

Introductory Topics in Linear Algebra for Data Analysis

Contributors:
Sugandha Roy
Hasnain Hossain
Rakibul Hasan Rajib
Amin Ahsan Ali
AI & ML Wing, CCDS

September 28, 2024

1 Matrix Multiplication

Multiplication of two matrices, two vectors, and a matrix and a vector can be written in many different ways.

- We will consider vectors and matrices with real elements unless otherwise stated.
- We will use lower case letters to denote n -vectors in column or row orientations. We will use, for example, \mathbf{a} to denote a column vector ($n \times 1$) and \mathbf{a}^* to denote a row vector ($1 \times n$) respectively.
- We will use upper case letters to denote a matrix. An $m \times n$ matrix A can be written in terms of its n columns (where each column \mathbf{a}_i is an m -vector) or its m rows (where each row \mathbf{a}_i^* is an n -vector):

$$\underbrace{\mathbf{A}}_{m \times n} = \left[\begin{array}{c|c|c|c} | & | & \dots & | \\ \mathbf{a}_1 & \mathbf{a}_2 & \dots & \mathbf{a}_n \\ | & | & \dots & | \end{array} \right] = \left[\begin{array}{c} - \mathbf{a}_1^* - \\ - \mathbf{a}_2^* - \\ \vdots \\ - \mathbf{a}_m^* - \end{array} \right]$$

1.1 Useful formulas:

1.

$$\underbrace{\mathbf{a}}_{\text{col}} \underbrace{\mathbf{b}^*}_{\text{row}} = \underbrace{\begin{bmatrix} a_1 \\ \vdots \\ a_m \end{bmatrix}}_{m \times 1} \underbrace{\begin{bmatrix} b_1 & \dots & b_p \end{bmatrix}}_{1 \times p} = \underbrace{\begin{bmatrix} | & | & \dots & | \\ b_1 \mathbf{a} & b_2 \mathbf{a} & \dots & b_p \mathbf{a} \\ | & | & \dots & | \end{bmatrix}}_{m \times p} \quad (\text{a rank 1 matrix})$$

2.

$$\underbrace{\mathbf{AB}}_{m \times p} = \underbrace{\begin{bmatrix} | & | & \dots & | \\ \mathbf{a}_1 & \mathbf{a}_2 & \dots & \mathbf{a}_n \\ | & | & \dots & | \end{bmatrix}}_{m \times n} \underbrace{\begin{bmatrix} - \mathbf{b}_1^* - \\ - \mathbf{b}_2^* - \\ \vdots \\ - \mathbf{b}_n^* - \end{bmatrix}}_{n \times p} \\ = \mathbf{a}_1 \mathbf{b}_1^* + \dots + \mathbf{a}_n \mathbf{b}_n^* \quad (\text{sum of rank 1 matrices})$$

3.

$$\mathbf{AB} = \mathbf{A} \left[\begin{array}{c|c|c|c} | & | & \dots & | \\ \mathbf{b}_1 & \mathbf{b}_2 & \dots & \mathbf{b}_p \\ | & | & \dots & | \end{array} \right] = [\mathbf{Ab}_1 \quad \mathbf{Ab}_2 \quad \dots \quad \mathbf{Ab}_p]$$

4. From 3

$$AB = \underbrace{\begin{bmatrix} -\mathbf{a}_1^* \\ -\mathbf{a}_2^* \\ \vdots \\ -\mathbf{a}_m^* \end{bmatrix}}_{m \times n} \underbrace{\begin{bmatrix} | & | & & | \\ \mathbf{b}_1 & \mathbf{b}_2 & \dots & \mathbf{b}_p \\ | & | & & | \end{bmatrix}}_{n \times p} = \begin{bmatrix} \mathbf{a}_1^* \mathbf{b}_1 & \mathbf{a}_1^* \mathbf{b}_2 & \dots & \mathbf{a}_1^* \mathbf{b}_p \\ \mathbf{a}_2^* \mathbf{b}_1 & \mathbf{a}_2^* \mathbf{b}_2 & \dots & \mathbf{a}_2^* \mathbf{b}_p \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{a}_m^* \mathbf{b}_1 & \mathbf{a}_m^* \mathbf{b}_2 & \dots & \mathbf{a}_m^* \mathbf{b}_p \end{bmatrix} = \underbrace{\begin{bmatrix} -\mathbf{a}_1^* B \\ -\mathbf{a}_2^* B \\ \vdots \\ -\mathbf{a}_m^* B \end{bmatrix}}_{m \times p}$$

5. From 2

$$\mathbf{A}\mathbf{x} = \underbrace{\begin{bmatrix} | & | & & | \\ \mathbf{a}_1 & \mathbf{a}_2 & \dots & \mathbf{a}_n \\ | & | & & | \end{bmatrix}}_{m \times n} \underbrace{\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}}_{n \times 1} = x_1 \mathbf{a}_1 + x_2 \mathbf{a}_2 + \dots + x_n \mathbf{a}_n$$

6. From 4

$$\mathbf{A}\mathbf{x} = \begin{bmatrix} -\mathbf{a}_1^* \mathbf{x} \\ -\mathbf{a}_2^* \mathbf{x} \\ \vdots \\ -\mathbf{a}_m^* \mathbf{x} \end{bmatrix} \quad \mathbf{x} = \begin{bmatrix} -\mathbf{a}_1^* \mathbf{x} \\ -\mathbf{a}_2^* \mathbf{x} \\ \vdots \\ -\mathbf{a}_m^* \mathbf{x} \end{bmatrix}$$

7. From 2

$$\mathbf{A} = \begin{bmatrix} | & | & & | \\ \mathbf{a}_1 & \mathbf{a}_2 & \dots & \mathbf{a}_n \\ | & | & & | \end{bmatrix}$$

$$\mathbf{A}\mathbf{A}^T = \sum_{i=1}^n a_i a_i^* \quad (\text{or } \sum_{i=1}^n a_i a_i^T)$$

8. From 4

$$\mathbf{A}^T \mathbf{A} = \begin{bmatrix} a_1^T a_1 & a_1^T a_2 & \dots & a_1^T a_n \\ a_2^T a_1 & a_2^T a_2 & \dots & a_2^T a_n \\ \vdots & \vdots & \ddots & \vdots \\ a_n^T a_1 & a_n^T a_2 & \dots & a_n^T a_n \end{bmatrix} = \begin{bmatrix} -\mathbf{a}_1^T \mathbf{A} \\ -\mathbf{a}_2^T \mathbf{A} \\ \vdots \\ -\mathbf{a}_n^T \mathbf{A} \end{bmatrix}$$

9. From 7

$$\sum_{i=1}^n (\mathbf{a}^T \mathbf{x}_i)^2 = \sum_{i=1}^n (\mathbf{a}^T \mathbf{x}_i)(\mathbf{x}_i^T \mathbf{a}) \quad (\mathbf{x}_i \text{'s are } n\text{-vectors})$$

$$= \sum_{i=1}^n \mathbf{a}^T (\mathbf{x}_i \mathbf{x}_i^T) \mathbf{a}$$

$$= \mathbf{a}^T \mathbf{X} \mathbf{X}^T \mathbf{a}$$

10. From 7

$$\sum_{i=1}^n (\mathbf{A}\mathbf{x}_i)(\mathbf{A}\mathbf{x}_i)^T = \mathbf{A} \left(\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T \right) \mathbf{A}^T = \mathbf{A} \mathbf{X} \mathbf{X}^T \mathbf{A}^T$$

11. Diagonal matrices

(a)

$$PD = \underbrace{\begin{bmatrix} | & | & & | \\ \mathbf{p}_1 & \mathbf{p}_2 & \dots & \mathbf{p}_m \\ | & | & & | \end{bmatrix}}_{m \times m} \underbrace{\begin{bmatrix} d_1 & 0 & \dots & 0 \\ 0 & d_2 & \dots & 0 \\ \vdots & \vdots & \dots & \vdots \\ 0 & 0 & \dots & d_m \end{bmatrix}}_{m \times m} = \begin{bmatrix} | & | & & | \\ d_1 \mathbf{p}_1 & d_2 \mathbf{p}_2 & \dots & d_m \mathbf{p}_m \\ | & | & & | \end{bmatrix}$$

$$DP = \underbrace{\begin{bmatrix} d_1 & 0 & \dots & 0 \\ 0 & d_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & d_m \end{bmatrix}}_{m \times m} \underbrace{\begin{bmatrix} | & | & & | \\ \mathbf{p}_1 & \mathbf{p}_2 & \dots & \mathbf{p}_m \\ | & | & & | \end{bmatrix}}_{m \times m} = \begin{bmatrix} - & d_1 \mathbf{p}_1^* & - \\ - & d_2 \mathbf{p}_2^* & - \\ \vdots & \vdots & \vdots \\ - & d_m \mathbf{p}_m^* & - \end{bmatrix}$$

(b)

$$U\Sigma = \underbrace{\begin{bmatrix} | & | & & | \\ \mathbf{u}_1 & \mathbf{u}_2 & \dots & \mathbf{u}_m \\ | & | & & | \end{bmatrix}}_{m \times m} \underbrace{\begin{bmatrix} \sigma_1 & 0 & 0 & \dots & 0 & 0 & 0 \\ 0 & \sigma_2 & 0 & \dots & 0 & 0 & 0 \\ 0 & 0 & \sigma_3 & \dots & 0 & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & \dots & \sigma_r & 0 & 0 \\ & & & \vdots & & & \\ 0 & 0 & 0 & \dots & 0 & 0 & 0 \end{bmatrix}}_{m \times n} = \underbrace{\begin{bmatrix} | & | & & | & | & | & | \\ \sigma_1 \mathbf{u}_1 & \sigma_2 \mathbf{u}_2 & \dots & \sigma_r \mathbf{u}_r & \mathbf{0} & \dots & \mathbf{0}_n \\ | & | & & | & | & | & | \end{bmatrix}}_{m \times n}$$

(c)

$$U\Sigma V^T = \underbrace{\begin{bmatrix} | & | & & | & | & | & | \\ \sigma_1 \mathbf{u}_1 & \sigma_2 \mathbf{u}_2 & \dots & \sigma_r \mathbf{u}_r & \mathbf{0} & \dots & \mathbf{0}_n \\ | & | & & | & | & | & | \end{bmatrix}}_{m \times n} \underbrace{\begin{bmatrix} - & \mathbf{v}_1^* & - \\ - & \mathbf{v}_2^* & - \\ \vdots & \vdots & \vdots \\ - & \mathbf{v}_n^* & - \end{bmatrix}}_{n \times n} \\ = \sigma_1 \mathbf{u}_1 \mathbf{v}_1^* + \sigma_2 \mathbf{u}_2 \mathbf{v}_2^* + \sigma_3 \mathbf{u}_3 \mathbf{v}_3^* + \dots + \sigma_r \mathbf{u}_r \mathbf{v}_r^* + \mathbf{0}$$

2 Matrix/Vector Derivatives

2.1 Layout Conventions

The derivative of a vector with respect to a vector, i.e. $\frac{\partial \mathbf{y}}{\partial \mathbf{x}}$, is often written in two competing ways. If the numerator \mathbf{y} is of size m and the denominator \mathbf{x} of size n , then the result can be laid out as either an $m \times n$ matrix or $n \times m$ matrix, i.e., the elements of \mathbf{y} laid out in columns and the elements of \mathbf{x} laid out in rows, or vice versa. This leads to the following possibilities:

1. Numerator layout, i.e. lay out according to \mathbf{y} and \mathbf{x}^T (i.e. contrarily to \mathbf{x}). This is sometimes known as the Jacobian formulation. This corresponds to the $m \times n$ layout.
2. Denominator layout, i.e. lay out according to \mathbf{y}^T and \mathbf{x} (i.e. contrarily to \mathbf{y}). This is sometimes known as the Hessian formulation. Some authors term this layout the gradient, in distinction to the Jacobian (numerator layout), which is its transpose. (However, gradient more commonly means the derivative $\frac{\partial \mathbf{y}}{\partial \mathbf{x}}$, regardless of layout.). This corresponds to the $n \times m$ layout in the previous example.
3. A third possibility sometimes seen is to insist on writing the derivative as $\frac{\partial \mathbf{y}}{\partial \mathbf{x}'}$, (i.e. the derivative is taken with respect to the transpose of \mathbf{x}) and follow the numerator layout. This makes it possible to claim that the matrix is laid out according to both numerator and denominator. In practice this produces results the same as the numerator layout.

2.1.1 Numerator-layout notation

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_1}{\partial x_2} & \cdots & \frac{\partial y_1}{\partial x_n} \\ \frac{\partial y_2}{\partial x_1} & \frac{\partial y_2}{\partial x_2} & \cdots & \frac{\partial y_2}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_m}{\partial x_1} & \frac{\partial y_m}{\partial x_2} & \cdots & \frac{\partial y_m}{\partial x_n} \end{bmatrix}$$

$$\frac{\partial \mathbf{y}}{\partial x} = \begin{bmatrix} \frac{\partial y_1}{\partial x} \\ \frac{\partial y_2}{\partial x} \\ \vdots \\ \frac{\partial y_m}{\partial x} \end{bmatrix}$$

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_1}{\partial x_2} & \cdots & \frac{\partial y_1}{\partial x_n} \\ \frac{\partial y_2}{\partial x_1} & \frac{\partial y_2}{\partial x_2} & \cdots & \frac{\partial y_2}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_m}{\partial x_1} & \frac{\partial y_m}{\partial x_2} & \cdots & \frac{\partial y_m}{\partial x_n} \end{bmatrix}$$

$$\frac{\partial \mathbf{y}}{\partial \mathbf{X}} = \begin{bmatrix} \frac{\partial y}{\partial x_{11}} & \frac{\partial y}{\partial x_{21}} & \cdots & \frac{\partial y}{\partial x_{p1}} \\ \frac{\partial y}{\partial x_{12}} & \frac{\partial y}{\partial x_{22}} & \cdots & \frac{\partial y}{\partial x_{p2}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y}{\partial x_{1q}} & \frac{\partial y}{\partial x_{2q}} & \cdots & \frac{\partial y}{\partial x_{pq}} \end{bmatrix}$$

Notice in the above \mathbf{X} is a $p \times q$ matrix.

In vector calculus, for a scalar valued function $f : R^n \rightarrow R$, $\nabla f = (\frac{\partial f}{\partial \mathbf{x}})^T$ is a column vector, called the **gradient** vector.

Also if \mathbf{y} is a vector-valued function ($R^n \rightarrow R^m$), then $\frac{\partial \mathbf{y}}{\partial \mathbf{x}}$ is a $m \times n$ matrix, called the **Jacobian** matrix.

