IOWA STATE UNIVERSITY

Department of Computer Science

COM S 573: Machine Learning

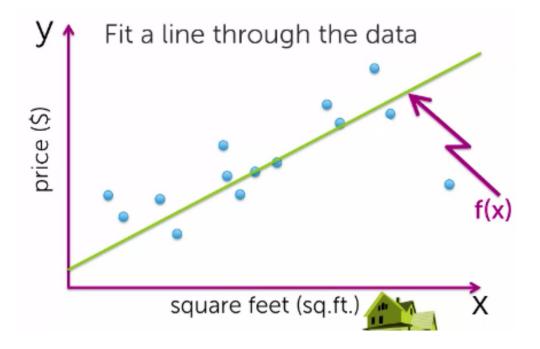
Lecture 3: Linear Regression

Linear Models

- Linear regression: predict a scalar
 - House price
 - Weight of a planet
- Linear perceptron: classifier
 - Predict an animal is a dog or not
 - Predict an image contains a square or not
- Logistic regression: classifier based on a probability
 - Predict how likely a team win
 - Predict how likely tomorrow is sunny

Linear Regression

Predict a scalar based on input features



What does "LINEAR" mean?

Linear Regression: Intuitions

- Given an input like house
- We extract some features from it:
 - For example: [size, #bedrooms, #floors, ...]

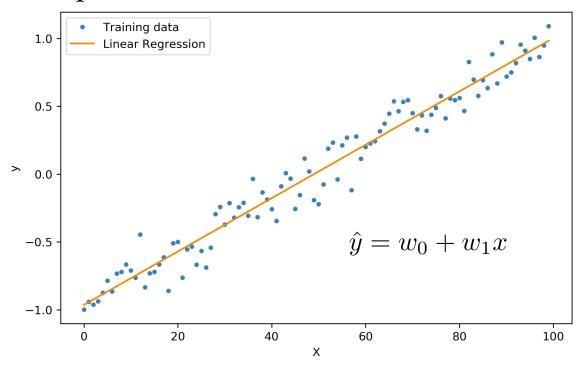
$$\boldsymbol{x} = [x_0, x_1, \cdots, x_d]^T \in \mathbb{R}^d$$

• We want to get an aggregation result by giving these features different weights.

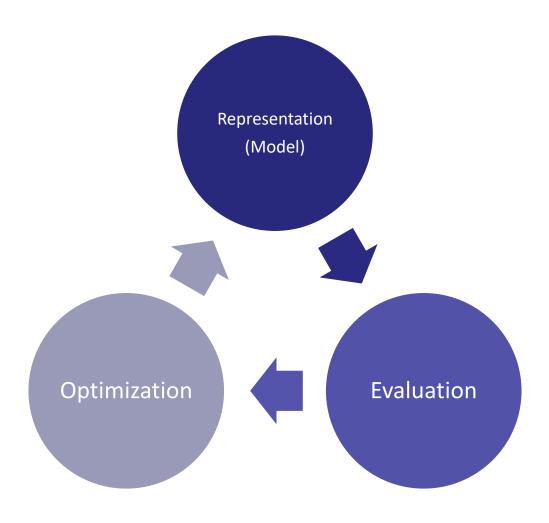
$$y = 1 + 3 * size + 4 * \#bedrooms + 5 * \#floors + \dots$$

Linear Regression

 Linear regression is a linear approach to modeling the relationship between a scalar response and one or more independent variables



Three Components of Learning



Linear Regression

- Given a training example $< m{x}, y>$
- Each \boldsymbol{x} has d features x_1, \dots, x_d
- The prediction is computed

$$\hat{y}_i = w_0 + w_1 x_{i,1} + w_2 x_{i,2} + \dots + w_d x_{i,d}$$

Linear Regression: Example

- Predict house price
 - x = [size, distance]
 - y = price

| Training Sample | Size (sq.ft.) | Distance (miles) | Price (\$) |
|-----------------|------------------|---------------------|---------------|
| 1 | 498 | 10 | 600 |
| 2 | 267 | 9 | 455 |
| 3 | 399 | 7.8 | 546 |
| | | ••• | |

Linear Regression: Representation

• For each training sample $< \boldsymbol{x}_i, y_i >$

•
$$\hat{y}_i = w_0 + w_1 x_{i,1} + w_2 x_{i,2} + \dots + w_d x_{i,d}$$

- Suppose $\boldsymbol{w} = [w_0, w_1, \cdots, w_d]^T$
- $\hat{y}_i = \mathbf{w}^T \mathbf{x}_i$

What is in \boldsymbol{x}_i ?

