

STAT24510 (W25): Statistical Theory and Methods IIa

Lecturer: Mei Wang

Notes by: Aden Chen

Tuesday 4th March, 2025

Contents

1	Introduction	3
2	Confidence Intervals	6
3	Change of Variable / Variable Transformation	10
4	Binomial and Poisson Distributions	12
5	Correlation	13
6	ANOVA (Analysis of Variance)	24
7	Two Way ANOVA	27
8	Comparison of Nested Models, Sequential ANOVA	30
9	Simple Linear Regression	34
10	Multiple Linear Regression Model	39
11	Simple Regression (A)	51

This set of notes is mostly unreadable.

1 Introduction

The goal of statistics is often to estimate a (population) parameter θ . From data, we may obtain point estimates $\hat{\theta}$ that depends on data, and construct confidence intervals to quantify the uncertainty of the estimate, with which we can conduct hypothesis testing.

In this course we will start from confidence intervals and hypothesis testing, and then move on to linear models. This course will be less structured; intend, it will be more like a collection of methods.

1.1 Examples of Construction: Review of Wald and Wilson

1.1.1 Pivotal method

Example 1.1 (Pivotal method, Normal, Known σ^2). Let $X_i \stackrel{\text{iid}}{\sim} \mathcal{N}(\mu, 5^2)$, $i = 1, \dots, n$. Our goal is to construct a $1 - \alpha$ CI for μ . We have the MLE estimator $\hat{\mu} = \bar{X} = n^{-1} \sum_{i=1}^n X_i$.

Note that

$$\frac{\bar{X} - \mathbb{E}[X]}{\sqrt{\text{Var}(X)}} = \frac{\bar{X} - \mu}{\sqrt{5^2/n}} \sim \mathcal{N}(0, 1).$$

In particular, note that the left side is a function of data and parameters, while the right side is free of parameters.

We thus have

$$\mathbb{P}\left(z_{\alpha/2} \leq \frac{\bar{X} - \mu}{\sqrt{5^2/n}} \leq z_{1-\alpha/2}\right) = 0.95,$$

using which we can construct the CI of μ : With $I := [\bar{X} - z_{1-\alpha/2}\sqrt{5^2/n}, \bar{X} + z_{1-\alpha/2}\sqrt{5^2/n}]$, we have $\mathbb{P}(\mu \in I) = 1 - \alpha$.

Notice that we obtain a probability statement of random interval containing a fixed quantity from a probability statement of a fixed interval containing a random quantity.

Example 1.2 (Pivotal method, Normal, Known μ). Let $X_i \stackrel{\text{iid}}{\sim} \mathcal{N}(4, \sigma^2)$, $i = 1, \dots, n$. The goal: CI for σ^2 , i.e., to find random variables L and U such that $\mathbb{P}(L \leq \sigma^2 \leq U) = 1 - \alpha$.

Note that

$$Y_i := \frac{X_i - 4}{\sigma} \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1)$$

and thus

$$T_n := \sum Y_i^2 = \sum \left(\frac{X_i - 4}{\sigma}\right)^2 \sim \chi_n^2.$$

Again, we obtained a function of data and parameters that follows a known distribution. From

$$\mathbb{P}\left(\chi_{n,\alpha/2}^2 \leq T_n \leq \chi_{n,1-\alpha/2}^2\right) = 1 - \alpha$$

we may again obtain the CI for σ^2 ,

$$\left[\frac{\sum (X_i - 4)^2}{\chi_{n,1-\alpha/2}^2}, \frac{\sum (X_i - 4)^2}{\chi_{n,\alpha/2}^2} \right].$$

Example 1.3 (Pivot failing). Let $X_i \stackrel{\text{iid}}{\sim} \text{Bernoulli}(p)$. The goal: CI for p . The MLE for p is $\hat{p} = \bar{X}$, thus we may be tempted to try

$$T_n := \frac{\hat{p} - p}{\sqrt{\text{Var}(\hat{p})}},$$

but the distribution of T_n depends on p . The method of pivots fail.

1.1.2 Asymptotic CI

I.e., when we have large sample size n .

Example 1.4 (Wald CI). Let X_i be iid with mean μ and variance σ^2 . From the CLT we have

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

that is,

$$\lim_{n \rightarrow \infty} \mathbb{P}(T_n \leq x) = P(Z \leq x) = \int_{-\infty}^x \frac{1}{\sqrt{2\pi}} e^{-t^2/2} dt.$$

Thus we have

$$\mathbb{P}\left(z_{\alpha/2} \leq \frac{\bar{X} - \mu}{\sqrt{\sigma^2/n}} \leq z_{1-\alpha/2}\right) \approx 1 - \alpha.$$

When σ^2 is known, we may derive an approximate CI for μ :

$$\mathbb{P}\left(\bar{X} - z_{1-\alpha/2}\sqrt{\sigma^2/n} \leq \mu \leq \bar{X} + z_{1-\alpha/2}\sqrt{\sigma^2/n}\right) \approx 1 - \alpha.$$

When σ^2 is unknown: If there exists random variables $U_n \rightarrow_p \sigma^2$ (that is, $\lim \mathbb{P}(U_n = \sigma^2) = 1$), then

$$T_n := \frac{\bar{X} - \mu}{\sqrt{U_n/n}} = \frac{\frac{\bar{X} - \mu}{\sqrt{\sigma^2/n}}}{\sqrt{U_n/\sigma^2}}$$

where $(\bar{X} - \mu)/\sqrt{\sigma^2/n} \xrightarrow{\mathcal{D}} \mathcal{N}(0, 1)$ and $\sqrt{U_n/\sigma^2} \rightarrow_p 1$, and thus by Slutsky's theorem we have

$$T_n \xrightarrow{\mathcal{D}} \mathcal{N}(0, 1),$$

using which we can again construct an approximate CI. Note that we used asymptotic approximation multiple times. This is called the Wald confidence interval.

Example 1.5 (Wald CI). Let $X_i \stackrel{\text{iid}}{\sim} \text{Poisson}(\lambda)$. The goal: asymptotic CI for λ . Note that we have the MLE of λ , $\hat{\lambda} = \bar{X}$, with $E[\hat{\lambda}] = \lambda$ and $\text{Var}[\hat{\lambda}] = \lambda/n$. We then have

$$\frac{\bar{X} - \lambda}{\sqrt{\lambda/n}} \approx Z \sim \mathcal{N}(0, 1).$$

We approximate a second time: $\frac{\bar{X} - \lambda}{\sqrt{\hat{\lambda}/n}} \approx Z$, from which we obtain the Wald CI for λ :

$$\left[\hat{\lambda} - z_{1-\alpha/2}\sqrt{\hat{\lambda}/n}, \hat{\lambda} + z_{1-\alpha/2}\sqrt{\hat{\lambda}/n} \right].$$

Example 1.6 (Wilson's method). Assume the same setup as above. Again, we use

$$\frac{\bar{X} - \lambda}{\sqrt{\lambda/n}} \approx Z \sim \mathcal{N}(0, 1),$$

which gives

$$\mathbb{P}\left(z_{\alpha/2} \leq \frac{\bar{X} - \lambda}{\sqrt{\lambda/n}} \leq z_{1-\alpha/2}\right) = \mathbb{P}\left(\left(\frac{\bar{X} - \lambda}{\sqrt{\lambda/n}}\right)^2 \leq z_{1-\alpha/2}^2\right) \approx 1 - \alpha.$$

Solving for λ in the middle expression gives the Wilson CI. We used one fewer approximation.

2 Confidence Intervals

2.1 Constructing CI

- Exact (i.e., coverage probability is exactly $1 - \alpha$) confidence intervals: **pivot method**:
 1. Find statistic $T_n = T(X, \theta)$ whose distribution is known and independent of θ . Such a statistic is called a pivot.
 2. Using knowledge on the distribution, find c_L and c_U such that $\mathbb{P}(c_L \leq T_n \leq c_U) = 1 - \alpha$.
 3. Convert the probability statement of the pivot to a probability statement of the parameter.
- Approximation methods, or asymptotic (sample size n is large) methods.
 - **Wald confidence interval** (the default CI produced by standard software): more than one approximations.
 - **Wilson's confidence interval** / score method / duality method: using CLT once.
 - * Wilson's CI is better in the sense that the actual coverage is closer to the desired coverage.
 - **Variance stabilization transformation (VST)** method.

Example 2.1 (Pivot, Normal). Assume $X_i \stackrel{\text{iid}}{\sim} \mathcal{N}(\mu, \sigma^2)$, $i = 1, \dots, n$. Goal: $(1 - \alpha)$ CI for μ .

1. Recall that

$$T_n = \frac{\bar{X} - \mu}{\sqrt{S^2/n}} \sim t_{n-1}.^1$$

2. We have

$$\mathbb{P}\left(t_{n-1, \alpha/2} \leq \frac{\bar{X} - \mu}{\sqrt{S^2/n}} \leq t_{n-1, 1-\alpha/2}\right) = 1 - \alpha.$$

3. Rearrangement gives

$$\mathbb{P}\left(\bar{X} - t_{n-1, 1-\alpha/2} \frac{S}{\sqrt{n}} \leq \mu \leq \bar{X} + t_{n-1, 1-\alpha/2} \frac{S}{\sqrt{n}}\right) = 1 - \alpha.$$

Example 2.2 (Wald, Bernoulli). Let $X_i \stackrel{\text{iid}}{\sim} \text{Bernoulli}(p)$, so that $\sum X_i \sim \text{Binomial}(n, p)$. Goal: $(1 - \alpha)$ CI for p . Note that $\hat{p} = \bar{X}$ is the MLE of p . CLT gives

$$\frac{\hat{p} - p}{\sqrt{\text{Var}(\hat{p})}} \approx \mathcal{N}(0, 1).$$

¹ t distribution has fatter tails than $\mathcal{N}(0, 1)$.

From Slutsky's theorem we have a **second approximation**

$$\text{Var}(\hat{p}) = \frac{p(1-p)}{n} \approx \frac{\hat{p}(1-\hat{p})}{n}.$$

This second approximation is characteristic of Wald's method. Using this second approximation we have

$$\frac{\bar{X} - p}{\sqrt{\hat{p}(1-\hat{p})/n}} \approx \mathcal{N}(0, 1),$$

with which we can construct the desired CI:

$$\left[\hat{p} - z_{1-\frac{\alpha}{2}} \sqrt{\hat{p}(1-\hat{p})/n}, \hat{p} + z_{1-\frac{\alpha}{2}} \sqrt{\hat{p}(1-\hat{p})/n} \right].$$

Example 2.3 (Wilson's CI; Bernoulli). We assume the same Bernoulli setup. Note that we obtained using just the CLT that

$$\mathbb{P} \left(z_{\frac{\alpha}{2}} \leq \frac{\hat{p} - p}{\sqrt{p(1-p)/n}} \leq z_{1-\frac{\alpha}{2}} \right) \approx 1 - \alpha.$$

We rewrite the middle expression as

$$n(\hat{p} - p)^2 \leq z_{1-\frac{\alpha}{2}}^2 p(1-p).$$

Solving a quadratic equation gives the desired Wilson's CI.

Remark 2.4.

- Note that the Wald CI is centered at \hat{p} , but Wilson's is not.
- Wilson's CI will always be contained in $[0, 1]$; the lower bound of Wald might be negative.

2.2 Variance Stabilization Transformation (VST) Method

Example 2.5 (VST, Poisson). Let $X_i \stackrel{\text{iid}}{\sim} \text{Poisson}(\lambda)$. Note that

$$\frac{\hat{\lambda} - \lambda}{\sqrt{\lambda/n}} \approx \mathcal{N}(0, 1).$$

and

$$\sqrt{n}(\hat{\lambda} - \lambda) \approx \mathcal{N}(0, \lambda).$$

Goal: find a transformation (usually smooth) g such that

$$\sqrt{n}(g(\hat{\lambda}) - g(\lambda)) \approx \mathcal{N}(0, 1).²$$

²Or just a normal distribution with fixed variance.

Tool: delta method (Taylor expansion for random variables). By Taylor expansion,

$$g(\hat{\lambda}) = g(\lambda) + g'(\lambda)(\hat{\lambda} - \lambda) + \frac{g''(\lambda)}{2}(\hat{\lambda} - \lambda)^2 + \dots$$

We thus have the approximation

$$g(\hat{\lambda}) \approx g(\lambda) + g'(\lambda)(\hat{\lambda} - \lambda) + O(n^{-1}),$$

where the last term follows from the fact that $(\hat{\lambda} - \lambda)/\sqrt{\lambda/n} \approx \mathcal{N}(0, 1)$. Then,

$$E[g(\hat{\lambda})] \approx g(\lambda) + g'(\lambda) E[\hat{\lambda} - \lambda].$$

Since $\hat{\lambda}$ is unbiased, we have

$$\sqrt{n}(g(\hat{\lambda}) - g(\lambda)) \approx g'(\lambda)\sqrt{n}(\hat{\lambda} - \lambda),$$

giving

$$\sqrt{n}(g(\hat{\lambda}) - g(\lambda)) \approx \mathcal{N}(0, [g'(\lambda)]^2 \lambda),$$

where we used the approximation

$$\text{Var}(g(\hat{\lambda})) = E[(g(\hat{\lambda}) - E[g(\hat{\lambda})])^2] \approx [g'(\lambda)]^2 \text{Var}(\hat{\lambda}).$$

To obtain $\sqrt{n}(g(\hat{\lambda}) - \lambda) \approx \mathcal{N}(0, 1)$, we need only $g'(\lambda)^2 = 1/\lambda$. $g(\lambda) = 2\sqrt{\lambda}$ will do.

That is, we have

$$\sqrt{n}(2\sqrt{\hat{\lambda}} - 2\sqrt{\lambda}) \approx \mathcal{N}(0, 1).$$

Using this we can obtain a CI for λ . Note that the left endpoint of the CI for $\sqrt{\lambda}$ may be negative, so in such cases when obtaining the CI for λ we need to use 0 instead.