The following definitions are only provided in numerator-layout notation:

$$\frac{\partial \mathbf{Y}}{\partial x} = \begin{bmatrix} \frac{\partial y_{11}}{\partial x} & \frac{\partial y_{12}}{\partial x} & \dots & \frac{\partial y_{1n}}{\partial x} \\ \frac{\partial y_{21}}{\partial x} & \frac{\partial y_{22}}{\partial x} & \dots & \frac{\partial y_{2n}}{\partial x} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_{m1}}{\partial x} & \frac{\partial y_{m2}}{\partial x} & \dots & \frac{\partial y_{mn}}{\partial x} \end{bmatrix}$$

$$d\mathbf{X} = \begin{bmatrix} dx_{11} & dx_{11} & \dots & dx_{1q} \\ dx_{21} & dx_{22} & \dots & dx_{2q} \\ \vdots & \vdots & \ddots & \vdots \\ dx_{p1} & dx_{p2} & \dots & dx_{pq} \end{bmatrix}$$

2.1.2 Denominator-layout notation

$$\frac{\partial y}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial y}{\partial x_1} \\ \frac{\partial y}{\partial x_2} \\ \vdots \\ \frac{\partial y}{\partial x_n} \end{bmatrix}$$

$$\frac{\partial \mathbf{y}}{\partial x} = \begin{bmatrix} \frac{\partial y_1}{\partial x} & \frac{\partial y_2}{\partial x} & \cdots & \frac{\partial y_m}{\partial x} \end{bmatrix}$$

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_2}{\partial x_1} & \cdots & \frac{\partial y_m}{\partial x_1} \\ \frac{\partial y_1}{\partial x_2} & \frac{\partial y_2}{\partial x_2} & \cdots & \frac{\partial y_m}{\partial x_2} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_1}{\partial x_n} & \frac{\partial y_2}{\partial x_n} & \cdots & \frac{\partial y_m}{\partial x_n} \end{bmatrix}$$

$$\frac{\partial y}{\partial \mathbf{X}} = \begin{bmatrix} \frac{\partial y}{\partial x_{11}} & \frac{\partial y}{\partial x_{12}} & \cdots & \frac{\partial y}{\partial x_{1q}} \\ \frac{\partial y}{\partial x_{21}} & \frac{\partial y}{\partial x_{22}} & \cdots & \frac{\partial y}{\partial x_{2q}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y}{\partial x_{p1}} & \frac{\partial y}{\partial x_{p2}} & \cdots & \frac{\partial y}{\partial x_{pq}} \end{bmatrix}$$

2.2 Useful formulas

1.

$$\frac{d}{d\mathbf{x}}(\mathbf{x}^T \mathbf{x}) = \begin{bmatrix} \frac{\partial}{\partial x_1}(\mathbf{x}^T \mathbf{x}) \\ \frac{\partial}{\partial x_2}(\mathbf{x}^T \mathbf{x}) \\ \vdots \\ \frac{\partial}{\partial x_n}(\mathbf{x}^T \mathbf{x}) \end{bmatrix} = \begin{bmatrix} 2x_1 \\ 2x_2 \\ \vdots \\ 2x_n \end{bmatrix} = 2\mathbf{x}$$

[since, $\frac{\partial}{\partial x_i}(\mathbf{x}^T \mathbf{x}) = \frac{\partial}{\partial x_i}(x_1^2 + x_2^2 + \dots + x_n^2) = 2x_i$]

2.

$$\frac{d}{d\mathbf{x}}(\mathbf{a}^T \mathbf{x}) = \frac{d}{d\mathbf{x}}(a_1 x_1 + a_2 x_2 + \dots + a_n x_n) = \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{bmatrix} = \mathbf{a}$$

3.

$$\frac{d}{d\mathbf{x}}(\mathbf{x}^T \mathbf{a}) = \mathbf{a}^T$$

4.

$$\frac{d}{d\mathbf{x}}(\mathbf{Ax}) = \mathbf{A}$$

5. Assume, \mathbf{A} is real and symmetric ($\mathbf{A}^T = \mathbf{A}$)

$$\begin{aligned}\frac{d}{d\mathbf{x}}(\mathbf{x}^T \mathbf{Ax}) &= 2 \sum_i \sum_j a_{ij} x_i x_j \\ \frac{\partial}{\partial x_i}(\mathbf{x}^T \mathbf{Ax}) &= 2 \sum_j a_{ij} x_j \\ &= 2(\mathbf{Ax})_i \\ \therefore \frac{d}{d\mathbf{x}} &= 2\mathbf{Ax}\end{aligned}$$

Example:

$$\begin{aligned}\Rightarrow [x_1 \quad x_2 \quad x_3] \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} &= \\ \Rightarrow a_{11} x_1^2 + a_{22} x_2^2 + a_{33} x_3^2 + 2 a_{12} x_1 x_2 + 2 a_{13} x_1 x_3 + 2 a_{23} x_2 x_3 & \\ \Rightarrow \frac{\partial}{\partial x_1}(\mathbf{x}^T \mathbf{Ax}) = 2(a_{11} x_1 + a_{12} x_2 + a_{13} x_3) & \\ \Rightarrow \frac{\partial}{\partial x_2}(\mathbf{x}^T \mathbf{Ax}) = 2(a_{21} x_1 + a_{22} x_2 + a_{23} x_3) & \\ \Rightarrow \frac{\partial}{\partial x_3}(\mathbf{x}^T \mathbf{Ax}) = 2(a_{31} x_1 + a_{32} x_2 + a_{33} x_3) &\end{aligned}$$

3 Orthogonal Projection and Least Square Approximation

3.1 The Linear Algebra Way

1. **Projection on a vector:** The orthogonal decomposition of \mathbf{x} on \mathbf{v} means is decomposing \mathbf{x} in the following manner:

$\mathbf{x} = \mathbf{p} + \mathbf{z}$ such that $\mathbf{p} = t\mathbf{v}$ (t is scalar) and $\mathbf{z} \perp \mathbf{v}$.
 $p = \text{proj}_{\mathbf{v}}x$ is called the orthogonal projection of x on v .

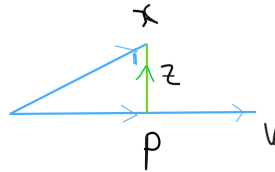
Show,

$$\mathbf{p} = \frac{\mathbf{x}^T \mathbf{v}}{\mathbf{v}^T \mathbf{v}} \mathbf{v}$$

and if \mathbf{v} is an unit vector, then $\mathbf{p} = (\mathbf{x}^T \mathbf{v})\mathbf{v}$.

Note that here we can consider $\mathbf{x}^T \mathbf{v}$ as the coordinate of \mathbf{p} in the space spanned by \mathbf{p} .

Solution:



Let $\mathbf{p} = t\mathbf{v}$, here t is a scalar

$$\mathbf{z} = \mathbf{x} - \mathbf{p} = \mathbf{x} - t\mathbf{v}$$

\mathbf{z} is orthogonal to \mathbf{v} if and only if

$$\begin{aligned} \mathbf{0} &= (\mathbf{x} - t\mathbf{v}) \cdot \mathbf{v} \\ \Rightarrow \mathbf{0} &= \mathbf{x} \cdot \mathbf{v} - (t\mathbf{v} \cdot \mathbf{v}) \\ \Rightarrow \mathbf{0} &= \mathbf{x} \cdot \mathbf{v} - t(\mathbf{v} \cdot \mathbf{v}) \\ \Rightarrow \mathbf{0} &= \mathbf{x}^T \mathbf{v} - t\mathbf{v}^T \mathbf{v} \\ \therefore t &= \frac{\mathbf{x}^T \mathbf{v}}{\mathbf{v}^T \mathbf{v}} \end{aligned}$$

Since $\mathbf{p} = t\mathbf{v}$, we can write

$$\mathbf{p} = \frac{\mathbf{x}^T \mathbf{v}}{\mathbf{v}^T \mathbf{v}} \mathbf{v}$$

If \mathbf{v} is an unit vector, then $\mathbf{v} \cdot \mathbf{v} = 1$

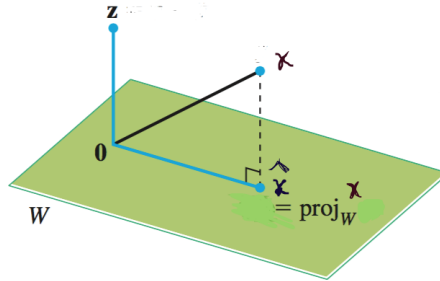
$$\therefore \mathbf{p} = (\mathbf{x}^T \mathbf{v})\mathbf{v}$$

2. **Projection on a subspace:** Projection of \mathbf{x} on a subspace W , such that $\mathbf{x} \notin W$, is given by

$$\mathbf{p} = \text{proj}_W \mathbf{x} = \hat{\mathbf{x}} = VV^T \mathbf{x}$$

here, W is spanned by orthonormal basis set $B = \{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k\}$ and $V = [\mathbf{v}_1 \mathbf{v}_2 \dots \mathbf{v}_k]$, the matrix with the \mathbf{v}_i 's as columns.

VV^T is called the projection matrix.



Solution:

$$\mathbf{x} = \mathbf{p} + \mathbf{z} \text{ s.t., } \mathbf{p} \in W \text{ or, } \hat{\mathbf{x}} \in W$$

$$\text{Then, } \hat{\mathbf{x}} = \mathbf{p} = \alpha_1 \mathbf{v}_1 + \dots + \alpha_k \mathbf{v}_k = V\alpha$$

$$\text{and } \mathbf{z} = \mathbf{x} - V\alpha$$

We also have, $\mathbf{z} \perp W$ which means $\mathbf{z} \perp \mathbf{w}$, for any $\mathbf{w} \in W$. Then $\mathbf{z} \perp \mathbf{v}_i$ because \mathbf{z} is in W^\perp and subspace W is spanned by the orthonormal basis vectors \mathbf{v}_i .

So,

$$\mathbf{z} \perp \mathbf{v}_1 \Rightarrow \mathbf{v}_1^T \mathbf{z} = 0 \Rightarrow \mathbf{v}_1^T (\mathbf{x} - V\alpha) = 0$$

$$\vdots$$

$$\mathbf{z} \perp \mathbf{v}_k \Rightarrow \mathbf{v}_k^T \mathbf{z} = 0 \Rightarrow \mathbf{v}_k^T (\mathbf{x} - V\alpha) = 0$$

Combining these equations, we can write

$$\begin{aligned} & \begin{pmatrix} \mathbf{v}_1^T \\ \mathbf{v}_2^T \\ \vdots \\ \mathbf{v}_k^T \end{pmatrix} (\mathbf{x} - V\alpha) = \mathbf{0} \\ & \Rightarrow V^T (\mathbf{x} - V\alpha) = \mathbf{0} \\ & \Rightarrow V^T \mathbf{x} - V^T V\alpha = \mathbf{0} \\ & \Rightarrow V^T \mathbf{x} - \alpha = \mathbf{0} \end{aligned}$$

Since $\mathbf{V}^T \mathbf{V} = \mathbf{I}$ (Property of orthonormal matrix and here \mathbf{V} is orthonormal matrix)

$$\therefore \boldsymbol{\alpha} = \mathbf{V}^T \mathbf{x}$$

We have $\hat{\mathbf{x}} = \mathbf{V} \boldsymbol{\alpha}$

$$\therefore \hat{\mathbf{x}} = \mathbf{V} \mathbf{V}^T \mathbf{x}$$

Note: Projection of \mathbf{x} in the direction of \mathbf{v}_i , for $i = 1, 2, \dots, k$:

$$proj_{\mathbf{v}_i} \mathbf{x} = (\mathbf{x}^T \mathbf{v}_i) \mathbf{v}_i$$

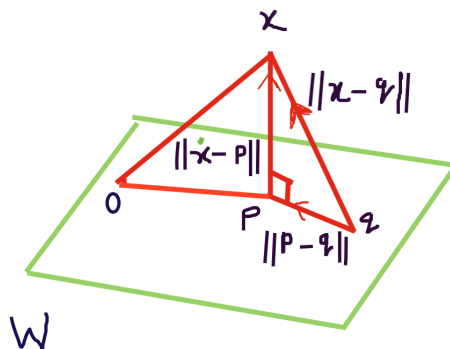
Therefore projection to the subspace can be expressed as,

$$\begin{aligned} \hat{\mathbf{x}} &= proj_W \mathbf{x} = \mathbf{V} \mathbf{V}^T \mathbf{x} \\ &= (\mathbf{x}^T \mathbf{v}_1) \mathbf{v}_1 + (\mathbf{x}^T \mathbf{v}_2) \mathbf{v}_2 + \dots + (\mathbf{x}^T \mathbf{v}_k) \mathbf{v}_k \\ &= proj_{\mathbf{v}_1} \mathbf{x} + proj_{\mathbf{v}_2} \mathbf{x} + \dots + proj_{\mathbf{v}_k} \mathbf{x} \end{aligned}$$

3. **Orthogonal Projection Gives the Best Approximation:** Using Pythagorean Theorem show that, $\|\mathbf{x} - \mathbf{p}\|^2 < \|\mathbf{x} - \mathbf{q}\|^2$, where $\mathbf{p} = proj_W \mathbf{x}$, \mathbf{q} is in W , and $\mathbf{p} \neq \mathbf{q}$. That is, \mathbf{p} is the best approximation of \mathbf{x} in the subspace W .

Solution:

Both \mathbf{p} and \mathbf{q} are in W and distinct from each other. Then $\mathbf{p} - \mathbf{q}$ is in W . $\mathbf{z} = \mathbf{x} - \mathbf{p}$ is



orthogonal to W . In particular, $\mathbf{x} - \mathbf{p}$ is orthogonal to $\mathbf{p} - \mathbf{q}$. Therefore,

$$\mathbf{x} - \mathbf{q} = (\mathbf{x} - \mathbf{p}) + (\mathbf{p} - \mathbf{q})$$

Using Pythagorean Theorem,

$$\|\mathbf{x} - \mathbf{q}\|^2 = \|\mathbf{x} - \mathbf{p}\|^2 + \|\mathbf{p} - \mathbf{q}\|^2$$

Now $\|\mathbf{p} - \mathbf{q}\|^2 > \mathbf{0}$ because $\mathbf{p} - \mathbf{q} \neq \mathbf{0}$ which implies $\|\mathbf{x} - \mathbf{p}\|^2 < \|\mathbf{x} - \mathbf{q}\|^2$

4. Change of Coordinates

$$\begin{aligned} \text{proj}_W \mathbf{x} &= \mathbf{V} \mathbf{V}^T \mathbf{x} \\ &= \begin{bmatrix} | & & | \\ \mathbf{v}_1 & \dots & \mathbf{v}_k \\ | & & | \end{bmatrix} \begin{bmatrix} \mathbf{v}_1^T \\ \mathbf{v}_2^T \\ \vdots \\ \mathbf{v}_k^T \end{bmatrix} \mathbf{x} \\ &= \begin{bmatrix} | & & | \\ \mathbf{v}_1 & \dots & \mathbf{v}_k \\ | & & | \end{bmatrix} \begin{bmatrix} \mathbf{x}^T \mathbf{v}_1 \\ \mathbf{x}^T \mathbf{v}_2 \\ \vdots \\ \mathbf{x}^T \mathbf{v}_k \end{bmatrix} \\ &= \begin{bmatrix} | & & | \\ \mathbf{v}_1 & \dots & \mathbf{v}_k \\ | & & | \end{bmatrix} \begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_k \end{bmatrix} \\ &\therefore [\hat{\mathbf{x}}]_S = V [\hat{\mathbf{x}}]_B \end{aligned}$$

Observe that

- $V^T \mathbf{x} = [\hat{\mathbf{x}}]_B$ gives the coordinates of the projection (using basis B) and
- then the coordinates of the projection are changed from the basis B to S (the standard basis).
- If \mathbf{x} is from a d -dimensional vector space, $[\hat{\mathbf{x}}]_S$ is a $d \times 1$ vector and $[\hat{\mathbf{x}}]_B$ is a $k \times 1$ vector.

