

LECTURE 2

- Readings: Sections 1.3-1.4

Lecture outline

- Review
- Conditional probability
- Three **important** tools:
 - Multiplication rule
 - Total probability theorem
 - Bayes' rule

Review of probability models

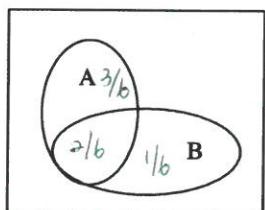
- **Sample space Ω**
 - Mutually exclusive
 - Collectively exhaustive
 - Right granularity
- **Event:** Subset of the sample space
- Allocation of probabilities to events

1. $P(A) \geq 0$
2. $P(\Omega) = 1$
3. If $A \cap B = \emptyset$,
then $P(A \cup B) = P(A) + P(B)$
- 3'. If A_1, A_2, \dots are disjoint events, then:
 $P(A_1 \cup A_2 \cup \dots) = P(A_1) + P(A_2) + \dots$

- Problem solving:
 - Specify sample space
 - Define probability law
 - Identify event of interest
 - Calculate...

* zero probability happens in continuous model !

Conditional probability



new universe
revision of belief
inference

- $P(A | B) =$ probability of A , given that B occurred
 - B is our new universe

- **Definition:** Assuming $P(B) \neq 0$,

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

$P(A | B)$ undefined if $P(B) = 0$

$$P(A \cap B) = P(B) \cdot P(A | B)$$

$$= P(A) \cdot P(B | A)$$

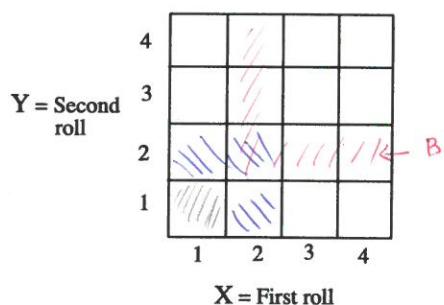
$$P(B | B) = 1$$

$B \rightarrow$ 일어났을 때 B 의 확률 | ($B \rightarrow X$, $B = \frac{1}{2}$)

$$\frac{2}{3} : \frac{1}{3}$$

총이 3이 B 는 일어나면 (1이므로)

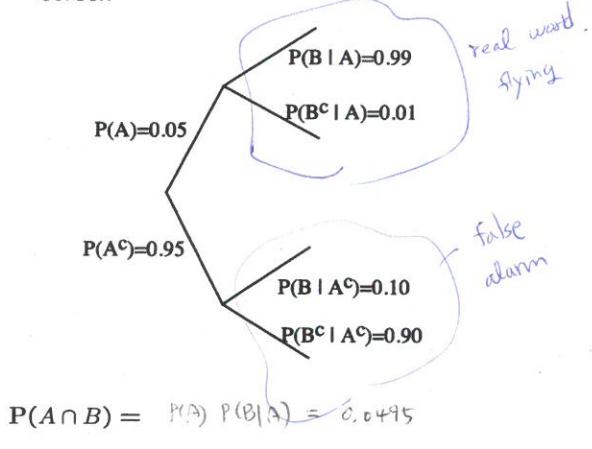
Die roll example



- Let B be the event: $\min(X, Y) = 2$
- Let $M = \max(X, Y)$
- $P(M = 1 | B) = 0$
- $P(M = 2 | B) = \frac{1}{5}$

Models based on conditional probabilities

- Event A: Airplane is flying above
- Event B: Something registers on radar screen



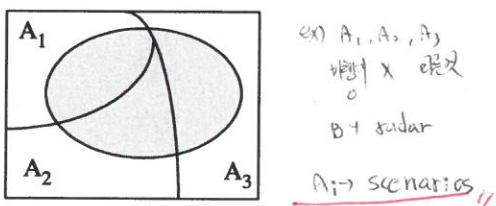
$$P(B) = 0.0495 + 0.95 \times 0.10 = 0.1445$$

$$P(A|B) = \frac{P(A \cap B)}{P(B)} = 0.34$$

(optimal) *Too low... because of false alarm happens frequently*

Total probability theorem

- Divide and conquer
- Partition of sample space into A_1, A_2, A_3
- Have $P(B|A_i)$, for every i



- One way of computing $P(B)$:

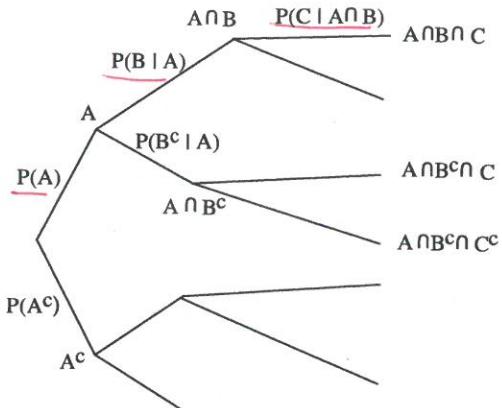
$$P(B) = P(A_1)P(B|A_1) + P(A_2)P(B|A_2) + P(A_3)P(B|A_3)$$

multiple rule로 증명함

add to 1

Multiplication rule

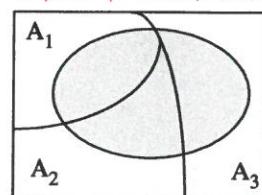
$$P(A \cap B \cap C) = P(A) \cdot P(B|A) \cdot P(C|A \cap B)$$



+ formal proof using definition of conditional prob
↳ 증명 가능

Bayes' rule

- "Prior" probabilities $P(A_i)$
 - initial "beliefs" 비행기 날고 있을 확률 ...
- We know $P(B|A_i)$ for each i
- Wish to compute $P(A_i|B)$
 - revise "beliefs", given that B occurred
비행기 날았을 때 비행기가 있을 확률 ...



$$\begin{aligned} P(A_i|B) &= \frac{P(A_i \cap B)}{P(B)} \\ &= \frac{P(A_i)P(B|A_i)}{P(B)} \\ &= \frac{P(A_i)P(B|A_i)}{\sum_j P(A_j)P(B|A_j)} \end{aligned}$$

LECTURE 3

- Readings: Section 1.5
- Review
- Independence of two events
- Independence of a collection of events

Review
reverse probability, sample space.

$$P(A | B) = \frac{P(A \cap B)}{P(B)}, \quad \text{assuming } P(B) > 0$$

- Multiplication rule:
- $$P(A \cap B) = P(B) \cdot P(A | B) = P(A) \cdot P(B | A)$$

- Total probability theorem:
- $$P(B) = P(A)P(B | A) + P(A^c)P(B | A^c)$$

- Bayes rule:

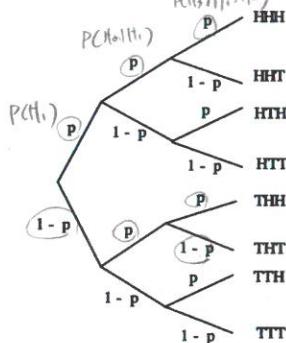
$$P(A_i | B) = \frac{P(A_i)P(B | A_i)}{P(B)}$$

cause effect
inference

Models based on conditional probabilities

- 3 tosses of a biased coin:

$$P(H) = p, P(T) = 1 - p$$



$$P(THT) = (1-p)p(1-p)$$

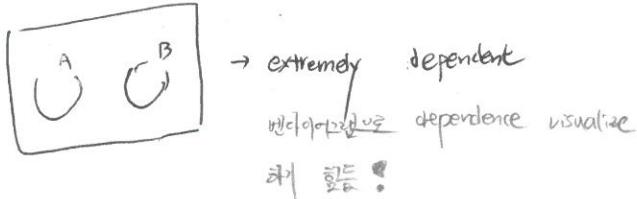
$$P(\text{1 head}) = 3p(1-p)$$

$$P(\text{first toss is H} | \text{1 head}) = \frac{P(\text{1st toss H} | \text{1 head})}{P(\text{1 head})}$$

$$= \frac{1}{3}$$

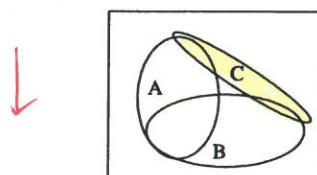
Independence of two events

- "Defn:" $P(B | A) = P(B)$
- "occurrence of A provides no information about B 's occurrence"
- initial belief is not changed though A happened
- Recall that $P(A \cap B) = P(A) \cdot P(B | A)$
- Defn: $P(A \cap B) = P(A) \cdot P(B)$
- Symmetric with respect to A and B
- applies even if $P(A) = 0$
- implies $P(A | B) = P(A)$



Conditioning may affect independence

- Conditional independence, given C , is defined as independence under probability law $P(\cdot | C)$
- Assume A and B are independent



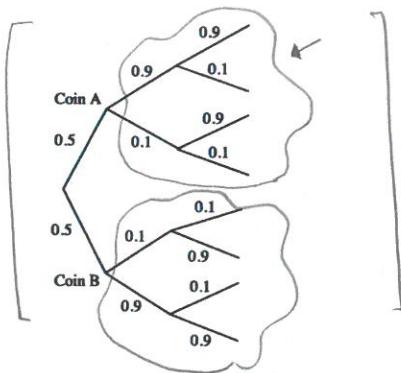
- If we are told that C occurred, are A and B independent?

$$P(A \cap B | C) = 0$$

\Rightarrow disjoint

Conditioning may affect independence

- Two unfair coins, A and B :
 $P(H \mid \text{coin } A) = 0.9, P(H \mid \text{coin } B) = 0.1$
choose either coin with equal probability



- Once we know it is coin A , are tosses independent? Yes, $(0.9, 0.9)$
- If we do not know which coin it is, are tosses independent?
 - Compare:
 $P(\text{toss 11} = H) = \frac{1}{2} \cdot 0.9 + \frac{1}{2} \cdot 0.1 = \frac{1}{2}$
 $P(\text{toss 11} = H \mid \text{first 10 tosses are heads})$
Inference - 10 heads ... $\rightarrow \text{coin } A \rightarrow \approx 0.9$

Independence of a collection of events

- Intuitive definition:
Information on some of the events tells us nothing about probabilities related to the remaining events

- E.g.: $\vdash \text{coin flip } \approx \text{ 1/2}$,

$$P(A_1 \cap (A_2^c \cup A_3) \mid A_5 \cap A_6^c) = P(A_1 \cap (A_2^c \cup A_3))$$

$\vdash 5, 6 \text{ 번째 토스 } \in \text{ 영상 } \quad \begin{matrix} 1, 2, 3 \\ 4, 5, 6 \end{matrix} \text{ 토스 } \}$

- Mathematical definition:

Events A_1, A_2, \dots, A_n

are called **independent** if:

$$P(A_i \cap A_j \cap \dots \cap A_q) = P(A_i)P(A_j) \dots P(A_q)$$

for any distinct indices i, j, \dots, q ,
(chosen from $\{1, \dots, n\}$)

X pairwise independence

$P(A_1 \cap A_2) = \sim$	different
$P(A_2 \cap A_3) = \sim$	
$P(A_1 \cap A_3) = \sim$	

$P(A_1 \cap A_2 \cap A_3 \dots)$

$\vdash \text{P(A}_1\text{)} \text{P(A}_2\text{)} \dots$

Independence vs. pairwise independence

- Two independent fair coin tosses
 - A : First toss is H
 - B : Second toss is H
 - $P(A) = P(B) = 1/2$

HH	HT
$\frac{1}{4}$	$\frac{1}{4}$
TH	TT
$\frac{1}{4}$	$\frac{1}{4}$

A
B
C

- C : First and second toss give same result

$$\begin{aligned}
 P(C) &= \frac{1}{2} \\
 P(C \cap A) &= \frac{1}{4}, \quad P(C \cap B) = \frac{1}{4} \\
 P(A \cap B \cap C) &= \frac{1}{4} \\
 P(C \mid A \cap B) &= 1 \neq \frac{1}{2}
 \end{aligned}$$

- Pairwise independence does not imply independence

$$\begin{aligned}
 P(C \cap A) &= P(C \cap B) = P(C) \cdot P(A) = P(C) \cdot P(B) \\
 \rightarrow P(C \mid A \cap B) &\neq P(C)
 \end{aligned}$$

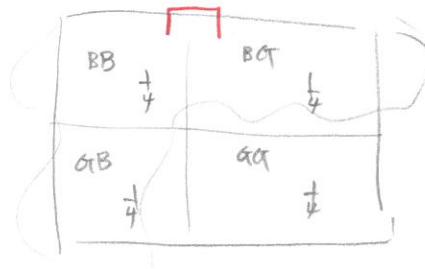
The king's sibling

- The king comes from a family of two children. What is the probability that his sibling is female?

$$\rightarrow \frac{1}{2} \text{ (naive)} \quad (\times)$$

boy female $\rightarrow \frac{2}{3}$

boy, girl $\rightarrow \frac{1}{2}$, independent.



LECTURE 4

- Readings: Section 1.6

Lecture outline

- Principles of counting
- Many examples
 - permutations
 - k -permutations
 - combinations
 - partitions
- Binomial probabilities

Discrete uniform law

- Let all sample points be equally likely

- Then,

$$P(A) = \frac{\text{number of elements of } A}{\text{total number of sample points}} = \frac{|A|}{|\Omega|}$$

- Just count...

Basic counting principle

- r stages
- n_i choices at stage i

$$n_1 = 4$$

$$n_2 = 3$$

$$n_3 = 2$$



- Number of choices is: $n_1 n_2 \cdots n_r$
- Number of license plates with 3 letters and 4 digits = $26^3 \cdot 10^4$
- ... if repetition is prohibited = $26 \cdot 25 \cdot 24 \cdot 10 \cdot 9 \cdot 8 \cdot 7$
- Permutations: Number of ways of ordering n elements is: $n \cdot (n-1) \cdots 1 = n!$
- Number of subsets of $\{1, \dots, n\} = 2^n$

Example

independent, fair

- Probability that six rolls of a six-sided die all give different numbers?
 - Number of outcomes that make the event happen: $6!$
 - Number of elements in the sample space: 6^6
 - Answer: $\frac{|A|}{|\Omega|} = \frac{6!}{6^6}$

Combinations

- $\binom{n}{k}$: number of k -element subsets of a given n -element set
- Two ways of constructing an ordered sequence of k distinct items:
 - Choose the k items one at a time:
 $n(n-1)\cdots(n-k+1) = \frac{n!}{(n-k)!}$ choices
 - Choose k items, then order them ($k!$ possible orders)
- Hence:

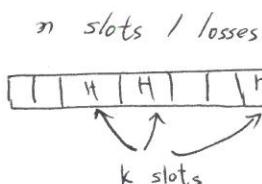
binomial coefficient \sim $\binom{n}{k} \cdot k! = \frac{n!}{(n-k)!}$

 $\sum_{k=0}^n \binom{n}{k} = 2^n$

Binomial probabilities

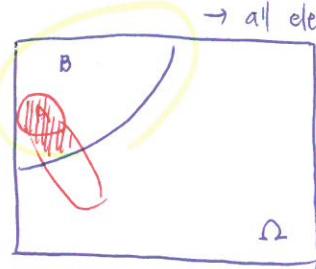
- n independent coin tosses
- $P(H) = p$
- $P(HTTHHH) = p^4(1-p)^2$
- $P(\text{sequence}) = p^{\# \text{ heads}}(1-p)^{\# \text{ tails}}$

$$\begin{aligned}
 P(k \text{ heads}) &= \sum_{\text{k-head seq.}} P(\text{seq.}) \\
 &= (\# \text{ of } k\text{-head seqs.}) \cdot p^k(1-p)^{n-k} \\
 &= \binom{n}{k} p^k (1-p)^{n-k} \quad k=0, 1, \dots, n
 \end{aligned}$$



Coin tossing problem

- event B : 3 out of 10 tosses were "heads".
- Given that B occurred, what is the (conditional) probability that the first 2 tosses were heads?
- All outcomes in set B are equally likely: probability $p^3(1-p)^7$
- Conditional probability law is uniform
- Number of outcomes in B : $10C_3$ or $\binom{10}{3}$
- Out of the outcomes in B , how many start with HH? 8



$$\frac{|A \cap B|}{|B|} = \left(\binom{10}{3} \right)^{-1}$$

* text book 27

Partitions

distribute ...

combinations

- 52-card deck, dealt to 4 players
- Find $P(\text{each gets an ace})$
- Outcome: a partition of the 52 cards
- number of outcomes: $\frac{52!}{13! 13! 13! 13!} \binom{52}{13} \binom{39}{13} \binom{26}{13} \binom{13}{13}$
- Count number of ways of distributing the four aces: $4 \cdot 3 \cdot 2$
- Count number of ways of dealing the remaining 48 cards

$$\frac{48!}{12! 12! 12! 12!}$$

- Answer:

$$\frac{4 \cdot 3 \cdot 2}{12! 12! 12! 12!} \frac{48!}{52!} \frac{1}{13! 13! 13! 13!}$$

more generally ...

$$(\binom{k}{n})^{n-k}$$

9 * $\frac{n!}{n_1! n_2! \cdots n_k!}$

LECTURE 5

- **Readings:** Sections 2.1-2.3, start 2.4

Lecture outline

- Random variables
 - Probability mass function (PMF)
 - Expectation
 - Variance

Random variables

- An assignment of a value (number) to every possible outcome

- Mathematically: A function from the sample space Ω to the real numbers
 - discrete or continuous values

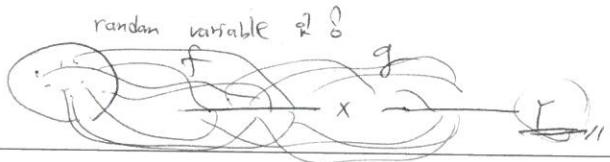
- Can have several random variables defined on the same sample space

height, weight
ex) \emptyset cm, m, feet
 oz kg

- Notation:

- random variable X function $\Omega \rightarrow \mathbb{R}$
 - numerical value $x \in \mathbb{R}$
number

