

MIT Probability System Analysis and Applied Probability

Lecture 1: Probability Models and Axioms

(1) Sample space Ω : list of possible outcomes

properties: * mutually exclusive

* collectively exhaustive

categories: discrete, continuous

(2) Axioms: (1) Nonnegativity $P(A) \geq 0$

(2) Normalization: $P(\Omega) = 1$

(3) Additivity: If $A \cap B = \emptyset$, then $P(A \cup B) = P(A + B)$

(3) Discrete uniform law: $P(A) = \frac{\# \text{ of elements of } A}{\# \text{ of total samples}}$
if all outcomes are equally likely.

(4) Continuous uniform law: Probability = Area

(5) Countable additivity axiom:

if A_1, A_2, \dots are disjoint events, then:

$$P(A_1 \cup A_2 \cup \dots) = P(A_1) + P(A_2) + \dots$$

Lecture 2: Conditioning and Bayes' Rule

(1) $P(A|B)$ = Probability of A , given B

$$\begin{cases} = \frac{P(A \cap B)}{P(B)} & \text{if } P(B) \neq 0 \\ \text{undefined} & \text{if } P(B) = 0 \end{cases}$$

a) Total probability theorem

sample space are divided into A_1, A_2, \dots, A_n

B is an event in the sample space

$$P(B) = P(B|A_1)P(A_1) + \dots + P(B|A_n)P(A_n)$$

(3) Bayes' rule

* prior probabilities $P(A_i)$ - initial beliefs

* know $P(B|A_i)$ for all i 's

* wish to know $P(A_i|B)$ - revise beliefs, given B occurred

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

$A_i \xrightarrow[\text{PCB|A}]{} B$

$A_i \xleftarrow[\text{P(Ai|B)}]{} B$

Lecture 3: Independence

(1) Definition of independence : $P(A \cap B) = P(A) \cdot P(B)$

(2) Conditioning may affect independence

$$P(A \cap B|C) \stackrel{?}{=} P(A|C) P(B|C)$$

(3) Independence of a collection of events : $P(A_1 \cap A_2 \cap \dots \cap A_q) = P(A_1) \cdot P(A_2) \cdots P(A_q)$

(4) If there are A, B, C events.

$$P(A \cap C) = P(A) P(C); P(A \cap B) = P(A) P(B); P(B \cap C) = P(B) P(C)$$

called pairwise independence

pairwise independence doesn't imply independence : i.e.

$$P(A \cap B \cap C) \stackrel{?}{=} P(A) P(B) P(C)$$

(5)

Lecture 4: Counting

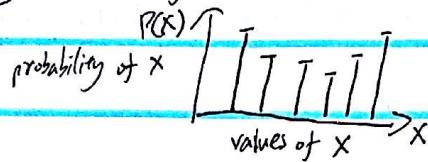
(1) Discrete uniform law

$$P(A) = \frac{\# \text{ of } A}{\text{total } \# \text{ of samples}} = \frac{|A|}{|S|}$$

Lecture 5: Discrete Random Variables; Probability Mass Functions; Expectations

(1) Random variables: functions from the sample space Ω to the real numbers. {discrete continuous}

(2) Probability mass function



(3) Expectation: $E[X] = \sum_x x P_X(x)$

(4) Properties of expectations:

* X is a r.v., $Y = g(X)$

$$E[Y] = \sum_y y P_Y(y) = \sum_x g(x) P_X(x)$$

* $E[\alpha] = \alpha$

* $E[\alpha X] = \alpha E[X]$

* $E[\alpha X + \beta] = \alpha E[X] + \beta$

* $E[g(X)] = g(E[X])$ if g is a linear function

(5) Second moment $E[X^2] = \sum_x x^2 P_X(x)$

(6) Variance: $\text{var}(X) = E[(X - E[X])^2]$

$$= \sum_x (x - E[X])^2 P_X(x) =$$

$$\begin{aligned} & \sum_x x^2 P_X(x) - \sum_x 2x E[X] P_X(x) \\ & + \sum_x E[X]^2 P_X(x) \end{aligned}$$

$$= E[X^2] - (E[X])^2 = E[X^2] - E[X]^2 = E[X^2] - 2E[X]^2 + E[X]$$

$$\text{var}(\alpha X + \beta) = \alpha^2 \text{var}(X)$$

Lecture 6 : Discrete Random Variable Examples; Joint PMFs

(1) Conditional expectation

$$E[X|A] = \sum_{\pi} \pi P_{X|A}(\pi)$$

$$\text{where } P_{X|A}(\pi) = P(X=\pi|A)$$

(2) Total Expectation theorem

$$P(B) = P(A_1)P(B|A_1) + \dots + P(A_n)P(B|A_n)$$

$$P_X(\pi) = P(A_1)P_{X|A_1}(\pi) + \dots + P(A_n)P_{X|A_n}(\pi)$$

$$E[X] = P(A_1)P_{X|A_1}(E[X|A_1]) + \dots + P(A_n)P_{X|A_n}(E[X|A_n])$$

(3) Joint PMF : $P_{X,Y}(\pi, y) = P(X=\pi \text{ and } Y=y)$

Lecture 7 : Multiple Discrete Random Variables: Expectations, Conditioning, Independence

