

STAT-4320 Notes

November 3, 2025

Contents

1 Class 1	3
1.1 The sample mean	3
1.2 Central Limit Theorem	3
2 Chapter 2	4
2.1 Breakdown Point and Efficiency	4
2.1.1 Sample and population median	4
2.1.2 Breakdown point	4
2.1.3 Efficiency	4
3 Class 3	6
3.1 Convergence of random variables	6
3.2 Slutsky's Lemma and Continuous Mapping Theorem	6
3.3 Weak Law of Large Numbers and Central Limit Theorem	6
3.4 Delta method	7
3.5 Multivariate Data	8
3.6 Multivariate Normal	9
4 Class 4	10
4.1 Linear Algebra Review	10
4.1.1 Positive definite matrices	10
4.2 Moment generating functions	11
4.3 Properties of the multivariate normal	11
4.4 Sample Variance	12
5 Class 5	13
5.1 Chi-squared distribution	13
5.2 Sample Variance, cont'd	13
5.3 Joint distribution of sample mean and sample variance	15
6 Class 6	16
6.1 Basic Framework of Statistical Estimation	16
6.2 Properties of Estimators	16
7 Class 7	20
7.1 Consistency	20
7.2 Method of moments estimation	20
8 Class 8	24
8.1 Maximum Likelihood Estimation	24
9 Class 9	27
9.1 Fisher Information	27
9.2 Asymptotic properties of the MLE estimator	27
9.3 Asymptotic properties of the MLE estimator for multiple parameters	30
10 Class 10	31
10.1 Cramer-Rao Lower Bound	31
10.2 Simple Linear Regression	32
10.2.1 Estimating β_0, β_1	32
10.2.2 Estimating σ^2	32
10.2.3 Facts about Least Square Estimates	33
11 Class 11	34
11.1 Least Absolute Deviation Line	34
12 Class 12	36
12.1 Confidence Intervals	36
12.1.1 Exact Confidence Intervals	36
12.1.2 Asymptotic Confidence Intervals	37

1 Class 1

1.1 The sample mean

Definition 1.1. (Sample mean): Given X_1, \dots, X_n iid from F , and $\mathbb{E}[X_i] = \mu$, the sample mean is the random variable defined as

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$$

Remark. The sample mean is random, so it has an expectation

$$\begin{aligned}\mathbb{E}[\bar{X}] &= \mathbb{E}\left[\frac{1}{n} \sum_{i=1}^n X_i\right] \\ &= \frac{1}{n} \sum_{i=1}^n \mathbb{E}[X_i] \quad \text{by linearity} \\ &= \mu\end{aligned}$$

The expectation of the sample mean is the population mean.

The sample mean is an unbiased estimator of the population mean.

Remark.

$$\begin{aligned}Var(\bar{X}) &= Var\left(\frac{1}{n} \sum_{i=1}^n X_i\right) \\ &= \frac{1}{n^2} \sum_{i=1}^n Var(X_i) \quad \text{since } \mathbb{E}[cX_i] = c\mathbb{E}[X_i], Var(cX_i) = c^2 Var(X_i), X_i \text{ indepedent} \\ &= \frac{\sigma^2}{n}\end{aligned}$$

Remark. We say that \bar{X} follows a sampling distribution. Think of this as a thought experiment. If a sample of size n is taken many times, we expect to see the sample mean exhibit the above expectation and variability.

1.2 Central Limit Theorem

Theorem 1.2. Theorem (Central Limit Theorem): Let X_1, \dots, X_n be iid from arbitrary distribution F wit mean μ and variance σ^2 , then as $n \rightarrow \infty$,

$$\frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \approx N(0, 1),$$

OR

$$\bar{X} \approx N\left(\mu, \frac{\sigma^2}{n}\right)$$

2 Chapter 2

2.1 Breakdown Point and Efficiency

Recall that sample mean is defined as

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$$

An issue with the sample mean is that corruption in 1 or few data points can make sample mean unstable. The sample median is more robust alternative.

2.1.1 Sample and population median

Definition 2.1. (Sample median): The sample median is the *middle value* when a list of numbers are sorted in non-decreasing order.

$$X_{med} = \begin{cases} X_{(\frac{n+1}{2})} & \text{if } n \text{ odd} \\ X_{(\frac{n}{2})} + X_{(\frac{n}{2}+1)} & \text{if } n \text{ even} \end{cases}$$

Where $X_{(i)}$ denotes the i -th smallest value in X_1, \dots, X_n

Definition 2.2. (population median): The population median of distribution F with density function f is the point m such that

$$\int_{\infty}^m f(x)dx = \int_m^{\infty} f(x)dx = \frac{1}{2}$$

2.1.2 Breakdown point

Definition 2.3. (Breakdown Point): The breakdown point of an estimate $\hat{\theta}_n$ based on data X_1, \dots, X_n is the fraction of data points that have to be moved to infinity for the estimate to also move to infinity.

Example. The breakdown point

- For sample mean $= \frac{1}{n}$
- For sample median $\approx \frac{1}{2}$

Remark. Note that this is **not** a direct consequence of CLT.

For example, $F = N(\mu, \sigma^2)$.

The sample mean follows **exactly** a normal distribution

$$\bar{X} \sim N\left(\mu, \frac{\sigma^2}{n}\right)$$

The sample median approximately follows

$$X_{med} \approx N\left(\mu, \frac{1}{4f(\mu)^2 n}\right)$$

Recall that

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2} \frac{(x-\mu)^2}{\sigma^2}}$$

Hence

$$f(\mu) = \frac{1}{\sqrt{2\pi}\sigma}$$

$$X_{med} \approx N\left(\mu, \frac{\pi\sigma^2}{2n}\right)$$

2.1.3 Efficiency

Definition 2.4. (Efficiency): The efficiency of two estimates is the ratio of their variances.

$$\text{Efficiency} \left(\tilde{X}_{med}, \bar{X} \right) = \frac{Var(\bar{X})}{Var(\tilde{X}_{med})}$$

Example. For sample mean and sample median, the efficnecy is $\frac{2}{\pi}$.

3 Class 3

3.1 Convergence of random variables

There are two kinds of convergence

- convergence in probability
- convergence in distribution

Definition 3.1. (Convergence in probability): We say a sequence of random variables $\{X_n\}_{n \geq 1}$ converges in probability to X if for any $\epsilon > 0$,

$$\lim_{n \rightarrow \infty} \mathbb{P}(|X_n - X| > \epsilon) = 0$$

Denoted

$$X_n \xrightarrow{p} X$$

Definition 3.2. (Convergence in distribution): We say a sequence of random variables $\{X_n\}_{n \geq 1}$ converges in distribution to X if

$$\mathbb{P}(X_n \leq x) \xrightarrow{n \rightarrow \infty} \mathbb{P}(X \leq x)$$

Which is equivalent to

$$F_n(x) \xrightarrow{n \rightarrow \infty} F(x)$$

Denoted

$$X_n \xrightarrow{d} X$$

3.2 Slutsky's Lemma and Continuous Mapping Theorem

Lemma 3.3. Lemma (Slutsky's Lemma): If $X_n \xrightarrow{d} X$ and $Y_n \xrightarrow{p} c$ for constant c , then the following hold

- $X_n + Y_n \xrightarrow{d} X + c$
- $X_n Y_n \xrightarrow{d} cX$
- $\frac{X_n}{Y_n} \xrightarrow{d} \frac{X}{c}$ if $c > 0$

Theorem 3.4. (Continuous mapping): If g is a continuous function, then

$$X_n \xrightarrow{p} X \implies g(X_n) \xrightarrow{p} g(X)$$

$$X_n \xrightarrow{d} X \implies g(X_n) \xrightarrow{d} g(X)$$

3.3 Weak Law of Large Numbers and Central Limit Theorem

The weak law of large numbers is an example of convergence in probability. The central limit theorem is an example of convergence in distribution.

Theorem 3.5. Weak law of large numbers: Suppose $X_1 \dots X_n$ iid from F with $\mathbb{E}[X_1] = \mu$ and $Var(X) = \sigma^2$, then

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i \xrightarrow{p} \mu$$

Proof.

$$\begin{aligned} \mathbb{P}(|\bar{X} - \mu| > \epsilon) &= \frac{Var(\bar{X})}{\epsilon^2} \quad \text{by Chebyshev's Inequality} \\ &= \frac{\sigma^2}{n\epsilon^2} \\ &\xrightarrow{n \rightarrow \infty} 0 \end{aligned}$$

✓

Theorem 3.6. Markov's Inequality: For any random variable X , and non-negative constant a ,

$$P(|X| \geq a) \leq \frac{\mathbb{E}[|X|]}{a}$$

Alternatively, for any non-negative random variable X ,

$$P(X \geq a) \leq \frac{\mathbb{E}[X]}{a}$$

Proof. We prove the general case, for any random variable X , let $Y = |X|$

$$\begin{aligned} \mathbb{E}[Y] &= \mathbb{E}[Y|Y \geq a] P(Y \geq a) + \mathbb{E}[Y|Y < a] P(Y < a) \text{ by Law of Total Expectation} \\ &\geq \mathbb{E}[Y|Y \geq a] P(Y \geq a) \\ &\geq a P(Y \geq a) \\ \implies P(Y \geq a) &\leq \frac{\mathbb{E}[Y]}{a} \end{aligned}$$

✓

Theorem 3.7. Chebyshev's Inequality: Let X be a random variable with mean μ and variance σ^2 , then for any $a > 0$

$$P(|X - \mu| \geq a) \leq \frac{\sigma^2}{a^2}$$

Proof.

$$\begin{aligned} P(|X - \mu| \geq a) &= P((X - \mu)^2 \geq a^2) \\ &\leq \frac{\mathbb{E}[(X - \mu)^2]}{a^2} \quad \text{by Markov's Inequality} \\ &= \frac{\text{Var}(X)}{a^2} \end{aligned}$$

✓

Theorem 3.8. Central limit theorem: Suppose $X_1 \dots X_n$ iid from F with $\mathbb{E}[X_1] = \mu$ and $\text{Var}(X) = \sigma^2$, then

$$\frac{\sqrt{n}(\bar{X} - \mu)}{\sigma} \xrightarrow{d} N(0, 1)$$

Alternatively, let

$$Z_n = \frac{\sqrt{n}(\bar{X} - \mu)}{\sigma}$$

$$\mathbb{P}(Z_n \leq t) = P(N(0, 1) \leq t) \text{ for all } t \in \mathbb{R}$$

Remark.

$$\bar{X}_n \approx N\left(\mu, \frac{\sigma^2}{n}\right)$$

Although this statement has no mathematical content. We are taking the limit of the distribution of \bar{X} as n gets large. $N\left(\mu, \frac{\sigma^2}{n}\right)$ cannot be a limit.

