

# Logistic Regression

## Cost Function & Gradient Descent

Rahul Singh  
rsingh@arrsingh.com

## Logistic Regression is used to predict a binary value

A model that returns a binary value (true / false)

The diagram shows the equation  $y = f(x)$  with two annotations. A green arrow points from a green-bordered box containing the text " $x$  is the independent variable" to the  $x$  in the equation. A red arrow points from a red-bordered box containing the text " $y$  (dependent variable) that returns 0 (false) or 1 (true)" to the  $y$  in the equation.

$$y = f(x)$$

$x$  is the independent variable

$y$  (dependent variable) that returns 0 (false) or 1 (true)

## Logistic Function

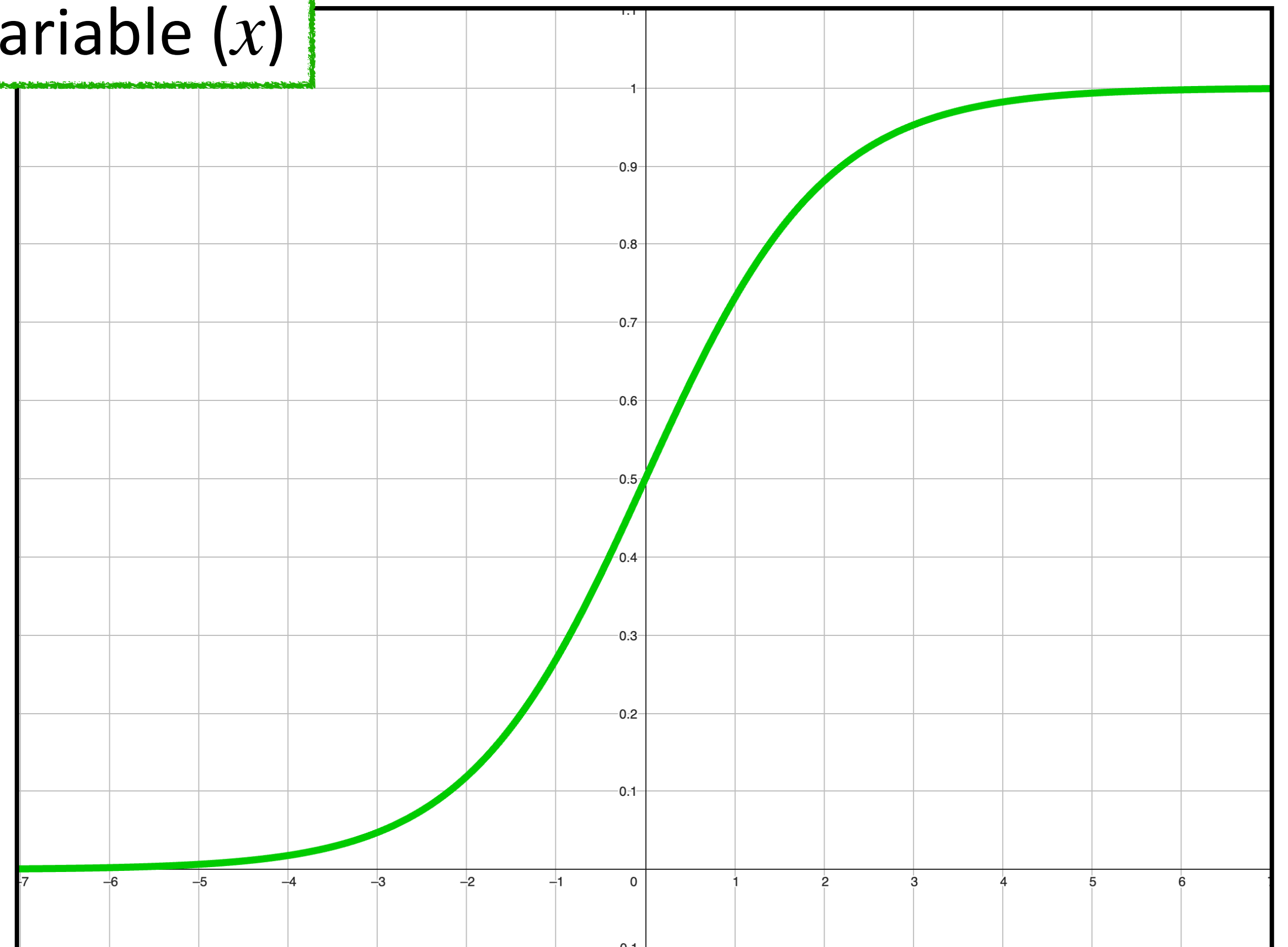
Logistic regression uses a sigmoid curve (Logistic Function) that converts a linear combination of input features (x values) into probabilities (y values)

$$y = \beta + \omega x$$

$$\sigma(y) = \frac{1}{1 + e^{-(\beta + \omega x)}}$$

1 independent variable (x)

1 dependent variable (y)



Generalizing this to  $k$  independent variables...

# Logistic Regression

## Logistic Function

The model is generalized to  $k$  independent variables and  $k + 1$  parameters

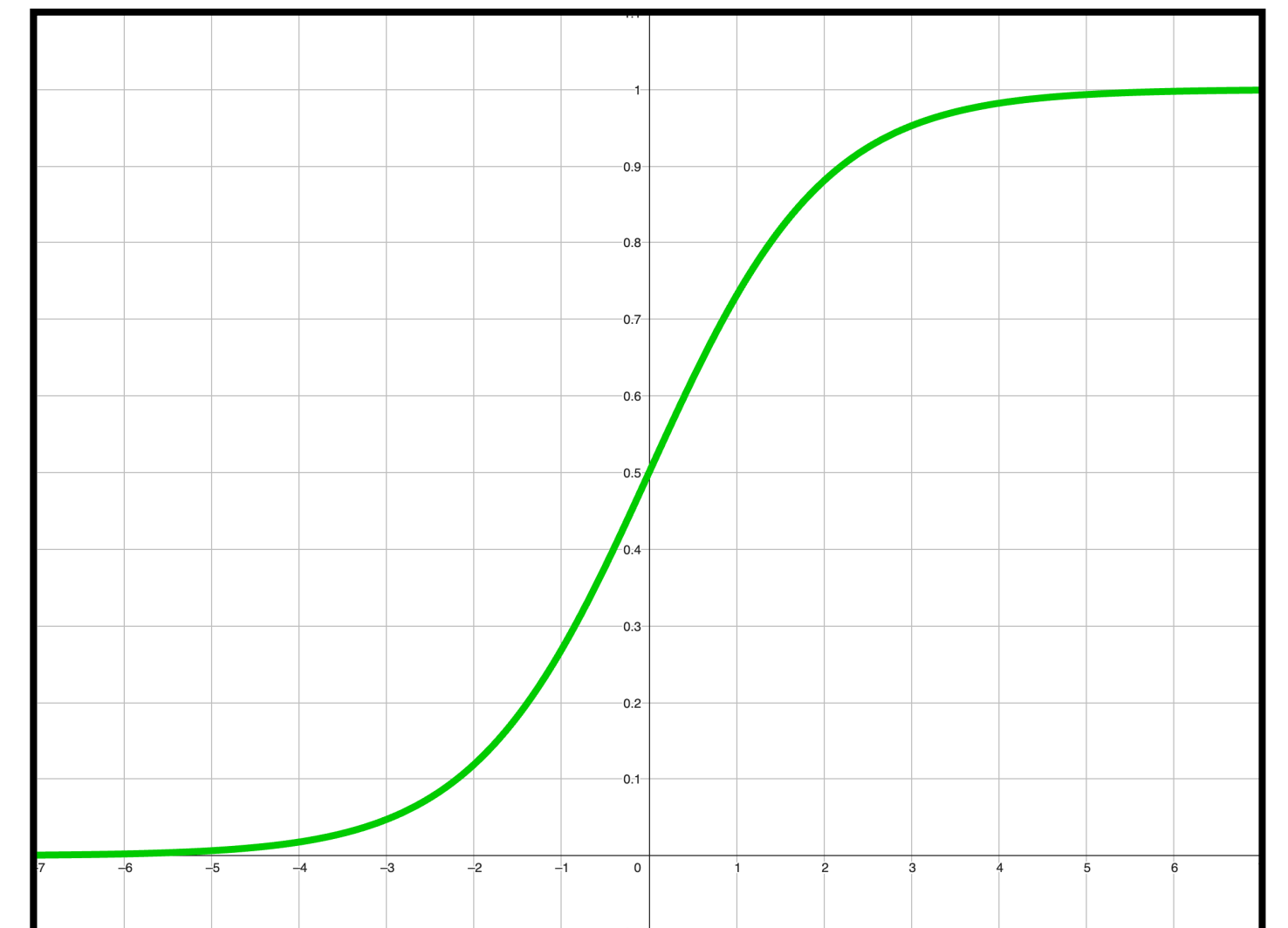
$$y = \beta + \omega_1 x_1 + \omega_2 x_2 + \omega_3 x_3 + \dots + \omega_k x_k$$

