

Introduction to Mathematical Statistics Summary

Introductions

Chapter 1 - what this module is about

- ① why bother: Stats is distinguishing signal from noise. This module introduces the basic grammar you need! Data \longrightarrow models \longrightarrow inference
 [parameters] \longrightarrow [about params]

Chapter 2 - Basic Concepts

- ② Probability Space: (Ω, \mathcal{F}, P)

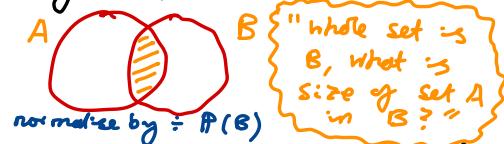
- Ω : $\omega \in \Omega$ are the points / elementary events
- \mathcal{F} : Subsets of Ω forming a σ -algebra [take measure theory!]
- P : probability measure, $P(E)$ assigned to each event E

- ③ Random variables: A function $X: \Omega \rightarrow \mathbb{R}$, measurement, $X(\omega)$ where ω is the outcome that actually happens. $[X \in B] = \{\omega \in \Omega : X(\omega) \in B\}$, the law is the function $B \mapsto P(X \in B)$ - its the distribution.

- Discrete: pmf, $p_X(x) = P(X=x)$
- Continuous: pdf, $f_X: \mathbb{R} \rightarrow [0, \infty)$ $P(X \in B) = \int_{B} f_X(x) dx$
- CDF captures all info about the distribution, $F_X(x) = P(X \leq x)$
 You can calculate expectation, variance, moments $E[X^n]$ ← n th moment
- X_1, \dots, X_n independent if $P(X_1 \leq a_1, \dots, X_n \leq a_n) = P(X_1 \leq a_1) \dots P(X_n \leq a_n)$
- If A, B events, $P(B) > 0$, then the conditional prob of A given B is

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

$$= M_X(u)$$



- ④ Moment Generating Functions: If $E[\exp(ux)] < \infty$ $\forall u \in (-c, c)$, $c > 0$ MGF ✓
 when the MGF exists, it characterizes the distribution

Thm 2.1: (Uniqueness thm for MGFs). X, Y RVs w/ $M_X(u) = M_Y(u)$ $\forall u \in (-c, c)$ for some $c > 0$. Then X and Y identically distributed

- ⑤ Conditional Expectations: Conditional pmf of Y given X

- If $X \geq 0 \Rightarrow E[X|Y] \geq 0$
- $E[aX + bZ|Y] = aE[X|Y] + bE[Z|Y]$
- $E[E[Y|X]] = E[Y]$
- $E[f(Y)X|Y] = f(Y)E[X|Y]$

$$P_{Y|X=x}(y) = P(Y=y | X=x)$$

[positivity]

[linearity]

[Tower Property]

[Tacking out what's known]

Statistical Models

Calculations to practice

Chapter 3 - A Glossary of Statistical Models

A statistical model is a parametrized family of distributions ($P_\theta : \theta \in \Theta$)

- ⑥ Bernoulli($i(p)$): $p \in (0, 1)$

failure \nwarrow success
 $X = \{0, 1\}$

- pmf: $P(X=1) = p = 1 - P(X=0)$
- models: wins/losses at roulette, success/fail
- $E[X] = p$, $\text{var}(X) = p(1-p)$, $M_X(u) = pe^u + 1 - p$ ($E[e^{uX}] = e^u \cdot p + e^0 \cdot (1-p)$)

- ⑦ Binomial(n, p): $n = \# \text{ trials}$, $p \in (0, 1)$ prob. success

- pmf: $P(X=k) = \binom{n}{k} p^k (1-p)^{n-k}$ $k = 0, 1, 2, \dots$

- $\mathbb{E}[X] = np$, $\text{Var}(X) = np(1-p)$, $M_x(u) = (1 + p(e^u - 1))^n$ pt: $X = \sum_{i=1}^n Y_i$ Bernoulli trials
 - Models: n iid coin flips

(8) Geometric(p): $N \sim \text{Geometric}(p)$ is counting time till the first success

- pmf: $P(N=k) = p(1-p)^{k-1}$ for $k=1, 2, 3, \dots$
- $P(N > n) = (1-p)^n$, $\mathbb{E}[N] = \frac{1}{p}$, $\text{Var}(X) = \frac{1-p}{p^2}$, $M_x(u) = \frac{pe^u}{1-(1-p)e^u}$
- Memoryless! $P(N=n+k | N > n) = \frac{p(1-p)^{n+k-1}}{(1-p)^n} = p(1-p)^{k-1} = P(X=k)$ [NO effect of previous event]
- Use: $\mathbb{E}[X] = \sum_{n=1}^{\infty} n P(X=n) = \sum_{k=0}^{\infty} P(X > n)$

(9) NegBinomial(r, p): n iid bernoulli, r th index $N_r = 1$ occurs.

- pmf: $P(N=k) = \binom{k-1}{r-1} p^r (1-p)^{k-r}$ for $k=r, r+1, \dots$
- $\mathbb{E}[N_r] = \frac{r}{p}$, $\text{Var}(N_r) = r \frac{(1-p)}{p^2}$, $M_x(u) = \left(\frac{pe^u}{1-(1-p)e^u}\right)$

Bernoulli (n, p) w/ $\lambda=np$ & large n

↳ # sequences length n w/ r ones & ends w/ a 1

[Independence + (8)]

(10) Poisson(λ): counts rare incidents, e.g. # typographical errors in long text

- pmf: $P(X=k) = \frac{\lambda^k}{k!} e^{-\lambda}$
- $\mathbb{E}[X] = \lambda$, $\text{Var}(X) = \lambda$, $M_x(u) = e^{\lambda(e^u - 1)}$

$$\mathbb{E}[e^{uX}] = \sum_{k=0}^{\infty} e^{uk} \frac{\lambda^k}{k!} e^{-\lambda} = e^{-\lambda} \sum_{k=0}^{\infty} \frac{(ue^u)^k}{k!}$$

(11) Categorical(p): Models questionnaire responses ...

- pmf: $P(X=i) = p_i$ for $i \in \{1, \dots, k\}$ $\sum p_i = 1$
- easy to fit to data, but many parameters. Occam's razor!!! Simplest explanation modn

↳ prob. simplex

(12) Multinomial(n, p): lots of categorical variables together. view $X \sim \text{Categorical}(p)$ as

- pmf: $P(M=\binom{m_1}{m_1, \dots, m_k}) = \binom{n}{m_1, \dots, m_k} \prod_{i=1}^k p_i^{m_i}$ ($\sum m_i = n$) so $X = e_i$ if $x = i$
- See hand stuff later...

$$\binom{n}{m_1, \dots, m_k} = \frac{n!}{m_1! \dots m_k!}$$

*note: measure theory
sets can't model
on all of \mathbb{R}*

(13) Uniform(a, b): model location of random objects on an interval

- $\mathbb{E}[u] = \frac{1}{2}(a+b)$, $\text{Var}(u) = \frac{1}{12}(b-a)^2$, CDF etc... all easy to calc...
- $\mathbb{E}[e^{uX}] = \frac{1}{b-a} \int_a^b e^{ux} dx = \frac{e^{bx} - e^{ax}}{b-a}$ can be a area/volume in high dim...

models radioactive decay

(14) Exponential(λ): continuous version of geometric. Memoryless / constant failure rate

- pdf: $f_T(x) = \lambda e^{-\lambda x}$ for $x > 0$, 0 o/w.
- $\mathbb{E}[T] = \frac{1}{\lambda}$, $\text{Var}(T) = \frac{1}{\lambda^2}$
- CDF: $F_T(x) = 1 - e^{-\lambda x}$ $x > 0$
- Memoryless: $P(T > u+t | T > u) = P(T > t)$
- MGF: $\mathbb{E}[e^{uT}] = \frac{\lambda}{\lambda-u}$ for $u < \lambda$

Deg: (The Gamma Function)

$$T(v) = \int_0^{\infty} x^{v-1} e^{-x} dx, v > 0$$

$$- T(v+1) = v T(v)$$

$$- T(n) = (n-1)!$$

*Think
Substitution
to get this
form*

(15) Gamma(v, λ): waiting times when NOT memoryless

- pdf: $f_X(x) = \frac{\lambda^v}{\Gamma(v)} x^{v-1} e^{-\lambda x}$, $x > 0$, 0 o/w.
- $\mathbb{E}[T] = \dots = \frac{v}{\lambda}$, $\text{Var}(T) = \frac{v}{\lambda^2}$
- $\text{Gamma}(1, \lambda) = \text{exponential}(\lambda) \Rightarrow \sum_{i=1}^n \text{iid exp}(\lambda) \sim \text{gamma}(n, \lambda)$
- $\text{Gamma}(v_1, \lambda) + \text{Gamma}(v_2, \lambda) = \text{Gamma}(v_1 + v_2, \lambda)$
- Mode: compute $\frac{d}{dx} \log f_X(x)$ & argues why max w/ limits

will discuss more distributions when they turn up

Chapter 4 - Transformations

(16) Univariate Transformation Theorem: $g: \mathbb{R} \rightarrow \mathbb{R}$ ctly diff, strictly increasing, X obs as P.V.

