

1 Back to Basics: Linear Algebra

Let $X \in \mathbb{R}^{n \times m}$. We introduce some important terms and notation

$$w \leftarrow Xv \quad v \in \mathbb{R}^m, \quad w \in \mathbb{R}^n$$

The **Columnspace**, also called the range, or span, of X is $\text{Range}(X) := \{y \mid y = Xv\}$.

The **Rowspace** is $\text{Row}(X) := \{y \mid y = X^T v\}$.

The **Nullspace**, or Kernel, of X is defined as $\mathcal{N}(X) := \{v \mid Xv = 0\}$.

The **Orthogonal Complement** of a subspace, U , is a subspace, U^\perp such that $u \in U, u' \in U^\perp \implies \langle u, u' \rangle = 0$

For this problem We do not assume that X has full rank.

$$A = \begin{bmatrix} 1 & 2 & 0 \\ 0 & 1 & 0 \end{bmatrix}$$

$$\text{Col}(A) = \mathbb{R}^2$$

$$\text{Row}(A) = \text{Col}(A^T) = \text{Span} \left\{ \begin{bmatrix} 1 \\ 2 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \right\}$$

(a) Check the following facts:

(i) The $\text{Row}(X) = \text{Range}(X^T)$
Rows of X are cols of X^T .

(ii) The $\mathcal{N}(X)^\perp = \text{Row}(X)$

$v \in \mathcal{N}(X) \implies Xv = 0 \implies X_i^T v = 0 \forall i \leq n$
 $\therefore v \perp$ all vectors in $\text{Row}(X) \implies \text{Row}(X) = \mathcal{N}(X)^\perp$

(iii) $\mathcal{N}(X^T X) = \mathcal{N}(X)$ Hint: if $v \in \mathcal{N}(X^T X)$, then $v^T X^T X v = 0$.

$v^T X^T X v = 0 \implies (Xv)^T (Xv) = 0 \implies \|Xv\|_2^2 = 0 \implies Xv = 0 \implies v \in \mathcal{N}(X)$

$v \in \mathcal{N}(X) \implies v \in \mathcal{N}(X^T X) \because X^T X v = X^T 0 = 0 \implies v \in \mathcal{N}(X)$

(iv) $\text{Row}(X^T X) = \text{Range}(X^T X) = \text{Row}(X)$ Hint: Use the relationship between nullspace and rowspace. $X^T X$ is sym $\because (X^T X)^T = X^T (X^T)^T = X^T X$

$\therefore \text{Col}(X^T X) = \text{Row}(X^T X)^T = \text{Row}(X^T X)$

From (iii) $\mathcal{N}(X)^\perp = \text{Row}(X)$ $\mathcal{N}(X^T X)^\perp = \mathcal{N}(X)^\perp \implies \text{Row}(X^T X) = \text{Row}(X)$

(b) We now prove an important result of linear algebra, the Rank-Nullity theorem. Let $\text{Rank}(X) = \dim(\text{Range}(X)) = \dim(\text{Row}(X))$ and $\text{Nullity}(X) = \dim(\mathcal{N}(X))$. The Rank nullity theorem says that for $X \in \mathbb{R}^{n \times m}$ we have

$$\text{Rank}(X) + \text{Nullity}(X) = m$$

Use the above results to prove this theorem. Hint: The complementary subspace theorem says that for a vector space V and subspace U , we can always find a complementary subspace U^\perp such that $U + U^\perp = V$

2 Probability Review

There are n archers all shooting at the same target (bulls-eye) of radius 1. Let the score for a particular archer be defined to be the distance away from the center (the lower the score, the better, and 0 is the optimal score). Each archer's score is independent of the others, and is distributed uniformly between 0 and 1. What is the expected value of the worst (highest) score?

(a) Define a random variable Z that corresponds with the worst (highest) score.

(b) Derive the Cumulative Distribution Function (CDF) of Z .