2.3 Delta Method

Consider an unbiased estimator $\hat{\theta}$. $E[\hat{\theta}] = \theta$. Assume that the variance is a function of the mean, $\text{Var}[\hat{\theta}] = v(\theta)$.³ Goal: find a function $g(t)$ such that $\text{Var}[g(\hat{\theta})]$ is constant. Taylor gives

$$g(\hat{\theta}) \approx g(\theta) + g'(\theta)(\hat{\theta} - \theta).$$

Note that we have

$$E[g(\hat{\theta})] \approx g(\theta) + g'(\theta) E[\hat{\theta} - \theta] = g(\theta).$$

and

$$\text{Var}[g(\hat{\theta})] = E[(g(\hat{\theta}) - g(\theta))^2] \approx E[g'(\theta)^2 (\hat{\theta} - \theta)^2] = g'(\theta)^2 \text{Var}[\hat{\theta}],$$

which we want to be constant, say 1. Thus we seek g such that

$$g'(\theta) = \frac{1}{\sqrt{v(\theta)}}.$$

A choice is of course $g(\theta) = \int 1/\sqrt{v(\theta)} d\theta$.

³This is a common case.

Example 2.6 (Delta Method, Exponential, distribution of MLE). Consider $X_i \stackrel{\text{iid}}{\sim} \text{Exponential}(\lambda)$. Recall that the MLE of λ is $\hat{\lambda} = 1/\bar{X}$. By the CLT we have for large n that

$$\frac{\bar{X} - \mathbb{E} X}{\sqrt{\text{Var } \bar{X}}} \approx \mathcal{N}(0, 1),$$

That is,

$$\frac{\bar{X} - \frac{1}{\lambda}}{\sqrt{\frac{1}{n\lambda^2}}} = \sqrt{n} \left(\bar{X} - \frac{1}{\lambda} \right) \approx \mathcal{N}(0, \lambda^{-2}).$$

Goal: find an approximate distribution of $\hat{\lambda} = 1/\bar{X}$ using Delta Method. Using $g(t) = t^{-1}$ we have

$$\sqrt{n} (g(\hat{\lambda}) - g(\lambda)) \approx \mathcal{N}(0, [g'(\lambda)]^2 \text{Var}(\hat{\lambda})) = \mathcal{N}(0, \lambda^2).$$

(Note that Fisher's theorem gives another way to derive the above.)

Example 2.7 (VST, Exponential, CI). Consider the same setup as above. We can then derive an approximate CI for λ using VST. Recall that from Delta Method, we have

$$\sqrt{n} (g(\hat{\lambda}) - g(\lambda)) \approx \mathcal{N}(0, [g'(\lambda)]^2 \text{Var}(\hat{\lambda})).$$

Thus we seek g such that $g'(\lambda)^2 \lambda^2$ is a constant, say 1, or $g'(\lambda) = 1/\lambda$. It is easy to see that $g = \log$ is such an option. Then,

$$\sqrt{n} (\log \hat{\lambda} - \log \lambda) \approx \mathcal{N}(0, 1),$$

using which we can obtain an approximate CI for $\log \hat{\lambda}$ and then $\hat{\lambda}$:

$$\left[\hat{\lambda} \exp \left(-\frac{z_{1-\alpha/2}}{\sqrt{n}} \right), \hat{\lambda} \exp \left(\frac{z_{1-\alpha/2}}{\sqrt{n}} \right) \right]$$

Theorem 2.8 (Delta Method). *If X_n is such that*

$$\sqrt{n}(X_n - \theta) \xrightarrow{\mathcal{D}} \mathcal{N}(0, \sigma^2)$$

and g is continuously differentiable, then

$$\sqrt{n}(g(X_n) - g(\theta)) \xrightarrow{\mathcal{D}} \mathcal{N}(0, \sigma^2 [g'(\theta)]^2).$$

Remark 2.9. Intuition: We can write

$$\sqrt{n}(g(X_n) - g(\theta)) = g'(\tilde{\theta}_n) \sqrt{n}(X_n - \theta), \quad \tilde{\theta}_n \in (x_n, \theta).$$

We know that $g'(\tilde{\theta}_n) \rightarrow_p g'(\theta)$ and $\sqrt{n}(X_n - \theta) \xrightarrow{\mathcal{D}} \mathcal{N}(0, \sigma^2)$, so Slutsky's gives the desired result.

3 Change of Variable / Variable Transformation

3.1 The Discrete Case

Let X and $Y = g(X)$ be discrete. Then,

$$\mathbb{P}(Y = k) = \sum_{X: g(X)=k} \mathbb{P}(X = k).$$

3.2 The Continuous, 1D Case

Consider a continuous random variable X and transformation $Y = g(X)$ with g smooth with derivative vanishing nowhere, that is, $g' \neq 0$. These conditions guarantee that Y is also continuous. Suppose further that g is smooth and $g' \neq 0$. We have then that

$$f_Y(y) = f_X(x) \frac{1}{|g'(x)|} = f_X(g^{-1}(y)) \left| \frac{d}{dy} g^{-1}(y) \right|.$$

Example 3.1. Let $g' > 0$. Note that

$$\mathbb{P}(Y \leq y) = \mathbb{P}(g(X) \leq y) = \mathbb{P}(X \leq g^{-1}(y)) = F_X(g^{-1}(y)).$$

Thus

$$f_Y(y) = \frac{d}{dy} \mathbb{P}(Y \leq y) = \frac{d}{dy} F_X(g^{-1}(y)) = f_X(g^{-1}(y)) \frac{d}{dy} g^{-1}(y).$$

Similarly, if $g' < 0$, then

$$\mathbb{P}(Y \leq y) = \mathbb{P}(g(X) \leq y) = \mathbb{P}(X \geq g^{-1}(y)) = 1 - F_X(g^{-1}(y)).$$

Then

$$f_Y(y) = \frac{d}{dy} \mathbb{P}(Y \leq y) = -f_X(g^{-1}(y)) \frac{d}{dy} g^{-1}(y) = f_X(g^{-1}(y)) \left| \frac{d}{dy} g^{-1}(y) \right|.$$

In general, if g is piecewise monotone, then

$$f_Y(y) = \sum_{\{x: g(x)=y\}} f_X(x) \frac{1}{|g'(x)|}.$$

Example 3.2. $Y = X^2$.

$$f_Y(y) = \sum_{\{x: x^2=y\}} f_X(x) \frac{1}{|2x|} = f_X(\sqrt{y}) \frac{1}{2\sqrt{y}} + f_X(-\sqrt{y}) \frac{1}{|-2\sqrt{y}|}.$$

3.3 The Continuous, 2D Case

Suppose X and Y are continuous with densities f_X and f_Y and joint density $f_{X,Y}$. Consider

$$\begin{cases} U = g_1(X, Y) \\ V = g_2(X, Y) \end{cases}$$

Suppose g_i 's are smooth. Question: what is the joint density of U and V ?

$$f_{U,V}(u, v) = f_{X,Y}(x, y)|J|,$$

where

$$\begin{aligned} J = \det \frac{\partial(x, y)}{\partial(u, v)} &= \det \begin{bmatrix} \partial x / \partial u & \partial x / \partial v \\ \partial y / \partial u & \partial y / \partial v \end{bmatrix} \\ &= \frac{\partial x}{\partial u} \frac{\partial y}{\partial v} - \frac{\partial x}{\partial v} \frac{\partial y}{\partial u} \end{aligned}$$

and we assume $J \neq 0$.

We have also

$$f_{X,Y}(x, y) = f_{U,V}(u, v)|J|^{-1},$$

where

$$J^{-1} = \det \frac{\partial(u, v)}{\partial(x, y)} = \frac{1}{J}.$$

Example 3.3 (2D Change of Variables, Quotient). Let $U = X/Y$ with $Y \neq 0$ and $V = Y$. We have then that

$$\begin{cases} X = UV \\ Y = V \end{cases}.$$

Then

$$\frac{\partial(x, y)}{\partial(u, v)} = \begin{bmatrix} v & u \\ 0 & 1 \end{bmatrix}$$

so $J = v$. Alternatively,

$$\det \frac{\partial(u, v)}{\partial(x, y)} = \det \begin{bmatrix} \frac{1}{y} & -\frac{x}{y^2} \\ 0 & 1 \end{bmatrix} = \frac{1}{y} = J^{-1}.$$

This gives

$$f_{U,V}(u, v) = f_{X,Y}(uv, v)|v|.$$

If we assume further that $X \perp\!\!\!\perp Y$, then

$$f_{U,V}(u, v) = f_X(uv)f_Y(v)|v|.$$

4 Binomial and Poisson Distributions

Consider Binomial(n, p_n) with $np_n \rightarrow \lambda$ (thus $p_n \rightarrow 0$). We have then that

$$\begin{aligned}\mathbb{P}(X_n = k) &= \binom{n}{k} p^k (1-p)^{n-k} \\ &= \frac{n!}{k!(n-k)!} \left(\frac{\lambda_n}{n}\right)^k \left(1 - \frac{\lambda_n}{n}\right)^{n-k} \\ &= \frac{\lambda^k}{k!} \frac{n}{n} \frac{n-1}{n} \dots \frac{n-k+1}{n} \left(1 - \frac{\lambda}{n}\right)^n \left(1 - \frac{\lambda}{n}\right)^{-k}.\end{aligned}$$

Since k is fixed, as $n \rightarrow \infty$, and recalling that $(1 + X/n)^n \rightarrow e^X$, we have

$$\mathbb{P}(X = k) \rightarrow \frac{\lambda^k}{k!} e^{-\lambda}.$$

That is, if $np \rightarrow \lambda$, we have Binomial(n, p) \rightarrow Poisson(λ). In this sense, the Poisson distribution is the limit of the binomial distribution, and we can thus use it as an approximation for the binomial distribution.

But note that for large n and np , since a Binomial is the sum of Bernoulli's, the CLT can be used to approximate the binomial distribution as $\mathcal{N}(np, np(1-p))$. Similarly, the Poisson distribution can be approximated as $\mathcal{N}(\lambda, \lambda)$.

Example 4.1. Let $X \sim \text{Binomial}(n, p)$. The CLT gives

$$\mathbb{P}\left(\frac{X - np}{\sqrt{np(1-p)}} \leq x\right) \approx \Phi(x)$$

and thus

$$\begin{aligned}\mathbb{P}(X = k) &= \mathbb{P}(k-1 < X \leq k) = \mathbb{P}\left(\frac{k-1-np}{\sqrt{np(1-p)}} < \frac{X-np}{\sqrt{np(1-p)}} \leq \frac{k-np}{\sqrt{np(1-p)}}\right) \\ &\approx \Phi\left(\frac{k-np}{\sqrt{np(1-p)}}\right) - \Phi\left(\frac{k-1-np}{\sqrt{np(1-p)}}\right).\end{aligned}$$

It turns out that the **continuity correction** almost always gives a better approximation:

$$\begin{aligned}\mathbb{P}(X = k) &= \mathbb{P}(k-0.5 < X \leq k+0.5) = \mathbb{P}\left(\frac{k-0.5-np}{\sqrt{np(1-p)}} < \frac{X-np}{\sqrt{np(1-p)}} \leq \frac{k+0.5-np}{\sqrt{np(1-p)}}\right) \\ &\approx \Phi\left(\frac{k+0.5-np}{\sqrt{np(1-p)}}\right) - \Phi\left(\frac{k-0.5-np}{\sqrt{np(1-p)}}\right).\end{aligned}$$

5 Correlation

5.1 Bivariate Normal

Let (X, Y) be bivariate normal. That is, X and Y have joint density

$$f_{XY}(x, y) = \frac{1}{2\pi\sigma_X\sigma_Y\sqrt{1-\rho^2}} \exp\left(-\frac{1}{2(1-\rho^2)}\left(\frac{(x-\mu_X)^2}{\sigma_X^2} - 2\rho\frac{(x-\mu_X)(y-\mu_Y)}{\sigma_X\sigma_Y} + \frac{(y-\mu_Y)^2}{\sigma_Y^2}\right)\right).$$

If $X \perp Y$, then $\rho = 0$ and the joint density simplifies to $f_X(x)f_Y(y)$.

5.2 Bivariate Normal, the Standard Case

The standard bivariate normal (U, V) is

$$f_{UV}(u, v) = \frac{1}{2\pi(1-\rho^2)} \exp\left(-\frac{1}{2(1-\rho^2)}(u^2 + v^2 - 2\rho uv)\right).$$

Note that

$$u^2 + v^2 - 2\rho uv = (v - \rho u)^2 + (1 - \rho^2)u^2.$$

Thus we have

$$\begin{aligned} f_V(v) &= \int_{-\infty}^{\infty} f_{UV}(u, v) \, dv = \frac{1}{2\pi(1-\rho^2)} \int_{-\infty}^{\infty} \exp\left(-\frac{(1-\rho^2)u^2}{2(1-\rho^2)}\right) \exp\left(-\frac{(v-\rho u)^2}{2(1-\rho^2)}\right) \, dv \\ &= \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{v^2}{2}\right) \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi(1-\rho^2)}} \exp\left(-\frac{(v-\rho u)^2}{2(1-\rho^2)}\right) \, dv. \end{aligned}$$

So the marginal distributions are also normal. One can verify similarly that $E[UV] = \rho$, giving

$$\text{Cov}(U, V) = E(UV) - E(U)E(V) = \rho$$

and

$$\text{Corr}(U, V) = \frac{\text{Cov}(U, V)}{\sqrt{\text{Var}(U)\text{Var}(V)}} = \rho.$$

Note that

$$\begin{aligned} f_{U|V}(u|v) &= \frac{f_{UV}(u, v)}{f_V(v)} = \frac{\frac{1}{2\pi\sqrt{1-\rho^2}} \exp\left(-\frac{1}{2(1-\rho^2)}(u^2 + v^2 - 2\rho uv)\right)}{\frac{1}{\sqrt{2\pi}} \exp\left(-\frac{v^2}{2}\right)} \\ &= \frac{1}{\sqrt{2\pi}\sqrt{1-\rho^2}} \exp\left(-\frac{1}{2(1-\rho^2)}(u^2 + v^2 - 2\rho uv - (1-\rho^2)v^2)\right) \\ &= \frac{1}{\sqrt{2\pi}\sqrt{1-\rho^2}} \exp\left(-\frac{(u-\rho v)^2}{2(1-\rho^2)}\right). \end{aligned}$$

Thus,

$$U|V = v \sim \mathcal{N}(\rho v, (1-\rho^2))$$

and similarly

$$V|U = u \sim \mathcal{N}(\rho u, (1 - \rho^2)).$$

Thus the conditional expectation functions are

$$U = \rho V$$

and

$$V = \rho U.$$

Remark 5.1. This demonstrates regression toward the mean, since $\rho \leq 1$.

