
Probability



Probability Axiom

Probability Space 由那 3 个元素组成

1. sample space Ω that contains all possible outcomes
2. σ -algebra F , a collection of subset of Ω
3. probability measure P , a function that assign non-negative value on every set of F

σ -field: F 是指 σ -field 为 collection of subset of Ω

1. $\emptyset \in F$
2. $A \in F \Rightarrow A^c \in F$
3. $A_i \in F \Rightarrow \bigcup_{i=1}^{\infty} A_i \in F$
 - A : F -measurable set / event
 - $\sigma(\omega)$: 包含 ω 的 F 的 σ -field.

probability measure: $P: F \rightarrow [0, 1]$ 满足 axiom 为 function

1. $P(\Omega) = 1$
2. $P(\emptyset) = 0$
3. $P(\bigcup A_i) = \sum P(A_i)$ countable additivity

Theorem: continuity of probability measure.

- \Updownarrow
- ① P is a probability measure
 - ② $A_i < A_{i+1}$ for all i and $A = \bigcup A_i \Rightarrow \lim_{n \rightarrow \infty} P(A_i) = P(A)$
 - ③ $A_i > A_{i+1}$ for all i and $A = \bigcap A_i \Rightarrow \lim_{n \rightarrow \infty} P(A_i) = P(A)$

Measurable Function

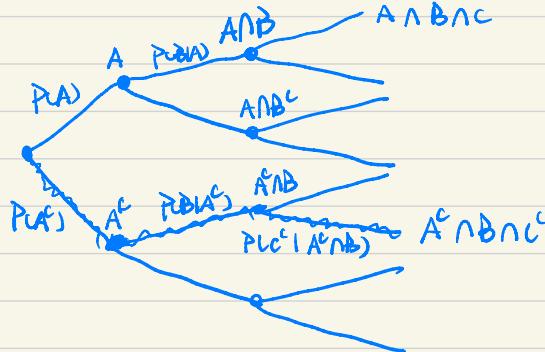
(Ω_1, F_1) & (Ω_2, F_2) be 2 measurable space. A function $f: \Omega_1 \rightarrow \Omega_2$ is called (F_1, F_2) measurable if $f^{-1}(B) \in F_1$ for every $B \in F_2$

Conditional probability

- Condition lead to revised (conditional) probabilities that take into account partial information.
- Useful in 'divide and conquer' for complex problem

1. Multiplication Rule

$$P(A_1 \cap A_2 \dots \cap A_n) = P(A_1) \prod_{i=2}^n P(A_i | A_1 \cap \dots \cap A_{i-1})$$



2. Bayes Rule

$$P(A_i | B) = \frac{P(A_i) P(A_i | B)}{\sum_j P(A_j) P(A_j | B)}$$

3. Independence

Events A_1, \dots, A_m are independent if

$$P(A_1 \cap A_2 \dots \cap A_m) = P(A_1) P(A_2) \dots P(A_m) \text{ for any distinct indices } i, j \neq k$$

pairwise independence

- Conditional PMF

$$P_{X|Y}(x|y) = \frac{P_{XY}(x,y)}{P_Y(y)} \text{ defined for } y \text{ such that } P_Y(y) > 0$$

- Conditional PDF

$$f_{X|Y}(x|y) = \frac{f_{XY}(x,y)}{f_{Y|y}} \text{ if } f_{Y|y} > 0$$

- Mixed

$$\checkmark P_{k|Y}(k|y) = \frac{P_k(k) f_{Y|k}(y|k)}{f_{Y|y}}$$

$$\checkmark f_{Y|k}(y|k) = \frac{f_k(k) P_{k|Y}(k|y)}{P_k(k)}$$

- Multiplication

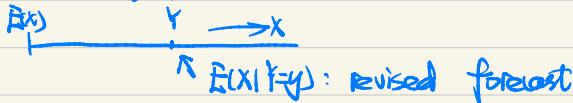
$$f_{X,Y,Z}(x,y,z) = f_Z(z) f_{Y|Z}(y|z) f_{X|Y,Z}(x|y,z)$$

- Moments

$$E(X|Y) = g(Y) \text{ a function of } Y$$

$$g(y) = E[X|Y=y] \Rightarrow E[E(X|Y)] = E(X)$$

i.e. Forecast Revision



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

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

- ① average variability within sectors
- ② variability between sectors

Counting Problem

- Permutation : # of ways of ordering n elements

$$n! = n(n-1)\cdots 1$$

- # of subsets of $\{1, \dots, n\}$

$$2^n$$

- # of k -element subsets of a given n -element set

$$\binom{n}{k} = \frac{n!}{k!(n-k)!} \quad \text{Binomial}$$

$$\sum_{k=0}^n \binom{n}{k} = \binom{n}{0} + \binom{n}{1} + \cdots + \binom{n}{n} = \# \text{ of all subsets} = 2^n$$

- Partitions

$$\# \text{ of partitions} = \frac{n!}{n_1! n_2! \cdots n_r!} \quad \text{multi-nomial}$$

Derived Distribution

- Discrete:

$$Y = g(X) \Rightarrow P(Y=y) = P(g(x)=y) \\ = \sum_{x: g(x)=y} P_X(x)$$

i.e. when $y = ax + b$ $P(Y=y) = P_X\left(\frac{y-b}{a}\right)$

- Continuous

$$Y = g(X)$$

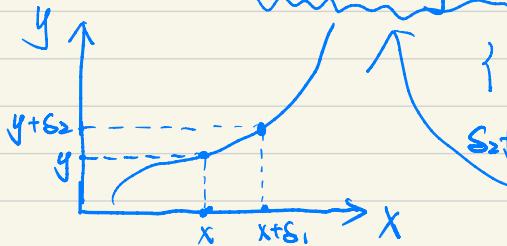
$$1. \text{ CDF: } F_Y(y) = P(Y \leq y) = P(g(x) \leq y)$$

$$2. \text{ Take derivative: } f_Y(y) = \frac{dF_Y}{dy}(y)$$

When g is monotonic, $x_1 > x_2 \Rightarrow g(x_1) > g(x_2)$

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

$$2. f_Y(y) = f_X(g^{-1}(y)) \left| \frac{d g^{-1}(y)}{dy} \right|$$



$$\begin{cases} y = g(x) \rightarrow S_2 \subset S_1, \frac{dg}{dx}(x) \\ x = g^{-1}(y) \rightarrow S_1 \subset S_2, \frac{dx}{dy}(g^{-1}(y)) \end{cases}$$

$$S_2 f_Y(y) \approx P(y \leq Y \leq y+s_2) = P(x \leq X \leq x+s_1)$$

$$\approx S_1 f_X(x) \approx S_2 \frac{dx}{dy}(g^{-1}(y)) f_X(x)$$

Sum of Random Variables

Convolution $Z = X + Y$ (X, Y independent)

1. discrete

$$P_Z(z) = \sum_x P_X(x) P_Y(z-x)$$

2. continuous

$$f_Z(z) = \int_{-\infty}^{\infty} f_X(x) \underbrace{f_Y(z-x)}_{f_{Z|X}(z|x)} dx$$

$$\text{Var}(X_1 + \dots + X_n) = \sum_{i=1}^n \text{Var}(X_i) + \sum_{i \neq j} \text{cov}(X_i, X_j)$$

$$\text{P}(X, Y) = \frac{\text{cov}(X, Y)}{\sigma_X \sigma_Y} = E\left[\frac{(X - \bar{X})}{\sigma_X} \frac{(Y - \bar{Y})}{\sigma_Y}\right]$$

- $P=1 \Rightarrow X - \bar{X} = C(Y - \bar{Y})$

- $\text{P}(ax+b, Y) = \text{sign}(a) P_{XY}$

Bayesian Inference frameworks

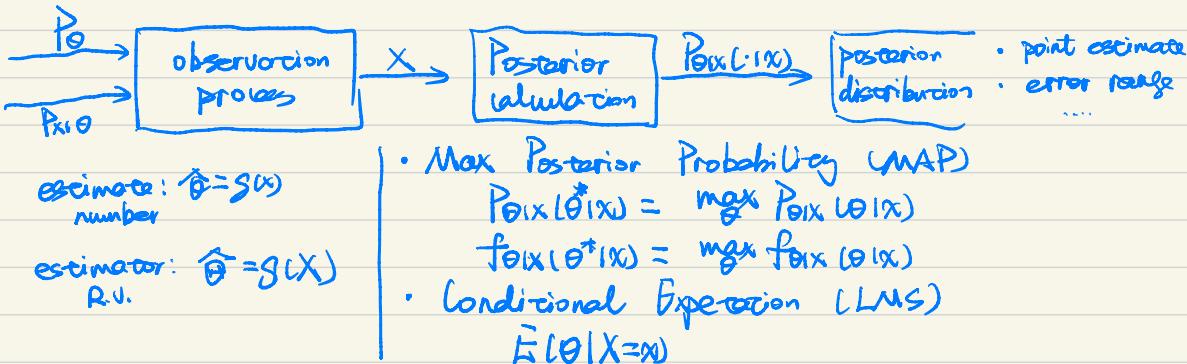
- Unknown Θ

 - treat as random variable

 - prior distribution P_θ or f_θ from
 - ① symmetry
 - ② known range
 - ③ earlier study
 - ④ arbitrary

- Observation X

 - Observation Model $P_{x|\theta}$ or $f_{x|\theta}$



Performance evaluation

1. conditional prob of error

$$P(\hat{\theta} \neq \theta | X=x) \text{ smallest under MAP}$$

2. Overall prob of error

$$\begin{aligned} P(\hat{\theta} \neq \theta) &= \sum P(\hat{\theta} \neq \theta | X=x) P_X(x) \\ &= \sum P(\hat{\theta} \neq \theta | \theta=\theta) P_\theta(\theta) \end{aligned}$$

3. Conditional mean square error

$$E[(\hat{\theta} - \theta)^2 | X=x]$$

4. Mean square error

$$E[(\hat{\theta} - \theta)^2]$$

LMS Estimation

- Without Observation

$$\min_{\hat{\theta}} E[(\theta - \hat{\theta})^2] \rightarrow \text{Mean square error}$$

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

$$\text{Also } E[(\theta - \hat{\theta})^2] = \text{Var}(\hat{\theta}) \leq E[(\theta - c)^2] \text{ for all } c$$

- LMS estimate of θ based on X

① unknown θ : prior $P(\theta)$

② observation X : model $P_{X|\theta}(x|\theta)$

$$\min_{\hat{\theta}} E[(\theta - \hat{\theta})^2 | X=x]$$

$$\Rightarrow \hat{\theta} = E[\theta | X=x]$$

$$\text{Also } E[(\theta - E[\theta | X=x])^2 | X=x] \leq E[(\theta - g(x))^2 | X=x]$$

$$\text{abstract } \downarrow E[(\theta - E[\theta | X])^2 | X] \leq E[(\theta - g(X))^2 | X]$$

$$\text{take expectation } E[(\theta - E[\theta | X])^2] \leq E[(\theta - g(X))^2]$$

$$\hat{\theta}_{\text{LMS}} = E[\theta | X] \text{ minimize } E[(\theta - g(X))^2] \text{ over all estimator } g$$

- Performance

1 LMS estimate: $\hat{\theta} = E[\theta | X=x]$

1 LMS estimator: $\hat{\theta} = E[\theta | X]$

① Expected performance once we have observation

$$MSE = E[(\theta - E[\theta | X=x])^2 | X=x]$$

$$= \text{Var}(\theta | X=x)$$

② Expected performance of the design (over all possible x)

$$MSE = E[(\theta - E[\theta | X])^2]$$

$$= E[\text{Var}(\theta | X)]$$

property

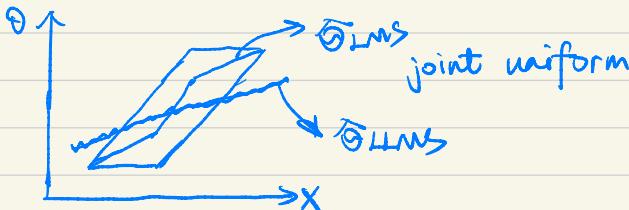
$$\text{① } \text{cov}(\hat{\theta}, \hat{\theta} - \theta) = 0$$

$$\text{② } \text{Var}(\hat{\theta}) = \text{Var}(\hat{\theta}) + \text{Var}(\hat{\theta} - \theta)$$

LLMS estimation

The conditional expectation may be hard to compute. In that case, we can restrict to estimators $\hat{\theta} = ax + b$ that minimize the mean square error.

