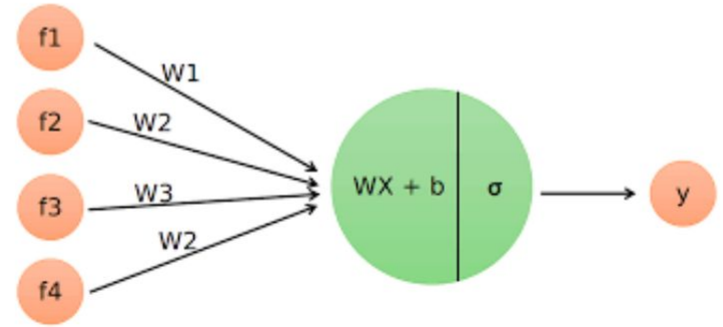
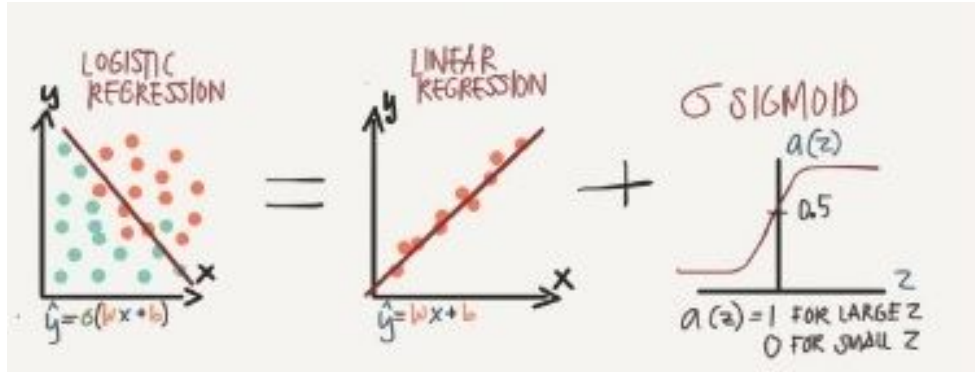


