

$$\begin{aligned}
 R_{\text{sq}}(\vec{w}) &= \|X\vec{w} - \vec{y}\|^2 \\
 \nabla_{\vec{w}} R_{\text{sq}}(\vec{w}) &= \frac{d}{d\vec{w}} R_{\text{sq}}(\vec{w}) \\
 &= 2X^T X \vec{w} - 2X^T \vec{y} \\
 (X^T X) \vec{w} &= X^T \vec{y}
 \end{aligned}$$

DSC 40A

Lecture 09
 Least Squares Regression, pt. IV

Last Time

- ▶ How do we make predictions using multiple features?
- ▶ Assume a linear decision rule:

$$H(\text{experience, GPA, \# internships}) =$$

$$w_0 + w_1 \times (\text{experience}) + w_2 \times (\text{GPA}) + w_3 \times (\#\text{ of internships})$$

- ▶ In general:

$$H(x_1, \dots, x_d) = w_0 + w_1 x_1 + w_2 x_2 + \dots + w_d x_d$$

Feature Vectors

- Nicer to pack into a **feature vector** and **parameter vector**:

$$\vec{x} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{pmatrix} \quad \vec{w} = \begin{pmatrix} w_1 \\ w_2 \\ \vdots \\ w_d \end{pmatrix}$$

- Then: $H(\vec{x}) = w_0 + \vec{w} \cdot \vec{x}$

Feature Vectors

- Nicer to pack into a **feature vector** and **parameter vector**:

$$\vec{x} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{pmatrix} \quad \vec{w} = \begin{pmatrix} w_1 \\ w_2 \\ \vdots \\ w_d \end{pmatrix}$$

- Then: $H(\vec{x}) = w_0 + \vec{w} \cdot \vec{x}$
- Actually, we should include w_0 in \vec{w} ...

Augmented Feature Vectors

- The **augmented feature vector** $\text{Aug}(\vec{x})$ is the vector obtained by adding a 1 to the front of \vec{x} :

$$\vec{x} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{pmatrix} \quad \text{Aug}(\vec{x}) = \begin{pmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_d \end{pmatrix} \quad \vec{w} = \begin{pmatrix} w_0 \\ w_1 \\ w_2 \\ \vdots \\ w_d \end{pmatrix}$$

- Then:

$$\begin{aligned} H(x_1, \dots, x_d) &= w_0 + w_1 x_1 + w_2 x_2 + \dots + w_d x_d \\ &= \text{Aug}(\vec{x}) \cdot \vec{w} \end{aligned}$$

Last Time

- ▶ We want to fit a decision rule of the form $H(\vec{x}) = \text{Aug}(\vec{x}) \cdot \vec{w}$.
- ▶ Minimize **mean squared error**:

$$R_{\text{sq}}(\vec{w}) = \frac{1}{n} \sum_{i=1}^n \left[(\vec{w} \cdot \text{Aug}(\vec{x}^{(i)})) - y_i \right]^2$$

Rewriting the Mean Squared Error

- ▶ Define the **design matrix**:

$$X = \begin{pmatrix} \text{Aug}(\vec{x}^{(1)}) \\ \text{Aug}(\vec{x}^{(2)}) \\ \vdots \\ \text{Aug}(\vec{x}^{(n)}) \end{pmatrix} = \begin{pmatrix} 1 & x_1^{(1)} & x_2^{(1)} & \dots & x_d^{(1)} \\ 1 & x_1^{(2)} & x_2^{(2)} & \dots & x_d^{(2)} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(n)} & x_2^{(n)} & \dots & x_d^{(n)} \end{pmatrix}$$

← 1 + # of features →

people

- ▶ And the vector of **observations**: $\vec{y} = (y_1, \dots, y_n)^T$

Rewriting the Mean Squared Error

- ▶ Then:

$$\begin{aligned} R_{\text{sq}}(\vec{w}) &= \frac{1}{n} \sum_{i=1}^n \left[(\vec{w} \cdot \text{Aug}(\vec{x}^{(i)})) - y_i \right]^2 \\ &= \frac{1}{n} \|X\vec{w} - \vec{y}\|^2 \end{aligned}$$

- ▶ Today's goal: find the \vec{w} that minimizes the MSE.