- Estimator $\hat{\theta} = g(x)$
- $\min_{\theta} E[(\hat{\theta} - \theta)^2] \Rightarrow \hat{\theta}_{LLMS} = E[\theta | x]$
- $\min_{a,b} E[(\theta - ax - b)^2] \Rightarrow \hat{\theta}_{LLMS}$
- If $E[\theta | x]$ is linear in $x \Rightarrow \hat{\theta}_{LLMS} = \hat{\theta}_{LMS}$



$$\begin{aligned}\hat{\theta}_{LLMS} &= E[\theta] + \frac{\text{cov}(\theta, x)}{\text{var}(x)} (x - E(x)) \\ &= E[\theta] + P \frac{\partial \theta}{\partial x} (x - E(x))\end{aligned}$$

- $P=0 \Rightarrow \hat{\theta}_{LLMS} = \bar{\theta}$
- $E[(\hat{\theta}_{LLMS} - \theta)^2] = (1-P^2) \text{var}(\theta)$

Multiple Observation Case :

$$\min_{a, b} \hat{\theta}[(a_1 x_1 + \dots + a_n x_n + b - \theta)^2]$$

Inequality

- bound $P(X \geq a)$ based on limited information of the Distribution
- Markov (based on mean)
- Chebyshew (based on mean and variance)

Markov Inequality

If $X \geq 0$ and $a > 0$, $P(X \geq a) \leq \frac{E(X)}{a}$

(If $E(X)$ small, X is unlikely to be large)

Chebyshew Inequality

X with finite mean and variance. $P(|X - \mu| \geq c) \leq \frac{\sigma^2}{c^2}$

(If variance small, X unlikely to be far away from mean)

Jensen's Inequality

g is convex $\rightarrow g(E(X), E(Y)) \leq E[g(X, Y)]$

R.U. convergence

WLLN:

X_1, X_2, \dots iid; finite mean and variance σ^2

• sample mean: $M_n = \frac{X_1 + \dots + X_n}{n}$

• $E(M_n) = \mu$.

• $\text{Var}(M_n) = \frac{\sigma^2}{n}$

$$P(|M_n - \mu| \geq \epsilon) \leq \frac{\text{Var}(M_n)}{\epsilon^2} = \frac{\sigma^2}{n\epsilon^2} \xrightarrow{n \rightarrow \infty} 0 \quad (\text{fix } \epsilon)$$

WLLN: $P(|M_n - \mu| \geq \epsilon)$

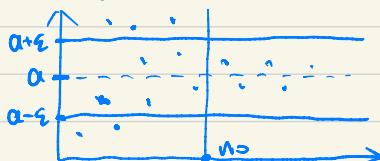
$$= P\left(\left|\frac{X_1 + \dots + X_n}{n} - \mu\right| \geq \epsilon\right) \xrightarrow{n \rightarrow \infty} 0$$

Interpretation: many observations with noise, the sample mean M_n is unlikely to be far from the true μ .

Convergence in probability

ordinary convergence

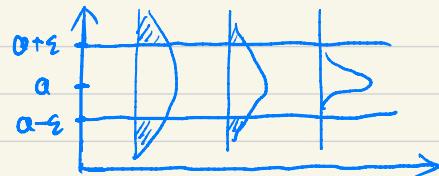
- sequence of number $a_n \rightarrow a$
 { a_n eventually gets and
 stay close to a }



- For every $\epsilon > 0$, there exists n_0 such that for every $n > n_0$
 $|a_n - a| < \epsilon$

convergence in probability

- sequence of R.V. $Y_n \rightarrow a$
 { almost all of the pdf/pmf
 of Y_n eventually get concentrated
 close to a }



for any $\epsilon > 0$ $P(|Y_n - a| \geq \epsilon) \rightarrow 0$

does not mean convergence in
 expectation (i.e.)

$$Y_n = \begin{cases} n^2, & \frac{1}{n} \\ 0, & 1 - \frac{1}{n} \end{cases} \quad \begin{matrix} Y_n \rightarrow 0 \\ E(Y_n) \rightarrow \infty \end{matrix}$$

Bernoulli Process

A sequence of Bernoulli trials

$$\begin{cases} P(X_i=1) = p \\ P(X_i=0) = 1-p \end{cases}$$

independence + Time-homogeneity

number of success/arrival

- $S = X_1 + \dots + X_n$
- $P(S=k) = \binom{n}{k} p^k (1-p)^{n-k}$ for $k=0, \dots, n$
- $E(S) = np$
- $\text{Var}(S) = np(1-p)$

Time until first arrival

- $T_1 = \min \{ i : X_i = 1 \}$
- $P(T_1=k) = (1-p)^{k-1} p$ for $k=1, \dots$
- $E(T_1) = 1/p$
- $\text{Var}(T_1) = 1/p^2$

Memoryless

- Fresh start after a random time N
 N : a causally determined (without future information)
The process X_{N+1}, X_{N+2}, \dots is
 - ① Bernoulli Process
 - ② independent of N, X_1, \dots, X_N

Time of the k-th arrival/success

T_k : time of k-th arrival

T_k : k-th inter-arrival time = $T_k - T_{k-1}$ ($k \geq 2$)

$T_K = T_1 + \dots + T_K$

T_2 is independent of T_1 \cap $\{X_1=0\}$

$$Y_k = T_1 + \dots + T_k$$

T_i are iid Geometric(p)

$$\bar{E}(Y_k) = k\bar{E}(T_i) = \frac{k}{n} \quad \text{Var}(Y_k) = \frac{ku-p}{p^2}$$

$$P_{Y_k}(t) = \binom{t-1}{k-1} p^k u^{-p} t^{-k} \quad t=k, k+1, \dots$$

Merging / Split

- $X_t \sim \text{Poisson}(p)$ $Y_t \sim \text{Poisson}(q)$

$$Z_t = g(X_t, Y_t)$$

- $X_t \sim \text{Poisson}(p)$
 - $Z_t \sim \text{Poisson}(pq)$
 - $Y_t \sim \text{Poisson}(pu-q)$

Poisson approximation to binomial

- when n large, p small, $\lambda = np$ moderate

- Number of arrivals S in n slots:

$$P_S(k) = \binom{n}{k} \left(\frac{\lambda}{n}\right)^k \left(1 - \frac{\lambda}{n}\right)^{n-k}$$

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

$$\xrightarrow{n \rightarrow \infty} 1 \cdot 1 \cdot \dots \cdot 1 \cdot \frac{\lambda^k}{k!} e^{-\lambda}$$

Poisson Process

- Prob of k arrivals in interval of duration S

$$P(k, S) = \begin{cases} 1 - \lambda S + o(S^2) & \text{if } k=0 \\ \lambda S + o(S^2) & \text{if } k=1 \\ o(S^2) & \text{if } k>1 \end{cases} \quad \lambda: \text{arrival rate}$$

- N_t : arrival in $[0, t]$

$$N_t = T/S = \text{interval} / \text{length of slot}$$

$N_t \sim \text{binomial}$

$$\rightarrow P = \lambda S + o(S^2)$$

$$np = \lambda T + o(S) \approx \lambda T$$

- $N_t \sim \text{Binomial}(n, p)$

$$n = T/S, \quad P = \lambda S + o(S^2)$$

$$P(k, t) = P(N_t = k) = \frac{(\lambda t)^k}{k!} e^{-\lambda t}$$

$$E(N_t) = \lambda t$$

$$\text{Var}(N_t) = \lambda t$$

T_1 : time until first arrival

$$P(T_1 < t) = 1 - P(T_1 \geq t) = 1 - P(0, t) = 1 - e^{-\lambda t}$$

$T_1 \sim \text{exponential}$

Y_k : time of the k -th arrival

$$P(Y_k \leq y) = \sum_{n=k}^{\infty} P(n, y)$$

- $f_{Y_k}(y) \leq n \quad P(y \leq Y_k \leq y+S)$

scenario: $\underbrace{k-1}_{\substack{y \\ \uparrow \\ k-1 \text{th}}} \quad \underbrace{k \text{-th}}_{\substack{y \\ \uparrow \\ k \text{th}}} \quad \underbrace{y+S}_{\substack{y \\ \uparrow \\ k+1 \text{th}}} \quad \underbrace{P(k-1, y) \lambda S}_{\substack{\text{kth in } [y, y+S]}}$

scenario: $k-2$ $\underbrace{\quad}_{\substack{y \\ \uparrow \\ k-2 \text{th}}} \quad \underbrace{k \text{-th}}_{\substack{y \\ \uparrow \\ k \text{th}}} \quad \underbrace{y+S}_{\substack{y \\ \uparrow \\ k+1 \text{th}}} \quad + P(k-2, y) o(S^2) \quad k-1, k \text{ in } [y, y+S]$
 $+ P(k-3, y) o(S^3) + \dots$

$$\xrightarrow{S \rightarrow \infty} \approx P(k-1, y) \lambda S$$

Explain: $f_{Y_k}(y) = \lambda P(k-1, y) = \lambda \frac{(xy)^{k-1} e^{-\lambda y}}{(k-1)!}$

$Y_k = T_1 + \dots + T_k \rightarrow$ a equivalent definition

- sum of iid exponentials
- $E(Y_k) = k/\lambda$
- $\text{Var}(Y_k) = k/\lambda^2$

Poisson v.s. Bernoulli

$$\lambda = \tau/\delta$$

$$P = \lambda \delta$$

$$np = \lambda \tau$$

	Poisson	Bernoulli
Time of arrival	continuous	discrete
Arrival rate	$\lambda/\text{unit time}$	$p/\text{per trial}$
PMF of # of arrivals	Poisson	Binomial
Interarrival time	Exponential	Geometric
Time to k -th arrival	Erlang	Pascal

Merge of 2 independent Poisson Proces

$M: \text{Poisson } (\mu) \rightarrow M+N: \text{Poisson } (\mu+\nu)$

$N: \text{Poisson } (\nu)$

$$\cdot P(\text{arrival is } m \mid \text{arrival at time } t) = \frac{\mu}{\mu+\nu}$$

$$\cdot P(k \text{ out of first } n \text{ arrival is } m) = \binom{n}{k} \left(\frac{\mu}{\mu+\nu} \right)^k \left(\frac{\nu}{\mu+\nu} \right)^{n-k}$$

Time of first / last arrivals of 3 independent Poissons

first arrival : $\min(X, Y, Z) \sim \text{Exp}(3\lambda)$

last arrival : $\max(X, Y, Z)$

$$E[\max(X, Y, Z)] = \frac{1}{3\lambda} + \frac{1}{2\lambda} + \frac{1}{\lambda}$$

Split

$\text{Poisson}(\lambda u) \leftarrow \begin{matrix} \text{Poisson}(\lambda u) \\ \text{Poisson}(\lambda(u-v)) \end{matrix}$ independent

Random Incidence

more likely to show up at a longer length of interval

$\text{Poisson}(\lambda) \quad \xrightarrow{\quad u \quad t^* \quad v \quad} \quad \rightarrow$

