# Preliminaries

Deep Learning (ENGG*6600*01)

School of Engineering,
University of Guelph, ON, Canada

Course Instructor: Benyamin Ghojogh
Summer 2023

**Dataset, Learning Model, and Learning Tasks**

# Dataset

- Consider the measurement of a quantity. This quantity can be:
  - personal health data, including blood pressure, blood sugar, and blood fat,
  - images from a specific scene but taken from different perspectives,
  - images from several categories of animals, such as cat, dog, frog, etc.,
  - medical images, such as digital pathology image patches, including both healthy and tumorous tissues,
  - or any other measured signal.
- The quantity can be multidimensional, i.e., a set of values, and therefore, every quantity can be considered a multidimensional data point in a Euclidean space.
- Let the dimensionality of this space be $d$, meaning that every quantity is a $d$-dimensional vector, or data point, in $\mathbb{R}^d$. The set of $d$ values for the quantity can be called **features** of the quantity.
- Multiple measurements of a quantity can exist, each of which is a $d$-dimensional data point. Therefore, there will be a set of $d$-dimensional data points, called a **dataset**.
- For example, the quantity can be an image, whose features are its pixels. The dataset can be a set of images from a specific scene but with different perspectives and angles.
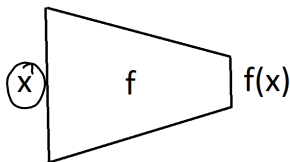
# Learning Model

- Consider a dataset of $n$ data points $\{x_i \in \mathbb{R}^d\}_{i=1}^{n}$ each of which is a $d$-dimensional vector in the $d$-dimensional Euclidean space. We can put these vectors column-wise in a matrix $X \in \mathbb{R}^{d \times n}$.

- Consider a learning model $f$ which is a map from data space to some output space:

$$
\begin{aligned}
f &: \mathbb{R}^d \to \mathbb{R}^p, \\
f &: x \mapsto f(x).
\end{aligned}
\tag{1}
$$

- Usually, $p \leq d$ but not necessarily.

# Learning Tasks

- Learning model is like a new-born baby (it first knows nothing and we should teach it)
- **Supervised**:
  - **Regression**: example of learning EMG signals for artificial leg, example of weather prediction

  $$f(\mathbf{x}) \in [0,1]^p \text{ or } f(\mathbf{x}) \in \mathbb{R}.$$

  - **Classification**: example of teaching apples and cucumbers to a baby

  $$f(\mathbf{x}) \in \{\ell_1, \ell_2, \ldots, \ell_m\}.$$

- **Unsupervised**:
  - **Clustering**: example of clustering apples and cucumbers by a baby

  $$f(\mathbf{x}) \in \{\ell_1, \ell_2, \ldots, \ell_m\}.$$

- **Environment (world)**:
  - **Reinforcement learning**: example of teaching a dog

  $$f(\mathbf{x}) = a \in \mathcal{A},$$

  where $\mathcal{A}$ is the set of possible actions.

# Learning Tasks

- **Dimensionality reduction (manifold learning)**: learning an embedding space

$$f : \mathbb{R}^d \to \mathbb{R}^p.$$

  where $p \leq d$, and usually $p \ll d$.

  - **Unsupervised** dimensionality reduction: embedding similar patterns close to each other
  - **Supervised** dimensionality reduction: Decreasing the intra-class variances and increase the inter-class variances
  - For more information on dimensionality reduction , you can see our textbook [1]: https://link.springer.com/book/10.1007/978-3-031-10602-6

- **Numerosity processing**:
  - **Outlier (anomaly) detection**: detecting outliers in data
  - **Prototype selection** [2]: selecting important instances
  - **Prototype generation** [3]: selecting and generating important instances
  - For more information on dimensionality reduction and numerosity processing, you can see my PhD thesis: [4]: https://uwspace.uwaterloo.ca/handle/10012/16813

# Other Fields of AI

Some other fields of Artificial Intelligence (AI):

- Soft computing:
    - Fuzzy logic and fuzzy control
    - Metaheuristic optimization and intelligent search
- Biological-inspired (third generation) neural networks - relation to neuroscience and cognitive science - Example: spiking neural network
- Feature engineering (pre-processing):
    - Feature selection
    - Feature extraction (dimensionality reduction)
- Application of AI in various fields of science and technology

