

LECTURE 1: Probability models and axioms

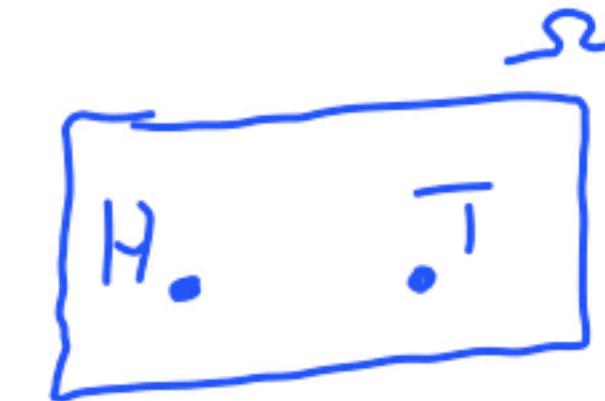
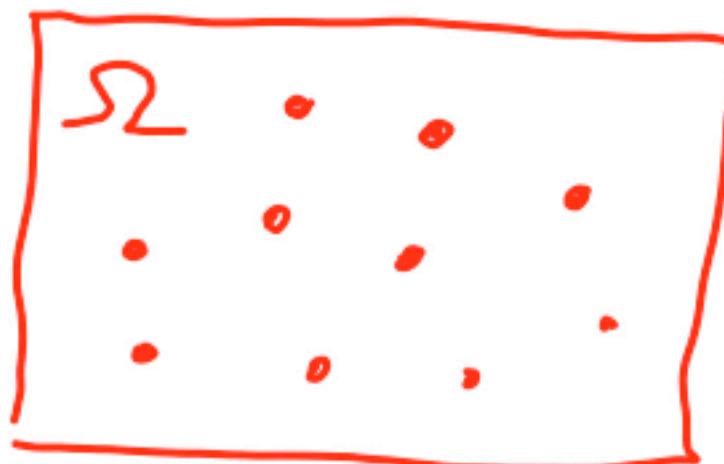
- Sample space
- Probability laws
 - Axioms
 - Properties that follow from the axioms
- Examples
 - Discrete
 - Continuous
- Discussion
 - Countable additivity
 - Mathematical subtleties
- Interpretations of probabilities

Sample space

- Two steps:
 - Describe possible outcomes
 - Describe beliefs about likelihood of outcomes

Sample space

- List (set) of possible outcomes, Ω
- List must be:
 - Mutually exclusive
 - Collectively exhaustive
 - At the “right” granularity

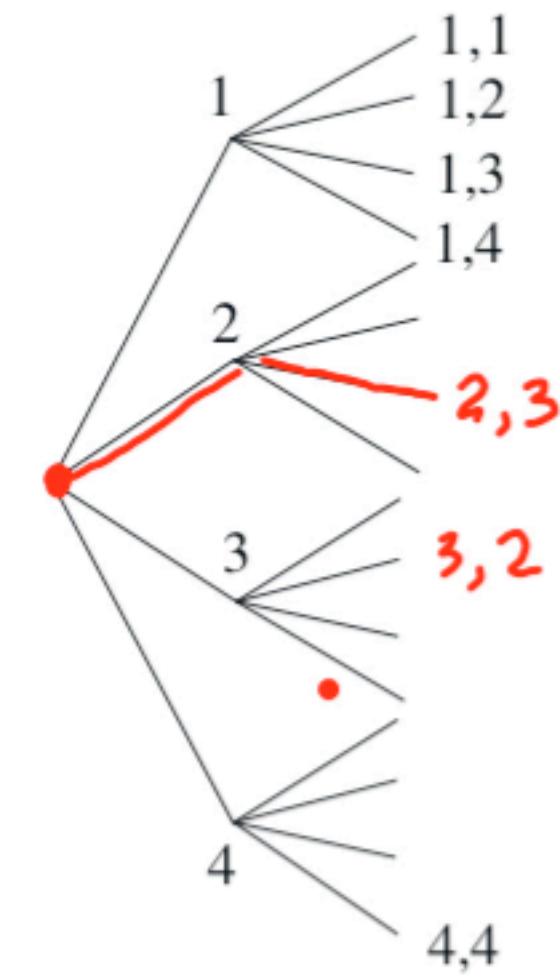
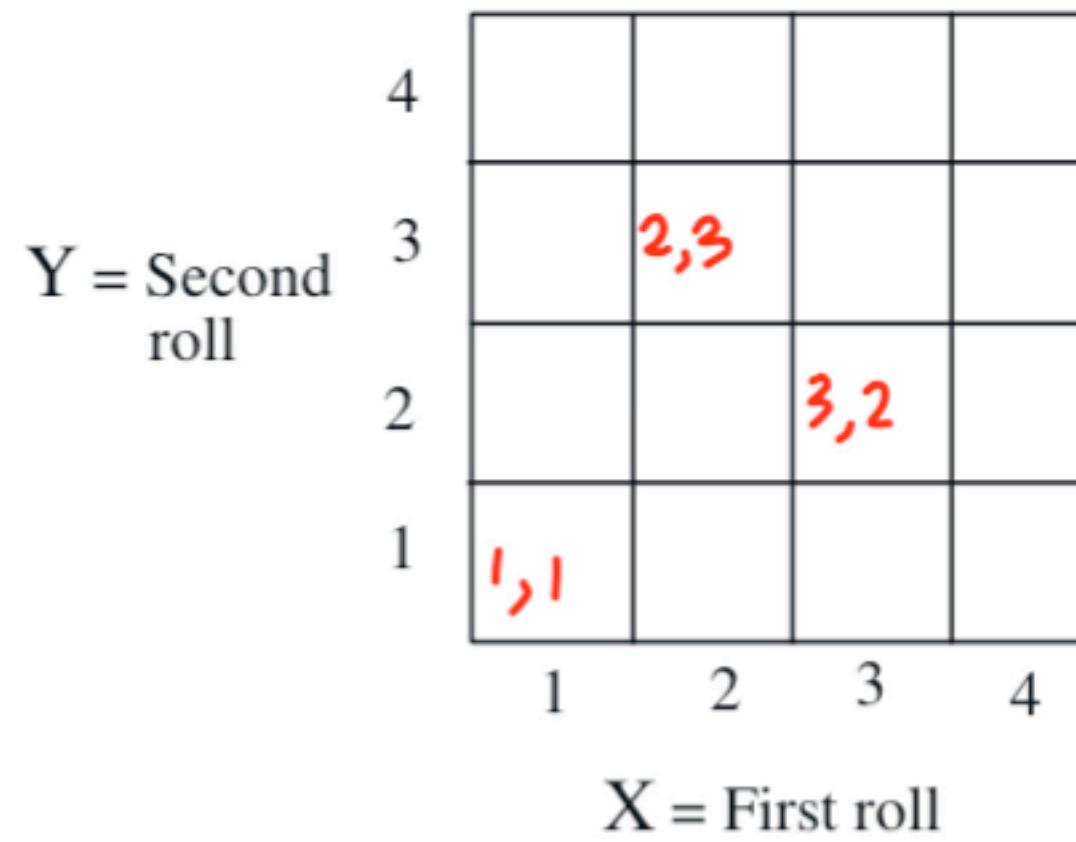


- H and rains Ω
- H and no rain
- T and rains
- T and no rain

Sample space: discrete/finite example

- Two rolls of a tetrahedral die

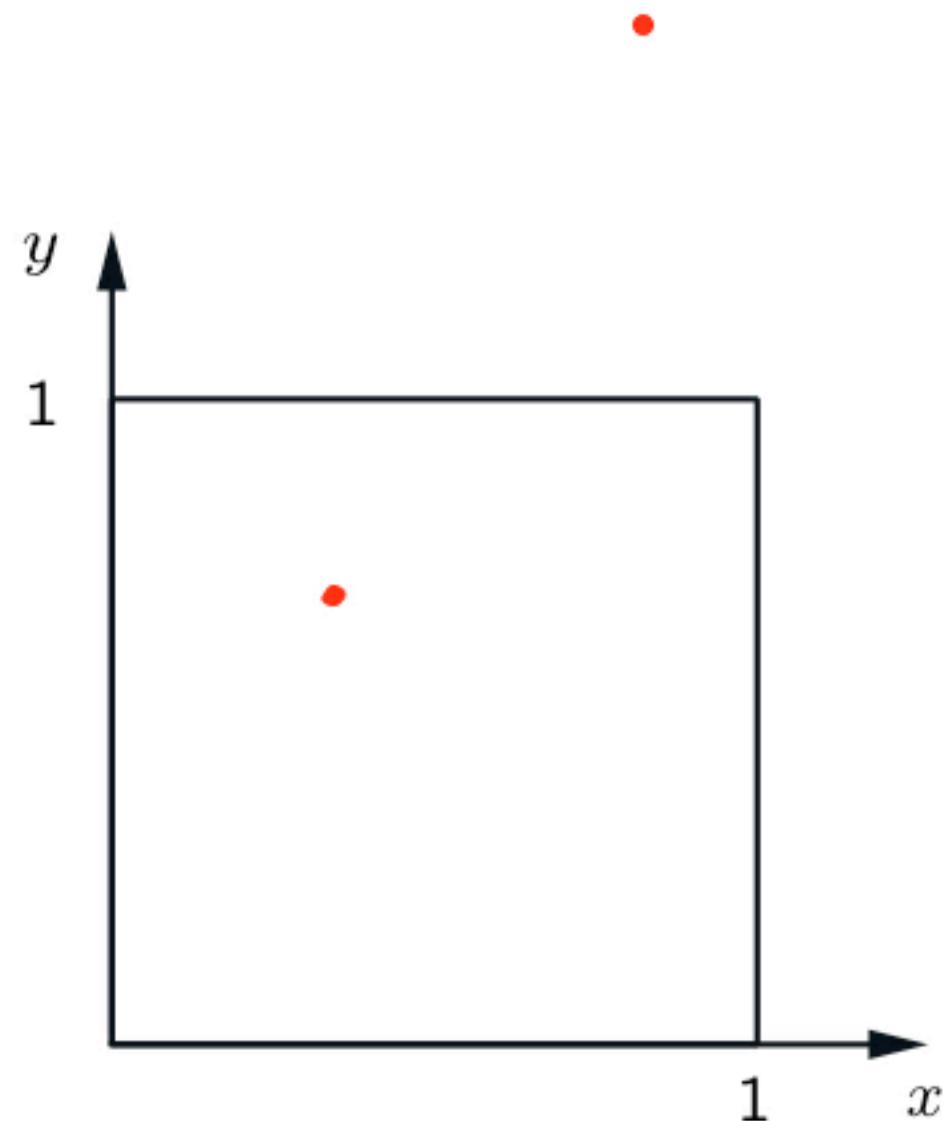
sequential description



Tree

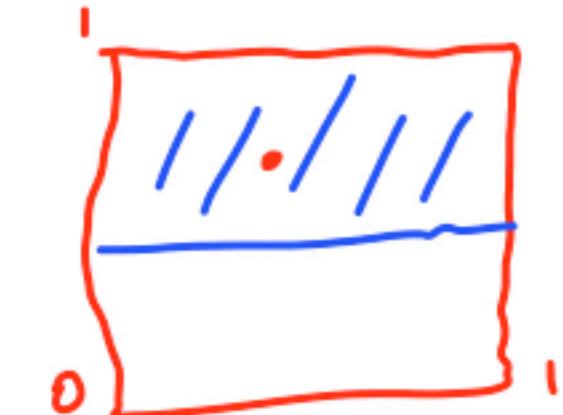
Sample space: continuous example

- (x, y) such that $0 \leq x, y \leq 1$



Probability axioms

- Event: a subset of the sample space
 - Probability is assigned to events



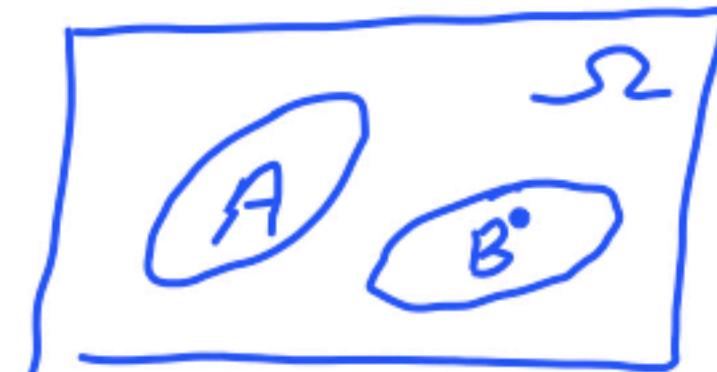
$\Omega(A)$



- **Axioms:**

- Nonnegativity: $P(A) \geq 0$
- Normalization: $P(\Omega) = 1$
- (Finite) additivity: (to be strengthened later)
If $A \cap B = \emptyset$, then $P(A \cup B) = P(A) + P(B)$

empty set



Some simple consequences of the axioms

Axioms

$$P(A) \geq 0$$

$$P(\Omega) = 1$$

For disjoint events:

$$P(A \cup B) = P(A) + P(B)$$

Consequences

$$P(A) \leq 1$$

$$P(\emptyset) = 0$$

$$P(A) + P(A^c) = 1$$

$$P(A \cup B \cup C) = P(A) + P(B) + P(C)$$

and similarly for k disjoint events

$$\begin{aligned} P(\{s_1, s_2, \dots, s_k\}) &= P(\{s_1\}) + \dots + P(\{s_k\}) \\ &= P(s_1) + \dots + P(s_k) \end{aligned}$$

Some simple consequences of the axioms

Axioms

(a) $P(A) \geq 0$



$$A \cup A^c = \Omega$$

(b) $P(\Omega) = 1$

$$1 \stackrel{(b)}{=} P(\Omega) = P(A \cup A^c)$$

For disjoint events:

(c) $P(A \cup B) = P(A) + P(B)$

$$\stackrel{(c)}{=} P(A) + P(A^c) \stackrel{(a)}{\leq} 1$$

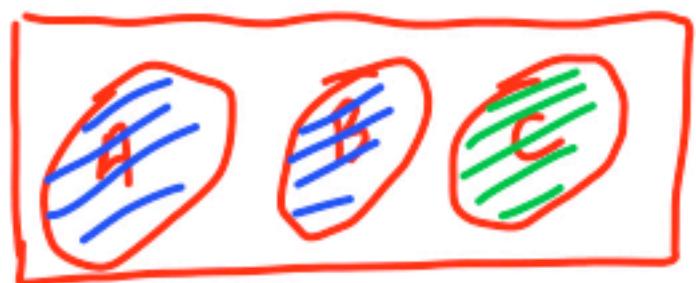
$$P(A) = 1 - \underbrace{P(A^c)}_{\leq 1} \leq 1$$

$$1 = P(\Omega) + P(\Omega^c)$$

$$1 = 1 + P(\emptyset) \Rightarrow P(\emptyset) = 0.$$

Some simple consequences of the axioms

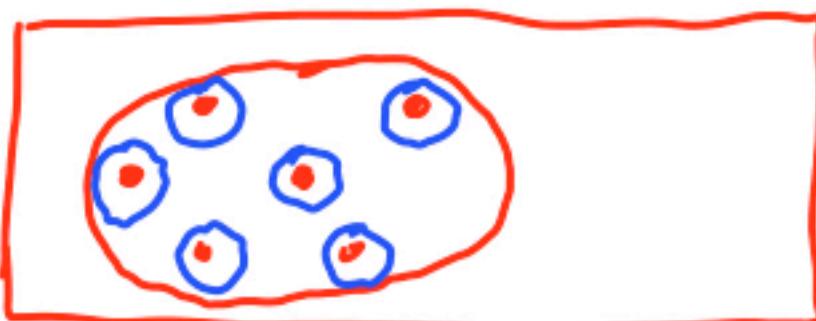
- A, B, C disjoint: $P(A \cup B \cup C) = P(A) + P(B) + P(C)$



$$\begin{aligned}P(A \cup B \cup C) &= P((A \cup B) \cup C) = P(A \cup B) + P(C) \\&= P(A) + P(B) + P(C)\end{aligned}$$

If A_1, \dots, A_k disjoint $\Rightarrow P(A_1 \cup \dots \cup A_k) = \sum_{i=1}^k P(A_i)$

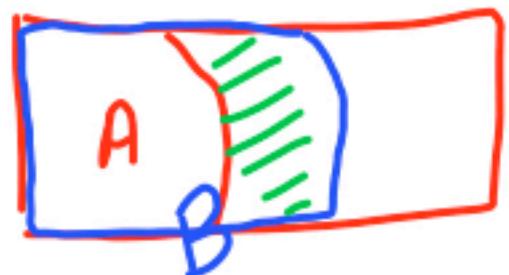
- $P(\{s_1, s_2, \dots, s_k\}) = P(\{s_1\} \cup \{s_2\} \cup \dots \cup \{s_k\})$



$$\begin{aligned}&= P(\{s_1\}) + \dots + P(\{s_k\}) \\&= P(s_1) + \dots + P(s_k)\end{aligned}$$

More consequences of the axioms

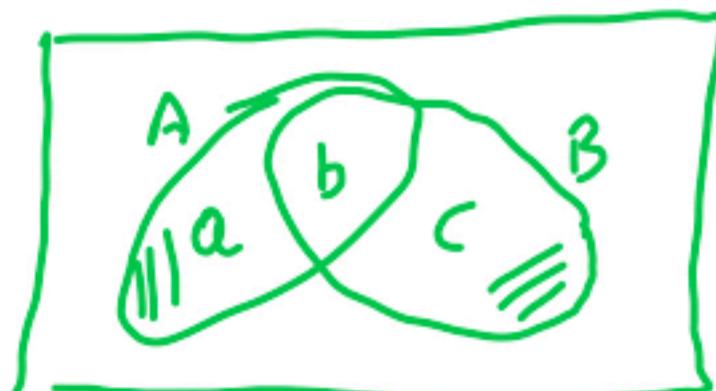
- If $A \subset B$, then $P(A) \leq P(B)$



$$B = A \cup (B \cap A^c)$$

$$P(B) = P(A) + \underline{P(B \cap A^c)} \geq P(A)$$

- $P(A \cup B) = P(A) + P(B) - \overbrace{P(A \cap B)}^{> 0}$



$$a = P(A \cap B^c) \quad b = P(A \cap B) \quad c = P(B \cap A^c)$$

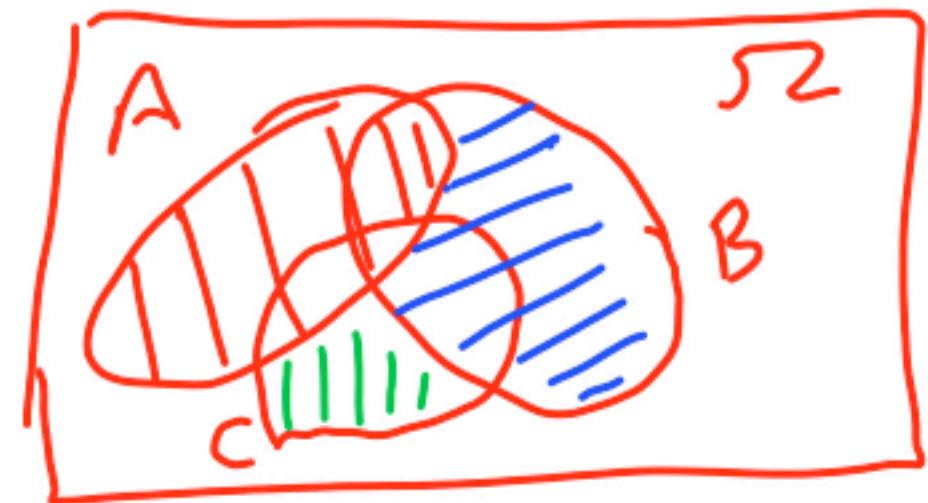
$$P(A \cup B) = a + b + c$$

$$P(A) + P(B) - P(A \cap B) = (a+b) + (\cancel{a+c}) - b$$

- $P(A \cup B) \leq P(A) + P(B)$ union bound = $a + b + c$

More consequences of the axioms

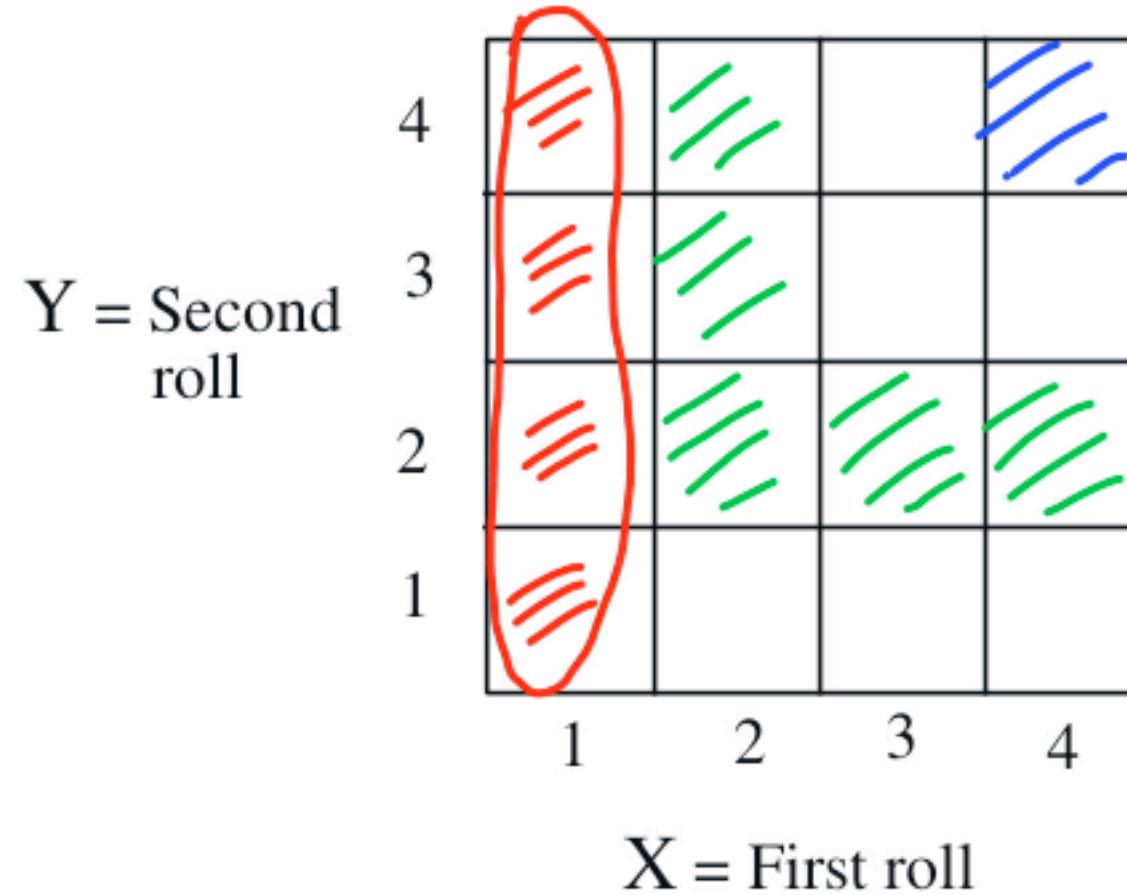
- $P(A \cup B \cup C) = P(A) + \underbrace{P(A^c \cap B)}_{\text{green}} + \underbrace{P(A^c \cap B^c \cap C)}_{\text{green}}$ • ↗



$$\begin{aligned} P(A \cup B \cup C) &= \\ &= A \cup \underbrace{(B \cap A^c)}_{\text{green}} \cup \underbrace{(C \cap A^c \cap B^c)}_{\text{green}} \end{aligned}$$

Probability calculation: discrete/finite example

- Two rolls of a tetrahedral die
- Let every possible outcome have probability $1/16$



