

Measure Theory

Sec1. Probability Space

Def. σ -field: i) $\Omega \in \mathcal{F}$; ii) if $A \in \mathcal{F}$, then $A^c \in \mathcal{F}$; iii) if $A_1, A_2, \dots \in \mathcal{F}$, then $\bigcup_{i=1}^{\infty} A_i \in \mathcal{F}$.
algebraic: if $A_1, \dots, A_n \in \mathcal{F}$, then $\cup_{i=1}^n A_i \in \mathcal{F}$.

Fact. If $F_i, i \in I$ are all σ -fields, $\cap_{i \in I} F_i$ is a σ -field.

Def. measure μ : i) $\mu(A) \geq 0, \forall A \in \mathcal{F}$; ii) $\mu(\emptyset) = 0$; iii) if $A_1, A_2, \dots \in \mathcal{F}$ disjoint, then $\mu(\cup_i A_i) = \sum_i \mu(A_i)$.

Thm 1.1.1. (i) continuity from below: if $A_i \downarrow A$, then $\mu(A_i) \uparrow \mu(A)$.
(ii) continuity from above: if $A_i \uparrow A$, then $\mu(A_i) \downarrow \mu(A)$.
Def. Borel σ -field: $\mathcal{B} = \sigma\{A, b : -\infty < a < b < \infty\}$.

Def. π -system \mathcal{P} : if $A, B \in \mathcal{P}$, then $A \cap B \in \mathcal{P}$.
 λ -system \mathcal{L} : i) $\Omega \in \mathcal{L}$; ii) if $A, B \in \mathcal{L}$ and $A \subset B$, then $B \setminus A \in \mathcal{L}$; iii) if $A_1, A_2, \dots \in \mathcal{L}$, and $A_i \uparrow A$, then $A \in \mathcal{L}$.

Ex * If \mathcal{F} is π -system and λ -system, then \mathcal{F} is σ -field.

Thm 2.1.2. If \mathcal{P} is π -system, \mathcal{L} is λ -system, $\mathcal{P} \subset \mathcal{L}$, then $(\mathcal{P}) \subset \mathcal{L}$.

Sec2-1. Measurable Function

Def. f is meas- if $\{\omega_1 \in \Omega_1 : f(\omega_1) \in A\} = f^{-1}(A) \in \mathcal{F}_1, \forall A \in \mathcal{F}_2$.

Fact. gen- σ -field by f : $\{f(\omega)\} = \{f^{-1}(A) : A \in \mathcal{F}_2\}$ is σ -field in Ω_1 , $\{A \subset \Omega_2 : f^{-1}(A) \in \mathcal{F}_1\}$ is a σ -field in Ω_2 .

Thm 1.3.1. If $\mathcal{F}_2 = \sigma(\mathcal{A}_2)$, and $f^{-1}(\mathcal{A}_2) \subset \mathcal{F}_1$, then f is meas-.

Thm 1.3.2. If f_1 and f_2 are meas-, then $f_2 \circ f_1$ is meas-.

Def. induced measure: $\mu_2(A) = \mu_1(f^{-1}(A))$.

Sec2-2. Random Variable

Thm 1.3.5. $\inf_n X_n, \sup_n X_n, \limsup X_n, \liminf X_n$ are r.v.-.

Ex 1.3.1. $\sigma(X^{-1}(A)) = X^{-1}(\sigma(A))$
Hint: $C = \{B \in \sigma(A) : X^{-1}(B) \in \sigma(X^{-1}(A))\}$.

Sec2-3. Distribution

Thm 1.2.2. $\Omega = (0, 1), \mathcal{F} = \mathcal{R}, P = \text{Lebesgue measure}$, then $X(\omega) = F^{-1}(\omega) = \inf\{\eta \in \mathbb{R} : F(\eta) \geq \omega\} = \sup\{\eta \in \mathbb{R} : F(\eta) < \omega\}$ with dist-
F.
Hint: $\{\omega : \omega \leq F(\omega)\} = \{\omega : X(\omega) \leq \bar{x}\}$, right-continuous of F .

Sec3. Expectation

Def. indicator \rightarrow simple \rightarrow non-negative \rightarrow arbitrary.
case3: $E[X] = \sup\{X(\omega) : 0 \leq Y \leq X, Y \text{ is simple}\}$.

Prop. a) monotonicity, b) linearity.
Hint: $Z_M^{(i)} = \frac{1}{M^i} [2^M Z], Z_M^{(i)} = \frac{1}{M^i} [2^M Z - 1]$ for truc- case.

Thm MC. X_n n-n seq- r.v.-, if $X_n \uparrow X$, then $E[X_n] \uparrow E[X]$.
Hint: $Y_e = \sum_i (b_i - 1) b_i$.

Thm Fatou's L. If $X_n \geq 0$, then $\liminf E[X_n] \geq E[\liminf X_n]$.
Thm DC. If $X_n \rightarrow X$, $|X_n| \leq Y$ with $E[Y] < \infty$, then $E[Y_n] \rightarrow E[Y]$.

Thm Jensen. If $E[X] < \infty$, φ is convex, then $E[\varphi(X)] \leq \varphi(E[X])$.

Thm Hölder. If $p, q \geq 1$ and $\frac{1}{p} + \frac{1}{q} = 1$, then $E[XY] \leq \|X\|_p \|Y\|_q$.
Hint: $x \cdot y \leq x^p/p + y^q/q$ via concav- of log.

Thm Minkowski. for $p \geq 1$, $\|X + Y\|_p \leq \|X\|_p + \|Y\|_p$.

Thm Markov. if r.v- $X \geq 0$ and $a > 0$, then $P(X \geq a) \leq \frac{1}{a} E[X]$.

Thm Chebyshev. if \exists var, then $P(|X - E[X]| \geq a) \leq \frac{1}{a^2} \text{Var}(X)$.

Law of Large Number

Sec1. Independence

Def. inde- events \rightarrow collections (σ -fields) \rightarrow random variables.

Thm 2.1.3. if sys- $\mathcal{A}_{i=1}^n$ are inde-, then $\sigma(\mathcal{A}_i)_{i=1}^n$ are indep-

Proof. It suffices to show that $P(B_1 \cdots B_n) = P(B_1) \cdots P(B_n), \forall B_i \in \sigma(\mathcal{A}_i)$.
Fix $B_i \in \mathcal{A}_i$ or $B_i = \Omega$ for $i = 1, \dots, n$. Define $L_1 := \{B_1 \in \sigma(\mathcal{A}_1) : P(B_1 \cdots B_n) = P(B_1) \cdots P(B_n)\}$.
Note that L_1 contains the π -system \mathcal{A}_1 by assumption. Moreover, it can be easily verified that L_1 is a λ -system (details omitted). Therefore, by the π - λ theorem, $P(A_1) \subset L_1$. We have proved that $\mathcal{A}_1, \mathcal{A}_2, \dots, \mathcal{A}_n$ are independent. **Repete the argument**, X_n

Thm 2.1.4. $R\text{-}V$ - $X_{i=1}^n$ are inde- if- f $P(\cap_i X_i \leq x_i) = \prod_i P(X_i \leq x_i)$.

Thm ~. If $X_{i=1}^n$ are inde-, then $\sigma(X_i : i \in I) \perp \sigma(X_j : j \in J)$.

Thm 2.1.5. If $X_{i=1}^n$ inde-, then $\{X_i, i \in I\} \perp \{X_j, j \in J'\}$, g, h meas-.

Thm 2.1.6. If X , Y inde-, $E[X] < \infty$ or ≥ 0 , then $E[XY] = E[X]E[Y]$.

Case 3. $X, Y \geq 0$. Choose a sequence of simple random variables $X_n \uparrow X$ and $\sigma(X_n) \subset \sigma(X)$. (X_M can be chosen as X_0^0 in Chapter 1, Section 3.2.) Choose also a sequence of simple random variables $Y_n \uparrow Y$ and $\sigma(Y_n) \subset \sigma(Y)$. From Case 2, we have $E(X_n Y_n) = E(X_n Y_0^0)$. Then the limit as $n \rightarrow \infty$ to obtain (1.1).

Thm Kolmogorov's 0-1 Law. If X_i 's inde-, tail σ -field $T = \cap_{n=1}^{\infty} \sigma(X_k, k \geq n)$, then $P(A) = 0, 1$ for $A \in T$.

Proof. Step 1: $\sigma(X_{i+1}, X_{i+2}, \dots)$ can be written as $\sigma(\cup_{j=i}^{\infty} \sigma(X_{i+1}, \dots, X_{j+1}))$. The union inside of $\sigma(X_{i+1}, \dots, X_{j+1})$ is independent of this σ -system by the condition. Therefore, $\sigma(X_{i+1}, \dots, X_{j+1})$ is independent of $\sigma(X_{i+1}, X_{i+2}, \dots)$.
Step 2: Write $\sigma(X_1, X_2, \dots) = \sigma(\cup_{j=1}^{\infty} \sigma(X_1, \dots, X_j))$ and argue similarly that $\sigma(X_1, X_2, \dots)$ is independent of T .
Step 3: By definition, $T \subset \sigma(X_1, X_2, \dots)$.
Therefore T is independent of itself. The theorem then follows from the previous fact.

Sec2-1. Weak Law of Large Number

Def. converges in prob: if $\forall \epsilon > 0$, $P(|Y_n - Y| > \epsilon) \rightarrow 0$, as $n \rightarrow \infty$. converges in L_p $E[Y_n - Y]^p \rightarrow 0$, as $n \rightarrow \infty$. $E[Y_n - Y]^p \leq \frac{C}{\epsilon^n}$

Theorem. If $(S_n)_{n=1}^{\infty}$ is a sequence of random variables with $\sigma_n^2 = \text{Var}(S_n)$ and $\sigma_n^2/b_n^2 \rightarrow 0$, **Thm SLLN**. i.e. $E[|X_i|] < \infty$, then $S_n/n \rightarrow \mu$ a.s.

[Step 1: Truncation.] Let $Y_k = X_{k \wedge \lfloor |X_k| \rfloor}$, $T_n = \sum_{i=1}^n Y_i$. Because $\sum_{k=1}^{\infty} P(|X_k| > k) \leq E|X_1| < \infty$, we have by B-C lemma (i)

$$P(X_k \neq Y_k \text{ i.o.}) = 0.$$

Therefore, to prove the theorem, it suffices to show $\frac{T_n}{b_n} \rightarrow \mu$ a.s..

[Step 2: 2nd moment calculation.] Fix $k > 1$. Choose $c(n) = [\alpha^n]$ (this means the integer part, and it is obviously $\leq \alpha^n$). We have, for any $\epsilon > 0$,

$$\sum_{n=1}^{\infty} P\left(\frac{T_{k(n)}}{b_n} > \epsilon\right) \leq 4(1 - \alpha^{-2})^{-1} \epsilon^{-2} \sum_{m=1}^{\infty} \frac{\text{Var}(Y_m)}{m^2}.$$

Since $\text{Var}(Y_m) \leq \frac{\text{Var}(X_m)}{m^2} = \frac{\sum_{j=1}^m j^2 P(|Y_j| > j) y_j dy}{m^2} \leq C \int_0^{\infty} P(|X_1| > y) dy = CE|X_1| < \infty$, we have

$$\sum_{n=1}^{\infty} P\left(\frac{T_{k(n)}}{b_n} > \epsilon\right) < \infty,$$

and by B-C lemma (i),

$$\frac{T_{k(n)}}{b_n} \rightarrow \mu \text{ a.s.}$$

Note also that, because $EY_m = EX_{1 \wedge \lfloor |X_1| \rfloor} \rightarrow \mu$ as $m \rightarrow \infty$, we have

$$\frac{ET_{k(n)}}{b_n^2} = \frac{\sum_{m=1}^{k(n)} EY_m}{b_n^2} \rightarrow \mu.$$

Therefore, $\frac{T_{k(n)}}{b_n} \rightarrow \mu$ a.s..

[Step 3: Subsequence method.] For $k(n) \leq m \leq k(n+1)$, we have

$$\frac{k(n)}{k(n+1)} \frac{T_{k(n)}}{b_n} \leq \frac{T_m}{b_m} \leq \frac{k(n+1)}{k(n+1)} \frac{T_{k(n+1)}}{b_{n+1}}.$$

The left-hand side tends to μ/a and the right-hand side tends to μ/a . Therefore, for any $\alpha > 1$,

$$\limsup_{m \rightarrow \infty} \frac{T_m}{m} \leq \mu/a, \quad \liminf_{m \rightarrow \infty} \frac{T_m}{m} \geq \mu/a.$$

The theorem follows by letting $\alpha \downarrow 1$. \square

Remark. for $0 < \epsilon < 1$, $E[|X|^{1-\epsilon}] < \infty$. $E[X] = \frac{1}{\epsilon} E[|X|^{1-\epsilon}] \rightarrow \mu$ in prob.

Lem 2.2.8. If $Y \geq 0$ and $p > 0$, then $E[Y^p] = \int_0^{\infty} py^{p-1} P(Y > y) dy$.

Thm 2.2.9. (WLNN without finite 1st moment)
i.i.d. $\mu = E[X_1]$. If $E[|X_1|] < \infty$, then $\frac{1}{n} S_n \rightarrow \mu$ in prob.

Proof of the Remark. Consider $\frac{X_n}{n} = \frac{S_n}{n} - \frac{S_{n-1}}{n-1}$. This tends to 0 a.s. because of the condition $\frac{S_n}{n} \rightarrow \mu$ a.s.. Therefore, $P(\left|\frac{X_n}{n}\right| > 1)$ is 0 and by the B-C lemma (ii),

$$\sum_{n=1}^{\infty} P\left(\left|\frac{X_n}{n}\right| > 1\right) < \infty.$$

Therefore, from the i.i.d. assumption,

$$E|X_1| \leq 1 + \sum_{n=1}^{\infty} P(|X_1| > n) < \infty.$$

Next result concerns the case of infinite mean.

Theorem. Let X_1, X_2, \dots be i.i.d. and $S_n = \sum_{i=1}^n X_i$. If $E(X_1^+) = \infty$ and $E(X_1^-) < \infty$, then $S_n \rightarrow \infty$ a.s..

Proof. Define $X_i^M = X_i \wedge M$. Then $E[X_i^M] < \infty$. By SLLN, for $S_0^M := \sum_{i=1}^n X_i^M$, we have

$$\frac{S_0^M}{n} \rightarrow EX_1^M \text{ a.s.}$$

Therefore,

$$\liminf_{n \rightarrow \infty} \frac{S_n}{n} \geq \lim_{n \rightarrow \infty} \frac{S_0^M}{n} = EX_1^M \rightarrow \infty, \text{ as } M \rightarrow \infty.$$

Next result concerns the case of infinite mean.

Theorem. Let X_1, X_2, \dots be i.i.d. and $S_n = \sum_{i=1}^n X_i$. If $E(X_1^+) = \infty$ and $E(X_1^-) < \infty$, then $S_n \rightarrow \infty$ a.s..

Proof. Define $X_i^M = X_i \wedge M$. Then $E[X_i^M] < \infty$. By SLLN, for $S_0^M := \sum_{i=1}^n X_i^M$, we have

$$\frac{S_0^M}{n} \rightarrow EX_1^M \text{ a.s.}$$

Prop. Define $Y_i = X_i \wedge 1$. Then $E[Y_i] < \infty$. By SLLN, for $S_0 = \sum_{i=1}^n Y_i$, we have

$$\frac{S_0}{n} \rightarrow EX_1 \text{ a.s.}$$

Thm 2.5.2. (Kolmogorov's Maximal Inequality)
indep., $E[X_i] = 0, E[X_i^2] < \infty$, then $P(\max_{k \leq n} |S_k| \geq x) \leq E[S_n^2]/x^2$.
Hint: $A_k := \{S_i : i < k, |S_i| \geq x\}$.

Thm 2.5.3. indep., $E[X_i] = 0, \sum_i E[X_i^2] < \infty$ \Rightarrow $\sum_i X_i$ converges a.s.
Hint: $\omega_M = \sup_{n>M} n/S_n \rightarrow 0$ a.s. as $M \rightarrow \infty$.
Proof. Let $S_n = \sum_{i=1}^n X_i$. It suffices to prove S_n is a Cauchy sequence, a.s., that is,
 $\omega_M := \sup_{n>M} |S_n - S_M| \rightarrow 0$ a.s. as $M \rightarrow \infty$.

To show this, for any $\epsilon > 0$, we write

$$P(\omega_M > 2\epsilon) = P(\omega_M > 2\epsilon, E[X] > \epsilon) + P(\omega_M > 2\epsilon, E[X] \leq \epsilon) \quad (\text{by monotonicity of } \omega_M)$$

$$= P(P(\omega_M > 2\epsilon) > 2\epsilon) \leq \frac{1}{2\epsilon} \sum_{i=1}^{\infty} P(|X_i| > \epsilon) \leq \frac{1}{2\epsilon} \sum_{i=1}^{\infty} E[X_i^2] \leq \frac{C}{\epsilon^2}.$$

Note that, by definition of ω_M and the union bound,

$$P(\omega_M > 2\epsilon) = P(\sup_{m > M} |S_m - S_M| > 2\epsilon) \leq P(\sup_{m > M} |S_m - S_{m-1}| > \epsilon) + P(\sup_{m > M} |S_{m-1} - S_M| > \epsilon) \leq 2P(\sup_{m > M} |S_m - S_{m-1}| > \epsilon).$$

Moreover, by the Kolmogorov's maximal inequality,

$$P(\sup_{m \geq M} |S_m - S_M| > \epsilon) \leq \frac{1}{\epsilon^2} \sum_{m=M+1}^{\infty} E[X_m^2] \rightarrow 0.$$

Combining the above arguments, we conclude that

$$P(\omega_M > 2\epsilon) = 0, \forall \epsilon > 0,$$

hence $\omega_M \rightarrow 0$. \square

Example. Let X_1, X_2, \dots be i.i.d. Rademacher variables, i.e., $P(X_i = 1) = P(X_i = -1) = \frac{1}{2}$. We have, for $\alpha > 1$, $\sum_{n=1}^{\infty} \frac{X_n}{n}$ converges a.s.

Thm 2.5.4. (Kolmogorov's three-series thm) X_i indep., $Y_i = X_{i \wedge \lfloor |X_i| \rfloor} \in A$,
i) $\sum_i P(|X_i| > A) < \infty$, ii) $\sum_i E|Y_i| < \infty$, iii) $\sum_i \text{Var}(Y_i) < \infty$, then $\sum_i X_i$ converges a.s.

Proof. WLOG, assume $\mu = 0$ (otherwise, consider $X_i - \mu$). We have

$$P\left(\frac{|S_n|}{n} > \epsilon\right) \leq \frac{E|S_n|^3}{\epsilon^3 n^2} \leq \frac{\sum_{i=1}^n E|X_i|^3}{\epsilon^3 n^2} \leq \frac{C}{\epsilon^3 n^2}$$

Therefore, $\sum_n X_n$ converges a.s. \square

Thm 2.5.5. (Kronecker's Lemma) If $a_n \uparrow \infty$ and $\sum_{i=1}^{\infty} \frac{x_i}{a_i}$ converges a.s., then $\frac{\sum_{i=1}^n x_i}{a_n} \rightarrow 0$.