CDF F_x , pdf f_x . Then $Y = g(X)$ abs ctz

$$F_Y(y) = F_X(g^{-1}(y)), \quad f_Y(y) = \frac{1}{|g'(g^{-1}(y))|} f_X(g^{-1}(y))$$

*be careful of
range of
functions*

Pl: $f_Y(y) = F'_Y(y) = \frac{d}{dy} (F_X(g^{-1}(y)))$, chain rule & $\frac{d}{dy} (g(g^{-1}(y))) = 1 \dots$

(17) Standard Normal Distribution: $Z \sim N(0,1)$, $f_Z(z) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}z^2}$ $\forall z \in \mathbb{R}$

recall double integral & polar coordinates trick...

- $E(Z) = 0$, $\text{var}(Z) = 1$. $F_Z(x) = \Phi(x)$

(18) General Normal Distribution: $X \sim N(\mu, \sigma^2)$, $f_X(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(x-\mu)^2}$ $\forall x \in \mathbb{R}$

- Note $X = \mu + \sigma Z$

- CDF $F_X(x) = \Phi\left(\frac{x-\mu}{\sigma}\right)$

- Linearity of expectation etc: $E[X] = \mu$, $\text{var}(X) = \sigma^2$

- MGF $M_X(t) = e^{\mu t + \frac{1}{2}\sigma^2 t^2}$ [complete the square]

- Independence: $X \sim N(\mu_1, \sigma_1^2)$, $Y \sim N(\mu_2, \sigma_2^2)$, $X \perp Y \Rightarrow X+Y \sim N(\mu_1 + \mu_2, \sigma_1^2 + \sigma_2^2)$

Pf: just use MGF, independence & uniqueness theorem for MGFs ■

asset prices/exchange rates/amount material

(19) Lognormal Distribution: To model non-negative data. $Y \sim \text{Log}N(\mu, \sigma^2)$ by

$Y = \exp(X)$ w/ $X \sim N(\mu, \sigma^2)$

- use transformation, then to get pdf/cdf.

- MGF does NOT exist! explodes... see notes...

- $E[Y^n] = E[e^{nX}] = M_X(n) = \exp(\mu n + \frac{1}{2}\sigma^2 n^2)$

(20) χ^2 Distribution: Central when assessing statistical procedure. $Y \sim \chi_d^2$ $Y = Z_1^2 + \dots + Z_d^2$ where $Z_i \sim N(0,1)$ iid. $d = \#$ degrees of freedom of chi-squared distribution

- If $d=1$, $Y = Z^2$, $F_Y(y) = \Phi(\sqrt{y}) - \Phi(-\sqrt{y})$, $\downarrow d$ independent gamma

- $Z_i^2 = \text{Gamma}(\frac{1}{2}, \frac{1}{2}) \Rightarrow \chi_d^2 = \text{Gamma}(\frac{d}{2}, \frac{1}{2})$ $\underbrace{\text{Gamma}(\frac{1}{2}, \frac{1}{2})}_{\text{Gamma}(\frac{d}{2}, \frac{1}{2})}$

Pf: $f_Y(y) = F'_Y(y) = \frac{1}{2\sqrt{y}} \Phi'(\sqrt{y}) + \frac{1}{2\sqrt{y}} \Phi'(-\sqrt{y}) = \frac{1}{\sqrt{\pi y}} \exp(-\frac{y}{2})$

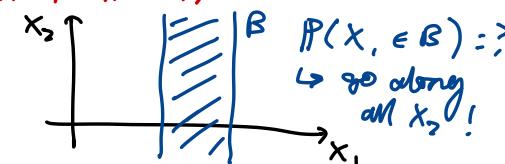
(21) Joint distributions: when R.V.s interact. $X \in \mathbb{R}^n$, then joint dist is $P(X \in A)$

- Multivariate distribution function $F_{x_1, \dots, x_n}(x_1, \dots, x_n) = P(X_1 \leq x_1, \dots, X_n \leq x_n)$

- $f_X(x) = \frac{\partial^n}{\partial x_1 \dots \partial x_n} F_X(x_1, \dots, x_n)$

- Multivariate MGF: $M(a) = E[\exp(a^T x)]$

- To get marginal distribution, just integrate out all else



(22) Covariance / Covariance matrices: How much do R.V.s depend on each other?

- $\text{Cov}(x, y) = E[(X - E[X])(Y - E[Y])] = E[XY] - E[X]E[Y]$. Note $X \perp Y \Rightarrow \text{Cov}(X, Y) = 0$

- Var(ax+by) = $a^2 \text{var}(x) + 2ab \text{cov}(x, y) + b^2 \text{var}(y)$

Pf: $b = E[(ax+by - E[ax+by])^2] = \dots$

- $-1 \leq \text{corr}(x, y) \leq 1$ \downarrow quadratic wrt $a \Rightarrow b^2 - 4ac \leq 0$ \downarrow

Pf: $\text{var}(ax+by) \geq 0 \Rightarrow (2b \text{cov}(x, y))^2 - 4b^2 \text{var}(x) \text{var}(y) \leq 0$ \circlearrowleft

- X_1, \dots, X_n independent wrt finite 2nd moments $\Rightarrow \text{var}(\sum a_i X_i) = \sum a_i^2 \text{var}(X_i)$

- If $X \in \mathbb{R}^n$, $\text{cov}(x)_{ij} = \text{cov}(X_i, X_j)$ [diagonals are variance]

$$\text{corr}(x, y) = \frac{\text{cov}(x, y)}{\sqrt{\text{var}(x)\text{var}(y)}}$$

!!! EXAM HINT ???

Most examples, writing going there =

Example: covariance matrix of $X \sim \text{multinomial}(n, p)$

$E[X] = (np_1, \dots, np_k)$, $\text{var}(x) = np_i(1-p_i) \therefore X_i \sim \text{Binomial}(n, p_i)$

$\text{cov}(X_i, X_j) @ i=j \Rightarrow \text{cov} = \text{var}(X_i) = p_i(1-p_i)$

@ $i \neq j \Rightarrow \text{cov} = E[X_i X_j] - E[X_i]E[X_j] = -p_i p_j$

so

$$\text{cov}(x) = \begin{pmatrix} np_1(1-p_1) & -np_1 p_2 & \cdots & -np_1 p_k \\ -np_2 p_1 & np_2(1-p_2) & & \vdots \\ \vdots & & \ddots & \\ -np_k p_1 & \cdots & & np_k(1-p_k) \end{pmatrix}$$

$X = \text{sum of } n \text{ iid categorical}(p)$

(23) Multivariate Transformation Theorem: $X \in \mathbb{R}^n$ abscts with density $f_x(x)$, $g: \mathbb{R}^n \rightarrow \mathbb{R}^n$ is 1:1 & onto & ctly diff with ctly diff inverse h , then $y = g(x)$ is abscts with density $f_y(y) = J_h(y) f_x(h(y))$, J_h is Jacobian $J_h = |\det(\frac{\partial h}{\partial y})| = \det \begin{pmatrix} \frac{\partial h_1}{\partial y_1} & \dots & \frac{\partial h_1}{\partial y_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial h_n}{\partial y_1} & \dots & \frac{\partial h_n}{\partial y_n} \end{pmatrix}$

Proof: Take $A \subset \mathbb{R}^n$, $P(Y \in A) = P(g(X) \in A) = P(X \in h(A))$

$$= \int_A f_x(x) dx = \int_{h(A)} f_x(h(y)) J_h(y) dy = \int_A f_y(y) dy$$

change of variable
formula for
multiple integrals

(24) Beta(a, b): Very flexible, can model whatever!

- pdf: $f_x(x) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} x^{a-1}(1-x)^{b-1}$, $x \in (0,1)$, $a, b > 0$

just norm constant

- If $X \sim \text{Gamma}(a, 1)$, $Y \sim \text{Gamma}(b, 1)$, $X \perp Y \Rightarrow \frac{X}{X+Y} \sim \text{Beta}(a, b)$

Pf: Define $Z = X+Y$, $B = \frac{X}{X+Y}$, $(B, Z) = g(x, y) = (\frac{X}{X+Y}, X+Y)$, use (23), find marginal π

(25) Quick Linear Algebra Review:

(a) Gram-Schmidt orthogonalization: Gives a basis $\{x_1, \dots, x_n\}$ of \mathbb{R}^n , orthogonalize by
 ① $u_1 = \frac{x_1}{\|x_1\|}$ ② $u_2 = x_2 - (x_2 \cdot u_1)u_1$, ③ normalize u_2 ④ repeat.

(b) $V \in \mathbb{R}^{n \times n}$ symmetric \Rightarrow real eigenvalues & orthonormal basis of eigenvectors. $Vu_i = \lambda_i u_i$

(c) If $U = [u_1 \ \dots \ u_n]$, Then $U^T U = I_n$, $A^T A = (UL^{\frac{1}{2}}U^T)^T(UL^{\frac{1}{2}}U^T)$

(d) $V = U \text{diag}(\lambda_1, \dots, \lambda_n) U^T$

(e) V positive-semidefinite $\Leftrightarrow x^T V x \geq 0 \quad \forall x \in \mathbb{R}^n$ [strict if true def]

(f) V has a sqrt $A = UL^{\frac{1}{2}}U^T$, $L = \text{diag}(\lambda_1^{\frac{1}{2}}, \dots, \lambda_n^{\frac{1}{2}})$ [avoids complex $\lambda_i^{\frac{1}{2}}$ w/ $i > 0$]

(g) If V symmetric positive semi-definite, Then so is A .

