

LECTURE NOTES ON LINEAR ALGEBRA

A. A. ERGÜR, O. ESKEW, AND D.M. GANDHI

Abstract. Lecture Notes on Linear Algebra given by Dr. A. A. Ergür on 23 January 2025 and 26 January 2025 respectively.

1. Vector Space over \mathbb{R} . A vector space over \mathbb{R} is a collection of objects that can be:

- Added to each other.
- Multiplied by a real number.

For example, in \mathbb{R}^2 , let $a = (a_1, a_2)$, $b = (b_1, b_2)$, and $3a = (3a_1, 3a_2)$. Then:

$$a + b = (a_1, a_2) + (b_1, b_2) = (a_1 + b_1, a_2 + b_2)$$

Example (see figure 1): Shapes in \mathbb{R}^2 that include the origin $(0,0)$.
A diagram illustrates:

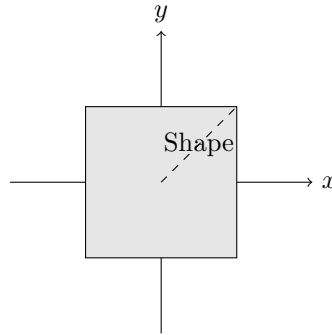


FIG. 1. Shapes in \mathbb{R}^2 including the origin[3]

- Shape K (a square with a dot at $(0,0)$) plus shape L (a circle with a dot at $(0,0)$) results in shape $K + L$ (a rounded square with a dot at $(0,0)$).
- Scalar multiplication: $2 \times K$ (a square with a dot at $(0,0)$) results in $2K$ (a larger square with a dot at $(0,0)$).

1.1. Basis of a Vector Space. In \mathbb{R}^2 , consider the standard basis:

$$e_1 = (1, 0), \quad e_2 = (0, 1)$$

Any vector $x = (x_1, x_2) \in \mathbb{R}^2$ can be written as:

$$x = x_1 e_1 + x_2 e_2$$

2. Norms. A norm is a function that attaches a number to each element x of a vector space, intended to measure its size. For $x \in \mathbb{R}^n$, where $x = (x_1, \dots, x_n)$:

*We thank Robbins family for supporting the Algorithmic Foundations of Data Science Course

- Euclidean norm (ℓ_2 -norm):

$$\|x\|_2 = (x_1^2 + x_2^2 + \cdots + x_n^2)^{1/2}$$

- ℓ_1 -norm:

$$\|x\|_1 = |x_1| + |x_2| + \cdots + |x_n|$$

- ℓ_p -norm ($1 \leq p < \infty$):

$$\|x\|_p = \left(\sum_{i=1}^n |x_i|^p \right)^{1/p}$$

- Infinity norm (ℓ_∞ -norm):

$$\|x\|_\infty = \max_{1 \leq i \leq n} |x_i|$$

- Example: Consider $x = \left(\frac{1}{\sqrt{n}}, \frac{1}{\sqrt{n}}, \dots, \frac{1}{\sqrt{n}} \right) \in \mathbb{R}^n$.

$$\|x\|_2 = \left(\frac{1}{n} + \frac{1}{n} + \cdots + \frac{1}{n} \right)^{1/2} = \left(n \cdot \frac{1}{n} \right)^{1/2} = 1$$

$$\|x\|_1 = \frac{1}{\sqrt{n}} + \frac{1}{\sqrt{n}} + \cdots + \frac{1}{\sqrt{n}} = n \cdot \frac{1}{\sqrt{n}} = \sqrt{n}$$

$$\|x\|_\infty = \frac{1}{\sqrt{n}}$$

Note: Images of ℓ_p -unit balls (Fig 2 and 3) are provided by Kayden Mimmace, with code available on GitHub.

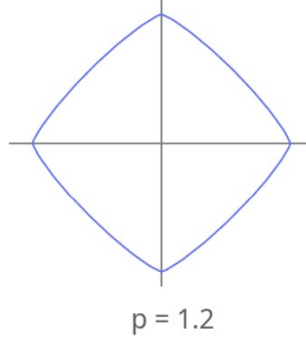


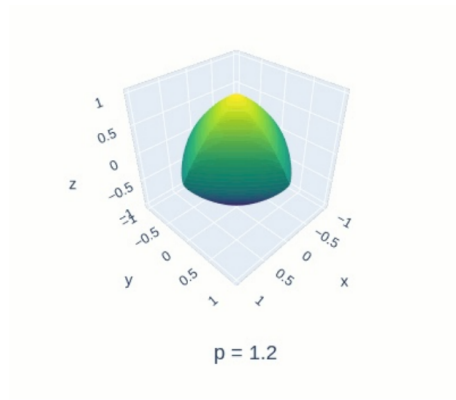
FIG. 2. ℓ_p -unit ball in \mathbb{R}^2 for $p = 1.2$ generated using the code [2].

2.1. ℓ_p -Unit Balls. The ℓ_p -unit ball in \mathbb{R}^n is defined as:

$$S_p := \{x \in \mathbb{R}^n : \|x\|_p \leq 1\}$$

Diagrams illustrate:

- For $p = 1$: A diamond shape in \mathbb{R}^2 .
- For $p = 1.2$: A rounded diamond in \mathbb{R}^2 and a 3D plot in \mathbb{R}^3 .

