the fact that for a scalar c , $\nabla_w c = 0$, and use the identity $\nabla_w a^\top w = a$ for a independent of w :

$$\nabla_w J(w) = \frac{1}{2m} \nabla_w [w^\top X^\top X w - 2y^\top X w + y^\top y]$$

Note: $y^\top y$ is a constant with respect to w .

Gradient Derivation: Step 3

We find the gradient of the quadratic form $w^\top X^\top X w$ and linear terms $y^\top X w$:

$$\nabla_w J(w) = \frac{1}{2m} [2X^\top X w - 2X^\top y]$$



Gradient Expression

$$x^T X w \approx x^T y$$

After simplifying, the gradient of the cost function with respect to w is:

$$\nabla_w J(w) = \frac{1}{m} X^T (Xw - y)$$

This is the direction of the steepest ascent in the cost function surface, and we need to move in the opposite direction to minimize the cost.

$$\nabla_w J = 0 \Rightarrow x^T x w - x^T y = 0$$
$$x^T x \underline{\underline{w}} = x^T \underline{\underline{y}}$$

Normal Eqs
 $\hat{w} = (x^T x)^{-1} x^T y$

Parameter Update Rule

In the gradient descent algorithm, we iteratively update w using the update rule:

$$w := w - \alpha \nabla_w J(w)$$

where α is the learning rate.

Conclusion

- We derived the expression for the gradient of the cost function in linear regression.
 - The gradient is used to update the parameters w in the direction that minimizes the cost.
 - Gradient descent is an iterative optimization algorithm used for finding the values of parameters w that minimize a cost function $J(w)$.
-

Gradient Descent Algorithm

- ① **Initialization:** Choose an initial vector of parameters \mathbf{w} and set learning rate α .
- ② **Repeat Until Convergence:**

- ① Compute the gradient:

$$\Delta \mathbf{w} = \frac{\partial J}{\partial \mathbf{w}} = \frac{1}{m} \mathbf{X}^T (\mathbf{X}\mathbf{w} - \mathbf{y})$$

- ② Update the weights:

$$\mathbf{w} = \mathbf{w} - \alpha \Delta \mathbf{w}$$

- ③ **Termination:** The algorithm converges when the cost function decreases by an amount smaller than a pre-defined threshold or after a pre-defined number of iterations.



General Supervised Learning Algorithm

- Step 1 : Data Collection – Collect Data in (x, y) pairs
- Step 2 : Parameter Initialization – Guess for w
- Step 3 : Forward Pass – For all data points in the dataset find \hat{y} using a hypothesis Model
- Step 4 : Gradient Descent – Improve w using gradient descent. $w_j = w_j - \alpha \frac{\partial J}{\partial w_j}$
- Step 5 : Check exit criterion.
 - If not satisfied, go to Step 3 with the updated parameters

$$\mathbf{x}^{(1)} = (x_1^{(1)}, x_2^{(1)}, \dots, x_n^{(1)})$$

no. of
data

$$\mathbf{y}^{(1)} = (y_1^{(1)}, y_2^{(1)}, \dots, y_K^{(1)})$$

Class

$$\hat{\mathbf{y}}^{(1)} = (\hat{y}_1^{(1)}, \hat{y}_2^{(1)}, \dots, \hat{y}_K^{(1)})$$

\mathbf{x}	\mathbf{y}	$\hat{\mathbf{y}}$	J
$\mathbf{x}^{(1)}$	$\mathbf{y}^{(1)}$	$\hat{\mathbf{y}}^{(1)}$	$J^{(1)}$
$\mathbf{x}^{(2)}$	$\mathbf{y}^{(2)}$	$\hat{\mathbf{y}}^{(2)}$	$J^{(2)}$
$\mathbf{x}^{(3)}$	$\mathbf{y}^{(3)}$	$\hat{\mathbf{y}}^{(3)}$	$J^{(3)}$
$\mathbf{x}^{(4)}$	$\mathbf{y}^{(4)}$	$\hat{\mathbf{y}}^{(4)}$	$J^{(4)}$
$\mathbf{x}^{(5)}$	$\mathbf{y}^{(5)}$	$\hat{\mathbf{y}}^{(5)}$	$J^{(5)}$
...
$\mathbf{x}^{(m)}$	$\mathbf{y}^{(m)}$	$\hat{\mathbf{y}}^{(m)}$	$J^{(m)}$
			$J = \sum J^{(i)}$

Linear Regression – Learning Algorithm

- Step 1 : Data Collection – Collect Data in (x, y) pairs
- Step 2 : Parameter Initialization – Guess for w
- Step 3 : Forward Pass – For all data points in the dataset
find \hat{y} using
 - The Linear Model $y^{(j)} = \sum_{i=0}^n w_i x_i^{(j)}$
- Step 4 : Gradient Descent – Improve w using gradient descent. $w_j = w_j - \alpha \frac{\partial J}{\partial w_j}$
 - Gradient $\frac{\partial J}{\partial w_j} = \frac{1}{m} \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)}) x_j^{(i)}$
- Step 5 : Check exit criterion.
 - If not satisfied, go to Step 3 with the updated parameters

x	y	\hat{y}	J
$x^{(1)}$	$y^{(1)}$	$\hat{y}^{(1)}$	$J^{(1)}$
$x^{(2)}$	$y^{(2)}$	$\hat{y}^{(2)}$	$J^{(2)}$
$x^{(3)}$	$y^{(3)}$	$\hat{y}^{(3)}$	$J^{(3)}$
$x^{(4)}$	$y^{(4)}$	$\hat{y}^{(4)}$	$J^{(4)}$
$x^{(5)}$	$y^{(5)}$	$\hat{y}^{(5)}$	$J^{(5)}$
...
$x^{(m)}$	$y^{(m)}$	$\hat{y}^{(m)}$	$J^{(m)}$

$J = \frac{1}{m} \sum J^{(i)}$

Practical Consideration – Batching Gradient Descent

Gradient Update

$$w_j = w_j - \alpha \frac{\partial J}{\partial w_j}$$

$m = 1000$ ↑
 $b = 50$, ↑
 20 updates

x	y	\hat{y}	J
$\mathbf{x}^{(1)}$	$\mathbf{y}^{(1)}$	$\hat{\mathbf{y}}^{(1)}$	$J^{(1)}$
$\mathbf{x}^{(2)}$	$\mathbf{y}^{(2)}$	$\hat{\mathbf{y}}^{(2)}$	$J^{(2)}$
$\mathbf{x}^{(3)}$	$\mathbf{y}^{(3)}$	$\hat{\mathbf{y}}^{(3)}$	$J^{(3)}$
$\mathbf{x}^{(4)}$	$\mathbf{y}^{(4)}$	$\hat{\mathbf{y}}^{(4)}$	$J^{(4)}$
$\mathbf{x}^{(5)}$	$\mathbf{y}^{(5)}$	$\hat{\mathbf{y}}^{(5)}$	$J^{(5)}$
...
$\mathbf{x}^{(m)}$	$\mathbf{y}^{(m)}$	$\hat{\mathbf{y}}^{(m)}$	$J^{(m)}$

$$\frac{\partial J}{\partial w_j} = \frac{1}{b} \sum \frac{\partial J^{(i)}}{\partial w_j}$$

b is the size of the batch

- Default option – Each update in gradient descent requires to see the whole dataset
 - That is, one update requires us to average the updates of every point in the whole dataset
 - This is called **Batch Gradient Descent** where $b = m$. **One epoch has one update**
 - Problems : (i) It is slow (ii) It is expensive for each update (iii) Data cannot be dynamic
- **Stochastic Gradient Descent** ↗
 - $b = 1$. **One epoch has m updates**. Many more updates per epoch
 - Can handle dynamic incoming data but loss convergence is more erratic
- **Mini-Batch Gradient Descent**
 - $1 < b < m$. **One epoch has $\frac{m}{b}$ updates**. More than 1 update per epoch
 - Loss convergence more erratic than Batch and less erratic than stochastic

Data Normalization and Its Impact on Gradient Descent for Linear Regression

Introduction

Data normalization and scaling, key preprocessing steps in machine learning, standardize the range of independent variables or features of the data. Especially for linear regression with gradient descent, these techniques can significantly enhance convergence speed and stability.

Machine Precision $\rightarrow 10^{-16}$

What is Data Normalization?

Normalization scales the dataset to have a mean of zero and a standard deviation of one. For a feature x , its normalized version x' is:

$$x' = \frac{x - \mu}{\sigma}$$

Batch norm in NN

Where:

- μ is the mean of the feature.
- σ is the standard deviation of the feature.



What is Min-Max Scaling?

Min-Max scaling transforms features to lie between 0 and 1. For a feature x , its scaled version x' is:

$$x' = \frac{x - x_{\min}}{x_{\max} - x_{\min}}$$

Where:

- x_{\min} is the minimum value of the feature.
- x_{\max} is the maximum value of the feature.

Why Normalize or Scale Data for Gradient Descent?

- **Faster Convergence:** Different feature scales result in skewed contour plots, leading to slower convergence.
- **Avoidance of Numerical Instability:** Large magnitude features can induce numerical issues.
- **Improved Model Interpretability:** Scaled features lead to more interpretable coefficients.
- **Consistent Regularization Impact:** Equal treatment of all features by regularization techniques.

Impact on Gradient Descent in Linear Regression

The cost function J in linear regression is:

$$J = \frac{1}{2m} \sum (\hat{y}^i - y^i)^2$$

Feature scales affect gradient magnitudes:

- Different scales lead to oscillatory behavior.
- Normalized or scaled data ensures balanced gradients and uniform weight updates.