(1) Random variables X, Y, Z are independent if

$$P_{X,Y,Z}(\pi, y, z) = p_X(\pi) \cdot p_Y(y) \cdot p_Z(z)$$

(2) if X and Y are independent random variables

$$E[XY] = \sum_{\pi} \sum_y \pi y P_{X,Y}(\pi, y)$$

$$= \sum_{\pi} \sum_y \pi y P_X(\pi) P_Y(y)$$

$$= \sum_{\pi} P_X(\pi) \sum_y y P_Y(y)$$

$$= E[X] E[Y]$$

$$E[g(x)h(y)] = \sum_{\pi} \sum_y g(\pi) h(y) P_{X,Y}(\pi, y)$$

$$= E[g(x)] E[h(y)]$$

(3) for general cases, $\text{Var}(aX) = a^2 \text{Var}(X)$

$$\text{Var}(X+a) = \text{Var}(X)$$

if X, Y are independent

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

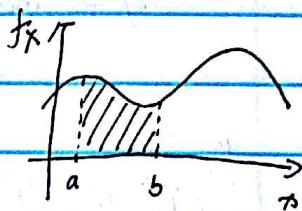
Lecture 8: Continuous Random Variables

(1) Probability density function f_X

$$P(a \leq X \leq b) = \int_a^b f_X(x) dx$$

$$\int_{-\infty}^{+\infty} f_X(x) dx = 1$$

$$P(X \in B) = \int_B f_X(x) dx$$



$$(2) E[X] = \int_{-\infty}^{+\infty} x f_X(x) dx$$

$$E[g(x)] = \int_{-\infty}^{+\infty} g(x) f_X(x) dx$$

$$\text{var}(X) = \sigma_X^2 = \int_{-\infty}^{+\infty} (x - E[X])^2 f_X(x) dx$$

(3) Cumulative distribution function (CDF)

$$F_X(x) = P(X \leq x) = \int_{-\infty}^x f_X(t) dt$$

$$\frac{dF_X(x)}{dx} = f_X(x)$$

(4) Gaussian (normal) PDF: $N(\mu, \sigma^2) : f_X(x) =$

$$N(\mu, \sigma^2) = f_X(x) = \frac{1}{\sigma \sqrt{2\pi}} e^{-(x-\mu)^2/2\sigma^2}$$

$$E[X] = \mu \quad \text{Var}[X] = \sigma^2$$

(5) No closed-form solution for CDF of normal distribution

$$\text{If } X \sim N(\mu, \sigma^2) \quad \frac{x-\mu}{\sigma} \sim N(0, 1)$$

Lecture 9: Multiple Continuous Random Variables

(1) Joint PDF: $P((X, Y) \in S) = \iint_S f_{X,Y}(x, y) dx dy$

$$P(x \leq X \leq x + \delta, y \leq Y \leq y + \delta) \approx f_{X,Y}(x, y) \cdot \delta^2$$

$$E[g(x, y)] = \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} g(x, y) f_{X,Y}(x, y) dx dy$$

$$\int_{-\infty}^{+\infty} f_{X,Y}(x, y) dx dy = f_X(x)$$

(2) X and Y are independent if $f_{X,Y}(x, y) = f_X(x) f_Y(y)$

(3) conditioning

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

if X and Y are independent: $f_{X,Y} = f_X f_Y$,
then $f_{X|Y}(x,y) = f_X(x)$

Lecture 10: Continuous Bayes' Rule; Derived Distributions

$$(1) \text{ The Bayes variations } \left\{ \begin{array}{l} P_{X|Y}(x|y) = \frac{P_{X,Y}(x,y)}{P_Y(y)} = \frac{P_X(x) P_{Y|X}(y|x)}{P_Y(y)} \\ f_{X|Y}(x|y) = \frac{f_{X,Y}(x,y)}{f_Y(y)} \end{array} \right.$$

discrete X , continuous Y

$$P_{X|Y}(x|y) = \frac{P_X(x) f_{Y|X}(y|x)}{f_Y(y)}$$

continuous X , discrete Y

$$f_{X|Y}(x|y) = \frac{f_X(x) P_{Y|X}(y|x)}{P_Y(y)}$$

(2) Derived distribution: a PMF or PDF of a function of one or more random variables with known probability law.

(3) Continuous derived distribution: cookbook procedure

* Get CDF of Y : $F_Y(y) = P(Y \leq y)$

* Differentiate to get $f_Y(y) = \frac{dF_Y}{dy}(y)$

(a) The PDF of $Y = ax + b$

$$f_Y(y) = \frac{1}{|a|} f_X\left(\frac{y-b}{a}\right)$$

Lecture 11: Derived Distributions; Convolution; Covariance and Correlation

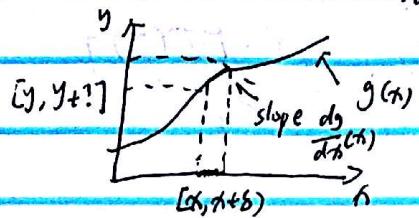
General formula for derived distributions.

