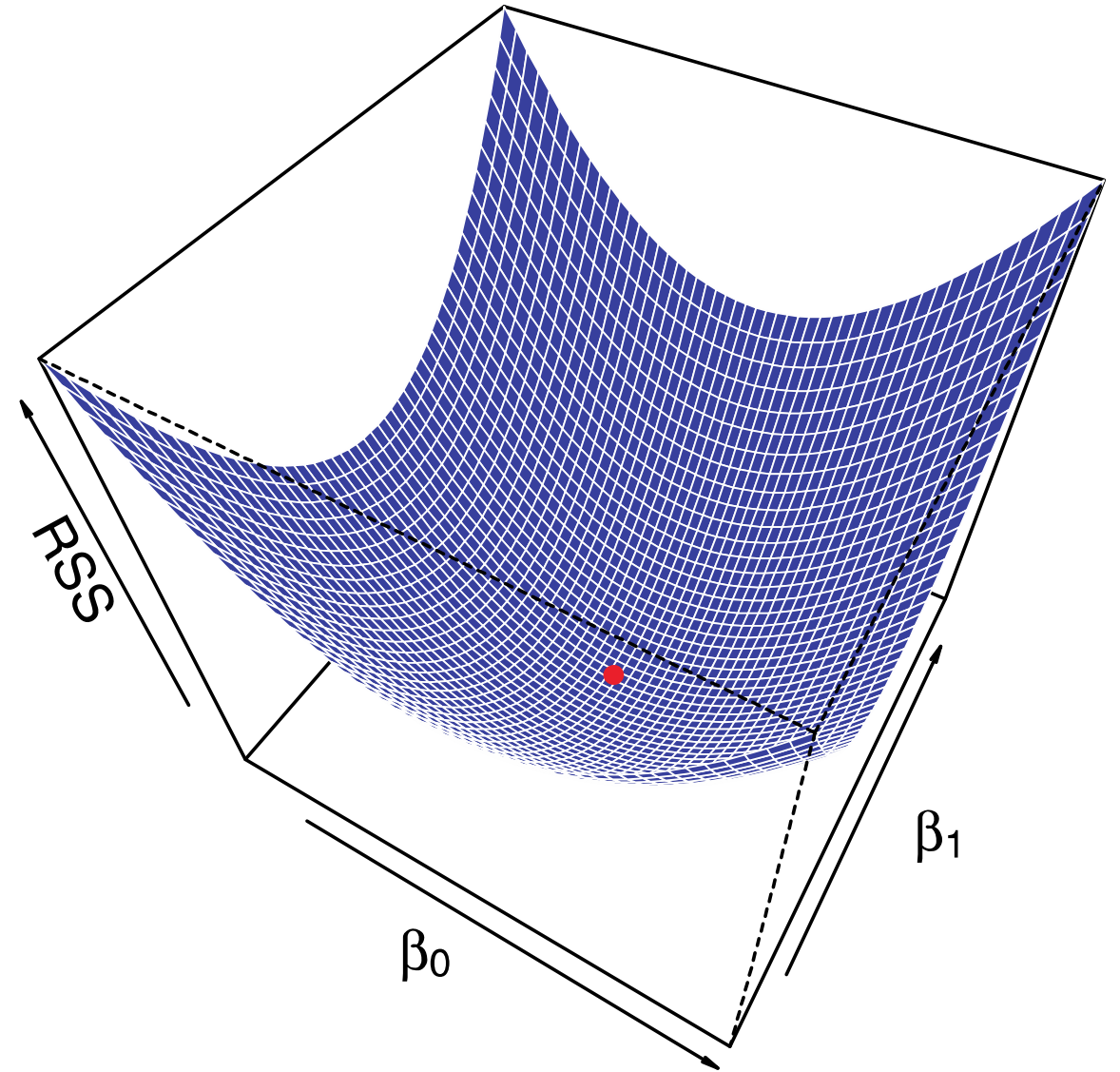
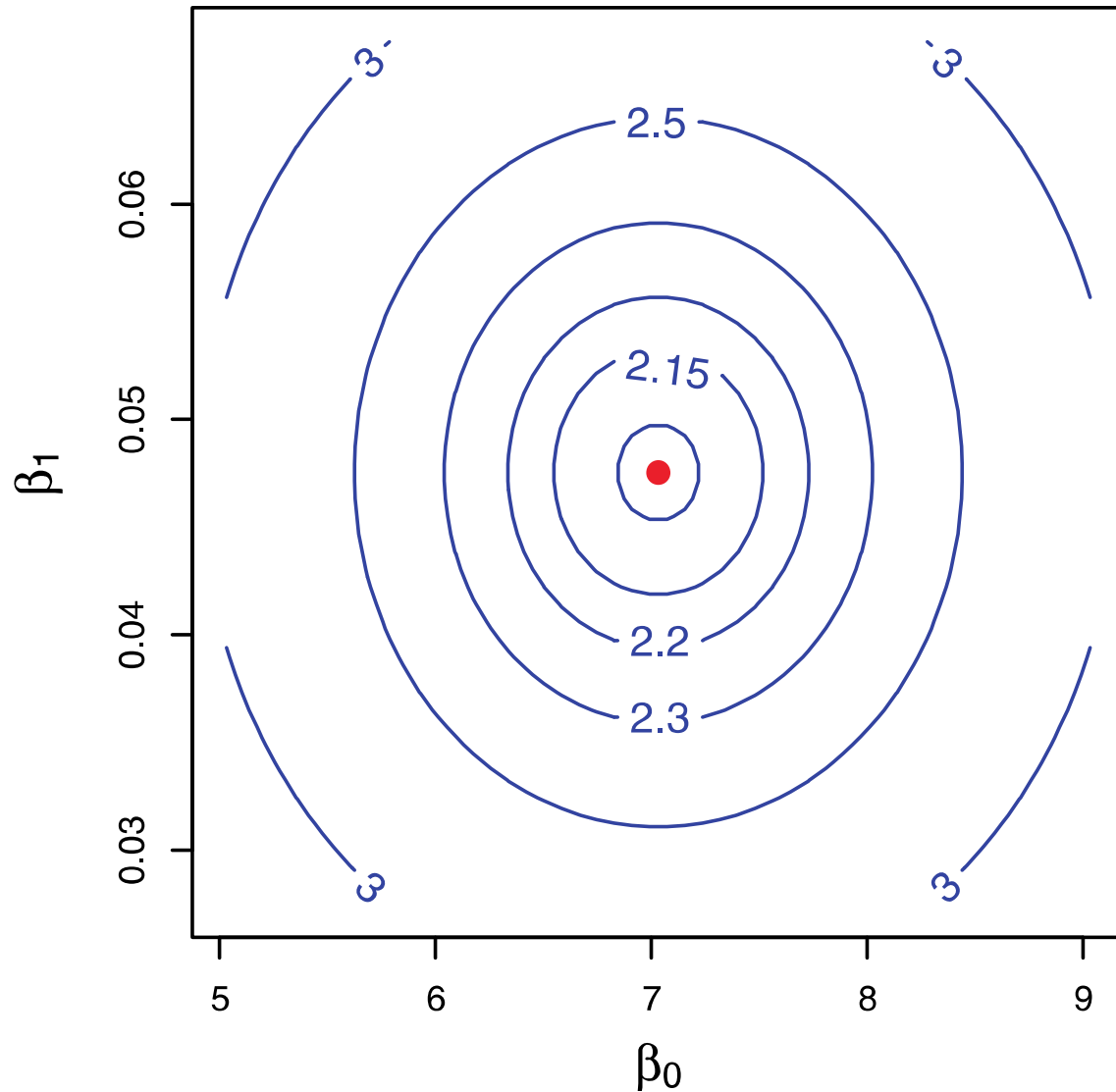


# Linear models I

## Lecture 06

# Quiz



# What is a linear model?

# Which of the following models are linear?

A  $y = w_0$

B  $y = w_0 + w_1x_1$

C  $y = w_0 + w_1x_1 + w_2x_2$

D  $y = w_0 + w_1x_1^2 + w_2x_2^{0.4}$

E  $y = w_0 + w_1x_1 + w_2x_2 + w_3x_1x_2$

F  $y = w_0 + w_1 \int \sqrt[3]{x_1} dx_1 + w_2 g(x_2) + w_3 \text{median}(x_1, x_2, x_3)$

# Which of the following models are linear?

A  $y = w_0$

B  $y = w_0 + w_1 x_1$

C  $y = w_0 + w_1 x_1 + w_2 x_2$

D  $y = w_0 + w_1 x_1^2 + w_2 x_2^{0.4}$

E  $y = w_0 + w_1 x_1 + w_2 x_2 + w_3 x_1 x_2$

F  $y = w_0 + w_1 \int \sqrt[3]{x_1} dx_1 + w_2 g(x_2) + w_3 \text{median}(x_1, x_2, x_3)$