Second Proof. Recall that WLOG, we can assume $X_1 \geq 0$. As in the first proof, we let $Y_i = X_{i \wedge \lfloor |X_i| \rfloor}$, $T_n = \sum_{i=1}^n Y_i$. We have argued that $P(X_1 \neq Y_1 \text{ i.o.}) = 0$ and it suffices to show $\frac{1}{a_n} \uparrow \mu$ a.s.. We have also argued that $\frac{EY_1}{a_1} \rightarrow \mu$. Therefore, we are left to show

$$\frac{n}{a_n} Y_1 - EY_1 \rightarrow 0 \text{ a.s.}$$

By Kronecker's lemma, we only N.T.S.

$$\sum_{k=1}^n \frac{Y_k - EY_k}{a_k} \text{ converges a.s.}$$

Taking $A = 1$ in Kolmogorov's Three-series Theorem, (i) and (ii) therein are automatic satisfied. For (iii), recall that we have verified in the first proof of SLLN that

$$\sum_{k=1}^{\infty} \frac{Y_k - EY_k}{k^2} \leq \sum_{k=1}^{\infty} \frac{E(Y_k^2)}{k^2} < \infty.$$

Then

$$W_n = \sum_{i=1}^n \frac{X_i - \mu}{\sqrt{n}}$$

Proof. Let $g: \mathbb{R} \rightarrow \mathbb{R}$ be any bounded continuous function with bounded and continuous derivatives up to the third order. By the above remark, it suffices to prove $Eg(W_n) \rightarrow Eg(Z)$. For each $i = 1, \dots, n$, let

$$\xi_i = \frac{X_i - \mu}{\sqrt{n}}$$

Define η_1, \dots, η_n on the same probability space such that $\{\xi_1, \dots, \xi_n, \eta_1, \dots, \eta_n\}$ are independent and $\eta_i \sim N(0, 1)$ for all i . Note that

$$Eg(W_n) = Eg(\eta_1) = E\eta_1^2 = \frac{1}{n} E[\xi_1^3] = E[\xi_1^3]/E[\xi_1^2] \rightarrow \frac{C}{n^2}. \quad (1.1)$$

Note also that $\sum_{i=1}^n \eta_i \sim N(0, 1)$. By using the telescoping sums of Taylor's expansion, we have

$$\begin{aligned} E(g(V_n) - Eg(Z)) \\ &= \sum_{k=1}^n \frac{1}{k!} \frac{\partial^k g}{\partial \xi^k}(0) \sum_{i=1}^k \frac{1}{i!(p-1)!} E[\xi_1^i | \cup_{j=1}^k \xi_j^i] p^k \\ &\stackrel{(1.1)}{=} -E \sum_{k=1}^n \frac{1}{k!} \frac{\partial^k g}{\partial \xi^k}(0) \sum_{i=1}^k \frac{1}{i!(p-1)!} E[\xi_1^i | \cup_{j=1}^k \xi_j^i] p^k \end{aligned}$$

where $V_k := \xi_1 + \dots + \xi_{k-1} + \eta_{k+1} + \dots + \eta_n$. Using independence and (1.1) for cancellation, we obtain

$$|E(g(V_n) - Eg(Z))| \leq Cn \frac{1}{n^2} \rightarrow 0.$$

Definition. the characteristic function (ch.f.) of a random variable X is defined to be

$$\varphi_X(t) := E[e^{itX}] = E[\cos(tX) + i \sin(tX)].$$

Properties. 1. $\varphi_X(0) = 1, |\varphi_X(t)| \leq 1$.
2. $\varphi_X(-t) = \varphi_X(t)$. (conjugate)
3. $|\varphi_X(t+h) - \varphi_X(t)| \leq |E[e^{ithX} - 1]| \rightarrow 0$, as $h \rightarrow 0$. (by DCT). That is, $\varphi_X(t)$ is uniformly continuous.
4. $\varphi_X(at+b) = a\varphi_X(t) + b\varphi_X(at)$.
5. If X_1 is independent of X_2 , then $\varphi_{X_1+X_2}(t) = \varphi_{X_1}(t) \varphi_{X_2}(t)$.

Theorem 3.3.8. If $E(X^2) < \infty$, then

$$\varphi_X(t) = 1 + t \cdot E[X] - \frac{t^2}{2} E(X^2) + o(t^2), \text{ as } t \rightarrow 0.$$

Ex 3.3.6. If X is a r.v. with $E[X] = 0$, $E[X^2] = 1$, then $\varphi_X(t) = e^{itX}$.

Proof. First note that

$$\sum_{i=1}^{\infty} E[Y_i^2] = \sum_{i=1}^{\infty} \frac{1}{i^2} E^2(X_{i \wedge \lfloor |X_i| \rfloor}) \leq \frac{1}{2} E^2(X) \frac{1}{1-e^{-2t}} e^{-2t} \rightarrow \frac{1}{2} E^2(X) e^{itX} = \varphi_X(t).$$

By Fubini's theorem, for any $b < a$, we have

$$\begin{aligned} \sum_{i=1}^{\infty} E[Y_i^2] &\leq \frac{1}{2} \int_0^T \int_0^t e^{-it(a-s)} ds dt = P(a < X < b) + \frac{1}{2} P(X = b) = \frac{1}{2} P(X = b). \end{aligned}$$

The existence of the limit is part of the statement of the theorem.

Central Limit Theorem

Sec1. Convergence in Distribution

Thm . (Stirling's Formula) $n! \sim n^{n+1/2} \sqrt{2\pi n}$ as $n \rightarrow \infty$.

Def. d.f.s $f_n \Rightarrow f$ weakly conv, if $f_n(y) \rightarrow f(y)$ at \forall cont-point y of f .
Fact 1.) If $X_n \rightarrow X$ in P , then $X_n \rightarrow X$ if $\forall x \in C$, if $X_n \rightarrow c$, then $X_n \rightarrow c$ in P .
Thm 3.2.2. (Skorokhod's Theorem) If $f_n \Rightarrow f$, X_n has d.f.s. f_n , and $Y_n \rightarrow Y$ a.s..
Proof. Let $\Omega = (0, 1)$, \mathcal{F} = (Borel sets), P = Lebesgue measure. Define $Y_n : \Omega \rightarrow \mathbb{R}$ to be

$$Y_n = F_n^{-1}(\omega),$$

where

$$F_n^{-1}(\omega) := \inf\{y : F_n(y) \geq \omega\} = \sup\{y : F_n(y) < \omega\}.$$

Recall from Chapter 1 that the d.f. of the above constructed F_n is Y_n . Let Ω_0 consists subsequential value limit is a distribution function. By the condition $\varphi_{Y_n}(t) \rightarrow \varphi_Y(t)$, Then Ω_0 is countable; hence Ω_0 has probability 1. Moreover, it follows from simple calculus that $\varphi_{Y_n}(t) \rightarrow \varphi_Y(t)$ for all $t \in \Omega$.
This implies the whole sequence converges weakly to the limit. \square

Theorem. Let X_1, X_2, \dots be a sequence of i.i.d. random variables such that $X_1 = \sum_{i=1}^{\infty} X_i/\sqrt{n}$. The random variables are defined on the same probability space. Then

$$X_n + Y_n \Rightarrow X + c, \quad X_n Y_n \Rightarrow c X.$$

Proof. Let x be a continuity point of F_X (so x is a continuity point of F_{X+n}). Then choose a decreasing sequence of $\epsilon > 0$ such that $x + \epsilon$ is a continuity point of F_X . We have

$$\begin{aligned} P(X_n + Y_n \leq x) &\leq P(Y_n \leq x - c) + P(X_n \leq x - c + \epsilon) \\ &\leq P(Y_n \leq x - c) + P(X_n \leq x + \epsilon) \rightarrow 0, \end{aligned}$$

Following similar arguments, we can prove the lower bound

$$P(X_n + Y_n \geq x) \leq P(X \geq x + c).$$

This shows $X_n + Y_n \Rightarrow X + c$ for the case $c > 0$, we instead, for $0 < c < \epsilon$,

$$P(X_n + Y_n \leq x) \leq P(Y_n \leq c - \epsilon) + P(X_n \leq x - \epsilon) \leq P(X_n \leq x - \epsilon) \rightarrow 0.$$

Then

$$W_n \xrightarrow{d} Z \sim N(0, 1).$$

Proof. N.T.S. $E e^{itW_n} \rightarrow e^{-t^2}$ for all $t \in \mathbb{R}$. We have, by the expression of W_n and independence,

$$\begin{aligned} E e^{itW_n} &= E \exp(it(\frac{X_1 - \mu}{\sqrt{n}} + \dots + \frac{X_n - \mu}{\sqrt{n}})) = \prod_{i=1}^n E e^{it(\frac{X_i - \mu}{\sqrt{n}})} \\ &= \prod_{i=1}^n \left[1 + i \frac{t}{\sigma \sqrt{n}} E(X_j - \mu) + \frac{t^2}{2\sigma^2 n} E(X_j - \mu)^2 + o(\frac{t^2}{n}) \right] \end{aligned}$$

From Theorem 3.3.8,

$$E e^{itW_n} = \frac{1}{n} \sum_{i=1}^n \frac{1}{i^2} E^2(X_j - \mu)^2 = \frac{1}{2\sigma^2 n} \sum_{i=1}^n i^2 \rightarrow \frac{1}{2\sigma^2}.$$

□

Random Walks

Sec1. Random Walks
Def. If r -v-s X_i i.i.d. $S_n = \sum_i X_i$, then $\{S_n\}$ is random walk with $S_0 = 0$.

Thm 4.1.1. (Hewitt-Savage 0-1 Law) If A is permutable (Not change under finite perm-), then $P(A) = 0$ or 1.

Thm 4.1.2. For random walk on \mathbb{R} , one of follow- has prob- 1:

i) $S_n = 0$ i.i.d. $S_n \rightarrow \infty$ iii) $S_n \rightarrow -\infty$ iv) $\liminf S_n < \limsup S_n = \infty$.

Proof. By the 0-1 law, $(\limsup_{n \rightarrow \infty} S_n) \geq c$ has probability 0 or 1. This mean with probability 1, $\limsup_{n \rightarrow \infty} S_n$ some value in $\mathbb{R} \cup (-\infty, \infty)$. Similarly, $\liminf_{n \rightarrow \infty} S_n$ equals some value in $\mathbb{R} \cup (-\infty, \infty)$. If $\limsup_{n \rightarrow \infty} S_n = c$ for some $c \in \mathbb{R}$, then from

$$\limsup S_n = \limsup_{n \rightarrow \infty} (S_{n+1} - X_1) = \limsup_{n \rightarrow \infty} S_{n+1} - X_1 = \limsup S_n - X_1,$$

we must have $X_1 = 0$ and it belongs to case (i). Other cases correspond to other combinations of possibilities of $\limsup_{n \rightarrow \infty} S_n$ and $\liminf_{n \rightarrow \infty} S_n$. \square

Def. $\mathcal{F}_0 = \{\emptyset, \Omega\}, \mathcal{F}_n = \sigma(X_1, \dots, X_n)$ seq- of incr- σ -fields is filtration.

Def. rand-time $\tau \in \mathbb{R}^+ \cup \{\infty\}$ is stopping time wrt $\{\mathcal{F}_n\}$ if $\{\tau = n\} \in \mathcal{F}_n$.

Fact. i) $\{\tau = n\} \in \mathcal{F}_n$ ii) $\{\tau \leq n\} \in \mathcal{F}_n$ iii) $\{\tau \geq n+1\} \in \mathcal{F}_n$ are equivalent.

Fact. i) $\tau_1 \vee \tau_2$ ii) $\tau_1 + \tau_2$ iii) $\tau_1 \wedge \tau_2$ are stopping times.

Thm 4.1.5. (Wald's eq) i.i.d. if $E[\mathbf{X}|\mathcal{F}_\tau] \leq \mathbf{X}[\tau]$, then $E[S_\tau] = E[\mathbf{X}_1]E[\tau]$.

Thm 4.1.6. (W's 2eq) i.i.d. $\sim (0, \sigma^2)$, if $E[\tau] < \infty$, then $E[S_\tau^2] = \sigma^2 E[\tau]$.

$E(S_\tau) = E(\sum_{i=1}^{\tau} X_i)$

$$= E(\sum_{i=1}^{\infty} X_i 1_{\{\tau \geq i\}})$$

$$= ES_{\tau|n}^2 = ES_{(n-1)+}^2 + E(X_{n+1} 1_{\{\tau \geq n\}}) + E(X_n^2 1_{\{\tau \geq n\}})$$

$$= ES_{\tau|n}^2 + \sigma^2 P(\tau \geq n)$$

$$= \dots$$

$$= \sum_{i=1}^{\infty} E(X_i) E(1_{\{\tau \geq i\}})$$

$$= E(X_1) E(\tau).$$

Example. Let X_1, X_2, \dots be i.i.d. $\sim \text{Uniform}(0, 1)$. $S_n = X_1 + \dots + X_n$. Let $\tau := \inf\{n : S_n > 1\}$. Then

$$E(\tau) = e, \quad E(S_\tau) = \frac{e}{2}.$$

Proof. We have

$$\begin{aligned} P(\tau > n) &= P(S_n \leq 1) \\ &= \int_0^1 \dots \int_{x_1+...+x_{n-1}} 1_{\{x_1+...+x_n \leq 1\}} dx_1 \dots dx_n \quad (\text{from the uniform distribution}) \\ &= \int_0^1 \int_0^{x_1} \dots \int_0^{x_{n-1}} dy_1 dy_2 \dots dy_n \quad (\text{by a change of variable}) \\ &= \frac{1}{n!}. \end{aligned}$$

This implies

$$E(\tau) = \sum_{n=0}^{\infty} P(\tau > n) = e.$$

By Wald's first equation, $E(S_\tau) = E(X_1)E(\tau) = \frac{e}{2}$.

Example. Let X_1, X_2, \dots be i.i.d. with $P(X_1 = 1) = P(X_1 = -1) = \frac{1}{2}$. Let $S_n = X_1 + \dots + X_n$ (SRW). Let a, b be two integers with $a < 0 < b$. Let

$$N := \inf\{n : S_n \notin \{a, b\}\}.$$

Then

- 1. $E(N) < \infty$,
- 2. $S_N = a$ or b ,
- 3. $P(S_N = a) = \frac{b}{b-a}$,
- 4. $E(N) = E(S_N^2) = -a/b$.

Proof. For any positive integer k , by dividing the interval $(0, k(b-a))$ into k subintervals of equal length and considering an extreme case behavior (keep going upwards) of the random walk within each subinterval, we obtain

$$P(N > k(b-a)) \leq (1 - \frac{1}{2^{b-a}})^k.$$

This implies 1.

- 2. is obvious.
- 3. follows from Wald's first equation.
- 4. follows from Wald's second equation. \square

Sec2. Recurrence v.s. Transience
 $\tau_1 = \inf\{m : S_m = 0\}$, $\tau_n = \inf\{m > \tau_{n-1} : S_m = 0\}$.

Theorem. The following are equivalent:
 (i) $P(\tau_1 < \infty) = 1$.
 (ii) $P(\tau_n < \infty) = 1, \forall n = 1, 2, 3, \dots$.
 (iii) $P(\tau_m = 0 \text{ i.o.}) = 1$.
 (iv) $\sum_{n=1}^{\infty} P(S_m = 0) = \infty$.

Proof. We have

$$\begin{aligned} P(\tau_2 < \infty) &= P(\tau_1 < \infty, \tau_2 - \tau_1 < \infty) \\ &= \sum_{m,n=1}^{\infty} P(\tau_1 = m, \tau_2 - \tau_1 = n) \\ &= \sum_{m,n=1}^{\infty} P(X_1 + \dots + X_m = 0, X_1 + \dots + X_n \neq 0, \forall 1 \leq u < m; \\ &\quad X_{m+1} + \dots + X_{m+n} = 0, X_{m+1} + \dots + X_{m+v} \neq 0, \forall 1 \leq v < n) \\ &= \sum_{m,n=1}^{\infty} P(\tau_1 = m)P(\tau_1 = n) \quad (\text{by i.i.d. assumption}) \\ &= (P(\tau_1 < \infty))^2. \end{aligned}$$

Similarly, we can prove $P(\tau_n < \infty) = (P(\tau_1 < \infty))^n$. (2.1)

$$\begin{aligned} \sum_{m=0}^{\infty} P(S_m = 0) &= \sum_{m=0}^{\infty} E[1_{\{S_m=0\}}] = E \sum_{m=0}^{\infty} 1_{\{S_m=0\}} = \sum_{n=0}^{\infty} (P(\tau_1 < \infty))^n. \end{aligned}$$

Theorem 4.2.3. SRW is recurrent in \mathbb{R}^1 and \mathbb{R}^2 and is transient in $\mathbb{R}^d, d \geq 3$.

Proof. In \mathbb{R}^1 ,

$$\begin{aligned} \sum_{m=1}^{\infty} P(S_m = 0) &= \sum_{n=1}^{\infty} P(S_{2n} = 0) \\ &\quad \left(\text{can only return to } 0 \right. \\ &\quad \left. \sum_{n=1}^{\infty} \binom{2n}{n} \cdot \frac{1}{2}^{2n} \right) \quad (\text{co}) \\ &\sim \sum_{n=1}^{\infty} \frac{\sqrt{2\pi n} \binom{2n}{n}}{(\sqrt{2\pi n})^n} \frac{1}{2^{2n}} \\ &= \sum_{n=1}^{\infty} \frac{1}{\sqrt{\pi n}} = \infty; \end{aligned}$$

hence recurrent by the previous theorem. By \mathbb{R}^2 , similar calculation yields

$$P(S_{2n} = 0) \asymp \frac{1}{n},$$

still sum to infinity; hence still recurrent.

In \mathbb{R}^3 , more complicated combinatorics give $P(S_{2n} = 0) \asymp \frac{1}{n^2}$ summing \mathbb{R}^3 . (Theorem (Reflection Principle). Let X_1, X_2, \dots be i.i.d. with $P(X_1 = 1) = P(X_1 = -1) = \frac{1}{2}$, $S_n = X_1 + \dots + X_n$. For any positive integer b , we have

$$P(\max_{1 \leq k \leq n} S_k \geq b) = 2P(S_n > b) + P(S_n = b).$$

Proof. We have

$$\begin{aligned} P(\max_{1 \leq k \leq n} S_k \geq b) &= P(\max_{1 \leq k \leq n} E(S_k) \geq b) \\ &= P(\max_{1 \leq k \leq n} E(S_k) \geq b, S_n > b) + P(\max_{1 \leq k \leq n} E(S_k) \geq b, S_n = b). \end{aligned}$$

The first two terms are equal by reflecting the random walk trajectory along the horizontal line $y = b$ after it first hits level b . They both equal $P(S_n > b)$. The third term is simply $P(S_n = b)$. This gives the result. \square

Martingale

Sec1. Conditional Expectation

Def. $(\Omega, \mathcal{F}, \mathbb{P}), \mathbb{E}[X] < \infty$, σ -field $\mathcal{A} \subset \mathcal{F}$, $\mathbb{E}[X|\mathcal{A}]$ cond-expectation, if i) $\mathbb{E}[X|\mathcal{A}]$ is \mathcal{A} -measurable, ii) $\forall A \in \mathcal{A}, \mathbb{E}[X 1_A] = \mathbb{E}[X|\mathcal{A}] 1_A$.

Prop. a) $\mathbb{E}[X|\mathcal{A}] = \mathbb{E}[X|\mathcal{A}|]$.

(c) If $X \in \mathcal{A}$, then $\mathbb{E}[X|\mathcal{A}] = X$.

(d) If X is independent of \mathcal{A} , then $\mathbb{E}[X|\mathcal{A}] = \mathbb{E}[X]$.

(e) (linearity) If $\mathbb{E}[X], \mathbb{E}[Y] < \infty$, then $\mathbb{E}[aX + bY|\mathcal{A}] = a\mathbb{E}[X|\mathcal{A}] + b\mathbb{E}[Y|\mathcal{A}]$.

(f) (monotonicity) If $X \leq Y$, then $\mathbb{E}[X|\mathcal{A}] \leq \mathbb{E}[Y|\mathcal{A}]$.

Hint: consider $A = \{\mathbb{E}[X|\mathcal{A}] \leq \mathbb{E}[Y|\mathcal{A}]\} \geq \emptyset$.

(g) (cmct) If $X_n \geq 0, X_n \uparrow X$, then $\mathbb{E}[X_n|\mathcal{A}] \uparrow \mathbb{E}[X|\mathcal{A}]$.