+ random variable or random variable Σ



Probability mass function (PMF)

- (“probability law”,
“probability distribution” of X)
 - Notation:

$$\begin{aligned} p_X(x) &= \mathbf{P}(X = x) \\ &= \mathbf{P}(\{\omega \in \Omega \text{ s.t. } X(\omega) = x\}) \end{aligned}$$

- $p_X(x) \geq 0$ $\sum_x p_X(x) = 1$
 - **Example:** X =number of coin tosses until first head

$$\begin{aligned} p_X(k) &= \mathbf{P}(X = k) \\ &= \mathbf{P}(TT \cdots TH) \\ &= (1-p)^{k-1}p, \quad k = 1, 2, \dots \end{aligned}$$

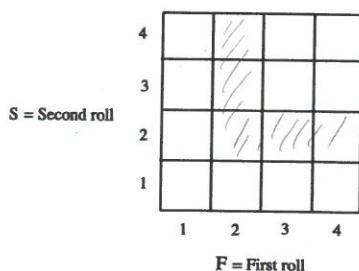
- geometric PMF

How to compute a PMF $p_X(x)$

- collect all possible outcomes for which X is equal to x
 - add their probabilities
 - repeat for all x

- **Example:** Two independent rolls of a fair tetrahedral die

F : outcome of first throw
 S : outcome of second throw
 $X = \min(F, S)$

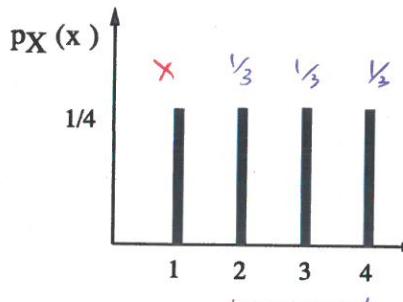


$$p_X(2) = \frac{5}{16}$$

X random var
 x value

Conditional PMF and expectation

- $p_{X|A}(x) = P(X = x | A)$
- new notation (definition) ↗
- $E[X | A] = \sum_x x p_{X|A}(x)$



- Let $A = \{X \geq 2\}$ new universe //

$$p_{X|A}(x) = \frac{1}{3} \quad (x=2, 3, 4)$$

$$E[X | A] = 3$$

$$E[g(x) | A] = \sum_x g(x) p_{X|A}(x)$$

example of random variables

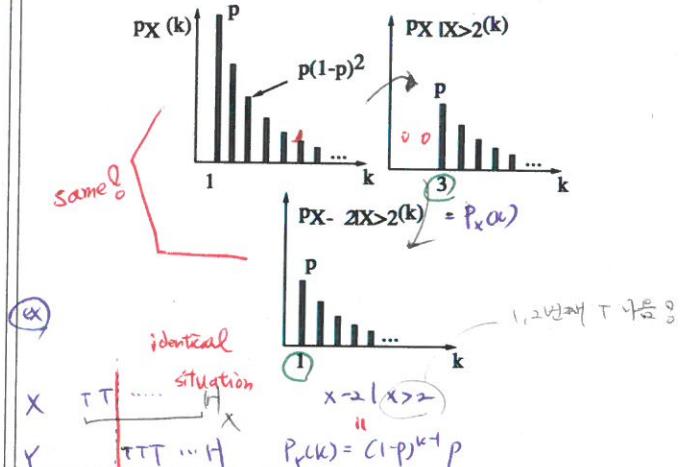
Geometric PMF

- X : number of independent coin tosses until first head

$$p_X(k) = (1-p)^{k-1}p, \quad k = 1, 2, \dots$$

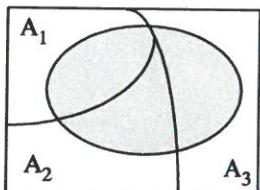
$$E[X] = \sum_{k=1}^{\infty} kp_X(k) = \sum_{k=1}^{\infty} k(1-p)^{k-1}p$$

- Memoryless property: Given that $X > 2$, the r.v. $X - 2$ has same geometric PMF



Total Expectation theorem

- Partition of sample space into disjoint events A_1, A_2, \dots, A_n



$$\begin{aligned} P(B) &= P(A_1)P(B | A_1) + \dots + P(A_n)P(B | A_n) \\ p_X(x) &= P(A_1)p_{X|A_1}(x) + \dots + P(A_n)p_{X|A_n}(x) \\ E[X] &= P(A_1)E[X | A_1] + \dots + P(A_n)E[X | A_n] \end{aligned}$$

- Geometric example:

$$A_1 : \{X = 1\}, \quad A_2 : \{X > 1\}$$

$$E[X] = P(X = 1)E[X | X = 1] + P(X > 1)E[X | X > 1]$$

- Solve to get $E[X] = 1/p$

$$\begin{aligned} E[X | X > 1] &= E[X-1 | X-1] + \\ &\quad \vdots \\ &= E[X] + 1 \end{aligned}$$

Joint PMFs

relate two random variables

- $p_{X,Y}(x, y) = P(X = x \text{ and } Y = y)$

	y			
	1	2	3	4
x	1/20	2/20	2/20	
1		1/20	3/20	1/20
2			1/20	
3				1/20
4				1/20

ex) $p_{X,Y}(2,3) = 3/20$
 $P_X(3) = 6/20$
 $y=2$

$$\sum_x \sum_y p_{X,Y}(x, y) = 1$$

$$p_X(x) = \sum_y p_{X,Y}(x, y) \quad \text{ex) } y=2 \text{ ↗ } \frac{1}{20} + \frac{3}{20} + \frac{1}{20} = \frac{5}{20}$$

$$p_{X|Y}(x | y) = P(X = x | Y = y) = \frac{p_{X,Y}(x, y)}{p_Y(y)}$$

$$\sum_x p_{X|Y}(x | y) = 1$$

$$\checkmark \quad 0 \quad \frac{1}{5} \quad \frac{3}{5} \quad \frac{1}{5} \quad (1:3:1 \text{ ratio})$$

Binomial PMF

- X : number of heads in n independent coin tosses
- $P(H) = p$
- Let $n = 4$

$$\begin{aligned} p_X(2) &= P(HHTT) + P(HTHT) + P(HTTH) \\ &\quad + P(THHT) + P(THTH) + P(TTHH) \\ &= 6p^2(1-p)^2 \\ &= \binom{4}{2}p^2(1-p)^2 \end{aligned}$$

In general:

$$p_X(k) = \binom{n}{k} p^k (1-p)^{n-k}, \quad k = 0, 1, \dots, n$$

(그림)



Expectation

- Definition:

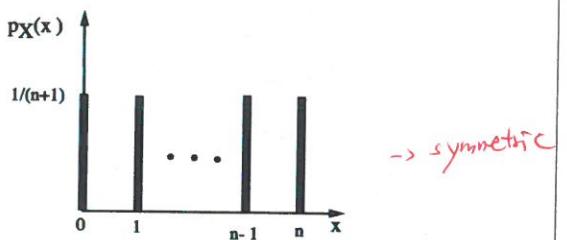
$$E[X] = \sum_x x p_X(x)$$

(평균의 정의)
(평균을 찾는다)

- Interpretations:

- Center of gravity of PMF
- Average in large number of repetitions of the experiment
(to be substantiated later in this course)

- Example: Uniform on $0, 1, \dots, n$



$$E[X] = 0 \times \frac{1}{n+1} + 1 \times \frac{1}{n+1} + \dots + n \times \frac{1}{n+1} = \frac{n}{2}$$

(by intuition)
gravity

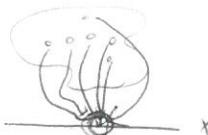
Properties of expectations

- Let X be a r.v. and let $Y = g(X)$
 - Hard: $E[Y] = \sum_y y p_Y(y)$
 - Easy: $E[Y] = \sum_x g(x) p_X(x)$ all possible x that happens
- Caution: In general, $E[g(X)] \neq g(E[X])$

= if g is linear

Properties: If α, β are constants, then:

- $E[\alpha] = \alpha$, α is a random variable
- $E[\alpha X] = \alpha E[X]$
- $E[\alpha X + \beta] = \alpha E[X] + \beta$



Variance

$$\text{Recall: } E[g(X)] = \sum_x g(x) p_X(x)$$

- Second moment: $E[X^2] = \sum_x x^2 p_X(x)$

Variance

$$\begin{aligned} \text{var}(X) &= E[(X - E[X])^2] \\ &= \sum_x (x - E[X])^2 p_X(x) \\ &= E[X^2] - (E[X])^2 \end{aligned}$$

Properties:

- $\text{var}(X) \geq 0$
- $\text{var}(\alpha X + \beta) = \alpha^2 \text{var}(X)$

(intuition은 어떤가요?)

$$g(x) = (x - E[x])^2$$

notation, practice

LECTURE 6

- Readings: Sections 2.4-2.6

Lecture outline

- Review: PMF, expectation, variance
- Conditional PMF
- Geometric PMF
- Total expectation theorem
- Joint PMF of two random variables

Review

- Random variable X : function from sample space to the real numbers
- PMF (for discrete random variables): $p_X(x) = P(X = x)$ numerical value ↗ how likely it's going to happen ↘ how likely they are
- Expectation: $E[X] = \sum_x x p_X(x)$

$$E[g(X)] = \sum_x g(x)p_X(x)$$

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

how far $\underbrace{X - E[X]}_{\text{random variable}}$

$$\bullet E[X - E[X]] = E[X] - E[E[X]] = 0$$

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

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

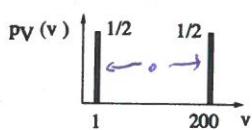
$$= E[X^2] - (E[X])^2$$

Standard deviation: $\sigma_X = \sqrt{\text{var}(X)}$

standard deviation
(%)

Random speed

- Traverse a 200 mile distance at constant but random speed V



$$\bullet d = 200, T = t(V) = 200/V$$

$$\bullet E[V] = \frac{1}{2} \cdot 1 + \frac{1}{2} \cdot 200 = 100.5$$

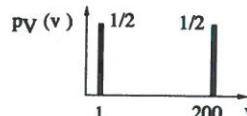
$$\bullet \text{var}(V) = \frac{1}{2} (1 - 100.5)^2 + \frac{1}{2} (200 - 100.5)^2 \approx 100^2$$

$$\bullet \sigma_V = 100$$

Cannot reason average on Non-linear things.

Average speed vs. average time

- Traverse a 200 mile distance at constant but random speed V



$$\bullet \text{time in hours} = T = t(V) =$$

$$\bullet E[T] = E[t(V)] = \sum_v t(v)p_V(v) = \frac{1}{2} \cdot 200 + \frac{1}{2} \cdot 1 \approx 100.5$$

$$\bullet E[TV] = 200 \neq E[T] \cdot E[V]$$

$$\bullet E[200/V] = E[T] \neq 200/E[V].$$

$\approx 100.5 \quad \approx 2$

expected time

LECTURE 7

- Readings: Finish Chapter 2

Lecture outline

- Multiple random variables
 - Joint PMF
 - Conditioning
 - Independence
- More on expectations
- Binomial distribution revisited
- A hat problem

discrete random variable

Review : notation *

$$p_X(x) = P(X = x)$$

$$p_{X,Y}(x,y) = P(X = x, Y = y)$$

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

condition that y is fixed value

$$p_X(x) = \sum_y p_{X,Y}(x,y)$$

$$p_{X,Y}(x,y) = p_X(x)p_{Y|X}(y | x)$$

$$\sum_x p_{X,Y}(x,y) = 1$$

$$* P(A \cap B) = P(A) P(B | A)$$

$$\{X=x\} \quad \{Y=Y\}$$

Independent random variables

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

\downarrow rewrite,

$$p_{X,Y,Z}(x,y,z) = p_X(x)p_{Y|X}(y | x)p_{Z|X,Y}(z | x,y)$$

$$\bullet P_X(x) = \sum_{y,z} p_{X,Y,Z}(y,z) : \text{marginal pmf of } X$$

- Random variables X, Y, Z are independent if:

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

for all x, y, z

$$\text{or } P_{X,Y}(x,y) = P_X(x)$$

$$\text{if } P_Y(y) > 0$$

4	1/20	2/20	2/20	
3	2/20	4/20	1/20	
2		1/20	3/20	
1		1/20		
	1	2	3	4

- Independent? No (by intuition)

- What if we condition on $X \leq 2$ and $Y \geq 3$?

Yes.

$$\begin{array}{c|cc} & \frac{1}{9} & \frac{2}{9} \\ \hline \frac{2}{9} & \frac{4}{9} \end{array} \quad P_{X|Y \geq 3, X \leq 2}(x | y) = P_{X|A}(x)$$

Expectations

$$E[X] = \sum_x x p_X(x)$$

$$E[g(X, Y)] = \sum_x \sum_y g(x, y) p_{X,Y}(x, y)$$

- In general: $E[g(X, Y)] \neq g(E[X], E[Y])$

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

$$\bullet E[X + Y + Z] = E[X] + E[Y] + E[Z]$$

- If X, Y are independent: $\rightarrow \star \star$
 $g(x) \cdot h(Y)$ also independent

$$- E[XY] = E[X]E[Y]$$

$$- E[g(X)h(Y)] = E[g(X)] \cdot E[h(Y)]$$

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

if X cannot tell anything about Y , so does $h(Y)$.

Variances

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

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

- Let $Z = X + Y$.

If X, Y are independent:

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

- Examples: $\underline{\text{What}}$

- If $X = Y$, $\text{Var}(X + Y) = 4\text{Var}(X)$

- If $X = -Y$, $\text{Var}(X + Y) = 0$

- If X, Y indep., and $Z = X - 3Y$,
 $\text{Var}(Z) = \text{Var}(X - 3Y)$

extremely
independent

$$= \text{Var}(X) + 9\text{Var}(Y)$$

Binomial mean and variance

- $X = \#$ of successes in n independent trials

— probability of success p

$$E[X] = \sum_{k=0}^n k \binom{n}{k} p^k (1-p)^{n-k}$$

↑ indicator ↓ binomial pmf

- $X_i = \begin{cases} 1, & \text{if success in trial } i, \\ 0, & \text{otherwise} \end{cases}$

$$\sum_i X_i = X$$

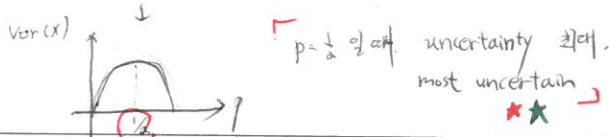
- $E[X_i] = 1 \cdot p + 0 \cdot (1-p) = p$

- $E[X] = np$

- $\text{Var}(X_i) = p(1-p)^2 + (1-p)(0-p)^2$ (by def)

$\text{or } E[X_i^2] - E[X_i]^2 = p - p^2 = p(1-p)$

- $\text{Var}(X) = np(1-p)$



The hat problem

- n people throw their hats in a box and then pick one at random.

- X : number of people who get their own hat

- Find $E[X]$

$$\frac{1}{n}$$

$$X_i = \begin{cases} 1, & \text{if } i \text{ selects own hat} \\ 0, & \text{otherwise.} \end{cases}$$

- $X = X_1 + X_2 + \dots + X_n$

- $P(X_i = 1) = \frac{1}{n}$

- $E[X_i] = 1 \cdot \frac{1}{n} + 0 \cdot (1 - \frac{1}{n}) = \frac{1}{n}$

- Are the X_i independent? No (why?) ok.

- $E[X] = \sum_i E[X_i] = n \cdot \frac{1}{n} = 1$

Variance in the hat problem

- $\text{Var}(X) = E[X^2] - (E[X])^2 = E[X^2] - 1$

$$X^2 = \sum_i X_i^2 + \sum_{i,j: i \neq j} X_i X_j$$

- $E[X_i^2] = \frac{1}{n}$

↳ going to be same
all i,j

$$P(X_1 X_2 = 1) = P(X_1 = 1) \cdot P(X_2 = 1 | X_1 = 1)$$

not independent
why?

$$= \frac{1}{n} \cdot \frac{1}{n-1}$$

- $E[X^2] = n \cdot \frac{1}{n} + \frac{1}{n} \cdot \frac{1}{n-1} \cdot \frac{n \cdot (n-1)}{n} = \frac{3+1}{2} = 2$

- $\text{Var}(X) =$

LECTURE 8

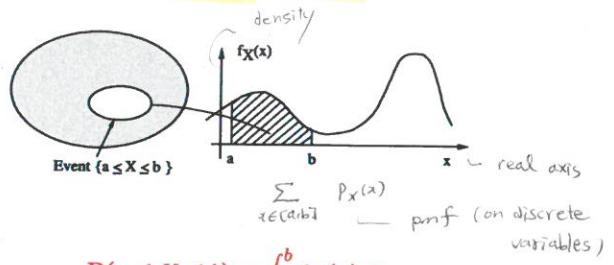
- Readings: Sections 3.1-3.3

Lecture outline

- Probability density functions
- Cumulative distribution functions
- Normal random variables

Continuous r.v.'s and pdf's

- A continuous r.v. is described by a probability density function f_X



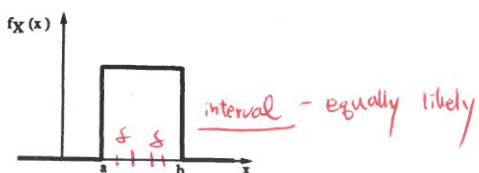
$$\int_{-\infty}^{\infty} f_X(x) dx = 1$$

$$P(x \leq X \leq x + \delta) = \int_x^{x+\delta} f_X(s) ds \approx f_X(x) \cdot \delta$$

$$P(X \in B) = \int_B f_X(x) dx, \text{ for "nice" sets } B$$

Means and variances

- $E[X] = \int_{-\infty}^{\infty} x f_X(x) dx$
- $E[g(X)] = \int_{-\infty}^{\infty} g(x) f_X(x) dx$
- $\text{var}(X) = \sigma_X^2 = \int_{-\infty}^{\infty} (x - E[X])^2 f_X(x) dx = E[X^2] - (E[X])^2$
- Continuous Uniform r.v.



