

# ITCS 6156/8156 Fall 2023

## Machine Learning

# Linear Regression

Instructor: Hongfei Xue

Email: [hongfei.xue@charlotte.edu](mailto:hongfei.xue@charlotte.edu)

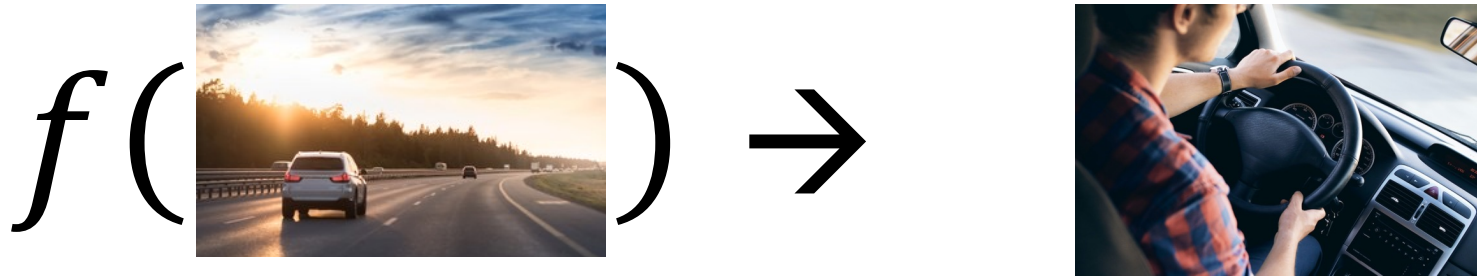
Class Meeting: Mon & Wed, 4:00 PM – 5:15 PM, CHHS 376



Some content in the slides is based on Dr. Razvan's lecture

# Machine Learning

- Function is everywhere!
  - Function  $\rightarrow f$  ; Input instance  $\rightarrow x$ ; Output Target  $\rightarrow y$



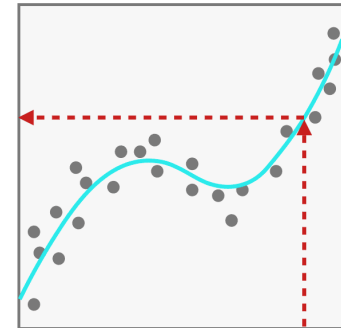
- Machine Learning Task:
  - learn an (unknown) function  $f: X \rightarrow Y$  that maps input instances  $x \in X$  to output targets  $f(x) \in Y$ .

# Classification vs. Regression

- Machine Learning Task:
  - learn an (unknown) function  $f: X \rightarrow Y$  that maps input instances  $x \in X$  to output targets  $f(x) \in Y$ .
- Linear Regression is a **Regression** algorithm.

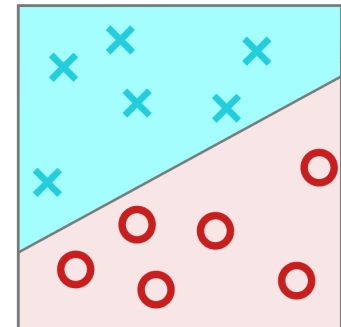
- **Regression:**

- The output targets  $f(x) \in Y$  is continuous, or has a continuous component.



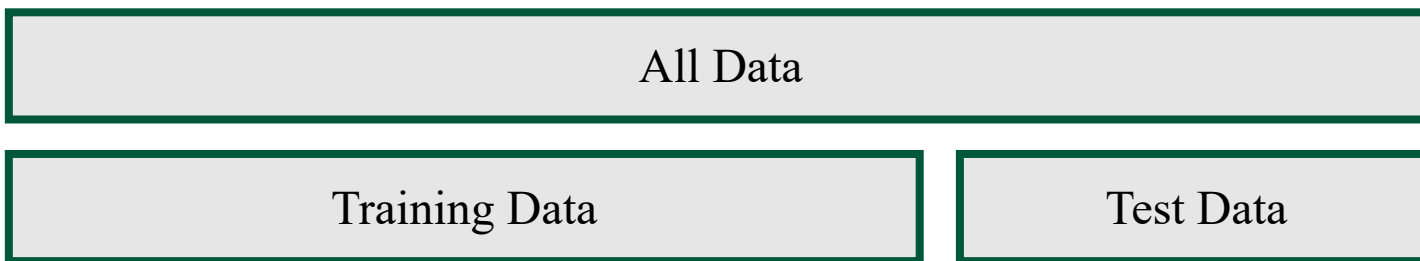
- **Classification:**

- The output targets  $f(x) \in Y$  is one of a finite set of discrete categories.