1 dependent variable  $y$

$k$  independent variables  $x_1 \dots x_k$

$$\sigma(y) = \frac{1}{1 + e^{-(\beta + \omega_1 x_1 + \omega_2 x_2 + \omega_3 x_3 + \dots + \omega_k x_k)}}$$

This model has  
 $k + 1$  parameters  
 $\beta, \omega_1, \omega_2 \dots \omega_k$



# Logistic Regression

## Logistic Function

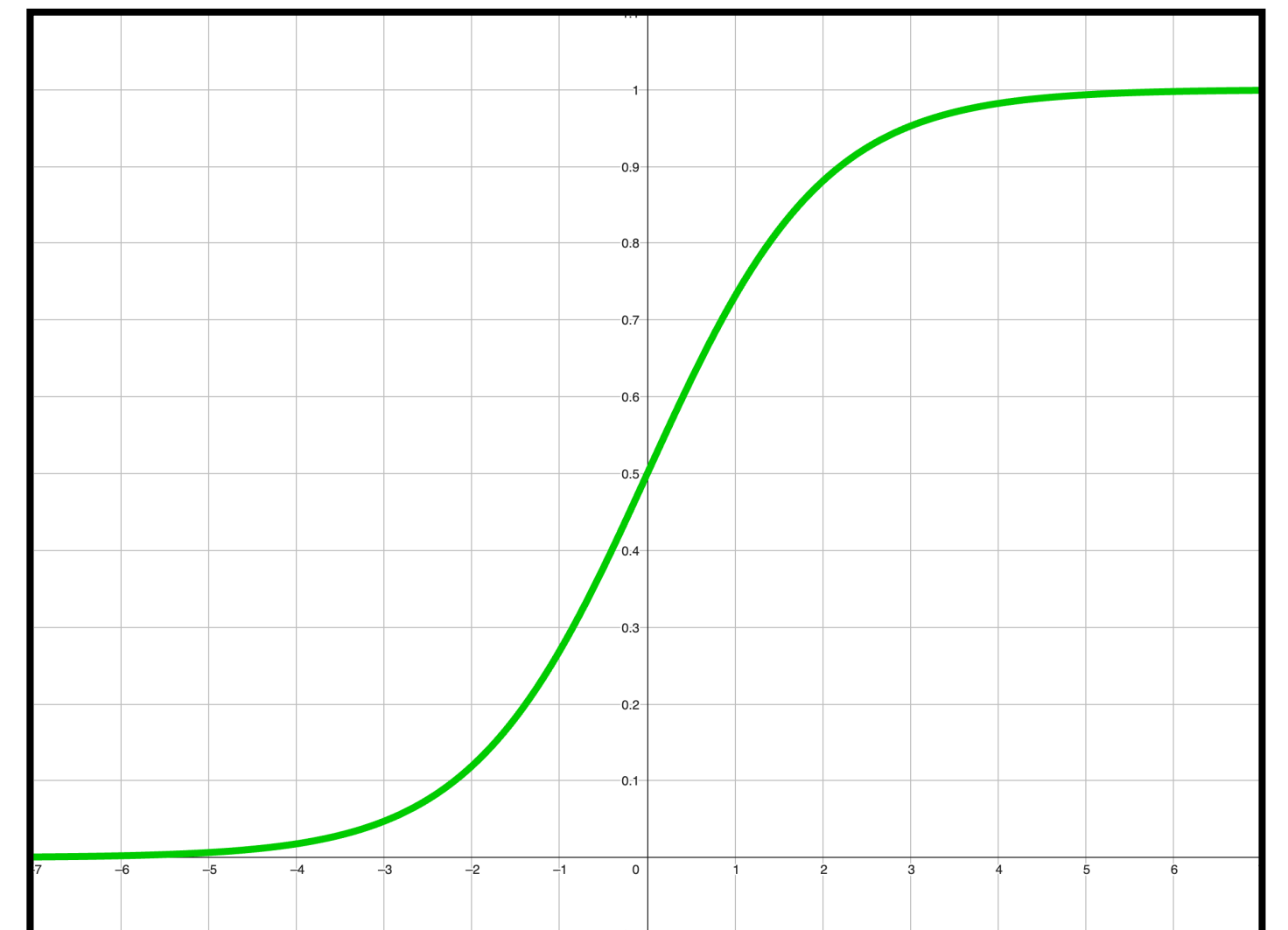
The model is generalized to  $k$  independent variables and  $k + 1$  parameters

$$y = \beta + \omega_1 x_1 + \omega_2 x_2 + \omega_3 x_3 + \dots + \omega_k x_k$$

General Matrix form:

$$\hat{Y} = \sigma(W^T X + \beta)$$

$W$  is a  $k \times 1$  vector of weights  $\omega_1, \omega_2, \omega_3 \dots \omega_k$   
 $X$  is a  $k \times n$  matrix of  $n$  observations  $x_1, x_2, x_3 \dots x_k$   
 $\beta$  is a scalar



# Logistic Regression

## Logistic Function

The model is generalized to  $k$  independent variables and  $k + 1$  parameters

$$y = \beta + \omega_1 x_1 + \omega_2 x_2 + \omega_3 x_3 + \dots + \omega_k x_k$$

General Matrix form:

