

STAT 24410

Lec 1A - 9/30

- Midterm in Nov.
- 6 or 7 HW

Example: Genome-Wide Assoc Studies (GWAS)

- Genetic Association

i) cases & controls

ii) disease risk increasing mutations more freq. in cases

$$\hookrightarrow P_m^{\text{cases}} > P_m^{\text{controls}}$$

iii) Single-nucleotide polymorphism (SNP)

$\hookrightarrow \frac{A}{G} \rightarrow$ look if $P_A^{\text{cases}} \neq P_A^{\text{control}}$, if so, close to mutation location

\hookrightarrow to do genome-wide, need millions of SNPs

iv) Say n cases, m controls

	A	G	count?
cases	n_A	n_G	$2n$
controls	m_A	m_G	$2m$

\hookrightarrow # chromosomes

v) Estimate P_A^{cases} ?

$$\hookrightarrow \hat{P}_A^{\text{cases}} = \frac{n_A}{2n}$$

\hookrightarrow why believe this? well here it's the natural estimator

\hookrightarrow in more complicated scenarios?

\hookrightarrow it's unbiased ($E[\hat{P}_A] = P_A$), good accuracy

\hookrightarrow we'll learn what makes a best estimator & r

vi) Testing. $H_0: P_A^{\text{cases}} = P_A^{\text{controls}}$, $H_a: P_A^{\text{cases}} \neq P_A^{\text{controls}}$

\hookrightarrow also need a test statistic & method

$$\hookrightarrow T = \frac{(\hat{p}_{cases} - \hat{p}_{controls})^2}{(\frac{1}{2n} + \frac{1}{2m}) \hat{p}_A (1 - \hat{p}_A)}, \quad \hat{p}_A = \frac{n_A + m_A}{2n + 2m}$$

freq. A in
combined data

+ we claim $T \stackrel{H_0}{\sim} \chi_1^2$. Why?

a) $n_A \sim \text{Bin}$ (sum of Bernoulli); LLT gives approx. \mathcal{N}

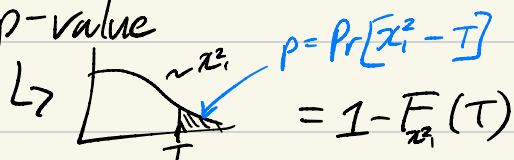
b) $Z \sim \mathcal{N}(0, 1) \Rightarrow Z^2 \sim \chi_1^2$

c) this is an approximation, best at the center, where our data is

- is this the "best" test stat.?

\hookrightarrow criteria to judge: power, type I error

vii) p-value



\hookrightarrow in fact, $p \stackrel{H_0}{\sim} \text{Unif}(0, 1)$ b/c $F_{\chi^2_1}(T) \sim \text{Unif}(0, 1)$

\hookrightarrow only an approx. b/c T was an approx.

viii) 1 million SNPs

\hookrightarrow now we have p_1, p_2, \dots, p_{1m}

\hookrightarrow w/ this many p-values, $p < 0.01$ isn't exciting (even under H_0)

- first, suppose we care abt the smallest p-val: $p_{(1)} < p_{(2)} < \dots < p_{(1m)}$

\hookrightarrow turns out $p_{(1)} \sim \text{Beta}(1, n)$, w/ EV $\frac{1}{n+1}$, meaning by chance, we expect to

see a $p = \frac{1}{1m+1}$

RVs

- Ω = sample space (e.g. \mathbb{R})

- subsets of Ω = events, elements = sample points

Def: σ -Field. \mathcal{F} is a collection of events s.t.

i) $\Omega \in \mathcal{F}$

ii) $A \in \mathcal{F} \Rightarrow A^c \in \mathcal{F}$

iii) $A_n \in \mathcal{F} \Rightarrow \bigcup_n A_n \in \mathcal{F}$

Eg: Borel σ -field.

$\hookrightarrow \Omega = \mathbb{R}$

$\hookrightarrow \mathcal{B}$ = smallest σ -field that contains open intervals (a, b)

$\hookrightarrow [(a, b)]$ also in \mathcal{B}

Def: Prob. Measure (P). Given (Ω, \mathcal{F}) , $P: \mathcal{F} \rightarrow [0, 1]$ s.t.

i) $P(\Omega) = 1$

ii) $P(\bigcup_n A_n) = \sum_n P(A_n)$ for $A_i \cap A_j = \emptyset$

Def: RV. on $\sim (\Omega, \mathcal{F}, P)$. X is an RV if $X: \mathcal{F} \rightarrow \mathbb{R}$ such that $\forall B \in \mathcal{B}$, $X^{-1}(B) \in \mathcal{F}$.

Eg: 2 RVs.

1) $(\Omega_1, \mathcal{F}_1, P_1)$ s.t. $\Omega_1 = \{H, T\}$, $\mathcal{F}_1 =$ all subsets, $P_1(\{H\}) = P_1(\{T\}) = \frac{1}{2}$

$$X_1(\omega) = \begin{cases} 1 & \text{if } H \\ 0 & \text{if } T \end{cases}$$

2) $\Omega_2 = [0, 1]$, $\mathcal{F}_2 = \mathcal{B}$, $P_2([b-a]) = b-a$. $X_2(\omega) = \begin{cases} 1 & \text{if } \omega \in [0, \frac{1}{2}] \\ 0 & \text{otherwise} \end{cases}$

we will only care
→ into the dists from here!

- even though these prob. spaces are really diff, the dists are the same

$$\hookrightarrow X_1 \stackrel{d}{=} X_2$$

Def. X_1 and X_2 are equal in dist. if $\Pr[X_1 \in A] = \Pr[X_2 \in A] \forall A$.

↳ basically an abstraction

Lec 1B - 10/2

- recall $X_1 \stackrel{D}{=} X_2 \Rightarrow \Pr[X_1 \in A] = \Pr[X_2 \in A]$

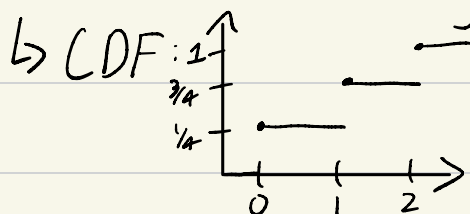
Def: Dist. Fn. X is RV. The CDF of X (denoted by F_X) is a fn.

$$F_X: \mathbb{R} \rightarrow [0, 1] \text{ s.t. } F_X(x) = \Pr[X \in (-\infty, x]] = \Pr[X \leq x]$$

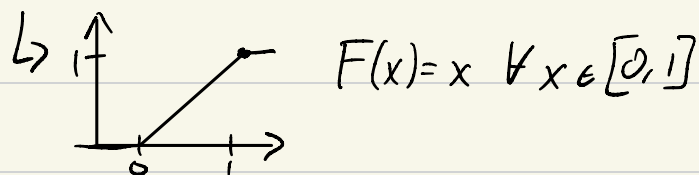
Lemma. $X \stackrel{D}{=} Y$ iff $F_X(t) = F_Y(t) \forall t \in \mathbb{R}$. No proof.

Eg.

1) Binomial $(2, \frac{1}{2})$ (tossing 2 coins)



2) Unit $(0, 1)$



CDF Properties.

1) $x \leq y \Rightarrow F(x) \leq F(y)$

2) right-continuous

3) $\lim_{x \rightarrow -\infty} F(x) = 0, \lim_{x \rightarrow \infty} F(x) = 1$

Questions

a) if $F: \mathbb{R} \rightarrow [0, 1]$ has properties 1-3, is F a cdf?

b) F is a cdf. How do we simulate RVs w/ this cdf?

Inverse of the Dist. Fn.

+ Case 1: F is continuous & strictly increasing.

↳ $F^{-1}: [0, 1] \rightarrow \mathbb{R}$ exists; continuous & strictly increasing

- Result 1: say X is RV w/ cdf F

↳ $F(X)$ is a RV

↳ $\Pr[F(X) \leq u] = \Pr[X \leq F^{-1}(u)] = F(F^{-1}(u)) = u \rightarrow \text{Unif}(0, 1)!$

⇒ $F(X) \sim \text{Unif}(0, 1)$

- this does not generalize; counterexample = discrete RVs

- Result 2: Let $U \sim \text{Unif}(0, 1)$. $F^{-1}(U)$ is a RV

↳ $\Pr[F^{-1}(U) \leq x] = \Pr[U \leq F(x)] = F(x) \Rightarrow F^{-1}(U) \stackrel{d}{=} X$

- thus, we can simulate a RV w/ F^{-1} by simulating $\text{Unif}(0, 1)$, $F^{-1}(U)$

- result 2 generalizes

↳ $F^{-1}(u) = \inf\{x : u \leq F(x)\}$ (Q2)

↳ it can be proved that $F^{-1}(U) \stackrel{d}{=} X$ (Q1)

- these results relevant for p-values (by simulation)

- and ordered statistics

↳ X_1, \dots, X_n iid F , $X_{(1)} \leq X_{(2)} \leq \dots \leq X_{(n)}$

↳ let U_1, \dots, U_n iid $U(0,1)$

↳ now $F^{-1}(U_1), \dots, F^{-1}(U_n)$ iid F , can show $F^{-1}(U_{(k)}) =_{\mathcal{D}} X_{(k)}$

↳ can learn dist. of $X_{(k)}$ from dist. of $U_{(k)} \sim \text{Beta}(k, n-k+1)$

RVs

Def: Discrete RV. Has a countable support: $S = \{x_1, x_2, \dots\}$, $\Pr[X \in S] = 1$

Notes

a) F_X is piecewise constant if S is the set of isolated points

ϵ ball, no other points w/in

b) Dist. is completely determined by $\Pr[X = x_i]$, called the prob. mass function.

→ pmf

Eg.

1) Binomial(n, p), $0 < p < 1$, $n \geq 1$

↳ # successes in n trials

↳ $S = \{0, 1, \dots, n\}$

↳ $\Pr[X = k] = \binom{n}{k} p^k (1-p)^{n-k}$

- Bin($1, p$) = Bernoulli(p)

↳ if $X_1, \dots, X_n \sim \text{Bernoulli}(p)$, $\sum_{i=1}^n X_i \sim \text{Bin}(n, p)$

2) Poisson(λ), $\lambda > 0$

↳ $S = \{0, 1, \dots\} = \mathbb{N}$

↳ memoryless arrivals

↳ $\Pr[X = k] = e^{-\lambda} \frac{\lambda^k}{k!}$

- if $\lim_{n \rightarrow \infty} np_n = \lambda$, $\lim_{n \rightarrow \infty} \binom{n}{k} p_n^k (1-p_n)^{n-k} = e^{-\lambda} \frac{\lambda^k}{k!}$

↳ Bin \rightarrow large, but np_n converges, approx. by poisson

Pf.

$$\binom{n}{k} p_n^k (1-p_n)^{n-k} = \frac{n!}{k!(n-k)!} p_n^k (1-p_n)^{n-k} = \frac{1}{k!} \cdot \frac{(n)(n-1)\dots(n-k+1)}{n^k} (np_n)^k (1-p_n)^{n-k} \left(1 - \frac{np_n}{n}\right)^n$$

in limit...
same
1
 λ^k
1
 $e^{-\lambda}$
 $\frac{a_n \rightarrow a}{(1 + \frac{a_n}{n})^n} \rightarrow e^a$

- $X \sim \text{Poi}(\lambda_1)$, $Y \sim \text{Poi}(\lambda_2) \Rightarrow X+Y \sim \text{Poi}(\lambda_1 + \lambda_2)$

↳ proof = exercise

3) Geometric (in notes)

+ Transformations

- X is discrete w/ support S , $f: \mathbb{R} \rightarrow \mathbb{R}$, $f(x) = ?$

↳ $f(X)$ is discrete

↳ $\Pr[f(X) = y] = \sum_{\substack{x \in S \\ f(x) = y}} \Pr[X = x]$

Def: Continuous RVs. RV X has a density if $\exists f: \mathbb{R} \rightarrow [0, \infty)$ s.t.

$\Pr[X \in A] = \int_A f(x) dx$. ($\forall A$ in Borel)

- $1 = \Pr[X \in \mathbb{R}] = \int_{\mathbb{R}} f(x) dx$

- density is not unique

↳ can modify a countable # of points w/o changing integral

$$- F(x) = P_r[X \in (-\infty, x]] = \int_{-\infty}^x f(t) dt$$

$$- f(x) = F'(x)$$

↳ remember $U(0,1)$; not differentiable everywhere

Eg.

1) $U(a,b)$

$$\hookrightarrow f(x) = \begin{cases} \frac{1}{b-a} & x \in [a,b] \\ 0 & \text{o.v.} \end{cases}, \quad F(x) = \begin{cases} 0 & x < a \\ \frac{x-a}{b-a} & x \in [a,b] \\ 1 & x > b \end{cases}$$

2) $N(\mu, \sigma^2)$

$$\hookrightarrow f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}, \quad x \in \mathbb{R}$$

↳ standard normal: $Z \sim N(0,1)$.

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

$$\hookrightarrow F(x) = \Phi(x) = \int_{-\infty}^x \frac{1}{\sqrt{2\pi}} e^{-\frac{t^2}{2}} dt$$

3) read a/b gamma, exp, beta

Lec 2A - 10/7

Transformations of Cont. RVs

- say X is a RV w/ density f_x

$$\hookrightarrow h: \mathbb{R} \rightarrow \mathbb{R}, Y = h(X)$$

\hookrightarrow what is f_y ?

+ Case 1: h strictly increasing & differentiable

$$\begin{aligned} F_Y(y) &= \Pr[\bar{Y} \leq y] = \Pr[\bar{h}(X) \leq y] = \Pr[X \leq h^{-1}(y)] = F_X(h^{-1}(y)) \\ \Rightarrow f_Y(y) &= F_X'(h^{-1}(y)) = f(h^{-1}(y)) \frac{\partial h^{-1}(y)}{\partial y} \end{aligned}$$

+ General Case: h is 1-to-1, differentiable

$$f_Y(y) = f_X(h^{-1}(y)) \left| \frac{\partial h^{-1}(y)}{\partial y} \right|$$

Eg. $X \sim \mathcal{U}(-\frac{\pi}{2}, \frac{\pi}{2})$, $h(x) = \tan(x)$

$$\Rightarrow f_X = \left\{ \frac{1}{\pi} \quad x \in \left[-\frac{\pi}{2}, \frac{\pi}{2}\right] \right.$$

\hookrightarrow note $h: (-\frac{\pi}{2}, \frac{\pi}{2}) \rightarrow \mathbb{R}$, $h^{-1}(y) = \arctan(y)$, $(h^{-1})'(y) = \frac{1}{1+y^2}$

$$\Rightarrow f_Y = \frac{1}{\pi} \frac{1}{1+y^2}, y \in \mathbb{R}$$

\hookrightarrow Fact: $Y \sim$ Cauchy

Eg. Many-to-one transformations.

$$X \sim \mathcal{N}(0,1), Y = X^2 \quad (Y \sim \chi^2_1)$$

$$\begin{aligned} \text{Look into } y > 0 &= F_Y(y) = \Pr[\bar{Y} \leq y] = \Pr[X^2 \leq y] \\ &= \Pr[-\sqrt{y} \leq X \leq \sqrt{y}] \end{aligned}$$

$$= \int_{-\sqrt{y}}^{\sqrt{y}} \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} dx = \frac{2}{\sqrt{2\pi}} \int_0^{\sqrt{y}} e^{-\frac{x^2}{2}} dx$$

Letting $x^2 = t$: $\frac{2}{\sqrt{2\pi}} \int_0^{\sqrt{y}} e^{-\frac{t}{2}} \frac{1}{2\sqrt{t}} dt$
 $= \int_0^y \frac{1}{\sqrt{2\pi}} e^{-\frac{t}{2}} \frac{1}{\sqrt{t}} dt$

Thus, $f_Y(y) = \begin{cases} \frac{1}{\sqrt{2\pi}} e^{-\frac{y}{2}} y^{-1/2} & : y > 0, \\ 0 & \text{o.w.} \end{cases}$

disjoint, h is 1-to-1 on each.

- "The trick": $\{X^2 \leq x\} = \{0 \leq x \leq \sqrt{x}\} \cup \{-\sqrt{x} \leq x \leq 0\}$

Expectation

Intuition: $E[X] = \sum x \Pr[X \in (x-\epsilon, x]]$
 $= \sum x (F(x) - F(x-\epsilon))$
 $\epsilon \rightarrow 0: \int_{-\infty}^{\infty} x dF(x)$

+ Case 1: $X \geq 0$ Riemann-Stieltjes

$$E[X] = \int_0^{\infty} x dF(x)$$

$$= \begin{cases} \int_0^{\infty} x f(x) dx & \text{if } X \text{ has density } f \\ \sum_x x \Pr[X=x] & X \text{ discrete} \end{cases}$$

+ General X

Let $X^+ = \max\{x, 0\}$, and $X^- = \max\{-x, 0\}$. Then $X = X^+ - X^-$.

