

# Day-85, feb 23, 2025 ( Falgun 11, 2081 B.S.)

- ① Regression with a Perceptron
- ② Regression with a Perceptron - loss function
- ③ Regression with a Perceptron - Gradient Descent
- ④ Classification with Perceptron
- ⑤ Classification with Perceptron - The Sigmoid function
- ⑥ Classification with Perceptron - Gradient Descent
- ⑦ Classification with  $\pi \rightarrow$  Calculating the derivatives

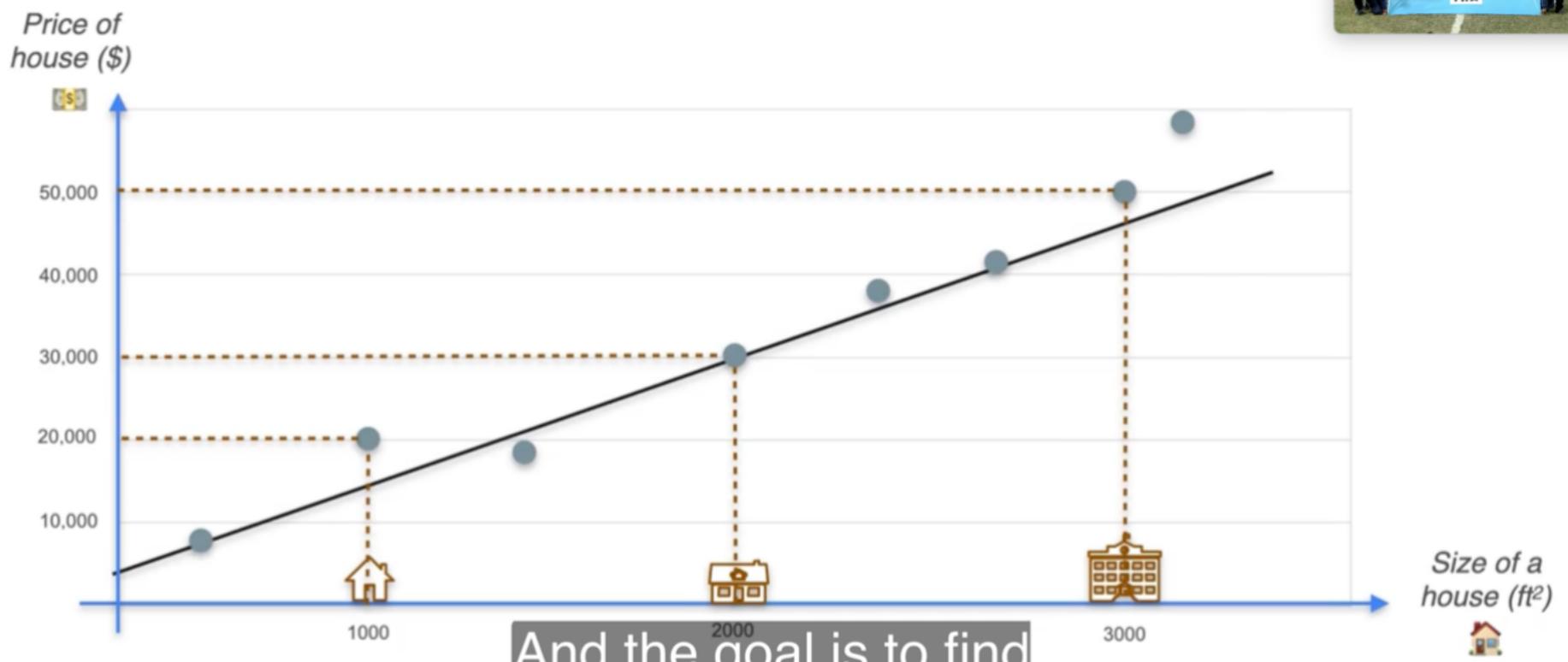
# # Regression with a Perceptron

## Regression With a Perceptron

	Size of a house (ft <sup>2</sup> )	Number of rooms	Price of house (\$)
House	1000 ft <sup>2</sup>	2	\$20,000
House	2000 ft <sup>2</sup>	4	\$30,000
House	3000 ft <sup>2</sup>	7	\$50,000

Well here's where the perception comes in.

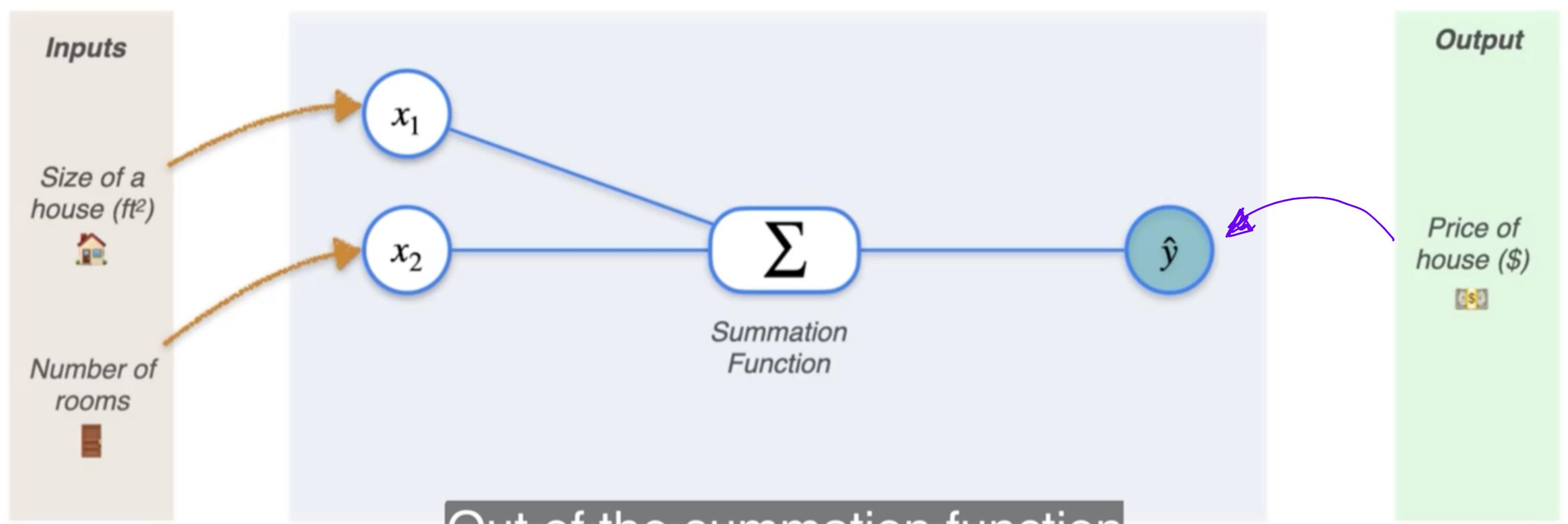
## Regression Problem Motivation



My goal: find weights and bias that will optimize the predictions.  
e.g. Reduce the errors in the predictions using loss function.