- $P(X = 1) = 4 \cdot \frac{1}{16} = \frac{1}{4}$

Let $Z = \min(X, Y)$

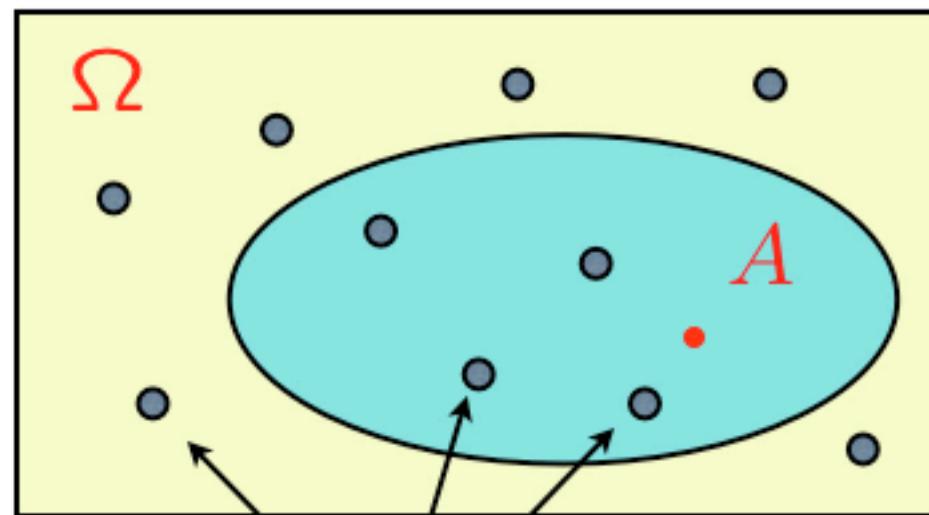
- $\underline{P(Z = 4)} = 1/16$
- $P(Z = 2) = 5 \cdot \frac{1}{16}.$

Discrete uniform law

✓finite

- Assume Ω consists of n equally likely elements
- Assume A consists of k elements

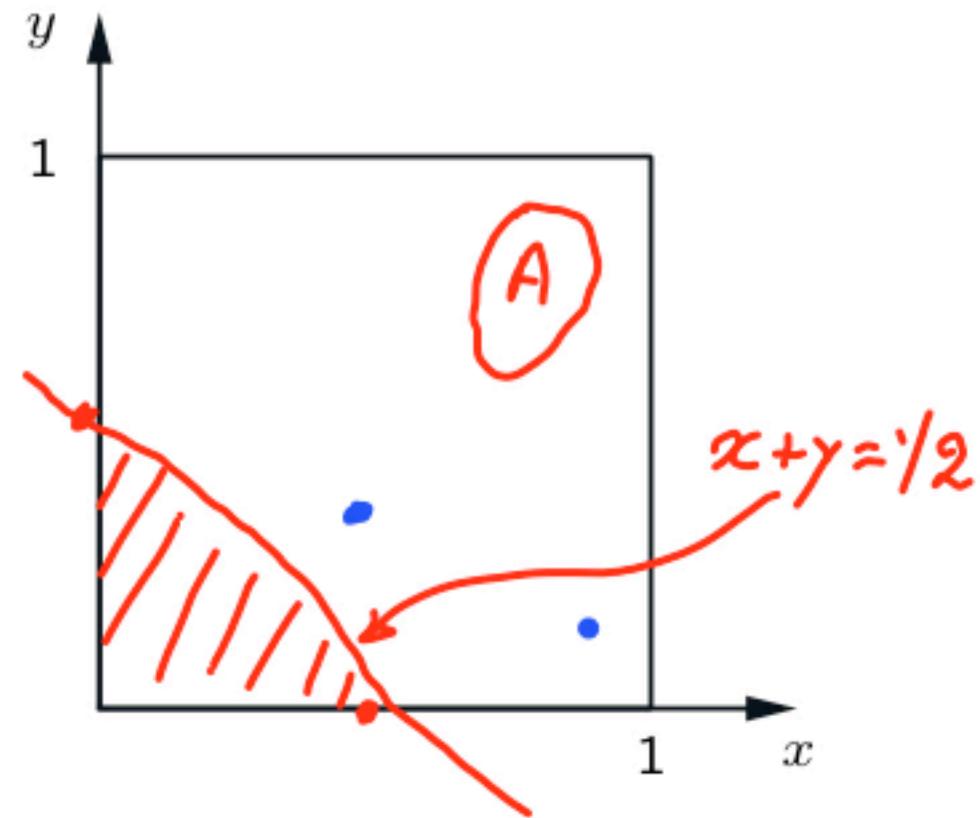
$$P(A) = k \cdot \frac{1}{n}$$



$$\text{prob} = \frac{1}{n}$$

Probability calculation: continuous example

- (x, y) such that $0 \leq x, y \leq 1$
- Uniform probability law: Probability = Area



$$P\left(\{(x, y) \mid x + y \leq 1/2\}\right) = \frac{1}{2} \cdot \frac{1}{2} \cdot \frac{1}{2} = \frac{1}{8}$$

$$P\left(\{(0.5, 0.3)\}\right) = 0$$

Probability calculation steps

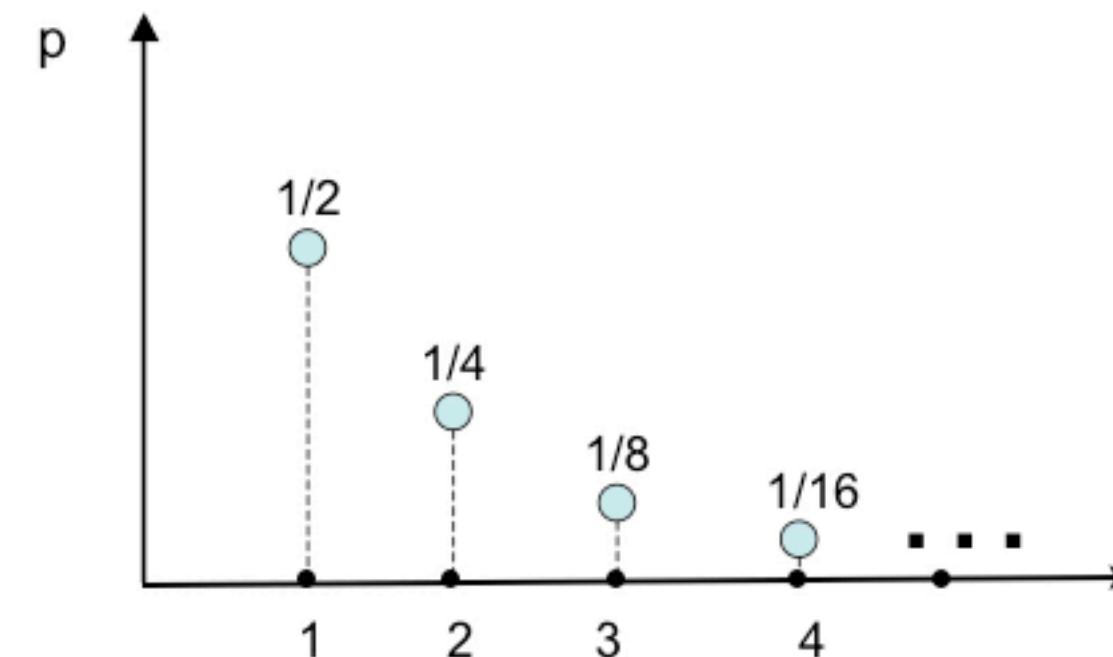
- Specify the sample space
- Specify a probability law
- Identify an event of interest
- Calculate...

Probability calculation: discrete but infinite sample space

- Sample space: $\{1, 2, \dots\}$

– We are given $P(n) = \frac{1}{2^n}, n = 1, 2, \dots$

$$\sum_{n=1}^{\infty} \frac{1}{2^n} = \frac{1}{2} \sum_{n=0}^{\infty} \frac{1}{2^n} = \frac{1}{2} \cdot \frac{1}{1 - (\frac{1}{2})} = 1$$



- $P(\text{outcome is even}) = P(\{2, 4, 6, \dots\})$

$$= P(\{2\} \cup \{4\} \cup \{6\} \cup \dots) \quad \textcircled{=} P(2) + P(4) + P(6) + \dots$$

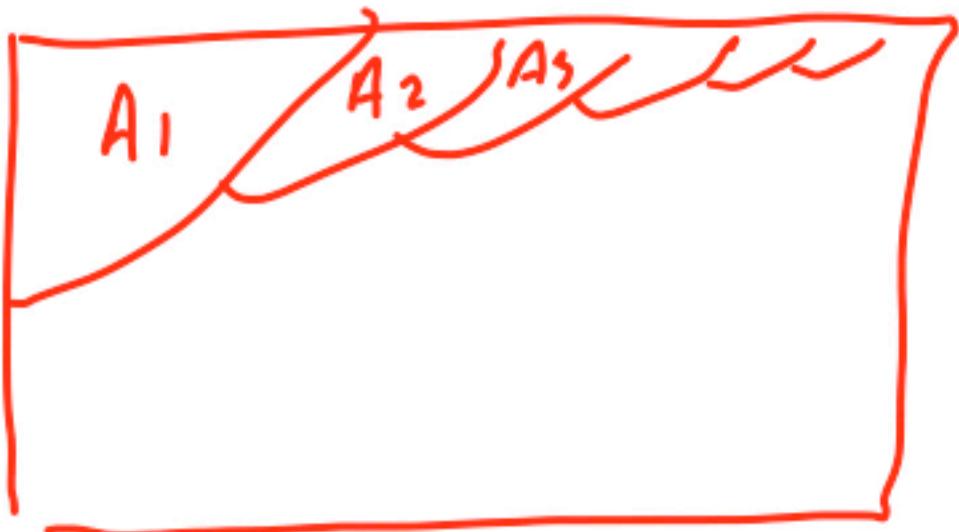
$$= \frac{1}{2^2} + \frac{1}{2^4} + \frac{1}{2^6} + \dots = \frac{1}{4} \left(1 + \frac{1}{4} + \frac{1}{4^2} + \dots \right) = \frac{1}{4} \cdot \frac{1}{1 - \frac{1}{4}} = \frac{1}{3}$$

Countable additivity axiom

- Strengthens the finite additivity axiom

Countable Additivity Axiom:

If A_1, A_2, A_3, \dots is an infinite sequence of disjoint events,
then $P(A_1 \cup A_2 \cup A_3 \cup \dots) = P(A_1) + P(A_2) + P(A_3) + \dots$



Mathematical subtleties

Countable Additivity Axiom:

If A_1, A_2, A_3, \dots is an infinite sequence of disjoint events,
then $P(A_1 \cup A_2 \cup A_3 \cup \dots) = P(A_1) + P(A_2) + P(A_3) + \dots$.


$$I = P(\Omega) = P\left(\bigcup \{(x,y)\}\right) \stackrel{?}{=} \sum P(\{(x,y)\}) = \sum 0 = 0$$

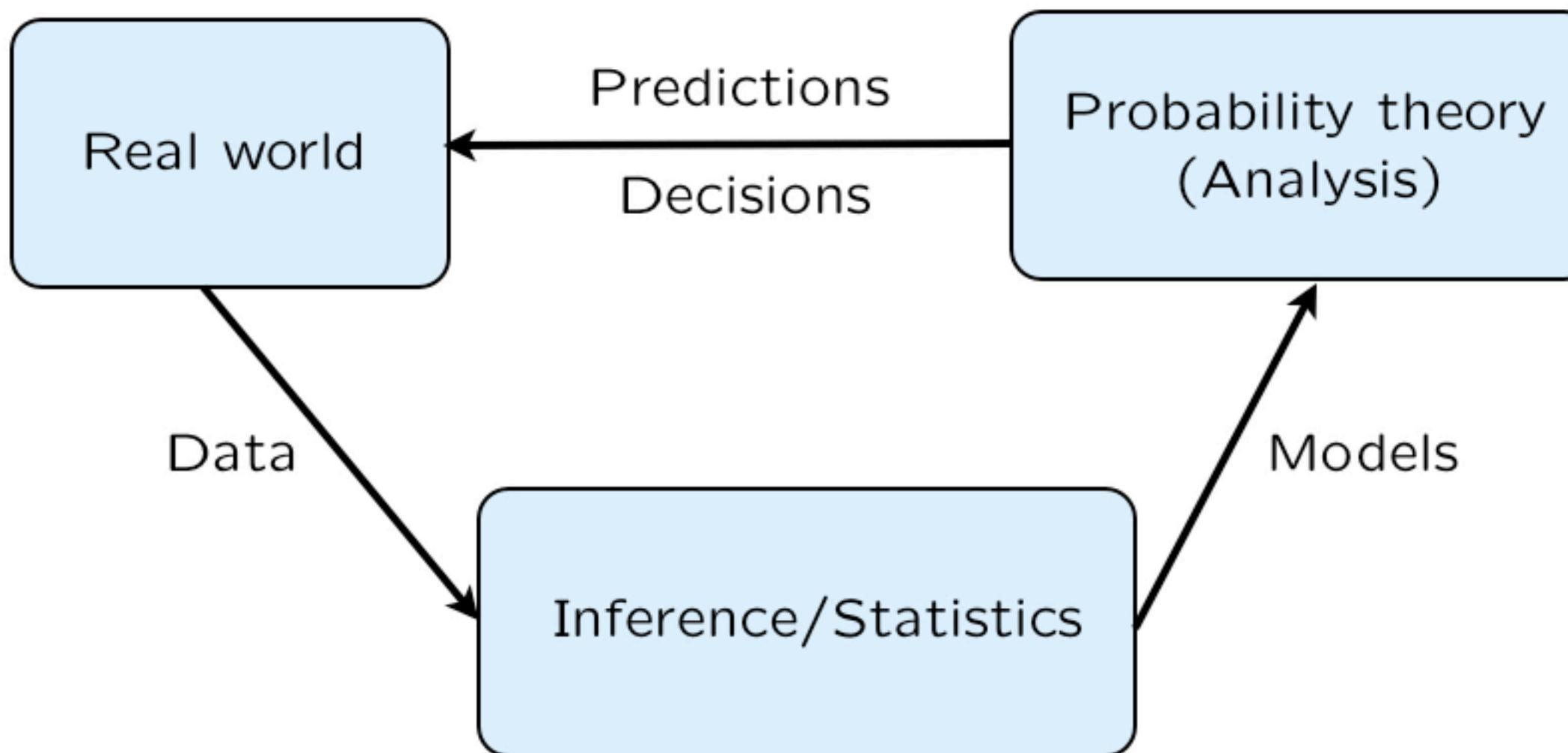
- Additivity holds only for “countable” sequences of events
- The unit square (similarly, the real line, etc.) is **not countable** (its elements cannot be arranged in a sequence)
- “Area” is a legitimate probability law on the unit square, as long as we do not try to assign probabilities/areas to “very strange” sets

Interpretations of probability theory

- A narrow view: a branch of math
 - Axioms \Rightarrow theorems “**Thm:**” “Frequency” of event A “is” $P(A)$
- Are probabilities frequencies?
 - $P(\text{coin toss yields heads}) = 1/2$
 - $P(\text{the president of ... will be reelected}) = 0.7$
- Probabilities are often interpreted as:
 - Description of beliefs
 - Betting preferences

The role of probability theory

- A framework for analyzing phenomena with uncertain outcomes
 - Rules for consistent reasoning
 - Used for predictions and decisions



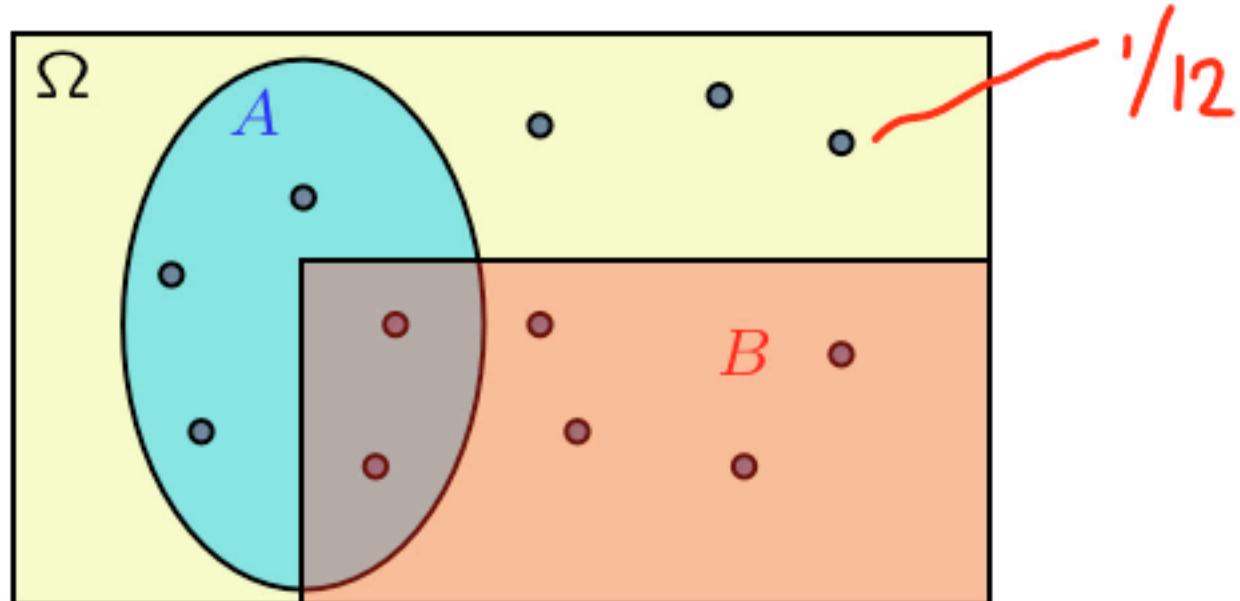
LECTURE 2: Conditioning and Bayes' rule

- Conditional probability
- Three **important** tools:
 - Multiplication rule
 - Total probability theorem
 - Bayes' rule (\rightarrow inference)

The idea of conditioning

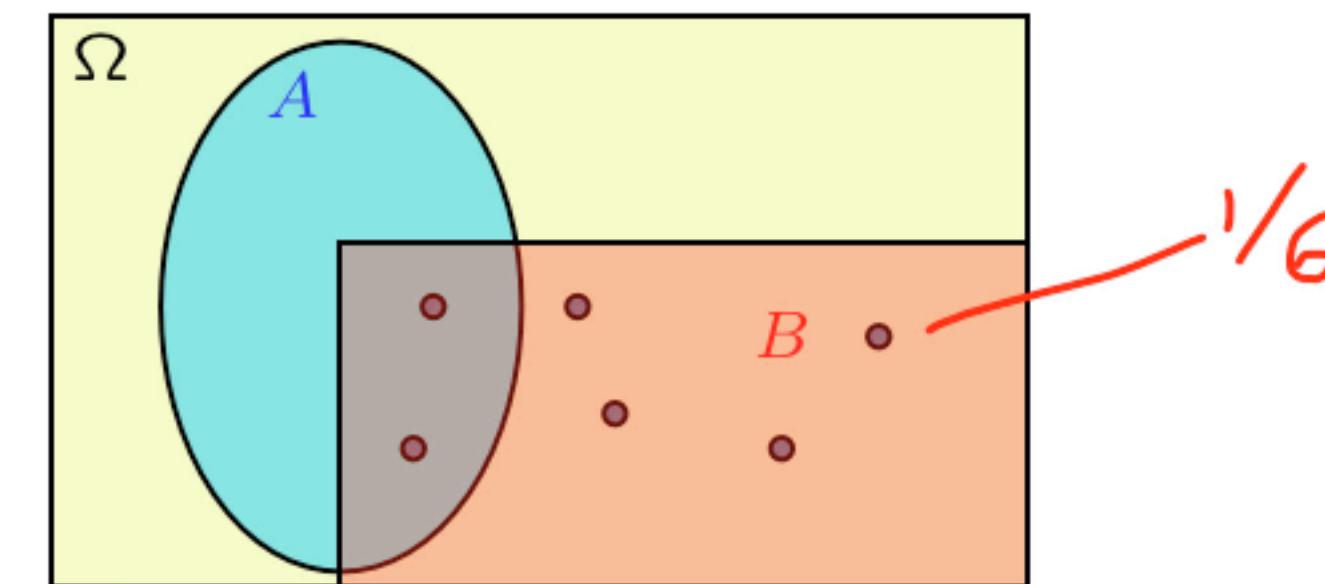
Use new information to revise a model

Assume 12 equally likely outcomes

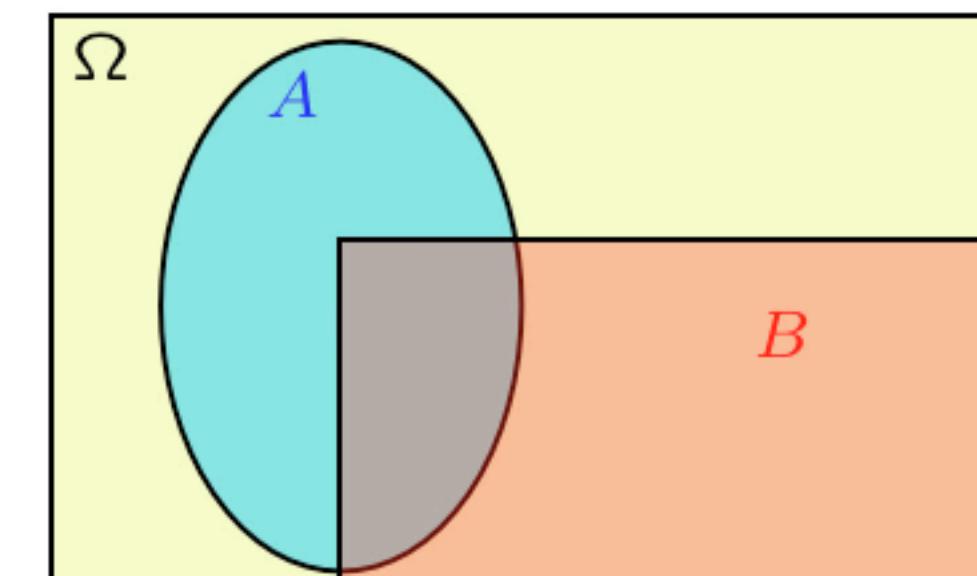
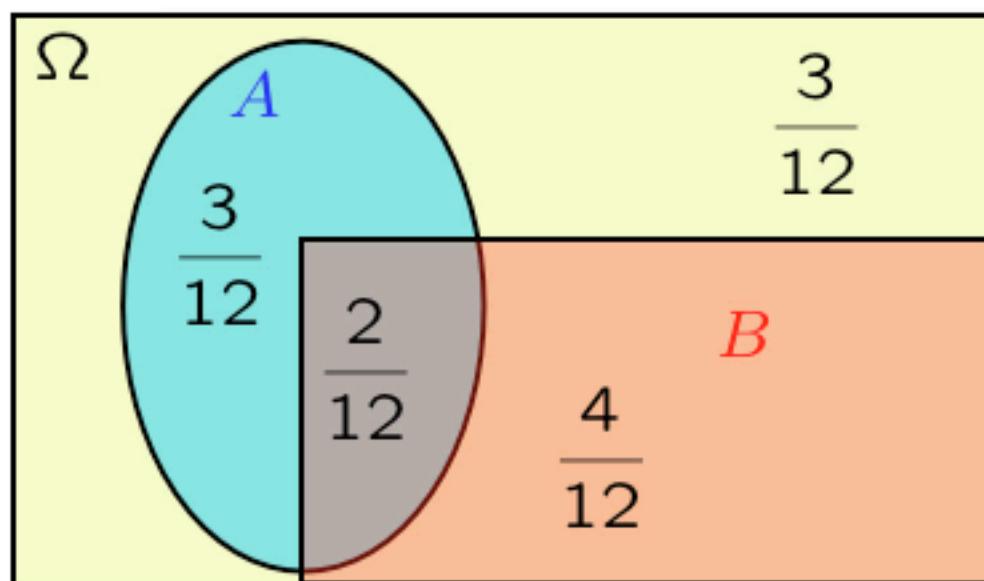


$$P(A) = \frac{5}{12} \quad P(B) = \frac{6}{12}$$

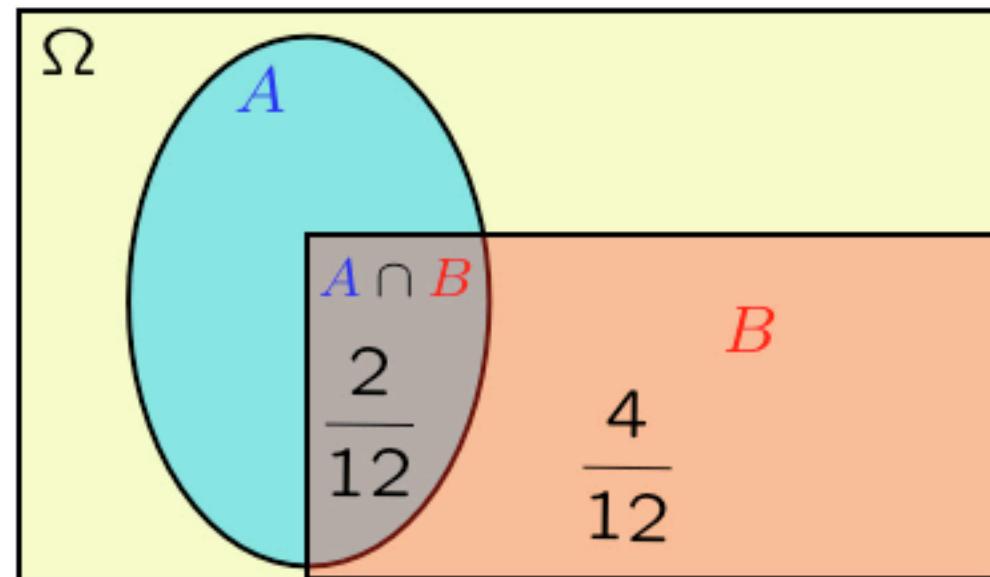
If told B occurred:



$$P(A | B) = \frac{2}{6} = \frac{1}{3} \quad P(B | B) = \underline{\underline{1}}$$