- u : last arrival
- v : next arrival
- t^* : arrival time

$$v-u = \frac{v-t^*}{\text{Exp}(\lambda)} + \frac{t^*-u}{\text{Exp}(\lambda)}$$

$$E(v-u) = E(v-t^*) + E(t^*-u) = 2/\lambda$$

sampling method matters

Markov Processes

$$\text{State}(t+1) = f(\text{State}(t), \text{noise})$$

- X_n : state after n transitions

- belong to a finite set

- initial state X_0 either given or random

- homogeneous transition probability

$$P_{ij} = P(X_{n+1}=j | X_n=i) = P(X_1=j | X_0=i)$$

- key assumption

Given current state, the past does not matter

$$P_{ij} = P(X_{n+1}=j | X_n=i)$$

$$= P(X_{n+1}=j | X_n=i, X_{n-1}, \dots, X_0)$$

- N -Step transition

$$P_{ij}(n) = P(X_n=j | X_0=i)$$

$$= P(X_{n+1}=j | X_0=i)$$

$$\text{Recursion: } P_{ij}(n) = \sum_k P_{ik}(n-1) P_{kj}$$

$$= \sum_k P_{ik} P_{kj}(n-1)$$

- Recurrent / transient

- State i is recurrent if starting from i , wherever you go, there is a way back to i

- not recurrent \rightarrow transient

- recurrent class: a collection of recurrent states communicating only between each other

- Steady state probability

- $\lim_{n \rightarrow \infty} P_{ij}(n) = \pi_j$ converge if

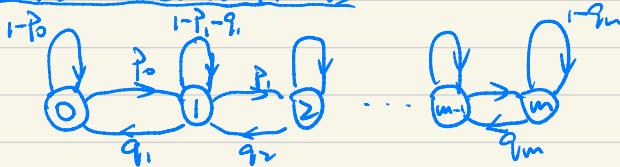
1. recurrent states are in one class

2. single recurrent class is not periodic

- calculate by solving linear system

$$\left\{ \begin{array}{l} \pi_j = \sum_k \pi_k P_{kj} \rightarrow \text{balance equation} \\ \sum \pi_j = 1 \end{array} \right.$$

Birth - Death Process



- $\pi_0 p_i = \pi_{i+1} q_i \Rightarrow \pi_0 + \pi_0 \frac{p_0}{q_1} + \pi_0 \frac{p_0 p_1}{q_1 q_2} + \dots = 1$
- $\sum \pi_i = 1$

$\Downarrow p_i, q_i$ fixed

- $\pi_0 = \pi_0 \left(\frac{p_0}{q_1}\right)^i$
- $p = q$ if $\Rightarrow \pi_i = \pi_0$
- $P < q$ if $\Rightarrow \pi_0 = 1 - \frac{P}{q}$ $E(X_n) = \frac{P/q}{1-P/q}$

Absorption

absorbing state: recurrent state/class k with $P_{kk} = 1$

• probability a_i that eventually settle if it start from i

$$a_j = \sum_{j=1}^m P_{ij} a_j$$

• Expected time to absorption

$$u_i = 1 + \sum_j P_{ij} u_j$$

Statistik



Introduction

Statistic: use data to gather insight and make decisions

- computational view

Data (sequence of numbers)



Algorithm (application of statistical principle)

sometimes will skip insights

- statistical view

Data: comes from a random process

Goal: learn how the process work and further predict

find determinant

Statistical Modeling

complicated process

""

simple process

+
random noise

) ① real randomness

② deterministic but too complex

Modeling consists in choosing) ③ plausible simple process
noise distribution

noise is what we do not understand. It is a way of
modeling lack of information

Probability Tool

- Law of large number (LLN)

When size n of the experiment becomes larger, sample mean is a good (consistent) estimator of expectation

- Central limit Theorem (CLT)

How good the estimator is

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

Convergence

Random Variable Convergence

- **Almost Sure** $T_n \xrightarrow{\text{a.s.}} T$ iff

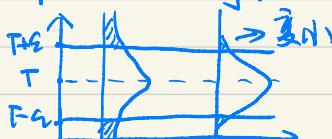
$$P[\{w : T_n(w) \xrightarrow{\text{a.s.}} T(w)\}] = 1$$

意味着 # of $|T_n - T| > \varepsilon$ 是 finite P.s., 除一些 prob 0 P.s. 外所有事件上都成立

converse, 也是 limit of number sequence

- **In Probability** $T_n \xrightarrow{\text{i.p.}} T$ iff

$$P[|T_n - T| > \varepsilon] \xrightarrow{n \rightarrow \infty} 0$$



- **In distribution** $T_n \xrightarrow{\text{d.}} T$ iff

$$\uparrow E[f(T_n)] \xrightarrow{n \rightarrow \infty} E[f(T)] \text{ for all continuous and bounded } f$$

$$\downarrow \lim_{n \rightarrow \infty} F_{T_n}(x) \rightarrow F_T(x) \text{ pointwise for all } x$$

Convergence operation

$$T_n \xrightarrow{\text{a.s./P}} T \quad U_n \xrightarrow{\text{a.s./P}} U$$

- $T_n + U_n \xrightarrow{\text{a.s./P}} T + U$
- $T_n U_n \xrightarrow{\text{a.s./P}} TU$
- if $U \neq 0$, $T_n/U_n \xrightarrow{\text{a.s./P}} T/U$
- In general does not apply to (d)

Suskey: partial operation hold for (d)

$$T_n \xrightarrow{\text{d.}} T \quad U_n \xrightarrow{\text{P}} u \text{ (deterministic)}$$

- $T_n + U_n = T + u$
- $T_n U_n = uT$
- If $u \neq 0$, $T_n/U_n = T/u$

Continuous Mapping

For continuous function f : $T_n \xrightarrow{\text{a.s./P/d.}} T \Rightarrow f(T_n) \xrightarrow{\text{a.s./P/d.}} f(T)$

Delta Method

| Let \hat{z}_n be asymptotical normal $\sqrt{n}(\hat{z}_n - \theta) \xrightarrow{d} N(0, \sigma^2)$

| Let g be continuous differentiable at the point θ

$g(\hat{z}_n)$ is also asymptotic normal:

$$\sqrt{n}(g(\hat{z}_n) - g(\theta)) \xrightarrow{d} N(0, g'(\theta)^2 \sigma^2)$$

Taylor Expansion approximation:

$$g(\hat{z}_n) - g(\theta) = (\hat{z}_n - \theta)g'(\theta) + \frac{1}{2}(\hat{z}_n - \theta)^2 g''(w) \quad w \in [\hat{z}_n, \theta]$$
$$\cong (\hat{z}_n - \theta)g'(\theta)$$

Multivariate:

$$\sqrt{n}(T_n - T) \rightarrow N(0, \Sigma) \quad \text{then}$$

$$\sqrt{n}(f(T_n) - f(T)) \rightarrow N(0, \nabla f(T) \Sigma \nabla f(T))$$

Statistical Inference

Goal: Use reasonable assumptions to find tractable models

1. estimation \rightarrow single number
2. confidence interval \rightarrow error bar around the number
3. Hypothesis test \rightarrow answer to yes/no question

Statistical Model

Let the observed outcome of a statistical experiment be a sample x_1, \dots, x_n of n i.i.d random variables in some measurable space E and common distribution P

$$(E, (P_\theta)_{\theta \in \Theta})$$

- E : sample space
- P : probability measure on E
- Θ : parameter set
 - parametrical: can be defined with finite unknowns
 - Nonparametrical: otherwise

Identifiability

Parameter θ is identified iff

$$\begin{aligned}\theta \neq \theta' &\Rightarrow P_\theta \neq P_{\theta'} \\ P_\theta = P_{\theta'} &\Rightarrow \theta = \theta'\end{aligned}\quad \Updownarrow$$

Parameter Estimation

Asymptotic normal: $\sqrt{n}(\hat{\theta}_n - \theta) \xrightarrow{n \rightarrow \infty} N(0, \sigma^2)$
 σ^2 is the asymptotic variance

Quadratic risk: $R(\hat{\theta}_n) = E[(\hat{\theta}_n - \theta)^2]$

$$R(\hat{\theta}_n) = E[(\hat{\theta}_n - E[\hat{\theta}_n])^2] + E[E(\hat{\theta}_n) - \theta]^2$$
$$\text{var}(\hat{\theta}_n) + \text{bias}^2$$

Confidence Interval

- Confidence Interval of level $1-\alpha$ for θ

Any random interval I (depend on X_1, \dots, X_n) whose boundary do not depend on θ such that

$$P_{\theta}[I \ni \theta] \geq 1-\alpha, \quad \forall \theta \in \Theta$$

- Asymptotic Variance

$$\lim_{n \rightarrow \infty} P_{\theta}[I \ni \theta] \geq 1-\alpha, \quad \forall \theta \in \Theta$$

Get Confidence Interval

- ① Find a pivotal statistic (a statistic without unknown)
 - Non-asymptotic CI
 - CLT / Delta Method
 - Asymptotic CI

- ② Solve for CI

- Two side symmetrical

$$P(\theta \in [\hat{\theta} - s, \hat{\theta} + s]) = 1-\alpha$$

$$\Rightarrow I = \hat{\theta} \pm q_{\alpha/2} \sigma_{\hat{\theta}} / \sqrt{n}$$

- One side

$$P(\theta \leq \hat{\theta} + s)$$

$$\Rightarrow I = (-\infty, \hat{\theta} + q_{\alpha} \sigma_{\hat{\theta}} / \sqrt{n})$$

Interpretation : probability of containing the true parameter
if you do many experiments.

Once plug in observations. It either contains the true parameter (100%) or not (0%)

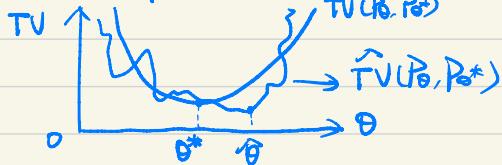
Estimation

Goal: Given x_1, \dots, x_n , find an estimator $\hat{\theta} = \hat{\theta}(x_1, \dots, x_n)$ such that $P_{\hat{\theta}}$ is close to P_{θ^*}
 $\Leftrightarrow |P_{\hat{\theta}}(A) - P_{\theta^*}(A)|$ is small for all $A \in \mathcal{E}$

Total Variation Distance

- $TV(P_{\theta}, P_{\theta'}) = \max_{A \in \mathcal{E}} |P_{\theta}(A) - P_{\theta'}(A)|$
 $= \frac{1}{2} \sum_{x \in E} |P_{\theta}(x) - P_{\theta'}(x)|$
 $= \frac{1}{2} \int_E |f_{\theta}(x) - f_{\theta'}(x)| dx$

- estimation strategy: Build an estimator $\hat{TV}(P_{\theta}, P_{\theta^*})$ for all θ .
Then find $\hat{\theta}$ that minimize the function $\theta \mapsto \hat{TV}(P_{\theta}, P_{\theta^*})$



- How to estimate $TV(P_{\theta}, P_{\theta^*})$. Nice if TV is expectation so that we can use LLN. TV is not KL

\leftarrow

Kullback-Leibler (KL) divergence

$$KL(P_{\theta}, P_{\theta^*}) = \begin{cases} \sum_{x \in E} P_{\theta}(x) \log \left(\frac{P_{\theta}(x)}{P_{\theta^*}(x)} \right) & \text{if discrete} \\ \int_E f_{\theta}(x) \log \left(\frac{f_{\theta}(x)}{f_{\theta^*}(x)} \right) dx & \text{if continuous} \end{cases}$$
$$= \mathbb{E}_{\theta} [\log(P_{\theta}(x)/P_{\theta^*}(x))]$$

Maximum Likelihood Estimation

KL Divergence to ML estimator

$$KL(P_{\theta^*}, P_{\theta}) = \mathbb{E}_{\theta^*} [\log \frac{P_{\theta^*}(x)}{P_{\theta}(x)}] = \mathbb{E}_{\theta^*} [\log P_{\theta^*}(x)] - \mathbb{E}_{\theta^*} [\log P_{\theta}(x)]$$

↓
constant

$$\hat{KL}(P_{\theta^*}, P_{\theta}) = \text{"constant"} - \frac{1}{n} \sum_{i=1}^n \log P_{\theta}(x_i)$$

$$\downarrow \min_{\theta \in \Theta} \hat{KL}(P_{\theta^*}, P_{\theta}) \iff \max_{\theta \in \Theta} \underbrace{\prod_{i=1}^n P_{\theta}(x_i)}$$

Likelihood

$$L_n = \begin{cases} L(x_1, \dots, x_n; \theta) & \rightarrow P_{\theta}[X_1=x_1, \dots, X_n=x_n] \quad \text{for discrete} \\ (x_1, \dots, x_n; \theta) & \rightarrow \prod_{i=1}^n f_{\theta}(x_i) \quad \text{for continuous} \end{cases}$$

ML Estimator

$$\hat{\theta}_n^{ML} = \underset{\theta \in \Theta}{\operatorname{argmax}} L(x_1, \dots, x_n; \theta)$$

Property under mild condition (continuous, identifiable, etc.)