FIG. 3. 3D ℓ_p -unit ball for $p = 1.2$ [2].

2.2. Hölder Inequality. For $x \in \mathbb{R}^n$, and $1 \leq p \leq q \leq \infty$:

$$\|x\|_q \leq \|x\|_p \leq n^{1/q-1/p} \|x\|_q$$

For $p = 1$, $q = 2$, and x as in the example above:

$$\|x\|_2 \leq \|x\|_1 \leq n^{1/2-1} \|x\|_2$$

3. Exercise. Show that for $x \in \mathbb{R}^n$, if $p > 2 \log n$, then:

$$\|x\|_\infty \leq \|x\|_p \leq c \|x\|_\infty$$

4. Inner Product. For a real vector space V , an inner product $\langle \cdot, \cdot \rangle$ satisfies:

- Symmetry: $\langle x, y \rangle = \langle y, x \rangle$
- Linearity: $\langle ax + by, z \rangle = a \langle x, z \rangle + b \langle y, z \rangle$
- Positive definiteness: $\langle x, x \rangle \geq 0$, and $\langle x, x \rangle = 0$ if and only if $x = 0$.

These properties hold for any $x, y, z \in V$, $a, b \in \mathbb{R}$. The inner product induces a norm:

$$\|x\| = \sqrt{\langle x, x \rangle}$$

4.1. Cauchy-Schwarz Inequality.

$$|\langle x, y \rangle| \leq \|x\| \cdot \|y\|$$

4.2. Angles. The angle θ between vectors x and y is given by:

$$\frac{\langle x, y \rangle}{\|x\| \cdot \|y\|} = \cos \theta$$

A diagram shows vectors x and y with an angle θ between them.

60 **5. Linear Maps and Matrices.** Every linear map is represented by a matrix. A linear map
 61 $f : V \rightarrow \mathbb{R}$ satisfies:

$$62 \quad f(ax + by) = af(x) + bf(y)$$

63 For example, consider:

$$64 \quad f(x) = 3x_1 + 2x_2, \quad x \in \mathbb{R}^2$$

$$65 \quad g(x) = 3x_1 + 2x_2 + 5$$

$$66 \quad h(x) = x_1^2 + 3x_2^2$$

67
 68
 69 A linear map $A : x \rightarrow Ax$ can be represented by a matrix. For instance:

$$70 \quad A = \begin{bmatrix} 3 & 1 \\ 5 & 2 \end{bmatrix}, \quad A \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 3x_1 + x_2 \\ 5x_1 + 2x_2 \end{bmatrix}$$

71 Define basis vectors:

$$72 \quad e_1 = \begin{bmatrix} 5 \\ 3 \end{bmatrix}, \quad e_2 = \begin{bmatrix} 3 \\ 2 \end{bmatrix}$$

73 The discussion involves the inner product $\langle \cdot, \cdot \rangle$. If A^T is the transpose of A , then:

$$74 \quad \langle Ax, y \rangle = \langle x, A^T y \rangle$$

75 for all x, y .

76 **Question:** How to see or hear what a matrix does to a vector?

77 **5.1. Eigenvalues and Eigenvectors.** For $A \in \mathbb{R}^{n \times n}$, $x \in \mathbb{R}^n$, and $\lambda \in \mathbb{R}$:

$$78 \quad Ax = \lambda x$$

79 Here, x is an eigenvector, and λ is an eigenvalue.

80 • Every $n \times n$ matrix A has n complex eigenvalues.

81 **5.2. Singular Value Decomposition (SVD).** Not every matrix A is diagonalizable. Con-
 82 sider using the eigenvalues of X , where:

$$83 \quad A^T A = X, \quad X^T A^T A = X X^T$$

84 **Theorem 4.22 (SVD Theorem):** Let $A \in \mathbb{R}^{m \times n}$ be a rectangular matrix of rank $r \in [0, \min(m, n)]$.
 85 The SVD of A is a decomposition of the form:

$$86 \quad A = U \Sigma V^T$$

87 where:

- 88 • $U \in \mathbb{R}^{m \times m}$ is an orthogonal matrix with column vectors u_i , $i = 1, \dots, m$, satisfying
- 89 $U^T U = I_m$.
- 90 • $V \in \mathbb{R}^{n \times n}$ is an orthogonal matrix with column vectors v_j , $j = 1, \dots, n$, satisfying $V^T V =$
- 91 I_n .
- 92 • $\Sigma \in \mathbb{R}^{m \times n}$ is a diagonal matrix with $\Sigma_{ii} = \sigma_i \geq 0$, and $\Sigma_{ij} = 0$ for $i \neq j$.

93 Thus:

$$94 \quad A = U\Sigma V^T$$

95 The σ_i are eigenvalues of $A^T A$.

$$96 \quad A \in \mathbb{R}^{m \times n}, \quad U \in \mathbb{R}^{m \times m}, \quad V \in \mathbb{R}^{n \times n}, \quad \Sigma \in \mathbb{R}^{m \times n}$$

$$97 \quad \Sigma = \begin{bmatrix} \sigma_1 & & & & & \\ & \sigma_2 & & & & \\ & & \ddots & & & \\ & & & \sigma_r & & \\ & & & & 0 & \\ & & & & & \ddots \\ & & & & & & 0 \end{bmatrix}$$

99 **5.3. Diagonalization.** In general, for any matrix A , if we have:

$$100 \quad P^T A P = D \quad \Rightarrow \quad A = P D P^{-1}$$

101 where $P \in \mathbb{R}^{n \times n}$ is invertible, and D is diagonal:

$$102 \quad D = \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{bmatrix}$$

103 we say A is diagonalizable.

