

DSC 40A

Theoretical Foundations of Data Science I

Last Time

- ▶ We used linear algebra to fit a prediction rule of the form

$$H(x) = w_0 + w_1x.$$

- ▶ Instead of using our formulas for w_0 and w_1 , we can find these parameters by solving the **normal equations**:

$$X^T X \vec{w} = X^T \vec{y}$$

where

$$X = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix}, \quad \vec{w} = \begin{bmatrix} w_0 \\ w_1 \end{bmatrix}, \quad \vec{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}.$$

In This Video

Can we change the form of the prediction rule to be nonlinear?

Recommended Reading

Course Notes: Chapter 2, Section 2

Review: Linear Algebra Textbook

The Hypothesis Vector

- ▶ The **hypothesis vector** is the vector $\vec{h} \in \mathbb{R}^n$ with components $H(x_i)$. This is the vector of predicted values.
- ▶ When our prediction rule is

$$H(x) = \underline{w_0 + w_1 x},$$

the hypothesis vector \vec{h} can be written

$$\vec{h} = \begin{bmatrix} H(x_1) \\ H(x_2) \\ \vdots \\ H(x_n) \end{bmatrix} = \begin{bmatrix} w_0 + w_1 x_1 \\ w_0 + w_1 x_2 \\ \vdots \\ w_0 + w_1 x_n \end{bmatrix} = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix} \begin{bmatrix} w_0 \\ w_1 \end{bmatrix}$$

(Handwritten blue annotations: a large 'X' is drawn over the matrix, and a blue arrow points from the matrix to the weight vector $\begin{bmatrix} w_0 \\ w_1 \end{bmatrix}$.)

The Hypothesis Vector

- ▶ The **hypothesis vector** is the vector $\vec{h} \in \mathbb{R}^n$ with components $H(x_i)$. This is the vector of predicted values.
- ▶ When our prediction rule is

$$H(x) = \underline{w_0 + w_1x + w_2x^2}$$

the hypothesis vector \vec{h} can be written

$$\vec{h} = \begin{bmatrix} H(x_1) \\ H(x_2) \\ \vdots \\ H(x_n) \end{bmatrix} = \begin{bmatrix} w_0 + w_1x_1 + w_2x_1^2 \\ w_0 + w_1x_2 + w_2x_2^2 \\ \vdots \\ w_0 + w_1x_n + w_2x_n^2 \end{bmatrix} = \begin{bmatrix} 1 & x_1 & x_1^2 \\ 1 & x_2 & x_2^2 \\ \vdots & \vdots & \vdots \\ 1 & x_n & x_n^2 \end{bmatrix} \begin{bmatrix} w_0 \\ w_1 \\ w_2 \end{bmatrix}$$

$\vec{h} = X \vec{w}$

The Hypothesis Vector

- ▶ The **hypothesis vector** is the vector $\vec{h} \in \mathbb{R}^n$ with components $H(x_i)$. This is the vector of predicted values.
- ▶ When our prediction rule is

$$\underline{H(x) = w_0 + w_1x + w_2x^2 + w_3x^3}$$

the hypothesis vector \vec{h} can be written

$$\vec{h} = \begin{bmatrix} H(x_1) \\ H(x_2) \\ \vdots \\ H(x_n) \end{bmatrix} = \begin{bmatrix} 1 & x_1 & x_1^2 & x_1^3 \\ 1 & x_2 & x_2^2 & x_2^3 \\ \vdots & \vdots & \vdots & \vdots \\ 1 & x_n & x_n^2 & x_n^3 \end{bmatrix} \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ w_3 \end{bmatrix}$$

\times \rightarrow W

The Hypothesis Vector

- ▶ The **hypothesis vector** is the vector $\vec{h} \in \mathbb{R}^n$ with components $H(x_i)$. This is the vector of predicted values.
- ▶ When our prediction rule is

$$H(x) = w_1 \frac{1}{x^2} + w_2 \sin x + w_3 e^x$$

the hypothesis vector \vec{h} can be written

$$\vec{h} = \begin{bmatrix} H(x_1) \\ H(x_2) \\ \vdots \\ H(x_n) \end{bmatrix} = \begin{bmatrix} \frac{1}{x_1^2} & \sin x_1 & e^{x_1} \\ \vdots & \vdots & \vdots \\ \frac{1}{x_n^2} & \sin x_n & e^{x_n} \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix}$$

$\vec{h} = X \vec{w}$

Minimizing the Mean Squared Error

- ▶ As long as the form of the prediction rule permits us to write $\underline{h} = X\vec{w}$ for some X and \vec{w} , the mean squared error is

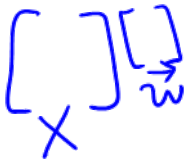
$$R_{\text{sq}}(\vec{w}) = \frac{1}{n} \|\vec{y} - X\vec{w}\|^2.$$

- ▶ Regardless of the values of X and \vec{w} ,

$$\begin{aligned} \frac{dR_{\text{sq}}}{d\vec{w}} &= 0 \\ \implies -2X^T\vec{y} + 2X^T X\vec{w} &= 0 \\ \implies X^T X\vec{w} &= X^T\vec{y}. \end{aligned}$$

- ▶ Solving the **normal equations** still works!