Definition of conditional probability



- $\underline{P(A | B)}$ = “probability of A , given that B occurred”

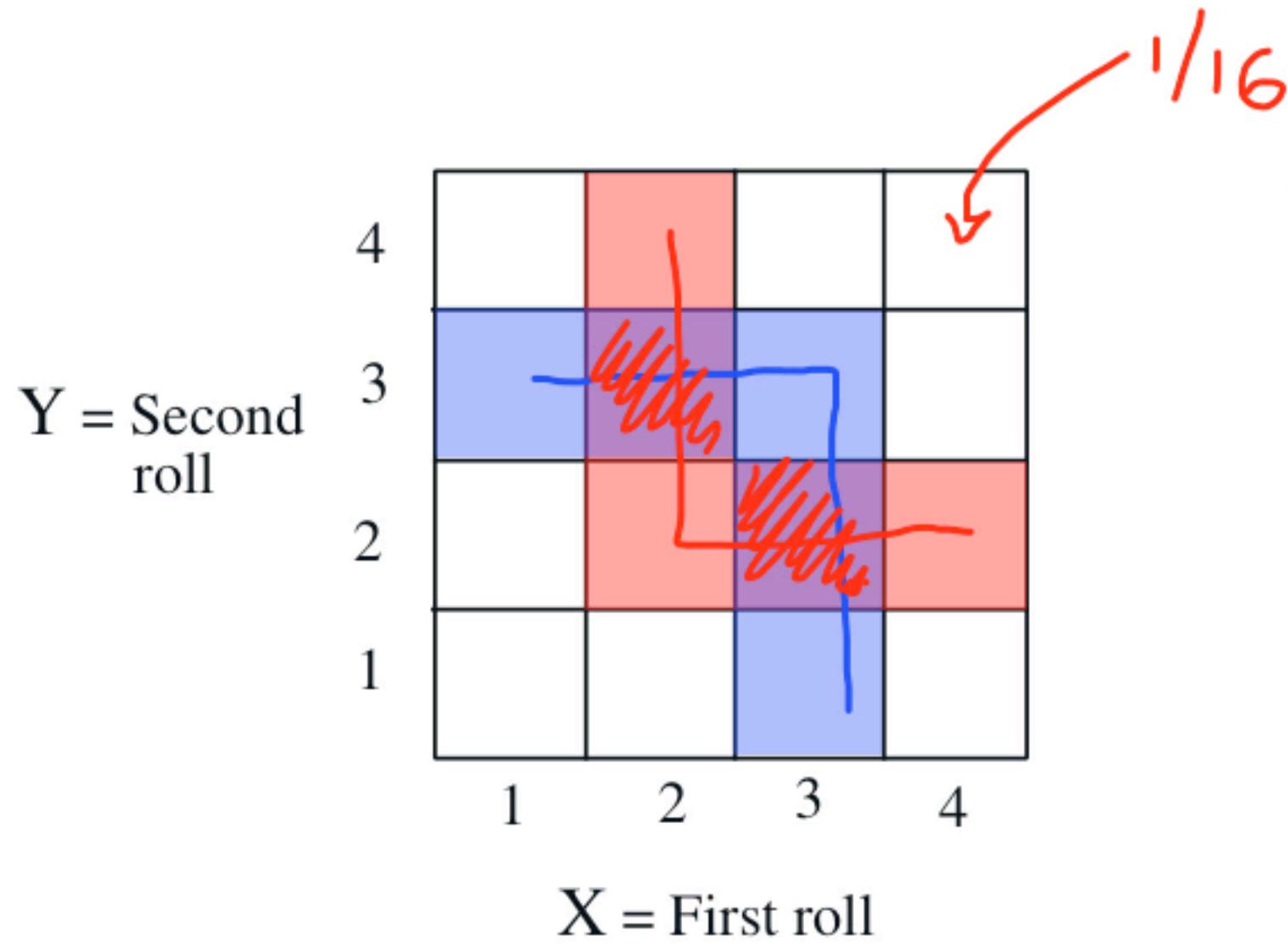
Def.

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

defined only when $P(B) > 0$

$$= \frac{2/12}{6/12} = \frac{1}{3}$$

Example: two rolls of a 4-sided die



- Let B be the event: $\min(X, Y) = 2$

Let $M = \max(X, Y)$

$$P(M = 1 | B) = 0$$

$$P(\underline{M = 3} | B) = \frac{P(M=3 \text{ and } B)}{P(B)}$$

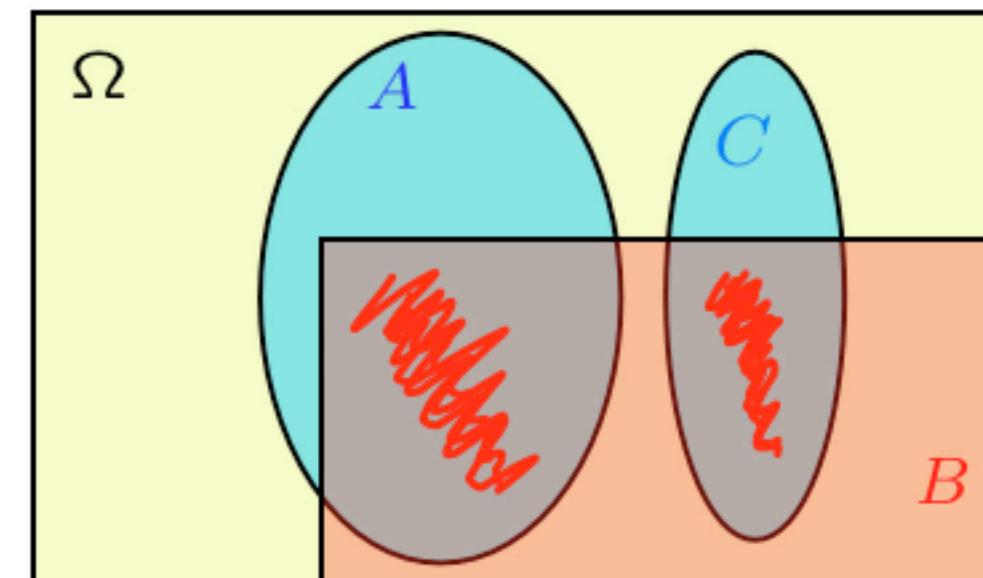
$$= \frac{2/16}{5/16} = \frac{2}{5}$$

Conditional probabilities share properties of ordinary probabilities

$$P(A | B) \geq 0 \quad \text{assuming } P(B) > 0$$

$$P(\Omega | B) = \frac{P(\Omega \cap B)}{P(B)} = \frac{P(B)}{P(B)} = 1$$

$$P(B | B) = \frac{P(B \cap B)}{P(B)} = 1$$



If $A \cap C = \emptyset$, then $P(A \cup C | B) = P(A | B) + P(C | B)$

$$= \frac{P((A \cup C) \cap B)}{P(B)} = \frac{P((A \cap B) \cup (C \cap B))}{P(B)} = \frac{P(A \cap B) + P(C \cap B)}{P(B)} =$$

$= P(A|B) + P(C|B)$ also finite
countable additivity

Models based on conditional probabilities

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

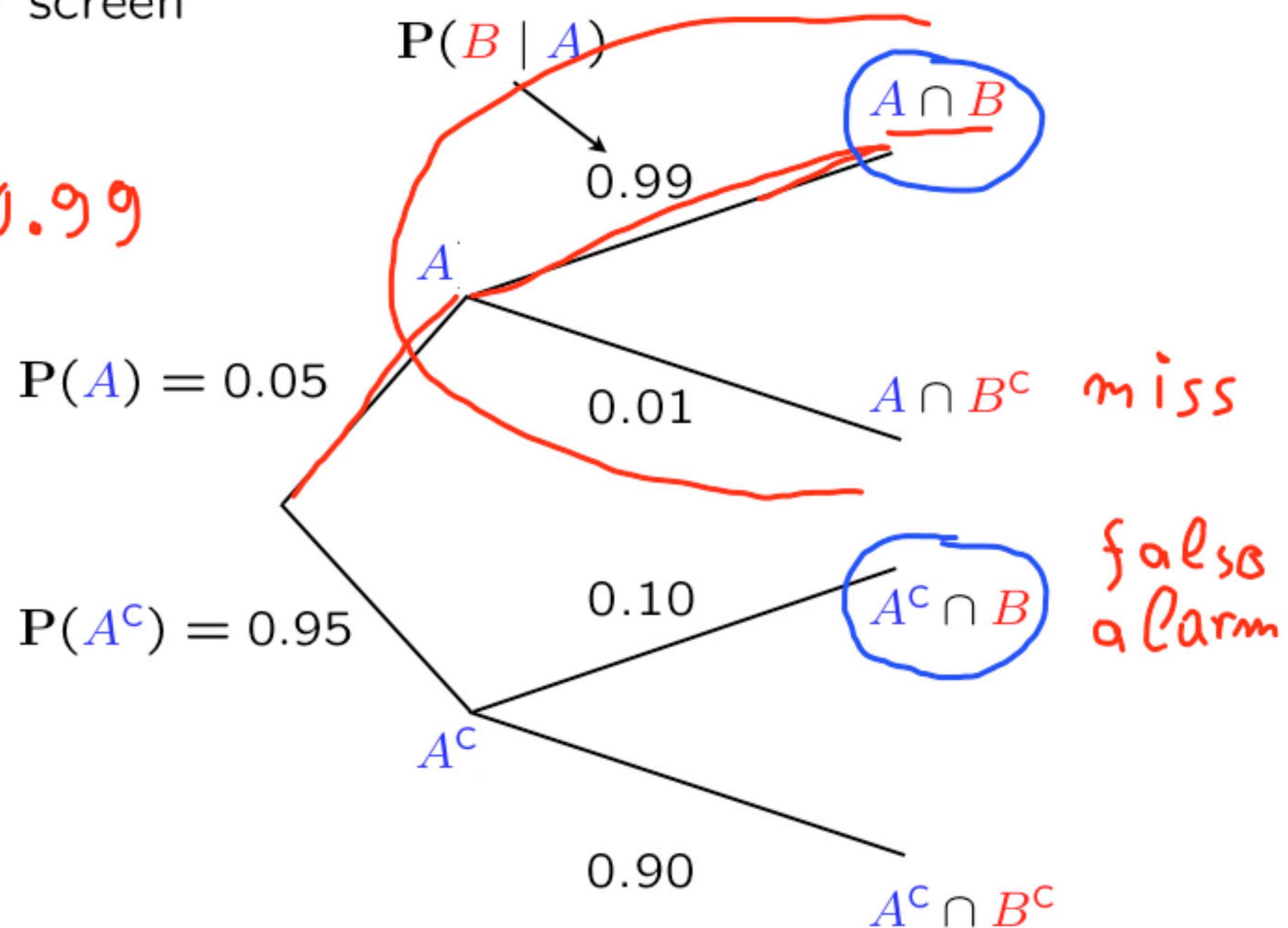
Event A : Airplane is flying above

Event B : Something registers on radar screen

- $P(A \cap B) = P(A) \cdot P(B | A) = 0.05 \cdot 0.99$

- $P(B) = 0.05 \cdot 0.99 + 0.95 \cdot 0.1 = 0.1445$

- $P(A | B) = \frac{0.05 \cdot 0.99}{0.1445} = 0.34$



The multiplication rule

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

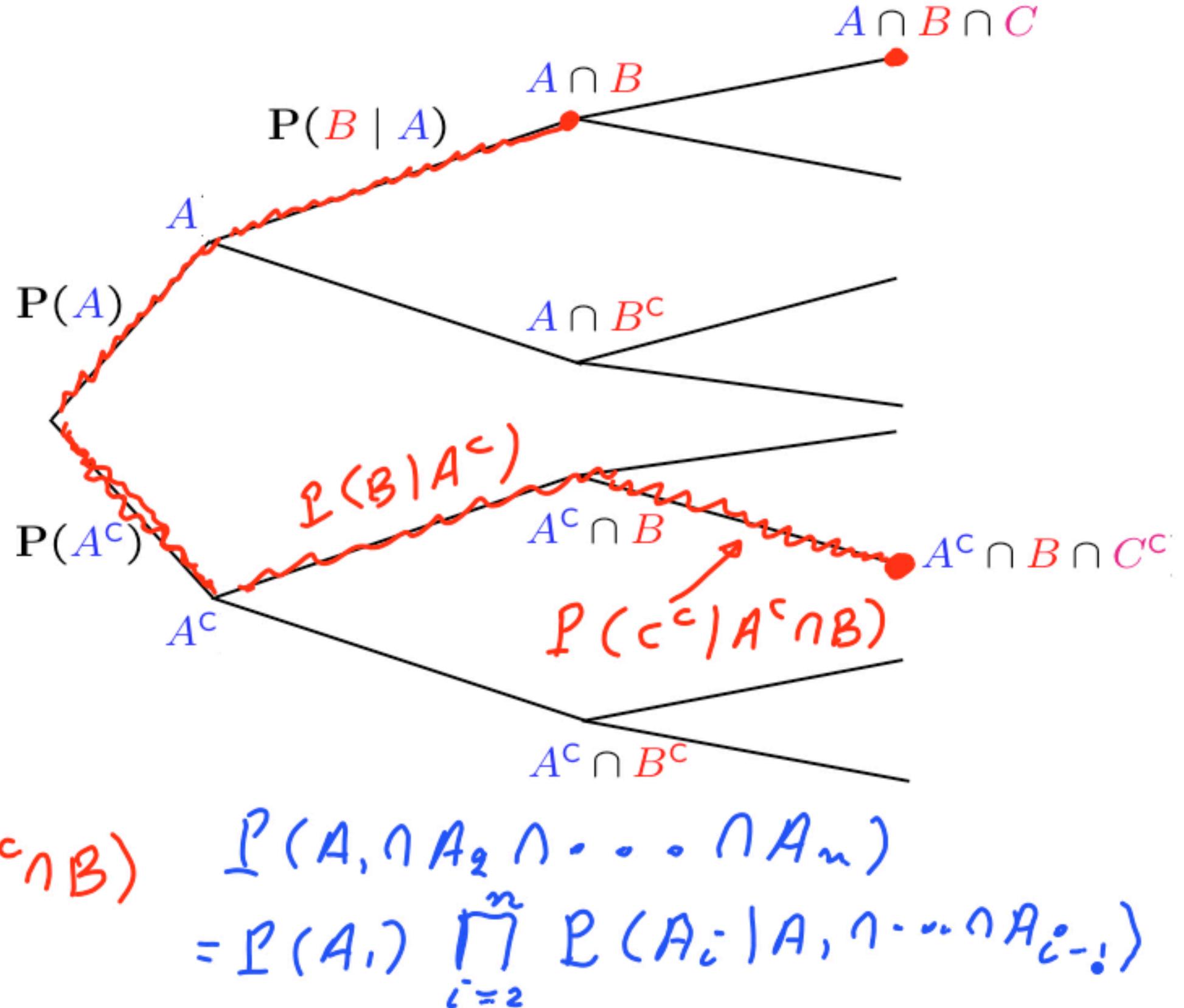
$$P(A \cap B) = P(B) P(A | B)$$

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

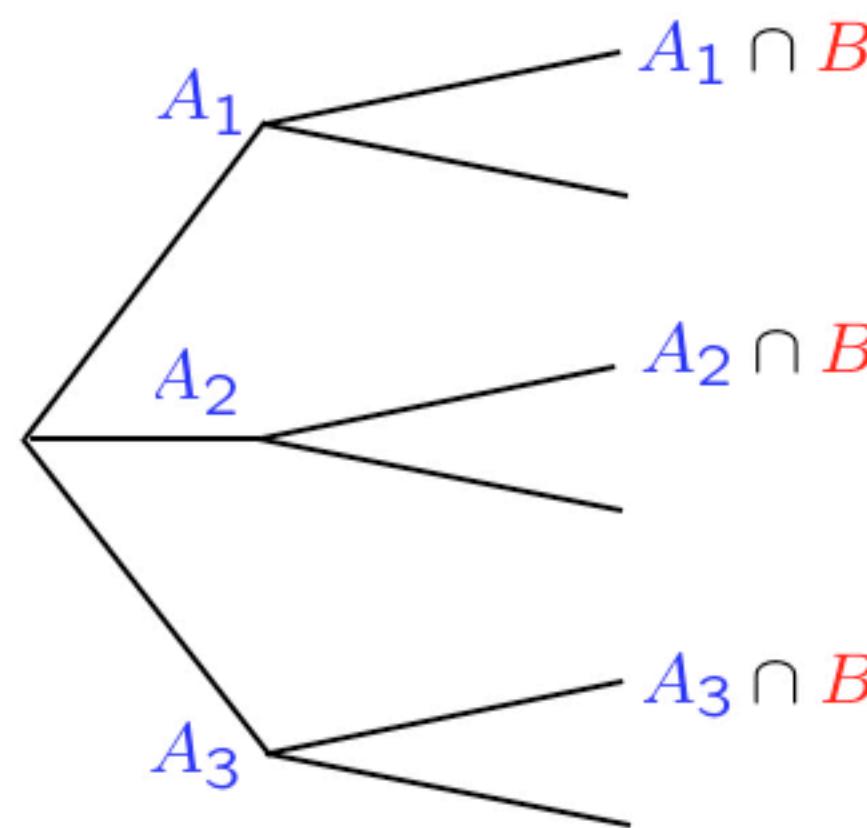
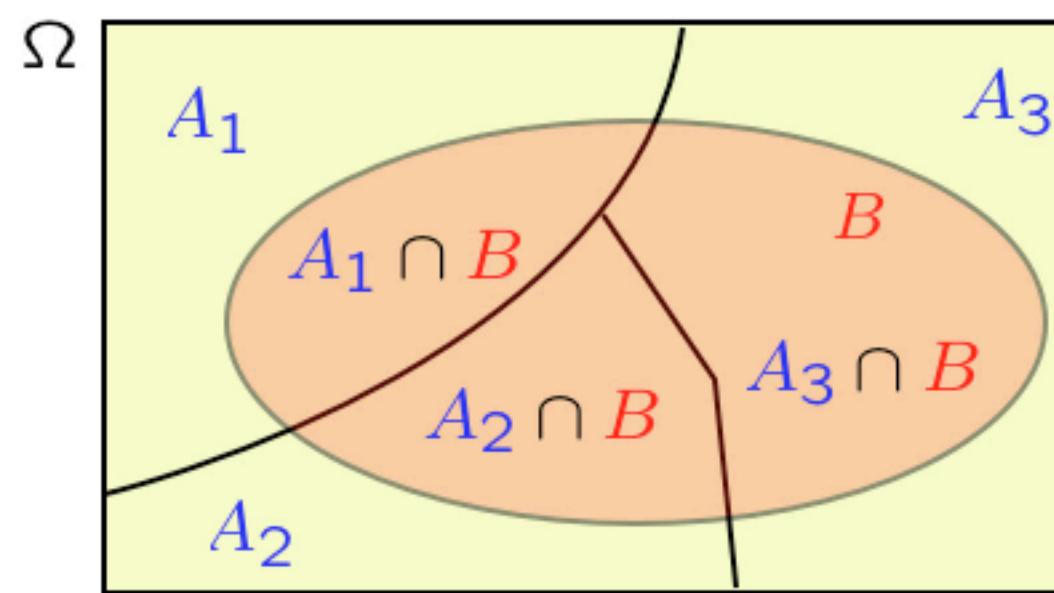
$$P(\underline{A^c \cap B} \cap \underline{C^c}) =$$

$$= P(A^c \cap B) P(C^c | A^c \cap B)$$

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



Total probability theorem



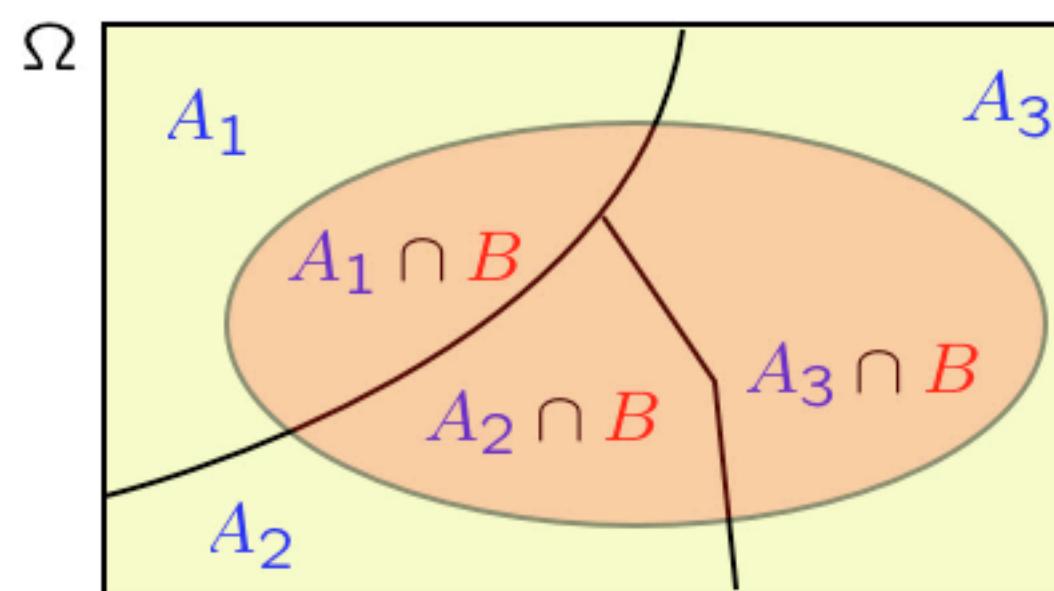
- Partition of sample space into A_1, A_2, A_3, \dots
- Have $P(A_i)$, for every i
- Have $P(B | A_i)$, for every i

$$\begin{aligned} P(B) &= P(B \cap A_1) + P(B \cap A_2) + P(B \cap A_3) \\ &= P(A_1)P(B | A_1) + \dots + \dots \end{aligned}$$