These are **ALL** linear in the **parameters,  $w$**

# Linear models are linear in the **parameters**

**Linear models are linear in the parameters**

They often model nonlinear relationships  
between features and targets

$$y_j = \sum_{i=0}^N w_i x_{i,j} + \epsilon$$

# Linear regression assumptions

1. Linear relationship between feature and target variables
2. Error is normally distributed
3. Features are not correlated with one another (no multicollinearity)
4. Assumes observations are independent from one another (no autocorrelation)
5. Variance of the error is constant (homoscedastic)

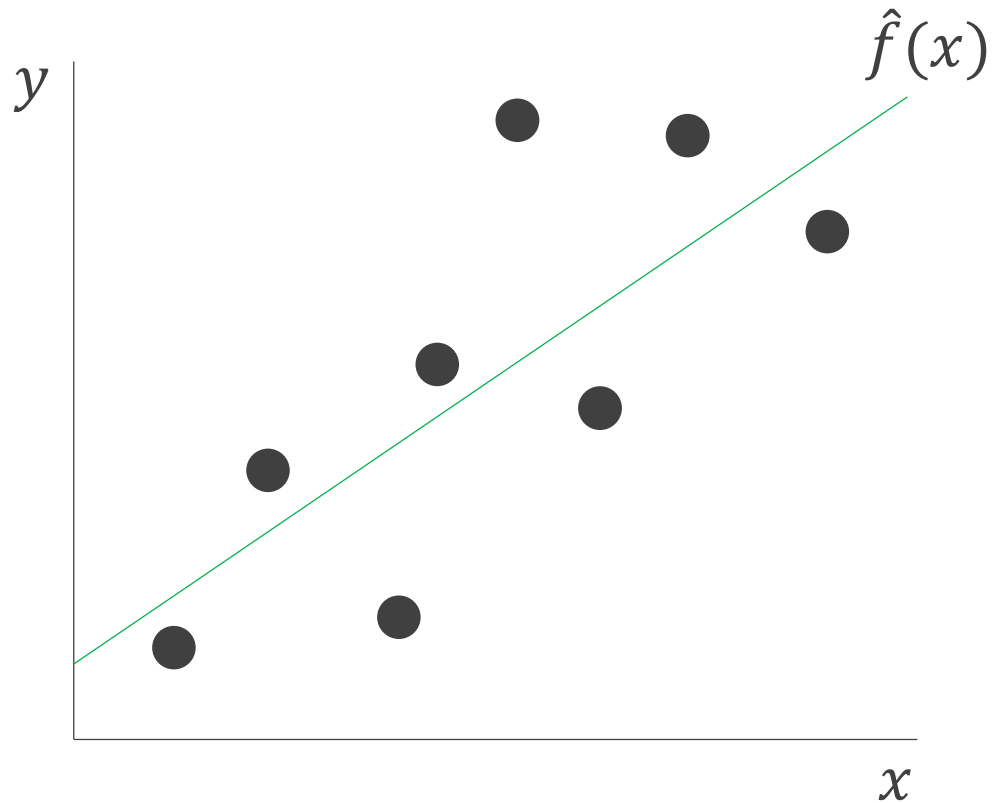


# Types of Linear Regression

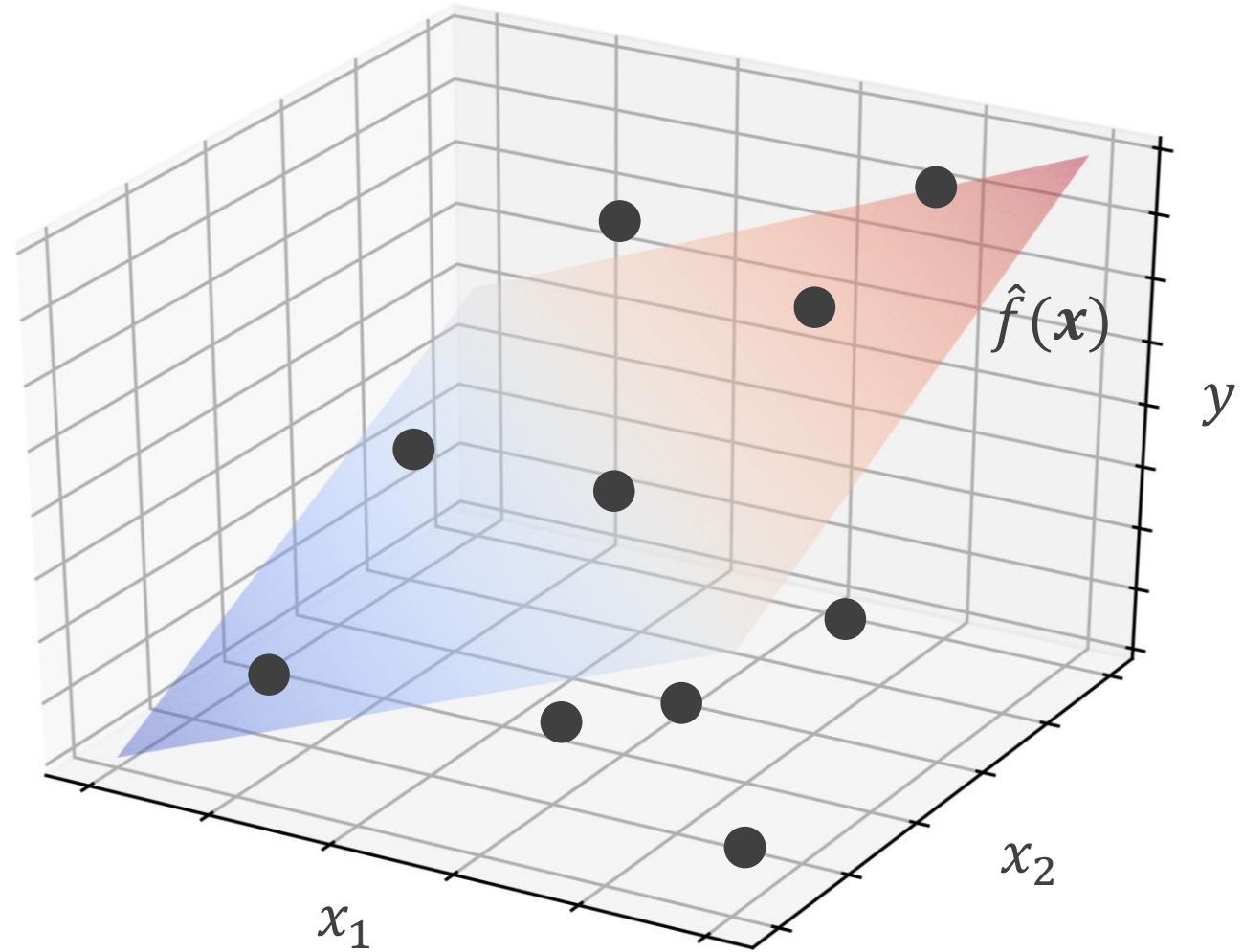
Here,  $x_0 \triangleq 1$

|                              | One feature variable  | One or more feature variables   |
|------------------------------|---|---|
| One target variable          | <b>Simple Linear Regression</b><br>$y = w_0 + w_1 x_1$  | <b>Multiple Linear Regression</b><br>$y = \sum_{i=0}^N w_i x_i \quad \text{or} \quad y = \mathbf{w}^T \mathbf{x}$ |
| One or more target variables | <b>Multivariate (Multiple) Linear Regression</b><br>$\mathbf{y} = \sum_{i=0}^p w_i \mathbf{x}_i \quad \text{or} \quad \mathbf{y} = \mathbf{X} \mathbf{w}$ |   |