3.2 Orthogonal Projection: The Calculus Way

1. **Projection on a vector** Using Calculus find the vector closest (in least square sense) to \mathbf{x} in the direction of \mathbf{v} . In other words, find the least square approximation of \mathbf{x} in the space spanned by \mathbf{v} .

Solution:

A vector \mathbf{p} in the direction of \mathbf{v} is given by $\mathbf{p} = t\mathbf{v}$, where t is a scalar. Therefore, we need to find t that minimizes

$$\begin{aligned} J(t) &= \|\mathbf{x} - t\mathbf{v}\|^2 \\ &= (\mathbf{x} - t\mathbf{v})^T (\mathbf{x} - t\mathbf{v}) = (\mathbf{x}^T - t\mathbf{v}^T)(\mathbf{x} - t\mathbf{v}) \\ &= \mathbf{x}^T \mathbf{x} - t\mathbf{v}^T \mathbf{x} - t\mathbf{x}^T \mathbf{v} + t^2(\mathbf{v}^T \mathbf{v}) \end{aligned}$$

$$\therefore J(t) = \|\mathbf{x}\|^2 - 2t\mathbf{x}^T \mathbf{v} + t^2\|\mathbf{v}\|^2$$

Differentiating the equation with respect to t , we get,

$$J'(t) = -2\mathbf{x}^T \mathbf{v} + 2t\|\mathbf{v}\|^2$$

Setting $J'(t) = 0$ and solving for t we get,

$$-2\mathbf{x}^T \mathbf{v} + 2t\|\mathbf{v}\|^2 = 0$$

$$\therefore t = \frac{\mathbf{x}^T \mathbf{v}}{\|\mathbf{v}\|^2}$$

Therefore, projection of \mathbf{x} on \mathbf{v} defined as \mathbf{p} can be written as,

$$\mathbf{p} = \frac{\mathbf{x}^T \mathbf{v}}{\|\mathbf{v}\|^2} \mathbf{v}$$

2. **Projection on a subspace** We seek the closest approximation of vector \mathbf{x} in the subspace \mathbf{W} which has dimension k . Assume, \mathbf{v}_i 's for $i = 1, 2, \dots, k$ form an orthonormal basis for \mathbf{W} . Find the α_i 's for $i = 1, 2, \dots, k$, s.t. the error J , given by

$$J = \|\mathbf{x} - \hat{\mathbf{x}}\|^2 = \left\| \mathbf{x} - \sum_{i=1}^k \alpha_i \mathbf{v}_i \right\|^2$$

is minimized.

Solution:

Any vector in \mathbf{W} can be written as $\sum_{i=1}^k \alpha_i \mathbf{v}_i$. Thus, \mathbf{x} will be represented by some vector in \mathbf{W} as $\sum_{i=1}^k \alpha_i \mathbf{v}_i$. To minimize the error J we need to take partial derivatives.

$$\begin{aligned}
J(\alpha_1, \dots, \alpha_k) &= \left\| \mathbf{x} - \sum_{i=1}^k \alpha_i \mathbf{v}_i \right\|^2 \\
&= \left(\mathbf{x} - \sum_{i=1}^k \alpha_i \mathbf{v}_i \right)^T \left(\mathbf{x} - \sum_{i=1}^k \alpha_i \mathbf{v}_i \right) \\
&= \left(\mathbf{x}^T - \sum_{i=1}^k \alpha_i \mathbf{v}_i^T \right) \left(\mathbf{x} - \sum_{i=1}^k \alpha_i \mathbf{v}_i \right) \\
&= \mathbf{x}^T \mathbf{x} - \mathbf{x} \sum_{i=1}^k \alpha_i \mathbf{v}_i^T - \mathbf{x}^T \sum_{i=1}^k \alpha_i \mathbf{v}_i + \left(\sum_{i=1}^k \alpha_i \mathbf{v}_i^T \right) \left(\sum_{i=1}^k \alpha_i \mathbf{v}_i \right) \\
&= \|\mathbf{x}\|^2 - 2 \sum_{i=1}^k \alpha_i \mathbf{x}^T \mathbf{v}_i + \sum_{i=1}^k \alpha_i^2 \left[\because \|\mathbf{v}_i\|^2 = 1 \text{ and } \mathbf{v}_i \text{'s are orthogonal} \right]
\end{aligned}$$

Then we take partial derivative with respect to α_i and set that to 0 for optimal value. We get,

$$\begin{aligned}
-2\mathbf{x}^T \mathbf{v}_i + 2\alpha_i &= 0 \\
\therefore \alpha_i &= \mathbf{x}^T \mathbf{v}_i
\end{aligned}$$

3.3 Ordinary Least Square Regression

Suppose we have n datapoints $(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(n)}, y^{(n)})$, where $\mathbf{x}^{(i)}$'s are d -vectors $[\mathbf{x}_1^{(i)}, \mathbf{x}_2^{(i)}, \dots, \mathbf{x}_d^{(i)}]$ (*i.e.*, x contains the values of d features or (independent) variables) and y is a real number (called the dependent variable).

We assume there is a function $f(\mathbf{x})$ such that $y = f(\mathbf{x})$. In linear regression, based on the data, we want to find a linear (affine) function \hat{f} that approximates f (in the least square sense).

Let,

$$\begin{aligned}
\hat{y} = \hat{f}(\mathbf{x}) &= \beta_0 + \beta_1 x_1 + \dots + \beta_d x_d \\
&= \beta_0 x_0 + \beta_1 x_1 + \dots + \beta_d x_d \quad [let, x_0 = 1]
\end{aligned}$$

That is, for each of the n datapoints $\mathbf{x}^{(i)}$:

$$\begin{aligned}
\beta_0 x_0^{(1)} + \beta_1 x_1^{(1)} + \dots + \beta_d x_d^{(1)} &= \hat{y}^{(1)} \\
\beta_0 x_0^{(2)} + \beta_1 x_1^{(2)} + \dots + \beta_d x_d^{(2)} &= \hat{y}^{(2)} \\
&\vdots \\
\beta_0 x_0^{(n)} + \beta_1 x_1^{(n)} + \dots + \beta_d x_d^{(n)} &= \hat{y}^{(n)}
\end{aligned}$$

$$\begin{aligned}
\begin{bmatrix} x_0^{(1)} & x_1^{(1)} & \dots & x_d^{(1)} \\ x_0^{(2)} & x_1^{(2)} & \dots & x_d^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ x_0^{(n)} & x_1^{(n)} & \dots & x_d^{(n)} \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_d \end{bmatrix} &= \begin{bmatrix} \hat{y}^{(1)} \\ \hat{y}^{(2)} \\ \vdots \\ \hat{y}^{(N)} \end{bmatrix} \\
\begin{bmatrix} 1 & x_1^{(1)} & \dots & x_d^{(1)} \\ 1 & x_1^{(2)} & \dots & x_d^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(N)} & \dots & x_d^{(N)} \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_d \end{bmatrix} &= \begin{bmatrix} \hat{y}^{(1)} \\ \hat{y}^{(2)} \\ \vdots \\ \hat{y}^{(N)} \end{bmatrix} \\
&\implies \mathbf{X}\boldsymbol{\beta} = \hat{\mathbf{y}}
\end{aligned}$$

That means, we need to find the $\beta_0, \beta_1, \dots, \beta_d$ coefficients that satisfy the above equations. We can view this problem as finding the solution to the system of linear equations ($\mathbf{X}\boldsymbol{\beta} = \hat{\mathbf{y}}$). However, this is an overdetermined system (more equations (or rows) than variables (or columns)). Therefore, we can only find the best $\boldsymbol{\beta}$ that approximately solves the system of linear equations.

Or, we can also view this as an optimization problem and find the $\boldsymbol{\beta}$ that minimizes the Mean Squared Error (MSE):

$$\frac{1}{n} \sum_{i=1}^N (y^{(i)} - \hat{y}^{(i)})^2 = \frac{1}{N} \|\mathbf{y} - \hat{\mathbf{y}}\|^2 = \frac{1}{N} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2$$

Note: Usually in linear regression the features or independent variables are transformed to create a new set of variables. This can be done through basis functions $\phi_j(\mathbf{x})$ that transform the data and create a datapoint in the transformed feature space, $z_j = \phi_j(\mathbf{x})$, $j = 1, \dots, p$. And then we do linear regression using the transformed datapoints. Note that the basis functions can be non-linear.

$$\begin{aligned}
\hat{f}(\mathbf{x}) &= \beta_0 + \beta_1 z_1 + \dots + \beta_p z_p \\
&= \beta_0 \phi_0(\mathbf{x}) + \beta_1 \phi_1(\mathbf{x}) + \dots + \beta_p \phi_p(\mathbf{x}) \quad [\phi_0(\mathbf{x}) = 1]
\end{aligned}$$

Again, the MSE can be written as

$$\frac{1}{n} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2$$

where,

$$\mathbf{X} = \begin{bmatrix} \phi_0(x^{(1)}) & \phi_1(x^{(1)}) & \dots & \phi_p(x^{(1)}) \\ \phi_0(x^{(2)}) & \phi_1(x^{(2)}) & \dots & \phi_p(x^{(2)}) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_0(x^{(N)}) & \phi_1(x^{(N)}) & \dots & \phi_p(x^{(N)}) \end{bmatrix}, \quad \boldsymbol{\beta} = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_N \end{bmatrix}$$

3.3.1 The Linear Algebra Way

Suppose we have a system of linear equations $A\boldsymbol{\beta} = \mathbf{y}$, where A is a $n \times d$ matrix.

$$\text{Here, } A = \begin{bmatrix} | & & | \\ \mathbf{a}_1 & \dots & \mathbf{a}_d \\ | & & | \end{bmatrix}, \boldsymbol{\beta} = \begin{bmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_d \end{bmatrix}, \text{ and } A\boldsymbol{\beta} = \beta_1\mathbf{a}_1 + \beta_2\mathbf{a}_2 + \dots + \beta_d\mathbf{a}_d.$$

If the system of linear equations does not have a solution, $\mathbf{y} \neq \beta_1\mathbf{a}_1 + \beta_2\mathbf{a}_2 + \dots + \beta_k\mathbf{a}_d$, i.e., $\mathbf{y} \notin \text{Col}(A) = \text{span}\{\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_d\}$.

Then least square solution is $\hat{\boldsymbol{\beta}}$, such that $\|\mathbf{y} - A\hat{\boldsymbol{\beta}}\|^2$ is minimum.

We observe, what we are asking for are the coordinates of $\hat{\mathbf{y}} = \text{proj}_{\text{Col}(A)}\mathbf{y}$. Now, we may not have an orthonormal basis of $\text{Col}(A)$, that is columns of A might not be orthonormal. Rather we have a basis $B' = \{\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_d\}$ of the $\text{Col}(A)$, assuming the columns of A are linearly independent.

Show,

$$\hat{\boldsymbol{\beta}} = (A^T A)^{-1} A^T \mathbf{y}$$

$$\hat{\mathbf{y}} = \text{proj}_{\text{Col}(A)}\mathbf{y} = A\hat{\boldsymbol{\beta}} = A(A^T A)^{-1} A^T \mathbf{y}$$

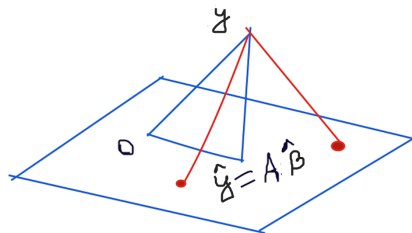
Note, $A(A^T A)^{-1} A^T$ is called the Projection matrix and $A^\dagger = (A^T A)^{-1} A^T$ is called the pseudo-inverse matrix.

Solution:

When a solution is demanded and none exists, in this scenario what we can do is to find a $\boldsymbol{\beta}$ that makes $A\boldsymbol{\beta}$ as close as possible to \mathbf{y} .

Here, A is $n \times d$ and \mathbf{y} is in \mathbb{R}^n .

Let $\hat{\mathbf{y}} = \text{proj}_{\text{Col}(A)}\mathbf{y}$



Because $\hat{\mathbf{y}}$ is in the column space of A , the equation $A\boldsymbol{\beta} = \hat{\mathbf{y}}$ is consistent and there is a $\hat{\boldsymbol{\beta}}$ in \mathbb{R}^k such that,

$$A\hat{\boldsymbol{\beta}} = \hat{\mathbf{y}}$$

The projection $\hat{\mathbf{y}}$ has the property that $\mathbf{y} - \hat{\mathbf{y}}$ is orthogonal to $Col(A)$, so $(\mathbf{y} - A\hat{\boldsymbol{\beta}})$ is orthogonal to each column of A .

If \mathbf{a}_j is any column of A , then $\mathbf{a}_j \cdot (\mathbf{y} - A\hat{\boldsymbol{\beta}}) = 0$ or $\mathbf{a}_j^T (\mathbf{y} - A\hat{\boldsymbol{\beta}}) = 0$.