# Regression With a Perceptron

Single Layer Neural Network Perceptron

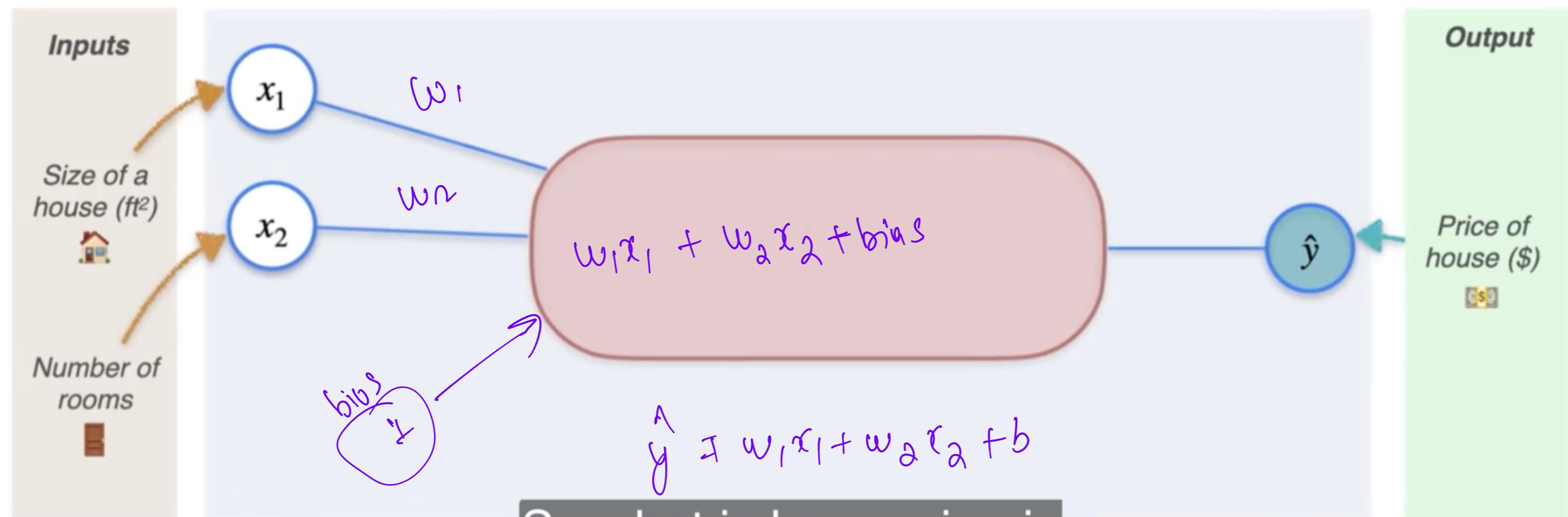


Out of the summation function  
comes the output  $y$  hat and

# Regression With a Perceptron



Single Layer Neural Network Perceptron

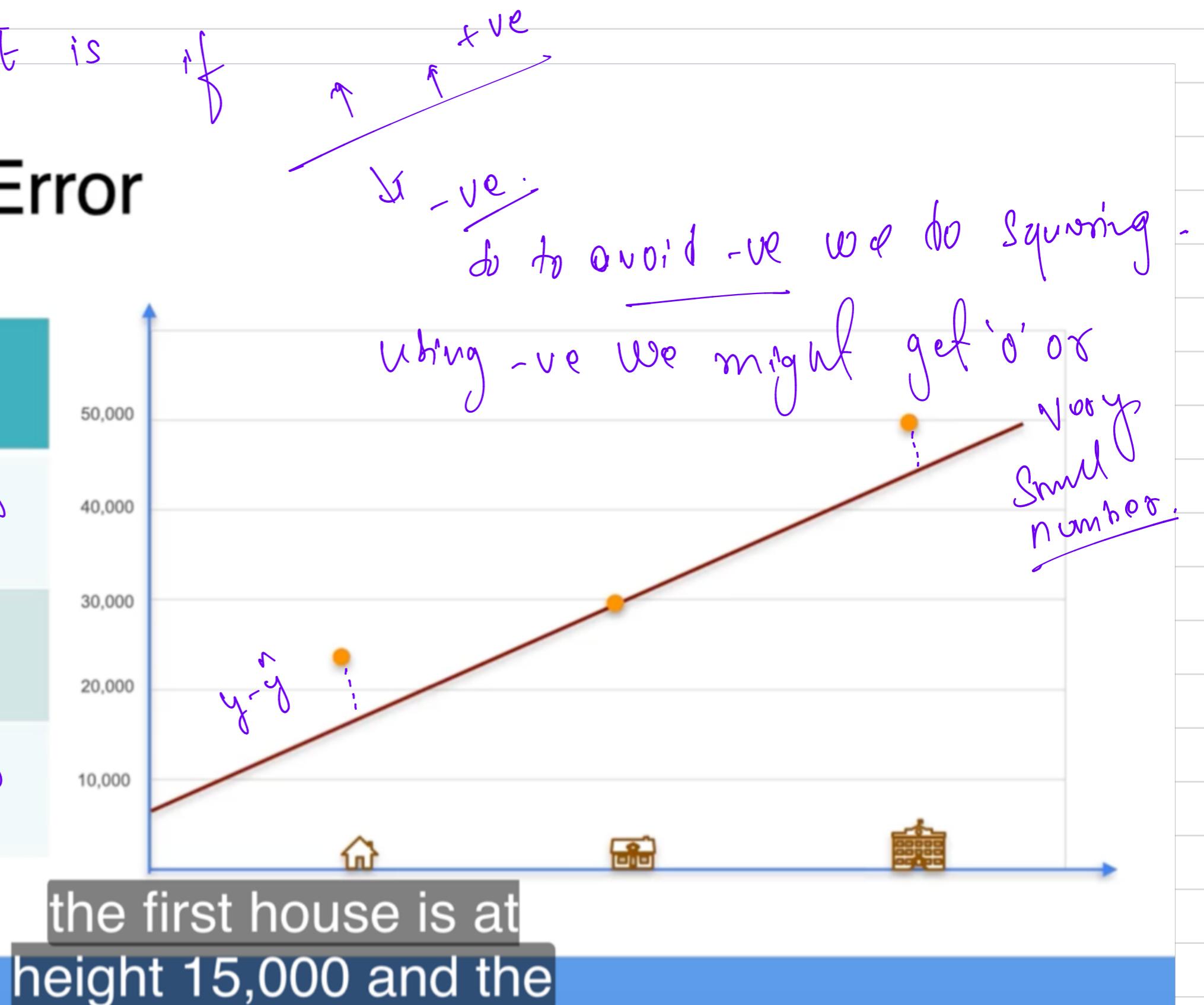


So what is happening in  
this summation step?

problem with MSE is if

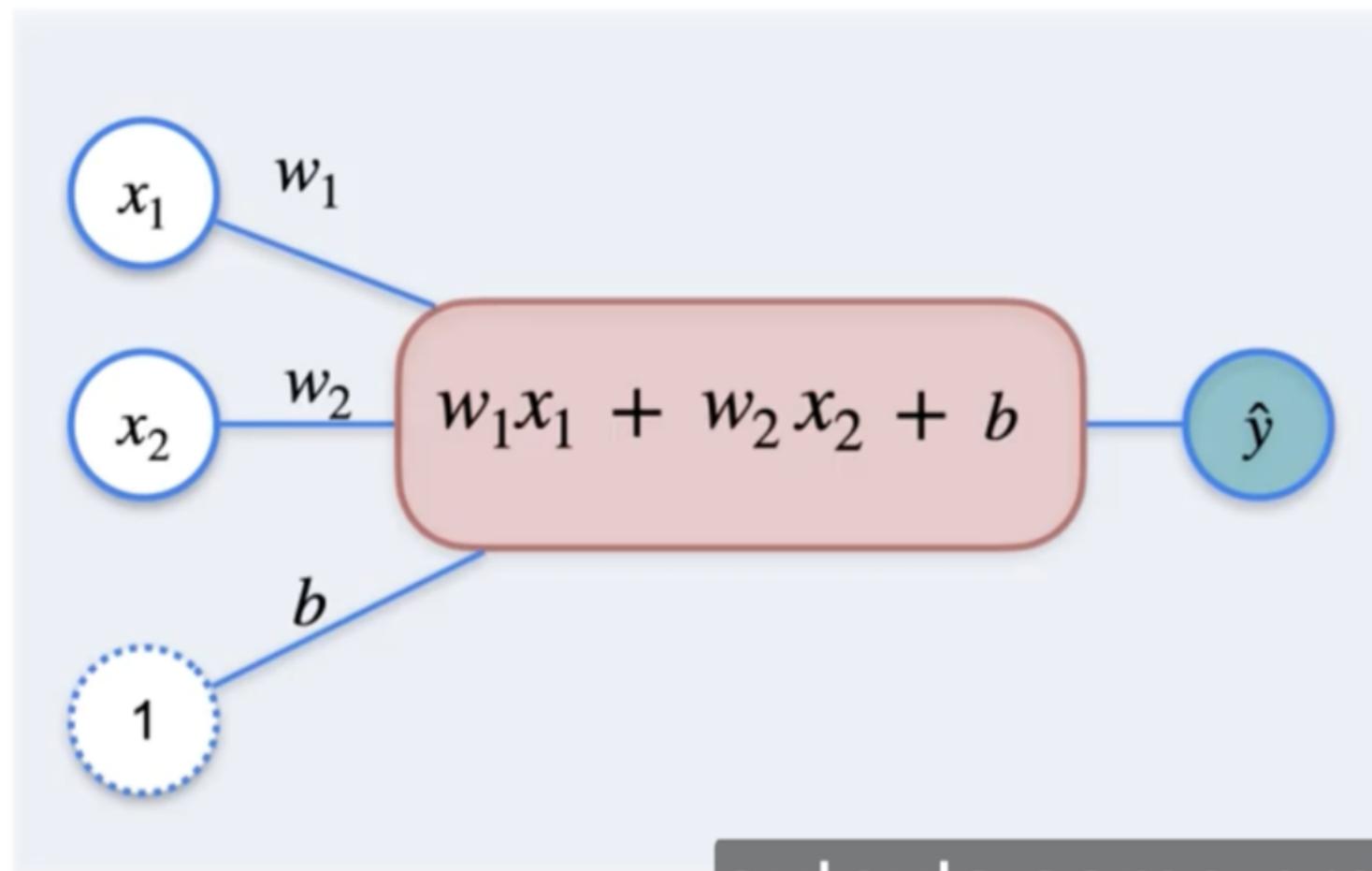
## Mean Squared Error

	$y$	$\hat{y}$	Error $y - \hat{y}$
House	\$20,000	\$15,000	\$5000
House	\$30,000	\$30,000	0
School	\$50,000	\$45,000	\$5000



# Regression With a Perceptron

Single Layer Neural Network Perceptron



outputs some prediction  $y$  and  
is given by this formula.

**Prediction Function:**

$$\hat{y} = w_1x_1 + w_2x_2 + b$$

loss function:  
$$L(y, \hat{y}) = \frac{1}{2} (y - \hat{y})^2$$

Min Goal:

find  $w_1, w_2, b$  that  
gives  $\hat{y}$  with  
the least error.

Best Model that makes smallest mistakes!

# Regression With a Perceptron

**Prediction Function:**

$$\hat{y} = w_1x_1 + w_2x_2 + b$$

**Loss Function:**

$$L(y, \hat{y}) = \frac{1}{2}(y - \hat{y})^2$$

**Main Goal:**

Find  $w_1, w_2, b$  that give  $\hat{y}$  with the least error

To find optimal values for:

$$w_1, w_2, b$$

You need gradient descent

$$w_1 \rightarrow w_1 - \alpha \cdot \frac{\partial L}{\partial w_1}$$

$$w_2 \rightarrow w_2 - \alpha \cdot \frac{\partial L}{\partial w_2}$$

$$b \rightarrow b - \alpha \cdot \frac{\partial L}{\partial b}$$

Now we know what the problem is.