# Linear models and error

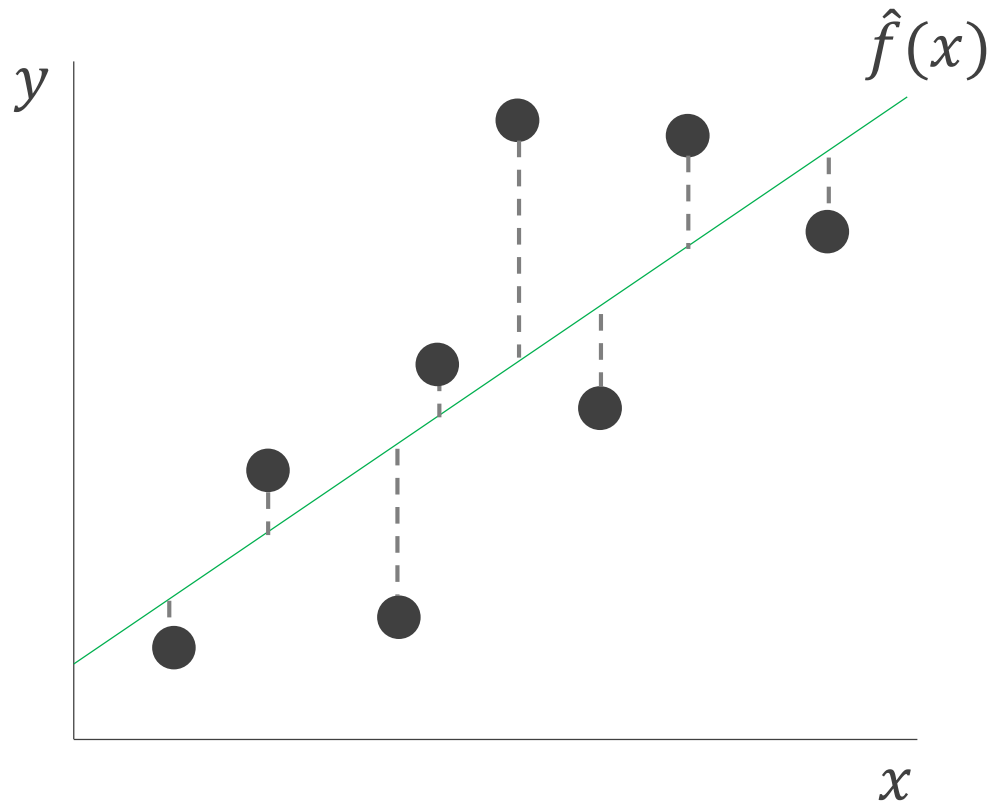


**simple linear regression**

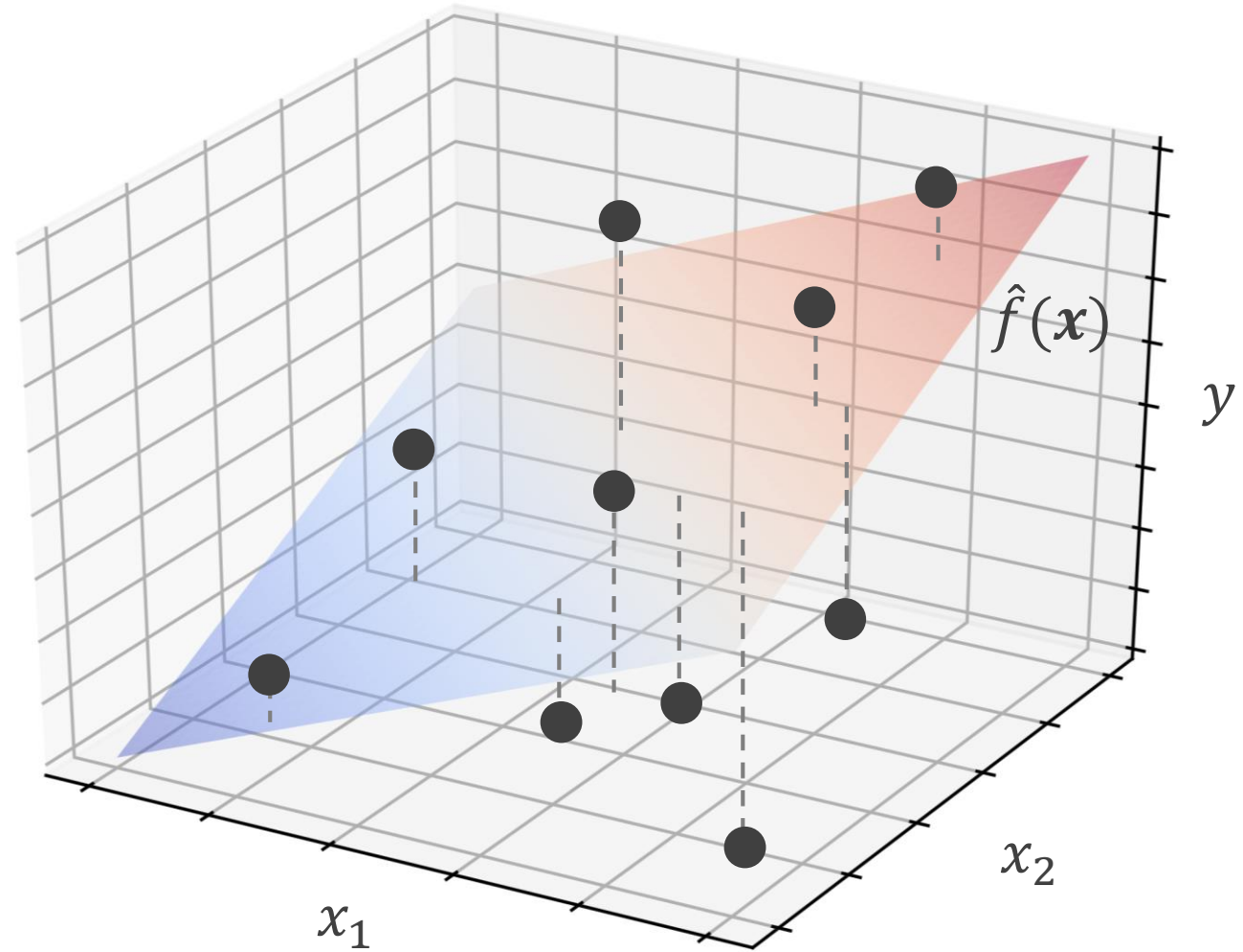


**multiple linear regression**

# Linear models and error

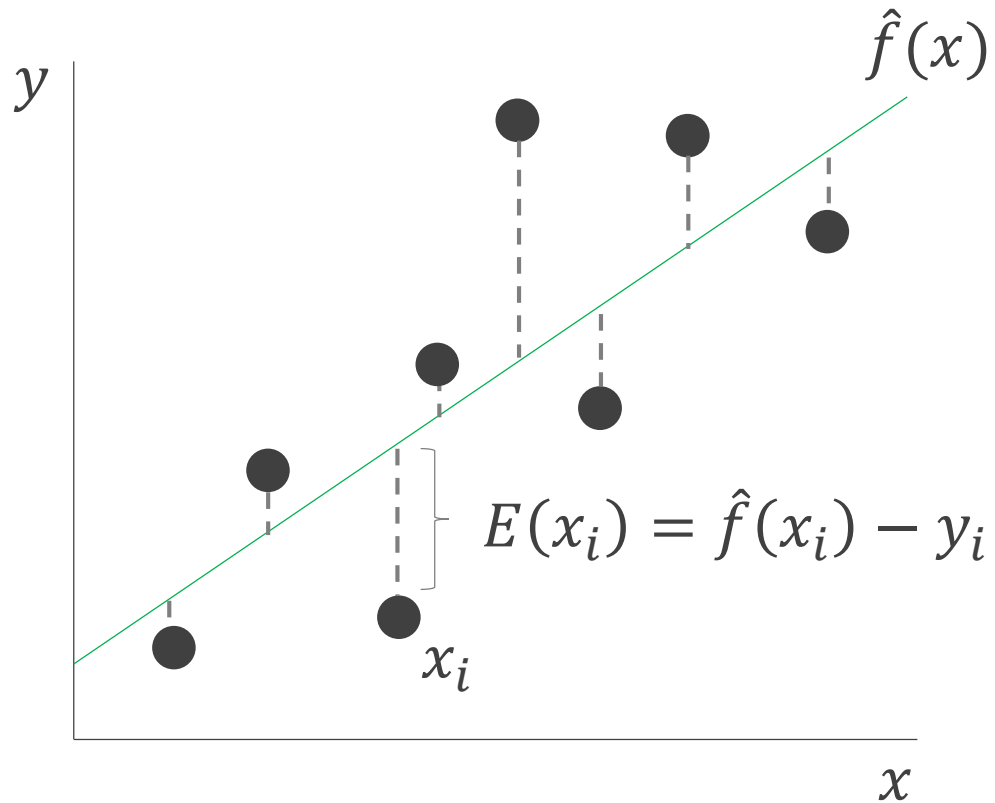


**simple linear regression**

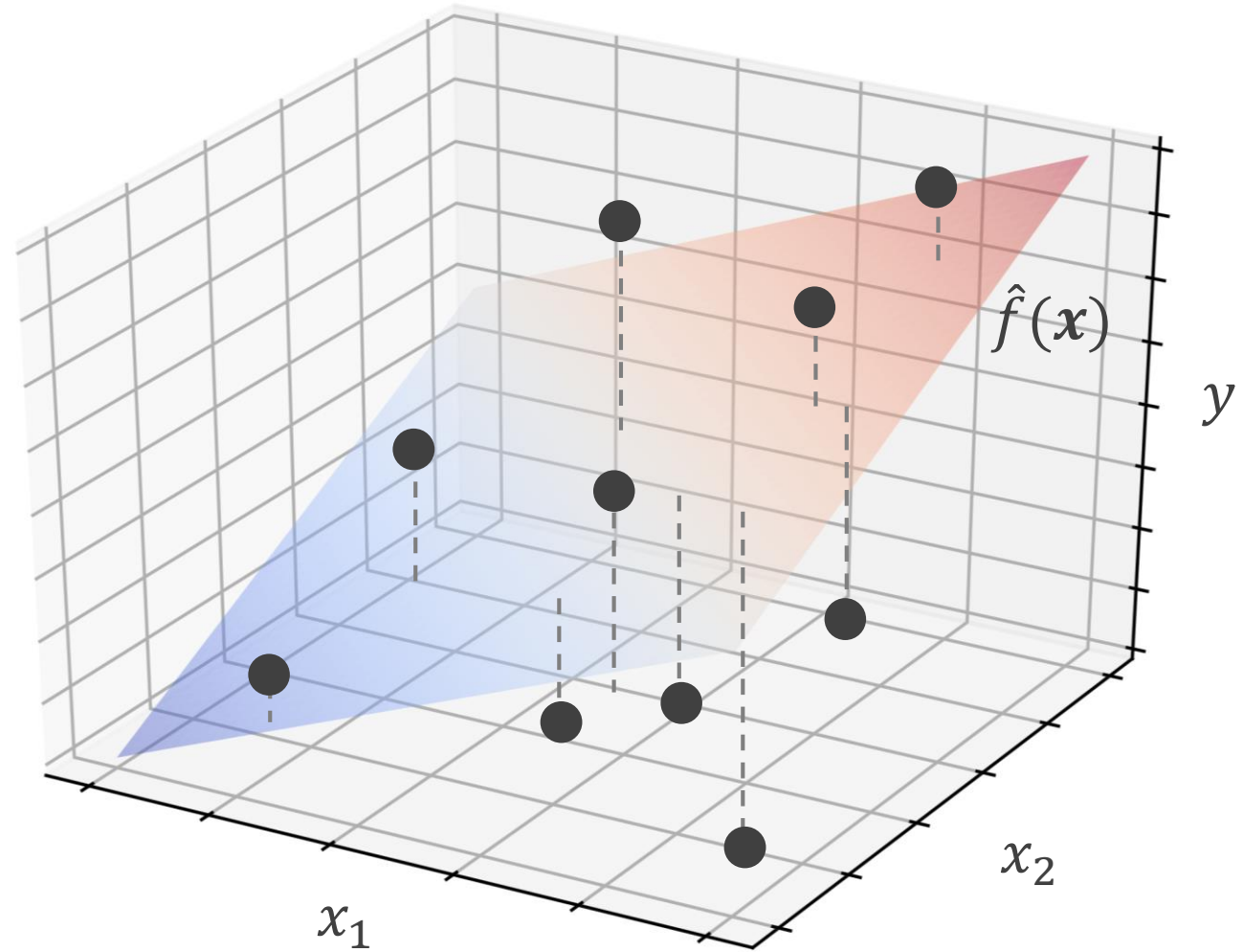


**multiple linear regression**

# Linear models and error



**simple linear regression**



**multiple linear regression**

# How do we fit a linear model to data?

We want the error between our estimates and predictions to be small

# How do we measure error?

How well does  $\hat{y} = \hat{f}(\mathbf{x}) = \sum_{i=0}^N w_i x_i$  approximate  $y$  ?

# How do we measure error?

How well does  $\hat{y} = \hat{f}(\mathbf{x}) = \sum_{i=0}^N w_i x_i$  approximate  $y$  ?

Error: difference between our estimate  $\hat{y}$  and our training data  $y$

$$\text{error} = \hat{y} - y$$

# How do we measure error?

How well does  $\hat{y} = \hat{f}(\mathbf{x}) = \sum_{i=0}^N w_i x_i$  approximate  $y$  ?

Error: difference between our estimate  $\hat{y}$  and our training data  $y$

$$\text{error} = \hat{y} - y$$

We use mean squared error to estimate training (in-sample) error:

**Training (in-sample) error:** 
$$E_{in}(\hat{f}) = \frac{1}{N} \sum_{n=1}^N (\hat{f}(\mathbf{x}_n) - y_n)^2$$



# How do we measure error?

How well does  $\hat{y} = \hat{f}(\mathbf{x}) = \sum_{i=0}^N w_i x_i$  approximate  $y$  ?

Error: difference between our estimate  $\hat{y}$  and our training data  $y$

$$\text{error} = \hat{y} - y$$

We use mean squared error to estimate training (in-sample) error:

**Training (in-sample) error:**  $E_{in}(\hat{f}) = \frac{1}{N} \sum_{n=1}^N (\hat{f}(\mathbf{x}_n) - y_n)^2$

We call this our **Cost Function**

**Cost Function:** 
$$E_{in}(\hat{f}) = \frac{1}{N} \sum_{n=1}^N (\hat{f}(x_n) - y_n)^2$$

**Cost Function:**  $E_{in}(\hat{f}) = \frac{1}{N} \sum_{n=1}^N (\hat{f}(\mathbf{x}_n) - y_n)^2$

Training error is a function of our **model** and the **training data**

**Cost Function:**  $E_{in}(\hat{f}) = \frac{1}{N} \sum_{n=1}^N (\hat{f}(\mathbf{x}_n) - y_n)^2$

Training error is a function of our **model** and the **training data**

We can't change the data, we must adjust our model to minimize cost

**Cost Function:**  $E_{in}(\hat{f}) = \frac{1}{N} \sum_{n=1}^N (\hat{f}(\mathbf{x}_n) - y_n)^2$

Training error is a function of our **model** and the **training data**

We can't change the data, we must adjust our model to minimize cost

We choose model **parameters** that minimize cost

**Cost Function:** 
$$E_{in}(\hat{f}) = \frac{1}{N} \sum_{n=1}^N (\hat{f}(\mathbf{x}_n) - y_n)^2$$

Training error is a function of our **model** and the **training data**

We can't change the data, we must adjust our model to minimize cost

We choose model **parameters** that minimize cost

This is an **optimization** problem

# How to fit our model to the training data?

Equivalently: how do we choose  $\mathbf{w}$  to minimize cost (error)

