

Lecture 1: Logistics and Preliminaries

Instructor: Yifan Chen

Scribes: Yujia Yin

Proof reader: Yifan Chen, Xiong Peng

Note: LaTeX template courtesy of UC Berkeley EECS dept.

Disclaimer: These notes have not been subjected to the usual scrutiny reserved for formal publications. They may be distributed outside this class only with the permission of the Instructor.

1.1 Gaussian Distribution

The Gaussian distribution, also known as the normal distribution, is a continuous probability distribution that is symmetric around its mean. It is characterized by its mean (μ) and covariance matrix (Σ). In this section, we study the matrix form of the distribution for the random variable $\mathbf{X} \in \mathbb{R}^n$ as follows:

$$p(\mathbf{X}; \mu, \Sigma) = \frac{1}{(\sqrt{2\pi})^n} \cdot \frac{1}{|\Sigma|^{1/2}} \cdot \exp\left(-\frac{1}{2}(\mathbf{X} - \mu)^\top \Sigma^{-1}(\mathbf{X} - \mu)\right),$$

here, μ is the mean vector of size $n \times 1$, and Σ is the covariance matrix, which describes the shape and orientation of the distribution. The term $|\Sigma|$ represents the determinant of the covariance matrix.

1.1.1 T-test

In a t-test, we consider two independent groups of observations. Suppose each group consists of n observations, and let \mathbf{X}_n^y and \mathbf{X}_n^b be two independent random vectors, following a multivariate normal distribution:

$$\mathbf{X}_n^y \sim N(\mu_y, \sigma_y^2 \mathbf{I}), \quad \mathbf{X}_n^b \sim N(\mu_b, \sigma_b^2 \mathbf{I}),$$

where μ_y and μ_b are both mean vectors, \mathbf{I} is the identity matrix.

Under the null hypothesis $H_0 : \mu_y = \mu_b$ and alternative hypothesis $H_1 : \mu_y > \mu_b$, we define a new random variable $\mathbf{X}_n^y - \mathbf{X}_n^b$. Since under H_0 both \mathbf{X}_n^y and \mathbf{X}_n^b are normally distributed and independent, their difference also follows a normal distribution:

$$\mathbf{X}_n^y - \mathbf{X}_n^b \sim N(\mathbf{0}, (\sigma_y^2 + \sigma_b^2) \mathbf{I}).$$

1.1.2 t distribution

The basic form of the t distribution is:

$$T = \frac{z}{\sqrt{s/d}},$$

where the random variables z and s satisfy the following conditions:

1. $z \sim N(0, 1)$,
2. $s \sim \chi^2(d)$, where d is the degrees of freedom,
3. z and s are independent.

In the case of estimating a population mean, the t distribution arises in the following way:

$$T = \frac{\bar{X} - \mu}{\hat{\sigma}/\sqrt{n}} \sim T(n-1), \quad \text{where } \hat{\sigma}^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{X})^2. \quad (1.1)$$

Proof. We can rewrite eq. (1.1) to match the basic form of the t distribution:

$$T = \frac{(\bar{X} - \mu)/\sqrt{\sigma^2/n}}{\hat{\sigma}/\sqrt{n}/\sqrt{\sigma^2/n}} \triangleq \frac{z}{\sqrt{s/(n-1)}}. \quad (1.2)$$

1. We can easily verify that condition item 1 holds because the numerator in eq. (1.2) follows a standard normal distribution.
2. Next, we prove condition item 2 holds in our case. we need to show that: $s = \frac{n-1}{\sigma^2} \hat{\sigma}^2 = \frac{1}{\sigma^2} \sum_{i=1}^n (x_i - \bar{X})^2$ follows a chi-square distribution with $n - 1$ degrees of freedom.

