# Advanced Machine Learning and Deep Learning

**Dr. Mohammad Pourhomayoun**
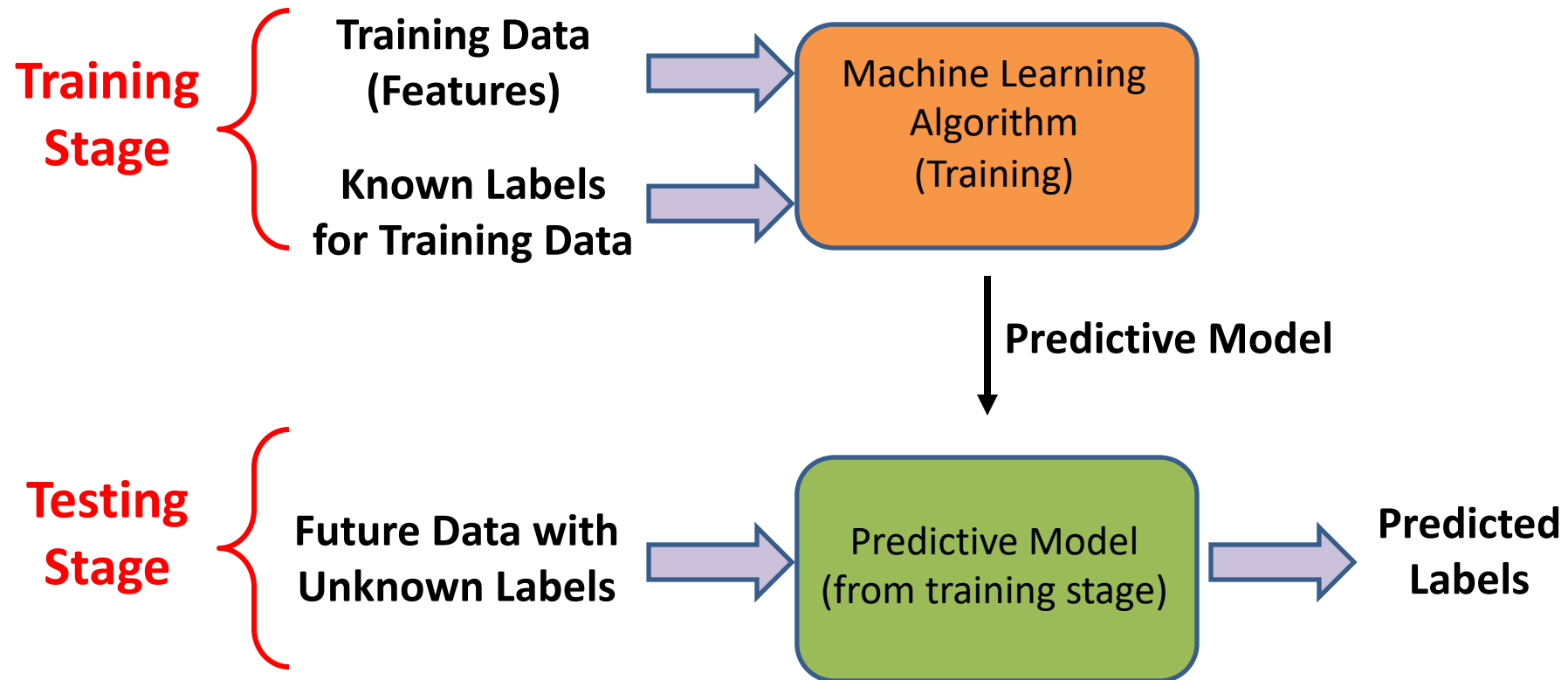