(26) Multivariate Normal: $X \sim MVN(\mu, V) \Leftrightarrow \exists \mu \in \mathbb{R}^n$, $A \in \mathbb{R}^{n \times m}$ s.t. $X = \mu + A\tilde{Z}$

$\tilde{Z} \in \mathbb{R}^m$ with $\tilde{Z}_i \sim N(0, 1)$ and $V = A A^T$

- $E[X_i] = E[\mu_i + \sum_j A_{ij} \tilde{Z}_j] = \mu_i$

- $\text{cov}(X_i, X_j) = E[(X_i - E[X_i])(X_j - E[X_j])] = E[(A\tilde{Z})_i (A\tilde{Z})_j] = E[(A\tilde{Z})(A\tilde{Z})^T]_{ij}$

- density: If V has full rank,

$$\text{use multi-dimensionally: } f_x(x) = \frac{1}{\sqrt{(2\pi)^n \det(V)}} \exp\left(-\frac{1}{2} (x - \mu)^T V^{-1} (x - \mu)\right)$$

Pf: $X = g(Z) = \mu + A\tilde{Z} \Rightarrow \tilde{Z} = h(x) = A^{-1}(x - \mu) \Rightarrow f_x(x) = J_h(x) f_{\tilde{Z}}(h(x))$

$$\frac{\partial h}{\partial x} = A^{-1} \Rightarrow J_h(x) = |\det(A^{-1})| = \frac{1}{\sqrt{\det(V)}} \quad \because V = A A^T, \text{ know } f_{\tilde{Z}}(z) = \frac{1}{\sqrt{(2\pi)^n}} e^{-\frac{1}{2} \|z\|^2}$$

$$\text{so } f_{\tilde{Z}}(h(x)) = \frac{1}{\sqrt{(2\pi)^n}} \exp\left(-\frac{1}{2} \|h(x)\|^2\right)$$

$$= \frac{1}{\sqrt{(2\pi)^n}} \exp\left(-\frac{1}{2} (A^{-1}(x - \mu))^T (A^{-1}(x - \mu))\right)$$

$$= \frac{1}{\sqrt{(2\pi)^n}} \exp\left(-\frac{1}{2} (x - \mu)^T V^{-1} (x - \mu)\right)$$

$$(A^{-1})^T A^{-1} = (A A^T)^{-1} = V^{-1}$$

- MGF: $M_x(u) = E[\exp(u^T x)] = \dots = \exp(u^T \mu + \frac{1}{2} u^T V u)$

(27) Fisher's Theorem: X_1, \dots, X_n independent $N(0, \sigma^2)$. $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$, $S_{xx}^2 = \sum_{i=1}^n (x_i - \bar{X}_n)^2$

Then, \bar{X}_n and S_{xx}^2 are independent & $\frac{\bar{X}_n}{\sigma/\sqrt{n}} \sim N(0, 1)$, $\frac{1}{\sigma^2} S_{xx}^2 \sim \chi_{n-1}^2$ independent of Y_2, \dots, Y_n

Pf: Take $u_i = (\frac{1}{\sqrt{n}}, \dots, \frac{1}{\sqrt{n}})^T$, Gram-Schmidt \Rightarrow orthonormal basis u_1, \dots, u_n of \mathbb{R}^n

Set $Y = U^T X$ with $U^T U = I_n$, $Y_i \sim \text{iid } N(0, 1)$. $Y_i = \sum \frac{1}{\sqrt{n}} x_i = \frac{\sqrt{n}}{\sigma} \sum x_i = \sqrt{n} \bar{X}_n$

$$\text{Note: } \sum Y_i^2 = \|U^T X\|^2 = X^T U U^T X = X^T X = \sum x_i^2$$

$$\text{so } S_{xx}^2 = \sum_{i=1}^n (x_i - \bar{X}_n)^2 = \sum_{i=1}^n x_i^2 - n\bar{X}_n^2 = \sum_{i=1}^n Y_i^2 - Y_1^2 = \sum_{i=2}^n Y_i^2 \sim \chi_{n-1}^2 \quad \text{independent of } Y_1 \text{ (rest formula)}$$

Corollary: $X_1, \dots, X_n \sim \text{iid } N(\mu, \sigma^2)$, Then $\frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} \sim N(0, 1) \Rightarrow \bar{X}_n \sim N(\mu, \frac{\sigma^2}{n})$

Pf: easy! Just write out...

(28) The F-distribution: $Y \sim F_{m,n}$ (m, n are degrees of freedom) if

$Y = \frac{X_m/m}{X_n/n}$ where $X_m \sim \chi_m^2$ and $X_n \sim \chi_n^2$

$$- \underline{\text{Lemma:}} \quad f_Y(y) = \frac{\Gamma(\frac{m+n}{2})}{\Gamma(\frac{m}{2}) \Gamma(\frac{n}{2})} \left(\frac{m}{n}\right)^{\frac{m}{2}} y^{\frac{m}{2}-1} \left(1 + \frac{m}{n}y\right)^{-\frac{1}{2}(m+n)}$$

Pf: Stated in notes...

used in statistical comparison of variance

Chapter 5 - Approximation Theorems

(29) Useful inequalities:

① Markov: $X \geq 0, \mathbb{P}(X \geq a) \leq \frac{\mathbb{E}[X]}{a}$ Pf: $Y = \begin{cases} 1 & \text{if } X \geq a \\ 0 & \text{o/w} \end{cases} \Rightarrow 0 \leq aY \leq X$

② Chebyshev: X R.V. $\mathbb{E}[X] = \mu, \text{Var}(X) = \sigma^2$, then $\mathbb{P}(|X - \mu| \geq k) \leq \frac{\sigma^2}{k^2}$ Pf: Markov on $|X - \mu|^2$

(30) Convergence in Random Variables: X, X_1, X_2, \dots R.V. defined on the same prob space (quadratic mean)

- Convergence in Quadratic mean: $X_n \xrightarrow{\text{in 2-mean}} X$ if $\mathbb{E}[|X_n - X|^2] \xrightarrow{n \rightarrow \infty} 0$
Pf: $\mathbb{P}(|X_n - X| > \varepsilon) = \mathbb{P}(|X_n - X|^2 > \varepsilon^2) \leq \varepsilon^{-2} \mathbb{E}[|X_n - X|^2] \xrightarrow{\text{counter}} 0$
- Convergence in Probability: $X_n \xrightarrow{P} X$ in prob if $\forall \varepsilon > 0 \quad \mathbb{P}(|X_n - X| > \varepsilon) \xrightarrow{n \rightarrow \infty} 0$
Pf: Omitted but not hard
counter: $U \sim \text{Uniform}(0,1), X_n = \begin{cases} n & U \leq \frac{1}{n} \\ 0 & \text{o/w} \end{cases}, \mathbb{E}[X_n] = n \times \frac{1}{n} = 1, X_n \xrightarrow{P} X$
- Convergence in Distribution: $X_n \xrightarrow{\text{"weakly" in law}} X$ in distribution if $\mathbb{E}[h(X_n)] \xrightarrow{n \rightarrow \infty} \mathbb{E}[h(X)]$
(for every bdd cts fn h)
 $\Rightarrow h(X_n) \xrightarrow{D} h(X)$

(31) Weak Law of Large Numbers: X_1, \dots, X_n iid, mean μ , var σ^2 . If $X_n = \frac{1}{n} \sum_{i=1}^n X_i$ is sample mean, then $\bar{X}_n \xrightarrow{P} \mu$ in prob.

Pf: $\mathbb{E}[\bar{X}_n] = \frac{1}{n} \sum_{i=1}^n \mathbb{E}[X_i] = \mu, \text{var}(\bar{X}_n) = \frac{1}{n^2} \text{var}(\sum_{i=1}^n X_i) = \frac{\text{var}(X_1)}{n} = \frac{\sigma^2}{n}$, pick $\varepsilon > 0$,
Chebyshev: $\mathbb{P}(|\bar{X}_n - \mu| > \varepsilon) \leq \frac{1}{\varepsilon^2} \frac{\sigma^2}{n} \xrightarrow{n \rightarrow \infty} 0$

Note: Convergence of MGFs can force weak convergence: Z_1, Z_2, \dots R.V. w/ MGF defined on $(-\infty, \infty)$, $c > 0$. If $M_{Z_n}(u) \xrightarrow{P} M_Z(u)$ for $-c < u < c$. THEN, $Z_n \xrightarrow{D} Z$ [weak convergence]

(32) Central Limit Theorem: X_1, X_2, \dots independent iid, mean ≥ 0 , $\text{var}(X_i) = \sigma^2 < \infty$, have a common MGF defined on $(-\infty, \infty) \ni 0$, $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$, Then $\sqrt{n} \bar{X}_n \xrightarrow{D} N(0, \sigma^2)$

Pf: $M(0) = \mathbb{E}[e^{0X}] = 1, M'(0) = 1, M''(0) = \sigma^2$, plug into pwr series: $M(u) = 1 + \frac{1}{2} \sigma^2 u^2 + O(u^3)$
 $S_n = \sum_{i=1}^n X_i, M_{S_n}(u) = (M(u))^n$ [iid], $Z_n = \sqrt{n} \bar{X}_n = \frac{1}{\sqrt{n}} S_n \Rightarrow M_{Z_n}(u) = M_{S_n}(\frac{u}{\sqrt{n}}) = (M(\frac{u}{\sqrt{n}}))^n$
 $M_{Z_n}(u) = \left(1 + \frac{\sigma^2 u^2}{2} + O(\frac{1}{n})\right)^n \xrightarrow{\text{MGF of } N(0, 1)} e^{\frac{\sigma^2 u^2}{2}}$ ← MGF of $N(0, 1)$, so rule $\Rightarrow \sqrt{n} \bar{X}_n \xrightarrow{D} N(0, \sigma^2)$