Let $\mathbf{X} = \{x_1, x_2, \dots, x_n\}^\top$ be the vector of observed data points. Using vector notation, we have:

$$\begin{aligned} s &= \frac{n-1}{\sigma^2} \hat{\sigma}^2 = \frac{1}{\sigma^2} \sum_{i=1}^n (x_i - \bar{X})^2 \\ &= \frac{1}{\sigma^2} (\mathbf{X} - \frac{1}{n} \cdot \mathbf{1} \cdot \mathbf{1}^\top \mathbf{X})^\top (\mathbf{X} - \frac{1}{n} \cdot \mathbf{1} \cdot \mathbf{1}^\top \mathbf{X}) \\ &= \frac{1}{\sigma^2} \left[(\mathbf{I} - \frac{1}{n} \cdot \mathbf{1} \cdot \mathbf{1}^\top) \mathbf{X} \right]^\top \left[(\mathbf{I} - \frac{1}{n} \cdot \mathbf{1} \cdot \mathbf{1}^\top) \mathbf{X} \right]. \end{aligned} \quad (1.3)$$

Let $\mathbf{P} := \mathbf{I} - \frac{1}{n} \cdot \mathbf{1} \cdot \mathbf{1}^\top$, which is a projection matrix satisfying $\mathbf{P}^2 = \mathbf{P}$ and $\mathbf{P} \cdot \mathbf{1} = 0$:

$$\mathbf{P} = \bar{\mathbf{U}} \begin{pmatrix} 1 & 0 & \cdots & 0 & 0 \\ 0 & 1 & \cdots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \cdots & 1 & 0 \\ 0 & 0 & \cdots & 0 & 0 \end{pmatrix} \bar{\mathbf{U}}^\top = \mathbf{U} \mathbf{U}^\top, \quad (1.4)$$

where $\mathbf{U} \in \mathbb{R}^{n \times (n-1)}$ is orthogonal. From eqs. (1.3) and (1.4), we obtain:

$$s = \frac{1}{\sigma^2} \mathbf{X}^\top \mathbf{P}^\top \mathbf{P} \mathbf{X} = \frac{1}{\sigma^2} \mathbf{X}^\top \mathbf{U} (\mathbf{U}^\top \mathbf{X}).$$

Let $\mathbf{Y}_{n-1} := \mathbf{U}^\top \mathbf{X} \sim N(\mu \cdot \mathbf{U}^\top \cdot \mathbf{1}, \sigma^2 \cdot \mathbf{U}^\top \mathbf{U})$, where $\mathbf{U}^\top \mathbf{U} = \mathbf{I}_{n-1}$. We can also infer that $\mathbf{P} \cdot \mathbf{1} = 0$ based on the properties of the projection matrix, which implies that $\mathbf{U}^\top \cdot \mathbf{1} = 0$. Finally, we have:

$$s = \frac{1}{\sigma^2} \mathbf{X}^\top \mathbf{U} (\mathbf{U}^\top \mathbf{X}) = \frac{1}{\sigma^2} \mathbf{Y}_{n-1}^\top \mathbf{Y}_{n-1} = \frac{1}{\sigma^2} \sum_{i=1}^{n-1} Y_i^2 \sim \chi^2(n-1).$$

Thus, condition item 2 holds.

3. To prove the independence between z and s , we calculate the covariance between \bar{X} and \mathbf{Y} .

$$\begin{aligned} \text{Cov}(\mathbf{Y}, \bar{X}) &= E[\mathbf{Y} \cdot \bar{X}] - E[\mathbf{Y}] \cdot E[\bar{X}] \\ &= E[\mathbf{U}^\top \mathbf{X} \cdot \frac{1}{n} \cdot \mathbf{X}^\top \cdot \mathbf{1}] \quad (E[\mathbf{Y}] = 0) \\ &= \frac{1}{n} \mathbf{U}^\top (E[\mathbf{X} \mathbf{X}^\top]) \cdot \mathbf{1} \quad (\text{linearity of expectation}) \\ &= \frac{1}{n} \mathbf{U}^\top \left[(\mu \cdot \mathbf{1})(\mu \cdot \mathbf{1})^\top + \sigma^2 \mathbf{I} \right] \cdot \mathbf{1} \quad (\text{covariance structure of } \mathbf{X}) \\ &= \frac{1}{n} \mathbf{U}^\top \mu^2 \cdot \mathbf{1} \cdot \mathbf{1}^\top \cdot \mathbf{1} + \frac{1}{n} \mathbf{U}^\top \sigma^2 \cdot \mathbf{I} \cdot \mathbf{1} = \mathbf{0}_{n-1}. \quad (\mathbf{U}^\top \mathbf{1} = 0) \end{aligned} \quad (1.5)$$