104 **Symmetric Matrices:** If $A^T = A$, then A has real eigenvalues $\lambda_1, \lambda_2, \dots, \lambda_n$, and eigenvectors
105 u_1, u_2, \dots, u_n , with:

$$106 \quad U = [u_1 \quad u_2 \quad \cdots \quad u_n], \quad \langle u_i, u_j \rangle = \begin{cases} 1 & \text{if } i = j \\ 0 & \text{if } i \neq j \end{cases}$$

107 If $U^T U = I_n$, then for all $x \in \mathbb{R}^n$:

$$108 \quad \|Ux\|_2 = \|x\|_2, \quad \langle Ux, Uy \rangle = \langle x, y \rangle$$

$$109 \quad x = yz, \quad x = z^T y^T$$

111 **5.4. Eigenbasis.** For $x = (x_1, x_2, \dots, x_n) \in \mathbb{R}^n$:

$$112 \quad x = x_1 u_1 + x_2 u_2 + \cdots + x_n u_n \quad (\text{in } U\text{-basis, } x = (x_1, x_2, \dots, x_n))$$

$$113 \quad Ax = (Ax_1 u_1 + Ax_2 u_2 + \cdots + Ax_n u_n)$$

$$114 \quad Ax = \lambda_1 x_1 u_1 + \lambda_2 x_2 u_2 + \cdots + \lambda_n x_n u_n$$

117 In the eigenbasis, A becomes:

$$118 \quad \Lambda = \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{bmatrix}$$

$$A = U\Lambda U^{-1}$$

This is undoing the change of basis to a diagonal form.

$$U^T = \begin{bmatrix} u_1 \\ u_2 \\ \vdots \\ u_n \end{bmatrix}, \quad U = [u_1 \quad u_2 \quad \cdots \quad u_n]$$

$$UU^T = \begin{bmatrix} 1 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 1 \end{bmatrix} = I_n$$

5.5. Symmetric Matrices and Eigenvalues.

- Some eigenvalues may be repeating.

- Example: For $A = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$, what is λ ? Compare with $\begin{bmatrix} 2 & \\ & 2 \end{bmatrix}$.

- If $A^T = A$, i.e., A is symmetric, then all eigenvalues are real.

Let λ_1, λ_2 be distinct eigenvalues of A with eigenvectors x and y :

$$Ax = \lambda_1 x, \quad Ay = \lambda_2 y$$

Suppose A is symmetric. Then:

$$\langle Ax, y \rangle = \langle x, A^T y \rangle = \langle x, Ay \rangle$$

$$\lambda_1 \langle x, y \rangle = \lambda_2 \langle x, y \rangle$$

Since $\lambda_1 \neq \lambda_2$:

$$\lambda_1 \langle x, y \rangle = \lambda_2 \langle x, y \rangle \quad \Rightarrow \quad \langle x, y \rangle = 0$$

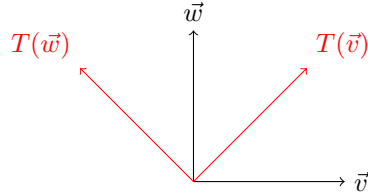


FIG. 4. Orthogonality and linear transformations[3]

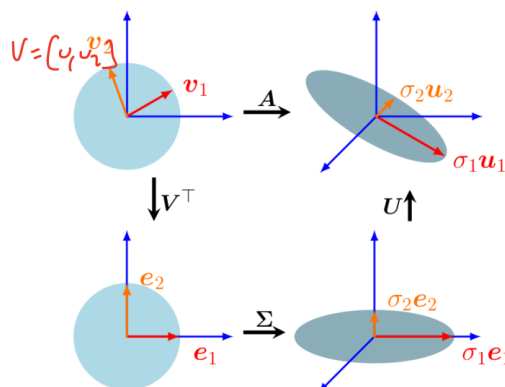


FIG. 5. SVD transformation of a unit sphere, adapted from [1]

5.6. SVD and Projections. If A is $n \times n$ and no x exists such that $Ax = 0$, consider the SVD:

$$A = U\Sigma V^T$$

Diagrams illustrate the transformation of a unit sphere under A :

- $V = \{v_1, v_2, v_3\}$, a sphere in \mathbb{R}^3 , transforms via Σ to an ellipsoid with axes $\sigma_1 v_1, \sigma_2 v_2, \sigma_3 v_3$, and then via U^T .
- V^T maps the ellipsoid back to a sphere with axes $\sigma_1 e_1, \sigma_2 e_2$, and U rotates it.

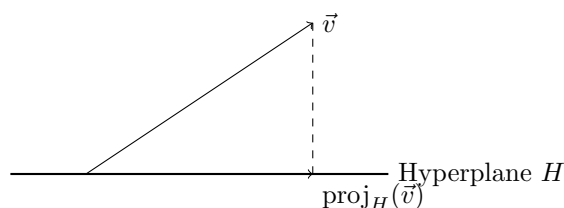
Question: How about projections?

