

Linear

# Model

Son Nguyen

• Given the data



ullet How are y and x related?

• Given the data



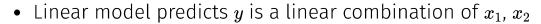
ullet Linear model predicts y is a linear combination of  $x_1$ ,  $x_2$ 

$$\hat{y} = w_0 + w_1 x_1 + w_2 x_2$$

Grant combration of the predictors

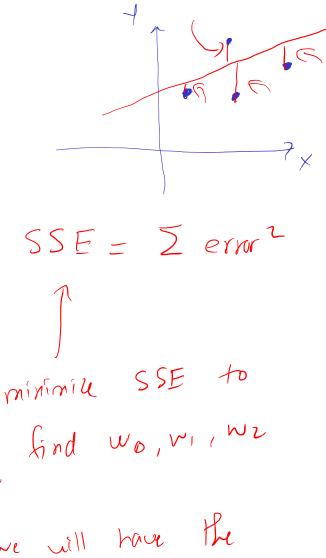
• Given the data





$$\hat{y} = w_0 + w_1 x_1 + w_2 x_2$$

- The goal of linear model is to solve for  $w_0$ ,  $w_1$  and  $w_2$
- To **train** a linear model is to find  $w_0$ ,  $w_1$  and  $w_2$



Solve this we will have the last squared solution.

$$f(w) = w^{2} + 2w + 5 - 50 = 6 \text{ for } \sqrt{f'(w)} = 0 = 0 \text{ for } \sqrt{f'(w)} = 0 \text{ for } \sqrt{f'(w)} = 0 = 0$$

Solve for 
$$\frac{\partial f}{\partial w_0} = 0$$
,  $\frac{\partial f}{\partial w_0} = 0$ 

Given the data

$x_1$	$x_2$	y
1	0	<b>-</b> 2
2	1	0
3	<del>-</del> 2	-1
4	3	1

what if we went to minimile I lerror

= 7, 17-91

• Linear model predicts y is a linear combination of  $x_1, x_2$ 

$$\hat{y} = w_0 + w_1 x_1 + w_2 x_2$$

- The goal of linear model is to solve for  $w_0$ ,  $w_1$  and  $w_2$
- To **train** a linear model is to find  $w_0$ ,  $w_1$  and  $w_2$

$$= \sum | t - w_0 - w_1 x_1 - w_2 x_1$$

Cannot use

Re "derivative" approach =

ble 1-1 is not differiable.

Solve this, we would a different solution.

• Given the data



• Linear model predicts y is a linear combination of  $x_1$ ,  $x_2$ 

$$\hat{y} = w_0 + w_1 x_1 + w_2 x_2$$

- ullet The goal of linear model is to solve for  $w_0,\,w_1$  and  $w_2$
- To **train** a linear model is to find  $w_0$ ,  $w_1$  and  $w_2$

• **Step 1**: Define the loss function  $l(y,\hat{y})$ 

,

- Step 1: Define the loss function  $l(y,\hat{y})$
- **Step 2**: Find w that minimizes the total loss function.

• Least Squared Method uses the **square loss** 

$$l(\hat{y},y) = (\hat{y}-y)^2$$

• Least Squared Method uses the **square loss** 

$$l(\hat{y},y) = (\hat{y}-y)^2$$

• We want to find  $w_0$ ,  $w_1$  and  $w_2$  that minimizes a loss function.

$x_1$	$x_2$	$oldsymbol{y}$	$\hat{y}=w_0+w_1x_1+w_2x_2$	$(\hat{y}-y)^2$
1	0	-2	$w_0 + w_1 \cdot 1 + w_2 \cdot 0$	$(w_0 + w_1 \cdot 1 + w_2 \cdot 0 + 2)^2$
2	1	0	$w_0+w_1\cdot 2+w_2\cdot 1$	$(w_0 + w_1 \cdot 2 + w_2 \cdot 1 - 0)^2$
3	-2	-1	$w_0+w_1\cdot 3+w_2\cdot -2$	$(w_0 + w_1 \cdot 3 + w_2 \cdot -2 + 1)^2$
4	3	1	$w_0+w_1\cdot 4+w_2\cdot 3$	$(w_0 + w_1 \cdot 4 + w_2 \cdot 3 - 1)^2$

• Least Squared Method uses the **square loss** 

$$l(\hat{y},y) = (\hat{y} - y)^2$$

• We want to find  $w_0$ ,  $w_1$  and  $w_2$  that minimizes a **loss function**.