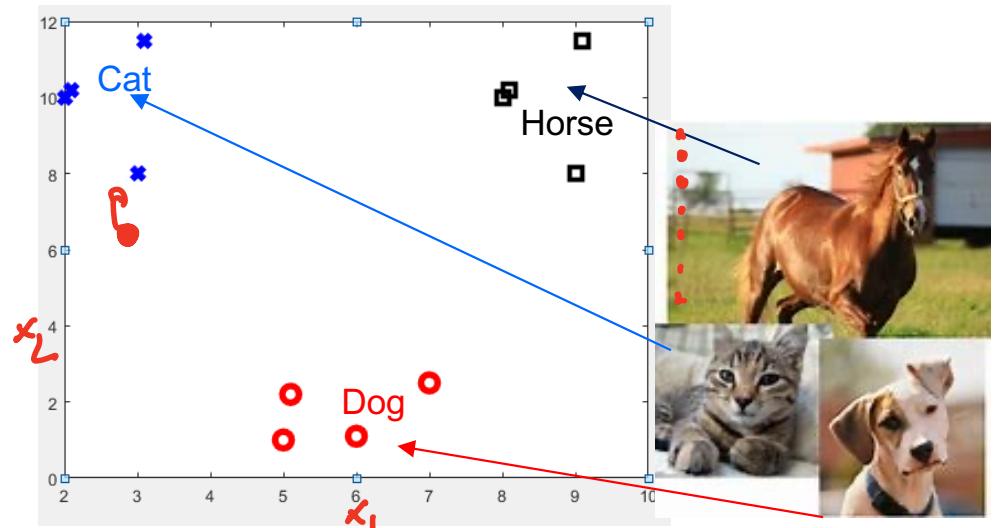
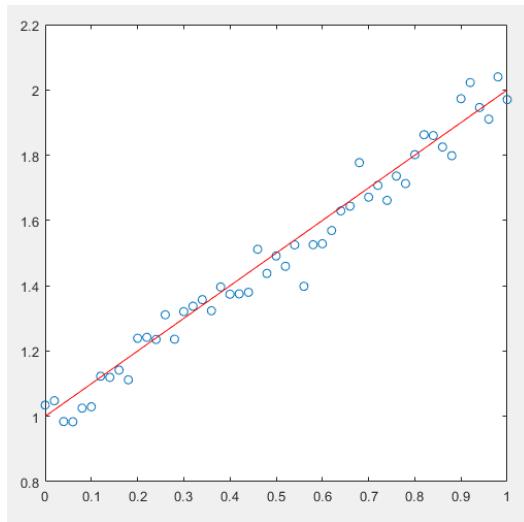
Conclusion

- Data normalization and scaling are crucial for efficient gradient descent in linear regression.
- Both techniques ensure faster and more reliable convergence.
- While multiple methods exist, normalization and Min-Max scaling are among the most popular for linear regression tasks.

Common non-linearities

Classification

Two problems in Supervised Learning



Regression	Classification
Fit it	Split it
Real number data	Discrete or Categorical data.
Has associated number	Has category associated
Example : Prediction of stock market	Example : Tumor classification
Model : Linear Model	Model : Logistic/Softmax

Binary Classification

Multiclass Classification

Definition of Machine Learning

- Simple Definition -- **Using Data** to answer questions
- Study of computer algorithms
 - that improve automatically
 - through experience.
- Formally, 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 .

- ML Recipe – Combine
 - ✓ Dataset representation
 - ✓ Cost function
 - ✓ Optimization procedure
 - ✓ Model

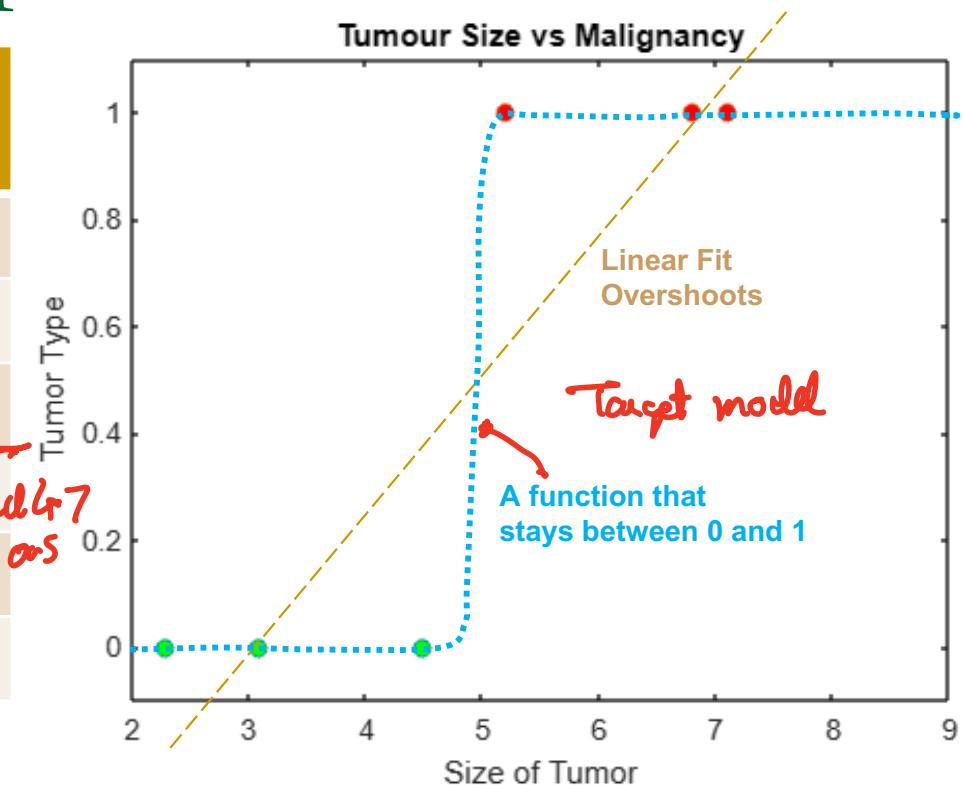
A simple binary classification problem

- Binary classification – Distinguishing between two classes
 - Benign/Malignant, Cat/Dog, Spam/Not Spam
- Imagine we are to characterize a tumor as benign or cancerous based on tumor size alone. Some prior data is given and we have to build a predictive model now → ✗
- Step 1 : Data Representation
 - Use numerical values for output – Benign = 0, Cancer = 1

Tumor Size (x in mm)	Outcome (y)
2.3	0
3.1	0
4.5	0
5.2	1
6.3	1
7.1	1

The Forward Model

Tumor Size (x in mm)	Outcome (y)
2.3	0
3.1	0
4.5	0
5.2	1
6.8	1
7.1	1

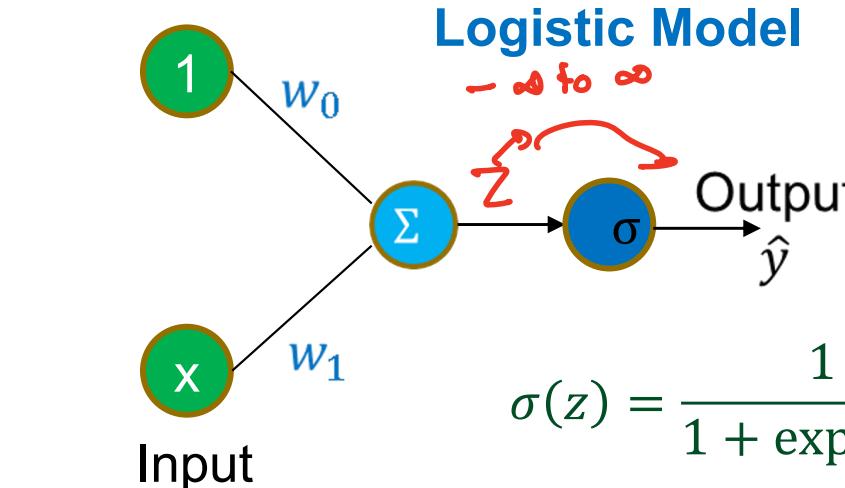
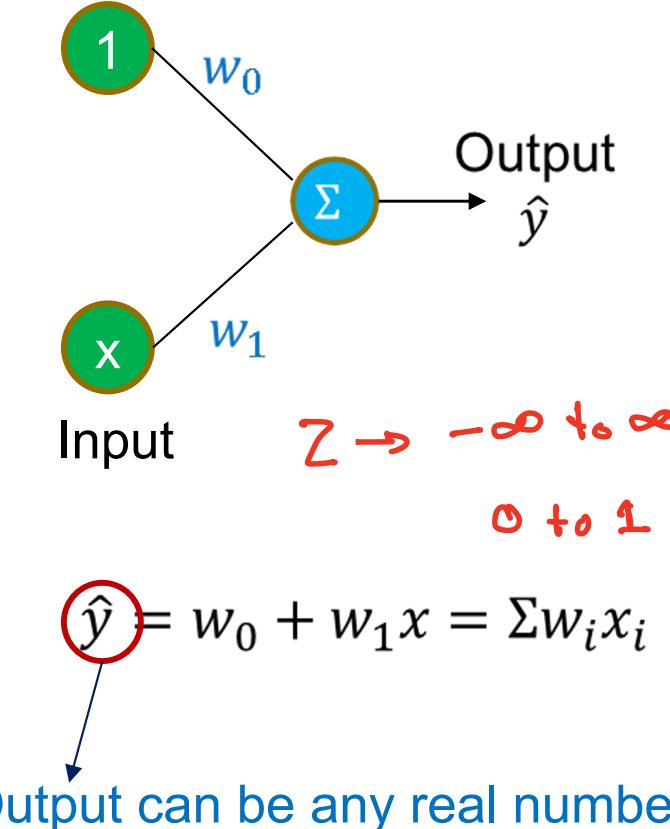


Step 2 : Create Forward Model

- We need a model that respects the data to be replicated
- We want the output to always be either 0 or 1
 - Or at least lie between 0 and 1
- A simple linear model $\hat{y} = w_0 + w_1 x$ will not obey this.
- So, we try the logistic model

From Linear to Logistic

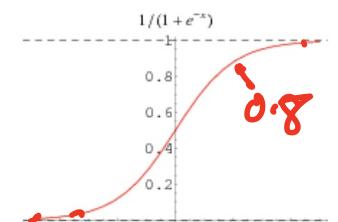
Linear Model



Logit/Linear Activation

$$z = w_0 + w_1 x = \sum w_i x_i$$

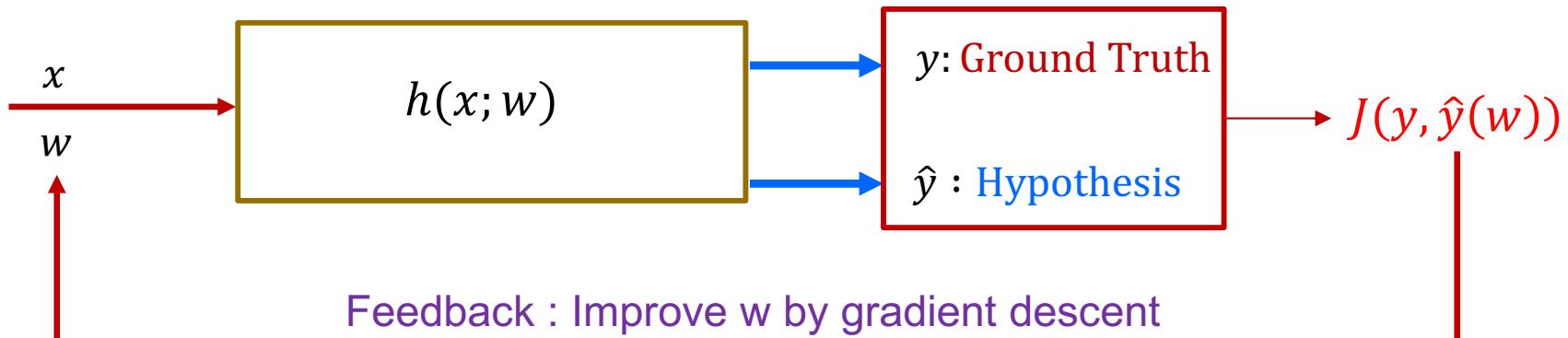
$\hat{y} = \sigma(z)$



$\hat{y} = \sigma(w_0 + w_1 x) \rightarrow$ Math model
Output will lie between 0 and 1 and can model probability $p(y=1 | x)$

