

Preliminaries

Statistical Machine Learning (ENGG*6600*02)

School of Engineering,
University of Guelph, ON, Canada

Course Instructor: Benyamin Ghojogh
Summer 2023

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

Dataset

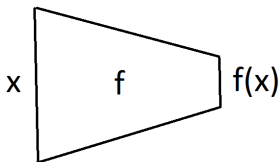
- Consider the measurement of a quantity. This quantity can be:
 - ▶ personal health data, including blood pressure, blood sugar, and blood fat,
 - ▶ images from a specific scene but taken from different perspectives,
 - ▶ images from several categories of animals, such as cat, dog, frog, etc.,
 - ▶ 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 $\{\mathbf{x}_i \in \mathbb{R}^d\}_{i=1}^n$, each of which is a d -dimensional vector in the d -dimensional Euclidean space. We can put these vectors column-wise in a matrix $\mathbf{X} \in \mathbb{R}^{d \times n}$.
- Consider a learning model f which is a map from data space to some output space:

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

- 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

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

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

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

- **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

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

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

**Expectation, Variance,
Covariance, Bias, MSE**

Random Variable, PDF, PMF, CDF

- **Random variable** X is a mathematical formalization of a quantity which depends on random events.
 - ▶ **Discrete** random variable: a countable set of possible values, e.g., $\{x_1, \dots, x_m\}$
 - ▶ **Continuous** random variable: an uncountable set of possible values, e.g., $(0, 1]$
- **Probability mass function (PMF)** shows the distribution of a discrete random variable:

$$f(x) = \mathbb{P}(X = x), \quad x \in \{x_1, \dots, x_m\}. \quad (2)$$

- **Probability density function (PDF)** shows the distribution of a continuous random variable:

$$f(x) = \lim_{\Delta x \rightarrow 0} \frac{\mathbb{P}(x < X < x + \Delta x)}{\Delta x}. \quad (3)$$

- **Cumulative distribution function (CDF)** shows the cumulative probabilities until that value:

$$F(x) = \mathbb{P}(X \leq x). \quad (4)$$

Expectation

- Let X be a random variable with probability mass/density function $f(x)$.
- **Expectation** or **expected value** is the mean of distribution.
- Expectation for discrete data, $x \in \{x_1, x_2, \dots, x_m\}$:

$$\mathbb{E}(X) = \sum_{i=1}^m x_i f(x_i) = x_1 f(x_1) + \dots + x_m f(x_m). \quad (5)$$

- Expectation for continuous data, $x \in \mathcal{D}$:

$$\mathbb{E}(X) = \int_{\mathcal{D}} x f(x) dx. \quad (6)$$

- If $h(x)$ is a function on X , then:

$$\mathbb{E}(h(x)) = \int_{\mathcal{D}} h(x) f(x) dx \quad (7)$$

- Expectation is a linear operator:

$$\mathbb{E}\left(\sum_{i=1}^k a_i X_i\right) = \sum_{i=1}^k a_i \mathbb{E}(X_i). \quad (8)$$

Variance

- Assume we have variable X and we estimate it. Let the random variable \hat{X} denote the estimate of X . Let $\mathbb{E}(\cdot)$ and $\mathbb{P}(\cdot)$ denote expectation and probability, respectively.
- The **variance** of estimating this random variable is defined as:

$$\text{Var}(\hat{X}) := \mathbb{E}((\hat{X} - \mathbb{E}(\hat{X}))^2), \quad (9)$$

which means average deviation of \hat{X} from the mean of our estimate, $\mathbb{E}(\hat{X})$, where the deviation is squared for symmetry of difference.

- This variance can be restated as:

$$\text{Var}(\hat{X}) = \mathbb{E}(\hat{X}^2) - (\mathbb{E}(\hat{X}))^2. \quad (10)$$

- Proof:

$$\begin{aligned} \text{Var}(\hat{X}) &= \mathbb{E}(\hat{X}^2 + (\mathbb{E}(\hat{X}))^2 - 2\hat{X}\mathbb{E}(\hat{X})) \\ &\stackrel{(a)}{=} \mathbb{E}(\hat{X}^2) + (\mathbb{E}(\hat{X}))^2 - 2\mathbb{E}(\hat{X})\mathbb{E}(\hat{X}) \\ &= \mathbb{E}(\hat{X}^2) - (\mathbb{E}(\hat{X}))^2, \end{aligned}$$

where (a) is because expectation is a linear operator and $\mathbb{E}(\hat{X})$ is not a random variable.

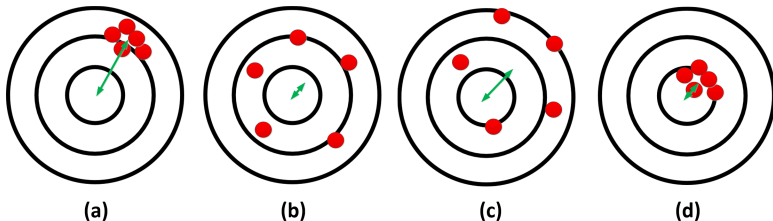
Bias

- Our estimation can have a bias. The **bias** of our estimate is defined as:

$$\text{Bias}(\hat{X}) := \mathbb{E}(\hat{X}) - X, \quad (11)$$

which means how much the mean of our estimate deviates from the original X .

- If the bias of an estimator is zero, i.e., $\mathbb{E}(\hat{X}) = X$, the estimator is **unbiased**. Otherwise, it is **biased**.
- The trade-off of bias and variance (dart example):



Mean Squared Error (MSE)

- The **Mean Squared Error (MSE)** of our estimate, \hat{X} , is defined as:

$$\text{MSE}(\hat{X}) := \mathbb{E}((\hat{X} - X)^2), \quad (12)$$

which means how much our estimate deviates from the original X .

- The relation of MSE, variance, and bias is as follows:

$$\text{MSE}(\hat{X}) = \text{Var}(\hat{X}) + (\text{Bias}(\hat{X}))^2. \quad (13)$$

- Proof:

$$\begin{aligned} \text{MSE}(\hat{X}) &= \mathbb{E}((\hat{X} - X)^2) = \mathbb{E}((\hat{X} - \mathbb{E}(\hat{X}) + \mathbb{E}(\hat{X}) - X)^2) \\ &= \mathbb{E}((\hat{X} - \mathbb{E}(\hat{X}))^2 + (\mathbb{E}(\hat{X}) - X)^2 + 2(\hat{X} - \mathbb{E}(\hat{X}))(\mathbb{E}(\hat{X}) - X)) \\ &\stackrel{(a)}{=} \mathbb{E}((\hat{X} - \mathbb{E}(\hat{X}))^2) + (\mathbb{E}(\hat{X}) - X)^2 + \underbrace{2(\mathbb{E}(\hat{X}) - \mathbb{E}(\hat{X}))(\mathbb{E}(\hat{X}) - X)}_0 \\ &\stackrel{(b)}{=} \text{Var}(\hat{X}) + (\text{Bias}(\hat{X}))^2, \end{aligned}$$

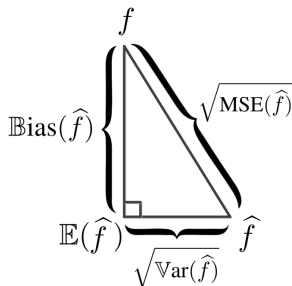
where (a) is because expectation is a linear operator and X and $\mathbb{E}(\hat{X})$ are not random, and (b) is because of Eqs. (9) and (11).