(h) (cfatou) If $X_n \geq 0, \mathbb{E}[|X_n|] < \infty$, then $\lim \mathbb{E}[X_n|\mathcal{A}] \geq \lim \mathbb{E}[X_n]$.

(i) (cdct) If $|X_n| \leq Y, |Y| < \infty$, then $\mathbb{E}[X_n|\mathcal{A}] \rightarrow \mathbb{E}[X|\mathcal{A}]$ a.s.

(j) (cn) If $\mathbb{E}[X|\mathcal{A}], \mathbb{E}[Y|\mathcal{A}] < \infty$, then $\mathbb{E}[\varphi(\mathbb{E}[X|\mathcal{A}])|\mathcal{A}] \leq \mathbb{E}[\varphi(\mathbb{E}[Y|\mathcal{A}])|\mathcal{A}]$.

(k) (chol) If $p, q \geq 1$ and $\frac{1}{p} + \frac{1}{q} = 1$, $\mathbb{E}[X^p|\mathcal{A}] \wedge \mathbb{E}[Y^q|\mathcal{A}] < \infty$,

$\mathbb{E}[XY|\mathcal{A}] \leq \mathbb{E}[X^p|\mathcal{A}]^{1/p} \mathbb{E}[Y^q|\mathcal{A}]^{1/q}$.

Hint: suppose $\mathbb{E}[X^p|\mathcal{A}] \geq \epsilon > 0$, take $|X'| = (|X|^p + \epsilon^p)^{1/p}$, DCT.

(l) (cMin) If $\mathbb{E}[X] \geq 0, \mathbb{E}[X^p|\mathcal{A}] \leq \mathbb{E}[X^p|\mathcal{A}]^p$,

then $(\mathbb{E}[X + Y|\mathcal{A}])^{1/p} \leq (\mathbb{E}[X^p|\mathcal{A}])^{1/p} + (\mathbb{E}[Y^p|\mathcal{A}])^{1/p}$.

(m) (cMar) If $X \geq a, a > 0$, then $\mathbb{P}(X \geq a|\mathcal{A}) \leq \mathbb{E}[X|\mathcal{A}/a]$.

(n) If $\mathbb{E}[X], \mathbb{E}[XY] < \infty$ and $X \in \mathcal{A}$, then $\mathbb{E}[XY|\mathcal{A}] = \mathbb{E}[X|\mathcal{A}]\mathbb{E}[Y|\mathcal{A}] < \infty$.

Proof. We proceed to verify that $\mathbb{E}[XY|\mathcal{A}]$ satisfies (i) and (ii) in the definition of conditional expectation. (i) follows from the assumption that $X \in \mathcal{A}$. To verify (ii), we N.T.S. for any $A \in \mathcal{A}$,

$$E(XY|A) = E(X|\mathcal{A})\mathbb{E}[Y|A].$$

To verify (2.1), we proceed by considering X being an 1. indicator variable, 2. simple variable, 3. positive variable and 4. general variable as before.

1. Suppose $X = 1_B, B \in \mathcal{A}$. Then

$$\begin{aligned} \text{RHHS(1.1)} &= E[1_B 1_A \mathbb{E}[Y|\mathcal{A}]] = E[Y 1_B 1_A] = E(YX 1_A) = LHS(1.1). \\ &\quad (\text{2.1}) \end{aligned}$$

2. Suppose $X = \sum_{i=1}^m b_i 1_{B_i}$. Then

$$\begin{aligned} E[\sum_{i=1}^m b_i 1_{B_i} Y|\mathcal{A}] &\stackrel{\text{(i)}}{=} \sum_{i=1}^m b_i E[1_{B_i} 1_A \mathbb{E}[Y|\mathcal{A}]] \\ &= \sum_{i=1}^m b_i 1_{B_i} E(Y|1_{B_i} \cap \mathcal{A}) = \sum_{i=1}^m b_i 1_{B_i} E(Y|\mathcal{A}). \end{aligned}$$

3. Use monotone convergence theorem of the conditional expectation above and approximate X by simple variables from below.

4. Write general $X \neq X' - X'$.

(o) (Tower Property) if $\mathbb{E}[X] < \infty, \forall i \in A, \mathbb{E}[X|A] = \mathbb{E}(X|A)$

 i) $\mathbb{E}[\mathbb{E}[X|\mathcal{A}_1]|\mathcal{A}_2] = \mathbb{E}[\mathbb{E}[X|\mathcal{A}_1]|\mathcal{A}_1] = \mathbb{E}[X|\mathcal{A}_1]$

$\forall A \in \mathcal{A}_1 : \mathbb{E}[\mathbb{E}[X|\mathcal{A}_1]|\mathcal{A}_1] = \mathbb{E}[\mathbb{E}(X|\mathcal{A}_1)|\mathcal{A}_1] = \mathbb{E}(X|\mathcal{A}_1)$

(p) (triangular eq) if $\mathbb{E}[X^2] < \infty$, then for $\forall Y \in \mathcal{A}$ with $\mathbb{E}[Y^2] < \infty$,

$\mathbb{E}[X(X-Y)|\mathcal{A}]^2 \leq \mathbb{E}[X-Y]^2$.

(q) Var $(X) \geq \mathbb{E}[\text{Var}(X|\mathcal{A})]$.

(r) $Z \perp\!\!\!\perp (X, Y)$, then $\mathbb{E}[X|Y, Z] = \mathbb{E}[X|\mathcal{A}]$.

Hint: $\mathcal{P} = \{B \cap C\}, \mathcal{L} = \{A \in \sigma(Y, Z) : \mathbb{E}[X|\mathcal{A}] = \mathbb{E}[\mathbb{E}[X|Y]|\mathcal{A}]\}$.

(s) If $\{T_n \geq 1\}_{n \in \mathbb{N}}$, $\forall n \in \mathbb{N}$. $\{\mathcal{F}_n\}$ is a submartingale, then

Thm 5.2.3-4. (martingale transformation)

i) if $\{S_n, \mathcal{F}_n\}$ mt, φ convex, $\mathbb{E}[\varphi(S_n)|\mathcal{F}_n] < \infty$, then $\{\varphi(S_n)\}$ is sub-mt.

ii) if $\{S_n\}$ sub-mt, φ convex \uparrow , $\mathbb{E}[\varphi(S_n)|\mathcal{F}_n] < \infty$, then $\{\varphi(S_n)\}$ is sub-mt.

For case (i), we use Jensen's inequality to obtain

$$\mathbb{E}[\varphi(S_n)|\mathcal{F}_n] \geq \varphi(\mathbb{E}[S_n|\mathcal{F}_n]) = \varphi(S_{n-1}).$$

verifying (iii).

For case (ii), we use Jensen's inequality and the additional assumption that φ is increasing to obtain

$$\mathbb{E}[\varphi(S_n)|\mathcal{F}_n] \geq \varphi(\mathbb{E}[S_n|\mathcal{F}_n]) = \varphi(S_{n-1}).$$

We have, from (2.7),

Theorem 5.4.2 (Doob's inequality). (S_n, \mathcal{F}_n) is a submartingale. Then for any $x > 0$, we have

$$P(\max_{1 \leq k \leq n} S_k \geq x) \leq \frac{1}{x} \mathbb{E}[S_n I_{\max_{1 \leq k \leq n} S_k \geq x}] \leq \frac{E(S_n^+)}{x}.$$

Exercise: Show that Kolmogorov's maximal inequality follows from the above inequality.

Proof of Theorem 5.4.2. Let $N = \inf\{k : S_k \geq x\}$. N is a stopping time. Let

$$A = \{ \max_{1 \leq k \leq n} S_k \geq x \} = \{N \leq n\}.$$

Then

$$\mathbb{E}[A] \leq \frac{S_{n+1}}{x} I_{\max_{1 \leq k \leq n} S_k \geq x}.$$

We have, from (2.7),

$$P(A) = \mathbb{E}[A] \leq \frac{1}{x} \mathbb{E}[S_{n+1} I_{\mathcal{A}}].$$

Exercise: Show $\mathbb{E}[S_n| \mathcal{F}_n] \geq S_{n-1}$.

Proof: $\mathcal{F}_0 = \{\emptyset, \Omega\}, \mathcal{F}_n = \sigma(X_1, \dots, X_n)$ seq- of incr- σ -fields is filtration.

Def. rand-time $\tau \in \mathbb{R}^+ \cup \{\infty\}$ is stopping time wrt $\{\mathcal{F}_n\}$ if $\{\tau = n\} \in \mathcal{F}_n$.

Fact. i) $\{\tau = n\} \in \mathcal{F}_n$ ii) $\{\tau \leq n\} \in \mathcal{F}_n$ iii) $\{\tau \geq n+1\} \in \mathcal{F}_n$ are equivalent.

Fact. i) $\tau_1 \vee \tau_2$ ii) $\tau_1 + \tau_2$ iii) $\tau_1 \wedge \tau_2$ are stopping times.

Thm 4.1.5. (Wald's eq) i.i.d. if $E[\mathbf{X}|\mathcal{F}_\tau] \leq \mathbf{X}[\tau]$, then $E[S_\tau] = E[\mathbf{X}_1]E[\tau]$.

Thm 4.1.6. (W's 2eq) i.i.d. $\sim (0, \sigma^2)$, if $E[\tau] < \infty$, then $E[S_\tau^2] = \sigma^2 E[\tau]$.

Thm 4.2.3. SRW is recurrent in \mathbb{R}^1 and \mathbb{R}^2 and is transient in $\mathbb{R}^d, d \geq 3$.

Proof. In \mathbb{R}^1 ,

$$\sum_{m=1}^{\infty} P(S_m = 0) = \sum_{n=1}^{\infty} P(S_{2n} = 0)$$

(can only return to 0 a

$$\sum_{n=1}^{\infty} \binom{2n}{n} \cdot \frac{1}{2}^{2n}$$

)

$$\approx \sum_{n=1}^{\infty} \frac{\sqrt{2\pi n} \binom{2n}{n}}{(\sqrt{2\pi n})^n} \frac{1}{2^{2n}} \\ = \sum_{n=1}^{\infty} \frac{1}{\sqrt{\pi n}} = \infty;$$

$$R(\theta, \delta) = E_{X \sim P(\theta)}(L(\theta, \delta(x)))$$

Data Reduction
Sufficient:

- def. conditional distribution $[X | T=t]$ doesn't depend on θ .
- B&C 6.2.2. if $p(x|\theta)/q(T(x)|\theta)$ is free of θ , then $T(X)$ is suff.
- NFFC. $T(X)$ is sufficient if.f. $p_\theta(x) = g_\theta(T(x))h(x)$.

Assuming γ is a counting measure as a countable set \mathcal{X} (i.e. X is discrete) Let the family of pmf's be given by $\{P_\theta : \theta \in \Theta\}$ discrete)
 \Leftarrow : Suppose $p_\theta(x) = g_\theta(T(x))h(x)$, $\forall x \in \mathcal{X}$. Need to show that T is sufficient.

$$P_\theta(X=x | T(X)=t) = \frac{P_\theta(X=x, T(X)=t)}{P(T(X)=t)}$$

$$\begin{aligned} &= \left\{ \begin{array}{ll} 0 & T(x) \neq t \\ \frac{P(T(X)=t)}{P(X) \neq t} & T(X)=t \end{array} \right. \\ &= \left\{ \begin{array}{ll} 0 & T(x) \neq t \\ g_\theta(t)h(x) & T(x)=t \end{array} \right. \\ &= \left\{ \begin{array}{ll} 0 & T(x) \neq t \\ h(x) & T(x)=t \\ \sum_{y \in x: T(y)=t} h(y) & \perp \theta \end{array} \right. \end{aligned}$$

$$\Rightarrow \text{Suppose is sufficient for } \theta, \text{ so } P_\theta(X=x) = P_\theta(X=x, T(X)=t) = P_\theta(X=x | T(X)=t) = p_\theta(x) \text{ as } P_\theta(X=x | T(X)=t) \text{ is free of } \theta \text{ by definition}$$

Minimal sufficient:

- def. sufficient T is min.suff. if T is function of any other suff. T' .
- if $p(x|\theta) = c_x y p(y|\theta) \Leftrightarrow T(x)=T(y)$, then T is min.suff.

We first prove that T is sufficient. Start with $T(X) = \{t : t = T(x) \text{ for some } x \in \mathcal{X}\}$

= range of T . For each $t \in T(X)$, we consider the preimage $A_t = \{x : T(x)=t\}$ and select an arbitrary representative x_t from each A_t . Then, for any $y \in \mathcal{X}$, we have $y \in A_{T(y)}$ and $X_{T(y)} \in A_{T(y)}$. By the definition of A_t , this implies that $T(y) = T(X_{T(y)})$. From the assumption of the theorem, $p(y|\theta) = c_y x_{T(y)} p(x_{T(y)}|\theta) = h(y)g_\theta(T(y))$ which yields sufficiency of T by the NFFC.

Consider another sufficient statistic T' . By NFFC, $p(x|\theta) = \tilde{g}_\theta(T'(x))\tilde{h}(x)$
Take any x, y such that $T'(x)=T'(y)$, then $p(x|\theta) = \tilde{g}_\theta(T'(x))\tilde{h}(x) = \tilde{g}_\theta(T'(y))\tilde{h}(y) = \frac{\tilde{h}(x)}{\tilde{h}(y)}$
Hence, $T(x)=T(y)$ by the assumption of the theorem. So $T'(x)=T'(y)$ implies $T(x)=T(y)$ for any sufficient statistic T' and any x and y . As a result, T is a minimal sufficient statistic.

Complete: (complete suff. \Rightarrow minimal suff.)

- def. V is ancillary if the distribution of V is free of θ .
- def. T is complete if $E_\theta[f(T)]=0$ for all θ implies $f(T)=0$ a.e.

Consider $X_1, \dots, X_n \stackrel{i.i.d.}{\sim} \text{CauchyLoc}(\theta)$. Then $A(X)=X(n)-X(1)$ is ancillary even though (X_1, \dots, X_n) is minimal sufficient. To see this, note that $X_i = Z_i + \theta$ for $Z_i \stackrel{i.i.d.}{\sim} \text{CauchyLoc}(\theta)$, we can see that $X_i = Z_i + \theta$ and $A(X)=A(Z) \perp \!\!\! \perp \theta$.

- def. V is first-order ancillary if $E_\theta[V]$ is free of θ .

Let $X_1, \dots, X_n \stackrel{i.i.d.}{\sim} \text{Bernoulli}(\theta)$, $\theta \in (0,1)$. Then $T(X) = \sum_{i=1}^n X_i$ is sufficient.

Suppose $E_\theta[f(T(X))]=0$ for all $\theta \in (0,1)$,

$$\sum_{j=0}^n f(j) \binom{n}{j} \theta^j (1-\theta)^{n-j} = 0, \quad \forall \theta \in (0,1)$$

Dividing both sides by θ^n and reparameterizing $\beta = \frac{\theta}{1-\theta}$ we can rewrite it as

$$\sum_{j=0}^n f(j) \binom{n}{j} \beta^j = 0, \quad \forall \beta > 0$$

if f is non-zero, then LHS is a polynomial of degree at most n . However, an n -th-degree polynomial has at most n roots. Hence, it is impossible for the LHS to be equal to 0 for every $\beta > 0$ unless $f=0$. Therefore, T is complete.

$$\begin{aligned} &\text{Bernoulli}(p) \{p, (1-p), (1-p)^2\} \rightarrow \text{Bernoulli}(np) \{n, (n-1)\} \\ &\text{Poiss}(n) \{n, \lambda, \exp(-\lambda), 1\}, P(X=x) = \frac{\lambda^x}{x!} e^{-\lambda} \\ &\text{Geometric}(p) \{1, p, 1-p, (1-p)^2\}, P(X=x) = p(1-p)^{x-1} \\ &\text{H(N,M,K)} \{1, p, q, pq, (pq)^2, \dots\} \\ &\text{NB}(np) \{1, np, (np)^2, (np)^3, \dots\} \\ &\text{Beta}(\alpha, \beta) \{1, \alpha/\beta, (\alpha/\beta)^2, (\alpha/\beta)^3, \dots\} \\ &\text{Gamma}(n, \sigma) \{1, n\sigma, (n\sigma)^2, (n\sigma)^3, \dots\} \\ &\text{Exp}(\mu, \sigma) \{1, \mu, \mu^2, \sigma^2/2, \dots\} \\ &\text{Unif}(a, b) \{1, (a+b)/2, (a+b)^2/3, (a+b)^3/4, \dots\} \\ &\text{Logit}(p, \sigma) \{1, p, p^2, (1-p)^2, (1-p)^3, \dots\} \\ &\text{N}(p, \sigma^2) \{1, p, \sigma^2, \exp(-\sigma^2/2), \dots\} \\ &t_n \cdot \frac{1}{n!} \cdot \frac{(n\sigma)^n}{(n\sigma)^2} \cdot \frac{1}{2} \cdot \frac{(n\sigma)^2}{2} \cdot \dots \cdot \frac{1}{(n-1)\sigma^2} \cdot \frac{1}{2} \cdot \frac{(n-1)\sigma^2}{2} \cdot \dots \cdot \frac{1}{\sigma^2} \end{aligned}$$

Let $X_1, \dots, X_n \stackrel{i.i.d.}{\sim} \text{N}(\theta, \sigma^2)$ with unknown θ and an known $\sigma^2 > 0$. Is $\bar{X}_n = n^{-1} \sum_{i=1}^n X_i$ complete for this model?

Consider the special case of $n=1$ and $\sigma=1$. $T(X)=X \sim \text{N}(\theta, 1)$

$$E_\theta[f(X)]=0, \quad \forall \theta \in \mathbb{R}$$

Multiplying both sides by $\sqrt{2\pi/\theta^2}$, we have $\int_{-\infty}^{\infty} f(x) \exp\left(-\frac{x^2}{\theta^2}\right) dx = 0, \quad \forall \theta \in \mathbb{R}$. We decompose f into its positive and negative parts as $f(x) = f_+(x) - f_-(x)$, where $f_+(x) = \max(f(x), 0)$, and $f_-(x) = \max(-f(x), 0)$. Note that $f_+ \geq 0$ and $f_- \geq 0$. For all $x \in \mathbb{R}$ $f_+(x) = f(x)$ if and only if $f_+(x) = f_-(x) = 0$. Suppose f_+ and f_- have non-zero components, and we may write

$$\int_{-\infty}^{\infty} f_+(x) e^{-\frac{x^2}{\theta^2}} dx = \int_{-\infty}^{\infty} f_-(x) e^{-\frac{x^2}{\theta^2}} dx = \int_{-\infty}^{\infty} f_-(x) e^{-\frac{x^2}{\theta^2}} dx$$

Note that

$$\begin{aligned} &f_+(x) e^{-\frac{x^2}{\theta^2}} \\ &- f_-(x) e^{-\frac{x^2}{\theta^2}} \end{aligned}$$

defines a probability density.

The equality of the mgfs implies equality of the densities, which in turn implies $f_+(x) = f_-(x)$ a.e.. Then $f_+(x) = f_-(x) = 0$ a.e., or in other words, $f(x) = 0$ a.e.. Hence T is complete.

Example: Let $X_1, \dots, X_n \stackrel{i.i.d.}{\sim} \text{CauchyLoc}(\theta)$ where both μ, σ^2 are unknown. Then, $\bar{X}_n = n^{-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2$ where

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

Fix any $\sigma > 0$ and consider sub-model $P_\sigma = \{N(\mu, \sigma^2)\}$. In each sub-model, \bar{X}_n is complete and sufficient, and $n^{-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2$ is ancillary. By Basu's theorem, $\bar{X}_n \perp \!\!\! \perp \sum_{i=1}^n (X_i - \bar{X}_n)^2$ under $N(\mu, \sigma^2)$ for any $\mu \in \mathbb{R}$. Since σ is arbitrary, the conclusion holds for the full model $P = \{N(\mu, \sigma^2) : \mu \in \mathbb{R}, \sigma > 0\}$.