(c) Let X be a non-negative random variable. The Tail-Sum formula states that

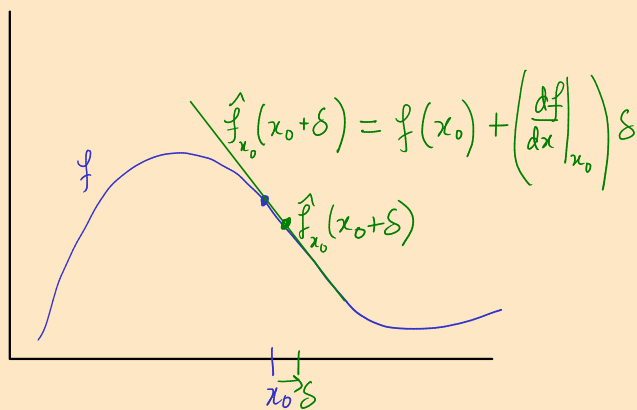
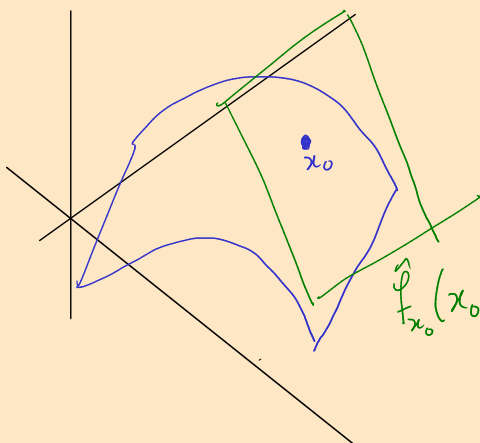
$$\mathbb{E}[X] = \int_0^{\infty} \mathbb{P}(X \geq t) dt$$

Using both the Tail-Sum formula and the CDF of Z derived above, calculate the expected value of Z *Hint: Write $\mathbb{P}(X \geq t)$ in terms of the CDF of X*

(d) Consider what happens to $\mathbb{E}[Z]$ as $n \rightarrow \infty$. Does this match your intuition?

$$f: \mathbb{R} \rightarrow \mathbb{R}$$

$$f: \mathbb{R}^n \rightarrow \mathbb{R}$$



$$\hat{f}_{x_0}(x_0 + \Delta) = f(x_0) + \underbrace{\left[\frac{\partial f}{\partial x} \right]}_{\text{deriv}} \Delta$$

$$\begin{bmatrix} \quad \end{bmatrix} \begin{bmatrix} \quad \end{bmatrix}$$

$$\nabla_x f \in \mathbb{R}^n$$

$$\frac{\partial f}{\partial x} : \mathbb{R}^n \rightarrow \mathbb{R}, \text{ linear}$$

$$f = x^2 + y^2$$

$$\frac{\partial f}{\partial x} = 2x$$

$$f: \mathbb{R}^n \rightarrow \mathbb{R}^m$$

$$\underbrace{\left[-\frac{\partial f}{\partial x_i} \right]}_n \text{ "derivative"}$$

$$\hat{f}_{x_0}(x_0 + \Delta) = f(x_0) + \left[\frac{\partial f}{\partial x} \right] \Delta$$

$$\begin{bmatrix} \equiv \end{bmatrix} \begin{bmatrix} \quad \end{bmatrix}$$

$$\begin{cases} \mathbb{R}^n \rightarrow \mathbb{R}^m \\ \text{linear} \\ \in \mathbb{R}^{m \times n} \end{cases}$$

$$\nabla_x f = \left(\frac{\partial f}{\partial x} \right)^T$$

$$j \rightarrow \begin{bmatrix} \text{"Jacobian"} \\ \frac{\partial f_j}{\partial x_i} \end{bmatrix}$$

$$f: \mathbb{R}^{m \times n} \rightarrow \mathbb{R}$$

$$\hat{f}_{x_0}(x_0 + \Delta) = f(x_0) + \underbrace{\text{trace}}_{n \times m} \left(\underbrace{\frac{\partial f}{\partial x}}_{A, B \in \mathbb{R}^{m \times n}} \Delta \right)$$

$$\sum_{ij} A_{ij} B_{ij}$$

$$\begin{aligned} \left\langle \begin{bmatrix} 1 & 2 \\ 0 & 1 \end{bmatrix}, \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \right\rangle_F &= 1+0+0+1=2 \\ &= \text{tr} \left(\begin{bmatrix} 1 & 2 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \right) \\ &= \text{tr} \left(\begin{bmatrix} 1 & 2 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \right) = 1+2+0=3 \end{aligned}$$