Interpretation

Completing the learning recipe



- ML Recipe for Logistic Regression for classification
 - ✓ Dataset representation – x, y Numerical representation
 - ✓ Model -- $\hat{y} = h(x; w)$
 - ✗ Cost function -- $J(y, \hat{y})$
 - ✗ Optimization procedure

$y \rightarrow \text{Either } 0 \text{ or } 1$

- The **cost function** $J(y, \hat{y})$ needs to satisfy the following properties

1. Correct Classifications should be as cheap as possible

That is, if $\hat{y} = y$ then $J = 0$

Least Square Loss satisfies this

2. Incorrect Classifications should be as expensive as possible

If $y = 1$ but $\hat{y} \rightarrow 0$ then $J \rightarrow \infty$

If $y = 0$ but $\hat{y} \rightarrow 1$ then $J \rightarrow \infty$

Least Square Loss
 $(y - \hat{y})^2$

We need these to
finish the recipe
Gradient descent

\hat{y} : Between
0 & 1

Least Square Loss does not
satisfy this

NOTE: Least square still works, despite this problem, but is slower compared to the Loss function we will use

Binary Cross-Entropy (BCE) Loss Function

- Recall that $J(y, \hat{y})$ the loss function quantifies the gap between the model output and the true output
- The Least Square error $J^{LSE}(y, \hat{y}) = (y - \hat{y})^2$
 - Can have a maximum value of only 1 for any given data point
 - Does not penalize incorrect classifications heavily enough
 - Derived from Gaussian distribution which is not accurate for binary data

- So, we introduce the binary cross entropy loss function

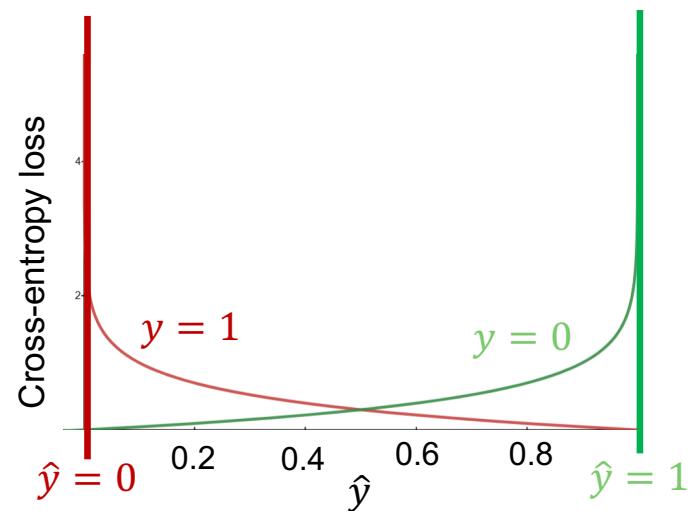
$$J^{BCE}(y, \hat{y}) = -y \ln \hat{y} - (1 - y) \ln (1 - \hat{y})$$

- This has the properties that
 - $J = 0$ if $y = \hat{y}$
 - $J \rightarrow \infty$ if y is totally misclassified
 - J is never negative
 - We will show these in the next slide

Binary Cross Entropy Properties

$$J^{BCE}(y, \hat{y}) = \begin{cases} -\log \hat{y} & \text{if } y = 1 \\ -\log(1 - \hat{y}) & \text{if } y = 0 \end{cases}$$

$$J^{BCE}(y, \hat{y}) = -y \log \hat{y} - (1 - y) \log(1 - \hat{y})$$

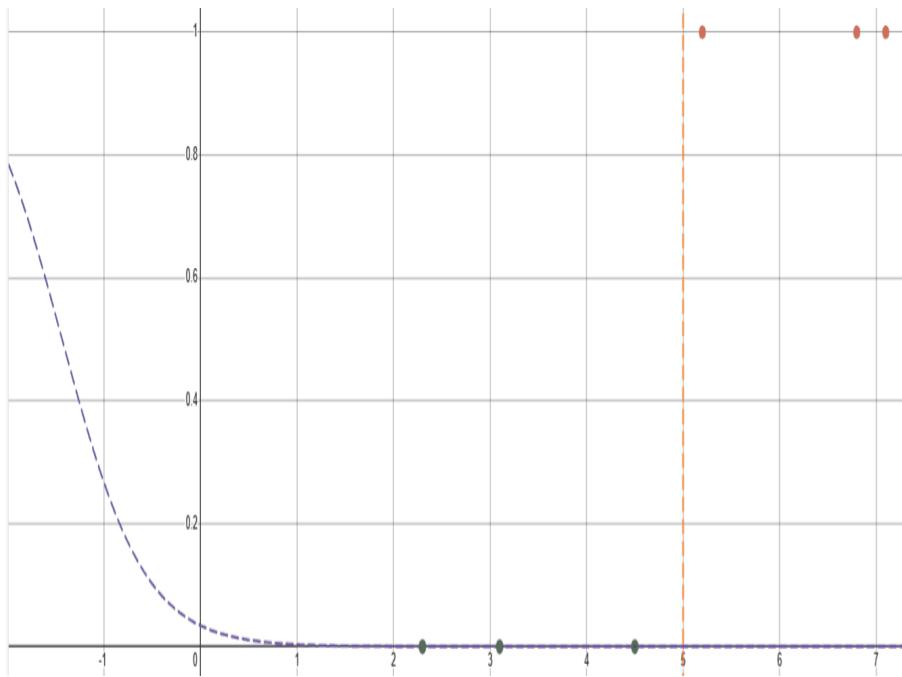
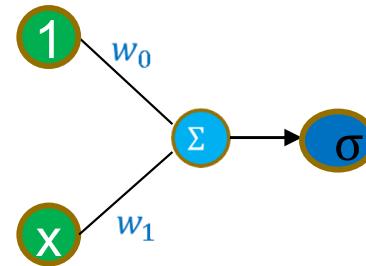


Properties of \ln

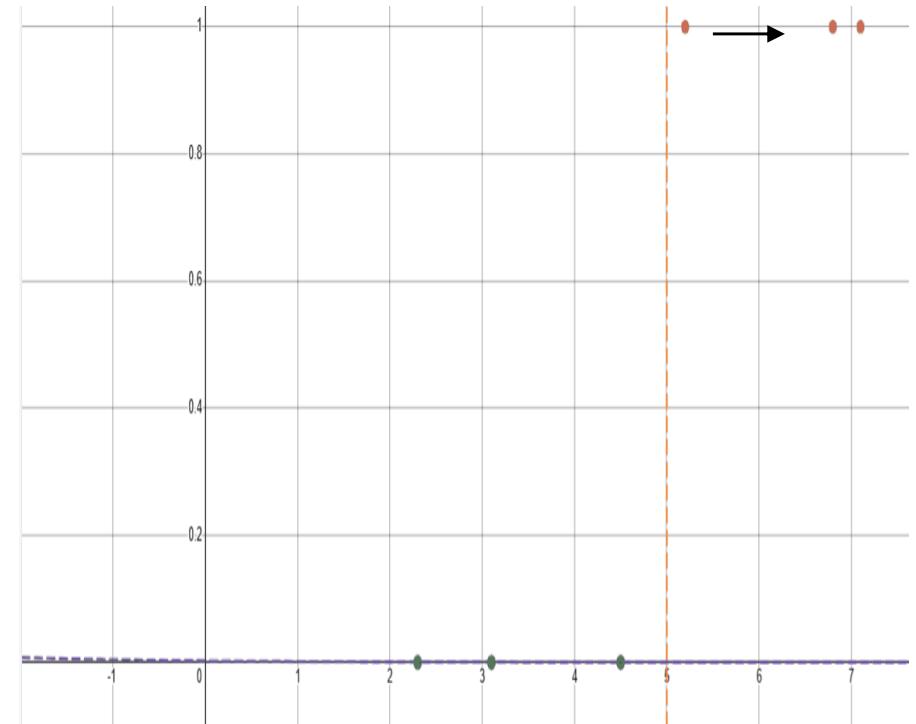
- Logarithm to the base e. Inverse of the exponential. That is, $\log(e^x) = 1$
- We can see that $\log(1) = 0$ and $\log(0) \rightarrow -\infty$

S.No	y	\hat{y}	$J(y, \hat{y}) = -y \log \hat{y} - (1 - y) \log(1 - \hat{y})$	Comments
1	1	≈ 1	$\approx -1 \times \log(1) - 0 \rightarrow 0$	Correct classification incurs no cost
2	0	≈ 0	$\approx -0 - 1 \times \log(1) \rightarrow 0$	Correct classification incurs no cost
3	1	≈ 0	$\approx -1 \times \log(0) - 0 \rightarrow \infty$	Incorrect classification incurs infinite cost
4	0	≈ 1	$\approx -0 - 1 \times \log(0) \rightarrow \infty$	Incorrect classification incurs infinite cost

Playing with the sigmoid – Varying the parameters



Changing w_0 with fixed w_1



Changing w_0 with fixed w_1

$$z = w_0 + w_1 x \quad \hat{y} = \sigma(z)$$

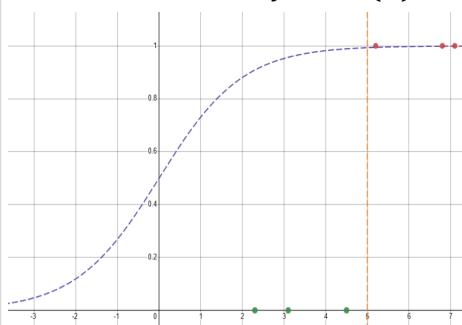
Changing w_0 seems to change the critical point and w_1 affects the sharpness