$x_1$	$x_2$	$oldsymbol{y}$	$\hat{y}=w_0+w_1x_1+w_2x_2$	$(\hat{y}-y)^2$
1	0	-2	$w_0 + w_1 \cdot 1 + w_2 \cdot 0$	$(w_0 + w_1 \cdot 1 + w_2 \cdot 0 + 2)^2$
2	1	0	$w_0 + w_1 \cdot 2 + w_2 \cdot 1$	$(w_0 + w_1 \cdot 2 + w_2 \cdot 1 - 0)^2$
3	-2	-1	$w_0+w_1\cdot 3+w_2\cdot -2$	$(w_0 + w_1 \cdot 3 + w_2 \cdot -2 + 1)^2$
4	3	1	$w_0+w_1\cdot 4+w_2\cdot 3$	$(w_0 + w_1 \cdot 4 + w_2 \cdot 3 - 1)^2$

• The total loss function:

$$L = L(w_0, w_1, w_2) = (w_0 + w_1 + 2)^2 + (w_0 + 2w_1 + w_2)^2 \ + (w_0 + 3w_1 - 2w_2 + 1)^2 + (w_0 + 4w_1 + 3w_2 - 1)^2$$

• Least Squared Method uses the **square loss** 

$$l(\hat{y},y) = (\hat{y} - y)^2$$

- We want to find  $w_0$ ,  $w_1$  and  $w_2$  that minimizes a **loss function**.
- The total loss function:

$$L = L(w_0, w_1, w_2) = (w_0 + w_1 + 2)^2 + (w_0 + 2w_1 + w_2)^2 \ + (w_0 + 3w_1 - 2w_2 + 1)^2 + (w_0 + 4w_1 + 3w_2 - 1)^2$$

• Solve for the partial derivatives equaling 0 to find  $w_0$ ,  $w_1$  and  $w_2$ .

#### How about other loss functions?

• Absolute loss:

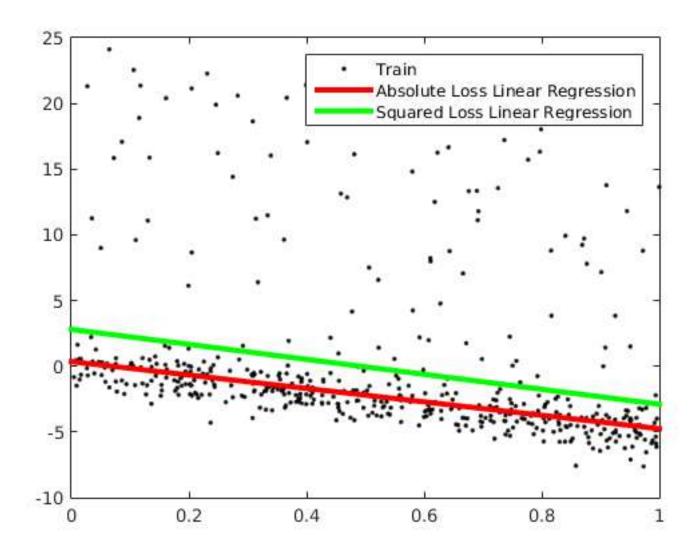
$$L(\hat{y},y) = |\hat{y} - y|$$

• The total loss function:

$$L = L(w_0, w_1, w_2) = |w_0 + w_1 + 2| + |w_0 + 2w_1 + w_2| \ + |w_0 + 3w_1 - 2w_2 + 1| + |w_0 + 4w_1 + 3w_2 - 1|$$

- Use Linear Programming to find  $w_0$ ,  $w_1$  and  $w_2$  that minimizes the total loss.
- Least absolute deviations regression

# Linear Models



#### How about other loss functions?

Ordinary least squares regression	Least absolute deviations regression
Not very robust	Robust
Stable solution	Unstable solution
Always one solution	Possibly multiple solutions