① Consistency

$$\hat{\theta}_n^{ML} \xrightarrow[n \rightarrow \infty]{\text{i.p.}} \theta^*$$

↑

(X axis converge)
↑ identifiable, etc..

$$\frac{1}{n} L(x_1, \dots, x_n; \theta) \xrightarrow[n \rightarrow \infty]{\text{i.p.}} KL(P_{\theta^*}, P_{\theta}) \quad (\text{y axis converge})$$

② Asymptotic normality

$$\sqrt{n} (\hat{\theta}_n^{ML} - \theta^*) \xrightarrow[n \rightarrow \infty]{\text{i.d.}} N(0, I(\theta^*)^{-1})$$

• Fisher Information

Define $l(\theta) = \log L(x, \theta)$ for one observation

$$I(\theta) = \mathbb{E}[\nabla l(\theta) \nabla l(\theta)^T] - \mathbb{E}[\nabla l(\theta)] \mathbb{E}[\nabla l(\theta)]^T = -\mathbb{E}[\nabla^2 l(\theta)]$$

variance of gradient

$$\text{cov}(\nabla l(\theta)) , \quad l(\theta) = \ln f_{\theta}(x)$$

EM Algorithm

Iterative process to find ML estimator in presence of latent variable

E: compute lower bound of log likelihood for current parameter

M: maximize the lower bound to get new parameter

Goal: Maximize likelihood

$$P(X|\theta) = \sum_z P(x, z|\theta)$$

X: observed variable

Z: latent variable

$$\ln P(X|\theta) = \ln q(\theta) + KL(q, p) \Leftarrow \ln P(X, Z|\theta) = \ln P(Z|X, \theta) + \ln P(X|\theta)$$
$$= \sum_z q(z) \ln \left(\frac{P(x, z|\theta)}{q(z)} \right) + \left[-\sum_z q(z) \ln \left(\frac{P(z|x, \theta)}{q(z)} \right) \right]$$

[for any $q(z)$]

$$\Downarrow KL \geq 0$$

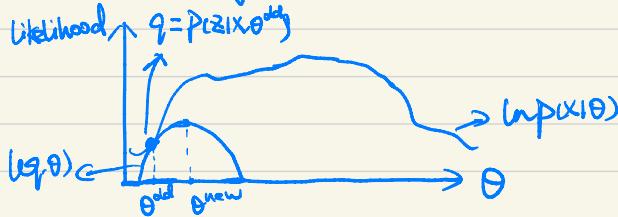
$\ln q(\theta) = \ln P(X|\theta)$ is a lower bound

Assumption

Optimization of $P(X|\theta)$ is difficult but optimization of $P(X, Z|\theta)$ is easier.

- ① E-step: maximizing $\ln q(\theta)$ holding θ^{old} constant
Since $\ln P(X|\theta)$ not dependent on $q(z)$
 $\max \ln q(\theta) \Leftrightarrow \min KL(q, p)$
 $\hat{q} = P(z|x, \theta^{\text{old}})$

- ② M-step: maximizing $\ln q(\theta)$ holding q , constant
Maximizing with complete data $\Rightarrow \theta^{\text{new}}$



$\ln q(\theta) \leq \ln P(x|\theta)$ to θ^{old} [by convex to $\frac{q}{\theta}$]

$$\ln P(x|\theta) = \ln \sum_z P(x, z|\theta) = \ln \sum_z q(z) \frac{P(x, z|\theta)}{q(z)} \geq \sum_z q(z) \ln \frac{P(x, z|\theta)}{q(z)}$$

(Jensen Inequality)

EM for Gaussian Mixture

1. Initialize μ_k , Σ_k and mix coefficient π_k

2. E-step

evaluate $P(z|x, \theta^{old}) = P_{j(i)}$ point i is generated by cluster j

$$P_{j(i)} = \frac{\pi_j N(x_i | \mu_j, \Sigma_j)}{\sum_{m=1}^k \pi_m N(x_i | \mu_m, \Sigma_m)}$$

M-step

$$\hat{\mu}_j = \frac{\sum_i^n P_{j(i)} x_i}{\sum_i^n P_{j(i)}}$$

$$\hat{\pi}_j = \frac{1}{n} \sum_i^n P_{j(i)}$$

$$\hat{\Sigma}_j = \frac{\sum_i^n P_{j(i)} (x_i - \hat{\mu}_j)(x_i - \hat{\mu}_j)^T}{\sum_i^n P_{j(i)}}$$

Evaluate for convergence

$$\ln p(x|\mu, \Sigma, \pi) = \sum_i^n \ln \left(\sum_{k=1}^k \pi_k N(x_i | \mu_k, \Sigma_k) \right)$$

M-Estimation

M-estimator

1. Loss function $p: \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ Possible value of u
 $Q(u) = E[p(x, u)]$ against true distribution
2. an estimate with $\hat{u} \approx p(x_i, u)$
3. choose p so $Q(u)$ is minimized at u^*
 $\hat{u} = \text{argmin}_u \frac{1}{n} \sum p(x_i, u)$

Asymptotic Normality

$$\begin{aligned} J(u) &= E[\nabla p] = E\left[\frac{\partial p}{\partial u} p(x, u)\right] \\ k(u) &= \text{Var}[\nabla p] = \text{Var}\left[\frac{\partial p}{\partial u}(x, u)\right] \end{aligned}$$

Assume:

1. u^* is the only minimizer of Q
2. $J(u)$ is invertible for all $u \in M$
3. A few technical conditions

we have:

$$\begin{cases} \hat{u} \xrightarrow{n \rightarrow \infty} u^* \\ \sqrt{n}(\hat{u} - u^*) \xrightarrow{n \rightarrow \infty} N(0, J(u^*)^{-1} k(u^*) J(u^*)) \end{cases}$$

• P~~¶~~ log-likelihood if $J(u) = k(u) = I(u)$

$M\vec{u}$ is also a m-estimator

$M\vec{u}$ tend to be best with smallest asymptotic variance
while $M\vec{u}$ need assumption about distribution

P-function

- $P(x, u) = \|x - u\|_2^2 \Rightarrow u^* = \bar{x}$
- $P(x, u) = |x - u| \Rightarrow u^* = \text{Median}$
- $P(x, u) = Q(x - u) \Rightarrow u^*: \alpha\text{-quantile}$

$$Q(x) = \begin{cases} -u + \alpha x & \text{if } x < 0 \\ \alpha x & \text{if } x \geq 0 \end{cases}$$

Method of Moment Estimation

Moments

X be random variable with P_θ , k -th moment is

$$m_k = m_{X(\theta)} = E_\theta(X^k)$$

Moment Generating Function (MGF)

$$M_X(t) = E[e^{tX}], t \in \mathbb{R}$$

$$m_k = M_X^{(k)}|_{t=0}$$

Method of Moments

sample mean

$$\hat{m}_k = \frac{1}{n} \sum_{i=1}^n x_i^{(k)} \xrightarrow{i.i.d.} m_k$$

Estimator

$$m_1(\hat{\theta}_n) = \hat{m}_1$$

$$m_2(\hat{\theta}_n) = \hat{m}_2 \quad \begin{matrix} \text{solve linear} \\ \text{system} \end{matrix}$$

$$m_d(\hat{\theta}_n) = \hat{m}_d$$

$$\Psi: \Theta \rightarrow (m_1(\theta), \dots, m_d(\theta))$$

↓ CLT + Delta Method

$$\sqrt{n} (\hat{\theta}_n - \theta) \xrightarrow{d} N(0, [\frac{\partial m_i}{\partial \theta}(m(\theta))]^T \Sigma(\theta) [\frac{\partial m_i}{\partial \theta}(m(\theta))])$$

When to use: don't know the function working with

Hypothesis Testing

Statistical Formulation

IID sample: x_1, \dots, x_n with statistical model $(\hat{E}, \{\hat{P}_\theta\}_{\theta \in \Theta})$

consider 2 disjoint hypotheses $\{$ Null $H_0: \theta = \theta_0$, Alternative $H_1: \theta \neq \theta_0\}$

If we believe that true θ is either in Θ_0 or Θ_1 , we can test H_0 against H_1 . H_0 is the status quo. Data can only be used to reject H_0 . (about specific method)

★ "Innocent till proven guilty" *

Test

A test is a statistic $\Psi \in \mathcal{F}_{0,1}$ that does not depend on unknowns

$\Psi = \{0 \Leftrightarrow \text{not reject } H_0, 1 \Leftrightarrow \text{reject } H_0\}$ can always be written in "reject region"
 $= \mathbb{R}(\Psi)$, $\mathbb{R}\Psi = \{x \in E^n : \Psi(x) = 1\}$

Errors

	Fail to reject H_0	Reject H_0
H_0 is true: $H_0 \neq H_1$	✓	Type I
H_1 is true: $H_0 \neq H_1$	Type II	✓

beta function $\beta(\varphi) = \dot{\varphi}(\varphi=0)$

① H_0 is true: $P(\psi) = P(\psi \text{ make type I error})$

③ H_0 is true: $1 - \beta(\alpha) = P(\text{not make type I error})$

The Neyman-Pearson paradigm

- Make sure $P(\text{Type I error}) \leq$ level of the test
 - Minimize $P(\text{Type II error})$ subject to the constraint

A test ψ has level α if

$$\max_{\theta \in \Theta} P(\Psi = 1) \leq \gamma$$

Power of a test

$$\bar{\pi}_\psi = \inf_{\theta \in \Theta} (\mu - p_\psi(\theta))$$

Build tests from C]

- ① $P_0(\theta \in [A, B]) \geq 1-\alpha$, confidence interval at level $1-\alpha$
- ② $H_0 : \theta = \theta_0$
- ③ $H_1 : \theta \neq \theta_0$
- ④ Natural candidate
 $\psi = \{ \theta_0 \notin [A, B] \}$
- ⑤ $P_{\theta_0}(\psi = 1) = P_{\theta_0}[\theta_0 \notin I] = 1 - P_{\theta_0}[\theta_0 \in I] \leq 1 - (1-\alpha) = \alpha$
 ψ is a test with level α .