Linear Regression: Representation

• For each training sample $< \boldsymbol{x}_i, y_i >$

•
$$\hat{y}_i = w_0 + w_1 x_{i,1} + w_2 x_{i,2} + \dots + w_d x_{i,d}$$

- Suppose $\boldsymbol{w} = [w_0, w_1, \cdots, w_d]^T$
- $\hat{y}_i = \mathbf{w}^T \mathbf{x}_i$

What is in
$$\boldsymbol{x_i}$$
 ? $\boldsymbol{x_i} = [1, x_1, \cdots, x_d]^T$

Linear Regression: Representation

• For all training sample $\langle x_1, y_1 \rangle, \cdots, \langle x_n, y_n \rangle$

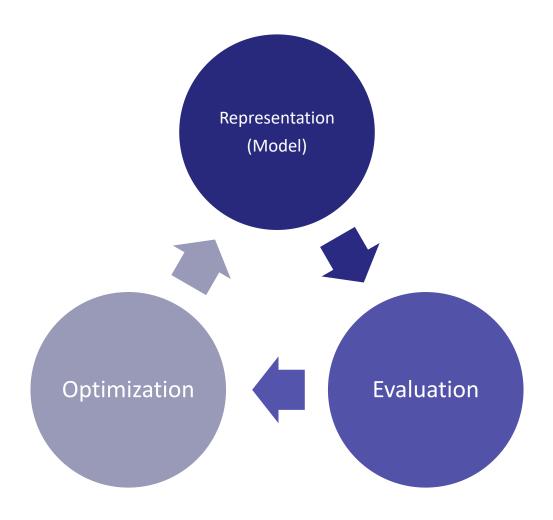
•
$$\hat{y}_i = w_0 + w_1 x_{i,1} + w_2 x_{i,2} + \dots + w_d x_{i,d}$$

$$egin{aligned} oldsymbol{\cdot} \begin{bmatrix} \hat{y_1} \\ \hat{y_2} \\ \vdots \\ \hat{y_n} \end{bmatrix} = \begin{bmatrix} 1 & x_{1,1} & x_{1,2} & \cdots & x_{1,d} \\ 1 & x_{2,1} & x_{2,2} & \cdots & x_{2,d} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n,1} & x_{n,2} & \cdots & x_{n,d} \end{bmatrix} oldsymbol{w}$$

•
$$\hat{y} = Xw$$

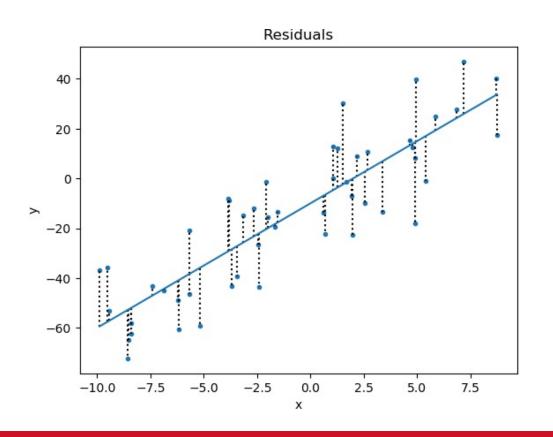
What's the dimension of X?