Minimizing the Mean Squared Error

- ▶ Our goal: minimize the function:

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

- ▶ Strategy:

1. Take partial derivatives,

$$\frac{\partial R_{\text{sq}}}{\partial w_0}(\vec{w}), \quad \frac{\partial R_{\text{sq}}}{\partial w_1}(\vec{w}), \quad \frac{\partial R_{\text{sq}}}{\partial w_2}(\vec{w}), \quad \dots \quad \frac{\partial R_{\text{sq}}}{\partial w_d}(\vec{w})$$

2. Set each equal to zero and solve for w_0, w_1, \dots, w_d .

Minimizing the MSE: Gradient Edition

- The vector of partial derivatives is called the **gradient**:

$$\left(\frac{\partial R_{\text{sq}}}{\partial w_0}(\vec{w}), \quad \frac{\partial R_{\text{sq}}}{\partial w_1}(\vec{w}), \quad \frac{\partial R_{\text{sq}}}{\partial w_2}(\vec{w}), \quad \dots, \quad \frac{\partial R_{\text{sq}}}{\partial w_d}(\vec{w}) \right)^T$$

- Written: $\nabla_{\vec{w}} R_{\text{sq}}(\vec{w})$ or $\frac{dR_{\text{sq}}}{d\vec{w}}(\vec{w})$
- Strategy:
 - Compute the gradient of $R_{\text{sq}}(\vec{w})$.
 - Set it to zero and solve for \vec{w} .

Gradients Review

Computing Gradients

When computing $\frac{df}{d\vec{x}}(\vec{x})$:

- ▶ Before: make sure that f takes in vectors, outputs scalars.

▶ **Example:** $\frac{d}{d\vec{x}} [A\vec{x}]$

▶ **Example:** $\frac{d}{d\vec{x}} [\vec{x} \cdot \vec{x}], \frac{d}{d\vec{x}} [\vec{x}^T A^T A \vec{x}]$

$$f(\vec{x}) = \vec{x} \cdot \vec{x}$$

- ▶ After: make sure your result is a vector.

$$(A\vec{x})^T = \vec{x}^T A^T$$

Finding the Gradient: Strategy #1 $\nabla \vec{a} \cdot \vec{x}$

Example: Find $\frac{d}{d\vec{x}} [\vec{a} \cdot \vec{x}]$ where \vec{x} and \vec{a} have d elements.

1. "Unpack" all matrix multiplications/dot products

$$\triangleright \vec{a} \cdot \vec{x} = a_1 x_1 + a_2 x_2 + \dots + a_d x_d$$

$$\vec{x} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{pmatrix} \quad \vec{a} = \begin{pmatrix} a_1 \\ a_2 \\ \vdots \\ a_d \end{pmatrix}$$

Finding the Gradient: Strategy #1

Example: Find $\frac{d}{d\vec{x}} [\vec{a} \cdot \vec{x}]$ where \vec{x} and \vec{a} have d elements.

1. "Unpack" all matrix multiplications/dot products

$$\blacktriangleright \vec{a} \cdot \vec{x} = a_1x_1 + a_2x_2 + \dots + a_dx_d = \sum_{i=1}^d a_i x_i$$

2. Take partial derivatives (perhaps with arbitrary index):

$$\frac{\partial}{\partial x_1} [a_1x_1 + a_2x_2 + \dots + a_dx_d] = a_1 \quad \frac{d}{dx_i} [a_1x_1 + a_2x_2 + \dots + a_dx_d] = a_i$$

$$\frac{\partial}{\partial x_2} [a_1x_1 + a_2x_2 + \dots + a_dx_d] = a_2 \quad = a_i$$

⋮

$$\frac{\partial}{\partial x_d} [a_1x_1 + a_2x_2 + \dots + a_dx_d] = a_d$$

Finding the Gradient: Strategy #1

3. Pack partial derivatives into a gradient vector:

$$\frac{d}{d\vec{x}} [\vec{a} \cdot \vec{x}] = (a_1, a_2, \dots, a_d)^T$$

4. Simplify:

$$(a_1, a_2, \dots, a_d)^T = \vec{a}$$

- ▶ So $\frac{d}{d\vec{x}} [\vec{a} \cdot \vec{x}] = \vec{a}$
- ▶ Check: **result is a vector.**

$$\frac{d}{dx} ax = a$$

Finding the Gradient: Strategy #1

- ▶ **Pro:** Always works, straightforward
- ▶ **Con:** Unpacking everything can get messy

Example

Show that $\frac{d}{d\vec{x}} [\vec{x}^T A^T A \vec{x}] = 2A^T A \vec{x}$, where A is $n \times d$ and \vec{x} is $n \times 1$.