5. (Rao-Blackwell) for conv. loss and suff. T , $R(\theta, E[\delta | T]) \leq R(\theta, \delta)$, $E(\delta | L(\eta|T), \delta | \eta|T)) \geq E(\delta | L(\eta|T), \eta|T))$

6. common steps:

- suppose $\int f(x)h(x)e^{\theta x} dx = 0$ for all $\theta \in \Omega$,
- decompose $f = f_+ - f_-$ with $f_+ + f_- \geq 0$,
- view f_+ and f_- as un-normalised densities p_+ and p_- ,
- argue that MGf of p_+ and p_- are equal, then $f_+ = f_-$ a.e.

$$\text{Example: } X \sim P(\theta); g(\theta) = e^{-\theta\theta}$$

UMRUE: $E_\theta(g(x)) = \sum_{x \in \mathcal{X}} g(x) \frac{e^{-\theta x}}{1-e^{-\theta}}$

$$\sum_{x \in \mathcal{X}} \frac{e^{-\theta x}}{1-e^{-\theta}} = e^{-\theta} \frac{1}{1-e^{-\theta}} = \frac{e^{-\theta}}{e^{-\theta} - 1} = \frac{1}{e^{-\theta}}$$

$$\text{Hence } g(x) = \frac{1}{e^{-\theta}}$$

- def. $R(\theta, \delta) \leq R(\theta, \delta')$ for $\forall \theta \in \Omega$ and \forall unbiased δ' .

- (Lehmann-Scheffe) if T is comp.suff. and $E_\theta[g(\cdot)] = g(\theta)$, then $h(T)$ is i) only unbiased fun. of T , ii) UMRUE under conv.loss. (unique UMVUE under Rao-Blackwellisation is very tedious. Instead, we can observe another fact that

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

And $S_*^2 \sim \sigma^2 \chi_{n-1}^2$. Hence, $E(S_*) = \sigma E(\chi_{n-1})$, which implies that

$$\frac{E(S_*)}{E(\chi_{n-1})} = \sigma$$

meaning that $\frac{S_*}{E(\chi_{n-1})}$ is unbiased for σ and hence UMVU.

For 3), Taking the expectation of the UMVUE for μ and squaring it, we obtain

$$E(\bar{X}_n^2) = \mu^2 + \frac{\sigma^2}{n}$$

So $\delta_n(X) = \bar{X}_n^2 - \frac{S_*^2}{n(n-1)}$

is the UMVUE. However, $\delta_n(X)$ can be negative even though the estimand is a non-negative quantity. The estimator is in fact inadmissible and dominated by the biased estimator $\max(0, \delta_n(X))$ since $E[\max(0, \delta_n(X))] \neq \max(E(\delta_n(X)), 0)$.

Example 5 (Semiparametric unbiased estimators).

For $n > 2$, let $X_1, \dots, X_n \stackrel{i.i.d.}{\sim} F$, where F is some unknown distribution on \mathbb{R} . Suppose that F is symmetric about some unknown point $\theta \in \mathbb{R}$. That is, suppose for all $X \sim F$, we have $X = d\theta - X$.

Consider the model $\mathcal{F} = \{\text{all distributions on } \mathbb{R} \text{ with finite variance symmetric about } \theta\}$.

Then there is no UMVUE for the point of symmetry θ . Proof. Suppose for sake of contradiction that the UMVUE $T(X)$ exists. Since \bar{X} is unbiased for the full model \mathcal{F} , $T(X)$ must have variance no larger than \bar{X} . However, we know that \bar{X} is the unique UMVUE for the Gaussian submodel, $\{N(\theta, 1) : \theta \in \mathbb{R}\}$, and so $T(X)$ must equal \bar{X} a.s. in the Gaussian submodel. This implies that $T(X) = \bar{X}$ a.s. under any continuous distribution on \mathbb{R} . In particular, $T(X) = \bar{X}$ a.s. under the uniform submodel $\{\text{Unif}(\theta-1, \theta+1) : \theta \in \mathbb{R}\}$. However, we learned in Homework 2, Problem 3 (b), that the distinct estimator $\frac{1}{2}(X_1 + X_n)$ has strictly better variance than \bar{X} in the uniform submodel. Since $\frac{1}{2}(X_1 + X_n)$ is also unbiased for the full model \mathcal{F} , $T(X)$ cannot be the UMVUE after all. We are forced to conclude that no UMVUE exists over the whole family.

Fisher Information:

$$\text{Var}(\frac{\partial \log f(x)}{\partial \mu})$$

1. $I(\theta) = E\left(\frac{\partial}{\partial \theta} \log f_\theta(x)\right)^2 = -E\left[\frac{\partial^2}{\partial \theta^2} \log f_\theta(x)\right]$.

2. Cramer-Rao lower bound: $\text{Var}(\delta) \geq [g'(\theta)]^2 / I(\theta)$.

$$\varphi(x) = \frac{\partial}{\partial \theta} \log f_\theta(x), E_\theta[\varphi(x)] = 0, E_\theta[\delta^2] < \infty, g'(\theta) = E_\theta[\delta \varphi]$$

3. If $X_1, \dots, X_n \stackrel{i.i.d.}{\sim} P_\theta$, $I_X(\theta) = I_{X_1}(\theta + \dots + I_{X_n}(\theta)) = nI_{X_1}(\theta)$.

Let \mathcal{P} be a one-parameter exponential family in canonical form and density p_θ given by $p_\theta = \exp(\eta(x) - A(\eta))h(x)$, then

$$\frac{\partial \log p_\theta(x)}{\partial \eta} = T(x) - A'(\eta).$$

By the previous results, we have $I(\eta) = \text{Var}_{p_\theta}(T(x) - A'(\eta)) = \text{Var}_{p_\theta}(T(x)) = A''(\eta)$.

If the family is parameterised instead by $\mu = A'(\eta) = \eta T(x)$, then

$$A''(\eta) = I(\mu)(A''(\eta))^2,$$

and so, because $A''(\eta) = \text{Var}(T)$, we have $I(\mu) = \frac{1}{\text{Var}(T)}$. Observe also that because T is UMVUE for μ , the lower bound variance

$\text{Var}(\delta) \geq (I(\mu))^{-1}$ for unbiased estimation δ of μ is sharp in this example.

Suppose X is ab. absolutely continuous random variable with density f . The family of distributions $\mathcal{P} = \{P_\theta : \theta \in \mathbb{R}\}$ with P_θ the distribution of $\theta + \epsilon$ is called a location family.

$$\begin{aligned} \int g(x)dP_\theta(x) &= E_\theta(g(X)) \\ &= E_\theta(g(\theta + \epsilon)) \\ &= \int g(x)f(x - \theta)d\epsilon \\ &= \int g(x)f(x)d\epsilon. \end{aligned}$$

So P_θ has the density $p_\theta = f(x - \theta)$. The corresponding Fisher information for the family is

$$\begin{aligned} I(\theta) &= E_\theta\left(\frac{\partial \log f(X - \theta)}{\partial \theta}\right)^2 \\ &= E_\theta\left(\frac{f'(X - \theta)}{f(X - \theta)}\right)^2 \\ &= E\left(\frac{f'(\epsilon)}{f(\epsilon)}\right)^2 \\ &= \int \frac{\{f'(x)\}^2}{f(x)}dx \perp \!\!\! \perp \theta. \end{aligned}$$

So, for the location families, $I(\theta)$ is constant with respect to θ .

Suppose $X_1, \dots, X_n \stackrel{i.i.d.}{\sim} N(\mu, \sigma^2), \mu \in \mathbb{R}$, and $\sigma^2 > 0$. We want to estimate $g_1(\mu, \sigma^2) = \mu$ and $g_2(\mu, \sigma^2) = \sigma^2$. Look at only unbiased estimators.

Claim 1: $\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2$ is UMVUE for $\text{Var}(X)$.

Proof. (a) $(\sum_{i=1}^n X_i, \sum_{i=1}^n X_i^2)$ is complete and sufficient.

(b) $\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2$ is a function of $(\sum_{i=1}^n X_i, \sum_{i=1}^n X_i^2)$. This is because $\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n) = \sum_{i=1}^n X_i^2 - n\bar{X}_n^2$.

(c) $\frac{1}{\sigma^2} \sum_{i=1}^n X_i^2 \rightarrow E\left(\frac{1}{\sigma^2} \sum_{i=1}^n X_i^2\right) = \frac{\sigma^2}{n-1} (n-1) = \sigma^2$

Claim 2: \bar{X}_n is UMVUE for μ . $E(\bar{X}_n - \mu)^2 = \sigma^2/n$. Note that $E\left(\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2 / (n-1) - \sigma^2\right)^2 = \text{Var}\left(\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2 / (n-1)\right) = 2\frac{(n-1)\sigma^4}{(n-1)^2} = \frac{2\sigma^4}{(n-1)}$.

$$\log p_{\mu, \sigma^2}(x) = -\frac{n}{2} \log \sigma^2 - \sum_{i=1}^n \frac{(x_i - \mu)^2}{2\sigma^2} + C_n, \quad \frac{\partial \log p_{\mu}(x)}{\partial \mu} = \sum_{i=1}^n \frac{x_i - \mu}{\sigma^2},$$

$$\frac{\partial^2 \log p_{\mu}(x)}{\partial \mu^2} = -\frac{1}{\sigma^2} + \sum_{i=1}^n \frac{1}{\sigma^2},$$

$$\frac{\partial \log p_{\mu}(x)}{\partial \sigma^2} = \frac{n}{\sigma^4} - \sum_{i=1}^n \frac{(x_i - \mu)^2}{\sigma^4}$$

$$\frac{\partial^2 \log p_{\mu}(x)}{\partial \mu \partial \sigma^2} = -\sum_{i=1}^n \frac{x_i - \mu}{\sigma^4}$$

Therefore, the Fisher information matrix is

$$I = \begin{pmatrix} \frac{n}{\sigma^2} & 0 \\ 0 & \frac{n}{\sigma^4} \end{pmatrix}$$

$$I_{22} = -\frac{n}{2\sigma^4} + \mathbb{E}\left(\sum_{i=1}^n \frac{(X_i - \mu)^2}{\sigma^6}\right) = -\frac{n}{2\sigma^4} + \frac{n\sigma^2}{\sigma^6} = \frac{n}{2\sigma^4}$$

$$\Rightarrow \text{CRLB for } \mu = (1 \ 0) I^{-1}(\frac{1}{\sigma^2}) = \frac{1}{11} = \frac{\sigma^2}{n}$$

$$\Rightarrow \text{CRLB for } \sigma^2 = (0 \ 1) I^{-1}(\frac{1}{\sigma^4}) = \frac{1}{22} = \frac{2\sigma^4}{n}$$

Exponential Family

general form: $p(x;\theta) = \exp\{\sum_{i=1}^n \eta_i(\theta) T_i(x) - B(\theta)\} h(x)$.

1. standardiser: $B(\theta) = \log f \exp\{\sum_{i=1}^n \eta_i(\theta) T_i(x)\} h(x)$
2. parameter space: $\Theta = \{\theta : B(\theta) < \infty\}$.
3. Moments Take $f(x)=1$, then

$$\begin{aligned} \frac{\partial A(\eta)}{\partial \eta_i} &= \int T_i(x) \exp\left(\sum_{j=1}^s \eta_j T_j(x) - A(\eta)\right) h(x) d\mu(x) \\ &= E_{\theta}[T_i(X)] \left(\text{harm. } \frac{1}{\sum_j \exp(A(j))} \right) \\ \frac{\partial^2 A(\eta)}{\partial \eta_i \partial \eta_j} &= \text{Cov}_{\eta}\left(T_i(X), T_j(X)\right) \end{aligned}$$

canonical form: $p(x;\eta) = \exp\{\sum_{i=1}^n \eta_i T_i(x) - A(\eta)\} h(x)$.

1. natural parameter η_i , and natural parameter space.
 2. def. canonical exp.fam. is minimal, if no affine T_i 's and η_i 's.
 $(\sum_i \lambda_i T_i(x) = \lambda_0$ implies $\lambda_i = 0$, similar for η_i 's.)
 3. def. min.exp.fam. is full-rank, if nat.par.space contain open rect.
 4. if exp.fam. is full-rank, then T is minimal sufficient and complete. $(\sum_i T_i(X) \sim \sum_i T_i(Y))$
- Hypothesis Testing
Basic definition:

1. test function $\phi(x)$: the prob. rejects H_0 given $X=x$,
2. power function: $\beta(\theta) = E_{\theta}[\phi(X)] = P_{\theta}(\text{rejects } H_0)$,
3. significant level $\alpha: \sup_{\theta \in \Omega_1} E_{\theta}[\phi(X)] \leq \alpha$,
4. level- α uniformly most powerful test ϕ : if $E_{\theta}[\phi(X)] \geq E_{\theta}[\phi^*(X)]$ for all $\theta \in \Omega_1$, for any other level- α test ϕ^* .
5. families with monotone likelihood ratio in $T(X)$:
 - (a) $\theta \neq \theta'$ implies $P_{\theta} \neq P_{\theta'}$, (identifiability)
 - (b) $\theta < \theta'$ implies the ratio $P_{\theta}(x)/P_{\theta'}(x)$ is a non-decreasing function of $T(X)$. (monotonicity)

Proof: For $\alpha=0$ and $\alpha=1$ the theorem is easily seen to be true provided the value $k=\pm\infty$ is admitted in (3.8) and $0\in\alpha$ is interpreted as 0 . Throughout the proof we shall therefore assume $0<\alpha<1$.
(i): Let $\alpha(c)=P_0\{p_1(X)>c p_0(X)\}$. Since the probability is computed under P_0 , the inequality need be considered only for the set where $p_0(x)>0$, so that $\alpha(c)$ is the probability that the random variable $p_1(X)/p_0(X)$ exceeds c . Thus $1-\alpha(c)$ is a cumulative distribution function, and $\alpha(c)$ is nonincreasing and continuous on the right, $\alpha(c-0)-\alpha(c)=P_0\{p_1(X)/p_0(X)=c\}$, $\alpha(-\infty)=1$, and $\alpha(\infty)=0$. Given any $0<\alpha<1$, let c_0 be such that $\alpha(c_0)\leq\alpha\leq\alpha(c_0-0)$, and consider the test ϕ defined by

$$\phi(x) = \begin{cases} 1 & \text{when } p_1(x) > c \\ 0 & \text{when } p_1(x) \leq c_0 \\ \frac{\alpha - \alpha(c_0)}{\alpha(c_0-0) - \alpha(c_0)} & \text{when } c_0 < p_1(x) < c \end{cases}$$

Here the middle expression is meaningful unless $\alpha(c_0)=\alpha(c_0-0)$; since then $P_0\{p_1(X)=c_0 p_0(X)\}=0$, ϕ is defined a.e.

The size of ϕ is

$$E_0 \phi(X) = P_0\left\{p_1(X) > c_0\right\} + \frac{\alpha - \alpha(c_0)}{\alpha(c_0-0) - \alpha(c_0)} P_0\left\{p_1(X) = c_0\right\} = \alpha,$$

so that c_0 can be taken as the k of the theorem.

(ii): Suppose that ϕ is a test satisfying (3.7) and (3.8) and that ϕ^* is any other test with

$E_0 \phi^*(X) \leq \alpha$. Denote by S^+ and S^- the sets in the sample space where $\phi(x) = \phi^*(x) > 0$ and < 0 respectively. If x is in S^+ , $\phi(x)$ must be > 0 and $p_1(x) \geq k p_0(x)$. In the same way $p_1(x) \leq k p_0(x)$ for all x in S^- , and hence $\int (\phi - \phi^*) (p_1 - k p_0) d\mu = \int S^+ + \int S^- (\phi - \phi^*) (p_1 - k p_0) d\mu \geq 0$

The difference in power between ϕ and ϕ^* therefore satisfies $\int (\phi - \phi^*) p_1 d\mu \geq \int (\phi - \phi^*) p_0 d\mu \geq 0$ as was to be proved.

(iii): Let ϕ^* be most powerful at level α for testing p_0 against p_1 , and let ϕ satisfy (3.7)

and (3.8). Let S be the intersection of the set $S^+ \cup S^-$, on which ϕ and ϕ^* differ, with the set $\{x : p_1(x) \neq k p_0(x)\}$, and suppose that $\mu(S) > 0$. Since $(\phi - \phi^*)(p_1 - k p_0)$ is positive on S , it follows from Problem 2.4 that $\int_{S^+ \cup S^-} (\phi - \phi^*)(p_1 - k p_0) d\mu = \int_S (\phi - \phi^*)(p_1 - k p_0) d\mu > 0$ the desired significance level is $\alpha = 50\%$. Let $\phi(x_1, x_2) = \begin{cases} 1 & \text{if } (x_1, x_2) = (1, 1) \\ 0 & \text{if } (x_1, x_2) = (0, 0) \end{cases}$, randomised when we observe $(0, 1)$ and $(1, 0)$, such that $E_{\theta_0}[\phi(X_1, X_2)] = \frac{1}{2}$. Example (One parameter exponential family). Consider the case where X_1, \dots, X_n i.i.d. $p_{\theta}(x) \propto h(x) \exp(\theta T(x))$, and we are interested in testing $H_0: \theta = \theta_0$ vs. $H_1: \theta = \theta_1$.

Example:

Let X_1, \dots, X_n i.i.d. $N(\mu, \sigma^2)$, with σ known. Consider the following two hypotheses:

$$H_0: \mu = 0 \text{ and } H_1: \mu = \mu_1$$

where μ_1 is known. Since this is a simple case (only two distributions), we calculate the likelihood ratio:

$$\begin{aligned} r(x) &= \frac{p_1(x)}{p_0(x)} = \frac{\prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma}} \exp\left(-(x_i - \mu_1)^2 / 2\sigma^2\right)}{\prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma}} \exp\left(-(x_i - \mu)^2 / 2\sigma^2\right)} \\ &= \exp\left(\frac{1}{\sigma^2} \mu_1 \sum_{i=1}^n x_i - \frac{n\mu_1^2}{2\sigma^2}\right). \end{aligned}$$

We observe that the likelihood ratio is only a function of the sufficient statistic (this is in general true by the Factorization Criterion). We have the following equivalences:

$$\begin{aligned} r(x) > k &\Leftrightarrow \mu_1 \sum_{i=1}^n x_i - \frac{n\mu_1^2}{2\sigma^2} > \log k \\ &\Leftrightarrow \mu_1 \sum_{i=1}^n x_i > k' \\ &\Leftrightarrow \begin{cases} \sum_{i=1}^n x_i > k'' & \text{if } \mu_1 > 0 \\ \sum_{i=1}^n x_i < k''' & \text{if } \mu_1 < 0 \end{cases} \end{aligned}$$

Let us focus on the first case where $\mu_1 > 0$ (but it is important to note that the two cases $\mu_1 > 0$ versus $\mu_1 < 0$ induce different rejection regions). We can rewrite the test in a different form so that the left hand side of the inequality has standard normal distribution under the null:

$$\Leftrightarrow \frac{\sqrt{n}\bar{x}}{\sigma} > k''''$$

Note: Here the sufficient statistic essentially determines a MP test. Also, observe that for a given level α , the constant involved in the LRT is uniquely determined by the constraint $E_{\theta_0}[\phi(X)] = \alpha$ and thus depends only on the distribution of the null hypothesis and not at all on μ_1 (provided that $\mu_1 > 0$).