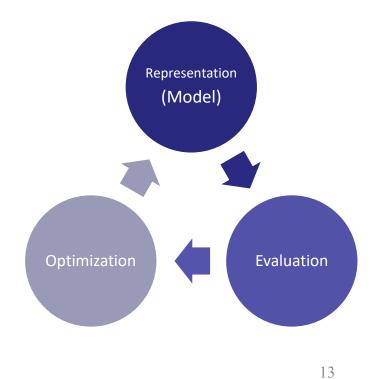
Three Components of Learning



Linear Regression: Evaluation

• Residual Squares $(y_i - {m w}_i^T {m x}_i)^2$





Linear Regression: Evaluation

Residual Sum of Squares (RSS)

$$RSS(\boldsymbol{w}) = \sum_{i=1}^{n} (y_i - \hat{y_i})^2$$

Equivalently

$$RSS(\boldsymbol{w}) = (\boldsymbol{y} - \boldsymbol{X}\boldsymbol{w})^T(\boldsymbol{y} - \boldsymbol{X}\boldsymbol{w})$$

Why?

•
$$RSS(\boldsymbol{w}) = (\boldsymbol{y} - \boldsymbol{X}\boldsymbol{w})^T(\boldsymbol{y} - \boldsymbol{X}\boldsymbol{w})$$

• Find the minimal $RSS(\boldsymbol{w})$

How to find it?

•
$$RSS(\boldsymbol{w}) = (\boldsymbol{y} - \boldsymbol{X}\boldsymbol{w})^T(\boldsymbol{y} - \boldsymbol{X}\boldsymbol{w})$$

- Find the minimal $RSS(\boldsymbol{w})$
- When the first derivative of a function equals zero, the minimum of a function is achieved.

$$\frac{\partial RSS(\boldsymbol{w})}{\partial \boldsymbol{w}} = 0$$

• The optimal $oldsymbol{w}$ is obtained by solving this equation.

•
$$RSS(\boldsymbol{w}) = (\boldsymbol{y} - \boldsymbol{X}\boldsymbol{w})^T(\boldsymbol{y} - \boldsymbol{X}\boldsymbol{w})$$

$$= \boldsymbol{y}^T \boldsymbol{y} - \underline{\boldsymbol{w}^T \boldsymbol{X}^T \boldsymbol{y} - \boldsymbol{y}^T \boldsymbol{X} \boldsymbol{w}} + \boldsymbol{w}^T \boldsymbol{X}^T \boldsymbol{X} \boldsymbol{w}$$

$$= \boldsymbol{y}^T \boldsymbol{y} - 2\boldsymbol{w}^T \boldsymbol{X}^T \boldsymbol{y} + \boldsymbol{w}^T \boldsymbol{X}^T \boldsymbol{X} \boldsymbol{w}$$

•
$$RSS(\boldsymbol{w}) = (\boldsymbol{y} - \boldsymbol{X}\boldsymbol{w})^T(\boldsymbol{y} - \boldsymbol{X}\boldsymbol{w})$$

 $= \boldsymbol{y}^T\boldsymbol{y} - \boldsymbol{w}^T\boldsymbol{X}^T\boldsymbol{y} - \boldsymbol{y}^T\boldsymbol{X}\boldsymbol{w} + \boldsymbol{w}^T\boldsymbol{X}^T\boldsymbol{X}\boldsymbol{w}$
 $= \boldsymbol{y}^T\boldsymbol{y} - 2\boldsymbol{w}^T\boldsymbol{X}^T\boldsymbol{y} + \boldsymbol{w}^T\boldsymbol{X}^T\boldsymbol{X}\boldsymbol{w}$

•
$$\frac{\partial RSS(\boldsymbol{w})}{\partial \boldsymbol{w}} = ?$$

•
$$RSS(\boldsymbol{w}) = (\boldsymbol{y} - \boldsymbol{X}\boldsymbol{w})^T(\boldsymbol{y} - \boldsymbol{X}\boldsymbol{w})$$

$$= \boldsymbol{y}^T \boldsymbol{y} - \boldsymbol{w}^T \boldsymbol{X}^T \boldsymbol{y} - \boldsymbol{y}^T \boldsymbol{X} \boldsymbol{w} + \boldsymbol{w}^T \boldsymbol{X}^T \boldsymbol{X} \boldsymbol{w}$$

$$= \boldsymbol{y}^T \boldsymbol{y} - 2\boldsymbol{w}^T \boldsymbol{X}^T \boldsymbol{y} + \boldsymbol{w}^T \boldsymbol{X}^T \boldsymbol{X} \boldsymbol{w}$$
• $\frac{\partial RSS(\boldsymbol{w})}{\partial \boldsymbol{w}} = ?$