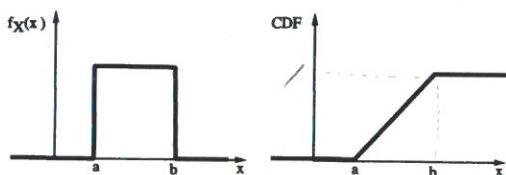
- $f_X(x) = \frac{1}{b-a}$ $a \leq x \leq b$
- $E[X] = \int_a^b x \frac{1}{b-a} dx = \frac{a+b}{2}$ or by intuition $\frac{a+b}{2}$
- $\sigma_X^2 = \int_a^b \left(x - \frac{a+b}{2}\right)^2 \frac{1}{b-a} dx = \frac{(b-a)^2}{12}$

$$\sigma_X = \sqrt{\frac{(b-a)^2}{12}}$$

Cumulative distribution function (CDF)

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

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

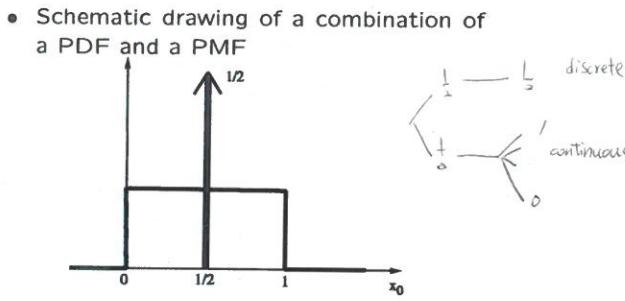


- Also for discrete r.v.'s:

$$F_X(x) = P(X \leq x) = \sum_{k \leq x} p_X(k)$$

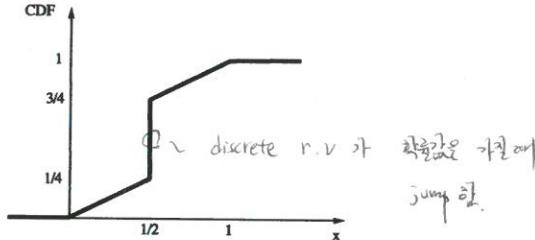


Mixed distributions

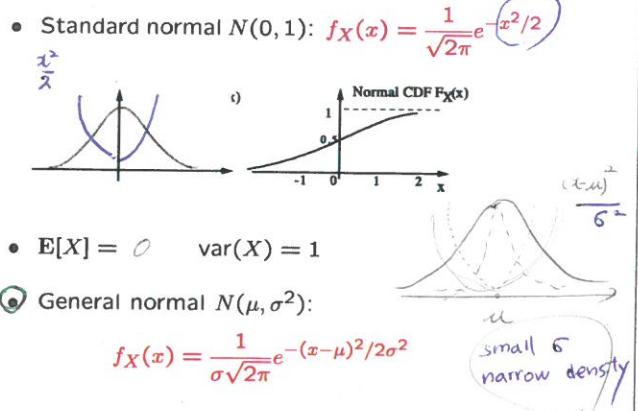


- The corresponding CDF:

$$F_X(x) = P(X \leq x)$$



Gaussian (normal) PDF



- It turns out that:
 $E[X] = \mu$ and $\text{Var}(X) = \sigma^2$.

- Let $Y = aX + b$
 - Then: $E[Y] = a\mu + b$ and $\text{Var}(Y) = a^2\sigma^2$
 - Fact: $Y \sim N(a\mu + b, a^2\sigma^2)$
- X is normal 이면 Y 도 normal

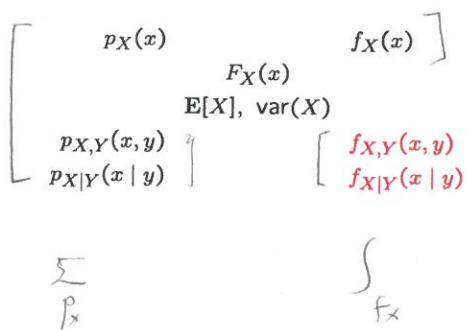
Calculating normal probabilities

- No closed form available for CDF
 - but there are tables (for standard normal)
- If $X \sim N(\mu, \sigma^2)$, then $\frac{X - \mu}{\sigma} \sim N(0, 1)$
- If $X \sim N(2, 16)$:
 $P(X \leq 3) = P\left(\frac{X - 2}{4} \leq \frac{3 - 2}{4}\right) = \text{CDF}(0.25)$

$\rightarrow N(0, 1)$

	.00	.01	.02	.03	.04	.05	.06	.07	.08	.09
0.0	.5000	.5040	.5080	.5120	.5160	.5199	.5239	.5279	.5319	.5350
0.1	.5398	.5438	.5478	.5517	.5557	.5596	.5636	.5675	.5714	.5753
0.2	.5793	.5832	.5871	.5910	.5948	.5987	.6026	.6064	.6103	.6141
0.3	.6179	.6217	.6255	.6293	.6331	.6368	.6406	.6443	.6480	.6517
0.4	.6554	.6591	.6628	.6664	.6700	.6736	.6772	.6808	.6844	.6879
0.5	.6915	.6950	.6985	.7019	.7054	.7088	.7123	.7157	.7190	.7224
0.6	.7257	.7291	.7324	.7357	.7389	.7422	.7454	.7486	.7517	.7549
0.7	.7580	.7611	.7642	.7673	.7704	.7734	.7764	.7794	.7823	.7852
0.8	.7881	.7910	.7939	.7967	.7995	.8023	.8051	.8078	.8106	.8133
0.9	.8159	.8184	.8212	.8238	.8264	.8289	.8315	.8340	.8365	.8389
1.0	.8413	.8436	.8461	.8485	.8508	.8531	.8554	.8577	.8599	.8621
1.1	.8643	.8665	.8686	.8708	.8729	.8749	.8770	.8790	.8810	.8830
1.2	.8849	.8869	.8888	.8907	.8925	.8944	.8962	.8980	.8997	.9015
1.3	.9032	.9049	.9066	.9082	.9099	.9115	.9131	.9147	.9162	.9177
1.4	.9192	.9207	.9222	.9236	.9251	.9265	.9279	.9292	.9306	.9319
1.5	.9332	.9345	.9357	.9370	.9382	.9394	.9406	.9418	.9429	.9441
1.6	.9452	.9463	.9474	.9484	.9495	.9505	.9515	.9525	.9535	.9545
1.7	.9554	.9564	.9573	.9582	.9591	.9599	.9608	.9616	.9625	.9633
1.8	.9641	.9649	.9656	.9664	.9671	.9678	.9686	.9693	.9699	.9706
1.9	.9713	.9719	.9726	.9732	.9738	.9744	.9750	.9756	.9761	.9767
2.0	.9772	.9778	.9783	.9788	.9793	.9798	.9803	.9808	.9812	.9817

The constellation of concepts



LECTURE 9

- Readings: Sections 3.4-3.5

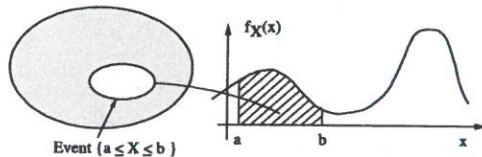
Outline

- PDF review
- Multiple random variables
 - conditioning
 - independence
- Examples

Summary of concepts

$p_X(x)$	$f_X(x)$
$F_X(x)$	
$\sum_x x p_X(x)$	$E[X] = \int x f_X(x) dx$
$\text{var}(X)$	
$p_{X,Y}(x,y)$	$f_{X,Y}(x,y)$ given
$p_{X A}(x)$	$f_{X A}(x)$ & certain event
$p_{X Y}(x y)$	$f_{X Y}(x y)$ (A) occurred

Continuous r.v.'s and pdf's



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

- $P(x \leq X \leq x + \delta) \approx f_X(x) \cdot \delta$ unit length
- $E[g(X)] = \int_{-\infty}^{\infty} g(x) f_X(x) dx$

* formal property
 $\lim_{\delta \rightarrow 0}, \int_{-\infty}^{\infty} f(x) dx = 1, P(X = a) = 0$

volume?
 volume? $\int_S f_{X,Y}(x,y) dx dy$

Joint PDF $f_{X,Y}(x,y)$
 def of joint pdf

$$P((X,Y) \in S) = \int_S f_{X,Y}(x,y) dx dy$$

- Interpretation:

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

probability for unit area
(δ^2)

- Expectations:

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

$\Sigma \Sigma$

- From the joint to the marginal:

$$f_X(x) \cdot \delta \approx P(x \leq X \leq x + \delta) =$$

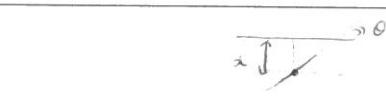
$$\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f_{X,Y}(\bar{x}, y) \cdot \delta f_{X,Y}(x,y) dy$$

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

- X and Y are called independent if

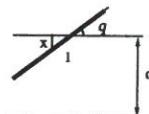
$$f_{X,Y}(x,y) = f_X(x)f_Y(y), \quad \text{for all } x, y$$

$$\Rightarrow P(X \in A, Y \in B) = P(X \in A) P(Y \in B)$$



Buffon's needle

- Parallel lines at distance d
- Needle of length ℓ (assume $\ell < d$)
- Find $P(\text{needle intersects one of the lines})$



- $X \in [0, d/2]$: distance of needle midpoint to nearest line

• Model: X, Θ uniform, independent

$$f_{X,\Theta}(x, \theta) = \begin{cases} \frac{1}{d/2} & 0 \leq x \leq d/2, 0 \leq \theta \leq \pi/2 \\ 0 & \text{otherwise} \end{cases}$$

$$f_X(x) f_\Theta(\theta) = \frac{2}{d} \cdot \frac{1}{\pi/2}$$

- Intersect if $X \leq \frac{\ell}{2} \sin \theta$
 (the line)

$$P\left(X \leq \frac{\ell}{2} \sin \theta\right) = \int \int_{x \leq \frac{\ell}{2} \sin \theta} f_X(x) f_\Theta(\theta) dx d\theta$$

$$= \frac{4}{\pi d} \int_0^{\pi/2} \int_0^{(\ell/2) \sin \theta} dx d\theta$$

$$= \frac{4}{\pi d} \int_0^{\pi/2} \frac{\ell}{2} \sin \theta d\theta = \frac{2\ell}{\pi d}$$

π 페인팅 방법으로도 예상

Conditioning

- Recall

$$P(x \leq X \leq x + \delta) \approx f_X(x) \cdot \delta$$

- By analogy, would like: $P(x \leq X \leq x + \delta | Y \approx y) \approx f_{X|Y}(x | y) \cdot \delta$ formal notation

- This leads us to the **definition**:

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

- For given y , conditional PDF is a (normalized) "section" of the joint PDF

- If independent, $f_{X,Y} = f_X f_Y$, we obtain

$$f_{X|Y}(x | y) = f_X(x)$$

Joint, Marginal and Conditional Densities

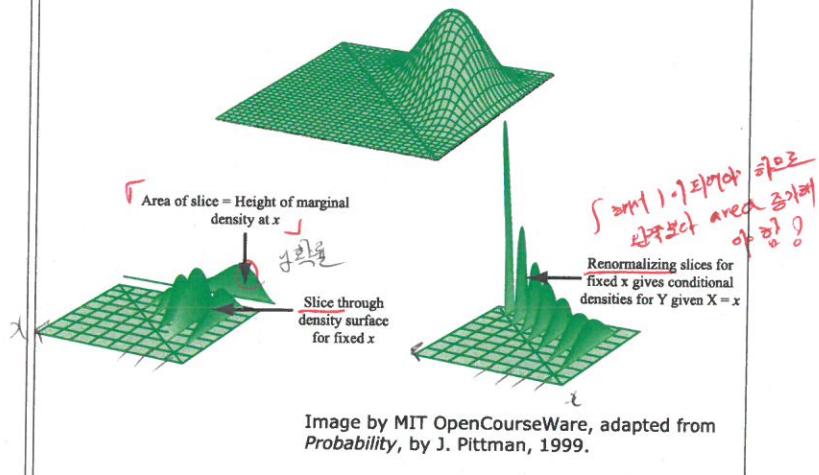
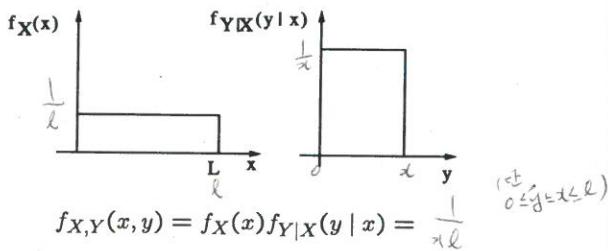


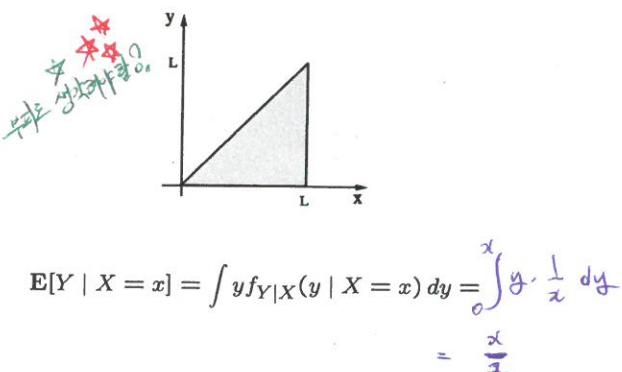
Image by MIT OpenCourseWare, adapted from Probability, by J. Pitman, 1999.

Stick-breaking example

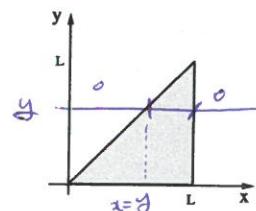
- Break a stick of length ℓ twice:
break at X : uniform in $[0, \ell]$; $0 \leq Y \leq X \leq \ell$
break again at Y , uniform in $[0, X]$



on the set:



$$f_{X,Y}(x, y) = \frac{1}{\ell x}, \quad 0 \leq y \leq x \leq \ell$$



$$\begin{aligned} f_Y(y) &= \int f_{X,Y}(x, y) dx \\ &= \int_y^\ell \frac{1}{\ell x} dx \\ &= \frac{1}{\ell} \log \frac{\ell}{y}, \quad 0 \leq y \leq \ell \end{aligned}$$

$$E[Y] = \int_0^\ell y f_Y(y) dy = \int_0^\ell y \frac{1}{\ell} \log \frac{\ell}{y} dy = \frac{\ell}{4}$$

LECTURE 10

Continuous Bayes rule; Derived distributions

- Readings:

Section 3.6; start Section 4.1

Review

$$p_{X|Y}(x|y) = \frac{p_{X,Y}(x,y)}{p_Y(y)}$$

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

$$p_X(x) = \sum_y p_{X,Y}(x,y) \quad f_X(x) = \int_{-\infty}^{\infty} f_{X,Y}(x,y) dy$$

$$F_X(x) = P(X \leq x)$$

$$E[X], \text{ var}(X)$$

$f \rightarrow \text{density}$

\curvearrowright continuous



The Bayes variations

$$p_{X|Y}(x|y) = \frac{p_{X,Y}(x,y)}{p_Y(y)} = \frac{p_X(x)p_{Y|X}(y|x)}{p_Y(y)}$$

$$p_Y(y) = \sum_x p_X(x)p_{Y|X}(y|x)$$

Example:

- $X = 1, 0$: airplane present/not present
- $Y = 1, 0$: something did/did not register on radar

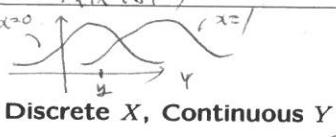
Continuous counterpart

$$f_{X|Y}(x|y) = \frac{f_{X,Y}(x,y)}{f_Y(y)} = \frac{f_X(x)f_{Y|X}(y|x)}{f_Y(y)}$$

$$f_Y(y) = \int_x f_X(x)f_{Y|X}(y|x) dx$$

Example: X : some signal; "prior" $f_X(x)$
 Y : noisy version of X measurement
 $f_{Y|X}(y|x)$: model of the noise

$X=0, 1$
 $Y=x+\eta$
 Gaussian noise



$$p_{X|Y}(x|y) = \frac{p_X(x)f_{Y|X}(y|x)}{f_Y(y)}$$

$$f_Y(y) = \sum_x p_X(x)f_{Y|X}(y|x)$$

Example:

- X : a discrete signal; "prior" $p_X(x)$
- Y : noisy version of X
- $f_{Y|X}(y|x)$: continuous noise model

Continuous X , Discrete Y

$$f_{X|Y}(x|y) = \frac{f_X(x)p_{Y|X}(y|x)}{p_Y(y)}$$

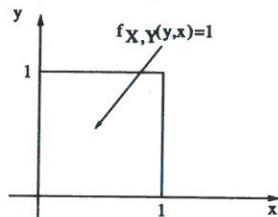
$$p_Y(y) = \int_x f_X(x)p_{Y|X}(y|x) dx$$

Example:

- X : a continuous signal; "prior" $f_X(x)$ (e.g., intensity of light beam);
- Y : discrete r.v. affected by X (e.g., photon count)
- $p_{Y|X}(y|x)$: model of the discrete r.v.

What is a derived distribution

- It is a PMF or PDF of a function of one or more random variables with known probability law. E.g.:



– Obtaining the PDF for

$$g(X, Y) = Y/X$$

involves deriving a distribution.