# Regression With a Perceptron

Prediction Function:

$$\hat{y} = w_1x_1 + w_2x_2 + b$$

Loss Function:

$$L(y, \hat{y}) = \frac{1}{2}(y - \hat{y})^2$$

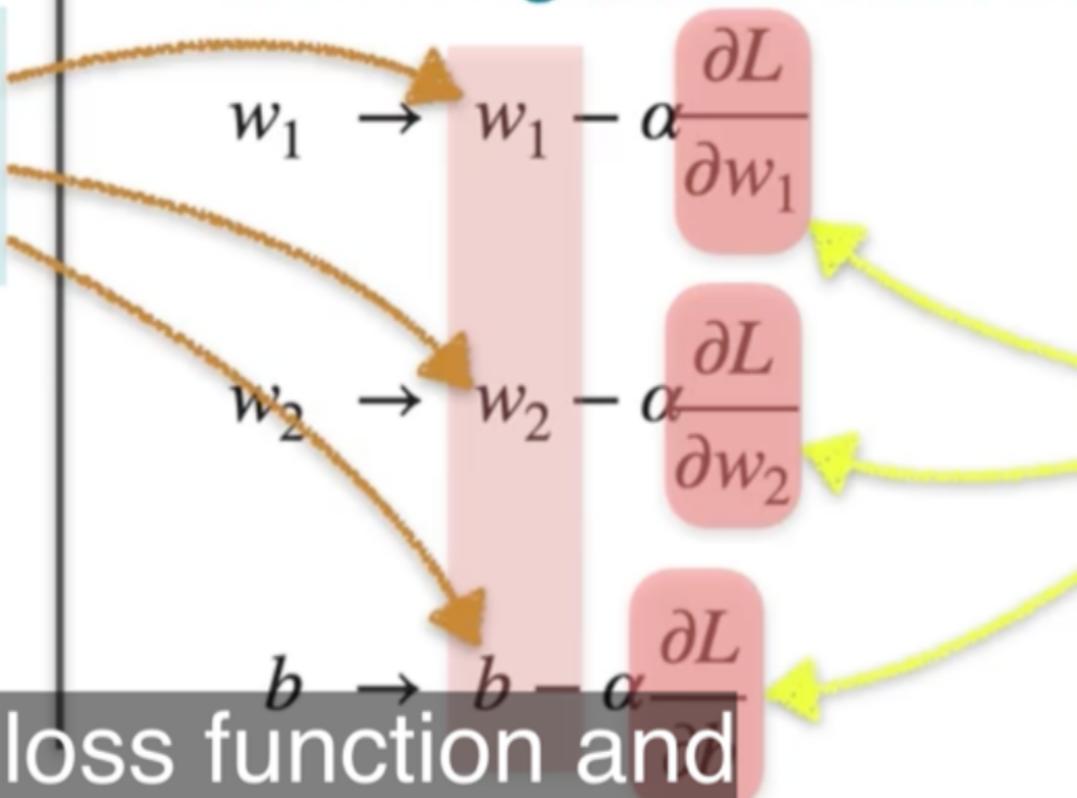
Main Goal:

Find  $w_1, w_2, b$  that give  $\hat{y}$  with the least error

To find optimal values for:  
 $w_1, w_2, b$

**You need gradient descent**

*Some initial starting values*



a very low loss function and  
that implies a good model.

# Regression With Perceptron

Prediction Function:

$$\hat{y} = w_1x_1 + w_2x_2 + b$$

Loss Function:

$$L(y, \hat{y}) = \frac{1}{2}(y - \hat{y})^2$$

depends on

Using chain rule:

$$\frac{\partial L}{\partial b} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial b}$$

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial w_1}$$

$$\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial b}$$

First of all, what  
is dL over db?

# Regression With Perceptron

**Prediction Function:**

$$\hat{y} = w_1x_1 + w_2x_2 + b$$

**Loss Function:**

$$L(y, \hat{y}) = \frac{1}{2}(y - \hat{y})^2$$

actually plugging  
in the whole thing.

DeepLearning.AI

Using chain rule:

$$\frac{\partial L}{\partial b} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial b}$$

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial w_1}$$

$$\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial w_2}$$

Let's find the  $\frac{\partial L}{\partial \hat{y}}$ ,  $\frac{\partial \hat{y}}{\partial w_1}$  and  
 $\frac{\partial \hat{y}}{\partial w_2}$  then we get,

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

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

$$\Rightarrow \frac{1}{2} \cdot 2(y - \hat{y}) \cdot (-1)$$

# Regression With Perceptron

**Prediction Function:**

$$\hat{y} = w_1x_1 + w_2x_2 + b$$

**Loss Function:**

$$L(y, \hat{y}) = \frac{1}{2}(y - \hat{y})^2$$

$\frac{\partial L}{\partial \hat{y}}$	$= -(y - \hat{y})$
$\frac{\partial \hat{y}}{\partial b}$	$= 1$
$\frac{\partial \hat{y}}{\partial w_1}$	$= x_1$
$\frac{\partial \hat{y}}{\partial w_2}$	$= x_2$

the constant accompany w\_2.

# Regression With Perceptron

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

$$\frac{\partial \hat{y}}{\partial b} = 1$$

$$\frac{\partial \hat{y}}{\partial w_1} = x_1$$

$$\frac{\partial \hat{y}}{\partial w_2} = x_2$$

$$\frac{\partial L}{\partial b} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial b}$$

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial w_1}$$

$$\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial w_2}$$

$$-(y - \hat{y}) \cdot 1$$

$$-(y - \hat{y}) \cdot x_1$$

$$-(y - \hat{y}) \cdot x_2$$

# Regression With a Perceptron

**Main Goal:**

Find  $w_1, w_2, b$  that give  $\hat{y}$  with the least error

**i.e. optimal values for:**

$w_1, w_2, b$

Perform Gradient Descent?

$$w_1 \rightarrow w_1 - \alpha \frac{\partial L}{\partial w_1}$$

$$\Rightarrow w_1 - \alpha (-x_1(y - \hat{y}))$$

$$w_2 \rightarrow w_2 - \alpha (-x_2(y - \hat{y}))$$

$$b \rightarrow b - \alpha (- (y - \hat{y}))$$

the predictions  $y$  had with  
the smallest possible error.

# Classification Problem Motivation

Sentence	Aack	Beep	Mood
Aack aack aack!	3	0	Happy
Beep beep!	0	2	Sad
Aack beep beep beep!	1	3	Sad
Aack beep aack!	2	1	Happy

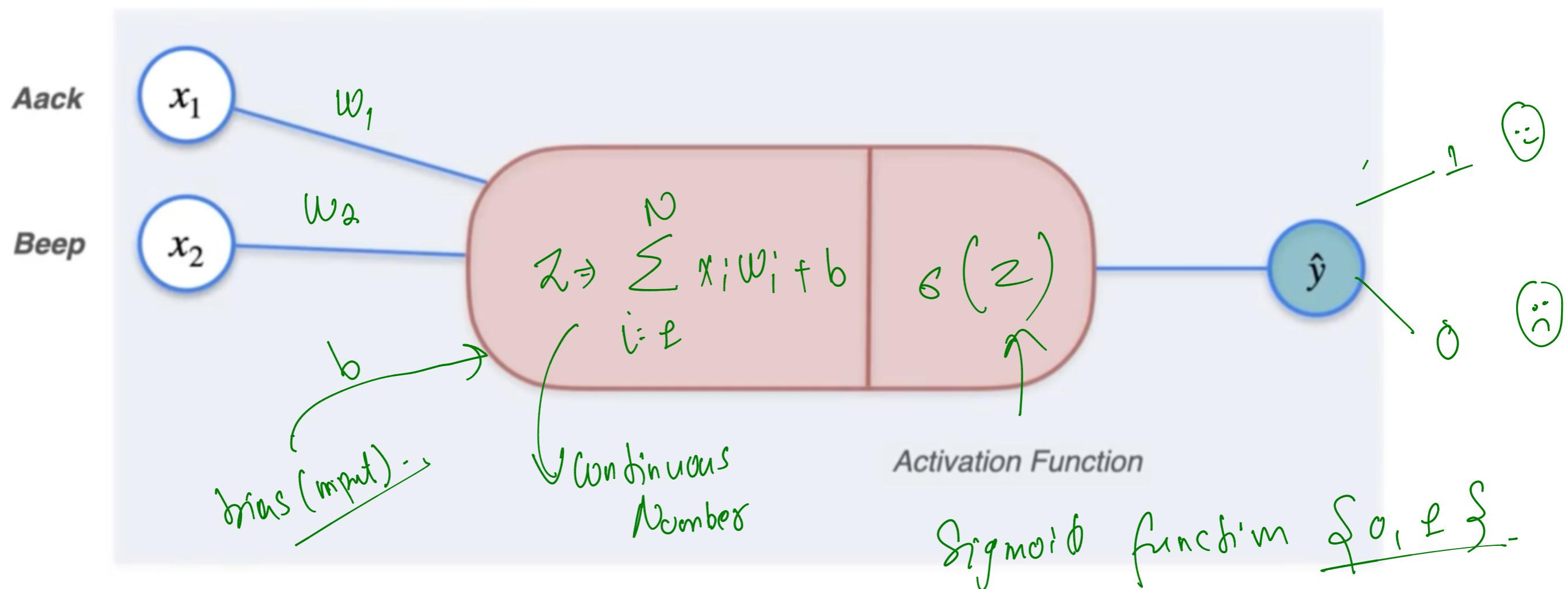
to tell the mood of  
the each sentence.