**Linear Projection**

# Column Space

- Consider $p$ basis vectors. We can define a $p$-dimensional Euclidean space by these $p$ basis vectors. For example, two vectors define a plane and three vectors define a 3D space.
- In terminology: The $p$ basis vectors **span** the $p$-dimensional Euclidean space. Or the $p$-dimensional Euclidean space is **spanned by** the $p$ basis vectors.
- Consider the $p$ vectors $\{u_1, \ldots, u_p\}$. These vectors can be stacked columnwise in matrix $U = [u_1, \ldots, u_p] \in \mathbb{R}^{d \times p}$.
- The space spanned by the columns of matrix $U$ is called the column space of matrix $U$, denoted by $\mathbb{C}ol(U)$:

$$\mathbb{C}ol(U) := \mathbf{span}\{u_1, \ldots, u_p\}. \tag{2}$$

- In other words, the space whose bases are the columns of matrix $U$ is called the column space of matrix $U$.
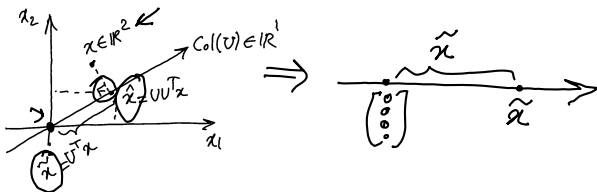
# Linear Projection

- Assume there is a data point $x \in \mathbb{R}^d$. The aim is to project this data point onto the vector space spanned by $p$ vectors $\{u_1, \ldots, u_p\}$, where each vector is $d$-dimensional and usually $p \ll d$.
- These vectors can be stacked columnwise in matrix $U = [u_1, \ldots, u_p] \in \mathbb{R}^{d \times p}$. In other words, the goal is to project $x$ onto the column space of $U$, denoted by $\mathbb{Col}(U)$.
- As $p < d$, this projection is projection onto a **subspace** because we are projecting from $d$-dimensional space onto a lower dimensional space.
- The **projection** of $x \in \mathbb{R}^d$ onto $\mathbb{Col}(U) \in \mathbb{R}^p$ is:

$$\mathbb{R}^p \ni \widetilde{x} := U^\top x. \tag{3}$$
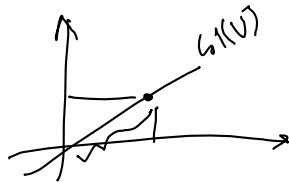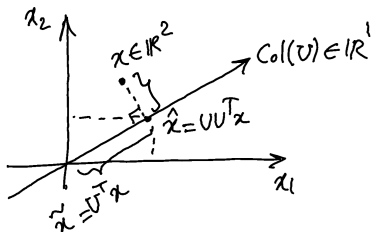
- The **reconstruction** of $\widetilde{x} \in \mathbb{R}^p$ in the $d$-dimensional space is:

$$\mathbb{R}^d \ni \widehat{x} := U\widetilde{x} = UU^\top x. \tag{4}$$

- Reconstruction is its representation in $\mathbb{R}^d$ again, but after projection.
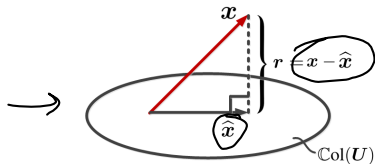
# Linear Projection



- **Reconstruction error**: There is a <u>residual/error</u> between the original data $x$ and its reconstruction (if the data point <u>is already in the column space</u>, this <u>residual is zero</u>):

$$r = x - \widehat{x} = x - UU^\top x. \tag{5}$$

**Norm**

# Inner product

## Definition (Inner product of vectors)

Consider two vectors $\boldsymbol{x} = [x_1, \ldots, x_d]^\top \in \mathbb{R}^d$ and $\boldsymbol{y} = [y_1, \ldots, y_d]^\top \in \mathbb{R}^d$. Their **inner product**, also called **dot product**, is:

$$\langle \boldsymbol{x}, \boldsymbol{y} \rangle = \boldsymbol{x}^\top \boldsymbol{y} = \sum_{i=1}^d x_i \, y_i.$$

$$x_1 y_1 + x_2 y_2 + \cdots + x_d y_d$$

## Definition (Inner product of matrices)

We also have inner product between matrices $\boldsymbol{X}, \boldsymbol{Y} \in \mathbb{R}^{d_1 \times d_2}$. Let $\boldsymbol{X}_{ij}$ denote the $(i, j)$-th element of matrix $\boldsymbol{X}$. The inner product of $\boldsymbol{X}$ and $\boldsymbol{Y}$ is:

$$\langle \boldsymbol{X}, \boldsymbol{Y} \rangle = \mathbf{tr}(\boldsymbol{X}^\top \boldsymbol{Y}) = \sum_{i=1}^{d_1} \sum_{j=1}^{d_2} \boldsymbol{X}_{i,j} \, \boldsymbol{Y}_{i,j},$$

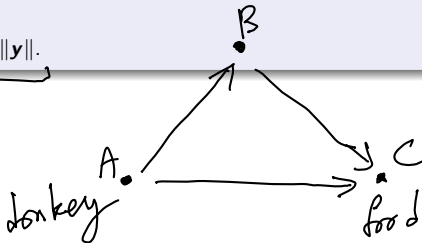where **tr**(.) denotes the trace of matrix.

# Norm

$\|x\|$



$\|x-y\|$

## Definition (Norm)

A function $\|\cdot\| : \mathbb{R}^d \to \mathbb{R}$, $\|\cdot\| : x \mapsto \|x\|$ is a **norm** if it satisfies:

1. $\|x\| \geq 0, \forall x$
2. $\|ax\| = |a|\,\|x\|, \forall x$ and all scalars $a$
3. $\|x\| = 0$ if and only if $x = 0$
4. Triangle inequality: $\|x + y\| \leq \|x\| + \|y\|$.



donkey A

B

C
food

# Important norms for vectors

Some important norms for a vector $\boldsymbol{x} = [x_1, \ldots, x_d]^\top$ are as follows.

- The $\ell_p$ **norm** is:

$$\|\boldsymbol{x}\|_p := \left(|x_1|^p + \cdots + |x_d|^p\right)^{1/p}, \quad = \sqrt[p]{|x_1|^p + \cdots + |x_d|^p}$$

where $p \geq 1$ and $|.|$ denotes the absolute value.

- Two well-known $\ell_p$ norms are $\ell_1$ **norm** and $\ell_2$ **norm** (also called the **Euclidean norm**) with $p = 1$ and $p = 2$, respectively:

$$\|\boldsymbol{x}\|_1 := |x_1| + \cdots + |x_d| = \sum_{i=1}^{d} |x_i|,$$

$$\|\boldsymbol{x}\|_2 := \sqrt{x_1^2 + \cdots + x_d^2} = \sqrt{\sum_{i=1}^{d} x_i^2},$$

# Important norms for matrices

Some important norms for a matrix $\boldsymbol{X} \in \mathbb{R}^{d_1 \times d_2}$ are as follows.

- The formulation of the **Frobenius norm** for a matrix is similar to the formulation of $\ell_2$ norm for a vector:

$$\|\boldsymbol{X}\|_F := \sqrt{\sum_{i=1}^{d_1} \sum_{j=1}^{d_2} \boldsymbol{X}_{i,j}^2},$$

where $\boldsymbol{X}_{ij}$ denotes the $(i,j)$-th element of $\boldsymbol{X}$.

## Quadratic forms using norms

$$\| (x - y)\|_2^2 = (x - y)^\top (x - y)$$

For $\boldsymbol{x} \in \mathbb{R}^d$ and $\boldsymbol{X} \in \mathbb{R}^{d_1 \times d_2}$, we have:

$$\star \quad \|\boldsymbol{x}\|_2^2 = \underbrace{\boldsymbol{x}^\top \boldsymbol{x}} = \langle \boldsymbol{x}, \boldsymbol{x} \rangle = \underbrace{\sum_{i=1}^d x_i^2},$$

$$\star \quad \|\boldsymbol{X}\|_F^2 = \underbrace{\mathbf{tr}(\boldsymbol{X}^\top \boldsymbol{X})} = \langle \boldsymbol{X}, \boldsymbol{X} \rangle = \underbrace{\sum_{i=1}^{d_1} \sum_{j=1}^{d_2} \boldsymbol{X}_{i,j}^2},$$

which are convex and in quadratic forms.