3.4 Delta method

CLT gives asymptotic distribution of \bar{X} . We want to get the asymptotic distribution of functions of \bar{X} . By Continuous Mapping Theorem, we get this for free

$$\frac{\sqrt{n}(\bar{X} - \mu)}{\sigma} \xrightarrow{d} N(0, 1) \implies g\left(\frac{\sqrt{n}(\bar{X} - \mu)}{\sigma}\right) \xrightarrow{d} g(N(0, 1))$$

However, we don't just want statements about $g(Z_n)$, we want statements about $g(\bar{X})$

Theorem 3.9. Delta Method: Suppose X_1, \dots, X_n iid F , with $\mathbb{E}[X_1] = \mu$ and $\text{Var}(X) = \sigma^2$ and g is a function

such that the derivative of $g'(\mu) \neq 0$. Then

$$\sqrt{n}(g(\bar{X}) - g(\mu)) \xrightarrow{d} N(0, \sigma^2(g'(\mu))^2)$$

Remark. Note: We know by Continuous Mapping Theorem that the in-probability limit of $g(\bar{X})$ is $g(\mu)$

$$g(\bar{X}) \xrightarrow{p} g(\mu)$$

Subtracting away the in-probability limit and taking the Z-score, delta method tells us that the z-score follows a normal distribution.

Proof.

Recall Taylor's Expansion

$$f(x) = f(a) + (x - a)f'(a) + \frac{1}{2}(x - a)^2 f''(a) + \dots$$

$$\begin{aligned} g(\bar{X}) - g(\mu) &= (\bar{X} - \mu)g'(\mu) + \text{error terms} \\ \sqrt{n}(g(\bar{X}) - g(\mu)) &= \sqrt{n}(\bar{X} - \mu)g'(\mu) + \text{error terms} \end{aligned}$$

By CLT, we know that

$$\sqrt{n}(\bar{X} - \mu) \xrightarrow{d} N(0, \sigma^2)$$

By Slutsky's,

$$\sqrt{n}(\bar{X} - \mu)g'(\mu) \xrightarrow{d} N(0, \sigma^2)g'(\mu) = N(0, \sigma^2(g'(\mu))^2)$$

Note: if $g'(\mu) = 0$, by taking higher orders in the Taylor expansion, we get

$$g(\bar{X}) - g(\mu) \approx \frac{1}{2}(\bar{X} - \mu)^2 g''(\mu)$$

Since

$$n(\bar{X} - \mu)^2 = (\sqrt{n}(\bar{X} - \mu))^2 \xrightarrow{d} (\sigma N(0, 1))^2 = \sigma^2 \chi_1^2$$

We get

$$n(g(\bar{X}) - g(\mu)) \xrightarrow{d} \frac{1}{2}\sigma^2 \chi_1^2 g''(\mu)$$

✓

3.5 Multivariate Data

For each unit of study, the number of measurements is greater than 1. For example, for n data points and p observed variables

$$\mathbf{X}_1 = \begin{pmatrix} X_{11} \\ X_{12} \\ \vdots \\ X_{1p} \end{pmatrix}, \mathbf{X}_2 = \begin{pmatrix} X_{21} \\ X_{22} \\ \vdots \\ X_{2p} \end{pmatrix}, \mathbf{X}_3 = \begin{pmatrix} X_{31} \\ X_{32} \\ \vdots \\ X_{3p} \end{pmatrix}$$

Where \mathbf{X}_i are iid p-dimentional observations with distribution F .

The mean vector is

$$\boldsymbol{\mu} = \mathbb{E}[\mathbf{X}_1] = \begin{pmatrix} \mathbb{E}[X_{11}] \\ \mathbb{E}[X_{12}] \\ \vdots \\ \mathbb{E}[X_{1p}] \end{pmatrix} \in \mathbb{R}^p$$

The covariance matrix is denoted

$$\boldsymbol{\Sigma} \in \mathbb{R}^{p \times p}$$

Where the $i - j$ -th element is

$$\boldsymbol{\Sigma}_{i,j} = Cov(X_{1i}, X_{1j}) = \mathbb{E}[X_{1i}X_{1j}] - \mathbb{E}[X_{1i}]\mathbb{E}[X_{1j}]$$

Hence, we can express $\boldsymbol{\Sigma}$ as the difference of two matrices

$$\boldsymbol{\Sigma} = \mathbb{E}[\mathbf{X}_1\mathbf{X}_1^T] - \mathbb{E}[\mathbf{X}]\mathbb{E}[\mathbf{X}]^T$$

3.6 Multivariate Normal

Definition 3.10. (Multivariate Normal): A p -dimensional random vector $\mathbf{X} = \begin{pmatrix} X_1 \\ X_2 \\ \vdots \\ X_p \end{pmatrix}$ is said to follow the multivariate normal distribution with mean vector $\boldsymbol{\mu}$ and positive definite (pd) covariance matrix $\boldsymbol{\Sigma}$ if it has a density function $f : \mathbb{R}^p \rightarrow \mathbb{R}$ of the form

$$f(\mathbf{x}) = \frac{1}{\sqrt{(2\pi)^p \det(\boldsymbol{\Sigma})}} e^{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}-\boldsymbol{\mu})}$$

Remark. This definition only works with a pd covariance matrix.

Remark. A multivariate normal can be defined without the pd matrix, using the linear combination definition.

Definition 3.11. (Multivariate CLT): Suppose $\mathbf{X}_1 \dots \mathbf{X}_n$ are iid p -dimensional random vectors with mean vector $\boldsymbol{\mu} \in \mathbb{R}^p$ and covariance matrix $\boldsymbol{\Sigma}$

The sample mean

$$\bar{\mathbf{X}} = \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i$$

has the following distribution

$$\sqrt{n} (\bar{\mathbf{X}} - \boldsymbol{\mu}) \xrightarrow{d} N_p(\mathbf{0}, \boldsymbol{\Sigma})$$

4 Class 4

4.1 Linear Algebra Review

4.1.1 Positive definite matrices

Definition 4.1. Definition (positive definite): A symmetric $p \times p$ matrix is said to be positive definite (pd) if for all $\mathbf{x} \in \mathbb{R}^p \setminus \{\mathbf{0}\}$,

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

Remark. All eigenvalues of a pd matrix are positive

Remark. By spectral decomposition, any pd matrix can be written as

$$A = P \Lambda P^T,$$

Where Λ is a diagonal matrix with the eigenvalues of A on the diagonals

$$\Lambda = \begin{bmatrix} \lambda_1 & \dots & \dots & 0 \\ 0 & \lambda_2 & \dots & 0 \\ \vdots & & & \\ 0 & \dots & \dots & \lambda_p \end{bmatrix}$$

and P is an orthogonal matrix

$$P^T P = P P^T = I_{p \times p}$$

Remark. A^{-1} exists and is given by

$$A^{-1} = P \Lambda^{-1} P^T$$

Proof.

$$A^{-1} A = P \Lambda^{-1} P^T P \Lambda P^T = P \Lambda^{-1} \Lambda P^T = P P^T = I$$

✓

Remark. A has a square root. Given a pd matrix A , we say that B is the square root of A if $BB = A$

$$B = P \begin{bmatrix} \sqrt{\lambda_1} & \dots & \dots & 0 \\ 0 & \sqrt{\lambda_2} & \dots & 0 \\ \vdots & & & \\ 0 & \dots & \dots & \sqrt{\lambda_p} \end{bmatrix} P^T$$

Remark. Sum of eigenvalues is the trace of A

$$\sum_{i=1}^n \lambda_i = \text{tr}(A)$$

Remark. What does it mean to assume that the covariance matrix is pd ?

Consider any $\mathbf{a} \in \mathbb{R}^p \setminus \{\mathbf{0}\}$. Recall that

$$\begin{aligned} \Sigma &= \mathbb{E} [\mathbf{X} \mathbf{X}^T] - \mathbb{E}[\mathbf{X}] \mathbb{E}[\mathbf{X}]^T \\ &= \mathbb{E} [(\mathbf{X} - \mathbb{E}[\mathbf{X}])(\mathbf{X} - \mathbb{E}[\mathbf{X}])^T] \end{aligned}$$

Consider $\mathbf{a}^T \Sigma \mathbf{a}$,

$$\begin{aligned} \mathbf{a}^T \Sigma \mathbf{a} &= \mathbb{E} [(\mathbf{a}^T (\mathbf{X} - \mathbb{E}[\mathbf{X}]))) (\mathbf{a}^T (\mathbf{X} - \mathbb{E}[\mathbf{X}]))^T] \text{ since } (AB)^T = B^T A^T \\ &= \mathbb{E} [(\mathbf{a}^T (\mathbf{X} - \mathbb{E}[\mathbf{X}]))^2] \\ &= \text{Var}(\mathbf{a}^T \mathbf{X}) \\ &> 0 \end{aligned}$$

i.e. projected onto any direction \mathbf{a} , the variance of \mathbf{X} is nonzero.

i.e. The RV is non-degenerate along every direction $\mathbf{a} \in \mathbb{R}^p \setminus \{\mathbf{0}\}$.