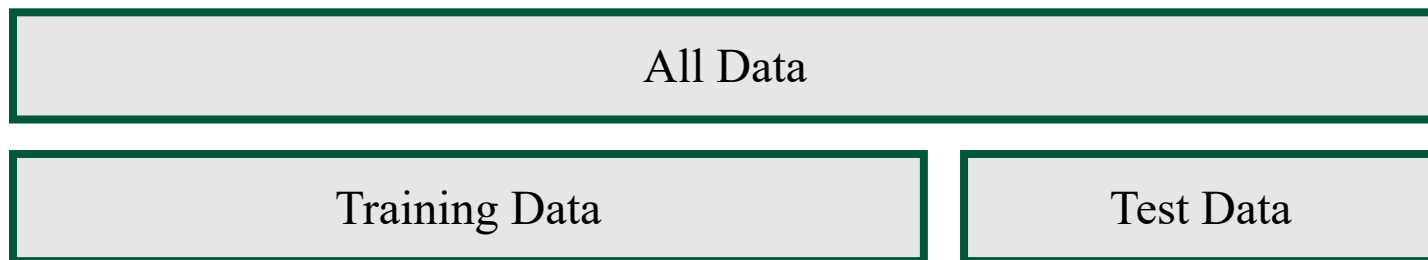
# Supervised Learning

- Machine Learning Task:
  - learn an (unknown) function  $f: X \rightarrow Y$  that maps input instances  $x \in X$  to output targets  $f(x) \in Y$ .
- Linear Regression is a Supervised Learning algorithm.
- Supervised Learning: The output targets are known in the set of training examples:
  - $(x_1, y_1), (x_2, y_2) \dots (x_i, y_i)$
- Training data vs. Test data



# Supervised Learning

- Training data vs. Test data



- The goal of machine learning algorithms is to build a function  $h(x)$  such that:
  - $h$  matches  $f$  well on the **training data**  $\rightarrow h$  is able to **fit data** that it has seen
  - $h$  also  $f$  well on the **test data**  $\rightarrow h$  is able to **generalize** to unseen data
- To achieve the goal, we want to choose  $h$  from a “nice” class of functions that depends on a vector of parameters  $w$ :
  - $h(x) \equiv h_w(x) \equiv h(w, x)$