**Preliminaries on Derivatives**

# Dimensionality of derivative

$$\frac{\delta\, tr(x^T x)}{\delta x} = 2X$$

- Consider a function $f : \mathbb{R}^{d_1} \to \mathbb{R}^{d_2}$, $f : \boldsymbol{x} \mapsto f(\boldsymbol{x})$.
- Derivative of function $f(\boldsymbol{x}) \in \mathbb{R}^{d_2}$ with respect to (w.r.t.) $\boldsymbol{x} \in \mathbb{R}^{d_1}$ has dimensionality $(d_1 \times d_2)$.
- This is because tweaking every element of $\boldsymbol{x} \in \mathbb{R}^{d_1}$ can change every element of $f(\boldsymbol{x}) \in \mathbb{R}^{d_2}$. The $(i, j)$-th element of the $(d_1 \times d_2)$-dimensional derivative states the amount of change in the $j$-th element of $f(\boldsymbol{x})$ resulted by changing the $i$-th element of $\boldsymbol{x}$.

## Examples

- The derivative of a scalar w.r.t. a scalar is a scalar.
- The derivative of a scalar w.r.t. a vector is a vector.
- The derivative of a scalar w.r.t. a matrix is a matrix.
- The derivative of a vector w.r.t. a vector is a matrix.
- The derivative of a vector w.r.t. a matrix is a rank-3 tensor.
- The derivative of a matrix w.r.t. a matrix is a rank-4 tensor.

# Dimensionality of derivative

In more details:

- If the function is $f : \mathbb{R} \to \mathbb{R}, f : x \mapsto f(x)$, the derivative $(\partial f(x)/\partial x) \in \mathbb{R}$ is a scalar because changing the scalar $x$ can change the scalar $f(x)$.
- If the function is $f : \mathbb{R}^d \to \mathbb{R}, f : \mathbf{x} \mapsto f(\mathbf{x})$, the derivative $(\partial f(\mathbf{x})/\partial \mathbf{x}) \in \mathbb{R}^d$ is a vector because changing every element of the vector $\mathbf{x}$ can change the scalar $f(\mathbf{x})$.
- If the function is $f : \mathbb{R}^{d_1 \times d_2} \to \mathbb{R}, f : \mathbf{X} \mapsto f(\mathbf{X})$, the derivative $(\partial f(\mathbf{X})/\partial \mathbf{X}) \in \mathbb{R}^{d_1 \times d_2}$ is a matrix because changing every element of the matrix $\mathbf{X}$ can change the scalar $f(\mathbf{X})$.
- If the function is $f : \mathbb{R}^{d_1} \to \mathbb{R}^{d_2}, f : \mathbf{x} \mapsto f(\mathbf{x})$, the derivative $(\partial f(\mathbf{x})/\partial \mathbf{x}) \in \mathbb{R}^{d_1 \times d_2}$ is a matrix because changing every element of the vector $\mathbf{x}$ can change every element of the vector $f(\mathbf{x})$.
- If the function is $f : \mathbb{R}^{d_1 \times d_2} \to \mathbb{R}^{d_3}, f : \mathbf{X} \mapsto f(\mathbf{X})$, the derivative $(\partial f(\mathbf{X})/\partial \mathbf{X})$ is a $(d_1 \times d_2 \times d_3)$-dimensional tensor because changing every element of the matrix $\mathbf{X}$ can change every element of the vector $f(\mathbf{X})$.
- If the function is $f : \mathbb{R}^{d_1 \times d_2} \to \mathbb{R}^{d_3 \times d_4}, f : \mathbf{X} \mapsto f(\mathbf{X})$, the derivative $(\partial f(\mathbf{X})/\partial \mathbf{X})$ is a $(d_1 \times d_2 \times d_3 \times d_4)$-dimensional tensor because changing every element of the matrix $\mathbf{X}$ can change every element of the matrix $f(\mathbf{X})$.