Def. X has an e.v. if at least one of $E[X^+]$ and $E[X^-]$ is finite. Then $E[X] = E[X^+] - E[X^-]$.

E.g. $X \sim \text{Cauchy}$; $f_x(x) = \frac{1}{\pi} \frac{1}{1+x^2}$, $x \in \mathbb{R}$

$$E x^+ = \int_0^{\infty} x \cdot \frac{1}{\pi(1+x^2)} dx$$

$$= \frac{1}{2\pi} \log(1+x^2) \Big|_0^{\infty} = +\infty$$

$$E x^- = -\infty$$

} no e.v.!

Properties.

1) X has a density f s.t. $\int_{-\infty}^{\infty} |x| f(x) dx < \infty \Rightarrow X$ has a finite expectation:
$$E[X] = \int_{-\infty}^{\infty} x f(x) dx$$

2) $E[X+Y] = E[X] + E[Y]$

3) $E[cX] = c E[X]$

4) $X \leq Y \Rightarrow E[X] \leq E[Y]$

5) Def. $X \preceq Y$ (Y stochastically dominates X) if $\Pr[X > t] \leq \Pr[Y > t]$.

$$X \preceq Y \Rightarrow E[X] \leq E[Y]$$

6) $E[\mathbb{1}_A(X)] = \Pr[X \in A]$

7) $X \perp Y \Rightarrow E[XY] = E[X] \cdot E[Y]$

8) $E[h(X)] = \int_{-\infty}^{\infty} h(x) f_x(x) dx$

9) $\text{Var}(X) = E[(X - E[X])^2]$ (if $E[X] < \infty$)

↳ measure of spread

a) $\text{Var}(X) = E[X^2] - E[X]^2$

b) $\text{Var}(X) = 0 \Leftrightarrow \Pr[X = E[X]] = 1$

$$c) \text{Var}(c \cdot X) = c^2 \text{Var}(X)$$

$$d) X \perp Y \Rightarrow \text{Var}(X+Y) = \text{Var}(X) + \text{Var}(Y)$$

10) Mixed Distributions

$$\hookrightarrow \text{eg: } Z \sim \mathcal{N}(0,1), Z^+ = \max(Z, 0)$$

$$- \Pr[Z=0] = \frac{1}{2} : \text{cdf}$$

- no density; no pmf

Lemma: $X \geq 0$ s.t. F_X is cont. differentiable on $[0, \infty) \setminus S$, where S is a countable set of isolated points $\{x_1, \dots, x_n\}$.

$$\Rightarrow E[X] = \sum_S x_i \Pr[X=x_i] + \int_0^\infty x F'_X(x) dx$$

\hookrightarrow back to Z^+

$$\begin{aligned} E[Z^+] &= 0 \cdot \Pr[Z^+=0] + \int_0^\infty x \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} dx \\ &= \frac{1}{\sqrt{2\pi}} - e^{-\frac{x^2}{2}} \Big|_0^\infty = \frac{1}{\sqrt{2\pi}} \end{aligned}$$

11) Probability Inequalities

a) Markov

$$X \geq 0, c > 0 \Rightarrow \Pr[X \geq c] \leq \frac{E[X]}{c}$$

\hookrightarrow equality if $\Pr[X=0, X=c] = 1$

b) Chebyshev

$$c > 0; \Pr[|X-\mu| \geq c] \leq \frac{\text{Var}(X)}{c^2}$$

c) Chernoff

$$c > 0, t \in \mathbb{R}; \Pr[X \geq c] \leq e^{-tc} E[e^{tx}]$$

$$\Pr[X \leq c] \leq e^{tc} E[e^{-tx}]$$

Proofs.

a) $Y = c$ if $X \geq c$, else 0 .

$$Y \leq X \Rightarrow E[Y] \leq E[X] \Rightarrow c \Pr[X \geq c] \leq E[X]$$

b) Markov for $(X - \mu)^2$

c) Markov for e^{tx} , e^{-tx}

Application 1. Weak Law of Large #s.

Suppose X_1, \dots, X_n iid RVs s.t. $E[X] = \mu < \infty$, $\text{Var}(X) = \sigma^2 < \infty$.

$$\text{Let } \bar{X}_n = \frac{X_1 + \dots + X_n}{n}$$

$$\Rightarrow \Pr[|\bar{X}_n - \mu| < \epsilon] \rightarrow 1 \text{ as } n \rightarrow \infty$$

\hookrightarrow "converges in probability": $\bar{X}_n \xrightarrow{p} \mu$

Pf. $E[\bar{X}_n] = \mu$, $\text{Var}(\bar{X}_n) = \frac{\sigma^2}{n}$. By Chebyshev, $\Pr[|\bar{X}_n - \mu| \geq \epsilon] \leq \frac{\sigma^2}{n\epsilon^2} \xrightarrow{n \rightarrow \infty} 0$.

Application 2: Large Deviations.

X_1, \dots, X_n iid RVs, finite expectations.

What is $\lim_{n \rightarrow \infty} \frac{1}{n} \log(\Pr[\bar{X}_n \geq c])$ ($c > \mu$)?

$$\hookrightarrow \text{Chernoff: } \Pr[\sum_i X_i \geq nc] \leq e^{-tnc} E[e^{t(\sum_i X_i)}] = e^{-tnc + nk(t)} \rightarrow k(t) = \log E[e^{tx}]$$

(proof not on midterm)
= cumulative generating fn.

$$\Rightarrow \frac{1}{n} \log \Pr[\bar{X}_n \geq c] \leq -[tc - k(t)] \forall t$$

$$\Rightarrow \leq -\sup_t [tc - k(t)] = -I(c) \quad (\text{in fact, the lim} = I(c)!) \quad \leftarrow \text{Legendre Transform}$$

$\Pr[\bar{X}_n \geq c] \rightarrow 0$, but how fast?

Lec 2B-10/9

CLT

Def: Convergence in Distribution. X_n converges in dist to X ($X_n \xrightarrow{d} X$)
iff $F_{X_n}(t) \rightarrow F_X(t) \forall t$.

Def: Moment Gen. Fn. $M: \mathbb{R} \rightarrow [0, +\infty]$, $M(t) = E[e^{tx}] \forall t$ (mgf).

Notes:

a) $e^{tx} = \sum_{n=0}^{\infty} \frac{(tx)^n}{n!}$. $E[e^{tx}] = \sum_{n=0}^{\infty} E[X^n] \frac{t^n}{n!}$ → only true if finite sum or Fubini Thm (not important)

↳ so $M^{(n)}(0) = E[X^n] = n$ th moment of X

b) $M_{ax}(t) = M_x(at)$

c) $X \perp Y$, $M_{X+Y}(t) = M_X(t) M_Y(t)$

d) X_1, \dots, X_n iid, $M_{\sum X_i}(t) = (M_X(t))^n$

e) $X_n \xrightarrow{d} X$ iff $M_{X_n}(t) \rightarrow M_X(t)$ in a neighborhood of 0

↳ can check convergence of distributions by conv. of fn.s!

→ powerful; easy to find moments of cumulation of iid RVs

CLT Proof.

X_1, \dots, X_n iid $E[X_i] = 0$, $\text{Var}(X_i) = \sigma^2$. WTS: $\frac{\sum X_i}{\sqrt{n}} \xrightarrow{d} \mathcal{N}(?, ?)$

↳ $M_{\frac{\sum X_i}{\sqrt{n}}}(t) = \left[M_X\left(\frac{t}{\sqrt{n}}\right) \right]^n$

↳ Taylor: $M_X\left(\frac{t}{\sqrt{n}}\right) = M_X(0) + \frac{t}{\sqrt{n}} M_X'(0) + \frac{t^2}{2n} M_X''(0) + \dots$

$\approx 1 + 0 + \frac{t^2}{2n} \sigma^2$

$\Rightarrow \left[1 + \frac{t^2 \sigma^2}{2n} \right]^n \xrightarrow{n \rightarrow \infty} e^{\frac{t^2 \sigma^2}{2}} = \text{mgf for } \mathcal{N}(0, \sigma^2)$

Complete proof uses characteristic fn: $E[e^{itx}]$
↳ $|e^{itx}| = 1$; nice integrals, easier math

Random Vectors

$$X = (X_1, \dots, X_n) : (\Omega, \mathcal{F}, Pr) \rightarrow \mathbb{R}^n$$

Def. X_1, \dots, X_n are indep. if $Pr[X_1 \in A_1, \dots, X_n \in A_n] = \prod_{i=1}^n Pr[X_i \in A_i]$

Def. cdf. $F_X : \mathbb{R}^n \rightarrow [0, 1]$, $F_X(x) = Pr[X_1 \leq x_1, \dots, X_n \leq x_n]$

Def: Discrete RVeas. $\exists S \subseteq \mathbb{R}^n$ countable s.t. $Pr[X \in S] = 1$.

Def. Cont. RVeas. X has a density $f : \mathbb{R}^n \rightarrow [0, \infty)$ s.t. $Pr[X \in A] = \int_A f(x) dx$

Notes

$$1) \int_{\mathbb{R}^n} f(x) dx = 1$$

$$2) F(x) = \int_{-\infty}^{x_1} \dots \int_{-\infty}^{x_n} f(y_1, \dots, y_n) dy_1 \dots dy_n$$

Thm. If f is cont. & differentiable, X has a density = $\frac{\partial^n F}{\partial x_1 \dots \partial x_n}$

If X_1, \dots, X_n ind, $F_X(x) = F_{X_1}(x_1) \dots F_{X_n}(x_n)$

\hookrightarrow if F differentiable: $f_X(x) = f_{X_1}(x_1) \dots f_{X_n}(x_n) \rightarrow$ converse also true; ^{can go} other way

Marginal Distributions

$X = (X_1, \dots, X_n)$ w/ density f . Let $I \subseteq \{1, 2, \dots, n\} \stackrel{= S_n}{}$ and $X_I = \{X_i\}_{i \in I}$.

Then, $f_{X_I}(x_I) = \int_{\mathbb{R}^{n-|I|}} f(x_{I_c}, \dots, x_{S_n \setminus I}) dx_{S_n \setminus I}$, $S_n = \{1, 2, \dots, n\}$

Transformations → open region

$$X: (\Omega, \mathcal{F}, P_r): H \subseteq \mathbb{R}^n$$

$h: H \rightarrow G \subseteq \mathbb{R}^n \rightarrow$ open region s.t. h is 1-to-1 differentiable and

$h^{-1}: G \rightarrow H$ is differentiable.

↳ X has a density f_x

↳ Let $Y = h(X)$. Y has a density → abs.

$$f_Y(y) = \begin{cases} f_X(h^{-1}(y)) |J(h^{-1}(y))| & y \in G \\ 0 & \text{o.w.} \end{cases}$$

Def: **Jacobian.** $g: G \rightarrow H \subseteq \mathbb{R}^n$. $J(g(y)) = \det\left(\frac{\partial g_i}{\partial y_j}\right)$

E.g.

$X \sim \mathcal{N}(0, 1) \perp Y \sim \chi_n^2$. $T = \frac{X}{\sqrt{Y/n}} \sim t_n$. Density?

↳ the defn. of the transform is $\mathbb{R}^n \rightarrow \mathbb{R}^n$, but this is $\mathbb{R}^2 \rightarrow \mathbb{R}^1$...

↳ Plan: $\begin{pmatrix} X \\ Y \end{pmatrix} \rightarrow \begin{pmatrix} T \\ Y \end{pmatrix} \xrightarrow[\text{out } Y]{\text{integrate}} T$ | $\Gamma(\lambda) = \int_0^\infty e^{-x} x^{\lambda-1} dx, \lambda > 0$

$$f_X(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}, x \in \mathbb{R}$$

$$\Gamma(1) = 1, \Gamma(1/2) = \sqrt{\pi}$$

$$f_Y(y) = \frac{1}{2^{n/2} \Gamma(n/2)} e^{-\frac{y}{2}} y^{\frac{n}{2}-1}, y > 0$$

$$\Gamma(\lambda+1) = (\lambda) \Gamma(\lambda)$$

$$h(x, y) = \left(\frac{\sqrt{y} x}{y}, y\right), h: \mathbb{R} \times (0, \infty) \rightarrow \mathbb{R} \times (0, \infty)$$

$$h^{-1}? \begin{cases} \frac{\sqrt{y} x}{y} = u \\ y = v \end{cases} \Rightarrow \begin{cases} x = \frac{\sqrt{v} u}{\sqrt{v}} \\ y = v \end{cases} \Rightarrow h^{-1}(u, v) = \left(\frac{\sqrt{v} u}{\sqrt{v}}, v\right)$$

$$\text{Jacobian: } \begin{vmatrix} \frac{\sqrt{v}}{\sqrt{v}} & \frac{u}{2\sqrt{v}} \\ 0 & 1 \end{vmatrix} = \frac{\sqrt{v}}{\sqrt{v}} \Rightarrow |J| = \frac{\sqrt{v}}{\sqrt{v}}$$

Density of $h(x, y)$ is

$$g(u, v) = f_x\left(\frac{\sqrt{v}u}{\sqrt{n}}\right) f_y(v) \frac{\sqrt{v}}{\sqrt{n}}$$

$$= \frac{1}{\sqrt{2\pi}} e^{-\frac{vu^2}{2n}} \frac{1}{2^{n/2} \Gamma(\frac{n}{2})} e^{-\frac{v}{2}} v^{\frac{n}{2}-1} \frac{\sqrt{v}}{\sqrt{n}}$$

$$\Rightarrow f_T(t) = \int_0^{\infty} g(u, v) dv$$

Lec 3A - 10/14

- recall that RVec transformations are $\begin{pmatrix} X \\ Y \end{pmatrix} \Rightarrow \begin{pmatrix} X \\ T \end{pmatrix} \Rightarrow T$

Mixed Distributions

Eg. $X \sim \mathcal{N}(0,1)$, $Y \sim \text{Bernoulli}(p)$

$\hookrightarrow (X, Y)$ not discrete: $\Pr[(X, Y) = (x, y)] = 0$

\hookrightarrow no density

Eg. (X, X^2) , $X \sim \mathcal{N}(0,1)$

\hookrightarrow only takes values on the (X, X^2) parabola

$\hookrightarrow \Rightarrow \Pr[(X, X^2) \in (a, b) \times (c, d)] = 0$

Cov & Corr

Def. (X, Y) RVec, $E[X] = \mu_x < \infty$, $E[Y] = \mu_y < \infty$. The covariance of X and Y is $\text{Cov}(X, Y) = E[(X - \mu_x)(Y - \mu_y)]$

Def. Correlation of X & Y : $\text{Corr}(X, Y) = \frac{\text{Cov}(X, Y)}{\sigma_x \sigma_y}$

$\hookrightarrow \sigma_x = \text{std. dev. of } X, \sigma_x < \infty$ (same for Y)

Properties

1) $\text{Cov}(a+bX, c+dY) = bd \cdot \text{Cov}(X, Y)$

2) $\text{Cov}(X, X) = \text{Var}(X)$

3) $\text{Var}(X+Y) = \text{Var}(X) + \text{Var}(Y) + 2 \text{Cov}(X, Y)$

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

5) $\text{Cov}(X, Y) = 0 \not\Rightarrow X \perp Y$

$\left\{ \begin{array}{l} Z \sim \mathcal{N}(0,1) \\ S, T \text{ random signs} \mid \text{all } \perp \\ X = SZ, Y = TZ \\ \text{Cov}(X, Y) = 0, \text{ but } X \not\perp Y \end{array} \right.$ \rightarrow knowing X gives info abt Y up to a sign

Conditional Expectation

+ the Plan

↳ (X, Y) RVec

↳ get dist. $F_{X|Y=y} \rightarrow E[X|Y=y] = \text{fn. of } y$

↳ so $E[X|Y] = f(y) = \text{a RV}$

Eg. Y discrete: $\Pr[X \in A, Y \in B] = \sum_{y \in B} \Pr[X \in A | Y=y] \Pr[Y=y]$
 $= \int_B \Pr[X \in A | Y=y] dF_Y(y)$

Def: Conditional Dist of $X|Y=y$. It is $q: \mathcal{F} \times \Omega \rightarrow [0, 1]$ s.t.

$$\Pr[X \in A, Y \in B] = \int_B q(A, y) dF_Y(y)$$

Def: Corresponding CDF. $F_{X|Y=y}(x) = q((-\infty, x], y) = \Pr[X \leq x | Y=y]$

Thm.

- a) If (X, Y) has density $f(x, y)$, $X|Y=y$ has a density $f_{X|Y=y}(x) = \frac{f(x, y)}{f_Y(y)}$.
- b) If (X, Y) discrete w/ pmf p , $X|Y=y$ is discrete w/ pmf $\frac{\Pr[X \in A, Y=y]}{\Pr[Y=y]}$

Recap: $E[X|Y=y]$ is a real #, $E[X|Y]$ is a RV.