P value

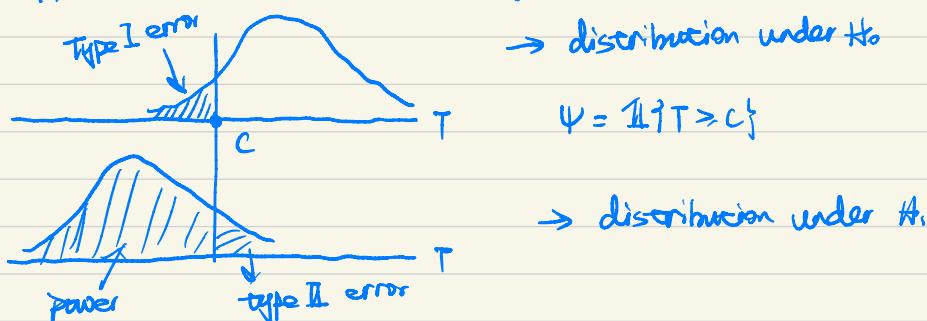
The (asymptotic) p-value of a test is the smallest (asymptotic) level α at which ψ rejects H_0 .

p-value $\leq \alpha \Leftrightarrow H_0$ is rejected by ψ at the level α

p-value $> \alpha \Leftrightarrow H_0$ is not rejected by ψ at the level α .

small p-value \Leftrightarrow strong evidence against H_0

Tradeoff between type I and type II



Also why usually one-sided test is more powerful than 2-sided test: the distribution under H_1 is usually bell shape not 2-peaked

Wald Test

For the estimator $\hat{\theta}$

We have $\frac{\hat{\theta} - \theta}{\sqrt{\text{Var}(\hat{\theta})}} \xrightarrow{n \rightarrow \infty} N(0, 1)$ Define $w := \frac{\hat{\theta} - \theta_0}{\sqrt{\text{Var}(\hat{\theta})}}$

	$H_0: \theta = \theta_0$ $H_1: \theta \neq \theta_0$	$H_0: \theta \leq \theta_0$ $H_1: \theta > \theta_0$	$H_0: \theta \geq \theta_0$ $H_1: \theta < \theta_0$
Wald Test	$\mathbb{I}\{ w > q_{\alpha/2}\}$	$\mathbb{I}\{w > q_{\alpha/2}\}$	$\mathbb{I}\{w < -q_{\alpha/2}\}$
P-value	$P(w > w^{\text{obs}})$	$P(w > w^{\text{obs}})$	$P(w < w^{\text{obs}})$
Asym P	$P(z > w^{\text{obs}})$	$P(z > w^{\text{obs}})$	$P(z < w^{\text{obs}})$

- $\text{Var}(\hat{\theta})$ can be any consistent variance estimator (Slusky)
- $\mathcal{Z} \sim N(0, 1)$
- For MLE estimator $\hat{\theta}^{\text{ML}}$, $\text{Var}(\hat{\theta}) = \frac{I(\hat{\theta})^{-1}}{n}$
 $w := \sqrt{n I(\hat{\theta}^{\text{ML}})} (\hat{\theta}^{\text{ML}} - \theta_0)$

Two Sample Wald Test

$$\begin{aligned} E(X) = \mu_1, \quad \text{Var}(X) = \sigma_1^2 &\Rightarrow H_0: \mu_1 = \mu_2, \quad \theta = \mu_1 - \mu_2 = 0 \\ E(Y) = \mu_2, \quad \text{Var}(Y) = \sigma_2^2 &\quad H_1: \mu_1 \neq \mu_2, \quad \theta = \mu_1 - \mu_2 \neq 0 \end{aligned}$$

$$\hat{\theta} = \bar{x}_n - \bar{t}_m \Rightarrow \frac{\hat{\theta} - \theta}{\sqrt{\text{Var}(\hat{\theta})}} \xrightarrow{n \rightarrow \infty, m \rightarrow \infty} N(0, 1)$$

$$\text{Var}(\hat{\theta}) = \text{Var}(\bar{x}_n) + \text{Var}(\bar{t}_m) = \frac{\sigma_1^2}{n} + \frac{\sigma_2^2}{m}$$

↓
Slusky: replace σ_1^2, σ_2^2 with $\hat{\sigma}_1^2, \hat{\sigma}_2^2$

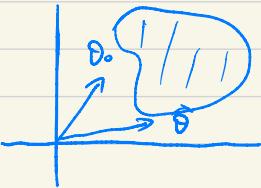
$$\frac{\hat{\theta} - \theta}{\sqrt{\frac{\hat{\sigma}_1^2}{n} + \frac{\hat{\sigma}_2^2}{m}}} \xrightarrow{n \rightarrow \infty, m \rightarrow \infty} N(0, 1)$$

Multi-variate Wald Test

$$\theta \in \mathbb{R}^d \quad H_0: \theta = \theta_0 \quad \text{Intuition: test if } \|\hat{\theta} - \theta_0\| \text{ is small enough}$$

$$H_1: \theta \neq \theta_0$$

$$W_n := n(\hat{\theta}_n^{\text{ML}} - \theta_0)^T I(\hat{\theta}_n^{\text{ML}})(\hat{\theta}_n^{\text{ML}} - \theta_0) \xrightarrow{n \rightarrow \infty} \chi_d^2$$



if $\|\theta_0 - \hat{\theta}\|$ small enough

Implicit hypotheses

Let $g: \mathbb{R}^k \rightarrow \mathbb{R}^k$ be continuous differentiable

Suppose $\sqrt{n}(\hat{\theta} - \theta^*) \xrightarrow{d} N(0, \Sigma(\theta^*))$, $\Sigma(\theta^*) \in \mathbb{R}^{k \times k}$

Multi-variate Delta

$$\sqrt{n}(g(\hat{\theta}_n) - g(\theta^*)) \xrightarrow{d} N_k(0, P(\theta))$$

$$P(\theta) = \nabla g(\theta) \Sigma(\theta) \nabla g(\theta) \in \mathbb{R}^{k \times k}$$

Assume $\Sigma(\theta)$ invertible and $\nabla g(\theta)$ has rank k

$$\sqrt{n}P(\theta)(g(\hat{\theta}_n) - g(\theta^*)) \xrightarrow{d} N_k(0, I_k)$$

Silvsky

$$\sqrt{n}P(\hat{\theta})(g(\hat{\theta}_n) - g(\theta^*)) \xrightarrow{d} N_k(0, I_k)$$

H₀: $g(\theta) = 0$ under null $g(\theta^*) = 0$

H₁: $g(\theta) \neq 0$

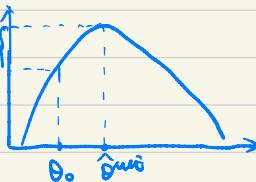
$$W_n := n g(\hat{\theta}_n)^T P^{-1}(\hat{\theta}_n) g(\hat{\theta}_n) \xrightarrow{d} \chi_k^2$$

Usually used with ML estimator

Likelihood Test

Basic Form

If large enough to
reject θ_0 .



Consider iid sample x_1, \dots, x_n with statistical model $(\tilde{E}, (\tilde{P}_\theta)_{\theta \in \Theta}), \Theta \subseteq \mathbb{R}^d$

$$H_0: (\theta_{\text{ref}}, \dots, \theta_d) = (\theta_{\text{ref}}^{(0)}, \dots, \theta_d^{(0)})$$

$$H_1: (\theta_{\text{ref}}, \dots, \theta_d) \neq (\theta_{\text{ref}}^{(0)}, \dots, \theta_d^{(0)})$$

$$\Downarrow \text{MLE: } \hat{\theta}_n = \underset{\theta \in \Theta}{\operatorname{argmax}} \ln(\theta)$$

$$\text{constrained MLE: } \hat{\theta}_n^c = \underset{\theta \in \Theta_0}{\operatorname{argmax}} \ln(\theta)$$

$$T_n = 2(\ln(\hat{\theta}_n) - \ln(\hat{\theta}_n^c))$$

↓ Wilk's Theorem

$$\left\{ \begin{array}{l} T_n \xrightarrow{n \rightarrow \infty} \chi_{d-r}^2 \\ \psi = \mathbb{I}\{T_n > q_\alpha\} \end{array} \right. \quad \begin{array}{l} \dim(\Theta \setminus \Theta_0) \\ \parallel \end{array} \quad d-r = \dim(\Theta) - \dim(\Theta_0)$$

q_α : $(1-\alpha)$ quantile of χ_{d-r}^2

Neyman-Pearson Lemma

Among all tests for the hypothesis $H_0: \theta = \theta_0, H_1: \theta \neq \theta_0$ at significant level α . Likelihood ratio test is the most powerful: highest probability to reject the null.

T-test

Cochran's Theorem

If $X_1, \dots, X_n \sim N(\mu, \sigma^2)$, $S_n = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$

1. \bar{X}_n is independent with S_n

2. $\frac{(n-1)S_n^2}{\sigma^2} \sim \chi_{n-1}^2$

↓ (lose 1 df to replace μ with \bar{X}_n)

$$\frac{\sqrt{n} \frac{\bar{X}_n - \mu}{\sigma}}{\sqrt{S_n^2 / \sigma^2}} \sim N(0, 1) \sim t_{n-1}$$

T test

- $X_1, \dots, X_n \sim N(\mu, \sigma^2)$
- $\sqrt{n} \frac{\bar{X}_n - \mu_0}{S_n} \sim t_{n-1}$ for all n

$T := \sqrt{n} \frac{\bar{X}_n - \mu_0}{S_n}$	$H_0: \mu = \mu_0$ $H_1: \mu \neq \mu_0$	$H_0: \mu \leq \mu_0$ $H_1: \mu > \mu_0$	$H_0: \mu \geq \mu_0$ $H_1: \mu < \mu_0$
T-test	$\{T \mid T > q_{\alpha/2}^{t_{n-1}}$	$\{T \mid T > q_{1-\alpha}^{t_{n-1}}$	$\{T \mid T < q_{\alpha}^{t_{n-1}}$
P-value	$P(T \mid T > T ^{obs})$	$P(T > T ^{obs})$	$P(T < T ^{obs})$

• $q_t > q_{\text{normal}}$: heavy tail

• Wald test p-value $\gtrsim n$, more likely to reject H_0

Two Sample Test

$X_1, \dots, X_n \sim N(\mu_1, \sigma_1^2)$

$H_0: \mu_1 - \mu_2 \leq 0$

$Y_1, \dots, Y_m \sim N(\mu_2, \sigma_2^2)$

$H_1: \mu_1 - \mu_2 > 0$

$$T = \frac{\bar{X}_n - \bar{Y}_m - (\mu_1 - \mu_2)}{\sqrt{S_1^2/n + S_2^2/m}} \sim t_{n+m-2}$$

$$N = \frac{(S_1^2/n + S_2^2/m)^2}{S_1^2/(n(n-1)) + S_2^2/(m(m-1))} \geq \min(n, m)$$

Welch-Satterthwaite (WS) formula

Multiple Hypothesis Test

Issue: try to make aggregate decision but only control error individually

	Not reject H_0 (non-significant)	Reject H_0 (significant)	
H_0 true	A (# of true non-significant)	V (# of type I error)	m
H_1 true	C (# of type II error)	D (# of true significant)	m_1
	N_0 (# of non-significant)	N_1 (# of significant)	m

FwFR: family wise error rate

- Probability of at least one false discovery (type I)
- $FwFR = P(V \geq 1) = P\left(\sum_{i=1}^m V_i \geq 1\right)$
 $= 1 - P(W=0) = 1 - (1-\alpha)^m \rightarrow 1$

Usually too strict for large m

FDR: False Discovery Rate

- Expected fraction of false significant results among all significant results
- $FDR = E[V/N_1]$
- FDR has higher power compared to FwFR. FwFR is stricter
 $E[V/N_1] \leq E[1_{W \geq 1}] = P(W \geq 1)$
 $\uparrow \quad \begin{cases} V/N_1 \leq 1_{W \geq 1} & \text{when } V \geq 1 \\ V/N_1 = 1_{W \geq 1} = 0 & \text{when } V = 0 \end{cases}$

$$\boxed{FDR \leq FwFR}$$

- Power of a series of tests with FDR are more powerful

Bonferroni Correction \rightarrow FWER

Use α/m instead of α to reject

- $P_i < \alpha/m \Leftrightarrow \text{FWER} < \alpha$
- $\text{FWER} = P(\cup_{i=1}^m \{\psi_i = 1\}) = \sum_{i=1}^m P(\psi_i = 1) = \sum_{i=1}^m \alpha/m = m\alpha/m \leq \alpha$
- $I_0 = \{i = 1, \dots, m \mid H_0 \text{-true}\}$