5.7. Projections. Define a projection:

$$A^2 = A \quad \Leftrightarrow \quad A \cdot (Ax) = Ax \quad \text{for all } x \in \mathbb{R}^n$$

Then A is a projection.

Example (Main): Let $y_1, y_2 \in \mathbb{R}^3$, and let H be the span of y_1, y_2 . A represents the projection of x onto H .

FIG. 6. Projection onto a hyperplane H [3]

A diagram shows $x \in \mathbb{R}^3$, H as a plane, and Ax as the projection of x onto H .

- Ax is the closest point to x in H .
- For $t \in H$, $x - Ax \perp t$, i.e., $\langle x - Ax, t \rangle = 0$.

154 For example:

$$155 \quad y_1 = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, \quad y_2 = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$$

$$156 \quad Y = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix} \in \mathbb{R}^{3 \times 2}$$

$$158 \quad Y = QR$$

160 Find A :

$$161 \quad A = QQ^T$$

$$162 \quad Q^T Q = I \Rightarrow QQ^T = I$$

164 For $A = U\Sigma V^T$:

$$165 \quad \langle Ax, Ax \rangle = \|Ax\|_2^2$$

$$166 \quad \|Ax\|_2^2 = \langle Ax, Ax \rangle = \langle x, A^T Ax \rangle$$

$$168 \quad \|Ax\|_2^2 = \lambda \langle x, x \rangle = \lambda \|x\|_2^2$$

170 So, all σ_i are non-negative.

$$171 \quad \Sigma = \begin{bmatrix} \sigma_1 & & & & & \\ & \sigma_2 & & & & \\ & & \ddots & & & \\ & & & \sigma_r & & \\ & & & & 0 & \\ & & & & \ddots & \\ & & & & & 0 \end{bmatrix} \in \mathbb{R}^{m \times n}$$

$$172 \quad \|Ax\|_2^2 = 0 \Rightarrow \|x\|_2 = 0 \Rightarrow Ax = 0$$

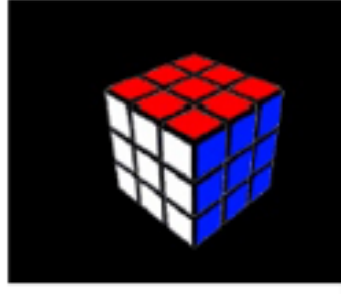
174 • The number of zeros is related to the kernel of A , i.e., $\{x : Ax = 0\}$.

$$175 \quad A = U\Sigma V^T$$

$$176 \quad Ax = U\Sigma V^T x = U \begin{bmatrix} \sigma_1 & & & \\ & \sigma_2 & & \\ & & \ddots & \\ & & & 0 \end{bmatrix} V^T x$$

178 **6. Tensors.** An order- d tensor with n variables is an $n \times n \times \cdots \times n$ (d times) data array.

- 179 • $d = 3$: Very common, e.g., $\mathbb{R}^{n \times n \times n}$.
- 180 • $d = 2$: Matrix, $\mathbb{R}^{n \times n}$.
- 181 • $d = 1$: Vector, \mathbb{R}^n .

FIG. 7. A $3 \times 3 \times 3$ tensor [3]

7. QR Decomposition. For $Y \in \mathbb{R}^{m \times n}$, $m \geq n$:

$$Y = QR, \quad Q \in \mathbb{R}^{m \times n}, \quad R \in \mathbb{R}^{n \times n}$$

where $Q^T Q = I_n$, Q is orthogonal, and R is upper triangular.

A diagram illustrates:

$$Y \in \mathbb{R}^{m \times n}, \quad Q \in \mathbb{R}^{m \times n}, \quad R \in \mathbb{R}^{n \times n}$$

with R having $n - \text{rank}(Y)$ zero rows.

- Suppose $\text{rank}(Y) = n$, $m \geq n$. Then $Y = QR$, $Y \in \mathbb{R}^{m \times n}$, $R \in \mathbb{R}^{n \times n}$ is invertible.

Goal: Given $x \in \mathbb{R}^m$, find $w \in \mathbb{R}^n$ such that:

$$\|x - Yw\|_2^2 = \min_{t \in \text{span}(Y)} \|x - t\|_2^2$$

$$Yw = QQ^T x, \quad w = R^{-1} Q^T x$$

$$Y \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \end{bmatrix} = w_1 y_1 + w_2 y_2 + \cdots + w_n y_n$$

Using Singular Values to Analyze A . Let A be an $m \times n$ matrix. The singular value decomposition (SVD) of A is given by:

$$A = U \begin{bmatrix} \delta_1 & & \\ & \ddots & \\ & & \delta_n \end{bmatrix} V^T$$

where U is $m \times m$, $U^T U = I_m$, V is $n \times n$, $V^T V = I_n$, and $\delta_1, \delta_2, \dots, \delta_n$ are the singular values of A .

- For $x \in \mathbb{R}^n$,

$$\|x\|_2 \cdot \delta_n(A) \leq \|Ax\|_2 \leq \delta_1(A) \cdot \|x\|_2$$

- $\delta_1(A)$ is called the **operator norm** of A , denoted by $\|A\|_2$ or $\|A\|_{\text{op}}$.