The eq. (1.5) shows that condition item 3 holds.

□

The t distribution is commonly used in hypothesis testing and in the construction of t-tests. It allows for inference about population means when the population standard deviation is unknown.

1.1.3 Maximum Likelihood Estimator

The Maximum Likelihood Estimator (MLE) of μ is the value that maximizes the likelihood function, i.e.,

$$\hat{\boldsymbol{\mu}} = \arg \max_{\boldsymbol{\mu}} \prod_{i=1}^n p(\mathbf{x}_i; \boldsymbol{\mu}, \boldsymbol{\Sigma}).$$

Next, we simplify the likelihood function by applying a logarithm, which converts the product into a sum. This transformation facilitates computation and differentiation. For example, taking the logarithm of an exponential function often leads to a more tractable additive form. In the case of MLE, this allows us to rewrite the optimization problem as follows:

$$\sum \ln p(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) \propto \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\mu}).$$

Next, we use matrix trace to simplify the problem:

$$\begin{aligned} & \min_{\boldsymbol{\mu}} \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\mu}) \\ \iff & \min \text{Tr} \left[(\mathbf{X} - \mathbf{1} \cdot \boldsymbol{\mu}^\top)^\top \boldsymbol{\Sigma}^{-1} (\mathbf{X} - \mathbf{1} \cdot \boldsymbol{\mu}^\top) \right] \\ := & f(\boldsymbol{\mu}). \end{aligned}$$

After defining $f(\boldsymbol{\mu})$, we calculate its derivative:

$$df = \text{Tr} \left[\left(\frac{\partial f}{\partial \boldsymbol{\mu}} \right)^\top d\boldsymbol{\mu} \right].$$

To compute the derivative of $f(\boldsymbol{\mu})$, let's recall how to differentiate the trace of a matrix. For matrices \mathbf{ABC} , the trace of the product is denoted as $\text{Tr}(\mathbf{ABC})$, and its derivative is:

$$df = \text{Tr} [\mathbf{dA} \cdot \mathbf{BC} + \mathbf{A} \cdot \mathbf{dB} \cdot \mathbf{C} + \mathbf{AB} \cdot \mathbf{dC}]. \quad (1.6)$$

Additionally, based on properties of matrix traces, we know that:

$$\text{Tr}(\mathbf{AB}) = \text{Tr}(\mathbf{BA}). \quad (1.7)$$

Using eqs. (1.6) and (1.7), the derivative of the MLE is:

$$df = \text{Tr} \left[d(\mathbf{X} - \mathbf{1} \cdot \boldsymbol{\mu}^\top)^\top \boldsymbol{\Sigma}^{-1} (\mathbf{X} - \mathbf{1} \cdot \boldsymbol{\mu}^\top) + \mathbf{0} + (\mathbf{X} - \mathbf{1} \cdot \boldsymbol{\mu}^\top)^\top \boldsymbol{\Sigma}^{-1} d(\mathbf{X} - \mathbf{1} \cdot \boldsymbol{\mu}^\top) \right],$$

let $\mathbf{A}^\top = \boldsymbol{\Sigma}^{-1}(\mathbf{X} - \mathbf{1} \cdot \boldsymbol{\mu}^\top)$, so we have:

$$\begin{aligned} df &= \text{Tr} \left[d(\mathbf{X} - \mathbf{1} \cdot \boldsymbol{\mu}^\top)^\top \cdot \mathbf{A}^\top + \mathbf{0} + \mathbf{A} \cdot d(\mathbf{X} - \mathbf{1} \cdot \boldsymbol{\mu}^\top) \right] \\ &= \text{Tr} \left[-\mathbf{1} \cdot d\boldsymbol{\mu}^\top \cdot \mathbf{A}^\top + \mathbf{A} \cdot (-d\boldsymbol{\mu}) \cdot \mathbf{1}^\top \right] \\ &= \text{Tr} \left[-2 \cdot -\mathbf{1}^\top \cdot \mathbf{A} \cdot d\boldsymbol{\mu} \right]. \end{aligned}$$

Thus, we derive the first derivative of $f(\boldsymbol{\mu})$. Since $f(\boldsymbol{\mu})$ is a convex function, the point at which the first derivative is zero corresponds to the global minimum. Therefore, we deduce:

$$\begin{aligned}\frac{\partial f}{\partial \mu} &= -2 \cdot \mathbf{A} \cdot \mathbf{1} = 0 \\ &\Rightarrow \Sigma^{-1}(\mathbf{X} - \mathbf{1} \cdot \boldsymbol{\mu}^\top) \cdot \mathbf{1} = 0 \\ &\Rightarrow \boldsymbol{\mu} = \frac{1}{n} \cdot \mathbf{X}^\top \cdot \mathbf{1}.\end{aligned}$$

Thus, we obtain the value of $\boldsymbol{\mu}$.

1.2 Linear Regression

1.2.1 Linear Model

Assume we have a linear model described by the equation:

$$\mathbf{Y} = (\mathbf{1}_n \quad \mathbf{X}) \cdot \begin{pmatrix} \boldsymbol{\beta}_0 \\ \boldsymbol{\beta}_1 \end{pmatrix} + \mathbf{e}, \quad \mathbf{X} \in \mathbb{R}^{n \times d}, \quad \boldsymbol{\beta}_1 \in \mathbb{R}^{d \times 1}. \quad (1.8)$$

In the eq. (1.8), \mathbf{X} is the design matrix of independent variables, and $\boldsymbol{\beta}_0, \boldsymbol{\beta}_1$ are the regression coefficients. The error term \mathbf{e} typically assumed to satisfy the Gaussian-Markov conditions:

1. Zero mean: $E[\mathbf{e}] = \mathbf{0}$,
2. Homoscedasticity and no autocorrelation: $\text{Var}[\mathbf{e}] = \sigma^2 \mathbf{I}_{n \times n}$.

1.2.2 Square Loss

To estimate $\bar{\boldsymbol{\beta}}$, we minimize the squared loss function, which is given by:

$$L(\bar{\boldsymbol{\beta}}) = \frac{1}{2n} (\mathbf{Y} - \bar{\mathbf{X}} \bar{\boldsymbol{\beta}})^\top (\mathbf{Y} - \bar{\mathbf{X}} \bar{\boldsymbol{\beta}}), \quad \text{where } \bar{\mathbf{X}} = (\mathbf{1}_n, \mathbf{X}), \quad \bar{\boldsymbol{\beta}} = \begin{pmatrix} \boldsymbol{\beta}_0 \\ \boldsymbol{\beta}_1 \end{pmatrix}.$$

Since the squared loss function is convex, we find its global minimum by setting its gradient to zero. Taking the derivative with respect to $\bar{\boldsymbol{\beta}}$ and setting it to zero, we obtain:

$$\nabla_{\bar{\boldsymbol{\beta}}} L(\bar{\boldsymbol{\beta}}) = -\frac{1}{n} \bar{\mathbf{X}}^\top (\mathbf{Y} - \bar{\mathbf{X}} \bar{\boldsymbol{\beta}}) = \mathbf{0}.$$