5.3 Bivariate Normal, the General Case

We may write

$$\begin{cases} X = \mu_X + \sigma_X U \sim \mathcal{N}(\mu_X, \sigma_X^2) \\ Y = \mu_Y + \sigma_Y V \sim \mathcal{N}(\mu_Y, \sigma_Y^2) \end{cases}.$$

We have

$$\text{Cov}(X, Y) = \text{Cov}(\mu_X + \sigma_X U, \mu_Y + \sigma_Y V) = \sigma_X \sigma_Y \text{Cov}(U, V)$$

and thus

$$\text{Corr}(X, Y) = \rho.$$

From

$$\frac{X - \mu_X}{\sigma_X} \mid \frac{Y - \mu_Y}{\sigma_Y} \sim \mathcal{N}\left(\rho \cdot \frac{Y - \mu_Y}{\sigma_Y}, 1 - \rho^2\right)$$

we know

$$X|Y = y \sim \mathcal{N}\left(\mu_X + \rho \frac{\sigma_X}{\sigma_Y}(y - \mu_Y), \sigma_X^2(1 - \rho^2)\right).$$

Similarly,

$$Y|X = x \sim \mathcal{N}\left(\mu_Y + \rho \frac{\sigma_Y}{\sigma_X}(x - \mu_X), \sigma_Y^2(1 - \rho^2)\right).$$

Note that the variance is smaller.

5.4 Bivariate Normal Data

Suppose (x_i, y_i) are iid bivariate normal with

$$(x_i, y_i) \stackrel{\text{iid}}{\sim} \mathcal{N}\left(\begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \begin{bmatrix} \sigma_1^2 & \rho \sigma_1 \sigma_2 \\ \rho \sigma_1 \sigma_2 & \sigma_2^2 \end{bmatrix}\right).$$

Then $X_i \sim \mathcal{N}(\mu_1, \sigma_1^2)$, $Y_i \sim \mathcal{N}(\mu_2, \sigma_2^2)$, and $\text{Cov}(X_i, Y_i) = \rho$.

5.5 MLE

We have

$$\begin{aligned}
L(\mu_X, \mu_Y, \sigma_X, \sigma_Y, \rho) &= \prod f(x_i, y_i) \\
&= \frac{1}{(2\pi)^2(1-\rho^2)^{n/2}} \\
&\quad \exp\left(\frac{1}{2(1-\rho^2)} \sum \left(\left(\frac{x_i - \mu_X}{\sigma_X}\right)^2 + \left(\frac{y_i - \mu_Y}{\sigma_Y}\right)^2 - 2\rho \left(\frac{x_i - \mu_X}{\sigma_X}\right) \left(\frac{y_i - \mu_Y}{\sigma_Y}\right) \right)\right).
\end{aligned}$$

We have⁴

$$\hat{\mu}_X = \bar{X}, \quad \hat{\mu}_Y = \bar{Y}, \quad \hat{\sigma}_X = n^{-1} \sum (X_i - \bar{X})^2, \quad \hat{\sigma}_Y = n^{-1} \sum (Y_i - \bar{Y})^2$$

and

$$R := \hat{\rho} = \frac{\sum (X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum (X_i - \bar{X})^2 \sum (Y_i - \bar{Y})^2}}.$$

Definition 5.2. R is called the **Pearson correlation coefficient**.

5.6 Distribution of the Pearson Correlation Coefficient

It is easy to see that $\text{supp } R = [-1, 1]$, but its distribution is quite complicated. We know however that

$$\mathbb{E}[\hat{\rho}] = \rho - \frac{\rho(1-\rho^2)}{2n} + \dots$$

and

$$\text{Var}[\hat{\rho}] = \frac{(1-\rho^2)^2}{n} + \dots$$

We have approximately that

$$\text{Var}[\hat{\rho}] \propto (1-\rho^2)^2 = v(\rho).$$

By Fisher we know that for large n we have

$$\hat{\rho} \sim \mathcal{N}\left(\rho, \frac{(1-\rho^2)^2}{n}\right).$$

To implement VST, we seek g such that $g'(\rho) = 1/(1-\rho^2)$. A choice is

$$g(\rho) = \int \frac{1}{(1+\rho)(1-\rho)} d\rho = \frac{1}{2} \int \frac{1}{1+\rho} + \frac{1}{1-\rho} d\rho = \frac{1}{2} \log \frac{1+\rho}{1-\rho}.$$

⁴These are *not* unbiased estimators. For the unbiased estimators, we need to divide by $n-1$ instead of n .

This is called the **Fisher transformation**. Note that it is monotone (and thus the construction of CIs is somewhat easy).

We have that

$$\frac{1}{2} \log \frac{1+\hat{\rho}}{1-\hat{\rho}} \sim \mathcal{N}\left(\frac{1}{2} \log \frac{1+\rho}{1-\rho}, \frac{1}{n-3}\right).$$

It turns out that $1/(n-3)$ is better than $1/n$ for $n > 5$.

Using this we can construct a CI for $g(\rho)$

$$\left[\frac{1}{2} \log \frac{1+\hat{\rho}}{1-\hat{\rho}} - z_{1-\frac{\alpha}{2}} \frac{1}{\sqrt{n-3}}, \frac{1}{2} \log \frac{1+\hat{\rho}}{1-\hat{\rho}} + z_{1-\frac{\alpha}{2}} \frac{1}{\sqrt{n-3}} \right] = [L, U]$$

and for then ρ : Note that

$$T = \frac{1}{2} \log \frac{1+R}{1-R} \iff e^{2T} = \frac{1+R}{1-R} \iff R = \frac{e^{2T}-1}{e^{2T}+1} = \tanh T.$$

So a CI for ρ is

$$\left[\frac{e^{2L}-1}{e^{2L}+1}, \frac{e^{2U}-1}{e^{2U}+1} \right].$$

Remark 5.3. This is the standard way to construct a CI for the correlation coefficient in software packages, regardless of the distribution of the original data. In our derivation, however, we assumed bivariate normality. For small n , be mindful of the many layers of approximation!

5.6.1 The Matrix Notation

For the standard bivariate normal, with

$$\begin{bmatrix} U \\ V \end{bmatrix} \sim \mathcal{N}\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix}\right),$$

we may write $[U, V]^\top \sim \mathcal{N}([0, 0]^\top, \Sigma)$, and then

$$f(u, v) = \frac{1}{2\pi|\Sigma|^{\frac{1}{2}}} \exp\left(\frac{1}{2}[uv]^\top \Sigma^{-1}[uv]\right),$$

where

$$\Sigma^{-1} = \frac{1}{1-\rho^2} \begin{bmatrix} 1 & -\rho \\ -\rho & 1 \end{bmatrix}.$$

For the more general case, with $X \sim \mathcal{N}(\mu, \Sigma)$, we have

$$f(x) = \frac{1}{2\pi|\Sigma|^{\frac{1}{2}}} \exp\left(\frac{1}{2}[x-\mu]^\top \Sigma^{-1}[x-\mu]\right),$$

5.6.2 Properties

Proposition 5.4. *For bivariate normally distributed variables, we have the variables are uncorrelated if and only if they are independence.*

Proposition 5.5. *If $X \sim \mathcal{N}_2(\mu, \Sigma)$ if and only if $aX + bY$ is of (univariate) normal distribution for any $a, b \in \mathbb{R}$.*

Example 5.6 (Non-example). Let $X \sim \mathcal{N}(0, 1)$ and

$$Y = \begin{cases} X & \text{with probability } 1/2 \\ -X & \text{with probability } 1/2 \end{cases}.$$

Then $Y \sim \mathcal{N}(0, 1)$, but $[XY]^\top$ is not of bivariate normal distribution. A way to see this is that the linear combination $X + Y$ has point mass at 0 with probability 1/2.

Proposition 5.7. *If $Z_1, Z_2 \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1)$, then*

$$\begin{bmatrix} X \\ Y \end{bmatrix} = A \begin{bmatrix} Z_1 \\ Z_2 \end{bmatrix} = \begin{bmatrix} a_{11}Z_1 + a_{12}Z_2 \\ a_{21}Z_1 + a_{22}Z_2 \end{bmatrix} \sim \mathcal{N}_2 \left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, AA^\top \right)$$

is bivariate normal.

Proposition 5.8. *If $[X_1 X_2]^\top \sim \mathcal{N}_2(\mu, \Sigma)$ and $[Y_1 Y_2]^\top = A[X_1 X_2]^\top$ is a linear transformation of $[X_1 X_2]^\top$ such that A^{-1} exists, then $[Y_1 Y_2]^\top \sim \mathcal{N}(A\mu, A\Sigma A')$. Similarly, $A[X_1 X_2]^\top + [b_1 b_2]^\top \sim \mathcal{N}_2([b_1 b_2]^\top + A\mu, A\Sigma A')$.*

5.7 Theorem of Sampling Distributions

5.7.1 The χ and t Distributions

χ^2 distribution. Let $Z_1, \dots, Z_k \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1)$. Then

$$Z_1^2 + \dots + Z_n^2 \sim \chi_k^2.$$

It turns out that $\chi_k^2 \sim \text{Gamma}(k/2, 1/2)$. If $Z \sim \mathcal{N}(0, 1)$ and $W \sim \chi_n^2$ with Z and W independent, then

$$\frac{Z}{\sqrt{W/k}} \sim t_k.$$

We have

$$f_{t_k}(t) = C_k \frac{1}{(1 + t^2/k)^{\frac{k+1}{2}}},$$

where C_k is a constant.

Remark 5.9.

- Note that as $k \rightarrow \infty$,

$$(1 + t^2/k)^{\frac{k+1}{2}} \longrightarrow e^{t^2/2}.$$

- $t_1 = Z_1/Z_2$ with $Z_1, Z_2 \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1)$ is also called the **Cauchy distribution**. It has density

$$f_1(t) = \frac{1}{\pi} \cdot \frac{1}{1+t^2}.$$

Note however that it does not have an expected value, since

$$\frac{1}{\pi} \int_0^\infty \frac{t}{1+t^2} dt \longrightarrow \infty,$$

and similarly the left side also diverges. The variance also does not exist. The law of large number does not apply!

Theorem 5.10. Suppose $X_i \stackrel{\text{iid}}{\sim} \mathcal{N}(\mu, \sigma)$. Then we have

- (i) $\frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} \sim \mathcal{N}(0, 1)$.
- (ii) $\frac{1}{\sigma^2} \sum (X_i - \bar{X}_n)^2 \sim \chi_{n-1}^2$.
- (iii) $\sum (X_i - \bar{X}_n)^2 \perp\!\!\!\perp \bar{X}_n$.
- (iv) $\frac{\bar{X}_n - \mu}{\sqrt{S^2/n}} \sim t_{n-1}$, where $S^2 = \frac{1}{n-1} \sum X_i$ is the sample variance.

These are exact distributions, not approximations.

Proof.

- (i) Follows from linearity of the normal distribution.

- (ii) & (iii) We prove by induction. Let $Z_i := (X_i - \mu)/\sigma$. Then $\bar{Z}_n = n^{-1} \sum Z_i = (\bar{X} - \mu)/\sigma$. Thus we have

$$\frac{1}{\sigma^2} \sum (X_i - \bar{X}_n)^2 = \sum \left(\frac{X_i - \bar{X}_n}{\sigma} \right)^2 = \sum (Z_i - \bar{Z})^2.$$

We need thus only show that $\sum (Z_i - \bar{Z})^2 \sim \chi_{n-1}^2$ and is independent of \bar{Z}_n .

For $n = 2$, we have

$$(Z_1 - \bar{Z})^2 + (Z_2 - \bar{Z})^2 = \left(Z_1 - \frac{Z_1 + Z_2}{2} \right)^2 + \left(Z_2 - \frac{Z_1 + Z_2}{2} \right)^2 = \frac{(Z_1 - Z_2)^2}{2}.$$

Recalling that $(Z_1 - Z_2)/\sqrt{2} \sim \mathcal{N}(0, 1)$, we have $(Z_1 - Z_2)^2/2 \sim \chi_1^2$.

Next, recall that for bivariate normal (X, Y) we have that $X \perp\!\!\!\perp Y$ if and only if $\rho = 0$. Thus, noting that $[Z_1 - Z_2 \ Z_1 + Z_2]^\top$ is bivariate normal and

$$\text{Cov}(Z_1 - Z_2, Z_1 + Z_2) = \text{Var}(Z_1) - \text{Var}(Z_2) = 0,$$

we know that $Z_1 - Z_2 \perp\!\!\!\perp Z_1 + Z_2$ and thus $(Z_1 + Z_2)^2/2 \perp\!\!\!\perp (Z_1 + Z_2)/2 = \bar{Z}_2$.

For the general case, suppose that $\sum_{i=1}^n (Z_i - \bar{Z})n^2$ is of χ_{n-1}^2 and is independent of \bar{Z}_n . Note that

$$\begin{aligned} \sum_{i=1}^n (Z_i - \bar{Z}_{n+1})^2 &= \sum_{i=1}^n (Z_i - \bar{Z}_n + \bar{Z}_n - \bar{Z}_{n+1})^2 + (Z_{n+1} - \bar{Z}_{n+1})^2 \\ &= \sum_{i=1}^n (Z_i - \bar{Z}_n)^2 + 2 \sum_{i=1}^n (Z_i - \bar{Z}_n)(\bar{Z}_n - \bar{Z}_{n+1}) \\ &\quad + \sum_{i=1}^n (\bar{Z}_n - \bar{Z}_{n+1})^2 + (Z_{n+1} - \bar{Z}_{n+1})^2. \end{aligned}$$

Note that the first term is χ_{n-1}^2 and the second term is 0. Thus we need only show that

$$J := n(\bar{Z}_n - \bar{Z}_{n+1})^2 + (Z_{n+1} - \bar{Z}_{n+1})^2$$

is of χ_1^2 distribution and is independent of the first term.

Observe now that

$$\bar{Z}_{n+1} = \frac{n}{n+1} \bar{Z}_n + \frac{Z_{n+1}}{n+1},$$

and thus

$$\begin{aligned} J &= n \left(\bar{Z}_n - \frac{n}{n+1} \bar{Z}_n - \frac{1}{n+1} Z_{n+1} \right)^2 + \left(Z_{n+1} - \frac{n}{n+1} \bar{Z}_n - \frac{1}{n+1} Z_{n+1} \right)^2 \\ &= \frac{n}{(n+1)^2} (\bar{Z}_n - Z_{n+1})^2 + \frac{n^2}{(n+1)^2} (Z_{n+1} - \bar{Z}_n)^2 \\ &= \frac{n+n^2}{(n+1)^2} (\bar{Z}_n - Z_{n+1})^2 = \frac{n}{n+1} (\bar{Z}_n - Z_{n+1})^2. \end{aligned}$$

Since $\bar{Z} - Z_{n+1} \sim \mathcal{N}(0, 1/n+1)$, we know that J is of χ_1^2 distribution. For claim (ii), it remains to show that $J = n/(n+1) \cdot (\bar{Z}_n - Z_{n+1})^2$ is independent to $\sum_{i=1}^n (Z_i - \bar{Z}_n)^2$. Since Z_{n+1} , \bar{Z}_n , and $\sum_{i=1}^n (Z_i - \bar{Z}_n)^2$ are pairwise independent, we obtain (ii). To complete the proof, note that a similar argument as before shows (iii).