- ▶ Check: **it is a scalar**

1. After unpacking: $\vec{x}^T A^T A \vec{x} = \sum_{i=1}^n \left(\sum_{j=1}^d A_{ij} x_j \right)^2$
2. Take partial derivatives:

$$\frac{\partial}{\partial x_1} \left[\sum_{i=1}^n \left(\sum_{j=1}^d A_{ij} x_j \right)^2 \right] = \sum_{i=1}^n \sum_{j=1}^d A_{i1} A_{ij} x_j$$

Example

3. Pack into a gradient vector:

$$\frac{d}{d\vec{x}} [\vec{x}^T A^T A \vec{x}] = \begin{pmatrix} \sum_{i=1}^n \sum_{j=1}^d A_{i1} A_{ij} x_j \\ \sum_{i=1}^n \sum_{j=1}^d A_{i2} A_{ij} x_j \\ \vdots \\ \sum_{i=1}^n \sum_{j=1}^d A_{id} A_{ij} x_j \end{pmatrix}$$

4. Somehow simplify this to $A^T A \vec{x}$...

Finding the Gradient: Strategy #2

Chain Rule: If $f : \mathbb{R} \rightarrow \mathbb{R}$, and $g : \mathbb{R}^d \rightarrow \mathbb{R}$, then:

$$\frac{d}{d\vec{x}} f(g(\vec{x})) = \frac{df}{dg} \frac{dg}{d\vec{x}}$$

Example: What is $\frac{d}{d\vec{x}} [(\vec{a} \cdot \vec{x})^2]?$

- ▶ $f(g) = g^2$
- ▶ $g(\vec{x}) = \vec{a} \cdot \vec{x}$
- ▶
$$\begin{aligned}\frac{d}{d\vec{x}} [(\vec{a} \cdot \vec{x})^2] &= 2g(\vec{x}) \vec{a} \\ &= 2(\vec{a} \cdot \vec{x}) \vec{a}\end{aligned}$$

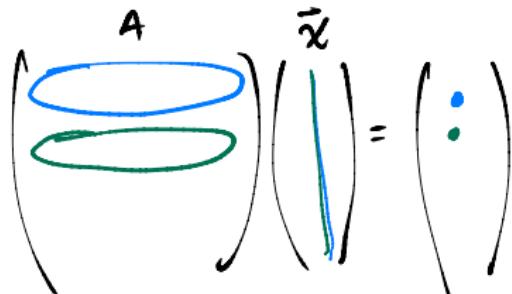
Finding the Gradient: Strategy #2

1. Unpack until we can use chain rule, but no more.
2. Use the chain rule.
3. Simplify.

Recall

Suppose A is $n \times d$.

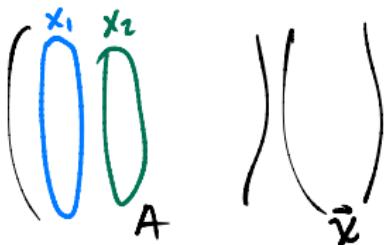
Let \vec{A}_{i*} denotes its i th row. Then:



$$A\vec{x} = \begin{pmatrix} \vec{A}_{1*} \cdot \vec{x} \\ \vec{A}_{2*} \cdot \vec{x} \\ \vdots \\ \vec{A}_{n*} \cdot \vec{x} \end{pmatrix}$$

Let \vec{A}_{*j} denotes its j th column, then:

$$A\vec{x} = \vec{A}_{*1}x_1 + \vec{A}_{*2}x_2 + \dots + \vec{A}_{*d}x_d$$



Finding the Gradient: Strategy #2

$$(A\vec{x})^T = \vec{x}^T A^T$$

Show that $\frac{d}{d\vec{x}} [\vec{x}^T A^T A \vec{x}] = 2A^T A \vec{x}$, where A is $n \times d$ and \vec{x} is $n \times 1$.