$$E_{in}(\hat{f}) = \frac{1}{N} \sum_{n=1}^N (\hat{f}(\mathbf{x}_n) - y_n)^2$$

# How to fit our model to the training data?

Equivalently: how do we choose  $\mathbf{w}$  to minimize cost (error)

$$E_{in}(\hat{f}) = \frac{1}{N} \sum_{n=1}^N (\hat{f}(\mathbf{x}_n) - y_n)^2$$

where  $\hat{f}(\mathbf{x}_n) = \mathbf{w}^T \mathbf{x}_n$



# How to fit our model to the training data?

Equivalently: how do we choose  $\mathbf{w}$  to minimize cost (error)

$$E_{in}(\hat{f}) = \frac{1}{N} \sum_{n=1}^N (\hat{f}(\mathbf{x}_n) - y_n)^2$$

where  $\hat{f}(\mathbf{x}_n) = \mathbf{w}^T \mathbf{x}_n$

So we want to minimize  $E_{in}(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n - y_n)^2$

# How to fit our model to the training data?

Equivalently: how do we choose  $\mathbf{w}$  to minimize cost (error)

$$E_{in}(\hat{f}) = \frac{1}{N} \sum_{n=1}^N (\hat{f}(\mathbf{x}_n) - y_n)^2$$

where  $\hat{f}(\mathbf{x}_n) = \mathbf{w}^T \mathbf{x}_n$

So we want to minimize  
...by varying  $\mathbf{w}$

$$E_{in}(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n - y_n)^2$$

# How to fit our model to the training data?

Take the derivative with respect to  $\mathbf{w}$ , set it to zero, and solve for  $\mathbf{w}$

$p$  = number of predictors  
 $N$  = number of data points

# How to fit our model to the training data?

Take the derivative with respect to  $\mathbf{w}$ , set it to zero, and solve for  $\mathbf{w}$

$$\nabla_{\mathbf{w}} E_{in}(\mathbf{w}) = \nabla_{\mathbf{w}} \left( \frac{1}{N} \sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n - y_n)^2 \right)$$

$p$  = number of predictors  
 $N$  = number of data points

# How to fit our model to the training data?

Take the derivative with respect to  $\mathbf{w}$ , set it to zero, and solve for  $\mathbf{w}$

$$\nabla_{\mathbf{w}} E_{in}(\mathbf{w}) = \nabla_{\mathbf{w}} \left( \frac{1}{N} \sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n - y_n)^2 \right)$$

$p$  = number of predictors  
 $N$  = number of data points

$$\nabla_{\mathbf{w}} E_{in}(\mathbf{w}) = \begin{bmatrix} \frac{\partial E_{in}}{\partial w_0} \\ \frac{\partial E_{in}}{\partial w_1} \\ \vdots \\ \frac{\partial E_{in}}{\partial w_p} \end{bmatrix} = \mathbf{0}$$

# How to fit our model to the training data?

Take the derivative with respect to  $\mathbf{w}$ , set it to zero, and solve for  $\mathbf{w}$

$$\nabla_{\mathbf{w}} E_{in}(\mathbf{w}) = \nabla_{\mathbf{w}} \left( \frac{1}{N} \sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n - y_n)^2 \right)$$

$p$  = number of predictors  
 $N$  = number of data points

$$\nabla_{\mathbf{w}} E_{in}(\mathbf{w}) = \begin{bmatrix} \frac{\partial E_{in}}{\partial w_0} \\ \frac{\partial E_{in}}{\partial w_1} \\ \vdots \\ \frac{\partial E_{in}}{\partial w_p} \end{bmatrix} = \mathbf{0}$$

Here we walk through the **ordinary least squares** (OLS) closed-form solution.