(1) Let $Y = g(X)$, g strictly monotonic

Event $x \leq X \leq x + \delta \Rightarrow g(x) \leq Y \leq g(x + \delta)$

$$\Rightarrow g(x) \leq Y \leq g(x) + \delta / |(dg/dx)(x)|$$

Hence, $s f_X(x) = s f_Y(y) \left| \frac{dy}{dx}(x) \right|$



(2) Discrete distribution of $X+Y$, $W=X+Y$

$$P_W(w) = P(X+Y=w) = \sum_{\pi} P(X=\pi) P(Y=w-\pi)$$

$$= \sum_{\pi} p_X(\pi) p_Y(w-\pi)$$

(3) Continuous distribution of $X+Y$, $W=X+Y$

$$f_W(w) = \int_{-\infty}^{\infty} f_X(x) f_Y(w-x) dx$$

(4) Sum of independent normal random variables is normal distributed.

$$X_i \sim N(\mu_i, \sigma_i^2) \quad i=1 \dots n$$

$$Z = \bar{X}_i \sim N(\bar{\mu}_i, \bar{\sigma}_i^2) \quad i=1 \dots n$$

(5) Covariance: $\text{Cov}(X, Y) = E[(X - E[X]) \cdot (Y - E[Y])]$

if $\text{Cov}(X, Y) > 0$, X and Y are positive related
 if $\text{Cov}(X, Y) < 0$, X and Y are negative related.

$$\text{Cov}(X, Y) = E[XY] - E[X]E[Y]$$

$$(6) \text{Var}\left(\sum_{i=1}^n X_i\right) = \sum_{i=1}^n \text{Var}(X_i) + \sum_{(i,j): i \neq j} \text{Cov}(X_i, X_j)$$

$$\therefore E[(X_1 + \dots + X_n)^2] = E[\sum_i X_i^2 + \sum_{i \neq j} X_i X_j]$$

$$\therefore E[(\bar{X}_i - \bar{E}[X_i])^2] = E[\sum_i X_i^2 + \sum_{i \neq j} X_i X_j - 2 \sum_i X_i E[X_i] - 2 \sum_i X_i \sum_j E[X_j] + \sum_i E[X_i]^2 + \sum_{i \neq j} E[X_i] E[X_j]]$$

$$= E[\sum_i X_i^2 - 2 \sum_i X_i E[X_i] + \sum_i E[X_i]^2] + E[\sum_{i \neq j} X_i X_j - 2 \sum_i X_i \sum_j E[X_j] + \sum_{i \neq j} E[X_i] E[X_j]]$$

$$= E[(X_i - E[X_i])^2] + \sum_{i \neq j} E[X_i X_j] - 2 E[\sum_i X_i] \sum_{i \neq j} E[X_i] E[X_j]$$

$$= \sum_i \text{Var}(X_i) + \sum_{i \neq j} \text{Cov}(X_i, X_j)$$

if X, Y independent $E[XY] = E[X]E[Y]$, $\text{Cov}(X, Y) = 0$

(1) Correlation coefficient (Dimensionless version of covariance)

$$\rho = E \left[\frac{(X - E[X])}{\sigma_X} \cdot \frac{(Y - E[Y])}{\sigma_Y} \right]$$
$$= \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y}$$

$$\rho \in [-1, 1]$$

if $|\rho| = 1 \Leftrightarrow (X - E[X]) = c(Y - E[Y])$ linear related

$\rho = 0 \Leftrightarrow$ independent. converse is not true

Lecture 12: Iterated Expectations; sum of a Random Number of Random Variables.

(1) Conditional expectations

$$E[X|Y=y] = \sum_x x P_{X|Y}(x|y)$$

$$E[X|Y=y] = \int_x x f_{X|Y}(x|y)$$

(2) Law of iterated expectations:

$$E[E[X|Y]] = \sum_y E[X|Y=y] P_Y(y) = E[X] \text{ Total Expectation Thm}$$

$$(3) \text{Var}(X|Y) = E[(X - E[X|Y=y])^2 | Y=y]$$

(4) Law of total variance:

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

variability of X
between Y groups

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

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

$$\Rightarrow E[\text{Var}(X|Y)] = E[X^2] - E[(E[X|Y])^2] \quad \textcircled{1}$$

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

$$= E[E[(E[X|Y])^2]] - E[X]^2 \quad \textcircled{2}$$

$$\textcircled{1} + \textcircled{2} \Rightarrow E[\text{Var}(X|Y)] + \text{Var}(E[X|Y]) = E[X^2] - E[X]^2 = \text{Var}(X)$$

Lecture 13: Bernoulli Process

memoryless, discrete time

- (1) Bernoulli process: a sequence of independent Bernoulli trials,
at each trial, i : $P(\text{success}) = P(X_i=1) = p$
 $P(\text{failure}) = P(X_i=0) = 1-p$

- (2) Interarrival times: T_i : number of trials until first success

$$P(T_i=t) = (1-p)^{t-1} p \quad \left. \begin{array}{l} \\ \end{array} \right\} \text{geometric distribution}$$
$$E[T_i] = \frac{1}{p}$$
$$\text{Var}(T_i) = \frac{(1-p)}{p^2}$$