4.2 Moment generating functions

Definition 4.2. (Moment generating function) For a random variable X , the *mgf* is a function $\mathbb{R} \rightarrow \mathbb{R}_{\geq 0}$,

$$\phi_X(t) = \mathbb{E}[e^{tX}]$$

Example. For $X \sim N(\mu, \sigma^2)$

$$\phi_X(t) = \exp\left(t\mu + \frac{1}{2}\sigma^2 t^2\right)$$

If X is a p-dimensional random variable, the *mgf* is a function $\mathbb{R}^p \rightarrow \mathbb{R}_{\geq 0}$

$$\phi_X(\mathbf{t}) = \exp\left(\mathbf{t}^T \boldsymbol{\mu} + \frac{1}{2} \mathbf{t}^T \boldsymbol{\Sigma} \mathbf{t}\right)$$

4.3 Properties of the multivariate normal

Proposition 4.3. If \mathbf{X} is a p-dimensional normal, $\mathbf{X} \sim N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, and \mathbf{A} is a $k \times p$ matrix such that $\text{rank}(\mathbf{A}) = k \leq p$ (i.e. full row rank) and $\mathbf{b} \in \mathbb{R}^k$ is a fixed vector,

$$\mathbf{AX} + \mathbf{b} \sim N_k(\mathbf{A}\boldsymbol{\mu}, \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^T)$$

Example. An important case of this is when $\mathbf{b} = \mathbf{0}$ and $k = 1$, then \mathbf{A} is a row vector $\mathbf{A} = [a_1, a_2, \dots, a_p]$. Take any $\mathbf{a} \in \mathbb{R}^p$, then

$$\mathbf{a}^T \mathbf{X} \sim N_1(\mathbf{a}^T \boldsymbol{\mu}, \mathbf{a}^T \boldsymbol{\Sigma} \mathbf{a})$$

This can also be expressed as

$$\mathbf{a}^T \mathbf{X} = \sum_{i=1}^p a_i X_i \sim N_1\left(\sum_{i=1}^p a_i \mu_i, \mathbf{a}^T \boldsymbol{\Sigma} \mathbf{a}\right)$$

Proof. *Exercise for enthusiasts:* Prove 1 using *mgfs*. ✓

Proposition 4.4. Suppose \mathbf{X} is a p_1 -dimensional RV and \mathbf{Y} is a p_2 -dimensional RV, such that

$$\begin{pmatrix} \mathbf{X} \\ \mathbf{Y} \end{pmatrix} \in \mathbb{R}^{p_1+p_2}, \quad \begin{pmatrix} \mathbf{X} \\ \mathbf{Y} \end{pmatrix} \sim N_{p_1+p_2}\left(\begin{pmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{pmatrix}, \begin{pmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{pmatrix}\right)$$

where $\boldsymbol{\mu}_1 \in \mathbb{R}^{p_1}$, $\boldsymbol{\mu}_2 \in \mathbb{R}^{p_2}$, and

$$\boldsymbol{\Sigma} = \begin{pmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{pmatrix} \in \mathbb{R}^{(p_1+p_2) \times (p_1+p_2)}$$

then we say that

$$\mathbf{X} \perp \mathbf{Y} \Leftrightarrow \boldsymbol{\Sigma}_{11} = \boldsymbol{\Sigma}_{21}^T = \mathbf{0}$$

Proof. The forward direction is obvious.

The converse is not true in general for a 1-dimensional normal, but it is true in multivariate normal. ✓

Example. Important case: $p_1 = p_2 = 1$, Suppose

$$\begin{pmatrix} X \\ Y \end{pmatrix} \sim N\left((\mu_1, \mu_2), \begin{bmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{12} & \sigma_2^2 \end{bmatrix}\right)$$

Then $X \perp Y$ if and only if $\sigma_{12} = 0$.

Remark. Relationship

- If $X \perp Y$, then $\text{Cov}(X, Y) = 0$. This is always true

$$X \perp Y \implies \text{Cov}(X, Y) = 0$$

- The converse is not true in general
- However, if (X, Y) are jointly normal or jointly bernoulli, then

$$\text{Cov}(X, Y) = 0 \implies X \perp Y$$

4.4 Sample Variance

The sample variance

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$$

Has a sampling distribution

$$s^2 \sim \frac{\sigma^2}{n-1} \chi_{n-1}^2$$

5 Class 5

5.1 Chi-squared distribution

Definition 5.1. The χ_d^2 is the sum of d iid standard normals squared.

$$Z_1, \dots, Z_d \sim N(0, 1), \quad \sum_{i=1}^d Z_i \sim \chi_d^2$$

Result 5.2. The *chi-squared* distribution has expectation

$$\begin{aligned}\mathbb{E}[\chi_d^2] &= d\mathbb{E}[Z_1^2] \\ &= d\end{aligned}$$

Proof.

$$\mathbb{E}[Z_1^2] = Var(Z_1) - \mathbb{E}[Z_1]^2$$

✓

Result 5.3. The *chi-squared* distribution has variance

$$\begin{aligned}Var(\chi_d^2) &= Var\left(\sum_{i=1}^d Z_i^2\right) \\ &= \sum_{i=1}^d Var(Z_i^2) \\ &= 2d\end{aligned}$$

Proof.

$$\begin{aligned}Var(Z_1^2) &= \mathbb{E}[Z_1^4] - \mathbb{E}[Z_1^2]^2 \\ &= \mathbb{E}[Z_1^4] - 1 \\ &= 3 - 1 \\ &= 2\end{aligned}$$

✓

5.2 Sample Variance, cont'd

Result 5.4. The sample variance follows a *chi-squared* distribution.

$$s^2 \sim \frac{\sigma^2}{n-1} \chi_{n-1}^2$$

Proof. Let

$$\mathbf{Y} = \begin{pmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{pmatrix} \begin{pmatrix} X_1 - \bar{X} \\ X_2 - \bar{X} \\ \vdots \\ X_n - \bar{X} \end{pmatrix}$$

We express $\mathbf{Y} = \mathbf{AX}$ for some matrix \mathbf{A} .

$$\mathbf{Y} = \begin{pmatrix} 1 - \frac{1}{n} & -\frac{1}{n} & \cdots & -\frac{1}{n} \\ -\frac{1}{n} & 1 - \frac{1}{n} & \cdots & -\frac{1}{n} \\ \vdots & \vdots & \ddots & \vdots \\ -\frac{1}{n} & -\frac{1}{n} & \cdots & 1 - \frac{1}{n} \end{pmatrix} \mathbf{X}$$

Since

$$\mathbf{X} \sim N_n(\mu \mathbf{1}, \sigma^2 \mathbf{I})$$

Hence

$$\begin{aligned}
\mathbf{Y} &\sim N_n(\mathbf{A}\boldsymbol{\mu}, \sigma^2 \mathbf{A}\mathbf{A}^T) \\
&= N_n(\mu \mathbf{A}\mathbf{1}, \sigma^2 \mathbf{A}\mathbf{A}^T) \\
&= N_n(\mathbf{0}, \sigma^2 \mathbf{A}\mathbf{A}^T) \text{ since } \mathbf{A}\mathbf{1} = \mathbf{0} \\
&= N_n(\mathbf{0}, \sigma^2 \mathbf{A}^2) \text{ since } \mathbf{A} \text{ symmetric}
\end{aligned}$$

Note that

$$\mathbf{A} = \mathbf{I} - \frac{1}{n} \mathbf{1}\mathbf{1}^T$$

Where $\mathbf{1}\mathbf{1}^T$ is a rank 1 symmetric matrix. Hence $\mathbf{1}\mathbf{1}^T$ has at most 1 non-zero eigenvalue. Since the sum of eigenvalues is the trace, and $\text{tr}(\mathbf{1}\mathbf{1}^T) = n$, the non-zero eigenvalue is n . $\mathbf{1}\mathbf{1}^T$ has eigenvalues $(n, 0, 0, \dots, 0)$. $\frac{1}{n} \mathbf{1}\mathbf{1}^T$ has eigenvalues $(1, 0, 0, \dots)$.

Therefore, $\mathbf{I} - \frac{1}{n} \mathbf{1}\mathbf{1}^T$ has eigenvalues $(1, 1, \dots, 1, 0)$ (This only works because of \mathbf{I}).

We can express s^2 in terms of \mathbf{Y} .

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2 = \frac{1}{n-1} \mathbf{Y}^T \mathbf{Y}$$

It suffices to show that

$$\mathbf{Y}^T \mathbf{Y} \sim \sigma^2 \chi_{n-1}^2$$

Note: \mathbf{Y} has elements that are normal, i.e. $Y_i = X_i - \bar{X}$ is a difference of normals. However, Y_i is not independent due to \bar{X} .

Fact: Denote $\Sigma = \mathbf{A}^2$, and we define $\mathbf{Z} = N_n(\mathbf{0}, \mathbf{I})$. Then,

$$Y \stackrel{d}{=} \sigma \Sigma^{\frac{1}{2}} \mathbf{Z} \sim N_n(\mathbf{0}, \sigma^2 \mathbf{A})$$

This is true because

$$\begin{aligned}
\sigma \Sigma^T \Sigma &\sim N_n \left(\mathbf{0}, \sigma^2 \left(\Sigma^{\frac{1}{2}} \right)^T \Sigma^{\frac{1}{2}} \right) \\
&= N_n \left(\mathbf{0}, \sigma^2 \mathbf{A}^T \mathbf{A} \right) \\
&= N_n \left(\mathbf{0}, \sigma^2 \mathbf{A}^2 \right)
\end{aligned}$$

Hence

$$\begin{aligned}
\mathbf{Y}^T \mathbf{Y} &= \sigma^2 \mathbf{Z}^T \left(\Sigma^{\frac{1}{2}} \right)^T \Sigma^{\frac{1}{2}} \mathbf{Z} \\
&= \sigma^2 \mathbf{Z} \mathbf{A}^T \mathbf{A} \mathbf{Z} \\
&= \sigma^2 \mathbf{Z}^T \mathbf{A} \mathbf{Z}
\end{aligned}$$

Hence it suffices to show that

$$\mathbf{Z}^T \mathbf{A} \mathbf{Z} \sim \chi_{n-1}^2$$

Note that if $\mathbf{A} = \mathbf{I}$, then $\mathbf{Z}^T \mathbf{Z} \sim \chi_n^2$.

By spectral decomposition,

$$\mathbf{Z} \mathbf{A}^T \mathbf{Z} = \mathbf{Z} \mathbf{P} \boldsymbol{\Lambda} \mathbf{P}^T \mathbf{Z}$$

Denote $\mathbf{P}^T \mathbf{Z} = \mathbf{W} \in \mathbb{R}^n$.

$$\mathbf{Z} \mathbf{A}^T \mathbf{Z} = \mathbf{Z} \mathbf{P} \boldsymbol{\Lambda} \mathbf{P}^T \mathbf{Z} = \mathbf{W}^T \boldsymbol{\Lambda} \mathbf{W}$$

Because of spectral decomposition, we know that \mathbf{P} and \mathbf{P}^T are orthogonal matrices whose product is the identity. Applying an orthogonal matrix to a multivariate standard normal, i.g. $\mathbf{A}\mathbf{Z}$, does not change the multivariate standard normal distribution.

$$\mathbf{W} \sim N_n(\mathbf{0}, \mathbf{P} \mathbf{P}^T) = N_n(\mathbf{0}, \mathbf{I})$$

Hence

$$\begin{aligned}
\mathbf{W}^T \boldsymbol{\Lambda} \mathbf{W} &= \mathbf{W}^T \begin{pmatrix} 1 & 0 & \dots & 0 & 0 \\ 0 & 1 & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & 0 \\ 0 & 0 & \dots & 1 & 0 \\ 0 & 0 & \dots & 0 & 0 \end{pmatrix} \mathbf{W} \\
&= \sum_{i=1}^{n-1} W_i^2 \\
&\sim \chi_{n-1}^2 \quad \checkmark
\end{aligned}$$

5.3 Joint distribution of sample mean and sample variance

We know the marginal distributions for the sample mean and variance

$$\bar{X} \sim N\left(\mu, \frac{\sigma^2}{n}\right), s^2 \sim \frac{\sigma^2}{n-1} \chi_{n-1}^2$$

Result 5.5. If X_1, X_2, \dots, X_n iid normal, \bar{X}, s^2 independent.