•
$$RSS(\boldsymbol{w}) = (\boldsymbol{y} - \boldsymbol{X}\boldsymbol{w})^T(\boldsymbol{y} - \boldsymbol{X}\boldsymbol{w})$$

$$= \boldsymbol{y}^T \boldsymbol{y} - \boldsymbol{w}^T \boldsymbol{X}^T \boldsymbol{y} - \boldsymbol{y}^T \boldsymbol{X} \boldsymbol{w} + \boldsymbol{w}^T \boldsymbol{X}^T \boldsymbol{X} \boldsymbol{w}$$

$$= \boldsymbol{y}^T \boldsymbol{y} - 2\boldsymbol{w}^T \boldsymbol{X}^T \boldsymbol{y} + \boldsymbol{w}^T \boldsymbol{X}^T \boldsymbol{X} \boldsymbol{w}$$
• $\frac{\partial RSS(\boldsymbol{w})}{\partial \boldsymbol{w}} = ?$

$$0$$
?

•
$$RSS(\boldsymbol{w}) = (\boldsymbol{y} - \boldsymbol{X}\boldsymbol{w})^T(\boldsymbol{y} - \boldsymbol{X}\boldsymbol{w})$$

$$= \boldsymbol{y}^T \boldsymbol{y} - \boldsymbol{w}^T \boldsymbol{X}^T \boldsymbol{y} - \boldsymbol{y}^T \boldsymbol{X} \boldsymbol{w} + \boldsymbol{w}^T \boldsymbol{X}^T \boldsymbol{X} \boldsymbol{w}$$

$$= \boldsymbol{y}^T \boldsymbol{y} - 2\boldsymbol{w}^T \boldsymbol{X}^T \boldsymbol{y} + \boldsymbol{w}^T \boldsymbol{X}^T \boldsymbol{X} \boldsymbol{w}$$
• $\frac{\partial RSS(\boldsymbol{w})}{\partial \boldsymbol{w}} = ?$

$$\boxed{-2\boldsymbol{X}^T \boldsymbol{y}}$$
?

•
$$RSS(\boldsymbol{w}) = (\boldsymbol{y} - \boldsymbol{X}\boldsymbol{w})^T(\boldsymbol{y} - \boldsymbol{X}\boldsymbol{w})$$

$$= \boldsymbol{y}^T\boldsymbol{y} - \boldsymbol{w}^T\boldsymbol{X}^T\boldsymbol{y} - \boldsymbol{y}^T\boldsymbol{X}\boldsymbol{w} + \boldsymbol{w}^T\boldsymbol{X}^T\boldsymbol{X}\boldsymbol{w}$$

$$= \boldsymbol{y}^T\boldsymbol{y} - 2\boldsymbol{w}^T\boldsymbol{X}^T\boldsymbol{y} + \boldsymbol{w}^T\boldsymbol{X}^T\boldsymbol{X}\boldsymbol{w}$$
• $\frac{\partial RSS(\boldsymbol{w})}{\partial \boldsymbol{w}} = ?$

$$\boxed{-2\boldsymbol{X}^T\boldsymbol{y}} \qquad 2\boldsymbol{X}^T\boldsymbol{X}\boldsymbol{w}$$

•
$$RSS(\boldsymbol{w}) = (\boldsymbol{y} - \boldsymbol{X}\boldsymbol{w})^T(\boldsymbol{y} - \boldsymbol{X}\boldsymbol{w})$$

 $= \boldsymbol{y}^T\boldsymbol{y} - \boldsymbol{w}^T\boldsymbol{X}^T\boldsymbol{y} - \boldsymbol{y}^T\boldsymbol{X}\boldsymbol{w} + \boldsymbol{w}^T\boldsymbol{X}^T\boldsymbol{X}\boldsymbol{w}$
 $= \boldsymbol{y}^T\boldsymbol{y} - 2\boldsymbol{w}^T\boldsymbol{X}^T\boldsymbol{y} + \boldsymbol{w}^T\boldsymbol{X}^T\boldsymbol{X}\boldsymbol{w}$