Playing with the sigmoid

$$z = w_0 + w_1 x \quad \hat{y} = \sigma(z)$$

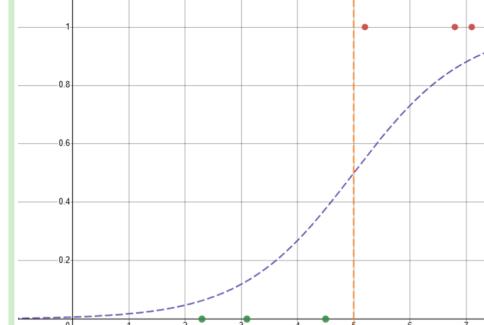
Let us try various values of the parameters to interpret what is happening

$$w_0 = 0, \quad w_1 = 1 \\ z = x \quad \hat{y} = \sigma(x)$$



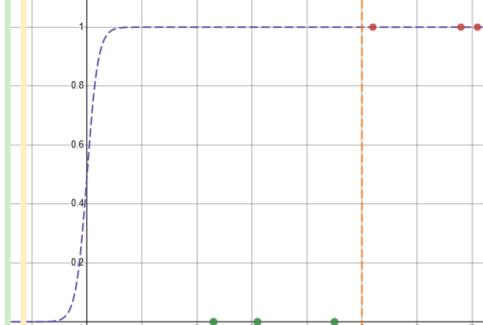
Case 1

$$w_0 = -5, \quad w_1 = 1 \\ z = x - 5 \quad \hat{y} = \sigma(x - 5)$$



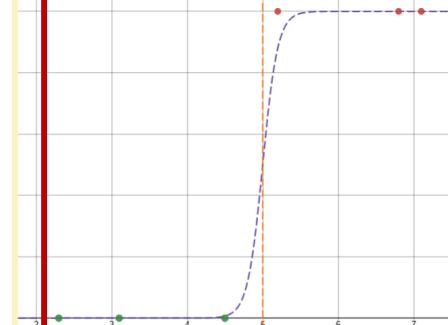
Case 2

$$w_0 = 0, \quad w_1 = 10 \\ z = 10x \quad \hat{y} = \sigma(10x)$$



Case 3

$$w_0 = -50, \quad w_1 = 10 \\ \sigma(10x - 50)$$



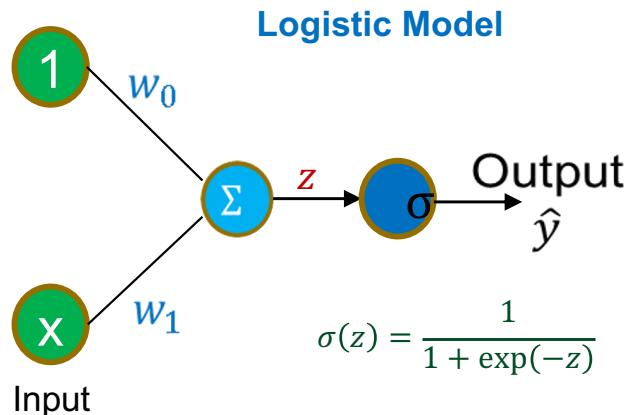
Case 4

Case #	w_0	w_1	z	Function \hat{y}	Comments
1	0	1	x	$\sigma(x)$	Does not capture the critical point and not sharp
2	-5	1	$x - 5$	$\sigma(x - 5)$	Captures critical point but not sharp
3	0	10	$10x$	$\sigma(10x)$	Sharp but underpredicts critical point
4	-50	10	$10x - 50$	$\sigma(10x - 50)$	Sharp and captures critical point reasonably

Let us now interpret the term $\sigma(z)$ and the logistic model so that we can systematize this

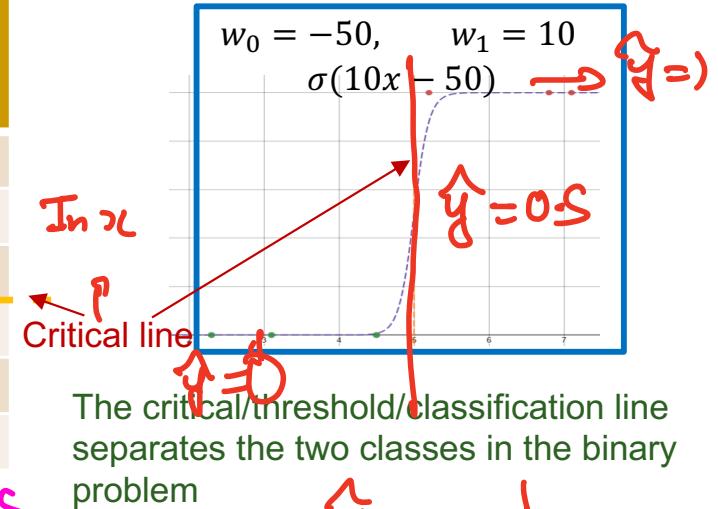
Interpreting the logistic model

$$z = w_0 + w_1 x \quad \hat{y} = \sigma(z)$$



Tumor Size (x in mm)	Outcome (y)
2.3	0
3.1	0
4.5	0
5.2	1
6.8	1
7.1	1

$$z_c = 0 \rightarrow \hat{y}_c = 0.5$$



$$\hat{y} = \frac{1}{1 + \exp(-z)}$$

0.5

$z_c = 0$

The criterion for the classifying line in terms of z (the network variable)

Far to the left of it we have $\hat{y} \rightarrow 0$ which means $z \rightarrow -\infty$

Far to the right of it we have $\hat{y} \rightarrow 1$ which means $z \rightarrow \infty$

- Right at the classifying line, we will have $\hat{y} \rightarrow 0.5$, i.e. $z_c = 0$

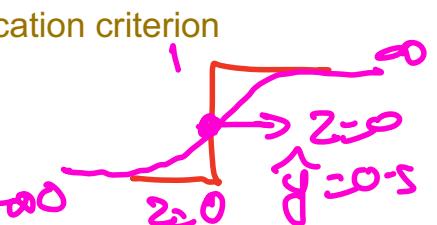
The criterion for the classifying line in terms of x (the physical variable)

$$z = w_0 + w_1 x \Rightarrow z_c = w_0 + w_1 x_c \Rightarrow x_c = (z_c - w_0)/w_1$$

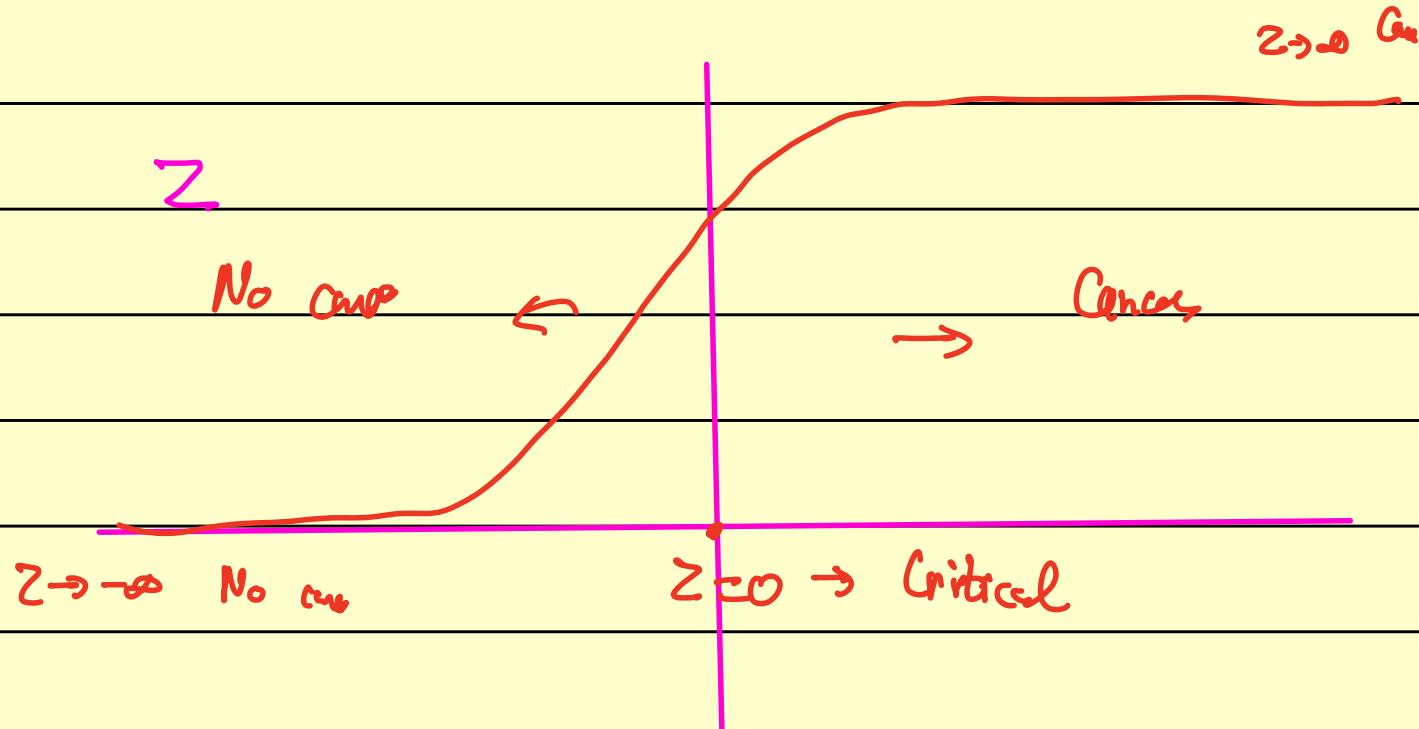
$$\text{Since } z_c = 0 \Rightarrow x_c = -w_0/w_1$$

We obtained $w_0 \approx -50, w_1 \approx 10$ manually $\Rightarrow x_c \approx -5$ which matches the data

Instead of obtaining these manually, can we automatically get these parameters via the “learning” process?



$$z = w_0 + w_1 x \rightarrow \text{Change of variable}$$



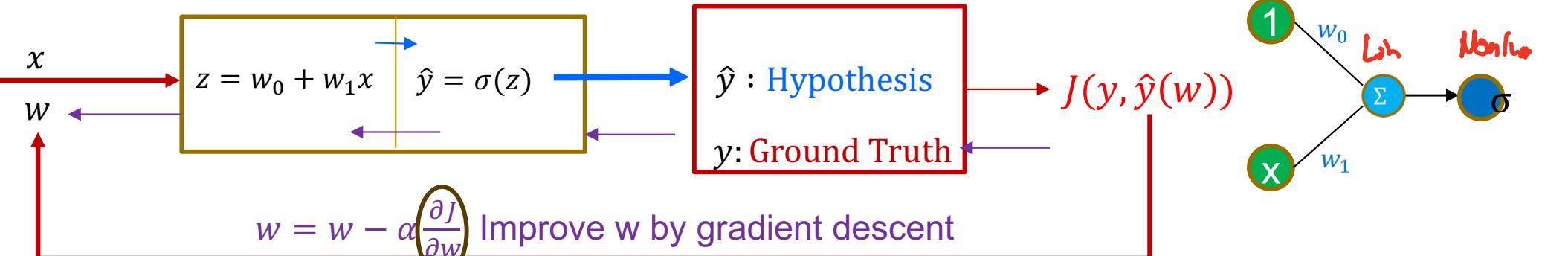
$w_0, w_1 \rightarrow$ Gradient Descent

$$x = 3 \text{ mm}$$

$$w_0 + w_1 x = \vec{z} = 1.1$$

$$\hat{y} = \sigma(z) = \sigma(1.1) \approx 0.83$$

Gradient Operator for logistic regression



We need $\frac{\partial J}{\partial w_0}$ and $\frac{\partial J}{\partial w_1}$.