### A general framework

- **Problem**: Given the data of  $x_1, x_2, \ldots, x_d, y$ , establish the *best* relation between y and  $x = [x_1, x_2, \ldots, x_d]$ .
- A solution framework:
  - $\circ$  Step 1: Assume the model function  $\hat{y} = f(x, w)$ , where w is a parameter vector.
  - $\circ$  Step 2: Define the loss function  $l(y,\hat{y})$

 $\circ$  Step 3: Find w that minimizes the loss function using gradient descent

approximation method to solve for  $\frac{\partial f}{\partial w_i} = 0$ 

#### LASSO

• Consider a linear model

$$y=100x_1+0.01x_2+50x_3-0.002x_4$$

- ullet  $x_2$  and  $x_4$  are not important because the coefficients are too small.
- ullet We want to get rid of  $x_2$  and  $x_4$

### LASSO - Principle

- LASSO forces the sum of the absolute value of the coefficients to be less than a fixed value.
- which forces certain coefficients (slopes) to be set to zero
- effectively making the model simpler

$$f(w) = (w - 3)^{2} + 4$$

$$w = 3 \quad \text{minimize} \quad f(w)$$

$$f(w) = (w-1)^{2} + 4$$

$$constrain: |w| |s|$$

### Linear Model vs. LASSO - Principle

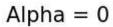
• Linear Model minimizes

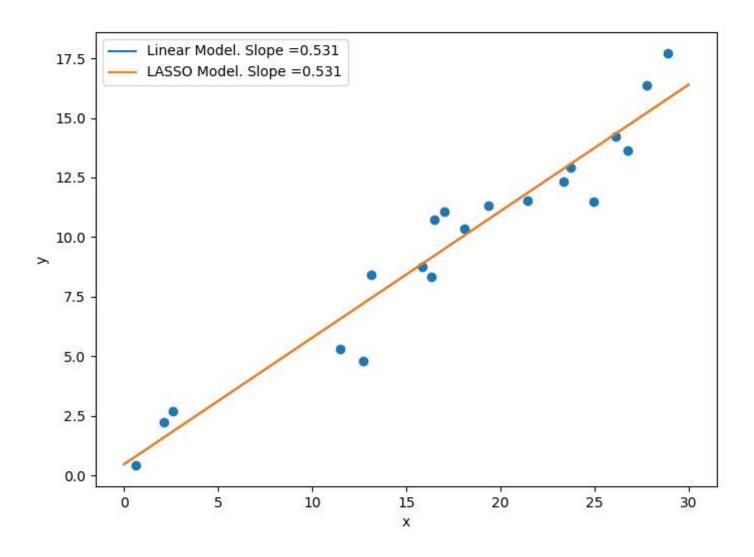
$$L(w_0,w_1,w_2)$$

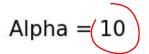
• LASSO minimizes

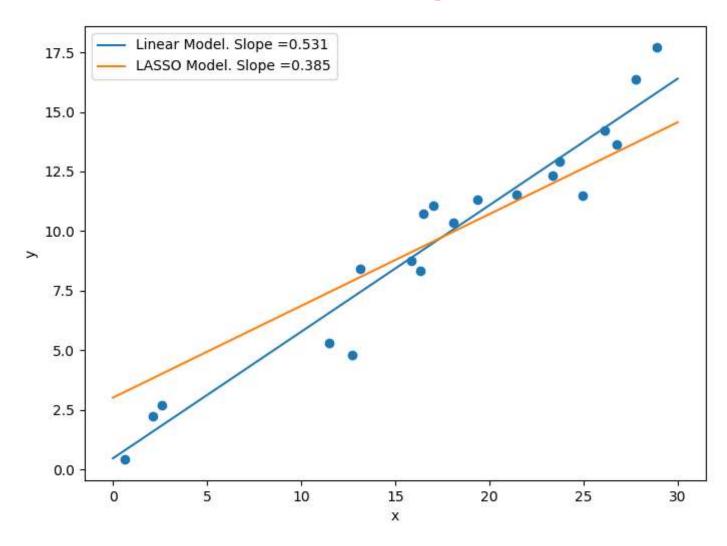
$$L(w_0,w_1,w_2)+lphaigg(|w_1|+|w_2|igg)$$

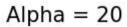
- The greater  $\alpha$ , the easier  $w_1$  and  $w_2$  will be zeros.
- When  $\alpha=0$ , LASSO is the linear model.