- (3) Time of the k -th arrival

$$P_{Y_k}(t) = P(Y_k = t) = P(k-1 \text{ arrivals in } [1, t-1], \text{ arrival at time } k)$$
$$= C_{t-1}^{k-1} p^{k-1} (1-p)^{t-k} \cdot p \Rightarrow \text{Pascal PMF}$$

$$E[Y_k] = \frac{k}{p} \Rightarrow Y_k = T_1 + T_2 + \dots + T_k$$

$$\text{Var}[Y_k] = \frac{kp(1-p)}{p^2} \quad \text{independent to each other}$$

$$\therefore E[Y_k] = \sum_{i=1}^k E[T_i] \quad \text{Var}[Y_k] = \sum_{i=1}^k \text{Var}[T_i]$$

- (4) Merging of independent Bernoulli Processes
yields a Bernoulli Process

Lecture 14: Poisson Process I

memoryless,
continuous time

- (1) Time homogeneity:

$P(k, \tau) = \text{Prob. of } k \text{ arrivals in } \tau \text{ interval of duration } \tau$

\times # of arrivals in disjoint time intervals are independent.

- (2) Small interval probabilities:

For very small δ :

$$P(k, \delta) \approx \begin{cases} 1 - \lambda \delta, & \text{if } k=0 \\ \lambda \delta, & \text{if } k=1 \\ 0, & \text{if } k>1 \end{cases}$$

where λ is "arrival rate"

(3) Approximate Poisson process by Bernoulli process

$$P(k, \tau) = \frac{(\lambda \tau)^k e^{-\lambda \tau}}{k!}, \quad k=0, 1, \dots$$

$$\left. \begin{aligned} E[N_t] &= \lambda t \\ \text{Var}[N_t] &= \lambda t \end{aligned} \right\} \text{based on Bernoulli}$$

$$\left. \begin{aligned} E[N_t] &= nP \\ \text{Var}[N_t] &= nP(1-P), \quad n=\frac{t}{\delta}, P=\lambda \delta \end{aligned} \right\} \text{let } \delta \rightarrow 0, P \rightarrow 0$$

(4) Interarrival times : time to the first arrival

* Y_k time of k th arrival

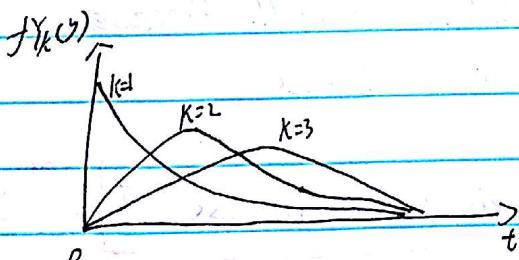
* Erlang distribution:

$$f_{Y_k}(y) = \frac{\lambda^k y^{k-1} e^{-\lambda y}}{(k-1)!}, \quad y \geq 0 \quad \text{depends on } k \text{ and } \lambda$$

$$P(Y_k \leq t) = P(t \leq Y_1 \leq t + \delta)$$

$$= P(k-1 \text{ arrivals in } [0, t]) \cdot \lambda \delta$$

$$= \frac{\lambda^{k-1} e^{-\lambda t}}{(k-1)!} \cdot \lambda \delta$$



when $k=1$

$$f_{Y_1}(y) = \lambda e^{-\lambda y}$$

exponential distribution

$$E[Y_k] = k \cdot E[Y_1] = k \cdot \int_0^\infty \lambda e^{-\lambda y} \cdot y dy$$

(5) Sum of independent Poisson random variables is Poisson.

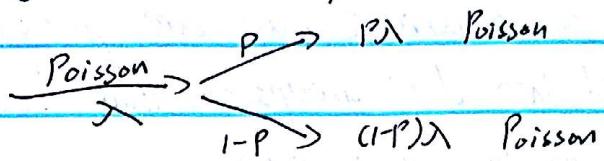
(6) Merging of independent Poisson processes is Poisson

$$\text{merged rate } \lambda = \sum_i \lambda_i$$

"Star Jima" & A. Asano

Lecture 15: Poisson Process II

(1) Splitting of Poisson process: still Poisson process



(2) Random incidence in "renewal process"

Lecture 16: Markov Chains I

(1) Finite state Markov chains

* X_n : state after n transitions

- belongs to a finite set, e.g., $\{1, \dots, m\}$

- X_0 is either given or random

* Markov property / assumption: (given current state, the past does not matter)

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

(2) State occupancy probabilities (given initial state i):

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

$$\text{key recursion: } r_{ij}(n) = \sum_{k=1}^m r_{ik}(n-1) P_{kj}$$

$$\text{with random initial state: } P(X_n = j) = \sum_{i=1}^m P(X_0 = i) r_{ij}(n)$$

(3) Recurrent and transient states

state i is recurrent if:

starting from i , and from whenever you can go, there is a way of returning to i .

If not recurrent, called transient

Lecture 17: Markov Chains II

- (1) Recurrent class: collection of recurrent states that "communicate" to each other
and to no other state
- (2) Periodic states

The states in a recurrent class are periodic if they can be grouped into $d > 1$ groups so that all transitions from one group lead to the next group.

(3) Steady-State Probabilities.