# How to fit our model to the training data?

Take the derivative with respect to  $\mathbf{w}$ , set it to zero, and solve for  $\mathbf{w}$

$$\nabla_{\mathbf{w}} E_{in}(\mathbf{w}) = \nabla_{\mathbf{w}} \left( \frac{1}{N} \sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n - y_n)^2 \right)$$

$p$  = number of predictors  
 $N$  = number of data points

$$\nabla_{\mathbf{w}} E_{in}(\mathbf{w}) = \begin{bmatrix} \frac{\partial E_{in}}{\partial w_0} \\ \frac{\partial E_{in}}{\partial w_1} \\ \vdots \\ \frac{\partial E_{in}}{\partial w_p} \end{bmatrix} = \mathbf{0}$$

Here we walk through the **ordinary least squares** (OLS) closed-form solution.

Could have used an iterative approach like **gradient descent**

# How to fit our model to the training data?

We can rewrite our objective function:

$$E_{in}(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n - y_n)^2$$

$p$  = number of predictors  
 $N$  = number of data points



# How to fit our model to the training data?

We can rewrite our objective function:

$$E_{in}(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n - y_n)^2$$

↑  
Scalar

$\mathbf{w}^T \in \mathbb{R}^{1 \times p+1}$   
 $\mathbf{x}_n \in \mathbb{R}^{p+1 \times 1}$

$p$  = number of predictors  
 $N$  = number of data points

# How to fit our model to the training data?

We can rewrite our objective function:

$$E_{in}(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N (\underbrace{\mathbf{w}^T \mathbf{x}_n}_{\text{Scalar}} - y_n)^2$$

$\mathbf{w}^T \in \mathbb{R}^{1 \times p+1}$   
 $\mathbf{x}_n \in \mathbb{R}^{p+1 \times 1}$

Convenient definitions:

$$\mathbf{y} \in \mathbb{R}^{N \times 1}$$
$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_N^T \end{bmatrix} \in \mathbb{R}^{N \times p+1}$$
$$\mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_p \end{bmatrix} \in \mathbb{R}^{p+1 \times 1}$$

$p$  = number of predictors  
 $N$  = number of data points

# How to fit our model to the training data?

We can rewrite our objective function:

$$E_{in}(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N (\underbrace{\mathbf{w}^T \mathbf{x}_n}_{\text{Scalar}} - y_n)^2$$

$\mathbf{w}^T \in \mathbb{R}^{1 \times p+1}$   
 $\mathbf{x}_n \in \mathbb{R}^{p+1 \times 1}$

We can rewrite our objective function:

$$E_{in}(\mathbf{w}) = \frac{1}{N} (\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y})$$

Convenient definitions:

$$\begin{aligned} \mathbf{y} &\in \mathbb{R}^{N \times 1} \\ \mathbf{X} &= \begin{bmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_N^T \end{bmatrix} \in \mathbb{R}^{N \times p+1} \\ \mathbf{w} &= \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_p \end{bmatrix} \in \mathbb{R}^{p+1 \times 1} \end{aligned}$$

$p$  = number of predictors  
 $N$  = number of data points

# How to fit our model to the training data?

$$E_{in}(\mathbf{w}) = \frac{1}{N} (\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y})$$

# How to fit our model to the training data?

$$E_{in}(\mathbf{w}) = \frac{1}{N} (\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y})$$

$$\nabla_{\mathbf{w}} E_{in}(\mathbf{w}) = \frac{2}{N} (\mathbf{X}^T \mathbf{X} \mathbf{w} - \mathbf{X}^T \mathbf{y}) = \mathbf{0}$$

# How to fit our model to the training data?

$$E_{in}(\mathbf{w}) = \frac{1}{N} (\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y})$$

$$\nabla_{\mathbf{w}} E_{in}(\mathbf{w}) = \frac{2}{N} (\mathbf{X}^T \mathbf{X}\mathbf{w} - \mathbf{X}^T \mathbf{y}) = \mathbf{0}$$

$$\mathbf{X}^T \mathbf{X}\mathbf{w} - \mathbf{X}^T \mathbf{y} = \mathbf{0}$$

# How to fit our model to the training data?

$$E_{in}(\mathbf{w}) = \frac{1}{N} (\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y})$$

$$\nabla_{\mathbf{w}} E_{in}(\mathbf{w}) = \frac{2}{N} (\mathbf{X}^T \mathbf{X}\mathbf{w} - \mathbf{X}^T \mathbf{y}) = \mathbf{0}$$

$$\mathbf{X}^T \mathbf{X}\mathbf{w} - \mathbf{X}^T \mathbf{y} = \mathbf{0}$$

$$\mathbf{X}^T \mathbf{X}\mathbf{w} = \mathbf{X}^T \mathbf{y} \quad (\text{normal equation})$$

# How to fit our model to the training data?

$$E_{in}(\mathbf{w}) = \frac{1}{N} (\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y})$$

$$\nabla_{\mathbf{w}} E_{in}(\mathbf{w}) = \frac{2}{N} (\mathbf{X}^T \mathbf{X}\mathbf{w} - \mathbf{X}^T \mathbf{y}) = \mathbf{0}$$

$$\mathbf{X}^T \mathbf{X}\mathbf{w} - \mathbf{X}^T \mathbf{y} = \mathbf{0}$$

$$\mathbf{X}^T \mathbf{X}\mathbf{w} = \mathbf{X}^T \mathbf{y} \quad (\text{normal equation})$$

$$\mathbf{w}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$



# How to fit our model to the training data?

$$E_{in}(\mathbf{w}) = \frac{1}{N} (\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y})$$

$$\nabla_{\mathbf{w}} E_{in}(\mathbf{w}) = \frac{2}{N} (\mathbf{X}^T \mathbf{X}\mathbf{w} - \mathbf{X}^T \mathbf{y}) = \mathbf{0}$$

$$\mathbf{X}^T \mathbf{X}\mathbf{w} - \mathbf{X}^T \mathbf{y} = \mathbf{0}$$

$$\mathbf{X}^T \mathbf{X}\mathbf{w} = \mathbf{X}^T \mathbf{y} \quad (\text{normal equation})$$

$$\mathbf{w}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

**Pseudoinverse**  $\mathbf{X}^\dagger = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$

# How to fit our model to the training data?

$$E_{in}(\mathbf{w}) = \frac{1}{N} (\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y})$$

$$\nabla_{\mathbf{w}} E_{in}(\mathbf{w}) = \frac{2}{N} (\mathbf{X}^T \mathbf{X}\mathbf{w} - \mathbf{X}^T \mathbf{y}) = \mathbf{0}$$

$$\mathbf{X}^T \mathbf{X}\mathbf{w} - \mathbf{X}^T \mathbf{y} = \mathbf{0}$$

$$\mathbf{X}^T \mathbf{X}\mathbf{w} = \mathbf{X}^T \mathbf{y} \quad (\text{normal equation})$$

$$\mathbf{w}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

**Pseudoinverse**  $\mathbf{X}^\dagger = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$

$$\mathbf{w}^* = \mathbf{X}^\dagger \mathbf{y}$$

# What is the pseudoinverse?

$$\begin{array}{c} \text{Features} \\ N \times p \end{array} \begin{array}{c} \text{Samples} \\ \left[ \begin{array}{c} \text{ } \\ \text{ } \\ \text{ } \end{array} \right] \end{array} \begin{array}{c} \left[ \begin{array}{c} \text{ } \\ \text{ } \\ \text{ } \end{array} \right] \\ p \times 1 \end{array} = \begin{array}{c} \left[ \begin{array}{c} \text{ } \\ \text{ } \\ \text{ } \end{array} \right] \\ N \times 1 \end{array}$$
$$X w = y$$

# What is the pseudoinverse?

$$\begin{array}{c} \text{Features} \\ N \times p \end{array} \begin{array}{c} \text{Samples} \\ \left[ \begin{array}{c} \\ \\ \end{array} \right] \end{array} \begin{array}{c} \left[ \begin{array}{c} \\ \\ \end{array} \right] \\ p \times 1 \end{array} = \begin{array}{c} \left[ \begin{array}{c} \\ \\ \end{array} \right] \\ N \times 1 \end{array}$$
$$X w = y$$