- 204 • The ratio $\frac{\delta_1(A)}{\delta_n(A)}$ is called the **condition number** of A , denoted by $\kappa(A)$. This is used in
- 205 LAPACK.
- 206 • The sum $\delta_1(A) + \delta_2(A) + \cdots + \delta_n(A)$ is called the **nuclear norm**, denoted by $\|A\|_*$.
- 207 • The sum $\delta_1(A)^2 + \delta_2(A)^2 + \cdots + \delta_n(A)^2$ is called the **Frobenius norm** or **Hilbert-Schmidt**
- 208 **norm**, denoted by $\|A\|_F$ or $\|A\|_{\text{HS}}$.

209 **Trace Norm or $\|A\|_2$.** This has a specific meaning. Consider a matrix A :

$$210 \quad A = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{bmatrix}$$

211 Then the Frobenius norm squared is:

$$212 \quad \|A\|_F^2 = \delta_1(A)^2 + \cdots + \delta_n(A)^2 = \sum_{1 \leq i, j \leq n} a_{ij}^2$$

213 **What's up with trace-norm naming?.** The trace of a matrix X is defined as:

$$214 \quad \text{Tr}(X) = x_{11} + x_{22} + \cdots + x_{nn} \quad \text{for} \quad X = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1n} \\ x_{21} & x_{22} & \cdots & x_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{nn} \end{bmatrix}$$

215 For the inner product on matrices $\langle A, B \rangle = \text{Tr}(B^T A)$, the Frobenius norm can be expressed as:

$$216 \quad \|A\|_F = \sqrt{\langle A, A \rangle} = \sqrt{\text{Tr}(A^T A)}$$

217 **When Do Eigenvalues and Singular Values Coincide?.** Let A be a symmetric matrix

218 with singular value decomposition (SVD) $A = U \Sigma V^T$. Then:

$$219 \quad A^T = V \Sigma U^T$$

220 Since A is symmetric, $A^T = A$, so:

$$221 \quad V \Sigma U^T = U \Sigma V^T$$

222 This implies $V = U$. Thus, the SVD becomes:

$$223 \quad A = U \Sigma U^T$$

224 Now consider the eigenvalue decomposition of A :

$$225 \quad A = V \Lambda V^T$$

226 Since $V = U$, we have:

$$227 \quad \Lambda = \Sigma$$

228 Thus, $\Lambda = \Sigma$, meaning the eigenvalues of A must coincide with the singular values. Additionally,

229 since Σ contains non-negative values, all eigenvalues of A are non-negative.

- 230 • If all eigenvalues are non-negative, A is positive semidefinite (PSD).
- 231 • If all eigenvalues are positive, A is positive definite (PD).

Cholesky Decomposition. If A is an $n \times n$ positive definite matrix, then it has a Cholesky decomposition:

$$A = RR^T$$

where R is a real, upper triangular matrix with positive diagonal entries.

If A is PD and has the form $A = RR^T$, then for any $x, y \in \mathbb{R}^n$:

$$\langle Ax, y \rangle = \langle x, Ay \rangle = \langle x, RR^T y \rangle$$

$$\langle Ax, y \rangle = \langle Rx, Ry \rangle$$

Thus, if we define a new inner product $\langle \cdot, \cdot \rangle_R$ such that:

$$\langle x, y \rangle_R = \langle Rx, Ry \rangle$$

we have:

$$\langle Ax, y \rangle = \langle x, y \rangle_R$$

This implies:

- All possible inner products on $\mathbb{R}^n \leftrightarrow$ all PD matrices A .
- All similarity measures using angles.

Linear Regression. Input: Labeled vectors (x_i, y_i) , $i = 1, 2, \dots, N$, where $x_i \in \mathbb{R}^n$ (vector with n coordinates) and $y_i \in \mathbb{R}$.

Goal: Develop a linear model to predict the output value y given $x = (x_1, \dots, x_n)$.

In other words, find a linear function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ that best fits the data.

What does "best" mean?

For now, define the residual sum of squares (RSS):

$$\text{RSS}(f) = \sum_{i=1}^N (y_i - f(x_i))^2$$

The goal is to minimize $\text{RSS}(f)$ among all linear functions f .

In this context, "linear" is used in a restrictive way to mean:

$$f(x) = w_1x_1 + w_2x_2 + \dots + w_nx_n + w_0$$

So the model is:

$$y_i \approx w_1x_{i1} + w_2x_{i2} + \dots + w_nx_{in} + w_0 \quad (\star)$$

Let $w = (w_1, w_2, \dots, w_n, w_0)$, and define:

$$\tilde{x}_i = (x_{i1}, x_{i2}, \dots, x_{in}, 1)$$

$$(\star) \quad \text{becomes} \quad \langle w, \tilde{x}_i \rangle \approx w^T \tilde{x}_i$$

The RSS can be written as:

$$\text{RSS}(f) = \sum_{i=1}^N (y_i - w^T \tilde{x}_i)^2$$

265 More concisely, define:

$$266 \quad y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}, \quad X = \begin{bmatrix} x_{11} & x_{21} & \cdots & x_{N1} \\ x_{12} & x_{22} & \cdots & x_{N2} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & \cdots & 1 \end{bmatrix}$$