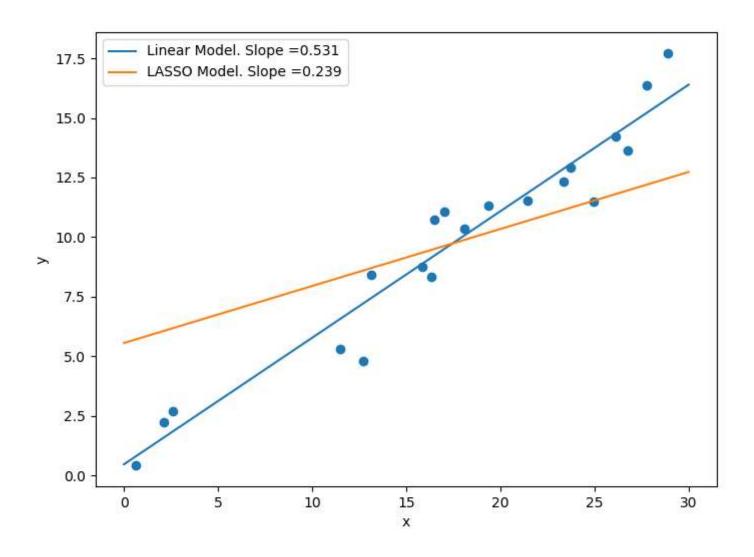


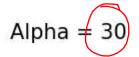


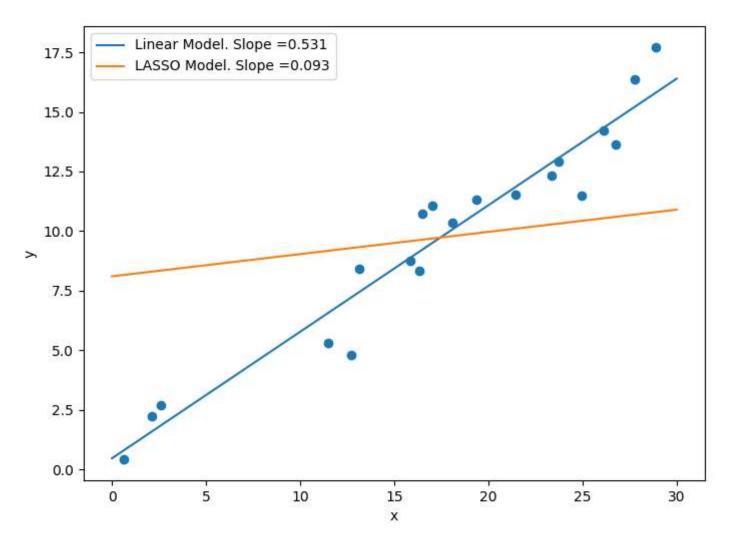


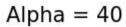


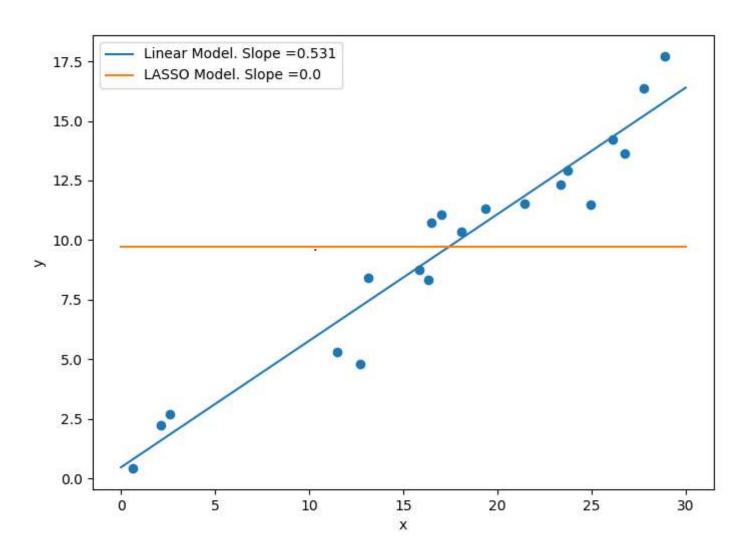


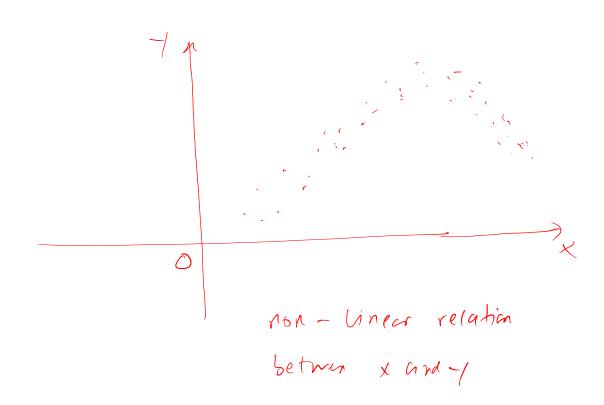








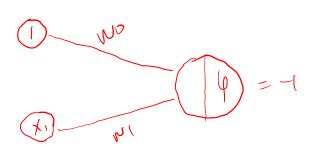




1) Graphical pre sentation of Unear model.

+=  $w_0+w_1X_1$ 

(algebraic presentation)



$$\varphi(t) = t$$

$$\frac{1}{1+e^{-(w_0+w_1y_1)}}$$

$$\varphi(t) = \frac{1}{1 + e^{-t}}$$

((t): actuation function

Example:

$$\begin{aligned}
Y &= SIn \left(N_0 + W_1 X_1\right) + W_2 \cos \left(W_3 + W_4 X_1\right) \\
+ W_5 &= \frac{W_6 X_1}{W_2} + W_4 \cos \left(W_3 + W_4 X_1\right) \\
&= \frac{W_6 X_1}{W_2 + W_3} + W_4 \cos \left(W_3 + W_4 X_1\right) \\
&= \frac{W_6 X_1}{W_2 + W_4 \times W_1} + W_4 \cos \left(W_3 + W_4 X_1\right) \\
&= \frac{W_6 X_1}{W_2 + W_4 \times W_1} + W_4 \cos \left(W_3 + W_4 X_1\right) \\
&= \frac{W_6 X_1}{W_2 + W_4 \times W_4} + W_4 \cos \left(W_3 + W_4 X_1\right) \\
&= \frac{W_6 X_1}{W_2 + W_4 \times W_4} + W_4 \cos \left(W_3 + W_4 X_1\right) \\
&= \frac{W_6 X_1}{W_2 + W_4 \times W_4} + W_4 \cos \left(W_3 + W_4 X_1\right) \\
&= \frac{W_6 X_1}{W_2 + W_4 \times W_4} + W_4 \cos \left(W_3 + W_4 X_1\right) \\
&= \frac{W_6 X_1}{W_4 + W_4 \times W_4} + W_4 \cos \left(W_3 + W_4 X_1\right) \\
&= \frac{W_6 X_1}{W_4 + W_4 \times W_4} + W_4 \cos \left(W_4 + W_4 X_1\right) \\
&= \frac{W_6 X_1}{W_4 + W_4 \times W_4} + W_4 \cos \left(W_4 + W_4 X_1\right) \\
&= \frac{W_6 X_1}{W_4 + W_4 \times W_4} + W_4 \cos \left(W_4 + W_4 X_1\right) \\
&= \frac{W_6 X_1}{W_4 + W_4 \times W_4} + W_4 \cos \left(W_4 + W_4 X_1\right) \\
&= \frac{W_6 X_1}{W_4 + W_4 \times W_4} + W_4 \cos \left(W_4 + W_4 X_1\right) \\
&= \frac{W_6 X_1}{W_4 + W_4 \times W_4} + W_4 \cos \left(W_4 + W_4 X_1\right) \\
&= \frac{W_6 X_1}{W_4 + W_4 \times W_4} + W_4 \cos \left(W_4 + W_4 X_1\right) \\
&= \frac{W_6 X_1}{W_4 + W_4 \times W_4} + W_4 \cos \left(W_4 + W_4 X_1\right) \\
&= \frac{W_6 X_1}{W_4 + W_4 \times W_4} + W_4 \cos \left(W_4 + W_4 X_1\right) \\
&= \frac{W_6 X_1}{W_4 + W_4 \times W_4} + W_4 \cos \left(W_4 + W_4 X_1\right) \\
&= \frac{W_6 X_1}{W_4 + W_4 \times W_4} + W_4 \cos \left(W_4 + W_4 X_1\right) \\
&= \frac{W_6 X_1}{W_4 + W_4 \times W_4} + W_4 \cos \left(W_4 + W_4 X_1\right) \\
&= \frac{W_6 X_1}{W_4 + W_4 \times W_4} + W_4 \cos \left(W_4 + W_4 X_1\right) \\
&= \frac{W_6 X_1}{W_4 + W_4 \times W_4} + W_4 \cos \left(W_4 + W_4 X_1\right) \\
&= \frac{W_6 X_1}{W_4 + W_4 \times W_4} + W_4 \cos \left(W_4 + W_4 X_1\right) \\
&= \frac{W_6 X_1}{W_4 + W_4 \times W_4} + W_4 \cos \left(W_4 + W_4 X_1\right) \\
&= \frac{W_6 X_1}{W_4 + W_4 \times W_4} + W_4 \cos \left(W_4 + W_4 X_1\right) \\
&= \frac{W_6 X_1}{W_4 + W_4 \times W_4} + W_4 \cos \left(W_4 + W_4 X_1\right) \\
&= \frac{W_6 X_1}{W_4 + W_4 \times W_4} + W_4 \cos \left(W_4 + W_4 X_1\right) \\
&= \frac{W_6 X_1}{W_4 + W_4 \times W_4} + W_4 \cos \left(W_4 + W_4 X_1\right) \\
&= \frac{W_6 X_1}{W_4 + W_4 \times W_4} + W_4 \cos \left(W_4 + W_4 X_1\right) \\
&= \frac{W_6 X_1}{W_4 + W_4 \times W_4} + W_4 \cos \left(W_4 + W_4 X_1\right) \\
&= \frac{W_6 X_1}{W_4 + W_4 \times W_4} + W_4 \cos \left(W_4 + W_4 X_1\right) \\
&= \frac{W_6 X_1}{W_4 + W_4 \times W_4} + W_4 \times W_4 \times W_4 \times W_4 \\
&= \frac{W_6 X_1}{W_4 + W_4 \times W_4} + W_4 \times W_4 \times W_4 \\
&=$$