Since each \mathbf{a}_j^T is a row of A^T ,

$$\begin{aligned} A^T(\mathbf{y} - A\hat{\boldsymbol{\beta}}) &= \mathbf{0} \\ \Rightarrow A^T\mathbf{y} - A^T A\hat{\boldsymbol{\beta}} &= \mathbf{0} \\ \Rightarrow A^T A\hat{\boldsymbol{\beta}} &= A^T\mathbf{y} \text{ [These are called the Normal Equations]} \\ \therefore \hat{\boldsymbol{\beta}} &= (A^T A)^{-1} A^T\mathbf{y} \end{aligned}$$

Finally we have, $\hat{\mathbf{y}} = \text{proj}_{Col(A)}\mathbf{y} = A\hat{\boldsymbol{\beta}} = A(A^T A)^{-1} A^T\mathbf{y}$

- Using QR decomposition:

Alternatively, $\hat{\mathbf{y}} = \mathbf{Q}\mathbf{Q}^T\mathbf{y}$ where $A = \mathbf{Q}\mathbf{R}$ is the QR decomposition of A

Here, the columns of \mathbf{Q} form an orthonormal basis for $Col(A)$ and \mathbf{R} is an upper triangular invertible matrix.

$$\begin{aligned} A\hat{\boldsymbol{\beta}} &= \hat{\mathbf{y}} \\ \Rightarrow \mathbf{Q}\mathbf{R}\hat{\boldsymbol{\beta}} &= \mathbf{Q}\mathbf{Q}^T\mathbf{y} \\ \therefore \hat{\boldsymbol{\beta}} &= \mathbf{R}^{-1}\mathbf{Q}^T\mathbf{y} \end{aligned}$$

Note, the pseudo-inverse of A , $A^\dagger = \mathbf{R}^{-1}\mathbf{Q}^T$

3.3.2 The Calculus Way

Using the derivative rules, as outlined in Section 2.2, we can use calculus to find the $\boldsymbol{\beta}$:

$$\begin{aligned} \|\mathbf{y} - \mathbf{A}\boldsymbol{\beta}\|^2 &= (\mathbf{y} - \mathbf{A}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{A}\boldsymbol{\beta}) \\ &= \|\mathbf{y}\|^2 - 2\mathbf{y}^T \mathbf{A}\boldsymbol{\beta} + \boldsymbol{\beta}^T \mathbf{A}^T \mathbf{A}\boldsymbol{\beta} \\ \frac{\partial}{\partial \boldsymbol{\beta}} \|\mathbf{y} - \mathbf{A}\boldsymbol{\beta}\|^2 &= 0 - 2\mathbf{A}^T \mathbf{y} + 2\mathbf{A}^T \mathbf{A}\boldsymbol{\beta} \\ \therefore \frac{\partial}{\partial \boldsymbol{\beta}} \|\mathbf{y} - \mathbf{A}\boldsymbol{\beta}\|^2 &= 0 \\ \Rightarrow \mathbf{A}^T \mathbf{A}\boldsymbol{\beta} &= \mathbf{A}^T \mathbf{y} \\ \therefore \hat{\boldsymbol{\beta}} &= (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{y} \end{aligned}$$

4 Principal Component Analysis

In the previous section, we were given a k dimensional subspace W or more specifically a basis of the subspace and we were asked to find a vector $\hat{\mathbf{x}}$ in that subspace that is the least square approximation (minimum reconstruction error) of a vector $\mathbf{x} \notin W$. We found that the $\hat{\mathbf{x}} = \text{proj}_W \mathbf{x}$, i.e., the orthogonal projection of \mathbf{x} on W .

In Principal Component Analysis (PCA) our objective is to find "best" the subspace or the orthonormal basis of the subspace for a given set of datapoints. PCA can be interpreted in the following two equivalent ways:

1. Least square reconstruction error minimization
2. Variance maximization

4.1 PCA as Least Square Reconstruction Error Minimization

In the least square reconstruction error minimization formulation, our goal is to find \mathbf{V} , where columns of V are orthonormal and they form the orthonormal basis of W such that, the reconstruction error/ projection error

$$J = \sum_{j=1}^n \|\mathbf{x}_j - \text{proj}_W \mathbf{x}_j\|^2$$

is minimum. Here, the data points $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ are d -vectors (i.e., there are d features).

Solution to PCA problem: using summation notation

Any vector in \mathbf{W} can be written as $\sum_{i=1}^k \alpha_i \mathbf{v}_i$. Thus \mathbf{x}_1 will be represented by some vector in \mathbf{W} as $\sum_{i=1}^k \alpha_{1i} \mathbf{v}_i$.

To minimize the error J , over all n datapoints, we need to take partial derivatives and enforce constraint that $\{v_1, \dots, v_k\}$ are orthogonal.

$$\begin{aligned} J(v_1, \dots, v_k, \alpha_{11}, \dots, \alpha_{nk}) &= \sum_{j=1}^n \left\| \mathbf{x}_j - \sum_{i=1}^k \alpha_{ji} \mathbf{v}_i \right\|^2 \\ &= \sum_{j=1}^n (\mathbf{x}_j - \sum_{i=1}^k \alpha_{ji} \mathbf{v}_i)^T (\mathbf{x}_j - \sum_{i=1}^k \alpha_{ji} \mathbf{v}_i) \\ &= \sum_{j=1}^n (\mathbf{x}_j^T - \sum_{i=1}^k \alpha_{ji} \mathbf{v}_i^T) (\mathbf{x}_j - \sum_{i=1}^k \alpha_{ji} \mathbf{v}_i) \\ &= \sum_{j=1}^n \left\{ \mathbf{x}_j^T \mathbf{x}_j - \mathbf{x}_j^T \sum_{i=1}^k \alpha_{ji} \mathbf{v}_i - \sum_{i=1}^k \alpha_{ji} \mathbf{v}_i^T \mathbf{x}_j + \sum_{i=1}^k \alpha_{ji} \mathbf{v}_i^T \sum_{i=1}^k \alpha_{ji} \mathbf{v}_i \right\} \\ &= \sum_{j=1}^n \|\mathbf{x}_j\|^2 - 2 \sum_{j=1}^n \sum_{i=1}^k \alpha_{ji} \mathbf{x}_j^T \mathbf{v}_i + \sum_{j=1}^n \sum_{i=1}^k \alpha_{ji}^2 \left[\because v_i \text{'s are orthonormal} \right] \end{aligned}$$

Then we take partial derivative with respect to some α_{ml} and set that to 0 for optimal value. We get,

$$-2\mathbf{x}_m^T \mathbf{v}_l + 2\alpha_{ml} = \mathbf{0}$$

$$\therefore \alpha_{ml} = \mathbf{x}_m^T \mathbf{v}_l$$

(But we already know this from the idea of orthogonal projection from the previous section)
Now we have to plug the optimal value for $\alpha_{ml} = \mathbf{x}_m^T \mathbf{v}_l$ back into J .

$$\begin{aligned} J(v_1, \dots, v_k) &= \sum_{j=1}^n \|\mathbf{x}_j\|^2 - 2 \sum_{j=1}^n \sum_{i=1}^k (\mathbf{x}_j^T \mathbf{v}_i) \mathbf{x}_j^T \mathbf{v}_i + \sum_{j=1}^n \sum_{i=1}^k (\mathbf{x}_j^T \mathbf{v}_i)^2 \\ &= \sum_{j=1}^n \|\mathbf{x}_j\|^2 - \sum_{j=1}^n \sum_{i=1}^k (\mathbf{x}_j^T \mathbf{v}_i)^2 \\ &= \text{const} - \sum_{i=1}^k \mathbf{v}_i^T \left(\sum_{j=1}^n \mathbf{x}_j \mathbf{x}_j^T \right) \mathbf{v}_i \\ &= \text{const} - \sum_{i=1}^k \mathbf{v}_i^T S \mathbf{v}_i \end{aligned}$$

In matrix notation:

Assume X is a $d \times n$ data matrix, where each d -dimensional datapoint \mathbf{x}_i stored as the columns of X . Then,

$$\begin{aligned} \sum_j \|\mathbf{x}_j - \text{proj}_W \mathbf{x}_j\|^2 &= \|X - VV^T X\|_{fro}^2 \\ &= \text{tr}((X - VV^T X)^T (X - VV^T X)) \quad \because \|A\|_{fro}^2 = \text{tr}(A^T A) \\ &= \text{tr}(X^T X - X^T VV^T X - X^T VV^T X + X^T VV^T VV^T X) \\ &= \text{tr}(X^T X - 2X^T VV^T X + X^T VV^T X) [\because V^T V = I] \\ &= \text{tr}(X^T X) - \text{tr}(X^T VV^T X) \\ &= \text{const} - \text{tr}(X^T VV^T X) \end{aligned}$$

Using the **cyclic property of trace**: $\text{tr}(ABC) = \text{tr}(CAB) = \text{tr}(BCA)$ we get,

$$\begin{aligned} \text{tr}(X^T VV^T X) &= \text{tr}(X X^T VV^T) \\ &= \text{tr}(V^T X X^T V) \\ &= \text{tr}(V^T S V) \\ &= \sum_i (v_i^T S v_i) \\ \therefore \sum_j \|\mathbf{x}_j - \text{proj}_W \mathbf{x}_j\|^2 &= \text{const} - \sum_i (v_i^T S v_i) \end{aligned}$$

Note: Minimizing the reconstruction error is equivalent to maximizing the sum of the quadratic forms $\mathbf{v}_i^T S \mathbf{v}_i$ where S is the Scatter matrix, if the data points are mean centered (more on this below). We will discuss about maximizing a quadratic form $Q(\mathbf{v}) = \mathbf{v}^T S \mathbf{v}$ in the next section.

4.2 Covariance and Scatter Matrix

The matrix S has a nice interpretation. Observe, each row i of the matrix X is a feature vector \mathbf{f}_i . If we subtract the row means (feature means, $\boldsymbol{\mu}_i$'s) of X from each element of the rows then we have made sure that each \mathbf{f}_i has zero mean. This is sometimes called centering the data. We will call mean subtracted data matrix, the demeaned data matrix \tilde{X} .

$$\begin{aligned}\tilde{X} &= X - \begin{pmatrix} -\boldsymbol{\mu}_1- \\ -\boldsymbol{\mu}_2- \\ \vdots \\ -\boldsymbol{\mu}_d- \end{pmatrix} = \begin{pmatrix} -\mathbf{f}_1- \\ -\mathbf{f}_2- \\ \vdots \\ -\mathbf{f}_d- \end{pmatrix} - \begin{pmatrix} -\boldsymbol{\mu}_1- \\ -\boldsymbol{\mu}_2- \\ \vdots \\ -\boldsymbol{\mu}_d- \end{pmatrix} \\ &= \begin{bmatrix} f_{11} - \mathbf{avg}(\mathbf{f}_1) & f_{12} - \mathbf{avg}(\mathbf{f}_1) & \dots & f_{1n} - \mathbf{avg}(\mathbf{f}_1) \\ f_{21} - \mathbf{avg}(\mathbf{f}_2) & f_{22} - \mathbf{avg}(\mathbf{f}_2) & \dots & f_{2n} - \mathbf{avg}(\mathbf{f}_2) \\ \vdots & \vdots & \ddots & \vdots \\ f_{d1} - \mathbf{avg}(\mathbf{f}_d) & f_{d2} - \mathbf{avg}(\mathbf{f}_d) & \dots & f_{dn} - \mathbf{avg}(\mathbf{f}_d) \end{bmatrix} \\ \tilde{X}^T &= \begin{bmatrix} f_{11} - \mathbf{avg}(\mathbf{f}_1) & f_{21} - \mathbf{avg}(\mathbf{f}_2) & \dots & f_{d1} - \mathbf{avg}(\mathbf{f}_d) \\ f_{12} - \mathbf{avg}(\mathbf{f}_1) & f_{22} - \mathbf{avg}(\mathbf{f}_2) & \dots & f_{d2} - \mathbf{avg}(\mathbf{f}_d) \\ \vdots & \vdots & \ddots & \vdots \\ f_{1n} - \mathbf{avg}(\mathbf{f}_1) & f_{2n} - \mathbf{avg}(\mathbf{f}_2) & \dots & f_{dn} - \mathbf{avg}(\mathbf{f}_d) \end{bmatrix}\end{aligned}$$

We can observe the ij -th element of $\tilde{X}\tilde{X}^T$ is,

$$\sum_{k=1}^n (f_{ik} - \mathbf{avg}(\mathbf{f}_i))(f_{jk} - \mathbf{avg}(\mathbf{f}_j)) = (n-1)s_{ij}$$

where, $s_{ij} = \mathit{cov}(f_i, f_j)$ is the sample covariance between the features f_i and f_j . Also observe, when $i = j$, $s_{ii} = \mathit{cov}(f_i, f_i) = \mathit{var}(f_i)$.

$$\begin{aligned}\therefore \tilde{X}\tilde{X}^T &= (n-1)\hat{\Sigma} \\ &= (n-1) \begin{bmatrix} s_{11}^2 & s_{12} & \dots & s_{1d} \\ s_{21} & s_{22}^2 & \dots & s_{2d} \\ \vdots & & \ddots & \\ s_{d1} & s_{d2} & \dots & s_{dd}^2 \end{bmatrix} \\ &= S\end{aligned}$$

Here, $\hat{\Sigma}$ is the Sample Covariance matrix, and S is the Scatter matrix.

Observation 1: What happens when data is stored as rows in X ?

$$X = \underbrace{\begin{pmatrix} -x_1- \\ -x_2- \\ \vdots \\ -x_n- \end{pmatrix}}_{n \times d} = \begin{pmatrix} | & | & & | \\ f_1 & f_2 & \dots & f_d \\ | & | & & | \end{pmatrix}$$

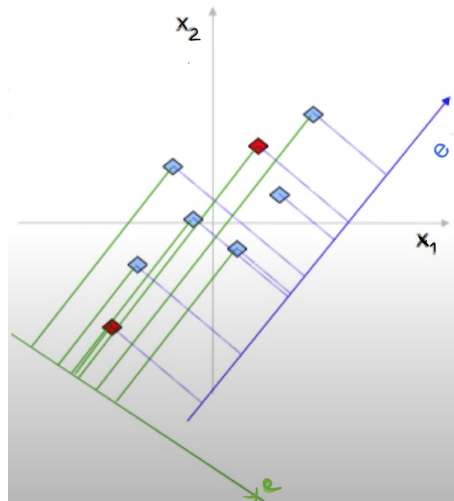
$$\therefore \tilde{X}^T \tilde{X} = (n-1) \hat{\Sigma}$$

We can easily show that we end up with the same maximization problem.

Observation 2: S is a real symmetric matrix.

4.3 PCA as Variance Maximization

We can think of the PCA as projecting \mathbf{x} into the subspace of dimension k so that we can capture maximum variance. That is, instead of thinking about reconstruction error, we can simply restrict attention to directions where scatter or variability of the data is the greatest.



Suppose we want to embed two dimensional points in a one dimensional space, a line e . We can see from the figure that the projection points are very spread out along the *blue* direction and they are very bunched up along the *green* direction. So the variance along the blue dimension is higher than the variance of the green one.

So why is good to have a high variance along the projection?

If we look at the two red points, they are pretty far away from each other in the 2-dimensional space. If we happen to project them to the green dimension, they will end up being on top of one another. So this dimension does not preserve distance of the original space. If points are far away in original space, we want them to remain that way in the lower dimensional space. But if we project them to blue line, they will stay apart from each other. (Though there will be some points which will end up being close when we project them. If we can choose the direction with maximum variability of the datapoints, we can reduce the number of such points).