Assistant Professor

Computer Science Department

California State University, Los Angeles
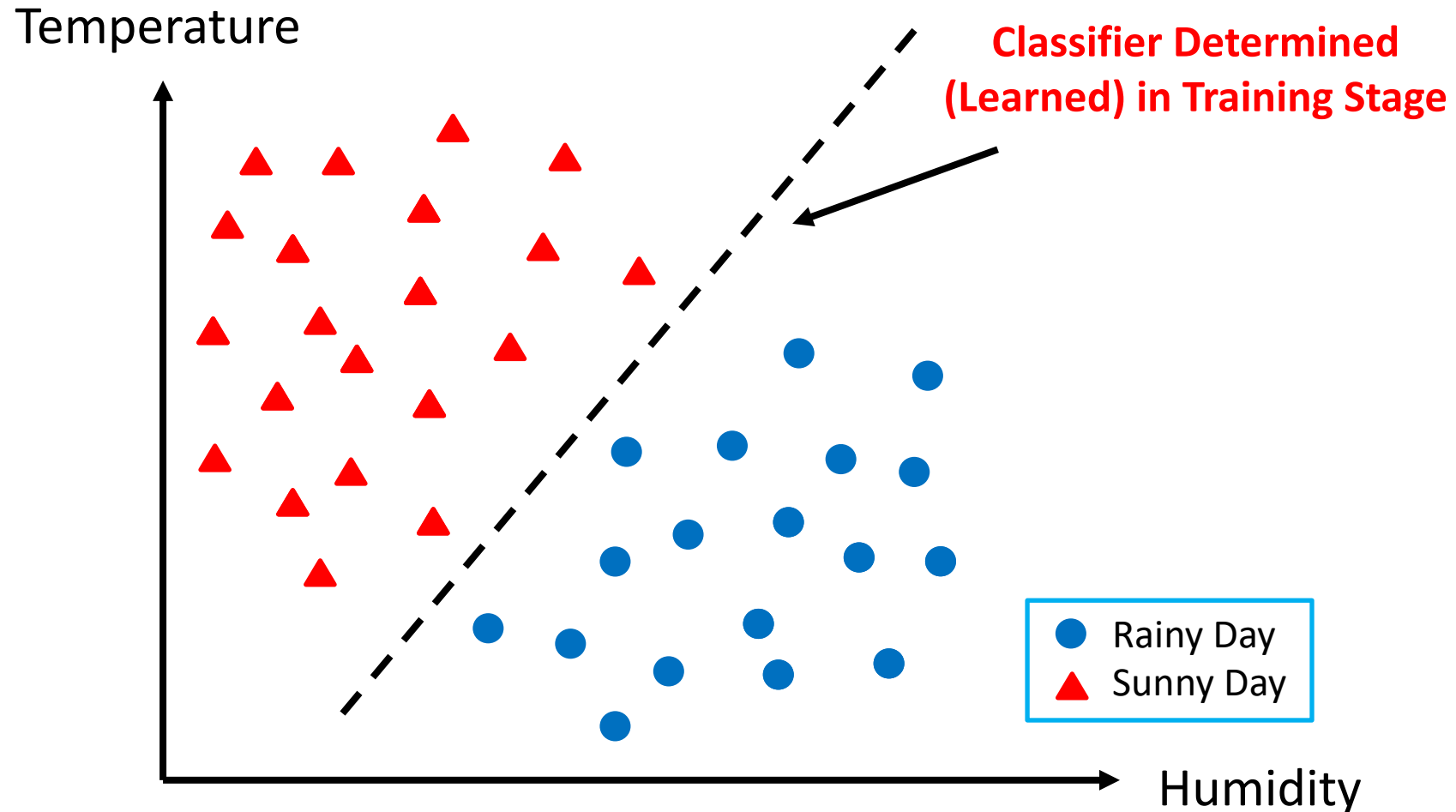
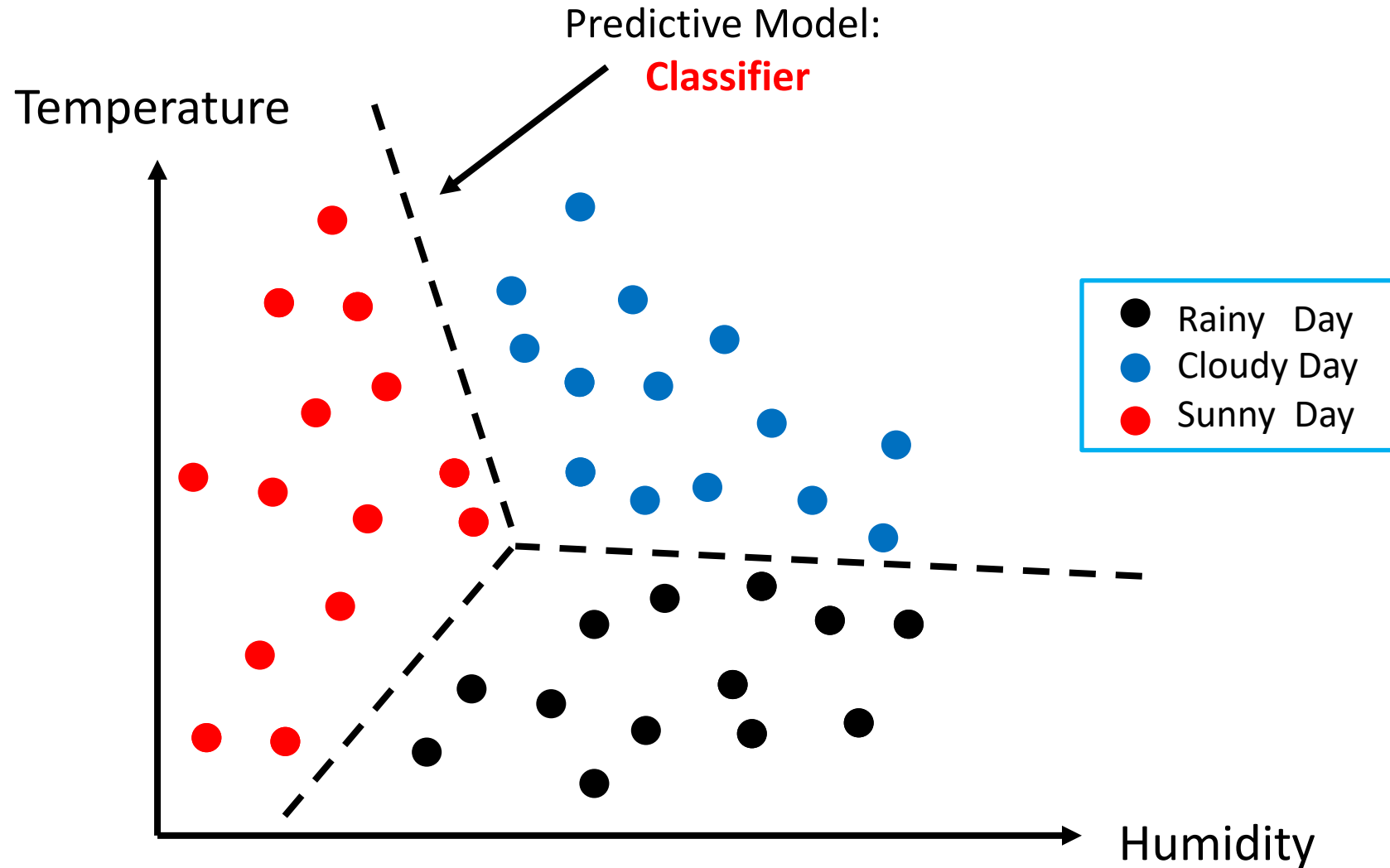# Supervised Learning:
# Learning from labeled Data