267 Then:

$$268 \quad \text{RSS}(f) = \|y - X^T w\|^2$$

269 This means projecting $y = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix}$ into the span of the rows of X .

270 If $X^T = QR$, where Q is $N \times (n+1)$ and R is $(n+1) \times (n+1)$, then:

$$271 \quad X^T w = QRw \quad \text{and} \quad w = R^{-1}Q^T y$$

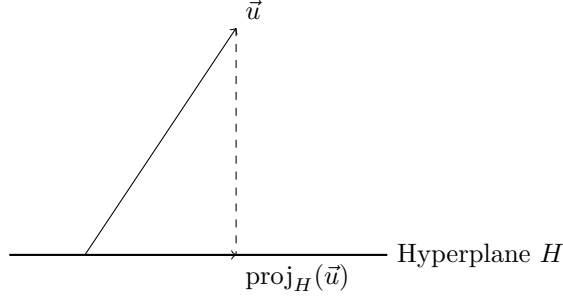


FIG. 8. Second example of a projection onto a hyperplane[3]

272

273 A diagram illustrates this: y (real labels) is projected onto the column span of X , with $X^T w$ being
274 the best linear approximation.

275 **How do you compute that Q ?**

276 **PCA (Principal Component Analysis).** For $d = 1$:

277 **Input:** $x_1, x_2, \dots, x_N \in \mathbb{R}^n$.

278 **Goal:** Find a d -dimensional vector space $L \subset \mathbb{R}^n$ such that:

$$279 \quad \sum_{i=1}^N \|x_i - P_L(x_i)\|^2$$

280 is minimized, where $P_L(x_i)$ is the projection of x_i onto L .

281 **• Simplification 1:** Define the mean of the data:

$$282 \quad \mu = \frac{1}{N} \sum_{i=1}^N x_i$$

So the centered data is:

$$(0, 0, \dots, 0) = \frac{1}{N} \sum_{i=1}^N (x_i - \mu)$$

We'll assume $\mu = (0, \dots, 0)$.

- **Simplification 2:** If $v_1, v_2, \dots, v_d \in \mathbb{R}^n$ is an orthonormal basis for the vector space L , and we set:

$$W = \begin{bmatrix} v_1 & v_2 & \dots & v_d \end{bmatrix} \quad (\text{an } n \times d \text{ matrix})$$

then the projection is:

$$P_L(x) = WW^T x$$

Goal: Find $W = \begin{bmatrix} v_1 & v_2 & \dots & v_d \end{bmatrix}$, an $n \times d$ matrix, such that $W^T W = I_d$, and:

$$\sum_{i=1}^N \|x_i - WW^T x_i\|^2 = \min_L \sum_{i=1}^N \|x_i - P_L(x_i)\|^2$$

Objective Function Restated.

$$\sum_{i=1}^N \|x_i - WW^T x_i\|^2 = \sum_{i=1}^N \tilde{x}_i^T (I - WW^T) \tilde{x}_i$$

where:

$$\tilde{x}_i = (x_i - \mu) \quad (\text{but we assumed } \mu = 0, \text{ so } \tilde{x}_i = x_i).$$

Define the data matrix:

$$X = \begin{bmatrix} x_1 & x_2 & \dots & x_N \end{bmatrix} \quad (\text{an } n \times N \text{ matrix}).$$

Then:

$$X^T (I - WW^T) X \quad (\text{an } N \times N \text{ matrix}).$$

The trace of the objective function is:

$$\text{Trace}(X^T (I - WW^T) X) = \sum_{i=1}^N \|x_i - WW^T x_i\|^2$$

So the optimization problem becomes:

$$\min_{W \in \mathbb{R}^{n \times d}, W^T W = I_d} \text{Trace}(X^T (I - WW^T) X)$$

This is equivalent to:

$$\min_{W \in \mathbb{R}^{n \times d}, W^T W = I_d} \text{Trace}(X^T X) - \text{Trace}(X^T W W^T X)$$

which is equivalent to:

$$\max_{W \in \mathbb{R}^{n \times d}, W^T W = I_d} \text{Trace}(X^T W W^T X) \quad \langle A, A \rangle \quad \text{where} \quad A^T = W^T X, \quad A = X^T W$$

308 Alternatively:

309
$$\max_{W \in \mathbb{R}^{n \times d}, W^T W = I_d} \|W^T X\|_F^2$$

310 where $W^T X$ is a $d \times N$ matrix, and X is an $n \times N$ matrix with d singular values.

311 This can also be written as:

312
$$\max_{W \in \mathbb{R}^{n \times d}, W^T W = I_d} \sum_{i=1}^d \delta_i (W^T X)^2$$

313 Given the SVD of X :

314
$$X = U \Sigma V^T \quad (\text{where } U \text{ is } n \times n, \Sigma \text{ is } n \times N, V \text{ is } N \times N),$$

315

316
$$\Sigma = \begin{bmatrix} \delta_1 & & & & \\ & \delta_2 & & & \\ & & \ddots & & \\ & & & \delta_r & \\ & & & 0 & \\ & & & & \ddots \\ & & & & & 0 \end{bmatrix},$$

317

318
$$U = [u_1 \quad u_2 \quad \cdots \quad u_n],$$

319 we pick the best d column vectors from U , i.e., $W^T = [u_1 \quad u_2 \quad \cdots \quad u_d]$.

Recap: Linear Regression.