Note: $g(X, Y)$ is a random variable

When not to find them

- Don't need PDF for $g(X, Y)$ if only want to compute expected value:

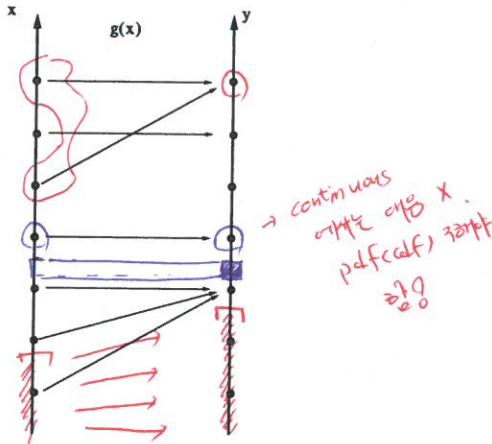
$$E[g(X, Y)] = \int \int g(x, y) f_{X,Y}(x, y) dx dy$$

How to find them

• Discrete case

- Obtain probability mass for each possible value of $Y = g(X)$

$$p_Y(y) = P(g(X) = y) = \sum_{x: g(x)=y} p_X(x)$$



The continuous case

• Two-step procedure:

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

- Differentiate to get

$$f_Y(y) = \frac{dF_Y}{dy}(y)$$

Example

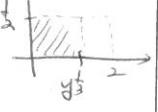
- X : uniform on $[0, 2]$
- Find PDF of $Y = X^3$ $Y \in [0, 8]$

Solution:

$$F_Y(y) = P(Y \leq y) = P(X^3 \leq y) = P(X \leq y^{1/3}) = \frac{1}{2}y^{1/3}$$

$$f_Y(y) = \frac{dF_Y}{dy}(y) = \frac{1}{6y^{2/3}}$$

$X \sim \text{Uniform}[0, 2] \dots P(X \leq x)$

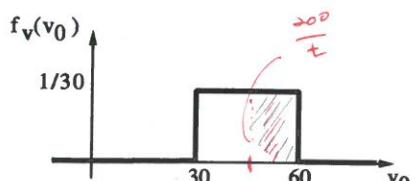


Example

- Joan is driving from Boston to New York. Her speed is uniformly distributed between 30 and 60 mph. What is the distribution of the duration of the trip?

- Let $T(V) = \frac{200}{V}$.

- Find $f_T(t)$



$$F_T(t) = P(T \leq t) \quad \text{replace } T \text{ with } v$$

$$= P\left(\frac{200}{v} \leq t\right)$$

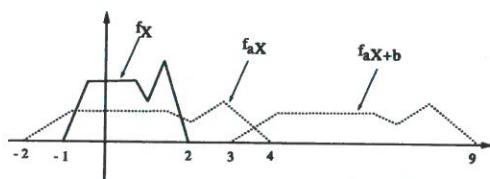
$$= P\left(v \geq \frac{200}{t}\right) = \frac{1}{30} \left(60 - \frac{200}{t}\right)$$

range: $30 \leq t \leq 60$

$$\frac{200}{60} \leq t \leq \frac{200}{30}$$

$$a=2 \quad b=5 \quad \text{The pdf of } Y=aX+b \quad f_{aX+b}(z) = \frac{1}{|a|} f_X\left(\frac{z-b}{a}\right)$$

$$Y = 2X + 5:$$



pdf of y curve. (sum=1)
scale down

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

- Use this to check that if X is normal, then $Y = aX + b$ is also normal.

$$F_Y(y) = P(Y \leq y)$$

$$= P(aX + b \leq y) = P\left(X \leq \frac{y-b}{a}\right)$$

$$= F_X\left(\frac{y-b}{a}\right)$$

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

$$E(Z) \neq \frac{E(Y)}{E(X)}$$

$$= E(Y) - E[\frac{1}{X}]$$

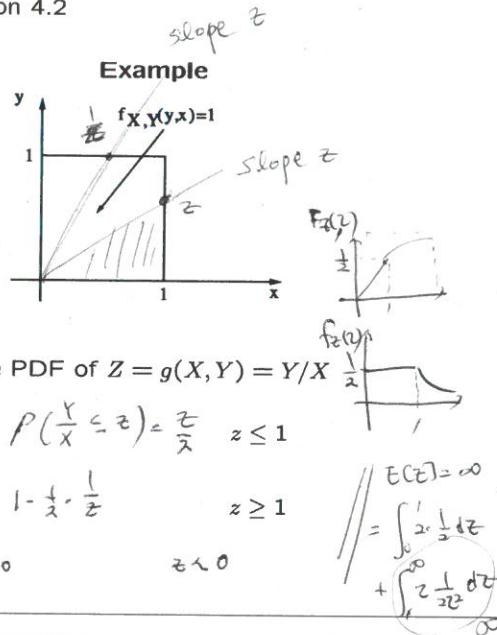
(x,y indep) $h(x,y)$ $g(x)$
indep

LECTURE 11

Derived distributions; convolution; covariance and correlation

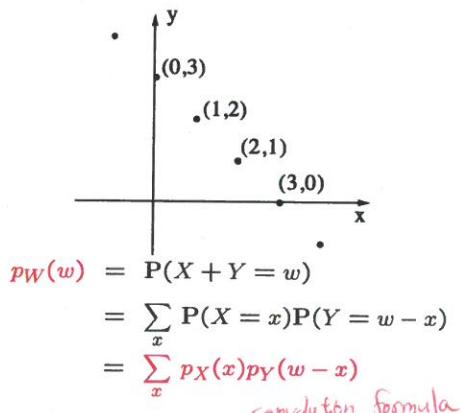
- Readings:**

Finish Section 4.1;
Section 4.2



The distribution of $X + Y$

- $W = X + Y$; X, Y independent



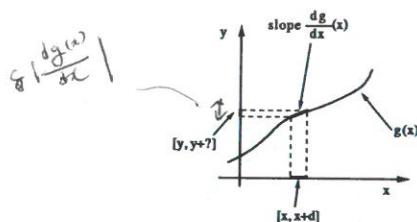
- Mechanics:

- Put the pmf's on top of each other
- Flip the pmf of Y
- Shift the flipped pmf by w (to the right if $w > 0$)
- Cross-multiply and add

$$\left. \begin{aligned} & P(Y = w - x) \\ & \text{cross} \end{aligned} \right\}$$

A general formula

- Let $Y = g(X)$
 g strictly monotonic.



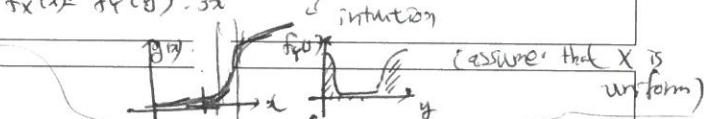
- Event $x \leq X \leq x + \delta$ is the same as $g(x) \leq Y \leq g(x + \delta)$ or (approximately) $g(x) \leq Y \leq g(x) + \delta |(dg/dx)(x)|$

- Hence,

$$\delta f_X(x) = \delta f_Y(y) \left| \frac{dg}{dx}(x) \right|$$

where $y = g(x)$ ✓ inversely proportional

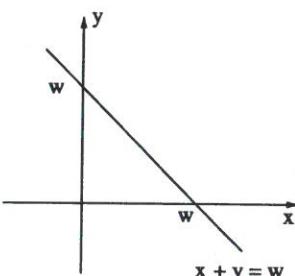
$$(b) Y = X^3, g(x) = x^3 = y, x = y^{1/3}, f_X(x) = f_Y(y) \cdot 3x^2$$



ii) conti

The continuous case

- $W = X + Y$; X, Y independent



- $f_{W|X}(w|x) = f_Y(w-x)$
- $f_{W,X}(w,x) = f_X(x)f_{W|X}(w|x)$
 $= f_X(x)f_Y(w-x)$
- $f_W(w) = \int_{-\infty}^{\infty} f_X(x)f_Y(w-x) dx$

注意到 $\int_{-\infty}^{\infty} f_X(x)f_Y(w-x) dx$

Two independent normal r.v.s

- $X \sim N(\mu_x, \sigma_x^2), Y \sim N(\mu_y, \sigma_y^2)$, independent

$$f_{X,Y}(x,y) = f_X(x)f_Y(y)$$

bell shape

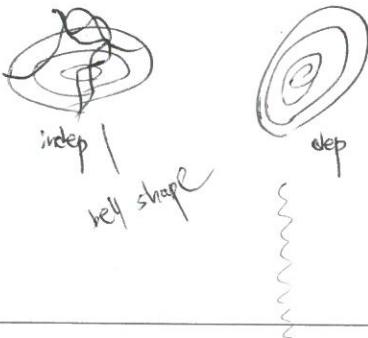
$$= \frac{1}{2\pi\sigma_x\sigma_y} \exp \left\{ -\frac{(x-\mu_x)^2}{2\sigma_x^2} - \frac{(y-\mu_y)^2}{2\sigma_y^2} \right\}$$

- PDF is constant on the ellipse where

$$\frac{(x-\mu_x)^2}{2\sigma_x^2} + \frac{(y-\mu_y)^2}{2\sigma_y^2}$$

is constant

- Ellipse is a circle when $\sigma_x = \sigma_y$



The sum of independent normal r.v.'s

- $X \sim N(0, \sigma_x^2), Y \sim N(0, \sigma_y^2)$, independent

- Let $W = X + Y$

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

$$= \frac{1}{2\pi\sigma_x\sigma_y} \int_{-\infty}^{\infty} e^{-x^2/2\sigma_x^2} e^{-(w-x)^2/2\sigma_y^2} dx$$

(algebra) $= ce^{-\gamma w^2}$

- Conclusion: W is normal

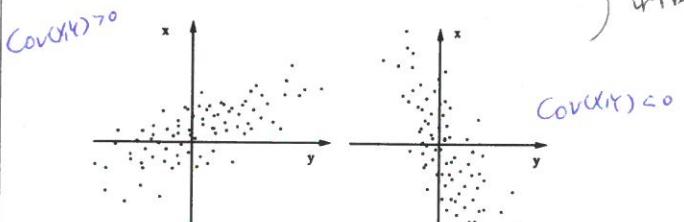
- mean = 0, variance $= \sigma_x^2 + \sigma_y^2$

- same argument for nonzero mean case

indep op. 3

Covariance

- $\text{cov}(X, Y) = E[(X - E[X])(Y - E[Y])]$
- Zero-mean case: $\text{cov}(X, Y) = E[XY]$



- $\text{cov}(X, Y) = E[XY] - E[X]E[Y]$ indep case.
- $\text{var}\left(\sum_{i=1}^n X_i\right) = \sum_{i=1}^n \text{var}(X_i) + \sum_{(i,j):i \neq j} \text{cov}(X_i, X_j)$
- independent $\Rightarrow \text{cov}(X, Y) = 0$
(converse is not true)

$$* \text{cov}(X|X) = \text{Var}(X)$$

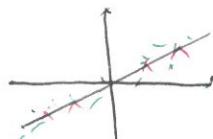
Correlation coefficient

- Dimensionless version of covariance:

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

- $-1 \leq \rho \leq 1$ ~ Exercise 9
- $|\rho| = 1 \Leftrightarrow (X - E[X]) = c(Y - E[Y])$ (linearly related)
- Independent $\Rightarrow \rho = 0$
(converse is not true)

... $\rho \neq 0$ \Rightarrow \exists strong relation
linearly related



4.3 / 4.5

* conditional Ex. Var \rightarrow rv $\in \mathbb{Q}$,**LECTURE 12**

- Readings:** Section 4.3;
parts of Section 4.5
(mean and variance only; no transforms)

Lecture outline

- Conditional expectation
 - Law of iterated expectations
 - Law of total variance
- Sum of a random number of independent r.v.'s
 - mean, variance

Conditional expectations

- Given the value y of a r.v. Y :

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

(integral in continuous case)



- Stick example: stick of length ℓ
break at uniformly chosen point Y
break again at uniformly chosen point X

- $E[X | Y = y] = \frac{y}{2}$ (number)
 $g(x)$

$$E[X | Y] = \frac{Y}{2} \quad (\text{r.v.}) \quad \text{function of rv of } y \rightarrow \text{rv}$$

- Law of iterated expectations:**

$$E[g(x)]$$

$$E[E[X | Y]] = \sum_y E[X | Y = y] p_Y(y) = E[X]$$

(total expectation theorem)

- In stick example:
 $E[X] = E[E[X | Y]] = E[Y/2] = \ell/4$ different scenarios

$X | E[X | Y = y]$] determined completely after we know the value of y
 $V[X | Y = y]$

var($X | Y$) and its expectation

- $\text{var}(X | Y = y) = E[(X - E[X | Y = y])^2 | Y = y]$
- $\text{var}(X | Y)$: a r.v.
with value $\text{var}(X | Y = y)$ when $Y = y$
- Law of total variance:**

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

Proof:

- (a) Recall: $\text{var}(X) = E[X^2] - (E[X])^2$
- (b) $\text{var}(X | Y) = E[X^2 | Y] - (E[X | Y])^2$
conditional universe Y
- (c) $E[\text{var}(X | Y)] = E[X^2] - E[(E[X | Y])^2]$
- (d) $\text{var}(E[X | Y]) = E[(E[X | Y])^2] - (E[X])^2$
rv $\in \mathbb{Q}$

Sum of right-hand sides of (c), (d):

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

intuition 증명! $\rightarrow E[X | Y]$: estimate of X in terms of Y

$\text{var}[X | Y]$: estimation error

(inference)

Section means and variances

Two sections:

 $y = 1$ (10 students); $y = 2$ (20 students)

$$y = 1 : \frac{1}{10} \sum_{i=1}^{10} x_i = 90 \quad y = 2 : \frac{1}{20} \sum_{i=11}^{30} x_i = 60$$

1st section 2nd section

$$E[X] = \frac{1}{30} \sum_{i=1}^{30} x_i = \frac{90 \cdot 10 + 60 \cdot 20}{30} = 70$$

$$E[X | Y = 1] = 90, \quad E[X | Y = 2] = 60$$

$E[X | Y]$ = $\begin{cases} 90, & \text{w.p. } 1/3 \\ 60, & \text{w.p. } 2/3 \end{cases}$

1 section
2nd section
picked

$$E[E[X | Y]] = \frac{1}{3} \cdot 90 + \frac{2}{3} \cdot 60 = 70 = E[X]$$

$$\sqrt{\text{var}(E[X | Y])} = \sqrt{\frac{1}{3}(90 - 70)^2 + \frac{2}{3}(60 - 70)^2} = \sqrt{\frac{600}{3}} = 20$$

$$\text{Var}[E[X | Y] - E[E[X | Y]]]$$

Section means and variances (ctd.)

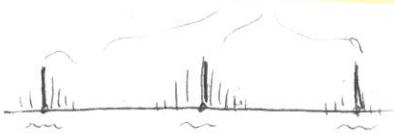
$$\frac{1}{10} \sum_{i=1}^{10} (x_i - 90)^2 = 10 \quad \frac{1}{20} \sum_{i=11}^{30} (x_i - 60)^2 = 20$$

$$\text{var}(X | Y = 1) = 10 \quad \text{var}(X | Y = 2) = 20$$

$$\text{var}(X | Y) = \begin{cases} 10, & \text{w.p. } 1/3 \\ 20, & \text{w.p. } 2/3 \end{cases}$$

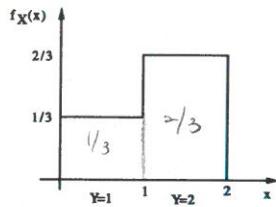
$$E[\text{var}(X | Y)] = \frac{1}{3} \cdot 10 + \frac{2}{3} \cdot 20 = \frac{50}{3}$$

$$\begin{aligned} \text{var}(X) &= E[\text{var}(X | Y)] + \text{var}(E[X | Y]) \\ &= \frac{50}{3} + 200 \\ &= (\text{average variability within sections}) \\ &\quad + (\text{variability between sections}) \end{aligned}$$



Example

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



$$E[X | Y = 1] = \frac{1}{2} \quad E[X | Y = 2] = \frac{3}{2}$$

$$\checkmark \quad \text{var}(X | Y = 1) = \frac{1}{12} \quad \text{var}(X | Y = 2) = \frac{1}{12}$$

$$E[X] = \frac{1}{3} \cdot \frac{1}{2} + \frac{2}{3} \cdot \frac{3}{2} = \frac{7}{6}$$

$$\checkmark \quad \text{var}(E[X | Y]) = \frac{1}{3} \left(\frac{1}{2} - \frac{7}{6} \right)^2 + \frac{2}{3} \left(\frac{3}{2} - \frac{7}{6} \right)^2$$

$$= \frac{1}{2}$$

$$E[E[X | Y] - E[E[X | Y]]]$$

Sum of a random number of independent r.v.'s

- N : number of stores visited (N is a nonnegative integer r.v.)