# Two Basic Approaches of <u>Supervised Learning</u>

- **Classification**: Predict a **<u>discrete</u>** <u>valued output</u> for each observation.
    - Labels are discrete (categorical)
    - Labels can be binary (e.g., rainy/sunny, spam/non-spam,) or non-binary (e.g., rainy/sunny/cloudy)

- **Regression:** Predict a **<u>continuous</u>** <u>valued output</u> for each observation.
    - Labels are continuous (numeric), e.g., stock price, housing price
    - Can define 'closeness' when comparing prediction with true values
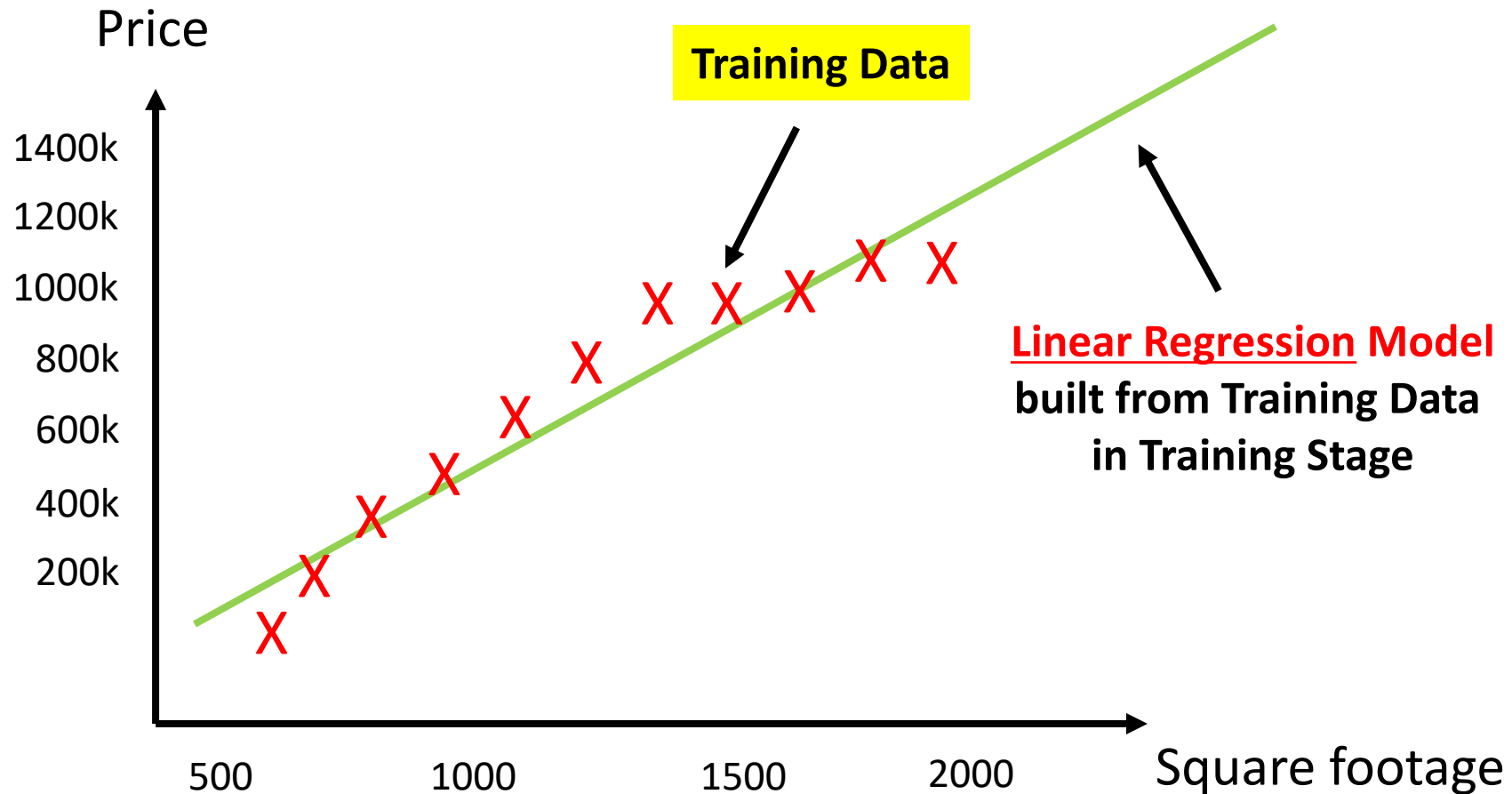
# Classification Example: Binary Label
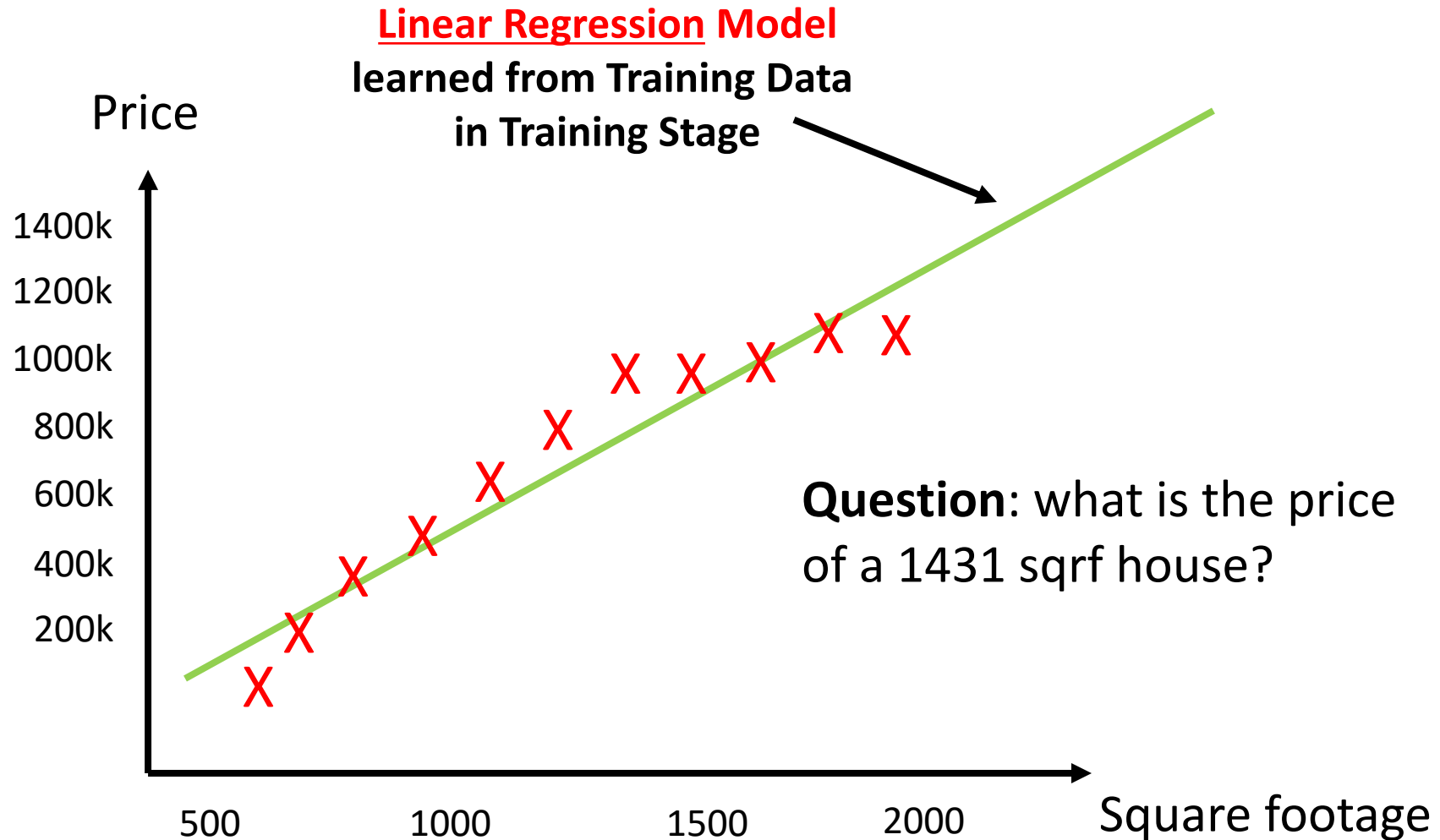
# Classification Example: Multiple Label