Thm. 1. If $X \perp Y$, $E[X|Y] = E[X]$ w/ $\Pr=1$

Thm. 2: LoTE/Tower Law. $E[E[X|Y]] = E[X]$

Thm. 3: Law of Total Variation. $\text{Var}(Y) = E[\text{Var}(Y|X)] + \text{Var}(E[Y|X])$

↳ $\text{Var}(Y|X) = E[Y^2|X] - E[Y|X]^2$

E.g. (Mossel's Dice)

Roll a fair die until you get a 6.

↳ Conditional on the event that all rolls are even, what is $E[\# \text{ rolls}]$?

$X = \# \text{ rolls until } 6; X \sim \text{Geom}(1/6)$

$$\rightarrow P_r[X=k] = (1-p)^{k-1} p, E[X] = \frac{1}{p}, \text{Var}(X) = \frac{1-p}{p^2}$$

Wrong answer: 3 (as if the die only had 3 faces)

Correct Answer

$N = \# \text{ times you roll } 2 \text{ or } 4 \text{ until you get something else } (1, 3, 5, 6)$

↳ $P = \text{value at the last roll}$

↳ $N \perp P$

$$\Rightarrow E[N] = E[N|P] = E[N|P=6]$$

$$N \sim \text{Geom}(2/3) \Rightarrow E[N] = \frac{3}{2} = E[N|P=6]$$

↳ soln is $\frac{3}{2}$!

↳ why smaller than 3? thought experiment w/ 1M sided die, still only 2 or 4 before 6

↳ most likely to hit 6 on first roll

New Prob: $Y = \# \text{ odd } \#s \text{ observed}$

$$Y|X=k \sim \text{Bin}(k-1, \frac{3}{5})$$

$$E[Y|X=k] = (k-1)\left(\frac{3}{5}\right), k=X, E[Y|X] = \frac{3}{5}(X-1)$$

$$E[Y] = E[E[Y|X]] = E\left[\frac{3}{5}(X-1)\right] = \frac{3}{5}(6-1) = 3$$

Now: $\text{Var}(Y)$

$$\hookrightarrow = E[\text{Var}(Y|X)] + \text{Var}(E[Y|X])$$

$$\hookrightarrow \text{Var}(Y|X) = ? \rightarrow \text{use binomial} \rightarrow = (X-1)\left(\frac{2}{5}\right)\left(\frac{2}{5}\right)$$

$$\hookrightarrow \text{Var}(E[Y|X]) = \text{Var}\left(\frac{2}{5}(X-1)\right) = \frac{4}{25}\text{Var}(X)$$

Bayesian inference:

$$\text{posterior} = \frac{\text{prior} \times \text{likelihood}}{\text{constant}}$$

don't always know

recall that $F^{-1}(U) \sim F$
but can't always do

Rejection Sampling

Goal: generate samples from dist. w/ density f_X

Assume:

1) $f(x) = \frac{g(x)}{N}$ \rightarrow known
 \rightarrow normalizing constant

2) I can sample from another dist. $h(x)$

3) $h(x) \geq c \cdot g(x)$, $c > 0$ (same support)

Algorithm

1) Sample Y from h

2) Draw $U \sim U[0, 1]$, $U \perp Y$

3) Accept Y and set $X = Y$ if $U \leq \frac{c \cdot g(Y)}{h(Y)}$

Justification

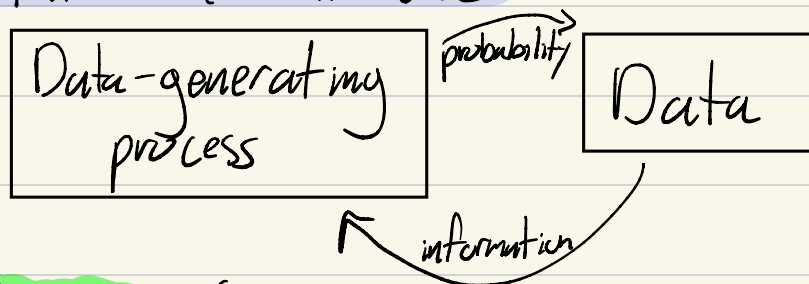
$$\Pr[X \in X \in x + \varepsilon] = \Pr[X \in Y \in x + \varepsilon \mid U \leq \frac{c \cdot g(Y)}{h(Y)}] = \frac{\Pr[X \in Y \in x + \varepsilon, U \leq \frac{c \cdot g(Y)}{h(Y)}]}{\Pr[U \leq \frac{c \cdot g(Y)}{h(Y)}]}$$

$$\begin{aligned} \hookrightarrow \text{Numerator} &\approx \Pr[X \in Y \in x + \varepsilon, U \leq \frac{c \cdot g(Y)}{h(Y)}] \\ &= \Pr[X \in Y \in x + \varepsilon] \Pr[U \leq \frac{c \cdot g(Y)}{h(Y)}] \\ &= \varepsilon \cdot h(x) \cdot \frac{c \cdot g(x)}{h(x)} = \varepsilon \cdot c \cdot g(x) \end{aligned}$$

$$\hookrightarrow \text{Denominator: } E[\mathbb{1}_{\{U \leq \frac{c \cdot g(Y)}{h(Y)}\}}] = E[E[\mathbb{1} | Y]] = E\left[\frac{c \cdot g(Y)}{h(Y)}\right] = \int \frac{c \cdot g(y)}{h(y)} h(y) dy = c \cdot N$$

Lec 3B - 10/16

Statistical Inference



Example (the big one)

A firm manufactures batteries & wants to learn their performance from data: $t_1, \dots, t_n =$ battery durations.

↳ Why do we care?

↳ eg. T RV models durations; $\Pr[T \geq 10]$?

Model

↳ We could approx. F_T w/ its empirical dist., as in HW 1.

T_1, \dots, T_n iid F_T . Goal: find F_T . ↳ not good for tail probabilities

Dists. Modeling Duration/Lifetime

$T: \Omega \rightarrow [0, \infty) =$ lifetime/duration of product

Def: Survival Fn. $S: [0, \infty) \rightarrow [0, 1]$, $S(x) = \Pr[T > x] = 1 - F_T(x)$

$$\Pr[T \leq x + \Delta x | T > x] = \frac{\Pr[x < T \leq x + \Delta x]}{\Pr[T > x]} \approx \frac{f(x) \cdot \Delta x}{S(x)}$$

Def: Failure Rate Fn. $h(x) = \frac{f(x)}{S(x)}$, $f =$ density

↳ h can be constant, increasing, decreasing

↳ no aging

↳ aging

↳ reverse aging

Q: If h known, what are f, F_T ?

$$\hookrightarrow h(x) = \frac{f(x)}{1-F(x)} = -\frac{(1-F)'(x)}{1-F(x)} = -\frac{d}{dx} \log(1-F(x))$$

$$\Rightarrow \log(1-F(x)) = -\int_0^x h(t) dt + C$$

$$\hookrightarrow F(0) = 0 \Rightarrow C = 0$$

$$\Rightarrow S(x) = \exp\left(-\int_0^x h(t) dt\right)$$

$$\Rightarrow f(x) = h(x) \exp\left(-\int_0^x h(t) dt\right).$$

E.g.

$$1) h(x) = \lambda x \quad \forall x \rightarrow x > 0$$

$$\Rightarrow f(x) = \lambda \cdot e^{-\lambda x} \Rightarrow T \sim \text{Exp}(\lambda)$$

\hookrightarrow Exp. dist. is memoryless: $\Pr[X > t+s | X > s] = \Pr[X > t]$

$$2) h(x) = \alpha + \beta x, \quad \alpha, \beta > 0$$

"Gompertz Family"

$$3) h(x) = \lambda \beta x^{\beta-1}$$

"Weibull"

Back 2 Example

Assume $T_1, \dots, T_n \stackrel{iid}{\sim} \text{Exp}(\lambda)$

$\hookrightarrow \lambda$ is our parameter of interest

Questions:

1) How to estimate λ from the data?

2) How good is the estimator?

↳ Metric?

3) UQ?

What do we know?

1) $T \sim \text{Exp}(\lambda)$

$$E[T] = \int_0^{\infty} x \cdot \lambda e^{-\lambda x} dx \rightarrow \lambda x = y, dx = \frac{1}{\lambda} dy$$

$$= \int_0^{\infty} y \cdot e^{-y} \frac{1}{\lambda} dy = \frac{1}{\lambda} \Gamma(2) = \frac{1}{\lambda}$$

$$E[T^2] = \int_0^{\infty} x^2 \lambda e^{-\lambda x} dx \rightarrow \lambda x = y$$

$$= \int_0^{\infty} \frac{1}{\lambda} y^2 \cdot e^{-y} \frac{1}{\lambda} dy = \frac{1}{\lambda^2} \Gamma(3) = \frac{2}{\lambda^2}$$

$$\text{Var}[T] = \frac{1}{\lambda^2}$$

2) $\bar{T}_n = \frac{T_1 + \dots + T_n}{n}$

$$E[\bar{T}_n] = \frac{1}{\lambda}, \text{Var}(\bar{T}_n) = \frac{1}{n\lambda^2}$$

$$\hookrightarrow \text{LLN: } \bar{T}_n \xrightarrow{P} \frac{1}{\lambda}$$

$$\text{CLT: } \sqrt{n}(\bar{T}_n - \frac{1}{\lambda}) \xrightarrow{D} \mathcal{N}(0, \frac{1}{\lambda^2})$$

Estimator for λ

$$E[\bar{T}_n] = \frac{1}{\lambda} \quad \begin{array}{l} \text{1st sample} \\ \text{moment} \end{array}$$

↳ LLN: $\bar{T}_n \sim \frac{1}{\lambda}$ for large n

$$\Rightarrow \hat{\lambda}_1 = \frac{1}{\bar{T}_n} = \text{Method of Moments}$$

Note: what ab 2nd moment?

$$\hookrightarrow E[T^2] = \frac{2}{\lambda^2}$$

$$\hookrightarrow \frac{\sum T_i^2}{n} \rightarrow \frac{2}{\lambda^2}, \text{ large } n$$

$$\Rightarrow \hat{\lambda}_2 = \frac{\sqrt{2n}}{\sqrt{\sum T_i^2}}$$

- now we have several estimators - which to use?

Properties of $\hat{\lambda}_1$

1) Consistency.

Thm (Continuous Mapping)

\mathcal{D}

$$a) Z_n \xrightarrow{p} b, g \text{ continuous} \Rightarrow g(Z_n) \xrightarrow{p} g(b)$$

$$b) Z_n \xrightarrow{D} Z, g \text{ continuous} \Rightarrow g(Z_n) \xrightarrow{D} g(Z)$$

$$\hookrightarrow \text{Thus, } T_n \xrightarrow{p} \frac{1}{\lambda} \Rightarrow \frac{1}{T_n} \xrightarrow{p} \lambda$$

Lemma. If $X_n \xrightarrow{D} X, Y_n \xrightarrow{D} C \Rightarrow X_n + Y_n \xrightarrow{D} X + C, X_n Y_n \xrightarrow{D} C X$

Thm (Delta Method). X_n s.t. $\sqrt{n}(X_n - \theta) \xrightarrow{D} \mathcal{N}(0, \sigma^2), g$ cont. differentiable
 $\Rightarrow \sqrt{n}(g(X_n) - g(\theta)) \xrightarrow{D} \mathcal{N}(0, \sigma^2 (g'(\theta))^2)$

Pf.

$$\text{Taylor: } g(X_n) = g(\theta) + g'(\tilde{\theta}_n)(X_n - \theta), \tilde{\theta}_n \in [\theta, X_n]$$

$$\sqrt{n}[g(X_n) - g(\theta)] = \underbrace{g'(\tilde{\theta}_n)}_{\xrightarrow{p} g'(\theta)} \underbrace{\sqrt{n}(X_n - \theta)}_{\xrightarrow{D} \mathcal{N}(0, \sigma^2)}$$

\Rightarrow the product converges

For us, $g(x) = \frac{1}{x}$, $g'(x) = \frac{1}{x^2}$, $\sqrt{n}(\bar{T}_n - \frac{1}{\lambda}) \xrightarrow{d} \mathcal{N}(0, \frac{1}{\lambda^2})$

↳ Delta Method: $\sqrt{n}(\frac{1}{\bar{T}_n} - \frac{1}{\lambda}) \xrightarrow{d} \mathcal{N}(0, \lambda^2)$

↳ λ_1 is asymptotically normal

For large n , $\hat{\lambda}_1 \approx \mathcal{N}(\lambda, \frac{\lambda^2}{n})$

↳ $\frac{\hat{\lambda}_1 - \lambda}{\lambda/\sqrt{n}} \approx \mathcal{N}(0, 1)$

$\Rightarrow \Pr[-1.96 \leq \frac{\hat{\lambda}_1 - \lambda}{\lambda/\sqrt{n}} \leq 1.96] \approx 0.95 \Leftrightarrow \Pr[-1.96 + \sqrt{n} \leq \frac{\hat{\lambda}_1}{\lambda} \leq 1.96 + \sqrt{n}] \approx 0.95$

$\Pr[\frac{\hat{\lambda}_1}{1.96 + \sqrt{n}} \leq \lambda \leq \frac{\hat{\lambda}_1}{-1.96 + \sqrt{n}}] \approx 0.95 \rightarrow 95\% \text{ CI}$

How about small n ?

Def. $X \sim \text{Gamma}(\alpha, \beta)$ if $f_X(x) = \frac{\beta^\alpha}{\Gamma(\alpha)} x^{\alpha-1} e^{-\beta x}$, $x > 0$

Properties

1) $E[X] = \frac{\alpha}{\beta}$, $\text{Var}[X] = \frac{\alpha}{\beta^2}$

2) $X_1, \dots, X_n \stackrel{iid}{\sim} \text{Gamma}(\alpha_i, \beta) \Rightarrow \sum X_i \sim \text{Gamma}(\sum \alpha_i, \beta)$

3) $E[\frac{1}{X}] = \frac{\beta}{\alpha-1}$, $\alpha > 1$

$$\sqrt{n}(\hat{\lambda}_n - \lambda) \xrightarrow{D} \mathcal{N}(0, \lambda^2)$$

(asymptotic normality)

Lec 4A- 10/21

Midterm

- cheat sheet (2-sided)

↳ can print, handwritten or typed

+ Topics

↳ Ch. 1- Ch. 3 (everything thru estimation)

- no coding problems, no calculator needed

- practice problems provided

↳ may be different (substantially) from HW, but not harder

↳ e.g. 2 probs similar to HW, 1 novel

↳ mix of computation & proof, but some proof

↗ or harder

↳ no explicitly tricky integrals, but e.g. $\int_0^{\infty} e^{-x}$ is expected

- problems: e.g. T_1, \dots, T_n

a) find MoM estimator for λ , b) use LoLN to show convergence

Estimator Selection (Metrics)

Say $X_1, \dots, X_n \stackrel{iid}{\sim} F_{\theta}$, w/ $\hat{\theta}$ an estimator of θ .

Def: Bias. = $E[\hat{\theta}] - \theta$.

↳ unbiased when bias = 0.

Def: MSE. $MSE(\hat{\theta}) = E[(\hat{\theta} - \theta)^2] = \text{bias}^2 + \text{variance}$

Small Sample Properties of $\hat{\lambda}_1$

Note: $\text{Exp}(\lambda) = \text{Gamma}(1, \lambda) \Rightarrow \sum_{i=1}^n T_i \sim \text{Gamma}(n, \lambda)$
 $\Rightarrow E[\hat{\lambda}_1] = n E[\sum T_i] = \frac{n}{n-1} \cdot \lambda$

$$E[\hat{\lambda}_1^2] = \frac{n^2}{(n-1)(n-2)} \lambda^2$$

$$\hookrightarrow \Rightarrow \text{Var}[\hat{\lambda}_1] = \frac{n^2}{(n-1)^2(n-2)} \cdot \lambda^2$$

- Twist: can I construct an unbiased estimator?

$$\hookrightarrow \tilde{\lambda} = \frac{n-1}{n} \cdot \hat{\lambda}_1 \Rightarrow E[\tilde{\lambda}] = \lambda$$

\hookrightarrow any tradeoff?

$$\hookrightarrow \text{Var}[\tilde{\lambda}] = \left(\frac{n-1}{n}\right)^2 \text{Var}(\hat{\lambda}_1)$$

\hookrightarrow so $\tilde{\lambda}$ has smaller bias & variance

\rightarrow doesn't always happen!

Recap

1) X_1, \dots, X_n iid

2) $\hat{\lambda} = f(X_1, \dots, X_n)$

3) Evaluate bias, var.

4) Look @ dist of $\hat{\lambda}$ (asymptotic, small sample)

5) CIs

Mom

$\hookrightarrow X_1, \dots, X_n$ iid F

$\hookrightarrow X \sim F; E X^k = k\text{th moment}$

$\hookrightarrow \frac{\sum X_i^k}{n} = k$ th sample moment

- to estimate k parameters, equate first k moments w/ first k sample moments

E.g.

$X_1, \dots, X_n \stackrel{iid}{\sim} U[0, \theta]$
 $E[X] = \frac{\theta}{2} \Rightarrow \hat{\theta} = 2\bar{X}$

E.g.