For ϕ we have by the NP lemma that the test: reject H_0 iff $\sqrt{n}\bar{x}/\sigma > k(\alpha)$ is MP where we pick $k(\alpha) = \Phi^{-1}(\alpha)$ (where $k(\alpha) \equiv z_{1-\alpha}$ is the $(1-\alpha)$ quantile of the standard normal) so that it equals the target level. This value is uniquely determined by the size constraint:

$$E_0 \phi(X) = P_0\left\{p_1(X) > c_0\right\} + \frac{\alpha - \alpha(c_0)}{\alpha(c_0-0) - \alpha(c_0)} P_0\left\{p_1(X) = c_0\right\} = \alpha,$$

Now an important observation: for any $\mu_1 > 0$ the MP test is the same, which means that it is actually a UMP at level α for testing:

$$H_0: \mu = 0 \text{ and } H_1: \mu > 0.$$

So we see no μ_1 dependence here and it's really nice that I can test against any $\mu_1 > 0$ at once. Similarly, we could derive a distinct UMP test for testing against $H_1: \mu < 0$. Unfortunately, no UMP test exists for testing

$$H_0: \mu = 0 \text{ and } H_1: \mu \neq 0$$

because the $\mu > 0$ test dominates the $\mu < 0$ test in the $\mu > 0$ scenario and vice versa i.e. the shape of the most powerful rejection (whether we reject for $\bar{x} > k$ or whether we reject for $\bar{x} < k$) is dependent on the sign of μ_1 . Example: Suppose X has a binomial distribution with success probability θ and $n=2$ trials. If we are interested in testing

$H_0: \theta = 1/2$ versus $H_1: \theta = 2/3$, then

$$L(x) = \frac{p_1(x)}{p_0(x)} = \frac{\binom{2}{x}}{\binom{2}{x}} \frac{(2/3)^x (1/3)^{2-x}}{(1/2)^x (1/2)^{2-x}} = \frac{2x}{9} \times 4.$$

and hence, we can restrict our attention to the sufficient statistics (Y, U) , where $Y = \bar{X}$ and $U = \sum_{i=1}^n (X_i - \bar{X})^2$. We know that $Y \sim N(\theta, \sigma^2/n)$, $U \sim \sigma^2 \chi_{n-1}^2$, and Y is independent of U by Basu's theorem. Thus, for Λ supported on $\sigma = \sigma_0$, we obtain the joint density of (Y, U) under H_A as

$$c_0 u^{\frac{n-3}{2}} \exp\left(-\frac{u}{2\sigma_0^2}\right) \int \exp\left(-\frac{y - \theta}{2\sigma_0^2}\right) d\Lambda(y)$$

and the joint density under alternative hypothesis (θ_1, σ_1) as

$$c_1 u^{\frac{n-3}{2}} \exp\left(-\frac{u}{2\sigma_1^2}\right) \exp\left(-\frac{y - \theta_1}{2\sigma_1^2}\right).$$

From the above observations, we see that the choice of Λ only affects the distribution of Y . To achieve minimal maximum power against the alternative (i.e., to be least favorable), we need to choose Λ such that the two distributions become as close as possible. Under the alternative hypothesis, $Y \sim N(\theta_1, \frac{\sigma_1^2}{n})$. Under H_A , the distribution of Y is in a convolution form, i.e., $Y = Z + \Theta$ for $Z \sim N(0, \frac{\sigma_0^2}{n})$, $\Theta \sim \Lambda$, where Z and Θ are independent.

Example 2 (Double exponential). Let $X \sim \text{DoubleExponential}(\theta)$, with density $p_{\theta}(x) = \frac{1}{2} \exp(-|x - \theta|)$. It is easy to see that the model is identifiable, so we need only check the second condition. Fix any $\theta' > \theta$ and consider the likelihood ratio

$$\frac{p_{\theta'}(x)}{p_{\theta}(x)} = \exp(|x - \theta| - |x - \theta'|).$$

Note that

$$|x - \theta| - |x - \theta'| = \begin{cases} \theta - \theta' & \text{if } x < \theta \\ 2x - \theta - \theta' & \text{if } \theta \leq x \leq \theta' \\ \theta' - \theta & \text{if } x > \theta' \end{cases}$$

which is non-decreasing in x . Therefore, the family has MLR in $T(x)=x$. Finally, we give an example of a model that does not exhibit MLR in $T(x)=x$.

Example 3 (Cauchy location model). Let X have density $p_{\theta}(x) = \frac{1}{\pi} \frac{1}{1+(x-\theta)^2}$. We find two points for which the MLR condition fails. For any fixed $\theta > 0$,

$$\frac{p_{\theta}(x)}{p_{\theta}(0)} = \frac{1+x^2}{1+(x-\theta)^2} \rightarrow 1 \quad \text{as } x \rightarrow \infty \text{ or } x \rightarrow -\infty,$$

but $p_{\theta}(0)/p_{\theta}(0) = 1/(1+\theta^2)$, which is strictly less than 1. Thus the ratio must increase at some values of x and decrease at others. In particular, it is not monotone in x . Here we have shown that the likelihood ratio in $T(x)=x$ is not MLR.

Theorem 1 (TSH 3.8.1). Suppose ϕ_{Λ} is a MP level- α test for testing H_{Λ} against g . If ϕ_{Λ} is level- α for the original hypothesis H_0 (i.e., $E_{\theta_0}[\phi_{\Lambda}(x)] \leq \alpha, \forall \theta_0 \in \Omega_0$), then 1. The test ϕ_{Λ} is MP for original $H_0: \theta \in \Omega_0$ vs. g . 2. The distribution Λ is least favorable.

Proof. 1. Let ϕ^* be any other level- α test of $H_0: \theta \in \Omega_0$ vs. g . Then ϕ^* is also a level- α test for H_{Λ} vs. g , which implies that

$$\int \phi^*(x) f_{\Lambda}(x) d\mu(x) = \int \phi^*(x) f_{\theta}(x) d\mu(x) \leq \alpha, \forall \theta \in \Omega_0.$$

Since ϕ_{Λ} is MP for H_{Λ} vs. g , we have

$$\int \phi_{\Lambda}(x) g(x) d\mu(x) \leq \int \phi_{\Lambda}(x) g(x) d\mu(x),$$

Hence ϕ_{Λ} is a MP test for H_0 vs. g , because ϕ_{Λ} is also level α . 2. Let Λ' be any distribution on Ω_0 . Since $E_{\theta}[\phi_{\Lambda}(x)] \leq \alpha, \forall \theta \in \Omega_0$, we know that ϕ_{Λ} must be level α for $H_{\Lambda'}$ vs. g . Thus $\beta_{\Lambda} \leq \beta_{\Lambda'}$, so Λ is the least favorable distribution.

Example 1 (Testing in the presence of nuisance parameters). Let X_1, \dots, X_n be i.i.d. $N(\theta, \sigma^2)$, where both θ, σ^2 are unknown. We consider testing $H_0: \theta \leq \theta_0$ against $H_1: \theta > \theta_0$. To find a UMP test, we follow the previously mentioned strategy:

1. First we fix a simple alternative (θ_1, σ_1) for some arbitrary θ_1 and $\sigma_1 > \theta_0$. 2. Second, we choose a prior distribution Λ to collapse our null hypothesis over. Intuitively, the least favorable prior should make the alternative hypothesis hard to distinguish. Hence, a rule of thumb consists in concentrating Λ on the boundary between H_1 and H_0 (i.e. the line $\{\sigma = \sigma_0\}$). Thus Λ will be a probability distribution over $\theta \in \mathbb{R}$ for the fixed $\sigma = \sigma_0$. Another useful observation is that, given any test function $\phi(x)$ and a sufficient statistic T , there exists a test function η that has the same power as ϕ but depends on x only through T :

$$\eta(T(x)) = E[\phi(x)|T(x)].$$

Here p_0 is a fixed quality parameter. Before we start our search for the UMP test, let us reparametrize the distribution \mathbb{P} as follows. Let \mathbb{P}^- and \mathbb{P}^+ be the conditional distributions of $X|X \leq u$ and $X|X > u$ respectively, and let $p = \mathbb{P}(X \leq u)$. Then, \mathbb{P} has a one-to-one correspondence with $(\mathbb{P}^+, \mathbb{P}^-)$. For any fixed \mathbb{P} , let \mathbb{P}^- and \mathbb{P}^+ be the conditional densities of \mathbb{P}^- and \mathbb{P}^+ with respect to some measure μ (existence of the densities and base measure can be justified, e.g. by Radon-Nikodym theorem in measure theory). The joint density of X_1, \dots, X_n at values

squared error loss considered here.
remark1 Remark 1. Recall that the Beta function can be evaluated as

$$\int_0^1 x^{k_1-1}(1-x)^{k_2-1}dx = \frac{\Gamma(k_1)\Gamma(k_2)}{\Gamma(k_1+k_2)}$$

whenever $k_1, k_2 > 0$. Therefore, if $Y \sim \text{Beta}(a, b)$, where $a, b > 0$, we can explicitly evaluate the expectation

$$\begin{aligned} E\left[\frac{1}{1-Y}\right] &= \int_0^1 \frac{1}{1-y} \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} y^{a-1}(1-y)^{b-1} dy \\ &= \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \cdot \frac{\Gamma(a)\Gamma(b-1)}{\Gamma(a+b-1)} \end{aligned}$$

where in the second step we require $b > 1$ in order to apply the relation 10.2. A similar argument yields

$$\begin{aligned} E\left[\frac{1}{Y(1-Y)}\right] &= \frac{(a+b-2)(a+b-1)}{(a-1)(b-1)}, \\ \text{ whenever } a > 1. \text{ Combining these identities, we have that, whenever } a, b > 1, \\ E\left[\frac{1}{Y(1-Y)}\right] &= \frac{a-1}{a+b-2}. \end{aligned}$$

10. Randomized Minimax Let $X \sim \text{Bin}(n, \theta)$, where $\theta \in [0, 1]$, and consider estimation of θ under the $0-1$ loss,

$$L(\theta, d) = \begin{cases} 0 & \text{if } |d-\theta| < \alpha \\ 1 & \text{otherwise} \end{cases}$$

First consider an arbitrary non-random estimator δ . Since X can take on only the $n+1$ values $\{0, 1, \dots, n\}$, the estimator $\delta(X)$ can take on only $n+1$ values, $\{\delta(0), \delta(1), \dots, \delta(n)\}$. If $\alpha < \frac{1}{2(n+1)}$, then we can always find θ_0 such that $|\delta(x) - \theta_0| \geq \alpha$ for every $x \in \{0, \dots, n\}$; see Figure 10.1. Hence, $R(\theta_0, \delta(X)) = 1$ is the maximum risk of any non-random δ .

Consider instead the estimator $\delta'(X, U) = U$, which is completely random and independent of the data X . Then, for any $\theta \in [0, 1]$,

$$\begin{aligned} R(\theta, \delta') &= E[L(\theta, \delta'(X, U))] \\ &= P(|U - \theta| \geq \alpha) \\ &= 1 - P(\theta - \alpha < U < \alpha + \theta) \\ &\leq 1 - \alpha < 1 \end{aligned}$$

and since $\alpha > 0$, the maximum risk of δ' is smaller than the maximum risk of any non-random δ . Hence, in this setting, there can be no deterministic minimax estimator.

- Lemma 1 (TPE 5.1.15). Suppose that δ is minimax for a submodel $\Omega_0 \subset \Omega$ and

$$\sup_{\theta \in \Omega_0} R(\theta, \delta) = \sup_{\theta \in \Omega} R(\theta, \delta)$$

Then, δ is minimax for the full model, $\theta \in \Omega$.

12. Example Let X_1, \dots, X_n be i.i.d $\mathcal{N}(\mu, \sigma^2)$, where both μ and σ^2 are unknown. Thus, our parameter vector, $\theta = (\mu, \sigma^2)$ and our parameter space $\Omega = \mathbb{R} \times \mathbb{R}^+$. Our task now is to estimate μ . Our loss function is the relative squared error loss, given by:

$$L((\mu, \sigma^2), d) = \frac{(\bar{d} - \mu)^2}{\sigma^2}$$

We consider the submodel where $\sigma^2 = 1$. That is, $\Omega_0 = \mathbb{R} \times \{1\}$, and our loss function simplifies to our usual squared error loss: $L((\mu, 1), d) = (\bar{d} - \mu)^2$. We saw in Example 1 of Lecture 10 that under this loss \bar{X} is minimax for Ω_0 . Moreover,

$$R((\mu, \sigma^2), \bar{X}) = \frac{1}{n} \quad \forall (\mu, \sigma^2) \in \Omega$$

Thus, the risk does not depend on σ^2 . Since $R((\mu, 1), \bar{X}) = R((\mu, \sigma^2), \bar{X})$, we have that the maximum risks are equal. (That is, $\sup_{\theta \in \Omega_0} R(\theta, \delta) = \sup_{\theta \in \Omega} R(\theta, \delta)$). Therefore, it follows from Lemma 1 that \bar{X} is minimax on Ω . Note that, thanks to our new loss function, we don't need to impose boundedness on our variance (like we did in our previous lecture) to establish minimaxity in a meaningful way.

13. Example 2 Suppose X_1, X_2, \dots, X_n are i.i.d with common CDF F , with mean $\mu(F) < \infty$, and variance $\sigma^2(F) < \infty$. Our goal is to find a minimax estimate of $\mu(F)$ under squared error loss. Constraint (a). Assume $\sigma^2(F) \leq B$.

Now, we've seen in the previous lecture that \bar{X} is minimax for the Gaussian submodel in this case. So a natural guess for us to make is that \bar{X} is minimax. We verify this by application of Lemma 1. First, we compute the supremum risk for the full model:

$$R(F, \bar{X}) = \frac{1}{n^2} \sum_i E(X_i - \mu(F))^2 = \frac{\sigma^2(F)}{n}$$

Since $\sigma^2(F) \in [0, B]$ by assumption, we get:

$$\begin{aligned} \sup_F R(F, \bar{X}) &= \frac{B}{n} \\ \text{Now we saw in Lecture 10 that for the submodel } \mathcal{F}_0 &= \mathcal{N}(\mu, \sigma^2) \text{ when } \sigma^2 \leq B, \bar{X} \text{ is minimax. Further, the supremum risk in this case is identical to that of the full model:} \\ E[(a\bar{X} + b - \theta)^2] &= E[(a(\bar{X} - \theta) + b + \theta(a-1))^2] \end{aligned}$$

Thus, using Lemma 1 we conclude that \bar{X} is minimax for the full model. (That is, the non-parametric model still constrained to have $\sigma^2(F) \leq B$.)

Constraint (b). Assume $F \in \mathcal{F}$ where \mathcal{F} is the set of all CDFs with support contained in $[0, 1]$, it turns out that \bar{X} isn't minimax. To show this, first consider the submodel, $\mathcal{F}_0 = \{\text{Ber}(\theta)\}_{\theta \in (0, 1)}$. Let $Y = \sum_{i=1}^n X_i$ so that $Y \sim \text{Bin}(n, \theta)$ and $X = Y/n$. Recall from Lecture 9 that the minimax estimator for $\mu(F) = \theta$, in the Binomial case, is:

$$\delta(X) = \frac{\sqrt{n}}{1+\sqrt{n}} \bar{X} + \frac{1}{2} \left(\frac{1}{1+\sqrt{n}} \right)$$

which has supremum risk $\frac{1}{4(1+\sqrt{n})^2}$. So

$$\sup_\theta R(\theta, \bar{X}) = \frac{1}{4n} > \frac{1}{4(1+\sqrt{n})^2} = \sup_\theta R(\theta, \delta)$$

Thus, \bar{X} has a higher worst-case risk than $\delta(X)$ as defined above, and hence, we have shown that \bar{X} is not minimax.

Now, let's get more ambitious, and try to see if we can find the minimax estimator under the full model. We know that this can't be \bar{X} , but it's possible that it could be $\delta(X)$. To examine this possibility, we conjecture that $\delta(X)$ is also minimax under the full model. If we are to establish this under the Lemma, we need to show that the supremum risk of $\delta(X)$ under the full model is no more than

$\frac{1}{4(1+\sqrt{n})^2}$ (which is the supremum risk for the binomial submodel). Let us compute:

$$E_F[\delta(X) - \mu(F)]^2$$

$$\begin{aligned} &= E_F \left[\left(\left(\frac{\sqrt{n}}{1+\sqrt{n}} \right) (\bar{X} - \mu(F)) + \frac{1}{1+\sqrt{n}} \right)^2 \right] \\ &= \left(\frac{1}{1+\sqrt{n}} \right)^2 \left[E[X_1^2] + \frac{1}{4} - \mu(F) \right] \end{aligned}$$

where the third step follows from the fact that $\text{Var}(X_1) = n \text{Var}(\bar{X}) = E[X_1^2] - (E[X_1])^2 = E[X_1^2] - (\mu(F))^2$. By assumption $X_1 \in [0, 1]$, so $X_1^2 \leq X_1$ and we can bound the risk:

$$\begin{aligned} E_F[\delta(X) - \mu(F)]^2 &\leq \left(\frac{1}{1+\sqrt{n}} \right)^2 \left[E[X_1] + \frac{1}{4} - \mu(F) \right] \\ &= \frac{1}{4(1+\sqrt{n})^2}. \end{aligned}$$

So, $\delta(X)$ is minimax for the Binomial submodel, and its worst-case risk is the same for the full model and for the Binomial submodel. Therefore, applying the Lemma, we conclude that $\delta(X)$ is minimax. Thus, we have found a minimax estimator.

Property:

- minimax esti. may not necessarily be Bayes esti.
- admissible with constant risk, implies minimax.
- minimality may not guarantee admissibility.
- Example Let $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \mathcal{N}(\theta, \sigma^2)$

where σ^2 is known, and θ is the estimand. Then the minimax estimator is \bar{X} under squared error loss, and we would like to determine whether \bar{X} is admissible. Instead of answering this directly, we answer a more general question: when is $a\bar{X} + b$, $a, b \in \mathbb{R}$, (basically, any affine function of \bar{X}) admissible?

sible? Case 1: $0 < a < 1$. $a\bar{X} + b$ is a convex combination of \bar{X} and b , it is a Bayes estimator with respect to some Gaussian prior on θ . Further, since we are using squared error loss, which is strictly convex, this Bayes estimator is unique. So, by Theorem 5.2.4, $a\bar{X} + b$ is admissible. Case 2: $a=0$. b is also a unique Bayes estimator with respect to a degenerate prior distribution with unit mass at $\theta=b$. So by Theorem 5.2.4, b is admissible. Case 3: $a=1, b \neq 0$. $\bar{X} + b$ is not admissible because it is dominated by \bar{X} . \bar{X} has the same variance as $\bar{X} + b$, but strictly smaller bias. In general, the risk of $a\bar{X} + b$ is:

$$\begin{aligned} E[(a\bar{X} + b - \theta)^2] &= E[(a(\bar{X} - \theta) + b + \theta(a-1))^2] \\ &= \frac{a^2 \sigma^2}{n} + (b + \theta(a-1))^2 \end{aligned}$$

Case 4: $a > 1$.

$$\mathbb{E}[(a\bar{X} + b - \theta)^2] \geq \frac{a^2 \sigma^2}{n} > \frac{\sigma^2}{n} = R(\theta, \bar{X}).$$

\bar{X} dominates $a\bar{X} + b$ when $a > 1$, and so in this case $a\bar{X} + b$ is inadmissible. Case 5: $a < 0$.

$$\mathbb{E}[(a\bar{X} + b - \theta)^2] > (b + \theta(a-1))^2$$

$$\begin{aligned} &= (a-1)^2 \left(\theta + \frac{b}{a-1} \right)^2 \\ &> \left(\theta + \frac{b}{a-1} \right)^2, \end{aligned}$$

and this is the risk of predicting the constant $-b/(a-1)$. So, $-b/(a-1)$ dominates $a\bar{X} + b$, and therefore, $a\bar{X} + b$ is again inadmissible. Case 6: $a=1, b=0$. Here, we use a limiting Bayes argument. Suppose \bar{X} is inadmissible. Then, assuming w.l.o.g that $\sigma^2=1$, we have:

$$R(6, \bar{X}) = \frac{1}{n}$$

By our hypothesis, there must exist an estimator δ' such that $R(\theta, \delta') \leq 1/n$ for all θ and $R(\theta, \delta') < 1/n$ for at least one $\theta' \in \Omega$. Because $R(\theta, \delta')$ is continuous in θ , there must exist $\epsilon > 0$ and an interval (θ_0, θ_1) containing θ' so that:

$$R(\theta, \delta') < \frac{1}{n} - \epsilon \quad \forall \theta \in (\theta_0, \theta_1).$$

Let r'_τ be the average risk of δ' with respect to the prior distribution $N(0, \tau^2)$ on θ . (Note that this is the exact same prior we used to show that \bar{X} is the limit of a Bayes estimator and hence minimax. We did this by letting $\tau \rightarrow \infty$, and therefore letting our prior tend to the improper prior $\pi(\theta) = 1/\theta$.) Let r_τ be the average risk of a Bayes estimator δ under the same prior.