How to find the direction of maximum variance?

$$\frac{1}{N} \sum_j (\mathbf{v}^T \mathbf{x}_j)^2 = \mathbf{v}^T \mathbf{X} \mathbf{X}^T \mathbf{v} = \mathbf{v}^T \mathbf{S} \mathbf{v}$$

From this, we can observe that the variance will be maximum when we project in the direction \mathbf{v} in which the quadratic form is maximized.

5 Maximization of Quadratic Form

There are two ways to maximize (or minimize) a quadratic form $Q(\mathbf{v}_i) = \mathbf{v}_i^T S \mathbf{v}_i$ subject to the constraint $\|\mathbf{v}_i\|^2 = \mathbf{1}$:

1. Calculus way: using the method of Lagrange Multipliers
2. Linear Algebra way: using Diagonalization method

5.1 Calculus way: Method of Lagrange Multipliers

Here, we will consider the sum of quadratic forms that we want to maximize in PCA. We will enforce constraints $\mathbf{v}_i^T \mathbf{v}_i = 1$ for all i and incorporate the constraints with undetermined $\lambda_1, \dots, \lambda_k$.

Now we will need to maximize a new function \hat{J} .

$$\hat{J}(\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k) = \sum_{i=1}^k \mathbf{v}_i^T S \mathbf{v}_i - \sum_{j=1}^k \lambda_j (\mathbf{v}_j^T \mathbf{v}_j - 1)$$

Computing the partial derivative with respect to \mathbf{v}_m ,

$$\frac{\partial \hat{J}(\mathbf{v}_1, \dots, \mathbf{v}_k)}{\partial \mathbf{v}_m} = 2S\mathbf{v}_m - 2\lambda_m \mathbf{v}_m = 0$$
$$S\mathbf{v}_m = \lambda_m \mathbf{v}_m$$

Therefore, λ_m 's and \mathbf{v}_m 's must be the eigenvalues and the corresponding eigenvectors of the Scatter matrix S .

5.2 Linear algebra way

In this section we deal with maximizing a single quadratic form. The maximization of the sum of quadratic form in PCA will follow from the solution to this problem. The problem of maximizing a quadratic form sometimes appears in engineering literature as maximizing the Rayleigh quotient.

Rayleigh Quotient problem: For a fixed symmetric matrix \mathbf{A} , the normalized quadratic form

$$\max_{\mathbf{x} \neq 0 \in \mathbb{R}^n} \frac{Q(\mathbf{x})}{\|\mathbf{x}\|^2} = \frac{\mathbf{x}^T \mathbf{A} \mathbf{x}}{\mathbf{x}^T \mathbf{x}}$$

is called Rayleigh Quotient.

Since the quotient is scaling invariant, we can write:

$$\max_{\mathbf{x} \in \mathbb{R}^n: \|\mathbf{x}\|=1} \mathbf{x}^T \mathbf{A} \mathbf{x}$$

Explanation: Let's consider $\mathbf{x} = c\mathbf{x}$. Then we get

$$\max_{c\mathbf{x} \neq 0 \in \mathbb{R}^n} \frac{(c\mathbf{x})^T \mathbf{A} (c\mathbf{x})}{(c\mathbf{x})^T (c\mathbf{x})} = \frac{c^2 \mathbf{x}^T \mathbf{A} \mathbf{x}}{c^2 \mathbf{x}^T \mathbf{x}} = \frac{\mathbf{x}^T \mathbf{A} \mathbf{x}}{\mathbf{x}^T \mathbf{x}}$$

It implies that we will get the same maximum value even when we choose a scalar multiple of \mathbf{x} . That is why we will choose to use unit vector, i.e., $\|\mathbf{x}\| = 1$.

At this point we can see that our problem: $\max \mathbf{v}^T \mathbf{S} \mathbf{v}$ subject to $\|\mathbf{v}\|^2 = 1$ is exactly the Rayleigh Quotient problem with the real symmetric matrix S .

Observe if S were a diagonal matrix, then solving this maximization problem is easy. For instance, consider the problem of finding the maximum value of

$$Q(\mathbf{x}) = 9x_1^2 + 4x_2^2 + 3x_3^2$$

subject to the constraint $\|\mathbf{x}\| = 1$.

In matrix form: $Q = [x_1 \ x_2 \ x_3] \begin{bmatrix} 9 & 0 & 0 \\ 0 & 4 & 0 \\ 0 & 0 & 3 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$

Solution: Since x_2^2 and x_3^2 are nonnegative, note that $4x_2^2 \leq 9x_2^2$ and $3x_3^2 \leq 9x_3^2$ and hence

$$\begin{aligned} Q(\mathbf{x}) &= 9x_1^2 + 4x_2^2 + 3x_3^2 \\ &\leq 9x_1^2 + 9x_2^2 + 9x_3^2 \\ &= 9(x_1^2 + x_2^2 + x_3^2) \\ &= 9 \end{aligned}$$

whenever $x_1^2 + x_2^2 + x_3^2 = 1$

So the maximum value of $Q(x)$ can not exceed 9 when x is a unit vector.

$\therefore Q(\mathbf{x}) = 9$ when $\mathbf{x} = (1, 0, 0)$

Thus 9 is the maximum value of $Q(\mathbf{x})$ for $\mathbf{x}^T \mathbf{x} = 1$.

Therefore, one way to solve the general quadratic form optimization problem will be to orthogonally diagonalize \mathbf{S} by changing of variable. In the following, we describe how to orthogonally diagonalize a matrix. We will show that for a real symmetric matrix such a diagonalization always exists.

5.2.1 Diagonalization, Similarity Transformation, and Eigendecomposition

Similar Matrices: If \mathbf{A} and \mathbf{B} are $n \times n$ matrices, then \mathbf{A} is similar to \mathbf{B} if there is an invertible matrix \mathbf{P} such that $\mathbf{P}^{-1} \mathbf{A} \mathbf{P} = \mathbf{B}$ or equivalently, $\mathbf{A} = \mathbf{P} \mathbf{B} \mathbf{P}^{-1}$.

Note: if \mathbf{A} and \mathbf{B} are **similar**, then they have the **same characteristic polynomial and hence the same eigenvalues**. They also have same number of independent eigenvectors. To show the claim:

First, if λ is an eigenvalue of \mathbf{A} , then $\mathbf{A} \mathbf{x} = \lambda \mathbf{x}$ and since \mathbf{A} and \mathbf{B} are similar, $\mathbf{P} \mathbf{B} \mathbf{P}^{-1} = \mathbf{A}$. so,

$$\begin{aligned} \mathbf{P} \mathbf{B} \mathbf{P}^{-1} \mathbf{x} &= \lambda \mathbf{x} \\ \mathbf{B}(\mathbf{P}^{-1} \mathbf{x}) &= \lambda \mathbf{P}^{-1} \mathbf{x} \end{aligned}$$

which is the eigenvalue problem formulation for \mathbf{B} . Therefore, we find, though the eigenvectors $\mathbf{P}^{-1} \mathbf{x} \neq \mathbf{x}$ but the eigenvalue λ is the same for \mathbf{A} and \mathbf{B} .

Another way to see this,

$$\begin{aligned}
\det(\mathbf{B} - \lambda\mathbf{I}) &= \det(\mathbf{P}^{-1}\mathbf{A}\mathbf{P} - \lambda\mathbf{I}) \\
&= \det(\mathbf{P}^{-1}\mathbf{A}\mathbf{P} - \lambda\mathbf{P}^{-1}\mathbf{P}) \\
&= \det(\mathbf{P}^{-1}(\mathbf{A} - \lambda\mathbf{I})\mathbf{P}) \\
&= \det(\mathbf{P}^{-1})\det(\mathbf{A} - \lambda\mathbf{I})\det(\mathbf{P}) \\
&= \det(\mathbf{P})^{-1}\det(\mathbf{A} - \lambda\mathbf{I})\det(\mathbf{P}) \\
&= \det(\mathbf{A} - \lambda\mathbf{I})
\end{aligned}$$

So, the matrices have the same characteristic equation.

Diagonalizable matrix: An $n \times n$ matrix \mathbf{A} is diagonalizable if and only if \mathbf{A} has n linearly independent eigenvectors. In fact, $\mathbf{A} = \mathbf{P}\mathbf{D}\mathbf{P}^{-1}$, with \mathbf{D} a diagonal matrix, if and only if the columns of \mathbf{P} are linearly independent eigenvectors of \mathbf{A} . In this case, the diagonal entries of \mathbf{D} are eigenvalues of \mathbf{A} that correspond, respectively, to the eigenvectors in \mathbf{P} .

Diagonalization Theorem

\mathbf{A} is diagonalizable iff \mathbf{A} has n linearly independent eigenvectors. In this case, we may construct \mathbf{P} by stacking the n eigenvectors and \mathbf{D} as a diagonal matrix with the corresponding eigenvalues.

Proof:

Consider the columns of $\mathbf{P} = (\mathbf{p}_1 \ \mathbf{p}_2 \ \dots \ \mathbf{p}_n)$ and $\mathbf{D} = \begin{pmatrix} d_1 & 0 & \dots & 0 \\ 0 & d_2 & \dots & 0 \\ \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & d_n \end{pmatrix}$

Let's assume that $\mathbf{A} = \mathbf{P}\mathbf{D}\mathbf{P}^{-1}$ and we multiply by \mathbf{P} on the right

$$\begin{aligned}
\mathbf{A}\mathbf{P} &= \mathbf{P}\mathbf{D} \\
(\mathbf{A}\mathbf{p}_1 \ \mathbf{A}\mathbf{p}_2 \ \dots \ \mathbf{A}\mathbf{p}_n) &= (d_1\mathbf{p}_1 \ d_2\mathbf{p}_2 \ \dots \ d_n\mathbf{p}_n)
\end{aligned}$$

This implies that

$$\begin{aligned}
\mathbf{A}\mathbf{p}_1 &= d_1\mathbf{p}_1 \\
\mathbf{A}\mathbf{p}_2 &= d_2\mathbf{p}_2 \\
&\dots \\
\mathbf{A}\mathbf{p}_n &= d_n\mathbf{p}_n
\end{aligned}$$

But this is the definition of eigenvector, so all the columns \mathbf{p}_i in \mathbf{P} must be eigenvectors of \mathbf{A} and d_i its corresponding eigenvalue. Since \mathbf{P} is invertible, its columns must be linearly independent.

Note that,

An $n \times n$ matrix \mathbf{A} is orthogonally diagonalizable if and only if \mathbf{A} is a symmetric matrix.

When is $\mathbf{A}_{n \times n}$ not diagonalizable?

While diagonalizing \mathbf{A} , if we come across fewer than n total vectors in all of the eigenspace bases, then the matrix is not diagonalizable. We can say, if the algebraic multiplicity of λ does not equal to the geometric multiplicity, then \mathbf{A} is not diagonalizable.

5.2.2 Eigenvalues and Eigenvectors of Real Symmetric Matrix:

Theorem: An $n \times n$ real symmetric matrix \mathbf{A} has real eigenvalues.

Proof: The conjugate transpose of a matrix \mathbf{A} is denoted by \mathbf{A}^H , then the Hermitian property can be written as

$$\begin{aligned}\mathbf{A} \text{ Hermitian} &\iff \mathbf{A} = \mathbf{A}^H \\ (\mathbf{x}^H \mathbf{A} \mathbf{x})^H &= \mathbf{x}^H \mathbf{A} \mathbf{x}\end{aligned}$$

$\therefore \mathbf{x}^H \mathbf{A} \mathbf{x}$ must be real.

$$\mathbf{A} \mathbf{x} = \lambda \mathbf{x}$$

Multiplying by \mathbf{x}^H from left on both sides,

$$\begin{aligned}\mathbf{x}^H \mathbf{A} \mathbf{x} &= \lambda \mathbf{x}^H \mathbf{x} = \lambda \|\mathbf{x}\|^2 \\ \lambda &= \frac{\mathbf{x}^H \mathbf{A} \mathbf{x}}{\mathbf{x}^H \mathbf{x}} = \text{real} \quad \mathbf{x} \neq 0\end{aligned}$$

Theorem: If \mathbf{A} is symmetric, then any two vectors from different eigenspaces are orthogonal.

Proof: Let \mathbf{v}_1 and \mathbf{v}_2 be eigenvectors that correspond to distinct eigenvalues, say λ_1 and λ_2 . To show that $\mathbf{v}_1 \cdot \mathbf{v}_2 = 0$, we have to compute

$$\begin{aligned}\lambda_1 \mathbf{v}_1 \cdot \mathbf{v}_2 &= (\lambda_1 \mathbf{v}_1)^T \mathbf{v}_2 \\ &= (\mathbf{A} \mathbf{v}_1)^T \mathbf{v}_2 \quad [\because \mathbf{v}_1 \text{ is an eigenvector}] \\ &= (\mathbf{v}_1^T \mathbf{A}^T) \mathbf{v}_2 \\ &= \mathbf{v}_1^T (\mathbf{A} \mathbf{v}_2) \quad [\because \mathbf{A}^T = \mathbf{A}] \\ &= \mathbf{v}_1^T (\lambda_2 \mathbf{v}_2) \quad [\because \mathbf{v}_2 \text{ is an eigenvector}] \\ &= \lambda_2 \mathbf{v}_1^T \mathbf{v}_2 \\ &= \lambda_2 \mathbf{v}_1 \cdot \mathbf{v}_2\end{aligned}$$

Hence $(\lambda_1 - \lambda_2) \mathbf{v}_1 \cdot \mathbf{v}_2 = 0$

But $(\lambda_1 - \lambda_2) \neq 0$, so $\mathbf{v}_1 \cdot \mathbf{v}_2 = 0$

5.2.3 Orthogonal Eigendecomposition and its Geometric Interpretation:

Suppose $\mathbf{A} = \mathbf{P} \mathbf{D} \mathbf{P}^{-1}$, where the columns of \mathbf{P} are orthonormal eigenvectors u_1, \dots, u_n of a real symmetric matrix \mathbf{A} and the corresponding eigenvalues $\lambda_1, \dots, \lambda_n$ are in the diagonal matrix \mathbf{D} . Then, since $\mathbf{P}^{-1} = \mathbf{P}^T$,

$$\begin{aligned}
\mathbf{A} &= \mathbf{P}\mathbf{D}\mathbf{P}^T \\
&= [\mathbf{u}_1 \quad \dots \quad \mathbf{u}_n] \begin{bmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_n \end{bmatrix} \begin{bmatrix} \mathbf{u}_1^T \\ \vdots \\ \mathbf{u}_n^T \end{bmatrix} \\
&= [\lambda_1 \mathbf{u}_1 \quad \dots \quad \lambda_n \mathbf{u}_n] \begin{bmatrix} \mathbf{u}_1^T \\ \vdots \\ \mathbf{u}_n^T \end{bmatrix}
\end{aligned}$$

We can write

$$A = \lambda_1 \mathbf{u}_1 \mathbf{u}_1^T + \lambda_2 \mathbf{u}_2 \mathbf{u}_2^T + \dots + \lambda_n \mathbf{u}_n \mathbf{u}_n^T$$

This representation of A is called a **spectral decomposition** of A because it breaks up A into pieces determined by the spectrum (eigenvalues) of A .