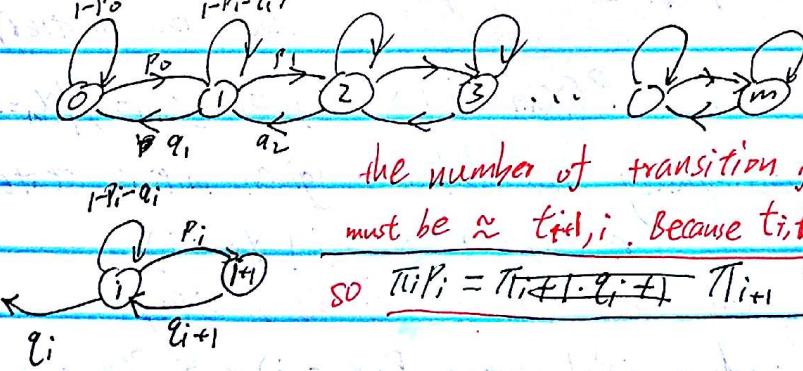
$\pi_{ij}(n)$ converge to π_{ij} (independent of the initial state)
if * recurrent states are all in a single class
and * single recurrent class is not periodic

(4) from $\pi_{ij}(n) = \sum_k \pi_{ik} (n-1) p_{kj}$

solve π_{ij} $\Rightarrow \pi_{ij} = \sum_k \pi_{ik} p_{kj}$, for all j Balance Equation

s solve π_{ij} and $\sum_j \pi_{ij} = 1$

(5) Birth-death process



the number of transition from (i) to (i+1), $t_{i,i+1}$
must be $\approx t_{i+1,i}$. Because $t_{i,i+1} = t_{i+1,i}$, or
 $\pi_i p_i = \pi_{i+1} q_{i+1}$ $t_{i,i+1} = t_{i+1,i} - 1$
or
 $t_{i,i+1} = t_{i+1,i} + 1$

Lecture 18: Markov chains III

(1) Absorption probability.

a_i is the probability that process eventually settles in transient states given the initial state i .

$$a_i = \sum_j a_{ij} P_{ij} \text{ for all other } i$$

(2) Expected time to absorption

μ_i is the expected number of transitions until reaching the absorbing state given the initial state i

$$\mu_i = 1 + \sum_j P_{ij} \mu_j$$

Lecture 19: Weak Law of Large Numbers

(1) Markov inequality (discrete variable)

$$X \geq 0, E[X] = \sum_{x \geq 0} x P_X(x) \geq \sum_{x \geq a} x P_X(x) \geq \sum_{x \geq a} a P_X(x) = a \sum_{x \geq a} P_X(x) = a P(X \geq a)$$

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

$$\mu = E[X] = P(|X - \mu| \geq a) a^2$$

(2) Chebychev's inequality (continuous variable)

$$\sigma^2 = \int (x - \mu)^2 f_X(x) dx \geq \int_{-\infty}^{\mu-c} (x - \mu)^2 f_X(x) dx + \int_{\mu+c}^{\infty} (x - \mu)^2 f_X(x) dx \geq c^2 \cdot P(|X - \mu| \geq c)$$

$$\Rightarrow P(|X - \mu| \geq c) \leq \frac{\sigma^2}{c^2}$$

$$\Rightarrow P(|X - \mu| \geq k\sigma) \leq \frac{1}{k^2}$$

(3) Convergence of the sample mean

X_1, X_2, \dots i.i.d. with finite mean μ and variance σ^2

$M_n = \frac{X_1 + \dots + X_n}{n}$ is a random variable

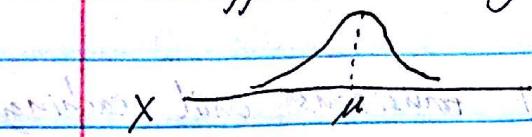
$$E[M_n] = \frac{E[X_1] + \dots + E[X_n]}{n} = \frac{n\mu}{n} = \mu$$

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

$$P(|M_n - \mu| \geq \epsilon) \leq \frac{\text{Var}(M_n)}{\epsilon^2} = \frac{\sigma^2}{n\epsilon^2}$$

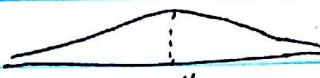
$$\text{for } \forall \epsilon > 0, P(|M_n - \mu| \geq \epsilon) \leq \frac{\sigma^2}{n\epsilon^2} \xrightarrow{n \rightarrow \infty} 0$$

(4) different scalings



differentiate by scale factors of σ

$$S_n = \sum_i x_i$$

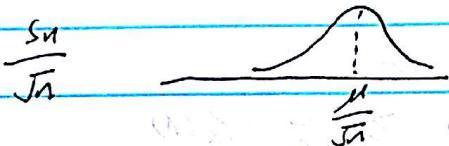


$$n\sigma$$

$$\frac{S_n}{n} = \frac{\sum_i x_i}{n}$$



$$\frac{\sigma}{\sqrt{n}}$$



$$\sigma$$

(5) The central limit theorem

"standardized": $S_n = X_1 + \dots + X_n$

$$\xrightarrow[\text{unit variance}]{\text{zero mean}} Z_n = \frac{S_n - E[S_n]}{\sigma_{S_n}} = \frac{S_n - nE[X]}{\sqrt{n}\sigma}$$

$$(S_n \sim \text{N}(0, \sigma^2) \quad \sigma^2 = n\sigma_{S_n}^2)$$

Theorem: for $\forall c \in \mathbb{R}$ $P(Z_n \leq c) \rightarrow P(\frac{Z}{\sqrt{n}} \leq c)$

Z is a standard normal r.v.