# Regression Example: Housing Price

# Regression Example: Housing Price

# Regression Example: Housing Price



**Linear Regression Model**
learned from Training Data
in Training Stage

Price

**Question**: what is the price of a 1431 sqrf house?

Square footage

# Regression Example: Housing Price

**Linear Regression Model**
**learned from Training Data**
**in Training Stage**

Price

1400k
1200k
1000k
800k
600k
400k
200k

**Question**: what is the price of a 1431 sqrf house?

1431

500    1000    1500    2000    Square footage

# Regression Example: Housing Price



Linear Regression Model
learned from Training Data
in Training Stage

Price

$848

Question: what is the price of a 1431 sqf house?

Answer: $848k

Square footage

# Regression Example: Housing Price



Quadratic Regression Model learned from Training Data in Training Stage

Price

Square footage

# Regression Example: Housing Price



**Quadratic Regression Model**
**learned from Training Data**
**in Training Stage**

Price

1400k
1200k
1000k
800k
600k
400k
200k

500        1000        1500        2000        Square footage

**Question**: what is the price of a 1431 sqf house?

**Answer**: $910k

# Regression Example: Housing Price



**In this case, quadratic model looks more accurate than linear model.**

Price

1400k
1200k
1000k
800k
600k
400k
200k

500    1000    1500    2000

Square footage

# Feature Table

- *Training dataset*: $\{(\textbf{\textit{x}}_1, y_1), (\textbf{\textit{x}}_2, y_2), \ldots, (\textbf{\textit{x}}_N, y_N)\}$ : $N$ data samples used for training.

| sepal length | sepal width | petal length | petal width | Label |
|---|---|---|---|---|
| 5.3 | 3.7 | 1.5 | 0.2 | setosa |
| 5 | 3 | 2 | 0.2 | setosa |
| 7.0 | 3.2 | 4.7 | 1.4 | versicolor |
| 6.4 | 3.2 | 4.5 | 1.5 | versicolor |
| 6.3 | 2.7 | 4.9 | 1.8 | virginica |
| 7.9 | 3.8 | 6.4 | 2 | virginica |

$\textbf{\textit{x}}_1$   $y_1$

$\textbf{\textit{x}}_2$   $y_2$

$\textbf{\textit{x}}_3$   $y_3$

- *Training dataset*: $\{(\pmb{x_1}, y_1), (\pmb{x_2}, y_2), \ldots, (\pmb{x_N}, y_N)\}$ with known label.
- Now, we have a new sample with unknown label: $(\pmb{x}, y=?)$

| sepal length | sepal width | petal length | petal width | Label |
|---|---|---|---|---|
| 5.3 | 3.7 | 1.5 | 0.2 | setosa |
| 5 | 3 | 2 | 0.2 | setosa |
| 7.0 | 3.2 | 4.7 | 1.4 | versicolor |
| 6.4 | 3.2 | 4.5 | 1.5 | versicolor |
| 6.3 | 2.7 | 4.9 | 1.8 | virginica |
| 7.9 | 3.8 | 6.4 | 2 | virginica |
| 7 | 3.9 | 5.9 | 1.3 | ??? |

$\pmb{x_1}$ → $y_1$

$\pmb{x_2}$ → $y_2$

: → :