Benjamini - Hochberg \rightarrow FDR

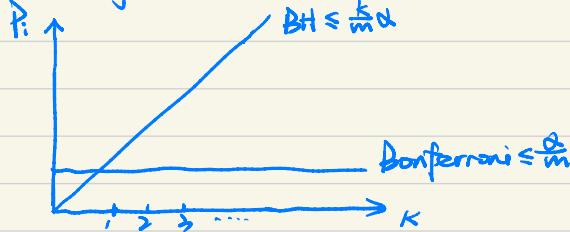
For a series of m independent test

1. Sort the p-value : $p^{(1)} \leq p^{(2)} \leq \dots \leq p^{(m)}$

2. Find the maximum k that

$$p^{(k)} \leq \frac{k}{m} \alpha$$

3. Reject $H_0^{(1)}, H_0^{(2)}, \dots, H_0^{(k)}$



Goodness of Fit - Discrete

x_1, \dots, x_n n discrete P.

| Sample space $\mathbb{E} = \{o_1, \dots, o_k\}$
| probability measure $P = [P_1, \dots, P_k]^T \quad \sum P_i = 1, P_i \geq 0$

GOF - Chi square

| $H_0: P = P^0 \Rightarrow T_n := n \sum_{j=1}^k \frac{(\hat{P}_j - P_j^0)^2}{P_j^0} \xrightarrow{n \rightarrow \infty} \chi^2_{k-1} \quad (\Psi = \mathbb{1}(T_n > \chi_{\alpha}^2))$
| $H_1: P \neq P^0$

$\Downarrow \text{Mö}$

$\hat{P}_j = \frac{n_j}{n}, \sqrt{n}(\hat{P}_j - P_j^0)$ is asymptotic normal in dim k-1

Degenerate since $\sqrt{n} \sum (\hat{P}_j - P_j^0) = 0$

$\sqrt{n} I^{\frac{1}{2}}(\hat{P}) (\hat{P} - P^0) \xrightarrow{n \rightarrow \infty} N_{k-1}(0, I_{k-1})$ no variance along 1 direction

GOF test on Family of Distribution

$H_0: P \in \{P_\theta\}_{\theta \in \mathbb{R}^d}$

$H_1: P \notin \{P_\theta\}_{\theta \in \mathbb{R}^d}$

• $\{o_1, \dots, o_k\}, f_\theta(o_i) = P(X=o_i)$

• Observe $x_1, \dots, x_n \Rightarrow \text{Mö } \hat{P} = [\frac{n_1}{n}, \dots, \frac{n_k}{n}]^T$

• $T_n := n \sum_{i=1}^k \frac{(\frac{n_i}{n} - f_\theta(o_i))^2}{f_\theta(o_i)} \xrightarrow{n \rightarrow \infty} \chi^2_{(k+1)-1-d}$

$\vee k+1$: the support size of P_θ (for all θ)

Goodness of Fit - Continuous

Convergence of Function

1. Convergence piecewise: for each $x \lim_{n \rightarrow \infty} g_n(x) = g(x)$
2. Convergence uniformly: for every $M > 0$, exists n_M such that $\sup_x |g_n(x) - g(x)| < M$ for all $n \geq n_M$
- Example of piecewise but not uniformly
$$g_n(x) = \frac{1}{x^n} \quad \sup_{x>0} g_n(x) = \sup_{x>0} \frac{1}{x^n} = 1$$

CDF

1. CDF : $F(t) = P(X \leq t) = E[\mathbb{1}(X \leq t)]$
empirical CDF : $F_n(t) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}(x_i \leq t)$

1. Convergence piecewise from LTP
 $F_n(t) \xrightarrow{a.s.} F(t)$ for all $t \in \mathbb{R}$
2. Convergence uniformly from GL (Glivenko-Cantelli) Theorem
 $\sup_{t \in \mathbb{R}} |F_n(t) - F(t)| \xrightarrow{a.s.} 0$
- Example of 1 but not 2
 $F_n(t) = F(t) + \frac{\epsilon}{n}$; $\forall t \quad F_n(t) = F(t)$
 $\sup_{t \in \mathbb{R}} |F_n(t) - F(t)| = \sup_{t \in \mathbb{R}} \left| \frac{\epsilon}{n} \right| = \infty$

$$\mathbb{1}\{X_i \leq t\} \sim \text{Ber}(F(t))$$

1. CLT : $\sqrt{n} (F_n(t) - F(t)) \xrightarrow{d} N(0, F(t)(1-F(t)))$

2. Donsker's Theorem

$$\sqrt{n} \sup_{t \in \mathbb{R}} |F_n(t) - F(t)| \xrightarrow{d} \sup_{t \in \mathbb{R}} B(t) \quad \text{Brownian bridge}$$

worst possible case A pivotal distribution that does not depend on $F(t)$
of all t

Goodness of continuous distribution

$X_1, \dots, X_n \sim \text{CDF } F$

$$H_0 : F = F^*$$

$$H_1 : F \neq F^*$$

Kolmogorov-Smirnov test

$$\begin{aligned} T_n &= \sup_{t \in \mathbb{R}} |F_n(t) - F^*(t)| \\ &= \sqrt{n} \sup_{t \in \mathbb{R}} \left| \left(\frac{1}{n} \sum_{i=1}^n \mathbb{1}(X_i \leq t) \right) - F^*(t) \right| \\ &\Downarrow \text{Let } \hat{F} = F^*, \quad Y_i \sim \text{Unif}[0, 1] \end{aligned}$$

$$= \sqrt{n} \sup_{t \in [0, 1]} \left| \left(\frac{1}{n} \sum_{i=1}^n \mathbb{1}(Y_i \leq \hat{t}) \right) - \hat{t} \right|$$

Calculation \Downarrow observation $x_0 \leq x_1 \leq \dots \leq x_n$

$$T_n = \sqrt{n} \max_{i=1, \dots, n} \{ \max(|\frac{i-1}{n} - \hat{F}(x_{i-1})|, |\frac{i}{n} - \hat{F}(x_i)|) \}$$

No analytical solution for quantile of Brownian

→ Simulation M : T_1, \dots, T_M . Use sample quantile as estimation

$$\uparrow \Psi = \mathbb{I}\{T_n > \hat{q}_{\alpha, M}\} \quad (\text{or table})$$

$$\text{P-value} \approx \frac{\#\{j=1, \dots, M : T_j > T_n\}}{M}$$

test based on function distance measure:

- sup-norm: $d(F_n, F) = \sup_t |F_n(t) - F(t)|$

- L2 distance: $d^2(F_n, F) = \int_R (F_n(t) - F(t))^2 dt$

- scaled L2: $\sqrt{R} (F_n(t) - F(t))^2 / (F(t)(1-F(t))) dt$

Kolmogorov-Lilliefors test

Check if from Gaussian distribution

$$\sup_{t \in \mathbb{R}} |F_n(t) - \Phi_{\text{std}}(t)| \rightarrow k-1 \text{ table}$$

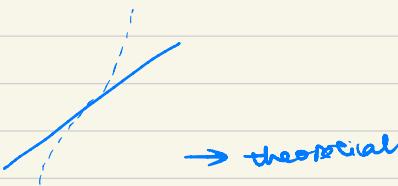
QQ plot

A visual way for GOF test

- $x_0 \leq \dots \leq x_n$

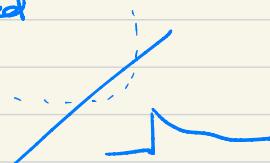
- $(F^{-1}(\frac{1}{n}), x_0), (F^{-1}(\frac{2}{n}), x_1), \dots, (F^{-1}(\frac{1}{n}), x_n), \dots, (F^{-1}(\frac{n-1}{n}), x_{n-1})$

① heavy tail
sample ↑



→ theoretical

② right skewed



③ left skewed



④ light tail



Bayesian Statistic

Most effective when experts have a lot of background knowledge
→ can bring a lot of prior information

Bayes Formula

Prior $\pi(\theta)$ } ✓ a function by which the likelihood function is weighted by in order to produce posterior
✓ is to be specified by the researcher in order to take into account previous knowledge about possible values of the parameter

$$\pi(\theta | x_1, \dots, x_n) = \frac{\pi(\theta) L(x_1, \dots, x_n | \theta)}{\int L(x_1, \dots, x_n | t) \pi(t) dt} \quad \forall \theta \in \Theta$$

$$\propto \pi(\theta) L(x_1, \dots, x_n | \theta) \quad \forall \theta \in \Theta$$

proportional, not depend on θ

Priors

1. Whether we can specify the parameters to approximate belief
2. Whether it is realistic
3. Whether computationally tractable
i.e. conjugate prior of the likelihood distribution
 \Leftrightarrow prior and posterior distribution from same family

We can still use Bayesian without prior information

Noninformative prior

$\pi(\theta) \propto 1$ } If Θ is bounded \Rightarrow uniform prior on Θ
If Θ is unbounded \Rightarrow improper prior $\Leftrightarrow \int \pi(\theta) d\theta = \infty$
can still use

Jeffreys Prior

Reparameterize Distribution

- $\pi(\theta)$, $\theta \in \mathbb{R}$
- \Downarrow reparameterize with $n = \phi(\theta)$ ϕ strict monotone
- $\pi(n)$

$$\int_{\theta_0}^{\theta_0 + \Delta\theta} \pi(\theta) d\theta = \int_{n_0}^{n_0 + \Delta n} \pi(n) dn$$

$\Downarrow \Delta\theta \rightarrow 0 \quad \Downarrow \Delta n \rightarrow 0$

$$\Delta\theta \pi(\theta_0) = \Delta n \pi(n_0)$$

$$\pi(n_0) = \pi(\theta_0) \left| \frac{\Delta\theta}{\Delta n} \right|_{\theta=\theta_0}$$

$$\Downarrow \theta = \phi(n) \quad \phi'(n) = dn/d\theta$$

$$\boxed{\pi(n) = \frac{\pi(\theta)}{|\phi'(\theta)|} = \frac{\pi(\phi(n))}{|\phi'(\phi(n))|}}$$

Jeffreys Prior

A non-informative prior. The probability assigned to a volume of the probability space is the same under all parameterization.

$$\pi_J(\theta) \propto \sqrt{\det(I(\theta))}$$

$I(\theta)^{-\frac{1}{2}}$: standard error of estimator

$|I(\theta)| \downarrow \Rightarrow$ more uncertainty \Rightarrow smaller weight

$|I(\theta)| \uparrow \Rightarrow$ less uncertainty \Rightarrow larger weight

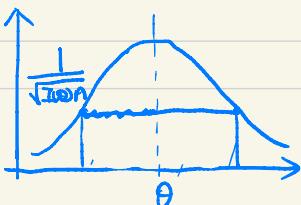
$$I(\theta) = E\left[\left(\frac{\partial}{\partial\theta} \ln L(x; \theta)\right)^2\right] = -E\left[\frac{\partial^2}{\partial\theta^2} \ln L(x; \theta)\right]$$

① marginal change of θ to log-likelihood

$$E\left[(\Delta \ln L(x; \theta))^2\right] = I(\theta) (\Delta\theta)^2$$

② curvature around θ

more curved \Leftrightarrow higher $I(\theta)$



Reparameterization Invariance of Jeffreys

taking sensitivity into account



$$\begin{cases} \ln L(\theta_i | \theta) = \ln L(\theta_i | n) \\ \ln L(\theta_i | \theta + \Delta\theta) = \ln L(\theta_i | n + sn) \end{cases} \Rightarrow \Delta\theta \frac{\partial}{\partial \theta} \ln L(\theta_i | \theta) = sn \frac{\partial}{\partial n} \ln L(\theta_i | n)$$

When $\phi'(\theta) > 1$, marginal change in log-likelihood is larger than marginal change caused by sn

$\pi_J(\theta)$ gives higher weight to

- 1. θ whose MLE estimate is more certain
- 2. θ where data has more information in deciding param
- 3. θ where marginal shifts have larger effect to x_i
 - $I(\theta)$ is proxy of how much, at a particular θ , would equivalent shifts to the parameter influence the data.

$$\begin{cases} \pi_J(\theta) \propto \sqrt{I(\theta)} \\ \pi_J(\theta) \propto \sqrt{I(\theta)} \end{cases} \xrightarrow{n=\phi(\theta)} \pi_J(n) \text{ is the reparameterization of } \pi_J(\theta) \text{ through } n=\phi(\theta)$$