Each term is a $n \times n$ matrix of rank 1. For example, every column of $\lambda_1 \mathbf{u}_1 \mathbf{u}_1^T$ is a multiple of \mathbf{u}_1 . The matrix-vector product $A\mathbf{x}$ is decomposed as

$$A\mathbf{x} = \sum_{j=1}^n \lambda_j u_j (u_j^T \mathbf{x})$$

- Each matrix $\mathbf{u}_j \mathbf{u}_j^T$ is a projection matrix in the sense that for each \mathbf{x} in \mathbb{R}^n , the vector $\mathbf{u}_j \mathbf{u}_j^T \mathbf{x}$ is the orthogonal projection of \mathbf{x} onto the subspace spanned by \mathbf{u}_j , i.e., $\mathbf{u}_j \mathbf{u}_j^T \mathbf{x} = \mathbf{u}_j (\mathbf{u}_j \cdot \mathbf{x}) = (\mathbf{u}_j \cdot \mathbf{x}) \mathbf{u}_j = \text{proj}_{\mathbf{u}_j} \mathbf{x}$
- $(u_1^T x, \dots, u_n^T x)$ are coordinates of x in the orthonormal basis $\{u_1, \dots, u_n\}$
- $(\lambda_1 u_1^T x, \dots, \lambda_n u_n^T x)$ are coordinates of Ax in the orthonormal basis $\{u_1, \dots, u_n\}$

Change of Variable in Quadratic Form and its Geometric View:

A quadratic form on \mathbb{R}^n is a function Q defined on \mathbb{R}^n whose value at a vector \mathbf{x} in \mathbb{R}^n can be computed by an expression of the form $Q(\mathbf{x}) = \mathbf{x}^T A \mathbf{x}$, where A is an $n \times n$ symmetric matrix.

Let $\mathbf{x} = \mathbf{P}\mathbf{y}$ or equivalently, $\mathbf{y} = \mathbf{P}^{-1}\mathbf{x}$

where \mathbf{P} is an invertible matrix and \mathbf{y} is a new variable vector in \mathbb{R}^n . Here \mathbf{y} is the coordinate vector of \mathbf{x} relative to the orthonormal basis of \mathbb{R}^n determined by the columns of \mathbf{P} .

If the change of variable is made in a quadratic form $\mathbf{x}^T A \mathbf{x}$, then

$$\mathbf{x}^T A \mathbf{x} = (\mathbf{P}\mathbf{y})^T A (\mathbf{P}\mathbf{y}) = \mathbf{y}^T (\mathbf{P}^T A \mathbf{P}) \mathbf{y}$$

and the new matrix of the quadratic form is $(\mathbf{P}^T A \mathbf{P})$. Since A is symmetric, it is guaranteed that there is an orthonormal matrix \mathbf{P} such that $(\mathbf{P}^T A \mathbf{P})$ is a diagonal matrix \mathbf{D} and we get $\mathbf{y}^T \mathbf{D} \mathbf{y}$ with no cross product term.

Geometric View

When A is not diagonal, then the graph of the equation $\mathbf{x}^T A \mathbf{x} = c$ is rotated out of standard

position. Finding the principal axes (columns of P) amounts to finding a new coordinate system with respect to which the graph is in standard position.

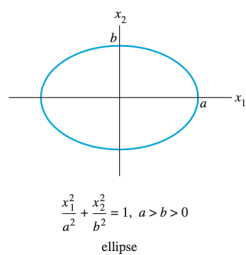


Figure 1: Ellipse in standard position

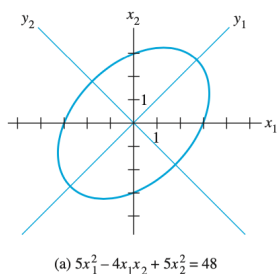


Figure 2: Ellipse *not* in standard position

5.2.4 Finally the Solution to the PCA Problem:

We approach our *maximization problem* by orthogonally diagonalizing \mathbf{S} as $\mathbf{S} = \mathbf{P}\mathbf{D}\mathbf{P}^T$ and making the change of variables $\mathbf{v} = \mathbf{P}\mathbf{y} \rightarrow \mathbf{y} = \mathbf{P}^T\mathbf{v}$. We know

$$Q(\mathbf{v}) = \mathbf{v}^T \mathbf{S} \mathbf{v} = \mathbf{y}^T \mathbf{D} \mathbf{y}$$

Additionally, $\|\mathbf{y}\|^2 = \|\mathbf{v}\|^2$ because

$$\|\mathbf{y}\|^2 = \mathbf{y}^T \mathbf{y} = (\mathbf{P}^T \mathbf{v})^T (\mathbf{P}^T \mathbf{v}) = \mathbf{v}^T \mathbf{v} = \|\mathbf{v}\|^2$$

Since \mathbf{D} is diagonal, we have

$$\mathbf{y}^T \mathbf{D} \mathbf{y} = \lambda_1 y_1^2 + \lambda_2 y_2^2 + \dots + \lambda_d y_d^2$$

Let's look for the maximum of these values subject to $\|\mathbf{y}\| = 1$. If we consider the maximum eigenvalue λ_{max} , then

$$\begin{aligned} \mathbf{y}^T \mathbf{D} \mathbf{y} &= \lambda_1 y_1^2 + \lambda_2 y_2^2 + \dots + \lambda_d y_d^2 \\ &\leq \lambda_{max} y_1^2 + \lambda_{max} y_2^2 + \dots + \lambda_{max} y_d^2 \\ &= \lambda_{max} (y_1^2 + y_2^2 + \dots + y_d^2) \\ &= \lambda_{max} \|\mathbf{y}\|^2 \\ &= \lambda_{max} \end{aligned}$$

the value λ_{max} is attained for $\mathbf{y}_{max} = (0 \ 0 \ \dots \ 0 \ 1 \ 0 \ \dots \ 0)$, where the 1 is at the location corresponding to λ_{max} . The corresponding \mathbf{v} is

$$\mathbf{v} = \mathbf{P}\mathbf{y} = \begin{pmatrix} u_1 & u_2 & \dots & u_{max-1} & u_{max} & u_{max+1} & \dots & u_d \end{pmatrix} \begin{pmatrix} 0 \\ 0 \\ \dots \\ 0 \\ 1 \\ 0 \\ \dots \\ 0 \end{pmatrix}$$

$$\therefore \mathbf{v} = \mathbf{u}_{max}$$

Therefore we can conclude that desired \mathbf{v} is \mathbf{u}_{max} which is the eigenvector associated to the largest eigenvalue λ_{max} .

Note that,

When we project on a subspace with orthonormal basis, the covariance of the projected data is zero.

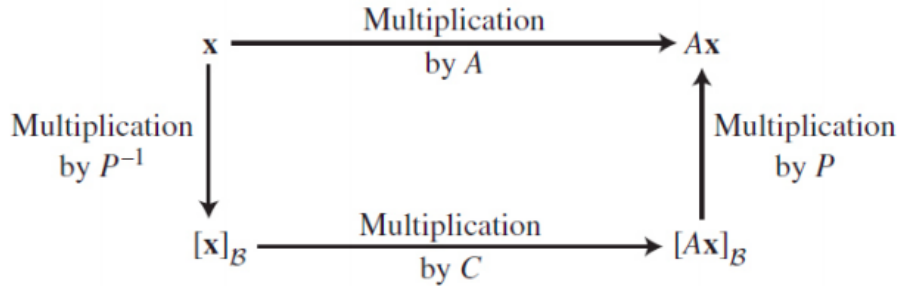
Additional Notes:

1. **Change of Basis and Similar Matrices:** To get more insight on the connection between similarity transformation and change of basis, let's see the following,

Suppose A is diagonalizable ($A = PDP^{-1}$). If B is the basis of \mathcal{R}^n formed by the columns of P , then

$$\begin{aligned}
 Ap_i \in \mathcal{R}^n & \quad \therefore Ap_i = Pc_i && \text{(for some } c_i) \\
 \implies [Ap_i]_B & = c_i \\
 \therefore A[p_1 \dots p_n] & = P[c_1 \dots c_n] \\
 \implies AP & = PC \\
 \implies P^{-1}AP & = C \\
 \implies A & = PCP^{-1} && \text{(x)}
 \end{aligned}$$

A and C are similar matrices if there exists another matrix P such that $A = PCP^{-1}$



$$\begin{aligned}
 A[x]_S & = PCP^{-1}[x]_S && \text{(under A)} \\
 & = PC[x]_B && \text{(under C)} \\
 & = P[Ax]_B && \text{(under C)} \\
 & = [Ax]_S && \text{() }
 \end{aligned}$$

C performs the same transformation, but in the coordinates defined by B . Now, if C is a simpler matrix (e.g. diagonal), then the transformation is easy. So, find a coordinate system (set of basis vectors), where the transformation T can be performed by a diagonal matrix.

Consider a linear transformation between two vectors spaces $T : V \rightarrow V$. Let $T(x) = t$. Let E be the basis of V and F be the basis of W

Let A be a matrix under T with respect to E , and B be a matrix under T with respect to F . Let S be the change of basis matrix from $F \rightarrow E$.

$$\begin{aligned}
[\mathbf{x}]_E &= \mathbf{S}_{F \rightarrow E} [\mathbf{x}]_F \\
[\mathbf{t}]_E &= \mathbf{A} [\mathbf{x}]_E \\
&\equiv [\mathbf{t}]_F = \mathbf{B} [\mathbf{x}]_F \\
\mathbf{S}^{-1} [\mathbf{t}]_E &= \mathbf{S}^{-1} \mathbf{A} [\mathbf{x}]_E \\
[\mathbf{t}]_F &= \mathbf{S}^{-1} \mathbf{A} \mathbf{S} [\mathbf{x}]_F \\
&\implies \mathbf{B} = \mathbf{S}^{-1} \mathbf{A} \mathbf{S} \\
&\implies \mathbf{A} = \mathbf{S} \mathbf{B} \mathbf{S}^{-1}
\end{aligned}$$

2. **Theorem:** A real symmetric matrix is Positive Definite (PD) if (i) $\mathbf{x}^T \mathbf{A} \mathbf{x} > 0 \quad \forall \mathbf{x} \neq 0$ and (ii) all the eigenvalues of \mathbf{A} satisfy $\lambda_i > 0$
 \mathbf{x} is eigenvector of \mathbf{A}

Proof: In the first step we will show that each eigenvalue will be positive.

If $\mathbf{x}^T \mathbf{A} \mathbf{x} > 0$:

$$\begin{aligned}
\mathbf{x}^T \lambda \mathbf{x} &= \lambda \|\mathbf{x}\|^2 \\
\mathbf{x}^T \mathbf{A} \mathbf{x} > 0 &\implies \lambda > 0
\end{aligned}$$

If $\lambda_i > 0$:

\mathbf{x} is any vector and $\mathbf{x} \neq 0$

$\mathbf{x} = c_1 \mathbf{x}_1 + c_2 \mathbf{x}_2 + \dots + c_n \mathbf{x}_n$

where $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ are eigenvectors.

$$\begin{aligned}
\therefore \mathbf{x}^T \mathbf{A} \mathbf{x} &= (c_1 \mathbf{x}_1^T + \dots + c_n \mathbf{x}_n^T) \mathbf{A} (c_1 \mathbf{x}_1 + \dots + c_n \mathbf{x}_n) \\
&= \sum_i c_i^2 \mathbf{x}_i^T \mathbf{A} \mathbf{x}_i \\
&= \sum_i c_i^2 \lambda_i \mathbf{x}_i^T \mathbf{x}_i > 0 \quad \because \lambda_i > 0, \quad \mathbf{x}_i^T \mathbf{x}_i > 0 \quad \forall \mathbf{x} \neq 0
\end{aligned}$$

3. The Transpose Trick PCA

Suppose the image database consists of $N = 30$ images. All images have the same dimension $m \times n = 100 \times 100$. The matrix \mathbf{P} has 30 columns and mn rows.

We are interested in the eigenvalues and eigenvectors of the Scatter matrix $\mathbf{P} \mathbf{P}^T$, but the matrix $\mathbf{P} \mathbf{P}^T$ has the dimensions 10000×10000 and computation of eigenvalues and eigenvectors become unfeasible.

The dimensions of $\mathbf{P}^T \mathbf{P}$ are only 30×30 and it is much more efficient to solve the eigenvalue problem for the matrix $\mathbf{P}^T \mathbf{P}$.

We will show the relation between the eigenvalues and eigenvectors of $\mathbf{P} \mathbf{P}^T$ and $\mathbf{P}^T \mathbf{P}$.

If \mathbf{x}_i is an eigenvector of $\mathbf{P}^T \mathbf{P}$ and λ_i is its corresponding eigenvalue, then we can write the following:

$$(\mathbf{P}^T \mathbf{P}) \mathbf{x}_i = \lambda_i \mathbf{x}_i$$

Multiplying by \mathbf{P} from the left,

$$\mathbf{P}(\mathbf{P}^T\mathbf{P})\mathbf{x}_i = \mathbf{P}\lambda\mathbf{x}_i$$

$$(\mathbf{P}\mathbf{P}^T)\mathbf{P}\mathbf{x}_i = \lambda(\mathbf{P}\mathbf{x}_i)$$

shows us that if \mathbf{x}_i is an eigenvector of $\mathbf{P}^T\mathbf{P}$, then $\mathbf{P}\mathbf{x}_i$ is the eigenvector of $\mathbf{P}\mathbf{P}^T$, with the same eigenvalue.

6 Linear Discriminant Analysis (LDA)

Fisher Linear Discriminant projects data to a line which preserves useful direction for *data classification*. Its main idea is to find projection to a line such that samples from different classes are well separated.

Suppose we have 2 classes and d -dimensional samples x_1, \dots, x_n where n_1 samples come from the first class (c_1) and n_2 samples come from the second class (c_2).

Let the line direction be given by unit vector \mathbf{v} . Thus the projection of sample x_i onto a line in direction \mathbf{v} is given by $\mathbf{v}^T x_i$. The scalar

$$y = \frac{\mathbf{v}^T \mathbf{x}}{\|\mathbf{v}\|}$$

is the projection of x along \mathbf{v} .

Measurement of separation between projections of different classes

Let $\tilde{\mu}_1$ and $\tilde{\mu}_2$ be the means of projections datapoints belonging to classes 1 and 2 respectively.