(iv) We have

$$\frac{\bar{X}_n - \mu}{\sqrt{S^2/n}} = \frac{\frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}}}{\sqrt{\frac{\sum (X_i - \bar{X})^2}{\sigma^2} / (n-1)}} = \frac{Z}{\sqrt{W/n-1}},$$

where $Z \sim \mathcal{N}(0, 1)$ and $W \sim \chi_n^2$.

□

5.8 Two-Sample

Suppose $X_i \stackrel{\text{iid}}{\sim} \mathcal{N}(\mu_1, \sigma_1^2)$ and $X_i \stackrel{\text{iid}}{\sim} \mathcal{N}(\mu_2, \sigma_2^2)$ and we want to estimate $\mu_1 - \mu_2$. A choice of an estimator is $\hat{\mu}_1 - \hat{\mu}_2 = \bar{X} - \bar{Y}$, but how can we obtain an CI?

Case 1: A “paired” sample. Suppose $\text{Cov}(X_1, Y_1) > 0$. We can use

$$\text{Var}(\bar{X} - \bar{Y}) = \text{Var}(\bar{X}) + \text{Var}(\bar{Y}) - 2\text{Cov}(\bar{X}, \bar{Y}) < \text{Var}(\bar{X}) + \text{Var}(\bar{Y})$$

to obtain a narrower CI.

Define $D_i := X_i - Y_i$ to reduce the problem to that of a one sample problem. We have in particular that

$$D_i \stackrel{\text{iid}}{\sim} \mathcal{N}(\mu_1 - \mu_2, \sigma_1^2 + \sigma_2^2 - 2\rho\sigma_1\sigma_2),$$

where $\rho := \text{Corr}(X_i, Y_i)$.

If σ is unknown, we may use $\hat{\text{Var}}[\bar{D}] = S_D^2/n$ to obtain

$$\frac{\bar{D} - (\mu_1 - \mu_2)}{\sqrt{S_D^2/n}} \sim t_{n-1}.$$

Case 2: $\text{Cov}(X_i, Y_i) = 0$ and $\sigma_1 = \sigma_2$. Note that the normality assumption gives $X_i \perp\!\!\!\perp Y_i$. Then, $\text{Var}[\bar{X} - \bar{Y}] = \text{Var}[\bar{X}] + \text{Var}[\bar{Y}] = \sigma^2/n + \sigma^2/m$. We may use the sample variances to estimate σ^2 :

$$S_X^2 = \frac{1}{n-1} \sum (X_i - \bar{X})^2, \quad S_Y^2 = \frac{1}{m-1} \sum (Y_i - \bar{Y})^2, \\ S_{\text{pool}}^2 = \frac{(n-1)S_X^2 + (m-1)S_Y^2}{n+m-2}.$$

We use the last term for higher power and narrower CI. We use

$$\hat{\text{Var}}[\bar{X} - \bar{Y}] = S_{\text{pool}}^2 \left(\frac{1}{n} + \frac{1}{m} \right)$$

to obtain

$$\frac{(\bar{X} - \bar{Y}) - (\mu_1 - \mu_2)}{S_{\text{pool}} \sqrt{\frac{1}{n} + \frac{1}{m}}} \sim t_{n+m-2}$$

and

$$\frac{S_{\text{pool}}^2}{\sigma^2} \sim \chi_{n+m-2}^2.$$

Case 3: $X_i \perp\!\!\!\perp Y_i$ and $\sigma_1 \neq \sigma_2$. We have $\text{Var}[\bar{X} - \bar{Y}] = \sigma_1^2/n + \sigma_2^2/m$. Then, approximately, we have

$$T := \frac{(\bar{X} - \bar{Y}) - (\mu_1 - \mu_2)}{\sqrt{\sigma_1^2/n + \sigma_2^2/m}} \sim \mathcal{N}(0, 1).$$

It turns out that for moderate sized n and m , a better approximation, called the **Welch–Satterthwaite** approximation is

$$T \sim t_\nu,$$

where ν is a function of sample variances and sample sizes. This is what is implemented in software packages.

5.8.1 Hypothesis Testing for Two-Sample

We define the likelihood ratio as

$$LR := \frac{\max_{H_0} \text{likelihood function}}{\max_{H_0 \cup H_a} \text{likelihood function}} \in (0, 1].$$

Note that if H_0 is true, LR is likely to be large (close to 1).

Remark 5.11. To maximize $f \in C^2$, we impose:

- (FOC) $\nabla f(x) = 0$
- (SOC) $H(x) = [\partial^2 f / \partial x_i \partial x_j]$ to be negative definite.

Definition 5.12. Let $A = [a_{ij}]_{p \times p}$ be a symmetric matrix. We say A is **negative definite** if one of the following equivalent statements hold:

- All eigenvalues of A are negative.
- Let M_k be the determinant of the k th leading principal minor. If $M_k < 0$ when k is odd, and $M_k > 0$ when k is even.

Proposition 5.13. Under H_0 , we have the following approximate asymptotic distribution of LR :

$$-2 \log LR \approx \chi_d^2,$$

where $d = \dim(H_a \cup H_0) - \dim H_0$.

Example 5.14. Suppose $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} \mathcal{N}(\mu_1, \sigma^2)$ and $Y_1, \dots, Y_m \stackrel{\text{iid}}{\sim} \mathcal{N}(\mu_2, \sigma^2)$ and $X_i \perp\!\!\!\perp Y_i$. Consider the hypothesis $H_0 : \mu_1 = \mu_2$ and $H_a : \mu_1 \neq \mu_2$. Under H_0 we have $\mu_1 = \mu_2 = \mu \in \mathbb{R}$ and $\sigma^2 \in \mathbb{R}$, so $\dim H_0 = 2$. Similarly we have $\dim(H_a \cup H_0) = 3$. Thus $d = 1$.

We have

$$\begin{aligned} L(\mu_1, \mu_2, \sigma^2) &= \prod f_X(x_i) f_Y(y_i) \\ &= \frac{1}{(2\pi\sigma^2)^{\frac{n+m}{2}}} \exp\left(-\frac{1}{2\sigma^2} \sum (x_i - \mu_1)^2 - \frac{1}{2\sigma^2} \sum (y_i - \mu_2)^2\right). \end{aligned}$$

The log-likelihood function is

$$\ell(\mu_1, \mu_2, \sigma^2) = -\frac{n+m}{2} \log(2\pi) - \frac{n+m}{2} \log \sigma^2 - \frac{1}{2\sigma^2} \sum (x_i - \mu_1)^2 - \frac{1}{2\sigma^2} \sum (y_i - \mu_2)^2.$$

Under H_0 we have $\mu_1 = \mu_2 = \mu_0$, and for large n ,

$$-2 \log LR \approx \chi_1^2.$$

Further, under H_0 ,

$$\frac{\partial \log L}{\partial \mu_0} = -\frac{1}{\sigma^2} \sum (x_i - \mu_0) - \frac{1}{\sigma^2} \sum (y_i - \mu_0).$$

Setting $\partial \log L / \partial \mu_0 = 0$, we get

$$\hat{\mu}_0 = \frac{\sum x_i + \sum y_i}{n + m}$$

for any σ^2 . Thus

$$\max_{H_0} L(\mu_1, \mu_2, \sigma^2) = L(\hat{\mu}_0).$$

Next, we have

$$\frac{\partial \log L}{\partial \sigma^2} = -\frac{n + m}{2\sigma^2} + \frac{1}{2(\sigma^2)^2} \sum (x_i - \mu_1)^2 + \frac{1}{2(\sigma^2)^2} \sum (y_i - \mu_2)^2.$$

Setting $\partial \log L / \partial \sigma^2 = 0$ at $\mu_1 = \mu_2 = \hat{\mu}_0$, we get

$$\hat{\sigma}_0^2 = \frac{1}{n + m} \left(\sum (x_i - \hat{\mu}_0)^2 + \sum (y_i - \hat{\mu}_0)^2 \right).$$

We have then that

$$\begin{aligned} \max_{H_0} L(\mu_1, \mu_2, \sigma^2) &= L(\hat{\mu}_0, \hat{\mu}_0, \hat{\sigma}_0^2) \\ &= \left(\frac{1}{\sqrt{\pi \hat{\sigma}_0^2}} \right)^{n+m} \exp \left(-\frac{\sum (x_i - \hat{\mu}_0)^2 + \sum (y_j - \hat{\mu}_0)^2}{2\hat{\sigma}_0^2} \right) \\ &= \left(\frac{1}{\sqrt{\pi \hat{\sigma}_0^2}} \right)^{n+m} \exp \left(-\frac{n + m}{2} \right). \end{aligned}$$

Under $H_0 \cup H_a$: we have

$$\frac{\partial \log L}{\partial \mu_1} = \frac{1}{\sigma^2} \sum (x_i - \mu_1), \quad \frac{\partial \log L}{\partial \mu_2} = \frac{1}{\sigma^2} \sum (y_i - \mu_2).$$

Setting the above terms to zero gives

$$\hat{\mu}_1 = \bar{X}, \quad \hat{\mu}_2 = \bar{Y}.$$

Again, setting $\partial \log L / \partial \sigma^2 = 0$ at $\mu_1 = \hat{\mu}_1$ and $\mu_2 = \hat{\mu}_2$, we get

$$\hat{\sigma}^2 = \frac{1}{n + m} \left(\sum (x_i - \hat{\mu}_1)^2 + \sum (y_i - \hat{\mu}_2)^2 \right).$$

We have then that

$$\begin{aligned}
\max_{H_0 \cup H_a} L(\mu_1, \mu_2, \sigma^2) &= L(\hat{\mu}_1, \hat{\mu}_2, \hat{\sigma}^2) \\
&= \left(\frac{1}{\sqrt{\pi \hat{\sigma}^2}} \right)^{n+m} \exp \left(-\frac{\sum (x_i - \hat{\mu}_1)^2 + \sum (y_j - \hat{\mu}_2)^2}{2\hat{\sigma}^2} \right) \\
&= \left(\frac{1}{\sqrt{\pi \hat{\sigma}^2}} \right)^{n+m} \exp \left(-\frac{n+m}{2} \right).
\end{aligned}$$

Then,

$$LR = \left(\frac{\hat{\sigma}^2}{\hat{\sigma}_0^2} \right)^{\frac{n+m}{2}} = \left(\frac{\sum (x_i - \bar{x})^2 + \sum (y_i - \bar{y})^2}{\sum (x_i - \hat{\mu}_0)^2 + \sum (y_i - \hat{\mu}_0)^2} \right)^{\frac{n+m}{2}}.$$

6 ANOVA (Analysis of Variance)

This chapter
is very wrong

For $i = 1, \dots, g$ with $g > 2$ suppose $Y_{i1}, \dots, Y_{in} \stackrel{\text{iid}}{\sim} \mathcal{N}(\mu_i, \sigma^2)$. Consider $H_0 : \mu_1 = \dots = \mu_g$ and $H_a : \mu_i \neq \mu_j$ for some $i \neq j$. (Consider for example the outcomes after different treatment levels i .)

Note that there are two types of variation: within group and across group. Denote the mean of the whole sample as $\bar{y} = y_{..}$. We may decompose the variance accordingly:

$$\begin{aligned} \sum_{i=1}^g \sum_{j=1}^n (y_{ij} - \bar{y})^2 &= \sum_{i=1}^g \sum_{j=1}^n [(y_{ij} - \bar{y}_i) - (\bar{y}_i - \bar{y})]^2 \\ &= \sum_{i=1}^g \sum_{j=1}^n (y_{ij} - \bar{y}_i)^2 + \sum_{i=1}^g n_i (\bar{y}_i - \bar{y})^2, \end{aligned}$$

where the first term is within group variation (the noise), and the second term is the across group variation (the treatment effect, the signal). Equivalently, we write this decomposition as

$$SS_{\text{total}} = SS_{\text{residual}} + SS_{\text{treatment}},$$

where SS stands for sum of squares.

6.0.1 ANOVA Table

Source	SS	df	MS	Var Ratio (F)
Treatment	SS_{trt}	$g - 1$	$SS_{\text{trt}}/(g - 1)$	MS_{trt}/MS_e
Residual (error)	SS_e	$n - g$	$SS_e/(n - g)$	
Total	SS_{total}	$n - 1$		

6.0.2 F distribution

If $W_1 \sim \chi_{k_1}^2$ and $W_2 \sim \chi_{k_2}^2$ with $W_1 \perp W_2$, then

$$\frac{W_1/k_1}{W_2/k_2} \sim F_{k_1, k_2}.$$

Note that $\text{supp } F \subset \mathbb{R}_+$.

