

Machine Learning: HW1

Brent Garey

February 3rd, 2025

Contents

1	1.	2
2	2.	3
3	3.	6
4	4.	9
5	5.	11

1. 1.

Show transformation of a vector by an orthonormal matrix (U) preserves the norm (magnitude) of the vector.

answer:

$$\text{Using class notes: } \|\theta\|_2 = \sqrt{\theta^T \theta} = \sqrt{\sum \theta_i^2} = \sqrt{\langle \theta, \theta \rangle} = \sqrt{\theta^T \theta}.$$

$$\sqrt{\theta^T \theta} = \|\theta\|_2 \rightarrow \text{square, so. } (\sqrt{\theta^T \theta})^2 = \|\theta\|_2^2$$

$$\rightarrow \theta^T \theta = \|\theta\|_2^2 \rightarrow \text{inner product equivalent.}$$

$$\text{so, } \|U\theta\|_2^2 = (U\theta)^T U\theta \rightarrow \theta^T U^T U \theta$$

• transpose of product

$$\theta^T \cancel{\sum} \theta$$

$$\begin{matrix} (\theta^T I_n) & & I_n \theta \\ \downarrow & & \downarrow \\ \theta^T & & \theta \end{matrix}$$

• Since Identity Matrix does not transform when doing multiplication

$$\text{so, } \rightarrow \theta^T \theta = \langle \theta, \theta \rangle = \|\theta\|_2^2 = \|\theta\|_2$$

2 2.

$$\text{Q. } 1* \left(\frac{\partial b^T X \theta}{\partial \theta} \right) = X^T b$$

$$1. \frac{\partial b^T X \theta}{\partial \theta} = \frac{\partial (b^T)(X\theta)}{\partial \theta} = \frac{\partial (X\theta)^T(b)}{\partial \theta}$$

$$\text{Since in class: } \frac{\partial (a^T x)}{\partial x} = \frac{\partial (x^T a)}{\partial x}$$

$$\text{so, } \frac{\partial (X\theta)^T(b)}{\partial \theta} = \frac{\partial \theta^T X^T(b)}{\partial \theta}$$

↳ multiplicative transpose property.

$$\text{in class: } \frac{\partial (x^T a)}{\partial x} = a, \text{ so }$$

$$\frac{\partial \theta^T X^T b}{\partial \theta} = \underline{\underline{X^T b}}$$

$$\frac{\partial \|\mathbf{X}\theta - \mathbf{b}\|_2^2}{\partial \theta} = 2\mathbf{X}^T(\mathbf{X}\theta - \mathbf{b})$$

set $\ell(\theta) = \|\mathbf{X}\theta - \mathbf{b}\|_2^2 = \left\| \begin{pmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_n^T \end{pmatrix} \begin{pmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{pmatrix} - \begin{pmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{pmatrix} \right\|_2^2$

transpose matrix \mathbf{X} to get all features of a single \mathbf{x}_n

$$\left\| \begin{pmatrix} \mathbf{x}_1^T \theta \\ \mathbf{x}_2^T \theta \\ \vdots \\ \mathbf{x}_n^T \theta \end{pmatrix} - \begin{pmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{pmatrix} \right\|_2^2 = \left\| \begin{pmatrix} \theta^T \mathbf{x}_1 \\ \theta^T \mathbf{x}_2 \\ \vdots \\ \theta^T \mathbf{x}_n \end{pmatrix} - \begin{pmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{pmatrix} \right\|_2^2$$

Since $\frac{\partial (\mathbf{x}^T \theta)}{\partial \theta} = \frac{\partial (\theta^T \mathbf{x})}{\partial \theta}$, transpose property rearranges order.

→ we now have vector operations, simplify:

$$\left\| \begin{pmatrix} \theta^T \mathbf{x}_1 - b_1 \\ \theta^T \mathbf{x}_2 - b_2 \\ \vdots \\ \theta^T \mathbf{x}_n - b_n \end{pmatrix} \right\|_2^2 \rightarrow \text{L2 norm definition "sum of squares"} \sum_{i=1}^N (\theta^T \mathbf{x}_i - b_i)^2$$

$$\frac{\partial}{\partial \theta} \sum_{i=1}^N (\theta^T x_i - b_i)^2;$$

↑
take derivative Using Chain rule.

$$\rightarrow n(f(x))^{n-1} \cdot \underbrace{f'(x)}_{\frac{\partial}{\partial \theta}}$$

Plug in: $\underbrace{2(\theta^T x_i - b_i) \cdot x_i}_{\text{since } x_i \text{ is features of index } i \text{ of } x, \text{ we can once again claim } x_i = X^T}$

Simplify: $2x_i(\theta^T x_i - b) = 2X^T(\theta^T X^T - b)$

$\xrightarrow{\text{transpose multiplicative property.}}$

$$= 2X^T((\theta X)^T - b) = 2X^T(X\theta - b)$$

3 3.