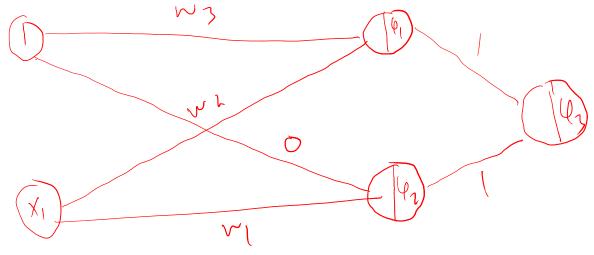
$$7 = w_1 \cdot x_1$$

$$\frac{1}{2} = (w_2 x_1 + w_3) \cdot \ln(w_1 x_1)$$

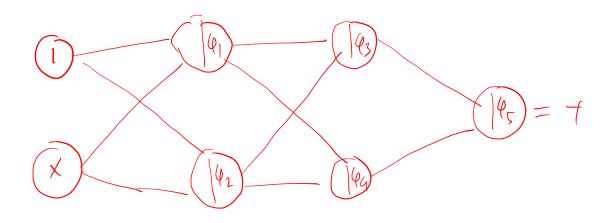
$$t = x_1 x_2$$

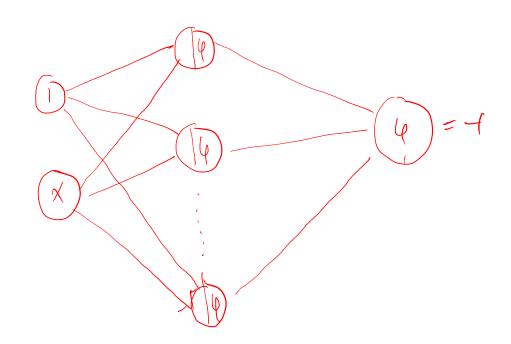
$$\ln \gamma = \ln x_1 + \ln x_2$$

$$ln(ln-l) = ln(w_2 x_1 + w_3) + ln(ln(v_1x_1))$$
  
 $4_1(t) = ln(t)$ 



$$(q_{3}(t) = 1n((n(t)))$$
 $(q_{3}(t) = t)$ 





All G, G2... 95

are liner activ.

Function

then the relation

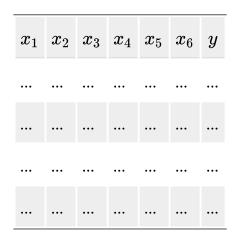
Setweren x and -1

Still linear.

This network uses
only 1 t-rpe of
actional Curction
(p(t), ron-linear.

#### LASSO for Variables Selection

• Data



• Assume that the truth relation between the input  $x_1, x_2, x_3, x_4, x_5, x_6$  and the output y is

$$y = 4x_2 + 3x_4 + 7x_6$$

- ullet We see that only  $x_2$ ,  $x_4$  and  $x_6$  impact y
- ullet LASSO can help to identify variables that have effect on y

#### LASSO for Variables Selection

• The result when training the linear model and the LASSO

	$w_1$	$w_2$	$w_3$	$w_4$	$w_5$	$w_6$
Truth	0	4	0	3	0	7
Linear Model	-0.244061	3.54013	0.221939	2.6042	0.0982158	6.83617
LASSO	-0	2.65623	0	1.84839	0	5.80624

- In Linear Model,  $x_1$ ,  $x_3$  and  $x_5$  have effect on y (which is WRONG!)
- In LASSO,  $x_1$ ,  $x_3$  and  $x_5$  have no effect on y (CORRECT!)
- LASSO can also be applied before another model.

# Logistic Regression

- classification,

$x_1$	$x_2$	y
1	0	1
2	1	0
3	<b>-</b> 2	0
4	3	1

• How are y and x related?

### Logistic Regression



• Logistic Regression models  $P(y=1|x)=\hat{y}$  as:

$$\hat{y} = rac{1}{1 + e^{-(w_0 + w_1 \cdot x_1 + w_2 \cdot x_2)}}$$

• OR,

$$\log\left(rac{\hat{y}}{1-\hat{y}}
ight) = w_0 + w_1 \cdot x_1 + w_2 \cdot x_2$$

where  $\hat{y}$  is the predicted value of the probability of y=1 given  $x_1$  and  $x_2$ .

### Logistic Regression



• Logistic Regression models  $P(y=1|x)=\hat{y}$  as:

$$\hat{y} = rac{1}{1 + e^{-(w_0 + w_1 \cdot x_1 + w_2 \cdot x_2)}}$$

• OR,

$$\log\left(rac{\hat{y}}{1-\hat{y}}
ight)=w_0+w_1\cdot x_1+w_2\cdot x_2$$

where  $\hat{y}$  is the predicted value of the probability of y=1 given  $x_1$  and  $x_2$ .

$$\log\left(rac{\hat{y}}{1-\hat{y}}
ight) = w_0 + w_1 \cdot x_1 + w_2 \cdot x_2$$

- $\left(\frac{\hat{y}}{1-\hat{y}}\right)$  is also called odd-ratio.
- Logistic Regression assumes that the log of the odd ratio is linear.