Note that $\delta \neq \delta'$ because $R(\theta, \delta_\tau) \rightarrow \infty$ as $\theta \rightarrow \infty$ which is not consistent with $R(\theta, \delta') \leq 1/n$ for all $\theta \in \mathbb{R}$. So, $r_\tau < r'_\tau$, because the Bayes estimator is unique almost surely with respect to the marginal distribution of θ . We will look at the following ratio, which is selected to simplify our algebra later. This ratio, which we will become arbitrarily large, which we will use to form a contradiction with $r_\tau < r'_\tau$. Using the form of the Bayes risk r_τ computed in a previous lecture (see TPE Example 5.1.14), we can write:

$$\frac{1}{n} - r'_\tau = \frac{1}{\sqrt{2\pi\tau}} \int_{-\infty}^{\infty} \left[\frac{1}{n} - R(\theta, \delta') \right] \exp\left(-\frac{\theta^2}{2\tau^2}\right) d\theta$$

Applying (11.1), we find:

$$\begin{aligned} \frac{1}{n} - r'_\tau &\geq \frac{1}{\sqrt{2\pi\tau}} \int_{-\infty}^{\theta_1} \frac{\theta_1}{n} \varepsilon^{\frac{-\theta^2}{2\tau^2}} d\theta \\ &= \frac{n(1+n\tau^2)}{\sqrt{2\pi\tau}} \varepsilon \int_{-\theta_1}^{\theta_1} e^{-\frac{\theta^2}{2\tau^2}} d\theta \end{aligned}$$

As $\tau \rightarrow \infty$, the first expression, $n(1+n\tau^2)/\sqrt{2\pi\tau} \rightarrow \infty$ and since the integrand converges monotonically to 1,

Lebesgue's monotone convergence theorem ensures that the integral approaches the positive quantity $\theta_1 - \theta_0$. So, for sufficiently large τ , we must have

$$\frac{1-r'_\tau}{n-r_\tau} > 1$$

This means that $r'_\tau < r_\tau$. However, this is a contradiction, because r_τ is the optimal average risk (since it is the Bayes risk). So our assumption that there was a dominating estimator was false, and in this case, $a\bar{X} + b = \bar{X}$ is admissible.

- 2) $X_1, \dots, X_n \sim N(\mu, \sigma^2)$ $Y_1, \dots, Y_m \sim N(\mu, \sigma^2)$
 σ^2 unknown
 (a) $H_0: \mu_2 \leq \mu_1 \quad H_1: \mu_2 > \mu_1$
 (b) show UMP

$$\begin{aligned} H_0: \text{ys} &\quad H_1: \text{y} > 0 \\ \text{pooled sample variance } S^2 &= \frac{1}{m+n-2} \left(\sum_{i=1}^m (X_i - \bar{Y})^2 + \sum_{j=1}^n (Y_j - \bar{Y})^2 \right) \\ \bar{X} &\sim N(\mu_1, \frac{\sigma^2}{m}) \quad \bar{Y} \sim N(\mu_2, \frac{\sigma^2}{n}) \\ \bar{Y} - \bar{X} &\sim N(\mu_2 - \mu_1, \sigma^2 \frac{1}{m+n-2}) \Rightarrow \bar{Y} - \bar{X} \sim N(0, 1) \end{aligned}$$

plug in the estimate S for the unknown σ^2 :

$$T = \frac{\bar{Y} - \bar{X}}{\sqrt{S}} \sim t_{m+n-2}$$

$$\phi(X|Y) = \begin{cases} 1, & \text{if } T < -t_{m+n-2}(1-\alpha) \\ 0, & \text{if } T \geq -t_{m+n-2}(1-\alpha) \end{cases}$$

where $t_{m+n-2}(1-\alpha) = -t_{m+n-2}(1-\alpha)$ is t_{m+n-2} -quantile of t_{m+n-2} -dist

$$\begin{aligned} X &\sim N(\theta) \\ E[X^2] &= E[X]X \\ E[X^2] &= \frac{1}{n} \sum_{i=1}^n E[X_i^2] \\ X_i &\sim \text{Beta}(1, n) \quad X_i \sim \text{Beta}(1, 1) \\ E[X_i^2] &= \frac{1}{1+n} \frac{1}{n} \sum_{i=1}^n E[X_i^2] \\ Y_i &\sim X_i \quad Y_i \sim X_i \quad Y_i \sim X_i \\ \sum_{i=1}^n X_i - \sum_{i=1}^n Y_i &\perp \text{I.I.} \quad Y_i \sim X_i \\ T_n &= \frac{\bar{Y} - \bar{X}}{\sqrt{S}} \perp \text{I.I.} \quad N(0, 1) \\ \text{let } \delta^* &= E[X|Y] \\ \text{by CLT, } \delta^* &\xrightarrow{D} N(0, 1) \\ \text{by SLN, } S &\xrightarrow{P} \delta^* \xrightarrow{P} N(0, 1) \\ \text{by Slutsky's theorem, } \frac{1}{\sqrt{S}} \frac{\delta^* - \bar{X}}{\bar{Y} - \bar{X}} &\xrightarrow{P} N(0, 1) \\ \text{or } P(\delta^* - \bar{X}) &= P(\delta^* - \bar{Y}) \xrightarrow{P} 0 \end{aligned}$$

Now,

$$\begin{aligned} P(X_1 \leq z_1, \dots, X_n \leq z_n) &= \left(\int_{-\infty}^{z_1} g_1(x_1) dx_1 \right) \cdots \left(\int_{-\infty}^{z_n} g_n(x_n) dx_n \right) \\ &= (c_2 \cdots c_n) \int_{-\infty}^{z_1} g_1(x_1) dx_1 \cdots \left(\int_{-\infty}^{z_n} g_n(x_n) dx_n \right) \\ &= \frac{1}{c_1} \int_{-\infty}^{z_1} g_1(x_1) dx_1 \cdots \left(\int_{-\infty}^{z_n} g_n(x_n) dx_n \right) \\ &= \prod_{i=1}^n P(X_i \leq z_i) \end{aligned}$$

By Lemma 2.1.5, A_i 's are independent. By Theorem 2.1.7, $\sigma(A_i)$'s are independent. Since $X_i^{-1}(A) \in \sigma(A)$ for all $A \in 2^{\mathbb{R}}$, $\sigma(X_i) \subset \sigma(A_i)$. Therefore, X_1, \dots, X_n are independent.

b) X_1, X_2, \dots are iid with mean 0 and variance 1. Prove that

$$\lim_{n \rightarrow \infty} \frac{1}{n^{1/2}} \mathbb{E} \left[\sum_{i=1}^n |X_i| \right] = \frac{2}{\pi}$$

$$(b) \text{ Let } Z \sim N(0, 1), \quad \mathbb{E}|Z| \approx \int_0^\infty z \frac{1}{\sqrt{2\pi}} e^{-\frac{z^2}{2}} dz = \sqrt{\frac{2}{\pi}}$$

$$\frac{1}{n} \sum_{i=1}^n |X_i| \xrightarrow{D} Z \Rightarrow \frac{1}{n} \mathbb{E} \left[\sum_{i=1}^n |X_i| \right] \rightarrow \mathbb{E}|Z|$$

Given a collection of set A_i derived from Ω , satisfies finite additivity and continuity.

For \forall disjoint A_1, \dots, A_n , take $A = \bigcup_{i=1}^n A_i$, $B_n = A - \bigcup_{i=1}^n A_i$, then $B_n \supset A_1 \supset \dots \supset A_n \supset \emptyset$ when $n \rightarrow \infty$ then according to the definition of continuity, $\lim_{n \rightarrow \infty} P(B_n) \rightarrow 0$ by B_n and $\bigcup_{i=1}^n A_i$ disjoint then $P(A) = P(B_n \cup \left(\bigcup_{i=1}^n A_i \right)) = P(B_n) + P(\bigcup_{i=1}^n A_i)$

$$= P(B_n) + \sum_{i=1}^n P(A_i) \rightarrow 0 + \sum_{i=1}^n P(A_i) \text{ when } n \rightarrow \infty.$$

\Rightarrow countable additivity

Generalized Inverse

decomposition:

$$\text{1. full rank: } A_{m \times n} = B_{m \times r} C_{r \times n}.$$

$$\text{2. } \exists \text{ nonsingular } P, Q \text{ s.t. } P \cdot A \cdot Q = (I_r)_{m \times n}.$$

$$\text{3. (symmetric) } \exists \text{ orthogonal } P, \text{ s.t. } P^T A P = D.$$

$$\text{idempotent: } A = A^2; \text{ Orthogonal: } Q^T Q = Q Q^T = I.$$

$$\text{1. all idempotent matrices (except } I\text{) are singular.}$$

$$\text{2. if } A \text{ is idempotent, then } \text{rank}(A) = \text{tr}(A).$$

$$\text{3. if } A \text{ is idempotent, then e.v. of } A \text{ are 1 or 0.}$$

$$\text{4. if } A \text{ is sym. with e.v. 1 and 0, then } A \text{ is idempotent.}$$

$$\text{Prf: (2) } \text{tr}(A) = \text{tr}(PQ) = \text{tr}(C) = \text{tr}(I_r) = r = \text{r}(A)$$

inverse:

$$\text{1. left inverse: } (r = n) A_{\text{left}}^{-1} = (A^T A)^{-1} A^T.$$

$$\text{2. right inverse: } (r = m) A_{\text{right}}^{-1} = A^T (AA^T)^{-1}.$$

$$\text{3. Moore-Penrose inverse:}$$

- (a) AA^+ and A^+A are symmetric,
- (b) $A = AA^+A$ and $A^+ = A^+AA^+$.

$$\text{4. generalized inverse: } AA^+A = A.$$

Moore-Penrose inverse:

$$\text{1. existence: } A^+ = C^T(CCT)^{-1}(B^T B)^{-1}B^T.$$

$$\text{2. uniqueness: } AA_1^+ = (AA_2^+AA_1^+)^T = AA_2^+.$$

$$\text{3. rank}(A) = \text{rank}(A^+).$$

$$\text{4. } (A^+)^T = (A^+)^T. \text{ and if } A = A^T, \text{ then } A^+ = (A^+)^T.$$

$$\text{5. } AA^+ + A^+A = AA^+ + I_m - A^+A \text{ are sym. and idempotent.}$$

$$\text{5. If } A \text{ is nonsingular, } A^{-1} = A^+.$$

$$\text{6. If } A \text{ is symmetric idempotent, } A^+ = A.$$

$$\text{7. If } r(A_{m \times n}) = m, \text{ then } A^+ = A^T(AA^T)^{-1}AA^+ = I_n.$$

$$\text{If } r(A_{m \times n}) = n, \text{ then } A^+ = (A^T A)^{-1}A^T, A^+A = I_m.$$

$$\text{generalized inverse: } AA^{-1}A = A.$$

$$\text{1. not unique.}$$

$$\text{2. Moore-Penrose inverse is also generalized inverse.}$$

$$\text{3. rank}(X^-) \geq r, \text{ rank}(X^-X) = \text{rank}(XX^-) = r.$$

$$\text{4. } X-X = I_n \text{ iff rank}(X) = n, XX^- = I_m \text{ iff rank}(X) = m. \text{ (Hint: full-rank -> nonsingular + (idempotent) -> identity.)}$$

$$\text{5. if } G \text{ is G-inv. of } X^T X, \text{ then } G^T \text{ is also G-inv. of } X^T X.$$

$$\text{6. } \star XGX^T X = X, \text{ i.e. } GX^T \text{ is G-inv. of } X.$$

$$\text{take } K = XGX^T, \text{ then } XK = X, X^T K = X^T$$

$$\text{7. } K = XGX^T \text{ is invariant w.r.t. the choice of } G.$$

$$\text{8. } K = XGX^T \text{ is symmetric and idempotent.}$$

$$\text{9. } K = XGX^T = XX^+.$$

Proof:

$$\text{(i) Since } G \text{ is a generalized inverse of } (X^T X), (X^T X)G(X^T X) = X^T X. \text{ Taking the transpose of both sides}$$

$$\begin{aligned}[X^T X]^T &= [(X^T X)G(X^T X)]^T \\ &= (X^T X)^T G^T (X^T X)^T\end{aligned}$$

$$\text{But } (X^T X)^T = X^T (X^T)^T = X^T X, \text{ hence } (X^T X)G^T (X^T X)^T = (X^T X).$$

$$\text{(ii) From (i) } (X^T X)G^T (X^T X) = (X^T X). \text{ Denote } (X^T X)G^T \text{ by } B. \text{ Then}$$

$$\begin{aligned}0 &= BX^T X - X^T X \\ &= (BX^T X - X^T X)(B^T - I) \\ &= BX^T XB^T - X^T XB^T - BX^T X - X^T X \\ &= (BX^T - X^T)(BX^T - X^T)^T\end{aligned}$$

$$\text{Hence, } 0 = BX^T - X^T$$

$$\Rightarrow BX^T = X^T$$

$$\Rightarrow X^T X G^T X^T = X^T$$

$$\text{Taking the transpose}$$

$$\begin{aligned}X &= (X^T X G^T X^T)^T \\ &= X G X^T\end{aligned}$$

$$\text{Hence, } GX^T \text{ is a generalized inverse for } X.$$

$$\text{matrix identity: }$$

$$\text{(iii) Suppose } F \text{ and } G \text{ are generalized inverses for } X^T X. \text{ Then, from (ii)}$$

$$\begin{aligned}X G X^T X &= X \\ \text{and} \\ X F X^T X &= X\end{aligned}$$

$$\text{It follows that}$$

$$\begin{aligned}0 &= X - X \\ &= (X G X^T X - X F X^T X) \\ &= (X G X^T X - X F X^T X)(G^T X^T - F^T X^T) \\ &= (X G X^T - X F X^T)X(G^T X^T - F^T X^T) \\ &= (X G X^T - X F X^T)(X G^T X^T - X F^T X^T) \\ &= (X G X^T - X F X^T)(X G X^T - X F X^T)^T\end{aligned}$$

$$\text{Since the (i,i) diagonal element of the result of multiplying a matrix by its transpose is the sum of the squared entries in the } i\text{-th row of the matrix, the diagonal elements of the product are all zero only if all entries are zero in every row of the matrix. Consequently,}$$

$$(X G X^T - X F X^T) = 0$$

$$\text{(iv) For any generalized inverse } G,$$

$$T = G X^T X G^T$$

$$\text{is a symmetric generalized inverse. Then}$$

$$X T X^T$$

$$\text{is symmetric and from (iii),}$$

$$X G X^T = X T X^T.$$

[Algorithm of Obtaining G-(m)] Let $A \in \mathbb{R}^{m \times n}$ be a matrix of rank r , and $A_{11} \in \mathbb{R}^{r \times r}$. If A_{11} is invertible, then $G = \begin{pmatrix} A_{11}^{-1} & 0 \\ 0 & 0 \end{pmatrix} \in \mathbb{R}^{m \times m}$ is a G -inv of A .

$$\text{quadratic form } \mathbf{x}' \mathbf{A} \mathbf{x} \text{ with } \mathbf{x} \sim \mathcal{N}(\mu, \Sigma):$$

$$\begin{aligned}1. \text{ m.g.f.: } &|\mathbf{I} - 2t\mathbf{A}\Sigma|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}t^2\mathbf{\mu}'[\mathbf{I} - (\mathbf{I} - 2t\mathbf{A}\Sigma)^{-1}]\mathbf{\mu}\right). \\ 2. \text{ If } \mathbf{A} \text{ & } \mathbf{V} \text{ sym., } \mathbf{V} \text{ p.d., then e.v. of } \mathbf{AV}, 0, 1 \rightarrow \mathbf{AV} \text{ idem.}\end{aligned}$$

Proof: Recall that if a symmetric matrix is positive definite if and only if it can be written as $P^T P$ for a nonsingular P . Then, $\mathbf{V} = P^T P$ for some nonsingular matrix P . If $|\mathbf{AV} - \lambda I| = 0$ has roots 0 and 1, then

$$\begin{aligned}|\mathbf{P}||\mathbf{AV}| |\mathbf{P}^{-1}| &= 0 \quad \text{has roots 0 & 1} \\ \Rightarrow |\mathbf{P} \mathbf{A} \mathbf{V}^T - \mathbf{A} \mathbf{P}^T \mathbf{P}^{-1}| &= 0 \quad \text{has roots 0 & 1} \\ \Rightarrow |\mathbf{P} \mathbf{A}^T \mathbf{P}^T \mathbf{P}^{-1} - \lambda I| &= 0 \quad \text{has roots 0 & 1} \\ \Rightarrow |\mathbf{P} \mathbf{A}^T - \lambda I| &= 0 \quad \text{has roots 0 & 1.} \end{aligned} \quad (1)$$

Thus, $\mathbf{P} \mathbf{A}^T \mathbf{P}^T$ has eigenvalues 0 and 1. But $\mathbf{P} \mathbf{A}^T \mathbf{P}^T$ is symmetric (because \mathbf{A} is symmetric). Hence, $\mathbf{P} \mathbf{A}^T \mathbf{P}^T$ is idempotent. That is

$$\mathbf{P} \mathbf{A}^T \mathbf{P}^T \mathbf{P}^T = \mathbf{P} \mathbf{A}^T \mathbf{P}^T \Rightarrow \mathbf{P} \mathbf{A} \mathbf{V} \mathbf{P}^T = \mathbf{P} \mathbf{A}^T \mathbf{P}^T.$$

As \mathbf{P} is nonsingular,

$$\mathbf{P}^{-1} \mathbf{P} \mathbf{A} \mathbf{V} \mathbf{P}^T \mathbf{P}^T = \mathbf{P}^{-1} \mathbf{P} \mathbf{A}^T \mathbf{P}^T \mathbf{P}^T$$

which implies

$$\mathbf{A} \mathbf{V} \mathbf{A}^T = \mathbf{A}^T \mathbf{V} \mathbf{A}$$

$$\text{3. If } \mathbf{A} \text{ sym. idem. with rank } r, \text{ then } \mathbf{A} \text{ has } r \text{ e.v.-1 and rest } -0.$$

$$\text{THEOREM 1. Let } \mathbf{x}_{px1} \sim \mathcal{N}(\mu, \Sigma) \text{ and let } \mathbf{A} \text{ be symmetric. Then, } q = \mathbf{x}' \mathbf{A} \mathbf{x} \text{ follows } X_{p,1}^2 \text{ with } r \text{ being the rank of } \mathbf{A} \text{ and } \lambda = \frac{\sigma^2 \mu^2}{2} \text{ if and only if } \mathbf{A} \mathbf{x} \text{ is idempotent.}$$

Corollaries:

$$\begin{aligned}1. \text{ If } \mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \text{ then } \mathbf{x}' \mathbf{A} \mathbf{x} &= X_{p,1}^2 \text{ if and only if } \mathbf{A} \text{ is idempotent of rank } r. \\ 2. \text{ If } \mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{V}), \text{ then } \mathbf{x}' \mathbf{A} \mathbf{x} &= X_{p,1}^2 \text{ if and only if } \mathbf{A} \mathbf{V} \text{ is idempotent of rank } r. \\ 3. \text{ If } \mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \sigma^2 \mathbf{I}), \text{ then } \frac{\mathbf{x}' \mathbf{x}}{\sigma^2} &\sim X_{(n-1)p,1}^2. \\ 4. \text{ If } \mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{I}), \text{ the } \mathbf{x}' \mathbf{A} \mathbf{x} &= X_{(n-p)_1}^2 \text{ if and only if } \mathbf{A} \text{ is idempotent of rank } r.\end{aligned}$$

Note

$$\begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} \begin{bmatrix} A_{11} & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix}.$$

By the definition of generalized inverse, it suffices to verify $A_{21} A_{11}^{-1} A_{12} = A_{22}$. Since the row rank of $[A_{11}, A_{12}]$ and \mathbf{A} are both r , then there exists a $(n-r) \times r$ -matrix \mathbf{P} such that $\mathbf{P}[A_{11}, A_{12}] = [A_{21}, A_{22}]$. And it can be found that $\mathbf{P} = A_{21} A_{11}^{-1}$, and $\mathbf{P} A_{12} = A_{22}$.