# Gradient and Hessian

## Definition (Gradient)

Consider a function $f : \mathbb{R}^d \to \mathbb{R}$, $f : \boldsymbol{x} \mapsto f(\boldsymbol{x})$. In optimizing the function $f$, the derivative of function w.r.t. its variable $\boldsymbol{x}$ is called the **gradient**, denoted by:

$$\nabla f(\boldsymbol{x}) := \frac{\partial f(\boldsymbol{x})}{\partial \boldsymbol{x}} \in \mathbb{R}^d.$$

## Definition (Hessian)

Consider a function $f : \mathbb{R}^d \to \mathbb{R}$, $f : \boldsymbol{x} \mapsto f(\boldsymbol{x})$. The second derivative of function w.r.t. to its derivative is called the **Hessian** matrix, denoted by:

$$\boldsymbol{B} = \nabla^2 f(\boldsymbol{x}) := \frac{\partial^2 f(\boldsymbol{x})}{\partial \boldsymbol{x}^2} \in \mathbb{R}^{d \times d}.$$

The Hessian matrix is symmetric. If the function is convex, its Hessian matrix is positive semi-definite.

# Chain rule

- When having composite functions (i.e., function of function), we use **chain rule** for derivative. Example:

$$f(x) = \sqrt{x^3 + x^2 - x + 10} = \sqrt{g(x)}, \quad g(x) = x^3 + x^2 - x + 10,$$

$$\frac{\partial f(x)}{\partial x} = \frac{\partial f(x)}{\partial g(x)} \times \frac{\partial g(x)}{\partial x} = \frac{1}{2\sqrt{g(x)}} \times (3x^2 + 2x - 1) = \frac{3x^2 + 2x - 1}{2\sqrt{x^3 + x^2 - x + 10}}$$

- The chain rule in matrix derivatives is usually stated right to left in matrix multiplications while transpose is used for matrices in multiplication.
- Let **vec**(.) denote vectorization of a $\mathbb{R}^{a \times b}$ matrix to a $\mathbb{R}^{ab}$ vector.
- Let **vec**$_{a \times b}^{-1}$(.) be de-vectorization of a $\mathbb{R}^{ab}$ vector to a $\mathbb{R}^{a \times b}$ matrix.

**Optimization**

# Optimization

- Lagrangian:

$$\underset{\boldsymbol{x}}{\text{minimize}} \quad f(\boldsymbol{x})$$
$$\text{subject to} \quad y_i(\boldsymbol{x}) \leq 0, \ i \in \{1, \ldots, m_1\},$$
$$h_i(\boldsymbol{x}) = 0, \ i \in \{1, \ldots, m_2\}.$$

$$\longrightarrow \mathcal{L}(\boldsymbol{x}, \boldsymbol{\lambda}, \boldsymbol{\nu}) := f(\boldsymbol{x}) + \sum_{i=1}^{m_1} \lambda_i y_i(\boldsymbol{x}) + \sum_{i=1}^{m_2} \nu_i h_i(\boldsymbol{x}) = f(\boldsymbol{x}) + \boldsymbol{\lambda}^\top \boldsymbol{y}(\boldsymbol{x}) + \boldsymbol{\nu}^\top \boldsymbol{h}(\boldsymbol{x}).$$

$$\frac{\partial \mathcal{L}}{\partial \boldsymbol{x}} \overset{\text{set}}{=} \boldsymbol{0} \implies \ldots$$

- Unconstrained optimization:

$$\underset{\boldsymbol{x}}{\text{minimize}} \quad f(\boldsymbol{x}),$$
$$\boldsymbol{x}^{(k+1)} := \boldsymbol{x}^{(k)} + (\Delta \boldsymbol{x})^{(k)},$$
$$\underline{\text{Gradient method:}} \ \boldsymbol{x}^{(k+1)} := \boldsymbol{x}^{(k)} - \eta^{(k)} \nabla f(\boldsymbol{x}^{(k)}),$$
$$\underline{\text{Newton's method:}} \ \boldsymbol{x}^{(k+1)} := \boldsymbol{x}^{(k)} - \eta^{(k)} \big( \nabla^2 f(\boldsymbol{x}^{(k)}) \big)^{-1} \nabla f(\boldsymbol{x}^{(k)}).$$

- Constrained optimization: We can use interior-point method or proximal methods, ...