$P(Z \leq c)$ is the standard normal CDF, available from the normal tables.

$$S_n = \sqrt{n}\sigma Z_n + \mu nE[X]$$

Lecture 20: Central Limit Theorem

(1) The polster's problem using the CLT

* f : fraction of population that "..."

* i th (randomly selected) person polled:

$$X_i = \begin{cases} 1 & \text{if yes} \\ 0 & \text{if no} \end{cases}$$

$$\star M_n = (X_1 + \dots + X_n)/n$$

suppose we want $P(|M_n - f| \geq 0.01) \leq 0.05$

$$|M_n - f| = \left| \frac{X_1 + \dots + X_n - nf}{n} \right| \geq 0.01$$

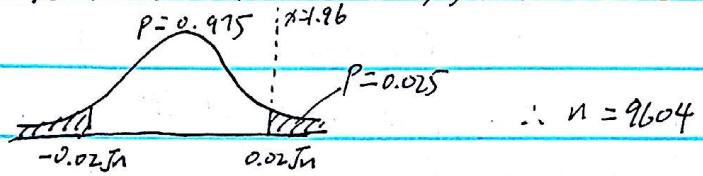
$$\left| \frac{X_1 + \dots + X_n - nf}{\sqrt{n}\sigma} \right| \geq \frac{0.01\sqrt{n}}{\sigma}$$

$$\therefore P(|M_n - f| \geq 0.01) \approx P(|Z| \geq 0.01/\sigma)$$

$$\therefore \sigma = \sqrt{f(1-f)} \quad \text{see Bernoulli variable}$$

$$\therefore \sigma \leq \frac{1}{2}$$

$\therefore P(|Z| \geq 0.01/\sigma) \leq P(|Z| \geq 0.02\sqrt{n})$, and then check table for ~~2P(Z > 0.025)~~



$$\therefore n = 9604$$

(2) CLT apply to Binomial

* Fix p , where $0 < p < 1$

* X_i : Bernoulli(p)

* $S_n = X_1 + \dots + X_n$: Binomial(n, p), with mean np ,

variance $np(1-p)$ ** Var(X) = $p(1-p)^2 + (1-p)(0-p)^2 = p(1-p)$, ~~Var(X)~~

$$\begin{aligned} \text{Var}(X_1 + \dots + X_n) \\ = n \text{Var}(X) \end{aligned}$$

* CDF of $\frac{S_n - np}{\sqrt{np(1-p)}}$ \rightarrow standard normal

(3) De Moivre-Laplace CLT (for binomial)

* using the $\frac{1}{2}$ correction for binomial approximation

e.g. $n=36$, $p=0.5$, mean np , variance $np(1-p)$

$$P(S_n = 19) = P(18.5 \leq S_n \leq 19.5)$$

$$\Leftrightarrow \frac{18.5 - np}{\sqrt{np(1-p)}} \leq \frac{S_n - np}{\sqrt{np(1-p)}} \leq \frac{19.5 - np}{\sqrt{np(1-p)}}$$

Lecture 21: Bayesian Statistical Inference I

(1) Bayesian inference: Use Bayes rule

* Hypothesis testing (discrete data)

- discrete data

$$P_{\theta|X}(\theta|x) = \frac{P_{\theta}(\theta) \cdot P_{X|\theta}(x|\theta)}{P_X(x)}$$

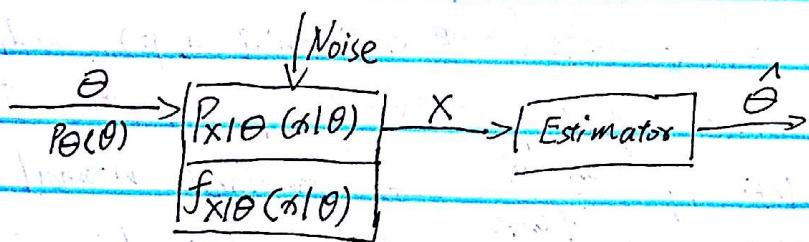
- continuous data

$$P_{\theta|X}(\theta|x) = \frac{P_{\theta}(\theta) \cdot f_{X|\theta}(x|\theta)}{f_X(x)}$$

* Estimation (continuous data)

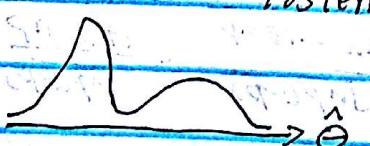
$$\text{- continuous } f_{\theta|X}(\theta|x) = \frac{f_{\theta}(\theta) f_{X|\theta}(x|\theta)}{f_X(x)}$$

$$\text{- discrete } f_{\theta|X}(\theta|x) = \frac{f_{\theta}(\theta) P_{X|\theta}(x|\theta)}{P_X(x)}$$

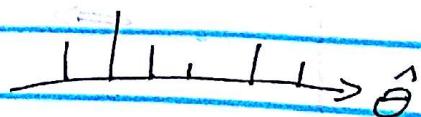


(2) * $f_{\theta}(\theta)$ or $P_{\theta}(\theta)$ is called prior distribution.