\checkmark X_i : money spent in store i
 – X_i assumed i.i.d. \sim independently identically distributed
 – independent of N

- Let $Y = X_1 + \dots + X_N$

$$\begin{aligned} E[Y | N = n] &= E[X_1 + X_2 + \dots + X_n | N = n] \\ &= E[X_1 + X_2 + \dots + X_n] \\ &= E[X_1] + E[X_2] + \dots + E[X_n] \\ &= n E[X] \end{aligned}$$

divide and conquer

- $E[Y | N] = N E[X]$

$$\begin{aligned} E[Y] &= E[E[Y | N]] \\ &= E[N E[X]] \\ &= E[N] E[X] \end{aligned}$$

Variance of sum of a random number of independent r.v.'s

$$\bullet \quad \text{var}(Y) = E[\text{var}(Y | N)] + \text{var}(E[Y | N])$$

$$\bullet \quad E[Y | N] = N E[X] \text{ const}$$

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

$$\bullet \quad \text{var}(Y | N = n) = n \text{var}(X) \quad \text{cov} = 0$$

$$\text{var}(Y | N) = N \text{var}(X) \text{ const}$$

$$E[\text{var}(Y | N)] = E[N] \text{var}(X)$$

$$\begin{aligned} \text{var}(Y) &= E[\text{var}(Y | N)] + \text{var}(E[Y | N]) \\ &= E[N] \text{var}(X) + (E[X])^2 \text{var}(N) \end{aligned}$$

LECTURE 13

lecturony
tutor 1
Problem 1

The Bernoulli process

- Readings: Section 6.1

Lecture outline

- Definition of Bernoulli process
- Random processes
- Basic properties of Bernoulli process
- Distribution of interarrival times
- The time of the k th success
- Merging and splitting

The Bernoulli process

- A sequence of independent Bernoulli trials
- At each trial, i :
 - $P(\text{success}) = P(X_i = 1) = p$ constant.
 - $P(\text{failure}) = P(X_i = 0) = 1 - p$
- Examples:
 - Sequence of lottery wins/losses
 - Sequence of ups and downs of the Dow Jones
 - Arrivals (each second) to a bank
 - Arrivals (at each time slot) to server

Random processes

- First view:
sequence of random variables X_1, X_2, \dots

$$\bullet E[X_t] = p \cdot 1 + (1-p) \cdot 0 = p$$

$$\bullet \text{Var}(X_t) = p(1-p)$$

- Second view:
what is the right sample space?

$$\bullet P(X_t = 1 \text{ for all } t) = \lim_{k \rightarrow \infty} p^k \quad (\text{if } p < 1) = 0$$

- Random processes we will study:

- Bernoulli process
(memoryless, discrete time)

- Poisson process
(memoryless, continuous time)

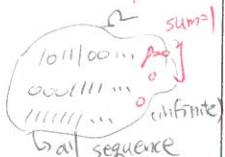
- Markov chains
(with memory/dependence across time)

distribution
of n
or sequence of

n .

continuous
process

sum



+ approximation

\rightarrow quiz

$\leq P(X_1, \dots, X_n)$

$p^{\oplus k}$

Bernoulli - time, num of success
(arrival)

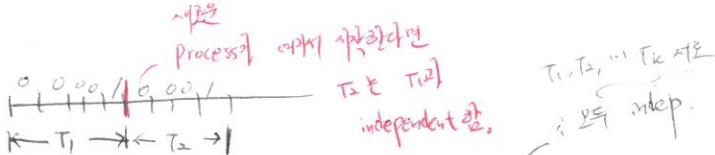
fixing time

Number of successes S in n time slots

$$\bullet P(S = k) = \binom{n}{k} p^k (1-p)^{n-k} \quad k=0, \dots, n$$

$$\bullet E[S] = np$$

$$\bullet \text{Var}(S) = np(1-p)$$



geometric param \rightarrow geometric same param
and dep with T_1)

Interarrival times

- T_1 : number of trials until first success

- $P(T_1 = t) = p \cdot (1-p)^{t-1}$

- Memoryless property

- $E[T_1] = 1/p$

- $\text{Var}(T_1) = (1-p)/p^2$

If you buy a lottery ticket every day, what is the distribution of the length of the first string of losing days?

L = geometric (p)

Time of the k th arrival

- Given that first arrival was at time t
i.e., $T_1 = t$:

additional time, T_2 , until next arrival

- has the same (geometric) distribution

- independent of T_1

- Y_k : number of trials to k th success

- $E[Y_k] = k \cdot \frac{1}{p} = \frac{k}{p}$

- $\text{Var}(Y_k) = k(1-p)/p^2$

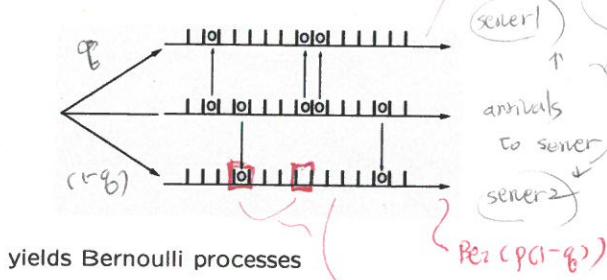
- $P(Y_k = t) = p \cdot \binom{t-1}{k-1} p^{k-1} (1-p)^{t-k}$
 $= \binom{t-1}{k-1} p^k (1-p)^{t-k}$

→ geometric opn^q .

$B_{\text{split}}(pq)$

Splitting of a Bernoulli Process

(using independent coin flips)



① 각 slot 들이

indep. \rightarrow

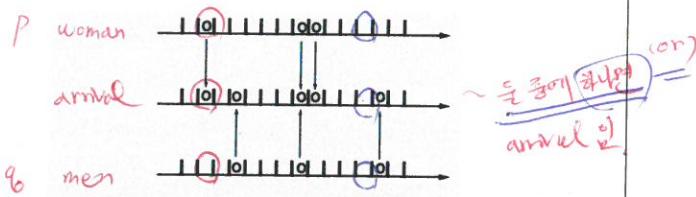
② 각 slot 확률

같음

* Splitting = indep function

remains indep.

Merging of Indep. Bernoulli Processes



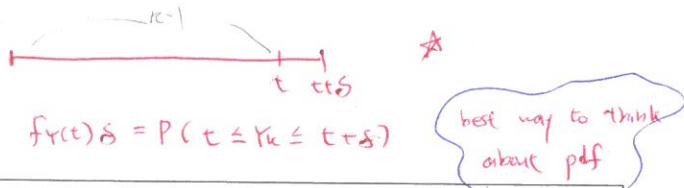
yields a Bernoulli process
(collisions are counted as one arrival)

$$P(\text{arrival}) = 1 - p(\text{no arrival})$$

$$= 1 - (1-p)(1-q)$$

$$= p + q - pq$$

Bernoulli \oplus (각 slot indep \Rightarrow arr)



Example

- You get email according to a Poisson process at a rate of $\lambda = 5$ messages per hour. You check your email every thirty minutes.
- Prob(no new messages) = $P(0, 0.5) = e^{-2.5} = 0.08$
- Prob(one new message) = $P(1, 0.5) = 2.5 e^{-2.5} = 0.20$

$$\lambda t = 5 \times 0.5 = 2.5$$

= prob arrival in [0, t] - λt

$$\therefore f_{Y_k}(y) = \frac{\lambda^y y^{k-1} e^{-\lambda y}}{(k-1)!}$$

Interarrival Times

- Y_k time of k th arrival

- Erlang distribution:

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

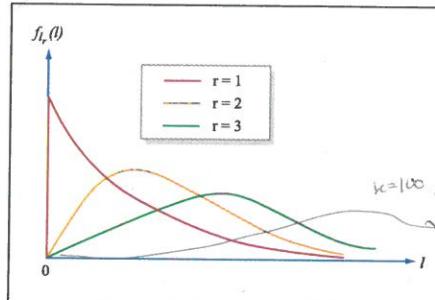


Image by MIT OpenCourseWare.

- Time of first arrival ($k = 1$):
exponential: $f_{Y_1}(y) = \lambda e^{-\lambda y}, \quad y \geq 0$
- Memoryless property: The time to the next arrival is independent of the past

$$Y_2 = T_1 + T_2$$

indep. exp(λ) same param

Bernoulli/Poisson Relation

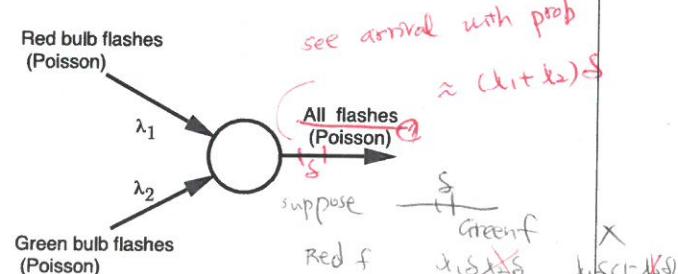


	POISSON	BERNOULLI
Times of Arrival	Continuous	Discrete
Arrival Rate	λ /unit time	p /per trial
PMF of # of Arrivals	Poisson	Binomial
Interarrival Time Distr.	Exponential	Geometric
Time to k -th arrival	Erlang	Pascal

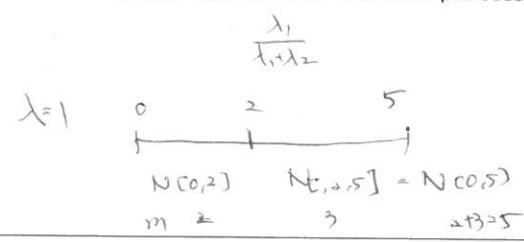
geo is discrete
version of exponential.
expo

Merging Poisson Processes

- Sum of independent Poisson random variables is Poisson
- Merging of independent Poisson processes is Poisson



- What is the probability that the next arrival comes from the first process?



LECTURE 15

Poisson process — II

- Readings: Finish Section 6.2.

- Review of Poisson process
- Merging and splitting
- Examples
- Random incidence

Review

- Defining characteristics
- Time homogeneity: $P(k, \tau)$
- Independence
- Small interval probabilities (small δ):

$$P(k, \delta) \approx \begin{cases} 1 - \lambda\delta, & \text{if } k = 0, \\ \lambda\delta, & \text{if } k = 1, \\ 0, & \text{if } k > 1. \end{cases} + o(\delta^2)$$

- N_τ is a Poisson r.v., with parameter $\lambda\tau$:

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

fixed
↓
 $E[N_\tau] = \text{var}(N_\tau) = \lambda\tau$
 $k \uparrow \text{prob} \downarrow$

$$E[N_\tau] = \text{var}(N_\tau) = \lambda\tau$$

- Interarrival times ($k = 1$): exponential:

$$f_{T_1}(t) = \lambda e^{-\lambda t}, \quad t \geq 0, \quad E[T_1] = 1/\lambda$$

- Time Y_k to k th arrival: Erlang(k):

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

λ, k : parameter of distribution



- Assume: Poisson, $\lambda = 0.6/\text{hour}$.
- Fish for two hours.
- if no catch, continue until first catch.

a) $P(\text{fish for more than two hours}) = P(0, 2)$

$$P(T_1 > 2) = \int_2^\infty f_{T_1}(t) dt = e^{-1.2}$$

b) $P(\text{fish for more than two and less than five hours}) = P(0, 2)(1 - P(0, 3)) = \int_2^3 f_{T_1}(t) dt$

c) $P(\text{catch at least two fish}) = \sum_{k=2}^{\infty} P(k, 2) = 1 - P(0, 2) - P(1, 2) = P(Y_2 \leq 2)$

d) $E[\text{number of fish}] =$

$$= 1.2 + P(0, 2) \cdot 1$$

e) $E[\text{future fishing time} | \text{fished for four hours}] =$

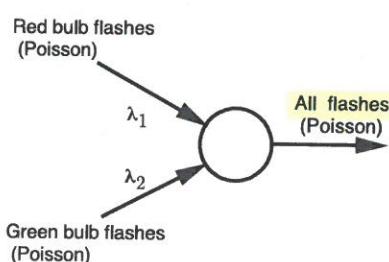
$$\lambda = 0.6$$

f) $E[\text{total fishing time}] =$

$$2 + P(0, 2) \cdot \frac{1}{0.6}$$



- Merging of independent Poisson processes is Poisson



- What is the probability that the next arrival comes from the first process?

$$\frac{\lambda_1}{\lambda_1 + \lambda_2}$$

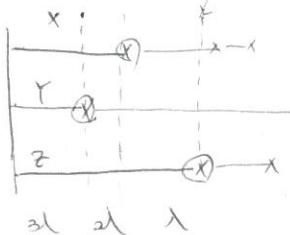
merging & Poisson

Light bulb example

- Each light bulb has independent, exponential(λ) lifetime \rightarrow first arrival time of poisson process
- Install three light bulbs.
Find expected time until last light bulb dies out.

$$E[Y_3] = \frac{1}{3\lambda} + \frac{1}{2\lambda} + \frac{1}{\lambda}$$

$T_1 \quad T_2 \quad T_3$



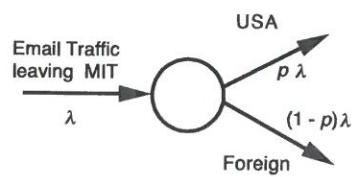
or

$$E[\max\{X_1, Y_1, Z_1\}]$$

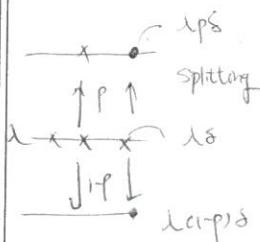
Thinning out

Splitting of Poisson processes

- Assume that email traffic through a server is a Poisson process. Destinations of different messages are independent.



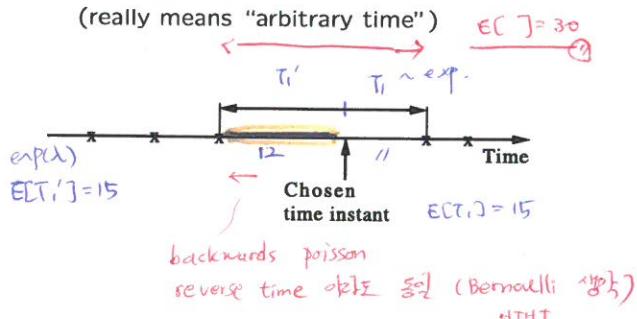
- Each output stream is Poisson.



$$\lambda = 4/\text{hours}, E[U] = 15$$

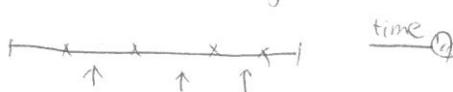
Random incidence for Poisson

- Poisson process that has been running forever
- Show up at some "random time" (really means "arbitrary time")



- What is the distribution of the length of the chosen interarrival interval? 30

biased in favor of longer interval.



Random incidence in "renewal processes"

- Series of successive arrivals
- i.i.d. interarrival times (but not necessarily exponential)

$$\begin{array}{c} 10 \text{ interval} \quad 5 \text{ interval} \\ \downarrow \quad \downarrow \\ \frac{2}{3} \cdot 10 + \frac{1}{3} \cdot 5 = 11.5 \end{array}$$

Example:

Bus interarrival times are equally likely to be 5 or 10 minutes

$$\begin{array}{c} 5 \quad 5 \quad 10 \quad 10 \dots \\ \downarrow \quad \downarrow \quad \downarrow \quad \downarrow \\ E[CT] = 7.5 \end{array}$$

- If you arrive at a "random time":

- what is the probability that you selected a 5 minute interarrival interval?
- what is the expected time to next arrival?

$$\frac{1}{3} \cdot 5 + \frac{2}{3} \cdot 10$$

ex)

Pick a bus

or

a person!

extremely crowded bus

and empty buses

LECTURE 16

Markov Processes – I

- Readings: Sections 7.1–7.2

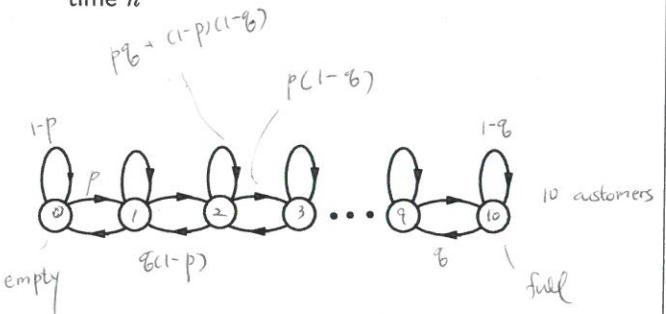
Lecture outline

- Checkout counter example
- Markov process definition
- n -step transition probabilities
- Classification of states

$$\text{new state} = f(\text{old state}, \text{noise})$$

Checkout counter model

- Discrete time $n = 0, 1, \dots$
- Customer arrivals: Bernoulli(p)
 - geometric interarrival times
- Customer service times: geometric(q)
 - departure
- "State" X_n : number of customers at time n



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

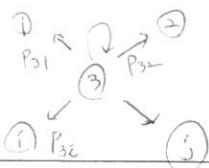
- **Markov property/assumption:**
(given current state, the past does not matter)

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

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

\hookrightarrow 이전에 도달했는지 여부 X
past 여부 여부.

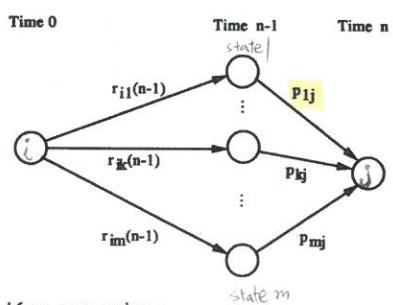
- Model specification:
 - identify the possible states
 - identify the possible transitions
 - identify the transition probabilities



n -step transition probabilities

- State occupancy probabilities, given initial state i :

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



- Key recursion:

$$r_{ij}(n) = \sum_{k=1}^m r_{ik}(n-1)p_{kj} \quad \forall i, j$$

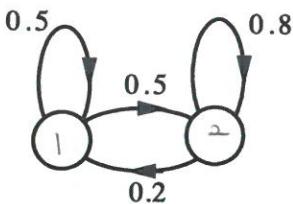
- With random initial state:

$$P(X_n = j) = \sum_{i=1}^m P(X_0 = i)r_{ij}(n)$$

Condition on initial state
(Total probability theorem)

$$\textcircled{X} Z_{ij}(n) = \sum_k P_{ik} r_{kj}(n-1)$$

Example



after longtime,
initial state doesn't
matter

	$n = 0$	$n = 1$	$n = 2$	$n = 100$	$n = 101$
$r_{11}(n)$	1	0.5	$\frac{0.5 \cdot 0.35}{0.2} = 0.35$	$\approx \frac{2}{7}$	$\frac{2}{7} \cdot \frac{1}{2} + \frac{5}{7} \cdot \frac{1}{5} = \frac{1}{2}$
$r_{12}(n)$	0	0.5	0.65	$\frac{5}{7}$	$\frac{5}{7}$
$r_{21}(n)$	0	0.2		$\approx \frac{2}{7}$	
$r_{22}(n)$	1	0.8		$\approx \frac{5}{7}$	

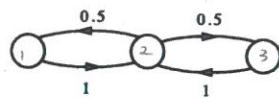
$$r_{11}(n) = r_{12}(n-1) \cdot 0.2 + r_{11}(n-1) \cdot 0.5$$

$$r_{12}(n) = 1 - r_{11}(n)$$

Whether $r_{ij}(n) \xrightarrow{?} r_j$

Generic convergence questions:

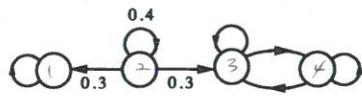
- Does $r_{ij}(n)$ converge to something?
(as $n \rightarrow \infty$)



$$\text{n odd: } r_{22}(n) = 0 \quad \text{n even: } r_{22}(n) = 1$$

convergence fails.