# Classification Problem Motivation



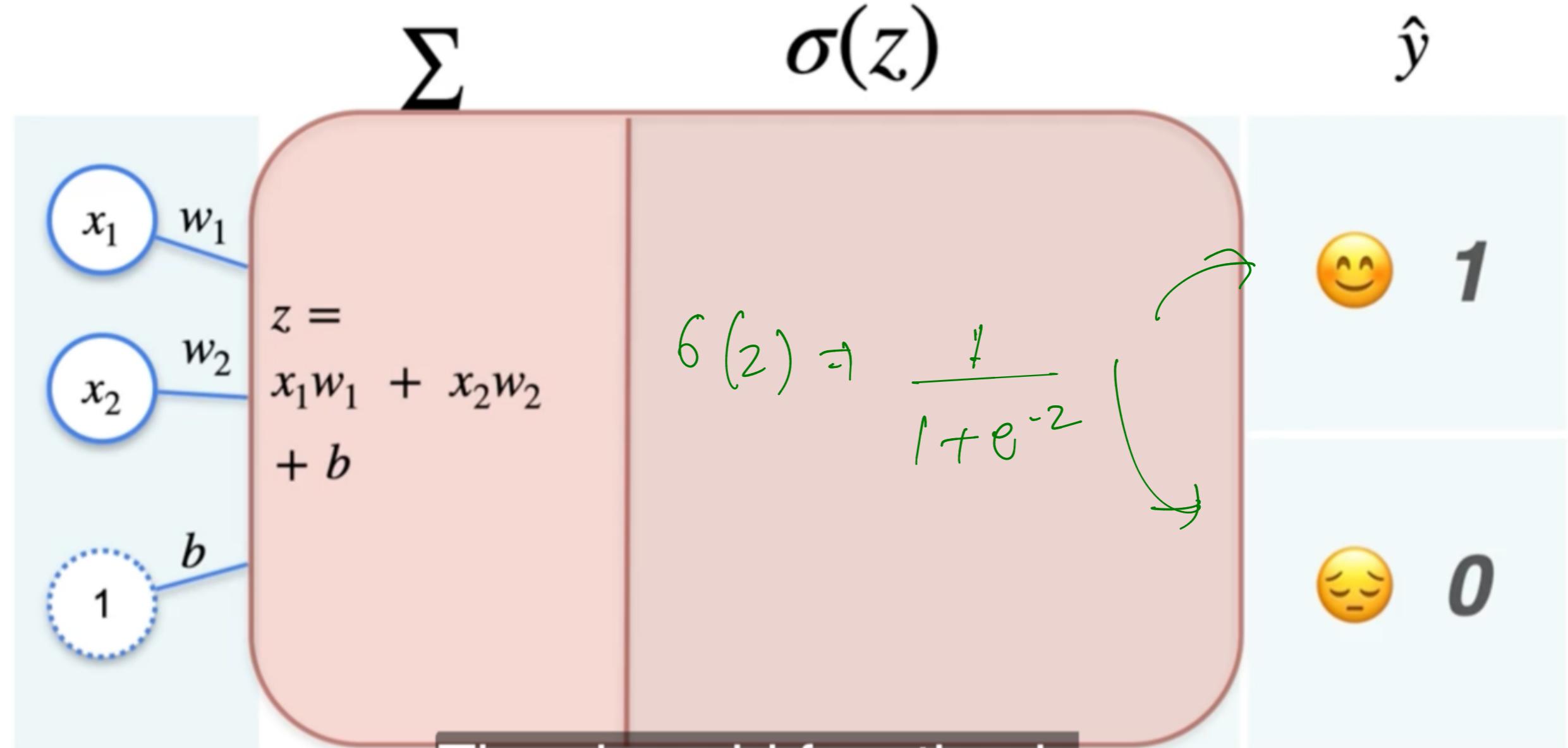
# Classification With a Perceptron

Single Layer Neural Network Perceptron





# Sigmoid Function



The sigmoid function is  
given by the formula

# Sigmoid Function

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

$$\delta(z) = (1 + e^{-z})^{-2}$$

$$\frac{d}{dz} \delta(z) = \frac{d}{dz} (1 + e^{-z})^{-2}$$



The graph shows the sigmoid function  $\sigma(z)$  plotted against  $z$ . The curve is a blue S-shape that passes through the point  $(0, 0.5)$ . The horizontal axis is labeled with  $z$  at the origin and  $1$  at the top. The vertical axis is labeled with  $0$  at the bottom and  $1$  at the top. A green arrow points from the text "the entire number line and" towards the left side of the graph, where the curve is very flat near the value of 0.

the entire number line and

$$\frac{d}{dz} \ln(z) = -1 (1 + e^{-z})^{-1-1} \left( \frac{d}{dz} (1 + e^{-z}) \right)$$

$$\Rightarrow -1 (1 + e^{-z})^{-2} \left( \frac{d}{dz} (1) + \frac{d}{dz} (e^{-z}) \right)$$

$$\Rightarrow -1 (1 + e^{-z})^{-2} \left( 0 + e^{-z} \left( \frac{d}{dz} (-z) \right) \right)$$

$$\left[ \frac{d}{dz} \ln(z) \Rightarrow -1 (1 + e^{-z})^{-2} (e^{-z}) (-1) \right]$$

$$\Rightarrow (1 + e^{-z})^{-2} (e^{-z})$$

$$\left[ \frac{-1}{(1 + e^{-z})^2} (e^{-z}) \right]$$

$$\Rightarrow \frac{e^{-z}}{(1 + e^{-z})^2}$$

$$\frac{d}{dz} f(z) \Rightarrow \frac{e^{-2} + 1 - 1}{(1+e^{-2})^2}$$

$$\Rightarrow \frac{1 + e^{-z} - 1}{(1+e^{-2})^2}$$

$$\Rightarrow \frac{1 + e^{-z}}{(1+e^{-2})^2} - \frac{1}{(1+e^{-2})^2}$$

$$\Rightarrow \frac{1}{(1+e^{-2})} - \left( \frac{1}{1+e^{-2}} \right) \left( \frac{1}{1+e^{-2}} \right)$$

$$\Rightarrow \frac{1}{1+e^{-2}} \left( 1 - \frac{1}{1+e^{-2}} \right)$$

$$\text{So, } \sigma(z) = \frac{1}{1 + e^{-z}}$$

$$\frac{d}{dz} \sigma(z) = \sigma(z) (1 - \sigma(z))$$



## Derivative of a Sigmoid Function

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

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

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

$$\frac{d}{dz} \sigma(z) = -1 (1 + e^{-z})^{-1-1} \left( \frac{d}{dz} (1 + e^{-z}) \right)$$

$$= -1 (1 + e^{-z})^{-2} \left( \frac{d}{dz} (1) + \frac{d}{dz} (e^{-z}) \right)$$

$$= -1 (1 + e^{-z})^{-2} (0 + e^{-z} \left( \frac{d}{dz} (-z) \right))$$

$$= -1 (1 + e^{-z})^{-2} (e^{-z}) (-1)$$



## Derivative of a Sigmoid Function

$$\frac{d}{dz} \sigma(z) = -1 (1 + e^{-z})^{-2} (e^{-z}) (-1)$$

$$= (1 + e^{-z})^{-2} (e^{-z})$$

$$= \frac{1}{(1 + e^{-z})^2} (e^{-z})$$

$$= \frac{e^{-z}}{(1 + e^{-z})^2}$$



# Derivative of a Sigmoid Function

$$\begin{aligned}\frac{d}{dz} \sigma(z) &= \frac{e^{-z}}{(1 + e^{-z})^2} + 1 - 1 \\&= \frac{1 + e^{-z} - 1}{(1 + e^{-z})^2} \\&= \frac{1 + e^{-z}}{(1 + e^{-z})^2} - \frac{1}{(1 + e^{-z})^2} \\&= \frac{1}{(1 + e^{-z})} - \frac{1}{(1 + e^{-z})^2}\end{aligned}$$

$$\begin{aligned}\frac{d}{dz} \sigma(z) &= \frac{1}{(1 + e^{-z})} - \left( \frac{1}{(1 + e^{-z})} \right) \left( \frac{1}{(1 + e^{-z})} \right) \\&= \frac{1}{(1 + e^{-z})} \left( 1 - \frac{1}{(1 + e^{-z})} \right)\end{aligned}$$

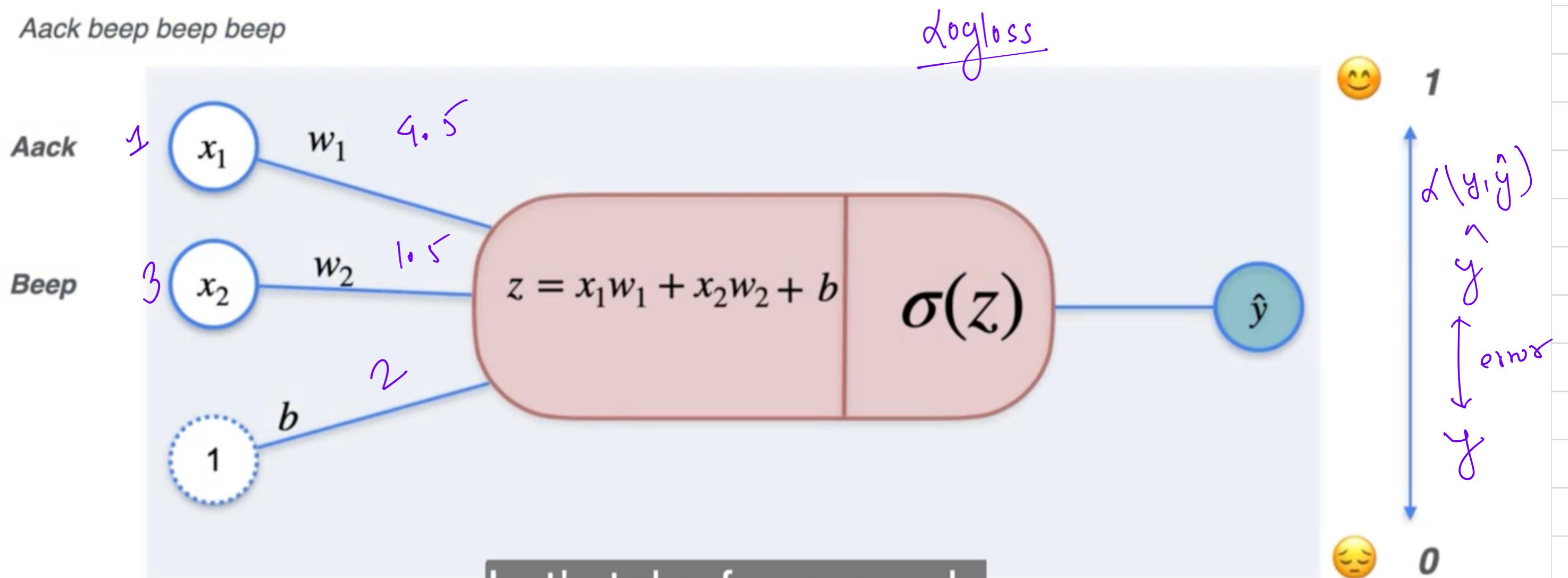
Recall that:  $\sigma(z) = \frac{1}{1 + e^{-z}}$