$X_1, \dots, X_n \stackrel{iid}{\sim} \text{Gamma}(\alpha, \beta)$
 $\hookrightarrow \frac{\alpha}{\beta} = \bar{X}$
 \hookrightarrow 2nd moment: $\frac{\alpha}{\beta^2} + \frac{\alpha^2}{\beta^2} = \frac{\sum X_i^2}{n}$

Note: for consistency, we use LLN and continuous mapping thm.
 \hookrightarrow For the asymptotic dist., we use CLT and delta method
 \hookrightarrow MoM only gives asymptotic guarantees

MLE

$X_1, \dots, X_n \stackrel{iid}{\sim} F_\theta, \theta \in \mathbb{R}^k$

\hookrightarrow let $f(x, \theta)$ be the density/pmf (or $f(x|\theta), f_\theta(x)$)

Goal: estimate θ

Def: Likelihood. The likelihood of θ given data X_1, \dots, X_n is the fn. $L(\theta) = L_n(\theta) = \prod_{i=1}^n f(x_i, \theta)$.

Def: MLE. $\hat{\theta} = \hat{\theta}_n = \operatorname{argmax}_{\theta} L_n(\theta)$ \rightarrow pick the θ that makes the data most likely

Def: Log-likelihood. $l(\theta) = \ln(\theta) = \log(L_n(\theta))$

E.g.

1) $X_1, \dots, X_n \stackrel{iid}{\sim} \text{Exp}(\lambda)$

$$f(x, \theta) = \lambda \exp(-\lambda x), \quad x > 0$$

$$L_n(\theta) = \lambda^n \exp(-\lambda \sum_{i=1}^n x_i), \quad x_i > 0$$

$$\ln(\theta) = n \ln(\lambda) - \lambda \sum_{i=1}^n x_i$$

$$l'_n(\theta) = \frac{n}{\lambda} - \sum_{i=1}^n x_i$$

$$\rightarrow l''(\lambda) = -\frac{n}{\lambda^2} < 0 \Rightarrow \text{global max}$$

$$\hookrightarrow l'_n(\theta) = 0 \Rightarrow \lambda = \frac{1}{\bar{x}_n}$$

2) $X_1, \dots, X_n \stackrel{iid}{\sim} U[0, \theta]$

$$f(x, \theta) = \frac{1}{\theta} : 0 < x < \theta$$

$$L_n(\theta) = \begin{cases} \frac{1}{\theta^n} & 0 < x_i < \theta \quad \forall i \in [n] \\ 0 & \text{o.w.} \end{cases}$$

$$\hookrightarrow 0 < x_i < \theta \quad \forall i \Leftrightarrow 0 < \max_i(x_i) < \theta$$

B/c $\frac{1}{\theta^n}$ decreasing, $\hat{\theta} = \max_i(x_i) = \text{leftmost val in range}$

3) (2 params) $X_1, \dots, X_n \stackrel{iid}{\sim} \mathcal{N}(\mu, \sigma^2) \Rightarrow \theta = (\mu, \sigma^2)$

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

$$\ln(\mu, \sigma^2) = \frac{1}{(2\pi)^{n/2}} \cdot \frac{1}{\sigma^2} \cdot \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2\right)$$

$$\hookrightarrow \ln(\mu, \sigma^2) = -\frac{n}{2} \ln(2\pi) - \frac{n}{2} \ln(\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2$$

$$\Rightarrow \hat{\mu} = \bar{X}$$

$$\hat{\sigma}^2 = \frac{1}{n} \sum (x_i - \bar{X})^2$$

Lec 4B - 10/23

Likelihood Estimation

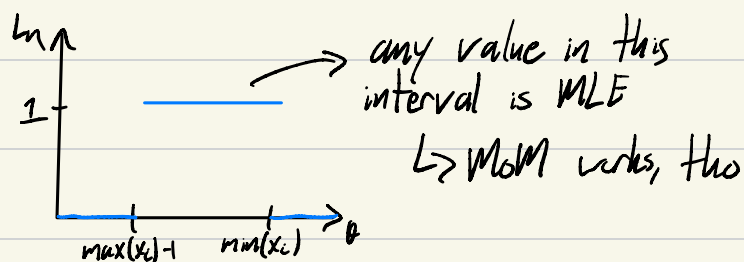
E.g.s

4) Non-Uniqueness

$$X_1, \dots, X_n \stackrel{iid}{\sim} \mathcal{U}[\theta, \theta+1]$$

$$f(x|\theta) = \begin{cases} 1 & \theta < x < \theta+1 \\ 0 & \text{o.w.} \end{cases}$$

$$L_n(\theta) = \begin{cases} 1 & \theta < x_i < \theta+1 \quad \forall i \in [n] \\ 0 & \text{o.w.} \end{cases} \Leftrightarrow \begin{cases} \theta < \min(x_i) \leq \max(x_i) < \theta+1 \\ \max(\theta)-1 < \theta < \min(x_i) \end{cases}$$



5) MLE DNE

$$X \text{ s.t. } X \sim \mathcal{N}(0,1) \text{ w/ } P_r = 1/2, \quad X \sim \mathcal{N}(\mu, \sigma^2) \text{ w/ } P_r = 1/2$$

$$f_X(x|\mu, \sigma^2) = \frac{1}{2} \left[\frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} \right] + \frac{1}{2} \left[\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \right]$$

$$X_1, \dots, X_n \stackrel{iid}{\sim} f_X$$

$$L_n(\mu, \sigma^2) = \prod_{i=1}^n f(x_i|\mu, \sigma^2)$$

\hookrightarrow Lets look at $\mu = X_1$

$$\hookrightarrow L_n(X_1, \sigma^2) \geq \frac{1}{2^n} \frac{1}{\sqrt{2\pi\sigma^2}} \cdot \left(\frac{1}{\sqrt{2\pi}} \right)^{n-1} \cdot e^{-\frac{\sum x_i^2}{2}}$$

$\hookrightarrow x = x_i$ in $\mathcal{N}(\mu, \sigma^2)$ \hookrightarrow take $\mathcal{N}(0,1)$ as lower bound instead

$$\sigma \rightarrow 0 \Rightarrow L_n(X_1, \sigma^2) \rightarrow \infty$$

\hookrightarrow but $\sigma = 0$ not permissible; no MLE

Why MLE?

+ under mild assumptions

- consistency
- asymptotic normality
- asymptotic variance known
- efficiency (optimal asymptotic variance)
- invariance

Theory (not for exams) (θ 1-dim)

Def: Score Fn. $\dot{\ell}_n(\theta | X_n) = \frac{\partial}{\partial \theta} \ell_n(\theta)$

Note: $\ell(\theta) = \sum \log f(x | \theta)$

$$\dot{\ell}(\theta) = \sum_{i=1}^n \frac{f'}{f}(x_i | \theta) \quad f' = \frac{\partial}{\partial \theta} f$$

Score equation: $\dot{\ell}_n(\theta | x) = 0$ gives MLE

Result 1. $E_{\theta_0} \left[\frac{\dot{\ell}_n(\theta)}{\dot{\ell}_n(\theta_0)} \right] = 1$

Pf.

$$E_{\theta_0} \left[\frac{\dot{\ell}_n(\theta)}{\dot{\ell}_n(\theta_0)} \right] = \int \frac{\dot{\ell}_n(\theta)}{\dot{\ell}_n(\theta_0)} f_n(x, \theta_0) dx = 1$$

↳ assumption (Common Support): $\{x : f(x, \theta) > 0\}$ does not depend on θ

↳ e.g. this is why the $U[\theta, \theta+1]$ didn't work!

Result 1 is also $E_{\theta_0} [e^{\ell_n(\theta) - \ell_n(\theta_0)}] = 1$

Result 2. a) $E_{\theta_0} [\dot{\ell}_n(\theta_0)] = 0$. b) $-E_{\theta_0} [\ddot{\ell}(\theta_0)] = E[\ddot{\ell}(\theta_0)]^2$ ^{2nd derivative}

Pf.

$$E_{\theta_0} [e^{\ell_n(\theta) - \ell_n(\theta_0)}] = 1$$

$$\hookrightarrow E_{\theta_0} \left[\frac{\partial}{\partial \theta} e^{\ell_n(\theta) - \ell_n(\theta_0)} \right] = 0$$

$$\Rightarrow E_{\theta_0} [\dot{\ell}_n(\theta) \cdot e^{\ell_n(\theta) - \ell_n(\theta_0)}] = 0$$

$$\hookrightarrow \theta = \theta_0 \Rightarrow @$$

For b) take $\frac{\partial}{\partial \theta} E[\dot{\ell}_n(\theta)]$

Assumption 2: Smoothness of densities. $\frac{\partial}{\partial \theta} E_{\theta_0} [g(x, \theta)] = E_{\theta_0} \left[\frac{\partial}{\partial \theta} g(x, \theta) \right]$

Def: Fisher Information. $I(\theta) = E_{\theta} [\dot{\ell}^2(\theta)] = E_{\theta} [-\ddot{\ell}(\theta)]$ ($n=1$)

Notes:

$$1) E_{\theta} [\ddot{\ell}(\theta)] = 0 \Rightarrow I(\theta) = \text{Var}(\dot{\ell}(\theta))$$

$$2) \text{Var}(\dot{\ell}_n(\theta)) = n \cdot I(\theta)$$

Thm: Cramer-Rao Inequality. $T(X_n)$ is unbiased for $g(\theta)$. Thus,

$$\text{Var}(T(X_n)) \geq \frac{(g'(\theta))^2}{n \cdot I(\theta)}$$

Note: if $g(\theta) = \theta$, bound is $\frac{1}{n \cdot I(\theta)}$. \rightarrow high information = low variance

Pf.

$$\text{Correlation}^2 \leq 1 \Rightarrow \frac{\text{Cov}^2}{\text{Var Var}} \leq 1$$

$$\text{Cov}_\theta^2(T(X_n), \ell_n) \leq \underbrace{\text{Var}(T(X_n))}_{\substack{\text{what we} \\ \text{want}}} \underbrace{\text{Var}(\ell_n)}_{= nI(\theta)}$$

$$\downarrow$$

$$\text{Cov}_\theta(T(X_n), \ell_n) = E_\theta[T(X_n) \cdot \ell] \rightarrow \text{b/c } E[\ell_n] = 0$$

$$= \int T(X_n) \ell_n \cdot f_n(\theta, x) dx$$

$$\hookrightarrow \text{Claim: } = \frac{\partial}{\partial \theta} f_n(\theta, x)$$

$$= \int T(X_n) \frac{\partial}{\partial \theta} f_n(\theta, x) dx$$

$$\stackrel{?}{=} \frac{\partial}{\partial \theta} \int T(X_n) f_n(\theta, x) dx = \frac{\partial}{\partial \theta} g(\theta)$$

↑ unbiased expectation

MLE Theory

Let $\hat{\theta}_n$ s.t. $0 = \ell_n(\hat{\theta}_n) \approx \ell(\theta_0) + (\hat{\theta}_n - \theta_0) \dot{\ell}(\theta_0) + \frac{1}{2}(\hat{\theta}_n - \theta_0)^2 \ddot{\ell}_n(\theta_0)$

$$\Rightarrow \sqrt{n}(\hat{\theta}_n - \theta_0) = \frac{1}{\sqrt{n}} \ell_n(\theta_0) - \frac{1}{n} \dot{\ell}_n(\theta_0) - \frac{1}{2n} (\hat{\theta}_n - \theta_0) \ddot{\ell}(\theta_0)$$

Notes:

a) $\frac{1}{\sqrt{n}} \ell_n(\theta_0) \xrightarrow{D} \mathcal{N}(0, I(\theta_0))$

b) $-\frac{1}{n} \dot{\ell}_n(\theta_0) \xrightarrow{P} I(\theta_0)$

c) Last term $\rightarrow 0$ as $n \rightarrow \infty$

\hookrightarrow need consistency

$$\Rightarrow \sqrt{n}(\hat{\theta}_n - \theta_0) \xrightarrow{D} \mathcal{N}(0, \frac{1}{I(\theta)})$$

\hookrightarrow gives asymptotic normality, asymptotic variance, asymptotic optimality

Thm. $\Pr_{\theta_0}[\ell_n(\theta_0) > \ell_n(\theta)] \xrightarrow{n \rightarrow \infty} 1$ (proof of consistency)

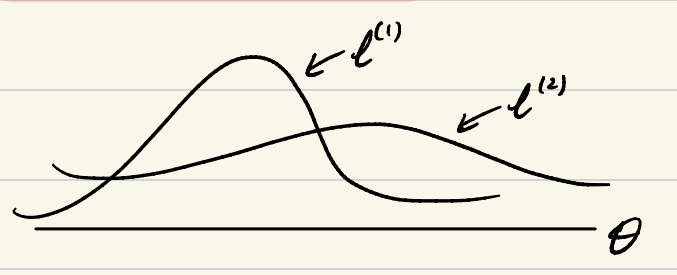
\hookrightarrow for any fixed θ

- proof requires Jensen

Assumption 3: $\theta_1 \neq \theta_2 \Rightarrow f(x, \theta_1) \neq f(x, \theta_2)$; distinct densities

MLE always consistent,
not necessarily unbiased

Lec 5A - 10/28



Which is better to learn about θ ?
 $\hookrightarrow l^{(1)}$ = peaked likelihood = good \Rightarrow large derivative
 $\hookrightarrow l^{(2)}$ = flat likelihood = bad
 $E[l^2]$ large = $I(\theta)$

$$0 = \dot{l}(\theta_0) \approx \dot{l}_n(\theta_0) + (\hat{\theta} - \theta_0) \ddot{l}_n(\theta_0)$$

$$\Rightarrow \sqrt{n} (\hat{\theta} - \theta_0) = \frac{\frac{1}{\sqrt{n}} \dot{l}_n(\theta_0)}{-\frac{1}{n} \ddot{l}_n(\theta_0)}$$

$$\Rightarrow \sqrt{n} (\hat{\theta} - \theta_0) \xrightarrow{D} \mathcal{N}(0, \frac{1}{I(\theta)})$$

\rightarrow gives asymptotic normality, optimality

E.g. $X_1, \dots, X_n \stackrel{iid}{\sim} \text{Bernoulli}(\theta)$

$$\Rightarrow f(x|\theta) = \theta^x (1-\theta)^{1-x}, \quad x = 0, 1$$

Then $l(\theta) = x \log \theta + (1-x) \log(1-\theta)$

$$\dot{l}(\theta) = \frac{x}{\theta} - \frac{1-x}{1-\theta}, \quad \ddot{l}(\theta) = -\frac{x}{\theta^2} - \frac{1-x}{(1-\theta)^2}$$

Score equation: $0 = \dot{l}_n(\theta) = \frac{\sum x_i}{\theta} - \frac{n - \sum x_i}{1-\theta}$

$$\Rightarrow \hat{\theta} = \bar{X}_n = \text{sample proportion}$$

$$I(\theta) = E[-\ddot{l}(\theta)] = E\left[\frac{x}{\theta^2} + \frac{1-x}{(1-\theta)^2}\right] = \frac{1}{\theta} + \frac{1}{1-\theta} = \frac{1}{\theta(1-\theta)}$$

We know $E[\hat{\theta}] = E[\bar{X}_n] = \theta$ (unbiased)

$$\hookrightarrow \text{Var}(\bar{X}_n) = \frac{1}{n} \text{Var}(X) = \frac{\theta(1-\theta)}{n} \quad (\text{smallest var from CR})$$

$\Rightarrow \hat{\theta} = \text{MVUE}$ (min. var. unbiased estimator)

But also, $\sqrt{n}(\hat{\theta} - \theta) \xrightarrow{D} \mathcal{N}(0, \theta(1-\theta))$

$$\Rightarrow \hat{\theta} \cong \mathcal{N}\left(\theta, \frac{\theta(1-\theta)}{n}\right) \text{ for large } n$$

$$\Rightarrow \Pr\left[-z_{\alpha/2} \leq \frac{\hat{\theta} - \theta}{\sqrt{\frac{\theta(1-\theta)}{n}}} \leq z_{\alpha/2}\right] \cong 1 - \alpha \text{ for } z_{\alpha/2} \text{ quantile of } \mathcal{N}(0, 1)$$

\hookrightarrow how do we isolate θ ? \rightarrow want $\Pr[A(x) < \theta < B(x)] = 1 - \alpha$

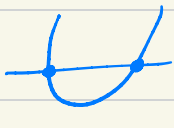
Two Confidence Intervals

$$1) \hat{\theta} \pm z_{\alpha/2} \sqrt{n \hat{\theta}(1-\hat{\theta})} \quad (\hat{\theta}(1-\hat{\theta}) \cong \theta(1-\theta))$$

↳ uses second approximation; the approx. gets worse

$$2) \Pr\left[\frac{(\hat{\theta}-\theta)^2}{\hat{\theta}(1-\hat{\theta})/n} < z_{\alpha/2}^2\right] = 1-\alpha$$

$$= \Pr\left[(\hat{\theta}-\theta)^2 < z_{\alpha/2}^2 \frac{\theta(1-\theta)}{n}\right]$$

↳ gives quadratic: 

Invariance of MLE

Thm. $\hat{\theta}_n$ is MLE of θ and $\tau = g(\theta)$, then $\hat{\tau} = g(\hat{\theta}_n)$ is the MLE of τ .