1. Unpack $\vec{x}^T A^T A \vec{x} = \vec{x}^T A^T \begin{pmatrix} \vec{A}_{1*} \cdot \vec{x} \\ \vec{A}_{2*} \cdot \vec{x} \\ \vdots \\ \vec{A}_{n*} \cdot \vec{x} \end{pmatrix}$

$$= (\vec{A}_{1*} \cdot \vec{x} \quad \vec{A}_{2*} \cdot \vec{x} \quad \dots \quad \vec{A}_{n*} \cdot \vec{x}) \begin{pmatrix} A_{1*} \cdot \vec{x} \\ A_{2*} \cdot \vec{x} \\ \vdots \\ A_{n*} \cdot \vec{x} \end{pmatrix}$$
$$= (A_{1*} \cdot \vec{x})^2 + (A_{2*} \cdot \vec{x})^2 + \dots + (A_{n*} \cdot \vec{x})^2$$

Finding the Gradient: Strategy #2

Show that $\frac{d}{d\vec{x}} [\vec{x}^T A^T A \vec{x}] = 2A^T A \vec{x}$, where A is $n \times d$ and \vec{x} is $n \times 1$.

2. Use chain rule:

$$\frac{d}{d\vec{x}} \left[(A_{1*} \cdot \vec{x})^2 + \dots + (A_{n*} \cdot \vec{x})^2 \right] = \frac{d}{d\vec{x}} (A_{1*} \cdot \vec{x})^2 + \dots + \frac{d}{d\vec{x}} (A_{n*} \cdot \vec{x})^2$$

Since $\frac{d}{d\vec{x}} (\vec{a} \cdot \vec{x})^2 = 2(\vec{a} \cdot \vec{x})\vec{a}$ (from the chain rule):

$$= 2(A_{1*} \cdot \vec{x})A_{1*} + \dots + 2(A_{n*} \cdot \vec{x})A_{n*}$$

Finding the Gradient: Strategy #2

Show that $\frac{d}{d\vec{x}} [\vec{x}^T A^T A \vec{x}] = 2A^T A \vec{x}$, where A is $n \times d$ and \vec{x} is $n \times 1$.

3. Show that this $= 2A^T A \vec{x}$.

$$\frac{d}{d\vec{x}} [\vec{x}^T A^T A \vec{x}] = 2(A_{1*} \cdot \vec{x}) A_{1*} + \dots + 2(A_{n*} \cdot \vec{x}) A_{n*}$$

Starting with $2A^T A \vec{x}$.

$$2A^T A \vec{x} = 2A^T \begin{pmatrix} A_{1*} \cdot \vec{x} \\ A_{2*} \cdot \vec{x} \\ \vdots \\ A_{n*} \cdot \vec{x} \end{pmatrix} = 2 \left[\begin{array}{l} \text{first col of } A^T \cdot A_{1*} \cdot \vec{x} \\ + \\ \vdots \\ + \\ n^{\text{th}} \text{ col of } A^T \cdot A_{n*} \cdot \vec{x} \end{array} \right]$$

Since the i^{th} col of $A^T = i^{\text{th}}$ row of A , i.e., A_{i*}

$$= 2 [A_{1*} A_{1*} \cdot \vec{x} + \dots + A_{n*} A_{n*} \cdot \vec{x}]$$

Back to Regression...

Minimizing the MSE

- We want to compute:

$$\frac{d}{d\vec{w}} \left[R_{\text{sq}}(\vec{w}) \right] = \frac{d}{d\vec{w}} \left[\|X\vec{w} - \vec{y}\|^2 \right]$$