3. Gradient Descent: $f(x_1, x_2) = (x_1 - x_2)^2 + (x_1 - 1)^2 + (x_2 - 1)^2$

partial Derivatives: $\begin{pmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \end{pmatrix} = \begin{pmatrix} 4x_1 - 2x_2 - 2 \\ -2x_1 + 4x_2 - 2 \end{pmatrix}$

1. Use learning rate $\rho = 1$.

iteration 1: $f(x_1, x_2) = (0, 0) = \theta^{(0)}$

$$\begin{pmatrix} 4(0) - 2(0) - 2 \\ -2(0) + 4(0) - 2 \end{pmatrix} = \begin{pmatrix} -2 \\ -2 \end{pmatrix} \rightarrow \text{gradient} \rightarrow \text{not } 0, \text{ so has NOT converged.}$$

↳ update position using step size & gradient:

$$(\theta^{(0)} - \rho(-2), \theta^{(0)} - \rho(-2)) \rightarrow (0 - 1(-2), 0 - 1(-2))$$

$= (2, 2)$ → updated position after 1 iteration.

iteration 2: $f(x_1, x_2) = (2, 2) = \theta^{(1)}$

$$\begin{pmatrix} 4(2) - 2(2) - 2 \\ -2(2) + 4(2) - 2 \end{pmatrix} = \begin{pmatrix} 2 \\ 2 \end{pmatrix} \rightarrow \text{gradient} \rightarrow \text{not converged!}$$

↳ update: $(\theta^{(1)} - 1(2), \theta^{(1)} - 1(2)) = (0, 0) = \underline{\text{updated position}}$

iteration 3: $f(x_1, x_2) = (0, 0) = \theta^{(2)}$

$$\begin{pmatrix} 4(0) - 2(0) - 2 \\ -2(0) + 4(0) - 2 \end{pmatrix} = \begin{pmatrix} -2 \\ -2 \end{pmatrix} \rightarrow \text{gradient} \rightarrow \text{Not converged!}$$

update: $(\theta^{(2)} - 1(-2), \theta^{(2)} - 1(-2)) = (2, 2)$

3 pt 2: $\rho = 0.5$,

iteration 1: $(x_1, x_2) = (0, 0) = \theta^0$.

$$\begin{pmatrix} 4(0) - 2(0) - 2 \\ -2(0) + 4(0) - 2 \end{pmatrix} = \begin{pmatrix} -2 \\ -2 \end{pmatrix} \rightarrow \text{not converged.}$$

update position: $(0 - 0.5(-2), 0 - 0.5(-2)) = \underline{(1, 1)}$

iteration 2: $(x_1, x_2) = (1, 1) = \theta^1$

$$\begin{pmatrix} 4(1) - 2(1) - 2 \\ -2(1) + 4(1) - 2 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix} \rightarrow \text{converged! (make no adjustment)}$$

update position yields: $(1 - 0.5(0), 1 - 0.5(0)) = \underline{(1, 1)}$

iteration 3: $(x_1, x_2) = (1, 1) = \theta^2$

$$\rightarrow \text{gradient} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

update: $(1 - 0.5(0), 1 - 0.5(0)) = \underline{(1, 1)}$.

3 pt 3: $\rho = 2$

iteration 1: $(x_1, x_2) = (0, 0) = \theta^0$

$$\begin{pmatrix} 4(0) - 2(0) - 2 \\ -2(0) + 4(0) - 2 \end{pmatrix} = \begin{pmatrix} -2 \\ -2 \end{pmatrix} \rightarrow \text{not converged.}$$

update position: $(0 - 2(-2), 0 - 2(-2)) = \underline{(4, 4)} = \theta^1$

iteration 2: $(x_1, x_2) = (4, 4) = \theta^2$

$$\begin{pmatrix} 4(4) - 2(4) - 2 \\ -2(4) + 4(4) - 2 \end{pmatrix} = \begin{pmatrix} 6 \\ 6 \end{pmatrix} \rightarrow \text{not converged.}$$

update position: $(4 - 2(6), 4 - 2(6)) = \underline{(-8, -8)} = \theta^2$

iteration 3: $(x_1, x_2) = (-8, -8) = \theta^3$

$$\begin{pmatrix} 4(-8) - 2(-8) - 2 \\ -2(-8) + 4(-8) - 2 \end{pmatrix} = \begin{pmatrix} -18 \\ -18 \end{pmatrix} \rightarrow \text{not converged.}$$

update position: $(-8 - 2(-18), -8 - 2(-18)) = \underline{(28, 28)} = \theta^3$

Analysis:

1. Learning rate = 1:

Gradient descent has NOT converged in 3 iterations. With this step size, GD has reached a local minima where it will constantly swap between (-2,-2) and (2,2) every other iteration, NEVER converging. The step size is slightly too large, where the ideal convergence criteria lays somewhere between the step size = 1 and 0 (0.5) since 1 step in either direction steps OVER the optimal convergence criteria of gradient = 0. The next iteration will always take the same step back OVER, resulting in a back and forth.