# LOGISTIC REGRESSION

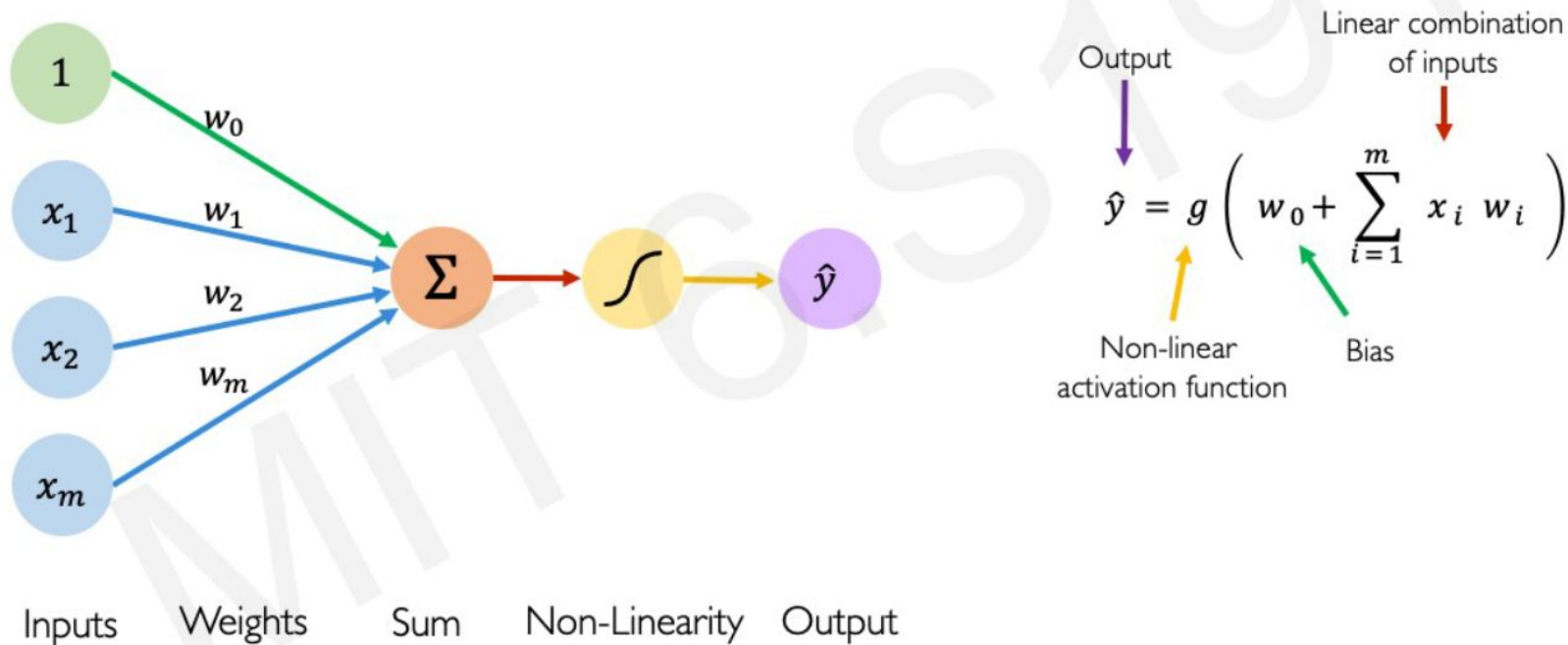


# LOGISTIC REGRESSION

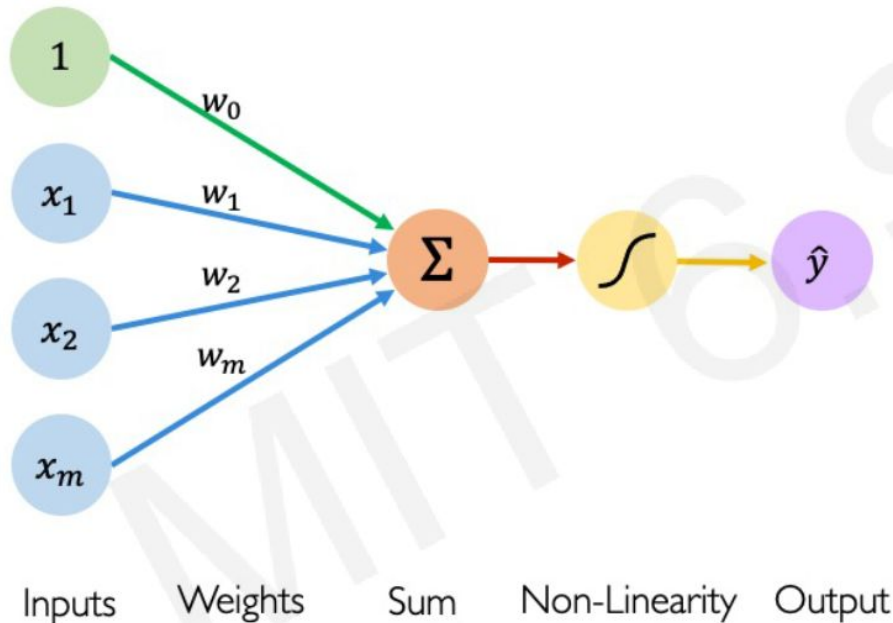
Logistic regression is a form of regression analysis in which the outcome variable is binary



# The Perceptron: Forward Propagation



# The Perceptron: Forward Propagation

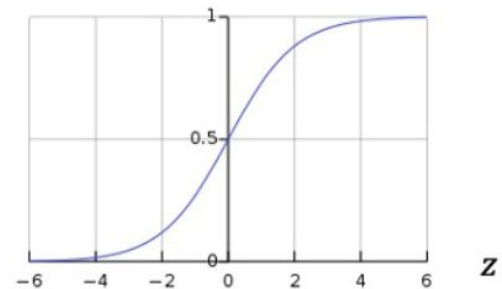


## Activation Functions

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

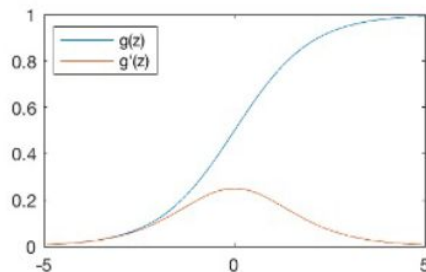
- Example: sigmoid function

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



# Common Activation Functions

Sigmoid Function



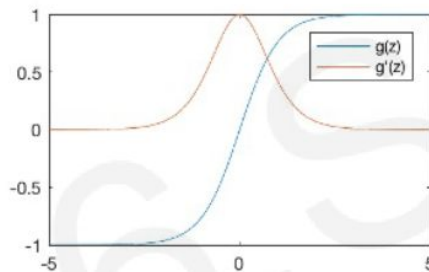
$$g(z) = \frac{1}{1 + e^{-z}}$$

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



`tf.math.sigmoid(z)`

Hyperbolic Tangent



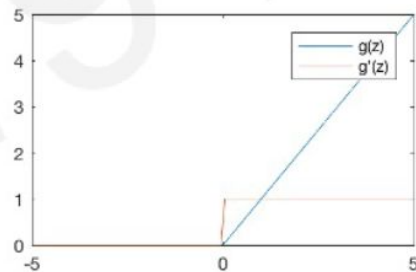
$$g(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

$$g'(z) = 1 - g(z)^2$$



`tf.math.tanh(z)`

Rectified Linear Unit (ReLU)



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

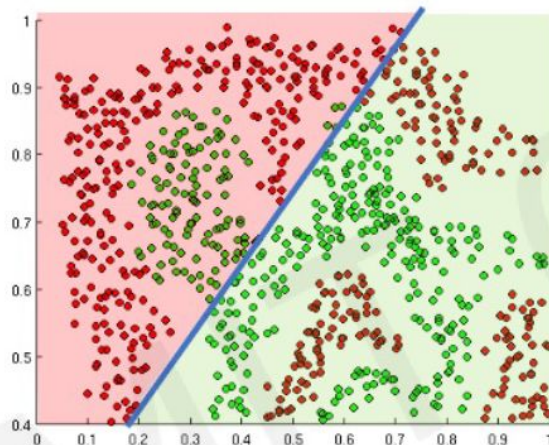
$$g'(z) = \begin{cases} 1, & z > 0 \\ 0, & \text{otherwise} \end{cases}$$



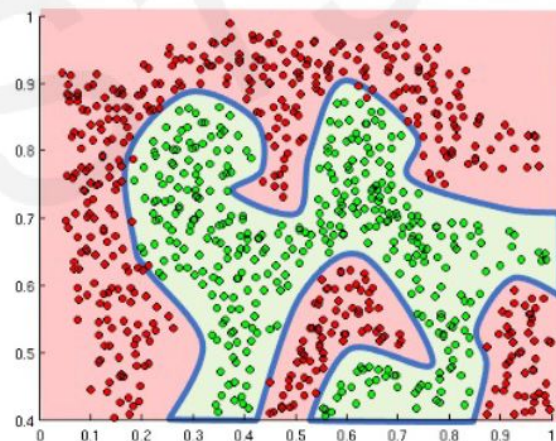
`tf.nn.relu(z)`

# Importance of Activation Functions

The purpose of activation functions is to **introduce non-linearities** into the network