- Step 1: Rewrite squared norm using dot product. Recall:

$$(A + B)^T = A^T + B^T$$

$$(AB)^T = B^T A^T$$

$$\vec{u} \cdot \vec{v} = \vec{v} \cdot \vec{u}$$

$$(\vec{u} + \vec{v}) \cdot (\vec{w} + \vec{z}) = \vec{u} \cdot \vec{w} + \vec{u} \cdot \vec{x} + \vec{v} \cdot \vec{w} + \vec{v} \cdot \vec{z}$$

$$\|\vec{u}\|^2 = \vec{u} \cdot \vec{u}$$

Step 1: Rewriting squared norm

$$\begin{aligned}\|X\vec{w} - \vec{y}\|^2 &= (\vec{X}\vec{w} - \vec{y})^\top (\vec{X}\vec{w} - \vec{y}) = ((\vec{X}\vec{w})^\top - \vec{y}^\top)(-\vec{w}) \\ &= (\vec{w}^\top \vec{X}^\top - \vec{y}^\top)(\vec{X}\vec{w} - \vec{y}) \\ &= \vec{w}^\top \vec{X}^\top \vec{X}\vec{w} - \vec{w}^\top \vec{X}^\top \vec{y} - \vec{y}^\top \vec{X}\vec{w} + \vec{y}^\top \vec{y} \\ &= \vec{w}^\top \vec{X}^\top \vec{X}\vec{w} - \vec{y}^\top \vec{X}\vec{w} + \vec{y}^\top \vec{y}\end{aligned}$$

Step 2: Take gradients

$$\begin{aligned}\frac{d}{d\vec{w}} [R_{\text{sq}}(\vec{w})] &= \frac{d}{d\vec{w}} [\vec{w}^T X^T X \vec{w} - 2\vec{y}^T X \vec{w} + \vec{y}^T \vec{y}] \\ &= \frac{d}{d\vec{w}} [\vec{w}^T X^T X \vec{w}] - 2 \frac{d}{d\vec{w}} [\vec{y}^T X \vec{w}] + \frac{d}{d\vec{w}} [\vec{y}^T \vec{y}] \\ &= 2X^T X \vec{w} - 2X^T \vec{y} + 0\end{aligned}$$

The Normal Equations

- ▶ To minimize $R_{\text{sq}}(\vec{w})$, set gradient to zero, solve for \vec{w} :

$$2X^T X \vec{w} - 2X^T \vec{y} = 0 \implies X^T X \vec{w} = X^T \vec{y}$$

- ▶ This is a system of equations in matrix form, called the **normal equations**.
- ▶ Solution¹: $\vec{w} = (X^T X)^{-1} X^T \vec{y}$.

¹Don't actually compute inverse! Use Gaussian elimination.

Regression with Multiple Features

- ▶ We want to find \vec{w} which minimizes $\|X\vec{w} - \vec{y}\|^2$.
- ▶ The answer: $\vec{w} = (X^T X)^{-1} X^T \vec{y}$.

Example: $\vec{x}^{(1)} = 2$ $\vec{x}^{(2)} = 5$ $\vec{x}^{(3)} = 7$ $\vec{x}^{(4)} = 8$
 $y_1 = 1$ $y_2 = 2$ $y_3 = 3$ $y_4 = 3$

Design matrix

$$X = \begin{pmatrix} 1 & 2 \\ 1 & 5 \\ 1 & 7 \\ 1 & 8 \end{pmatrix} \quad \vec{y} = \begin{pmatrix} 1 \\ 2 \\ 3 \\ 3 \end{pmatrix}$$

Compute $\vec{w} = (X^T X)^{-1} X^T \vec{y}$

$$X^T \vec{y} = \begin{pmatrix} 1 & 1 & 1 & 1 \\ 2 & 5 & 7 & 8 \end{pmatrix} \begin{pmatrix} 1 \\ 2 \\ 3 \\ 3 \end{pmatrix} = \begin{pmatrix} 9 \\ 57 \end{pmatrix}$$

$$X^T X = \begin{pmatrix} 1 & 1 & 1 & 1 \\ 2 & 5 & 7 & 8 \end{pmatrix} \begin{pmatrix} 1 & 2 \\ 1 & 5 \\ 1 & 7 \\ 1 & 8 \end{pmatrix} = \begin{pmatrix} 4 & 25 \\ 25 & 142 \end{pmatrix}$$

$\vec{w} = \boxed{\begin{pmatrix} 4 & 25 \\ 25 & 142 \end{pmatrix}^{-1} \begin{pmatrix} 9 \\ 57 \end{pmatrix}}$