Proof. Define

$$\begin{aligned} \mathbf{Z} &= \begin{pmatrix} X_1 - \bar{X} \\ X_2 - \bar{X} \\ \vdots \\ X_n - \bar{X} \\ \bar{X} \end{pmatrix} \\ \mathbf{Z} &= \mathbf{B}\mathbf{X}, \mathbf{B} \in \mathbb{R}^{(n+1) \times n}, \mathbf{B} = \begin{pmatrix} & & & & \mathbf{A} \\ & & & & \\ & & & & \\ \frac{1}{n} & \frac{1}{n} & \dots & \frac{1}{n} & \end{pmatrix} \end{aligned}$$

Hence \mathbf{Z} is $(n+1)$ -dimensional normal.

If we show that $X_1 - \bar{X} \perp \bar{X}, X_2 - \bar{X} \perp \bar{X} \dots X_n - \bar{X} \perp \bar{X}$, then

$$\begin{pmatrix} X_1 - \bar{X} \\ X_2 - \bar{X} \\ \vdots \\ X_n - \bar{X} \end{pmatrix} \perp \bar{X} \implies s^2 \perp \bar{X}$$

Since $\mathbf{Z} = \mathbf{B}\mathbf{X}$ is multivariate normal, it suffices to check that the covariance is 0.

To show $Cov(\bar{X}, X_i - \bar{X}) = 0$ for all $1 \leq i \leq n$.

Note that

$$\begin{aligned} \bar{X} &= \frac{1}{n} \sum_{i=1}^n X_i = \sum_{i=1}^n a_i X_i, \quad a_1 = a_2 = \dots = a_n = \frac{1}{n} \\ X_1 - \bar{X} &= \sum_{i=1}^n b_i X_i, \quad b_1 = 1 - \frac{1}{n}, b_2 = b_3 = \dots = b_n = -\frac{1}{n} \end{aligned}$$

$$\begin{aligned} Cov\left(\sum_{i=1}^n a_i X_i, \sum_{i=1}^n b_i X_i\right) &= \sum_{i=1}^n \sum_{j=1}^n a_i b_j Cov(X_i X_j) \\ &= \sum_{i=1}^n a_i b_i Var(X_i) \\ &= \sigma^2 \sum_{i=1}^n a_i b_i \\ &= 0 \end{aligned}$$

Remark. As a general strategy, to show independence, we can show in two steps

1. show jointly normal
2. show 0 covariance

6 Class 6

6.1 Basic Framework of Statistical Estimation

Given iid samples X_1, X_2, \dots, X_n , how do we infer / estimate parameters of F ?

Example. (Bernoulli): Given X_1, X_2, \dots, X_n iid $Ber(p)$,

$$X_i = \begin{cases} 1 & \text{with probability } p \\ 0 & \text{with probability } 1 - p \end{cases}$$

An estimate of p is

$$\hat{p} = \frac{1}{n} \sum_{i=1}^n X_i = \frac{T}{n}$$

Where T is the number of heads and $T \sim Binom(n, p)$.

Example. (Normal): X_1, X_2, \dots, X_n iid $N(\mu, \sigma^2)$. Estimates for parameters are

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n X_i = \bar{X}$$

$$\hat{\sigma}^2 = s^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$$

6.2 Properties of Estimators

Definition 6.1. (Unbiasedness): An estimate $\hat{\theta}$ is said to be unbiased for a parameter θ if

$$\mathbb{E}[\hat{\theta}] = \theta$$

for all values of the population.

Example. (Bernoulli):

$$\mathbb{E}[\hat{p}] = \frac{\mathbb{E}[T]}{n} = \frac{np}{n} = p$$

Hence \hat{p} is an unbiased estimate of p .

Note that \hat{p}^2 is not an unbiased estimate of p^2 .

$$\begin{aligned} \mathbb{E}[\hat{p}^2] &= \mathbb{E}[T^2] \left(\frac{1}{n^2} \right) \\ &= (np(1-p) + n^2 p^2) \left(\frac{1}{n^2} \right) \\ &= p^2 + \frac{p(1-p)}{n} \\ &\neq p^2 \end{aligned}$$

Note that we can rearrange terms to get

$$\begin{aligned} \mathbb{E}[T^2] &= n^2 p^2 + np(1-p) \\ &= (n^2 - n)p^2 + np \\ \frac{\mathbb{E}[T^2]}{n(n-1)} &= p^2 + \frac{p}{n-1} \end{aligned}$$

We try estimating $\frac{p}{n-1}$ by $\frac{T}{(n-1)n}$.

Consider the estimate

$$\tilde{p} = \frac{T^2}{n(n-1)} - \frac{T}{n(n-1)} = \frac{T(T-1)}{n(n-1)}$$

The expectation is

$$\mathbb{E}[\tilde{p}] = p^2 + \frac{p}{n-1} - \frac{p}{n-1} = p^2$$

Note that (Proof left as exercise)

$$\mathbb{E}\left[\frac{T(T-1)}{n(n-1)}\right] = \sum_{r=2}^n \frac{r(r-1)}{n(n-1)} \binom{n}{r} p^r (1-p)^{n-r} = p^2$$

In general, an unbiased estimate for p^k is

$$\frac{T(T-1)(T-2)\dots(T-k+1)}{n(n-1)(n-2)\dots(n-k+1)}$$

An unbiased estimate of $2p^2 + 5p^3$ is

$$2\frac{T(T-1)}{n(n-1)} + 5\frac{T(T-1)(T-2)}{n(n-1)(n-2)}$$

Example. (Normal): Estimate σ^2 with

$$\hat{\sigma}^2 = s^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2 \sim \chi_{n-1}^2 \left(\frac{\sigma^2}{n-1} \right)$$

The expectation of sample variance is

$$\begin{aligned} \mathbb{E}[s^2] &= \frac{\sigma^2}{n-1} \mathbb{E}[\chi_{n-1}^2] \\ &= \frac{\sigma^2}{n-1} (n-1) \\ &= \sigma^2 \end{aligned}$$

Example. (Network sampling)

A *network* refers to a graph of vertices that are connected if they have an interaction.

Network sampling helps with understanding *how networks look* by studying a small section of the network, and understanding *features of a large unobserved network* from a sample subgraph.

Motifs refer to patterns of small subgraphs, such as an edge, or a triangle.

Motif estimation refers to estimating the number of a particular type of motif, based on an observed sample subgraph.

Let G_n be a population graph on n -vertices. Our subgraph sampling model involves sampling each vertex of G_n with probability $p \in (0, 1)$ independently, and then observing the subgraph on the set of sampled vertices.

Goal: estimate the number of edges in G_n based on the observed graph.

- initial guess: count the number of edges in the observed graph
- Denote $\hat{E}(G_n)$ as the number of edges in observed graph
- Denote $E(G_n)$ as number of edges in population graph

Result:

$$\mathbb{E}[\hat{E}(G_n)] = p^2$$

Hence,

$$\frac{\hat{E}(G_n)}{p^2} = \frac{\# \text{ edges in observed graph}}{p^2}$$

is an unbiased estimate of the number of edges in the population.

proof: Note that

$$E(G_n) = \sum_{1 \leq i \leq j \leq n} a_{ij}$$

Where

$$a_{ij} = \begin{cases} 1 & \text{if } (i, j) \text{ edge in } G_n \\ 0 & \text{otherwise} \end{cases}$$

Denote S the set of sampled vertices.

$$\begin{aligned}
\hat{E}(G_n) &= \sum_{1 \leq i \leq j \leq n, i \in S, j \in S} a_{ij} \\
&= \sum_{1 \leq i \leq j \leq n} a_{ij} \mathbf{1}[i \in S] \mathbf{1}[j \in S] \\
\mathbb{E}[\hat{E}(G_n)] &= a_{ij} \sum_{1 \leq i \leq j \leq n} \mathbb{E}[\mathbf{1}[i \in S] \mathbf{1}[j \in S]] \\
&= \sum_{1 \leq i \leq j \leq n} a_{ij} \mathbb{P}(i \in S) \mathbb{P}(j \in S) \\
&= p^2 \sum_{1 \leq i \leq j \leq n} a_{ij} \\
&= p^2 E(G_n)
\end{aligned}$$

Definition 6.2. (Variance of an estimate): The variance of an estimate $\hat{\theta}$ is

$$Var(\hat{\theta}) = \mathbb{E} \left[(\hat{\theta} - \mathbb{E}[\hat{\theta}])^2 \right]$$

Example. (Bernoulli):

$$Var(\hat{p}) = Var \left(\frac{Binom(n, p)}{n} \right) = \frac{np(1-p)}{n^2} = \frac{p(1-p)}{n}$$

Example. (Normal):

$$\begin{aligned}
Var(\hat{\mu}) &= Var(\bar{X}) = \frac{\sigma^2}{n} \\
Var(\hat{\sigma}^2) &= Var(s^2) \\
&= Var \left(\frac{\sigma^2}{n-1} \chi_{n-1}^2 \right) \\
&= \left(\frac{\sigma^2}{n-1} \right)^2 Var(\chi_{n-1}^2) \\
&= \frac{\sigma^4}{(n-1)^2} 2(n-1) \\
&= \frac{2\sigma^4}{n-1}
\end{aligned}$$

Definition 6.3. Definition (Mean Squared Error): The mean squared error of an estimate $\hat{\theta}$ is

$$MSE(\hat{\theta}) = \mathbb{E} \left[(\hat{\theta} - \theta)^2 \right]$$

Result 6.4.

$$MSE(\hat{\theta}) = (Bias(\hat{\theta}))^2 + Var(\hat{\theta})$$

Proof.