Linear activation functions produce linear decisions no matter the network size



Non-linearities allow us to approximate arbitrarily complex functions

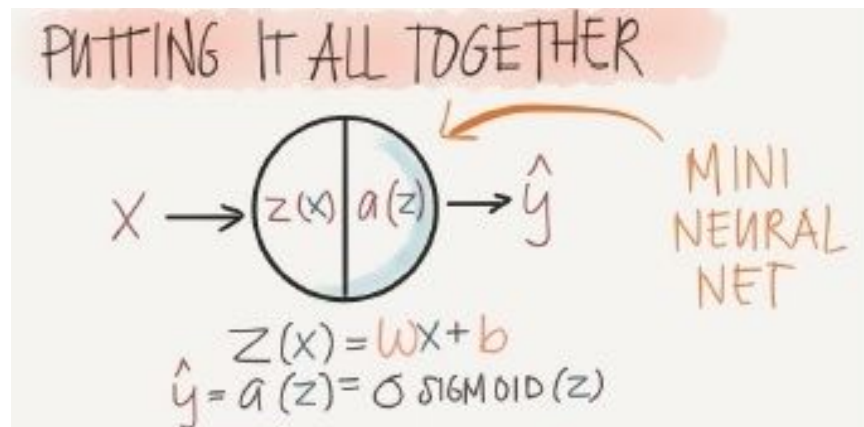
# LOGISTIC REGRESSION

$$z = \left( \sum_{i=1}^n w_i x_i \right) + b$$

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

$$\begin{aligned} P(y=1) &= \sigma(w \cdot x + b) \\ &= \frac{1}{1 + e^{-(w \cdot x + b)}} \end{aligned}$$

$$\begin{aligned} P(y=0) &= 1 - \sigma(w \cdot x + b) \\ &= 1 - \frac{1}{1 + e^{-(w \cdot x + b)}} \\ &= \frac{e^{-(w \cdot x + b)}}{1 + e^{-(w \cdot x + b)}} \end{aligned}$$





## Logistic Regression

$$z = b + a_1x_1 + a_2x_2 + a_3x_3$$
$$p = 1.0 / (1.0 + e^{-z})$$

Ex:

$$\begin{array}{ll} x_1 = 1.0 & a_1 = 0.01 \\ x_2 = 2.0 & a_2 = 0.02 \\ x_3 = 3.0 & a_3 = 0.03 \\ & b = 0.05 \end{array}$$

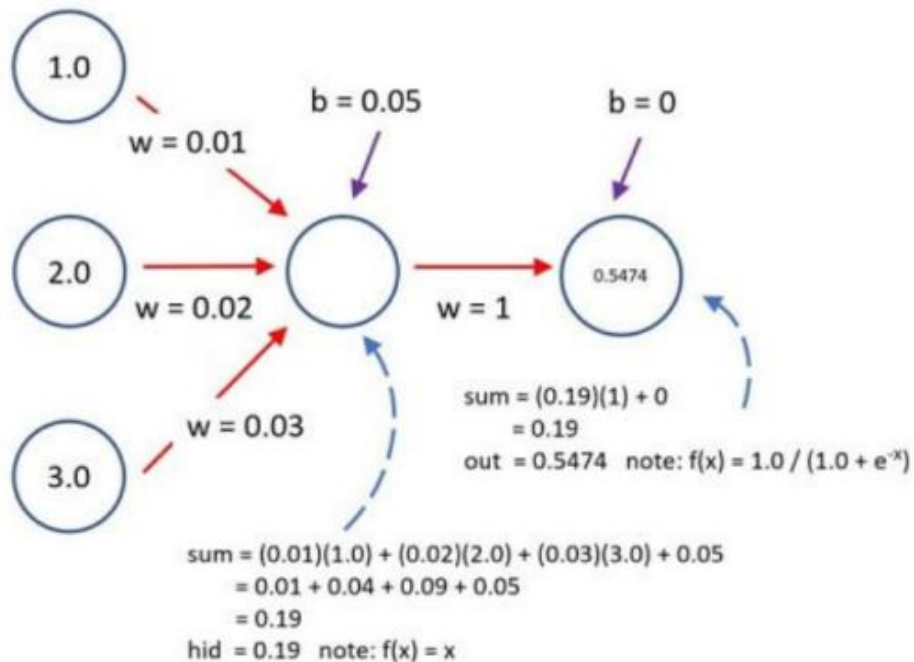
$$\begin{aligned} z &= (0.05) + (0.01)(1.0) + \\ &\quad (0.02)(2.0) + (0.03)(3.0) \\ &= 0.05 + 0.01 + 0.04 + 0.09 \\ &= 0.19 \end{aligned}$$

$$\begin{aligned} p &= 1.0 / (1.0 + e^{-0.19}) \\ &= 0.5474 \text{ (predicted class = 1)} \end{aligned}$$

## Neural Network

single hidden layer, identity activation  $f(x) = x$

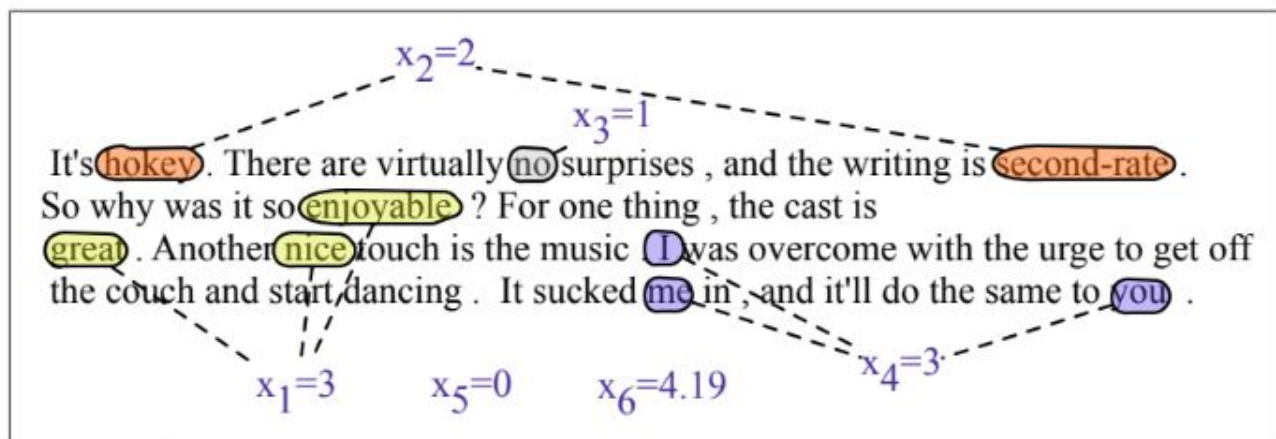
single output node, logistic sigmoid activation  $f(x) = 1 / (1 + e^{-x})$