$\pmb{x_N}$ → $y_N$

$\pmb{x}$ → $y=?$
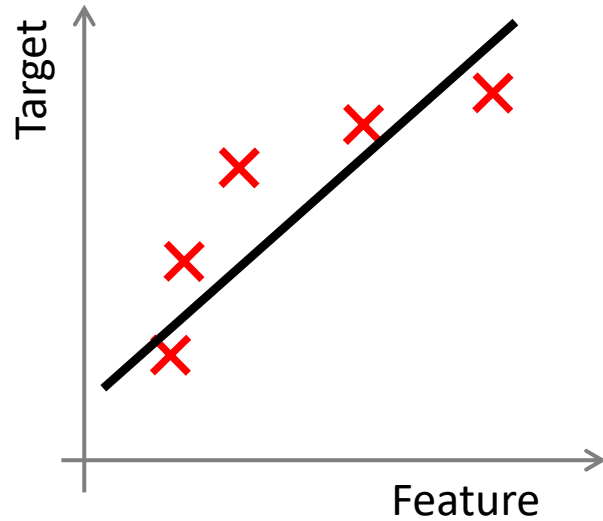
# The Problem of Overfitting

# The Problem of Overfitting

- **Overfitting** happens when the predictive model (classification model or regression model) **fits too much** with the **training samples** so that it starts capturing, learning, and representing the **noise and randomness or outlier samples** of the training dataset.


- Overfitting provides excellent accuracy for training data, but poor results for future data samples (testing set)!
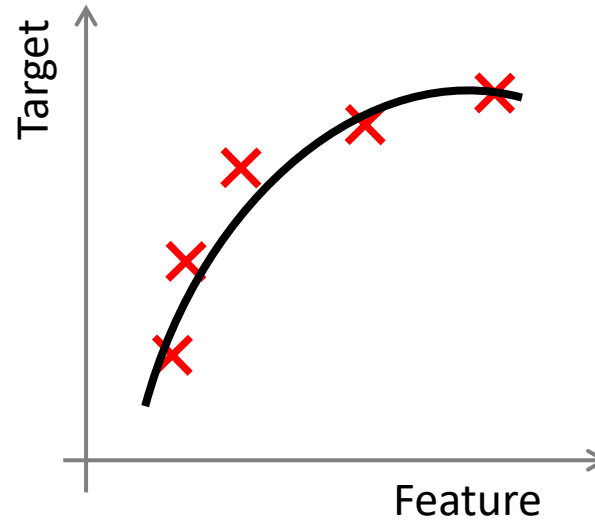
# The Problem of Overfitting

- **Overfitting** occurs when a model is excessively complex. The two main reasons that makes a model too complex are:

    1. having **too many features**.
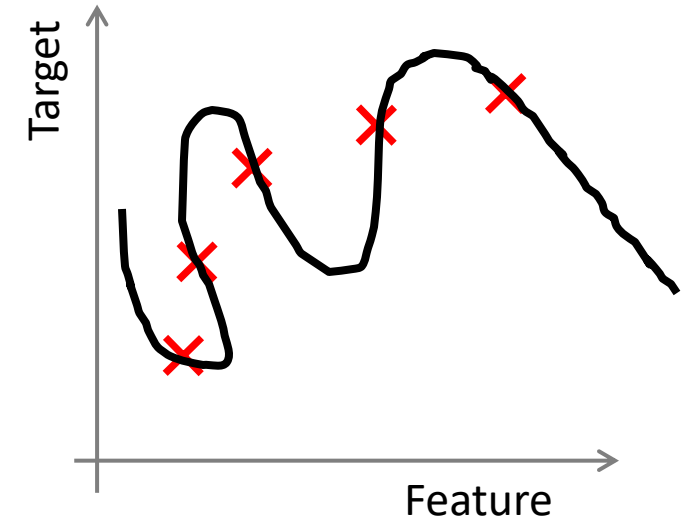    2. having a **complex model with very high order**.