320
$$x_1, x_2, \dots, x_N \in \mathbb{R}^n, \quad y_1, y_2, \dots, y_N \in \mathbb{R}, \quad (y_1, y_2, \dots, y_N) \in \mathbb{R}^N$$

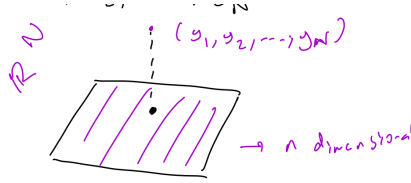


FIG. 9. Basis of a subspace and projection onto it[3].

321

322 A diagram illustrates \mathbb{R}^N with the n -dimensional row span of X :

323
$$X = [x_1 \quad x_2 \quad \cdots \quad x_N] \quad (\text{an } n \times N \text{ matrix}).$$

324 Another diagram shows $x_1, x_2, \dots, x_N \in \mathbb{R}^n$ projected onto an n -dimensional subspace (PCA):

325
$$X = [x_1 \quad x_2 \quad \cdots \quad x_N] \quad (\text{an } n \times N \text{ matrix}).$$

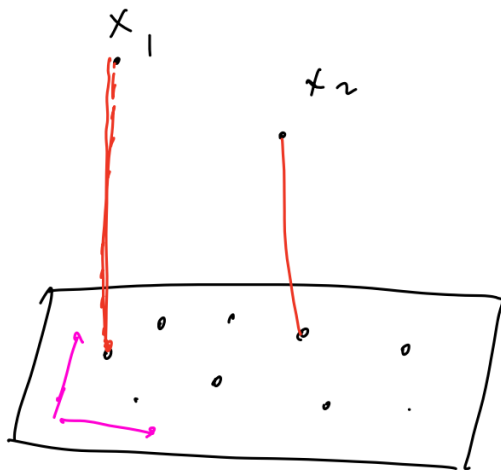


FIG. 10. Projection of vectors onto a plane[3].

326 The SVD of X :

$$327 \quad X = U \begin{bmatrix} \delta_1 & & & & \\ & \delta_2 & & & \\ & & \ddots & & \\ & & & \delta_r & \\ & & & & 0 \\ & & & & & \ddots \\ & & & & & & 0 \end{bmatrix} V^T$$

328 where U is $n \times n$, and we pick the first r columns corresponding to non-zero singular values.

329 Exercises.

330 **Problem 1.** Let A be a symmetric matrix ($A^T = A$), and let $\langle \cdot, \cdot \rangle$ be an inner product. Show
 331 that $\langle x, Ay \rangle = \langle Ax, y \rangle$.

332 Since A is symmetric, $A^T = A$. Using the standard inner product $\langle x, y \rangle = x^T y$, we have:

$$333 \quad \langle x, Ay \rangle = x^T (Ay) = x^T A y$$

334

$$335 \quad \langle Ax, y \rangle = (Ax)^T y = x^T A^T y = x^T A y \quad (\text{since } A^T = A)$$

336 Thus:

$$337 \quad \langle x, Ay \rangle = x^T A y = \langle Ax, y \rangle$$

338

339

340 **Problem 2.** Define the function $f : \mathbb{R}^2 \rightarrow \mathbb{R}^2$, where $f(x)$ is obtained by turning x counter-
 341 clockwise by 45° . Find the matrix that represents this function using the standard basis $e_1 = (1, 0)$
 342 and $e_2 = (0, 1)$.
 343 A counter-clockwise rotation by 45° in \mathbb{R}^2 is represented by the matrix:

$$344 \quad R = \begin{pmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{pmatrix}$$

345 For $\theta = 45^\circ$, we have $\cos 45^\circ = \sin 45^\circ = \frac{\sqrt{2}}{2}$, so:

$$346 \quad R = \begin{pmatrix} \frac{\sqrt{2}}{2} & -\frac{\sqrt{2}}{2} \\ \frac{\sqrt{2}}{2} & \frac{\sqrt{2}}{2} \end{pmatrix}$$

347 **Problem 3.** Let $Q \in \mathbb{R}^{n \times n}$ be a matrix such that $Q^T Q = I_n$. Show that for any $x, y \in \mathbb{R}^n$,
 348 the angle between x and y is the same as the angle between Qx and Qy .
 349 The angle θ between x and y is given by:

$$350 \quad \cos \theta = \frac{\langle x, y \rangle}{\|x\| \|y\|} = \frac{x^T y}{\|x\| \|y\|}$$

351 For Qx and Qy , compute the inner product:

$$352 \quad \langle Qx, Qy \rangle = (Qx)^T (Qy) = x^T Q^T Q y = x^T I_n y = x^T y = \langle x, y \rangle$$

353 The norms are:

$$354 \quad \|Qx\| = \sqrt{(Qx)^T (Qx)} = \sqrt{x^T Q^T Q x} = \sqrt{x^T x} = \|x\|$$

355 Similarly, $\|Qy\| = \|y\|$. Thus:

$$356 \quad \cos \theta' = \frac{\langle Qx, Qy \rangle}{\|Qx\| \|Qy\|} = \frac{x^T y}{\|x\| \|y\|} = \cos \theta$$

357 The angles are the same.

358 **Problem 4.** Recall that for $1 \leq p < \infty$ and $x \in \mathbb{R}^n$, we define $\|x\|_p = (\sum_{i=1}^n |x_i|^p)^{1/p}$.

359 **Part a.** Let $x \in \mathbb{R}^n$ be a vector with 100 non-zero entries. Show that $\frac{1}{10} \leq \frac{\|x\|_1}{\|x\|_\infty}$.