Proof:

$$\begin{aligned} MSE(\hat{\theta}) &= \mathbb{E} \left[(\hat{\theta} - \theta)^2 \right] \\ &= \mathbb{E} \left[(\hat{\theta} - \mathbb{E}[\hat{\theta}] + \mathbb{E}[\hat{\theta}] - \theta)^2 \right] \\ &= \mathbb{E} \left[(\hat{\theta} - \mathbb{E}[\hat{\theta}])^2 + 2(\hat{\theta} - \mathbb{E}[\hat{\theta}])(\mathbb{E}[\hat{\theta}] - \theta) + (\mathbb{E}[\hat{\theta}] - \theta)^2 \right] \\ &= \mathbb{E} \left[(\hat{\theta} - \mathbb{E}[\hat{\theta}])^2 \right] + 2(\mathbb{E}[\hat{\theta}] - \theta) \mathbb{E} [\hat{\theta} - \mathbb{E}[\hat{\theta}]] + (\mathbb{E}[\hat{\theta}] - \theta)^2 \\ &= \text{Var}(\hat{\theta}) + 0 + (\mathbb{E}[\hat{\theta}] - \theta)^2 \\ &= \text{Var}(\hat{\theta}) + (\text{Bias}(\hat{\theta}, \theta))^2 \end{aligned}$$

✓

Corollary 6.5. If $\hat{\theta}$ unbiased, then $MSE(\hat{\theta}) = Var(\hat{\theta})$

7 Class 7

7.1 Consistency

Definition 7.1. (Consistency): Suppose $\hat{\theta}_n$ is an estimate of θ based on n iid samples. Then, $\hat{\theta}_n$ is said to be consistent for θ if

$$\hat{\theta}_n \xrightarrow{P} \theta$$

as $n \rightarrow \infty$

Result 7.2. MSE converging to 0 in the limit implies consistency

$$MSE(\hat{\theta}_n) \rightarrow 0 \implies \hat{\theta}_n \xrightarrow{P} \theta$$

Proof. By Markov's Inequality, for any $\epsilon > 0$

$$P(|\hat{\theta}_n - \theta| > \epsilon) \leq \frac{\mathbb{E}[(\hat{\theta}_n - \theta)^2]}{\epsilon^2}$$

✓

Remark. To prove consistency, we can

- Show MSE vanishes
- Use Continuous Mapping
- Use Slutsky's Lemma

Example. (Bernoulli):

$$Bias(\hat{p}) = 0, Var(\hat{p}) = \frac{p(1-p)}{n}$$

As $n \rightarrow \infty$

$$MSE(\hat{p}) = Var(\hat{p}) = \frac{p(1-p)}{n} \rightarrow 0$$

Example. (Normal):

$$\hat{\mu} = \bar{X} \implies bias(\hat{\mu}) = 0, Var(\hat{\mu}) = Var(\bar{X}) = \frac{\sigma^2}{n}$$

As $n \rightarrow \infty$

$$MSE(\bar{X}) = \frac{\sigma^2}{n} \rightarrow 0$$

$$\hat{\sigma}^2 = s^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$$

$$bias(\hat{\sigma}^2) = 0, Var(\hat{\sigma}^2) = Var\left(\sigma^2 \frac{\chi_{n-1}^2}{n-1}\right) = \frac{2\sigma^4}{n-1}$$

As $n \rightarrow \infty$

$$MSE(\hat{\sigma}^2) = Var(\hat{\sigma}^2) \rightarrow 0$$

7.2 Method of moments estimation

Definition 7.3. (Method of Moments): Suppose X_1, X_2, \dots, X_n iid from distribution F which has k unknown parameters $(\theta_1, \theta_2, \dots, \theta_k)$.

Denote $\boldsymbol{\theta} = (\theta_1, \theta_2, \dots, \theta_k)$.

For each $1 \leq j \leq k$, compute the j -th moment of F

$$\alpha_j(\boldsymbol{\theta}) = \mathbb{E}[X^j] = \int X_j f(x) dx$$

Define the j -th sample moment as

$$\hat{\alpha}_j = \frac{1}{n} \sum_{i=1}^n X_i^j$$

Set up k equations

$$\begin{cases} \alpha_1 = \hat{\alpha}_1 \\ \alpha_2 = \hat{\alpha}_2 \\ \vdots \\ \alpha_k = \hat{\alpha}_k \end{cases}$$

The method of moments estimate $\hat{\theta}_n$ is the solution to the system of equations.

Remark. Note that

- The above system is a system of k equations in k unknowns
- The above system may have 0, 1 or multiple solutions. If solution is unique, then the solution is the method of moments estimate of $\hat{\theta}$

Example. (Bernoulli): $X_1, \dots, X_n \sim Ber(p)$

$$\begin{aligned} \alpha_1 &= \alpha_1(p) = \mathbb{E}[X] = p \\ \hat{\alpha}_1 &= \frac{1}{n} \sum_{i=1}^n X_i \end{aligned}$$

Solving for p :

$$\begin{aligned} \hat{\alpha}_1 &= \alpha_1 \\ \implies \hat{p} &= \frac{1}{n} \sum_{i=1}^n X_i \end{aligned}$$

Example. (Normal) $X_1, \dots, X_n \sim N(\mu, \sigma^2)$

The population moments are

$$\begin{aligned} \alpha_1 &= \mu \\ \alpha_2 &= \mathbb{E}[X^2] = \sigma^2 + \mu^2 \end{aligned}$$

Equating to sample moments

$$\begin{aligned} \mu &= \alpha_1 = \hat{\alpha}_1 = \frac{1}{n} \sum_{i=1}^n X_i \\ \mu^2 + \sigma^2 &= \alpha_2 = \hat{\alpha}_2 = \frac{1}{n} \sum_{i=1}^n X_i^2 \end{aligned}$$

Therefore,

$$\begin{aligned} \hat{\mu} &= \sum_{i=1}^n X_i \\ \hat{\sigma}^2 &= \frac{1}{n} \sum_{i=1}^n X_i^2 - \left(\frac{1}{n} \sum_{i=1}^n X_i \right)^2 \\ &= \frac{1}{n} \sum_{i=1}^n X_i^2 - \bar{X}^2 \\ &= \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2 \end{aligned}$$

The last equality follows because

$$\begin{aligned}
& \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2 \\
&= \frac{1}{n} \sum_{i=1}^n (X_i^2 - 2X_i \bar{X} + \bar{X}^2) \\
&= \frac{1}{n} \sum_{i=1}^n X_i^2 - \frac{2}{n} \bar{X} \sum_{i=1}^n X_i + \bar{X}^2 \\
&= \frac{1}{n} \sum_{i=1}^n X_i^2 - \bar{X}^2
\end{aligned}$$

Note: $\hat{\sigma}_{MLE}^2$ is a biased estimator of σ^2 .

$$\mathbb{E} [\hat{\sigma}^2] = \frac{n-1}{n} \mathbb{E} [s^2] = \frac{n-1}{n} \sigma^2$$

As $n \rightarrow \infty$

$$\begin{aligned}
Bias(\hat{\theta}^2) &= \frac{n-1}{n} \sigma^2 - \sigma^2 \rightarrow 0 \\
Var(\hat{\theta}^2) &= Var \left(\frac{n-1}{n} s^2 \right) \\
&= \left(\frac{n-1}{n} \right)^2 \left(\frac{2\sigma^4}{n-1} \right) \\
&= \frac{n-1}{n^2} 2\sigma^4 \\
&\rightarrow 0
\end{aligned}$$

Hence $MSE(\hat{\sigma}) \rightarrow 0$, and $\hat{\sigma}_{MLE}^2$ is consistent despite being biased.

Example. Ex (Pearson, 1984).

Contributions to the mathematical theory of evolution (1894). Weldon (1893) crab data: body length. The histogram of the data did not look like the bell-shaped Gaussian curve.

Conjecture: maybe two species, each Gaussian? (Gaussian Mixture Model.)

Pearson conjectured that the data consists of two different kinds of crabs, each with its own Gaussian distribution.

Gaussian Mixture Model (GMM)

Let W be a random variable such that

$$W \sim \begin{cases} N(\mu_1, \sigma_1^2), & \text{with prob. } p, \\ N(\mu_2, \sigma_2^2), & \text{with prob. } (1-p). \end{cases}$$

Distribution Function

$$P(W \leq t) = P(W \leq t, Z = 1) + P(W \leq t, Z = 2),$$

where $Z = \begin{cases} 1 & \text{if group 1 is sampled,} \\ 2 & \text{if group 2 is sampled.} \end{cases}$

So,

$$P(W \leq t) = p P(N(\mu_1, \sigma_1^2) \leq t) + (1-p) P(N(\mu_2, \sigma_2^2) \leq t).$$

Note: GMM is *not* adding two normals. Adding two normals gets another normal, but GMM is not a sum. It is a mixture.

Density Function

$$f(t) = \frac{d}{dt} P(W \leq t) = p \phi \left(\frac{t - \mu_1}{\sigma_1} \right) + (1-p) \phi \left(\frac{t - \mu_2}{\sigma_2} \right),$$

where $\phi(t) = \frac{1}{\sqrt{2\pi}} \exp \left(-\frac{t^2}{2} \right)$.

Estimation of Parameters of GMM

1. Compute sample moments:

$$\hat{m}_r = \frac{1}{n} \sum_{i=1}^n X_i^r, \quad 1 \leq r \leq 5.$$

2. Compute population moments:

$$m_r = \mathbb{E}[X^r], \quad 1 \leq r \leq 5.$$

3. Solve for $\mu_1, \mu_2, \sigma_1^2, \sigma_2^2, p$ from

$$\hat{m}_r = m_r, \quad r = 1, \dots, 5.$$

Pearson's Sixth-Moment Test

In general, the system of equations can have multiple roots. To choose a root, Pearson's approach: choose the root $(\hat{\mu}_1, \hat{\mu}_2, \hat{\sigma}_1^2, \hat{\sigma}_2^2, \hat{p})$ that is closest to the sixth sample moment.

That is, the feasible root with the smallest value of

$$|\hat{m}_6 - m_6| = \left| \frac{1}{n} \sum_{i=1}^n X_i^6 - \mathbb{E}[X^6] \right|.$$

Theorem (Kalai, Moitra, Valiant, 2010): the sixth-moment method gives consistent estimates of the parameters of GMM and can be computed efficiently.