$$\frac{d}{dz} \sigma(z) = \sigma(z) (1 - \sigma(z))$$

The derivative of sigmoid is

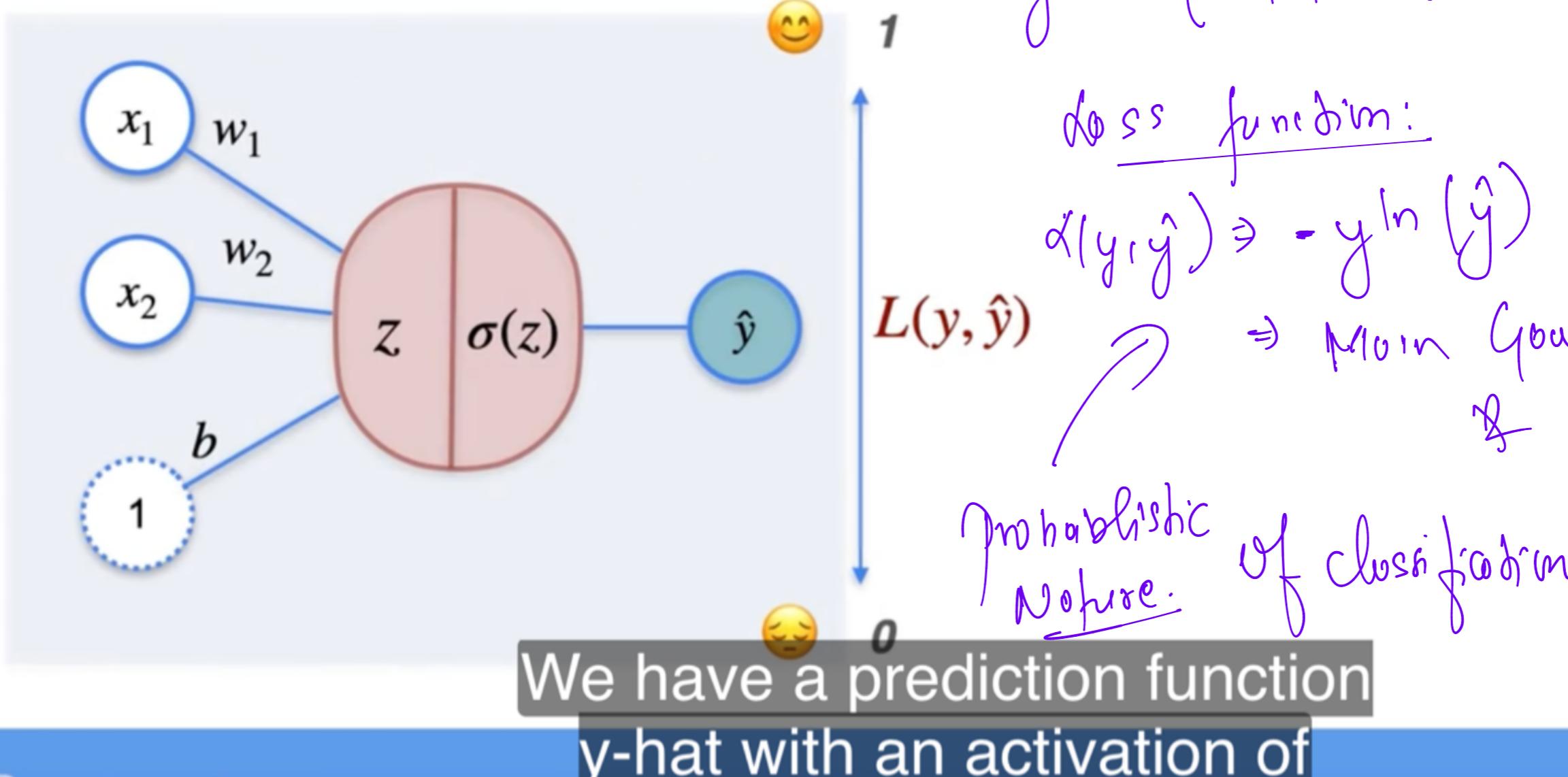
# Classification With a Perceptron

Aack beep beep beep



Let's take for example  
the sentence Aack,

# Classification With a Perceptron



$$\hat{y} = \sigma(w_1x_1 + w_2x_2 + b)$$

loss function:

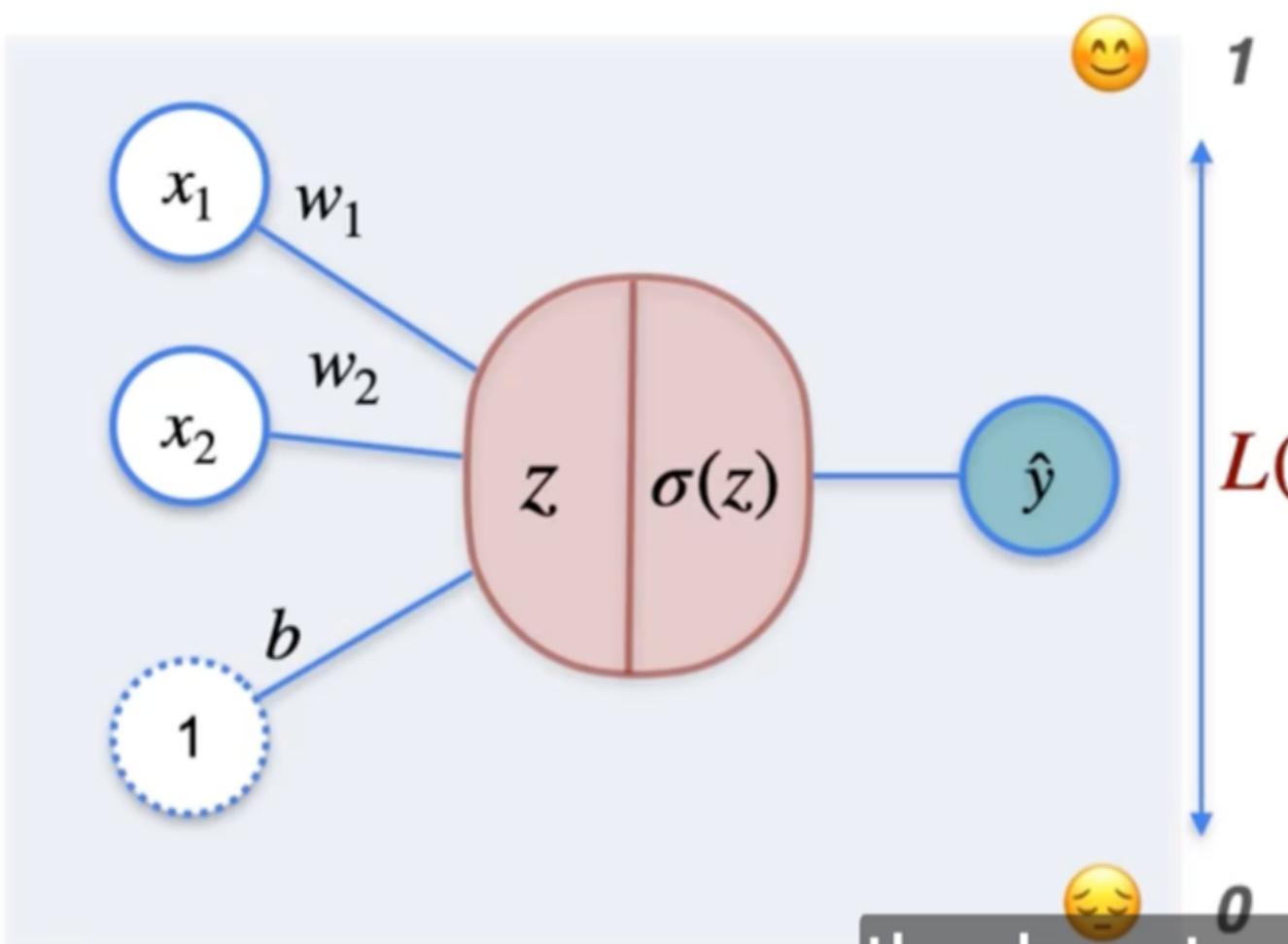
$$L(y, \hat{y}) = -y \ln(\hat{y}) - (1-y) \ln(1-\hat{y})$$

⇒ Main Goal  
find  $w_1, w_2$

if  $b$  to minimize  
loss

function.

# Classification With a Perceptron



**Prediction Function:**

$$\hat{y} = \sigma(w_1x_1 + w_2x_2 + b)$$

**Loss Function:**

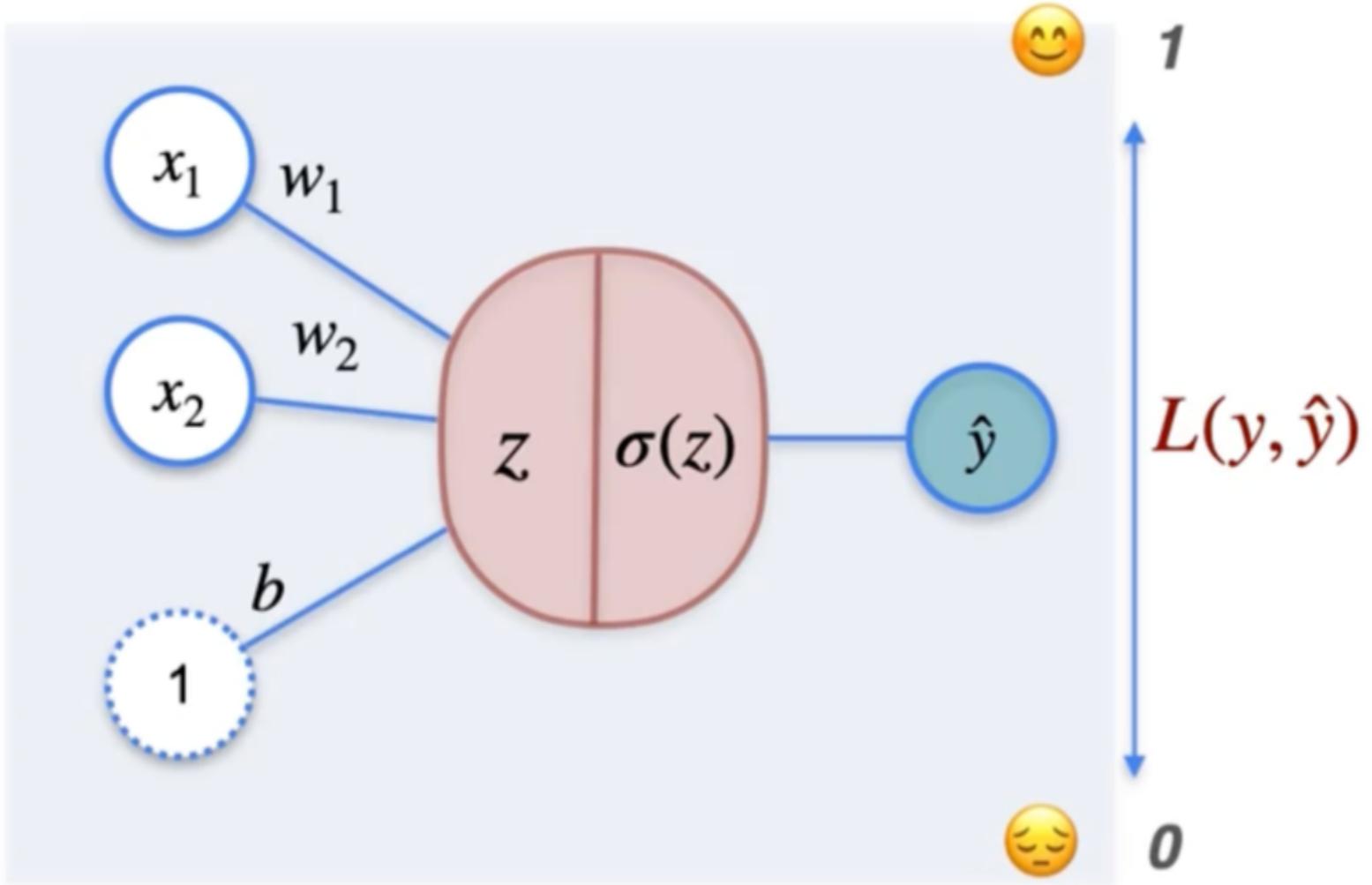
$$L(y, \hat{y}) \quad L(y, \hat{y}) = -y \ln(\hat{y}) - (1 - y) \ln(1 - \hat{y})$$

**Main Goal:**

Find  $w_1, w_2, b$  that give  $\hat{y}$  with the least error

the least amount of log loss  $L(y, \hat{y})$  error.