(33) Continuous Mapping Theorem: $X_i \in \mathbb{R}^k$ R.V. $X_n \xrightarrow{P} X, g: \mathbb{R}^k \rightarrow \mathbb{R}^m$ cts. Then $g(X_n) \xrightarrow{P} g(X)$
Pf: $X_n \xrightarrow{P} X \Rightarrow f(X_n) \xrightarrow{P} f(X)$ & bdd cts $f: \mathbb{R}^k \rightarrow \mathbb{R}$. Pick $h: \mathbb{R}^m \rightarrow \mathbb{R}$ bdd cts, then
 $h \circ g: \mathbb{R}^k \rightarrow \mathbb{R}$ bdd cts $\Rightarrow (X_n \xrightarrow{P} X) \Rightarrow h(g(X_n)) \xrightarrow{P} h(g(X))$ so \Rightarrow

(34) Slutsky's Thm: $X_n \xrightarrow{P} X, Y_n \xrightarrow{P} c$, then $X_n + Y_n \xrightarrow{P} X + c, X_n Y_n \xrightarrow{P} cX, \frac{X_n}{Y_n} \xrightarrow{P} \frac{X}{c} (c \neq 0)$
Pf: omitted

Statistical Inference

Heart of module

Chapter 6 - Likelihood

- Process
- ① Get data $\{x_1, \dots, x_n\}$
 - ② Model via $\{X_1, \dots, X_n; \theta\}$
 - ③ estimate $\hat{\theta}$, accurate? confidence intervals
- hypothesis tests consistent w/ data?

- ⑤ Observed Likelihood Function: Data x_1, \dots, x_n modelled as observed values of R.V.s w/ joint density/pdf $f_{x_1, \dots, x_n}(x_1, \dots, x_n; \theta)$, Observed likelihood function is $L(\theta) = f_{x_1, \dots, x_n}(x_1, \dots, x_n; \theta)$
- Not pdf/density so no need to integrate to 1, function of θ , data is fixed!
- (MLE)

- ⑥ Maximum Likelihood Estimate: MLE $\hat{\theta} \rightarrow$ the value of θ that maximises $L(\theta)$ assuming it exists & is unique: $\hat{\theta} = \arg \max_{\theta} L(\theta) = \arg \max_{\theta} \ell(\theta)$ [same :: log monotonic]
- MLE {(i) makes our data more probable
(ii) value of θ that has most support from data} log-likelihood

Warning: correlation $\not\Rightarrow$ causation. Rains CAUSES ppl to stay indoors (makes sense) Staying indoors CAUSES rain [≠ make sense]

$$\ell(\theta) = \log L(\theta)$$

beware of confounding variables

- ⑦ Example MLE Calculations: CORE SKILL OF THE MODULE!!!

- ⑧ MLE for sequence of Bernoulli R.V.s
- Likelihood: $L(\theta) = \theta^{\sum x_i} (1-\theta)^{n-\sum x_i} \Rightarrow \ell(\theta) = \log \theta (\sum x_i) + \log (1-\theta) (n-\sum x_i)$ only records $\sum x_i$ & n so good for privacy!
 - Diff: $\frac{\partial \ell}{\partial \theta} = \frac{1}{\theta} \sum x_i - \frac{1}{1-\theta} (n-\sum x_i) + \log (1-\theta) (n-\sum x_i)$
 - $\frac{\partial \ell}{\partial \theta} = 0 \Leftrightarrow \dots \Leftrightarrow \hat{\theta}_{MLE} = \frac{1}{n} \sum x_i$
 - JU: $\ell(\theta) \rightarrow 0$ when $\theta \rightarrow 0$ or $\theta \rightarrow 1 \rightarrow \text{MAX}$

Assumption: INDEPENDENCE!!

↳ Q: How were populations sampled?

- ① If you sample from internet - you self-select for nerdy internet type...
- ② Phone for political survey: NOT EVERYONE PICKS UP THE PHONE \Rightarrow NOT a random sample of pop. of interest...
- ③ Snowball sampling: got friends to also answer \Rightarrow Health/Health Qs \Rightarrow very bad!

Getting responses $>$ good quality data
dangerous!!!

- ⑨ MLE for data modelled by binomial RV

- Model: $Y = \sum x_i$ w/ x_i binomial & remember n , $Y \sim \text{Binomial}(n, \theta)$
- Likelihood: $L(\theta) = P_{\theta}(Y=y) = \binom{n}{y} \theta^y (1-\theta)^{n-y} \Rightarrow \ell(\theta) = \log \binom{n}{y} + y \log \theta + (n-y) \log (1-\theta)$
- Diff: $\frac{\partial \ell}{\partial \theta} = \frac{y}{\theta} - \frac{n-y}{1-\theta}, \frac{\partial \ell}{\partial \theta} = 0 \Leftrightarrow \hat{\theta}_{MLE} = \frac{y}{n}$
- JU: $\theta \rightarrow 1, \theta \rightarrow 0, \ell(\theta) \xrightarrow{\theta \rightarrow 1} 0, \ell(\theta) \xrightarrow{\theta \rightarrow 0} 0 \Rightarrow \text{MAX}$

$$-\frac{1}{2} \sigma^{-2} \rightarrow \frac{1}{2} \sigma^{-1}$$

- ⑩ MLE for data modelled by Normal RVs (I)

- Model: $x_1, \dots, x_n \sim n \text{ iid } N(\mu, \sigma^2)$ $\theta = (\mu, \sigma) \in \mathbb{R} \times (0, \infty)$
- $L(\theta) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(x_i - \mu)^2} \Rightarrow \ell(\theta) = -\frac{n}{2} \log(2\pi\sigma^2) - \sum_{i=1}^n \frac{(x_i - \mu)^2}{2\sigma^2}$: log easier to diff
- Diff: $\frac{\partial \ell}{\partial \mu} = 0 \text{ & } \frac{\partial \ell}{\partial \sigma} = 0 \Rightarrow \frac{\partial \ell}{\partial \mu} = \sum_{i=1}^n \frac{(x_i - \mu)}{\sigma^2} \Rightarrow \mu = \frac{1}{n} \sum_{i=1}^n x_i, \frac{\partial \ell}{\partial \sigma} = \frac{-n}{2\pi\sigma^2} - \frac{1}{2\sigma} \sum (x_i - \mu)^2$
- JU: Hessian has strictly negative eigenvalues / $L(\mu, \sigma^2) \rightarrow 0$ on boundary of $\mathbb{R} \times (0, \infty) \rightarrow \text{MAX}$
- So can say that $\hat{\mu}_{MLE} = \frac{1}{n} \sum_{i=1}^n x_i = \bar{x}, \hat{\sigma}_{MLE}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$

- ⑪ MLE for data modelled by discrete uniform RVs

- Model: bag has unknown # balls you to estimate. Select k at random: data x_1, \dots, x_k each #s \hookrightarrow use x_1, \dots, x_k as independent Uniform($\{1, \dots, n\}$) to model this data.
- $L(n) = \prod_{i=1}^k \frac{1}{n} \mathbb{1}\{x_i \in \{1, \dots, n\}\} = \frac{1}{n^k} \mathbb{1}\{n > \max \{x_1, \dots, x_k\}\}$
- $L(n)$ monotone $\because \frac{1}{n^k} \Rightarrow L(n) \downarrow$ for $n > \max \{x_1, \dots, x_k\} \Rightarrow \hat{n}_{MLE} = \max \{x_1, \dots, x_k\}$

"linear regression w/ normally distributed errors"

- ⑫ MLE for simple linear regression [Normal RVs (II)] params random influenced by cov

- Model: $x_i = \text{age}$, $y_i = \text{blood pressure} \Rightarrow Y_i = \alpha + \beta X_i + \sigma \varepsilon_i ; \varepsilon_i \sim \text{iid } N(0, 1)$
- Fe-parametrise: $\bar{x} = \frac{1}{n} \sum x_i, x'_i = x_i - \bar{x}$ [centered], $Y_i = \alpha + \beta x'_i + \sigma \varepsilon_i$ [$\alpha = \bar{y} + \beta \bar{x}$]
- switch: $\varepsilon_i = \frac{1}{\sigma} (y_i - \alpha - \beta x'_i)$ use formula $\Rightarrow L(\alpha, \beta, \sigma) = -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - (\alpha + \beta x'_i))^2$
- Diff: $\frac{\partial \ell}{\partial \alpha} = 0 \Leftrightarrow \hat{\alpha}_{MLE} = \bar{y}, \frac{\partial \ell}{\partial \beta} = 0, \frac{\partial \ell}{\partial \sigma} = 0 \Leftrightarrow \hat{\beta}_{MLE} = \frac{\sum (x'_i - \bar{x})(y_i - \bar{y})}{\sum (x'_i - \bar{x})^2}, \hat{\sigma}_{MLE}^2 = \frac{1}{n} \sum (y_i - (\bar{y} + \hat{\beta} x'_i))^2$
- JU: Can $L(\alpha, \beta, \sigma) \rightarrow \infty$ as $(\alpha, \beta, \sigma) \rightarrow \partial(\mathbb{R} \times \mathbb{R} \times (0, \infty))$? Informed, cheats LN. Generalises to multiple covariates. VERY important! GLMs...