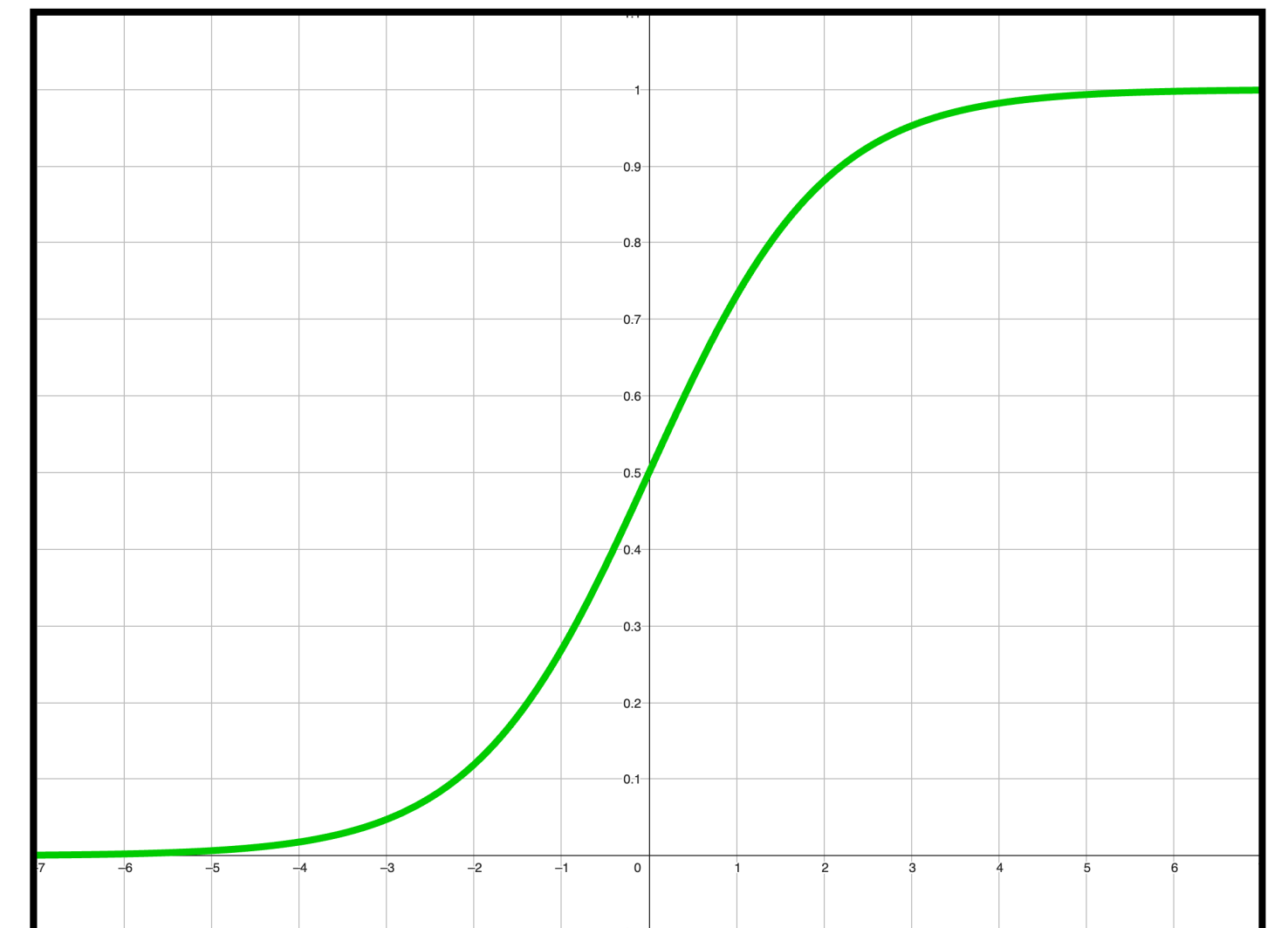
$$\hat{Y} = \sigma(W^T X + \beta)$$

$\hat{Y}$  is the vector of predicted values from the model.  
 $\hat{Y}$  is a vector of probabilities each between 0 and 1

A Threshold converts a given probability to a binary value

if  $\hat{y}_i \geq 0.5$  then 1

if  $\hat{y}_i < 0.5$  then 0





## Logistic Function

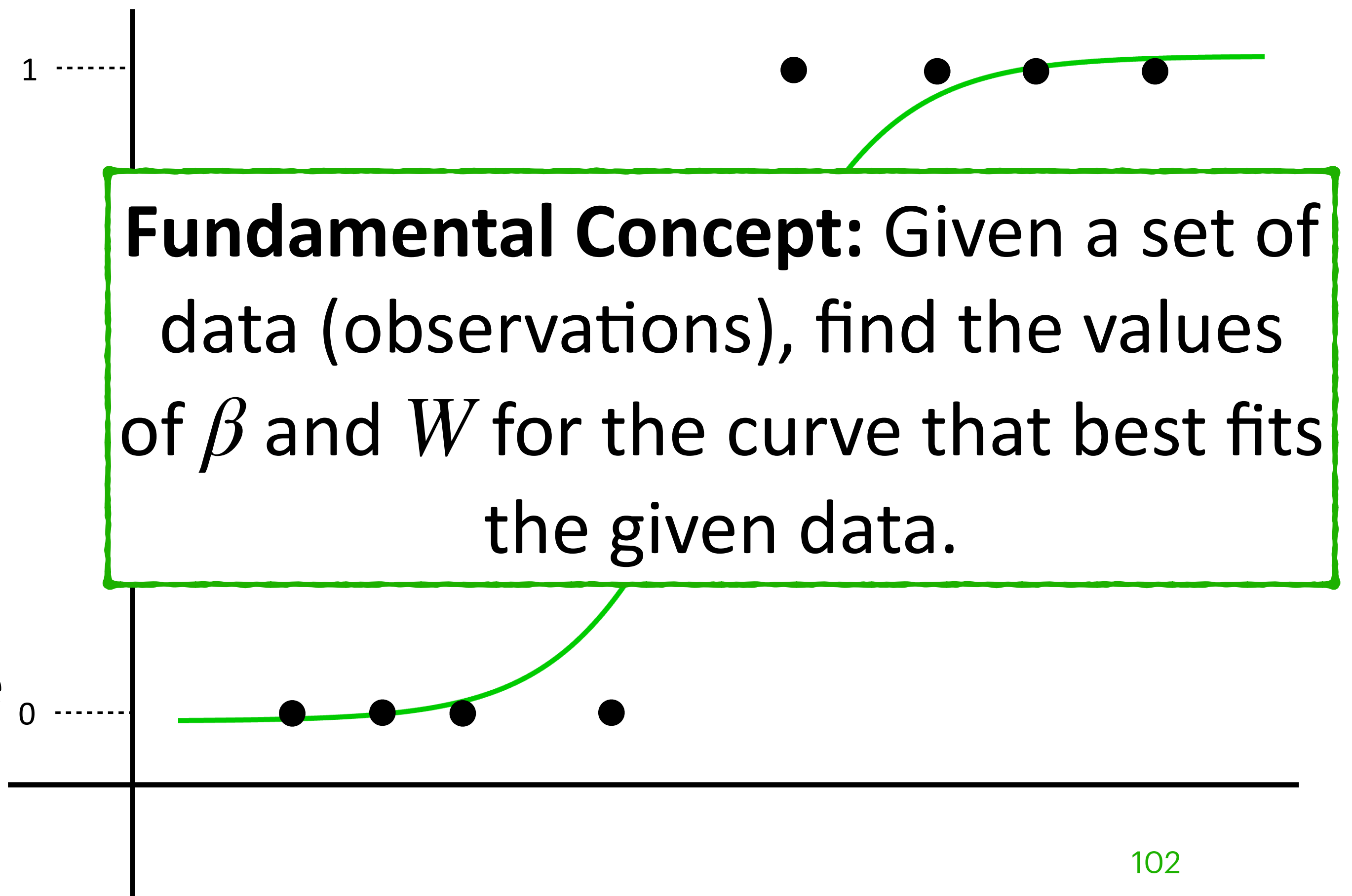
Logistic regression uses a sigmoid curve (Logistic Function) that converts a linear combination of input features (x values) into probabilities (y values)

Curve of best fit is...

$$\hat{Y} = \sigma(W^T X + \beta)$$

$$\Rightarrow \hat{y}_i = \frac{1}{1 + e^{-(\omega_i x_i + \beta)}}$$

We can use **Gradient Descent** to find the optimal values of  $\beta$  and  $W$





# Recap: Gradient Descent and Linear Regression

# Multiple Regression

Linear Model in  
 $k + 1$  Dimensions

$$\hat{y} = \beta_0 + \beta_1 \hat{x}_1 + \beta_2 \hat{x}_2 + \beta_3 \hat{x}_3 + \dots + \beta_k \hat{x}_k$$

$$\hat{Y} = \hat{X}\beta$$

The Mean Squared Error (MSE):

$$\frac{1}{2n} \| Y - X\beta \|^2$$

$$\begin{bmatrix} \hat{y}_0 \\ \hat{y}_1 \\ \hat{y}_2 \\ \hat{Y} \\ \cdot \\ \cdot \\ \hat{y}_n \end{bmatrix} = \begin{bmatrix} 1 & \hat{x}_{10} & \hat{x}_{20} & \hat{x}_{30} & \cdot & \cdot & \cdot & \hat{x}_{k0} \\ 1 & \hat{x}_{11} & \hat{x}_{21} & \hat{x}_{21} & \cdot & \cdot & \cdot & \hat{x}_{k1} \\ 1 & \hat{x}_{12} & \hat{x}_{22} & \hat{x}_{22} & \cdot & \cdot & \cdot & \hat{x}_{k2} \\ 1 & \hat{x}_{13} & \hat{x}_{23} & \hat{X} & \cdot & \cdot & \cdot & \hat{x}_{k3} \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ 1 & \hat{x}_{1n} & \hat{x}_{2n} & \hat{x}_{3n} & \cdot & \cdot & \cdot & \hat{x}_{kn} \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \beta \\ \cdot \\ \cdot \\ \beta_k \end{bmatrix}$$

# Multiple Regression

$$\hat{Y} = \hat{X}\beta$$

Linear Model in  
 $k + 1$  Dimensions

The **Mean Squared Error (MSE)**:

$$\frac{1}{2n} \| Y - X\beta \|^2$$

Cost Function

Gradient descent can be used to minimize the cost function.