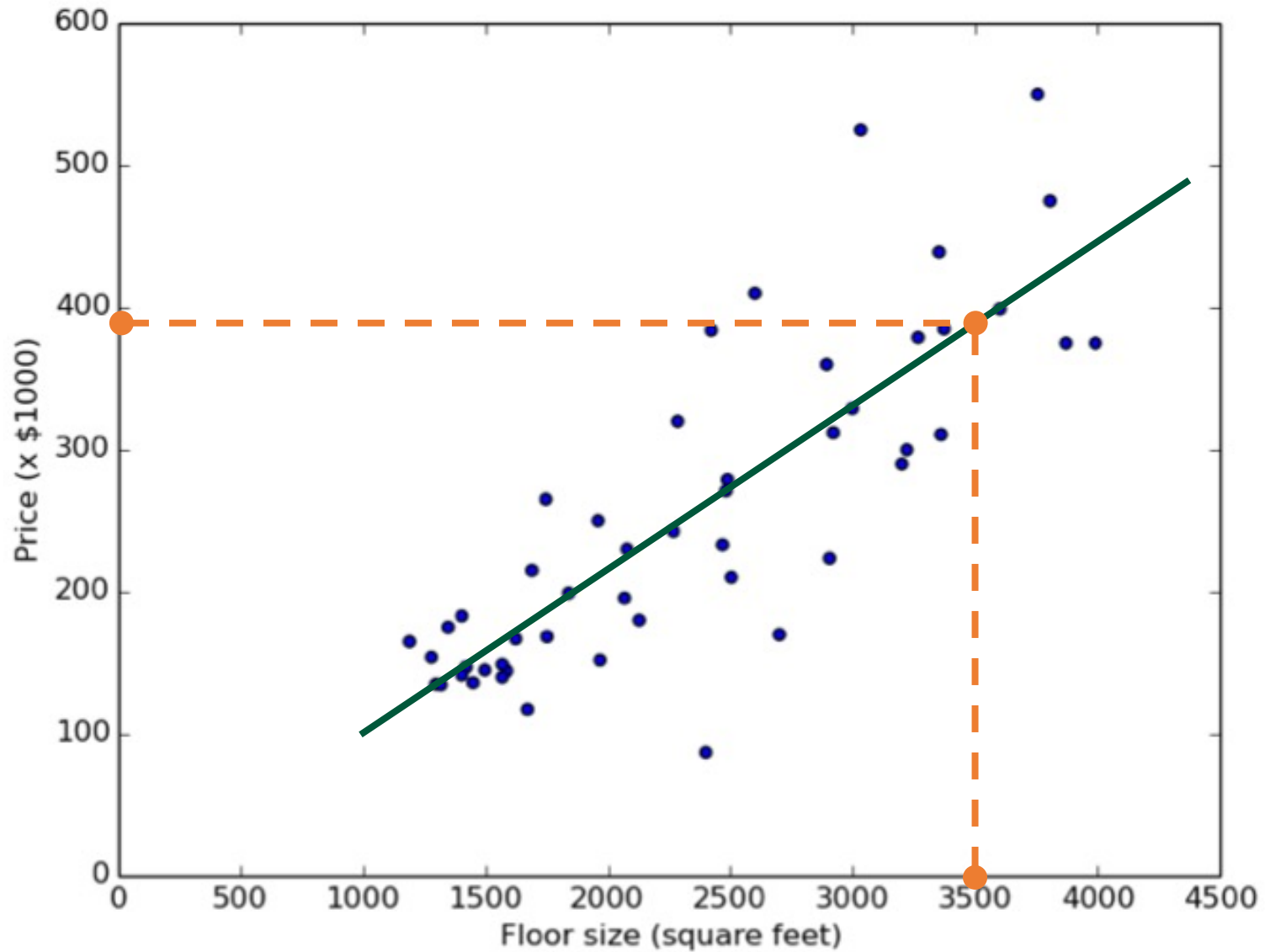


# House Price Prediction

- Given the floor size in square feet, predict the selling price:
  - Input  $x$ : the floor size of the house
  - Output  $y$ : the selling price of the house
  - Need to learn a function  $h$  such that  $h(x) \approx f(x)$
- Is this a classification or regression task?
  - **Regression**, because the house price is real-valued.
  - (Simple) linear regression, because only one input value.
  - Would a problem with only two labels  $y_1 = 0.5$  and  $y_2 = 1.0$  still be regression?



# House Price Prediction



# Linear Regression

- Use a linear function to approximate the real (unknown) function:
  - $h_{\mathbf{w}}(X) = \mathbf{w}^T X = [w_0, w_1]^T [1, x] = w_1 x + w_0$
  - In our case,  $h_{\mathbf{w}}(X)$  is a straight line
    - $w_0$  is the intercept (or the bias term)
    - $w_1$  controls the slope
- Actually, the floor size of the house is not the only factor determining its sell price. There are many factors:  $(x_1, \dots, x_d)$ .
  - $$h_{\mathbf{w}}(X) = \mathbf{w}^T X = [w_0, w_1, \dots, w_d]^T [1, x_1, \dots, x_d]$$
$$= w_0 + w_1 x_1 + w_2 x_2 + \dots + w_d x_d$$



# Error Measurement

- Our linear approximation function:
  - $$h_{\mathbf{w}}(X) = \mathbf{w}^T X = [w_0, w_1, \dots, w_d]^T [1, x, \dots, x_d]$$
$$= w_0 + w_1 x_1 + w_2 x_2 + \dots + w_d x_d$$
- Error Measurement: Find  $\mathbf{w}$  that obtains the best fit on the training data, i.e. find  $\mathbf{w}$  that minimizes the **sum of square errors**:

$$J(\mathbf{w}) = \frac{1}{2N} \sum_{n=1}^N (h_{\mathbf{w}}(X_n) - y_n)^2$$
$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} J(\mathbf{w})$$