E.g. $X_1, \dots, X_n \stackrel{iid}{\sim} \mathcal{N}(\mu, 1)$.

↳ MLE of μ is \bar{X}_n .

$$\Rightarrow \text{MLE } \mu^2 = (\bar{X}_n)^2$$

$$\text{MLE } e^{\mu} = e^{\bar{X}_n}$$

Pf.

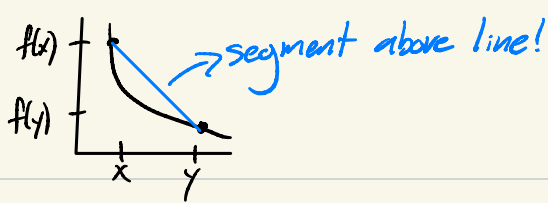
$$\hat{\theta}_n = \operatorname{argmax}_{\theta} L_n(\theta|x)$$

Case 1: g bijective.

$$\hat{\tau}_n = g(\hat{\theta}_n) \Rightarrow \hat{\theta}_n = g^{-1}(\hat{\tau}_n)$$

$$L(\tau) = L(\theta) \leq L(\hat{\theta}_n) = L(\hat{\tau}_n)$$

Don't care abt case 2.



Consistency

Def: Convex Fn. f is convex if $\forall x, y, 0 < \alpha < 1,$

$$f(\alpha x + (1-\alpha)y) \leq \alpha f(x) + (1-\alpha)f(y)$$

↳ if this is " $<$ ", we have strict convexity

Note

Z is a RV s.t. $\begin{cases} x & Pr = \alpha \\ y & Pr = 1-\alpha \end{cases} \Rightarrow f(E[Z]) \leq E[f(Z)] \rightarrow$ holds for all RVs!

Thm: Jensen's Inequality. X RV s.t. $E[|X|] < \infty, f$ convex, then $f(E[X]) \leq E[f(X)].$
we also assume differentiable

Pf.

Tangent has slope $f'(y)$ and is below $f.$

$$\Rightarrow f(x) \geq f(y) + f'(y)(x-y) \quad (\text{first-order Taylor})$$

Let $\mu = E[X].$ Then $f(x) \geq f(\mu) + f'(\mu)(x-\mu)$

$$E[f(x)] \geq f(\mu) = f(E[X])$$

Notes

a) If f strictly convex & $Pr[X=A] < 1 \forall A,$ we get a strict ineq.

b) f convex $\Rightarrow -f$ concave

↳ so if f concave, $f(EX) \geq E f(x)$

c) $f(x) = x^2$ convex $\Rightarrow (E[X])^2 \leq E[X^2]$

$g(x) = \log x$ concave $\Rightarrow \log E[X] \geq E[\log X]$

d) $x_1, \dots, x_n > 0, p_i \geq 0$ s.t. $\sum p_i = 1$

$$\Rightarrow \underbrace{\prod_{i=1}^n x_i^{p_i}}_{\text{geometric mean}} \leq \underbrace{\sum_{i=1}^n p_i x_i}_{\text{arithmetic mean}}$$

Pf: take log, X is RV s.t. $X = x_i$ w/ $Pr = p_i$. Then $E \log X \leq \log EX$

e) **Holder's Inequality.**

$0 \leq X, Y$ RVs, $\frac{1}{p} + \frac{1}{q} = 1$
 $\Rightarrow E[XY] \leq (E[X^p])^{1/p} (E[Y^q])^{1/q}$

Pf.

Case 1.

$$(E[X^p])^{1/p} = 1 = (E[Y^q])^{1/q}$$

$$XY = (X^p)^{1/p} (Y^q)^{1/q} \stackrel{d)}{\leq} \frac{1}{p} X^p + \frac{1}{q} Y^q$$

$$\Rightarrow E[XY] \leq \underbrace{\frac{1}{p} E[X^p]}_{\leq 1} + \underbrace{\frac{1}{q} E[Y^q]}_{\leq 1} = 1$$

Case 2.

Use case 1 for $X = \frac{X}{(E[X^p])^{1/p}}, Y = \frac{Y}{(E[Y^q])^{1/q}}$

f) **Likelihood Consistency.**

$$Pr_{\theta_0}[\ell_n(\theta_0) > \ell_n(\theta)] \xrightarrow{n \rightarrow \infty} 1$$

Pf.

$$\frac{1}{n} (\ell_n(\theta_0) - \ell_n(\theta)) = \frac{1}{n} \sum_{i=1}^n \log \frac{f(x_i | \theta_0)}{f(x_i | \theta)}$$

$\xrightarrow{Pr_{\theta_0}} \frac{1}{p} E_{\theta_0} \left[\log \frac{f(x | \theta_0)}{f(x | \theta)} \right]$

$< \log E_{\theta_0} \left[\frac{f(x | \theta_0)}{f(x | \theta_0)} \right] = \log(1) = 0$

Annotations:
 - just log properties (pointing to the log)
 - iid RVs (pointing to the sum)
 - whole quantity is their avg (pointing to the sum)
 - Jensen w/ strictly concave log (pointing to the inequality)
 - nontrivial b/c diff density assumption! (pointing to the expectation)
 - not always a single val (pointing to the expectation)
 - $\int \frac{f(x|\theta)}{f(x|\theta_0)} f(x|\theta_0) dx = 1$ (pointing to the expectation)

Lec 5B - 10/30

Multivariate Normal

Def. $X = (X_1, \dots, X_k)$ has a multivariate normal distribution if $\forall a \in \mathbb{R}^k$, $a^T X$ is normal (any LC is normal).

Notation: $\mu = E[X] = (E[X_1] \dots E[X_k])^T$
 $\Sigma = \text{Var}(X) = E[(X-\mu)(X-\mu)^T] \in \mathbb{R}^{k \times k}$

Notes

a) if Σ is positive-definite $\Rightarrow X$ has density

$$f_X(x) = \frac{1}{(2\pi)^{k/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu)\right) \rightarrow |\Sigma| = \det(\Sigma)$$

E.g. $Z \sim \mathcal{N}(0, 1)$, $(Z, 1-Z)$ has no density

\hookrightarrow support is on a line: $x+y=1$; $\text{Cov}\begin{pmatrix} Z \\ 1-Z \end{pmatrix} = \begin{pmatrix} 1 & -1 \\ -1 & 1 \end{pmatrix} = \text{rank } 1$

b) (X_1, X_2) bivariate normal & $\text{Cov}\begin{pmatrix} X_1 \\ X_2 \end{pmatrix} = 0 \Rightarrow X_1 \perp X_2$

c) $U \sim \mathcal{N}_k(\mu, \Sigma)$

\hookrightarrow let $V = a + BU$, $a \in \mathbb{R}^p$, $B \in \mathbb{R}^{p \times k}$

$\hookrightarrow V \sim \mathcal{N}_p(a + B\mu, B\Sigma B^T)$

Sample Mean & Variance

$X_1, \dots, X_n \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1)$

$$\hookrightarrow \bar{X}_n = \frac{X_1 + \dots + X_n}{n}$$

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

Notes

a) $E[\bar{X}_n] = \mu, \text{Var}(S^2) = \sigma^2$

b) $\bar{X} \perp S^2$

Pf. ($n=2$)

$$\bar{X} = \frac{x_1 + x_2}{2}, \quad S^2 = \frac{(x_1 - x_2)^2}{2}$$

If $x_1 + x_2 \perp x_1 - x_2$, we're done.

↳ note $\begin{pmatrix} x_1 + x_2 \\ x_1 - x_2 \end{pmatrix}$ is bivariate normal

$$\text{Cov}(x_1 + x_2, x_1 - x_2) = 0 \Rightarrow (x_1 + x_2) \perp (x_1 - x_2) \quad \square$$

c) $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} \mathcal{N}(\mu, \sigma^2)$

$$\Rightarrow V = (n-1) \frac{S^2}{\sigma^2} \sim \chi_{n-1}^2$$

$$T = \sqrt{n} \frac{\bar{X} - \mu}{S} \sim t_{n-1}$$

Confidence Intervals

$$X_1, \dots, X_n \text{ iid } F_\theta$$

CI's = probabilistic statements on data of the form

$$P_\theta [A(X_1, \dots, X_n) \leq \theta \leq B(X_1, \dots, X_n)] = \alpha$$

↳ Then $[A(X_1, \dots, X_n), B(X_1, \dots, X_n)]$ is a α -CI

↳ 0.95, 0.99, etc.

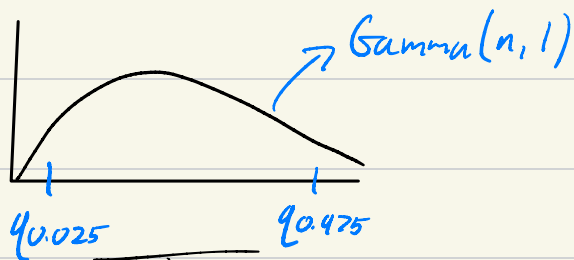
Notes

a) ends of the CI are RVs

b) Interpretation: long-term frequency of intervals that cover θ is α .

E.g. X_1, \dots, X_n iid $\text{Exp}(\lambda)$. Recall $X_i \sim \text{Gamma}(1, \lambda)$
 $\Rightarrow \sum X_i \sim \text{Gamma}(n, \lambda)$

Further, $\lambda \sum X_i \sim \text{Gamma}(n, 1)$



quantiles := $\Pr[Y \leq q_\alpha] = \alpha$

$$\Rightarrow \Pr[q_{0.025} \leq \lambda \sum X_i \leq q_{0.975}] = 0.95 \Rightarrow \Pr\left[\frac{q_{0.025}}{\sum X_i} \leq \lambda \leq \frac{q_{0.975}}{\sum X_i}\right] = 0.95$$

$\hookrightarrow \left[\frac{q_{0.025}}{\sum X_i}, \frac{q_{0.975}}{\sum X_i}\right]$ is a 95% CI for λ

+ the key: $\lambda \sum X_i$ has a dist. \perp of λ .

Def: Pivot Statistic. X_1, \dots, X_n iid F_θ . $g(X_1, \dots, X_n, \theta)$ is called pivot if its dist. does not depend on θ .

If $g(X_1, \dots, X_n, \theta) \sim F^*$, and L, U quantiles s.t.
 $\Pr[L < g(X_1, \dots, X_n, \theta) < U]$, we can solve for θ (may not give interval).

Shield AI

E.g.s

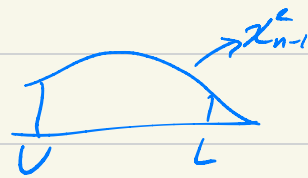
a) X_1, \dots, X_n iid $\mathcal{N}(\mu, \sigma^2)$

pivot: $g = T = \sqrt{n} \frac{\bar{x} - \mu}{s} \sim t_{n-1}$

gives $\bar{x} \pm t_{n-1}^{\alpha/2} \frac{s}{\sqrt{n}}$

b) Normal variance.

pivot $g: V = (n-1) \frac{s^2}{\sigma^2} \sim \chi_{n-1}^2 \rightarrow$



c) Large sample mean.

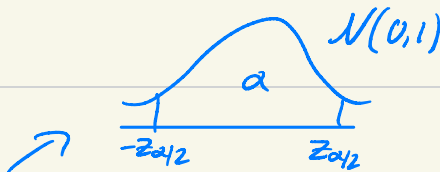
X_1, \dots, X_n iid F , $E[X_i] = \mu$, $\text{Var}(X_i) = \sigma^2$

$\Rightarrow \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} \xrightarrow{d} \mathcal{N}(0, 1)$ (CLT)

\Rightarrow for n large, $\frac{\bar{x} - \mu}{\sigma/\sqrt{n}} \approx \mathcal{N}(0, 1)$

If σ known

\Rightarrow a CI: $\left[\bar{x} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}}, \bar{x} + z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \right]$ (approximation)



σ unknown:

$\left[\bar{x} - z_{\alpha/2} \frac{s}{\sqrt{n}}, \bar{x} + z_{\alpha/2} \frac{s}{\sqrt{n}} \right] \rightarrow$ double approximation

d) MLE. X_1, \dots, X_n iid F_θ . $\hat{\theta}$ is MLE (valid assumptions)

$\Rightarrow \sqrt{n}(\hat{\theta} - \theta) \approx \mathcal{N}(0, \frac{1}{I(\theta)})$

$\approx \mathcal{N}(0, \frac{1}{I(\hat{\theta})}) \rightarrow$ recall Bernoulli last lec.

e) 2 samples means.

X_1, \dots, X_n iid $\mathcal{N}(\mu_1, \sigma_1^2) \perp Y_1, \dots, Y_m$ iid $\mathcal{N}(\mu_2, \sigma_2^2)$

\hookrightarrow want CI for $\mu_1 + \mu_2$

$\bar{X} - \bar{Y} \sim \mathcal{N}(\mu_1 - \mu_2, \frac{\sigma_1^2}{n} + \frac{\sigma_2^2}{m})$

$\Rightarrow (\bar{X} - \bar{Y}) - (\mu_1 - \mu_2) \sim \mathcal{N}(0, \frac{\sigma_1^2}{n} + \frac{\sigma_2^2}{m})$

★ Ch. 4 on
midterm (CIs)

What if we don't have the above scenarios?

↳ Case 1: pivot $g \sim F^*$, F^* unknown

↳ Case 2: T_n is unbiased for θ , but $\text{Var}(T_n)$ unknown

Bootstrap

Def: Empirical Dist. F_n . Given X_1, \dots, X_n iid F , the EDF \hat{F}_n is the cdf that puts mass $\frac{1}{n}$ at each X_i : $\hat{F}_n(x) = \frac{1}{n} \sum \mathbb{1}_{\{X_i \leq x\}}$

Notes

- \hat{F}_n for fixed x is sample proportion for $\text{Bern}(F(x))$
- $E[\hat{F}_n(x)] = F(x)$
- $\text{Var}(\hat{F}_n(x)) = \frac{1}{n} [F(x)(1-F(x))]$
- LLN: $\hat{F}_n(x) \xrightarrow{P} F(x)$
- Thm (Glivenko-Catelli): $\sup_x |\hat{F}_n(x) - F(x)| \rightarrow 0$ a.s.

Let $\theta = T(F)$

↳ mean $\int x dF(x)$

var $\int (x-\mu)^2 dF(x)$

median $F^{-1}(\frac{1}{2})$

Plug in estimator of θ : $\hat{\theta} = T(\hat{F}_n)$ (mean \rightarrow sample mean)

Lec 6A - 11/4

Midterm

- minimal proofs (similar to HW)
- difficulty on par w/ HW
 - ↳ but no copied HW problems

CI's

Pivot Stats: X_1, \dots, X_n iid F_θ

$g(X_1, \dots, X_n, \theta) \sim F^*$ is known, ind. of θ

↳ $\Pr[L \leq g(X_1, \dots, X_n, \theta) \leq R] = 0.95 \rightarrow$ solve for θ

Bootstrap

- relies on empirical distribution f_n $\hat{F}_n(x) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{X_i \leq x\}}$

$\theta = T(F)$ mean $\int x dF(x)$

$\hat{\theta} = T(\hat{F}_n)$ mean $\int x d\hat{F}_n(x) = \bar{X}_n$

E.g. Variance Estimation

X_1, \dots, X_n iid F , $T_n = g(X_1, \dots, X_n)$ a statistic, need $\text{Var}[T_n]$

Case 0: Can simulate from F

① Take k samples of size n from F .

② For each sample, calculate $g(\cdot)$

↳ $T_{n,1}, T_{n,2}, \dots, T_{n,k}$

③ Estimate variance w/ $\frac{1}{k} \sum_{j=1}^k (T_{n,j} - \bar{T})^2$

General Case

$$\text{Var}_F[T_n] \approx \text{Var}_{\hat{F}_n}[\bar{T}_n]$$

↳ simulate from \hat{F}_n (bootstrap samples \equiv sampling from X_1, \dots, X_n w/ replacement)

Bootstrap CIs

1. (The Normal Interval) X_1, \dots, X_n iid F_θ . Let $\hat{\theta} = T_n(X_1, \dots, X_n)$ be unbiased for θ . $\hat{\sigma}_{\text{boot}}$ = bootstrap estimate of sd of θ

The α CI is $[\hat{T}_n - z_{\alpha/2} \hat{\sigma}, \hat{T}_n + z_{\alpha/2} \hat{\sigma}]$

↳ need T_n to be close to normal

2. (Pivot Statistic). X_1, \dots, X_n iid F_θ , $g(X_1, \dots, X_n, \theta) \sim F^*$ unknown is our pivot.