Solving for $\bar{\boldsymbol{\beta}}$, we derive the closed-form solution:

$$\bar{\mathbf{X}}^\top \bar{\mathbf{X}} \bar{\boldsymbol{\beta}} = \bar{\mathbf{X}}^\top \mathbf{Y}. \quad (1.9)$$

Since eq. (1.9) uniquely determines the optimal estimate of $\bar{\boldsymbol{\beta}}$, we denote the solution as $\hat{\bar{\boldsymbol{\beta}}}$:

$$\hat{\bar{\boldsymbol{\beta}}} = (\bar{\mathbf{X}}^\top \bar{\mathbf{X}})^{-1} \bar{\mathbf{X}}^\top \mathbf{Y}.$$

Next, we analyze the properties of $\hat{\bar{\boldsymbol{\beta}}}$, focusing on its unbiasedness and variance.

Lemma 1.1. *The estimator $\hat{\bar{\boldsymbol{\beta}}}$ is unbiased.*

Proof. For $\mathbf{Y} = \bar{\mathbf{X}} \bar{\boldsymbol{\beta}} + \mathbf{e}$, we can deduce:

$$\begin{aligned}E[\hat{\bar{\boldsymbol{\beta}}}] &= E[(\bar{\mathbf{X}}^\top \bar{\mathbf{X}})^{-1} \bar{\mathbf{X}}^\top (\bar{\mathbf{X}} \bar{\boldsymbol{\beta}} + \mathbf{e})] \\ &= (\bar{\mathbf{X}}^\top \bar{\mathbf{X}})^{-1} (\bar{\mathbf{X}}^\top \bar{\mathbf{X}}) \bar{\boldsymbol{\beta}} + (\bar{\mathbf{X}}^\top \bar{\mathbf{X}})^{-1} (\bar{\mathbf{X}}^\top \bar{\mathbf{X}}) E[\mathbf{e}].\end{aligned}$$

Since $E[\mathbf{e}] = \mathbf{0}$, we have:

$$E[\hat{\bar{\boldsymbol{\beta}}}] = \bar{\boldsymbol{\beta}}.$$

□

Lemma 1.2. The variance of $\hat{\beta}$ is given by $\text{Var}(\hat{\beta}) = \sigma^2(\bar{X}^\top \bar{X})^{-1}$.

Proof. Using the formula for the variance of a linear transformation, we deduce:

$$\begin{aligned}\text{Var}(\hat{\beta}) &= (\bar{X}^\top \bar{X})^{-1} \bar{X}^\top \cdot \text{Var}(Y) \cdot \bar{X} (\bar{X}^\top \bar{X})^{-1} \\ &= \sigma^2 (\bar{X}^\top \bar{X})^{-1} \bar{X}^\top \bar{X} (\bar{X}^\top \bar{X})^{-1} \\ &= \sigma^2 (\bar{X}^\top \bar{X})^{-1}.\end{aligned}$$

□

While we obtained $\hat{\beta}$ by minimizing the squared loss, we now show that this solution also corresponds to the Maximum Likelihood Estimator (MLE) under the assumption of Gaussian errors.

Lemma 1.3. $\hat{\beta}$ is the MLE of β under the Gaussian error term e .

Proof. From the linear model in eq. (1.8), we can easily get that Y is following a normal distribution. The likelihood function of Y given \bar{X} is:

$$L(\bar{\beta}, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{n/2}} \exp\left(-\frac{1}{2\sigma^2}(\mathbf{Y} - \bar{X}\bar{\beta})^\top(\mathbf{Y} - \bar{X}\bar{\beta})\right).$$

Taking the logarithm, we have:

$$\log L(\bar{\beta}, \sigma^2) = -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2}(\mathbf{Y} - \bar{X}\bar{\beta})^\top(\mathbf{Y} - \bar{X}\bar{\beta}).$$