"all models are wrong but some are useful"

calculations easier

"Think political opinion poll: # individuals \rightarrow categories"

⑥ MLE for multinomial Distributions

- Setup: Sample n individuals from pop. w/ k types. n_i of type i , want p_i , $\sum_{i=1}^k n_i = n$
- Model: $N_1, \dots, N_k \sim \text{Multinomial}(n, p)$, $p \in \mathbb{R}^k$. $\sum p_i = 1$, $p_i > 0 \quad \forall i \in \{1, \dots, k\}$, $\sum p_i = 1 - p_{k+1} - \dots - p_m$
- Likelihood: $L(p) = \prod_{i=1}^k p_i^{n_i}$, $n \in \mathbb{R}^k$, max p over parameter space (Lagrange multipliers)
- maximizer: $\ell(p) = \log(\hat{n}) + \sum_{i=1}^k n_i \log p_i$, For $i = \{1, \dots, k-1\}$, Split ... $\Rightarrow \hat{p}_{i, \text{MLE}} = \frac{n_i}{n}$
- JV: When $n_i > 0 \forall i$, $\ell(p) \rightarrow 0$ as $p_i \rightarrow 0$

major tool in modern data science

"when we raise money it's AI, when we hire it's machine learning, when we do the work, it's logistic regression."

blurry line between stats, ML, data science etc... offers data science sounds shiny & gets funding, but substance is often...

⑦ MLE for Logistic regression

- Scenario: $x_i = \text{BMI}$, $y_i = \begin{cases} 1 & \text{type II diabetes} \\ 0 & \text{otherwise} \end{cases}$
- Model: $y_i \sim \text{Bernoulli}(f(x_i; \beta_0, \beta_1))$ where $f(x; \beta_0, \beta_1) = \frac{\exp(\beta_0 + \beta_1 x)}{1 + \exp(\beta_0 + \beta_1 x)}$
- Attractive for interpretation: $1 - f(x; \beta_0, \beta_1) = \frac{1}{1 + \exp(\beta_0 + \beta_1 x)} \Rightarrow \log\left(\frac{P(Y_i = 1)}{P(Y_i = 0)}\right) = \beta_0 + \beta_1 x$
- Likelihood: $L(\beta_0, \beta_1) = \prod_{i=1}^n f(x_i; \beta_0, \beta_1)^{y_i} (1 - f(x_i; \beta_0, \beta_1))^{1-y_i}$
- MLE: compute numerically.

Chapter 7 - Repeated Sampling Principle

How accurate are our MLE estimates?

Does NOT depend on true model parameters!!!

e.g. MLE

- ⑧ Statistic: A statistic is simply a function $t(x_1, \dots, x_n)$ of the observed data x_1, \dots, x_n .

- ⑨ Repeated Sampling Principle: Statistical decision procedures should be evaluated on the basis of their behaviors in hypothetical repetitions of the experiment that generated the original data x_1, \dots, x_n . If we observe data that \neq fit model \Rightarrow reject model.

Note: if data x_1, \dots, x_n modelled by X_1, \dots, X_n w/ param θ , if $t(x_1, \dots, x_n)$ estimates model parameter θ , then $T = t(X_1, \dots, X_n)$ is an estimator for θ .

want small BUT large size & we don't care...

$$E_\theta[T] = \theta$$

- ⑩ Bias of an Estimator: Bias of the estimator $T = T(x_1, \dots, x_n)$ for θ is $E_\theta[T - \theta]$

- ⑪ Mean-Square-Error of an estimator: MSE of T for θ is $E_\theta[(T - \theta)^2]$

- ⑫ Standard error of an estimator: S.E. \approx RMSE \approx RMSD of T for θ is $\sqrt{E_\theta[(T - \theta)^2]}$

- ⑬ Bias-variance decomposition of MSE: $MSE = E_\theta[(T - \theta)^2] = \text{Var}_\theta(T) + (E_\theta(T) - \theta)^2$

$$\text{pf: } E[(T - \theta)^2] = E[(T - E_\theta[T] + E_\theta[T] - \theta)^2] = \underbrace{E[(T - E_\theta[T])^2]}_{\text{variance}} + 2E[(T - E_\theta[T])(E_\theta[T] - \theta)] + \underbrace{(E_\theta[T] - \theta)^2}_{\text{bias}^2}$$

- ⑭ Consistency of Estimators: A sequence of estimators T_1, T_2, \dots is consistent for θ if $T_n \xrightarrow{n \rightarrow \infty} \theta$ in probability.

Note: $T_n = T_n(x_1, \dots, x_n)$ estimates θ . $MSE(T_n) = E_\theta[(T_n - \theta)^2] \xrightarrow{n \rightarrow \infty} 0$

(*) $\Rightarrow T_n \xrightarrow{n \rightarrow \infty} \theta$ in probability

Note: MLE could be biased but consistent.

- Data $x_1, \dots, x_n \sim \text{iid Uniform}(0, \theta)$, θ unknown
- MLE $\hat{\theta} = M_n = \max\{x_1, \dots, x_n\}$, $\ell(\theta) = \frac{1}{\theta^n}$ for $\theta > \max\{x_1, \dots, x_n\}$
- CDF of M_n : $F_{M_n}(m) = \left(\frac{m}{\theta}\right)^n$ for $0 < m < \theta$, $P(M_n \leq m) = (P(X_1 \leq m))^n$
- PDF of M_n : $f_M(m) = \frac{nm^{n-1}}{\theta^n}$ for $0 < m < \theta$ [diff]
- $E_\theta[M_n] = \int m \cdot \frac{nm^{n-1}}{\theta^n} dm = \frac{n}{n+1} \theta \Rightarrow \text{bias} = E_\theta[M_n] - \theta = -\frac{1}{n+1} \theta$
- $E_\theta[M_n^2] = \dots = \frac{n\theta^2}{n+2}$, $MSE = E_\theta[(M_n - \theta)^2] = \dots = \frac{2}{(n+1)(n+2)} \theta^2$

more cells shows biased but not consistent...

- ⑮ MLE General Picture: Under 'suitable regularity conditions', sequences of MLE estimators are asymptotically unbiased, iff sequences of MLEs have MSE which are asymptotically as small as possible. Hence MLE consistent (* * *)

{ Sequences of MLEs $\not\rightarrow$ boundary of param space
data $\xrightarrow{n \rightarrow \infty} \infty$
dim parameters θ cannot increase for sequences of MLEs }

(46) MLE may not yield optimal MSE: idea is that intuitively natural estimators can actually be biased. Need to adjust to make them unbiased. *Say we wanna estimate μ, σ^2 ...*

$$\text{MSE} = \text{Variance} + \text{bias}^2$$

Pick your poison...

This is straightforward

(a) Estimation of mean using sample mean

- Take X_1, \dots, X_n iid mean μ , variance σ^2 . If $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$ is sample mean
- $E[\bar{X}_n] = \frac{1}{n} \sum E[X_i] = \frac{n\mu}{n} = \mu \Rightarrow \text{Bias: } E[\bar{X}_n] - \mu = 0 \Rightarrow \text{unbiased}$
- $MSE = \text{var} + \text{bias}^2 = \text{var}(\bar{X}_n) = \frac{1}{n^2} \sum \text{var}(X_i) = \frac{\sigma^2}{n} \leftarrow \text{MSE}$

*MLE variance is biased
sample variance is unbiased but P variance*

(b) Estimation of variance using sample variance

- MLE is biased! $\hat{\sigma}_{MLE}^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2$.
- Sample variance is unbiased! $S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2 = nE[X_i^2] - nE[\bar{X}_n^2]$
- why? $E[S^2] = \frac{1}{n-1} E\left[\sum_{i=1}^n X_i^2 - 2 \sum_{i=1}^n \bar{X}_n X_i + \sum_{i=1}^n \bar{X}_n^2\right] = \frac{1}{n-1} (nE[X_i^2] - 2nE[\bar{X}_n^2] + nE[\bar{X}_n^2])$
 $\text{So bias} = E[S^2] - \sigma^2 = 0$
- Let $S_{n,a} = \frac{1}{a} \sum_{i=1}^n (X_i - \bar{X}_n)^2$
- $\arg\min(a \cdot \text{MSE}(S_{n,a})) = a^* = n+1$

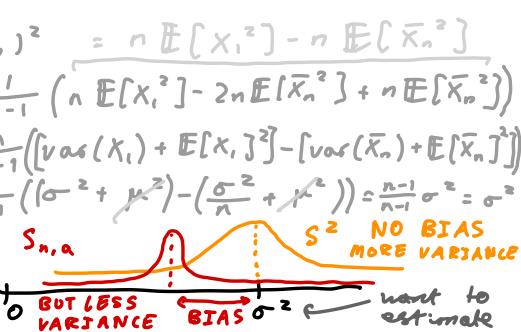
$$\text{MSE}(\hat{\sigma}_{MLE}^2) < \text{MSE}(S^2)$$

biased unbiased
better MSE worse MSE

Trade off between bias & variance

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

$$\Leftrightarrow E[z^2] = \text{Var}(z) + E[z]^2$$



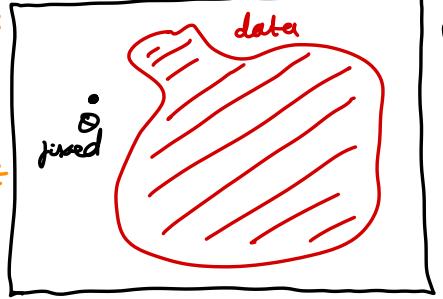
Note: method of moments is an alternative to MLE. 200 years old. MLE is gold standard!