Proof:

$$\text{he hat matrix } \mathbf{H} \text{ is symmetric idempotent;}$$

Proof: Note that $(\mathbf{X}^T \mathbf{X})$ is symmetric and $(\mathbf{X}^T \mathbf{X})^{-1}$ is also symmetric. Then,

$$\mathbf{I} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T = \mathbf{H}^T, \quad \mathbf{HH} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T = \mathbf{H}.$$

$$2. \mathbf{X}' \hat{\mathbf{\beta}} = \mathbf{0}; \text{ (This holds because of } \mathbf{X}^T \mathbf{H} = \mathbf{X}^T, \mathbf{H} \mathbf{X} = \mathbf{X} \text{ and } \mathbf{X}^T(I - \mathbf{H}) = 0, (\mathbf{I} - \mathbf{H}) \mathbf{X} = 0\text{.)}$$

$$\text{Proof: Since}$$

$$\begin{aligned}\mathbf{X}^T \mathbf{H} &= \mathbf{X}^T \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T = \mathbf{X}^T, \\ \mathbf{H} \mathbf{X} &= \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T = \mathbf{X},\end{aligned}$$

then,

$$\mathbf{X}^T(I - \mathbf{H}) = \mathbf{X}^T - \mathbf{X}^T \mathbf{H} = \mathbf{X}^T - \mathbf{X}^T = \mathbf{0},$$

$$(\mathbf{I} - \mathbf{H}) \mathbf{X} = \mathbf{X} - \mathbf{H} \mathbf{X} = \mathbf{X} - \mathbf{X} = \mathbf{0}.$$

Clearly,

$$\mathbf{X}^T \hat{\mathbf{\beta}} = \mathbf{X}^T(I - \mathbf{H}) \mathbf{Y} = 0.$$

3. $\hat{\mathbf{\beta}}^T \hat{\mathbf{\beta}} = 0;$

Proof: Write

$$\begin{aligned}(\mathbf{HY})^T (\mathbf{I} - \mathbf{H}) \mathbf{Y} &= \mathbf{Y}^T \mathbf{H}^T (\mathbf{I} - \mathbf{H}) \mathbf{Y} = \mathbf{Y}^T \mathbf{H} (\mathbf{I} - \mathbf{H}) \mathbf{Y} \\ &= \mathbf{Y}^T \mathbf{O} \mathbf{Y} = 0.\end{aligned}$$

4. $\mathbf{I} - \mathbf{H}$ is symmetric idempotent;

5. $E(\hat{\beta}) = \beta_0$ (unbiased estimate);

Proof:

$$E(\hat{\beta}) = E((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T E(\mathbf{Y}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{B}_0 = \beta_0.$$

Consequently,

$$\frac{SSE}{\sigma^2} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i^T \mathbf{B}_0 - \mathbf{x}_i^T \mathbf{B})^2.$$

Therefore, the noncentrality parameter is 0. In addition,

$$r(\frac{I - H}{\sigma^2}) = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i^T \mathbf{B})^2 = n - r(\mathbf{X}).$$

$$\text{Consequently, } \frac{SSE}{\sigma^2} = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i^T \mathbf{B})^2 - \frac{n}{n-r(\mathbf{X})} \sum_{i=1}^n (\mathbf{x}_i^T \mathbf{B})^2.$$

(Lemma) If $M = \begin{bmatrix} X & Z \\ Z & W \end{bmatrix}$ where $W = (D - B^T A^{-1} B)^{-1}$

$$\begin{aligned}Cov(\hat{\beta}) &= Cov((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T Cov(\mathbf{Y}) \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \\ &= \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{X}^T \mathbf{X})^{-1} = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}.\end{aligned}$$

7. $tr(I_n - \mathbf{H}) = n - p - \frac{1}{\sigma^2} \sum_{i=1}^n Var(\hat{\beta}_i) = \frac{1}{\sigma^2} \mathbf{P}^T \sigma^2 \mathbf{I} \mathbf{P} = \mathbf{P}^T \mathbf{I} \mathbf{P} = (S(V))$

Proof: Note that

$$tr(\mathbf{H}) = tr(\mathbf{X}^T \mathbf{X}^{-1} \mathbf{X}^T) = tr(\mathbf{X}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1}) = tr(\mathbf{I}_p) = p.$$

Then,

$$tr(I_n - \mathbf{H}) = tr(I_n) - tr(\mathbf{H}) = n - p.$$

8. $\hat{\mathbf{\beta}}^T \hat{\mathbf{\beta}} = tr(\mathbf{Y} \mathbf{Y}^T (\mathbf{I} - \mathbf{H}))$

Proof: We can easily show that

$$\begin{aligned}\hat{\mathbf{\beta}}^T \hat{\mathbf{\beta}} &= \mathbf{Y}^T (\mathbf{I} - \mathbf{H})^T (\mathbf{I} - \mathbf{H}) \mathbf{Y} = \mathbf{Y}^T \mathbf{H} (\mathbf{I} - \mathbf{H}) \mathbf{Y} \\ &= \mathbf{Y}^T \mathbf{O} \mathbf{Y} = 0.\end{aligned}$$

Woodbury matrix identity:

$$(A + BCD)^{-1} = A^{-1} - A^{-1}B(C^{-1} + DA^{-1}B)^{-1}DA^{-1}.$$

Prf: $\mathbf{F}(\mathbf{R}) = \frac{\hat{\mathbf{\beta}}^T \hat{\mathbf{\beta}}}{\sigma^2} = \frac{SSE}{\sigma^2} = \frac{1}{n-p} = \mathbf{F}(\mathbf{R})$

$$\text{9. } E(\mathbf{Y} \mathbf{Y}^T) = \sigma^2 \mathbf{I} + \mathbf{X} \mathbf{B} \mathbf{B}^T \mathbf{X}^T;$$

$$10. \hat{\mathbf{\beta}}^T \hat{\mathbf{\beta}}/(n-p) \text{ is an unbiased estimate of } \sigma^2, \text{ that is}$$

$$E(\frac{\hat{\mathbf{\beta}}^T \hat{\mathbf{\beta}}}{n-p}) = \sigma^2.$$

Proof: Write

$$E(\hat{\mathbf{\beta}}^T \hat{\mathbf{\beta}}) = E(\mathbf{Y}^T (\mathbf{I} - \mathbf{H}) (\mathbf{I} - \mathbf{H}) \mathbf{Y}) = E(\mathbf{Y}^T (\mathbf{I} - \mathbf{H}) \mathbf{Y})$$

$$= tr((\mathbf{I} - \mathbf{H}) \mathbf{Y}) + \mathbf{B}^T \mathbf{X}^T (\mathbf{I} - \mathbf{H}) \mathbf{Y}$$

$$= \sigma^2 tr(\mathbf{I} - \mathbf{H}) = \sigma^2(n - p).$$

weighted least squares estimation

1. assumption: Σ is known.

2. loss function: $S(\beta) = (\mathbf{Y} - \mathbf{X} \beta)^T (\mathbf{Y} - \mathbf{X} \beta)$

3. WLSE: $\hat{\beta} = (\mathbf{X}' \Sigma^{-1} \mathbf{X})^{-1} \mathbf{X}' \Sigma^{-1} \mathbf{Y}$

best linear unbiased estimator:

1. minimizes $\text{Var}(\hat{\mathbf{\beta}}') = \Sigma \mathbf{A} \mathbf{A}$

s.t. $\mathbf{E}[\hat{\mathbf{\beta}}'] = \mathbf{B}$

Hypothesis Testing $H_0: K^\top \beta = m$.

$$1. (\mathbf{K}'\beta - m) \sim N(\mathbf{K}'\beta - m, \sigma^2 \mathbf{K}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{K}).$$

$$2. Q = (\mathbf{K}'\beta - m)' [\mathbf{K}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{K}]^{-1} (\mathbf{K}'\beta - m) \sim \sigma^2 \chi_{(s, n)}^2.$$

3. Q and SSE are independent.

$$\text{Hint: } \mathbf{K}'\beta - m = \mathbf{K}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' [\mathbf{Y} - \mathbf{X}\mathbf{K}(\mathbf{K}'\mathbf{K})^{-1}\mathbf{m}].$$

$$4. F(H) = \frac{\mathbf{Q}/s}{\text{SSE}/(n - r(\mathbf{X}))} \stackrel{H_0}{\sim} F_{s, n-r(\mathbf{X})}.$$

$\therefore \text{Var}:$

Estimation under constraint (Null hypothesis):

1. use 2θ as Lagrange multiplier

2. minimize $W(\theta, \beta) = (\mathbf{Y} - \mathbf{X}\beta)'(\mathbf{Y} - \mathbf{X}\beta) + 2\theta'(\mathbf{K}'\beta - m)$,

3. get $\hat{\beta} = \bar{\beta} - (\mathbf{X}'\mathbf{X})^{-1}\mathbf{K}[\mathbf{K}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{K}]^{-1}(\mathbf{K}'\beta - m)$.

$$\therefore \text{SSE}_{\text{H}(\beta)}(\hat{\beta}) = \text{SSE}(\hat{\beta}) + Q \geq \text{SSE}(\hat{\beta}) = (\mathbf{Y} - \hat{\beta})'(\mathbf{Y} - \hat{\beta})$$

(Properties of $\hat{\beta}$) (1) $E(\hat{\beta}) = G\mathbf{X}^\top \mathbf{X}\beta$. (2)

$$\text{Var}(\hat{\beta}) = G\mathbf{X}^\top \mathbf{X}G^2. \quad (3) \hat{y} = \mathbf{X}\hat{\beta} = \mathbf{X}\mathbf{G}\mathbf{X}^\top \mathbf{y}, \quad (4)$$

$$E(\hat{y}) = \mathbf{X}\mathbf{G}\mathbf{X}^\top \mathbf{X}\beta = \mathbf{X}\beta. \quad (5) \text{SSE} = \mathbf{y}^\top (\mathbf{I} - \mathbf{X}\mathbf{G}\mathbf{X}^\top) \mathbf{y}. \quad (6)$$

$$\text{SSR} = \mathbf{y}^\top \mathbf{X}\mathbf{G}\mathbf{X}^\top \mathbf{y} = (\hat{\beta})^\top \mathbf{X}^\top \mathbf{y}. \quad (7) \hat{\sigma}^2 = \frac{\text{SSR}}{n-r(\mathbf{X})}$$

estimator of σ^2 . (8) $\text{SSR}_0 = \mathbf{y}^\top (\mathbf{X}\mathbf{G}\mathbf{X}^\top - \frac{1}{n}\mathbf{I}) \mathbf{y}$.

(Distributional Properties) (1) $\hat{\beta} \sim N(\mathbf{G}\mathbf{X}^\top \mathbf{X}\beta, \mathbf{G}\mathbf{X}^\top \mathbf{X}G^2)$. The covariance matrix is singular. (2) $\frac{\text{SSR}}{\text{SSR}_0} \sim \chi^2_{n-r(\mathbf{X})}$.

$$\frac{\text{SSR}}{\text{SSR}_0} \sim \chi^2(n), \quad \frac{\text{SSR}}{\text{SSR}_0} \sim F(r(\mathbf{X}), n-r(\mathbf{X})), \quad \frac{\text{SSR}}{\text{SSR}_0} \sim \chi^2_{n-r(\mathbf{X})}$$

(Identifiability) Formally, the parameter β is identifiable if $f(\beta_1) = f(\beta_2)$ implies that $\beta_1 = \beta_2$ for any β_1 and β_2 . More generally, the vector-valued function $g(\beta)$ is identifiable if $g(\beta_1) = g(\beta_2)$ implies that $g(\beta_1) = g(\beta_2)$.

[Proposition] In a linear model for which \mathbf{X} is of full rank, β is identifiable.

[Proposition] A function $g(\beta)$ is identifiable if and only if $\mathcal{M}_v(g(\beta)) = \mathcal{M}_v(f(\beta))$ for some function f .

(Estimable functions) If a vector t exists such that $t^\top \mathbf{E}(y) = t^\top \beta$, then $t^\top \beta$ is said to be estimable. Linear combinations of estimable functions are estimable.

• $E(t^\top y) = t^\top \mathbf{X}\beta = q^\top \beta \Rightarrow t^\top \mathbf{X} = q^\top t$ for some t . This is equivalent to saying that q is in the row space of \mathbf{X} .

Theorem 1. (Gauss-Markov Theorem) The best linear unbiased estimator of the estimable function $q^\top \beta$ is $\hat{\beta}$.

Proof. We prove this theorem from the following three aspects:

(i) $q^\top \beta^0 = q^\top \mathbf{G}\mathbf{X}^\top \mathbf{y} \iff$ linear function of y_i .

(ii)

$$\begin{aligned} E(q^\top \beta^0) &= q^\top E(\beta^0) \\ &= q^\top \mathbf{G}\mathbf{X}^\top \mathbf{E}(y) \\ &= q^\top \mathbf{G}\mathbf{X}^\top \mathbf{X}\beta \\ &= t^\top \mathbf{X}\mathbf{G}\mathbf{X}^\top \mathbf{X}^\top \mathbf{t}^\top \\ &= t^\top \mathbf{X}\mathbf{G}\mathbf{X}^\top \mathbf{t}^\top \\ &= q^\top \mathbf{G}\mathbf{t}^\top. \end{aligned}$$

$$\text{var}(q^\top \beta^0) = q^\top \text{var}(\beta^0)q$$

$$= q^\top \mathbf{G}\mathbf{X}^\top \mathbf{X}G^2 q^2$$

$$= t^\top \mathbf{X}\mathbf{G}\mathbf{X}^\top \mathbf{X}^\top \mathbf{t}^\top q^2$$

$$= t^\top \mathbf{X}\mathbf{G}\mathbf{X}^\top \mathbf{t}^\top q^2$$

$$= q^\top \mathbf{G}\mathbf{t}^\top q^2.$$

$$\text{var}(q^\top \beta^0 - k^\top \beta) = \text{var}(q^\top \beta^0) - k^\top \text{var}(\beta)k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2\text{cov}(q^\top \beta^0, k^\top \beta)$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0 - k^\top \beta)$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$= \text{var}(q^\top \beta^0) + \text{var}(k^\top \beta) - 2q^\top \mathbf{G}\mathbf{t}^\top k$$

$$F(y) = P(Y \leq y).$$

is called the τ -th quantile of Y . The median $F^{-1}(1/2)$ plays the central role.

3.6.2 Linear quantile regression.

Given a covariate vector $X \in \mathbb{R}^{p+1}$, the τ -th conditional quantile function of $Y | R$ is modeled as

$$F^-(\tau) = \inf_{y < \tau} \{y | F(y) \geq \tau\}.$$

$Q_\tau(X^\top)X^{-1}$.

for certain specific $\tau \in (0, 1)$ of interest, and $\beta_\tau = (\beta_0^\top, \beta_\tau^\top)$ —vector usually including an intercept. For model (1) with complete data, a classical estimate of β_τ , denoted as $\hat{\beta}_\tau$, is obtained by minimizing

$$\begin{aligned} R(\beta_\tau) &= \sum_{i=1}^n \rho_\tau(y_i - x_i^\top \beta_\tau), \\ \text{over } \beta_\tau, \text{ where } (y_i, x_i) \text{ are iid copies of } (Y, X), i = 1, \dots, n, \\ \hat{\beta}_\tau(\mathbf{y}) &= u(\tau - I(u < 0)), \end{aligned}$$

(2)

directional derivative of R in direction w is given by

$$\begin{aligned} \nabla R(\beta_\tau, w) &\equiv \frac{d}{dt} R(\beta_\tau + tw)|_{t=0} \\ &= \frac{d}{dt} \sum_{i=1}^n [y_i - x_i^\top (\beta_\tau + tw)]^\top [\tau - I(y_i - x_i^\top (\beta_\tau + tw) < 0)]|_{t=0} \\ &= - \sum_{i=1}^n \psi_\tau(y_i - x_i^\top \beta_\tau, -x_i^\top w)x_i^\top w, \end{aligned}$$

where

$$\psi_\tau(u, v) = \begin{cases} \tau - I(u < 0) & \text{if } u \neq 0; \\ \tau - I(v < 0) & \text{if } u = 0. \end{cases}$$

3.6.3 Asymptotic properties of quantile regression.

The conditional distribution function of $Y | R$ will be written as

$$P(Y_i < y | x_i) = F_{Y_i}(y | x_i) = F(y)$$

and so

$$\begin{aligned} Q_\tau(F)(\tau) &= F_{Y_i}^{-1}(\tau) = \xi(\tau), \\ \text{we will employ the following regularity conditions to explore the asymptotic behavior of the estimator } \hat{\beta}_\tau(\tau). \end{aligned}$$

- Condition (A1): The distribution functions $\{F_i\}$ are absolutely continuous, with continuous densities $f_i(x_i)$ uniformly bounded away from zero and at the point $\xi_i(\tau)$, $i = 1, \dots, n$.
- Condition (A2(i)): There exists positive definite matrices D_0 and D_1 such that

$$(i) \lim_{n \rightarrow \infty} n^{-1} \sum_i x_i x_i^\top = D_0;$$

$$(ii) \lim_{n \rightarrow \infty} n^{-1} \sum_i f_i(x_i) \hat{\beta}_\tau = D_1 F(\tau);$$

$$(iii) \max_{1 \leq i \leq n} \|x_i\|^2 < \tau.$$

- Condition (A2(ii)): $\int f_i(x_i) \hat{\beta}_\tau dF_i(x_i) = 0$.
- Condition (A3): μ is $N(0, V)$ and Σ is V^{-1} .
- Condition (A4): $\Delta \hat{\beta}_\tau(\tau) \overset{d}{\rightarrow} N(0, \Omega(\tau))$, $\tau \in (0, 1)$.

Theorem 4.1 Under Conditions A1 and A2,

$$\sqrt{n}[\hat{\beta}_\tau(\tau) - \beta_0(\tau)] \rightarrow N(0, \tau D_1^{-1} D_0 D_1^{-1})$$

Conditions (A2(i) and A2(ii)) are familiar throughout the literature on M-estimators for regression models; some variant of them is necessary to ensure that a Δ identifies condition A2(ii) is really a matter of notational convenience and could be deduced from Condition A2(i) and a slightly strengthened version of Condition A4.

Theorem 4.2 Under Conditions A1 and A2,

$$\sqrt{n}[\hat{\beta}_{\tau0}(\tau) - \beta_{\tau0}] \xrightarrow{d} (0, V(\tau))$$

in distribution as $n \rightarrow \infty$, where $V(\tau) = (1 - \tau)/F(\tau)$ and $V = E(I(\bar{X} < \tau))$.
Remark: The asymptotic properties above were established under fixed design, which can be extended to handle random design without further difficulties.

Theorem 4.3 Under Conditions A1 and A2,

$$\sqrt{n}[\hat{\beta}_\tau(\tau) - \beta_0(\tau)] \rightarrow N(0, \Omega(\tau))$$

In the random model (random design) $X_i \sim z_i^\top \beta + u_i$, $i = 1, \dots, n$, where z_i have a common distribution function F_z , in particular, a special case of quantile regression when $\tau = 0.5$. Its estimate $\hat{\beta}_\tau(\tau)$ is defined to be the minimizer of

$$\min_{\beta} \sum_{i=1}^n V_\tau(y_i - z_i^\top \beta),$$

Hence, the following theorem can be established in a straightforward fashion under regularity conditions.