To find MLE of $\bar{\beta}$, we maximize $\log L(\bar{\beta}, \sigma^2)$. Since the first term is independent of $\bar{\beta}$, we only need to minimize $(\mathbf{Y} - \bar{X}\bar{\beta})^\top(\mathbf{Y} - \bar{X}\bar{\beta})$. We take its derivative with respect to $\bar{\beta}$ and set it to $\mathbf{0}$:

$$\begin{aligned}\nabla_{\bar{\beta}}(\mathbf{Y} - \bar{X}\bar{\beta})^\top(\mathbf{Y} - \bar{X}\bar{\beta}) &= \mathbf{0} \\ \Rightarrow -2\bar{X}^\top(\mathbf{Y} - \bar{X}\hat{\beta}) &= \mathbf{0}\end{aligned}$$

Thus, we obtain the solution:

$$\hat{\beta} = (\bar{X}^\top \bar{X})^{-1} \bar{X}^\top \mathbf{Y}.$$

□

1.2.3 Population Risk

In this section, we investigate the population risk of our linear model. We remark that during the training process, we are minimizing the *empirical risk* on the fixed design matrix \bar{X} .

Specifically, we assume a new random sample \mathbf{x} and the corresponding label $y = \mathbf{x}^\top \beta + \epsilon$; the risk is expressed as below:

$$E_{\mathbf{x}, y, \hat{\beta}} (y - \mathbf{x}^\top \hat{\beta})^2 = E_{\hat{\beta}} \left[E_{\mathbf{x}, y} (y - \mathbf{x}^\top \hat{\beta})^2 \mid \hat{\beta} \right].$$

Considering that the noise term ϵ is a random variable independent of the feature vector \mathbf{x} , and that its expectation is zero, the population risk depends on: $E_{\mathbf{x}, \hat{\beta}} [\mathbf{x}^\top (\beta - \hat{\beta})(\beta - \hat{\beta})^\top \mathbf{x}]$. For simplicity, we refer to this as the prediction error in the well-specified case, where the model is correctly specified. We are equivalently minimizing:

$$\begin{aligned}\min E_{\mathbf{x}, \hat{\beta}} &\left[\mathbf{x}^\top (\beta - \hat{\beta})(\beta - \hat{\beta})^\top \mathbf{x} \right] \\ &= E_{\hat{\beta}} \left[E_{\mathbf{x}} \text{Tr}(\mathbf{x} \mathbf{x}^\top (\beta - \hat{\beta})(\beta - \hat{\beta})^\top \mid \hat{\beta}) \right] \\ &= E_{\hat{\beta}} \left[\text{Tr}(E[\mathbf{x} \mathbf{x}^\top](\beta - \hat{\beta})(\beta - \hat{\beta})^\top) \mid \hat{\beta} \right] \\ &= \text{Tr}(E[\mathbf{x} \mathbf{x}^\top] \cdot E[(\beta - \hat{\beta})(\beta - \hat{\beta})^\top]) \\ &= \sigma^2 \text{Tr}(E[\mathbf{x} \mathbf{x}^\top] (\bar{X}^\top \bar{X})^{-1}).\end{aligned}\tag{1.10}$$

In eq. (1.10), we recall that $\bar{\mathbf{X}}$ is the fixed design matrix in the training process. In other words, during the training, we can manipulate $\bar{\mathbf{X}}$ to minimize the population risk, while other quantities remain fixed. The expectation $E[\mathbf{x}\mathbf{x}^\top]$ represents the second-moment matrix of the feature distribution and remains constant across realizations of \mathbf{x} , whereas individual samples \mathbf{x} are random. Assuming $E[\mathbf{x}\mathbf{x}^\top] = \mathbf{I}$, the preceding expression simplifies to: $\sigma^2 \text{Tr}[(\bar{\mathbf{X}}^\top \bar{\mathbf{X}})^{-1}] = \text{Tr}[\text{Var}(\hat{\boldsymbol{\beta}})]$, indicating the connection between population risk and parameter variance in the case of linear regression.