Reverse the order of calculations to calculate these derivatives

Forward Prop $\rightarrow x; w \rightarrow z \rightarrow \hat{y} \rightarrow J$
 Backward Prop $\rightarrow \Delta J \rightarrow \Delta \hat{y} \rightarrow \Delta z \rightarrow \Delta w$

Calculating $\frac{\partial J}{\partial \hat{y}}$

$$\text{Use } \frac{d}{dx} \log(x) = \frac{1}{x}$$

$$J = -y \log \hat{y} - (1-y) \log(1-\hat{y})$$

$$\frac{\partial J}{\partial \hat{y}} = -\frac{y}{\hat{y}} + \frac{(1-y)}{1-\hat{y}}$$

$$\frac{\partial J}{\partial w_0} = \frac{\partial J}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z} \frac{\partial z}{\partial w_0}$$

$$\frac{\partial J}{\partial w_1} = \frac{\partial J}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z} \frac{\partial z}{\partial w_1}$$

Common term $\frac{\partial J}{\partial z}$

$$\Rightarrow \frac{\partial J}{\partial z} = \frac{\partial J}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z} = \hat{y} - y$$

$$\frac{\partial J}{\partial w_0} = \frac{\partial J}{\partial z} \frac{\partial z}{\partial w_0} = (\hat{y} - y) \times 1$$

$$\frac{\partial J}{\partial w_1} = \frac{\partial J}{\partial z} \frac{\partial z}{\partial w_1} = (\hat{y} - y) \times x$$

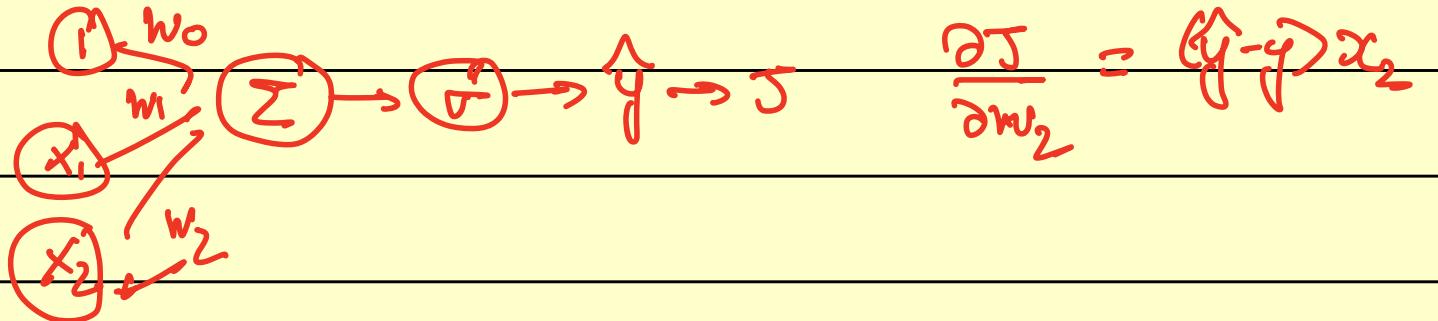
Calculating $\frac{\partial \hat{y}}{\partial z} \Rightarrow \hat{y} = \frac{1}{1 + \exp(-z)}$

$$\Rightarrow \frac{d\hat{y}}{dz} = \frac{\exp(-z)}{(1 + \exp(-z))^2}$$

$$\Rightarrow \frac{d\hat{y}}{dz} = \frac{\exp(-z)}{(1 + \exp(-z))} \times \frac{1}{(1 + \exp(-z))}$$

$$\Rightarrow \frac{d\hat{y}}{dz} = \hat{y} \times (1 - \hat{y})$$

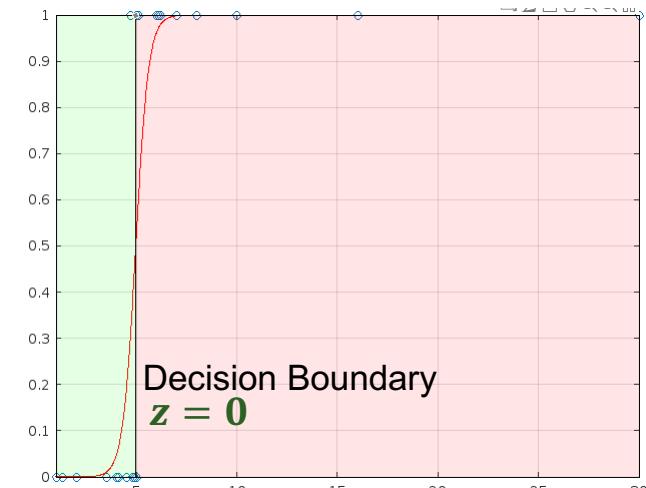
DELT A RULE



Binary Classification Algorithm

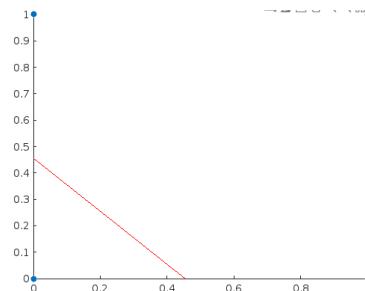
- Step 1 : Data Collection – Collect Data in (x, y) pairs
 - Denote y only as 0 or 1 (Binary)
- Step 2 : Parameter Initialization – Guess for w
 - For a 1 feature model, this is w_0 and w_1
- Step 3 : Forward Pass – For all data points in the dataset find \hat{y}
 - For a 1 feature model, this is $\hat{y}_i = \sigma(w_0 + w_1 x_i)$
- Step 4 : Gradient Descent – Improve w using gradient descent
 - $w_j = w_j - \alpha \frac{\partial J}{\partial w_j}$
 - For a 1 feature model, $\frac{\partial J}{\partial w_j} = (\hat{y} - y) \times x_j$
- Step 5 : Check exit criterion.
 - If not satisfied, go to Step 3 with the updated parameters

Tum or Size (x in mm)	True Outcome (y)	Model Outcome \hat{y}	Data Loss J_i
2.3	0		
3.1	0		
4.5	0		
5.2	1		
6.8	1		
7.1	1		



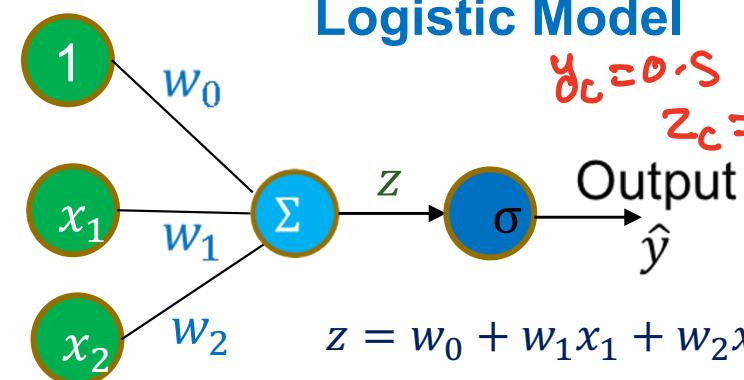
Binary Classification for more than one feature

	x_1	x_2	y
1	0	0	0
2	1	0	1
3	0	1	1
4	1	1	1



Logistic Model

$$y_c = 0 \cdot S \\ z_c = 0$$



$$z = w_0 + w_1 x_1 + w_2 x_2 = \sum w_i x_i$$

OR gate as a classification problem

- The technique for modeling more than 1 features remains the same.

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

- For two features this becomes

$$z = w_0 + w_1 x_1 + w_2 x_2 \quad \hat{y} = \sigma(z)$$

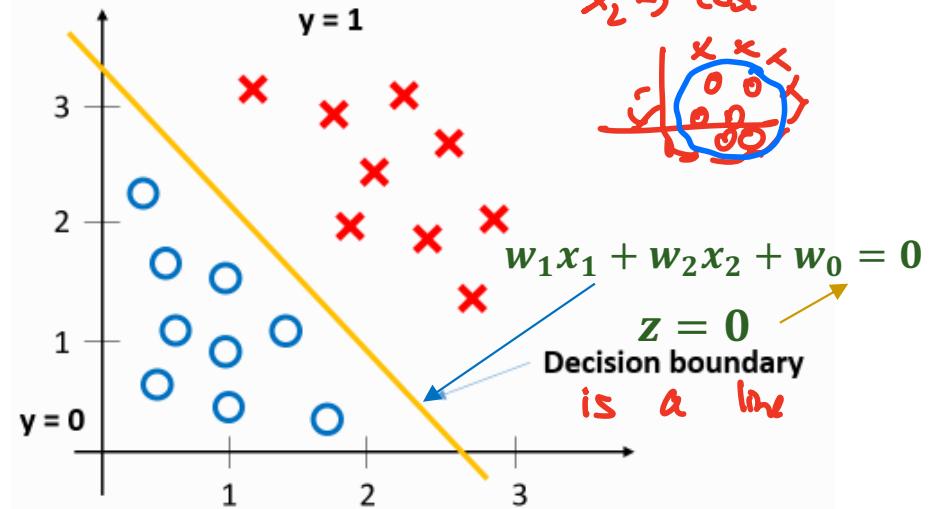
- The decision boundary remains $z = 0$
 - This means that the **boundary is always a line**
- Gradients follow the “Delta Rule”