I.F : $X \sim \text{Ber}(q^{10})$

$$\begin{aligned} \textcircled{1} \quad I(q) &= -E\left[\frac{\partial^2}{\partial q^2} \ln L(x_i | q)\right] = -E\left[\frac{\partial^2}{\partial q^2} \ln(q^{10x} (1-q)^{10(1-x)})\right] \\ &= 100q^8/(1-q)^2 \end{aligned}$$

$$\textcircled{2} \quad P = q^{10} \quad \phi(p) = p^{1/10} \quad \phi'(p) = \frac{1}{10} p^{-9/10} \quad \phi''(p) = p^{-10}$$

$$\pi_J(p) \propto \frac{1}{p^{1/10}}$$

$$\frac{\pi_J(p)}{|\phi'(\phi(p))|} = \sqrt{\frac{100p^8}{1-p^{10}}}$$

Linear Regression

Goal: Given data point $\{(x_i, y_i)\}$, fit a model to predict Y given a value for X

Model Set Up

(x_i, y_i) , $i=1, \dots, n$ are iid from unknown joint distribution P
Plan be entirely described by:

① joint PDF $h(x, y)$ or

② marginal of X : $h_X(x) = \int_y h(x, y) dy$

conditional: $h_{Y|X}(y|x) = h(x, y)/h_X(x)$

can also describe the distribution partially i.e: expectation
regression function of Y with respect to X :

$$\mu(x) = E[Y|X=x]$$

$$= \int_y y f(y|x) dy$$

linear regression function (\geq representation)

$$\mu(x) = E[Y|X=x] = a + bx$$

$$\underset{a, b}{\text{argmin}} E[(Y - a - bx)^2] \Rightarrow$$

$$b^* = \frac{\text{cov}(x, Y)}{\text{var}(x)}, a^* = E(Y) - \frac{\text{cov}(x, Y)}{\text{var}(x)} E(x)$$

$\epsilon = Y - (a^* + b^*x)$ is noise

with $E(\epsilon) = 0$, $\text{cov}(x, \epsilon) = 0$

↓ replace expectation with average

LSE of (a^*, b^*) is $\hat{\beta} = \underset{\beta}{\text{argmin}} \sum_i (Y_i - x_i^\top \beta)^2$

Multi-variate Linear Model

$$Y_i = X_i^T \beta + \varepsilon_i \quad | \quad Y_n = X_n^T \beta + \varepsilon_n \quad \Leftrightarrow \quad \begin{pmatrix} Y_1 \\ \vdots \\ Y_n \end{pmatrix} = \begin{pmatrix} X_1^T \\ \vdots \\ X_n^T \end{pmatrix}_{n \times p} \begin{pmatrix} \beta \\ \vdots \\ \beta \end{pmatrix} + \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{pmatrix}$$

$$Y = X \beta + \varepsilon$$

$Y = X^T \beta + \varepsilon$ is an assumption on the model. Equivalently, we can assume that the regression function is linear: $u(x) = E(Y|X=x) = X^T \beta$ with the understanding that $E(\varepsilon) = 0$
 (can always do linear regression even when model is misspecified)

↓ linear regression

Find the best $\hat{\beta}$

$$\underset{\beta}{\text{LS}} = \underset{\beta}{\operatorname{argmin}} \|Y - X\beta\|^2 \quad \xrightarrow{\text{1st order}} \quad X^T X \hat{\beta} = X^T Y$$

$$= (X^T X)^{-1} X^T Y \quad \xleftarrow{\text{rank}(X) = p \Rightarrow \text{unique } \hat{\beta}}$$

$\hat{Y} = X \hat{\beta} = X(X^T X)^{-1} X^T Y$ is linear combination of column of X . The orthogonal projection of Y on the space spanned by column of X

Fisher Information

$$f(y; \beta) = \frac{1}{\sigma^2} \exp(-\frac{1}{2\sigma^2} (Y - X\beta)^2)$$

$$\Rightarrow \nabla_{\beta} l(y; \beta) = -\frac{1}{\sigma^2} (-2YX + 2X^T X\beta)$$

$$= \frac{1}{\sigma^2} (YX - X^T X\beta)$$

$$\Rightarrow H_{\beta} l(y; \beta) = -X X^T / \sigma^2$$

$$I(\beta) = -\sum_i H_{\beta} l(y_i; \theta) = -\frac{n}{\sigma^2} \text{Tr}[-\frac{1}{\sigma^2} X^T X]$$

$$= \frac{1}{\sigma^2} X^T X$$

LSE property under deterministic X and Gaussian ε

$$\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} \quad \text{assume } \begin{cases} 1. \mathbf{X} \text{ is deterministic, } \text{rank}(\mathbf{X}) = p \\ 2. \text{homoskedastic, } \varepsilon_1, \dots, \varepsilon_n \text{ iid} \\ 3. \varepsilon \sim N(0, \sigma^2 I_n) \end{cases}$$

$$1. LSE = \text{ML}$$

$$Y_i \sim N(X_i \beta^*, \sigma^2)$$

$$\log L(Y_1, \dots, Y_n, \beta) = -n \ln(\sigma \sqrt{2\pi}) - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - X_i \beta)^2$$

$$2. \hat{\beta} \sim N_p(\beta^*, \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1})$$

$$\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{X} \beta^* + \varepsilon)$$

$$= \beta^* + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \varepsilon \quad \text{N}(0, \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1})$$

$$= \beta^* + N(0, \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1})$$

• $\mathbf{X}^T \mathbf{X}$ more spread $\rightarrow \text{Var}(\hat{\beta})$ smaller

$$3. \text{Quadratic risk } E[\|\hat{\beta} - \beta^*\|_2^2] = \sigma^2 \text{tr}(\mathbf{X}^T \mathbf{X})$$

$$E[\|\hat{\beta} - \beta^*\|_2^2] \\ = E[\text{tr}((\hat{\beta} - \beta^*)(\hat{\beta} - \beta^*)^T)]$$

$$\Downarrow \quad \|\mathbf{x}\|_2^2 = \text{tr}(\mathbf{x}^T \mathbf{x}) = \text{tr}(\mathbf{x} \mathbf{x}^T)$$

$$= \text{tr} E[(\hat{\beta} - \beta^*)(\hat{\beta} - \beta^*)^T] = \sigma^2 \text{tr}((\mathbf{X}^T \mathbf{X})^{-1})$$

$$4. \text{Prediction Error } E[\|Y - \mathbf{X} \hat{\beta}\|_2^2] = \sigma^2(n-p)$$

$$\|Y - \mathbf{X} \hat{\beta}\|_2^2 = \sum \hat{\varepsilon}_i^2$$

$$= \|Y - P\mathbf{Y}\|_2^2 = \|(\mathbf{I}_n - P)\mathbf{Y}\|_2^2 = \|P^\perp \mathbf{Y}\|_2^2$$

$$5. \text{Unbiased estimator of } \hat{\sigma}^2 = \frac{1}{n-p} \|Y - \mathbf{X} \hat{\beta}\|_2^2$$

$$\hat{\sigma}^2 = \frac{1}{n-p} \|Y - \mathbf{X} \hat{\beta}\|_2^2 = \frac{1}{n-p} \sum \hat{\varepsilon}_i^2$$

\Downarrow Cochran's Theorem

$$(n-p) \frac{\hat{\sigma}^2}{\sigma^2} \sim \chi_{n-p}^2$$

$$\hat{\beta} \perp \hat{\sigma}^2$$

Hypothesis Test

$$\text{Test: } H_0: u^\top \beta^* \leq 0 \quad u \in \mathbb{R}^p$$

$$H_1: u^\top \beta^* > 0$$

Assume X full rank \Rightarrow identifiable

1. MLE estimator

$$\hat{P}_{\beta, \sigma^2}(Y_1, \dots, Y_n) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(Y_i - X_i^\top \beta)^2\right)$$

$$\begin{aligned} \hat{L}_{\beta, \sigma^2} &= \log \hat{P}_{\beta, \sigma^2}(Y_1, \dots, Y_n) \\ &= -\frac{n}{2} \log(2\pi) - \frac{n}{2} \log \sigma^2 - \frac{1}{2\sigma^2} \|Y - X\beta\|_2^2 \end{aligned}$$

$$\begin{cases} \frac{\partial \hat{L}}{\partial \beta} = 0 \Rightarrow X^\top Y - X^\top X \beta = 0 & \Rightarrow \boxed{\hat{\beta} = (X^\top X)^{-1} X^\top Y} \\ \frac{\partial \hat{L}}{\partial \sigma^2} = -\frac{n}{2} - \frac{1}{2\sigma^2} \|Y - X\beta\|_2^2 = 0 & \Rightarrow \boxed{\hat{\sigma}_{\text{MLE}}^2 = \frac{1}{n} \|Y - X\beta\|_2^2} \end{cases}$$

2. Pivotal Statistic

$$\begin{aligned} u^\top \hat{\beta} &= u^\top (X^\top X)^{-1} X^\top Y = u^\top (X^\top X)^{-1} X^\top (X\beta^* + \varepsilon) \\ &= u^\top \beta^* + u^\top (X^\top X)^{-1} X^\top \varepsilon \end{aligned}$$

$$\Rightarrow \boxed{u^\top \hat{\beta} - u^\top \beta^* \sim N(0, \hat{\sigma}^2 u^\top (X^\top X)^{-1} u)}$$

$$\begin{aligned} \|Y - X\hat{\beta}\|_2^2 &= \|Y - X(X^\top X)^{-1} X^\top Y\|_2^2 = \|\varepsilon - X(X^\top X)^{-1} X \varepsilon\|_2^2 \\ &= \|(I_n - H)\varepsilon\|_2^2 \end{aligned}$$

$$\Downarrow \begin{aligned} &\exists \text{ orthogonal } V \text{ and diagonal } \Lambda \text{ such that} \\ &H = V \Lambda V^\top = V \begin{bmatrix} I_p & 0 \\ 0 & 0 \end{bmatrix} V^\top = \begin{bmatrix} V_1 & V_2 \\ \vdots & \vdots \\ 0 & 0 \end{bmatrix} \begin{bmatrix} I_p & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} V_1^\top \\ V_2^\top \end{bmatrix} \end{aligned}$$

$$\begin{aligned} I_n - H &= I_n - V \Lambda V^\top = V V^\top - V \Lambda V^\top = V_2 V_2^\top \\ &= \|V_2^\top \varepsilon\|_2^2 \end{aligned}$$

$$\Rightarrow \boxed{\|Y - X\hat{\beta}\|_2^2 / (\hat{\sigma}^2 n-p) \sim \chi^2_{n-p}}$$

$$\boxed{\frac{u^\top \hat{\beta} - u^\top \beta^*}{\hat{\sigma} \sqrt{u^\top (X^\top X)^{-1} u}} \sim N(0, 1)}$$

$$\boxed{\frac{u^\top \hat{\beta} - u^\top \beta^*}{\hat{\sigma} \sqrt{u^\top (X^\top X)^{-1} u}} \sim t_{n-p}}$$