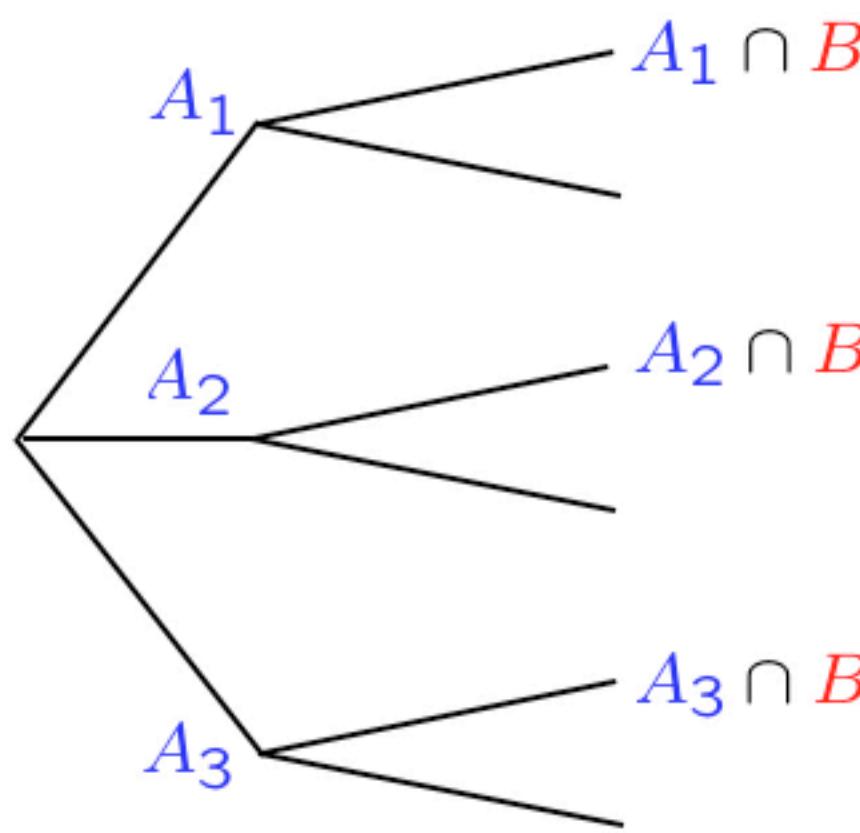
$$P(B) = \sum_i P(A_i)P(B | A_i)$$

$\sum_i P(A_i) = 1$ weights
weighted average
of $P(B | A_i)$

Bayes' rule



- Partition of sample space into A_1, A_2, A_3
 - Have $P(A_i)$, for every i initial “beliefs”
 - Have $P(B | A_i)$, for every i
- revised “beliefs,” given that B occurred:



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

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

Bayes' rule and inference

- Thomas Bayes, presbyterian minister (c. 1701-1761)
- “Bayes’ theorem,” published posthumously
- systematic approach for incorporating new evidence
- Bayesian inference
 - initial beliefs $P(A_i)$ on possible causes of an observed event B
 - model of the world under each A_i : $P(B | A_i)$

$$A_i \xrightarrow{\text{model}} B$$
$$P(B | A_i)$$

- draw conclusions about causes

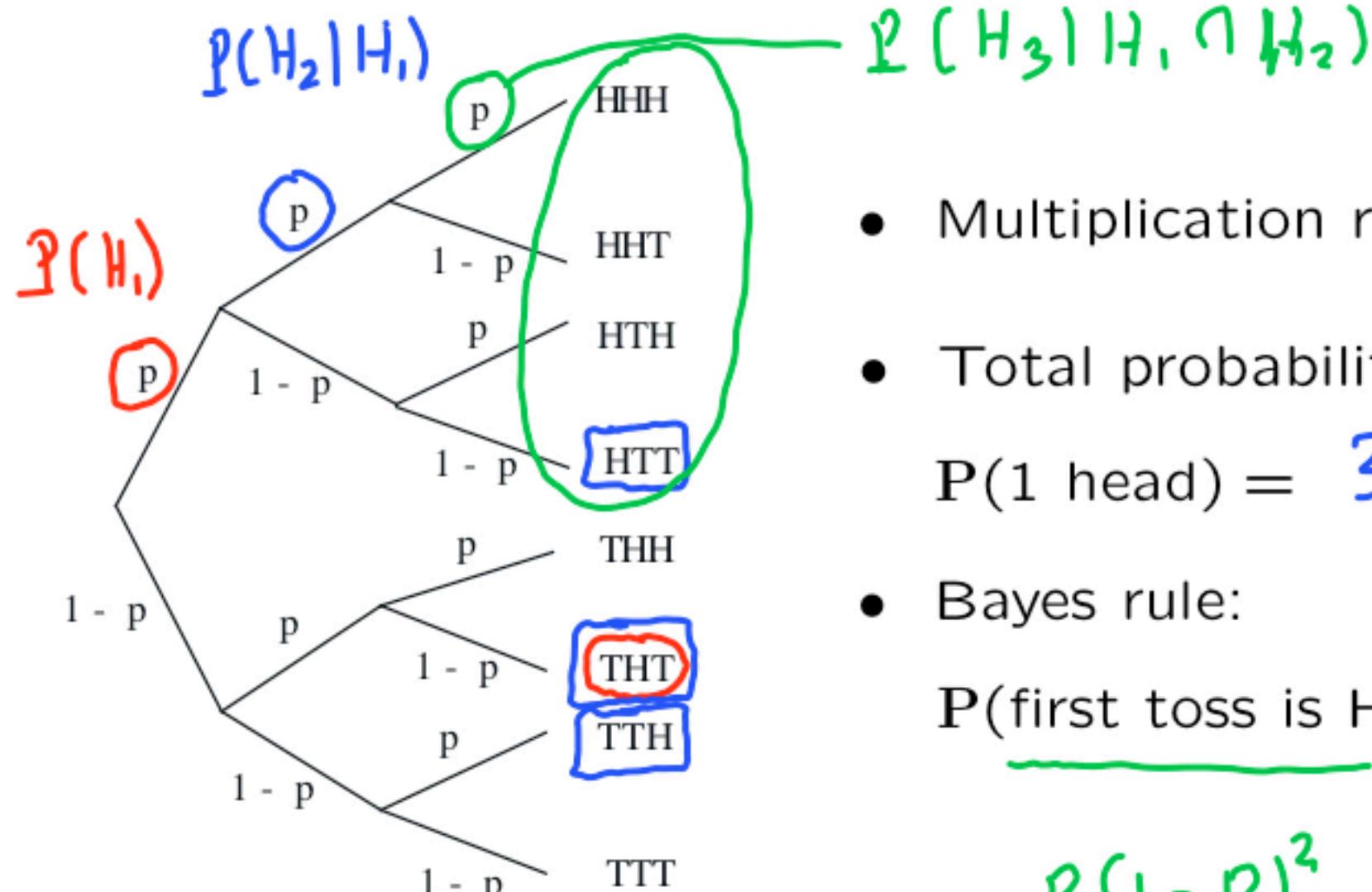
$$B \xrightarrow{\text{inference}} A_i$$
$$P(A_i | B)$$

LECTURE 3: Independence

- Independence of two events
- Conditional independence
- Independence of a collection of events
- Pairwise independence
- Reliability
- The king's sibling puzzle

A model based on conditional probabilities

- 3 tosses of a biased coin: $P(H) = p$, $P(T) = 1 - p$



$$\begin{aligned}
 P(H_2 | H_1) &= p = P(H_2 | T_1) \\
 P(H_2) &= P(H_1) P(H_2 | H_1) \\
 &\quad + P(T_1) P(H_2 | T_1) \\
 &= p
 \end{aligned}$$

- Multiplication rule: $P(\underline{THT}) = \underline{(1-p)p(1-p)}$

- Total probability:

$$P(1 \text{ head}) = \underline{3} p(1-p)^2$$

- Bayes rule:

$$\begin{aligned}
 P(\text{first toss is } H | 1 \text{ head}) &= \frac{P(H_1 \cap 1 \text{ head})}{P(1 \text{ head})} \\
 &= \frac{p(1-p)^2}{3 p(1-p)^2} = \frac{1}{3}
 \end{aligned}$$

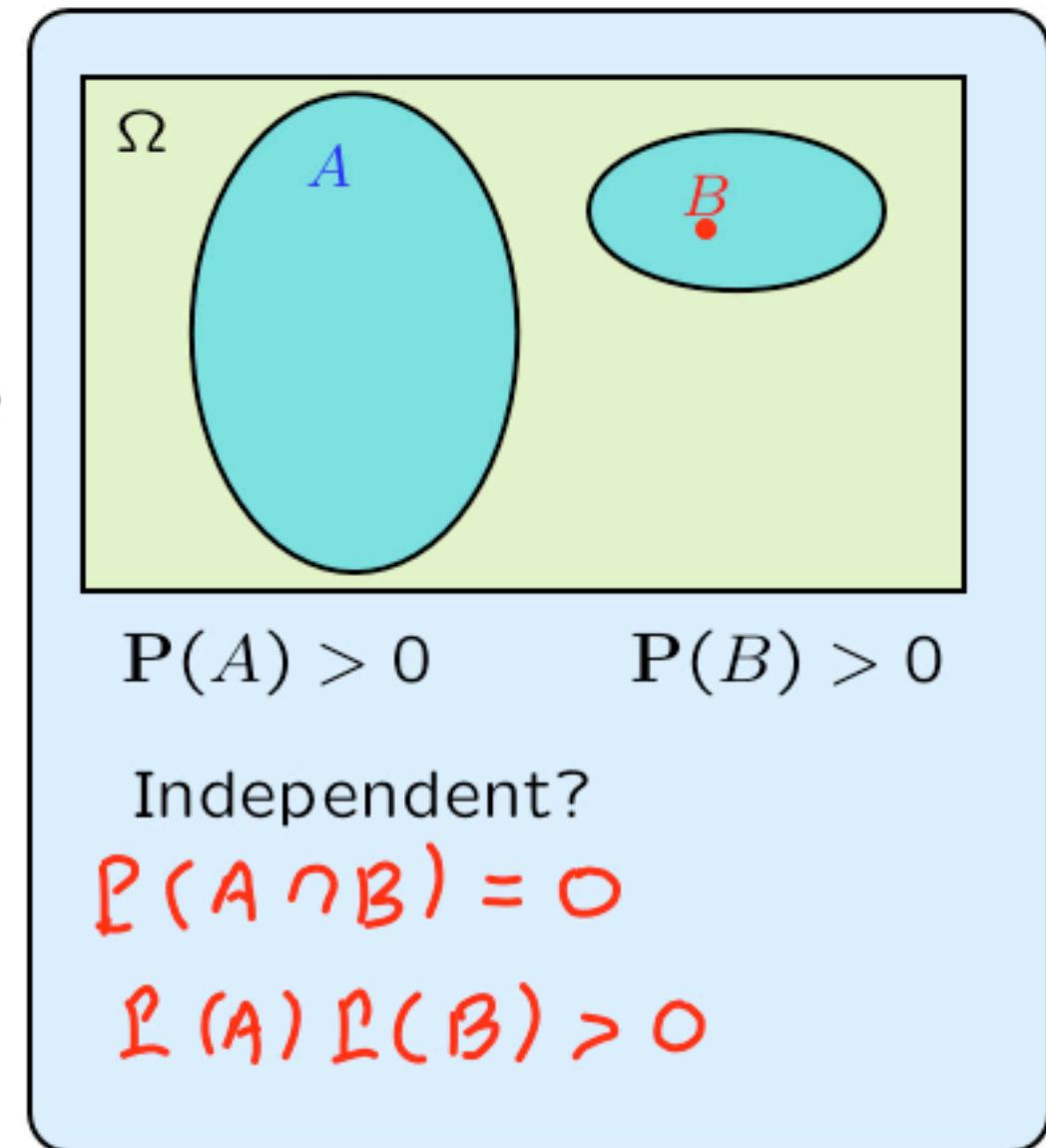
Independence of two events

- Intuitive “definition”: $P(B | A) = P(B)$
 - occurrence of A provides no new information about B

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

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

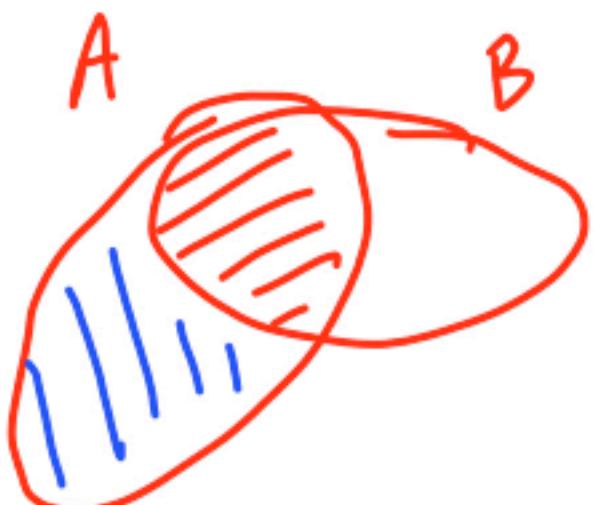
- Symmetric with respect to A and B
- implies $P(A | B) = P(A)$
- applies even if $P(A) = 0$



Independence of event complements

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

- If A and B are independent, then A and B^c are independent.
 - Intuitive argument
 - Formal proof



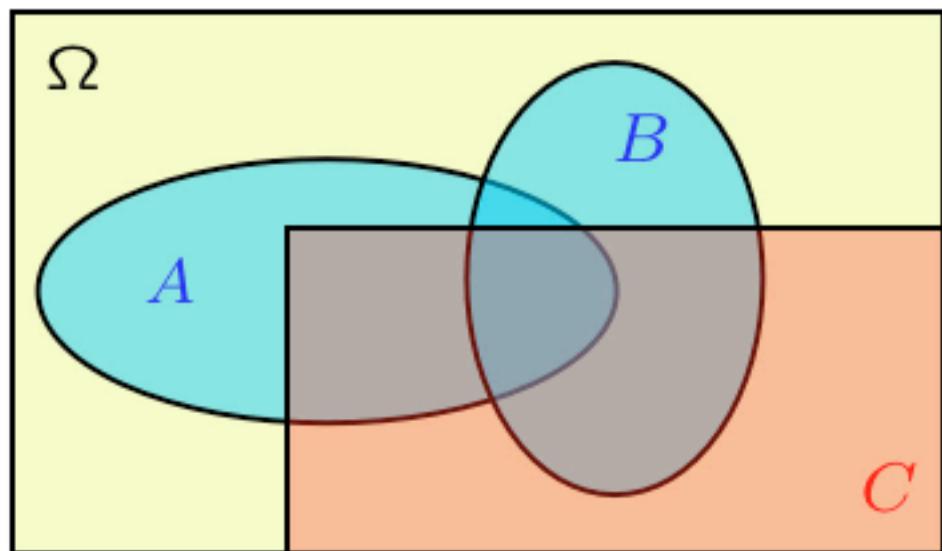
$$A = (A \cap B) \cup (A \cap B^c)$$

$$\begin{aligned}P(A) &= P(A \cap B) + P(A \cap B^c) \\&= P(A)P(B) + P(A \cap B^c)\end{aligned}$$

$$\begin{aligned}P(A \cap B^c) &= P(A) - P(A)P(B) = P(A)(1 - P(B)) \\&= P(A)P(B^c)\end{aligned}$$

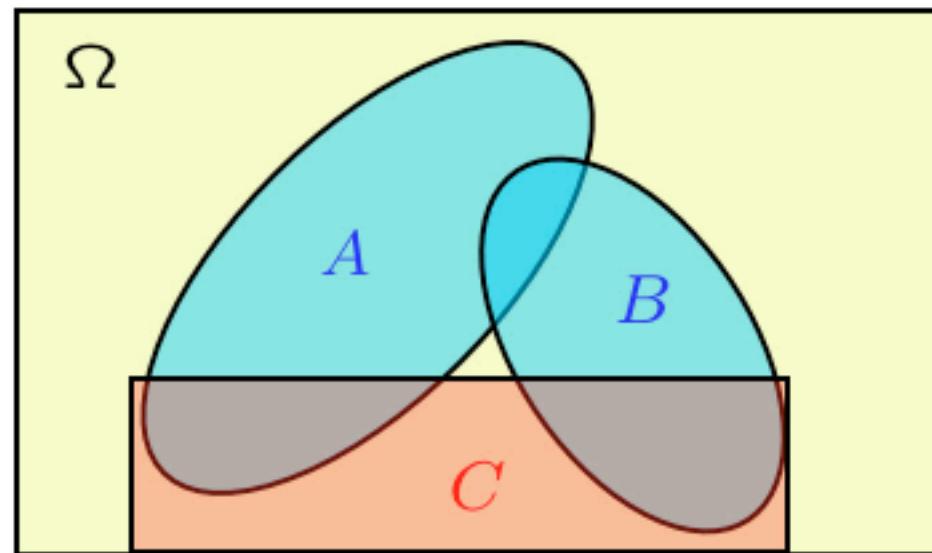
Conditional independence

- Conditional independence, given C , is defined as independence under the probability law $\mathbf{P}(\cdot | C)$



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

Assume A and B are independent



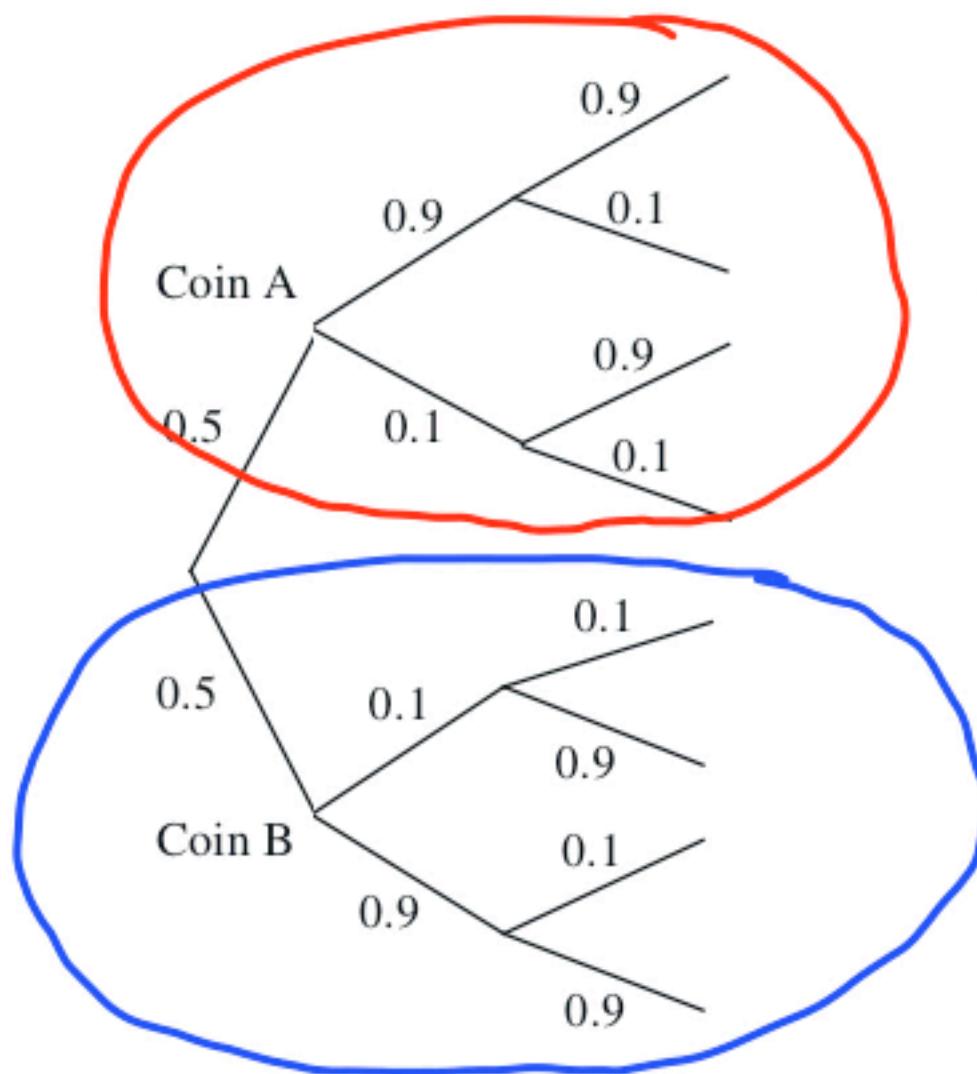
- If we are told that C occurred, are A and B independent? **No**

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



- Compare:

$$\begin{aligned} P(\text{toss } 11 = H) &= P(A)P(H_{11} | A) + P(B)P(H_{11} | B) \\ &= 0.5 \times 0.9 + 0.5 \times 0.1 = 0.5 \end{aligned}$$

$$P(\text{toss } 11 = H \mid \text{first 10 tosses are heads})$$

$$\approx P(H_{11} | A) = 0.9$$

given or coin:
independent tosses

- Are coin tosses independent?

No!