$$\frac{\partial J}{\partial w_j} = (\hat{y} - y) \times x_j$$

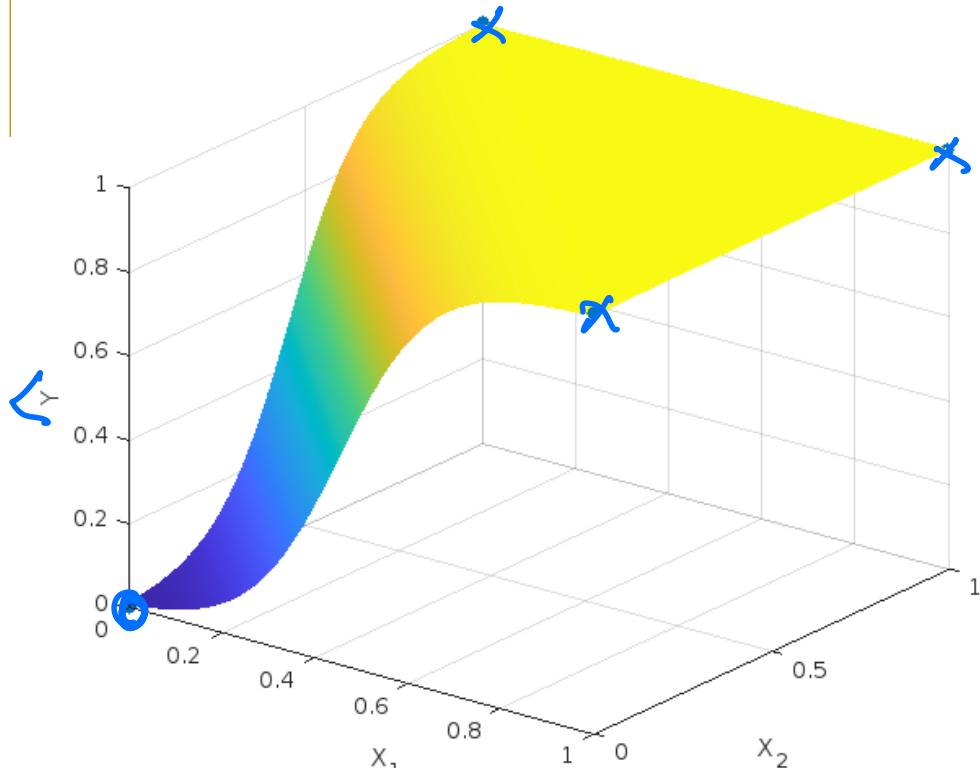
$$\text{Input } Z_C \quad \hat{y} = \sigma(z)$$

$$\underline{w_0 + w_1 x_1 + w_2 x_2 = 0} \quad \text{Locus of pts}$$

$x_1 \rightarrow \text{years}$
 $x_2 \rightarrow \text{cost}$

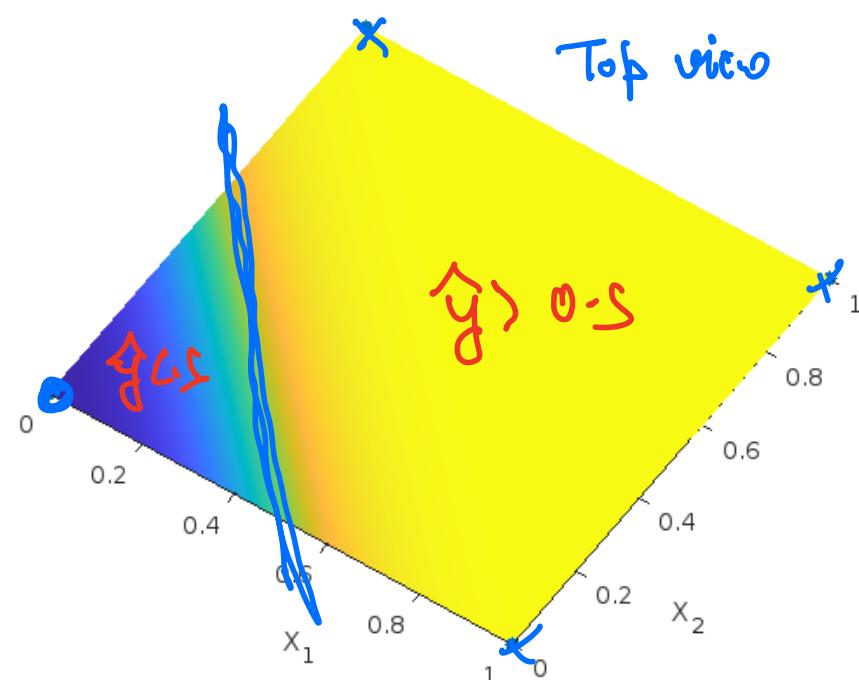


Log. Reg work only if data is linearly separable



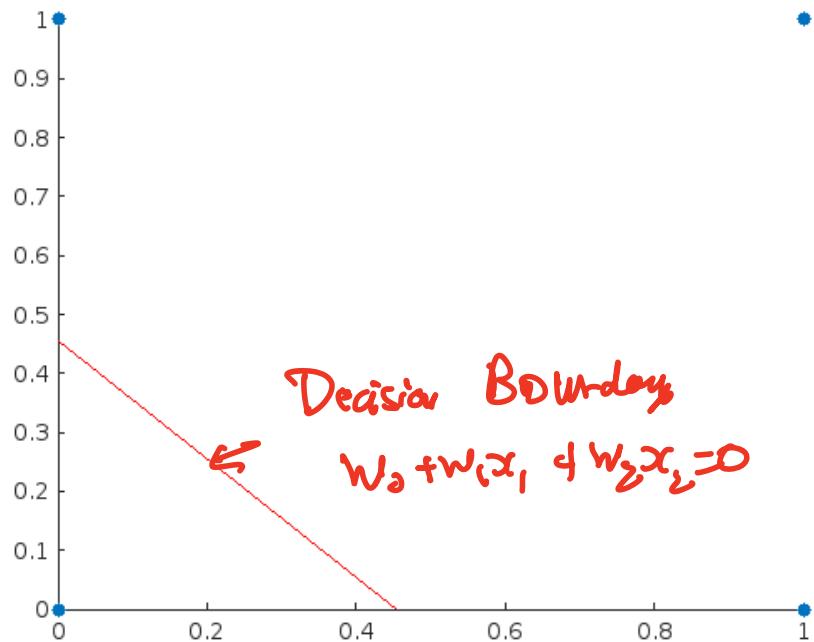
$$m=4$$

	x_1	x_2	y
1	0	0	0
2	1	0	1
3	0	1	1
4	1	1	1



$$z = w_0 + w_1x_1 + w_2x_2 = \sum w_i x_i$$

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

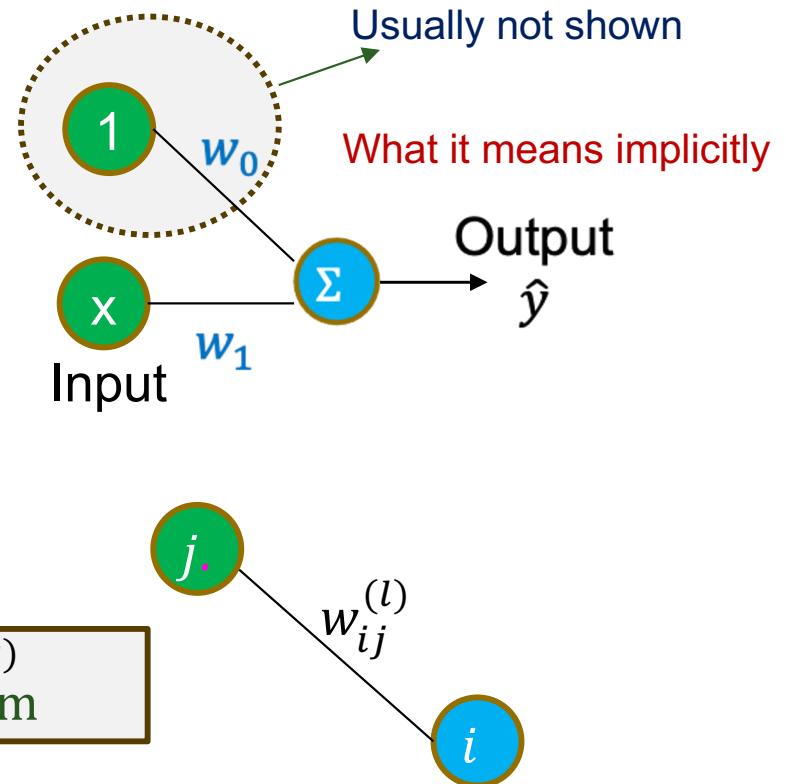
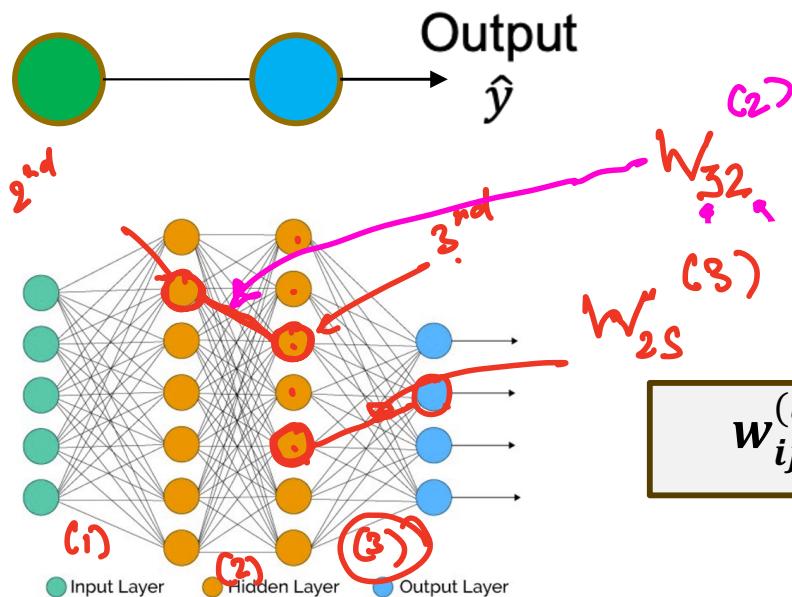


MULTICLASS CLASSIFICATION

(Multinomial Classification)

Weights notation

What is shown explicitly



- Constant neurons representing a “feature” with 1 are usually not depicted in the pictorial representation.
- The parameters going from the constant neuron is called a bias and the neuron is called a “bias unit”
- The other parameters are called “weights”

Multiclass classification



$K = 2$ Binary Classification



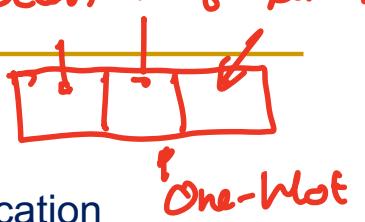
$K = 3$ Multiclass Classification

- **K : Number of classes.**
 - $K = 2$ Binary classification
 - $K > 2$ Multiclass (Multinomial) classification
- **ML Recipe – Requires the following**
 - **Data Representation** – How do we represent the output (it cannot be 0 or 1)
 - **Model** – What is the nonlinearity instead of sigmoid? (That was used to lie between 0 and 1)
 - **Loss** – What is the Loss function?
 - **Optimization** – Gradient Descent but What is the gradient for the new model?