- $N$ : Total number of samples in training set
  - $J(\mathbf{w})$ : Error function
  - $\hat{\mathbf{w}}$ : Optimal  $\mathbf{w}$  that minimize  $J(\mathbf{w})$
- **Why do we use the square errors?**

# Inductive Learning Hypothesis

- Learning = finding the “right” parameters  $\mathbf{w}^T = [w_0, w_1, \dots, w_d]$ 
  - Find  $\mathbf{w}$  that minimizes an error function  $J(\mathbf{w})$  which measures the misfit between  $h(\mathbf{x}_i, \mathbf{w})$  and  $t_i$ .
  - Expect that  $h(\mathbf{x}, \mathbf{w})$  performing well on training examples  $\mathbf{x}_i \Rightarrow h(\mathbf{x}, \mathbf{w})$  will perform well on arbitrary test examples  $\mathbf{x}_j \in \mathbf{X}$ .



Inductive Learning Hypothesis

# Matrix Notation

- **Linear Regression Learning Task**
  - learn  $\mathbf{w}$  given training examples  $\langle \mathbf{X}, \mathbf{y} \rangle$ .
  - The training data is denoted as  $\langle \mathbf{X}, \mathbf{y} \rangle$ , where  $\mathbf{X}$  is a  $N \times D$  data matrix consisting of  $N$  data examples such that each data example is a  $D$  dimensional vector.  $\mathbf{y}$  is a  $N \times 1$  vector consisting of corresponding target values for the examples in  $\mathbf{X}$ .
- The derivation of the least squares estimate can be done by first converting the expression of the squared loss into **matrix notation**, i.e.,

$$J(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N (y_i - \mathbf{w}^T \mathbf{x}_i)^2 = \frac{1}{2} (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w})$$

# Analytical Solution

- To minimize the error, we first compute its derivative with respect to  $\mathbf{w}$ :

$$\frac{\partial LL(\mathbf{w})}{\partial \mathbf{w}} = \frac{1}{2} \frac{\partial}{\partial \mathbf{w}} (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w})$$

- Note that, we use the fact that  $(\mathbf{X}\mathbf{w})^T \mathbf{y} = \mathbf{y}^T \mathbf{X}\mathbf{w}$ , since both quantities are scalars, and the transpose of a scalar is equal to itself. Continuing with the derivative:

$$\frac{\partial LL(\mathbf{w})}{\partial \mathbf{w}} = \frac{1}{2} (2\mathbf{w}^T \mathbf{X}^T \mathbf{X} - 2\mathbf{y}^T \mathbf{X})$$

# Analytical Solution

- Setting the derivation to 0, we get:

$$2\mathbf{w}^\top \mathbf{X}^\top \mathbf{X} - 2\mathbf{y}^\top \mathbf{X} = 0$$

$$\mathbf{w}^\top \mathbf{X}^\top \mathbf{X} = \mathbf{y}^\top \mathbf{X}$$

$$(\mathbf{X}^\top \mathbf{X})^\top \mathbf{w} = \mathbf{X}^\top \mathbf{y} \text{ (Taking transpose both sides)}$$

$$(\mathbf{X}^\top \mathbf{X}) \mathbf{w} = \mathbf{X}^\top \mathbf{y}$$

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



The Moore-Penrose pseudo-inverse

# Summary

- Our **linear approximation function**:

- $$h_{\mathbf{w}}(X) = \mathbf{w}^T X = [w_0, w_1, \dots, w_d]^T [1, x, \dots, x_d]$$
$$= w_0 + w_1 x_1 + w_2 x_2 + \dots + w_d x_d$$

- Error Measurement: Find  $\mathbf{w}$  that minimizes the **sum of square errors**:

$$J(\mathbf{w}) = \frac{1}{2N} \sum_{n=1}^N (h_{\mathbf{w}}(X_n) - y_n)^2$$
$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} J(\mathbf{w})$$

- Analytical Solution:

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



# Questions?