# Classification With a Perceptron



To find optimal values for:  
 $w_1, w_2, b$

We need gradient descent:

$$w_1 \rightarrow w_1 - \lambda \cdot \frac{\partial L}{\partial w_1}$$

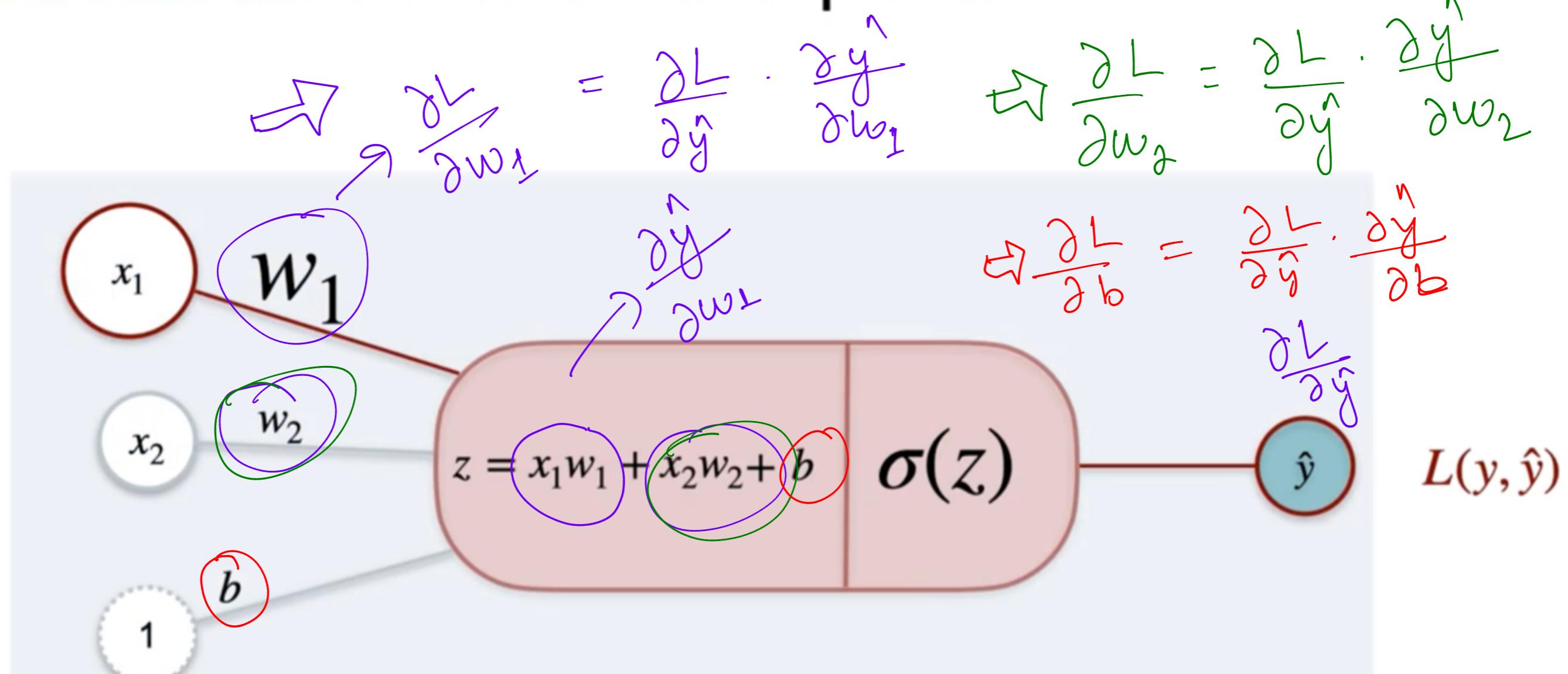
$$w_2 \rightarrow w_2 - \lambda \cdot \frac{\partial L}{\partial w_2}$$

$$b \rightarrow b - \lambda \cdot \frac{\partial L}{\partial b}$$

In order to be able to do this,

# Classification with a Perceptron

## Classification With a Perceptron



all those derivatives  
and be able to do

# Classification With a Perceptron

$$\frac{\partial L}{\partial b} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial b}$$

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial w_1}$$

$$\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial w_2}$$

$$\frac{\partial L}{\partial \hat{y}} =$$

$$\frac{\partial \hat{y}}{\partial b} =$$

$$\frac{\partial \hat{y}}{\partial w_1} = ?$$

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

$$\hat{y} = \sigma(w_1x_1 + w_2x_2 + b)$$

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

dy Hat is there many times. What are these ones?

# Classification With a Perceptron

$$\frac{\partial L}{\partial \hat{y}}$$

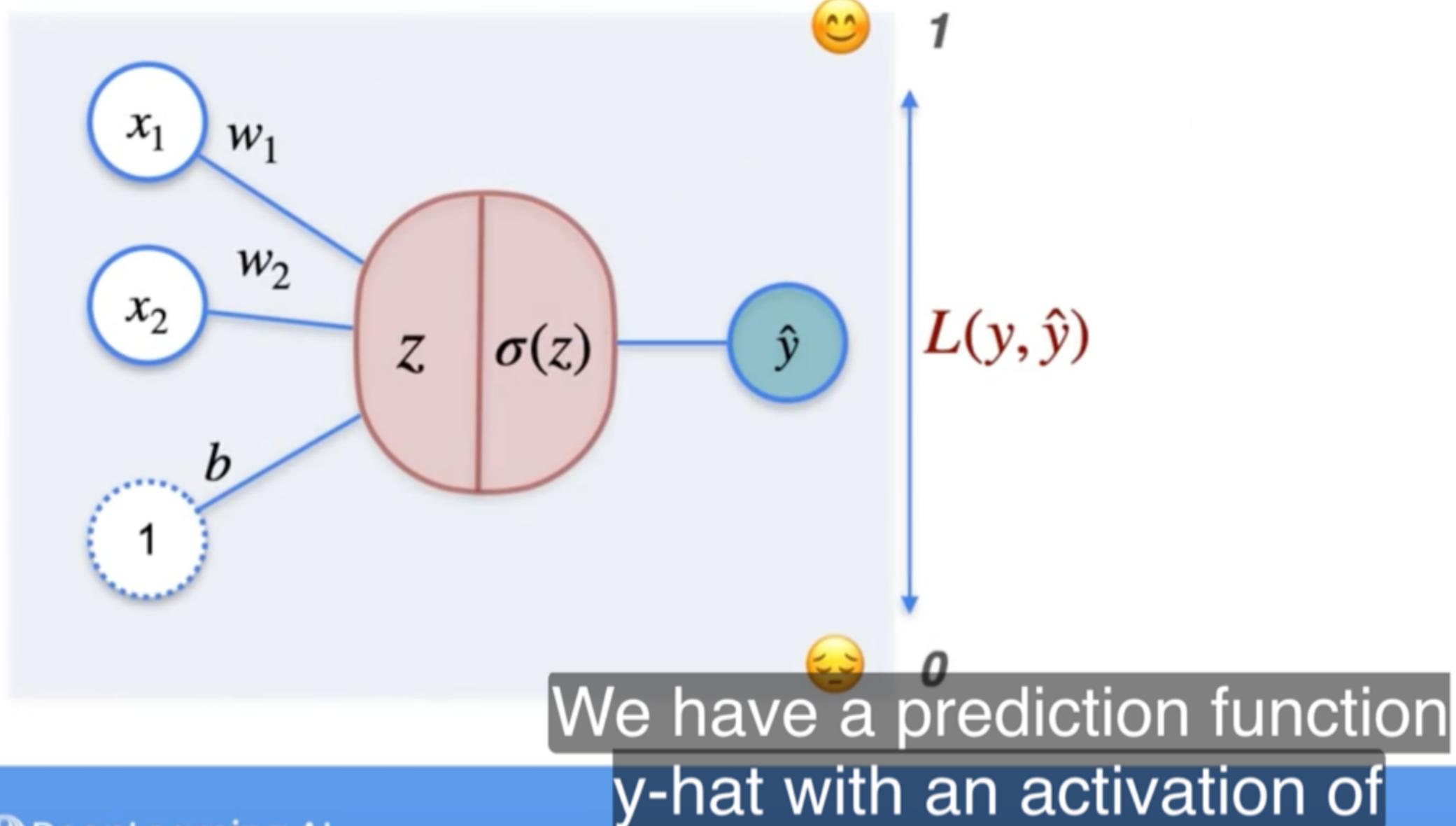
$$L(y, \hat{y}) = -y \ln(\hat{y}) - (1 - y) \ln(1 - \hat{y})$$

$$\frac{\partial L}{\partial \hat{y}} = \frac{-y}{\hat{y}} + \frac{1-y}{1-\hat{y}}$$