Let us look at these issues one by one

$b(\text{Cat}) \quad p(\text{Dog}) \quad b(\text{Horse})$

Data Representation – One Hot Vector



$K = 2$ Binary Classification



$y = 0$

Binary Representation

Scalar

$y = 1$

$p(\text{Dog})$

$K = 3$ Multiclass Classification



1	0	0
---	---	---

Ground Truth y

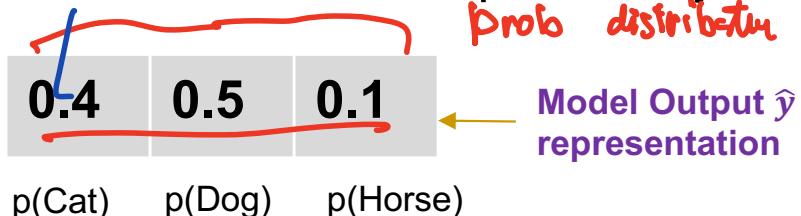
0	1	0
---	---	---

One Hot Representation

0	0	1
---	---	---

Image Not

- “One-Hot Representation” – Represent the output as a vector instead of a scalar
- Each output class is represented by a 1 in one location and 0 in others – One Hot
- We can interpret the entries in each box as the probability of the class

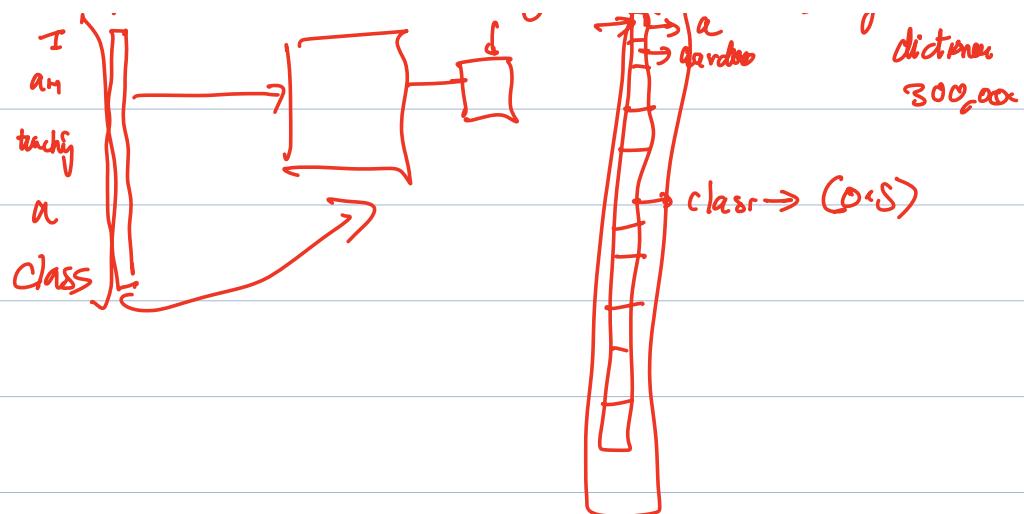


For both y and \hat{y} , every entry represents $p(y_i = 1|x)$

x

y

$K = ?$ Size of



$$z_1 = w_1 x_1 + w_2 x_2 + w_3 x_3 + w_4 x_4 \quad \hat{y}_1 = \sigma(z_1)$$

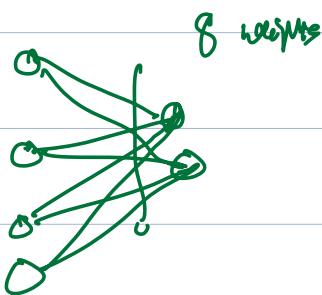
$$\hat{y}_1 \in \{0, 1\}$$

$$z_2 = w_5 x_1 + \dots + w_8 x_4 \quad \hat{y}_2 = \sigma(z_2)$$

$$\hat{y}_2 \in \{0, 1\}$$

$$z_3 = w_9 x_1 + w_{10} x_4 \quad \hat{y}_3 = \sigma(z_3)$$

$$\hat{y}_3 \in \{0, 1\}$$



$$\mathcal{J}(\hat{y}) =$$

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

Model Nonlinearity -- Softmax



0.4 0.5 0.1

p(Cat) p(Dog) p(Horse)

$$\hat{y} = [0.4 \quad 0.5 \quad 0.1]$$

$$= [\hat{y}_1 \quad \hat{y}_2 \quad \hat{y}_3]$$

	z	$\sigma(z)$	Softmax(z)
z_1	1	0.7311	0.1191.
z_2	3	0.9526	0.8801.
z_3	-4	0.0180	0.0008.
SUM Σ		1.7016	1.0000

Not = 1

- We can define $z = \sum w_i x_i$ as before
- However, we cannot use $\hat{y} = \sigma(z)$ since we would not have $\sum \hat{y}_i = 1$
- To guarantee that the output can be interpreted as a probability we introduce softmax. $sf_1 = \frac{e^{z_1}}{D}, sf_2 = \frac{e^{z_2}}{D}, sf_3 = \frac{e^{z_3}}{D}$ $sf_1 + sf_2 + sf_3 = \frac{D}{D} = 1$
- This is defined as

$$\text{softmax}(z_1) = \frac{\exp(z_1)}{(\exp(z_1) + \exp(z_2) + \exp(z_3))}$$

Normalizing const

- This has all outputs between 0 and 1 as well as $\sum \hat{y}_i = 1$

Categorical Cross-Entropy (CCE) Loss Function

- Recall the binary cross entropy loss function

$$J^{BCE}(y, \hat{y}) = -y \ln \hat{y} - (1 - y) \ln (1 - \hat{y})$$

- We use a generalization of this for the multiclass case

$$J^{CCE}(y, \hat{y}) = -\sum_{i=1}^K (y_i \ln \hat{y}_i)$$

Categorical Cross Entropy

- What happens if apply CCE to the binary case with one-hot?

!



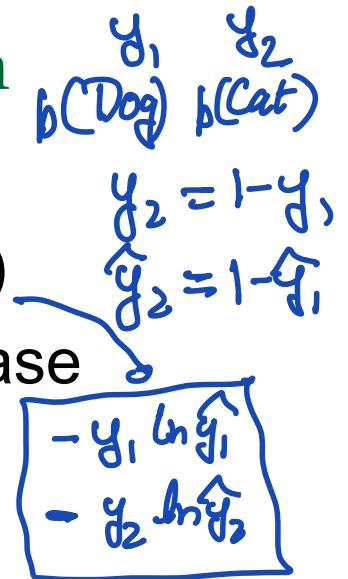
- $J^{CCE}(y, \hat{y}) = -\sum_{i=1}^K (y_i \ln \hat{y}_i) = -y_1 \ln \hat{y}_1 - y_2 \ln \hat{y}_2$
 - In the binary case $\hat{y}_2 = (1 - \hat{y}_1)$ and $y_2 = (1 - y_1)$
- $$\Rightarrow J^{CCE}(y, \hat{y}) = -y_1 \ln \hat{y}_1 - (1 - y_1) \ln (1 - \hat{y}_1)$$
- This is identical to the binary cross entropy loss!

1	Binary	0
1	One-hot	0 1

$$-\sum_{i=1}^2 y_i \ln \hat{y}_i = -[y_1 \ln \hat{y}_1 + y_2 \ln \hat{y}_2]$$

Multiclass

$$\begin{aligned} y_1 &= 1 - \hat{y}_2; & y_2 &= 1 - \hat{y}_1 \\ \hat{y}_1 &= 1 - \hat{y}_2; & \hat{y}_2 &= 1 - \hat{y}_1 \end{aligned}$$



$$= -[y_1 \ln \hat{p}_1 + (1-y_1) \ln (1-\hat{p}_1)]$$

$$= -y_1 \ln \hat{p}_1 - (1-y_1) \ln (1-\hat{p}_1)$$

BCE

$$-y_1 \ln \hat{p}_1 - y_2 \ln \hat{p}_2 - y_3 \ln \hat{p}_3$$

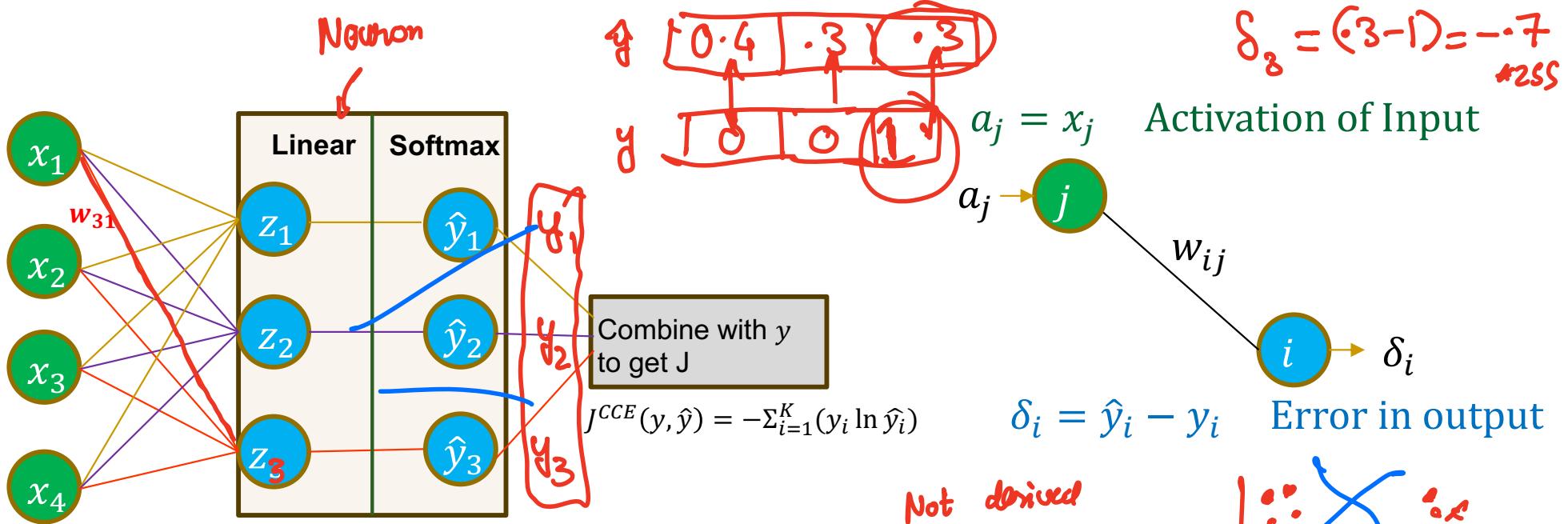
$\underbrace{[1-(y_1+y_2)]}$ $\underbrace{[1-\hat{p}_1-\hat{p}_2]}$

$$0 \leq Sf(z) \leq 1$$

$$\frac{e^{z_i}}{\sum e^{z_i}}$$

$$\frac{\partial J}{\partial w_j} = \frac{\partial J}{\partial z} \frac{\partial z}{\partial z} \frac{\partial z}{\partial w_j}$$