- Does the limit depend on initial state?



$$r_{11}(n) = 1, \forall n$$

$$r_{31}(n) = 0, \forall n$$

$$\cancel{r_{21}(n) = \frac{1}{2} \text{ (as } n \rightarrow \infty)}$$

Classify state into two types

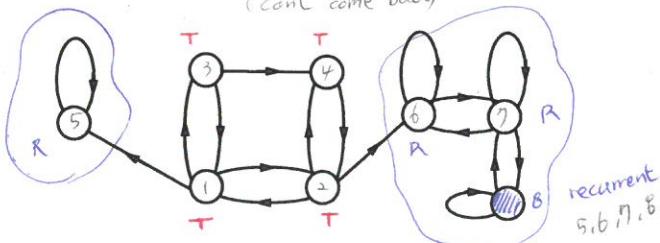
Recurrent and transient states

- State i is **recurrent** if:

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

- If not recurrent, called **transient** ~~transient~~

(can't come back)



- i transient:

$P(X_n = i) \rightarrow 0$,
 i visited finite number of times

initial state
matters?

- **Recurrent class:**

collection of recurrent states that
"communicate" with each other
and with no other state

LECTURE 17

Markov Processes – II

- Readings: Section 7.3

Lecture outline

- Review
- Steady-State behavior
 - Steady-state convergence theorem
 - Balance equations
- Birth-death processes

Review

- Discrete state, discrete time, time-homogeneous
 - Transition probabilities p_{ij} 1 step
 - Markov property

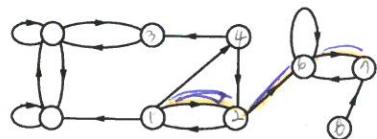
- $r_{ij}(n) = P(X_n = j | X_0 = i)$: n step
 $r_{ij}(1) = p_{ij}$

- Key recursion:

$$r_{ij}(n) = \sum_k r_{ik}(n-1)p_{kj}$$

summing over k
 over

Warmup



$$P(X_1 = 2, X_2 = 6, X_3 = 7 | X_0 = 1) = P_{12} P_{26} P_{67}$$

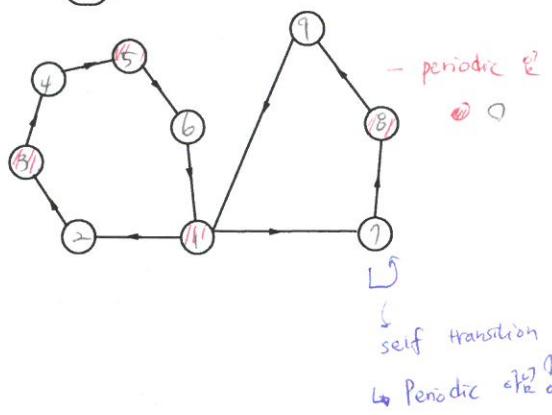
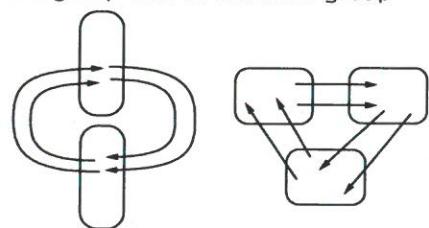
$$P(X_4 = 7 | X_0 = 2) = P_{26} P_{67} P_{76} P_{67} + P_{26} P_{66} P_{66} P_{67} + P_{21} P_{12} P_{66} P_{67}$$

Recurrent and transient states

- State i is **recurrent** if:
starting from i ,
and from wherever you can go,
there is a way of returning to i
- If not recurrent, called **transient**
- **Recurrent class:**
collection of recurrent states that
“communicate” to each other
and to no other state

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



Steady-State Probabilities

- Do the $r_{ij}(n)$ converge to some π_j ? (independent of the initial state i)
- Yes, if:
 - intuition: collision 후 두 개의 확률은 같음
 - recurrent states are all in a single class, and
 - single recurrent class is not periodic

- Assuming "yes," start from key recursion

$$r_{ij}(n) = \sum_k r_{ik}(n-1)p_{kj}$$

- take the limit as $n \rightarrow \infty$

$$\pi_j = \sum_k \pi_k p_{kj}, \quad \text{for all } j$$

] unique solution

Additional equation:

$$\sum_j \pi_j = 1$$

Visit frequency interpretation

Balance equations

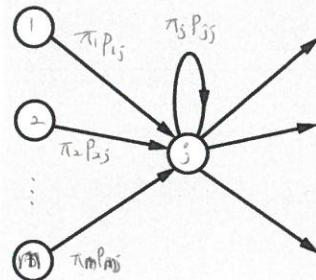
$$\pi_j = \sum_k \pi_k p_{kj}$$

- (Long run) frequency of being in j : π_j

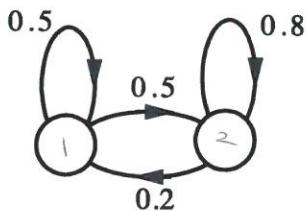
- Frequency of transitions $k \rightarrow j$: $\pi_k p_{kj}$

(total)

- Frequency of transitions into j : $\sum_k \pi_k p_{kj}$



Example



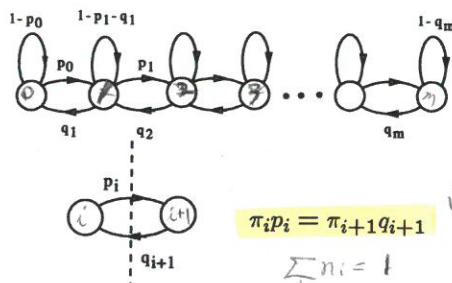
$$\begin{aligned} \pi_1 &= \pi_1 \cdot 0.5 + \pi_2 \cdot 0.2 \\ \pi_2 &= \pi_1 \cdot 0.5 + \pi_2 \cdot 0.8 \end{aligned} \quad \text{linearity dependent}$$

$$0.5\pi_1 = 0.2\pi_2 \Rightarrow \pi_1 + \pi_2 = 1$$

$$\pi_1 = \frac{2}{7}, \pi_2 = \frac{5}{7}$$

Probability settle to state 1
convergence to probability

Birth-death processes



'insight'

supermarket example

- Special case: $p_i = p$ and $q_i = q$ for all i
 $\rho = p/q = \text{load factor}$

$$\pi_{i+1} = \pi_i \frac{p}{q} = \pi_i \rho$$

$\rho > 1$, heavy load ↑

$$\pi_i = \pi_0 \rho^i, \quad i = 0, 1, \dots, m$$

$\rho < 1$, drain down ↓

- Assume $p < q$ and $m \approx \infty$

$$\pi_0 = 1 - \rho \quad \pi_i = (1-\rho) \rho^i \quad (i=0, 1, \dots)$$

$p=q$ symmetric random walk

$$\checkmark \quad \mathbb{E}[X_n] = \frac{\rho}{1-\rho} \quad (\text{in steady-state})$$

LECTURE 18

Markov Processes – III

Readings: Section 7.4

Lecture outline

- Review of steady-state behavior
- Probability of blocked phone calls
- Calculating absorption probabilities
- Calculating expected time to absorption

Review

- Assume a single class of recurrent states, aperiodic; plus transient states. Then,

$$\lim_{n \rightarrow \infty} r_{ij}(n) = \pi_j \quad \text{steady state}$$

where π_j does not depend on the initial conditions:

$$\lim_{n \rightarrow \infty} P(X_n = j | X_0 = i) = \pi_j$$

$$\lim_{n \rightarrow \infty} P(X_n = j) = \pi_j$$

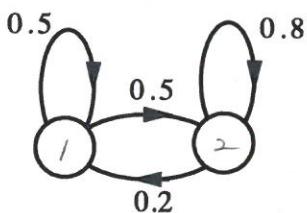
- π_1, \dots, π_m can be found as the unique solution to the balance equations

$$\pi_j = \sum_k \pi_k p_{kj}, \quad j = 1, \dots, m,$$

together with

$$\sum_j \pi_j = 1$$

Example



$$\pi_1 = 2/7, \pi_2 = 5/7$$

- Assume process starts at state 1.

$$P(X_0=1)$$

- $P(X_1 = 1, \text{ and } X_{100} = 1) =$

$$P(X_1 = 1 | X_0 = 1) P(X_{100} = 1 | X_0 = 1) \approx P_1 \cdot \pi_1^{99}$$

- $P(X_{100} = 1 \text{ and } X_{101} = 2) =$

$$= \pi_1 \cdot P_{12} \quad ($$

$$\approx \pi_1 \cdot P_{12}$$

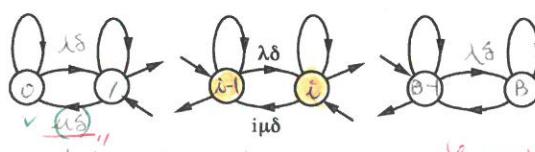
$$P(X_{100} = 1, X_{200} = 1)$$

$$\approx \pi_1 \cdot \pi_1 \cdot \pi_1 \approx \pi_1^3$$

The phone company problem

- Calls originate as a Poisson process, rate λ calls/min
- Each call duration is exponentially distributed (parameter μ) rate /min
- B lines available

- Discrete time intervals of (small) length δ



- Balance equations: $\lambda \pi_{i-1} = i \mu \pi_i$

$$\pi_i = \pi_0 \frac{\lambda^i}{\mu^i i!} \quad \pi_0 = 1 / \sum_{i=0}^B \frac{\lambda^i}{\mu^i i!}$$

• steady-state probability (= frequency)

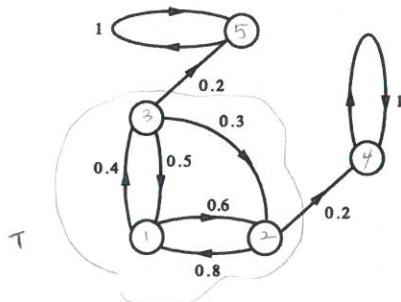
$$\pi_B = p(\text{busy})$$

$$\approx 1\% \rightarrow B = 106,$$

$$\text{margin: } 106 - 10 = 167$$

Calculating absorption probabilities

- What is the probability a_i that process eventually settles in state 4, given that the initial state is i ?



For $i = 4$, $a_i = 1$

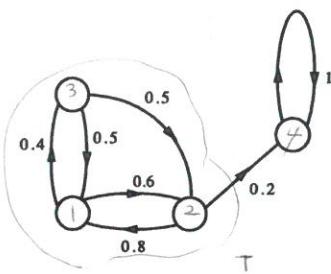
For $i = 5$, $a_i = 0$

$$a_i = \sum_j p_{ij} a_j, \text{ for all other } i$$

ex) $a_2 = 0.2 + 0.8 a_1$

- unique solution

Expected time to absorption



- Find expected number of transitions μ_i , until reaching the absorbing state, given that the initial state is i ?

$$\mu_i = 0 \text{ for } i = 7$$

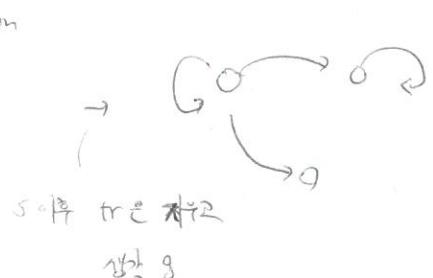
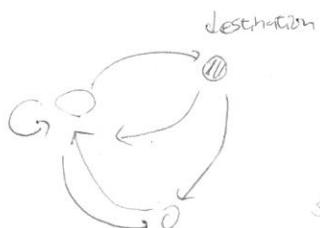
$$\text{For all other } i: \mu_i = 1 + \sum_j p_{ij} \mu_j$$

- unique solution

$$i=1, \mu_1 = 1 + 0.6 \mu_2 + 0.4 \mu_3$$

first transition

* 'or' destination multiple to single one



Mean first passage and recurrence times

- Chain with one recurrent class; fix s recurrent
- Mean first passage time from i to s :**

$$t_i = E[\min\{n \geq 0 \text{ such that } X_n = s\} | X_0 = i]$$

- t_1, t_2, \dots, t_m are the unique solution to

$$t_s = 0, \\ t_i = 1 + \sum_j p_{ij} t_j, \quad \text{for all } i \neq s$$

- Mean recurrence time of s :**

$$t_s^* = E[\min\{n \geq 1 \text{ such that } X_n = s\} | X_0 = s]$$

$$t_s^* = 1 + \sum_j p_{sj} t_j$$

LECTURE 19

Limit theorems – I

- Readings:** Sections 5.1-5.3; start Section 5.4

- X_1, \dots, X_n i.i.d.

$$M_n = \frac{X_1 + \dots + X_n}{n}$$

What happens as $n \rightarrow \infty$?

$$E[X]$$

sample mean
random variable

- Why bother?
- A tool: Chebyshev's inequality
- Convergence "in probability"
- Convergence of M_n (weak law of large numbers)

Markov inequality: mean - probability relationship

$$\begin{aligned} \text{① } X \geq 0 \quad E[X] &= \sum x P(x) \geq \sum_{x \geq a} x P(x) \\ &\geq \sum_{x \geq a} a P(x) = a P(X \geq a) \\ E[(X-\mu)^2] &\geq P((X-\mu)^2 \geq a^2) a^2 \end{aligned}$$

$$\text{② } \text{var}(X) \geq P(|X-\mu| \geq a) a$$

Deterministic limits

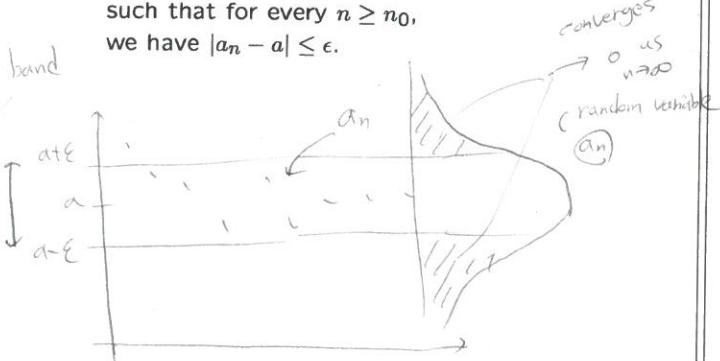
- Sequence a_n
Number a

- a_n converges to a

$$\lim_{n \rightarrow \infty} a_n = a$$

" a_n eventually gets and stays (arbitrarily) close to a "

- For every $\epsilon > 0$, there exists n_0 , such that for every $n \geq n_0$, we have $|a_n - a| \leq \epsilon$.



Chebyshev's inequality

- Random variable (X) continuous (with finite mean μ and variance σ^2)

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

intuition

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

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

Convergence "in probability"

- Sequence of random variables Y_n
- converges in probability to a number a : "(almost all) of the PMF/PDF of Y_n , eventually gets concentrated (arbitrarily) close to a "

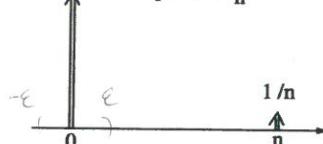
- For every $\epsilon > 0$,

$$\lim_{n \rightarrow \infty} P(|Y_n - a| \geq \epsilon) = 0$$

for $\epsilon > 0$ $\exists n_0$ ~~such that~~ outside of band

$$\forall n \geq n_0 : P(|Y_n - a| \geq \epsilon) \leq \epsilon$$

pmf of Y_n



Does Y_n converge? Yes.

$$Y_n \xrightarrow{\text{ip.}} 0$$

$$\begin{aligned} E[Y_n] &= / \\ E[Y_n^2] &= n \rightarrow \infty \end{aligned}$$

Convergence of the sample mean (Weak law of large numbers)

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

$$M_n = \frac{X_1 + \dots + X_n}{n}$$

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

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

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

- M_n converges in probability to μ

\checkmark $E[2]: P(\cdot) \leq \frac{\sigma^2}{n\epsilon^2} \rightarrow 0$

130pm

Different scalings of M_n

- X_1, \dots, X_n i.i.d.
finite variance σ^2
- Look at three variants of their sum:

• $S_n = X_1 + \dots + X_n$ variance $n\sigma^2$

• $M_n = \frac{S_n}{n}$ variance σ^2/n
converges "in probability" to $E[X]$ (WLLN)

intermediate scaling $\frac{S_n}{\sqrt{n}}$ constant variance σ^2

- Asymptotic shape? variance doesn't change!

$$\text{var}\left(\frac{S_n}{\sqrt{n}}\right) = \frac{1}{n} \text{var}(S_n) = \frac{n\sigma^2}{n} = \sigma^2$$

The pollster's problem

true fraction

- f : fraction of population that "..."
- i th (randomly selected) person polled:

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

- $M_n = (X_1 + \dots + X_n)/n \rightarrow \hat{f}$: fraction of "yes" in our sample

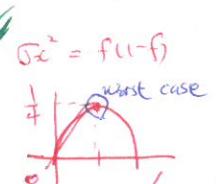
- Goal: 95% confidence of $\leq 1\%$ error

$$P(|M_n - f| \geq .01) \leq .05$$

- Use Chebyshev's inequality:

$$P(|M_n - f| \geq .01) \leq \frac{\sigma_{M_n}^2}{(0.01)^2} = \frac{\sigma_x^2}{n(0.01)^2} \leq \frac{1}{4n(0.01)^2}$$

- If $n = 50,000$, then $P(|M_n - f| \geq .01) \leq .05$ (conservative)



The central limit theorem

- "Standardized" $S_n = X_1 + \dots + X_n$:

$$Z_n = \frac{S_n - E[S_n]}{\sigma_{S_n}} = \frac{S_n - nE[X]}{\sqrt{n}\sigma}$$

- zero mean
- unit variance

- Let Z be a standard normal r.v. (zero mean, unit variance)