Relation of Variance, Bias, and MSE

- The relation of MSE, variance, and bias is as follows:

$$\text{MSE}(\hat{X}) = \text{Var}(\hat{X}) + (\text{Bias}(\hat{X}))^2.$$

- According to Pythagorean theorem:



Covariance

- **Covariance** is defined as:

$$\mathbb{Cov}(\hat{X}, \hat{Y}) := \mathbb{E}\left((\hat{X} - \mathbb{E}(\hat{X}))(\hat{Y} - \mathbb{E}(\hat{Y}))\right). \quad (14)$$

- We can restate it as:

$$\mathbb{Cov}(\hat{X}, \hat{Y}) = \mathbb{E}(\hat{X}\hat{Y}) - \mathbb{E}(\hat{X})\mathbb{E}(\hat{Y}). \quad (15)$$

- Proof:

$$\begin{aligned} \mathbb{Cov}(\hat{X}, \hat{Y}) &= \mathbb{E}\left((\hat{X} - \mathbb{E}(\hat{X}))(\hat{Y} - \mathbb{E}(\hat{Y}))\right) \\ &= \mathbb{E}\left(\hat{X}\hat{Y} - \hat{X}\mathbb{E}(\hat{Y}) - \mathbb{E}(\hat{X})\hat{Y} + \mathbb{E}(\hat{X})\mathbb{E}(\hat{Y})\right) \\ &= \mathbb{E}(\hat{X}\hat{Y}) - \mathbb{E}(\hat{X}\mathbb{E}(\hat{Y})) - \mathbb{E}(\mathbb{E}(\hat{X})\hat{Y}) + \mathbb{E}(\mathbb{E}(\hat{X})\mathbb{E}(\hat{Y})) \\ &= \mathbb{E}(\hat{X}\hat{Y}) - \mathbb{E}(\hat{X})\mathbb{E}(\hat{Y}) - \mathbb{E}(\hat{X})\mathbb{E}(\hat{Y}) + \mathbb{E}(\hat{X})\mathbb{E}(\hat{Y}) \\ &= \mathbb{E}(\hat{X}\hat{Y}) - \mathbb{E}(\hat{X})\mathbb{E}(\hat{Y}). \end{aligned}$$

Variance and Covariance

- If we have two random variables \hat{X} and \hat{Y} , we can say:

$$\begin{aligned}\text{Var}(a\hat{X} + b\hat{Y}) &\stackrel{(10)}{=} \mathbb{E}((a\hat{X} + b\hat{Y})^2) - (\mathbb{E}(a\hat{X} + b\hat{Y}))^2 \\ &\stackrel{(a)}{=} a^2 \mathbb{E}(\hat{X}^2) + b^2 \mathbb{E}(\hat{Y}^2) + 2ab \mathbb{E}(\hat{X}\hat{Y}) - a^2 (\mathbb{E}(\hat{X}))^2 - b^2 (\mathbb{E}(\hat{Y}))^2 - 2ab \mathbb{E}(\hat{Y})\mathbb{E}(\hat{X}) \\ &\stackrel{(10)}{=} a^2 \text{Var}(\hat{X}) + b^2 \text{Var}(\hat{Y}) + 2ab \text{Cov}(\hat{X}, \hat{Y}),\end{aligned}\tag{16}$$

where (a) is for linearity of expectation.

- If the two random variables are independent, i.e., $X \perp\!\!\!\perp Y$, we have:

$$\begin{aligned}\mathbb{E}(\hat{X}\hat{Y}) &\stackrel{(a)}{=} \iint \hat{x}\hat{y}f(\hat{x}, \hat{y})d\hat{x}d\hat{y} \stackrel{\perp\!\!\!\perp}{=} \iint \hat{x}\hat{y}f(\hat{x})f(\hat{y})d\hat{x}d\hat{y} \\ &= \int \hat{y}f(\hat{y}) \underbrace{\int \hat{x}f(\hat{x})d\hat{x}}_{\mathbb{E}(\hat{X})} d\hat{y} = \mathbb{E}(\hat{X}) \underbrace{\int \hat{y}f(\hat{y})d\hat{y}}_{\mathbb{E}(\hat{Y})} = \mathbb{E}(\hat{X})\mathbb{E}(\hat{Y}) \implies \text{Cov}(\hat{X}, \hat{Y}) = 0,\end{aligned}\tag{17}$$

where (a) is according to definition of expectation. Note that Eq. (17) is not true for the reverse implication (we can prove by counterexample).

- So, if two random variables \hat{X} and \hat{Y} are independent, then:

$$X \perp\!\!\!\perp Y \implies \mathbb{E}(\hat{X}\hat{Y}) = \mathbb{E}(\hat{X})\mathbb{E}(\hat{Y}) \implies \text{Cov}(\hat{X}, \hat{Y}) = 0.\tag{18}$$

Variance and Covariance

- We can extend Eqs. (16) and (15) to multiple random variables:

$$\mathbb{V}\text{ar}\left(\sum_{i=1}^k a_i X_i\right) = \sum_{i=1}^k a_i^2 \mathbb{V}\text{ar}(X_i) + \sum_{i=1}^k \sum_{j=1, j \neq i}^k a_i a_j \mathbb{C}\text{ov}(X_i, X_j), \quad (19)$$

$$\mathbb{C}\text{ov}\left(\sum_{i=1}^{k_1} a_i X_i, \sum_{j=1}^{k_2} b_j Y_j\right) = \sum_{i=1}^{k_1} \sum_{j=1}^{k_2} a_i b_j \mathbb{C}\text{ov}(X_i, Y_j), \quad (20)$$

where a_i 's and b_j 's are not random.

- If X_i 's are independent, we have:

$$\mathbb{V}\text{ar}\left(\sum_{i=1}^k a_i X_i\right) = \sum_{i=1}^k a_i^2 \mathbb{V}\text{ar}(X_i). \quad (21)$$

Dependence and Correlation

- **Dependence** is any dependence of two random variables.
- **Correlation** is linear dependence.
- **Pearson correlation coefficient** is defined as:

$$\mathbb{C}\text{orr}(X, Y) := \frac{\mathbb{C}\text{ov}(X, Y)}{\sqrt{\mathbb{V}\text{ar}(X)\mathbb{V}\text{ar}(Y)}}.$$

$$-1 \leq \mathbb{C}\text{orr}(X, Y) \leq 1.$$

Conditional Expectation and Variance

- Conditional probability:

$$\mathbb{P}(X, Y) = \mathbb{P}(X|Y) \mathbb{P}(Y) = \mathbb{P}(Y|X) \mathbb{P}(X). \quad (22)$$

- Bayes' rule:

$$\mathbb{P}(X|Y) \mathbb{P}(Y) = \mathbb{P}(Y|X) \mathbb{P}(X) \implies \mathbb{P}(X|Y) = \frac{\mathbb{P}(Y|X) \mathbb{P}(X)}{\mathbb{P}(Y)}. \quad (23)$$

- We can have:

$$\mathbb{P}(X|Y) = \frac{\mathbb{P}(Y, X)}{\mathbb{P}(Y)}. \quad (24)$$

- **Law of total probability** or **marginalization**:

$$f(y) = \sum_x f(y|x) f(x), \quad (25)$$

$$f(y) = \int_x f(y|x) f(x) dx. \quad (26)$$

Conditional Expectation and Variance

- Conditional expectation:

$$\mathbb{E}(Y|X) := \sum_y y f(y|x) \quad (27)$$

- **Law of total expectation or Adam's law:**

$$\mathbb{E}(Y) = \mathbb{E}(\mathbb{E}(Y|X)) \quad (28)$$

- Proof:

$$\begin{aligned} \mathbb{E}(Y) &= \sum_y y f(y) = \sum_y y \sum_x f(y|x) f(x) = \sum_x \left[\sum_y y f(y|x) \right] f(x) \\ &= \sum_x \mathbb{E}(Y|X) f(x) = \mathbb{E}(\mathbb{E}(Y|X)). \end{aligned}$$