Proposition 6.1. *We have under H_0 that*

$$\frac{MS_{\text{trt}}}{MS_e} \sim F_{g-1, n-g}.$$

Perhaps some intuition:

Note that

$$\begin{aligned} \mathbb{E}[SS_e] &= \mathbb{E} \left[\sum_i (n_i - 1) \sum \frac{(Y_{ij} - \bar{Y}_i)^2}{n_i - 1} \right] \\ &= \sum_i (n_i - 1) \mathbb{E}[S_i^2] = \sum_i (n_i - 1) \sigma^2 = (n - g) \sigma^2. \end{aligned}$$

Recalling that

$$\frac{1}{\sigma^2} \sum_j (Y_{ij} - \bar{Y}_i)^2 \sim \chi_{n_i-1}^2,$$

we have

$$\frac{1}{\sigma^2} SS_e \sim \chi_{n-g}^2.$$

Next,

$$\begin{aligned} E[SS_{\text{trt}}] &= E \left[\sum n_i (\bar{Y}_i - \bar{Y})^2 \right] = \sum n_i E [(\bar{Y}_i - \bar{Y})^2] \\ &= \sum n_i [\text{Var}[\bar{Y}_i - \bar{Y}] + [E(\bar{Y}_i - \bar{Y})]^2]. \end{aligned}$$

Now note that

$$E[\bar{Y}_i - \bar{Y}] = E \left[\frac{1}{n_i} \sum_j Y_{ij} \right] - E \left[\frac{1}{n} \sum_i \sum_j Y_{ij} \right] = \mu_i - \frac{1}{n} \sum_i n_i \mu_i =: \mu_i - \mu.$$

Noting that $\text{Cov}(Y_{ij}, Y_{kl}) = 0$ for $i \neq k$ and $\text{Cov}(Y_{ij}, Y_{ik}) = 0$ for $j \neq k$, we have

$$\text{Cov}(\bar{Y}_i, \bar{Y}) = \text{Cov} \left(\frac{1}{n_i} \sum_j Y_{ij}, \frac{1}{n} \sum_i \sum_j Y_{ij} \right) = \frac{1}{n} \frac{1}{n_i} \sum_i \text{Var}(Y_{ij}) = \frac{\sigma^2}{n}.$$

We have thus that

$$\text{Var}(\bar{Y}_i - \bar{Y}) = \text{Var}(\bar{Y}_i) + \text{Var}(\bar{Y}) - 2 \text{Cov}(\bar{Y}_i, \bar{Y}) = \frac{\sigma^2}{n_i} + \frac{\sigma^2}{n} - 2 \frac{\sigma^2}{n} = \frac{\sigma^2}{n_i} - \frac{\sigma^2}{n}.$$

Plugging back,

$$E[SS_{\text{trt}}] = \sum n_i \left[\left(\frac{\sigma^2}{n_i} - \frac{\sigma^2}{n} \right) + (\mu_i - \mu)^2 \right] = (g-1)\sigma^2 + \sum_i n_i (\mu_i - \mu)^2.$$

Then,

$$E[MS_{\text{trt}}] = \frac{(g-1)\sigma^2}{g-1} + \sum_i n_i (\mu_i - \mu)^2 = \sigma^2 + \frac{1}{g-1} \left(\sum_i (\mu_i - \mu)^2 \right).$$

Under H_0 , there holds $E[MS_{\text{trt}}] = \sigma^2$. Moreover,

$$\frac{\sum \sum (Y_{ij} - \bar{Y})^2}{\sigma^2} = \frac{\sum \sum (Y_{ij} - \bar{Y}_i)^2}{\sigma^2} + \frac{\sum n_i (\bar{Y}_i - \bar{Y})^2}{\sigma^2} \sim \chi_{n-1}^2.$$

6.0.3 Relation to Linear Models

We may rewrite the assumptions for one-way ANOVA as: $Y_{ij} = \mu_i + \epsilon_{ij}$ with $\epsilon_{ij} \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma^2)$. (Alternatively, this may be viewed as an assumption of the data generation process.)

Proposition 6.2. *The MLE of μ_i is \bar{y}_i .*

Another way to parameterize the model is

$$Y_{ij} = \mu + \alpha_i + \epsilon_{ij}$$

where $\epsilon_{ij} \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma^2)$ for $i = 1, \dots, g$ and $j = 1, \dots, n_i$. We impose constraint on the parameters of the form $\sum \alpha_i = 0$, or $\alpha_0 = 0$, or $\alpha_g = 0$.

These are constraints that we can freely impose without loss of generality; they may be viewed as a way to define μ . But we have always that $\hat{\mu}_i = \hat{\mu} + \hat{\alpha}_i$.

Example 6.3. Use $\partial \log L / \partial \alpha_i = 0$ for $i = 0, \dots, g - 1$ and use constraint $\sum \alpha_i = 0$ to derive MLE for α_i and μ . It turns out that $\hat{\mu} = \sum \bar{y}_i / g$.

7 Two Way ANOVA

This chapter
is very wrong

	SS	df		SS	df
trt _A		$I - 1$	trt _B		$J - 1$
Res		$n - I$	Res		$n - J$

Consider y_{ijk} , where $i = 1, \dots, I$ is treatment level of A, $j = 1, \dots, J$ is treatment level of B, and $k = 1, \dots, K_{ij}$ is the replicate. Consider first the case where $K_{ij} = K$. We denote the means as:

- $\bar{y}_{i..} := \frac{1}{IK} \sum_{j=1}^J \sum_{k=1}^K y_{ijk}$.
- $\bar{y}_{.j.} := \frac{1}{JK} \sum_{i=1}^I \sum_{k=1}^K y_{ijk}$.

Source	SS	df	MS	F
Treatment A	$SS_A = JK \sum_i (\bar{y}_{i..} - \bar{y})^2$	$I - 1$	$SS_A / (I - 1)$	MS_A / MS_e
Treatment B	$SS_B = IK \sum_j (\bar{y}_{.j.} - \bar{y})^2$	$J - 1$	$SS_B / (J - 1)$	MS_B / MS_e
A × B	SS_{AB}	$(I - 1)(J - 1)$		MS_{AB} / MS_e
Residue	$SS_e = \sum_{ijk} (y_{ijk} - \bar{y}_{ij.})^2$	df_e		
Total	$SS_{\text{total}} = \sum_{ijk} (y_{ijk} - \bar{y})^2$	$IK - 1$		

When there is no interactions, we have the partition

$$SS_{\text{total}} = \sum_{ijk} (y_{ijk} - \bar{y})^2 = SS_A + SS_B + SS_e.$$

7.1 Relation to Linear Models

We may write the model as

$$y_{ijk} \stackrel{\text{iid}}{\sim} \mathcal{N}(\mu_{ij}, \sigma^2)$$

or

$$y_{ijk} \stackrel{\text{iid}}{\sim} \mu_{ij} + \epsilon_{ijk}$$

with $\epsilon_{ijk} \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma^2)$ and $\mu_{ij} = \mu + \alpha_i + \beta_j + \gamma_{ij}$. That is,

$$y_{ijk} = \mu + \alpha_i + \beta_j + \gamma_{ij} + \epsilon_{ijk}, \quad \epsilon_{ijk} \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma^2).$$

A typical constraint is $\sum \alpha_i = \sum \beta_j = \sum_j \gamma_{1j} = \sum_i \gamma_{i1} = 0$.

Proposition 7.1. *We have the MLE*

$$\begin{aligned} \hat{\mu} &= \bar{y}, \quad \hat{\alpha}_i = \bar{y}_{i..} - \bar{y}, \quad \hat{\beta}_j = \bar{y}_{.j.} - \bar{y}, \\ \hat{\gamma}_{ij} &= \bar{y}_{ij.} - (\hat{\mu} - \hat{\alpha}_i - \hat{\beta}_j) = \bar{y}_{ij.} - \bar{y}_{i..} - \bar{y}_{.j.} + \bar{y}. \end{aligned}$$

7.2 Hypothesis testing

- $H_A : \alpha_i = 0, \forall i$. We have $\frac{MS_A}{MS_e} \sim F_{I-1, df_e}$.
- $H_B : \beta_i = 0, \forall j$. We have $\frac{MS_B}{MS_e} \sim F_{J-1, df_e}$.
- $H_{AB} : \gamma_{ij} = 0, \forall i, j$. We have $\frac{MS_{AB}}{MS_e} \sim F_{df_{AB}, df_e}$.

Example 7.2. Consider the case of no replicates.

$$Y_{ij} = \mu + \alpha_i + \beta_j + \epsilon_{ij}, \quad \epsilon_{ij} \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma^2).$$

$y_{ij} = \hat{\mu}_{ijk}$	B-1	B-2	B-3	B-4	row sum	row mean
A-1	3	2	1	2	8	2
A-2	8	7	8	5	28	7
A-3	7	6	9	2	24	6
col sum	18	15	18	9	60	
col mean	6	5	6	3		

We have then that $\hat{\mu} = 60/12 = 5$. We subtract the mean.

$\hat{\mu}_{ijk} - \hat{\mu}$	B-1	B-2	B-3	B-4	row sum	$\hat{\alpha}_i$
A-1	-2	-3	-4	-3	12	-3
A-2	3	2	3	0	8	2
A-3	2	1	4	-3	4	1
col sum	3	0	3	-6		
$\hat{\beta}_j$	1	0	1	-2		

We then subtract the row and column means to obtain the interactions (or residual if no interaction in model).

$\hat{\mu}_{ijk} - \hat{\mu} - \hat{\alpha}_i - \hat{\beta}_j$	B-1	B-2	B-3	B-4
A-1	0	0	-2	2
A-2	0	0	0	0
A-3	0	0	2	-2

We have the following ANOVA table:

Source	SS	df
A	SS_A	$I - 1$
B	SS_B	$J - 1$
Res	SS_E	$IJK - 1$

Definition 7.3 (Nested Models). Consider models

$$M_1 : \beta_0 + \beta_1 x_1 + \cdots + \beta_p x_p + \cdots + \beta_{p+q} x_{pq} + \epsilon$$

and

$$M_2 : \beta_0 + \beta_{p+1}x_{p+1} + \cdots + \beta_{p+q}x_{p+q} + \epsilon.$$

We say M_2 is nested in M_1 and write $M_2 \subset M_1$.

Proposition 7.4. *Let $M_{small} \subset M_{large}$ and denote the residual sum of squares SS_e as RSS . We have then that $RSS_{large} \leq RSS_{small}$. Moreover, we have*

$$RSS_{small} = RSS_{large} + \sum (fitted\ value_{small} - fitted\ value_{large})^2.$$

8 Comparison of Nested Models, Sequential ANOVA

Suppose

$$Y_{ijk} \stackrel{\text{iid}}{\sim} \mathcal{N}(\mu_{ij}, \sigma^2), \quad i = 1, \dots, I, \quad j = 1, \dots, J, \quad k = 1, \dots, K_{ij}.$$

Consider the model $\mu_{ij} = \mu + \alpha_i + \beta_j$ with constraints $\sum \alpha_i = \sum \beta_j = 0$.

Proposition 8.1. *We have the following MLE:*

- $\hat{\mu} = \bar{y}$
- $\hat{\alpha}_i = \bar{y}_{i.} - \bar{y}$
- $\hat{\beta}_j = \bar{y}_{.j} - \bar{y}$

If $K_{ij} \equiv K$, by writting

$$y_{ijk} - \bar{y} = y_{ijk} - \bar{y} - (\bar{y}_{i.} - \bar{y}) - (\bar{y}_{.j} - \bar{y}) + \bar{y} - \bar{y}_{i.} + \bar{y}_{.j} - \bar{y}$$

and recalling

$$(a + b + c)^2 = a^2 + b^2 + c^2 + 2ab + 2ac + 2bc,$$

we have

$$\begin{aligned} \sum_{ijk} (y_{ijk} - \bar{y})^2 &= \sum_{ijk} (y_{ijk} - \bar{y}_{i.} - \bar{y}_{.j} + \bar{y})^2 + \sum_{ijk} (\bar{y}_{i.} - \bar{y})^2 + \sum_{ijk} (\bar{y}_{.j} - \bar{y})^2 \\ &= SS_e + SS_A + SS_B, \end{aligned}$$

where the remaining three terms are 0 by the constraints and $K_{ij} \equiv K$.

8.1 Residual Sum of Squares

Example 8.2. Consider the two sample case with

$$H_0 : \mu_1 = \mu_2, \quad H_a : \mu_1 \neq \mu_2.$$

Consider the models

$$\begin{aligned} M_0 : X_i &= \mu_0 + \epsilon_i, \quad Y_j = \mu_0 + \epsilon_j, \quad \epsilon \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma^2). \\ M_1 : X_i &= \mu_1 + \epsilon_i, \quad Y_j = \mu_2 + \epsilon_j, \quad \epsilon \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma^2). \end{aligned}$$

Note that M_0 corresponds to the null hypothesis, and M_1 includes both the null and the alternative hypotheses.

We have

$$\hat{\mu}_0 = \frac{\sum X_i + \sum Y_j}{n + m}, \quad \hat{\mu}_1 = \bar{X}, \quad \hat{\mu}_2 = \bar{Y}.$$

The residual sum of squares are

$$RSS_0 = \sum (X_i - \hat{\mu}_0)^2 + \sum (Y_j - \hat{\mu}_0)^2, \quad RSS_1 = \sum (X_i - \hat{\mu}_1)^2 + \sum (Y_j - \hat{\mu}_2)^2.$$

We have always that $RSS_0 \geq RSS_1$.

We write

$$\sum (X_i - \hat{\mu}_0)^2 = \sum (X_i - \hat{\mu}_1 + \hat{\mu}_1 - \hat{\mu}_0)^2 = \sum (X_i - \hat{\mu}_1)^2 + \sum (\hat{\mu}_1 - \hat{\mu}_0)^2$$

and similarly for Y_j . Then,

$$\begin{aligned} RSS_0 &= RSS_1 + \sum (\hat{\mu}_1 - \hat{\mu}_0)^2 + \sum (\hat{\mu}_2 - \hat{\mu}_0)^2 \\ &= RSS_1 + \sum (\text{fitted value}_1 - \text{fitted value}_0)^2. \end{aligned}$$

Example 8.3. In one way ANOVA we have

$$\sum (y_{ij} - \bar{y})^2 = \sum (\bar{y}_i - \bar{y})^2 + \sum (y_{ij} - \bar{y}_i)^2.$$

That is,

$$SS_{\text{total}} = SS_{\text{treatment}} + SS_{\text{residual}}.$$

Consider

$$H_0 : \mu_1 = \cdots = \mu_g, \quad H_a : \mu_i \neq \mu_j \text{ for some } i \neq j.$$

The corresponding models are

$$M_0 : Y_{ij} = \mu_0 + \epsilon_{ij}, \quad \epsilon_{ij} \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma^2),$$

$$M_1 : Y_{ij} = \mu_i + \epsilon_{ij}, \quad \epsilon_{ij} \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma^2).$$

We have then that $\hat{\mu}_0 = \bar{y}$ and $\hat{\mu}_i = \bar{y}_i$.

$$SS_{\text{total}} = RSS_0, \quad SS_{\text{treatment}} = \sum (\text{fitted value}_1 - \text{fitted value}_0)^2, \quad SS_{\text{residual}} = RSS_1.$$

8.2 Sequential ANOVA

Consider models

$$M_0 : Y_{ijk} = \mu + \epsilon_{ijk}$$

$$M_1 : Y_{ijk} = \mu + \alpha_i + \epsilon_{ij}$$

$$M_1^* : Y_{ijk} = \mu + \beta_j + \epsilon_{ij}$$

$$M_2 : Y_{ijk} = \mu + \alpha_i + \beta_j + \epsilon_{ij} \quad M_3 : Y_{ijk} = \mu + \alpha_i + \beta_j + \gamma_{ij} + \epsilon_{ij}$$

Source	SS	df	Source	SS	df
A	$SS_A = \sum (\bar{y}_i - \bar{y})^2$		B	$SS_B = \sum (\bar{y}_j - \bar{y})^2$	
Residual	$SS_E = \sum (y - \bar{y}_i)$		Residual	SS_E	
Source	SS	df	Source	SS	df
A	SS_A		A	SS_A	
B	SS_B		B	SS_B	
Residual	SS_E		A × B	SS_{AB}	
			Residual	SS_E	

In general, consider a string of models

$$M_0 \subset M_1 \subset \cdots \subset M_n.$$

The parameter space gets larger:

$$\Theta_0 \subset \Theta_1 \subset \cdots \subset \Theta_n.$$

Sequential ANOVA SS columns: we are effectively partitioning the residual sum of squares as follows

$$\begin{aligned} RSS_0 &= RSS_0 - RSS_1 \\ &\quad + RSS_1 - RSS_2 \\ &\quad \vdots \\ &\quad + RSS_{n-1} - RSS_n \\ &= RSS_n. \end{aligned}$$

Remark 8.4. The effect in interested should be conducted in the last step.