8 Class 8

8.1 Maximum Likelihood Estimation

Definition 8.1. (Maximum Likelihood Estimation): Let X_1, \dots, X_n iid from F with pdf or pmf f_θ . Consider the joint pdf or pmf

$$L(\theta|x_1, x_2, \dots, x_n) = \prod_{i=1}^n f_\theta(x_i)$$

The maximum likelihood estimator is given by

$$\hat{\theta}_n = \arg \max_{\Theta} L(\theta|x_1, x_2, \dots, x_n)$$

Equivalently, define

$$\begin{aligned} l(\theta|x_1, \dots, x_n) &:= \log L(\theta|x_1, \dots, x_n) \\ &= \sum_{i=1}^n \log f_\theta(x_i) \end{aligned}$$

$$\hat{\theta}_n = \arg \max_{\Theta} l(\theta|x_1, x_2, \dots, x_n)$$

Example. (Bernoulli) $X_1, \dots, X_n \sim Ber(p)$.

$$\begin{aligned} f_p(x) &= P(X = x) \\ &= \begin{cases} p & \text{if } x = 1 \\ 1 - p & \text{if } x = 0 \end{cases} \\ &= p^x (1 - p)^{1-x} \end{aligned}$$

The likelihood function is

$$\begin{aligned} L(p|X_1, \dots, X_n) &= \prod_{i=1}^n f_p(X_i) \\ &= \prod_{i=1}^n p^{X_i} (1 - p)^{1-X_i} \\ &= p^{\sum_{i=1}^n X_i} (1 - p)^{n - \sum_{i=1}^n X_i} \\ l(p|X_1, X_2, \dots, X_n) &= \left(\sum_{i=1}^n X_i \right) \log p + \left(n - \sum_{i=1}^n X_i \right) \log(1 - p) \end{aligned}$$

Define $T = \sum_{i=1}^n X_i$,

$$\begin{aligned} l(p|X_1, \dots, X_n) &= T \log p + (n - T) \log(1 - p) \\ \frac{d}{dp} l(p|X_1, \dots, X_n) &= \frac{T}{p} - \frac{n - T}{1 - p} \\ &= 0 \end{aligned}$$

Hence

$$\hat{p}_{MLE} = \frac{T}{n} = \frac{1}{n} \sum_{i=1}^n X_i$$

To show that this is indeed a maximum,

$$\frac{d^2}{dp^2} l(p|X_1, \dots, X_n) = -\frac{T}{p^2} - \frac{T}{(1-p)^2} < 0$$

Hence $l(p|X_1, \dots, X_n)$ is concave and $\hat{p} = \frac{T}{n}$ is the unique maximum.

By LLN,

$$\hat{p} \xrightarrow{p} p$$

By CLT,

$$\sqrt{n}(\hat{p} - p) \xrightarrow{d} N(0, p(1-p))$$

Example. (Normal): X_1, \dots, X_n iid $N(\mu, \sigma^2)$

$$\begin{aligned} L(\mu, \sigma^2 | X_1, \dots, X_n) &= \prod_{i=1}^n f_{\mu, \sigma^2}(X_i) \\ &= \prod_{i=1}^n \left(\frac{1}{\sqrt{2\pi\sigma^2}} \right) \exp \left(-\frac{1}{2\sigma^2}(X_i - \mu)^2 \right) \\ &= \left(\frac{1}{\sqrt{2\pi\sigma^2}} \right)^n \exp \left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (X_i - \mu)^2 \right) \end{aligned}$$

The log likelihood is

$$\begin{aligned} l(\mu, \sigma^2 | X_1, \dots, X_n) &= \log(L(\mu, \sigma^2 | X_1, \dots, X_n)) \\ &= n \log \frac{1}{\sqrt{2\pi\sigma^2}} - \frac{1}{2\sigma^2} \sum_{i=1}^n (X_i - \mu)^2 \\ &= -\frac{n}{2} \log 2\pi\sigma^2 - \frac{n}{2} \log \sigma^2 - \frac{1}{2\sigma^2} \sum_{i=1}^n (X_i - \mu)^2 \end{aligned}$$

Define

$$g(\mu, \sigma^2) := -\frac{n}{2} \log \sigma^2 - \frac{1}{2\sigma^2} \sum_{i=1}^n (X_i - \mu)^2$$

Setting first order partials to zero, we get

$$\begin{aligned} \frac{\partial g(\mu, \sigma^2)}{\partial \mu} &= \frac{1}{\sigma^2} \sum_{i=1}^n (X_i - \mu) = 0 \\ \frac{\partial g(\mu, \sigma^2)}{\partial \sigma^2} &= -\frac{n}{2} \left(\frac{1}{\sigma^2} \right) + \frac{1}{2\sigma^4} \sum_{i=1}^n (X_i - \mu)^2 = 0 \end{aligned}$$

Hence

$$\begin{aligned} \hat{\mu}_{MLE} &= \frac{1}{n} \sum_{i=1}^n X_i = \bar{X} \\ \hat{\sigma}_{MLE}^2 &= \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2 \end{aligned}$$

To show that this is the global max, we can either

- Find the Hessian and show that it is negative definite, i.e. the function is concave
- Argue that the derivative is 0 at only one point, and show that for extreme values of μ and σ^2 (i.e. $\mu \rightarrow -\infty$, $\mu \rightarrow \infty$, $\sigma^2 \rightarrow 0$, $\sigma^2 \rightarrow \infty$)

$$l(\mu, \sigma^2) \rightarrow -\infty$$

Note

1. $\hat{\mu} \rightarrow \mu$ by LLN
2. $\sqrt{n}(\hat{\mu} - \mu) \xrightarrow{d} N(0, \sigma^2)$
3. $\hat{\sigma}^2 \xrightarrow{p} \sigma^2$ by computing MSE.
4. $\sqrt{n}(\hat{\sigma}^2 - \sigma^2) \xrightarrow{d} N(0, 2\sigma^4)$

Proof of 4:

$$\begin{aligned} \sqrt{n}(\hat{\sigma}^2 - \sigma^2) &= \sqrt{n} \left(\frac{\sigma^2}{n-1} \chi_{n-1}^2 - \sigma^2 \right) \\ &= \sigma^2 \sqrt{n} \left(\frac{\chi_{n-1}^2}{n-1} - 1 \right) \end{aligned}$$

Note that

$$\sqrt{n} \left(\frac{\chi_n^2}{n} - 1 \right) \xrightarrow{d} N(0, 2)$$

Hence

$$\sqrt{n} (\hat{\sigma}^2 - \sigma^2) \xrightarrow{d} N(0, 2\sigma^4)$$

REVISIT: check with the professor regarding this proof on asymptotic normality of $\hat{\sigma}_{MLE}^2$ without using Fisher Information.

9 Class 9

9.1 Fisher Information

Definition 9.1. (Score Function): The score function is defined as

$$s(\theta) = \frac{d}{d\theta} \log f_\theta(X)$$

Definition 9.2. Definition (Fisher Information): The Fisher Information is defined as

$$I(\theta) = \mathbb{E}[(s(\theta))^2] = \mathbb{E}\left[\left(\frac{d}{d\theta} \log f_\theta(X)\right)^2\right]$$

Lemma 9.3. Lemma (Properties of the score function):

1. Zero expectation

$$\mathbb{E}[s(\theta)] = \mathbb{E}\left[\frac{d}{d\theta} \log f_\theta(X)\right] = 0$$

2. The variance is Fisher Information

$$\text{Var}\left(\frac{d}{d\theta} \log f_\theta(x)\right) = I(\theta)$$

3. The expectation of the second derivative is negative Fisher information

$$I(\theta) = -\mathbb{E}\left[\frac{d^2}{d\theta^2} \log f_\theta(X)\right]$$

Proof. (1):

$$\begin{aligned}\mathbb{E}\left[\frac{d}{d\theta} \log f_\theta(X)\right] &= \int \frac{d}{d\theta} \log f_\theta(x) f_\theta(x) dx \\ &= \int \frac{\frac{d}{d\theta} f_\theta(x)}{f_\theta(x)} f_\theta(x) dx \\ &= \int \frac{d}{d\theta} f_\theta(x) dx \\ &= \frac{d}{d\theta} \int f_\theta(x) dx \\ &= \frac{d}{d\theta} 1 \\ &= 0\end{aligned}$$

(2): Proof is immediate

$$\begin{aligned}\text{Var}(s(\theta)) &= \mathbb{E}[(s(\theta))^2] - \mathbb{E}[s(\theta)]^2 \\ &= \mathbb{E}[(s(\theta))^2] \\ &= I(\theta)\end{aligned}$$

✓

9.2 Asymptotic properties of the MLE estimator

Result 9.4. (Asymptotic properties of MLE): Suppose X_1, \dots, X_n iid with pdf or pmf $f_{\theta_0}(x), \theta_0 \in \Omega$ where θ_0 is the true parameter and Ω is the parameter space.

Denote by $\hat{\theta}$ the MLE estimate based on X_1, \dots, X_n .

Suppose the following assumptions

- The density / pmf f_θ has the same support for all $\theta \in \Omega$, i.e.

$$\{x : f_\theta(x) > 0\}$$

- does not depend on θ
- θ_0 is an interior point of Ω
- The log-likelihood function $l_n(\theta) = l(\theta|X_1, \dots, X_n)$ is differentiable in θ
- $\hat{\theta}$ is the unique value of $\theta \in \Omega$ such that

$$l'_n(\theta) = 0$$

Then the following hold

- $\hat{\theta}$ consistent for θ_0 , i.e. as $n \rightarrow \infty$

$$\hat{\theta} \rightarrow \theta_0$$

- 2.

and

$$\sqrt{n} (\hat{\theta} - \theta_0) \xrightarrow{d} N\left(0, \frac{1}{I(\theta_0)}\right)$$

Proof. Normality implies consistency: if we know that the asymptotic distribution is normal with variance being the inverse of Fisher Information, then

$$(\hat{\theta} - \theta_0) = \frac{\sqrt{n} (\hat{\theta} - \theta)}{\sqrt{n}} \xrightarrow{p} 0 \text{ by Slutsky's}$$

Asymptotic normality:

$$l'_n(\hat{\theta}) = 0$$

By Taylor Expansion around θ_0 ,

$$\begin{aligned} l'_n(\hat{\theta}) &\approx l'_n(\theta_0) + (\hat{\theta} - \theta_0) l''_n(\theta_0) \\ \implies (\hat{\theta} - \theta_0) &\approx \frac{-l'_n(\theta_0)}{l''_n(\theta_0)} \\ \implies \sqrt{n} (\hat{\theta} - \theta_0) &\approx \frac{-\frac{1}{\sqrt{n}} l'_n(\theta_0)}{\frac{1}{n} l''_n(\theta_0)} \end{aligned}$$

Consider the denominator

$$\begin{aligned} &\frac{1}{n} l''_n(\theta_0) \\ &= \frac{1}{n} \sum_{i=1}^n \frac{\partial^2}{\partial \theta^2} \log f_\theta(X_i) \Big|_{\theta=\theta_0} \\ &\xrightarrow{p} \mathbb{E} \left[\frac{\partial^2}{\partial \theta^2} \log f_\theta(X_i) \Big|_{\theta=\theta_0} \right] \text{ by law of large numbers} \\ &= -I(\theta_0) \end{aligned}$$

Consider the numerator

$$\begin{aligned} &\frac{1}{\sqrt{n}} l'_n(\theta_0) \\ &= \frac{1}{\sqrt{n}} \sum_{i=1}^n \frac{\partial}{\partial \theta} \log f_\theta(X_i) \Big|_{\theta=\theta_0} \\ &= \sqrt{n} \left(\frac{1}{n} \sum_{i=1}^n \frac{\partial}{\partial \theta} \log f_\theta(X_i) \Big|_{\theta=\theta_0} \right) \\ &\xrightarrow{d} N(0, I(\theta_0)) \text{ by CLT} \end{aligned}$$