We compute the partial derivative of the cost function w.r.t the parameters  $\beta$

Partial Derivative w.r.t  $\beta$ :

$$\frac{\partial}{\partial \beta} \frac{1}{2n} \| Y - X\beta \|^2 = -\frac{1}{n} X^T (Y - X\beta)$$

Partial Derivative w.r.t  $\beta$

A linear model in  $k + 1$  dimensions...

$$\hat{Y} = \hat{X}\beta$$

Cost Function (Mean Squared Error (MSE)):

$$\frac{1}{2n} \| Y - X\beta \|^2$$

Gradient Vector (Partial Derivative w.r.t  $\beta$ ):

$$\frac{\partial}{\partial \beta} \frac{1}{2n} \| Y - X\beta \|^2 = -\frac{1}{n} X^T (Y - X\beta)$$

# Gradient Descent

## Gradient Descent: Basic Concept

**Step 1:** Start with random values for  $\beta$

**Step 2:** Compute the partial derivative of the cost function w.r.t  $\beta$

**Step 3:** Calculate a step size that is proportional to the slope

**Step 4:** Calculate new values for  $\beta$  by subtracting the step size

**Step 5:** Go to step 2 and repeat

We need a cost function for Logistic Regression

# Logistic Regression

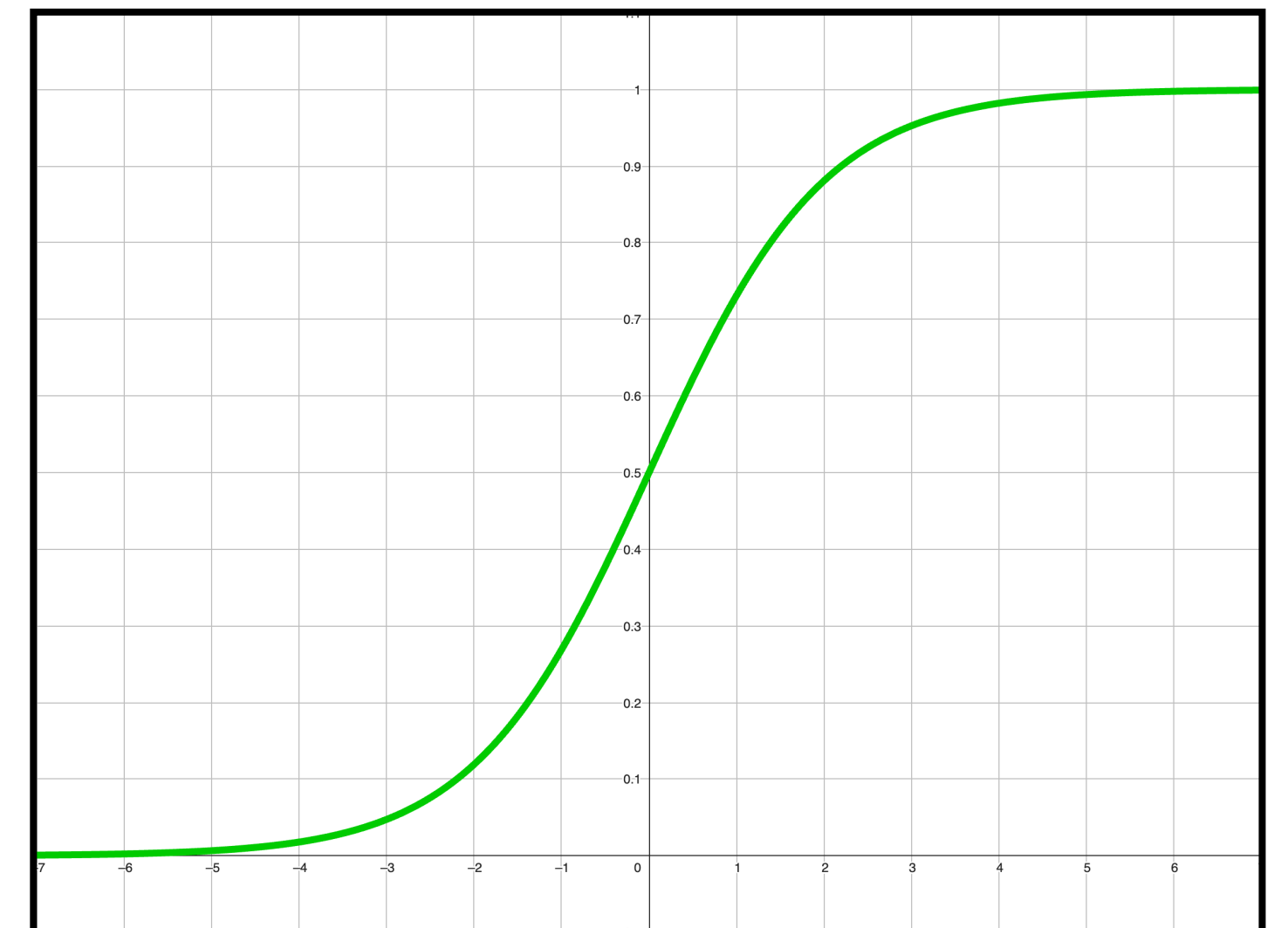
# Logistic Regression

$$\hat{Y} = \sigma(W^T X + \beta)$$

$W$  is a  $k \times 1$  vector of weights  $\omega_1, \omega_2, \omega_3 \dots \omega_k$   
 $X$  is a  $k \times n$  matrix of observations  $x_1, x_2, x_3 \dots x_k$   
 $\beta$  is a scalar

$\hat{Y}$  is the predicted values for the model with parameters  $W$  and  $\beta$

Logistic Regression uses Maximum Likelihood Estimation to find the parameters  $W$  and  $\beta$



## Logistic Regression

Logistic Regression uses Maximum Likelihood Estimation to find the parameters  $W$  and  $\beta$

$$p(y | x; W, \beta)$$

Likelihood of a value  $y$  given a value  $x$  for a model with parameters  $W$  and  $\beta$