Let μ_1 and μ_2 be the means of classes 1 and 2.

$$\tilde{\mu}_1 = \frac{1}{n_1} \sum \mathbf{v}^T x_i = \mathbf{v}^T \left(\frac{1}{n_1} \sum_{x_i \in C_1} x_i \right) = \mathbf{v}^T \mu_1$$

Similarly, $\tilde{\mu}_2 = \mathbf{v}^T \mu_2$

The larger $|\tilde{\mu}_1 - \tilde{\mu}_2|$, the better is the expected separation. The problem with $|\tilde{\mu}_1 - \tilde{\mu}_2|$ is that it does not consider the variance of the classes.

We have to normalize $|\tilde{\mu}_1 - \tilde{\mu}_2|$ by scatter. y_i 's are the projected samples.

Scatter for projected samples of class 1 is: $\sigma_1^2 = \sum (y_i - \tilde{\mu}_1)^2$

Scatter for projected samples of class 2 is: $\sigma_2^2 = \sum (y_i - \tilde{\mu}_2)^2$

Thus Fisher Linear Discriminant is to project on line in the direction \mathbf{v} which maximizes:

$$J(\mathbf{v}) = \frac{(\tilde{\mu}_1 - \tilde{\mu}_2)^2}{\sigma_1^2 + \sigma_2^2}$$

To maximize J we want $(\tilde{\mu}_1 - \tilde{\mu}_2)^2$ to be large, i.e., the projected means to be far from each other.

And we also want σ_1^2 , the scatter in class 1, to be as small as possible, i.e. samples of class 1 cluster around the mean $\tilde{\mu}_1$. The same goes for class 2.

If we find \mathbf{v} which makes $J(\mathbf{v})$ large, we are guaranteed that the classes are well separated.

Assuming the classes to be equiprobable, it can be shown that:

$$(\tilde{\mu}_1 - \tilde{\mu}_2)^2 = \mathbf{v}^T (\mu_1 - \mu_2) (\mu_1 - \mu_2)^T \mathbf{v} \propto \mathbf{v}^T \mathbf{S}_b \mathbf{v}$$

where \mathbf{S}_b is the *Between-class scatter matrix* and $tr(\mathbf{S}_b)$ is the measure of the average (over all classes) distance of the mean of each individual class from the respective mean.

We have,

$$\sigma_i^2 = E[(y - \mu_i)^2] = E[\mathbf{v}^T (x - \mu_i) (x - \mu_i)^T \mathbf{v}] = \mathbf{v}^T \boldsymbol{\Sigma}_i \mathbf{v}$$

where for each $i = 1, 2$, samples $y(x)$ from the respective class v_i have been used.

$\boldsymbol{\Sigma}_i$ is the covariance matrix corresponding to the data of the class v_i in the d dimensional space. Using the definition of \mathbf{S}_w we get,

$$\sigma_1^2 + \sigma_2^2 \propto \mathbf{v}^T \mathbf{S}_w \mathbf{v}$$

where \mathbf{S}_w is the *Within-class scatter matrix*. $tr(\mathbf{S}_w)$ is a measure of the average, over all classes, variance of the features.

Combining we end up that the optimal direction is obtained by maximizing Fisher's criterion:

$$J(\mathbf{v}) = \frac{\mathbf{v}^T \mathbf{S}_b \mathbf{v}}{\mathbf{v}^T \mathbf{S}_w \mathbf{v}}$$

This is called the *Generalized Rayleigh quotient*.

Since \mathbf{S}_w is PD, then $\mathbf{S}_w = \mathbf{R}^T \mathbf{R}$

Let $\mathbf{y} = \mathbf{R} \mathbf{v}$

So, $\mathbf{v} = \mathbf{R}^{-1} \mathbf{y} = \mathbf{C} \mathbf{y}$

$$\therefore \frac{\mathbf{v}^T \mathbf{S}_b \mathbf{v}}{\mathbf{v}^T \mathbf{S}_w \mathbf{v}} = \frac{\mathbf{y}^T \mathbf{C}^T \mathbf{S}_b \mathbf{C} \mathbf{y}}{\mathbf{v}^T \mathbf{R}^T \mathbf{R} \mathbf{v}} = \frac{\mathbf{y}^T \mathbf{C}^T \mathbf{S}_b \mathbf{C} \mathbf{y}}{(\mathbf{R} \mathbf{v})^T \mathbf{R} \mathbf{v}} = \frac{\mathbf{y}^T \mathbf{C}^T \mathbf{S}_b \mathbf{C} \mathbf{y}}{\mathbf{y}^T \mathbf{y}}$$

Now it has turned into the aforementioned Rayleigh Quotient.

The Generalized Rayleigh quotient is maximized if \mathbf{v} is chosen such that $\mathbf{S}_b \mathbf{v} = \lambda \mathbf{S}_w \mathbf{v}$. This is the generalized symmetric eigenvalue problem. Since the extreme values λ of the Generalized Rayleigh quotient satisfy

$$\begin{aligned} \mathbf{S}_b \mathbf{v} &= \lambda \mathbf{S}_w \mathbf{v} \\ \Rightarrow \mathbf{S}_b \mathbf{v} &= \lambda \mathbf{R}^T \mathbf{R} \mathbf{v} \\ \Rightarrow \mathbf{S}_b \mathbf{C} \mathbf{y} &= \lambda \mathbf{R}^T \mathbf{y} \\ \Rightarrow \mathbf{R}^{T^{-1}} \mathbf{S}_b \mathbf{C} \mathbf{y} &= \lambda \mathbf{y} \\ \therefore \mathbf{C}^T \mathbf{S}_b \mathbf{C} \mathbf{y} &= \lambda \mathbf{y} \end{aligned}$$

The top eigenvector \mathbf{y}_1 of $\mathbf{C}^T \mathbf{S}_b \mathbf{C} \mathbf{y}$:

$$\max \frac{\mathbf{y}^T \mathbf{C}^T \mathbf{S}_b \mathbf{C} \mathbf{y}}{\mathbf{y}^T \mathbf{y}} = \lambda_1$$

when

$$\mathbf{C}^T \mathbf{S}_b \mathbf{C} \mathbf{y}_1 = \lambda_1 \mathbf{y}_1$$

Then

$$\mathbf{S}_b \mathbf{v}_1 = \lambda \mathbf{S}_w \mathbf{v}_1 \text{ for } \mathbf{v}_1 = \mathbf{R}^{-1} \mathbf{y}_1 = \mathbf{S}_w^{-\frac{1}{2}} \mathbf{y}_1$$

The eigenvalues of $\mathbf{S}_b \mathbf{v} = \lambda \mathbf{S}_w \mathbf{v}$ are real because $\mathbf{C}^T \mathbf{S}_b \mathbf{C}$ is a real symmetric matrix.

The orthogonality condition of two eigenvectors being orthogonal (in the case of two distinct eigenvalues) for a symmetric matrix extends to $\mathbf{S}_b \mathbf{v}_1 = \lambda \mathbf{S}_w \mathbf{v}_1$ with two symmetric matrices. For this we have to assume \mathbf{S}_w is positive definite and we have to change from $\mathbf{v}_1^T \mathbf{v}_2 = 0$ to "M-orthogonality" of \mathbf{v}_1 and \mathbf{v}_2 .

Two vectors are M-orthogonal if $\mathbf{v}_1^T \mathbf{S}_w \mathbf{v}_2 = \mathbf{v}_1^T \mathbf{R}^T \mathbf{R} \mathbf{v}_2 = \mathbf{y}_1^T \mathbf{y}_2 = 0$ when $\mathbf{S}_b \mathbf{v}_1 = \lambda_1 \mathbf{S}_w \mathbf{v}_1$, $\mathbf{S}_b \mathbf{v}_2 = \lambda_2 \mathbf{S}_w \mathbf{v}_2$ and $\lambda_1 \neq \lambda_2$

7 Singular Value Decomposition

A rectangular $m \times n$ matrix \mathbf{A} can not be decomposed in the previous manner. Singular Value Decomposition fills this gap. Now we need two sets of singular vectors the \mathbf{u} 's and the \mathbf{v} 's. The connection between \mathbf{u} 's and \mathbf{v} 's is not $\mathbf{Ax} = \lambda\mathbf{x}$, but rather $\mathbf{AV} = \mathbf{U}\Sigma$ or, $\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^T$. The matrices on the R.H.S are, respectively, (orthogonal), (diagonal) and (orthogonal).

$$\mathbf{A}_{m \times n} \underbrace{\begin{bmatrix} | & | & & | \\ \mathbf{v}_1 & \mathbf{v}_2 & \dots & \mathbf{v}_n \\ | & | & & | \end{bmatrix}}_{n \times n} = \underbrace{\begin{bmatrix} | & | & & | \\ \mathbf{u}_1 & \mathbf{u}_2 & \dots & \mathbf{u}_m \\ | & | & & | \end{bmatrix}}_{m \times m} \underbrace{\begin{bmatrix} \sigma_1 & 0 & 0 & \dots & 0 & 0 & 0 \\ 0 & \sigma_2 & 0 & \dots & 0 & 0 & 0 \\ 0 & 0 & \sigma_3 & \dots & 0 & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & \dots & \sigma_r & 0 & 0 \\ & & & & \vdots & & \\ 0 & 0 & 0 & \dots & 0 & 0 & 0 \end{bmatrix}}_{m \times n} \quad (1)$$

$$\Rightarrow \underbrace{\begin{bmatrix} | & | & & | \\ \mathbf{Av}_1 & \mathbf{Av}_2 & \dots & \mathbf{Av}_n \\ | & | & & | \end{bmatrix}}_{m \times n} = \underbrace{\begin{bmatrix} | & | & & | \\ \mathbf{u}_1 & \mathbf{u}_2 & \dots & \mathbf{u}_m \\ | & | & & | \end{bmatrix}}_{m \times m} \underbrace{\begin{bmatrix} \sigma_1 & 0 & 0 & \dots & 0 & 0 & 0 \\ 0 & \sigma_2 & 0 & \dots & 0 & 0 & 0 \\ 0 & 0 & \sigma_3 & \dots & 0 & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & \dots & \sigma_r & 0 & 0 \\ & & & & \vdots & & \\ 0 & 0 & 0 & \dots & 0 & 0 & 0 \end{bmatrix}}_{m \times n} \quad (2)$$

Here,

\mathbf{U} is $m \times m$ orthogonal matrix, u_1, \dots, u_m are called m left singular vectors in \mathbb{R}^m .

\mathbf{V} is $n \times n$ orthogonal matrix. v_1, \dots, v_n are called n right singular vectors in \mathbb{R}^n

$\Sigma_{m \times n}$ has positive entries $\sigma_1, \dots, \sigma_r$ which are in descending order $\sigma_1 \geq \sigma_2 \geq \dots > 0$. They are called the singular values of \mathbf{A} . They fill the first r places on the main diagonal of Σ - when \mathbf{A} has rank r . The rest of Σ is zero. We can rewrite (2) as:

$$\begin{aligned}
\mathbf{A}\mathbf{v}_1 &= \sigma_1 \mathbf{u}_1 \\
\mathbf{A}\mathbf{v}_2 &= \sigma_2 \mathbf{u}_2 \\
&\vdots \\
\mathbf{A}\mathbf{v}_r &= \sigma_r \mathbf{u}_r \\
&\vdots \\
\mathbf{A}\mathbf{v}_{r+k} &= 0
\end{aligned}$$

Observations:

1. You will see later on why the rank r of matrix \mathbf{A} and subsequently the ranks of $\mathbf{A}^T\mathbf{A}$ and $\mathbf{A}\mathbf{A}^T$ are important. Essentially, we only need r \mathbf{u} and \mathbf{v} vectors and that is enough for the decomposition of matrix \mathbf{A} , which will be achieved by using $\mathbf{A}^T\mathbf{A}$ and $\mathbf{A}\mathbf{A}^T$.
2. It is also interesting to note that $\mathbf{A}^T\mathbf{A}$ and $\mathbf{A}\mathbf{A}^T$ are square matrices. So, we will prove later on that these square matrices have exactly r non-zero eigenvalues (with repetitions) and $n-r$ zero eigenvalues. These r non-zero eigenvalues will yield the r eigenvectors of \mathbf{V} and the rest of the $n-r$ eigenvectors will come from the eigenspace of the $n-r$ zero eigenvalues.
3. The ranks are therefore important because we need to be certain that the matrix \mathbf{A} will be decomposed correctly into \mathbf{U} , $\mathbf{\Sigma}$ and \mathbf{V}

7.1 Finding the singular vectors

Our goal is to find two sets of singular vectors (which are orthonormal)- the \mathbf{u}' s and \mathbf{v}' s. **We begin with finding \mathbf{V} .**

One way to find \mathbf{V} is to form the symmetric matrix $\mathbf{A}^T\mathbf{A}$.

$$\mathbf{A}^T\mathbf{A} = (\mathbf{V}\mathbf{\Sigma}^T\mathbf{U}^T)(\mathbf{U}\mathbf{\Sigma}\mathbf{V}^T) = \mathbf{V}\mathbf{\Sigma}^T\mathbf{\Sigma}\mathbf{V}^T$$

The right side has the special form \mathbf{PDP}^T .

Eigenvalues are in $\mathbf{D} = \mathbf{\Sigma}^T\mathbf{\Sigma}$

So now we know how \mathbf{V} connects to the symmetric matrix $\mathbf{A}^T\mathbf{A}$. \mathbf{V} contains orthonormal eigenvectors of $\mathbf{A}^T\mathbf{A}$.

Before we prove Observation (2) above, we will first prove that if \mathbf{A} has rank $r = \dim(\text{colspace}(\mathbf{A})) = \dim(\text{rowspace}(\mathbf{A}))$, then $\mathbf{A}^T\mathbf{A}$ and $\mathbf{A}\mathbf{A}^T$ both will have rank $= r$.

Proof: The set of vectors which satisfy $\mathbf{A}^T\mathbf{A}x = 0$ is called the nullspace of the matrix $\mathbf{A}^T\mathbf{A}$. So, $x \in N(\mathbf{A}^T\mathbf{A})$

$$\mathbf{A}^T\mathbf{A}x = 0$$

$$\begin{aligned}
x^T A^T A x &= 0 \\
(Ax)^T (Ax) &= 0 \\
\therefore Ax &= 0 \quad [:\cdot x^T x = 0 \Leftrightarrow x = 0] \\
\therefore x &\in N(A)
\end{aligned}$$

$\therefore N(A^T A) \subseteq N(A)$ Again, $x \in N(A)$

$$Ax = 0$$

$$A^T Ax = 0$$

$$\therefore x \in N(A^T A)$$

$$\therefore N(A) \subseteq N(A^T A)$$

So it is obvious that $N(A^T A) = N(A) \Rightarrow \dim(N(A^T A)) = \dim(N(A))$

According to the Rank-Nullity Theorem

$$\dim(\text{colspace}(A)) + \dim(N(A)) = n = \dim(\text{colspace}(A^T A)) + \dim(N(A^T A))$$

$$\text{It gives, } r = \dim(A) = \dim(A^T A)$$

$$\text{We can also get, } \dim(N(A)) = \dim(N(A^T A)) = n - r$$

Again,

$$\begin{aligned}
AA^T y &= 0 \\
y^T AA^T y &= 0 \\
(A^T y)^T (A^T y) &= 0 \\
\therefore A^T y &= 0 \quad [:\cdot x^T x = 0 \Leftrightarrow x = 0]
\end{aligned}$$

The left nullspace is the space of all vectors y such that $A^T y = 0$ or equivalently $y^T A = 0$.