2. Learning rate = 0.5:

GD has converged in 2 iterations because the gradient becomes (0,0) on iteration 2. I wrote down iteration 3 just in case, but from iteration 3 and beyond the gradient will always evaluate to (0,0), so GD will never move the point from (1,1).

3. Learning rate = 2:

GD has NOT converged in 3 iterations. The step size is too large, so we achieve an oscillating pattern that will forever increase the absolute values of the gradient, with every other iteration changing signs from positive to negative. This will never converge.

4. 4.

1. Closed form solution:

4. 1. Closed Form Solution

$$\sum_{i=1,3,5,\dots} (y_i - \theta^T x_i)^2 + \sum_{i=2,4,6,\dots} 2(y_i - \theta^T x_i)^2$$

$$\|y - \theta^T x\|_2^2 + \|2(y - \theta^T x)\|_2^2$$

$$\underbrace{\cancel{2}(y - \theta^T x)^T (y - \theta^T x)}_{\text{odds}} + \underbrace{\cancel{2}(y - \theta^T x)^T (y - \theta^T x)}_{\text{evens}}$$

Similar to Q1: multiply by a (weighted) matrix where the diagonals correspond to the weight of the index:

$$\text{so, } \rightarrow \|w(y - \theta^T x)\|_2^2 \text{ where } w = \begin{cases} \frac{1}{\sqrt{2}} & \text{on odd diagonals} \\ \sqrt{2} & \text{on even diagonals} \\ 0 & \text{elsewhere...} \end{cases}$$

$$\rightarrow \left\| \begin{pmatrix} (1)(y_1 - \theta^T x_1) \\ (\sqrt{2})(y_2 - \theta^T x_2) \\ (1)(y_3 - \theta^T x_3) \\ (\sqrt{2})(y_4 - \theta^T x_4) \end{pmatrix} \right\|_2^2$$

$$\begin{matrix} 1 & 0 & 0 & 0 \\ 0 & \sqrt{2} & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & \sqrt{2} \end{matrix}$$

$$w \cdot y = y'$$

$$\text{so, } \sum_{i=1}^n w_i (y_i - \theta^T x_i)^2$$

$$\theta^* = (X^T(WX))^{-1} X^T W y \text{ (on the bottom of next page)}$$

Now, $\|(\omega(y - \theta^T x))\|_2^2 \xrightarrow{\text{w} = y, x' \text{ to simplify notation}} \|y - \theta^T x'\|_2^2$

$$\frac{\partial}{\partial \theta} \rightarrow 2x'^T (y' - x'^T \theta)$$

$$\rightarrow 2x'^T y' - 2x'^T \underbrace{x'^T \theta}_{x'} = 2x'^T y' - 2x'^T \theta$$

$$\rightarrow x'^T y' - x'^T \theta = 0 \quad \cdot \text{ multiply by } x'^T x$$

$$\rightarrow x'^T y' = x'^T \theta$$

$$\hookrightarrow \theta^* = (x'^T x')^{-1} x'^T y'$$

in terms of ω, x, y

$$\rightarrow \theta^* = (X^T (W X))^{-1} X^T W Y$$

2. Steps for gradient descent:

Step 1: Initialize θ^0 either randomly, or as 0's.

Step 2: Determine if index is odd, or even (if $\% 2 == 0$) to determine weight of gradient in next step.

Step 3: Calculate gradient ($\nabla_{\theta} = 2X^T W(X\theta - y)$):

Check for convergence! (gradient either 0 or some arbitrarily small number)

Step 4: Use gradient and some arbitrary step size to update position (find θ^{t+1})
 $\theta^{t+1} = \theta^t - \alpha \nabla_{\theta}(\theta)$

5 5.

Q5: x_i , sample i with d features $\begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$

• y vector becomes covariance measures.
↳ multivariate regression $\rightarrow \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \begin{bmatrix} \theta_1 \\ \vdots \\ \theta_n \end{bmatrix}$

• extra sum?
↳ similar to Frobenius Norm?

where $\|X\|_F = \left(\sum_{i=1}^m \sum_{j=1}^n x_{ij}^2 \right)^{1/2}$

now, matrices

$$\frac{\partial \|X\theta - Y\|_F^2}{\partial \theta} = \text{tr}((X\theta - Y)(X\theta - Y)^T) =$$
$$= \text{tr}(B^T X^T X\theta - Y^T X\theta - \theta^T X^T Y + Y^T Y)$$
$$\frac{\partial}{\partial \theta} = 2X^T X\theta - 2X^T Y = \underline{2X^T(X\theta - Y)}$$

likewise: $\sum_{i=1}^n (y_i - \theta_0 - \sum_{j=1}^d \theta_j x_{ij})^2 \rightarrow \underline{\theta^* = (X^T X)^{-1} X^T Y}$