## Logistic Regression Cost Function

if  $y = 1$  then  $p(y | x; W, \beta) = \hat{y}$

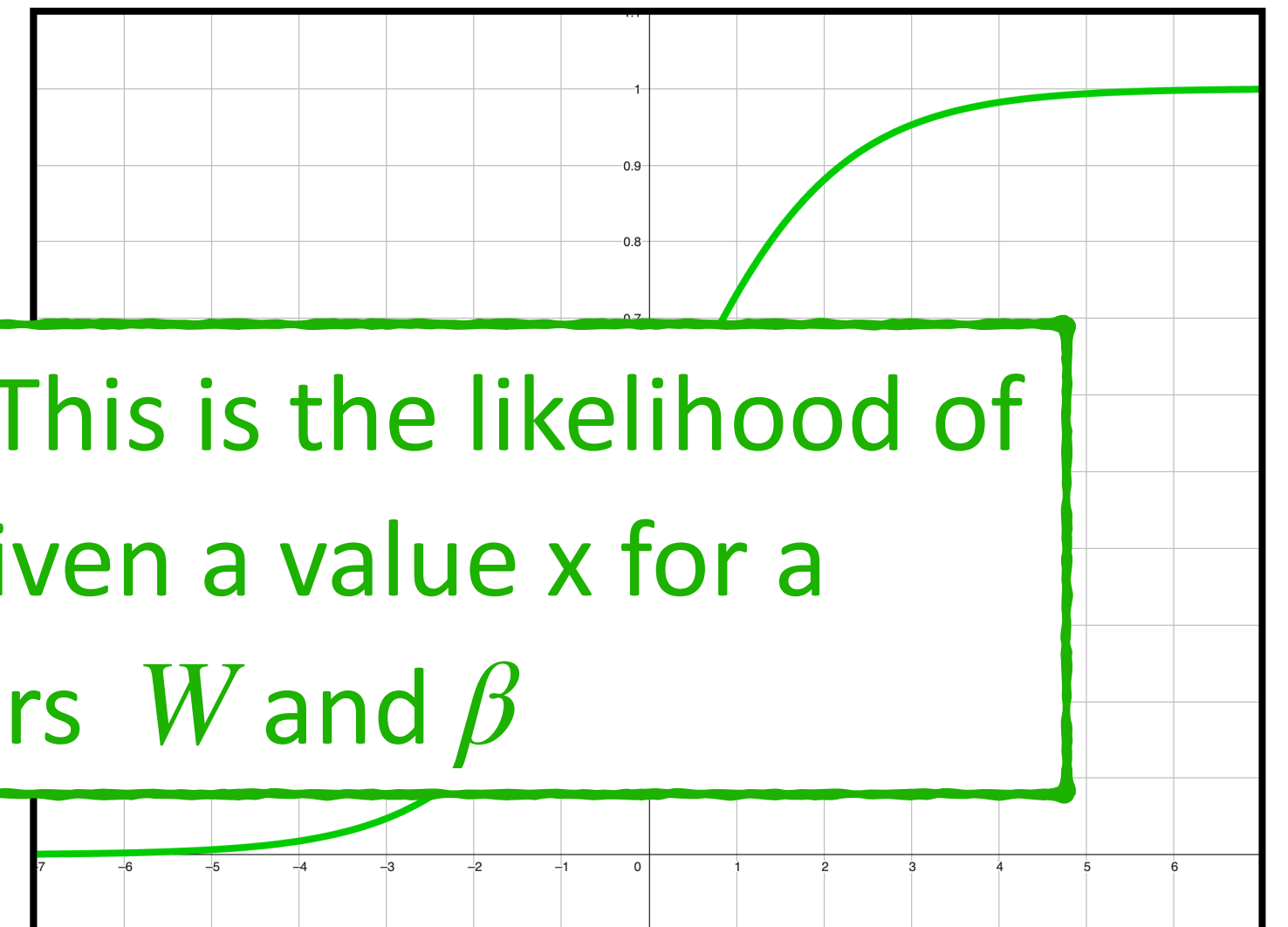
if  $y = 0$  then  $p(y | x; W, \beta) = 1 - \hat{y}$

## Logistic Regression Cost Function

A compact representation of the cost function

$$p(y | x; W, \beta) = \hat{y}^y (1 - \hat{y})^{(1-y)}$$

**Likelihood Function:** This is the likelihood of predicting a value  $y$  given a value  $x$  for a model with parameters  $W$  and  $\beta$





## Logistic Regression

Logistic Regression uses Maximum Likelihood Estimation to find the parameters  $W$  and  $\beta$

### Logistic Regression Cost Function

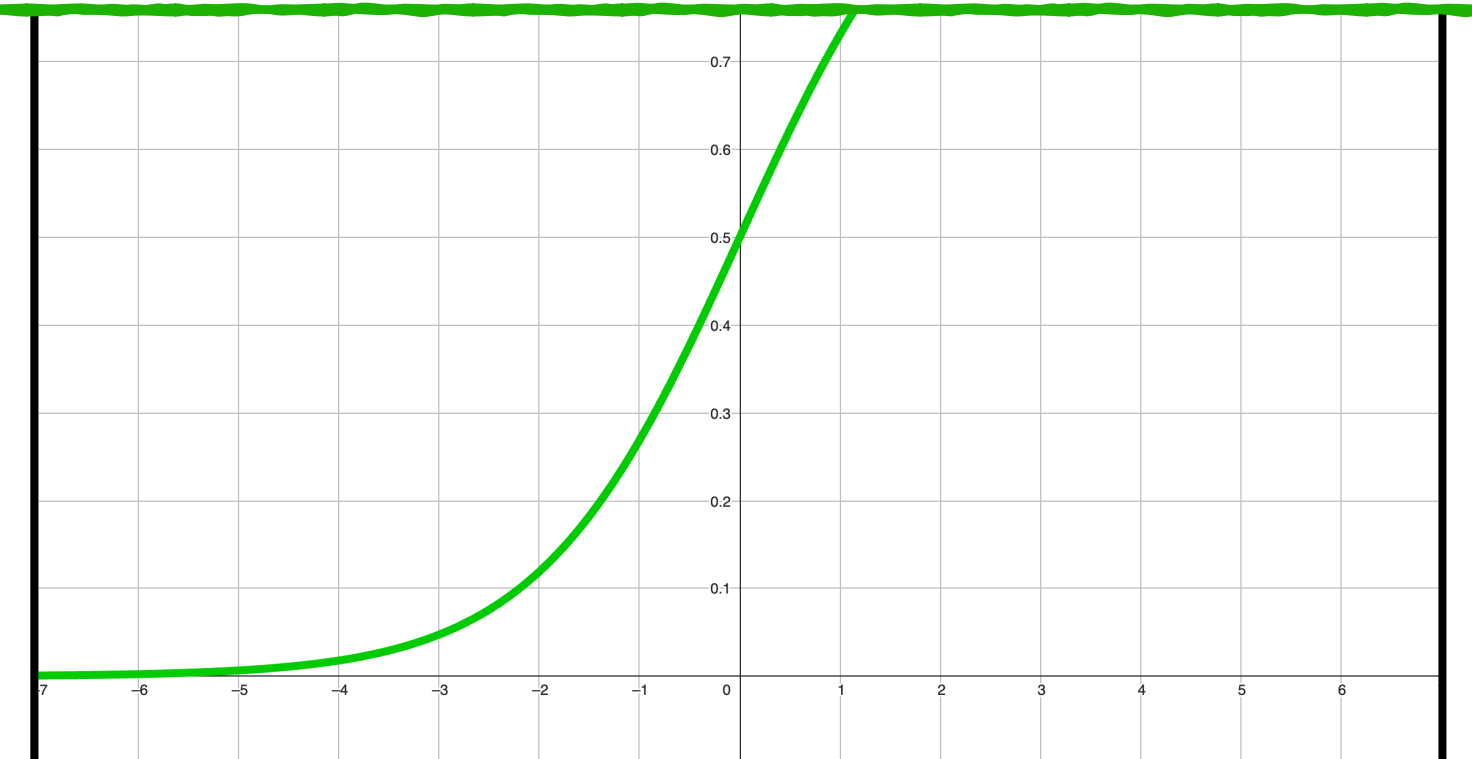
A compact representation of the cost function

$$p(y | x; W, \beta) = \hat{y}^y (1 - \hat{y})^{(1-y)}$$

Maximizing the Likelihood is the same as maximizing the log likelihood

**Likelihood Function:** This is the likelihood of predicting a value  $y$  given a value  $x$  for a model with parameters  $W$  and  $\beta$

We want to find the values of  $W$  and  $\beta$  that maximize this function. Hence Maximum Likelihood Estimation



## Logistic Regression

Logistic Regression uses Maximum Likelihood Estimation to find the parameters  $W$  and  $\beta$