Optimization – Gradients in the Multiclass case



- It can be shown that $\frac{\partial J}{\partial z_1} = \hat{y}_1 - y_1$. Similarly, $\frac{\partial J}{\partial z_i} = \hat{y}_i - y_i$
- Same $\frac{\partial J}{\partial z}$ expression as linear and logistic
- $\frac{\partial J}{\partial w_{31}} =$ Error in Output 3 \times Activation of Input 1 $= \delta_3 \times a_1 = (\hat{y}_3 - y_3) \times x_1$

General Delta Rule

$$\frac{\partial J}{\partial w_{ij}} = \delta_i a_j$$

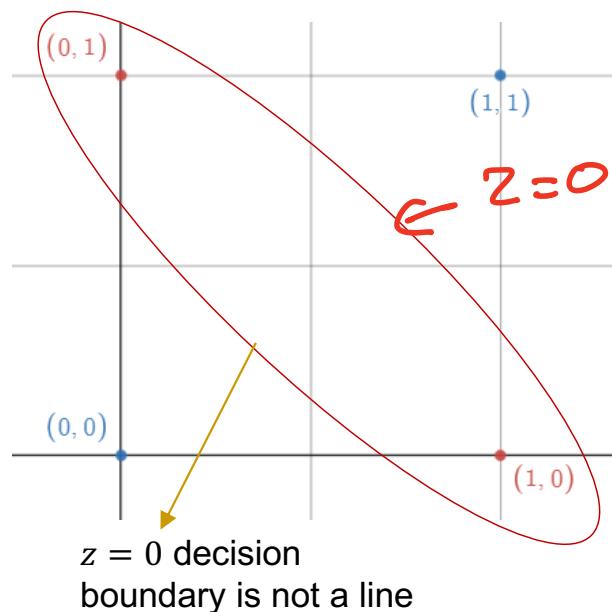
Error in output
Activation of
corresponding input

No hidden layers \rightarrow No nonlinear decision bds

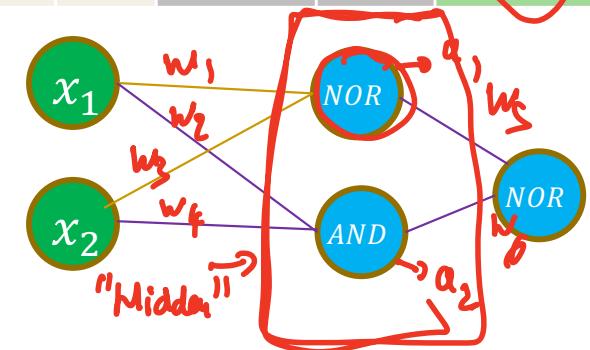
The XOR problem – Need for deep networks

	x_1	x_2	y
1	0	0	0
2	1	0	1
3	0	1	1
4	1	1	0

$$\hat{y} = x_1 + x_2 - 2x_1x_2$$



x_1	x_2	NOR a_1	AND a_2	$NOR(NOR, AND)$
0	0	1	0	0
1	0	0	0	1
0	1	0	0	1
1	1	0	1	0



- In all cases so far, we map the input to $z = w_0 + w_1x_1 + \dots$
- The decision boundary, $z = 0$ will, therefore, always be a line
 - If we directly map $x \rightarrow z$ like so far, it only works for **linearly separable data**
- Two solutions
 - Use nonlinear features. Example -- $\hat{y} = w_0 + w_1x_1 + w_2x_2 + w_3x_1x_2$
 - Use “deep” networks.
- We can use $XOR(x_1, x_2) = NOR(NOR(x_1, x_2), AND(x_1, x_2))$ for example
- But, can we automatically come up with this representation via learning?

$$x_3 = x_1x_2 \text{ Nonlinear feature}$$

Two approaches

(1) Fix the input → Nonlinear features

$$x_1, x_2, x_3$$

SVM → "Kernel" trick

$$x_1 x_2 \ x_2 x_3 \ x_3 x_1$$

Feature Engineering

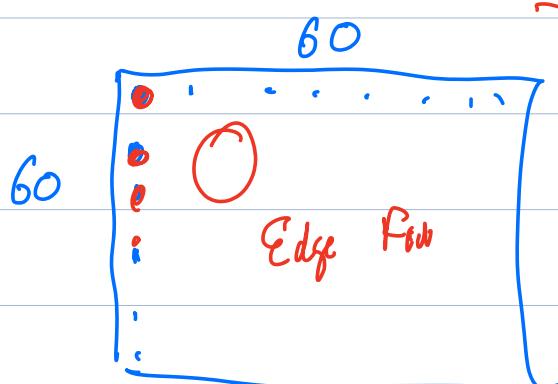
$$x_1^2 \ x_2^2 \ x_3^2$$

Curse of dimensionality

(2) Fix the process → Deep networks

→ Add more "intermediate"

Feature learning



Computato
60x60
3600 weights

3 classes

not linearly separable

will not work
(3600)
1000

1 extra weight

1 feature $w_0 + w_1 x_1 + w_2 x_2$
↪ Quadrat.,

$N^D \rightarrow 3 \rightarrow 6$

3 features

$$x_1^2, x_1 x_2, x_2^2, x_1 x_3, x_2 x_3, x_3^2$$

2 factors $x_1, x_2 \rightarrow w_0 + w_1 x_1 + w_2 x_2$

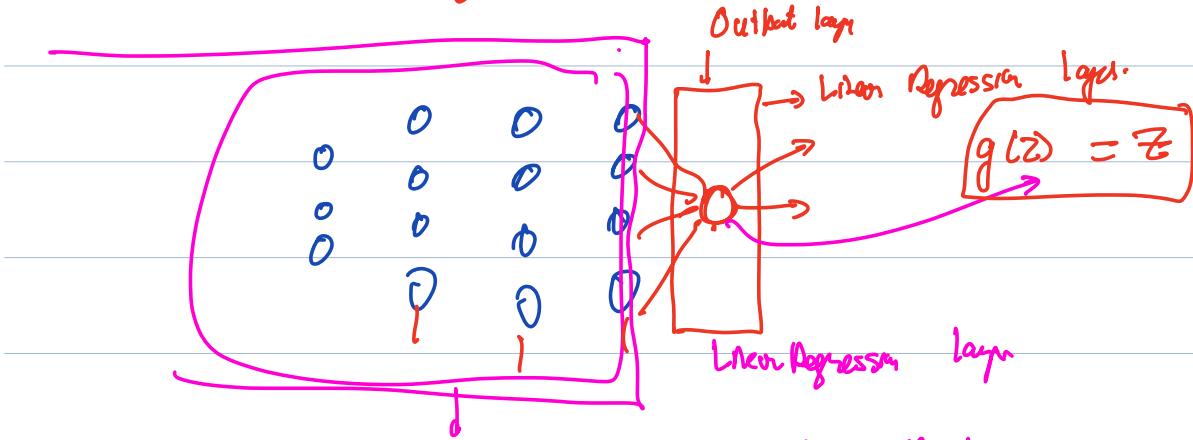
3 for quadrat.

Quadrat.

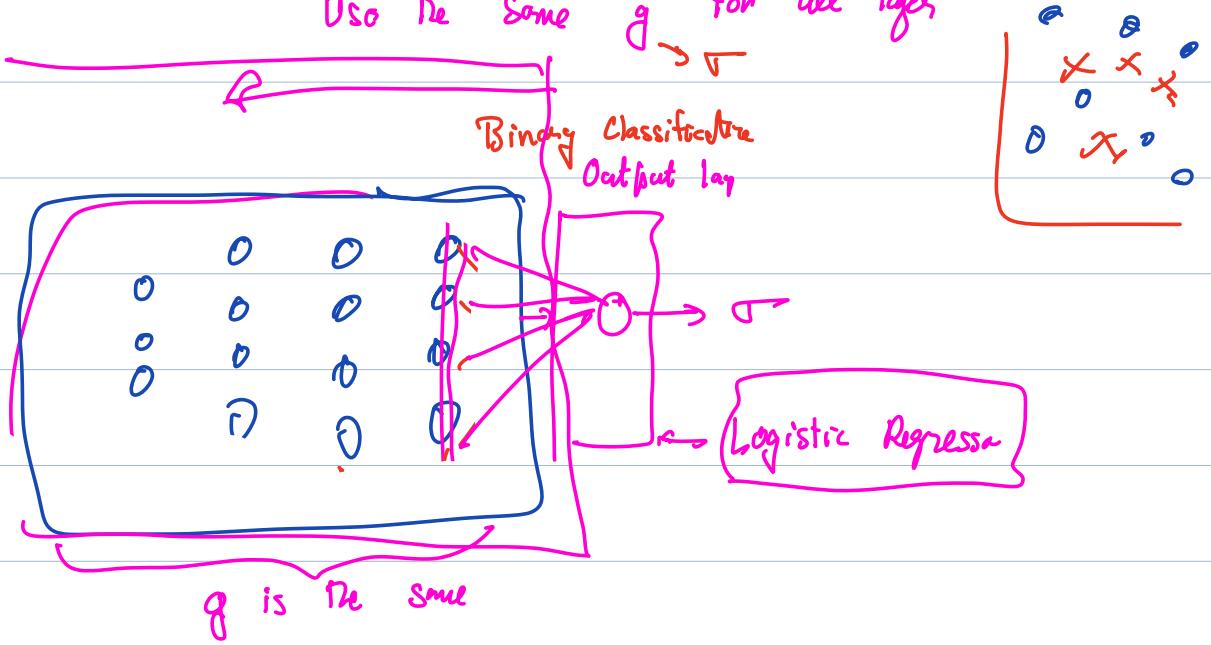
$$w_0 + w_1 x_1 + w_2 x_2$$

$$+ w_3 x_1^2 + w_4 x_2^2 + w_5 x_1 x_2$$

Regression, no limits on \hat{y} .



Use the same g for all layers



g is the same