* pdf $f_{\theta|X}(\cdot|x)$ or pmf $P_{\theta|X}(\cdot|x)$ is called Posterior distribution.



or



(3) Maximum a posteriori probability (MAP):

$$P_{\theta|X}(\theta^*|x) = \max_{\theta} P_{\theta|X}(\theta|x)$$

$$f_{\theta|X}(\theta^*|x) = \max_{\theta} f_{\theta|X}(\theta|x) \text{ or}$$

Conditional expectation: $E[\theta | X = x] = \int \theta f_{\theta|X}(\theta|x) d\theta$

(4) Least Mean Squares Estimation (LMS)

Knowing distribution of θ , find estimation estimate c , to minimize $E[(\theta - c)^2]$

$$E[(\theta - c)^2] = E[\theta^2] - 2cE[\theta] + c^2$$

$$\frac{dE[(\theta - c)^2]}{dc} = -2E[\theta] + 2c = 0$$

$$\therefore c = E[\theta]$$

Optimal mean squared error: $E[(\theta - c)^2] = \text{Var}(\theta)$

(5) Conditional LMS Estimation

* Two r.v.'s θ, X

* observe that $X = x$

* $E[(\theta - c)^2 | X = x]$ is minimized by

$$c = E[\theta | X = x]$$

* ~~$E[c]$~~ assume there is another estimator $g(x)$
we can have that: $E[(\theta - E[\theta | X = x])^2]$

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

because x is random, so if take average for all x

$$E[(\theta - E[\theta | X])^2 | X] \leq E[(\theta - g(x))^2 | X]$$

take expectation of two sides

$$E[(\theta - E[\theta | X])^2] \leq E[(\theta - g(x))^2]$$

$\therefore E[\theta | X]$ minimizes $E[(\theta - g(x))^2]$ over all estimators $g(x)$

Lecture 22: Bayesian Statistical Inference II

Q) Some properties of LMS estimation

- Estimator: $\hat{\theta} = E[\theta | X]$

- Estimation error: $\hat{\theta} = \hat{\theta} - \theta$

* $E[\hat{\theta}] = 0$, $E[\hat{\theta}|x=x] = 0$

proof: $E[\hat{\theta}|x] = E[\hat{\theta} - \theta|x]$

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

$$= \hat{\theta} - \theta = 0$$

* $E[\hat{\theta} h(x)] = 0$, for any function h

proof: $E[\hat{\theta} h(x)|x] = h(x)E[\hat{\theta}|x] = 0$

take expectation to both sides: $E[\hat{\theta} h(x)] = 0$

* $\text{Cov}(\hat{\theta}, \hat{\theta}) = 0$

proof: $\text{Cov}(\hat{\theta}, \hat{\theta}) = E[\hat{\theta} \hat{\theta}] - \underbrace{E[\hat{\theta}]}_0 \underbrace{E[\hat{\theta}]}_0 = 0$

* $\text{Var}(\theta) = \text{Var}(\hat{\theta}) + \text{Var}(\tilde{\theta})$

proof: $\theta = \hat{\theta} - \tilde{\theta}$

(2) Linear LMS

$$\hat{\theta} = aX + b$$

$$\text{minimize } E[(\theta - aX - b)^2]$$

$$\hat{\theta}_L = E[\theta] + \frac{\text{Cov}(X, \theta)}{\text{Var}(X)} (X - E[X])$$

$$\text{optimal } a = \frac{\text{Cov}(X, \theta)}{\text{Var}(X)}$$

$$b = -\frac{\text{Cov}(X, \theta)}{\text{Var}(X)} [E[X] + E[\theta]]$$

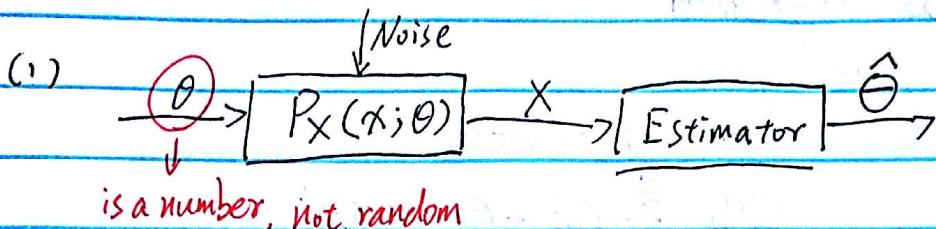
* properties:

$$E[(\hat{\theta}_L - \theta)^2] = (1 - \rho^2) \text{Var}(\theta)$$

correlation coefficient

$$\rho \text{Corr}(X, \theta)$$

Lecture 23: Classical Statistical Inference I



(2) Maximum Likelihood Estimation

* Model with unknown parameter(s):

$$X \sim P_x(x; \theta) \quad \text{or} \quad X \sim f_x(x; \theta)$$

* Pick θ that "makes data most likely"

$$\hat{\theta}_{\text{ML}} = \arg \max_{\theta} P_x(x; \theta) \quad \text{or} = \arg \max_{\theta} f_x(x; \theta)$$

* Compare to Bayesian MAP estimation:

$$\hat{\theta}_{\text{MAP}} = \arg \max_{\theta} P_{\text{prior}}(\theta) P_{\text{data}}(x|\theta)$$

$P_{\text{prior}}(\theta)$ constant
same as $\hat{\theta}_{\text{ML}}$

∴ ML estimation is MAP with uniform distribution of prior ~~prior~~ distribution.