Var	Definition	Value in Fig. 5.2
$x_1$	count(positive lexicon) $\in$ doc)	3
$x_2$	count(negative lexicon) $\in$ doc)	2
$x_3$	$\begin{cases} 1 & \text{if "no"} \in \text{doc} \\ 0 & \text{otherwise} \end{cases}$	1
$x_4$	count(1st and 2nd pronouns $\in$ doc)	3
$x_5$	$\begin{cases} 1 & \text{if "!"} \in \text{doc} \\ 0 & \text{otherwise} \end{cases}$	0
$x_6$	log(word count of doc)	$\ln(66) = 4.19$

Let's assume for the moment that we've already learned a real-valued weight for each of these features, and that the 6 weights corresponding to the 6 features are  $[2.5, -5.0, -1.2, 0.5, 2.0, 0.7]$ , while  $b = 0.1$ . (We'll discuss in the next section how



**Figure 5.2** A sample mini test document showing the extracted features in the vector  $x$ .

Given these 6 features and the input review  $x$ ,  $P(+|x)$  and  $P(-|x)$  can be computed using Eq. 5.5:

$$\begin{aligned}
 p(+|x) &= P(Y = 1|x) = \sigma(w \cdot x + b) \\
 &= \sigma([2.5, -5.0, -1.2, 0.5, 2.0, 0.7] \cdot [3, 2, 1, 3, 0, 4.19] + 0.1) \\
 &= \sigma(.833) \\
 &= 0.70
 \end{aligned} \tag{5.6}$$

$$\begin{aligned}
 p(-|x) &= P(Y = 0|x) = 1 - \sigma(w \cdot x + b) \\
 &= 0.30
 \end{aligned}$$

## Algorithm

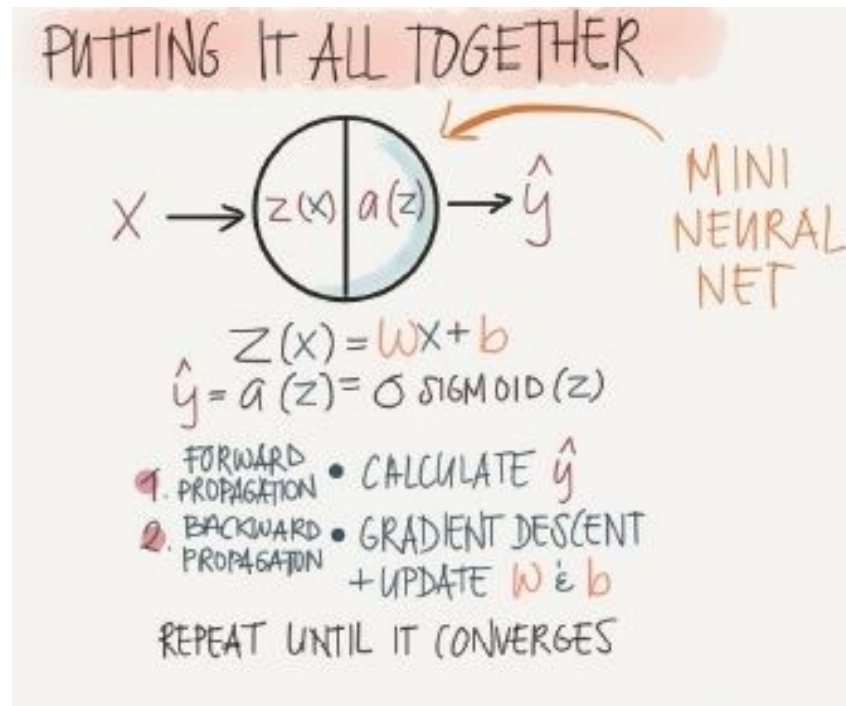
- Calculate the prediction  $\hat{y}$  using the current parameters ( $W$  and  $b$ )
- Calculate the loss of the current values
- Calculate the gradients of the loss function with respect to the parameters
- Adjust the weights (optimize) using the gradients
- Repeat for the number of epochs, i.e. the number of times to go through the provided examples (dataset)

# LEARNING IN LOGISTIC REGRESSION

How are the parameters of the model, the weights  $w$  and bias  $b$ , learned?

Two components required:

1. **Loss function or Cost function:** metric for how close the current label ( $\hat{y}$ ) is to the true gold label  $y$ .
2. **Optimization algorithm:** for iteratively updating the weights



# LOSS FUNCTION OR COST FUNCTION

For an observation  $x$ , how close the classifier output ( $\hat{y} = \sigma(w \cdot x + b)$ ) is to the correct output ( $y$ , which is 0 or 1)

$L(\hat{y}, y)$  = How much  $\hat{y}$  differs from the true  $y$

# LOSS FUNCTION OR COST FUNCTION

Conditional maximum likelihood estimation prefers the correct class labels of the training examples to be more likely.

We choose the parameters  $w$ ,  $b$  that maximize the log probability of the true  $y$  labels in the training data given the observations  $x$ .