We need L, U quantiles for dist. of g

Bootstrap Way:

① Gen bootstrap sample

② Calc. g w/ $\hat{\theta} = \hat{\theta}_{\text{boot}}$, $\theta = \hat{\theta}$

③ Repeat and get \hat{L}, \hat{U}

e.g. if $g = \hat{\theta} - \theta$, we'd
calc. $\hat{\theta}_{\text{boot}} - \hat{\theta}$

* bootstrap not on midterm

Hypothesis Testing

+ "Is this data consistent w/ a certain data-generating mechanism?"

E.g. X_1, \dots, X_n iid F ; is F the cdf of a normal (aka is the data normally distributed)?

Ideas:

t, cauchy
↑

1) Plot histogram. Flaw: subjective, other dist.s look similar

2) Quantiles (e.g. median, lower/upper quartile)

- need a comparison; $\mathcal{N}(0,1)$ quantiles are known

- standardize the data: $\frac{X_i - \bar{X}}{s_x}$

- QQ plot = standard normal vs. sample quantiles, if a line then sample is normal

- still visual!

3) Quantiles, pt. Z \rightarrow 3rd quartile

$$\hookrightarrow \text{take } \frac{(X_{(\frac{3}{4}n)} - X_{(\frac{1}{4}n)})^2}{\sum (X_i - \bar{X})^2}$$

Generalizes to Shapiro-Wilks: $\frac{(\sum a_i X_{(i)})^2}{\sum (X_i - \bar{X})^2}$

4) EDF

$\hookrightarrow \hat{F}_n$ of $\frac{X_i - \bar{X}}{s_x}$, $F = \text{cdf of } \mathcal{N}(0,1)$

\hookrightarrow Kolmogorov-Smirnov: $\sup_x |\hat{F}_n(x) - F(x)|$ (works for any dist.)

Parametric Models

X_1, \dots, X_n iid F_θ , $\theta \in \Omega$

+ Hypotheses:

$$\left. \begin{array}{l} \text{Null: } \theta \in \Omega_0 \subset \Omega \\ \text{Alternative: } \theta \in \Omega_1 \subset \Omega \end{array} \right\} \begin{array}{l} \Omega_0 \cap \Omega_1 = \emptyset \\ \Omega_0 \cup \Omega_1 = \Omega \end{array}$$

Notes:

a) H_0 = null hypothesis

H_a / H_1 = alternative

b) In our normal example, H_0 = data was normal, H_a = wasn't normal

c) Asymmetry btwn H_0 , H_a

↳ data is used to investigate if the data is consistent w/ the null, not alternative

Def. S is the set of all possible values for (X_1, \dots, X_n) . Also,

$$S = S_0 \cup S_1, \quad S_0 \cap S_1 = \emptyset$$

↳ S_0 = values for which we do not reject the null (acceptance region)

S_1 = vals for which we reject the null (rejection region)

Normally we define S_1 using a statistic, e.g., $\{x: \sup_x |F_n(x) - F(x)|\}$

$\rightarrow R_1 \subseteq \mathbb{R}$

Def: Test Statistic. T a statistic s.t. $S_1 = \{x : T(x) \in R_1\}$

E.g. X_1, \dots, X_n iid $\mathcal{N}(\mu, 1)$.

$H_0: \mu = \mu_0 \equiv \Omega_0 = \{\mu_0\}$

$H_a: \mu \neq \mu_0 \equiv \Omega_1 = \mathbb{R} \setminus \{\mu_0\}$

$T(x) = |\bar{X}_n - \mu_0|$

$S_1 = \{x : |\bar{X}_n - \mu_0| \geq c\} \rightarrow$ how to pick c ?

Error Types

	$X \in S_0$	$X \in S_1$
$\theta \in \Omega_0$ (H_0)	Correct	Type I Error
$\theta \in \Omega_1$ (H_a)	Type II Error	Correct

E.g. X_1, \dots, X_{10} iid $\text{Bern}(\theta)$

$H_0: \theta = 1/2$

$H_a: \theta > 1/2$

$T = \sum X_i; S_1 = \{\sum X_i \geq 7\}$

+ Type I Error

$$Pr_{\theta=1/2}[S_1] = Pr_{\theta=1/2}[\sum X_i \geq 7] = \sum_{k=7}^{10} \binom{10}{k} (1/2)^k (1/2)^{n-k} = 0.17$$

+ Type II Error

\hookrightarrow for which θ ? next time

always easier to calc
 $\max < x$ than $\max > x$!

Lec 6B - 11/6

Error Types

	S_0	S_1
H_0	✓	Type I
H_a	Type II	✓

Def: Power Function. X_1, \dots, X_n iid F_θ . Power fn. is $\pi: \Omega \rightarrow [0,1]$ s.t.
 $\pi(\theta) = \Pr_\theta[x \in S_1]$ ($\Pr_\theta \Rightarrow$ as if data came from F_θ).

\hookrightarrow Type I error governed by $\pi(\theta)$, $\theta \in \Omega_0$
Type II error " " $1 - \pi(\theta)$, $\theta \in \Omega_1$

E.g. X_1, \dots, X_n iid $U[0, \theta]$, $\theta > 0$ ($\Omega = \mathbb{R}_+$)

$$H_0: 3 \geq \theta \geq 4$$

$$H_a: \theta < 3 \text{ or } \theta > 4$$

MLE of θ is $X_{(n)} \in (0, \theta)$

$$\text{Rejection region: } S_1 = \{X_{(n)} \geq 4\} \cup \{X_{(n)} \leq 2.9\}$$

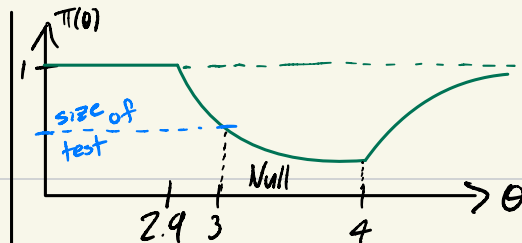
Power fn:

$$\pi(\theta) = \Pr_\theta(S_1) = \Pr_\theta(X_{(n)} \geq 4) + \Pr_\theta(X_{(n)} \leq 2.9)$$

$$\text{Case } \theta \leq 2.9: \pi(\theta) = 0 + 1$$

$$\text{Case } 2.9 < \theta < 4: \pi(\theta) = 0 + \Pr_\theta[X_{(n)} \leq 2.9] = \Pr_\theta[X_1 \leq 2.9, \dots, X_n \leq 2.9] \\ = \left(\frac{2.9}{\theta}\right)^n$$

$$\text{Case } \theta \geq 4: \pi(\theta) = [1 - \Pr_\theta(X_{(n)} < 4)] + \Pr_\theta[X_{(n)} \leq 2.9] \\ = 1 - \left(\frac{4}{\theta}\right)^n + \left(\frac{2.9}{\theta}\right)^n$$



→ So the errors aren't clean #s like in elementary statistics

Ideally, $\pi(\theta) = \begin{cases} 0 & \theta \in \Omega_0 \\ 1 & \theta \in \Omega_1 \end{cases}$ or $\pi(\theta) = \begin{cases} \text{low} & \theta \in \Omega_0 \\ \text{high} & \theta \in \Omega_1 \end{cases}$

Solution: make $\pi(\theta)$ low on $\Omega_0 \equiv$ find S_1 s.t.
 $\pi(\theta) \leq \alpha_0 \forall \theta \in \Omega_0$ (e.g., $\alpha_0 = 0.1$)

Def: Size of Test. $\sup_{\theta \in \Omega_0} \pi(\theta)$

Def. Test is a level α test if $\text{size} = \sup_{\theta \in \Omega_0} \pi(\theta) \leq \alpha$.

Note: size is easy to evaluate when $\Omega_0 = \{\theta_0\}$, $\text{size} = \text{Pr}_{\theta_0}[X \in S_1]$

If $S_1 = \{T = T(x_1, \dots, x_n) \geq c\}$, then a level- α test is s.t.
 $\sup_{\theta \in \Omega_0} \text{Pr}_{\theta}[T \geq c] \leq \alpha$.

E.g. X_1, \dots, X_n iid F_{μ} , $E[X_i] = \mu$, $\text{Var}(X_i) = \sigma^2$ is known.

$H_0: \mu = \mu_0$, $H_a: \mu > \mu_0$

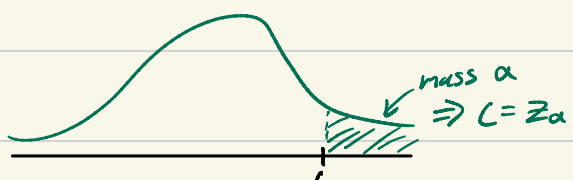
Test Statistic

$\hookrightarrow S_1 = \{\bar{x}_n - \mu_0 \geq c\} \rightarrow$ but can't calc probabilities

$S_1 = \{\sqrt{n} \frac{\bar{x}_n - \mu_0}{\sigma} \geq c\} \rightarrow$ approx $\mathcal{N}(0,1)$ by CLT

$\Rightarrow T = \sqrt{n} \frac{\bar{x}_n - \mu_0}{\sigma} \approx \mathcal{N}(0,1)$

$\text{Pr}_{\mu_0}[X \in S_1] = \text{Pr}_{\mu_0}[T \geq c] \rightarrow$



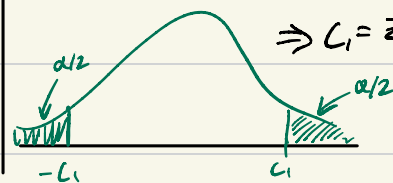
$$C = z_\alpha \rightarrow \Pr[\bar{Z} \sim \mathcal{N}(0,1) > z_\alpha] = \alpha$$

$$N_{\alpha,1}, S_1 = \left\{ x : \bar{X}_n \geq \mu_0 + z_\alpha \cdot \frac{\sigma}{\sqrt{n}} \right\}$$

E.g. Same setting.

$$H_0: \mu = \mu_0, H_a: \mu \neq \mu_0, T = \sqrt{n} \frac{\bar{X}_n - \mu_0}{\sigma}$$

$$S_1 = \{|T| \geq c_1\}$$



$$\Rightarrow c_1 = z_{\alpha/2}$$

$$S_1 = \{|T| \geq z_{\alpha/2}\}$$

But neither of these examples say anything a/b the type II error!

Power of the Above Test

$$\Pi(\mu) = \Pr_\mu[S_1] = \Pr_\mu\left[\left|\sqrt{n} \frac{\bar{X}_n - \mu_0}{\sigma}\right| \geq z_{\alpha/2}\right]$$

$$= \Pr_\mu\left[\left|\sqrt{n} \frac{\bar{X}_n - \mu}{\sigma} + \frac{\mu - \mu_0}{\sigma}\right| \geq z_{\alpha/2}\right] \rightarrow z = \sqrt{n} \frac{\bar{X}_n - \mu}{\sigma} \approx \mathcal{N}(0,1) \quad (\mu = \mu_0)$$

$$= \Pr_\mu\left[\left|Z + \sqrt{n} \frac{\mu - \mu_0}{\sigma}\right| \geq z_{\alpha/2}\right]$$

$$= \Pr_\mu\left[Z + \sqrt{n} \frac{\mu - \mu_0}{\sigma} \geq z_{\alpha/2}\right] + \Pr_\mu\left[Z + \sqrt{n} \frac{\mu - \mu_0}{\sigma} \leq -z_{\alpha/2}\right]$$

$$= \Pr_\mu\left[Z \geq z_{\alpha/2} - \sqrt{n} \frac{\mu - \mu_0}{\sigma}\right] + \Pr_\mu\left[Z \leq -z_{\alpha/2} - \sqrt{n} \frac{\mu - \mu_0}{\sigma}\right]$$

Power depends on

- sample size (n): $\Pi \xrightarrow{n \rightarrow \infty} 1$ more samples, \rightarrow more power
- signal ($\mu - \mu_0$): $\Pi \xrightarrow{(\mu - \mu_0) \rightarrow \infty} 1$ more signal \rightarrow more power

E.g. Same setting, σ unknown

$$T = \sqrt{n} \frac{\bar{X}_n - \mu_0}{s} \approx t_{n-1}$$

\hookrightarrow sample
std. dev.

E.g. Variance of Normal

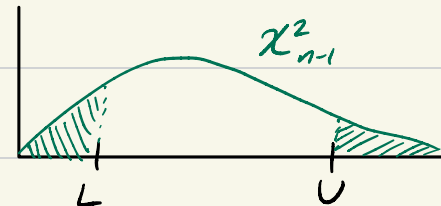
$$X_1, \dots, X_n \text{ iid } \mathcal{N}(\mu, \sigma^2) \quad \Omega_0 = \mathbb{R} \times \{\sigma_0\}$$

$$H_0: \sigma = \sigma_0, H_a: \sigma \neq \sigma_0 \quad \Omega_1 = \mathbb{R} \times \mathbb{R}_+ \setminus \{\sigma_0\}$$

$$T = (n-1) \frac{S^2}{\sigma_0^2} \stackrel{H_0}{\sim} \chi_{n-1}^2 \quad (\text{midterm mention?})$$

↳ null says $\approx (n-1)$, reject if smaller/greater

$$\hookrightarrow S_1 = \{T < L, T > U\}$$



Need L, U s.t. $\Pr_{\sigma_0}[S_1] = \alpha$

↳ L, U quantiles of χ_{n-1}^2

Power

$$\Pr_{\sigma_0}[S] = \Pr_{\sigma_0}[T < L, T > U]$$

$$\text{Recall } T = (n-1) \frac{S^2}{\sigma_0^2} = \underbrace{(n-1) \frac{S^2}{\sigma^2}}_{\chi_{n-1}^2} \cdot \frac{\sigma^2}{\sigma_0^2}$$

$$\Rightarrow \Pr_{\sigma_0} \left[\chi_{n-1}^2 < L \frac{\sigma_0^2}{\sigma^2}, \chi_{n-1}^2 > U \frac{\sigma_0^2}{\sigma^2} \right]$$

Lec 7B - 11/13

Recall hypothesis testing easy when $\pi(\theta)$ constant over Ω_0
↳ if $S_1 = \{T(x) \geq c\}$, $\Rightarrow \mathcal{L}(T(x))$ same for all $\theta \in \Omega_0$

P-Values

Def. p-value is the smallest level α for which we reject the null w/ the observed data.

↳ means we do not need to specify α

- eg. $p = 0.008$

↳ $\alpha = 0.01 > p \rightarrow$ reject at this level

↳ $\alpha = 0.005 < p \rightarrow$ do not reject

E.g. X_1, \dots, X_n iid F , $E[X] = \mu$, $\text{Var}[X] = \sigma^2$ known

$H_0: \mu = \mu_0$, $H_A: \mu \geq \mu_0$

$$T(x) = \sqrt{n} \frac{\bar{X}_n - \mu_0}{\sigma} \stackrel{H_0}{\sim} \mathcal{N}(0,1)$$

Rejection region: $\{T(x) \geq z_\alpha\} \rightarrow$ smallest α at which we reject is $T_{\text{obs}} = z_\alpha$
= p-value

$$\Rightarrow \text{if } Z \sim \mathcal{N}(0,1), p = \Pr[Z \geq T_{\text{obs}}]$$

E.g. Same as above, but $H_A: \mu \neq \mu_0$

$$S_1 = \{|T(x)| \geq z_{\alpha/2}\}$$

$$\Rightarrow p\text{-val is } \alpha \text{ s.t. } \Pr[\bar{Z} \geq |T|] = \frac{\alpha}{2}$$

$$= 2 \Pr[\bar{Z} \geq |T|] \quad \text{under } H_0, U[0,1]$$

Note: $p\text{-val} = \Pr[\bar{Z} > T_{\text{obs}}] = 1 - \Phi(T) \stackrel{H_0}{\sim} U[0,1]$

↳ helps w/ interpretation: under H_0 , p should behave like $\sim U[0,1]$

E.g.: Comparing Variances. X_1, \dots, X_n iid $\mathcal{N}(\mu_1, \sigma_1^2)$, Y_1, \dots, Y_m iid $\mathcal{N}(\mu_2, \sigma_2^2)$

$$H_0: \sigma_1^2 = \sigma_2^2, \quad H_A: \sigma_1^2 \neq \sigma_2^2$$

Notice that $\Omega_0 = \{(\mu_1, \mu_2, \sigma_1, \sigma_2) \in \mathbb{R} : \sigma_1 = \sigma_2\}$ → ugly; we want dist. of T same $\forall \theta \in \Omega_0!$

But recall $\frac{(n-1)S_x^2}{\sigma_1^2} \sim \chi_{n-1}^2$, $\frac{(m-1)S_y^2}{\sigma_2^2} \sim \chi_{m-1}^2$

Note: $U \sim \chi_m^2$, $V \sim \chi_n^2 \Rightarrow \frac{U/m}{V/n} \sim F_{m,n}$ ("only find the dist if you are very bored")