If  $N = p$ , then there are the same number of features as samples

# What is the pseudoinverse?

$$\begin{array}{c} \text{Features} \\ N \times p \end{array} \begin{array}{c} \text{Samples} \\ \left[ \begin{array}{c} \text{ } \end{array} \right] \\ p \times 1 \end{array} = \begin{array}{c} \left[ \begin{array}{c} \text{ } \end{array} \right] \\ N \times 1 \end{array}$$

$X \quad w = y$

If  $N = p$ , then there are the same number of features as samples

If  $N > p$ , then the system of equations is overdetermined (more samples than features)

# What is the pseudoinverse?

Consider the case when  $N = 3$ ,  $p = 2$

# What is the pseudoinverse?

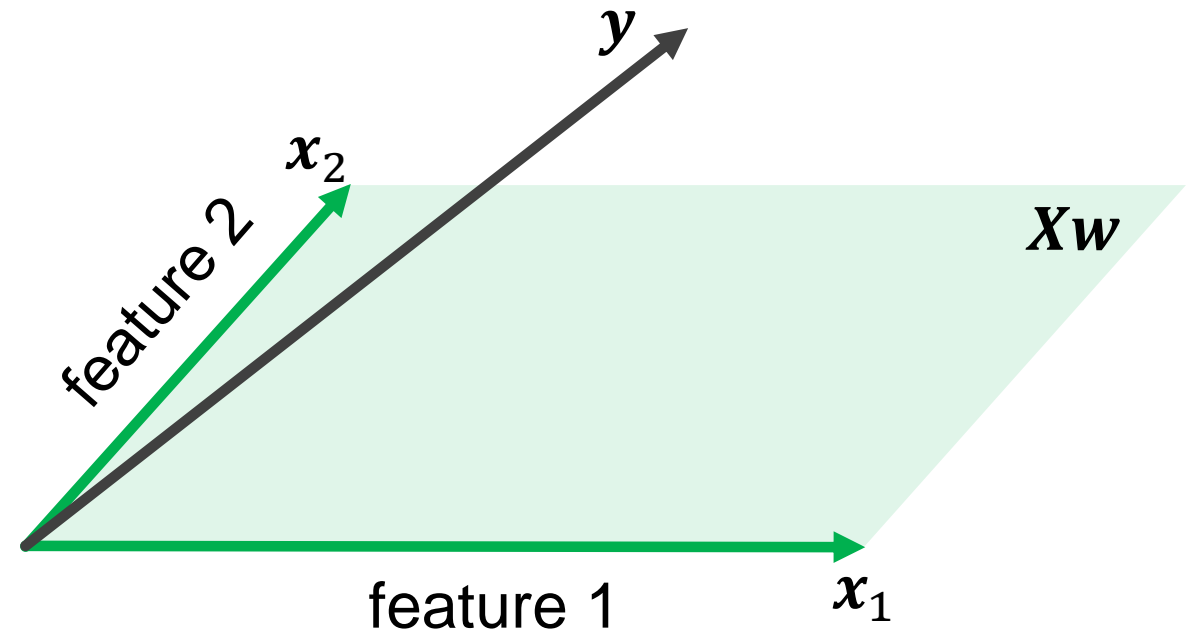
Consider the case when  $N = 3$ ,  $p = 2$

$$\begin{array}{c} \text{Samples} \end{array} \begin{array}{c} \text{Features} \\ \begin{bmatrix} x_{1,1} & x_{1,2} \\ x_{2,1} & x_{2,2} \\ x_{3,1} & x_{3,2} \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix} \end{array}$$

# What is the pseudoinverse?

Consider the case when  $N = 3$ ,  $p = 2$

$$\begin{array}{c} \text{Samples} \end{array} \begin{array}{c} \text{Features} \\ \begin{bmatrix} x_{1,1} & x_{1,2} \\ x_{2,1} & x_{2,2} \\ x_{3,1} & x_{3,2} \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix} \end{array}$$



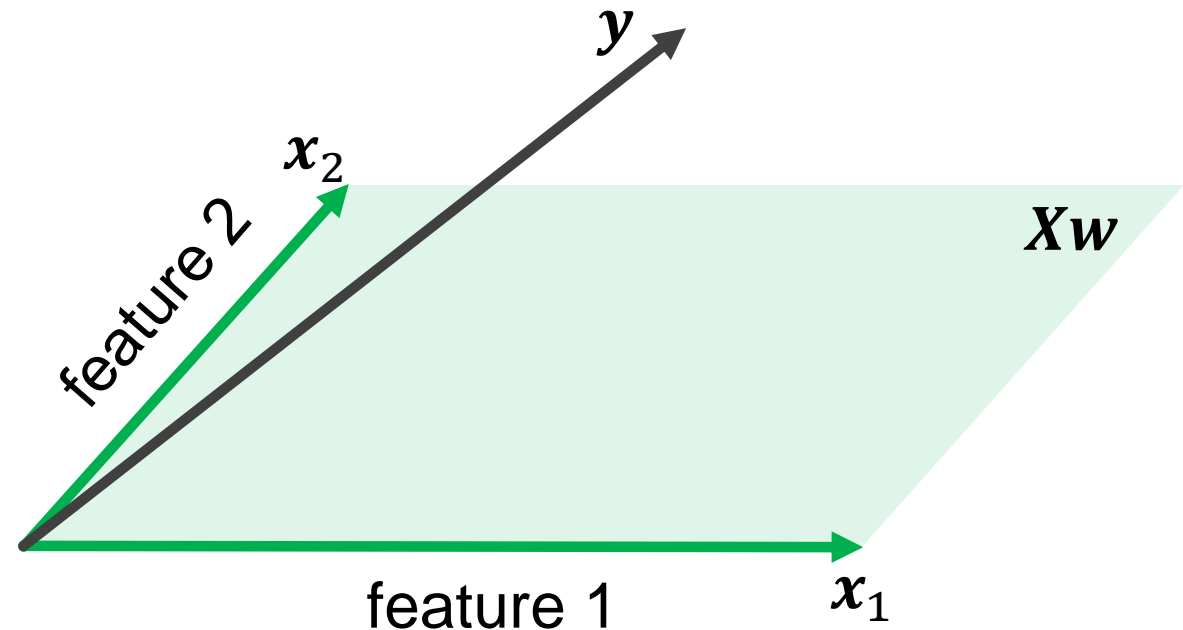


# What is the pseudoinverse?

Consider the case when  $N = 3, p = 2$

$$\begin{array}{c} \text{Samples} \\ \begin{bmatrix} x_{1,1} & x_{1,2} \\ x_{2,1} & x_{2,2} \\ x_{3,1} & x_{3,2} \end{bmatrix} \end{array} \begin{array}{c} \text{Features} \\ \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} \end{array} = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix}$$

$X \quad w \neq y$



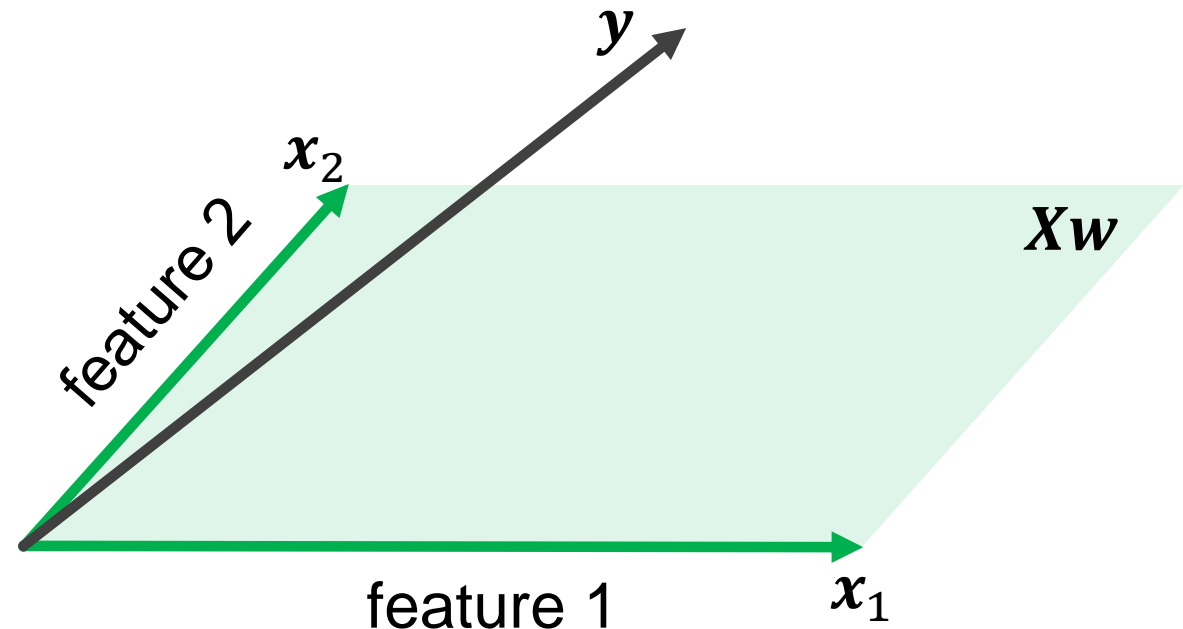
# What is the pseudoinverse?