We have then that for $H_1 : \alpha_o = i \equiv 0$,

$$\frac{(RSS_0 - RSS_1)/(df_0 - df_1)}{RSS_n/df_n} \sim F_{df_0 - df_1, df_n}$$

and so on.

8.3 Design Matrix

Consider the model

$$y_{ij} = \mu + \alpha_i + \epsilon_{ij}.$$

We may write this as

$$\begin{bmatrix} y_{11} \\ \vdots \\ y_{1J} \\ \vdots \\ y_{I1} \\ \vdots \\ y_{IJ} \end{bmatrix} = \begin{bmatrix} 1 & 1 & 0 & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & 0 & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & 0 & 0 & \cdots & 1 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & 0 & 0 & \cdots & 1 \end{bmatrix} \begin{bmatrix} \mu \\ \alpha_1 \\ \vdots \\ \alpha_I \end{bmatrix}$$

or

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$$

Note in particular that we have $\mathbf{Y} \in \mathbb{R}^{n \times 1}$, $\mathbf{X} \in \mathbb{R}^{n \times (g+1)}$, $\boldsymbol{\beta} \in \mathbb{R}^{(g+1) \times 1}$, and $\boldsymbol{\epsilon} \in \mathbb{R}^{n \times 1}$. Note that the rank of X is g , since one column can be written as a linear combination of the others.

To make the design matrix full rank, we need parameter constraints such as $\sum \alpha_i = 0$. Under this constraint, we have $\alpha_I = -\sum_{i=1}^{I-1} \alpha_i$. We may then write the design matrix as

$$\begin{bmatrix} y_{11} \\ \vdots \\ y_{1J} \\ \vdots \\ y_{I1} \\ \vdots \\ y_{IJ} \end{bmatrix} = \begin{bmatrix} 1 & 1 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & -1 & -1 & \dots & 1 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & -1 & -1 & \dots & 1 \end{bmatrix} \begin{bmatrix} \mu \\ \alpha_1 \\ \vdots \\ \alpha_{I-1} \end{bmatrix}.$$

The design matrix is now full rank.

9 Simple Linear Regression

Recall that with (x_i, y_i) bivariate normal we have the following MLE:

$$\hat{\mu}_x = \bar{x}, \quad \hat{\mu}_y = \bar{y}, \quad \hat{\sigma}_x^2 = \frac{1}{n} \sum (x_i - \bar{x})^2, \quad \hat{\sigma}_y^2 = \frac{1}{n} \sum (y_i - \bar{y})^2$$

and

$$\hat{\rho} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2 \sum (y_i - \bar{y})^2}}.$$

We have also

$$Y|X = x \sim \mathcal{N}\left(\mu_y + \frac{\sigma_y}{\sigma_x} \rho (x - \mu_x), \sigma_y^2 (1 - \rho^2)\right).$$

We consider now an alternative model which focuses on the distribution of one variable conditioning on a specific value of the other variables. Suppose we have data $(x_1, y_1), \dots, (x_n, y_n)$. Writing

$$E[Y|X = x] = \beta_0 + \beta_1 x,$$

we have comparing to the previous results for bivariate normal variables that

$$\beta_0 = \mu_y - \rho \frac{\sigma_y}{\sigma_x} \mu_x, \quad \beta_1 = \rho \frac{\sigma_y}{\sigma_x}.$$

Consider thus the model

$$Y_i = \beta_0 + \beta_1 x_i + \epsilon_i, \quad \epsilon_i \stackrel{\text{iid}}{\sim} (0, \sigma^2).$$

Definition 9.1 (Least Square Estimator). The least square estimators of $\{\beta_0, \beta_1\}$ are the values that minimize

$$\sum (y_i - (\beta_0 + \beta_1 x_i))^2 := \text{RSS}.$$

That is,

$$\{\hat{\beta}_0, \hat{\beta}_1\} := \arg \min \text{RSS}.$$

Note that

$$\frac{\partial \text{RSS}}{\partial \beta_0} = -2 \sum (y_i - \beta_0 - \beta_1 x_i), \quad \frac{\partial \text{RSS}}{\partial \beta_1} = -2 \sum x_i (y_i - \beta_0 - \beta_1 x_i).$$

Setting the above to 0 gives

$$n^{-1} \sum y_i - \beta_0 - \beta_1 n^{-1} \sum x_i = 0, \quad n^{-1} \sum x_i y_i - n^{-1} \sum x_i \beta_0 - n^{-1} \sum \beta_1 x_i^2 = 0$$

and thus

$$\hat{\beta}_1 = \frac{\frac{1}{n} \sum x_i y_i - (\frac{1}{n} \sum x_i)(\frac{1}{n} \sum y_i)}{\frac{1}{n} \sum x_i^2 - (\frac{1}{n} \sum x_i)^2} = \frac{\text{sample covariance of } x_i \text{ and } y_i}{\text{sample variance of } x_i}.$$

9.1 Properties of least square estimators

Note that $\hat{\beta}_0$ and $\hat{\beta}_1$ are random. Think of the x_i 's the formula for $\hat{\beta}_1$ as fixed, and only y_i 's are random:

$$\hat{\beta}_1 = \frac{\frac{1}{n} \sum x_i y_i - (\frac{1}{n} \sum x_i)(\frac{1}{n} \sum y_i)}{\frac{1}{n} \sum x_i^2 - (\frac{1}{n} \sum x_i)^2}.$$

We may rewrite the above as

$$\hat{\beta}_1 = \frac{\sum x_i y_i - \bar{x} \sum y_i}{\sum x_i^2 - \bar{x} \sum x_i} = \frac{\sum (x_i - \bar{x}) y_i}{\sum (x_i - \bar{x}) x_i} = \frac{\sum (x_i - \bar{x}) y_i}{\sum (x_i - \bar{x})^2} = \sum c_i y_i,$$

a linear combination of the y_i 's.

Alternatively, we can write

$$\hat{\beta}_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2},$$

the familiar covariance over variance form.

From the linear combination form, we have

$$E[\hat{\beta}_1] = E\left[\sum c_i Y_i\right] = \sum c_i E[Y_i] = \sum c_i (\beta_0 + \beta_1 x) = \left(\sum c_i\right) \beta_0 + \left(\sum c_i x_i\right) \beta_1.$$

Recalling that

$$c_i = \frac{x_i - \bar{x}}{\sum (x_i - \bar{x})^2},$$

we see that

$$\sum c_i = 0, \quad \sum c_i x_i = \sum \frac{(x_i - \bar{x})(x_i - \bar{x})}{\sum (x_i - \bar{x})^2} = 1.$$

Therefore $E[\hat{\beta}_1] = \beta_1$. Moreover, we have

$$E[\hat{\beta}_0] = E[\bar{Y} - \hat{\beta}_1 \bar{x}] = E[\bar{Y}] - \bar{x} E[\hat{\beta}_1] = (\beta_0 + \beta_1 \bar{x}) - \bar{x} \beta_1 = \beta_0.$$

That is, we have:

Proposition 9.2. *The least square estimators are unbiased.*

We consider next the variance. We have

$$\begin{aligned} \text{Var}[\hat{\beta}_1] &= \text{Var}\left[\sum c_i Y_i\right] = \sum c_i^2 \text{Var}[Y_i] \\ &= \left(\sum c_i^2\right) \sigma^2 = \sigma^2 \frac{1}{\sum (x_i - \bar{x})^2} = \frac{\sigma^2}{n} \frac{1}{n^{-1} \sum (x_i - \bar{x})^2}, \end{aligned}$$

where we used the fact that

$$\sum c_i^2 = \sum \frac{(x_i - \bar{x})^2}{(\sum (x_i - \bar{x})^2)^2} = \frac{1}{\sum (x_i - \bar{x})^2}.$$

Remark 9.3. The variance is large if either the noise σ^2 is large or x_i 's are concentrated ($n^{-1} \sum (x_i - \bar{x})^2$ is small).

For $\hat{\beta}_0$, we can write

$$\begin{aligned}\text{Var}[\hat{\beta}_0] &= \text{Var}[\bar{Y} - \hat{\beta}_1 \bar{x}] = \text{Var}[\bar{Y}] + \bar{x}^2 \text{Var}[\hat{\beta}_1] - 2\bar{x} \text{Cov}(\bar{Y}, \hat{\beta}_1) \\ &= \frac{\sigma^2}{n} + \frac{\bar{x}^2 \sigma^2}{\sum (x_i - \bar{x})^2},\end{aligned}$$

where the last inequality comes from

$$\text{Cov}(\bar{y}, \hat{\beta}_1) = \frac{1}{n} \text{Cov}\left(\sum Y_i, \sum c_i Y_i\right) = \frac{1}{n} \sum c_i \text{Cov}(Y_i, Y_i) = 0.$$

$$\text{Cov}(\hat{\beta}_0, \hat{\beta}_1) = \text{Cov}(\bar{Y} - \hat{\beta}_1 \bar{x}, \hat{\beta}_1) = \text{Cov}(\bar{Y}, \hat{\beta}_1) - \bar{x} \text{Var}(\hat{\beta}_1) = -\frac{\bar{x} \sigma^2}{\sum (x_i - \bar{x})^2}$$

9.2 Distributional results

We now consider the model

$$\begin{aligned}Y_i &= \beta_0 + \beta_1 x_i + \epsilon_i, \quad \epsilon_i \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma^2), \\ \iff Y_i &= \beta_0 + \beta_1 x_i + \sigma^2 Z_i, \quad Z_i \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1), \\ \iff Y_i &\text{ independent, } Y_i \sim \mathcal{N}(\beta_0 + \beta_1 x_i, \sigma^2).\end{aligned}$$

Proposition 9.4. *Under the assumption above, the MLE of β_0 and β_1 are the least square estimators.*

With the added bivariate normality assumption, we have $[\hat{\beta}_0 \hat{\beta}_1]^\top$ is bivariate normal as a linear combination of the data. Specifically, with the results above, we have

$$\begin{bmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \end{bmatrix} \sim \mathcal{N}\left(\begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix}, \sigma^2 \begin{bmatrix} \frac{1}{n} + \frac{\bar{x}^2}{\sum (x_i - \bar{x})^2} & -\frac{\bar{x}}{\sum (x_i - \bar{x})^2} \\ -\frac{\bar{x}}{\sum (x_i - \bar{x})^2} & \frac{1}{\sum (x_i - \bar{x})^2} \end{bmatrix}\right)$$

9.2.1 Prediction

Consider a new observation of x , say x_0 . We have

$$\mathbb{E}[Y_0] = \beta_0 + \beta_1 x_0, \quad \hat{y}_0 := \mathbb{E}[\hat{Y}_0] = \hat{\beta}_0 + \hat{\beta}_1 x_0$$

and

$$\begin{aligned}\text{Var}(\hat{\beta}_0 + \hat{\beta}_1 x_0) &= \frac{\sigma^2}{n} + \frac{\sigma^2 \bar{x}^2}{\sum (x_i - \bar{x})^2} + \frac{\sigma^2 x_0^2}{\sum (x_i - \bar{x})^2} + \frac{-2x_0 \bar{x} \sigma^2}{\sum (x_i - \bar{x})^2} \\ &= \frac{\sigma^2}{n} + \sigma^2 \frac{(x_0 - \bar{x})^2}{\sum (x_i - \bar{x})^2}.\end{aligned}$$

We have

$$\frac{\hat{\beta}_0 + \hat{\beta}_1 x_0 - (\beta_0 + \beta_1 x_0)}{\sigma \sqrt{\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{\sum (x_i - \bar{x})^2}}} \sim \mathcal{N}(0, 1).$$

We may then construct a CI at $\hat{y}_0 = \hat{\beta}_0 + \hat{\beta}_1 x_0$:

$$(\hat{\beta}_0 + \hat{\beta}_1 x_0) \pm z_{1-\frac{\alpha}{2}} \sqrt{\frac{\sigma^2}{n} + \sigma^2 \frac{(x_0 - \bar{x})^2}{\sum (x_i - \bar{x})^2}},$$

We use the estimator

$$\hat{\sigma}^2 := \frac{\text{RSS}}{n-2}.$$

Remark 9.5. The CI is wider when x_0 is far from the center \bar{x} .

It turns out that

$$\frac{\text{RSS}}{\hat{\sigma}^2} \sim \chi_{n-2}^2,$$

where the degree of freedom is $n-2$ because we are estimating two parameters, and

$$\frac{\hat{\beta}_0 + \hat{\beta}_1 x_0 - (\beta_0 + \beta_1 x_0)}{\sigma \sqrt{\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{\sum (x_i - \bar{x})^2}}} \cdot \frac{1}{\sqrt{\frac{\text{RSS}}{\sigma^2} / (n-2)}} \sim t_{n-2},$$

since the first term is $\mathcal{N}(0, 1)$ and the second is $1/\sqrt{\chi_{n-2}^2/(n-2)}$. We then have the CI at $(x_0, \hat{E}(Y_0)) = (x_0, \hat{\beta}_0 + \hat{\beta}_1 x_0)$ for $\beta_0 + \beta_1 x_0$ is

$$\hat{\beta}_0 + \hat{\beta}_1 x_0 \pm t_{n-2, 1-\frac{\alpha}{2}} \sqrt{\frac{\text{RSS}}{n-2} \left(\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{\sum (x_i - \bar{x})^2} \right)}$$

Consider now

$$H_0 : \beta_1 = 0, \quad H_a : \beta_1 \neq 0.$$

Consider

$$Y_i = \beta_0 + \beta_1 x_i + \epsilon_i, \quad \epsilon_i \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1).$$

We have the LSE/MLE for $\{b_0, b_1\}$

$$\begin{cases} \hat{\beta}_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} \\ \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \end{cases}$$

and the MLE for σ^2

$$\hat{\sigma}_{ML}^2 = \frac{\text{RSS}}{\sigma^2}.$$

Recall that we can write $\hat{\beta}_0$ and $\hat{\beta}_1$ as linear combinations of Y_i . They are thus jointly normal.