Hence

$$\begin{aligned}
& \sqrt{n} (\hat{\theta} - \theta_0) \\
& \approx -\sqrt{n} \left(\frac{l_n'(\theta_0)}{l_n''(\theta_0)} \right) \\
& \approx -\frac{\frac{1}{\sqrt{n}} l_n'(\theta_0)}{\frac{1}{n} l_n''(\theta_0)} \\
& \xrightarrow{d} \frac{N(0, I(\theta_0))}{I(\theta_0)} \text{ by Slutsky's} \\
& = N\left(0, \frac{1}{I(\theta_0)}\right)
\end{aligned}$$

✓

Example. (Bernoulli): $X_1, X_2, \dots, X_n \sim Ber(P)$. Then,

$$\begin{aligned}
f_p(x) &= p^x (1-p)^{1-x} \\
\log f_p(x) &= x \log p + (1-x) \log(1-p) \\
\frac{\partial}{\partial p} \log f_p(x) &= \frac{x}{p} - \frac{1-x}{1-p} \\
-\frac{\partial^2}{\partial p^2} \log f_p(x) &= \frac{x}{p^2} + \frac{1-x}{(1-p)^2}
\end{aligned}$$

Therefore

$$\begin{aligned}
I(p) &= -\mathbb{E} \left[\left(\frac{\partial^2}{\partial p^2} \log f_p(x) \right) \right] \\
&= \frac{p}{p^2} + \frac{1-p}{(1-p)^2} \\
&= \frac{1}{p} + \frac{1}{1-p} \\
&= \frac{1}{p(1-p)}
\end{aligned}$$

Hence

$$\begin{aligned}
\sqrt{n} (\hat{p}_{MLE} - p) &\xrightarrow{p} N\left(0, \frac{1}{I(p)}\right) \\
&= N(0, p(1-p))
\end{aligned}$$

Example. (Normal): $X_1, X_2, \dots, X_n \sim N(\mu, \sigma^2)$. The pdf of X is

$$\begin{aligned}
f_{\mu, \sigma^2}(x) &= \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(x-\mu)^2} \\
\log f_{\mu, \sigma^2}(x) &= -\frac{1}{2} \log 2\pi - \frac{1}{2} \log \sigma^2 - \frac{1}{2\sigma^2}(x-\mu)^2
\end{aligned}$$

Fisher information for μ :

$$\begin{aligned}
\frac{\partial}{\partial \mu} \log f(x) &= \frac{1}{\sigma^2}(x-\mu) \\
-\frac{\partial^2}{\partial \mu^2} \log f(x) &= \frac{1}{\sigma^2}
\end{aligned}$$

Therefore,

$$I(\mu) = \frac{1}{\sigma^2}$$

Hence

$$\sqrt{n}(\hat{\mu}_{MLE} - \mu) \xrightarrow{d} N(0, \sigma^2)$$

Fisher information for σ^2 :

$$\begin{aligned}
\frac{\partial}{\partial \sigma^2} \log f(x) &= -\frac{1}{2\sigma^2} + \frac{1}{2\sigma^4}(x-\mu)^2 \\
-\frac{\partial^2}{(\partial \sigma^2)^2} \log f(x) &= -\frac{1}{2\sigma^4} + \frac{1}{\sigma^6}(x-\mu)^2
\end{aligned}$$

Therefore,

$$I(\sigma^2) = -\mathbb{E} \left[\frac{\partial^2}{(\partial \sigma^2)^2} \right] = -\frac{1}{2\sigma^4} + \frac{1}{\sigma^4} = \frac{1}{2\sigma^4}$$

And

$$\begin{aligned} \sqrt{n} (\hat{\sigma}_{MLE} - \sigma^2) &\xrightarrow{d} N \left(0, \frac{1}{I(\sigma^2)} \right) \\ &= N(0, 2\sigma^4) \end{aligned}$$

Remark. As a heuristic

$$\hat{\theta} \approx N \left(\theta_0, \frac{1}{nI(\theta)} \right)$$

9.3 Asymptotic properties of the MLE estimator for multiple parameters

Let $\{f_{\theta}(x) : \theta \in \Omega\}$ be a family of distributions where $\theta = (\theta_1, \theta_2, \dots, \theta_k) \in \mathbb{R}^k$ is a k -dimensional parameter.

Theorem 9.5. Suppose X_1, \dots, X_n iid from distribution with density / pmf f_{θ} . Let $\hat{\theta}_n$ be the MLE based on X_1, X_2, \dots, X_n .

Let the $k \times k$ Fisher Information matrix be defined as

$$(I(\theta))_{ij} = Cov \left(\frac{\partial}{\partial \theta_i} \log f_{\theta}(x), \frac{\partial}{\partial \theta_j} \log f_{\theta}(x) \right)$$

Under the same conditions as the univariate normal,

$$\sqrt{n} (\hat{\theta}_n - \theta) \xrightarrow{d} N_k(\mathbf{0}, I(\theta)^{-1})$$

Proof. We WTS that

$$Cov \left(\frac{\partial}{\partial \theta_i} \log f_{\theta}(x), \frac{\partial}{\partial \theta_j} \log f_{\theta}(x) \right) = -\mathbb{E} \left[\frac{\partial^2}{\partial \theta_i \partial \theta_j} \log f_{\theta}(x) \right]$$

✓

Example. (Normal):

$$\begin{aligned} -\frac{\partial^2}{\partial \theta_i \partial \theta_j} l_n(\mu, \sigma^2) &= \frac{1}{\sigma^4} (x_i - \mu) \\ \mathbb{E} \left[-\frac{\partial^2}{\partial \theta_i \partial \theta_j} l_n(\mu, \sigma^2) \right] &= 0 \end{aligned}$$

The Fisher Information matrix is

$$I(\mu, \sigma^2) = \begin{pmatrix} \frac{1}{\sigma^2} & 0 \\ 0 & \frac{1}{2\sigma^4} \end{pmatrix}$$

Hence

$$I(\mu, \sigma^2)^{-1} = \begin{pmatrix} \sigma^2 & 0 \\ 0 & 2\sigma^4 \end{pmatrix}$$

and

$$\sqrt{n} \begin{pmatrix} \hat{\mu}_{MLE} - \mu \\ \hat{\sigma}_{MLE}^2 - \sigma^2 \end{pmatrix} = N_2 \left(\mathbf{0}, \begin{pmatrix} \sigma^2 & 0 \\ 0 & 2\sigma^4 \end{pmatrix} \right)$$

Example. (Gamma distribution)

10 Class 10

10.1 Cramer-Rao Lower Bound

Theorem 10.1. (Cramer-Rao Lower Bound) Consider a parametric model of distributions $\{f_\theta(x), \theta \in \Omega\}$ satisfying certain *mild regularity conditions*, and T is any unbiased estimator θ based on X_1, X_2, \dots, X_n iid. Then

$$Var(T) \geq \frac{1}{nI(\theta)}$$

Proof. Denote the score function

$$s(\theta, x) = \frac{\partial}{\partial \theta} \log f_\theta(x) = \frac{\frac{\partial f_\theta(x)}{\partial \theta}}{f_\theta(x)}$$

Let

$$S = s(X_1, X_2, \dots, X_n) = \sum_{i=1}^n s(\theta, X_i)$$

For any unbiased estimator T ,

$$\begin{aligned} (Corr(S, T))^2 &\leq 1 \\ \implies Cov(S, T)^2 &\leq Var(S) \cdot Var(T) \\ \implies Var(T) &\geq \frac{(Clv(S, T))^2}{Var(S)} \end{aligned}$$

Where

$$\begin{aligned} Var(S) &= nVar(s) \\ &= nVar\left(\frac{\partial}{\partial \theta} \log f_\theta(X_1)\right) \\ &= nI(\theta) \end{aligned}$$

Hence

$$Var(T) \geq \frac{(Cov(S, T))^2}{nI(\theta)}$$

To show that $Cov(S, T) = 1$, by unbiasedness of T

$$\begin{aligned} \theta &= \mathbb{E}[(T)] \\ &= \int_{\mathbb{R}^n} T(X_1, X_2, \dots, X_n) f_\theta(X_1) f_\theta(X_2) \dots f_\theta(X_n) dX_1 dX_2 \dots dX_n \end{aligned}$$

Taking the derivative wrt θ on both sides,

$$\begin{aligned} 1 &= \int_{\mathbb{R}^n} T(X_1, X_2, \dots, X_n) \left[\sum_{i=1}^n \left(\frac{\partial}{\partial \theta} f_\theta(x_i) \prod_{j \in [1, n], j \neq i} f_\theta(x_j) \right) \right] dX_1 dX_2 \dots dX_n \\ &= \int_{\mathbb{R}^n} T(X_1, \dots, X_n) S(X_1, \dots, X_n) f_\theta(X_1) f_\theta(X_2) \dots f_\theta(X_n) dX_1 dX_2 \dots dX_n \\ &= \mathbb{E}_\theta[TS] \end{aligned}$$

Since the score function has zero expectation,

$$1 = \mathbb{E}[TS] - \mathbb{E}[T]\mathbb{E}[S] = Cov(S, T)$$

✓

Remark. An unbiased estimator is said to be efficient if its variance is

$$\frac{1}{nI(\theta)}$$

Remark. Since the MLE is asymptotically unbiased and the variance of the MLE attains the Cramer-Rao Lower Bound asymptotically, the MLE is said to be **asymptotically efficient**.

10.2 Simple Linear Regression

10.2.1 Estimating β_0, β_1

Definition 10.2. (Simple Linear Regression) Suppose we have data $(X_1, Y_1), (X_2, Y_2), \dots, (X_n, Y_n)$ related by model

$$Y_i = \beta_0 + \beta_1 X_i + \epsilon_i$$

The **least squares estimates** are obtained via

$$(\hat{\beta}_0, \hat{\beta}_1) = \arg \min_{\beta_0, \beta_1} \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 X_i)^2$$

The **fitted values** are defined as

$$\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i$$

The **residuals** are defined as

$$E_i = Y_i - \hat{Y}_i$$

Remark. Define

$$f(\beta_0, \beta_1) = \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 X_i)^2$$

The LS estimates $(\hat{\beta}_0, \hat{\beta}_1)$ are solved via

$$\begin{cases} \frac{\partial f}{\partial \beta_0} = 0 \\ \frac{\partial f}{\partial \beta_1} = 0 \end{cases}$$

Result 10.3. The LS line always passes through the point (\bar{X}, \bar{Y})

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sum_{i=1}^n (X_i - \bar{X})^2}$$

$$\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X}$$

Remark. Note that the results thus far assumes nothing about the ϵ_i 's.