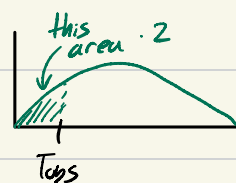
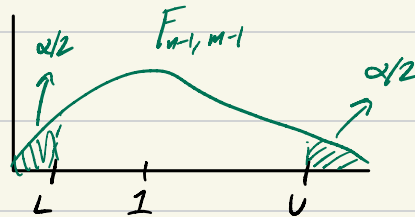
↳ use this w/ above to get $\frac{S_x^2/\sigma_1^2}{S_y^2/\sigma_2^2} \sim F_{n-1, m-1}$

$$\Rightarrow T(x) = \frac{S_x^2}{S_y^2} \stackrel{H_0}{\sim} F_{n-1, m-1}$$

↳ under H_0 , ≈ 1 .

$$\Rightarrow S_1 = \{T > U \text{ or } T < L\}$$

↳ $p = \text{take area of tail of } T_{\text{obs}} \cdot 2$



Likelihood Ratio Tests

X_1, \dots, X_n iid F_θ , $f(x|\theta)$

Case 1: Simple Hypotheses

$$H_0: \theta = \theta_0, \quad H_A: \theta = \theta_1$$

$$LR(x) = \frac{L(\theta_0|x)}{L(\theta_1|x)}$$

$$\Rightarrow S_1 = \{LR(\theta) \leq c\}$$

E.g. $H_0: X_1, \dots, X_n$ iid $\mathcal{N}(1, 1)$, $H_A: X_1, \dots, X_n$ iid $\mathcal{N}(2, 2)$.

$$\theta = (\mu, \sigma^2)$$

$$\begin{aligned} LR(X) &= \frac{\left(\frac{1}{\sqrt{2\pi}}\right)^n \exp\left(-\sum_{i=1}^n \frac{(x_i-1)^2}{2}\right)}{\left(\frac{1}{2\sqrt{2\pi}}\right)^n \exp\left(-\sum_{i=1}^n \frac{(x_i-2)^2}{4}\right)} = (\sqrt{2})^n \exp\left(-\sum_{i=1}^n \left[\frac{x_i^2}{2} - x_i + \frac{1}{2} - \frac{x_i^2}{4} + x_i - 1\right]\right) \\ &= e^{\frac{n}{2}} \exp\left(-\sum_{i=1}^n \frac{x_i^2}{4}\right) (\sqrt{2})^n \Rightarrow S_c = \left\{ \sum x_i^2 \geq c \right\} \end{aligned}$$

Neyman-Pearson Lemma. H_0, H_A simple hypotheses, $S_c = \{LR \leq c\}$, $\alpha =$ Type I error, $\beta =$ Type II. Then, any other test w/ the same α as the LRT has a larger β .

General LRT

$$H_0: \theta \in \Omega_0, H_A: \theta \in \Omega_1, \Omega_0 \cup \Omega_1 = \Omega$$

Likelihood Ratio $L(x) = \frac{\sup_{\theta \in \Omega_0} L(\theta|x)}{\sup_{\theta \in \Omega} L(\theta|x)}$

Thm. $\Omega \in \mathbb{R}^p$ open. Ω_0 is obtained by fixing k coordinates. If MLE assumptions are satisfied,

$$-2 \log L(x) \stackrel{H_0}{\sim} \chi_k^2$$

E.g. X_1, \dots, X_n iid $\mathcal{N}(\mu, \sigma^2)$ (2 params $\Rightarrow p=2$)

$$H_0: \mu = \mu_0 \Rightarrow k=1$$

$$H_0: \mu = \mu_0, \sigma^2 = 1 \Rightarrow k=2$$

Pf.

For simplicity, $\theta \in \Omega \subset \mathbb{R} \Rightarrow p=1, k=1, w/ H_0: \theta = \theta_0$.

Let $\hat{\theta}$ be the MLE.

$$\Rightarrow -2 \log \Lambda(x) = -2(\ell(\theta_0) - \ell(\hat{\theta}))$$

$$\text{Taylor: } \ell(\theta_0) \approx \ell(\hat{\theta}) + \underbrace{(\theta_0 - \hat{\theta})}_{0} \dot{\ell}(\hat{\theta}) + \frac{1}{2} (\theta_0 - \hat{\theta})^2 \ddot{\ell}(\hat{\theta})$$

$$\Rightarrow -2(\ell(\theta_0) - \ell(\hat{\theta})) = \underbrace{(\theta_0 - \hat{\theta})^2}_{\text{normal square}} \underbrace{[-\ddot{\ell}(\hat{\theta})]}_{I(\theta_0)} \rightarrow \text{cancel to get } \mathcal{N}(0,1)$$

Lec 8A - 11/18

- how to do better on final?

↳ remember midterm mistakes

↳ practice a lot

Tests vs. CIs

E.g. X_1, \dots, X_n iid $\mathcal{N}(\mu, 1)$. $H_0: \mu = \mu_0$, $H_A: \mu \neq \mu_0$

$T = \sqrt{n} \left| \frac{\bar{X}_n - \mu_0}{1} \right|$, $S_\alpha = \left\{ \sqrt{n} \left| \frac{\bar{X}_n - \mu_0}{1} \right| \geq z_{\alpha/2} \right\} = \text{rejection region @ level } \alpha$

$$S_0 = \left\{ x : \mu_0 - \frac{z_{\alpha/2}}{\sqrt{n}} \leq \bar{X}_n \leq \mu_0 + \frac{z_{\alpha/2}}{\sqrt{n}} \right\}$$

$$= \left\{ x : \bar{X}_n - \frac{z_{\alpha/2}}{\sqrt{n}} \leq \mu_0 \leq \bar{X}_n + \frac{z_{\alpha/2}}{\sqrt{n}} \right\}$$

⇒ Given data X_1, \dots, X_n , any $\mu_0 \in \bar{X}_n \pm \frac{z_{\alpha/2}}{\sqrt{n}}$ is not rejected by our test!

↳ notice that range is a $1-\alpha$ CI

General Case X_1, \dots, X_n iid F_θ .

Let $A(\theta_0)$ denote the acceptance region of $H_0: \theta = \theta_0$,

$H_A: \theta \neq \theta_0$ at level α .

Let $S(x) = \{ \theta_0 : x \in A(\theta_0) \}$, the set of parameters that are not rejected by our data.

↳ then $S(x)$ is a confidence set

Pf. $P_{\theta_0}[\theta_0 \in S(x)] = P_{\theta_0}[x \in A(\theta_0)] \geq 1 - \alpha$

Note: the converse is true

Let $S(x)$ be a $1-\alpha$ confidence set. Let $A(\theta_0) = \{x: \theta_0 \in S(x)\}$.

Then $A(\theta_0)$ is the acceptance region for a level α test of $H_0: \theta = \theta_0$.

↳ proof is just other way as above

Multiple Testing

Recall the GWAS example: 1M SNPs, 1M tests, 1M p-values

↳ under H_0 , $p \sim U[0,1]$.

E.g.

Suppose we have m tests w/ p-values p_1, \dots, p_m . Assume p_1, \dots, p_m independent.

↳ if all H_0 true, $\Pr[\exists p_i < \gamma] = \Pr[p_{(1)} \leq \gamma] = 1 - \Pr[p_1 > \gamma, \dots, p_m > \gamma] = 1 - (1-\gamma)^m$

↳ w/ $m=1000$, $\gamma=0.001$, $\Pr=0.63$ — high chance to see small p

Could try to control this at level α : ↗ Sidek Correction

$$\Rightarrow 1 - (1-\gamma)^m \leq \alpha \Rightarrow (1-\alpha)^{1/m} \leq 1-\gamma \Rightarrow \gamma \leq 1 - (1-\alpha)^{1/m}$$

$$\begin{aligned} \text{↳ b/c } e^x \approx 1+x \text{ for small } x, \quad 1 - (1-\alpha)^{1/m} &\approx 1 - e^{-\alpha/m} \approx 1 - \left(1 - \frac{\alpha}{m}\right) \\ &= \frac{\alpha}{m} \end{aligned}$$

Bonferroni Correction: $\gamma < \frac{\alpha}{m}$

+ What do we control?

	Reject	Non-reject	
H_0 true	V	V	m_0
H_0 false	T	S	$m - m_0$
	$m - R$	R	m

unknown

all letters are non-neg. integers

- if I reject when $p < \gamma$

↳ V = false positives

↳ T = false negatives

Error Rates

a) Family-wide error rates (FWER)

↳ $\Pr[V > 0]$ ↳ control w/ Sidak, Bonferroni

↳ e.g. $\alpha = 0.05$, $m = 10^6 \rightarrow$ Bonferroni: $\gamma < \frac{\alpha}{m} = 5e-8$

quite small!
b/c no false positives

b) Per-Family Error Rate (PFER)

↳ $E[V]$ (less conservative)

c) False Discovery Rate (FDR)

↳ $E[V/R]$

1995 paper

Benjamini & Hochberg. Control FDR at level α .

1. Order p-values $p_{(1)} \leq \dots \leq p_{(m)}$
2. Find largest j s.t. $p_{(j)} \leq \frac{\alpha \cdot j}{m}$
3. Reject all H_0 corresponding to $p_{(1)}, \dots, p_{(j)}$.

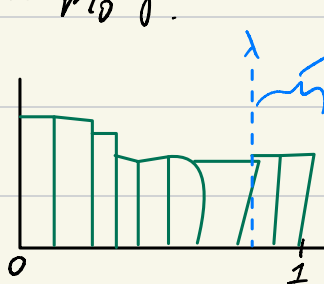
Note: Estimating FDR.

Given a rejection region $[0, \gamma]$, $E[V] \approx m_0 \cdot \gamma$.

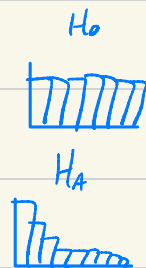
↳ need to know m_0 or $\pi_0 = \frac{m_0}{m}$

Trick: histogram of p-values:

$$\Rightarrow \#(p > \lambda) \approx m_0(1-\lambda)$$



likely from the null:



$$\Rightarrow m_0 \approx \frac{\#(p > \lambda)}{(1-\lambda)} \rightarrow FDR \approx \frac{m_0}{R} = \frac{\#(p > \lambda) \gamma}{R(1-\lambda)}$$

- bias-variance tradeoff: large $\lambda \Rightarrow$ pure signal, but low sample size

Pf (BH).

Let $N \subseteq [m]$ be the set of indices for H_0

$$\hookrightarrow |N| = m_0$$

$$\hookrightarrow \text{Let } \alpha_R = \frac{\alpha \cdot R}{m}, R \in [m]$$

$$\text{Then } E\left[\frac{V}{R}\right] = E\left[\frac{\sum_{k \in N} \mathbb{1}\{p_k \leq \alpha_R\}}{R}\right] \rightarrow \alpha_R \text{ is the threshold, counting qualifying p-vals under } H_0$$

$$= \sum_{k \in N} \sum_{r=1}^m \frac{1}{r} \Pr[p_k \leq \alpha_r, R=r]$$

↳ R_k can be R or $R+1$

Trick: $R_k = \#$ discoveries when doing BH at level α w/ $p_k = 0$.

$$\text{Then, } \Pr[p_k \leq \alpha_r, R_k = r] = \Pr[p_k \leq \alpha_r, R_k = r]$$

$$\hookrightarrow \text{by independence, } = \Pr[p_k \leq \alpha_r] \Pr[R_k = r]$$

$$\Rightarrow E[\bar{V}_R] = \sum_{k \in N} \sum_{r=1}^n \frac{1}{r} \alpha_r \Pr[R_k = r]$$

$$= \sum_{k \in N} \frac{\alpha}{m} \underbrace{\sum_{r=1}^m \Pr[R_k = r]}_{=1}$$

$$\Rightarrow E[\bar{V}_R] = \alpha \cdot \frac{m_0}{m} \leq \alpha$$

Lec 8B - 11/20

Bayesian Statistics

Def: Bayes Formula. $\Pr[A|B] = \frac{\Pr[B|A] \Pr[A]}{\Pr[B]}$.

H_1, \dots, H_k partition of sample space S , D s.t. $\Pr[D] > 0$.

$$\hookrightarrow \Pr[H_i | D] = \frac{\Pr[D|H_i] \Pr[H_i]}{\sum_i \Pr[D|H_i] \Pr[H_i]} \rightarrow \text{how our beliefs are "updated" in light of new information}$$

new info ↙

E.g. CAD = coronary artery disease, PAD = periphery ...

$H_1 = \{\text{CAD}+, \text{PAD}+\}$, $H_2 = \{\text{CAD}+, \text{PAD}-\}$, $H_3 = \{\text{CAD}-, \text{PAD}+\}$,

$H_4 = \{\text{CAD}-, \text{PAD}-\}$

$D = \text{high cholesterol}$

$\Pr[H_i]$ = prior probability, $\Pr[D|H_i]$ = likelihood, $\Pr[H_i|D]$ = posterior

Def: Third Bayes $f(\theta|x) = \frac{f(\theta)f(x|\theta)}{h(x)}$

prior = $f(\theta) \rightarrow \theta$ is RV

likelihood = $f(x|\theta)$

posterior = $f(\theta|x)$

marginal dist. of data

$$= h(x) = \int f(\theta) f(x|\theta) d\theta$$

Bayesian vs. Frequentist

All starts w/ interp of probability.

↳ coin toss $p = 0.5$

↳ frequentist gets 0.5 from $n \rightarrow \infty$ repeats

↳ Bayesians get 0.5 from subjective understanding

Inference

↳ frequentist: θ fixed

↳ Bayesian: θ RV

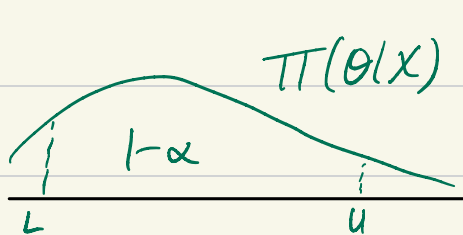
Bayesian Inference

Setup: X_1, \dots, X_n iid $f(x|\theta) \rightarrow$ likelihood

↳ prior: $\theta \sim f_\theta(\theta) = \pi(\theta) \rightarrow$ proportional

↳ posterior: $\pi(\theta|x) = f(\theta|x) \propto \pi(\theta) f(x|\theta)$

+ How do we use the posterior?



Estimation \rightarrow this is it!

1) posterior mean

3) posterior median

2) posterior mode

- Credible Intervals (equiv of CIs)

↳ pick L, U s.t. $\Pr_{\theta|x} [L \leq \theta \leq U] = 1 - \alpha$

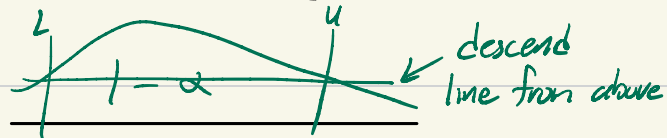
How to pick L, U ?

Def: Equal-Tailed Interval. Let $F_{\theta|x}$ be the posterior cdf. Then
| pick $F_{\theta|x}^{-1}(\alpha/2) \leq \theta \leq F_{\theta|x}^{-1}(1-\alpha/2)$.



Def: High-Posterior Interval. $I = \{\theta : f(\theta|x) \geq c\}$ s.t.

| $\Pr_{\theta|x}(I) = 1 - \alpha$.



+ Hyp. Testing

Let $H_0: \theta \in \Omega_0$, $H_1: \theta \in \Omega_1$.

↳ take ratio of probabilities: $\frac{\Pr_{\theta|x}[\theta \in \Omega_0]}{\Pr_{\theta|x}[\theta \in \Omega_1]}$

E.g: Coin Tossing.

θ = pr of heads. Let X_1, \dots, X_n iid Bern(θ).

We need a prior for θ .

↳ let's start w/ $\theta \sim U[0, 1]$ → basically no prior information

Likelihood

↳ $X = \sum_i X_i \rightarrow f(x|\theta) = \theta^x (1-\theta)^{n-x}$

Posterior

↳ $\pi(\theta|x) \propto f(x|\theta) \pi(\theta) = \theta^x (1-\theta)^{n-x}$, $\theta \in (0, 1)$

Def: Beta Dist. Beta(α, β) has density $f(x) = \frac{1}{\beta(\alpha, \beta)} x^{\alpha-1} (1-x)^{\beta-1}$,


where $\beta(\alpha, \beta) = \frac{1}{\Gamma(\alpha+\beta) \Gamma(\alpha) \Gamma(\beta)}$.


↳ $x \in (0, 1)$

Beta Facts

a) $\text{Beta}(1, 1) = \mathcal{U}[0, 1]$

b) w/ $\text{Beta}(\alpha, \alpha)$

↳ if $\alpha < 1$: 

↳ if $\alpha > 1$: 

c) If $X \sim \text{Beta}(\alpha, \beta)$, $E[X] = \frac{\alpha}{\alpha + \beta}$ and $\text{Var}[X] = \frac{\alpha}{\alpha + \beta} \cdot \frac{\beta}{\alpha + \beta} \cdot \frac{1}{\alpha + \beta + 1}$

d) If $\alpha, \beta > 1$, the mode is $\frac{\alpha - 1}{\alpha + \beta - 2}$.