**Eigenvalue and Singular Value Decomposition**

# Eigenvalue Problem

- Eigenvalue and generalized eigenvalue problems play important roles in different fields of science, including machine learning, physics, statistics, and mathematics.

- In the eigenvalue problem, the eigenvectors of a matrix represent the most important and informative directions of that matrix. For example, if the matrix is a covariance matrix of data, the eigenvectors represent the directions of the spread or variance of data and the corresponding eigenvalues are the magnitude of the spread in these directions [5].

- These directions are impacted by another matrix in the generalized eigenvalue problem. If the other matrix is the identity matrix, this impact is cancelled and the eigenvalue problem captures the directions of the maximum spread.

- The **eigenvalue problem** [6, 7] of a symmetric matrix $\boldsymbol{A} \in \mathbb{R}^{d \times d}$ is defined as:

$$\boldsymbol{A}\boldsymbol{\phi}_i = \lambda_i \boldsymbol{\phi}_i, \quad \forall i \in \{1, \ldots, d\}, \tag{6}$$

and in matrix form, it is:

$$\boldsymbol{A}\boldsymbol{\Phi} = \boldsymbol{\Phi}\boldsymbol{\Lambda}, \tag{7}$$

where the columns of $\mathbb{R}^{d \times d} \ni \boldsymbol{\Phi} := [\boldsymbol{\phi}_1, \ldots, \boldsymbol{\phi}_d]$ are the eigenvectors and diagonal elements of $\mathbb{R}^{d \times d} \ni \boldsymbol{\Lambda} := \mathbf{diag}([\lambda_1, \ldots, \lambda_d]^\top)$ are the eigenvalues. Note that $\boldsymbol{\phi}_i \in \mathbb{R}^d$ and $\lambda_i \in \mathbb{R}$.

# Eigenvalue Problem

- For the eigenvalue problem, the matrix $A$ can be nonsymmetric. If the matrix is symmetric, its eigenvectors are orthogonal/orthonormal and if it is nonsymmetric, its eigenvectors are not orthogonal/orthonormal.

- Equation (7) can be restated as:

$$A\Phi = \Phi\Lambda \implies A\Phi\Phi^\top = \Phi\Lambda\Phi^\top \implies A = \Phi\Lambda\Phi^\top = \Phi\Lambda\Phi^{-1}, \qquad (8)$$

where $\Phi^\top = \Phi^{-1}$ because $\Phi$ is an orthogonal matrix.

- There is always $\Phi^\top\Phi = I$ for orthogonal $\Phi$, but there is only $\Phi\Phi^\top = I$ if "all" columns of the orthogonal $\Phi$ exist (it is not truncated, i.e., it is a square matrix). Equation (8) is referred to as "**eigenvalue decomposition**", "eigen-decomposition", or "spectral decomposition".

$$\Phi = \begin{bmatrix} | & | & | & | \\ & & & \\ | & | & | & | \end{bmatrix} \rightarrow \Phi^\top\Phi = \Phi\Phi^\top \leq I$$

$$\begin{bmatrix} | & | & | \end{bmatrix} \rightarrow \Phi^\top\Phi \leq I$$

$$\Phi\Phi^\top \neq I$$

# Generalized Eigenvalue Problem

- The **generalized eigenvalue problem** [8, 7] of two symmetric matrices $\boldsymbol{A} \in \mathbb{R}^{d \times d}$ and $\boldsymbol{B} \in \mathbb{R}^{d \times d}$ is defined as:

$$\boxed{\boldsymbol{A}\phi_i = \lambda_i \boldsymbol{B}\phi_i,} \quad \forall i \in \{1, \ldots, d\}, \tag{9}$$

and in matrix form, it is:

$$\boxed{\boldsymbol{A}\boldsymbol{\Phi} = \boldsymbol{B}\boldsymbol{\Phi}\boldsymbol{\Lambda},} \tag{10}$$

where the columns of $\mathbb{R}^{d \times d} \ni \boldsymbol{\Phi} := [\phi_1, \ldots, \phi_d]$ are the eigenvectors and diagonal elements of $\mathbb{R}^{d \times d} \ni \boldsymbol{\Lambda} := \mathbf{diag}([\lambda_1, \ldots, \lambda_d]^\top)$ are the eigenvalues. Note that $\phi_i \in \mathbb{R}^d$ and $\lambda_i \in \mathbb{R}$.

- The generalized eigenvalue problem of Eq. (9) or (10) is denoted by $(\boldsymbol{A}, \boldsymbol{B})$.
- The $(\boldsymbol{A}, \boldsymbol{B})$ is called a "pair" or "pencil" [8], and the order in the pair matters, according to Eq. (10).
- The $\boldsymbol{\Phi}$ and $\boldsymbol{\Lambda}$ are called the generalized eigenvectors and eigenvalues of $(\boldsymbol{A}, \boldsymbol{B})$.
- The $(\boldsymbol{\Phi}, \boldsymbol{\Lambda})$ or $(\phi_i, \lambda_i)$ is called the "eigenpair" of the pair $(\boldsymbol{A}, \boldsymbol{B})$ in the literature [8].
- Comparing Eqs. (6) and (9) or Eqs. (7) and (10) demonstrates that the eigenvalue problem is a special case of the generalized eigenvalue problem where $\boldsymbol{B} = \boldsymbol{I}$.

# Singular Value Decomposition

- **Singular Value Decomposition (SVD)** [9] is one of the most well-known and effective matrix decomposition methods. There are different methods for obtaining this decomposition, one of which is Jordan's algorithm [9].
- SVD has two different forms, i.e., complete and incomplete.
- Consider a matrix $A \in \mathbb{R}^{\alpha \times \beta}$. The **complete SVD** decomposes the matrix as:

$$\mathbb{R}^{\alpha \times \beta} \ni A = U \Sigma V^\top, \tag{11}$$
$$U \in \mathbb{R}^{\alpha \times \alpha}, \quad V \in \mathbb{R}^{\beta \times \beta}, \quad \Sigma \in \mathbb{R}^{\alpha \times \beta},$$

where the columns of $U$ and the columns of $V$ are called *left singular vectors* and *right singular vectors*, respectively.

- In complete SVD, $\Sigma$ is a *rectangular diagonal matrix* whose main diagonal includes the *singular values*. In the cases with $\alpha > \beta$ and $\alpha < \beta$, this matrix is in the following forms:

$$\Sigma = \begin{bmatrix} \sigma_1 & 0 & 0 \\ \vdots & \ddots & \vdots \\ 0 & 0 & \sigma_\beta \\ 0 & 0 & 0 \\ \vdots & \vdots & \vdots \\ 0 & 0 & 0 \end{bmatrix} \text{ and } \begin{bmatrix} \sigma_1 & 0 & 0 & 0 & \cdots & 0 \\ \vdots & \ddots & & 0 & \cdots & 0 \\ 0 & 0 & \sigma_\alpha & 0 & \cdots & 0 \end{bmatrix},$$

respectively. In other words, the number of singular values is $\min(\alpha, \beta)$.

# Singular Value Decomposition

- The **incomplete SVD** decomposes the matrix as:

$$\mathbb{R}^{\alpha \times \beta} \ni \boldsymbol{A} = \boldsymbol{U} \boldsymbol{\Sigma} \boldsymbol{V}^{\top}, \tag{12}$$

$$\underbrace{\boldsymbol{U} \in \mathbb{R}^{\alpha \times k}}, \quad \underbrace{\boldsymbol{V} \in \mathbb{R}^{\beta \times k}}, \quad \underbrace{\boldsymbol{\Sigma} \in \mathbb{R}^{k \times k}},$$

where [10]:

$$k := \min(\alpha, \beta), \tag{13}$$

and the columns of $\boldsymbol{U}$ and the columns of $\boldsymbol{V}$ are called *left singular vectors* and *right singular vectors*, respectively.