### Logistic Regression Cost Function

A compact representation of the cost function

$$p(y | x; W, \beta) = \hat{y}^y (1 - \hat{y})^{(1-y)}$$

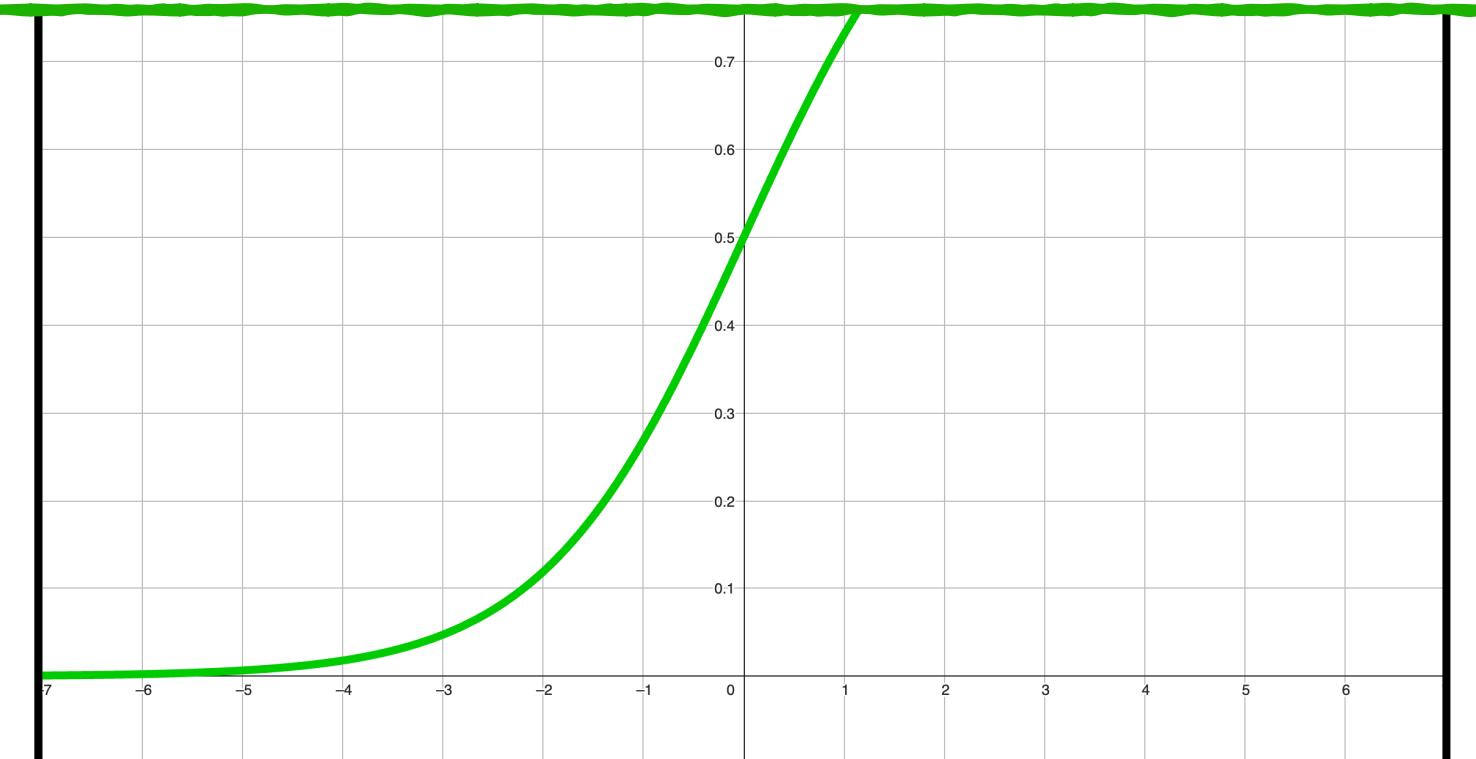
Maximizing the Likelihood is the same as maximizing the log likelihood

$$\begin{aligned} \log_e p(y | x; W, \beta) &= \log_e \hat{y}^y (1 - \hat{y})^{(1-y)} \\ &= y \log_e \hat{y} + (1 - y) \log_e (1 - \hat{y}) \end{aligned}$$

Next we average the cost across all observations...

**Likelihood Function:** This is the likelihood of predicting a value  $y$  given a value  $x$  for a model with parameters  $W$  and  $\beta$

We want to find the values of  $W$  and  $\beta$  that maximize this function. Hence Maximum Likelihood Estimation

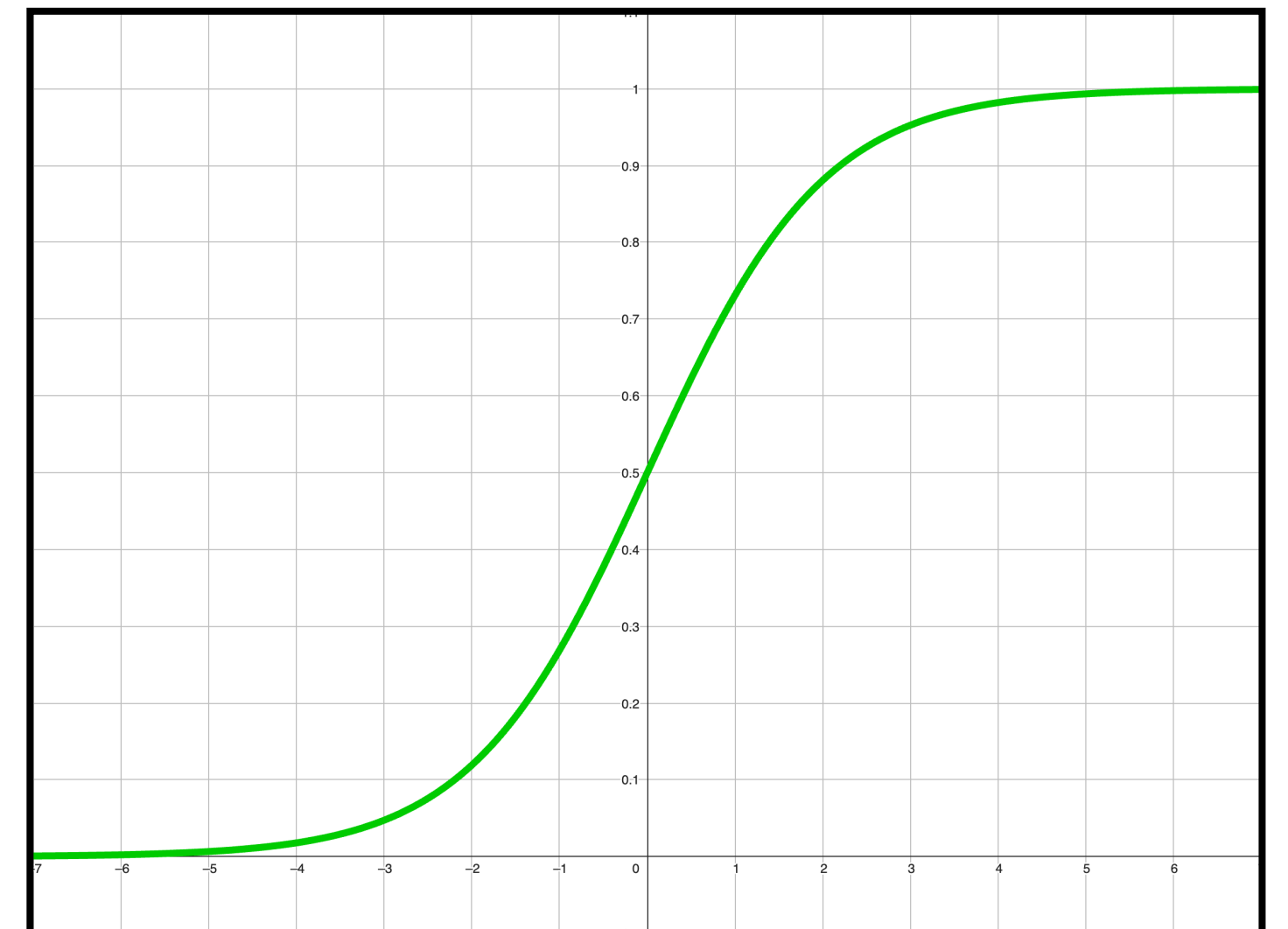


# Logistic Regression

## Logistic Regression Cost Function

$$\hat{Y} = \sigma(W^T X + \beta)$$

$$L(W, \beta) = -\frac{1}{n} \sum_{i=0}^n y \log_e \hat{y} + (1 - y) \log_e (1 - \hat{y})$$