Conditional Expectation and Variance

- **Law of total variance or Eve's law:**

$$\text{Var}(Y) = \mathbb{E}(\text{Var}(Y|X)) + \text{Var}(\mathbb{E}(Y|X)). \quad (29)$$

- **Proof:**

$$\text{Var}(Y) = \mathbb{E}(Y^2) - \mathbb{E}(Y)^2 \implies \mathbb{E}(Y^2) = \text{Var}(Y) + \mathbb{E}(Y)^2.$$

By law of total expectation:

$$\begin{aligned} \mathbb{E}(Y^2) &= \mathbb{E}(\text{Var}(Y|X) + \mathbb{E}(Y|X)^2) \\ &\implies \mathbb{E}(Y^2) - \mathbb{E}(Y)^2 = \mathbb{E}(\text{Var}(Y|X) + \mathbb{E}(Y|X)^2) - \mathbb{E}(Y)^2 \\ &= \mathbb{E}(\text{Var}(Y|X) + \mathbb{E}(Y|X)^2) - \mathbb{E}(\mathbb{E}(Y|X))^2 \\ &= \mathbb{E}(\text{Var}(Y|X)) + (\mathbb{E}(\mathbb{E}(Y|X)^2) - \mathbb{E}(\mathbb{E}(Y|X))^2) = \mathbb{E}(\text{Var}(Y|X)) + \text{Var}(\mathbb{E}(Y|X)). \end{aligned}$$

Monte Carlo Approximation

- Suppose we are considering some d -dimensional data $x \in \mathbb{R}^d$. Let $f(x)$ be the Probability Density Function (PDF) of data. Consider $h(x)$ is a function over the data x . According to definition, the expectation of function $h(x)$ over the distribution $f(x)$ and the probability of function $h(x)$ belonging to a set \mathcal{A} are:

$$\mathbb{E}(h(x)) = \int h(x) f(x) dx, \quad (30)$$

$$\mathbb{P}(h(x) \in \mathcal{A}) = \int_{h(x) \in \mathcal{A}} f(x) dx, \quad (31)$$

respectively.

- Using a sample of size n from distribution $f(x)$ (i.e., $\{x_1, \dots, x_n\} \sim f(x)$), we can approximate Eqs. (30) and (31) by (**Monte Carlo approximation**):

$$\mathbb{E}(h(x)) \approx \frac{1}{n} \sum_{i=1}^n h(x_i), \quad (32)$$

$$\mathbb{P}(h(x) \in \mathcal{A}) \approx \frac{1}{n} \sum_{i=1}^n \mathbb{I}(h(x_i) \in \mathcal{A}), \quad (33)$$

where $\mathbb{I}(\cdot)$ denotes the indicator function which is one and zero when its condition is and is not satisfied, respectively.

- As the above definition states, the MC approximation generates many samples from the distribution in order to approximate the expectation by mean (or average) of the samples. Obviously, the more the n is, the better the approximation becomes.

Sample Mean and Covariance

- Assume we have n data points $\{\mathbf{x}_i \in \mathbb{R}^d\}_{i=1}^n$.
- Sample mean:

$$\mathbb{R}^d \ni \boldsymbol{\mu} := \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i. \quad (34)$$

- Sample covariance matrix:

$$\mathbb{R}^{d \times d} \ni \boldsymbol{\Sigma} := \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^\top. \quad (35)$$

- For this to be an unbiased sample covariance, n should be $n - 1$ in the denominator:

$$\mathbb{R}^{d \times d} \ni \boldsymbol{\Sigma} := \frac{1}{n-1} \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^\top. \quad (36)$$

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 $\{\mathbf{u}_1, \dots, \mathbf{u}_p\}$. These vectors can be stacked columnwise in matrix $\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_p] \in \mathbb{R}^{d \times p}$.
- The space spanned by the columns of matrix \mathbf{U} is called the column space of matrix \mathbf{U} , denoted by $\text{Col}(\mathbf{U})$:

$$\text{Col}(\mathbf{U}) := \text{span}\{\mathbf{u}_1, \dots, \mathbf{u}_p\}. \quad (37)$$

- In other words, the space whose bases are the columns of matrix \mathbf{U} is called the column space of matrix \mathbf{U} .

Linear Projection

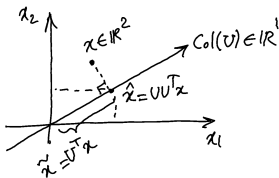
- Assume there is a data point $\mathbf{x} \in \mathbb{R}^d$. The aim is to project this data point onto the vector space spanned by p vectors $\{\mathbf{u}_1, \dots, \mathbf{u}_p\}$, where each vector is d -dimensional and usually $p \ll d$.
- These vectors can be stacked columnwise in matrix $\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_p] \in \mathbb{R}^{d \times p}$. In other words, the goal is to project \mathbf{x} onto the column space of \mathbf{U} , denoted by $\text{Col}(\mathbf{U})$.
- 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 $\mathbf{x} \in \mathbb{R}^d$ onto $\text{Col}(\mathbf{U}) \in \mathbb{R}^p$ is:

$$\mathbb{R}^p \ni \tilde{\mathbf{x}} := \mathbf{U}^\top \mathbf{x}. \quad (38)$$

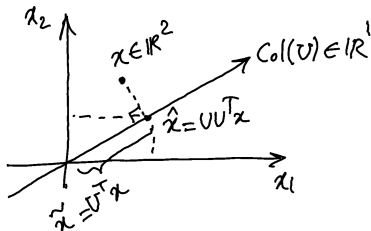
- The **reconstruction** of $\tilde{\mathbf{x}} \in \mathbb{R}^p$ in the d -dimensional space is:

$$\mathbb{R}^d \ni \hat{\mathbf{x}} := \mathbf{U}\tilde{\mathbf{x}} = \mathbf{U}\mathbf{U}^\top \mathbf{x}. \quad (39)$$

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

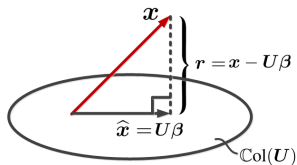


Linear Projection



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

$$r = x - \hat{x} = x - UU^T x. \quad (40)$$



Centering Matrix

Centering Matrix

- Consider a matrix $\mathbf{A} \in \mathbb{R}^{\alpha \times \beta}$, which is represented by its rows, $\mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_\alpha]^\top$ or by its columns, $\mathbf{A} = [\mathbf{b}_1, \dots, \mathbf{b}_\beta]$, where \mathbf{a}_i and \mathbf{b}_j denote the i -th row and j -th column of \mathbf{A} , respectively. Note that the vectors are column vectors.
- The **left centering matrix** is defined as:

$$\mathbb{R}^{\alpha \times \alpha} \ni \mathbf{H} := \mathbf{I} - (1/\alpha)\mathbf{1}\mathbf{1}^\top, \quad (41)$$

where $\mathbf{1} = [1, \dots, 1]^\top \in \mathbb{R}^\alpha$ and $\mathbf{I} \in \mathbb{R}^{\alpha \times \alpha}$ is the identity matrix.