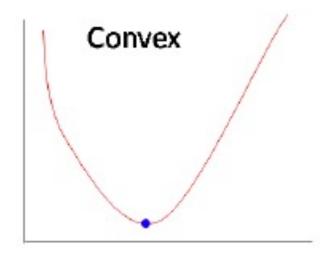
•
$$\frac{\partial RSS(\boldsymbol{w})}{\partial \boldsymbol{w}} = -2\boldsymbol{X^T}\boldsymbol{y} + 2\boldsymbol{X^T}\boldsymbol{X}\boldsymbol{w} = 0$$

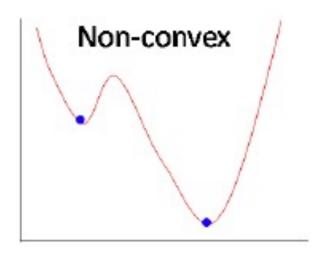
•
$$w^* = (X^T X)^{-1} X^T y$$

Any protentional issue?

• w^* are global optima?

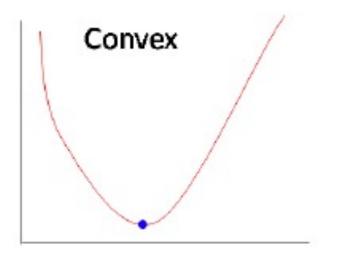
- $oldsymbol{w}^*$ are global optima? Yes
- RSS(w) is a convex function

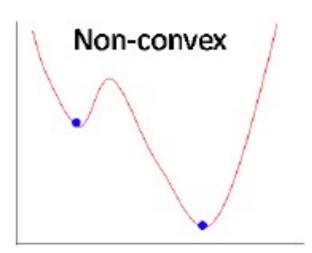




• What is convex?

- w^* are global optima? Yes
- RSS(w) is a convex function

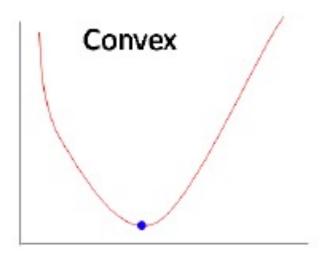


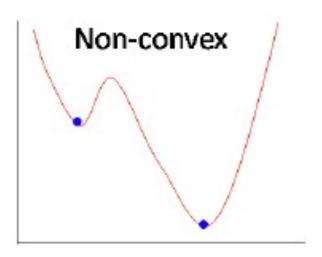


- What is convex?
- Convex is a property that a line joining any two points on its graph lies on or above the graph.

26

- w^* are global optima? Yes
- RSS(w) is a convex function





- What is convex?
- How to prove a function is convex?

- w^* are global optima? Yes
- RSS(w) is a convex function

$$H(\boldsymbol{w}) = \frac{\partial^2 RSS(\boldsymbol{w})}{\partial \boldsymbol{w}^2}$$

• For every $oldsymbol{u} \in \mathbb{R}^d$, we have

$$\boldsymbol{u}^T H(\boldsymbol{w}) \boldsymbol{u} >= 0$$

Linear Regression: Summary

Representation

$$\hat{y_i} = oldsymbol{w} oldsymbol{x}_i^T$$

Predict a continuous scalar.

Evaluation

$$RSS(\boldsymbol{w}) = (\boldsymbol{y} - \boldsymbol{X}\boldsymbol{w})^T(\boldsymbol{y} - \boldsymbol{X}\boldsymbol{w})^T$$

Optimization

$$\frac{\partial RSS(\boldsymbol{w})}{\partial \boldsymbol{w}} = 0 \to \boldsymbol{w}^* = (\boldsymbol{X}^T\boldsymbol{X})^{-1}\boldsymbol{X}^T\boldsymbol{y}$$

Linear Models: Next

- Linear regression: predict a scalar
 - House price
 - Weight of a planet
- Linear perceptron: classifier of discrete prediction
 - Predict an animal is a dog or not
 - Predict an image contains a square or not
- Logistic regression: classifier based on a probability
 - Predict how likely a team win
 - Predict how likely tomorrow is sunny