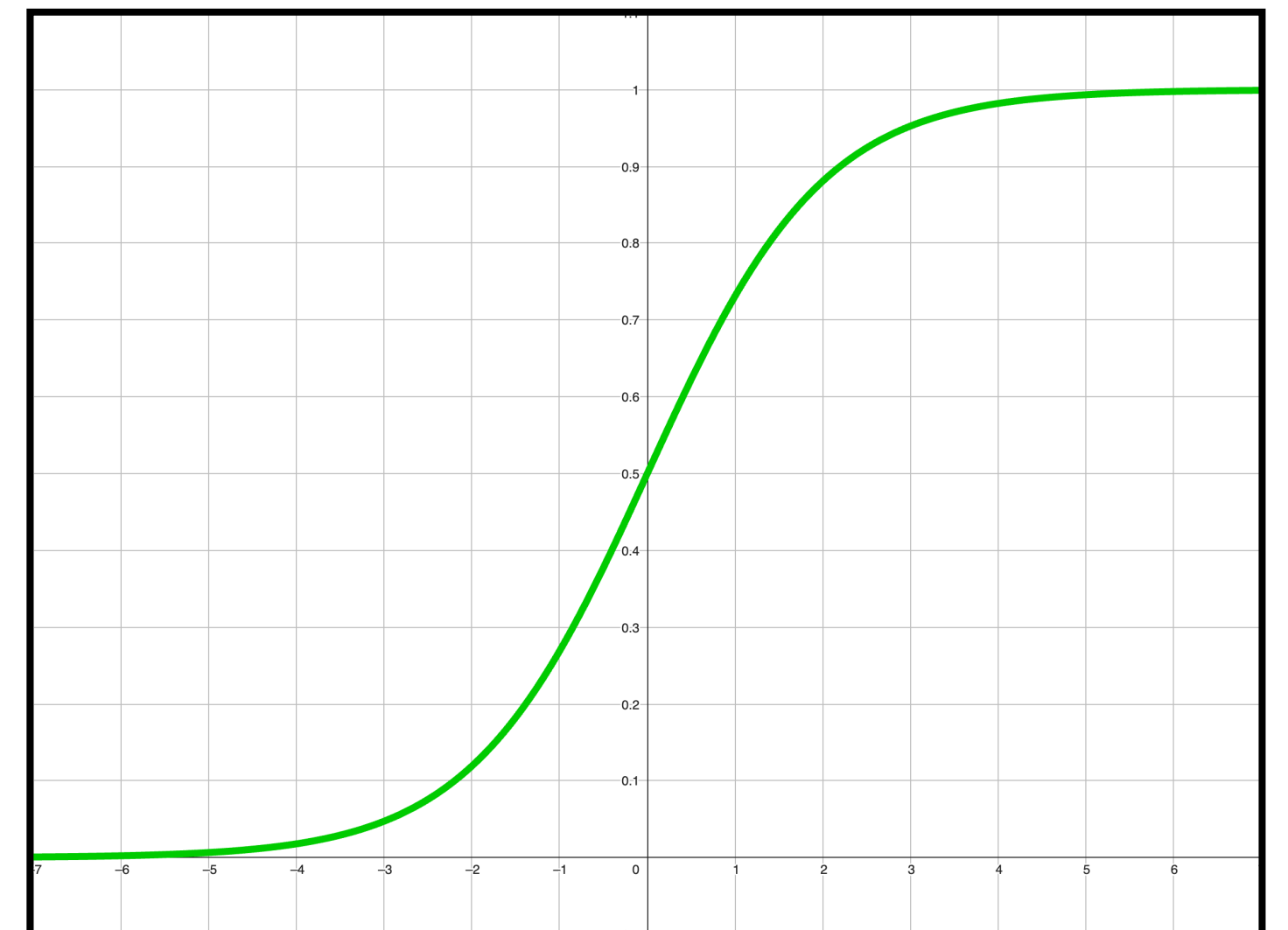


## Logistic Regression Cost Function

$$\hat{Y} = \sigma(W^T X + \beta)$$

$$L(W, \beta) = -\frac{1}{n} \sum_{i=0}^n y \log_e \hat{y} + (1 - y) \log_e (1 - \hat{y})$$

The negative sign indicates that we want to minimize the cost (aka maximize the likelihood)



# Logistic Regression

## Logistic Regression Cost Function

$$\hat{Y} = \sigma(W^T X + \beta)$$

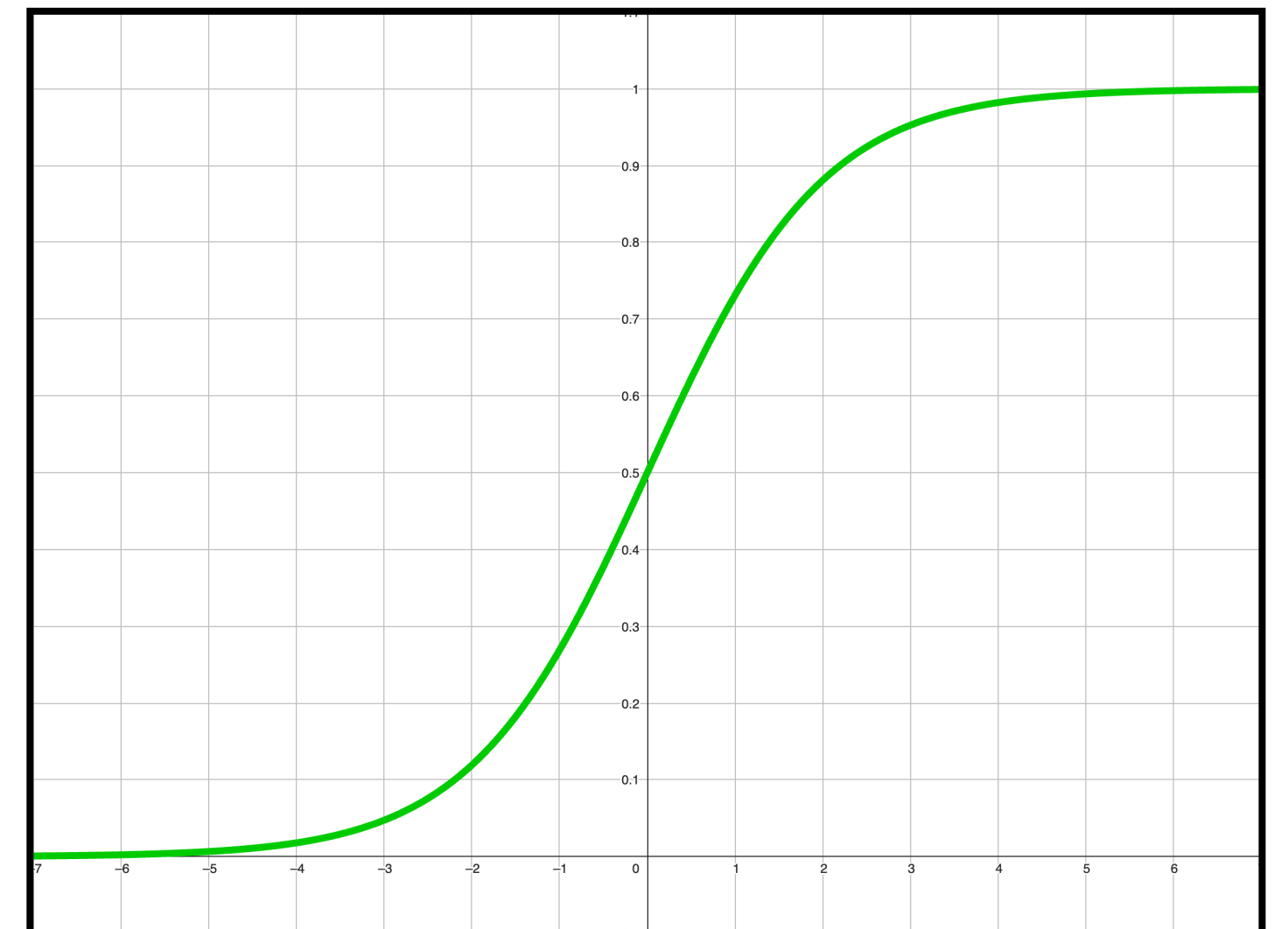
$$L(W, \beta) = -\frac{1}{n} \sum_{i=0}^n y \log_e \hat{y} + (1 - y) \log_e (1 - \hat{y})$$

## Partial Derivatives of the Cost Function w.r.t $W$ and $\beta$

$$\frac{\partial}{\partial W} L(W, \beta) = (\hat{Y} - Y) X$$

$$\frac{\partial}{\partial \beta} L(W, \beta) = (\hat{Y} - Y)$$

With the cost function and the first derivatives above we can use gradient descent to estimate the values of  $W$  and  $\beta$  that minimize the cost



# Logistic Regression

## Logistic Regression Cost Function

$$\hat{Y} = \sigma(W^T X + \beta)$$

$$L(W, \beta) = -\frac{1}{n} \sum_{i=0}^n y \log_e \hat{y} + (1 - y) \log_e (1 - \hat{y})$$