3 Vector Calculus

1

Below, $\mathbf{x} \in \mathbb{R}^d$ means that \mathbf{x} is a $d \times 1$ column vector with real-valued entries. Likewise, $\mathbf{A} \in \mathbb{R}^{d \times d}$ means that \mathbf{A} is a $d \times d$ matrix with real-valued entries. In this course, we will by convention consider vectors to be column vectors.

Consider $\mathbf{x}, \mathbf{w} \in \mathbb{R}^d$ and $\mathbf{A} \in \mathbb{R}^{d \times d}$. In the following questions, $\nabla_{\mathbf{x}}$ denotes the gradient with respect to \mathbf{x} , which, by convention, is a column vector.

Calculate the following derivatives.

- (a) $\nabla_{\mathbf{x}}(\mathbf{w}^T \mathbf{x})$ $f(x) = \sum_{i=1}^d w_i x_i$ $\frac{\partial f}{\partial x_i} = w_i$ $\begin{bmatrix} \partial f / \partial x_1 \\ \vdots \\ \partial f / \partial x_d \end{bmatrix} = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_d \end{bmatrix} = \mathbf{w}$
 $f(x) = \underline{ax}$ $\frac{\partial f}{\partial x} = w^T$
- (b) $\nabla_{\mathbf{x}}(\mathbf{w}^T \mathbf{A} \mathbf{x})$ $\frac{\partial f}{\partial \mathbf{x}} = \mathbf{w}^T \mathbf{A}$ $\text{tr}(\mathbf{w}^T \mathbf{A} \mathbf{x}) = \text{tr}(\mathbf{x} \mathbf{w}^T \mathbf{A})$
- (c) $\nabla_{\mathbf{A}}(\mathbf{w}^T \mathbf{A} \mathbf{x})$ $f(A) = \mathbf{w}^T \mathbf{A} \mathbf{x}$ $f(A + \Delta) = \mathbf{w}^T (\mathbf{A} + \Delta) \mathbf{x} = \mathbf{w}^T \mathbf{A} \mathbf{x} + \mathbf{w}^T \Delta \mathbf{x}$
 $\text{tr}(ABC) = \text{tr}(CAB)$ $= f(A) + \text{tr}(\mathbf{x} \mathbf{w}^T \Delta)$
 $= \text{tr}(BCA)$ $\frac{\partial f}{\partial A}$
- (d) $\nabla_{\mathbf{x}}(\mathbf{x}^T \mathbf{A} \mathbf{x})$ $\sum w_j A_{ji} x_i$ $\frac{\partial f}{\partial A_{ij}} = w_j x_i$ $\mathbf{w} \mathbf{x}^T$
- (e) $\nabla_{\mathbf{x}}^2(\mathbf{x}^T \mathbf{A} \mathbf{x})$

Now let's apply our identities derived above to a practical problem. Given a design matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$ and label vector $\mathbf{Y} \in \mathbb{R}^n$, the Ordinary least squares regression problem becomes

$$\mathbf{w}^* = \min_{\mathbf{w}} \frac{1}{2} \|\mathbf{X} \mathbf{w} - \mathbf{Y}\|_2^2$$

- (f) Using parts (a) - (e), derive a necessary condition for \mathbf{w}^* . *Note: We do not necessarily assume \mathbf{X} is full rank!*

¹Good resources for matrix calculus are:

- The Matrix Cookbook: <https://www.math.uwaterloo.ca/~hwolkowi/matrixcookbook.pdf>
- Wikipedia: https://en.wikipedia.org/wiki/Matrix_calculus
- Khan Academy: <https://www.khanacademy.org/math/multivariable-calculus/multivariable-derivatives>
- YouTube: <https://www.youtube.com/playlist?list=PLSQL0a2vh4HC5feHa6Rc5c0wbRTx56nF7>.