360 We have $\|x\|_1 = \sum_{i=1}^n |x_i|$ and $\|x\|_\infty = \max_i |x_i|$. Let x have 100 non-zero entries, say $|x_i| = a_i$ for
 361 $i = 1$ to 100, and $x_i = 0$ otherwise. Then:

$$362 \quad \|x\|_1 = \sum_{i=1}^{100} a_i, \quad \|x\|_\infty = \max_{i=1}^{100} a_i = M$$

363 If all non-zero entries are equal to M , then:

$$364 \quad \|x\|_1 = 100M, \quad \|x\|_\infty = M \implies \frac{\|x\|_1}{\|x\|_\infty} = 100 \geq \frac{1}{10}$$

365 In the minimal case (e.g., one entry is M , others smaller), $\|x\|_1 \geq M$, so:

$$366 \quad \frac{\|x\|_1}{\|x\|_\infty} \geq 1 \geq \frac{1}{10}$$

367 The inequality holds.

Part b. Let $x \in \mathbb{R}^{8000}$. Show that $\ell \cdot \|x\|_q \leq \|x\|_\infty \leq \|x\|_q$, where e denotes the natural base.
 (Note: Assuming ℓ is a typo or constant; interpreting as a norm comparison.) For $\|x\|_\infty = \max_i |x_i|$
 and $\|x\|_q = \left(\sum_{i=1}^{8000} |x_i|^q \right)^{1/q}$, we have:

$$\|x\|_\infty \leq \|x\|_q \leq 8000^{1/q} \|x\|_\infty$$

The exact role of ℓ or e is unclear, but the standard norm comparison holds as shown.

Exercises (First Set).

Exercise 1. Let L be the line spanned by the vector $(-1, 1, 0) \in \mathbb{R}^3$. Let A be the matrix that represents the projection onto this line. Compute A .
 The vector $v = (-1, 1, 0)$ spans the line. The projection matrix is:

$$A = \frac{vv^T}{v^T v}$$

$$v^T v = (-1)^2 + 1^2 + 0^2 = 2$$

$$vv^T = \begin{pmatrix} -1 \\ 1 \\ 0 \end{pmatrix} \begin{pmatrix} -1 & 1 & 0 \end{pmatrix} = \begin{pmatrix} 1 & -1 & 0 \\ -1 & 1 & 0 \\ 0 & 0 & 0 \end{pmatrix}$$

$$A = \frac{1}{2} \begin{pmatrix} 1 & -1 & 0 \\ -1 & 1 & 0 \\ 0 & 0 & 0 \end{pmatrix} = \begin{pmatrix} \frac{1}{2} & -\frac{1}{2} & 0 \\ -\frac{1}{2} & \frac{1}{2} & 0 \\ 0 & 0 & 0 \end{pmatrix}$$

Exercise 2. Generate 100 random Gaussian five-dimensional vectors. Compute the matrix that represents the projection onto the span of $x_1, x_2, \dots, x_{100} \in \mathbb{R}^5$ using QR decomposition.
 Using Python with NumPy:

```
import numpy as np
np.random.seed(42)
X = np.random.randn(5, 100) # 5x100 matrix
Q, R = np.linalg.qr(X)
A = Q @ Q.T # Projection matrix
```


 The matrix A is 5×5 and projects onto the span of the columns of X .

Exercises (Second Set).

Exercise 2. We define the Hilbert-Schmidt norm as $\|A\|_{HS} = \left(\sum_{i,j} a_{ij}^2 \right)^{1/2}$, and the trace inner product as $\langle A, B \rangle = \text{Trace}(B^T A)$, with norm $\|A\|_2 = \sqrt{\langle A, A \rangle}$. Show $\|A\|_{HS} = \|A\|_2$.

$$\|A\|_{HS} = \left(\sum_{i,j} a_{ij}^2 \right)^{1/2}$$

$$\|A\|_2 = \sqrt{\text{Trace}(A^T A)}$$

399 The (i, i) -th entry of $A^T A$ is $\sum_j a_{ji}^2$, so:

$$400 \quad \text{Trace}(A^T A) = \sum_i \sum_j a_{ji}^2 = \sum_{i,j} a_{ij}^2$$

$$401 \quad \|A\|_2 = \sqrt{\text{Trace}(A^T A)} = \sqrt{\sum_{i,j} a_{ij}^2} = \|A\|_{HS}$$

402
403 For further study, we recommend the textbooks by Murphy [3] and Streil [4].

404 REFERENCES

- 405 [1] M. P. DEISENROTH, A. A. FAISAL, AND C. S. ONG, *Mathematics for machine learning*, Cambridge University
406 Press, 2020.
- 407 [2] K. MIMMACE, *Unit ball visualizations*. <https://github.com/kaydenmimmace/unit-ball-visualizations>, 2025. Ac-
408 cessed: 2025-05-14.
- 409 [3] K. P. MURPHY, *Probabilistic Machine Learning: An Introduction*, 2023, [https://probml.github.io/pml-book/](https://probml.github.io/pml-book/book1.html)
410 [book1.html](https://probml.github.io/pml-book/book1.html). Available online.
- 411 [4] S. STREIL, *Linear algebra done wrong*, 2023, <https://www.math.brown.edu/streil/papers/LADW/LADW.html>.
412 Available online.