- Left multiplying this matrix by \mathbf{A} , i.e., $\mathbf{H}\mathbf{A}$, removes the mean of rows of \mathbf{A} from all of its rows:

$$\mathbf{H}\mathbf{A} \stackrel{(41)}{=} \mathbf{A} - (1/\alpha)\mathbf{1}\mathbf{1}^\top \mathbf{A} = (\mathbf{A}^\top - \boldsymbol{\mu}_{\text{rows}})^\top, \quad (42)$$

where the column vector $\boldsymbol{\mu}_{\text{rows}} \in \mathbb{R}^\beta$ is the mean of the rows of \mathbf{A} .

- The **right centering matrix** is defined as:

$$\mathbb{R}^{\beta \times \beta} \ni \mathbf{H} := \mathbf{I} - (1/\beta)\mathbf{1}\mathbf{1}^\top, \quad (43)$$

where $\mathbf{1} = [1, \dots, 1]^\top \in \mathbb{R}^\beta$ and $\mathbf{I} \in \mathbb{R}^{\beta \times \beta}$ is the identity matrix.

- Right multiplying this matrix to \mathbf{A} , i.e., $\mathbf{A}\mathbf{H}$, removes the mean of the columns of \mathbf{A} from all of its columns:

$$\mathbf{A}\mathbf{H} \stackrel{(43)}{=} \mathbf{A} - (1/\beta)\mathbf{A}\mathbf{1}\mathbf{1}^\top = \mathbf{A} - \boldsymbol{\mu}_{\text{cols}}, \quad (44)$$

where the column vector $\boldsymbol{\mu}_{\text{cols}} \in \mathbb{R}^\alpha$ is the mean of the columns of \mathbf{A} .

Centering Matrix

- Both left and right centering matrices can be used at the same time to have a **double-centered** matrix \mathbf{A} :

$$\begin{aligned} \mathbf{H}\mathbf{A}\mathbf{H} &= (\mathbf{I}_\alpha - (1/\alpha)\mathbf{1}_\alpha\mathbf{1}_\alpha^\top)\mathbf{A}(\mathbf{I}_\beta - (1/\beta)\mathbf{1}_\beta\mathbf{1}_\beta^\top) \\ &= (\mathbf{A} - (1/\alpha)\mathbf{1}_\alpha\mathbf{1}_\alpha^\top\mathbf{A})(\mathbf{I}_\beta - (1/\beta)\mathbf{1}_\beta\mathbf{1}_\beta^\top) \\ &= \mathbf{A} - (1/\alpha)\mathbf{1}_\alpha\mathbf{1}_\alpha^\top\mathbf{A} - (1/\beta)\mathbf{A}\mathbf{1}_\beta\mathbf{1}_\beta^\top + (1/(\alpha\beta))\mathbf{1}_\alpha\mathbf{1}_\alpha^\top\mathbf{A}\mathbf{1}_\beta\mathbf{1}_\beta^\top. \end{aligned} \quad (45)$$

- The second term removes the mean of the rows of \mathbf{A} according to Eq. (42) and the third term removes the mean of the columns of \mathbf{A} according to Eq. (44). The last term, however, adds the overall mean of \mathbf{A} back to it.
- In summary:

$$\begin{aligned} \mathbf{H}\mathbf{A} &\approx (\mathbf{A}^\top - \mu_{\text{rows}})^\top, \\ \mathbf{A}\mathbf{H} &\approx \mathbf{A} - \mu_{\text{cols}}, \\ \mathbf{H}\mathbf{A}\mathbf{H} &\approx (\mathbf{A}^\top - \mu_{\text{rows}})^\top - \mu_{\text{cols}} + \mu_{\text{all}}. \end{aligned}$$

Centering Matrix

- The centering matrix is symmetric because:

$$\mathbf{H}^\top = (\mathbf{I} - (1/\alpha)\mathbf{1}\mathbf{1}^\top)^\top = \mathbf{I}^\top - (1/\alpha)(\mathbf{1}\mathbf{1}^\top)^\top = \mathbf{I} - (1/\alpha)\mathbf{1}\mathbf{1}^\top \stackrel{(41)}{=} \mathbf{H}. \quad (46)$$

- The centering matrix is also idempotent:

$$\mathbf{H}^k = \underbrace{\mathbf{H}\mathbf{H}\cdots\mathbf{H}}_{k \text{ times}} = \mathbf{H}, \quad (47)$$

where k is a positive integer. Proof:

$$\begin{aligned} \mathbf{H}\mathbf{H} &= (\mathbf{I} - (1/\alpha)\mathbf{1}\mathbf{1}^\top)(\mathbf{I} - (1/\alpha)\mathbf{1}\mathbf{1}^\top) \\ &= \mathbf{I} - (1/\alpha)\mathbf{1}\mathbf{1}^\top - (1/\alpha)\mathbf{1}\mathbf{1}^\top + (1/\alpha^2)\underbrace{\mathbf{1}\mathbf{1}^\top\mathbf{1}\mathbf{1}^\top}_{\alpha} \\ &= \mathbf{I} - (1/\alpha)\mathbf{1}\mathbf{1}^\top - (1/\alpha)\mathbf{1}\mathbf{1}^\top + (1/\alpha)\mathbf{1}\mathbf{1}^\top = \mathbf{I} - (1/\alpha)\mathbf{1}\mathbf{1}^\top \stackrel{(41)}{=} \mathbf{H}. \end{aligned}$$

Therefore:

$$\mathbf{H}^k = (\underbrace{\underbrace{\underbrace{\mathbf{H}(\mathbf{H}\mathbf{H}))}_{\mathbf{H}}}_{\mathbf{H}})_{\mathbf{H}} = \mathbf{H}.$$

Norm

Inner product

Definition (Inner product of vectors)

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

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

Definition (Inner product of matrices)

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

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

where $\text{tr}(\cdot)$ denotes the trace of matrix.

Norm

Definition (Norm)

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

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

Important norms for vectors

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

- The ℓ_p **norm** is:

$$\|\mathbf{x}\|_p := (|x_1|^p + \dots + |x_d|^p)^{1/p},$$

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

- Two well-known ℓ_p norms are ℓ_1 **norm** and ℓ_2 **norm** (also called the **Euclidean norm**) with $p = 1$ and $p = 2$, respectively:

$$\begin{aligned}\|\mathbf{x}\|_1 &:= |x_1| + \dots + |x_d| = \sum_{i=1}^d |x_i|, \\ \|\mathbf{x}\|_2 &:= \sqrt{x_1^2 + \dots + x_d^2} = \sqrt{\sum_{i=1}^d x_i^2},\end{aligned}$$

Important norms for matrices

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

- The formulation of the **Frobenius norm** for a matrix is similar to the formulation of ℓ_2 norm for a vector:

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

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

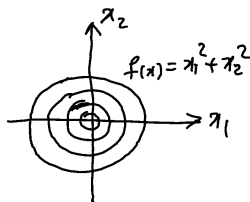
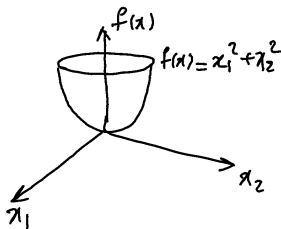
Quadratic forms using norms

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

$$\|\mathbf{x}\|_2^2 = \mathbf{x}^\top \mathbf{x} = \langle \mathbf{x}, \mathbf{x} \rangle = \sum_{i=1}^d x_i^2,$$

$$\|\mathbf{X}\|_F^2 = \text{tr}(\mathbf{X}^\top \mathbf{X}) = \langle \mathbf{X}, \mathbf{X} \rangle = \sum_{i=1}^{d_1} \sum_{j=1}^{d_2} x_{i,j}^2,$$

which are convex and in quadratic forms.