AIM: construct intervals using data \Rightarrow where unknown θ lies...

Chapter 8 - Confidence Intervals

(47) Confidence Interval: Data x_1, \dots, x_n modelled by X_1, \dots, X_n w/ param $\theta \in \Theta$. A pair of statistics $l(x_1, \dots, x_n), r(x_1, \dots, x_n)$ specifying an interval $[l(x_1, \dots, x_n), r(x_1, \dots, x_n)]$ which is called a $100(1-\alpha)\%$ CI for $\phi(\theta)$ if for any $\theta \in \Theta$, $P_\theta[l(X_1, \dots, X_n) \leq \phi(\theta) \leq r(X_1, \dots, X_n)] = 1 - \alpha$

"Throw data at wall;
100(1- α)% chance of hitting θ in your wet sponge of data"



(48) Gaussian percentiles: For $\alpha \in (0, 1)$ set the 100(1- α) percentile of the standard Normal distribution to be z_α s.t. if $Z \sim N(0, 1)$, then $P(Z \leq z_\alpha) = 1 - \alpha$

notes:

- $\alpha \cdot P(Z \leq z_\alpha) = 1 - \alpha$
- $\Leftrightarrow \Phi(z_\alpha) = 1 - \alpha$
- $\Leftrightarrow z_\alpha = \Phi^{-1}(1 - \alpha)$
- $\Phi(P(Z > -z_\alpha)) = P(Z < z_\alpha) = 1 - \alpha$

all integrates to one!

(49) Normal CI with known variance, unknown mean: data x_1, \dots, x_n modelled by $X_i \sim N(\mu, \sigma^2)$ unknown μ , KNOWN $\sigma^2 > 0$. Choose MLE estimate for μ : $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$; $[\bar{X}_n - z_{\alpha/2} \frac{\sigma}{\sqrt{n}}, \bar{X}_n + z_{\alpha/2} \frac{\sigma}{\sqrt{n}}]$ $\bar{X}_n \sim N(\mu, \frac{\sigma^2}{n}) \Rightarrow \frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} \sim N(0, 1) \Rightarrow P_{\mu, \sigma^2}[\alpha \leq \frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} \leq b] = \Phi(b) - \Phi(a) \Rightarrow$ has prob. $1 - \alpha$ of containing μ

(50) Pivot: is a random variable whose distribution \neq depend on params of underlying model

Note: If we can find a pivot of simple form \Rightarrow easy to build a CI!! *n degrees of freedom*

(51) Student's t-distribution: $Z \sim N(0, 1)$, $V \sim \chi^2_n$ are independent R.V.s $T = \frac{Z}{\sqrt{V/n}} \sim t_n$

(52) More Normal confidence Intervals: Data x_1, \dots, x_n modelled by $X_i \sim N(\mu, \sigma^2)$

(a) Estimation of variance (ignore mean)

- MLE: $\hat{\mu} = \bar{x}$, $\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{n-1}{n} S^2$ where $S^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$
- Fischer's thm $\Rightarrow \hat{\mu}$ independent from $\hat{\sigma}^2$ and $\frac{n-1}{\sigma^2} S^2 = \frac{1}{\sigma^2} \sum_{i=1}^n (x_i - \bar{x})^2 \sim \chi^2_{n-1}$
- Pivot: $\alpha \in (0, 1) \Rightarrow P_{\mu, \sigma^2}(\chi^2_{n-1, 1-\alpha} \leq \frac{n-1}{\sigma^2} S^2 \leq \chi^2_{n-1, \alpha}) = 1 - \alpha$

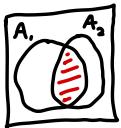
$$\Leftrightarrow P_{\mu, \sigma^2}(\frac{(n-1)S^2}{\chi^2_{n-1, 1-\alpha}} \leq \sigma^2 \leq \frac{(n-1)S^2}{\chi^2_{n-1, \alpha}}) = 1 - \alpha$$

unbiased sample variance

pivot!

$\chi^2_{n-1, \alpha}$ is the 100(1- α) percentile of $Y \sim \chi^2_{n-1}$
 $P(Y \leq \chi^2_{n-1, \alpha}) = 1 - \alpha$

- (b) Estimation of mean w/ variance unknown
- Setup: $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i \sim N(\mu, \frac{\sigma^2}{n}) \Rightarrow \frac{\bar{X}_n - \mu}{\frac{\sigma}{\sqrt{n}}} \sim N(0, 1)$, $S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2$
 - Fisher's Thm: $\frac{n-1}{\sigma^2} S^2 \sim \chi^2_{n-1}$ independent of \bar{X}_n $\Rightarrow T = \frac{\bar{Z}}{\sqrt{\frac{S^2}{n-1}}} = \frac{\bar{X}_n - \mu}{\frac{\sigma}{\sqrt{n-1}}} \sim t_{n-1}$
 - CI: pivot is $T = \frac{\bar{X}_n - \mu}{\frac{S}{\sqrt{n}}} \Rightarrow P(-t_{n-1, \alpha/2} \leq T \leq t_{n-1, \alpha/2}) \sqrt{\frac{S^2}{n-1}} = P(\bar{X}_n - t_{n-1, \alpha/2} \frac{S}{\sqrt{n}} \leq \mu \leq \bar{X}_n + t_{n-1, \alpha/2} \frac{S}{\sqrt{n}}) = 1 - \alpha$
- Note:** widths of CI for μ & σ^2 are not independent...
- $P(A_1 \cap A_2) > 1 - P(A_1^c) - P(A_2^c)$
- $\mu \in CI_{\sigma^2} \quad \sigma^2 \in CI_{\sigma^2}$ w/ prob $1 - 2\alpha$
- $\mu, \sigma^2 \in$ their CI



(c) Regression case (use fisher's thm for linear regression) ...

- (3) CI for Binomial Distribution: Unknown proportion p of pop. carry genetic marker
- $\hat{p} = \frac{x}{n}$ where x carry marker in sample size of n , model as $X = \sum_{i=1}^n B_i$, $B_i \stackrel{iid}{\sim} \text{Bern}(p)$
 - CLT: $x \sim \text{Sum iid R.V.} \xrightarrow{CLT} X \stackrel{d}{\sim} N(np, np(1-p)) \Rightarrow \frac{(X - np)\sqrt{n}}{\sqrt{np(1-p)}} \stackrel{d}{\approx} N(0, 1)$ (doesn't involve p)
 - WLLN: $\frac{X}{n} \xrightarrow{P} p$
 - CI mapping thm: $\sqrt{\frac{X}{n}(1-\frac{X}{n})} \xrightarrow{d} \sqrt{p(1-p)}$
 - Slutsky's thm: $\sqrt{n}(\frac{X}{n} - p) \xrightarrow{d} N(0, 1) \Rightarrow P_p[-z_{\alpha/2} \leq \frac{\sqrt{n}(\frac{X}{n} - p)}{\sqrt{\frac{X}{n}(1-\frac{X}{n})}} \leq z_{\alpha/2}] = 1 - \alpha$
 - **approximate CI for $p \rightarrow$ need big n**
- approximate CI for $p \rightarrow$ need big n**

- (4) Multinomial CI: Lead of one political party over another
- Setup: data n_1, \dots, n_k ; $\sum n_i = n$, modelled by multinomial(n, p). Estimate $p_1 - p_2$
 - MLE: $p_1 - p_2 \hat{s} \hat{p}_1 - \hat{p}_2 = \frac{n_1}{n} - \frac{n_2}{n}$
 - Mean: $E[\hat{p}_1 - \hat{p}_2] = \dots = p_1 - p_2$, $\text{Var}(\hat{p}_1 - \hat{p}_2) = \text{Var}(\hat{p}_1) - 2\text{cov}[\hat{p}_1, \hat{p}_2] + \text{Var}(\hat{p}_2) = \dots$
 - Approx. Pivot: $\frac{\frac{n_1 - n_2}{n} - (p_1 - p_2)}{\sqrt{\frac{(n_1 - n_2)^2}{n} - (\hat{p}_1 - \hat{p}_2)^2}} \approx N(0, 1) \Rightarrow \hat{p}_1 - \hat{p}_2 \pm \sqrt{\frac{\hat{p}_1 - \hat{p}_2 - (\hat{p}_1 - \hat{p}_2)^2}{n}} \times z_{\alpha/2}$
- replace some params w/ their estimators via WLLN / CLT**

(5) Median instead of mean:

- $X_1, \dots, X_n \stackrel{iid}{\sim}$ continuous $\Rightarrow \frac{1}{n} \sum X_i \sim$ Cauchy \Rightarrow cannot estimate θ based on \bar{X}_n
- Trick: compute median instead of mean! Define $Y_i = \begin{cases} 1 & \text{if } X_i \leq \theta \\ 0 & \text{or w/o } Y_i \sim \text{Bern}(\frac{1}{2}) \end{cases}$
- $P(X_{(k)} \leq \theta \leq X_{(n-k+1)}) = P(\sum_{i=1}^k Y_i \in (k, n-k+1)) = A_k$
- Pict: k s.t. $P(A_k) > 1 - \alpha$.
- Pivot: $\sum_{i=1}^k Y_i \sim \text{Binomial}(n, \frac{1}{2})$
- Have another look... a little confusing

Chapter 9 - Statistical Tests

From estimating \mapsto Testing!

possible param space

- (6) Testing setup: Given data x_1, \dots, x_n modelled by a joint distribution depending on $\theta \in \Theta$. The statistical model \rightarrow likelihood $L(\theta; x_1, \dots, x_n)$. Define:
- (i) null hypothesis $\theta_0 \in \Theta_0 \subset \Theta$ depends on data
 - (ii) alternative hypothesis $\theta \in \Theta_1 \subset \Theta$
- intuition:** $L(\theta) \approx$ how much does data support param value θ ?
- (7) Likelihood Ratio: $W(x_1, \dots, x_n) = W(\theta_0, \theta_1) = \frac{\sup_{\theta \in \Theta_1} \{L(\theta)\}}{\sup_{\theta \in \Theta_0} \{L(\theta)\}}$
- note:** Large value of $W \Rightarrow$ "evidence vs null hypothesis"
- \hookrightarrow There are $\theta \in \Theta$, that are much more strongly supported by the data than any $\theta \in \Theta_0$

best case scenario under alt. hypothesis

best case scenario under null hypothesis

measures the maximum probability of seeing evidence against the null hypothesis when null hypothesis is actually true.