Note that the regression line always pass through (\bar{x}, \bar{y}) .

10 Multiple Linear Regression Model

Consider

$$(x_{i1}, x_{i2}, \dots, x_{ir}, y_i), \quad i = 1, \dots, n.$$

Model assumption:

$$Y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_r x_{ir} + \epsilon_i, \quad \epsilon_i \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma^2)$$

or equivalently

$$Y_i \sim \mathcal{N}(\beta_0 + \beta_1 x_{i1} + \dots + \beta_r x_{ir}, \sigma^2).$$

Think of x_{i1}, \dots, x_{ir} as fixed.

We have

$$\begin{aligned} L(\beta, \sigma^2; x, y) &= \prod_{i=1}^n f(y_i; \beta, x, \sigma^2) \\ &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y_i - \beta_0 - \beta_1 x_{i1} - \dots - \beta_r x_{ir})^2}{2\sigma^2}\right) \\ &= (2\pi\sigma^2)^{-n/2} \exp\left(-\frac{1}{2\sigma^2} \sum (y_i - (\beta_0 + \beta_1 x_{i1} + \dots + \beta_r x_{ir}))^2\right). \end{aligned}$$

We write

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}.$$

Note that $\mathbf{y}, \boldsymbol{\epsilon} \in \mathbb{R}^{n \times 1}$, $\mathbf{X} \in \mathbb{R}^{n \times p}$ where $p = r + 1$, $\boldsymbol{\beta} \in \mathbb{R}^{p \times 1}$. Using this notation we can write

$$\sum (y_i - (\beta_0 + \beta_1 x_{i1} + \dots + \beta_r x_{ir}))^2 = \boldsymbol{\epsilon}^\top \boldsymbol{\epsilon} = \langle \boldsymbol{\epsilon}, \boldsymbol{\epsilon} \rangle = \|\boldsymbol{\epsilon}\|^2 = \|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}\|^2.$$

The log likelihood is then

$$\ell = -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{\sigma^2} \|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}\|^2.$$

We see thus that the MLE is attained when $\|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}\|$ is minimized. The MLE coincides with the least square estimator. Specifically, we have

$$\frac{\partial \ell}{\partial \boldsymbol{\beta}} = -\frac{1}{2\sigma^2} \frac{\partial}{\partial \boldsymbol{\beta}} \|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}\|^2.$$

Noting that

$$\|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}\|^2 = (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}) = \mathbf{Y}^\top \mathbf{Y} - 2\boldsymbol{\beta}^\top \mathbf{X}^\top \mathbf{Y} + \boldsymbol{\beta}^\top \mathbf{X}^\top \mathbf{X} \boldsymbol{\beta},$$

we have

$$\frac{\partial \ell}{\partial \boldsymbol{\beta}} = -\frac{1}{2\sigma^2} (-2\mathbf{X}^\top \mathbf{Y} + 2\mathbf{X}^\top \mathbf{X} \boldsymbol{\beta}) = -\frac{1}{\sigma^2} (\mathbf{X}^\top \mathbf{Y} - \mathbf{X}^\top \mathbf{X} \boldsymbol{\beta}).$$

Setting it to 0 gives $\mathbf{X}^\top \mathbf{X} \boldsymbol{\beta} = \mathbf{X}^\top \mathbf{Y}$. Thus if \mathbf{X} has rank $p \leq n$, we have

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

That is, we have

Proposition 10.1. *The MLE and the least square estimator coincide, and are given by*

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

Proof. Note that

$$\|Y - X\beta\|^2 = (Y - X\beta)^\top (Y - X\beta) = Y^\top Y - 2\beta^\top X^\top Y + \beta^\top X^\top X\beta,$$

Taking the derivative with respect to β gives

$$\frac{\partial \|Y - X\beta\|^2}{\partial \beta} = -2X^\top Y + 2X^\top X\beta,$$

where the second term comes from noting that

$$\begin{aligned} \frac{\partial}{\partial x_k} (x^\top Ax) &= \frac{\partial}{\partial x_k} \left(\sum_i a_{ii} x_i^2 + \sum_{i \neq j} a_{ij} x_i x_j \right) \\ &= 2a_{kk} x_k + \sum_{i \neq k} a_{ik} x_i + \sum_{j \neq k} a_{kj} x_j \\ &= \sum_i a_{ik} x_i + \sum_j a_{kj} x_j \\ &= k\text{th entry of } Ax + k\text{th entry of } A^\top x. \end{aligned}$$

Thus $\partial(x^\top Ax)/\partial x_k = Ax + A^\top x$ and when A is symmetric, we have $\partial(x^\top Ax)/\partial x = 2Ax$. \square

Then

$$E[\hat{\beta}] = (X^\top X)^{-1} X^\top E[Y] = (X^\top X)^{-1} X^\top (X\beta) = \beta.$$

The least square estimator is unbiased.

The variance is

$$\text{Var}[\hat{\beta}] = \text{Var}[(X^\top X)^{-1} X^\top Y] = (X^\top X)^{-1} X^\top \text{Var}[Y] X (X^\top X)^{-1},$$

where we recall

$$\text{Cov}(AY, BY) = A \text{Cov}(Y, Y) B^\top.$$

Since X and β are fixed, we have

$$\text{Var}[Y] = \text{Var}[\epsilon] = \sigma^2 I_n.$$

Then,

$$\text{Var}[\hat{\beta}] = \sigma^2 (X^\top X)^{-1}.$$

So $\hat{\beta}$ is normally distributed with

$$\hat{\beta} \sim \mathcal{N}(\beta, \sigma^2 (X^\top X)^{-1}).$$

The RSS is

$$\text{RSS} = \|Y - X\hat{\beta}\|^2$$

with

$$\frac{\text{RSS}}{\sigma^2} \sim \chi_{n-p}^2,$$

where, recall, $p = r - 1$.

We have the following

Proposition 10.2.

- $\hat{\beta} \sim \mathcal{N}(\beta, \sigma^2(\mathbf{X}^\top \mathbf{X})^{-1})$.
- $\text{RSS}/\sigma^2 \sim \chi_{n-p}^2$.
- $\hat{\beta} \perp \text{RSS}$.

10.1 Projection Matrix

We write

$$X = \begin{bmatrix} 1 & X_1 & \dots & X_{p-1} \end{bmatrix}.$$

Then,

$$X\beta = \beta_0 1 + \beta_1 X_1 + \dots + \beta_{p-1} X_{p-1}$$

is a linear combination of the columns of X . In particular, $X\beta$ is in the column space of X , and the least squares condition dictates that it is the orthogonal projection of Y onto the column space of X .

$$X\hat{\beta} = X(X^\top X)^{-1}X^\top Y = P_X Y =: \hat{Y},$$

where P_X or H denotes the projection matrix.

Definition 10.3. P is a **projection matrix** if $P = P^\top$ and $P = P^2$.

Example 10.4.

- I_n , the identity matrix.
- $H = P_X := X(X^\top X)^{-1}X^\top$ is a projection matrix onto the column space of X .
- Let v be a vector with norm $\|v\| = v^\top v = 1$. Then $P = vv^\top$ is a projection matrix onto the line spanned by v (a “rank-1” matrix). This is a special case of the previous example.
- Let $P \in \mathbb{R}^{n \times n}$ be projection matrix. Then $I_n - P$ is also a projection matrix.

Proposition 10.5. *The only two possible eigenvalues of a projection matrix P is 1 and 0.*

Proof. Let v_1, \dots, v_n be a set of orthonormal eigenvectors of P with $Pv_i = \lambda_i v_i$ for $i = 1, \dots, n$ and $\lambda_1 \geq \dots \geq \lambda_n$. Note that since P is symmetric, v_i 's exist and $\lambda_i \in \mathbb{R}$ are real. We have then that

$$P \begin{bmatrix} v_1 & \dots & v_n \end{bmatrix} = \begin{bmatrix} \lambda_1 v_1 & \dots & \lambda_n v_n \end{bmatrix} = \begin{bmatrix} v_1 & \dots & v_n \end{bmatrix} \begin{bmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & \lambda_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \lambda_n \end{bmatrix}.$$

Denote the last two terms by V and Λ . Since v_1, \dots, v_n is orthonormal, we have $V^\top V = I_n$. Then from

$$PV = V\Lambda$$

we have by multiplying each side by V^\top that

$$P = V\Lambda V^\top.$$

Since P is a orthogonal projection, we have

$$P^2 = V\Lambda V^\top V\Lambda V^\top = V\Lambda^2 V^\top = P = V\Lambda V^\top.$$

In particular,

$$V\Lambda^2 V^\top = V\Lambda V^\top.$$

Multiplying both sides by V^\top on the left and V on the right gives

$$\Lambda^2 = \Lambda,$$

which gives the desired result. □

Note that we obtained $P = V\Lambda V^\top$, which gives

Corollary 10.6. *The rank of P is equal to the number of eigenvalues equal to 1. Thus, $\text{rank } P = \sum \lambda_i$.*

Corollary 10.7. $\text{tr } P = \text{rank } P$

Proof.

$$\text{tr } P = \text{tr}(V\Lambda V^\top) = \text{tr}(V^\top V\Lambda) = \text{tr } \Lambda = \text{rank } P.$$

□

Note that we used the following:

Proposition 10.8.

- $\text{tr}(AB) = \text{tr}(BA)$,
- $\text{tr}(A + B) = \text{tr } A + \text{tr } B$.

An application:

Proposition 10.9. Let P be a projection matrix and $Z \in \mathcal{N}(0, I_n)$. We have

$$\|PZ\|^2 \sim \chi_k^2,$$

where $k := \text{rank } P$.

Proof. Let V and Λ be as in the previous proof (note that λ_i 's are in decreasing order). Then

$$\begin{aligned} \|PZ\|^2 &= (PZ)^\top (PZ) = (V\Lambda V^\top Z)^\top (V\Lambda V^\top Z) \\ &= Z^\top V\Lambda^2 V^\top Z \\ &= Z^\top V\Lambda V^\top Z = Y^\top \Lambda Y, \end{aligned}$$

where we define $Y := Z^\top V$.

Note that Y is a linear combination of Z 's and thus of multi-normal distribution. In particular we have

$$E[Y] = V^\top E[Z] = V^\top 0 = 0$$

and

$$\text{Var}[Y] = \text{Cov}(V^\top Z) = V^\top \text{Var}[Z]V = V^\top I_n V = I_n.$$

Then,

$$\|PZ\|^2 = Y^\top \Lambda Y = Y_1^2 + \dots + Y_k^2 \sim \chi_k^2.$$

□

Note that we used the following:

Proposition 10.10. Let A and B be matrices and X and Y be random vectors. Then, $\text{Cov}(AX, BY) = A \text{Cov}(X, Y) B^\top$.

10.2 Statistical Properties of the Least Squares Estimators

Recall that for $Y_i \stackrel{\text{iid}}{\sim} \mathcal{N}(\mu, \sigma^2)$ and the hypothesis $H_0 : \mu = \mu_0$, we have under null that

- $\frac{\bar{Y} - \mu_0}{\sqrt{\sigma^2/n}} \sim \mathcal{N}(0, 1),$
- $\frac{\sum (Y_i - \bar{Y})^2}{\sigma} \sim \chi_{n-1}^2,$
- $S^2 \perp \hat{\mu},$

which gives

$$\frac{\bar{Y} - \mu_0}{\sqrt{S^2/n}} \sim t_{n-1},$$

using which we can construct a CI for μ . This is a simple case of the following multiple regression model:

$$Y_i = \beta_0 + \beta_1 x_{i1} + \cdots + \beta_r x_{i(p-1)} + \epsilon_i, \quad \epsilon_i \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma^2),$$

or

$$Y = X\beta + \epsilon, \quad \epsilon \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma^2 I_n).$$

Recall that we obtained the MLE/LSE of β as $\hat{\beta} = (X^\top X)^{-1} X^\top Y$. Thus,

$$\hat{\beta} \sim \mathcal{N}(\beta, \sigma^2 (X^\top X)^{-1}).$$

Definition 10.11. The standard error is defined by

$$\begin{aligned} \text{se } \hat{\beta}_k &= \sqrt{\text{Var}(\hat{\beta}_k)} \\ &= \text{the } (k+1)\text{th diagonal entry of the covariance matrix } \sigma^2 (X^\top X)^{-1} \end{aligned}$$

We will prove the following:

Proposition 10.12. *We have the following:*

- $\frac{\text{RSS}}{\sigma^2} = \frac{\|Y - X\hat{\beta}\|^2}{\sigma^2} = \chi_{n-p}^2$.
- $\text{RSS} \perp \hat{\beta}$.

This similarly will enable us to do all kinds of t tests.

Proof.