# MAXIMUM LIKELIHOOD ESTIMATION

$$\begin{aligned} P(y=1) &= \sigma(w \cdot x + b) \\ &= \frac{1}{1 + e^{-(w \cdot x + b)}} \end{aligned}$$

For samples labeled as '1', we try to estimate (w and b) such that the product of all probability  $p(x)$  is as close to 1 as possible.

And for samples labeled as '0', we try to estimate (w and b) such that the product of all probability is as close to 0 as possible in other words  $(1 - p(x))$  should be as close to 1 as possible.

$$\text{for samples labelled as 1 : } \prod_{s \text{ in } y_i = 1} p(x_i)$$

$$\text{for samples labelled as 0 : } \prod_{s \text{ in } y_i = 0} (1 - p(x_i))$$



# MAXIMUM LIKELIHOOD ESTIMATION

On combining both , (w and b) parameters should be such that the product of both of these products is maximum over all elements of the dataset.

$$L(\beta) = \prod_{s \text{ in } y_i = 1} p(x_i) * \prod_{s \text{ in } y_i = 0} (1 - p(x_i))$$

This function is the one we need to optimize and is called the **likelihood function**.

# Loss Function

$$\text{If } y = 1: \quad p(y|x) = \hat{y}$$

$$\text{If } y = 0: \quad p(y|x) = 1 - \hat{y}$$

$$P(y|x) = \hat{y}^y (1-\hat{y})^{(1-y)}$$

$$\text{If } y = 1 \Rightarrow \hat{y}$$

$$\text{If } y = 0 \Rightarrow 1-\hat{y}$$

# LOG LIKELIHOOD

Since there are only two discrete outcomes (1 or 0), this is a Bernoulli distribution.

$p(y|x)$  for one observation:

$$p(y|x) = \hat{y}^y (1 - \hat{y})^{1-y}$$

$$\begin{aligned}\log p(y|x) &= \log [\hat{y}^y (1 - \hat{y})^{1-y}] \\ &= y \log \hat{y} + (1 - y) \log(1 - \hat{y})\end{aligned}$$

whatever values maximize a probability will also maximize the log of the probability.

# LOG LIKELIHOOD

Since there are only two discrete outcomes (1 or 0), this is a Bernoulli distribution.

$p(y|x)$  for  $n$  observation:

$$L(\beta) = \prod_s \hat{y}^y (1 - \hat{y})^{1-y}$$

$$l(\beta) = \sum_{i=1}^n y \log \hat{y} + (1 - y) \log(1 - \hat{y})$$

where,  $l(\beta)$  is log-likelihood

# CROSS ENTROPY

Log likelihood should be maximized.

In order to turn this into loss function (something that we need to minimize), we'll just flip the sign.

negative log likelihood loss = cross-entropy loss.

$$L_{CE}(\hat{y}, y) = -\log p(y|x) = -[y \log \hat{y} + (1 - y) \log (1 - \hat{y})]$$

$$L_{CE}(w, b) = -[y \log \sigma(w \cdot x + b) + (1 - y) \log (1 - \sigma(w \cdot x + b))]$$

# BINARY CROSS ENTROPY

negative log likelihood loss = cross-entropy loss.

Since we are working on 2 classes it becomes

## Binary Cross Entropy Loss

$$L_{CE}(\hat{y}, y) = -\log p(y|x) = -[y \log \hat{y} + (1 - y) \log (1 - \hat{y})]$$

$$L_{CE}(w, b) = -[y \log \sigma(w \cdot x + b) + (1 - y) \log (1 - \sigma(w \cdot x + b))]$$

# BINARY CROSS ENTROPY LOSS

$$L_{CE}(\hat{y}, y) = -\log p(y|x) = -[y \log \hat{y} + (1 - y) \log(1 - \hat{y})]$$

Binary cross-entropy is often calculated as the average cross-entropy across all data examples.

$$L = -\frac{1}{N} \left[ \sum_{j=1}^N [y \log \hat{y} + (1 - y) \log(1 - \hat{y})] \right]$$



# WHY NOT MSE?



$\log(p(x)) < 0$  for all  $p(x)$  in  $(0,1)$ .

Actual label: "1" Predicted: "0"

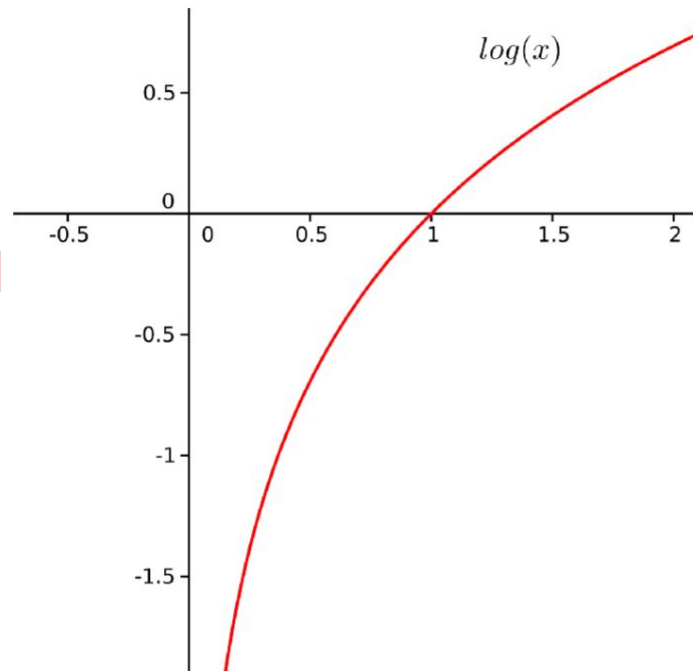
MSE:  $(1 - 0)^2 = 1$

Log loss:  $-(1 * \log(0) + 0 * \log(1)) = \text{infinity!!}$

Actual label: "1" Predicted: "1"