[Likelihood Ratio Test - LRT]

\hookrightarrow unlikely if null hypothesis true \Rightarrow small $P \Rightarrow$ reject θ_0

measures the maximum probability of seeing evidence against the null hypothesis when null hypothesis is actually true.

- (8) P-value of a Test Statistic: Suppose $W(x_1, \dots, x_n)$ is a statistic to test θ_0 against θ_1 . $W(x_1, \dots, x_n)$ is P.V.
- The p-value for $W(x_1, \dots, x_n)$ is $p(x_1, \dots, x_n) = \sup_{\theta \in \Theta_1} P_{\theta} (W(x_1, \dots, x_n) \geq W(x_1, \dots, x_n))$
- ↳ small p-value \Rightarrow strong evidence vs null.**

How often will test statistic be wrong?

- (59) Decision Theoretic Interpretation: Null hypothesis simple so $\Theta_0 = \{\theta_0\} \subset \Theta$, [nested]
- (a) Pick significance level $\alpha \in [0, 1)$ [Based on contextual evaluation of risk of making mistakes]
 - (b) Model the scenario: obtain data x_1, \dots, x_n & model w/ R.V.s \rightarrow likelihood $L(\theta; x_1, \dots, x_n)$
 - (c) Test null hypothesis vs. alt. hypothesis using LRT (or other chosen test)
 - (d) Consider the resulting p-value:
 - $p(x_1, \dots, x_n) \leq \alpha \Rightarrow$ reject null hypothesis $\theta = \theta_0$ at significance level α
 - $p(x_1, \dots, x_n) > \alpha \Rightarrow$ do not reject θ_0 at S.I. α [$\not\Rightarrow$ accept θ_0 , - could be worse!]

- (60) Uniform Distribution of p-value under a simple null hypothesis:
- Suppose $\Theta_0 = \{\theta_0\}$, $w(x)$ has cts dist for $\theta = \theta_0$, then, under null, P_{θ_0} , p-value $P(X) \sim \text{Uniform}(0, 1)$
- PF: Let Z, X be R.V. $Z \perp X$, distributed P_{θ_0} , by test statistic def distribution of $w(X)$ when $\theta = \theta_0$
- $$p(X) = P_{\theta_0}(w(Z) \geq w(X) | X) = 1 - P_{\theta_0}(w(Z) < w(X) | X) = 1 - F_w(w(X); \theta_0) \quad 1 - (1 - \alpha) = \alpha$$
- $$P_{\theta_0}(p(X) \leq \alpha) = P_{\theta_0}(1 - F_w(w(X); \theta_0) \leq \alpha) = P_{\theta_0}(F_w(w(X); \theta_0) \geq 1 - \alpha) = 1 - P_{\theta_0}(F_w(w(X); \theta_0) \leq 1 - \alpha)$$
- (The t-test)

- (61) One Sample t-test:
- Setup: data x_1, \dots, x_n modeled by $X_1, \dots, X_n \stackrel{iid}{\sim} N(\mu, \sigma^2)$, $\Theta_0 = \{(\mu, \sigma) : \mu = \theta_0, \sigma > 0\}$, $\Theta_1 = \{(\mu, \sigma) : \mu \in \mathbb{R}, \sigma > 0\}$
 - Likelihood Ratio: $W(\theta_0, \theta_1) = \frac{L(\bar{x}, \hat{\sigma}^2)}{L(\theta_0, \hat{\sigma}^2)} = \dots = \left(1 + \frac{n\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2}\right)^{-\frac{n}{2}} = \left(1 + \frac{t^2}{n-1}\right)^{-\frac{n}{2}}$ where $t = \frac{\bar{x}}{\sqrt{\frac{s^2}{n-1}}}$, $s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$
 - Fischer's thm: under $\mu = 0$, $t(X) \sim t_{n-1}$ [30]
 - $w(x)$ T function of $t^2 \Rightarrow p(x) = P_{\theta_0, \theta_1}(w(X) > w(x)) = P_{\theta_0, \theta_1}(t(X)^2 > t(x)^2) = 2(1 - F_{t_{n-1}}(|t(x)|))$
 - Reject null hypothesis if $|t(x)| \geq t_{n-1, \alpha/2}$
- Note: rejecting null hypothesis $\Leftrightarrow \theta \notin 100(1-\alpha)\%$ confidence interval for μ . LWt! be very careful what α value you choose!

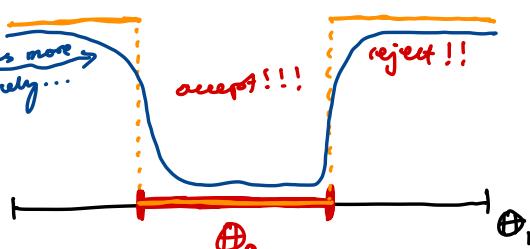
- (62) Significance & power of a test: Θ -parameter space, $\mathcal{X}_0 : \theta \in \Theta_0$, $\mathcal{X}_1 : \theta \in \Theta_1$, $\Theta_0 \subset \Theta_1 \subset \Theta$
- A test has significance level α if it has prob $\leq \alpha$ of rejecting \mathcal{X}_0 when it is true
 - A test \rightarrow test statistic $w(x)$: reject \mathcal{X}_0 if $w(x) \geq c$ where c chosen to achieve sig.level α .
 - Set $R_\alpha = \{x : w(x) \geq c\}$. Reject \mathcal{X}_0 if $x \in R_\alpha$ [rejection region]
 - Type I error: rejecting \mathcal{X}_0 when \mathcal{X}_0 is true.
 - Type II error: not rejecting \mathcal{X}_0 when \mathcal{X}_0 false.
 - In brief:
 - Power function of a test: $\beta(\theta) = P_{\theta}(X \notin R_\alpha)$ for $\theta \in \Theta_1 \setminus \Theta_0$
 - Significance level / size of test: $\alpha = \sup_{\theta \in \Theta_0} \beta(\theta)$

Note: Ideally $\beta(\theta) \approx 1$ for $\theta \in \Theta_1 \setminus \Theta_0$.

$\alpha > \beta(\theta) \approx 0$ for $\theta \in \Theta_0$

\hookrightarrow Not possible $\because \beta(\theta)$ is continuous... *this is more likely...*

Power = $P(\text{reject } \mathcal{X}_0)$ when \mathcal{X}_0 false
Size = max $P(\text{reject } \mathcal{X}_0)$ when \mathcal{X}_0 true



Booster note to recall:

- (A) MGFs Specifying Distributions: If X, Y are two random variables on \mathbb{R} with $M_X(t) = M_Y(t)$ on some interval $(-\epsilon, \epsilon)$ with $\epsilon \in (-\epsilon, \epsilon)$, Then X and Y have the same distributions
 X abs ctRV.
- (B) Univariate Transformation: If $Y = g(X)$ where $g: \mathbb{R} \rightarrow \mathbb{R}$ is increasing, bijective, continuously differentiable & positive, then $F_Y(y) = F_X(g^{-1}(y))$ and $f_Y(y) = \frac{1}{g'(g^{-1}(y))} f_X(g^{-1}(y))$
- (C) $X \sim MVN$: $X \in \mathbb{R}^n$ has MVN distribution if $X = \mu + A\mathbf{Z}$ where $\mu \in \mathbb{R}^n$ is the mean vector, $\mathbf{Z} \in \mathbb{R}^m$ is an m -vector of iid $N(0, 1)$ random numbers and $A \in \mathbb{R}^{n \times m}$ is a matrix such that $V = A A^T \in \mathbb{R}^{n \times n}$ is the variance-covariance matrix, so that

$$X \sim MVN(\mu, V)$$
- (D) A distribution is memoryless if $P(X > u+t | X > t) = P(X > u)$
- (E) Fisher's theorem: If X_1, \dots, X_n are a sequence of $N(\mu, \sigma^2)$ random variables, Then setting $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$ and $S^2 = \sum_{i=1}^n (X_i - \bar{X}_n)^2$, Then the following are independent