Independence of a collection of events

- **Intuitive “definition”:** Information on some of the events does not change probabilities related to the remaining events

$$A_1, A_2, \dots, \text{indep} \Rightarrow P(A_3 \cap A_4^c) = P(A_3 \cap A_4^c | A_1 \cup (A_2 \cap A_5^c)) .$$

$$P(A_3) = P(A_3 | A_1 \cap A_2) = P(A_3 | A_1 \cap A_2^c) = P(A_3 | A_1^c \cap A_2)$$

Definition: Events A_1, A_2, \dots, A_n are called **independent** if:

$$P(A_i \cap A_j \cap \dots \cap A_m) = P(A_i)P(A_j) \dots P(A_m) \quad \text{for any distinct indices } i, j, \dots, m$$

$n = 3$:

$$\left. \begin{array}{l} P(A_1 \cap A_2) = P(A_1) \cdot P(A_2) \\ P(A_1 \cap A_3) = P(A_1) \cdot P(A_3) \\ P(A_2 \cap A_3) = P(A_2) \cdot P(A_3) \end{array} \right\} \text{pairwise independence}$$

$$P(A_1 \cap A_2 \cap A_3) = P(A_1) \cdot P(A_2) \cdot P(A_3)$$

Independence vs. pairwise independence

- Two independent fair coin tosses

- H_1 : First toss is H
- H_2 : Second toss is H

$$P(H_1) = P(H_2) = 1/2$$

- C : the two tosses had the same result $= \{HH, TT\}$

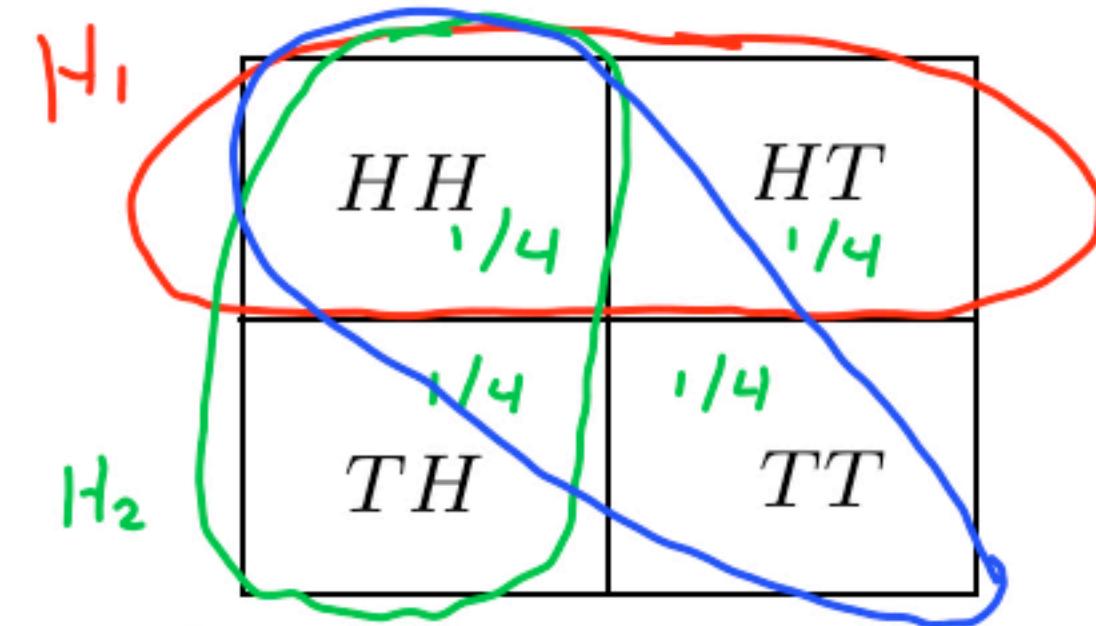
$$P(H_1 \cap C) = P(H_1 \cap H_2) = 1/4 \quad P(H_1) P(C) = \frac{1}{2} \cdot \frac{1}{2} = \frac{1}{4} \quad H_1, C: \text{indep.}$$

$$P(H_1 \cap H_2 \cap C) = P(HH) = 1/4 \quad P(H_1) P(H_2) P(C) = 1/8 \quad \text{diff.}$$

$$P(C | H_1) = P(H_2 | H_1) = P(H_2) = 1/2 = P(C)$$

$$P(C | H_1 \cap H_2) = 1 \neq P(C) = 1/2$$

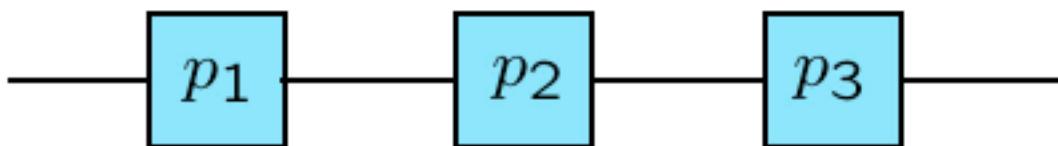
H_1 , H_2 , and C are pairwise independent, but not independent



$H_1, C: \text{indep.}$
 $H_2, C: \text{indep.}$

Reliability

p_i : probability that unit i is “up”
independent units

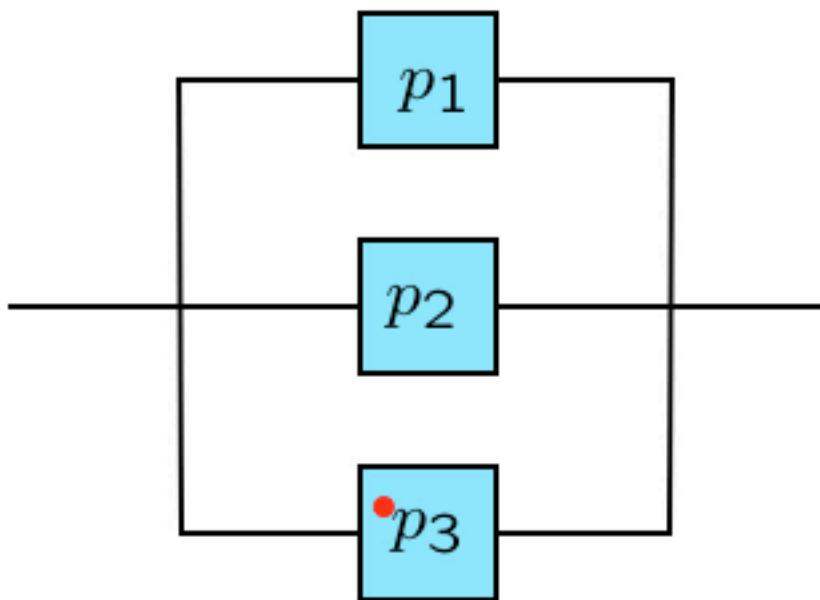


U_i : i th unit up
 U_1, U_2, \dots, U_m independent

F_i : i th unit down
 $\Rightarrow F_i$ independent

probability that system is “up”?

$$\begin{aligned} P(\text{system up}) &= P(U_1 \cap U_2 \cap U_3) \\ &= P(U_1) P(U_2) P(U_3) = p_1 p_2 p_3 \end{aligned}$$



$$\begin{aligned} P(\text{system is up}) &= P(U_1 \cup U_2 \cup U_3) \\ &= 1 - P(F_1 \cap F_2 \cap F_3) \\ &= 1 - P(F_1) P(F_2) P(F_3) \\ &\approx 1 - (1-p_1)(1-p_2)(1-p_3) \end{aligned}$$

The king's sibling

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

~~1/2~~ ?

boy have precedence

$$P(\text{boy}) = P(\text{girl}) = 1/2$$

independent

BB $\frac{1}{4}$	BG $\frac{1}{4}$
GB $\frac{1}{4}$	GG $\frac{1}{4}$

2/3

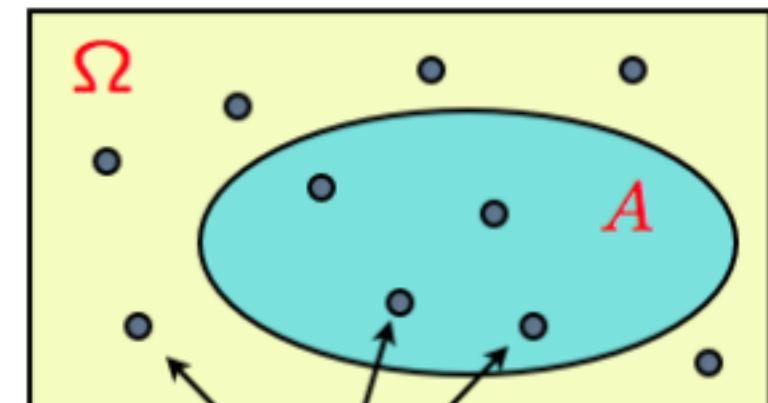
- till 1 boy $\Rightarrow P(G) = 1$
- till 2 boys $\Rightarrow P(G) = 0$

LECTURE 4: Counting

Discrete uniform law

- Assume Ω consists of n equally likely elements
- Assume A consists of k elements

Then : $P(A) = \frac{\text{number of elements of } A}{\text{number of elements of } \Omega} = \frac{k}{n}$



- Basic counting principle
- Applications

permutations
combinations
partitions

number of subsets
binomial probabilities

Basic counting principle

4 shirts

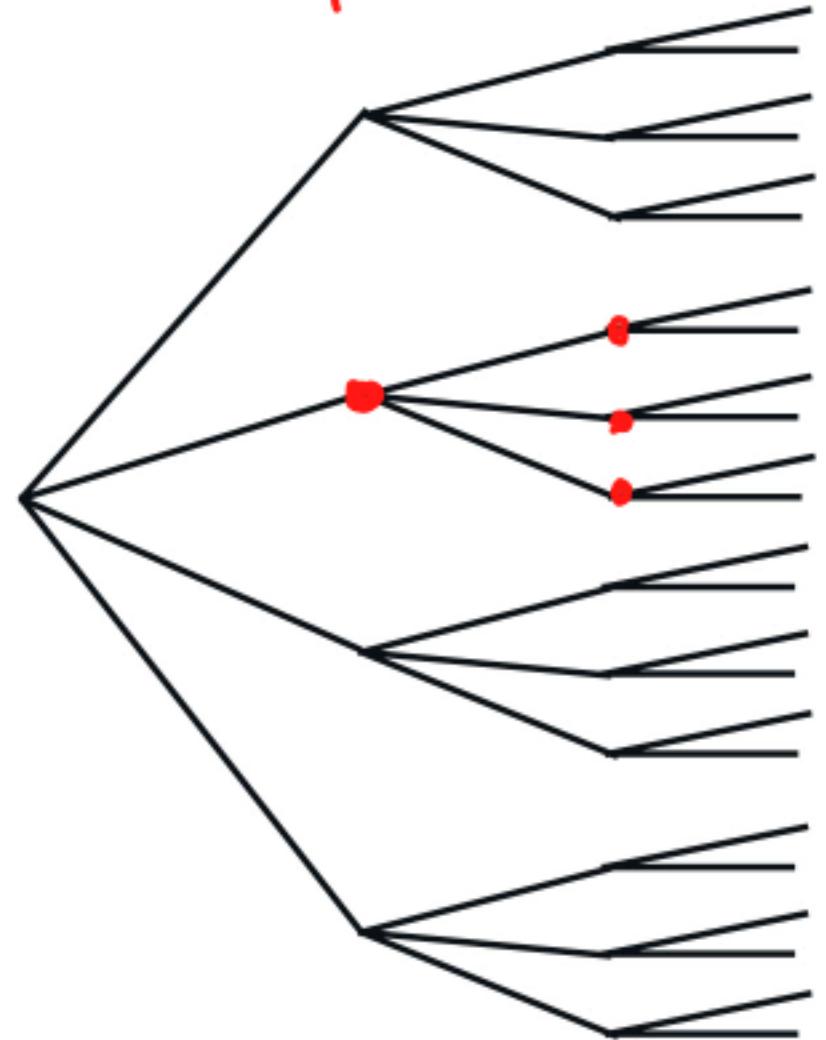
3 ties

2 jackets

Number of possible attires?

- r stages
- n_i choices at stage i

$$4 \quad 12 \quad 24 = 4 \cdot 3 \cdot 2$$



$$r = 3$$

$$n_1 = 4$$

$$n_2 = 3$$

$$n_3 = 2$$

Number of choices is: $n_1 \cdot n_2 \cdots n_r$

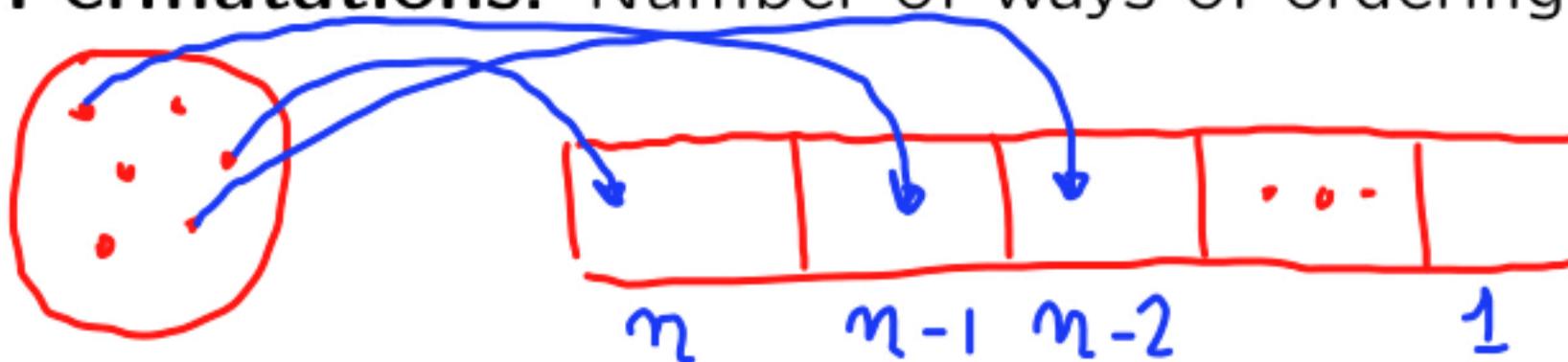
Basic counting principle examples

- Number of license plates with 2 letters followed by 3 digits:

$$26 \cdot 26 \cdot 10 \cdot 10 \cdot 10$$

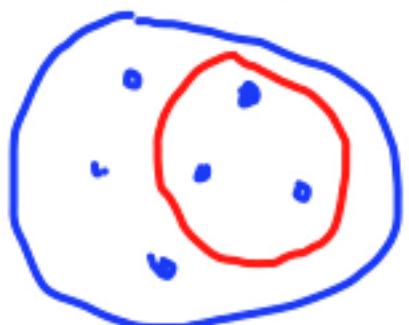
- ... if repetition is prohibited: $26 \cdot 25 \cdot 10 \cdot 9 \cdot 8$

- Permutations:** Number of ways of ordering n elements:



$$n \cdot (n-1) \cdot (n-2) \cdots 1 = n!$$

- Number of subsets of $\{1, \dots, n\}$:



$$2 \cdot 2 \cdots 2 = 2^n$$

$$\begin{array}{ccc} n=1 & \{\} & 2^1 = 2 \\ \{\} & \emptyset & \end{array}$$

Example

- Find the probability that:
six rolls of a (six-sided) die all give different numbers.

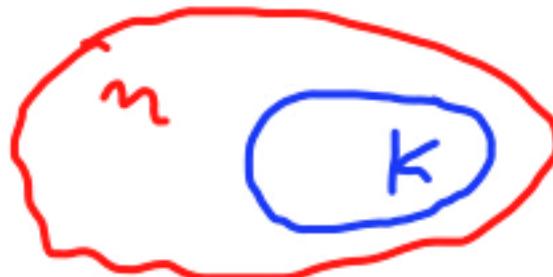
(Assume all outcomes equally likely.)

typical outcome $P(2,3,4,3,6,2) = 1/6^6$

" element of A : $(2,3,4,1,6,5) = 6!$

$$P(A) = \frac{\# \text{ in } A}{\# \text{ possible outcomes}} = \frac{6!}{6^6}$$

Combinations

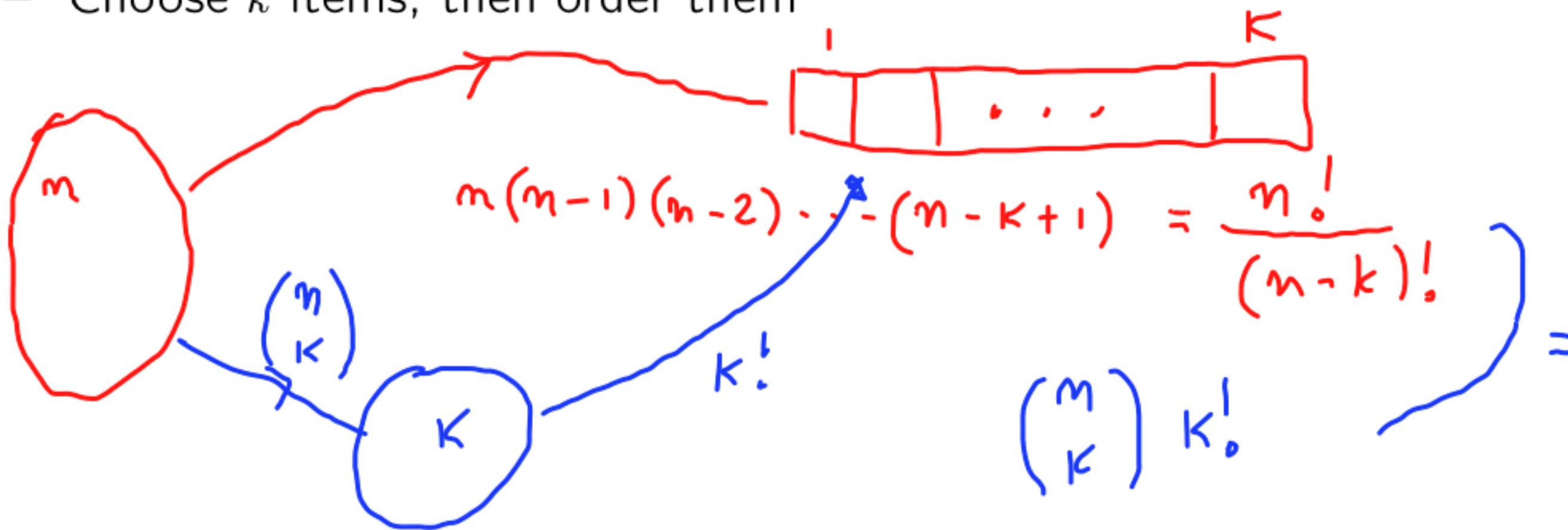


- Definition: $\binom{n}{k}$: number of k -element subsets of a given n -element set

$$= \frac{n!}{k!(n-k)!}$$

$n = 0, 1, 2, \dots$

- Two ways of constructing an **ordered** sequence of k **distinct** items: $k = 0, 1, \dots, n$
 - Choose the k items one at a time
 - Choose k items, then order them



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

$$\binom{n}{n} = 1 \quad \frac{n!}{n! 0!}$$

$$\binom{n}{0} = \frac{n!}{0! n!} = 1 \quad \phi$$

$$0! = 1 \quad \text{convention}$$

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

Binomial coefficient $\binom{n}{k}$ → Binomial probabilities

- $n \geq 1$ independent coin tosses; $P(H) = p$

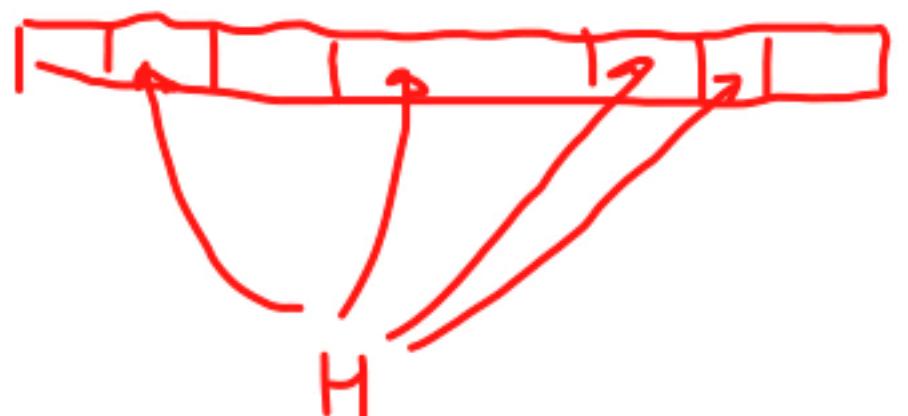
$n=6$

$$\bullet P(HTTHHH) = P((1-p)(1-p)p p p p) = p^4(1-p)^2$$

$$\bullet P(\text{particular sequence}) = P^{\#\text{heads}}(1-p)^{\#\text{tails}}$$

$$\bullet P(\text{particular } k\text{-head sequence}) = p^k(1-p)^{n-k}$$

$$P(k \text{ heads}) = p^k(1-p)^{n-k} \cdot (\# k\text{-head sequences})$$



$$\binom{n}{k}$$

$$P(k \text{ heads}) = \binom{n}{k} p^k (1-p)^{n-k}$$

A coin tossing problem

- Given that there were 3 heads in 10 tosses, what is the probability that the first two tosses were heads?
 - event A : the first 2 tosses were heads
 - event B : 3 out of 10 tosses were heads

Assumptions:

- independence
- $P(H) = p$

$$P(k \text{ heads}) = \binom{n}{k} p^k (1-p)^{n-k}$$

- First solution:

$$P(A | B) = \frac{P(A \cap B)}{P(B)} = \frac{\cancel{P(H_1 H_2 \text{ and one } H \text{ in tosses } 3, \dots, 10)}}{\cancel{P(B)}}$$

$$= \frac{p^2 \cdot \binom{8}{1} p^1 \cdot (1-p)^7}{\binom{10}{3} p^3 (1-p)^7} = \frac{\binom{8}{1}}{\binom{10}{3}} = \frac{8}{\binom{10}{3}}.$$

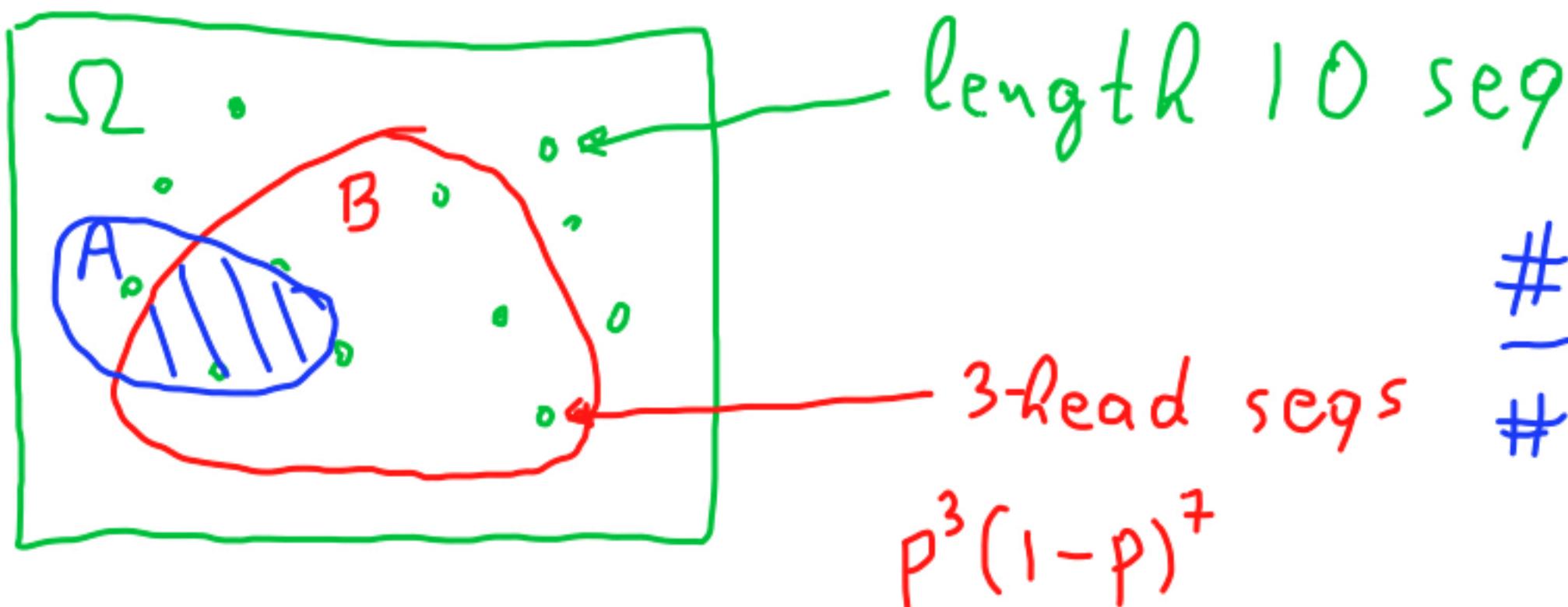
A coin tossing problem

- Given that there were 3 heads in 10 tosses, what is the probability that the first two tosses were heads?
 - event A : the first 2 tosses were heads
 - event B : 3 out of 10 tosses were heads
- Second solution: Conditional probability law (on B) is uniform

Assumptions:

- independence
- $P(H) = p$

$$P(k \text{ heads}) = \binom{n}{k} p^k (1-p)^{n-k}$$

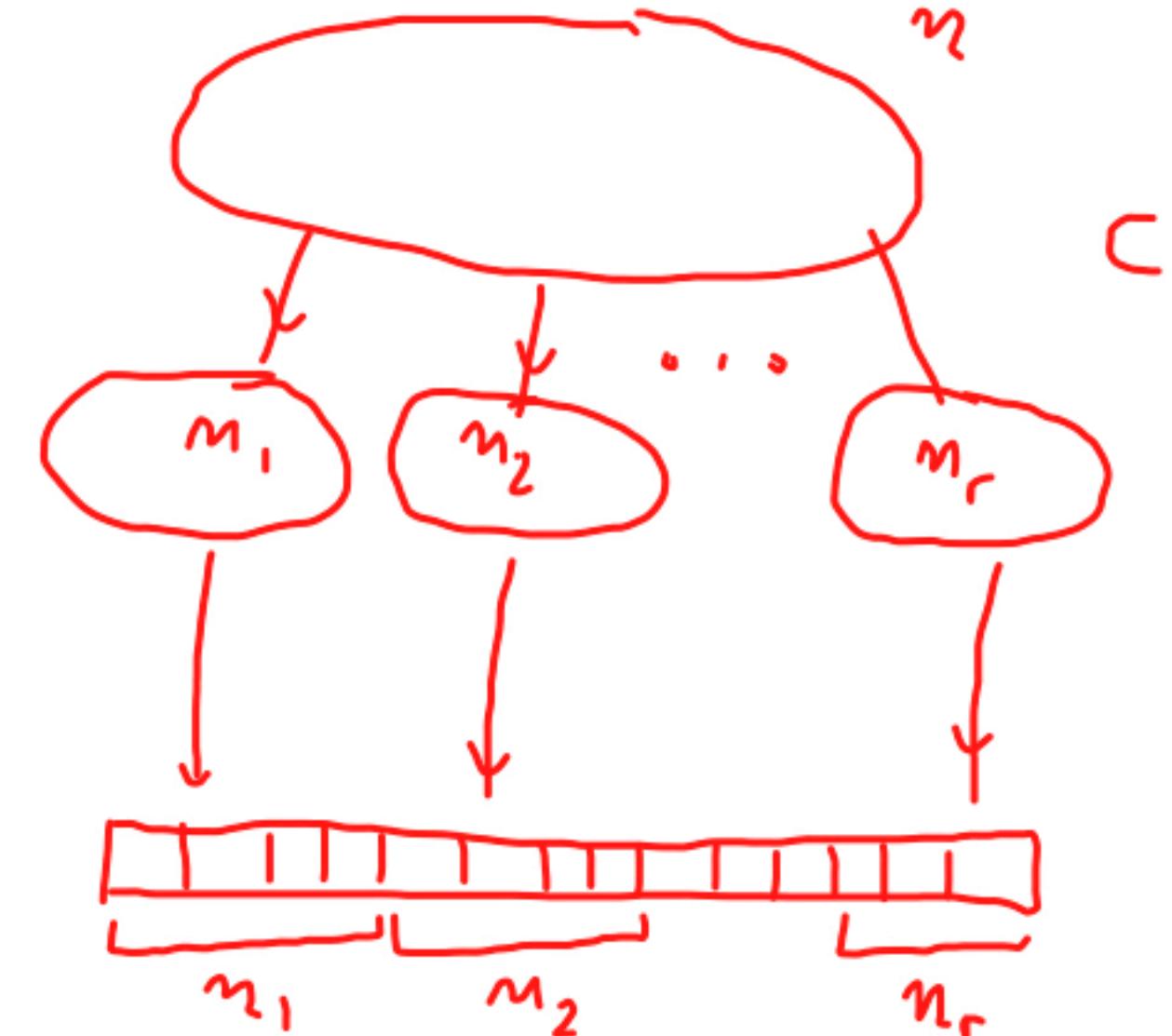


$$\frac{\#\text{in}(A \cap B)}{\#\text{in } B} = \frac{8}{\binom{10}{3}}$$

Partitions

- $n \geq 1$ distinct items; $r \geq 1$ persons give n_i items to person i
 - here n_1, \dots, n_r are given nonnegative integers
 - with $n_1 + \dots + n_r = n$
- Ordering n items: $n!$
 - Deal n_i to each person i , and then order

$$n_1! n_2! \dots n_r! = n!$$



$$r=2 \quad m_1=k \quad m_2=n-k$$

number of partitions = $\frac{n!}{n_1! n_2! \dots n_r!}$ • (multinomial coefficient)

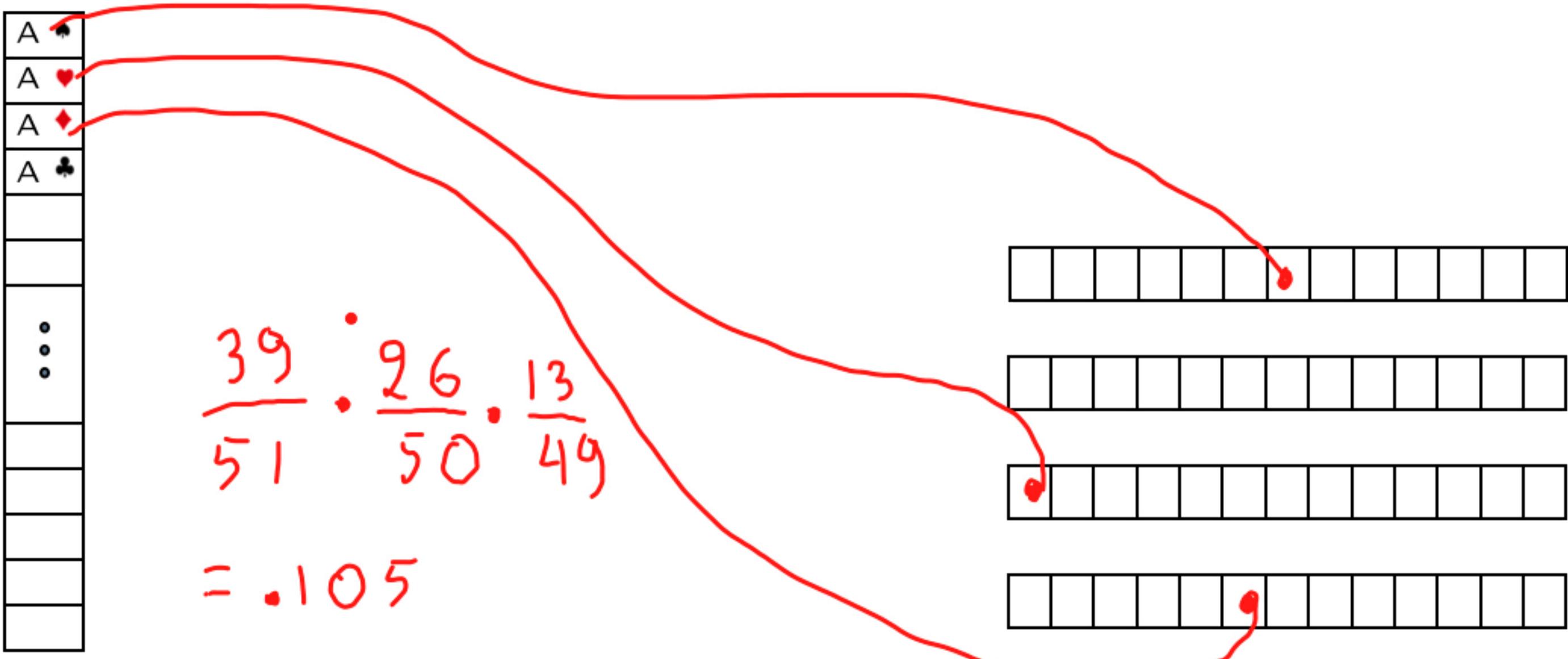
Example: 52-card deck, dealt (fairly) to four players.
Find $P(\text{each player gets an ace})$

- Outcomes are: **partition equally likely**
 - number of outcomes: $\frac{52!}{13! 13! 13! 13!} \cdot$
- Constructing an outcome with one ace for each person:
 - distribute the aces $4 \cdot 3 \cdot 2 \cdot 1$
 - distribute the remaining 48 cards $\frac{48!}{12! 12! 12! 12!}$
- Answer:
$$\frac{4 \cdot 3 \cdot 2 \cdot \frac{48!}{12! 12! 12! 12!}}{\frac{52!}{13! 13! 13! 13!}}$$

Example: 52-card deck, dealt (fairly) to four players.
Find $P(\text{each player gets an ace})$

A smart solution

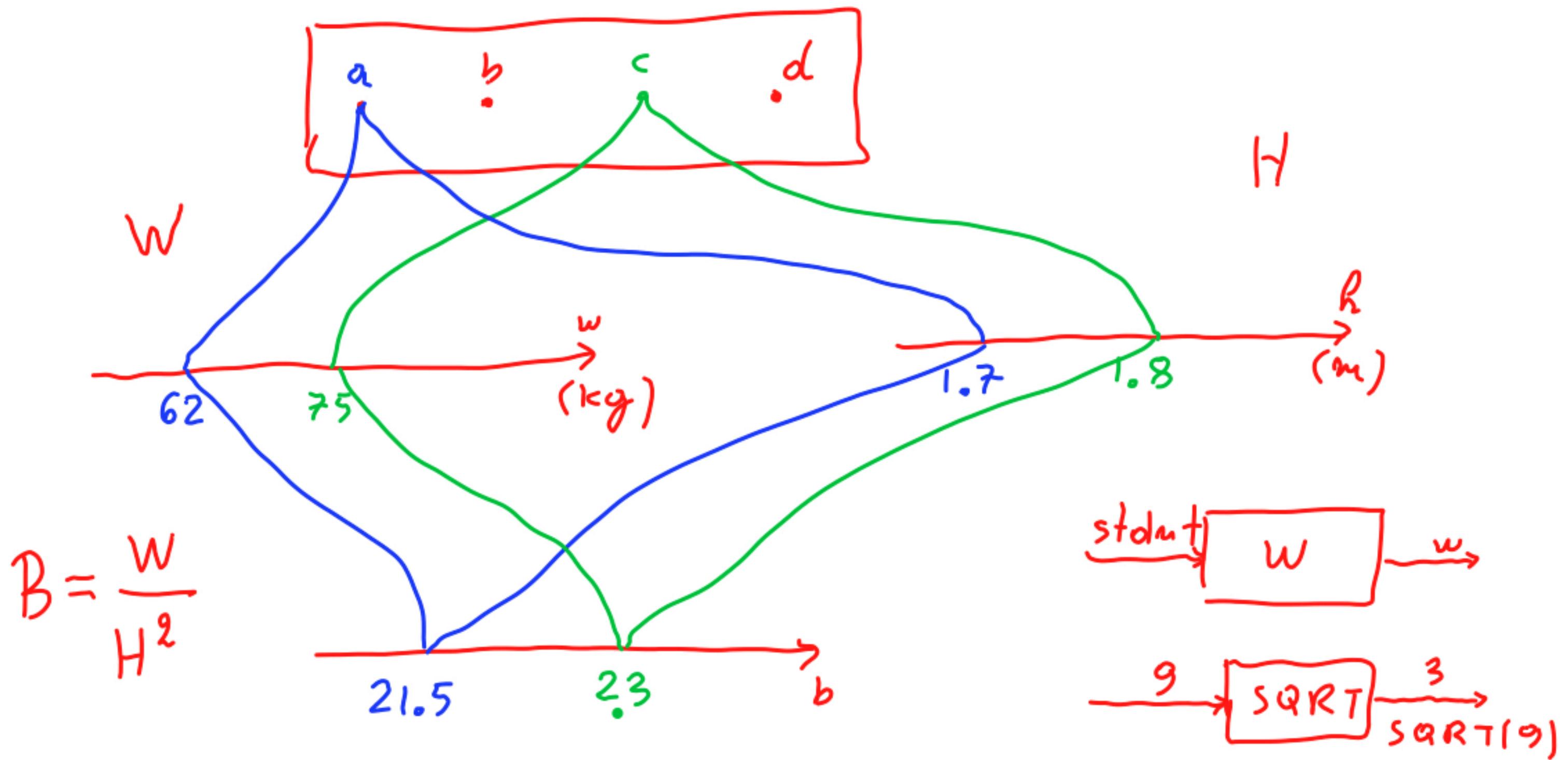
Stack the deck, aces on top



LECTURE 5: Discrete random variables: probability mass functions and expectations

- Random variables: the idea and the definition
 - **Discrete:** take values in finite or countable set
- Probability mass function (PMF)
- Random variable examples
 - Bernoulli
 - Uniform
 - Binomial
 - Geometric
- Expectation (mean) and its properties
 - The expected value rule
 - Linearity

Random variables: the idea



Random variables: the formalism

- A random variable (“r.v.”) associates a value (a number) to every possible outcome
- Mathematically: A function from the sample space Ω to the real numbers
- It can take discrete or continuous values

Notation: random variable X numerical value x

- We can have several random variables defined on the same sample space
- A function of one or several random variables is also a random variable
 - meaning of $X + Y$: r.v takes value $x+y$, when X takes value x , Y takes value y

Probability mass function (PMF) of a discrete r.v. X

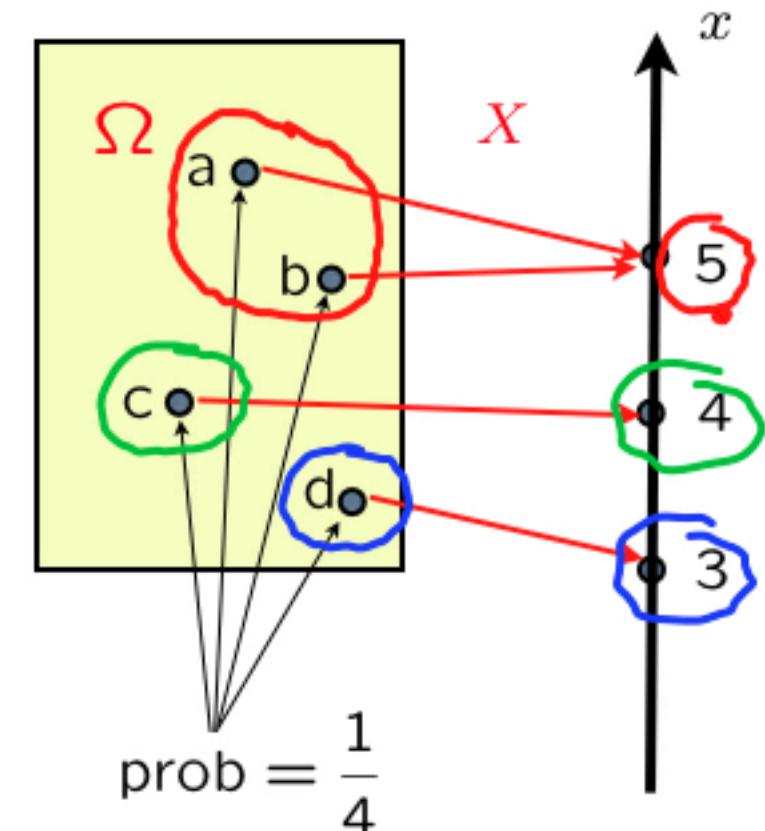
- It is the “probability law” or “probability distribution” of X
- If we fix some x , then “ $X = x$ ” is an event

$$x=5 \quad X=5 \quad \{\omega : X(\omega)=5\} = \{a, b\}$$

$$P_X(5) = 1/2$$

$$p_X(x) = P(X=x) = P(\{\omega \in \Omega \text{ s.t. } X(\omega)=x\})$$

• Properties: $p_X(x) \geq 0$ $\sum_x p_X(x) = 1$



$$P_Y(y)$$



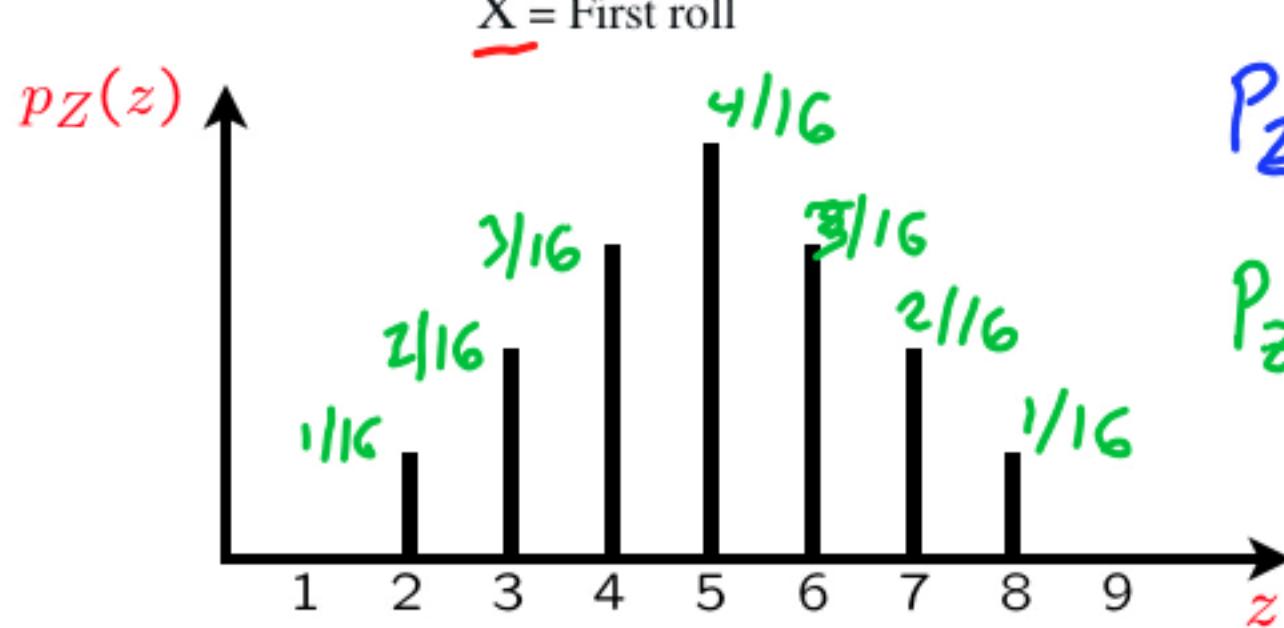
PMF calculation

- Two rolls of a tetrahedral die
- Let every possible outcome have probability $1/16$

	5	6	7	8
4	4	5	6	7
3	3	4	5	6
2	2	3	4	5
1	1	2	3	4

$\text{Y} = \text{Second roll}$

$\text{X} = \text{First roll}$



$$Z = X + Y$$

Find $p_Z(z)$ for all z

- repeat for all z :
 - collect all possible outcomes for which Z is equal to z
 - add their probabilities

$$P_Z(2) = P(Z=2) = 1/16$$

$$P_Z(3) = P(Z=3) = 2/16$$

$$P_Z(4) = P(Z=4) = 3/16$$

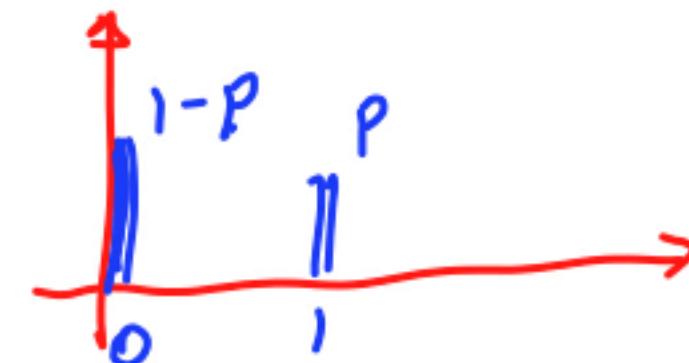
•
•
•

The simplest random variable: Bernoulli with parameter $p \in [0, 1]$

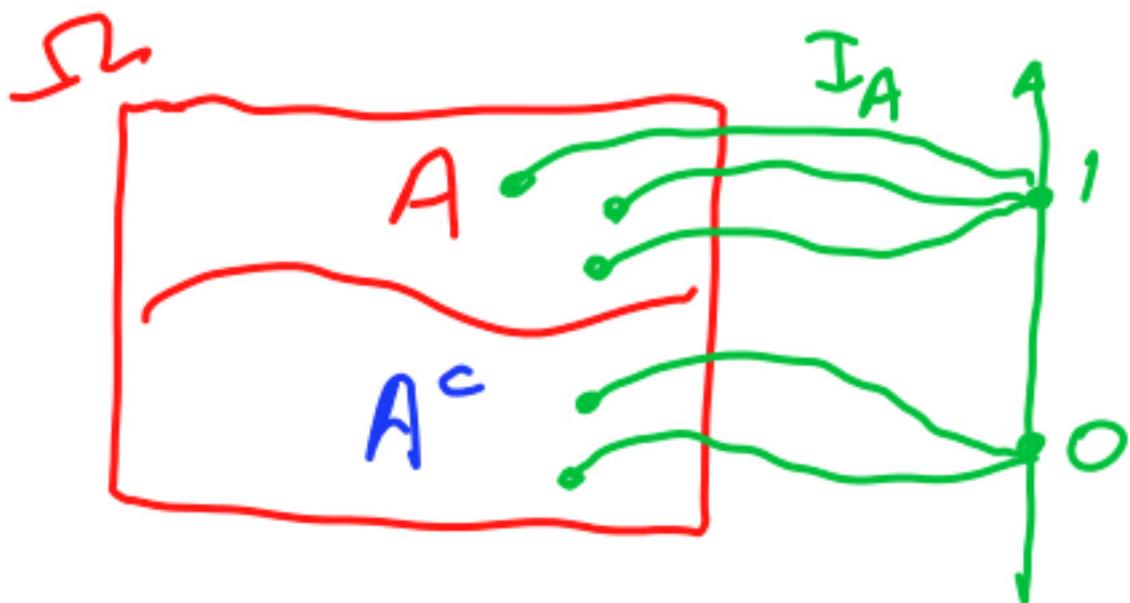
$$X = \begin{cases} 1, & \text{w.p. } p \\ 0, & \text{w.p. } 1 - p \end{cases}$$

$$P_X(0) = 1 - p$$

$$P_X(1) = p$$



- Models a trial that results in success/failure, Heads/Tails, etc.
- Indicator r.v. of an event A : $I_A = 1$ iff A occurs



$$P_{I_A}(1) = P(I_A = 1) = P(A)$$



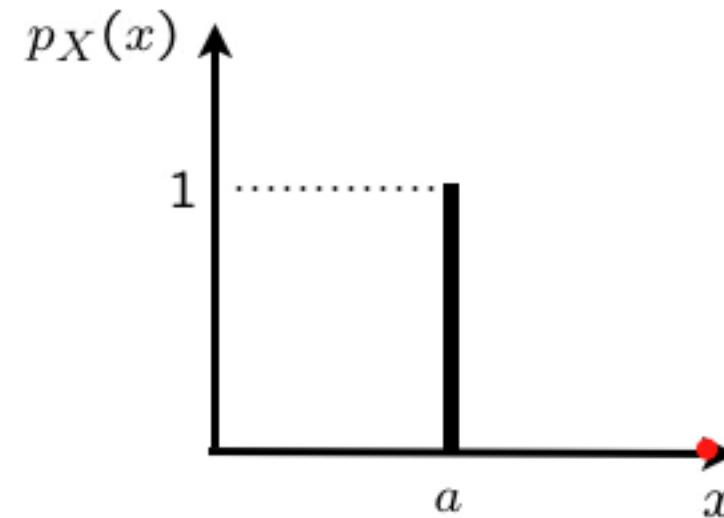
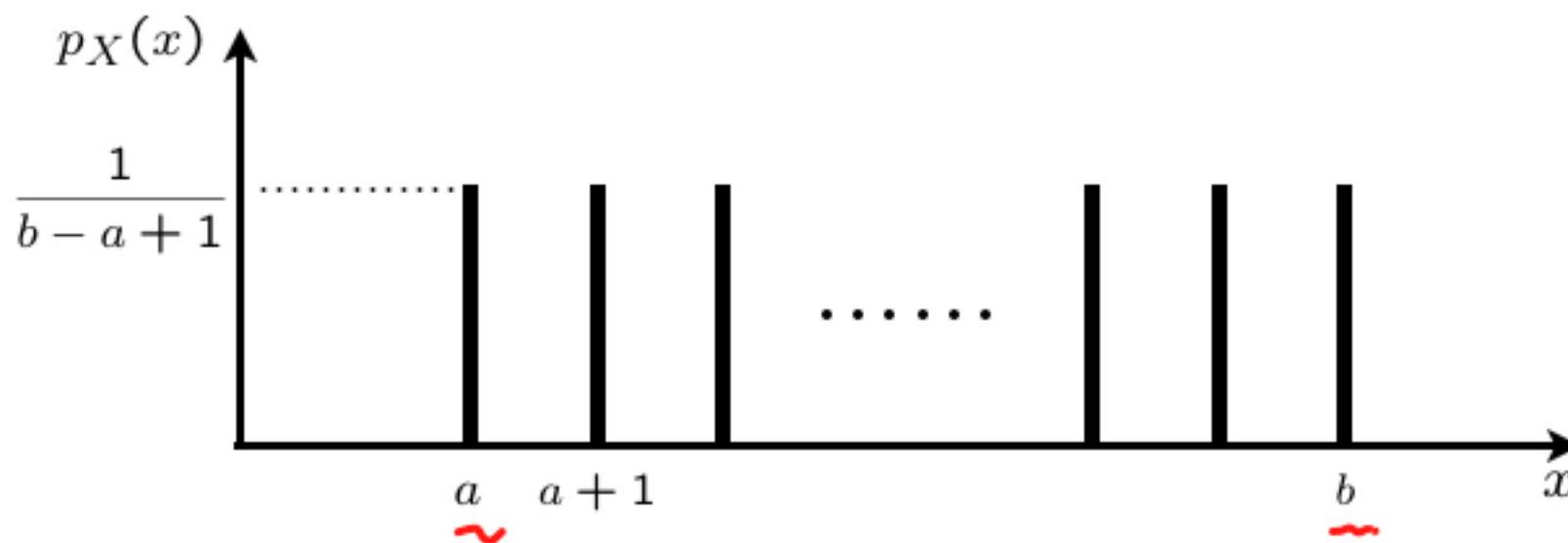
Discrete uniform random variable; parameters a, b