Consider the case when  $N = 3, p = 2$

$$\begin{array}{c} \text{Features} \\ \text{Samples} \end{array} \begin{bmatrix} x_{1,1} & x_{1,2} \\ x_{2,1} & x_{2,2} \\ x_{3,1} & x_{3,2} \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix}$$

$$X w \neq y$$

We CAN solve for the least squares solution:  $X^\dagger w^* = y$



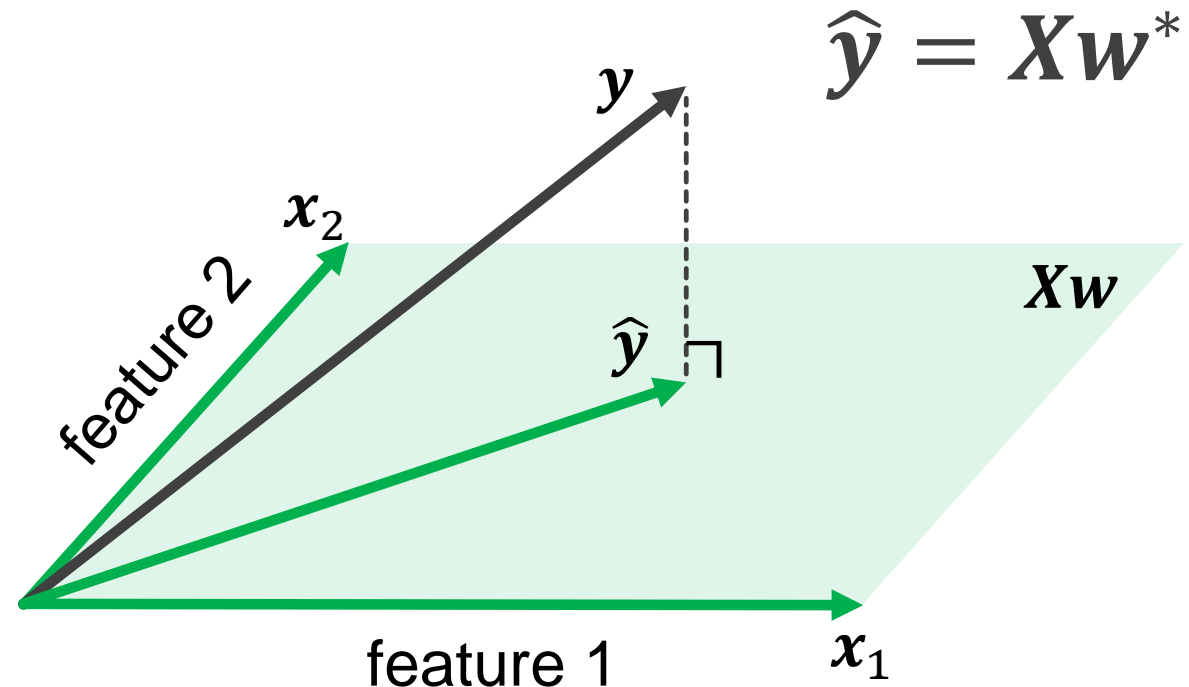
# What is the pseudoinverse?

Consider the case when  $N = 3, p = 2$

$$\begin{array}{c} \text{Features} \\ \text{Samples} \end{array} \begin{bmatrix} x_{1,1} & x_{1,2} \\ x_{2,1} & x_{2,2} \\ x_{3,1} & x_{3,2} \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix}$$

$$X w \neq y$$

We CAN solve for the least squares solution:  $X^\dagger w^* = y$



# What is the pseudoinverse?

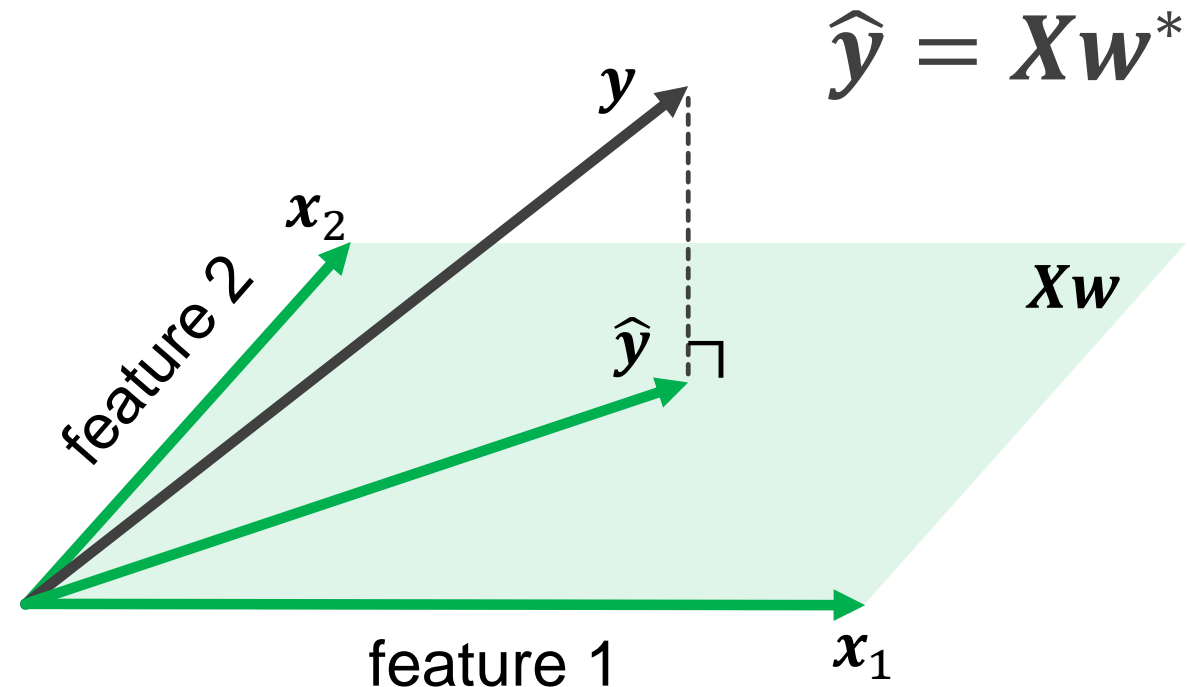
Consider the case when  $N = 3, p = 2$

The least squares solution is the best we can do given  $N > p$

$$\begin{array}{c} \text{Features} \\ \text{Samples} \end{array} \begin{bmatrix} x_{1,1} & x_{1,2} \\ x_{2,1} & x_{2,2} \\ x_{3,1} & x_{3,2} \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix}$$

$$X w \neq y$$

We CAN solve for the least squares solution:  $X^\dagger w^* = y$



# Common paradigm for model fitting

# Common paradigm for model fitting

1. Choose a **hypothesis set of models** to train  
(e.g. linear regression with 4 predictor variables)

# Common paradigm for model fitting

1. Choose a **hypothesis set of models** to train  
(e.g. linear regression with 4 predictor variables)
2. Identify a **cost function** to measure the model fit to the training data  
(e.g. mean square error)

# Common paradigm for model fitting

1. Choose a **hypothesis set of models** to train  
(e.g. linear regression with 4 predictor variables)
2. Identify a **cost function** to measure the model fit to the training data  
(e.g. mean square error)
3. **Optimize model parameters** to minimize cost  
(e.g. closed form solution using the normal equations for OLS)