MSE:  $(1 - 1)^2 = 0$

Log loss:  $-(1 * \log(1) + 0 * \log(0)) = 0$

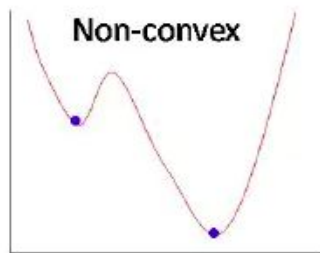
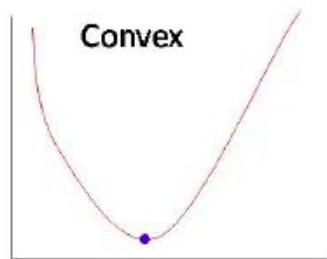


# WHY NOT MSE?



Mean squared error is a non-convex function.

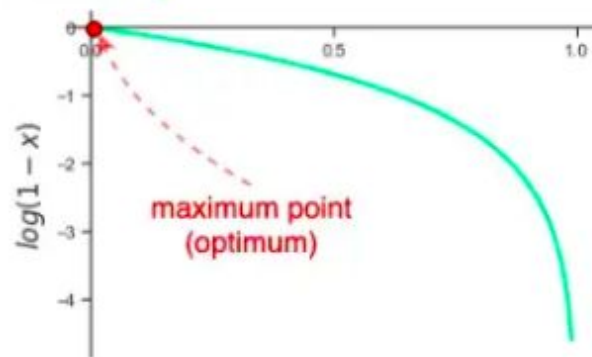
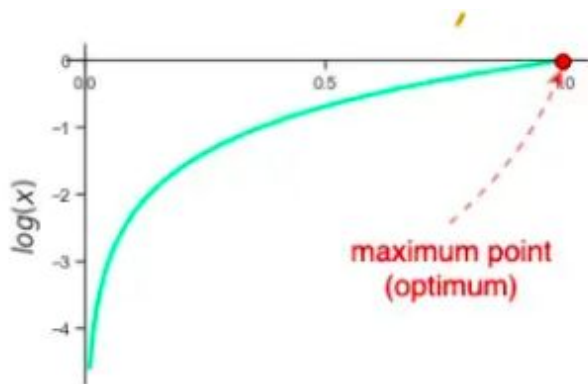
Cross entropy is a convex function.



Gradient descent will converge into global minimum only if the function is convex.

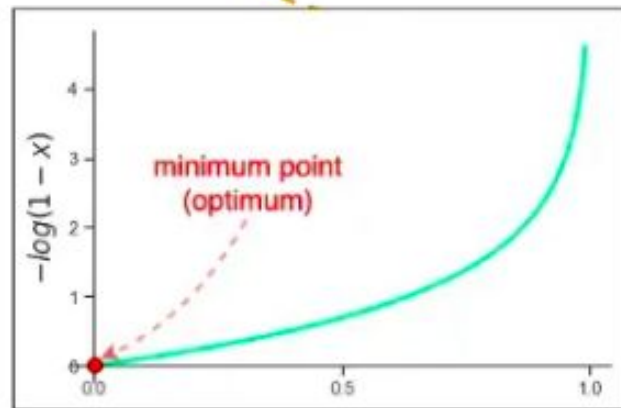
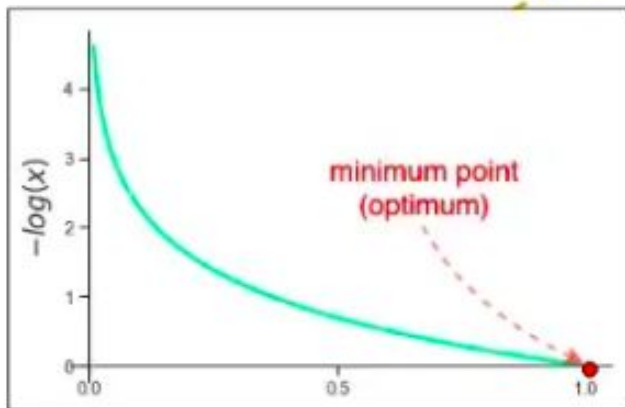
# BINARY CROSS ENTROPY LOSS

$$L_{CE}(\hat{y}, y) = -\log p(y|x) = -[y \log \hat{y} + (1 - y) \log(1 - \hat{y})]$$



# BINARY CROSS ENTROPY LOSS

$$L_{CE}(\hat{y}, y) = -\log p(y|x) = -[y \log \hat{y} + (1 - y) \log(1 - \hat{y})]$$

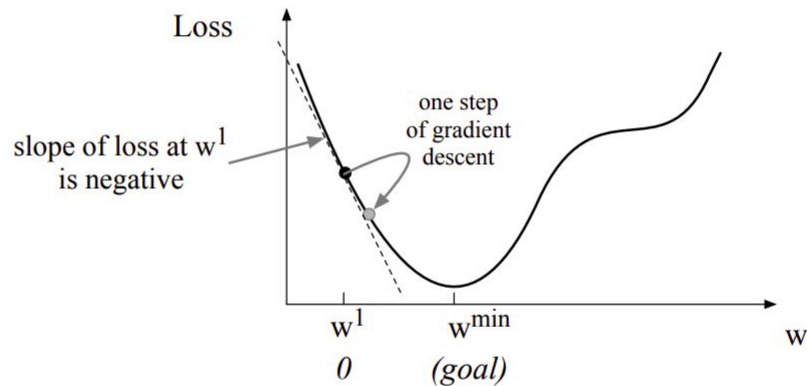


# GRADIENT DESCENT ALGORITHM

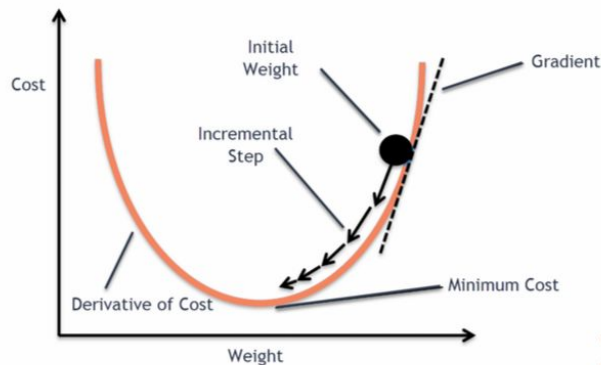
Goal:

- To find the optimal weights (in the case of logistic regression  $\theta = w, b$ )
- Minimize the loss function

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \frac{1}{m} \sum_{i=1}^m L_{CE}(y^{(i)}, x^{(i)}; \theta)$$



# GRADIENT DESCENT ALGORITHM



Gradient Descent steps are:

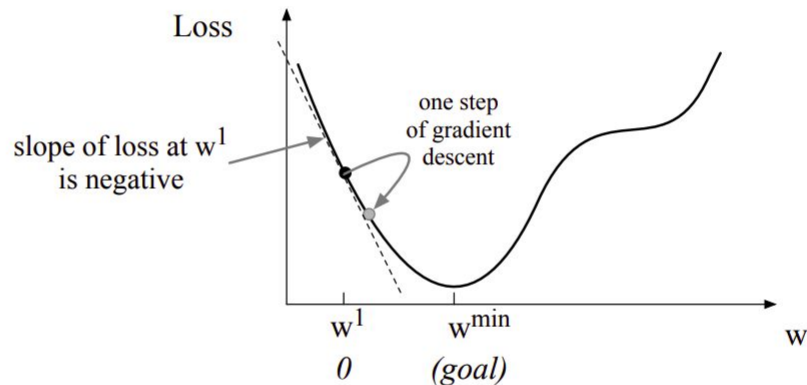
1. choose a starting point (initialisation)
2. calculate gradient at this point
3. make a scaled step in the opposite direction to the gradient (objective: minimise)
4. repeat points 2 and 3 until one of the criteria is met:
  - maximum number of iterations reached
  - step size is smaller than the tolerance (due to scaling or a small gradient).

# GRADIENT DESCENT ALGORITHM

Goal:

- To find the optimal weights (in the case of logistic regression  $\theta = w, b$ )
- Minimize the loss function

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \frac{1}{m} \sum_{i=1}^m L_{CE}(y^{(i)}, x^{(i)}; \theta)$$





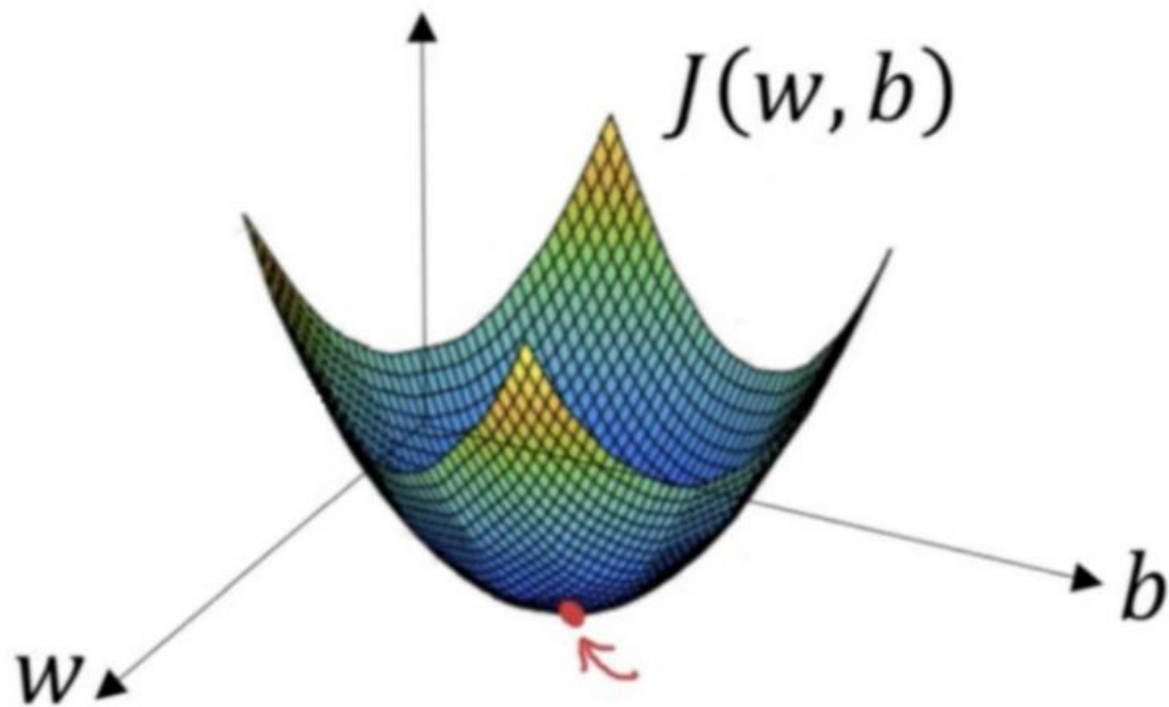
# Gradient Descent

$$J(w, b) = \frac{\sum_1^m L(\hat{y}^i, y^i)}{m} = \frac{\sum_1^m y^i \log \hat{y}^i + (1 - y^i) \log(1 - \hat{y}^i)}{m}$$

In order to minimize the cost function for minimal error across the training data set to find  $w$  and  $b$

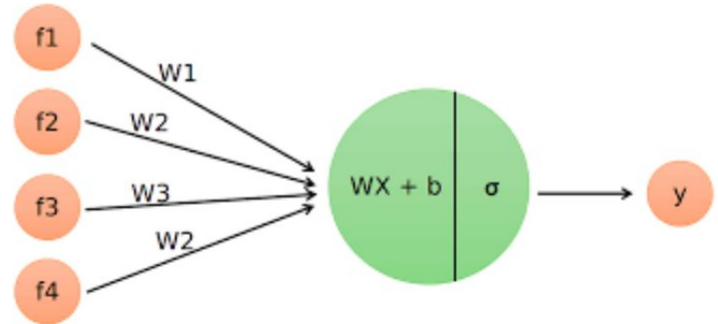
The value of the parameters can be achieved using gradient descent technique.