#### How to find $w_0, w_1, w_2$ ?

- **Step 1**: Define the loss function  $l(\hat{y}, y)$
- Step 2: Find  $\boldsymbol{w}$  that minimizes the total loss function

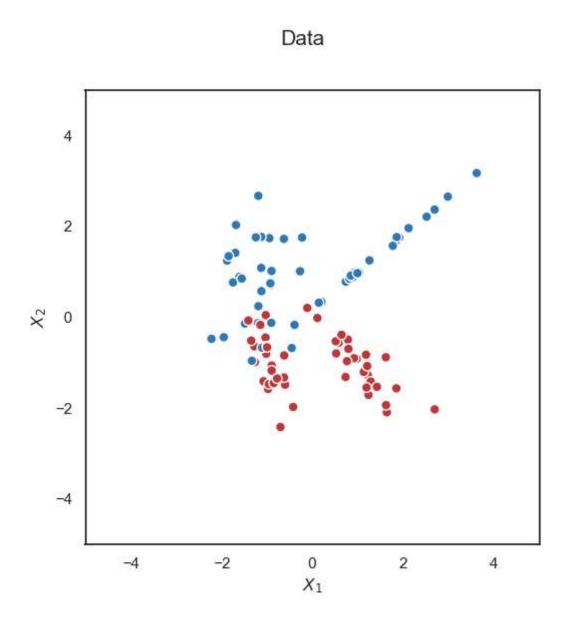
• Define the loss function: We use the log-loss or cross-entropy loss function

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

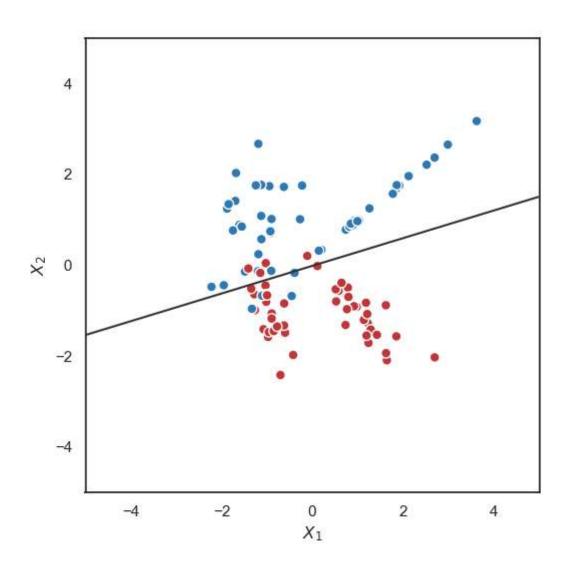
• Total Loss:

$$egin{split} L(w_0,w_1,w_2) &= -\log\left(rac{1}{1+e^{-w_0-w_1}}
ight) \ &-\log\left(1-rac{1}{1+e^{-w_0-2w_1-w_2}}
ight) \ &-\log\left(1-rac{1}{1+e^{-w_0-3w_1+w_2}}
ight) \ &-\log\left(rac{1}{1+e^{-w_0-4w_1-3w_2}}
ight) \end{split}$$

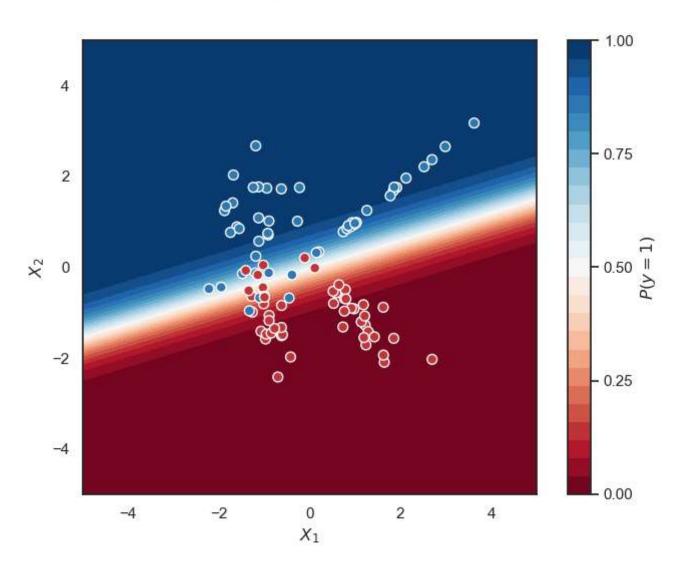
• We need to find  $w_0, w_1, w_2$  that minimizes the total loss











• The idea is the same as for linear model

.