Preliminaries on Derivatives

Dimensionality of derivative

- Consider a function $f : \mathbb{R}^{d_1} \rightarrow \mathbb{R}^{d_2}$, $f : \mathbf{x} \mapsto f(\mathbf{x})$.
- Derivative of function $f(\mathbf{x}) \in \mathbb{R}^{d_2}$ with respect to (w.r.t.) $\mathbf{x} \in \mathbb{R}^{d_1}$ has dimensionality $(d_1 \times d_2)$.
- This is because tweaking every element of $\mathbf{x} \in \mathbb{R}^{d_1}$ can change every element of $f(\mathbf{x}) \in \mathbb{R}^{d_2}$. 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 $f(\mathbf{x})$ resulted by changing the i -th element of \mathbf{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 $f : \mathbb{R} \rightarrow \mathbb{R}, f : x \mapsto f(x)$, 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 $f : \mathbb{R}^d \rightarrow \mathbb{R}, f : \mathbf{x} \mapsto f(\mathbf{x})$, 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 $f : \mathbb{R}^{d_1 \times d_2} \rightarrow \mathbb{R}, f : \mathbf{X} \mapsto f(\mathbf{X})$, 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 $f : \mathbb{R}^{d_1} \rightarrow \mathbb{R}^{d_2}, f : \mathbf{x} \mapsto f(\mathbf{x})$, 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 $f : \mathbb{R}^{d_1 \times d_2} \rightarrow \mathbb{R}^{d_3}, f : \mathbf{X} \mapsto f(\mathbf{X})$, 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 $f : \mathbb{R}^{d_1 \times d_2} \rightarrow \mathbb{R}^{d_3 \times d_4}, f : \mathbf{X} \mapsto f(\mathbf{X})$, 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 $f : \mathbb{R}^d \rightarrow \mathbb{R}$, $f : \mathbf{x} \mapsto f(\mathbf{x})$. In optimizing the function f , the derivative of function w.r.t. its variable \mathbf{x} is called the **gradient**, denoted by:

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

Definition (Hessian)

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

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

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 $\text{vec}(\cdot)$ denote vectorization of a $\mathbb{R}^{a \times b}$ matrix to a \mathbb{R}^{ab} vector.
- Let $\text{vec}_{a \times b}^{-1}(\cdot)$ be de-vectorization of a \mathbb{R}^{ab} vector to a $\mathbb{R}^{a \times b}$ matrix.

Optimization

Optimization

- Lagrangian:

$$\begin{aligned} & \underset{\mathbf{x}}{\text{minimize}} && f(\mathbf{x}) \\ & \text{subject to} && y_i(\mathbf{x}) \leq 0, \quad i \in \{1, \dots, m_1\}, \\ & && h_i(\mathbf{x}) = 0, \quad i \in \{1, \dots, m_2\}. \end{aligned}$$

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

$$\frac{\partial \mathcal{L}}{\partial \mathbf{x}} \stackrel{\text{set}}{=} \mathbf{0} \implies \dots$$

- Unconstrained optimization:

$$\underset{\mathbf{x}}{\text{minimize}} \quad f(\mathbf{x}),$$

$$\mathbf{x}^{(k+1)} := \mathbf{x}^{(k)} + (\Delta \mathbf{x})^{(k)},$$

$$\text{Gradient method: } \mathbf{x}^{(k+1)} := \mathbf{x}^{(k)} - \eta^{(k)} \nabla f(\mathbf{x}^{(k)}),$$

$$\text{Newton's method: } \mathbf{x}^{(k+1)} := \mathbf{x}^{(k)} - \eta^{(k)} (\nabla^2 f(\mathbf{x}^{(k)}))^{-1} \nabla f(\mathbf{x}^{(k)}).$$

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

Rayleigh-Ritz Quotient

Rayleigh-Ritz Quotient

- The **Rayleigh-Ritz quotient** or **Rayleigh quotient** is defined as [5, 6]:

$$\mathbb{R} \ni R(\mathbf{A}, \mathbf{x}) := \frac{\mathbf{x}^\top \mathbf{A} \mathbf{x}}{\mathbf{x}^\top \mathbf{x}}, \quad (48)$$

where \mathbf{A} is a symmetric matrix and \mathbf{x} is a nonzero vector:

$$\mathbf{A} = \mathbf{A}^\top, \quad \mathbf{x} \neq \mathbf{0}. \quad (49)$$

- One of the properties of the Rayleigh-Ritz quotient is:

$$R(\mathbf{A}, c\mathbf{x}) = R(\mathbf{A}, \mathbf{x}), \quad (50)$$

where c is a scalar. Proof:

$$R(\mathbf{A}, c\mathbf{x}) = \frac{(c\mathbf{x})^\top \mathbf{A} c\mathbf{x}}{(c\mathbf{x})^\top c\mathbf{x}} \stackrel{(a)}{=} \frac{c\mathbf{x}^\top \mathbf{A} c\mathbf{x}}{c\mathbf{x}^\top c\mathbf{x}} \stackrel{(b)}{=} \frac{c^2}{c^2} \times \frac{\mathbf{x}^\top \mathbf{A} \mathbf{x}}{\mathbf{x}^\top \mathbf{x}} \stackrel{(48)}{=} R(\mathbf{A}, \mathbf{x}),$$

where (a) and (b) are because c is a scalar.

Rayleigh-Ritz Quotient

- Consider the optimization problem of the Rayleigh-Ritz quotient:

$$\underset{\mathbf{x}}{\text{minimize/maximize}} \quad R(\mathbf{A}, \mathbf{x}). \quad (51)$$

- According to Eq. (50), this is equivalent to the following problem [6]:

$$\begin{aligned} &\underset{\mathbf{x}}{\text{minimize/maximize}} \quad R(\mathbf{A}, \mathbf{x}) \\ &\text{subject to} \quad \|\mathbf{x}\|_2 = 1. \end{aligned} \quad (52)$$

- Proof: Let $\mathbf{y} := (1/\|\mathbf{x}\|_2) \mathbf{x}$. The Rayleigh-Ritz quotient is:

$$R(\mathbf{A}, \mathbf{y}) = \frac{\mathbf{y}^\top \mathbf{A} \mathbf{y}}{\mathbf{y}^\top \mathbf{y}} = \frac{1/\|\mathbf{x}\|_2^2}{1/\|\mathbf{x}\|_2^2} \times \frac{\mathbf{x}^\top \mathbf{A} \mathbf{x}}{\mathbf{x}^\top \mathbf{x}} = R(\mathbf{A}, \mathbf{x}).$$

Due to the following:

$$\|\mathbf{y}\|_2^2 = \frac{1}{\|\mathbf{x}\|_2^2} \times \|\mathbf{x}\|_2^2 = 1 \implies \|\mathbf{y}\|_2 = 1,$$

$R(\mathbf{A}, \mathbf{y})$ should be optimized subject to $\|\mathbf{y}\|_2 = 1$. Changing the dummy variable \mathbf{y} to \mathbf{x} gives Eq. (52).

Rayleigh-Ritz Quotient

- Another equivalent problem for Eq. (51) is:

$$\begin{aligned} & \underset{\mathbf{x}}{\text{minimize/maximize}} && \mathbf{x}^\top \mathbf{A} \mathbf{x} \\ & \text{subject to} && \|\mathbf{x}\|_2 = 1, \end{aligned} \tag{53}$$

obtained by inserting the constraint $\|\mathbf{x}\|_2^2 = \mathbf{x}^\top \mathbf{x} = 1$ in Eqs. (48) and (52).