10.2.2 Estimating σ^2

Now assume that

$$\epsilon_i \stackrel{iid}{\sim} N(0, \sigma^2)$$

By definition

$$\begin{aligned} \epsilon_i &= Y_i - \beta_0 - \beta_1 X_i \\ \|\epsilon\|^2 &= \sum_{i=1}^n \epsilon_i^2 = \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 X_i)^2 \end{aligned}$$

If β_0, β_1 known, then

$$\mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n \epsilon_i^2 \right] = \sigma^2$$

and an unbiased estimator of σ^2 is

$$\frac{\|\epsilon\|^2}{n}$$

Definition 10.4. (Sum of Squared Errors, SSE): The Sum of squared errors is defined as

$$SSE := \sum_{i=1}^n E_i^2 = \sum_{i=1}^n (Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_i)^2$$

Result 10.5. $\hat{\sigma}^2$ defined below is an unbiased estimate of σ^2

$$SSE \sim \sigma^2 \chi_{n-2}^2 \implies \hat{\sigma}^2 = \frac{SSE}{n-2}$$

10.2.3 Facts about Least Square Estimates

Result 10.6. The LS estimates are unbiased.

Result 10.7. The LS estimates are normally distributed.

Result 10.8. The LS estimates find the vector $\hat{\mathbf{v}}$ in the plane spanned by the vectors $\mathbf{1}, \mathbf{X}$ that is the closest to \mathbf{Y} , where

$$\hat{\mathbf{v}} = \hat{\beta}_0 \mathbf{1} + \hat{\beta}_1 \mathbf{X}$$

Result 10.9. $\hat{\beta}_0, \hat{\beta}_1$ are independent of $\hat{\sigma}^2$

Result 10.10. $\hat{\beta}_0, \hat{\beta}_1$ coincides with the MLE estimates.

11 Class 11

11.1 Least Absolute Deviation Line

Definition 11.1. (Least Absolute Deviation) Suppose we have data $(X_1, Y_1), (X_2, Y_2), \dots, (X_n, Y_n)$ related by model

$$Y_i = \beta_0 + \beta_1 X_i + \epsilon_i$$

The **least absolute deviation line** is obtained by minimizing

$$\min_{\beta_0, \beta_1} \sum_{i=1}^n |Y_i - \beta_0 - \beta_1 X_i|$$

Result 11.2.

$$\arg \min_m \sum_{i=1}^n |X_i - m| = \text{median}\{X_1, \dots, X_n\}$$

Result 11.3. The LAD line passes through a pair of points $(X_i, Y_i), (X_j, Y_j)$.

Proof. For simplicity, assume n odd.

Begin by fixing β_1 and defining

$$Z_i = Y_i - \beta_1 X_i$$

Then

$$\begin{aligned} \arg \min_{\beta_0} \sum_{i=1}^n |Z_i - \beta_0| &= \arg \min_{\beta_0} \sum_{i=1}^n |Z_i - \beta_0| \\ &= \text{median}\{Z_1, \dots, Z_n\} \\ &= Z_{i_0} \text{ for some } i_0 \\ &= Y_{i_0} - \beta_1 X_{i_0} \end{aligned}$$

This implies the LAD line for fixed β_1 passes through some point (X_{i_0}, Y_{i_0}) .

by shifting the coordinate system, we can assume that the LAD line passes through the origin

$$Y - Y_{i_0} = \beta_1 (X - X_{i_0})$$

We solve for β_1 with

$$\begin{aligned} \min_{\beta_1} \sum_{i=1}^n |(Y_i - Y_{i_0}) - \beta_1 (X_i - X_{i_0})| \\ = \min_{\beta_1} \sum_{i=1}^n |Y_i - \beta_1 X_i| \end{aligned}$$

For simplicity, we just write Y_i, X_i

$$\begin{aligned} \min_{\beta_1} \sum_{i=1}^n |Y_i - \beta_1 X_i| \\ = \min_{\beta_1} \sum_{i=1}^n |X_i| \left| \frac{Y_i}{X_i} - \beta_1 \right| \end{aligned}$$

This is a weighted median problem, the sum of absolute deviations is piecewise linear between each data point and convex. Hence, the minimum is attained at some i_* , i.e.

$$\hat{\beta}_1 = \frac{Y_{i_*}}{X_{i_*}}$$

The LAD line is

$$Y = Y_{i_0} + \frac{Y_{i_*} - Y_{i_0}}{X_{i_*} - X_{i_0}} (X - X_{i_0})$$

✓

Remark. The LAD line can be computed by checking over all the $\binom{n}{2}$ pairwise lines determined by the data points. These lines are called the **elemental lines**.

Remarks. Note that

1. The LAD line is one of the elemental lines
2. The LAD estimates are the MLEs where the errors have the Laplace / Double Exponential distribution

$$Ae^{\frac{-|x|}{B}}$$

3. The slope of the LS line is a weighted average of the slopes of the elemental lines.

12 Class 12

12.1 Confidence Intervals

Definition 12.1. For an unknown parameter θ and a sample X_1, X_2, \dots, X_n , A $100(1 - \alpha)\%$ confidence interval for θ is a random interval

$$[L(X_1, X_2, \dots, X_n), U(X_1, X_2, \dots, X_n)]$$

such that

$$P(L \leq \theta \leq U) = 1 - \alpha$$

Remark. Interpretation of CI: If the experiment is repeated, $100(1 - \alpha)\%$ of intervals will contain true θ .

12.1.1 Exact Confidence Intervals

Definition 12.2. (t-distribution). If $U \sim N(0, 1)$ and $V \sim \chi_d^2$ and U, V independent, then

$$\frac{U}{\sqrt{\frac{V}{d}}} \sim t_d$$

i.e. the t -distribution with d degrees of freedom.

Example. (CI for μ, σ^2 , normal data) Suppose $X_1, X_2, \dots, X_n \sim iid N(\mu, \sigma^2)$ where μ, σ^2 unknown.

Confidence interval for μ :

1. Estimate μ with $\hat{\mu} = \bar{X}$
2. Find distribution of estimate

$$\begin{aligned}\bar{X} &\sim N\left(\mu, \frac{\sigma^2}{n}\right) \\ \implies Z &= \frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}} \sim N(0, 1) \\ \implies P\left(\left|\frac{\bar{X} - \mu}{\sigma/\sqrt{n}}\right| \leq z_{\alpha/2}\right) &= 1 - \alpha \\ \implies P\left(\left|\bar{X} - \mu\right| \leq \frac{\sigma}{\sqrt{n}}z_{\alpha/2}\right) &= 1 - \alpha \\ \implies P\left(-\frac{\sigma}{\sqrt{n}}z_{\alpha/2} \leq \mu \leq \frac{\sigma}{\sqrt{n}}z_{\alpha/2}\right) &= 1 - \alpha \\ \implies P\left(\bar{X} - \frac{\sigma}{\sqrt{n}}z_{\alpha/2} \leq \mu \leq \bar{X} + \frac{\sigma}{\sqrt{n}}z_{\alpha/2}\right) &= 1 - \alpha\end{aligned}$$

3. Estimate σ^2 :

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2 \sim \frac{\sigma^2}{n-1} \chi_{n-1}^2$$

Then, define

$$\begin{aligned}\hat{Z} &:= \frac{\bar{X} - \mu}{s/\sqrt{n}} \\ &= \frac{\sqrt{n}(\bar{X} - \mu)/\sigma}{\sqrt{s^2/\sigma^2}} \\ &= \frac{N(0, 1)}{\sqrt{\chi_{n-1}^2/(n-1)}} \\ &\sim t_{n-1}\end{aligned}$$

Hence

$$\begin{aligned}
P(\left| \hat{Z} \right| \leq t_{n-1,\alpha/2}) &= 1 - \alpha \\
\implies P\left(-t_{n-1,\alpha/2} \leq \left| \frac{\bar{X} - \mu}{s/\sqrt{n}} \right| \leq t_{n-1,\alpha/2} \right) &= 1 - \alpha \\
\implies P\left(\bar{X} - t_{n-1,\alpha/2} \frac{s}{\sqrt{n}} \leq \mu \leq \bar{X} + t_{n-1,\alpha/2} \frac{s}{\sqrt{n}} \right)
\end{aligned}$$

Confidence interval for σ^2

1. Estimate σ^2 with

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$$

2. Find the distribution of the estimate

$$s^2 \sim \frac{\sigma^2}{n-1} \chi_{n-1}^2 \implies \frac{s^2}{\sigma^2} \sim \frac{\chi_{n-1}^2}{n-1}$$

3. Since the distribution is free of unknown parameters, we can immediately find a confidence interval. Define $\chi_{n-1,\alpha}$ be the $(1-\alpha)$ -th percentile of the χ_{n-1}^2 distribution.

$$\begin{aligned}
P\left(\chi_{n-1,1-\alpha/2}^2 \leq \frac{(n-1)s^2}{\sigma^2} \leq \chi_{n-1,\alpha/2}^2 \right) &= 1 - \alpha \\
\implies P\left(\frac{(n-1)s^2}{\chi_{n-1,\alpha/2}} \leq \sigma^2 \leq \frac{(n-1)s^2}{\chi_{n-1,1-\alpha/2}} \right) &= 1 - \alpha
\end{aligned}$$

12.1.2 Asymptotic Confidence Intervals

Example. (Confidence interval for μ , arbitrary distribution): Suppose $X_1, X_2, \dots, X_n \sim iid F$ with $\mathbb{E}[X_1] = \mu, Var(X_1) = \sigma^2$.

Confidence interval for μ :

1. By CLT

$$\frac{\sqrt{n}(\bar{X} - \mu)}{\sigma} \xrightarrow{d} N(0, 1)$$

2. By LLN,

$$s^2 \xrightarrow{p} \sigma^2 \implies s \xrightarrow{p} \sigma$$

3. By Slutsky's

$$\frac{\sqrt{n}(\bar{X} - \mu)}{s} \xrightarrow{d} N(0, 1)$$

4. Find the CI

$$\begin{aligned}
P\left(\mu \in \left[\bar{X} - z_{\alpha/2} \frac{s}{\sqrt{n}}, \bar{X} + z_{\alpha/2} \frac{s}{\sqrt{n}} \right] \right) &\xrightarrow{n \rightarrow \infty} 1 - \alpha \\
\implies \left[\bar{X} - z_{\alpha/2} \frac{s}{\sqrt{n}}, \bar{X} + z_{\alpha/2} \frac{s}{\sqrt{n}} \right] \text{ is an asymptotically } 100(1 - \alpha) \% \text{ confidence interval for } \mu
\end{aligned}$$