# Gradient Descent



# GRADIENT— BACKWARD PROPAGATION

$$L_{CE}(w, b) = -[y \log \sigma(w \cdot x + b) + (1 - y) \log (1 - \sigma(w \cdot x + b))]$$



## Gradient

$$z = w_1 x_1 + w_2 x_2 + b \rightarrow \hat{y} = a = \sigma(z) \rightarrow L(\hat{y}, y)$$

$$\Leftrightarrow a = \hat{y}$$

$$\boxed{w_1} \Rightarrow \frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial a} \cdot \frac{\partial a}{\partial z} \cdot \frac{\partial z}{\partial w_1}$$

$$\begin{aligned} \frac{\partial L}{\partial a} &= \frac{\partial}{\partial a} (-y \log a - (1-y) \log(1-a)) \\ &= -y \left( \frac{1}{a} \right) - (-1) \frac{(1-y)}{(1-a)} \end{aligned}$$

$$\boxed{\frac{\partial L}{\partial a} = \left( \frac{-y}{a} \right) + \left( \frac{1-y}{1-a} \right)}$$

$$\boxed{\frac{\partial a}{\partial z} = a(1-a)}$$

$$\boxed{\frac{\partial z}{\partial w_1} = x_1}$$

$$\frac{\partial L}{\partial w_1} = \left( \left( \frac{-y}{a} + \frac{(1-y)}{1-a} \right) \cdot (a)(1-a) \right) \cdot x_1$$

$$= (a-y) \cdot x_1$$

update for  $w_1$ ,

$$\frac{\partial L}{\partial w_1} = (a-y) \cdot x_1$$

$$\text{Here, } (a-y) = \frac{\partial L}{\partial z}$$

$$\boxed{w_1 = w_1 - \alpha \frac{\partial L}{\partial w_1}}$$

Similarly, for all parameters

$$\boxed{w_i = w_i - \alpha \frac{\partial L}{\partial w_i}}$$

$$i = 1, 2, \dots, m$$

$m$  = no. of parameters

$$\boxed{b = b - \alpha \frac{\partial L}{\partial b}}$$

$$\text{where, } \frac{\partial L}{\partial b} = (a-y)$$

# LOGISTIC REGRESSION - LOSS CALCULATION

Observation	True Label ( $y$ )	Predicted Probability ( $p$ )
1	1	0.9
2	0	0.2
3	1	0.8
4	0	0.4

Total Loss =  $0.1054 + 0.2231 + 0.2231 + 0.5108 = 1.0624$

Average Loss (BCE) =  $1.06244 / 4 = 0.2656$

# LOGISTIC REGRESSION - EXAMPLE DONE IN CLASS

$$w_1 = 0.5$$

$$w_2 = 0.31$$

$$b = 0.09$$

$$x_1 = 5$$

$$x_2 = 3$$

$$y = 1$$

$$\text{Learning Rate} = 0.01$$



# LOGISTIC REGRESSION - EXAMPLE DONE IN CLASS

Forward propagation:

$z=3.52$ ,  $a=0.97$

$L=0.305$  (Binary Cross Entropy Loss)

$dL/dw_1=-0.15$ ,  $dL/dw_2=-0.09$ ,  $dL/dwb=-0.03$

$w_{1\_new}=0.5015$

$w_{2\_new}=0.3109$

$b\_new=0.0903$

2. Consider a classification problem where we are building a logistic regression classifier. The task is to predict if a given city has a risk of a disease epidemic or not. The data is defined using two input

## Practice Problem 6

- Construct a logistic regression model with two predictors for the riding mower example with  $\beta_0 = -25.9382$ ,  $\beta_1 = 0.1109$ ,  $\beta_2 = 0.9638$ , where  $\beta_1$  and  $\beta_2$  are for the "Income" and "Lot\_Size" variables, respectively.
- Using the logistic regression model with probability cutoff = 0.75, classify the following 6 customers as "Owner" or "Nonowner" : if  $p \geq 0.75$  then the case as a "Owner". Present the results in a classification matrix.

Customer#	Income	Lot_Size	Ownership
1	60.0	18.4	Owner
2	64.8	21.6	Owner
3	84.0	17.6	Nonowner
4	59.4	16.0	Nonowner
5	108.0	17.6	Owner
6	75	19.6	Nonowner

# REFERENCES

- <https://towardsdatascience.com/logistic-regression-detailed-overview-46c4da4303bc>
- <https://arunaddagatla.medium.com/maximum-likelihood-estimation-in-logistic-regression-f86ff1627b67>
- <https://towardsdatascience.com/cross-entropy-loss-function-f38c4ec8643e>
- <https://www.marktechpost.com/2021/04/08/logistic-regression-with-keras/>
- <https://www.datasciencecentral.com/logistic-regression-as-a-neural-network/>
- <https://github.com/SSaishruthi/LogisticRegressionVectorizedImplementation/blob/master/LogisticRegression.ipynb>
- <https://towardsdatascience.com/where-did-the-binary-cross-entropy-loss-function-come-from-ac3de349a715>
- <https://www.youtube.com/watch?v=nzNp05AyBM8>
-