# Example of Overfitting for Regression
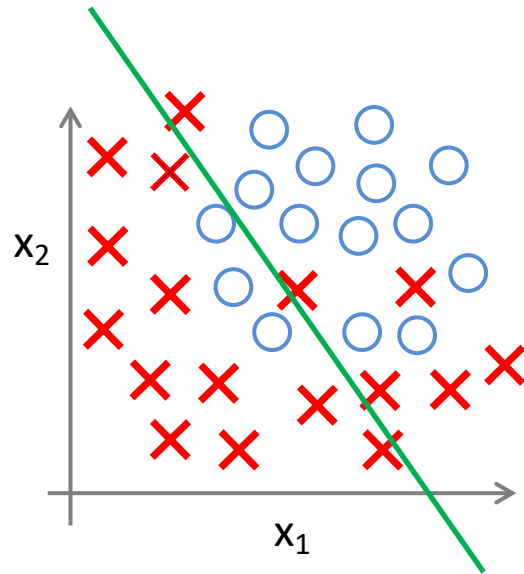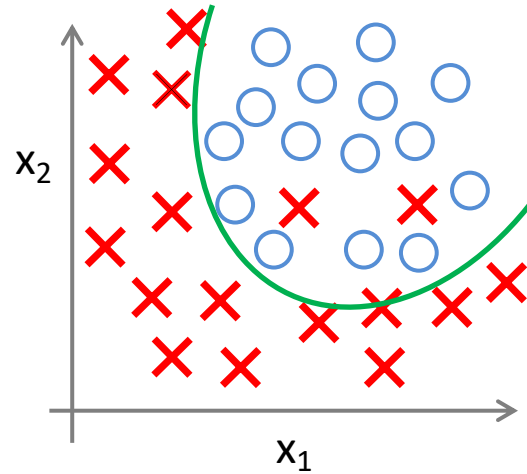


**Under-fit
(High Bias)**

**Ideal fit**

**Over-fit
(High Variance)**

# Example of Overfitting for Classification



Under-fit
(High Bias)

Ideal fit

Over-fit
(High Variance)

# Addressing the Overfitting Problem: Approach 1: Dimensionality Reduction

- **Approach 1**: <span style="color:red">**Dimensionality Reduction:**</span>

  - Reduce the number of features $x$ (e.g. rather using 20 features for prediction, use only the best 3 features)

$$\theta_0 + \theta_1 x_1 + \theta_2 x_2 + ... + \theta_{20} x_{20} \quad \rightarrow \quad \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3$$

  We can:

  a) <u>Manually</u> select which features to keep.

  b) Detecting the **best features** using <u>automated</u> **Feature Selection** and/or **Dimensionality Reduction** algorithms (will be covered later).

# Feature Selection

- **Feature selection** is an important step in machine learning. The classic feature selection algorithms usually focus on specific metrics to **quantify the relevance and/or redundancy of each feature** with the goal of finding **the smallest subset of features that provides the maximum amount of useful information** for prediction.

- Thus, the **main goal of feature selection algorithms is to eliminate redundant or irrelevant features** in a given feature set.

- Applying an effective feature selection algorithm not only decreases the complexity of the system by reducing the dimensionality, but also increases the performance of the classifier by avoiding overfitting and also removing irrelevant and confusing features.

# Addressing the Overfitting Problem: Approach 2: Regularization

- **Approach 2**: <span style="color:red">**Regularization**</span>:
  - Keep all features, but reduce the magnitude/values of parameters of the model ($\theta_j$) to simplify the model.

$$\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \theta_4 x_1^2 + \theta_5 x_2^2 + \theta_6 x_2 x_3 + \theta_6 x_2 x_3^2 + \ldots$$

$$\rightarrow \quad \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \theta_4 x_1^2 + \theta_6 x_2 x_3$$

# Thank You!

**Questions?**