- **Parameters:** integers a, b ; $a \leq b$
- **Experiment:** Pick one of $a, a+1, \dots, b$ at random; all equally likely
- **Sample space:** $\{a, a+1, \dots, b\}$
- **Random variable X :** $X(\omega) = \omega$
- **Model of:** complete ignorance

$b - a + 1$ possible values

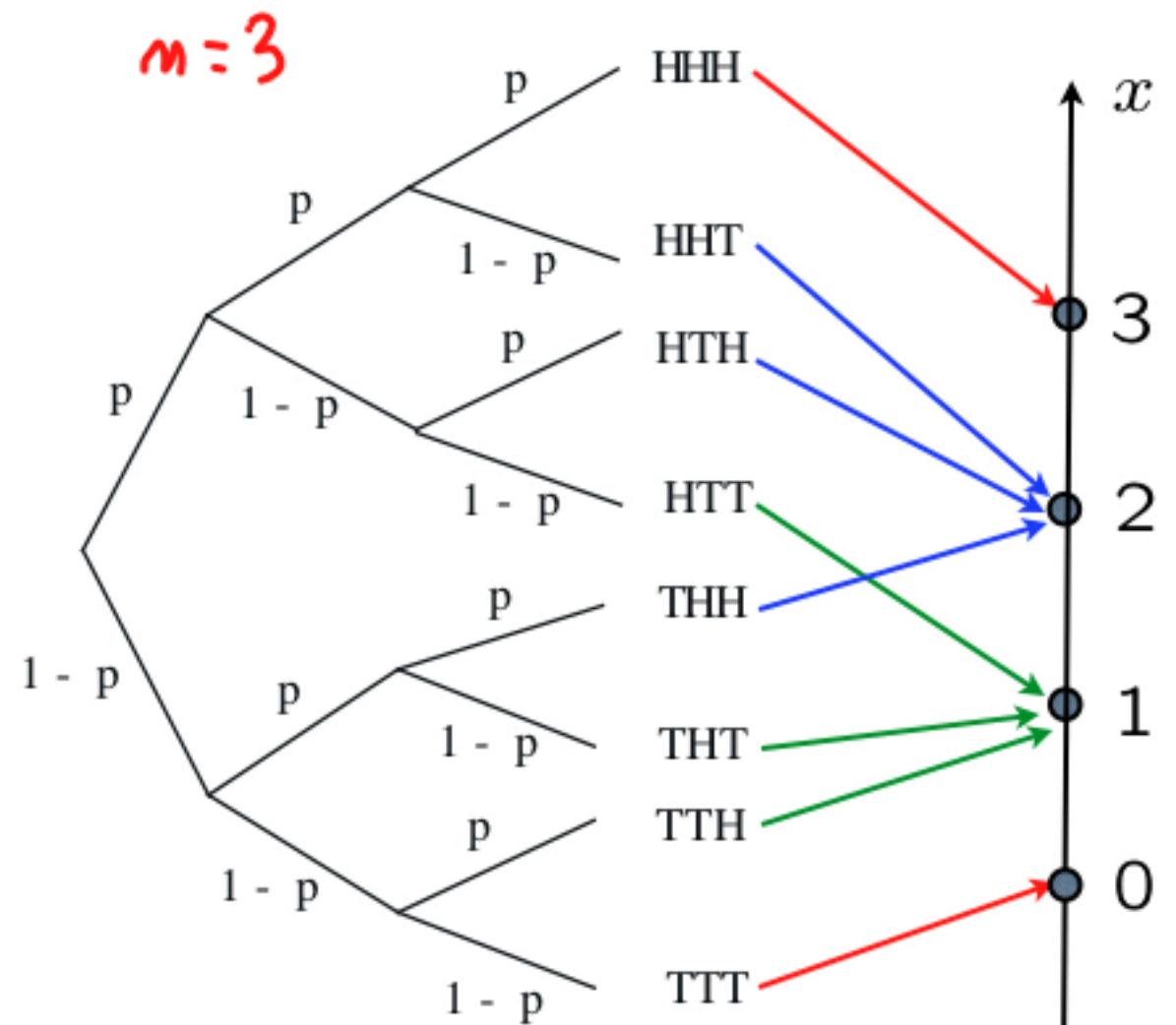
11:52:26 $\{0, 1, \dots, 59\}$

Special case: $a = b$
constant/deterministic r.v.



Binomial random variable; **parameters:** positive integer n ; $p \in [0, 1]$

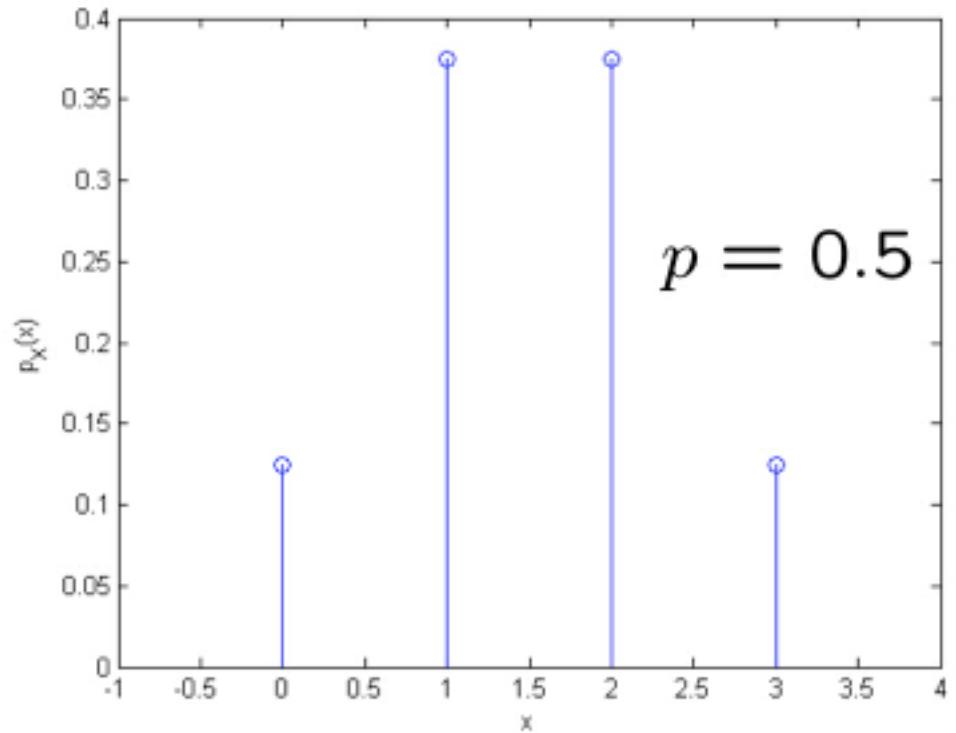
- **Experiment:** n independent tosses of a coin with $P(\text{Heads}) = p$
- **Sample space:** Set of sequences of H and T, of length n
- **Random variable X :** number of Heads observed
- **Model of:** number of successes in a given number of independent trials



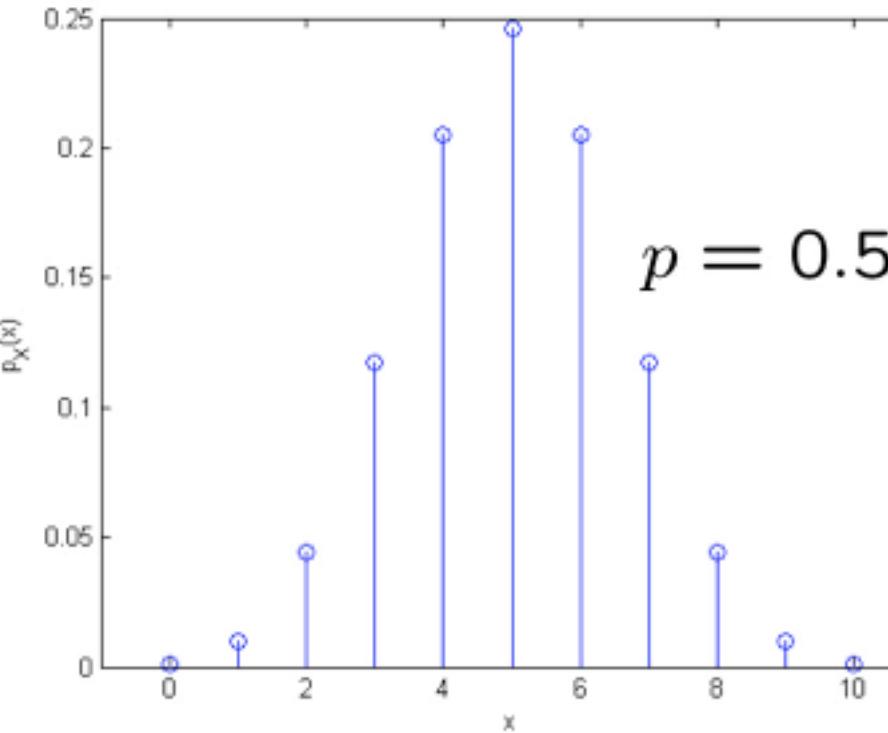
$$\begin{aligned} P_X(2) &= P(X=2) \\ &= P(HHT) + P(HTH) + P(THH) \\ &= 3p^2(1-p) \approx \binom{3}{2} p^2 (1-p) \end{aligned}$$

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

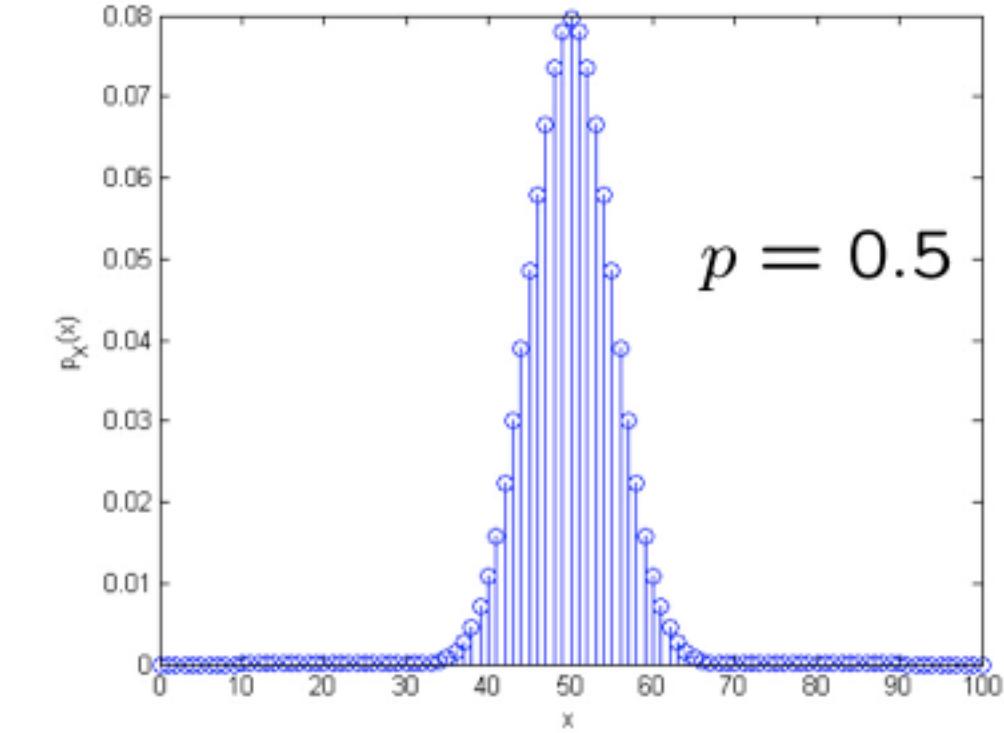
$n = 3$



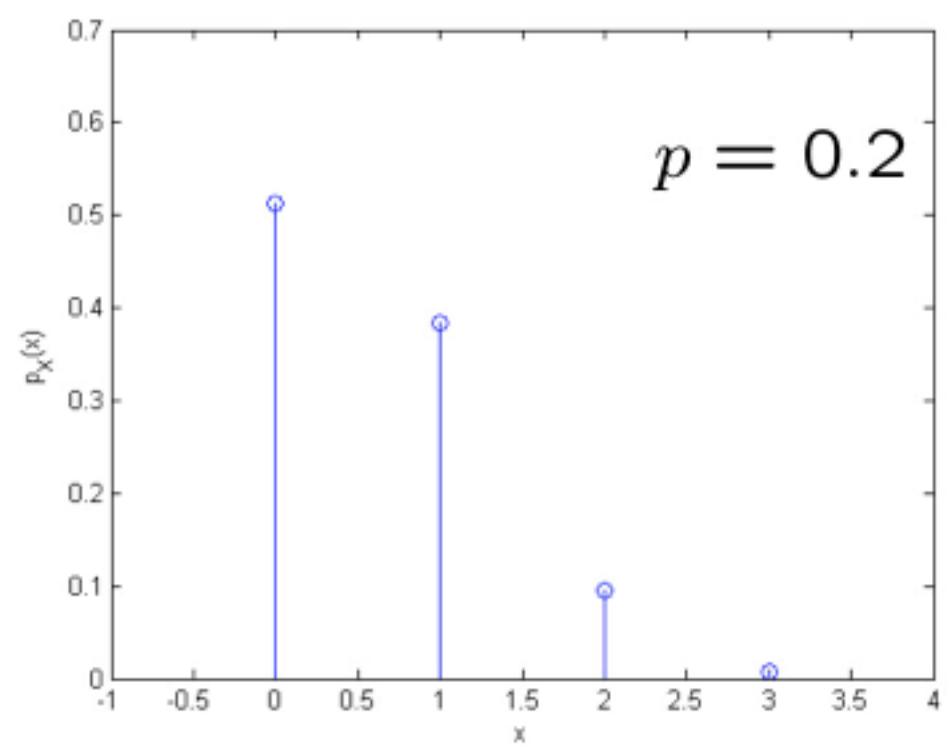
$n = 10$



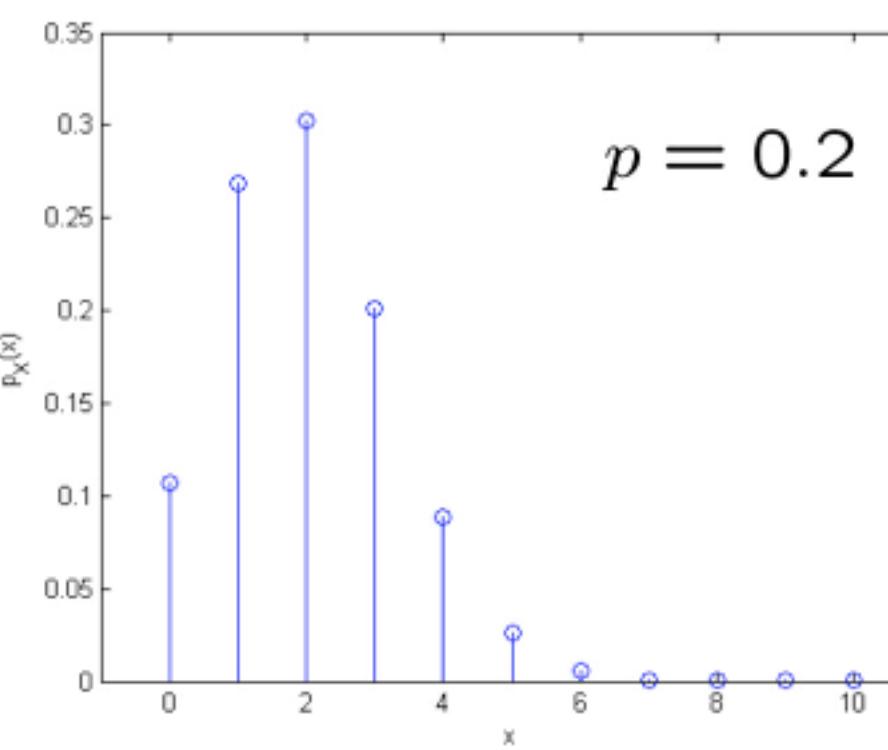
$n = 100$



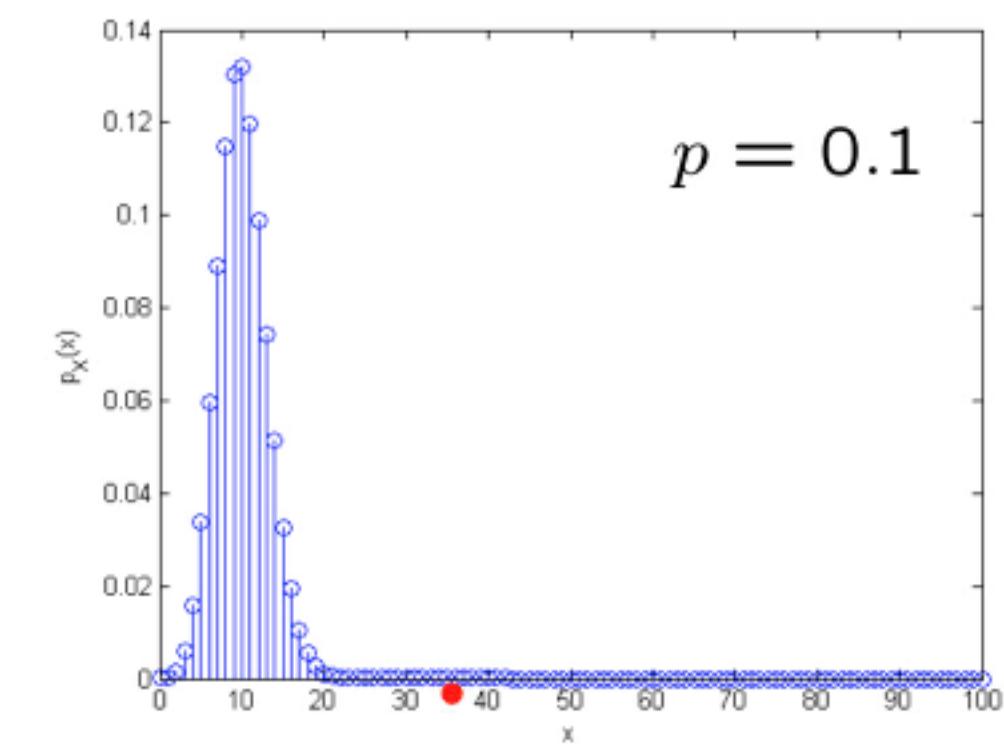
$p = 0.2$



$p = 0.2$



$p = 0.1$



Geometric random variable; parameter p : $0 < p \leq 1$

- **Experiment:** infinitely many independent tosses of a coin; $P(\text{Heads}) = p$
- **Sample space:** Set of infinite sequences of H and T TTTTHHT...
- **Random variable X :** number of tosses until the first Heads $X = 5$

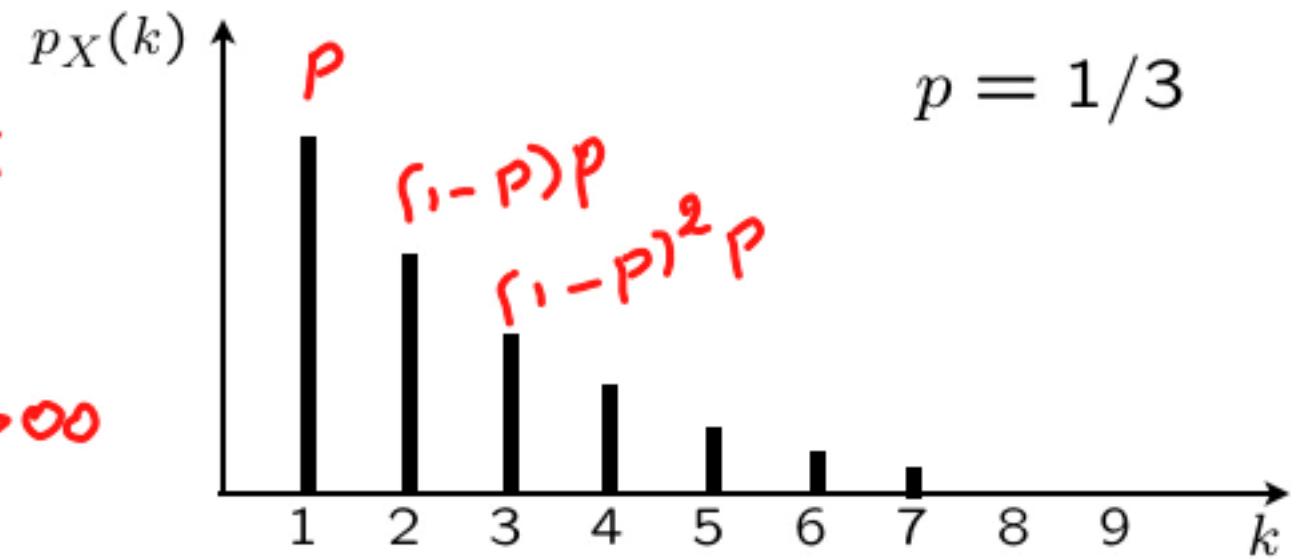
- **Model of:** waiting times; number of trials until a success

$$p_X(k) = P(X=k) = P\left(T\underbrace{\dots}_{k-1}H\right) = (1-p)^{k-1}p \quad k=1, 2, 3, \dots$$

$P(\text{no Heads ever}) \leq P\left(T\underbrace{\dots}_{k}T\right) = (1-p)^k$

$\underbrace{(TTT\dots)}_{\text{"X=oo"}} = 0$

$\downarrow k \rightarrow \infty$



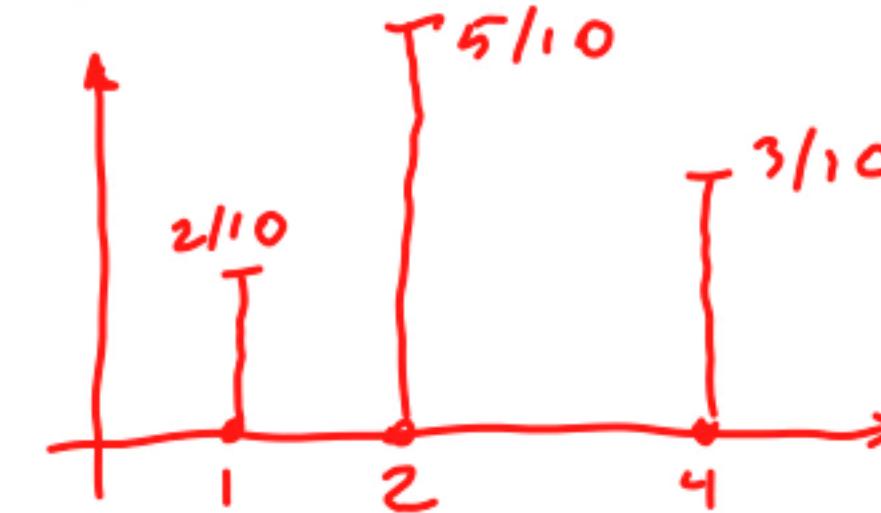
Expectation/mean of a random variable

- **Motivation:** Play a game 1000 times.
Random gain at each play described by:
- “Average” gain:

$$\begin{aligned} & \frac{1 \cdot 200 + 2 \cdot 500 + 4 \cdot 300}{1000} \\ &= 1 \cdot \frac{2}{10} + 2 \cdot \frac{5}{10} + 4 \cdot \frac{3}{10} \end{aligned}$$

$$X = \begin{cases} 1, & \text{w.p. } 2/10 \\ 2, & \text{w.p. } 5/10 \\ 4, & \text{w.p. } 3/10 \end{cases}$$

~ 200
 ~ 500
 ~ 300



- **Definition:** $E[X] = \sum_x x p_X(x)$

- **Caution:** If we have an infinite sum, it needs to be well-defined.
We assume $\sum_x |x| p_X(x) < \infty$
- **Interpretation:** Average in large number of independent repetitions of the experiment

Expectation of a Bernoulli r.v.

$$X = \begin{cases} 1, & \text{w.p. } p \\ 0, & \text{w.p. } 1 - p \end{cases}$$

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

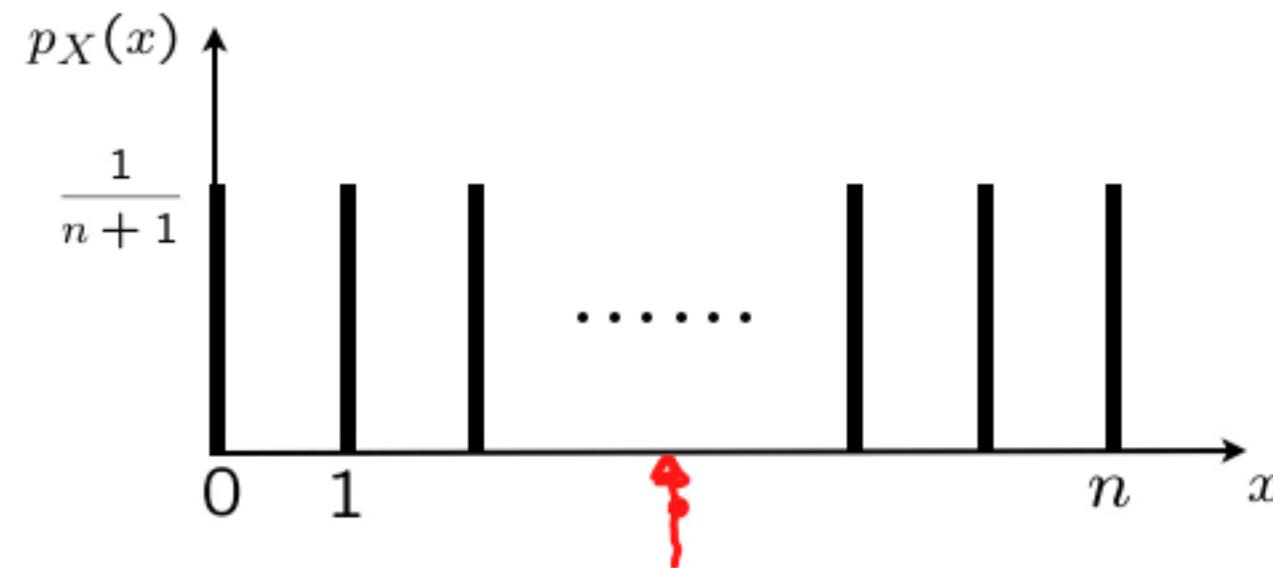
If X is the indicator of an event A , $X = I_A$:

$$X = 1 \quad \text{iff } A \text{ occurs} \quad p = P(A)$$

$$E[I_A] = P(A)$$

Expectation of a uniform r.v.