It is obvious that $N(AA^T) = N(A^T) \Rightarrow \dim(N(AA^T)) = \dim(N(A^T))$

According to the Rank-Nullity Theorem

$$\dim(\text{colspace}(A^T)) + \dim(N(A^T)) = m = \dim(\text{colspace}(AA^T)) + \dim(N(AA^T))$$

$$\text{Since } \dim(A) = \dim(A^T), \text{ it gives, } r = \dim(A) = \dim(AA^T)$$

$$\text{We can also get, } \dim(N(A^T)) = \dim(N(AA^T)) = m - r$$

Now we will prove Observation 2.

Lemma: The rank of any square matrix equals the number of nonzero eigen-values (with repetitions), so the square matrix $A^T A_{n \times n}$ of rank r will have r numbers of eigenvalues and $n - r$ number of zero eigenvalues. We already have $A^T A = V \Sigma^T \Sigma V^T$,

So $A^T A$ and $S (= \Sigma^T \Sigma)$ are similar matrices. We have to show that they have the same rank r .

Proof: Let $Y = A^T A$

We have $Y = V S V^T$

$\implies V^{-1} Y V = S \implies Y V = V S$

Let $S u = 0$, when $u \in N(S)$

$V S u = 0, \therefore u \in N(V S)$

$\therefore N(S) \subseteq N(V S)$

When V is invertible

$V S x = 0$, if $x \in N(V S)$

$S x = V^{-1} 0 = 0, \therefore x \in N(S)$

$\therefore N(V S) \subseteq N(S)$

$\therefore N(V S) = N(S)$

From Rank-Nullity Theorem

$\text{rank}(V S) = \text{rank}(S)$ when V is invertible

Similarly, $\text{rank}(Y V) = \text{rank}(Y)$ when V is invertible

$\therefore \text{rank}(Y) = \text{rank}(S)$

If $A^T A$ is diagonalizable and $A^T A = V S V^T$ then $A^T A$ and S are similar and of course S contains the eigenvalues of $A^T A$. And since the ranks are same and equal to r , then S must contain r nonzero rows and $n - r$ rows of zeros which implies $A^T A$ has r nonzero rows and $n - r$ rows of zeros.

So the eigenvalues of $A^T A$ are all nonnegative. We may assume that the eigenvalues are rearranged so that $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n \geq 0$

The **singular values** of A are the square roots of the eigenvalues of $A^T A$, denoted by $\sigma_1, \dots, \sigma_r$.

$$A^T A = (U \Sigma V^T)^T (U \Sigma V^T) = V \Sigma^T U^T U \Sigma V^T = V_r \text{diag}(\sigma_1^2, \dots, \sigma_r^2) V_r^T$$

$$A A^T = (U \Sigma V^T) (U \Sigma V^T)^T = U \Sigma V^T V \Sigma^T U^T = U_r \text{diag}(\sigma_1^2, \dots, \sigma_r^2) U_r^T$$

That is, $\sigma_i = \sqrt{\lambda_i}$ for $1 \leq i \leq r$

Observation:

1. The singular values of A are the lengths of the vectors $A v_1, \dots, A v_r$.

It can be seen in the following way: Let $\lambda_1, \dots, \lambda_r$ be the associated eigenvalues of $A^T A_{n \times n}$.

Then for $1 \leq i \leq r$,

$$\begin{aligned}\|Av_i\|^2 &= (Av_i)^T Av_i \\ &= v_i^T A^T Av_i \\ &= v_i^T (\lambda_i v_i) \quad [\because v_i \text{ is an eigen vector of } A^T A] \\ &= \lambda_i \quad [\because v_i \text{ is a unit vector}]\end{aligned}$$

2. The eigenvalues found from $A^T A$ and AA^T are incidentally the same. The eigenvectors are different.

So, we have found the vectors in \mathbf{V} . Earlier, we saw that SVD requires $Av_k = \sigma_k u_k$. It connects each right singular vector v_k to a left singular vector u_k , for $k = 1, \dots, r$. The \mathbf{u} vectors obtained from this will be valid if and only if:

1. They are eigenvectors of AA^T
2. They are orthogonal to each other

So the unit eigenvector is $u_k = \frac{Av_k}{\sigma_k}$ for $k = 1, \dots, r$

Check that these u 's are eigenvectors of AA^T :

$$AA^T u_k = AA^T \left(\frac{Av_k}{\sigma_k} \right) = A \left(\frac{A^T Av_k}{\sigma_k} \right) = A \frac{\sigma_k^2 v_k}{\sigma_k}$$

Using $u_k = \frac{Av_k}{\sigma_k} \implies \sigma_k = \frac{Av_k}{u_k}$ we get,

$$AA^T u_k = \sigma_k^2 u_k$$

The v 's were chosen to be orthonormal. Now we will check that u 's are also orthonormal:

$$u_j^T u_k = \left(\frac{Av_j}{\sigma_j} \right)^T \left(\frac{Av_k}{\sigma_k} \right) = \frac{v_j^T (A^T Av_k)}{\sigma_j \sigma_k} = \frac{\sigma_k}{\sigma_j} v_j^T v_k$$

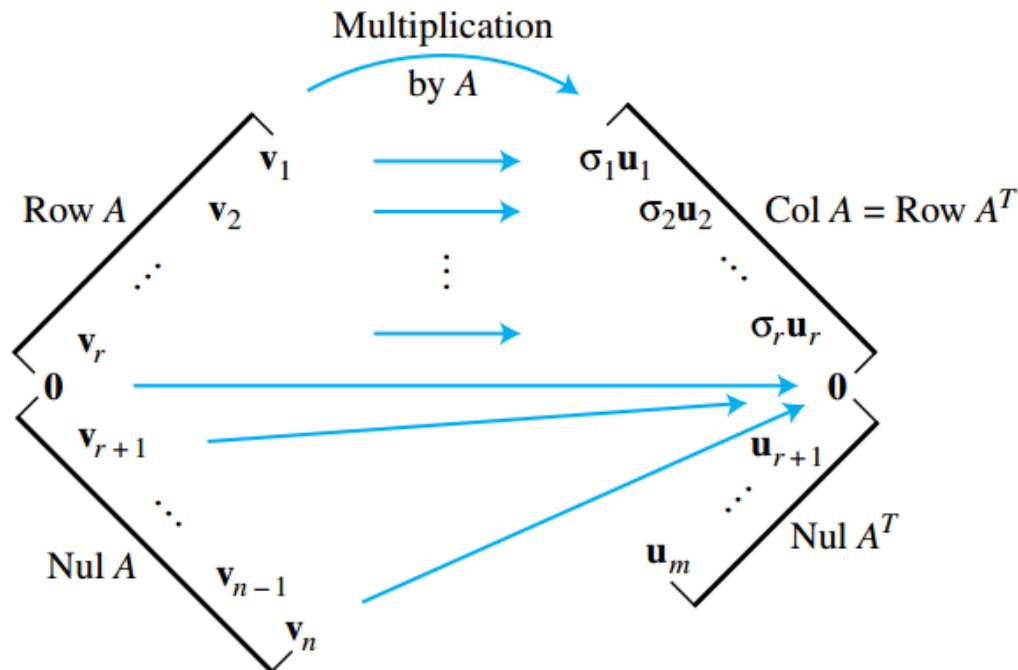
$$\therefore u_j^T u_k = \frac{\sigma_k}{\sigma_j} v_j^T v_k = \begin{cases} 1 & \text{if } j = k \\ 0 & \text{if } j \neq k \end{cases}$$

Finally we have to choose the last $n - r$ vectors v_{r+1} to v_n and the last $m - r$ vectors u_{r+1} to u_m , so that they form the orthogonal bases for \mathbb{R}^n and \mathbb{R}^m respectively.

Observations

1. **These v 's and u 's are in the nullspaces of A and A^T .** We can choose any orthonormal bases for those nullspaces.

2. It can therefore be noted that obtaining the eigenvectors from the eigenspace of the zero eigenvalues is equivalent to obtaining the vectors from the nullspace of A .
3. The last $n - r$ vectors of V and the last $m - r$ vectors of U will automatically be orthogonal to the first v 's in the row space of A and the first u 's in the column space of A^T respectively because the whole spaces are orthogonal: $N(A) \perp C(A^T)$ and $N(A^T) \perp C(A)$



We will prove that u_1, \dots, u_r provide an orthogonal basis for Col A .

Proof: We saw earlier $u_i^T u_j = 0$ for $i \neq j$, so $(Av_i)^T (Av_j) = 0$

Thus $\{Av_1, \dots, Av_n\}$ is an orthogonal set. Since there are r nonzero singular values, $Av_i \neq 0$ if and only if $1 \leq i \leq r$, so Av_1, \dots, Av_r are linearly independent vectors and they are in Col A . Finally

for any y in Col A say $y = Ax$ we can write $x = c_1 v_1 + \dots + c_n v_n$

and $y = Ax = c_1 Av_1 + \dots + c_r Av_r + c_{r+1} Av_{r+1} \dots + c_n Av_n$

$\therefore y = c_1 Av_1 + \dots + c_r Av_r + 0 + \dots + 0$

Thus y is in Span $\{Av_1, \dots, Av_r\}$ which shows that $\{u_1, \dots, u_r\}$ is an orthogonal basis for Col A and rank $= r$. **[Proved]**

As we know, $(\text{Col } A)^\perp = \text{Nul } A^T$. Hence u_{r+1}, \dots, u_m is an orthonormal basis for $\text{Nul } A^T$.

Since $\|Av_i\| = \sigma_i$ for $1 \leq i \leq n$, and σ_i is 0 if and only if $i > r$, the vectors v_{r+1}, \dots, v_n span a subspace of $\text{Nul } A$ of dimensions $n - r$.

By the Rank theorem, $\dim \text{Nul } A = n - \text{rank } A$. It follows that v_{r+1}, \dots, v_n is an orthonormal basis for $\text{Nul } A$.

The orthogonal complement of $\text{Nul } A^T$ is $\text{Col } A$. Interchanging A and A^T , note that $(\text{Nul } A)^\perp = \text{Col } A^T = \text{Row } A$. **Hence v_1, \dots, v_r is an orthonormal basis for Row A .**

(Additional) **Observation 1.** Eigenvectors of AA^T must go into the columns of U

$$AA^T = (U\Sigma V^T)(V\Sigma^T U^T) = U\Sigma\Sigma^T U^T$$

As we already saw earlier, U contains the orthogonal eigenvectors of AA^T . The m by m diagonal matrix is in the middle $\Sigma\Sigma^T$ with the eigenvalues $\sigma_1^2, \dots, \sigma_r^2$ on the diagonal.

Observation 2. For positive definite matrices, Σ is D and $U\Sigma V^T$ is identical to PDP^T . For other symmetric matrices, any negative eigenvalues in D become positive in Σ .

For complex matrices, Σ remains real but U and V become unitary (the complex version of orthogonal). We take complex conjugates in $U^H U = I$ and $V^H V = I$ and $A = U\Sigma V^H$

7.2 Best Matrix approximation

SVD separates the matrix into rank one pieces. A special property of the SVD is that **those pieces come in order of importance**. The first piece $\sigma_1 u_1 v_1^T$ is the closest rank-one matrix to A . A rank- k approximation is obtained by keeping the leading k singular values and vectors and discarding the rest.

A truncated SVD basis (and the resulting approximation matrix A_k) will be denoted by $A_k = \tilde{U}\tilde{\Sigma}\tilde{V}^T$

Since Σ is diagonal, the rank- k SVD approximation is given by the sum of k distinct rank-1 matrices:

$$A_k = \sum_{j=1}^k \sigma_j u_j v_j^T = \sigma_1 u_1 v_1^T + \sigma_2 u_2 v_2^T + \dots + \sigma_k u_k v_k^T$$

For a given k , there is no better approximation for A , in the ℓ_2 sense, than the truncated SVD approximation.

Eckart-Young Theorem: If B has rank k , then $\|A - B\|_F \geq \|A - A_k\|_F$

Proof: $A_k = U \text{diag}(\sigma_1, \sigma_2, \dots, \sigma_k, 0, \dots, 0) V^T$, $\text{rank}(A_k) = k$

$$A_k = \sum_{i=1}^k \sigma_i u_i v_i^T$$

$$\|A - A_k\|_2 = \left\| \sum_{i=k+1}^n \sigma_i u_i v_i^T \right\|_2 = \sigma_{k+1}$$

The whole proof of $\|A - B\|_2 \geq \sigma_{k+1}$ depend on a good choice of the vector x in computing the norm $\|A - B\|$:

We will choose $x \neq 0$ so that $Bx = 0$ and $x = \sum_1^{k+1} c_i v_i$

First, the nullspace B has dimension $\geq n - k$, because B has rank $\leq k$.

Second, the combinations of v_1 to v_{k+1} produce a subspace of dimension $k + 1$. Those two subspaces must intersect.

When dimensions add to $(n - k) + (k + 1) = n + 1$, the subspaces must share a line (at least).

Choose a nonzero vector x on this line. Use that x to estimate the norm $A - B$. As $Bx = 0$ and $Av_i = \sigma_i u_i$:

$$\|(A - B)x\|^2 = \|Ax\|^2 = \left\| \sum \frac{\sigma_i u_i}{v_i} c_i v_i \right\|^2 = \|c_i \sigma_i u_i\|^2 = \sum_1^{k+1} c_i^2 \sigma_i^2$$

That sum is at least as large as $(\sum c_i^2) \sigma_{k+1}^2$, which is exactly $\|x\|^2 \sigma_{k+1}^2$.

So $\|(A - B)x\| \geq \sigma_{k+1} \|x\|$

$$\|A - B\| \geq \sigma_{k+1} = \|A - A_k\|$$

We still have $V_k^T V_k = I_k$ and $U_k^T U_k = I_k$ from orthogonal unit vectors v 's and u 's.

But when V_k and U_k are not square, we can no longer have two-sided inverses: $V_k V_k^T = (n \times k)(k \times n) \neq I_n$ and $U_k U_k^T = (m \times k)(k \times m) \neq I_m$