- We have

$$\text{RSS} = \|Y - X\hat{\beta}\|^2 = \|Y - HY\|^2 = \|(I - H)Y\|^2.$$

Note that

$$\begin{aligned} (I - H)Y &= (I - H)(X\beta + \epsilon) \\ &= X\beta - H\beta + (I - H)\epsilon = (I - H)\epsilon. \end{aligned}$$

Plugging back in, we obtain

$$\text{RSS} = \sigma^2 \|(I - H)\epsilon\|^2,$$

where, note, $I - H$ is a projection matrix. Since $\epsilon =: \sigma Z \sim \mathcal{N}(0, \sigma^2 I)$. We have

$$\text{rank}(I - H) = \text{tr}(I - H) = \text{tr}(I) - \text{tr}(H) = n - p,$$

where we used $\text{tr}(P) = \text{tr}(V\Lambda V^{-1}) = \text{tr}(\Lambda V^\top V) = \text{tr}(\Lambda) = p$. Using the result from the previous subsection, we have

$$\frac{\text{RSS}}{\sigma^2} = \frac{\|(I - H)Z\|^2}{\sigma^2} \sim \chi_{n-p}^2.$$

- In light of normality, we need only check

$$\text{Cov}(Y - X\hat{\beta}, \hat{\beta}) = 0.$$

We have

$$\begin{aligned} \text{Cov}(Y - X\hat{\beta}, \hat{\beta}) &= \text{Cov}(Y - X(X^\top X)^{-1}X^\top Y, (X^\top X)^{-1}X^\top Y) \\ &= \text{Cov}\left([I - X(X^\top X)^{-1}X^\top]Y, (X^\top X)^{-1}X^\top Y\right) \\ &= [I - X(X^\top X)^{-1}X^\top] \text{Cov}(Y, Y) [(X^\top X)^{-1}X^\top]^\top \\ &= [I - X(X^\top X)^{-1}X^\top] \left(\sigma^2 I_n\right) X(X^\top X)^{-1} \\ &= \sigma^2 [I - X(X^\top X)^{-1}X^\top] X(X^\top X)^{-1} = 0. \end{aligned}$$

□

Note that from the first result we can obtain an unbiased estimator for σ^2 :

$$E\left(\frac{\text{RSS}}{n-p}\right) = \sigma^2.$$

We summarize the results we obtained so far as follows:

Proposition 10.13. *Consider the linear model*

$$Y = \mathcal{N}_n(X\beta, \sigma^2 I_n).$$

We have

- $$\hat{\beta} = (X^\top X)^{-1}X^\top Y \sim \mathcal{N}_p(\beta, \sigma^2(X^\top X)^{-1}).$$

- $$\frac{\text{RSS}}{\sigma^2} = \frac{\|Y - X\hat{\beta}\|^2}{\sigma^2} \sim \chi_{n-p}^2.$$

- $$\hat{\beta} \perp \text{RSS} \quad (\text{thus } \text{MSS} \perp \text{RSS}).$$

- Under the null hypothesis that $\beta_k = 0$, we have

$$\frac{\hat{\beta}_k - 0}{\text{se } \hat{\beta}_k} \sim t_{n-p},$$

where $\text{se } \hat{\beta}_k = \sqrt{\text{Var}(\hat{\beta})}$ and $\text{Var}[\hat{\beta}_k]$ is the i th diagonal entry of $\text{Cov}(\hat{\beta}) = \sigma^2(X^\top X)^{-1}$.

10.3 Mode Comparison

Consider the following two models:

$$\begin{aligned} M_0 : \quad Y_i &= \beta_0 + \epsilon_i, & \epsilon_i &\stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma_0^2) \\ M_1 : \quad Y_i &= \beta_0 + \beta_1 x_{i,1} + \cdots + \beta_{p-1} x_{i,p-1} + \epsilon_i, & \epsilon_i &\stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma^2). \end{aligned}$$

We can rewrite the models as

$$\begin{aligned} M_0 : \quad Y_i &\sim \mathcal{N}(\beta_0 \mathbb{1}_n, \sigma_0^2 I_n), \\ M_1 : \quad Y_i &\sim \mathcal{N}(X\beta, \sigma^2 I_n), \end{aligned}$$

where $\mathbb{1}_n$ is a vector of 1's. The model estimates are

$$\begin{aligned} M_0 : \quad \hat{Y}_i &= \bar{Y} = J_n Y, \\ M_1 : \quad \hat{Y}_i &= X\hat{\beta} = P_X Y, \end{aligned}$$

where $J_n := n^{-1} \mathbb{1}_n \mathbb{1}_n^\top$ is easily checked to be a projection matrix of rank 1 (it's column space is just a line). We have for both models that

Proposition 10.14.

$$\|Y - \bar{Y}\|^2 = \|Y - \hat{Y}\|^2 + \|\hat{Y} - \bar{Y}\|^2.$$

Or,

$$SS_{total} = SS_e + SS_{model}.$$

Since the column space of $\mathbb{1}_n$ is clearly a subspace of the column space of X , we have

$$\|Y - \hat{Y}_0\|^2 \leq \|Y - \bar{Y}_1\|^2.$$

We have moreover that

Proposition 10.15.

$$\|Y - J_n Y\|^2 = \|Y - HY\|^2 + \|HY - J_n Y\|^2.$$

Proof. We have

$$\begin{aligned} \|Y - J_n Y\|^2 &= \|Y - HY + HY - J_n Y\|^2 \\ &= \|Y - HY\|^2 + \|HY - J_n Y\|^2 + 2(Y - HY)^\top (HY - J_n Y). \end{aligned}$$

Since

$$\begin{aligned} (Y - HY)^\top (HY - J_n Y) &= Y^\top (I_n - H)(H - J_n)Y \\ &= Y^\top (H - J_n - H^2 + HJ_n)Y. \end{aligned}$$

Since the column space of $\mathbb{1}_n$ is a subspace of the column space of X , we know that $HJ_n = J_n$, thus completing the proof. Alternatively, from $HX = X$ and the fact that the first column of X is $\mathbb{1}_n$, we know that $H\mathbb{1}_n = \mathbb{1}_n$, from which we have similarly that $HJ_n = J_n$. \square

Remark 10.16. Remember the picture!

We summarize the results as follows:

Proposition 10.17. *We have*

- $SS_{total} = SS_e + SS_{model}$.
- $SS_e/\sigma^2 = RSS/\sigma^2 \sim \chi_{n-p}^2$.
- $SS_e \perp\!\!\!\perp SS_{model}$.
- *Under the null hypothesis that M_0 is true, we have*

$$\frac{SS_{total}}{\sigma^2} \sim \chi_{n-p}^2, \quad \frac{SS_{model}}{\sigma^2} \sim \chi_{p-1}^2.$$

Thus,

$$\frac{(SS_{total} - SS_e)/[(n-1) - (n-p)]}{SS_e/(n-p)} = \frac{SS_{model}/(p-1)}{SS_e/(n-p)} \sim F_{p-1, n-p}.$$

More generally,

$$\frac{(SS_{small} - SS_{large})/[(n-p_s) - (n-p_l)]}{SS_{large}/(n-p_l)} \sim F_{p_s-p_l, n-p_l}.$$

10.4 LR and ANOVA

Example 10.18. Consider the case $Y_i = \beta_1 x_i + \epsilon_i$, $i = 1, \dots, n$. We have $X^\top X = \|X\|$ and so $P_X = X(X^\top X)^{-1}X^\top = XX^\top/\|X\|$.

Example 10.19. Consider the case $Y_i = \beta_1 x_{i1} + \beta_2 x_{i2} + \epsilon_i$, $i = 1, \dots, n$. We have $X^\top X = \|X\|$ and so $P_X = X(X^\top X)^{-1}X^\top = XX^\top/\|X\|$. Thus

$$X^\top X = \begin{bmatrix} x_1^\top x_1 & x_1^\top x_2 \\ x_2^\top x_1 & x_2^\top x_2 \end{bmatrix}.$$

If in addition $x_1^\top x_2 = 0$, then

$$x^\top x = \begin{bmatrix} \|x_1\|^2 & 0 \\ 0 & \|x_2\|^2 \end{bmatrix}$$

and thus

$$P_X = X(X^\top X)^{-1}X^\top = \begin{pmatrix} x_1 & x_2 \end{pmatrix} \begin{bmatrix} \frac{1}{\|x_1\|^2} & 0 \\ 0 & \frac{1}{\|x_2\|^2} \end{bmatrix} \begin{bmatrix} x_1^\top \\ x_2^\top \end{bmatrix} = \begin{bmatrix} \frac{x_1 x_1^\top}{\|x_1\|^2} & 0 \\ 0 & \frac{x_2 x_2^\top}{\|x_2\|^2} \end{bmatrix}$$

This shows that $P_{x_1} \perp P_{x_2}$, which gives $P_X = P_{x_1} + P_{x_2}$. We have

$$P_X Y = P_{x_1} Y + P_{x_2} Y.$$

From $Y = P_X Y + (I - P_X)Y$ we have

$$\|Y\|^2 = \|P_{x_1} Y\|^2 + \|P_{x_2} Y\|^2 + \|(I - P_X)Y\|^2,$$

the two way ANOVA partition. We see that the order does not matter for sequential ANOVA if $x_1^\top x_2 = 0$.

Example 10.20. Consider

$$Y_{ijk} = \mu + \alpha_i + \beta_j + \epsilon_{ijk}.$$

We may write it as

$$Y_i = \mu + \beta_1 x_{1i} + \beta_2 x_{2i} + \epsilon_i,$$

where x_{1i} is a dummy for receiving treatment A at level i , and x_{2i} a dummy for B at level i . Consider the usual constraints $\sum \alpha_i = \sum \beta_j = 0$. The usual adjustment of the design matrix is needed to obtain a full rank matrix. Writing $X = [X_0 \ X_A \ X_B \ X_{AB}]$ we see that $X_A \perp X_B$ if the design is balanced.

10.5 Q-Q Plot

$$RSS = SS_e = \sum \hat{\epsilon}_i^2 = \|Y - X\hat{\beta}\|^2.$$

$\hat{\epsilon}_i$'s should be close to iid $\mathcal{N}(0, \hat{\sigma}^2)$. To test this hypothesis, we may use a Q-Q plot (quantile-quantile plot).

Suppose we have data X_1, \dots, X_n (e.g., the residuals of a linear model). Order the X_i 's as $X_{(1)} \leq \dots \leq X_{(n)}$. Are X_i 's approximately iid of some normal distribution $\mathcal{N}(\mu, \sigma^2)$. Pick $Z_i \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1)$ and order them as $Z_{(1)} \leq \dots \leq Z_{(n)}$. We plot $(Z_{(i)}, X_{(i)})$ for $i = 1, \dots, n$.

Claim 10.21. If the X_i 's are approximately iid of some normal distribution, the points should approximately lie on the line $X = \mu + \sigma Z$.

In practice, to make the process deterministic, we often plot $(E[Z_{(i)}], X_{(i)})$, with

$$E[Z_{(i)}] = \Phi^{-1}(\Phi(E[Z_{(i)}])) \approx \Phi^{-1}(E[\Phi(Z_{(i)})]) = \Phi\left(\frac{i}{n+1}\right).$$

The last equality comes from the discussion below.

10.5.1 Order Statistics

Let X_1, \dots, X_n be iid with CDF F and density f , both continuous. Order them as $X_{(1)} \leq \dots \leq X_{(n)}$.

Example 10.22. We have

$$F_{(1)}(x) = \mathbb{P}(X_{(1)} \leq x) = 1 - \mathbb{P}(X_{(1)} \geq x) = 1 - \mathbb{P}(X_i > x, \forall i) = 1 - (1 - F(x))^n.$$

$$F_{(n)}(x) = \mathbb{P}(X_{(n)} \leq x) = \mathbb{P}(X_i < x, \forall i) = F(x)^n.$$

More generally,

$$\begin{aligned} F_{(k)}(x) &= \mathbb{P}(X_{(k)} \leq x) = \mathbb{P}(\text{as least } k \text{ of } X_i \leq x) \\ &= \sum_{i=k}^n \binom{n}{i} F(x)^i (1 - F(x))^{n-i}. \end{aligned}$$

Hence,

$$\begin{aligned} f_{(k)} &= F'_{(k)}(x) = \sum_{i=k}^n \frac{n!}{(n-i)!i!} i F(x)^{i-1} f(x) (1 - F(x))^{n-i} \\ &\quad - \sum_{i=k}^n \frac{n!}{(n-i)!i!} F(x)^i (n-i) (1 - F(x))^{n-i-1} f(x) \\ &= n f(x) \left[\sum_{i=k}^n \frac{(n-1)!}{(n-i)!(i-1)!} F(x)^{i-1} (1 - F(x))^{n-i} + \sum_{i=k}^n \frac{(n-1)!}{i!(n-i-1)!} F(x)^i (1 - F(x))^{n-i-1} \right] \\ &= n f(x) \frac{(n-1)!}{(k-1)![(n-1)-(k-1)]!} F(x)^{k-1} (1 - F(x))^{(n-1)-(k-1)} \\ &= n f(x) \binom{n-1}{k-1} F(x)^{k-1} (1 - F(x))^{n-k}. \end{aligned}$$

That is, we have

Proposition 10.23. *Let X_1, \dots, X_n be iid with CDF F and density f , both continuous. Order them as $X_{(1)} \leq \dots \leq X_{(n)}$. We have then that*

$$F_{(k)} = \sum_{i=k}^n \binom{n}{i} F(x)^i (1 - F(x))^{n-i}$$

and

$$f_{(k)} = n f(x) \binom{n-1}{k-1} F(x)^{k-1} (1 - F(x))^{n-k}.$$

Remark 10.24 (Sidenote). Let $W = \Phi(Z)$ and $Z \sim \mathcal{N}(0, 1)$. We have

$$\mathbb{P}(W \leq w) = \mathbb{P}(\Phi(Z) \leq w) = \mathbb{P}(Z \leq \Phi^{-1}(w)) = \Phi(\Phi^{-1}(w)) = w$$

and thus

$$W \sim \text{Uniform}[0, 1].$$

This is true in fact for all random variables. That is, we have

Proposition 10.25. *Let X be a random variable with distribution F . Then*

$$F(X) \sim \text{Uniform}[0, 1], \quad X = F^{-1}(U),$$

where

$$F^{-1}(u) := \inf \{y : f(y) \geq u\}.$$

The latter result is useful for simulation.

We can now apply this to the discussion of Q-Q plots above: For $U \sim \mathcal{N}(0, 1)$, we have

$$f_{(k)}(x) = n f(x) \binom{n-1}{k-1} F(x)^{k-1} (1 - F(x))^{n-k} \sim \text{Beta}(k, n - k + 1).$$

Thus it has mean

$$\mathbb{E}[f_{(k)}] = \frac{k}{k + (n - k + 1)} = \frac{k}{n + 1}.$$

11 Simple Regression (A)

Proposition 11.1. *The least square estimators are*

$$\hat{\beta}_1 = \frac{\overline{(xy)} - \bar{x} \cdot \bar{y}}{\overline{x^2} - \bar{x}^2}, \quad \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}.$$

We can also write

$$\bar{\beta}_1 = \sum c_i Y_i, \quad c_i := \frac{x_i - \bar{x}}{\sum (x_i - \bar{x})^2}.$$

Note that $\sum c_i = 0$, $\sum c_i x_i = 1$, $\sum c_i^2 = 1/[\sum (x_i - \bar{x})^2]$. We can use these to calculate the mean and variance of the least square estimators.

Thus we have:

Proposition 11.2. *The least square estimators are unbiased.*