**Much of machine learning is  
optimizing a cost function**

# What about classification?

# Moving from regression to classification

**Regression**

$$y = \sum_{i=0}^N w_i x_i$$

**Classification**  
(perceptron model)

$$y = \sum_{i=0}^N w_i x_i$$

Source: Abu-Mostafa, Learning from Data, Caltech

# Moving from regression to classification

**Regression**

$$y = \sum_{i=0}^N w_i x_i$$

**Classification**  
(perceptron model)

$$y = \text{sign} \left( \sum_{i=0}^N w_i x_i \right)$$

Source: Abu-Mostafa, Learning from Data, Caltech

# Moving from regression to classification

**Regression**

$$y = \sum_{i=0}^N w_i x_i$$

**Classification**  
(perceptron model)

$$y = \text{sign} \left( \sum_{i=0}^N w_i x_i \right)$$

where

$$\text{sign}(x) = \begin{cases} 1 & x > 0 \\ -1 & \text{else} \end{cases}$$

Source: Abu-Mostafa, Learning from Data, Caltech

# Moving from regression to classification

**Regression**

$$y = \sum_{i=0}^N w_i x_i$$

**Classification**  
(perceptron model)

$$y = \text{sign} \left( \sum_{i=0}^N w_i x_i \right)$$

$$y = \begin{cases} 1 & \sum_{i=0}^N w_i x_i > 0 \\ -1 & \text{else} \end{cases}$$

where

$$\text{sign}(x) = \begin{cases} 1 & x > 0 \\ -1 & \text{else} \end{cases}$$

Source: Abu-Mostafa, Learning from Data, Caltech

# Moving from regression to classification

## Linear Regression

$$\hat{f}(\mathbf{x}) = \sum_{i=0}^N w_i x_i$$

## Linear Classification (perceptron)

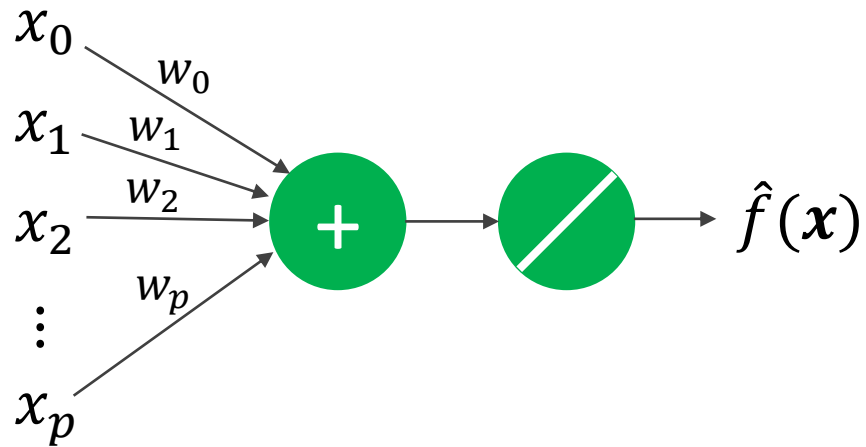
$$\hat{f}(\mathbf{x}) = \text{sign} \left( \sum_{i=0}^N w_i x_i \right)$$

Source: Abu-Mostafa, Learning from Data, Caltech

# Moving from regression to classification

## Linear Regression

$$\hat{f}(\mathbf{x}) = \sum_{i=0}^N w_i x_i$$



## Linear Classification (perceptron)

$$\hat{f}(\mathbf{x}) = \text{sign} \left( \sum_{i=0}^N w_i x_i \right)$$

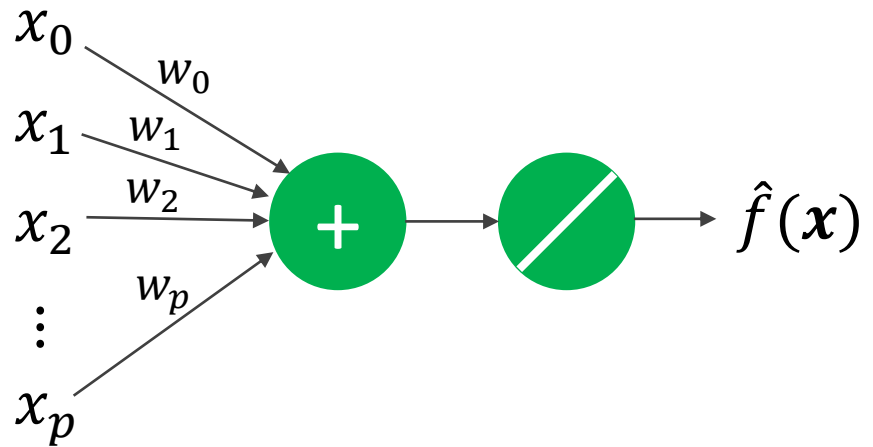
Source: Abu-Mostafa, Learning from Data, Caltech



# Moving from regression to classification

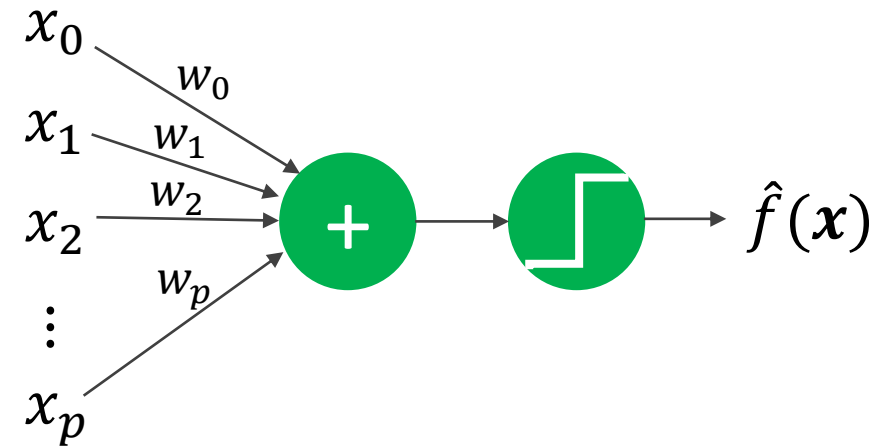
## Linear Regression

$$\hat{f}(\mathbf{x}) = \sum_{i=0}^N w_i x_i$$



## Linear Classification (perceptron)

$$\hat{f}(\mathbf{x}) = \text{sign} \left( \sum_{i=0}^N w_i x_i \right)$$

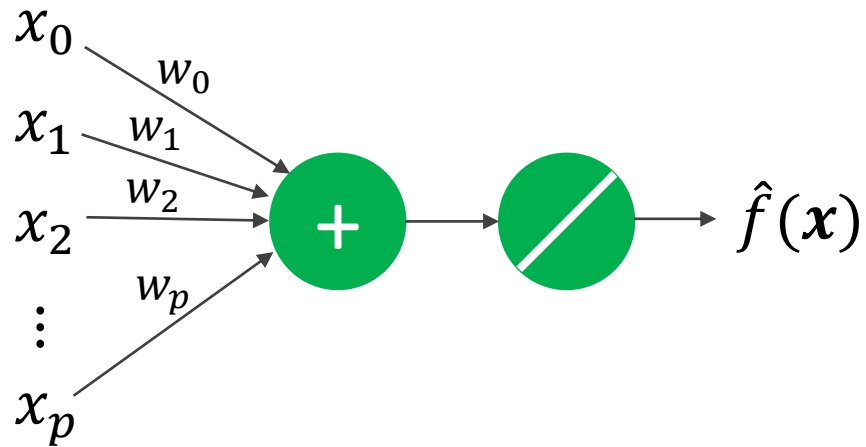


Source: Abu-Mostafa, Learning from Data, Caltech

# Moving from regression to classification

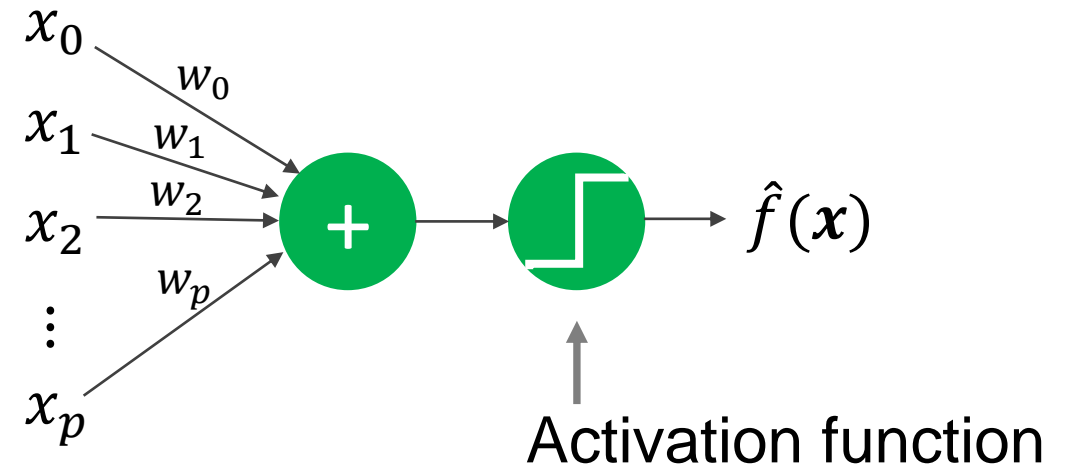
## Linear Regression

$$\hat{f}(\mathbf{x}) = \sum_{i=0}^N w_i x_i$$



## Linear Classification (perceptron)

$$\hat{f}(\mathbf{x}) = \text{sign} \left( \sum_{i=0}^N w_i x_i \right)$$



Source: Abu-Mostafa, Learning from Data, Caltech

# Takeaways

Linear models are **linear in the weights**

Linear models can be used for **both regression and classification**

Model fitting/training (valid beyond linear models):

- Choose a hypothesis set of models to train
- Identify a cost function
- **Optimize the cost function** by adjusting model parameters