- The constraint in Eqs. (52) and (53) can be equal to any constant (1 is a constant and disappears in derivative of Lagrangian).
- The **generalized Rayleigh-Ritz quotient** or **generalized Rayleigh quotient** is defined as [5]:

$$\mathbb{R} \ni R(\mathbf{A}, \mathbf{B}; \mathbf{x}) := \frac{\mathbf{x}^\top \mathbf{A} \mathbf{x}}{\mathbf{x}^\top \mathbf{B} \mathbf{x}}, \tag{54}$$

where \mathbf{A} and \mathbf{B} are symmetric matrices and \mathbf{x} is a nonzero vector:

$$\mathbf{A} = \mathbf{A}^\top, \quad \mathbf{B} = \mathbf{B}^\top, \quad \mathbf{x} \neq \mathbf{0}. \tag{55}$$

Rayleigh-Ritz Quotient

- If the symmetric \mathbf{B} is positive definite:

$$\mathbf{B} \succ 0, \quad (56)$$

it has a Cholesky decomposition:

$$\mathbf{B} = \mathbf{C}\mathbf{C}^\top, \quad (57)$$

where \mathbf{C} is a lower triangular matrix.

- In case $\mathbf{B} \succ 0$, the generalized Rayleigh-Ritz quotient can be converted to a Rayleigh-Ritz quotient:

$$R(\mathbf{A}, \mathbf{B}; \mathbf{x}) = R(\mathbf{D}, \mathbf{C}^\top \mathbf{x}), \quad (58)$$

where:

$$\mathbf{D} := \mathbf{C}^{-1} \mathbf{A} \mathbf{C}^{-\top}. \quad (59)$$

- Proof:

$$\begin{aligned} \text{RHS} &= R(\mathbf{D}, \mathbf{C}^\top \mathbf{x}) \stackrel{(48)}{=} \frac{(\mathbf{C}^\top \mathbf{x})^\top \mathbf{D} (\mathbf{C}^\top \mathbf{x})}{(\mathbf{C}^\top \mathbf{x})^\top (\mathbf{C}^\top \mathbf{x})} \\ &\stackrel{(59)}{=} \frac{\mathbf{x}^\top \mathbf{C} \mathbf{C}^{-1} \mathbf{A} (\mathbf{C} \mathbf{C}^{-1})^\top \mathbf{x}}{\mathbf{x}^\top (\mathbf{C} \mathbf{C}^\top) \mathbf{x}} \stackrel{(a)}{=} \frac{\mathbf{x}^\top \mathbf{A} \mathbf{x}}{\mathbf{x}^\top \mathbf{B} \mathbf{x}} \stackrel{(54)}{=} R(\mathbf{A}, \mathbf{B}; \mathbf{x}) = \text{LHS}, \end{aligned}$$

where RHS and LHS are short for right and left hand sides and (a) is because of Eq. (57) and $\mathbf{C}\mathbf{C}^{-1} = \mathbf{I}$ because \mathbf{C} is a square matrix.

Rayleigh-Ritz Quotient

- Similarly, one of the properties of the generalized Rayleigh-Ritz quotient is:

$$R(\mathbf{A}, \mathbf{B}; c\mathbf{x}) = R(\mathbf{A}, \mathbf{B}; \mathbf{x}), \quad (60)$$

where c is a scalar. Proof:

$$R(\mathbf{A}, \mathbf{B}; c\mathbf{x}) = \frac{(c\mathbf{x})^\top \mathbf{A} c\mathbf{x}}{(c\mathbf{x})^\top \mathbf{B} c\mathbf{x}} \stackrel{(a)}{=} \frac{c\mathbf{x}^\top \mathbf{A} c\mathbf{x}}{c\mathbf{x}^\top \mathbf{B} c\mathbf{x}} \stackrel{(b)}{=} \frac{c^2}{c^2} \times \frac{\mathbf{x}^\top \mathbf{A} \mathbf{x}}{\mathbf{x}^\top \mathbf{B} \mathbf{x}} \stackrel{(54)}{=} R(\mathbf{A}, \mathbf{B}; \mathbf{x}),$$

where (a) and (b) are because c is a scalar.

- Consider the optimization problem of the generalized Rayleigh-Ritz quotient:

$$\underset{\mathbf{x}}{\text{minimize/maximize}} \quad R(\mathbf{A}, \mathbf{B}; \mathbf{x}). \quad (61)$$

According to Eq. (60), it has an equivalent form:

$$\begin{aligned} &\underset{\mathbf{x}}{\text{minimize/maximize}} \quad \mathbf{x}^\top \mathbf{A} \mathbf{x} \\ &\text{subject to} \quad \mathbf{x}^\top \mathbf{B} \mathbf{x} = 1, \end{aligned} \quad (62)$$

for the same reason as the Rayleigh-Ritz quotient. The constraint can be equal to any constant because in the derivative of Lagrangian, the constant will be dropped.

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 [7].
- 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** [8, 9] of a symmetric matrix $\mathbf{A} \in \mathbb{R}^{d \times d}$ is defined as:

$$\mathbf{A}\phi_i = \lambda_i\phi_i, \quad \forall i \in \{1, \dots, d\}, \quad (63)$$

and in matrix form, it is:

$$\mathbf{A}\Phi = \Phi\Lambda, \quad (64)$$

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

Eigenvalue Problem

- For the eigenvalue problem, the matrix \mathbf{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 (64) can be restated as:

$$\mathbf{A}\Phi = \Phi\Lambda \implies \mathbf{A}\underbrace{\Phi\Phi^\top}_I = \Phi\Lambda\Phi^\top \implies \mathbf{A} = \Phi\Lambda\Phi^\top = \Phi\Lambda\Phi^{-1}, \quad (65)$$

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

- There is always $\Phi^\top\Phi = I$ for orthogonal Φ , but there is only $\Phi\Phi^\top = I$ if “all” columns of the orthogonal Φ exist (it is not truncated, i.e., it is a square matrix). Equation (65) is referred to as “**eigenvalue decomposition**”, “eigen-decomposition”, or “spectral decomposition”.

Generalized Eigenvalue Problem

- The **generalized eigenvalue problem** [5, 9] of two symmetric matrices $\mathbf{A} \in \mathbb{R}^{d \times d}$ and $\mathbf{B} \in \mathbb{R}^{d \times d}$ is defined as:

$$\mathbf{A}\phi_i = \lambda_i \mathbf{B}\phi_i, \quad \forall i \in \{1, \dots, d\}, \quad (66)$$

and in matrix form, it is:

$$\mathbf{A}\Phi = \mathbf{B}\Phi\mathbf{\Lambda}, \quad (67)$$

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

- The generalized eigenvalue problem of Eq. (66) or (67) is denoted by (\mathbf{A}, \mathbf{B}) .
- The (\mathbf{A}, \mathbf{B}) is called a “pair” or “pencil” [5], and the order in the pair matters, according to Eq. (67).
- The Φ and $\mathbf{\Lambda}$ are called the generalized eigenvectors and eigenvalues of (\mathbf{A}, \mathbf{B}) .
- The $(\Phi, \mathbf{\Lambda})$ or (ϕ_i, λ_i) is called the “eigenpair” of the pair (\mathbf{A}, \mathbf{B}) in the literature [5].
- Comparing Eqs. (63) and (66) or Eqs. (64) and (67) demonstrates that the eigenvalue problem is a special case of the generalized eigenvalue problem where $\mathbf{B} = \mathbf{I}$.