Example $M_1: Y_g = \mu + \Delta \log + Eg$ $\log \sim N(0, \Omega(\tau))$
 $M_2: Y_g = \mu + \Delta^2 \log + Eg$
 $\log \sim N(0, \Omega(\tau))$ $\eta_j g_j = T_j(Y_j + \Gamma_j(F(g_j)) + Eg_j)$
 $M_{2a}: Y_{gj} = \mu_j + t^2 \Delta \log + \Delta \log + Eg_j$

link model $M_{2b}: Y_{gj} = \mu_j + \Delta \log + Eg_j + \Delta_j \log + Eg_j$
 $\log \sim N(0, \Omega(\tau))$ $\eta_j g_j = T_j(Y_j + \Gamma_j(F(g_j)) + Eg_j)$

Hierarchical Data Come from a number of diff groups with known hierarchical structure.
L1 (within group) $y_{ij} = y_j + \Delta_j \log j + Eg_j$ $y_j \sim N(0, \Omega_j)$
L2 (between group) $y_g = \mu_g + \Delta_{2g} \log g + Eg_g$ $y_g \sim N(0, \Omega_{2g})$
 $\rightarrow U_{ij} = \mu_j + \Delta_{2j} \log g + Eg_j + \Delta_j \log y_j + Eg_j$
 $\eta_j g_j = T_j(Y_j + \Gamma_j(F(y_j)) + Eg_j)$
 $\eta_g g_g = T_g(Y_g + \Gamma_g(F(y_g)) + Eg_g)$
Assumption within group errors are invariant/shared. Coeff & errors are distinct.
(Random Permutation) Let $\psi = (\Omega, W, \theta)$, the permutation sampler for generating ψ from the posterior $p(\psi | \mathbf{Y})$ is implemented as follows:
(1) Generate ψ from the unconstrained posterior $p(\psi | \mathbf{Y})$ using standard Gibbs sampling steps;
(2) Select some permutation $\pi(1), \dots, \pi(K)$ and define $\psi = \psi(\pi)$ from ψ by reordering the labeling through this permutation: $(\pi(1), \dots, \pi(K)) \rightarrow (\pi(\psi_1), \dots, \pi(\psi_K))$, and
Model Comparison $p(\mathbf{y} | \mathbf{x}, \theta) = p(\mathbf{w}_1, \dots, \mathbf{w}_m)$ $= p(\mathbf{w}_1, \dots, \mathbf{w}_m)$ scale label switching

Consider the following competing models:
 $M_1: \mathbf{Y} | \theta, \tau \stackrel{D}{=} \sum_{k=1}^K \pi_k f_k(\mathbf{Y} | \mu_k, \Sigma_k)$
 $M_0: \mathbf{Y} | \theta, \tau \stackrel{D}{=} \sum_{k=1}^K \pi_k^* f_k(\mathbf{Y} | \mu_k, \Sigma_k)$,

where $1 \leq c < K$. They can be linked up by a variable M_c as follows:

$$\begin{aligned} M_c: & \mathbf{Y} | \theta, \tau \stackrel{D}{=} \frac{\pi_1}{c} (1 - t_1) f_1(\mathbf{y}_{c+1}, \dots, \mathbf{y}_c) f_1(\mathbf{y}_1, \Sigma_1) + \dots + \\ & + \frac{\pi_c}{c} (1 - t_c) f_c(\mathbf{y}_{c+1}, \dots, \mathbf{y}_c) f_c(\mathbf{y}_1, \Sigma_1) + \dots + t_{c+1} \pi_{c+1} f_{c+1}(\mathbf{y}_{c+2}, \Sigma_{c+1}) + \dots + t_m \pi_m f_m(\mathbf{y}_{c+1}, \Sigma_m). \end{aligned}$$

$$\log p(\mathbf{y}, \Omega, \Sigma, \theta) = \sum_{i=1}^n \log \left[\sum_{k=1}^c \pi_k (1 - t_k) \sum_{h=c+1}^K \pi_h f_h(\mathbf{y}_i, \omega_i | \mu_k, \Sigma_k) \right]$$

$$\text{Let } \Lambda_a^{(k)} \text{ and } \Lambda_b^{(k)} \text{ be the } k\text{-th row of } \Lambda_a \text{ and } \Lambda_b^{(k)}, \text{ respectively. The prior distributions of } \pi_k^{(k)} | \mu_k \text{ are given by:}$$

$$(i) \text{ If } \Lambda_a^{(1)} = \dots = \Lambda_a^{(c)} = \Lambda_a, \quad \Lambda_a \stackrel{D}{=} N(\Lambda_{0a}, H_{0a}), \text{ constraint} \\ (ii) \text{ If } \Lambda_a^{(1)} \neq \dots \neq \Lambda_a^{(c)} = \Lambda_a^{(c)}, \quad \Lambda_a^{(c)} \stackrel{D}{=} N(\Lambda_{0a}, \mu_a^{(c)}), \text{ non-constraint.}$$

$$\text{Similarly, let } \psi_k^{(k)} \text{ be the } k\text{-th diagonal element of } \Psi^{(k)}. \text{ The prior distributions of } \psi_k^{(k)} | \mu_k \text{ are:}$$

$$\text{Gamma}_k(\alpha_k^{(1)}, \beta_k^{(1)}), \quad \mu_k^{(1)} \stackrel{D}{=} N(\mu_k^{(1)}, \rho_k^{(1)}), \quad \Psi_k^{(1)} = B_k(\mu_k^{(1)}, \rho_k^{(1)}).$$

$$\text{Let } M_0 \text{ and } M_1 \text{ be two different sets of non-nested competing models (hypotheses). } M_0 \text{ and } M_1 \text{ can then be compared using the Bayes factor:}$$

$$B(O_1, O_2) = \frac{p(O_1 | \mathbf{y}, \omega_1, \mu_1, \Sigma_1)}{p(O_2 | \mathbf{y}, \omega_2, \mu_2, \Sigma_2)}$$

$$\text{In computing } \log \log_B \text{ through path sampling, searching for a good path to link } M_0 \text{ and } M_1 \text{ is crucial. As an illustrative example, suppose that the competing models } M_1 \text{ and } M_2 \text{ are defined as follows: For } g = 1, 2, \dots, N_g,$$

$$\sum_{i=1}^n \sum_{h=c+1}^K \pi_h \sum_{k=1}^c \pi_k \delta_k(\mathbf{y}_i, \omega_i | \mu_k, \Sigma_k) + \sum_{h=c+1}^K \pi_h \delta_h(\mathbf{y}_i, \omega_i | \mu_k, \Sigma_k)$$

$$\text{where } \delta_k(\cdot) \text{ is written as follows:}$$

$$\delta_k(\mathbf{y}_i, \omega_i | \mu_k, \Sigma_k) = (2\pi)^{-d/2} |\Psi_k|^{-1/2} \times \exp \left[\frac{1}{2} (\mathbf{y}_i - \mu_k - \Delta_k \omega_i)^\top \Psi_k^{-1} (\mathbf{y}_i - \mu_k - \Delta_k \omega_i) \right]$$

$$\times \exp \left[-\frac{1}{2} (\mathbf{y}_i - \mu_k - \Delta_k \omega_i)^\top \Pi_k \Psi_k^{-1} (\mathbf{y}_i - \mu_k - \Delta_k \omega_i) \right] \times \exp \left[\frac{1}{2} (\eta - \mu_k - \Delta_k \omega_i)^\top \Psi_k^{-1} (\eta - \mu_k - \Delta_k \omega_i) \right]$$

$$\times (2\pi)^{-d/2} |\Psi_k|^{-1/2} \exp \left[-\frac{1}{2} (\eta - \mu_k - \Delta_k \omega_i)^\top \Psi_k^{-1} \eta \right].$$

and $\Delta_{ik} = (\Pi_k, \Gamma_k)$. Thus, the Bayes factor can be estimated with

$$\tilde{U}_{ij} = \sum_{k=1}^J U_{ijk}(\mathbf{y}_i, \Omega^{(1)}, \theta^{(1)}),$$

Posterior distributions of mixture model.

Step (a): Generate $W^{(1)}, \Omega^{(1)}$ from $p(W, \Omega | \mathbf{Y}; \theta^{(0)})$.

Step (b): Generate $\Omega^{(2)} \sim p(\Omega^{(1)} | \mathbf{Y}; \mathbf{W}^{(1)}, \Omega^{(1)}, \theta^{(1)})$.

Step (c): Resample the label through the permutation sampler to achieve the identifiability.

As $p(\mathbf{W}, \Omega | \mathbf{Y}, \theta) = p(\mathbf{Y}, \Omega | \mathbf{W}, \theta)$, Step (a) can be further decomposed into the following two steps:

Step (a1): Generate $W^{(1)} \sim p(\mathbf{Y}, \Omega^{(1)}, \theta^{(0)})$.

Step (a2): Generate $\Omega^{(1)} \sim p(\mathbf{Y}, W^{(1)}, \Omega^{(1)}, \theta^{(0)})$.

The link model M_{02} for linking M_0 and M_2 is defined as follows:

$$M_{02}: \quad \begin{aligned} \eta^{(1)} &= \mu^{(1)} + \Lambda_{02}^{(1)}(\omega^{(1)} - \epsilon^{(1)}), \quad \text{Hence } \Delta_{02} = \Delta_{02}^{(1)} \\ \eta^{(2)} &= \mu^{(2)} + \Lambda_{02}^{(2)}(\omega^{(2)} + \epsilon^{(2)}), \end{aligned}$$

$$M_2: \quad \begin{aligned} \eta^{(2)} &= \mu^{(2)} + \mu^{(1)} + \Lambda_{22}^{(2)}(\omega^{(2)} + \epsilon^{(2)}) + \epsilon^{(2)}, \quad \text{Hence } \Delta_{22} = \Delta_{22}^{(2)}, \quad \epsilon^{(2)} = \epsilon^{(1)} \\ \eta^{(2)} &= \mu^{(2)} + \Lambda_{22}^{(2)}(\omega^{(2)} + \epsilon^{(2)}) \end{aligned}$$

When $t = 1$, M_0 reduces to M_1 , and when $t = 0$, M_0 reduces to M_2 . The parameter vector in M_0 contains $\mu^{(1)}, \mu^{(2)}, \Lambda_{01}^{(1)}, \Lambda_{02}^{(1)}, \Psi_1^{(1)}, \Gamma_1^{(1)}, \Omega^{(1)}, \Phi^{(1)}, \Psi_2^{(1)}, \Gamma_2^{(1)}$.

The link model M_{02} for linking M_0 and M_2 is defined as follows:

$$M_{02}: \quad \begin{aligned} \eta^{(2)} &= \mu^{(1)} + \Lambda_{02}^{(1)}(\omega^{(1)} - \epsilon^{(1)}) + \epsilon^{(2)}, \quad \text{and } \Delta_{02} = \Delta_{02}^{(1)}, \\ \eta^{(2)} &= \mu^{(2)} + \mu^{(1)} + \Lambda_{22}^{(2)}(\omega^{(2)} + \epsilon^{(2)}) + \epsilon^{(2)}, \quad \text{and } \Delta_{22} = \Delta_{22}^{(2)}. \end{aligned}$$

Moreover,

$$p(\mathbf{y}_i | \mathbf{x}_i, \theta) = \frac{(t-1) + (1-t)}{2} p(\mathbf{y}_i | \mathbf{x}_i, \theta^{(1)}) + \frac{1-t}{2} p(\mathbf{y}_i | \mathbf{x}_i, \theta^{(2)}).$$

Clearly, when $t = 1$, M_0 reduces to M_1 , and when $t = 0$, M_0 reduces to M_2 . The parameter vector in M_0 contains $\mu^{(1)}, \mu^{(2)}, \Lambda_{01}^{(1)}, \Lambda_{02}^{(1)}, \Psi_1^{(1)}, \Gamma_1^{(1)}, \Omega^{(1)}, \Phi^{(1)}, \Psi_2^{(1)}, \Gamma_2^{(1)}$.

We first compute $\log B_{12}$ and $\log B_{21}$, and then obtain $\log B_{12}$ via the following equations:

$$\log B_{12} = \log \left(\frac{p(M_1 | \mathbf{y}, \omega_1, \mu_1, \Sigma_1)}{p(M_2 | \mathbf{y}, \omega_2, \mu_2, \Sigma_2)} \right) = \log B_{12} - \log B_{21}.$$

where $\mathbf{Y}_1 = \sum_{i=1}^n (\mathbf{x}_i - \Delta_{01}^{(1)}) \omega_i + \mu_1$, with $\sum_{i=1}^n$ denotes the summation with respect to those i such that $x_i = 1$, and

$$\log B_{12} = \frac{1}{2} \sum_{i=1}^n \sum_{h=1}^K \sum_{k=1}^c \pi_h \delta_k(\mathbf{y}_i, \omega_i | \mu_k, \Sigma_k),$$

$\delta_k(\cdot) = \sum_{l=1}^L \pi_l \delta_l(\mathbf{y}_i, \omega_i | \mu_k, \Sigma_k)$. $\delta_k(\mathbf{y}_i, \omega_i | \mu_k, \Sigma_k) = \delta_k(\mathbf{y}_i, \omega_i | \mu_k, \Sigma_k, \Delta_{01}^{(1)}, \Delta_{02}^{(1)})$.

$$\delta_k(\mathbf{y}_i, \omega_i | \mu_k, \Sigma_k, \Delta_{01}^{(1)}, \Delta_{02}^{(1)}) = \frac{1}{2} \log \left(\frac{p(\mathbf{y}_i | \mathbf{x}_i, \theta^{(1)})}{p(\mathbf{y}_i | \mathbf{x}_i, \theta^{(2)})} \right),$$

$\Delta_{01}^{(1)} = \Delta_{02}^{(1)} = \Delta_{02}^{(1)}(\mathbf{y}_i, \omega_i | \mu_k, \Sigma_k, \Delta_{01}^{(1)})$.

$\Delta_{02}^{(1)} = \Delta_{02}^{(1)}(\mathbf{y}_i, \omega_i | \mu_k, \Sigma_k, \Delta_{01}^{(1)})$.

$\Delta_{01}^{(1)} = \Delta_{01}^{(1)}(\mathbf{y}_i, \omega_i | \mu_k, \Sigma_k, \Delta_{01}^{(1)})$.

$\Delta_{02}^{(1)} = \Delta_{02}^{(1)}(\mathbf{y}_i, \omega_i | \mu_k, \Sigma_k, \Delta_{01}^{(1)})$.

$\Delta_{01}^{(1)} = \Delta_{01}^{(1)}(\mathbf{y}_i, \omega_i | \mu_k, \Sigma_k, \Delta_{01}^{(1)})$.

$\Delta_{02}^{(1)} = \Delta_{02}^{(1)}(\mathbf{y}_i, \omega_i | \mu_k, \Sigma_k, \Delta_{01}^{(1)})$.

$\Delta_{01}^{(1)} = \Delta_{01}^{(1)}(\mathbf{y}_i, \omega_i | \mu_k, \Sigma_k, \Delta_{01}^{(1)})$.

$\Delta_{02}^{(1)} = \Delta_{02}^{(1)}(\mathbf{y}_i, \omega_i | \mu_k, \Sigma_k, \Delta_{01}^{(1)})$.

$\Delta_{01}^{(1)} = \Delta_{01}^{(1)}(\mathbf{y}_i, \omega_i | \mu_k, \Sigma_k, \Delta_{01}^{(1)})$.

$\Delta_{02}^{(1)} = \Delta_{02}^{(1)}(\mathbf{y}_i, \omega_i | \mu_k, \Sigma_k, \Delta_{01}^{(1)})$.

$\Delta_{01}^{(1)} = \Delta_{01}^{(1)}(\mathbf{y}_i, \omega_i | \mu_k, \Sigma_k, \Delta_{01}^{(1)})$.

$\Delta_{02}^{(1)} = \Delta_{02}^{(1)}(\mathbf{y}_i, \omega_i | \mu_k, \Sigma_k, \Delta_{01}^{(1)})$.

$\Delta_{01}^{(1)} = \Delta_{01}^{(1)}(\mathbf{y}_i, \omega_i | \mu_k, \Sigma_k, \Delta_{01}^{(1)})$.

$\Delta_{02}^{(1)} = \Delta_{02}^{(1)}(\mathbf{y}_i, \omega_i | \mu_k, \Sigma_k, \Delta_{01}^{(1)})$.

$\Delta_{01}^{(1)} = \Delta_{01}^{(1)}(\mathbf{y}_i, \omega_i | \mu_k, \Sigma_k, \Delta_{01}^{(1)})$.

$\Delta_{02}^{(1)} = \Delta_{02}^{(1)}(\mathbf{y}_i, \omega_i | \mu_k, \Sigma_k, \Delta_{01}^{(1)})$.

$\Delta_{01}^{(1)} = \Delta_{01}^{(1)}(\mathbf{y}_i, \omega_i | \mu_k, \Sigma_k, \Delta_{01}^{(1)})$.

$\Delta_{02}^{(1)} = \Delta_{02}^{(1)}(\mathbf{y}_i, \omega_i | \mu_k, \Sigma_k, \Delta_{01}^{(1)})$.

$\Delta_{01}^{(1)} = \Delta_{01}^{(1)}(\mathbf{y}_i, \omega_i | \mu_k, \Sigma_k, \Delta_{01}^{(1)})$.

$\Delta_{02}^{(1)} = \Delta_{02}^{(1)}(\mathbf{y}_i, \omega_i | \mu_k, \Sigma_k, \Delta_{01}^{(1)})$.

$\Delta_{01}^{(1)} = \Delta_{01}^{(1)}(\mathbf{y}_i, \omega_i | \mu_k, \Sigma_k, \Delta_{01}^{(1)})$.

$\Delta_{02}^{(1)} = \Delta_{02}^{(1)}(\mathbf{y}_i, \omega_i | \mu_k, \Sigma_k, \Delta_{01}^{(1)})$.

$\Delta_{01}^{(1)} = \Delta_{01}^{(1)}(\mathbf{y}_i, \omega_i | \mu_k, \Sigma_k, \Delta_{01}^{(1)})$.

$\Delta_{02}^{(1)} = \Delta_{02}^{(1)}(\mathbf{y}_i, \omega_i | \mu_k, \Sigma_k, \Delta_{01}^{(1)})$.

$\Delta_{01}^{(1)} = \Delta_{01}^{(1)}(\mathbf{y}_i, \omega_i | \mu_k, \Sigma_k, \Delta_{01}^{(1)})$.

$\Delta_{02}^{(1)} = \Delta_{02}^{(1)}(\mathbf{y}_i, \omega_i | \mu_k, \Sigma_k, \Delta_{01}^{(1)})$.

$\Delta_{01}^{(1)} = \Delta_{01}^{(1)}(\mathbf{y}_i, \omega_i | \mu_k, \Sigma_k, \Delta_{01}^{(1)})$.

$\Delta_{02}^{(1)} = \Delta_{02}^{(1)}(\mathbf{y}_i, \omega_i | \mu_k, \Sigma_k, \Delta_{01}^{(1)})$.

$\Delta_{01}^{(1)} = \Delta_{01}^{(1)}(\mathbf{y}_i, \omega_i | \mu_k, \Sigma_k, \Delta_{01}^{(1)})$.

$\Delta_{02}^{(1)} = \Delta_{02}^{(1)}(\mathbf{y}_i, \omega_i | \mu_k, \Sigma_k, \Delta_{01}^{(1)})$.

$\Delta_{01}^{(1)} = \Delta_{01}^{(1)}(\mathbf{y}_i, \omega_i | \mu_k, \Sigma_k, \Delta_{01}^{(1)})$.