(3) Desirable properties of estimators

* Unbiased: $E[\hat{\theta}_n] = \theta$

* Consistent: $\hat{\theta}_n \rightarrow \theta$ (converge in probability)

* small mean squared error

$$\begin{aligned} E[(\hat{\theta} - \theta)^2] &= \text{var}(\hat{\theta} - \theta) + (E[\hat{\theta} - \theta])^2 \\ &= \text{var}(\hat{\theta}) + (\text{bias})^2 \end{aligned}$$

$$\because E[X^2] = E[X]^2 + \text{Var}(X)$$

(4) Confidence intervals (CIs)

An $1-\alpha$ confidence interval is a (random) interval

$$[\hat{\theta}_L, \hat{\theta}_U], \text{ s.t. } P(\hat{\theta}_L \leq \theta \leq \hat{\theta}_U) \geq 1-\alpha, \forall \theta$$

Lecture 24: Classical Inference II

(1) Linear Regression

* Least squares

$$\text{with } Y = \theta_0 + \theta_1 X$$

$$\min_{\theta_0, \theta_1} \sum_{i=1}^n (y_i - \theta_0 - \theta_1 x_i)^2$$

* Maximum Likelihood

$$Y = \theta_0 + \theta_1 X + W; W \sim N(0, \sigma^2) \text{ i.i.d}$$

$$f_W(w) = ce^{-\frac{w^2}{2\sigma^2}}$$

$$\therefore \text{likelihood function } f_{X,Y|O}(x, y; \theta) = c \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \theta_0 - \theta_1 x_i)^2 \right\}$$

take logs, same as least squares.

Least squares \leftrightarrow pretend W i.i.d. normal

(2) Solution to $Y = \theta_0 + \theta_1 X$

$$\min_{\theta_0, \theta_1} \sum_{i=1}^n (y_i - \theta_0 - \theta_1 x_i)^2$$

$$\text{is } \hat{\theta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$\hat{\theta}_0 = \bar{y} - \hat{\theta}_1 \bar{x}$$

interpretation: $Y = \theta_0 + \theta_1 X + W$, multiply X to both sides

$$\rightarrow YX = \theta_0 X + \theta_1 X^2 + WX, \text{ take expectation to both sides}$$

$$\rightarrow E[YX] = \theta_0 E[X] + \theta_1 E[X^2] + E[W]E[X]$$

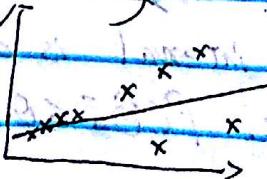
$$\text{plug in } \theta_0 = E[Y] - \theta_1 E[X]$$

$$\rightarrow E[YX] = E[X]E[Y] - \theta_1 E[X]^2 + \theta_1 E[X^2] + E[W]E[X]$$

$$E[W] = 0 \rightarrow \text{Cov}(X, Y) = \theta_1 \text{Var}(X)$$

$$\hat{\theta}_1 = \frac{\text{Cov}(X, Y)}{\text{Var}(X)}$$

(3) Heteroskedasticity



(4) Hypothesis testing

two possible hypotheses : H_0 and H_1

$$X \sim P_X(x; H_0) \text{ [or } f_X(x; H_0)]$$

or

$$X \sim P_X(x; H_1) \text{ [or } f_X(x; H_1)]$$

* Bayesian case : choose H_1 if (likelihood ratio test)

$$P(H_1 | X=x) > P(H_0 | X=x)$$

$$\Rightarrow \frac{P(X=x | H_1) P(H_1)}{P(X=x)} > \frac{P(X=x | H_0) P(H_0)}{P(X=x)}$$

$$\Rightarrow \frac{P(X=x | H_1)}{P(X=x | H_0)} > \frac{P(H_0)}{P(H_1)}$$

* Nonbayesian version: choose H_1 if

$$\frac{P(X=x; H_1)}{P(X=x; H_0)} > \xi \quad \text{or} \quad \frac{\mathbb{P} f_X(x; H_1)}{f_X(x; H_0)} > \xi$$

trade off FN and FP

Lecture 25 : Classical Inference III ; Course Overview

(1) Chi-square test :

$$\sum_i \frac{(N_i - np_i)^2}{np_i} > \xi$$

(2) Test the correctness of PDF

- partition the range into bins

- np_i : expected incidence of bin i from the PDF

- N_i : observed incidence of bin i

- use chi-square test

(3) Kolmogorov-Smirnov test :

form empirical CDF, F_x , from data

$$D_n = \max_x |F_x(a) - \hat{F}_x(a)|$$

$$P(|J_n D_n| \geq 1.36) \approx 0.05$$