Optimization Forms of (Generalized) Eigenvalue Problem

- **Optimization Form 1:**

$$\begin{aligned} & \underset{\phi}{\text{maximize}} && \phi^\top \mathbf{A} \phi, \\ & \text{subject to} && \phi^\top \mathbf{B} \phi = 1, \end{aligned} \tag{68}$$

where $\mathbf{A} \in \mathbb{R}^{d \times d}$ and $\mathbf{B} \in \mathbb{R}^{d \times d}$. The Lagrangian for Eq. (68) is:

$$\mathcal{L} = \phi^\top \mathbf{A} \phi - \lambda (\phi^\top \mathbf{B} \phi - 1),$$

where $\lambda \in \mathbb{R}$ is the Lagrange multiplier. Equating the derivative of Lagrangian to zero gives us:

$$\mathbb{R}^d \ni \frac{\partial \mathcal{L}}{\partial \phi} = 2\mathbf{A}\phi - 2\lambda\mathbf{B}\phi \stackrel{\text{set}}{=} 0 \implies \mathbf{A}\phi = \lambda\mathbf{B}\phi,$$

which is a generalized eigenvalue problem (\mathbf{A}, \mathbf{B}) according to Eq. (66), where ϕ is the eigenvector and λ is the eigenvalue.

- As Eq. (68) is a *maximization* problem, the eigenvector is the one having the largest eigenvalue. If Eq. (68) is a *minimization* problem, the eigenvector is the one having the smallest eigenvalue.

Optimization Forms of (Generalized) Eigenvalue Problem

- **Optimization Form 2:**

$$\begin{aligned} & \underset{\Phi}{\text{maximize}} && \text{tr}(\Phi^\top \mathbf{A} \Phi), \\ & \text{subject to} && \Phi^\top \mathbf{B} \Phi = \mathbf{I}, \end{aligned} \tag{69}$$

where $\mathbf{A} \in \mathbb{R}^{d \times d}$ and $\mathbf{B} \in \mathbb{R}^{d \times d}$.

The Lagrangian for Eq. (69) is:

$$\mathcal{L} = \text{tr}(\Phi^\top \mathbf{A} \Phi) - \text{tr}(\Lambda^\top (\Phi^\top \mathbf{B} \Phi - \mathbf{I})),$$

where $\Lambda \in \mathbb{R}^{d \times d}$ is a diagonal matrix whose entries are the Lagrange multipliers. Equating the derivative of \mathcal{L} to zero gives us:

$$\mathbb{R}^{d \times d} \ni \frac{\partial \mathcal{L}}{\partial \Phi} = 2 \mathbf{A} \Phi - 2 \mathbf{B} \Phi \Lambda \stackrel{\text{set}}{=} \mathbf{0} \implies \mathbf{A} \Phi = \mathbf{B} \Phi \Lambda,$$

which is an eigenvalue problem (\mathbf{A}, \mathbf{B}) according to Eq. (67).

Optimization Forms of (Generalized) Eigenvalue Problem

- **Optimization Form 3:** Consider the following optimization problem with the variable $\phi \in \mathbb{R}^d$:

$$\begin{aligned} & \underset{\phi}{\text{minimize}} && ||\mathbf{X} - \phi \phi^\top \mathbf{X}||_F^2, \\ & \text{subject to} && \phi^\top \phi = 1, \end{aligned} \tag{70}$$

where $\mathbf{X} \in \mathbb{R}^{d \times n}$. The objective function in Eq. (70) is simplified as:

$$\begin{aligned} ||\mathbf{X} - \phi \phi^\top \mathbf{X}||_F^2 &= \text{tr}((\mathbf{X} - \phi \phi^\top \mathbf{X})^\top (\mathbf{X} - \phi \phi^\top \mathbf{X})) \\ &= \text{tr}((\mathbf{X}^\top - \mathbf{X}^\top \phi \phi^\top)(\mathbf{X} - \phi \phi^\top \mathbf{X})) = \text{tr}(\mathbf{X}^\top \mathbf{X} - 2\mathbf{X}^\top \phi \phi^\top \mathbf{X} + \underbrace{\mathbf{X}^\top \phi \phi^\top \phi \phi^\top \mathbf{X}}_1) \\ &= \text{tr}(\mathbf{X}^\top \mathbf{X} - \mathbf{X}^\top \phi \phi^\top \mathbf{X}) = \text{tr}(\mathbf{X}^\top \mathbf{X}) - \text{tr}(\mathbf{X}^\top \phi \phi^\top \mathbf{X}) \\ &= \text{tr}(\mathbf{X}^\top \mathbf{X}) - \text{tr}(\mathbf{X} \mathbf{X}^\top \phi \phi^\top) = \text{tr}(\mathbf{X}^\top \mathbf{X} - \mathbf{X} \mathbf{X}^\top \phi \phi^\top). \end{aligned}$$

The Lagrangian is:

$$\mathcal{L} = \text{tr}(\mathbf{X}^\top \mathbf{X}) - \text{tr}(\mathbf{X} \mathbf{X}^\top \phi \phi^\top) - \lambda(\phi^\top \phi - 1),$$

where λ is the Lagrange multiplier. Equating the derivative of \mathcal{L} to zero gives:

$$\mathbb{R}^d \ni \frac{\partial \mathcal{L}}{\partial \phi} = 2\mathbf{X} \mathbf{X}^\top \phi - 2\lambda \phi \stackrel{\text{set}}{=} \mathbf{0} \implies \mathbf{X} \mathbf{X}^\top \phi = \lambda \phi \implies \mathbf{A} \phi = \lambda \phi,$$

which is an eigenvalue problem for \mathbf{A} according to Eq. (66), where ϕ is the eigenvector and λ is the eigenvalue.

Optimization Forms of (Generalized) Eigenvalue Problem

- **Optimization Form 4:** Consider the following optimization problem with the variable $\Phi \in \mathbb{R}^{d \times d}$:

$$\begin{aligned} & \underset{\Phi}{\text{minimize}} && ||\mathbf{X} - \Phi \Phi^\top \mathbf{X}||_F^2, \\ & \text{subject to} && \Phi^\top \Phi = \mathbf{I}, \end{aligned} \tag{71}$$

where $\mathbf{X} \in \mathbb{R}^{d \times n}$. Similar to Eq. (70), the objective function in Eq. (71) is simplified as:

$$||\mathbf{X} - \Phi \Phi^\top \mathbf{X}||_F^2 = \text{tr}(\mathbf{X}^\top \mathbf{X} - \mathbf{X} \mathbf{X}^\top \Phi \Phi^\top)$$

The Lagrangian is:

$$\mathcal{L} = \text{tr}(\mathbf{X}^\top \mathbf{X}) - \text{tr}(\mathbf{X} \mathbf{X}^\top \Phi \Phi^\top) - \text{tr}(\Lambda^\top (\Phi^\top \Phi - \mathbf{I})),$$

$$\mathbb{R}^{d \times d} \ni \frac{\partial \mathcal{L}}{\partial \Phi} = 2 \mathbf{X} \mathbf{X}^\top \Phi - 2 \Phi \Lambda \stackrel{\text{set}}{=} \mathbf{0} \implies \mathbf{X} \mathbf{X}^\top \Phi = \Phi \Lambda \implies \mathbf{A} \Phi = \Phi \Lambda,$$

which is an eigenvalue problem for \mathbf{A} according to Eq. (67). The columns of Φ are the eigenvectors of \mathbf{A} and the diagonal elements of Λ are the eigenvalues.