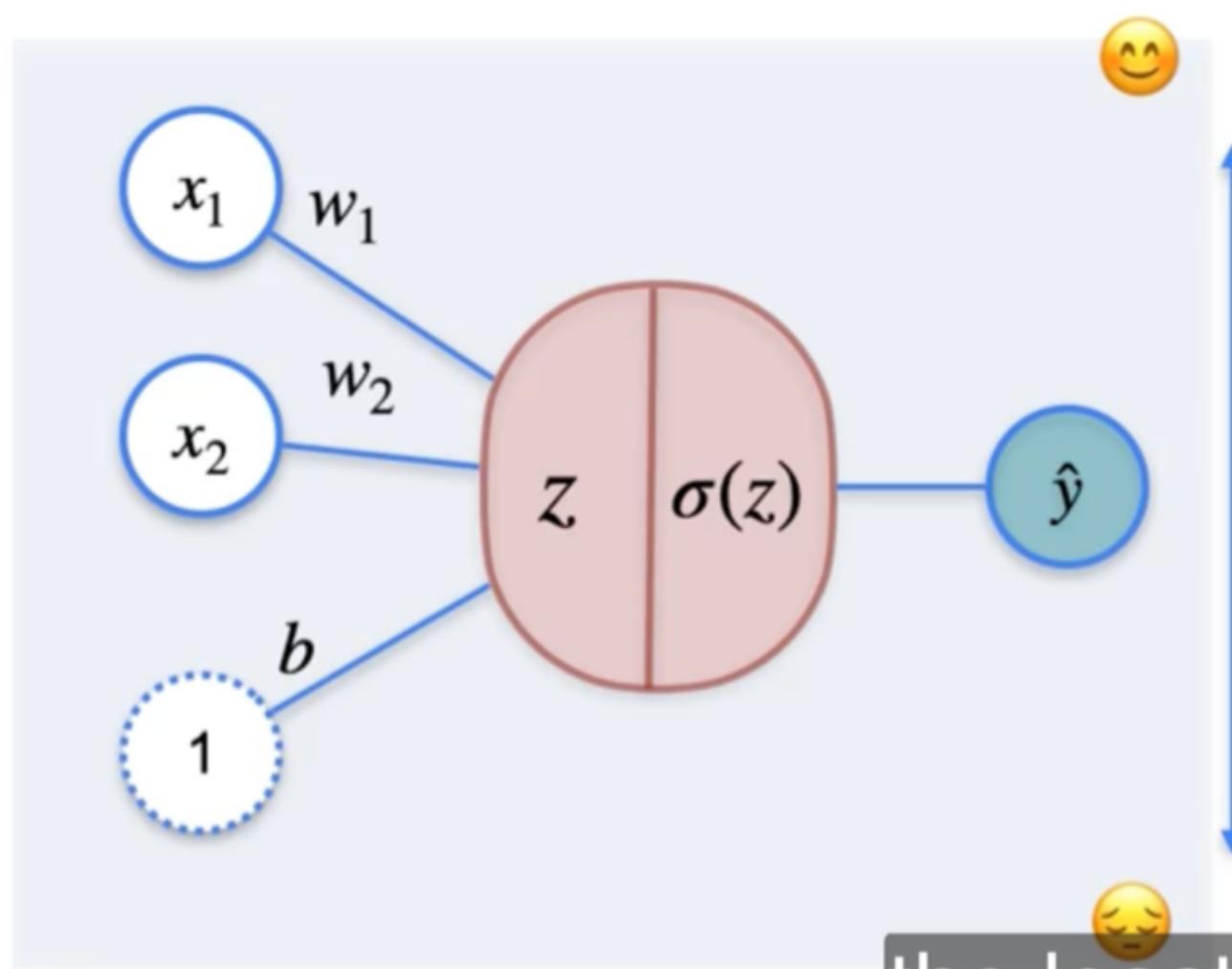
$$\begin{aligned} &= \frac{-y + y\hat{y}^{\wedge} + \hat{y} - \hat{y}\hat{y}^{\wedge}}{\hat{y}^{\wedge}(1 - \hat{y}^{\wedge})} \\ &= 1 - y(\hat{y} - \hat{y}^{\wedge}) \end{aligned}$$

Notice that the logarithm of

# Classification With a Perceptron



# Classification With a Perceptron



**Prediction Function:**

$$\hat{y} = \sigma(w_1x_1 + w_2x_2 + b)$$

**Loss Function:**

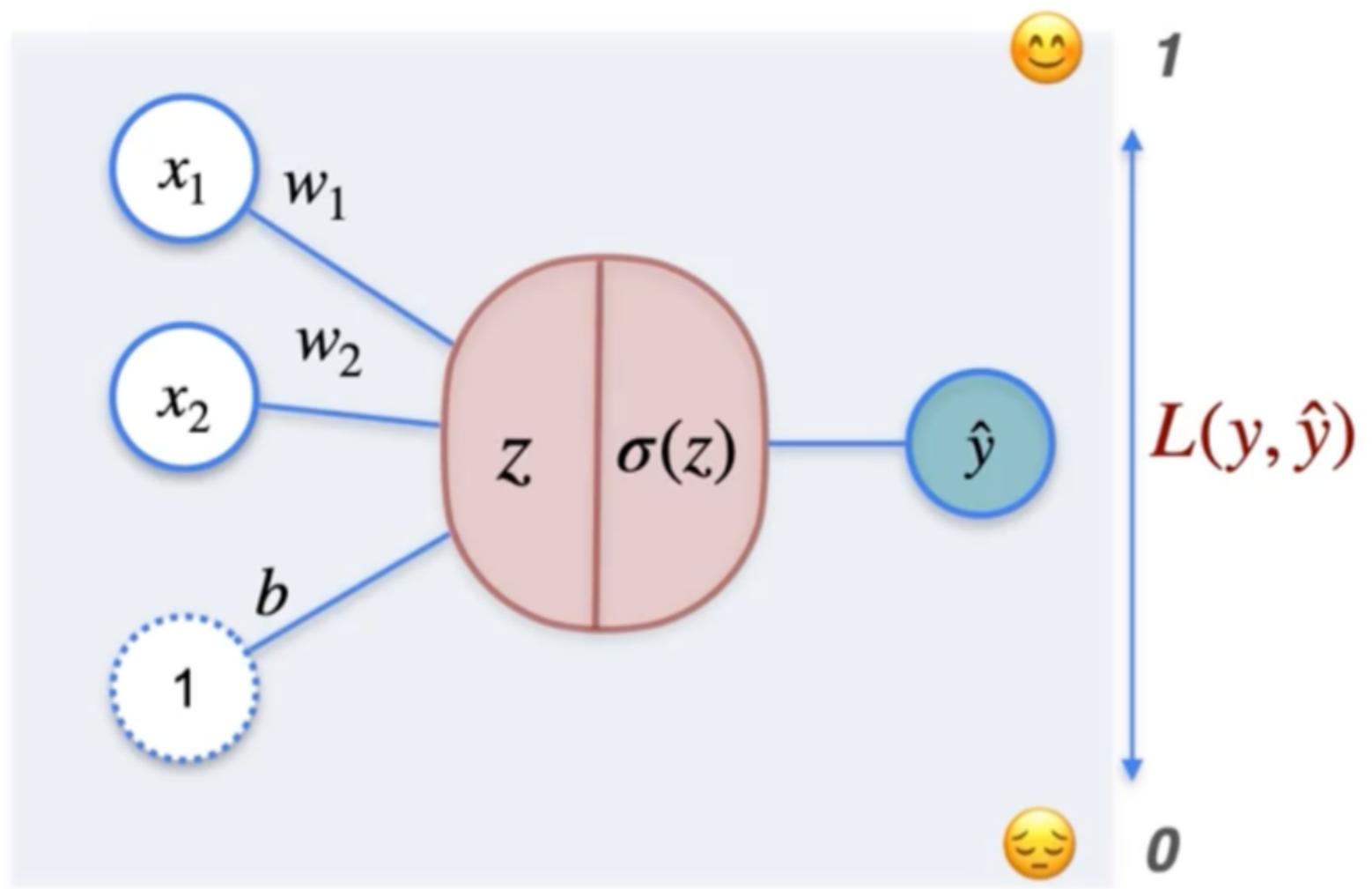
$$L(y, \hat{y}) \quad L(y, \hat{y}) = -y \ln(\hat{y}) - (1 - y) \ln(1 - \hat{y})$$

**Main Goal:**

Find  $w_1, w_2, b$  that give  $\hat{y}$  with the least error

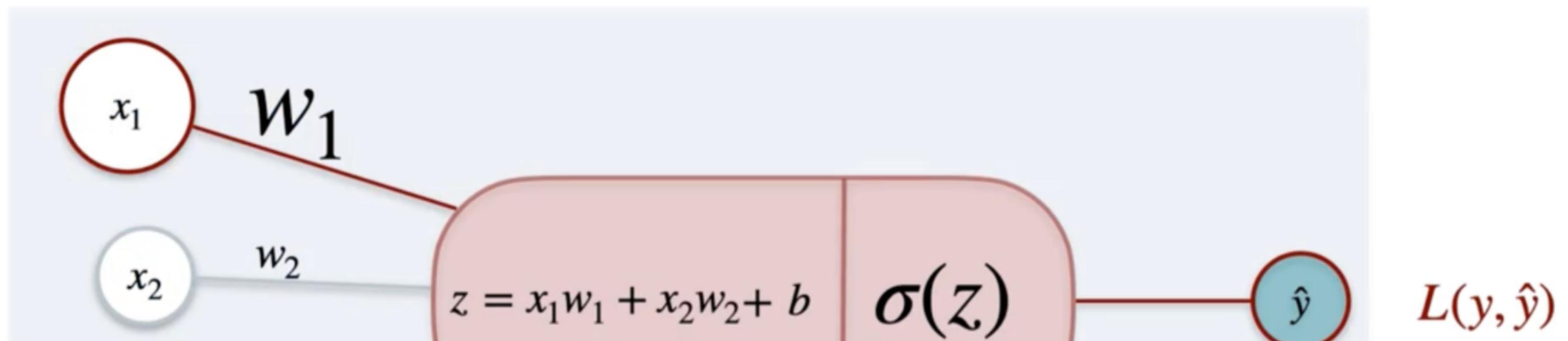
the least amount of log loss  $L(y, \hat{y})$  error.

# Classification With a Perceptron



In order to be able to do this,

# Classification With a Perceptron



all those derivatives  
and be able to do

# Classification With a Perceptron

$$\frac{\partial L}{\partial b} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial b}$$

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial w_1}$$

$$\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial w_2}$$



dy Hat is there many times. What are these ones?