Multiple Hypothesis Test

$$Y = X\beta + \varepsilon \quad , \quad \beta, \sigma^2 \text{ unknown}$$

Test $H_0: 0 \leq \beta_i \leq \beta_i \quad \Leftrightarrow \quad H_0: 0 \leq \beta_i \text{ and } \beta_i \leq \beta_i$
 $H_1: \text{not } 0 \leq \beta_i \leq \beta_i \quad H_1: 0 > \beta_i \text{ or } \beta_i > \beta_i$

$$H_0 = \begin{cases} H_0^{(1)}: 0 \leq \beta_i \Rightarrow \psi_{0,i}^{(1)} \Rightarrow \psi_0 = \psi_{0,1}^{(1)} \text{ or } \psi_{0,2}^{(1)} \\ H_0^{(2)}: \beta_i = \beta_i \Rightarrow \psi_{0,i}^{(2)} \end{cases} \quad \max(\psi_0^{(1)}, \psi_0^{(2)})$$

$$P(\psi \text{ reject}) = P(\psi^{(1)} \text{ reject or } \psi^{(2)} \text{ reject}) \geq P(\psi^{(1)} \text{ reject})$$

|| Union bound

$$P(A_1 \cup \dots \cup A_n) \leq \sum_{i=1}^k P(A_i)$$

Fisher test: Test whether a group of explanatory variable is significant in linear regression

$$\psi_{SA} = \bigcup_{j \in S} \psi_{j, \text{SA}}$$

$$\left. \begin{array}{l} \textcircled{1} \quad \hat{\beta}^{\text{MB}} = (X^T X)^{-1} X^T Y \sim N(\beta, \sigma^2 (X^T X)^{-1}) \\ \hat{\beta}_i = e_i^T \hat{\beta}^{\text{MB}}, \quad e_i = (1, 0, \dots, 0)^T \\ \hat{\beta}_i \sim N(e_i^T \beta, e_i^T [\sigma^2 (X^T X)^{-1}] e_i) \\ T^{(1)} = \frac{e_i^T \hat{\beta}}{\sqrt{\hat{\sigma}^2 M_{ii}}} \Rightarrow \text{If } T^{(1)} < q_{1-\alpha}^{t\text{ap}} \} \\ \textcircled{2} \quad \hat{\beta}_2 - \hat{\beta}_1 = U^T \hat{\beta}^{\text{MB}} \quad U = (1, -1, 0, \dots, 0)^T \\ \sim N(U^T \beta, U^T [\sigma^2 (X^T X)^{-1}] U) \\ \hat{\beta}_2 - \hat{\beta}_1 \sim N(0, \sigma^2 (M_{11} + M_{22} - 2M_{12})) \\ \Rightarrow T^{(2)} = \frac{\hat{\beta}_2 - \hat{\beta}_1}{\sqrt{\hat{\sigma}^2 (M_{11} + M_{22} - 2M_{12})}} \Rightarrow \text{If } T^{(2)} < q_{1-\alpha}^{t\text{ap}} \} \end{array} \right.$$

Ridge Regression

$$y_i = x_i^T \beta + \varepsilon_i, \quad \varepsilon_i \stackrel{iid}{\sim} N(0, \sigma^2)$$

$$Y = X \beta + \varepsilon, \quad \varepsilon \sim N(0, \sigma^2 I_n)$$

$n \approx p \approx n$

$$\text{OLS: } \hat{\beta}^{\text{OLS}} = \arg \min \|Y - X\beta\|_2^2 = (X^T X)^{-1} X^T Y$$

1. Bayesian - Assume $\pi \sim N(0, \tau^2 I_p)$

$$\pi(\beta | y) = \frac{\pi(\beta) P(Y|\beta)}{P(Y)} \propto \pi(\beta) P(Y|\beta)$$

$$= \frac{1}{\sqrt{2\pi}^p} \frac{1}{\sqrt{2\pi}^n} \exp \left[-\frac{1}{2\sigma^2} (\|Y - X\beta\|_2^2 + \frac{\tau^2}{2} \|\beta\|_2^2) \right]$$

$$\Downarrow \text{Let } \lambda = \frac{\tau^2}{2\sigma^2}$$

$$\star = (Y - X\beta)^T (Y - X\beta) + \lambda \beta^T \beta$$

$$= \sigma^2 (\beta - u)^T \Sigma^{-1} (\beta - u) + Y^T Y - \sigma^2 u^T \Sigma u$$

$$\pi(\beta | y) \sim N(u, \Sigma) \Rightarrow \begin{cases} u = (X^T X + \lambda I_p)^{-1} X^T Y \\ \Sigma = \sigma^2 (X^T X + \lambda I_p)^{-1} \end{cases}$$

2. $\hat{\beta}^{(R)} = \arg \min \|Y - X\beta\|_2^2 + \lambda \|\beta\|_2^2$ R: ridge regression
with $\tau^2 = \frac{\sigma^2}{2}$, $\hat{\beta}^{(R)} = \hat{\beta}^{(2)}$

3. Special case: assume $p=n$, $X^T X = I_p$

$$\hat{\beta}^{(R)} = (X^T X + \lambda I_p)^{-1} X^T Y = (u + \lambda I_p)^{-1} X^T Y$$

$$= X^T Y / (u + \lambda)$$

$$E[\|\hat{\beta}^{(R)} - \beta\|_2^2] = \frac{\lambda^2}{1+\lambda} \|\beta\|_2^2 + \frac{\sigma^2}{1+\lambda}$$

Remark

Especially useful when X is rank deficient ($p > n$). When X deficient, it looks flat, there's no unique solution or hard to get. $\lambda \|\beta\|_2^2$ add some curvature at β so quadratic risk looks like "ridge"

Generalized Linear Model

GLM - relax the linear constraint

1. Random component

$Y|X \sim \text{some distribution}$

2. Regression function

$$g(u(x)) = X^T \beta \quad g \circ u \text{ is linear}$$

} g : link function

u : regression function

In general no closed form

Exponential Family Distribution

$\{f_{\theta}(y) : \theta \in \Theta\}$, $\Theta \subset \mathbb{R}^k$ is k -parameter exponential family on \mathbb{R}^q iff

$$f_{\theta}(y) = h(y) \exp \left[\sum_{j=1}^k \eta_j(\theta) T_j(y) - b(\theta) \right]$$

$$\begin{matrix} k \rightarrow k \\ \underbrace{\eta_1(\theta) \dots \eta_k(\theta)}_{T(y)} \end{matrix} \left| \begin{matrix} T_j(y) \\ T(y) \end{matrix} \right. \quad q \rightarrow k$$

One-parameter canonical exponential family - another parametrization

$$f_{\theta}(y) = \exp \left(\frac{y\theta - b(\theta)}{\phi} + c(y, \phi) \right) \quad \begin{cases} b(\theta) : \text{log-partition function} \\ \phi : \text{dispersion parameter} \end{cases}$$

Likelihood: $l(\theta) = \log f(\gamma)$ meets the following identity

$$\begin{cases} \text{① First Identity} & E\left(\frac{\partial l}{\partial \theta}\right) = 0 \\ & \int f_{\theta}(y) dy = 1 \Rightarrow \int \frac{\partial}{\partial \theta} f_{\theta}(y) dy = 0 \\ \text{② Second Identity} & E\left(\frac{\partial^2 l}{\partial \theta^2}\right) + E\left(\frac{\partial l}{\partial \theta}\right)^2 = 0 \\ \quad \left[\begin{matrix} E(Y) = b'(\theta) \\ \text{Var}(Y) = b''(\theta) \phi \end{matrix} \right] \end{cases}$$

Link Function

Map range of expectation of Y to \mathbb{R} (range of $X^T \beta$)

$$\text{if } Y \sim \text{Bernoulli } p \Rightarrow \text{logit: } \log \left(\frac{u(x)}{1-u(x)} \right) = X^T \beta$$

$$\text{Probit: } \Phi^{-1}(u(x)) = X^T \beta$$

Canonical link

The function g that links the mean μ to the canonical parameter θ is call canonical link.

$$g(\mu(x)) = \theta = (b')^T \mu$$

I.E.: Bernoulli Y.

$$\text{Pmf} = p^y (1-p)^{1-y}$$

$$= \exp(y \frac{\theta}{\phi}) + \frac{\log(1-\theta)}{-\phi}$$

$$\Rightarrow p = \frac{e^\theta}{1+e^\theta}, \text{ canonical link is the logit link}$$

Parameter the model with θ

$$Y_i | X_i = x_i \sim \text{for}(y_i) = \exp((y_i \beta - b(x_i)) / \phi + \text{const.} / \phi)$$

$$\mu_i = b'(\theta_i)$$

$$\theta_i = g(\mu_i) = x_i \beta$$

$$\theta_i = (b')^{-1}(\mu_i)$$

$$= (b')^{-1}(g'(x_i^T \beta)) \equiv h(x_i^T \beta)$$

$$h = (b')^{-1} \circ g^{-1} = \underbrace{(g \circ b')^{-1}}_{\text{for canonical link}}$$

$$\downarrow \quad g = (b')^{-1} \Rightarrow h = I$$

$$\theta_i = x_i^T \beta$$

↓ Log-Likelihood

$$\log L(x_1, y_1; \dots x_n, y_n; \theta_1, \dots \theta_n)$$

$$\downarrow f_\theta(y_i)$$

$$= \sum_i (Y_i \theta_i - b(\theta_i)) / \phi + \text{const}$$

$$= \frac{1}{\phi} \sum_i (Y_i h(x_i^T \beta) - b(h(x_i^T \beta))) + \text{const}$$

↓ Canonical link $h = I$

$$\ln(L(x, \beta)) = \frac{1}{\phi} \sum_i (Y_i x_i^T \beta - b(x_i^T \beta))$$

1. Given exponential family there exist param that canonical \rightarrow invertible

2. $L(x, \beta)$ is strictly concave if $\phi > 0$

3. Gaussian $\Rightarrow \hat{\beta}^{LS} = \hat{\beta}^{ML} = (X^T X)^{-1} X^T Y$

Asymptotic Normality

(X_i, Y_i) be iid with a distribution from exponential family
 $g(u_i) = x_i^T \beta$

where $u_i = E(Y_i | X_i)$, $u_i = b(\theta_i)$ $g(u)$ is link function. θ_i be canonical parameter for each i . $\theta_i = h(x_i^T \beta^*)$

$$\downarrow \text{MC conditions}$$

$$\sqrt{n}(\hat{\beta}_n - \beta^*) \xrightarrow{d} N(0, I(\beta))$$

Hypothesis Test for Logit

$y_i \in \{0, 1\}$, $x_i \in \mathbb{R}^p$, (X_i, Y_i) independent, $Y_i \sim \text{Ber}(p_i)$

$$P(Y_i = y_i | P_i) = \begin{cases} p_i & \text{if } y_i = 1 \\ 1-p_i & \text{if } y_i = 0 \end{cases} = p_i^{y_i} (1-p_i)^{1-y_i}$$

$$= \exp(y_i \log \frac{p_i}{1-p_i} + \log(1-p_i))$$

$$= \theta_i = x_i^T \beta \Leftrightarrow p_i = \frac{e^{\beta^T x_i}}{1+e^{\beta^T x_i}}$$

$$\Rightarrow P(Y_i = y_i | x_i^T \beta) = \exp(y_i x_i^T \beta - \log(1+e^{x_i^T \beta}))$$

$$\downarrow \begin{array}{l} x_i \text{ iid distribution} \\ Y_i | x_i^T \beta \sim \text{Ber}\left(\frac{e^{x_i^T \beta}}{1+e^{x_i^T \beta}}\right) \end{array}$$

$$\ln(\beta | y, x) = \frac{1}{n} [\ln(x_i^T \beta) - \ln(1+e^{x_i^T \beta})]$$

$$\hat{\beta} = \arg \max \ln(\beta | y, x)$$

$$\text{Test : } H_0: \hat{\beta}_j = \beta_j^* \\ H_1: \hat{\beta}_j \neq \beta_j^*$$

Wald Test

$$\widehat{W}_n = n(\hat{\theta} - \theta^*) I(\hat{\theta})(\hat{\theta} - \theta^*) \xrightarrow{d} \chi_d^2 \text{ under } H_0$$

$$= -E[\hat{\beta}^T \ln(x_i^T \beta)]$$

$$\text{MLE : } \bar{n}(\hat{\theta} - \theta) \xrightarrow{d} N(0, I(\theta))$$

$$\Rightarrow n(\hat{\theta} - \theta) I(\theta)(\hat{\theta} - \theta) \xrightarrow{d} \chi_d^2$$

$$\text{Slusky : } \Rightarrow n(\hat{\theta} - \theta) I(\hat{\theta})(\hat{\theta} - \theta) \xrightarrow{d} \chi_d^2 \quad (\hat{\theta} \text{ is consistent})$$

$$\Psi = \mathbb{I}\{T_n > \chi_{\alpha}^2\}$$