Optimization Forms of (Generalized) Eigenvalue Problem

- **Optimization Form 5:** Consider the following optimization problem [5] with the variable $\phi \in \mathbb{R}^d$:

$$\underset{\phi}{\text{maximize}} \quad \frac{\phi^\top \mathbf{A} \phi}{\phi^\top \mathbf{B} \phi}. \quad (72)$$

According to the generalized Rayleigh-Ritz quotient method [6], this optimization problem can be restated as:

$$\begin{aligned} \underset{\phi}{\text{maximize}} \quad & \phi^\top \mathbf{A} \phi, \\ \text{subject to} \quad & \phi^\top \mathbf{B} \phi = 1, \end{aligned} \quad (73)$$

The Lagrangian is:

$$\mathcal{L} = \phi^\top \mathbf{A} \phi - \lambda(\phi^\top \mathbf{B} \phi - 1),$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}} = 2 \mathbf{A} \phi - 2 \lambda \mathbf{B} \phi \stackrel{\text{set}}{=} \mathbf{0} \implies 2 \mathbf{A} \phi = 2 \lambda \mathbf{B} \phi \implies \mathbf{A} \phi = \lambda \mathbf{B} \phi,$$

which is a generalized eigenvalue problem (\mathbf{A}, \mathbf{B}) according to Eq. (66).

Singular Value Decomposition

- **Singular Value Decomposition (SVD)** [10] 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 [10].
- SVD has two different forms, i.e., complete and incomplete.
- Consider a matrix $\mathbf{A} \in \mathbb{R}^{\alpha \times \beta}$. The **complete SVD** decomposes the matrix as:

$$\mathbb{R}^{\alpha \times \beta} \ni \mathbf{A} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T, \quad (74)$$
$$\mathbf{U} \in \mathbb{R}^{\alpha \times \alpha}, \quad \mathbf{V} \in \mathbb{R}^{\beta \times \beta}, \quad \mathbf{\Sigma} \in \mathbb{R}^{\alpha \times \beta},$$

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

- In complete SVD, $\mathbf{\Sigma}$ 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:

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

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

Singular Value Decomposition

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

$$\begin{aligned}\mathbb{R}^{\alpha \times \beta} \ni \mathbf{A} &= \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T, \\ \mathbf{U} \in \mathbb{R}^{\alpha \times k}, \quad \mathbf{V} \in \mathbb{R}^{\beta \times k}, \quad \mathbf{\Sigma} \in \mathbb{R}^{k \times k},\end{aligned}\tag{75}$$

where [11]:

$$k := \min(\alpha, \beta),\tag{76}$$

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

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

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

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, \mathbf{U} and \mathbf{V} are both orthogonal matrices so:

$$\mathbf{U}^\top \mathbf{U} = \mathbf{I}, \quad (77)$$

$$\mathbf{V}^\top \mathbf{V} = \mathbf{I}. \quad (78)$$

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

$$\mathbf{U}\mathbf{U}^\top = \mathbf{I}, \quad (79)$$

$$\mathbf{V}\mathbf{V}^\top = \mathbf{I}. \quad (80)$$

Relation of SVD and EVD

- In both complete and incomplete SVD of matrix \mathbf{A} , the left and right singular vectors are the eigenvectors of $\mathbf{A}\mathbf{A}^\top$ and $\mathbf{A}^\top\mathbf{A}$, respectively, and the singular values are the square root of eigenvalues of either $\mathbf{A}\mathbf{A}^\top$ or $\mathbf{A}^\top\mathbf{A}$.
- Proof: There is:

$$\mathbf{A}\mathbf{A}^\top = (\mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top)(\mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top)^\top = \mathbf{U}\mathbf{\Sigma}\underbrace{\mathbf{V}^\top\mathbf{V}}_{\mathbf{I}}\mathbf{\Sigma}\mathbf{U}^\top = \mathbf{U}\mathbf{\Sigma}\mathbf{\Sigma}\mathbf{U}^\top = \mathbf{U}\mathbf{\Sigma}^2\mathbf{U}^\top,$$

which is the eigen-decomposition [12] of $\mathbf{A}\mathbf{A}^\top$ where the columns of \mathbf{U} are the eigenvectors and the diagonal of $\mathbf{\Sigma}^2$ are the eigenvalues so the diagonal of $\mathbf{\Sigma}$ are the square root of eigenvalues.

- Also:

$$\mathbf{A}^\top\mathbf{A} = (\mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top)^\top(\mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top) = \mathbf{V}\mathbf{\Sigma}\underbrace{\mathbf{U}^\top\mathbf{U}}_{\mathbf{I}}\mathbf{\Sigma}\mathbf{V}^\top = \mathbf{V}\mathbf{\Sigma}\mathbf{\Sigma}\mathbf{V}^\top = \mathbf{V}\mathbf{\Sigma}^2\mathbf{V}^\top,$$

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

Acknowledgement

- Some slides are based on our textbook: “Elements of Dimensionality Reduction and Manifold Learning” [1]
- Some slides of this slide deck (expectation parts) are inspired by the lectures of Prof. Mu Zhu at University of Waterloo, Department of Statistics and Actuarial Science.
- For more information on optimization, refer to my course “Optimization Techniques” in University of Guelph. Link in my YouTube channel:
<https://www.youtube.com/playlist?list=PLPrxGIUWsQP3ZBM4Zy5YqfCh1BqM5sJov>
- More information on Rayleigh-Ritz quotient, eigenvalue problem, and SVD: see our tutorial paper “Eigenvalue and generalized eigenvalue problems: Tutorial” [12]
- More information on expectation: see our two tutorial papers “The Theory Behind Overfitting, Cross Validation, Regularization, Bagging, and Boosting: Tutorial” [13] and “Sampling algorithms, from survey sampling to Monte Carlo methods: Tutorial and literature review” [14]

References

- [1] B. Ghojogh, M. Crowley, F. Karay, 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] B. N. Parlett, “The symmetric eigenvalue problem,” *Classics in Applied Mathematics*, vol. 20, 1998.
- [6] E. Croot, “The Rayleigh principle for finding eigenvalues,” tech. rep., Georgia Institute of Technology, School of Mathematics, 2005.
Online, Accessed: March 2019.
- [7] I. Jolliffe, *Principal component analysis*. Springer, 2011.

References (cont.)

- [8] J. H. Wilkinson, *The algebraic eigenvalue problem*, vol. 662. Oxford Clarendon, 1965.
- [9] G. H. Golub and C. F. Van Loan, *Matrix computations*, vol. 3. The Johns Hopkins University Press, 2012.
- [10] G. W. Stewart, “On the early history of the singular value decomposition,” *SIAM review*, vol. 35, no. 4, pp. 551–566, 1993.
- [11] G. H. Golub and C. Reinsch, “Singular value decomposition and least squares solutions,” *Numerische mathematik*, vol. 14, no. 5, pp. 403–420, 1970.
- [12] B. Ghojogh, F. Karray, and M. Crowley, “Eigenvalue and generalized eigenvalue problems: Tutorial,” *arXiv preprint arXiv:1903.11240*, 2019.
- [13] B. Ghojogh and M. Crowley, “The theory behind overfitting, cross validation, regularization, bagging, and boosting: tutorial,” *arXiv preprint arXiv:1905.12787*, 2019.
- [14] B. Ghojogh, H. Nekoei, A. Ghojogh, F. Karray, and M. Crowley, “Sampling algorithms, from survey sampling to Monte Carlo methods: Tutorial and literature review,” *arXiv preprint arXiv:2011.00901*, 2020.