## Partial Derivatives of the Cost Function w.r.t $W$ and $\beta$

$$\frac{\partial}{\partial W} L(W, \beta) = (\hat{Y} - Y) X$$

$$\frac{\partial}{\partial \beta} L(W, \beta) = (\hat{Y} - Y)$$

## Gradient Descent for Logistic Regression

**Step 1:** Start with random values for  $W$  and  $\beta$



# Logistic Regression

## Logistic Regression Cost Function

$$\hat{Y} = \sigma(W^T X + \beta)$$

$$L(W, \beta) = -\frac{1}{n} \sum_{i=0}^n y \log_e \hat{y} + (1 - y) \log_e (1 - \hat{y})$$

## Partial Derivatives of the Cost Function w.r.t $W$ and $\beta$

$$\frac{\partial}{\partial W} L(W, \beta) = (\hat{Y} - Y) X$$

$$\frac{\partial}{\partial \beta} L(W, \beta) = (\hat{Y} - Y)$$

## Gradient Descent for Logistic Regression

Step 1: Start with random values for  $W$  and  $\beta$

**Step 2:** Compute the partial derivative of the cost function w.r.t  $W$  and  $\beta$

$$\frac{\partial}{\partial W} L(W, \beta) = (\hat{Y} - Y) X$$

$$\frac{\partial}{\partial \beta} L(W, \beta) = (\hat{Y} - Y)$$



# Logistic Regression

## Logistic Regression Cost Function

$$\hat{Y} = \sigma(W^T X + \beta)$$

$$L(W, \beta) = -\frac{1}{n} \sum_{i=0}^n y \log_e \hat{y} + (1 - y) \log_e (1 - \hat{y})$$

## Partial Derivatives of the Cost Function w.r.t $W$ and $\beta$

$$\frac{\partial}{\partial W} L(W, \beta) = (\hat{Y} - Y) X$$

$$\frac{\partial}{\partial \beta} L(W, \beta) = (\hat{Y} - Y)$$

## Gradient Descent for Logistic Regression

Step 1: Start with random values for  $W$  and  $\beta$

Step 2: Compute the partial derivative of the cost function w.r.t  $W$  and  $\beta$

$$\frac{\partial}{\partial W} L(W, \beta) = (\hat{Y} - Y) X$$

$$\frac{\partial}{\partial \beta} L(W, \beta) = (\hat{Y} - Y)$$

**Step 3:** Calculate a step size that is proportional to the slope

$$step\_size_W = \frac{\partial}{\partial W} L(W, \beta) \times learning\_rate$$

$$step\_size_\beta = \frac{\partial}{\partial \beta} L(W, \beta) \times learning\_rate$$

# Logistic Regression

## Logistic Regression Cost Function

$$\hat{Y} = \sigma(W^T X + \beta)$$

$$L(W, \beta) = -\frac{1}{n} \sum_{i=0}^n y \log_e \hat{y} + (1 - y) \log_e (1 - \hat{y})$$

## Partial Derivatives of the Cost Function w.r.t $W$ and $\beta$

$$\frac{\partial}{\partial W} L(W, \beta) = (\hat{Y} - Y) X$$

$$\frac{\partial}{\partial \beta} L(W, \beta) = (\hat{Y} - Y)$$

## Gradient Descent for Logistic Regression

Step 1: Start with random values for  $W$  and  $\beta$

Step 2: Compute the partial derivative of the cost function w.r.t  $W$  and  $\beta$

$$\frac{\partial}{\partial W} L(W, \beta) = (\hat{Y} - Y) X$$

$$\frac{\partial}{\partial \beta} L(W, \beta) = (\hat{Y} - Y)$$

Step 3: Calculate a step size that is proportional to the slope

$$step\_size_W = \frac{\partial}{\partial W} L(W, \beta) \times learning\_rate$$

$$step\_size_\beta = \frac{\partial}{\partial \beta} L(W, \beta) \times learning\_rate$$

**Step 4:** Calculate new values for  $W$  and  $\beta$  by subtracting the step size

$$W = W - step\_size_W$$

$$\beta = \beta - step\_size_\beta$$

# Logistic Regression

## Logistic Regression Cost Function

$$\hat{Y} = \sigma(W^T X + \beta)$$

$$L(W, \beta) = -\frac{1}{n} \sum_{i=0}^n y \log_e \hat{y} + (1 - y) \log_e (1 - \hat{y})$$

## Partial Derivatives of the Cost Function w.r.t $W$ and $\beta$

$$\frac{\partial}{\partial W} L(W, \beta) = (\hat{Y} - Y) X$$

$$\frac{\partial}{\partial \beta} L(W, \beta) = (\hat{Y} - Y)$$

## Gradient Descent for Logistic Regression

Step 1: Start with random values for  $W$  and  $\beta$

Step 2: Compute the partial derivative of the cost function w.r.t  $W$  and  $\beta$

$$\frac{\partial}{\partial W} L(W, \beta) = (\hat{Y} - Y) X$$

$$\frac{\partial}{\partial \beta} L(W, \beta) = (\hat{Y} - Y)$$

Step 3: Calculate a step size that is proportional to the slope

$$step\_size_W = \frac{\partial}{\partial W} L(W, \beta) \times learning\_rate$$

$$step\_size_\beta = \frac{\partial}{\partial \beta} L(W, \beta) \times learning\_rate$$

Step 4: Calculate new values for  $W$  and  $\beta$  by subtracting the step size

$$W = W - step\_size$$

$$\beta = \beta - step\_size$$

**Step 5: Go to step 2 and repeat**

# Logistic Regression

## Logistic Regression Cost Function

$$\hat{Y} = \sigma(W^T X + \beta)$$

$$L(W, \beta) = -\frac{1}{n} \sum_{i=0}^n y \log_e \hat{y} + (1 - y) \log_e (1 - \hat{y})$$

## Partial Derivatives of the Cost Function w.r.t $W$ and $\beta$

$$\frac{\partial}{\partial W} L(W, \beta) = (\hat{Y} - Y) X$$

$$\frac{\partial}{\partial \beta} L(W, \beta) = (\hat{Y} - Y)$$

## Gradient Descent for Logistic Regression

Step 1: Start with random values for  $W$  and  $\beta$

Step 2: Compute the partial derivative of the cost function w.r.t  $W$  and  $\beta$

**Gradient Descent continues in this manner until the step size is close to zero or a fixed number of iterations**

$$step\_size_{\beta} = \frac{\partial}{\partial \beta} L(W, \beta) \times learning\_rate$$

Step 4: Calculate new values for  $W$  and  $\beta$  by subtracting the step size

$$W = W - step\_size$$

$$\beta = \beta - step\_size$$

Step 5: Go to step 2 and repeat

# Related Tutorials & Textbooks

## Logistic Regression ↗

An introduction to Logistic Regression. A Logistic Regression model use used to predict a binary value (the dependent variable) for one or more independent variables using a threshold to classify a probability.

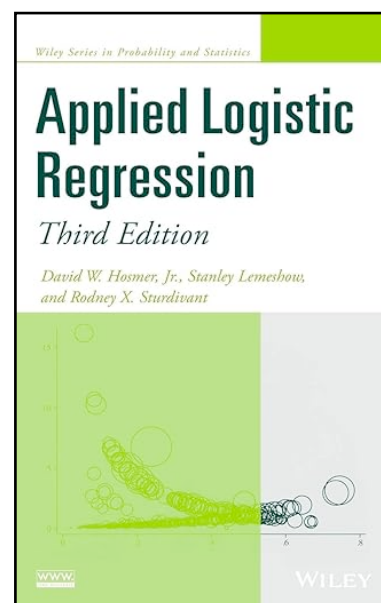
## Multiple Regression ↗

Multiple regression extends the two dimensional linear model introduced in Simple Linear Regression to  $k + 1$  dimensions with one dependent variable,  $k$  independent variables and  $k+1$  parameters.

## Gradient Descent for Multiple Regression ↗

Gradient Descent algorithm for multiple regression and how it can be used to optimize  $k + 1$  parameters for a Linear model in multiple dimensions.

## Recommended Textbooks



### **Applied Logistic Regression**

by David W. Hosmer Jr., Stanley Lemeshow, Rodney X. Sturdivant

**For a complete list of tutorials see:**

<https://arrsingh.com/ai-tutorials>