Linear in the Parameters



- ▶ We can fit rules like:

$$\underline{w_0} + \underline{w_1}x + \underline{w_2}x^2 \quad \underline{w_1}e^{-x^2} + \underline{w_2} \cos(x + \pi) + \underline{w_3} \frac{\log 2x}{x}$$

- ▶ This includes arbitrary polynomials.
- ▶ We can't fit rules like:

$$e^{\underline{w_1}x} + w_0 \quad \sin(\underline{w_1}x + \underline{w_0})$$

- ▶ We can have any number of parameters, but must be **linear in the parameters**.

Determining Function Form

- ▶ How do we know what form our prediction rule should take?
- ▶ Sometimes, we know from *theory*, using knowledge about what the variables represent and how they should be related.
- ▶ Other times, we make a guess based on the data.
- ▶ Generally, start with simpler functions first.

Question

Suppose you collect data on the height, or position, of a freefalling object at various times t_i . Which form should your prediction rule take to best fit the data?

a) constant, $H(t) = w_0$

b) linear, $H(t) = w_0 + w_1 t$

c) quadratic, $H(t) = w_0 + w_1 t + w_2 t^2$

d) no way to know without plotting the data

$$a(t) = -9.8$$

$$v(t) = -9.8t + C_0$$

$$s(t) = -\frac{9.8t^2}{2} + C_0 t + C_1$$

theory

Example: Amdahl's Law

- ▶ Amdahl's Law relates the runtime of a program on p processors to the time to do the sequential and nonsequential parts on one processor.

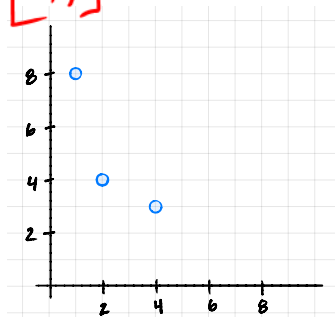
$$H(p) = \frac{t_{NS}}{p} + t_S$$

- ▶ Collect data by timing a program with varying numbers of processor:

Processors	Time (Hours)
1	8
2	4
4	3

$$\vec{w} = \begin{bmatrix} 1 \\ 4/7 \end{bmatrix}$$

Example: fitting $H(x) = \underbrace{w_1}_{t_{NS}} \cdot \frac{1}{x_i} + \underbrace{w_0}_{t_S}$ theory



$$\begin{bmatrix} H(x_1) \\ \vdots \\ H(x_n) \end{bmatrix} = \begin{bmatrix} 1 & \frac{1}{x_1} \\ \vdots & \vdots \\ 1 & \frac{1}{x_n} \end{bmatrix} \begin{bmatrix} w_0 \\ w_1 \end{bmatrix}$$

\vec{w}

normal eqns $X^T X \vec{w} = X^T y$

$$\begin{bmatrix} 1 & 1 & 1 \\ 1 & 1/2 & 1/4 \end{bmatrix} \begin{bmatrix} 1 & 1/2 \\ 1 & 1/4 \end{bmatrix} \begin{bmatrix} w_0 \\ w_1 \end{bmatrix} =$$

$$(X^T X) \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1/2 & 1/4 \end{bmatrix} \begin{bmatrix} 8 \\ 4 \\ 3 \end{bmatrix}$$

x_i	y_i
1	8
2	4
4	3

$$y = \begin{bmatrix} 8 \\ 4 \\ 3 \end{bmatrix}$$

Example: Amdahl's Law

- ▶ We found: $t_{NS} = \frac{48}{7} \approx 6.88$, $t_S = 1$
- ▶ Therefore our prediction rule is:

$$\begin{aligned} H(p) &= \frac{t_{NS}}{p} + t_S \\ &= \frac{6.88}{p} + 1 \end{aligned}$$

Demo

Summary

- ▶ Whenever our prediction rule is **linear in the parameters**, we can use the **normal equations**

$$X^T X \vec{w} = X^T \vec{y}$$

to find the parameters that minimize the mean squared error.

- ▶ This means we can use regression to fit prediction rules of many forms, including arbitrary polynomials.
- ▶ **Next time:** We'll use the normal equations to make predictions based on multiple features.