- Theorem:** For every c :

$$P(Z_n \leq c) \rightarrow P(Z \leq c)$$

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

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

LECTURE 20

THE CENTRAL LIMIT THEOREM

- Readings: Section 5.4
- X_1, \dots, X_n i.i.d., finite variance σ^2
- "Standardized" $S_n = X_1 + \dots + X_n$:

$$Z_n = \frac{S_n - E[S_n]}{\sigma_{S_n}} = \frac{S_n - nE[X]}{\sqrt{n}\sigma}$$

- $E[Z_n] = 0, \quad \text{var}(Z_n) = 1$

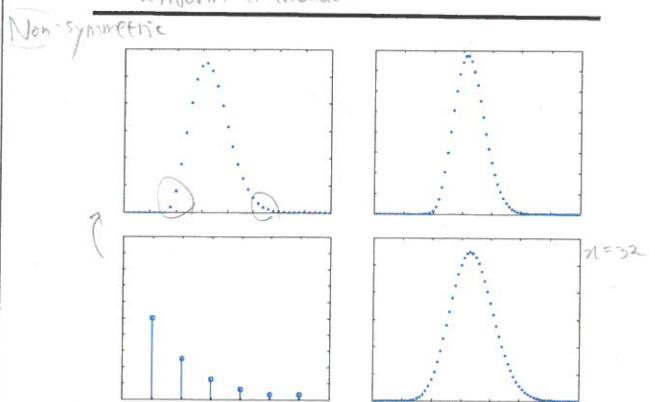
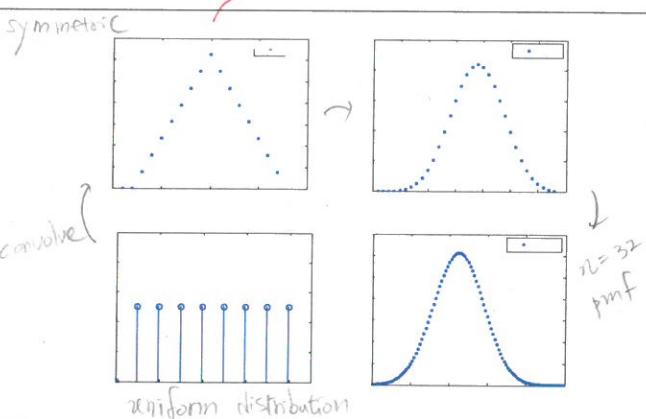
- Let Z be a standard normal r.v. (zero mean, unit variance)

- **Theorem:** For every c :

$$P(Z_n \leq c) \rightarrow P(Z \leq c)$$

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

distribution fixed,



Usefulness

- universal; only means, variances matter otherwise
- accurate computational shortcut (convolving n times)
- justification of normal models (random motion particles...)

What exactly does it say?

- **CDF** of Z_n converges to normal **CDF**
 - not a statement about convergence of PDFs or PMFs

Normal approximation

- Treat Z_n as if normal
 - also treat S_n as if normal (linear transformation)

→ limit theorem

Can we use it when n is "moderate"?

- Yes, but no nice theorems to this effect
- Symmetry helps a lot

The pollster's problem using the CLT

- f : fraction of population that "..."
- i th (randomly selected) person polled:

$$X_i = \begin{cases} 1, & \text{if yes, } w.p. f \\ 0, & \text{if no. } 1-f \end{cases}$$

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

- Suppose we want:

$$P(|M_n - f| \geq .01) \leq .05$$

- Event of interest: $|M_n - f| \geq .01$

$$\left| \frac{X_1 + \dots + X_n - nf}{n} \right| \geq .01$$

(CLT app)

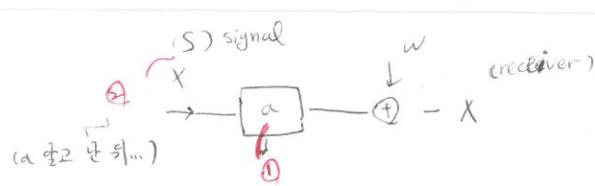
$$\left| \frac{X_1 + \dots + X_n - nf}{\sqrt{n}\sigma} \right| \geq \frac{.01\sqrt{n}}{\sigma}$$

$$P(|M_n - f| \geq .01) \approx P(|Z| \geq .01\sqrt{n}/\sigma) \leq P(|Z| \geq .02\sqrt{n})$$

$$\sigma = \sqrt{f(1-f)}$$

* Probability (통계학적) Study 한 것)

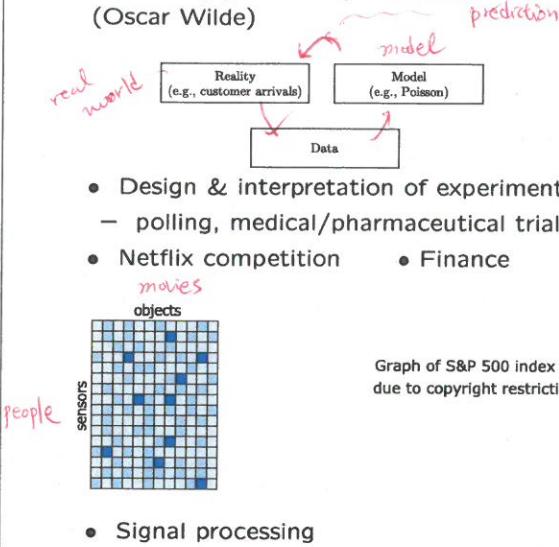
학부의 통계학적 lec 21 ~ .



LECTURE 21

• Readings: Sections 8.1-8.2

"It is the mark of truly educated people to be deeply moved by **statistics**."
(Oscar Wilde)



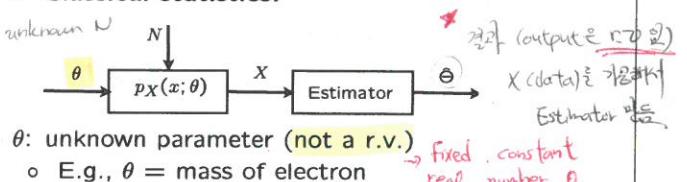
Probability는 서로 다른 개념이 아니라 같은 예제를 배울 것임.

probability - unique solution
statistics, inference - 다른 것들...

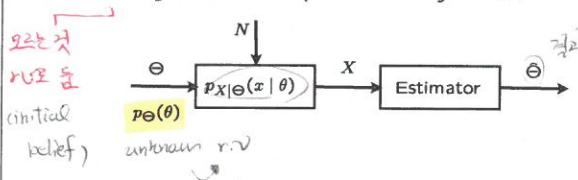
Types of Inference models/approaches

- Model building versus inferring unknown variables. E.g., assume $X = aS + W$
 - Model building:
know "signal" S , observe X , infer a
 - Estimation in the presence of noise:
know a , observe X , estimate S .
- Hypothesis testing: unknown takes one of few possible values; aim at small probability of incorrect decision
- Estimation: aim at a small estimation error

• Classical statistics:



• Bayesian: Use priors & Bayes rule



Bayesian inference: Use Bayes rule

• Hypothesis testing prior (initial belief)

- discrete data

$$p_{\Theta|X}(\theta | x) = \frac{p_\theta(\theta) p_{X|\Theta}(x | \theta)}{p_X(x)}$$

model of experiment

- continuous data

$$p_{\Theta|X}(\theta | x) = \frac{p_\theta(\theta) f_{X|\Theta}(x | \theta)}{f_X(x)}$$

θ : continuous

• Estimation; continuous data

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

$$\begin{aligned} Z_t &= \Theta_0 + t\Theta_1 + t^2\Theta_2 \\ X_t &= Z_t + W_t, \quad t = 1, 2, \dots, n \end{aligned}$$

Bayes rule gives:

$$f_{\Theta_0, \Theta_1, \Theta_2 | X_1, \dots, X_n}(\theta_0, \theta_1, \theta_2 | x_1, \dots, x_n)$$

i, j, k : multidimensional θ_i & x_i .

Estimation with discrete data

$$f_{\Theta|X}(\theta | x) = \frac{f_\theta(\theta) p_{X|\Theta}(x | \theta)}{p_X(x)}$$

$$p_X(x) = \int f_\theta(\theta) p_{X|\Theta}(x | \theta) d\theta$$

• Example:

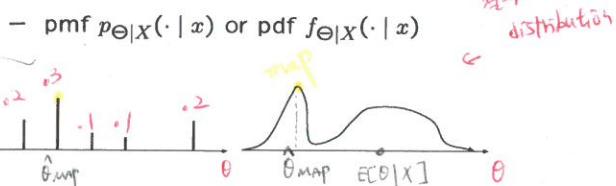
- Coin with unknown parameter θ (e.g. $\hat{\theta}_m = \frac{X}{n}$)
- Observe X heads in n tosses
- What is the Bayesian approach?
 - Want to find $f_{\Theta|X}(\theta | x)$ Bayesian
 - Assume a prior on Θ (e.g., uniform) (prior)

assume distribution of θ .
Uniform normal

most likely to be true

Output of Bayesian Inference

- Posterior distribution:



- If interested in a single answer:

- Maximum a posteriori probability (MAP):

- $p_{\Theta|X}(\theta^*|x) = \max_{\theta} p_{\Theta|X}(\theta|x)$ minimizes probability of error; often used in hypothesis testing

- $f_{\Theta|X}(\theta^*|x) = \max_{\theta} f_{\Theta|X}(\theta|x)$

- Conditional expectation:

$$E[\Theta | X = y] = \int \theta f_{\Theta|X}(\theta|x) d\theta \quad E[\Theta|X=y]$$

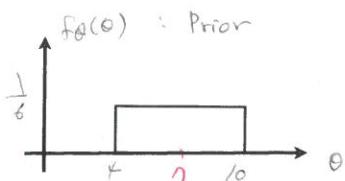
- Single answers can be misleading!

↳ MAP [Point estimation, given data, $\theta \sim \text{prior}$, report $E[\Theta|X=x]$]

1) 2) 3)

Least Mean Squares Estimation

- Estimation in the absence of information



- find estimate c , to:

$$\text{minimize } E[(\Theta - c)^2] \quad \begin{matrix} \text{Avg} \\ \text{squared error} \end{matrix}$$

- Optimal estimate: $c = E[\Theta]$ (derivative = 0)

- Optimal mean squared error:

$$E[(\Theta - E[\Theta])^2] = \text{Var}(\Theta) \quad \begin{matrix} \text{Avg size of} \\ \text{estimation error} \end{matrix}$$

LMS Estimation of Θ based on X

- Two r.v.'s Θ, X

- we observe that $X = x$

- new universe: condition on $X = x$

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

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

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

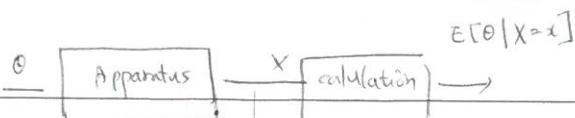
$$\leq E[(\Theta - g(x))^2 | X = x]$$

\hat{c}_{LMS} estimator

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

iterated expectation
optimal estimator

$E[\Theta | X]$ minimizes $E[(\Theta - g(X))^2]$
over all estimators $g(\cdot)$



43

Somebody's choices $\rightarrow g(x)$ other estimators ... (alternative)

Conditional $E[\Theta | X = x]$ ok? ... x real life ↓

LMS Estimation w. several measurements

- Unknown r.v. Θ

several data in real world

- Observe values of r.v.'s X_1, \dots, X_n

- Best estimator: $E[\Theta | X_1, \dots, X_n]$ - (prior Θ)

- Can be hard to compute/implement

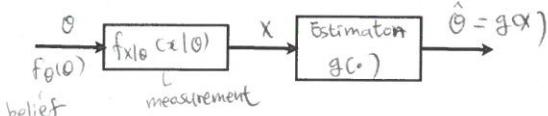
- involves multi-dimensional integrals, etc.

LECTURE 22

- Readings: pp. 225-226; Sections 8.3-8.4

Topics

- (Bayesian) Least means squares (LMS) estimation
- (Bayesian) Linear LMS estimation



- MAP estimate: $\hat{\theta}_{MAP}$ maximizes $f_{\theta|X}(\theta|x)$

- LMS estimation:

- $\hat{\theta} = E[\theta | X]$ minimizes $E[(\theta - g(X))^2]$ over all estimators $g(\cdot)$
- for any x , $\hat{\theta} = E[\theta | X = x]$ minimizes $E[(\theta - \hat{\theta})^2 | X = x]$ over all estimates $\hat{\theta}$

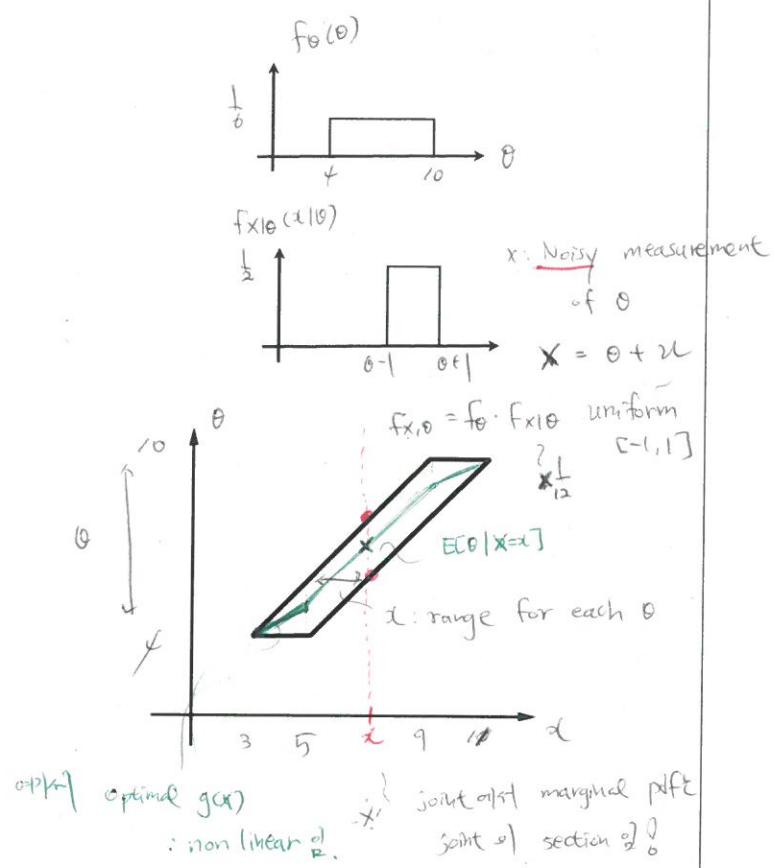
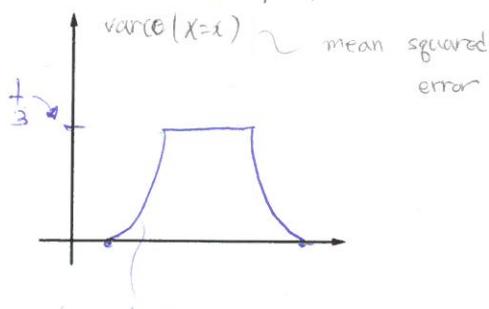
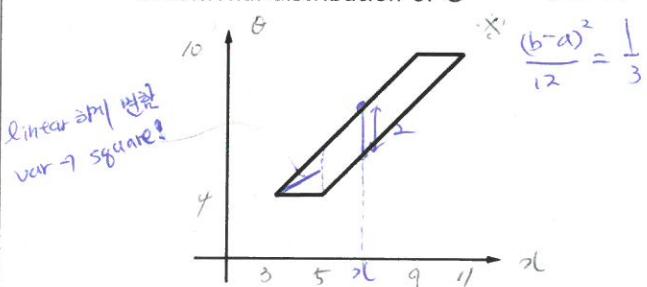
specific x happens...

MAP



Conditional mean squared error

- $E[(\theta - E[\theta | X])^2 | X = x]$
- same as $\text{Var}(\theta | X = x)$: variance of the conditional distribution of θ



4.3 정리

$\hat{\theta}$ unbiased

Some properties of LMS estimation

- Estimator: $\hat{\theta} = E[\theta | X]$
- Estimation error: $\tilde{\theta} = \hat{\theta} - \theta$
imply that this is true for all X !
- $E[\tilde{\theta}] = 0 \rightarrow E[\tilde{\theta} | X = x] = 0$

• $E[\tilde{\theta}h(X)] = 0$, for any function h

derived from $\text{cov}(\tilde{\theta}, \hat{\theta}) = 0$

• Since $\theta = \hat{\theta} - \tilde{\theta}$:
 $\text{var}(\theta) = \text{var}(\hat{\theta}) + \text{var}(\tilde{\theta})$

$\text{cov}(\hat{\theta}, h(x)) = 0 \rightarrow \text{cov}(\tilde{\theta}, \hat{\theta}) = 0$
any function of X . $h(x) = \hat{\theta} = E[\theta | X]$

function of X

$\tilde{\theta} \rightarrow \hat{\theta}$

X 일 때 $\hat{\theta}$ 자동으로 알 수 있다...