Back to our example: $\pi(\theta|x) \propto \theta^x (1 - \theta)^{n-x}$

$$\Rightarrow \theta|x \sim \text{Beta}(x+1, n-x+1)$$

Now let's estimate

↳ mean = $\frac{x+1}{n+2}$

↳ mode = $\frac{x}{n}$ → also MLE

Case 2: New Prior.

Now let's say $\theta \sim \text{Beta}(\alpha, \beta)$

$$\Rightarrow \pi(\theta|x) \propto \underbrace{(\theta^{\alpha-1} (1-\theta)^{\beta-1})}_{\text{prior}} \underbrace{(\theta^x (1-\theta)^{n-x})}_{\text{likelihood}} = \theta^{\alpha+x-1} (1-\theta)^{\beta+n-x-1}$$

↳ $\sim \text{Beta}(\alpha+x, \beta+n-x)$

New posterior depends on the data (x) & prior (α, β).

↳ mean = $\frac{\alpha+x}{\alpha+\beta+n}$ → w/ $\alpha = \# \text{ heads}$
↳ posterior "sample size"

Also think of weighted mean:

$$\frac{\alpha}{\alpha+\beta} \left(\frac{\alpha+\beta}{\alpha+\beta+n} \right) + \frac{X}{n} \left(\frac{n}{\alpha+\beta+n} \right) \rightarrow \text{sample weight}$$

\leftarrow prior mean \leftarrow prior weight \leftarrow sample mean

Concrete Example.

$$\alpha = \beta = 3, n = 10, X = 3$$

$$\hookrightarrow \text{we get mean} = \frac{6}{16} = \frac{3}{8} = 0.375 > \text{frequentist } \frac{3}{10} = 0.3$$

b/c the prior was $\frac{1}{2}$

- but how do we pick priors? \exists a few camps:

1) use your judgement

2) assume no information

3) pick a convenient prior

\hookrightarrow Conjugate Family \rightarrow $\left. \begin{array}{l} \theta \quad \text{beta} \\ X \quad \text{binomial} \\ \theta|X \quad \text{beta} \end{array} \right\} \rightarrow \text{subjective, but easy to work w/}$

E.g. Travel to a city. Based on the first bus you see, numbered T , and assuming n busses $\in [N]$, what is n ?

+ Frequentist

$$\text{Likelihood: } f(T|N) = \frac{1}{N} \text{ (discrete uniform)}$$

$$\hookrightarrow \text{mean} = \frac{N}{2} \rightarrow \text{MoM} \Rightarrow \hat{N} = 2T$$

+ Bayesian

Need prior on $N \in \{1, 2, 3, \dots\} \rightarrow$ what to use?

Lec 9A- 12/2

Bayesian Inference

X_1, \dots, X_n iid $f(x|\theta) \rightarrow$ likelihood

$\theta \sim \pi(\theta) \rightarrow$ prior

$\pi(\theta|X) \propto \pi(\theta) \prod f(x_i|\theta) \rightarrow$ posterior

\hookrightarrow prior \hookrightarrow likelihood

E.g. (exponential)

X_1, \dots, X_n iid $\text{Exp}(\lambda)$

Prior: $\lambda \sim \text{Gamma}(k, r)$ (conjugate prior)

Posterior: $\pi(\lambda|X) \propto \pi(\lambda) f(X|\lambda)$

$$= \left(\frac{r^k}{\Gamma(k)} \lambda^{k-1} e^{-r\lambda} \right) \left(\lambda^n e^{-\lambda \sum x_i} \right)$$

$$\propto \lambda^{n+k-1} \exp(-\lambda(r + \sum x_i)) \rightarrow \text{can drop terms w/o}$$

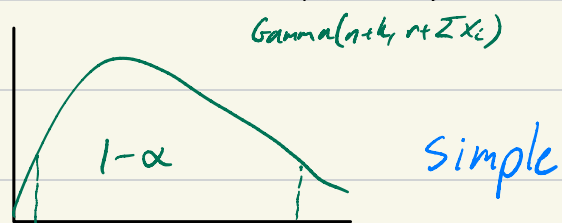
$$\Rightarrow \lambda|X \sim \text{Gamma}(n+k, r + \sum x_i) \text{ params}$$

Estimation

- posterior mean: $\frac{n+k}{r + \sum x_i}$

\hookrightarrow recall the MLE is $\frac{1}{\bar{x}}$; same as bayes if $k=r=0$

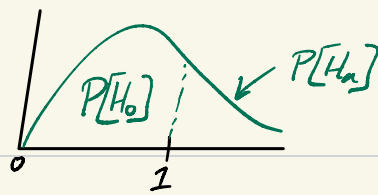
Credible Intervals



Hyp. Testing.

- $H_0: \lambda \in (0, 1), H_a: \lambda \geq 1$

- again look @ posterior:



E.g. Normal Mean.

X_1, \dots, X_n iid $\mathcal{N}(\mu, \sigma^2)$, σ^2 is known

Prior: $\mu \sim \mathcal{N}(\mu_0, \nu^2)$

↳ can be shown that $\mu | X \sim \mathcal{N}\left(\bar{X} \frac{n/\sigma^2}{n/\sigma^2 + 1/\nu^2} + \mu_0 \frac{1/\nu^2}{n/\sigma^2 + 1/\nu^2}, \frac{1}{n/\sigma^2 + 1/\nu^2}\right)$

↳ $\frac{n}{\sigma^2}$ = information in the data; n large = more info

↳ $\frac{1}{\nu^2}$ = information in the prior; ν^2 small = lots of info

Selecting the Priors

Most critical & criticized part of Bayesian inference.

Subjective Determination. $\theta \in \Omega$

a) Say Ω is discrete. Then just use past experience

b) If Ω is an interval, two options:

1) discretize (histogram approach)

2) matching a parametric family

E.g. $X \sim \text{Bin}(n, p)$, p = pr. of getting A in 24410

From previous years, p has mean 0.7 & var 0.1

If we restrict to Beta priors, $p \sim \text{Beta}(\alpha, \beta)$

↳ want $\frac{\alpha}{\alpha + \beta} = 0.7$, $\frac{\alpha\beta}{(\alpha + \beta)^2(\alpha + \beta + 1)} = 0.1 \Rightarrow \alpha = 0.77, \beta = 0.33$

$$\Rightarrow p \sim \text{Beta}(0.77, 0.33)$$

Conjugate Priors

Def: A family F of distributions is said to be closed under sampling for a model $f(x|\theta)$ if $\forall f \in F$, the posterior $f(\theta|x) \propto f(\theta)f(x|\theta) \in F$.

E.g.s

<u>Model</u>	<u>Conjugate Prior</u>
Bernoulli	Beta
Normal	Normal
Exponential	Gamma
Poisson	Gamma

- in the previous example problem, we combined subjective & conjugate priors

Noninformative Priors: favor no particular value.

Case 1: $|\Omega| = a$ (finite)

↳ prior = $\frac{1}{a}$ for each value

Case 1b: Ω is countable (e.g. \mathbb{Z})

↳ sol'n: constant, improper prior (invalid dist.)

Case 2: Ω interval (e.g. $\Omega = [0, 1]$)

↳ Is $U[0, 1]$ noninformative?

↳ Reparameterizing to odds:

$$\eta = \frac{\theta}{1-\theta} \sim \frac{1}{(1+\eta)^2} \text{ on } \mathbb{R}_+ \text{ (can be shown)}$$

↳ Unit is not invariant under reparameterization

Case 3: $\Omega = \mathbb{R}$.

↳ sol'n: flat/constant improper prior (can't integrate to 1)

↳ usable as long as posterior is proper

Jeffrey's Prior

X_1, \dots, X_n iid $f(x|\theta)$ w/ Fisher information $I(\theta)$

Def: Jeffrey's Prior: $\pi_J(\theta) \propto I^{1/2}(\theta)$

Thm: Invariance. Let $\theta \sim \pi_J(\theta)$, and $\eta = g(\theta)$. Then η is distributed according to Jeffrey's prior of $g(\theta)$.

Pf: Lec. notes.

E.g. $X \sim \text{Bin}(n, \theta)$.

$$f(x|\theta) = \binom{n}{x} \theta^x (1-\theta)^{n-x}$$

$$\hookrightarrow I(\theta) = \frac{n}{\theta(1-\theta)}$$

$$\Rightarrow \pi_J(\theta) \propto \theta^{-1/2} (1-\theta)^{-1/2} = \text{Beta}(\frac{1}{2}, \frac{1}{2}) \text{ (proper distribution!)}$$

Note: not all Jeffrey priors are proper.

Lec 9B - 12/4

Statistical Decision Theory

↳ framework for making "decisions"

↳ e.g. choosing estimators: $\theta \in \Omega$ = "truth", \mathcal{D} = collection of decisions

Def: Loss Fn. $L: \Omega \times \mathcal{D} \rightarrow [0, \infty)$. $L(\theta, d)$ = loss associated w/ decision $d \in \mathcal{D}$.

E.g.

$\mathcal{D} = \{d_1, d_2\}$ (say d_1 = football game, d_2 = movie)

$\Omega = \{\theta_1, \theta_2\}$ (θ_1 = rain, θ_2 = no rain)

Loss:

θ_1	d_1	1
θ_1	d_2	$\frac{1}{4}$
θ_2	d_1	0
θ_2	d_2	$\frac{4}{10}$

Assume a 40% chance of rain.

Expected loss:

↳ choose d_1 : $1 \cdot \frac{4}{10} + 0 \cdot (1 - \frac{4}{10}) = \frac{4}{10}$

↳ choose d_2 : $\frac{1}{4} \cdot \frac{4}{10} + \frac{4}{10} (1 - \frac{4}{10}) = \frac{34}{100}$

So the expected loss favors d_2

Inference as a Decision Problem

Problem: X_1, \dots, X_n iid $f(x|\theta)$, $\theta \in \Omega$ w/ prior $\pi(\theta)$.

Decision: choice of estimator $\hat{\theta} = \hat{\theta}(X_1, \dots, X_n)$

↳ our loss fn. is $L: \Omega \times \Omega \rightarrow [0, \infty)$

Types of Loss Functions

- 1) Squared error: $L(\theta, \hat{\theta}) = (\theta - \hat{\theta})^2$
- 2) Absolute error: $L(\theta, \hat{\theta}) = |\theta - \hat{\theta}|$
- 3) Zero-One loss: $L(\theta, \hat{\theta}) = \mathbb{1}_{\{\theta \neq \hat{\theta}\}}$

Def: (frequentist) risk. $R_{\hat{\theta}}(\theta) = R(\theta, \hat{\theta}) = E_{x^n} [L(\theta, \hat{\theta})] = \int L(\theta, \hat{\theta}) f(x^n | \theta) dx^n$

Notes

- a) this is a pre-data measure of performance
- b) if we use squared error loss, $R(\theta, \hat{\theta}) = \text{MSE}(\hat{\theta})$
- c) for a given estimator, risk is a fu. of θ

E.g.s

a) X_1, X_2 iid $\begin{cases} \theta-1 & p=1/2 \\ \theta+1 & p=1/2 \end{cases}$

Let's use zero-one loss, w/ estimators

$$\hat{\theta}_1 = \frac{X_1 + X_2}{2} \rightarrow R(\theta, \hat{\theta}_1) = \Pr_{\theta}[X_1 = X_2] = 1/2$$

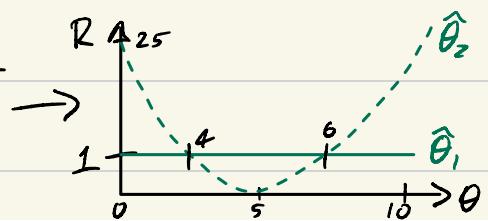
$$\hat{\theta}_2 = X_1 + 1 \rightarrow R(\theta, \hat{\theta}_2) = \Pr_{\theta}[X_1 = \theta + 1] = 1/2$$

↳ neither risks depend on θ , & risk is equal

b) $X \sim \mathcal{N}(\theta, 1)$, $\theta \in (0, 10) = \Omega$ w/ squared error loss

$$\hat{\theta}_1 = X \rightarrow R(\theta, \hat{\theta}_1) = E[(X - \theta)^2] = \text{Var}[X] = 1$$

$$\hat{\theta}_2 = 5 \rightarrow R(\theta, \hat{\theta}_2) = E[(5 - \theta)^2] = (5 - \theta)^2$$



Neither estimator always wins
↳ how to choose?

+ Frequentist Approach

Def: Maximum Risk. $\bar{R}(\hat{\theta}) = \max_{\theta} R(\theta, \hat{\theta})$

Def: Minimax Estimator. $\hat{\theta} \text{ s.t. } \bar{R}(\hat{\theta}) = \inf_{\tilde{\theta}} \bar{R}(\tilde{\theta})$.

+ Bayesian Approach

Def: Bayes Risk. $r(\hat{\theta}) = \int R(\theta, \hat{\theta}) \pi(\theta) d\theta$.

Def: Bayes Estimator. associated w/ a loss fn. & prior π is $\hat{\theta} \text{ s.t. } r(\hat{\theta}) = \min_{\tilde{\theta}} r(\tilde{\theta})$.

E.g. $X \sim \text{Bin}(n, \theta)$

Prior: $\pi(\theta) = U[0, 1]$, loss = square loss

Let $\hat{\theta}_1 = \frac{X}{n}$ and $\hat{\theta}_2 = \frac{\alpha + X}{\alpha + \beta + X}$ (posterior mean for Beta(α, β) prior)

↳ for $\hat{\theta}_1$: $R(\theta, \hat{\theta}_1) = E[(\hat{\theta}_1 - \theta)^2] = \text{Var}[\hat{\theta}_1] = \frac{\theta(1-\theta)}{n}$

$$\bar{R}(\hat{\theta}_1) \stackrel{\theta=1/2}{=} \frac{1}{4n}$$

$$r(\hat{\theta}_1) = \int_0^1 \frac{\theta(1-\theta)}{n} \cdot 1 d\theta = \frac{1}{6n}$$

↳ for $\hat{\theta}_2$: $R(\theta, \hat{\theta}_2) = \frac{n\theta(1-\theta) + (\alpha - (\alpha + \beta)\theta)^2}{(\alpha + \beta + n)^2} = \hat{\theta}_2^{\alpha, \beta}$

↳ turns out that $\hat{\theta}_2^{\frac{\sqrt{n}, \sqrt{n}}{2}, \frac{\sqrt{n}}{2}} = \frac{n}{4(n + \sqrt{n})^2} \rightarrow$ does not depend on θ !

$\bar{R}(\hat{\theta}_2^{\frac{\sqrt{n}, \sqrt{n}}{2}, \frac{\sqrt{n}}{2}}) = \frac{n}{4(n + \sqrt{n})^2} \rightarrow \hat{\theta}_2^{\frac{\sqrt{n}, \sqrt{n}}{2}, \frac{\sqrt{n}}{2}}$ wins! maybe b/c of small part of sample space

$r(\hat{\theta}_2^{\frac{\sqrt{n}, \sqrt{n}}{2}, \frac{\sqrt{n}}{2}}) = \frac{n}{4(n + \sqrt{n})^2} \rightarrow$ loses here (for large n)

The best Bayes estimator here is $r(\hat{\theta}_2^{1,1}) = \frac{1}{6(n+2)} \rightarrow \hat{\theta}_2^{1,1} =$ posterior mean

Def: Posterior Risk. $r(\hat{\theta} | X^n) = \int L(\theta, \hat{\theta}) \pi(\theta | X^n) d\theta$

Thm. Let $\hat{\theta} = \hat{\theta}(X^n)$ be the value that minimizes the posterior risk. Then $\hat{\theta}$ is the Bayes estimator.

Pf.

$$\begin{aligned} r(\hat{\theta}) &= \int R(\theta, \hat{\theta}) \pi(\theta) d\theta = \iint L(\theta, \hat{\theta}) f(x|\theta) \pi(\theta) dx d\theta \\ &\text{by Fubini: } = \iint L(\theta, \hat{\theta}) \pi(\theta|x) f(x) d\theta dx \\ &= \int r(\hat{\theta} | X^n) f(X^n) dX^n. \end{aligned}$$

Thm. Bayes estimators are:

- square loss \rightarrow posterior mean
- absolute loss \rightarrow posterior median
- zero-one loss \rightarrow posterior mode

Pf. (of ①)

Let X be a RV w/ mean μ .

$$\hookrightarrow E[(X-c)^2] = E[(X-\mu + \mu-c)^2]$$

$$= E[(X-\mu)^2] + (\mu-c)^2 \rightarrow \text{minimized when } c = \mu$$

$$r(\hat{\theta} | X) = \int (\theta - \hat{\theta})^2 \pi(\theta | X^n) d\theta$$

$\uparrow_x \quad \uparrow_c \Rightarrow$ minimized for $\hat{\theta} =$ posterior mean