- In incomplete SVD, $\boldsymbol{\Sigma}$ is a *square* diagonal matrix whose main diagonal includes the *singular values*. The matrix $\boldsymbol{\Sigma}$ is in the form:

$$\boldsymbol{\Sigma} = \begin{bmatrix} \sigma_1 & 0 & 0 \\ \vdots & \ddots & \vdots \\ 0 & 0 & \sigma_k \end{bmatrix}.$$

# Singular Value Decomposition

- Note that in both complete and incomplete SVD, the left singular vectors are orthonormal and the right singular vectors are also orthonormal; therefore, $U$ and $V$ are both orthogonal matrices so:

$$U^\top U = I, \tag{14}$$

$$V^\top V = I. \tag{15}$$

If these orthogonal matrices are not truncated and thus are square matrices, e.g. for complete SVD, there are also:

$$U U^\top = I, \tag{16}$$

$$V V^\top = I. \tag{17}$$

# Relation of SVD and EVD

$$(xy)^\top = y^\top x^\top$$

- In both complete and incomplete SVD of matrix $A$, the left and right singular vectors are the eigenvectors of $AA^\top$ and $A^\top A$, respectively, and the singular values are the square root of eigenvalues of either $AA^\top$ or $A^\top A$.

- Proof: There is:

$$AA^\top = (U\Sigma V^\top)(U\Sigma V^\top)^\top = U\Sigma V^\top V\Sigma U^\top = U\Sigma\Sigma U^\top = U\Sigma^2 U^\top,$$

which is the eigen-decomposition [11] of $AA^\top$ where the columns of $U$ are the eigenvectors and the diagonal of $\Sigma^2$ are the eigenvalues so the diagonal of $\Sigma$ are the square root of eigenvalues.

- Also:

$$A^\top A = (U\Sigma V^\top)^\top(U\Sigma V^\top) = V\Sigma U^\top U\Sigma V^\top = V\Sigma\Sigma V^\top = V\Sigma^2 V^\top,$$

which is the eigenvalue decomposition of $A^\top A$ where the columns of $V$ are the eigenvectors and the diagonal of $\Sigma^2$ are the eigenvalues, so the diagonal of $\Sigma$ are the square root of eigenvalues.

# Acknowledgement

# References

[1] B. Ghojogh, M. Crowley, F. Karray, and A. Ghodsi, *Elements of Dimensionality Reduction and Manifold Learning*.
Springer Nature, 2023.

[2] S. Garcia, J. Derrac, J. Cano, and F. Herrera, "Prototype selection for nearest neighbor classification: Taxonomy and empirical study," *IEEE transactions on pattern analysis and machine intelligence*, vol. 34, no. 3, pp. 417–435, 2012.

[3] I. Triguero, J. Derrac, S. Garcia, and F. Herrera, "A taxonomy and experimental study on prototype generation for nearest neighbor classification," *IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews)*, vol. 42, no. 1, pp. 86–100, 2011.

[4] B. Ghojogh, "Data reduction algorithms in machine learning and data science," 2021.

[5] I. Jolliffe, *Principal component analysis*.
Springer, 2011.

[6] J. H. Wilkinson, *The algebraic eigenvalue problem*, vol. 662.
Oxford Clarendon, 1965.

[7] G. H. Golub and C. F. Van Loan, *Matrix computations*, vol. 3.
The Johns Hopkins University Press, 2012.

[8] B. N. Parlett, "The symmetric eigenvalue problem," *Classics in Applied Mathematics*, vol. 20, 1998.

# References (cont.)

[9]  G. W. Stewart, "On the early history of the singular value decomposition," *SIAM review*, vol. 35, no. 4, pp. 551–566, 1993.

[10] G. H. Golub and C. Reinsch, "Singular value decomposition and least squares solutions," *Numerische mathematik*, vol. 14, no. 5, pp. 403–420, 1970.

[11] B. Ghojogh, F. Karray, and M. Crowley, "Eigenvalue and generalized eigenvalue problems: Tutorial," *arXiv preprint arXiv:1903.11240*, 2019.