# Classification With a Perceptron

$$\hat{y} = \sigma(w_1x_1 + w_2x_2 + b)$$

$$\frac{\partial \hat{y}}{\partial w_1} = \hat{y}(1 - \hat{y})x_1$$

$$\frac{\partial \hat{y}}{\partial w_2} = \hat{y}(1 - \hat{y})x_2$$

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

of something is sigmoid  
times one minus sigmoid.

# Classification With a Perceptron

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

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

$$\frac{\partial \hat{y}}{\partial w_1} = \hat{y}(1 - \hat{y})x_1$$

$$\frac{\partial \hat{y}}{\partial w_2} = \hat{y}(1 - \hat{y})x_2$$

$$\frac{\partial L}{\partial b} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial b}$$

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial w_1}$$

$$\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial w_2}$$

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

$$= \frac{-(y - \hat{y})}{\hat{y}(1 - \hat{y})} \hat{y}(1 - \hat{y})x_1$$

$$= \frac{-(y - \hat{y})}{\hat{y}(1 - \hat{y})} \hat{y}(1 - \hat{y})x_2$$

Now we need to cancel  
out a bunch of things.

# Classification With a Perceptron

$$\frac{\partial L}{\partial b} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial b} = \frac{-(y - \hat{y})}{\hat{y}(1 - \hat{y})}$$

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial w_1} = \frac{-(y - \hat{y})}{\hat{y}(1 - \hat{y})} \hat{y}(1 - \hat{y})x_1$$

$$\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial w_2} = \frac{-(y - \hat{y})}{\hat{y}(1 - \hat{y})} \hat{y}(1 - \hat{y})x_2$$

cancel out on the denominator  
and the numerator.

# Classification With a Perceptron

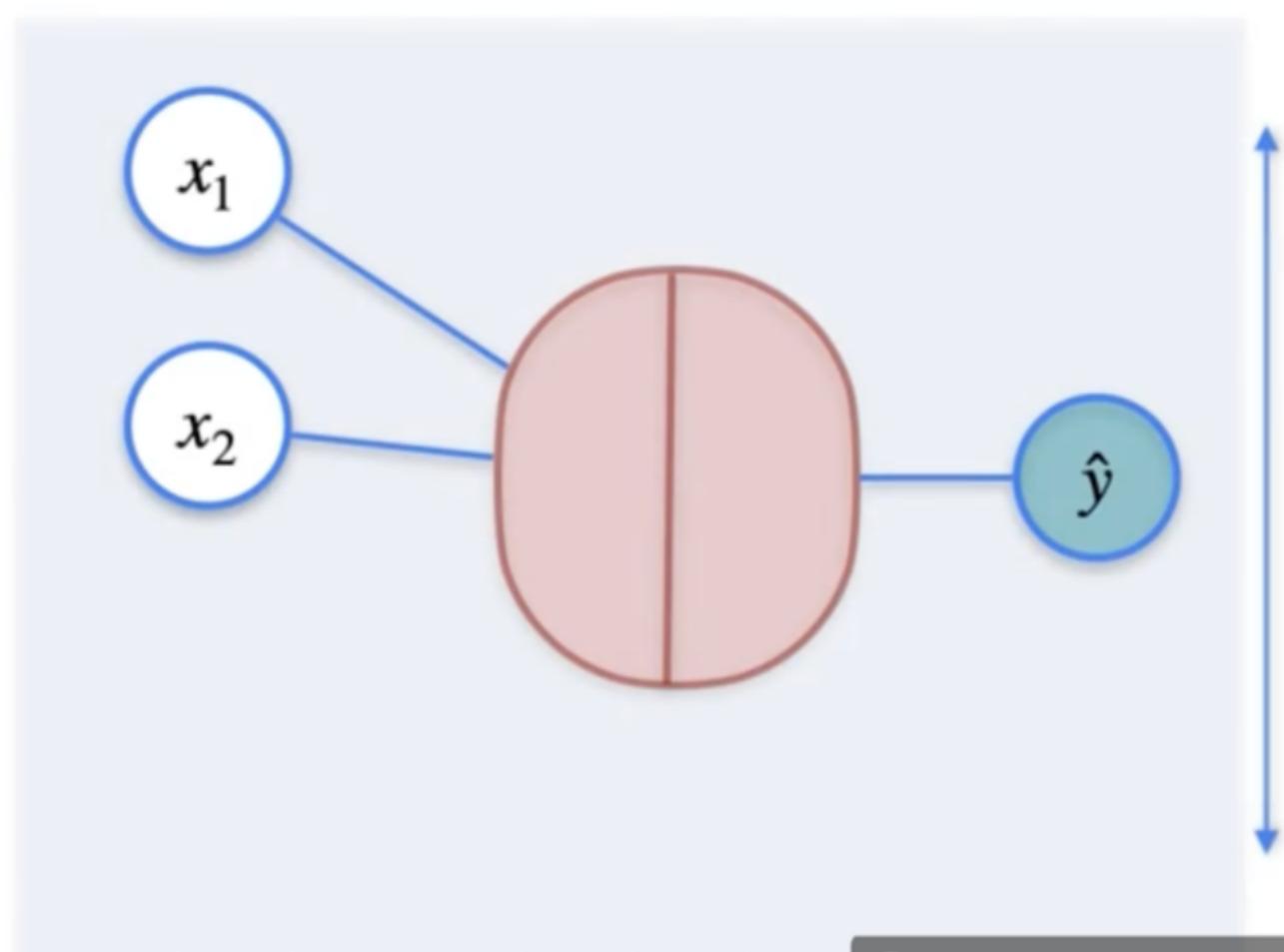
$$\frac{\partial L}{\partial b} = -(y - \hat{y})$$

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

$$\frac{\partial L}{\partial w_2} = -(y - \hat{y})x_2$$

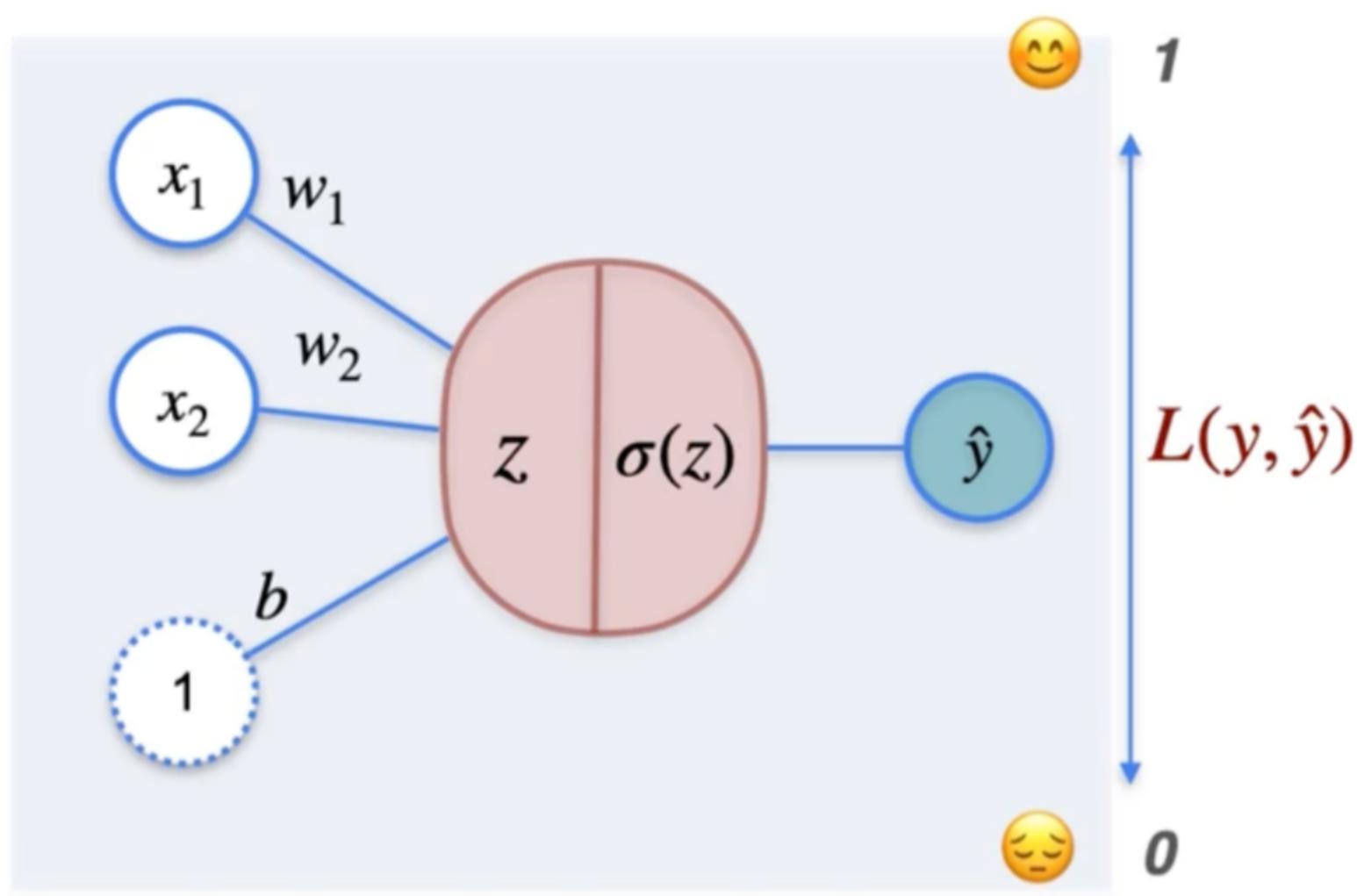
Therefore, you need a large

# Classification With a Perceptron



Our main goal was to find  
the optimal values for  $W_1$ ,

# Classification With a Perceptron



To find optimal values for:  
 $w_1, w_2, b$

*You need gradient descent*

$$w_1 \rightarrow w_1 - \alpha(-x_1(y - \hat{y}))$$

$$w_2 \rightarrow w_2 - \alpha(-x_2(y - \hat{y}))$$

$$b \rightarrow b - \alpha(-(y - \hat{y}))$$

You start with some  $W1$ ,