$$\frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} \sim N(0, 1), \quad \frac{1}{\sigma^2} S^2 \sim \chi_{n-1}^2$$
- (F) Convergence:
- $X_n \xrightarrow{\text{in 2-mean}} X$ if $\mathbb{E}[(X_i - X)^2] \xrightarrow{n \rightarrow \infty} 0$
 - \downarrow
 - $X_n \xrightarrow{\text{in prob}} X$ if $\forall \epsilon > 0 \quad P(|X_n - \mu| > \epsilon) \xrightarrow{n \rightarrow \infty} 0$
 - \downarrow
 - $X_n \xrightarrow{\text{in distribution}} X$ if $\mathbb{E}[h(X_n)] \xrightarrow{n \rightarrow \infty} \mathbb{E}[h(x)]$ bdd $h: \mathbb{R} \rightarrow \mathbb{R}$
- (G) WLLN: X_1, X_2, \dots sequence of random variables iid with finite mean μ and variance σ^2 . If $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$ is the sample mean, Then $\bar{X}_n \xrightarrow{\text{prob}} \mu$
Pf: $P(|\bar{X}_n - \mu| > \epsilon) \leq \frac{\text{var}(\bar{X}_n)}{\epsilon^2} = \frac{\sigma^2}{n\epsilon^2} \xrightarrow{n \rightarrow \infty} 0 \quad \forall \epsilon > 0$
- (H) CLT: X_1, X_2, \dots sequence of iid random variables w/ zero mean and finite variance, and have a common MGF defined on some interval $(-\epsilon, \epsilon)$ with $\epsilon \in (-\epsilon, \epsilon)$. Then $\sqrt{n} \bar{X}_n \xrightarrow{\text{dist}} N(0, \sigma^2)$
- (I) Cts mapping Thm: If $X_n \xrightarrow{\text{in } \mathbb{R}^n}$ and $g: \mathbb{R}^n \rightarrow \mathbb{R}^m$ is continuous then $g(X_n) \xrightarrow{\text{in } \mathbb{R}^m}$ Shatry's & cts mapping are both convergence in dist.
- (J) Shatry's Thm: If $X_n \xrightarrow{\text{in }} X$, $Y_n \xrightarrow{\text{in }} c$, Then (i) $X_n + Y_n \xrightarrow{\text{in }} X + c$ (ii) $X_n Y_n \xrightarrow{\text{in }} Xc$ and (iii) $\frac{X_n}{Y_n} \xrightarrow{\text{in }} \frac{X}{c}$ if $c \neq 0$
- (K) An estimator is consistent if $T_n \xrightarrow{\text{P}} \theta$
- (L) p-value: $p(x_1, \dots, x_n) = \sup_{\theta \in \Theta_0} P_\theta(w(X) \geq w(x))$. Max probability of seeing evidence against the null hypothesis when θ_0 is actually true.
- (M) Type I: $P(\text{reject } \theta_0 \text{ when } \theta_0 \text{ true})$
Type II: $P(\text{don't reject } \theta_0 \text{ when } \theta_0 \text{ false})$

Rejection region $R_\alpha = \{x : w(x) > c\}$ where c chosen to achieve $\text{s.e. } \alpha$.
Reject H_0 if $x \in R_\alpha$.

$$\sup_{\theta \in \Theta_0} P_\theta(w(x) > c) = \sup_{\theta \in \Theta} P_\theta(x \in R_\alpha) \leq \alpha$$

Power $\beta(\theta) = P_\theta(H_0 \text{ rejected})$

$$\beta = \sup_{\theta \in \Theta} \beta(\theta)$$

Result check:

$$\forall t \in (-c, c)$$

- If $\mathbb{E}[e^{tx}] < \infty$, then MGF of X $m_x(t) = \mathbb{E}[e^{tx}] = 1 + t\mathbb{E}[x] + \frac{1}{2}t^2\mathbb{E}[x^2] + \dots$
 $\forall t \in (-c, c)$ where $c > 0$ some constant.
- $\int_x^\infty \lambda e^{-\lambda x} dx = [-e^{-\lambda x}]_{-\infty}^\infty = \infty e^{-\lambda x} \leftarrow \text{need}$

List of Sneaky Tricks

- X_1, \dots, X_n w/ $X_i \sim \text{Exponential}(\lambda)$
 - $\sum_{i=1}^n X_i \sim \text{Gamma}(n, \lambda)$ [Sum of n iid exp is gamma]
 - $P(X > y) = e^{-\lambda y}$ which is easier to work w/!
 - $E\left[\frac{1}{\sum_{i=1}^n X_i}\right] = E\left[\frac{1}{Y}\right] = \int_0^\infty \frac{1}{y} \frac{1}{\Gamma(n)} \lambda^n y^{n-1} e^{-\lambda y} dy$
 $Y \sim \text{Gamma}(n, \lambda)$
 - Memorizes $P(X > u+t | X > u) = P(X > t)$ $\forall u, t > 0$
 - $\text{Gamma}(n, \frac{1}{2}) \sim \chi_n^2$
- Uniform distribution on a weird area: $\frac{1}{\text{Area}(A)} \int_A \prod_{(x,y) \in A} dx dy$
 ↳ use symmetry $\text{Corr}(X, Y) = \text{Corr}(X, -Y)$
- Transformations then requires STRICTLY T functions
- State WHEN you use a named theorem.
 - ↳ sum iid Random Variables \Rightarrow CLT applies
 - ↳ exp cts \Rightarrow apply CMT
- Fisher's Thm for sequence X_1, \dots, X_n but μ, σ^2 unknown.
 - ↳ $\frac{1}{\sigma^2} \sum_{i=1}^n (X_i - \bar{X})^2 \sim \chi_{n-1}^2$ so $Z \sim N(0, 1)$ R.V.
 - so $E\left[\sum_{i=1}^n (X_i - \bar{X})^2\right] = \sigma^2 E\left[Z^2 + \dots + Z_{n-1}^2\right] = (n-1)\sigma^2$
- Be VERY careful!!! NOT $X_1 \sim \text{Exp}(\lambda), X_2 \sim \text{Exp}(\lambda)$
- Justify MLE makes.
- NOTE: SHIFTED DISTRIBUTIONS $\Rightarrow X = Z + \gamma$ (not odd factors)
- DEFINE ALL YOUR VARIABLES. (leave nothing unexplained.)
- CLT WLLN
 - iid
 - common MGF
 - iid
 - finite μ, σ^2
- Say test statistic $u(x) \downarrow$, defined on $x > 0$, then will be biggest on $[0, c]$, so that is the rejection region.
 ↳ use \uparrow / \downarrow of test statistic to determine shape of R_α

- $X \sim N(0, 1) \quad Y = X^2$
- $\hookrightarrow F_Y(y) = P(Y \leq y) = P(X^2 \leq y) = \int_{-\infty}^{\infty} \mathbb{I}\{x^2 \leq y\} dx = \int_{-\sqrt{y}}^{\sqrt{y}} f_X(x) dx$
- $\hookrightarrow f_Y(y) = F'_Y(y)$
- $Z \in \mathbb{R}^n \quad \mathbb{E}[Z] = \mu$
- $V_{ij} = \mathbb{E}[(Z_i - \mathbb{E}[Z_i])(Z_j - \mathbb{E}[Z_j])]$
- $= \mathbb{E}[Z_i Z_j] - \mu_i \mu_j \Rightarrow V = ZZ^T + \mu \mu^T$
- conditional expectations:
 - Take out what's known
 - Tower property (can sneak an extra $\mathbb{E}[\cdot | X]$ in!)
- Sum n iid R.V. \Rightarrow apply CLT \Rightarrow $\log N$.
- Several Normals $\rightarrow MVN(\mu, V)$ where $V \in \mathbb{R}^{n \times n}$
 - \hookrightarrow Specifying rows/columns of matrix carefully. $AA^T = V$
 - \hookrightarrow watch out for X_{n+2}^2 at any point \hookrightarrow ALL RELATED!
- $Y_i = \gamma + \beta x_i + \varepsilon_i \Rightarrow Y_i \sim N(\gamma + \beta x_i, \sigma^2) \rightarrow$ PDF!
- Lagrange multipliers...
- For $\lambda \in \mathbb{I}$
 - $\hookrightarrow \mu, \sigma^2$ unknowns \Rightarrow Fishers, t_{n-1} create it!
 - $\hookrightarrow \sigma^2$ known \Rightarrow simpler χ^2_{n-1} .
- Likelihood ratio statistic \Rightarrow PI for values of parameters in diff scenarios. E.g. $\mu = \frac{n\bar{x} + m\bar{y}}{n+m}$
- NLLN proof: $H \varepsilon > 0$

$$P(|\bar{X}_n - \mu| > \varepsilon) = P((\bar{X}_n - \mu)^2 > \varepsilon^2) \leq \frac{\text{Var}(\bar{X}_n - \mu)}{\varepsilon^2} = \frac{n\sigma^2}{n^2\varepsilon^2} \xrightarrow{n \rightarrow \infty} 0$$
- Choose questions CAREFULLY!
 - \hookrightarrow make sure you can do all parts/ have ideas for how to attack.
- If you can't JV a logical step, LEAVE IT & come back.
 - \hookrightarrow all double, just can't go wrong.



USE 11 questions
to identify ranges

$$\boxed{+} + \boxed{-} = \boxed{\square}$$

$$\int_{-\infty}^{\infty} \mathbb{I}\{x^2 \leq y\} dx = \int_{-\sqrt{y}}^{\sqrt{y}} f_X(x) dx$$