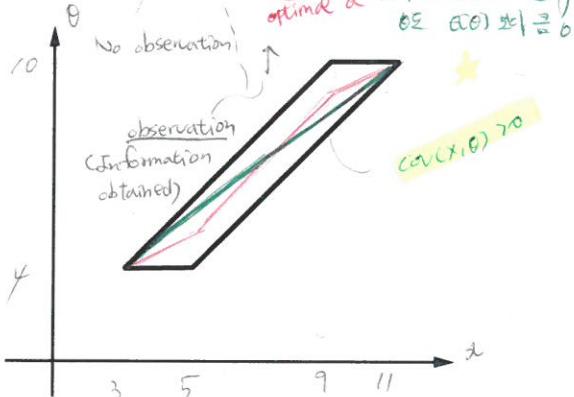
variance \equiv uncertainty \equiv variance

\downarrow
variance
variance

approximation of (x_1, x_2, \dots, x_n)

Linear LMS

- Consider estimators of Θ , of the form $\hat{\Theta} = aX + b$ choose best linear function
- Minimize $E[(\Theta - aX - b)^2] = h(a, b)$; quad
- Best choice of a, b ; best linear estimator: $\hat{\Theta}_L = E[\Theta] + \frac{\text{Cov}(X, \Theta)}{\text{var}(X)}(X - E[X])$ intuition



Linear LMS properties

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

$\rightarrow E[(\hat{\Theta}_L - \Theta)^2] = (1 - \rho^2)\sigma_\theta^2$ $\rho \uparrow$ error ↓
if $\rho = 0$, observation does not help. can't reduce uncertainty

Linear LMS with multiple data

- Consider estimators of the form: $E[\Theta | x_1, \dots, x_n]$

$$\hat{\Theta} = a_1x_1 + \dots + a_nx_n + b$$

- Find best choices of a_1, \dots, a_n, b
- Minimize:

$$E[(a_1x_1 + \dots + a_nx_n + b - \Theta)^2]$$

- Set derivatives to zero linear system in b and the a_i
- Only means, variances, covariances matter

$$a_1^2 E[x_1^2] + 2a_1 a_2 E[x_1 x_2] + \dots$$

observation from experiment, x_1, x_2, \dots, x_n
measurement. data points,

The cleanest linear LMS example

$$X_i = \Theta + W_i, \quad \Theta, W_1, \dots, W_n \text{ independent}$$

$$\Theta \sim \mu, \sigma_\theta^2 \quad W_i \sim 0, \sigma_i^2 \quad \text{Variance } \downarrow \text{ noise } \uparrow$$

$$\text{Prior mean } \mu/\sigma_\theta^2 + \sum_{i=1}^n X_i/\sigma_i^2$$

$$\hat{\Theta}_L = \frac{\sum_{i=0}^n 1/\sigma_i^2}{\sum_{i=0}^n 1/\sigma_i^2}$$

(weighted average of μ, X_1, \dots, X_n)

If all normal, $\hat{\Theta}_L = E[\Theta | X_1, \dots, X_n]$

LMS = linear LMS

in Bayesian Inference

Big picture

Standard examples:

- X_i uniform on $[0, \theta]$; uniform prior on θ
- X_i Bernoulli(p); uniform (or Beta) prior on p
- X_i normal with mean θ , known variance σ^2 ; normal prior on θ ; $X_i = \Theta + W_i$

Estimation methods:

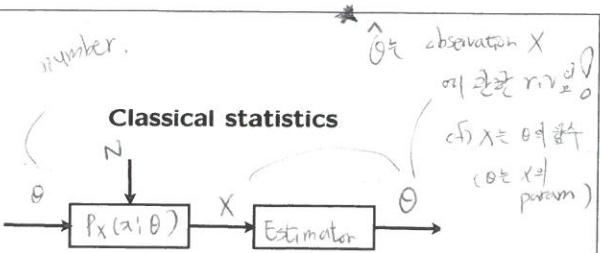
- MAP
- MSE
- Linear MSE

Choosing X_i in linear LMS

- $E[\Theta | X]$ is the same as $E[\Theta | X^3]$
- Linear LMS is different:
 - $\hat{\Theta} = aX + b$ versus $\hat{\Theta} = aX^3 + b$
 - Also consider $\hat{\Theta} = a_1X + a_2X^2 + a_3X^3 + b$

LECTURE 23

- Readings:** Section 9.1
(not responsible for t -based confidence intervals, in pp. 471-473)
- Outline**
 - Classical statistics
 - Maximum likelihood (ML) estimation
 - Estimating a sample mean
 - Confidence intervals (CIs)
 - CIs using an estimated variance



- also for vectors X and θ :
 $p_{X_1, \dots, X_n}(x_1, \dots, x_n; \theta_1, \dots, \theta_m)$
- These are NOT conditional probabilities;**
 θ is NOT random (constant) (x̄ param of θ)
- mathematically: many models, one for each possible value of θ
- Problem types:**
 - Hypothesis testing:
 $H_0: \theta = 1/2$ versus $H_1: \theta = 3/4$
 - Composite hypotheses:
 $H_0: \theta = 1/2$ versus $H_1: \theta \neq 1/2$
 - Estimation: design an estimator $\hat{\theta}$, to keep estimation error $\hat{\theta} - \theta$ small

Maximum Likelihood Estimation

- Model, with unknown parameter(s):
 $X \sim p_X(x; \theta)$
- Pick θ that "makes data most likely"

$$\hat{\theta}_{ML} = \arg \max_{\theta} p_X(x; \theta)$$

fx

- Compare to Bayesian MAP estimation:

$$\hat{\theta}_{MAP} = \arg \max_{\theta} p_{\theta|X}(\theta | x)$$

constant of $p_{\theta|X}(\theta | x)$

$$\hat{\theta}_{MAP} = \arg \max_{\theta} \frac{p_X(x|\theta)p_{\theta}(\theta)}{p_X(x)}$$

MAP = ML

- Example: X_1, \dots, X_n : i.i.d., exponential(θ)

$$\max_{\theta} \prod_{i=1}^n \theta e^{-\theta x_i}$$

$$\max_{\theta} \left(n \log \theta - \theta \sum_{i=1}^n x_i \right)$$

r.v. of observation X

$$\hat{\theta}_{ML} = \frac{n}{x_1 + \dots + x_n}$$

$$\hat{\theta} = \frac{n}{X_1 + \dots + X_n}$$

function of observation X

ML Estimator = 특징값

는 X 의 r.v.
는 θ 의 param!

Desirable properties of estimators (should hold FOR ALL θ !!!)

- Unbiased:** $E[\hat{\theta}_n] = \theta$ $\star \int f_{\hat{\theta}_n}(\hat{\theta}; \theta) \hat{\theta} d\theta = E[\hat{\theta}_n]$
 - exponential example, with $n = 1$:
 $E[1/X_1] = \infty \neq \theta$
 (biased)
- Consistent:** $\hat{\theta}_n \rightarrow \theta$ (in probability)
 - exponential example: un P. 269
 - $(X_1 + \dots + X_n)/n \rightarrow E[X] = 1/\theta$
 - can use this to show that:
 $\hat{\theta}_n = n/(X_1 + \dots + X_n) \rightarrow 1/E[X] = \theta$
- "Small" mean squared error (MSE)**

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

$$= \text{var}(\hat{\theta}) + (\text{bias})^2$$

θ : constant of

ex) $X \sim N(\theta, 1)$

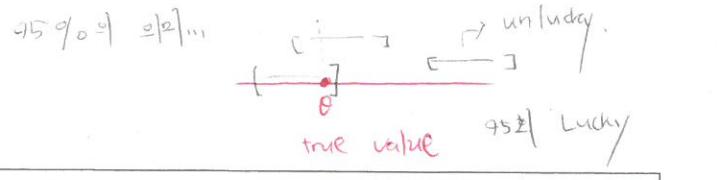
$$\hat{\theta} = 100, \text{ bias} = 100 - \theta$$

$\text{var}(\hat{\theta})$ 작자인

偏差 + 1 (Bias + 1은 정답을 찾을 때 중요)

정답은 θ 에 대한 예상치가 아닙니다

ex) $\theta = 0, \dots$



Estimate a mean using sample mean

- X_1, \dots, X_n : i.i.d., mean θ , variance σ^2

$$X_i = \theta + W_i$$

W_i : i.i.d., mean 0, variance σ^2

ex) $\hat{\Theta}_n = \text{sample mean} = M_n = \frac{X_1 + \dots + X_n}{n}$
estimator, sample mean

Properties:

- $E[\hat{\Theta}_n] = \theta$ (unbiased)
- WLLN: $\hat{\Theta}_n \rightarrow \theta$ (consistency) P. 269
- MSE: σ^2/n $E[(\hat{\Theta}_n - \theta)^2]$
- Sample mean often turns out to also be the ML estimate.
E.g., if $X_i \sim N(\theta, \sigma^2)$, i.i.d.

Confidence intervals (CIs)

- An estimate $\hat{\Theta}_n$ may not be informative enough (single value vs. range)

- An $1 - \alpha$ confidence interval is a (random) interval $[\hat{\Theta}_n^- , \hat{\Theta}_n^+]$,

s.t. $P(\hat{\Theta}_n^- \leq \theta \leq \hat{\Theta}_n^+) \geq 1 - \alpha, \forall \theta$

- often $\alpha = 0.05$, or 0.25, or 0.01
- interpretation is subtle

- CI in estimation of the mean
 $\hat{\Theta}_n = (X_1 + \dots + X_n)/n$ (sample mean)

normal tables: $\Phi(1.96) = 1 - 0.05/2$

$$P\left(\frac{|\hat{\Theta}_n - \theta|}{\sigma/\sqrt{n}} \leq 1.96\right) \approx 0.95 \quad (\text{CLT})$$

$$P\left(\hat{\Theta}_n - \frac{1.96\sigma}{\sqrt{n}} \leq \theta \leq \hat{\Theta}_n + \frac{1.96\sigma}{\sqrt{n}}\right) \approx 0.95$$

More generally: let z be s.t. $\Phi(z) = 1 - \alpha/2$

$$P\left(\hat{\Theta}_n - \frac{z\sigma}{\sqrt{n}} \leq \theta \leq \hat{\Theta}_n + \frac{z\sigma}{\sqrt{n}}\right) \approx 1 - \alpha$$

interval of σ 을 알 수 있다.

→ estimate σ .

The case of unknown σ

- Option 1: use upper bound on σ
 - if X_i Bernoulli: $\sigma \leq 1/2$
- Option 2: use ad hoc estimate of σ
 - if X_i Bernoulli(θ): $\hat{\sigma} = \sqrt{\hat{\theta}(1 - \hat{\theta})}$
↳ $\hat{\theta} = P(\text{observed})$
- Option 3: Use generic estimate of the variance
 - Start from $\sigma^2 = E[(X_i - \theta)^2]$ ↳ how to compute in if n is large.

$$\hat{\sigma}_n^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \theta)^2 \rightarrow \sigma^2$$

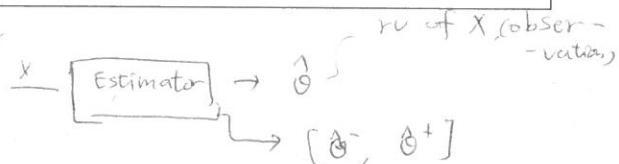
(but do not know θ) ↳ n large → converges to σ^2

$$\hat{S}_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \hat{\Theta}_n)^2 \rightarrow \sigma^2$$

(unbiased: $E[\hat{S}_n^2] = \sigma^2$) ↳ plug in the estimate of mean.

잘 했어!

잘 했어!



LECTURE 24

- Reference: Section 9.3
- Course Evaluations (until 12/16)
<http://web.mit.edu/subjectevaluation>

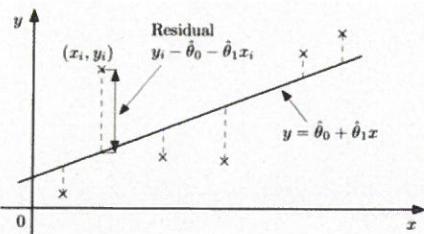
Outline

- Review
 - Maximum likelihood estimation
 - Confidence intervals
- Linear regression
- Binary hypothesis testing
 - Types of error
 - Likelihood ratio test (LRT)

Review

- Maximum likelihood estimation
 - Have model with unknown parameters:
 $X \sim p_X(x; \theta)$ ↳ multiple models
 - Pick $\hat{\theta}$ that "makes data most likely"
 $\max_{\theta} p_X(x; \theta)$
- Compare to Bayesian MAP estimation:
 $\max_{\theta} p_{\Theta|X}(\theta | x)$ or $\max_{\theta} \frac{p_X(x|\theta)p_{\Theta}(\theta)}{p_Y(y)}$ Prior flat ↳ $p_{\Theta}(\theta) = 1$
- Sample mean estimate of $\theta = E[X]$
 $\hat{\theta}_n = (X_1 + \dots + X_n)/n$
- $1 - \alpha$ confidence interval ↳ Prob with respect to Interval
 $P(\hat{\theta}_n^- \leq \theta \leq \hat{\theta}_n^+) \geq 1 - \alpha, \forall \theta$
 ↳ random interval itself to capture θ
- confidence interval for sample mean
 - let z be s.t. $\Phi(z) = 1 - \alpha/2$
 - $P\left(\hat{\theta}_n - \frac{z\sigma}{\sqrt{n}} \leq \theta \leq \hat{\theta}_n + \frac{z\sigma}{\sqrt{n}}\right) \approx 1 - \alpha$

Regression



- Data: $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$
- Model: $y \approx \theta_0 + \theta_1 x$ (hypothesize ...)

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

Probabilistic justification

- One interpretation:
 $Y_i = \theta_0 + \theta_1 x_i + W_i, W_i \sim N(0, \sigma^2)$, i.i.d.

– Likelihood function $f_{X,Y|\theta}(x, y; \theta)$ is:

$$W_i = Y_i - \theta_0 - \theta_1 x_i \quad \text{c} \cdot \exp\left\{-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \theta_0 - \theta_1 x_i)^2\right\}$$

- Take logs, same as (*)
- Least sq. ↔ pretend W_i i.i.d. normal

$$E[Y] = E[\theta_0 + \theta_1 X + W] = \theta_0 + \theta_1 E[X]$$

$$E[Y^2] = E[\theta_0^2 + 2\theta_0 \theta_1 X + \theta_1^2 X^2 + W^2] = \theta_0^2 + 2\theta_0 \theta_1 E[X] + \theta_1^2 E[X^2] + \sigma^2$$

$$\text{cov}(X, Y) = \theta_1 \text{var}(X)$$

Linear regression

- Model $y \approx \theta_0 + \theta_1 x$

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

- Solution (set derivatives to zero):

$$\bar{x} = \frac{x_1 + \dots + x_n}{n}, \quad \bar{y} = \frac{y_1 + \dots + y_n}{n}$$

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

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

- Interpretation of the form of the solution

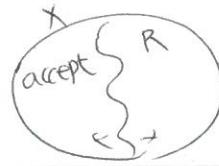
- Assume a model $Y = \theta_0 + \theta_1 X + W$
 W independent of X , with zero mean
- Check that

$$\hat{\theta}_1 = \frac{\text{cov}(X, Y)}{\text{var}(X)} = \frac{E[(X - E[X])(Y - E[Y])]}{E[(X - E[X])^2]}$$

- Solution formula for $\hat{\theta}_1$ uses natural estimates of the variance and covariance

$$\theta_0 = E[Y] - \hat{\theta}_1 E[X] \quad (\bar{Y}, \bar{X} \in \text{data})$$

$$\hat{\theta}_0 = \bar{Y} - \hat{\theta}_1 \bar{X} \quad (\text{est. value})$$



- i) shape of boundary
- ii) move \leftrightarrow

LECTURE 25 Outline

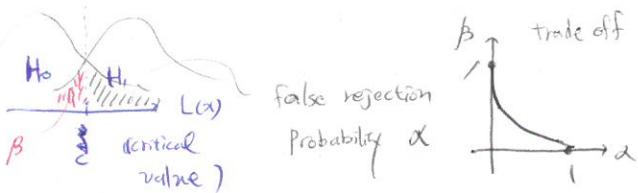
- Reference: Section 9.4
- Course Evaluations (until 12/16)
<http://web.mit.edu/subjectevaluation>
- Review of simple binary hypothesis tests
 - examples
- Testing composite hypotheses
 - is my coin fair?
 - is my die fair?
 - goodness of fit tests

Simple binary hypothesis testing

- null hypothesis H_0 : \sim (default)
- $X \sim p_X(x; H_0)$ [or $f_X(x; H_0)$]
- alternative hypothesis H_1 :
 $X \sim p_X(x; H_1)$ [or $f_X(x; H_1)$]
- Choose a **rejection region** R ;
reject H_0 iff data $\in R$
- Likelihood ratio test: reject H_0 if

$$L(x) = \frac{p_X(x; H_1)}{p_X(x; H_0)} > \xi \quad \text{or} \quad \frac{f_X(x; H_1)}{f_X(x; H_0)} > \xi$$
- ⊖ fix false rejection probability α
(e.g., $\alpha = 0.05$)
- choose ξ so that $P(\text{reject } H_0; H_0) = \alpha$

Likelihood ratio test gives α, β ,



Example (test on normal mean)

- n data points, i.i.d.
 $H_0: X_i \sim N(0, 1)$
 $H_1: X_i \sim N(1, 1)$
- Likelihood ratio test; rejection region:

$$\frac{(1/\sqrt{2\pi})^n \exp\{-\sum_i(X_i - 1)^2/2\}}{(1/\sqrt{2\pi})^n \exp\{-\sum_i X_i^2/2\}} > \xi$$

simplify
- algebra: reject H_0 if: $\sum_i X_i > \xi'$

equivalent
- Find ξ' such that

$$P\left(\sum_{i=1}^n X_i > \xi'; H_0\right) = \alpha$$

ex) $\alpha = 0.05$
- use normal tables

$$\sum_{i=1}^n X_i \sim \text{Normal}(n, n)$$

Example (test on normal variance)

- n data points, i.i.d.
 $H_0: X_i \sim N(0, 1)$
 $H_1: X_i \sim N(0, 4)$
- Likelihood ratio test; rejection region:

$$\frac{(1/2\sqrt{2\pi})^n \exp\{-\sum_i X_i^2/(2 \cdot 4)\}}{(1/\sqrt{2\pi})^n \exp\{-\sum_i X_i^2/2\}} > \xi$$

simplify
- algebra: reject H_0 if $\sum_i X_i^2 > \xi'$
- Find ξ' such that

$$P\left(\sum_{i=1}^n X_i^2 > \xi'; H_0\right) = \alpha$$
- the distribution of $\sum_i X_i^2$ is known
(derived distribution problem)
 - ex) or
 - convolution
- "chi-square" distribution;
tables are available