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



• **Definition:** $E[X] = \sum_x x p_X(x)$

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

$$= \frac{1}{n+1} (0 + 1 + \dots + n) = \frac{1}{n+1} \cdot \frac{n(n+1)}{2} = \frac{n}{2}$$

Expectation as a population average

- n students
- Weight of i th student: x_i
- Experiment: pick a student at random, all equally likely
- Random variable X : weight of selected student
 - assume the x_i are distinct

$$p_X(x_i) = \frac{1}{n}$$

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

Elementary properties of expectations

- If $X \geq 0$, then $E[X] \geq 0$

for all $\omega: X(\omega) \geq 0$

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

≥ 0 ≥ 0 ≥ 0

- If $a \leq X \leq b$, then $a \leq E[X] \leq b$

for all $\omega: a \leq X(\omega) \leq \bar{b}$

$$\begin{aligned} E[X] &= \sum_x x p_X(x) \geq \sum_x a p_X(x) \\ &= a \sum_x p_X(x) = a \cdot 1 = a \end{aligned}$$

- If c is a constant, $E[c] = c$

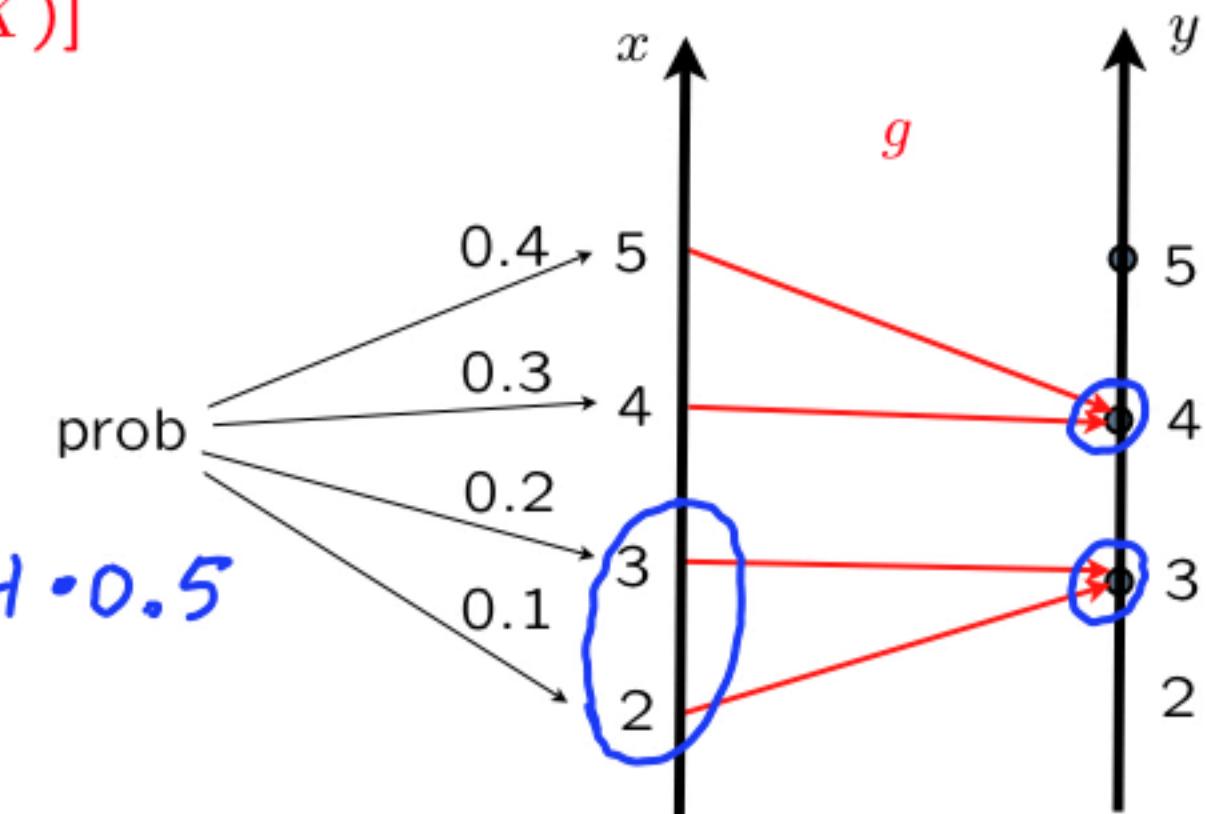


$$E[c] = c \cdot p(c) = c$$

The expected value rule, for calculating $E[g(X)]$

- Let X be a r.v. and let $Y = g(X)$
- Averaging over y : $E[Y] = \sum_y y p_Y(y)$
 $3 \cdot (0.1 + 0.2) + 4 \cdot (0.3 + 0.4)$
- Averaging over x : $3 \cdot 0.1 + 3 \cdot 0.2 + 4 \cdot 0.3 + 4 \cdot 0.5$

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



Proof:

$$\sum_y \sum_{x: g(x)=y} g(x) p_x(x)$$

$$= \sum_y \sum_{x: g(x)=y} y p_x(x) = \sum_y \sum_{x: g(x)=y} p_x(x)$$

$$= \sum_y y p_y(y) = E[Y]$$

- $E[X^2] = \sum_x x^2 p_x(x)$
 $g(x) = x^2$
- Caution: In general, $E[g(X)] \neq g(E[X])$

$$E[X^2] \neq (E[X])^2$$

Linearity of expectation: $E[aX + b] = aE[X] + b$

$X = \text{salary}$ $E[X] = \text{average salary}$

$Y = \text{new salary} = 2X + 100$ $E[Y] = E[2X + 100] = 2E[X] + 100$

- Intuitive
- **Derivation**, based on the expected value rule:

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

$$= \sum_x (ax + b) p_x(x) = a \sum_x x p_x(x) + b \underbrace{\sum_x p_x(x)}_1$$

$$E[g(x)] = g(E[x]) = aE[X] + b$$

exceptional g

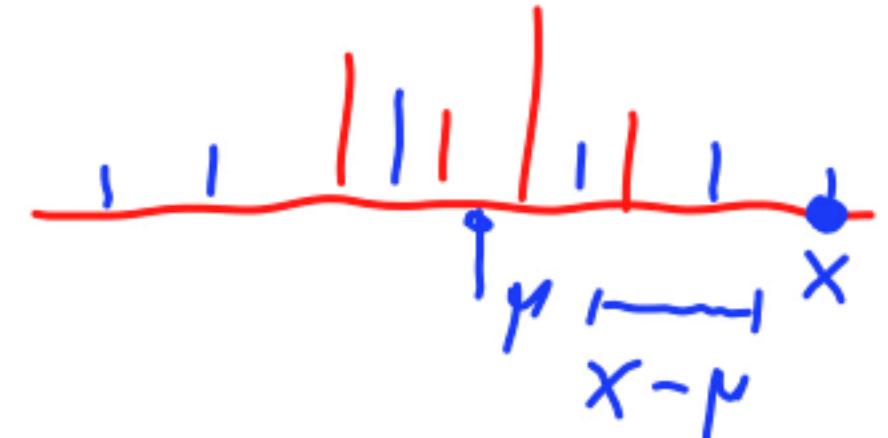
LECTURE 6: Variance; Conditioning on an event; Multiple random variables

- Variance and its properties
 - Variance of the Bernoulli and uniform PMFs
- Conditioning a r.v. on an event
 - Conditional PMF, mean, variance
 - Total expectation theorem
- Geometric PMF
 - Memorylessness
 - Mean value
- Multiple random variables
 - Joint and marginal PMFs
 - Expected value rule
 - Linearity of expectations
- The mean of the binomial PMF

Variance — a measure of the spread of a PMF

- Random variable X , with mean $\mu = E[X]$
- Distance from the mean: $X - \mu$
- Average distance from the mean?

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



- **Definition of variance:** $\text{var}(X) = E[(X - \mu)^2] \geq 0$
- Calculation, using the expected value rule, $E[g(X)] = \sum_x g(x)p_X(x)$
$$g(x) = (x - \mu)^2 \quad \text{var}(X) = E[g(X)] = \sum_x (x - \mu)^2 p_X(x)$$

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

Properties of the variance

- Notation: $\mu = E[X]$

$$\text{var}(aX + b) = a^2 \text{var}(X)$$

- Let $Y = X + b$ $\gamma = E[Y] = \mu + b$
 $\text{var}(Y) = E[(Y - \gamma)^2] = E[(X + b - (\mu + b))^2] = E[(X - \mu)^2] = \text{var}(X)$
- Let $Y = aX$ $\gamma = E[Y] = a\mu$
 $\text{var}(Y) = E[(aX - a\mu)^2] = E[a^2(X - \mu)^2] = a^2 E[(X - \mu)^2] = a^2 \text{var}(X)$

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

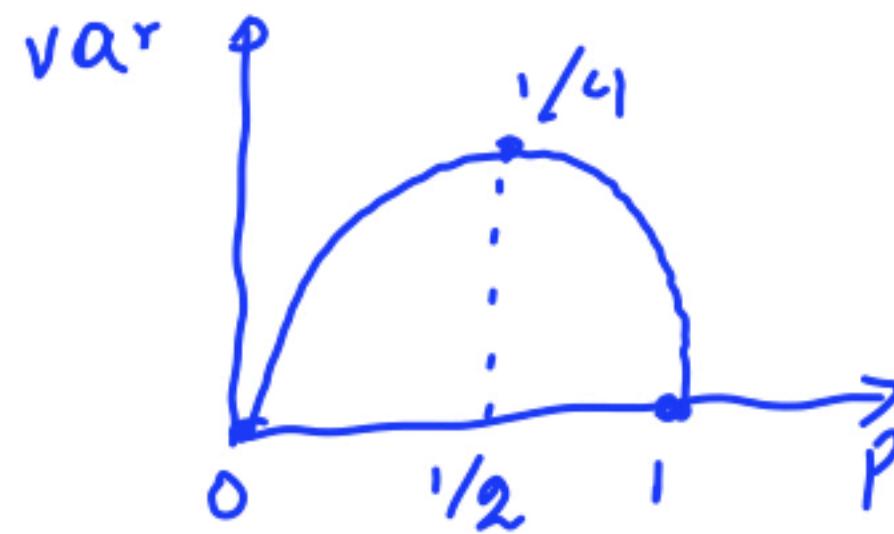
$$\begin{aligned}\text{var}(X) &= E[(X - \mu)^2] = E[X^2 - 2\mu X + \mu^2] \\ &= E[X^2] - 2\mu E[X] + \mu^2 = E[X^2] - (E[X])^2\end{aligned}$$

$$\begin{aligned}\text{var}(3 - 4X) &= (-4)^2 \text{var}(X) \\ &= 16 \text{var}(X)\end{aligned}$$

Variance of the Bernoulli

$$X = \begin{cases} 1, & \text{w.p. } p \\ 0, & \text{w.p. } 1 - p \end{cases}$$

$$E[X] = p$$



$$\begin{aligned} \text{var}(X) &= \sum_x (x - E[X])^2 p_X(x) = (1-p)^2 p + (0-p)^2 \cdot (1-p) \\ &= p - 2p^2 + p^2 + p^2 - p^3 = p - p^2 = p(1-p) \end{aligned}$$

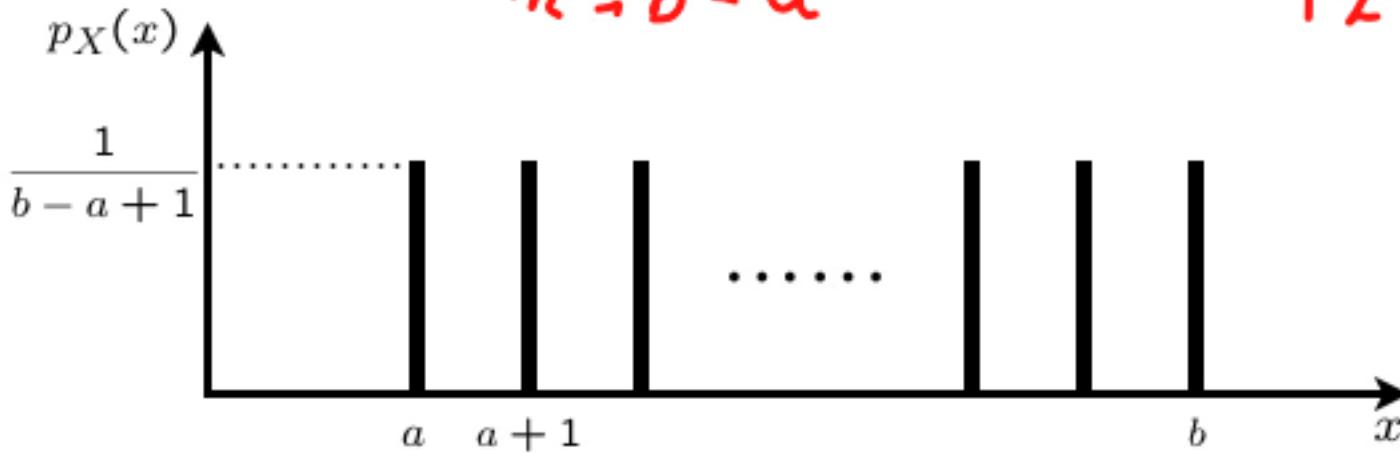
$$\begin{aligned} \text{var}(X) &= E[X^2] - (E[X])^2 = E[X] - (E[X])^2 = p - p^2 = \boxed{p(1-p)} \\ X^2 &= X \end{aligned}$$

Variance of the uniform



$$\text{var}(x) = E[x^2] - (E[x])^2 = \frac{1}{n+1} (0^2 + 1^2 + 2^2 + \dots + n^2) - \left(\frac{n}{2}\right)^2$$

$$= \frac{1}{12} n(n+2)$$



$$\text{Var}(x) = \frac{1}{12} (b-a)(b-a+2)$$

Conditional PMF and expectation, given an event

- Condition on an event $A \Rightarrow$ use conditional probabilities

$$p_X(x) = \mathbf{P}(X = x)$$

$$\sum_x p_X(x) = 1$$

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

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

$$\underline{p_{X|A}(x)} = \underline{\mathbf{P}(X = x | A)}$$

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

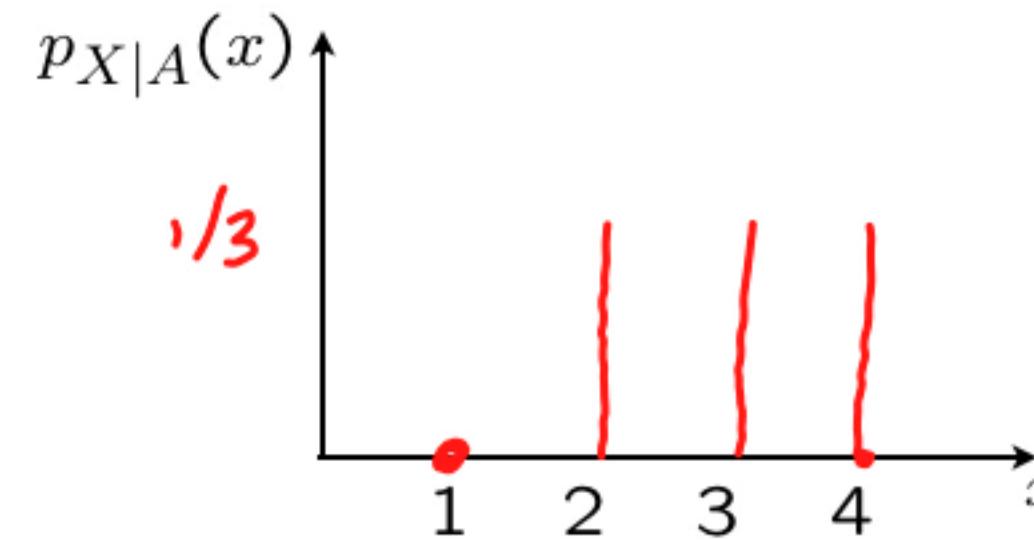
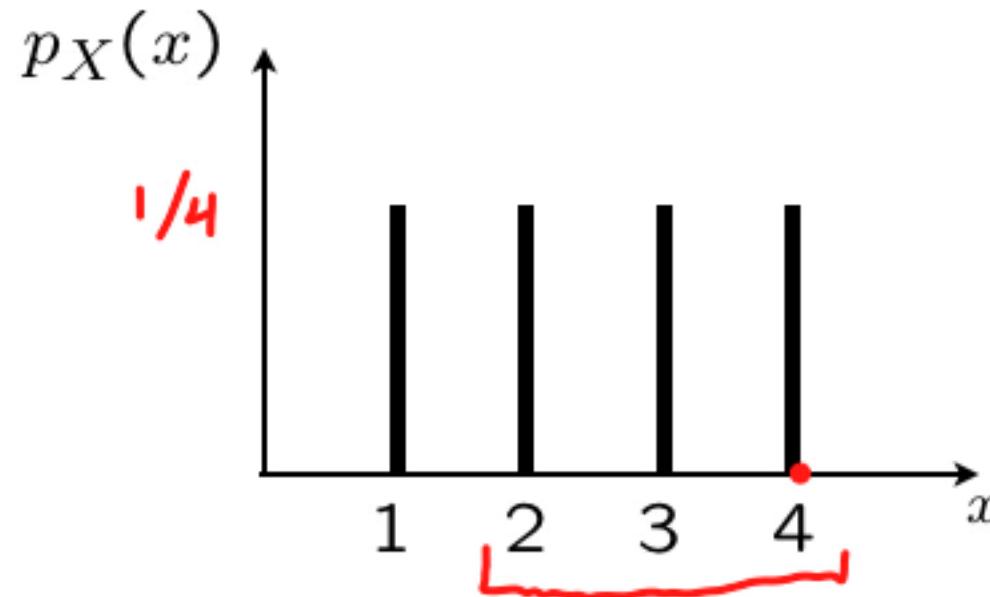
$$\mathbf{E}[X | A] = \sum_x x p_{X|A}(x)$$

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

assume
 $\mathbf{P}(A) > 0$

Example of conditioning

- Let $A = \{X \geq 2\}$



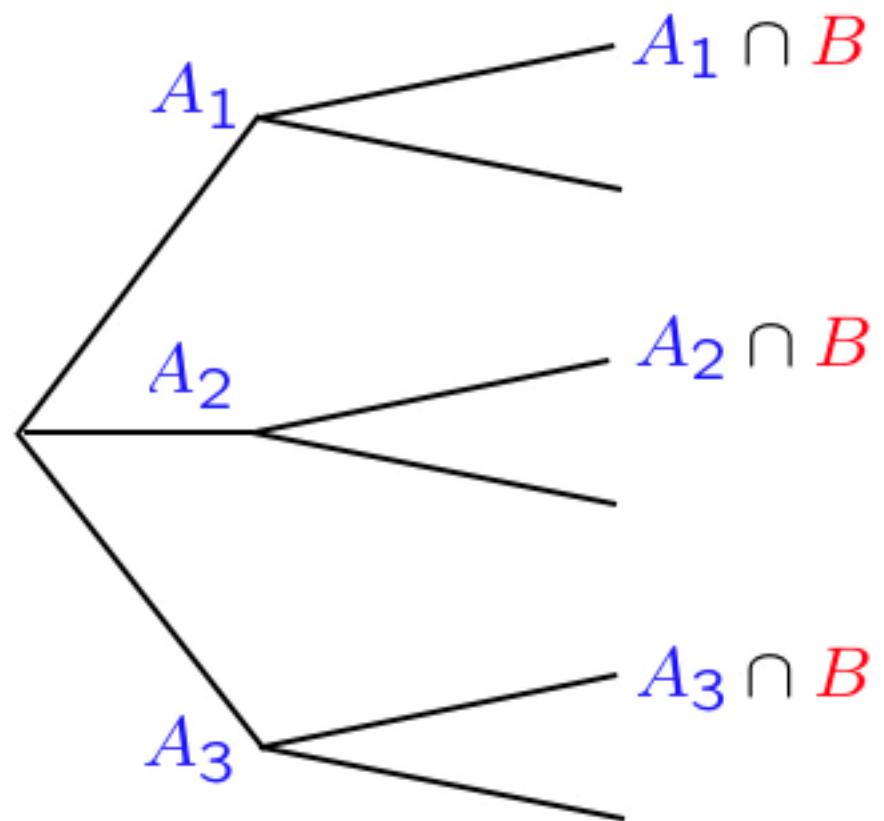
$$E[X] = 2.5$$

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

$$\begin{aligned}\text{var}(X) &= \frac{1}{12}(b-a)(b-a+2) \\ &= \frac{1}{12} 3 \cdot 5 = \frac{5}{4}\end{aligned}$$

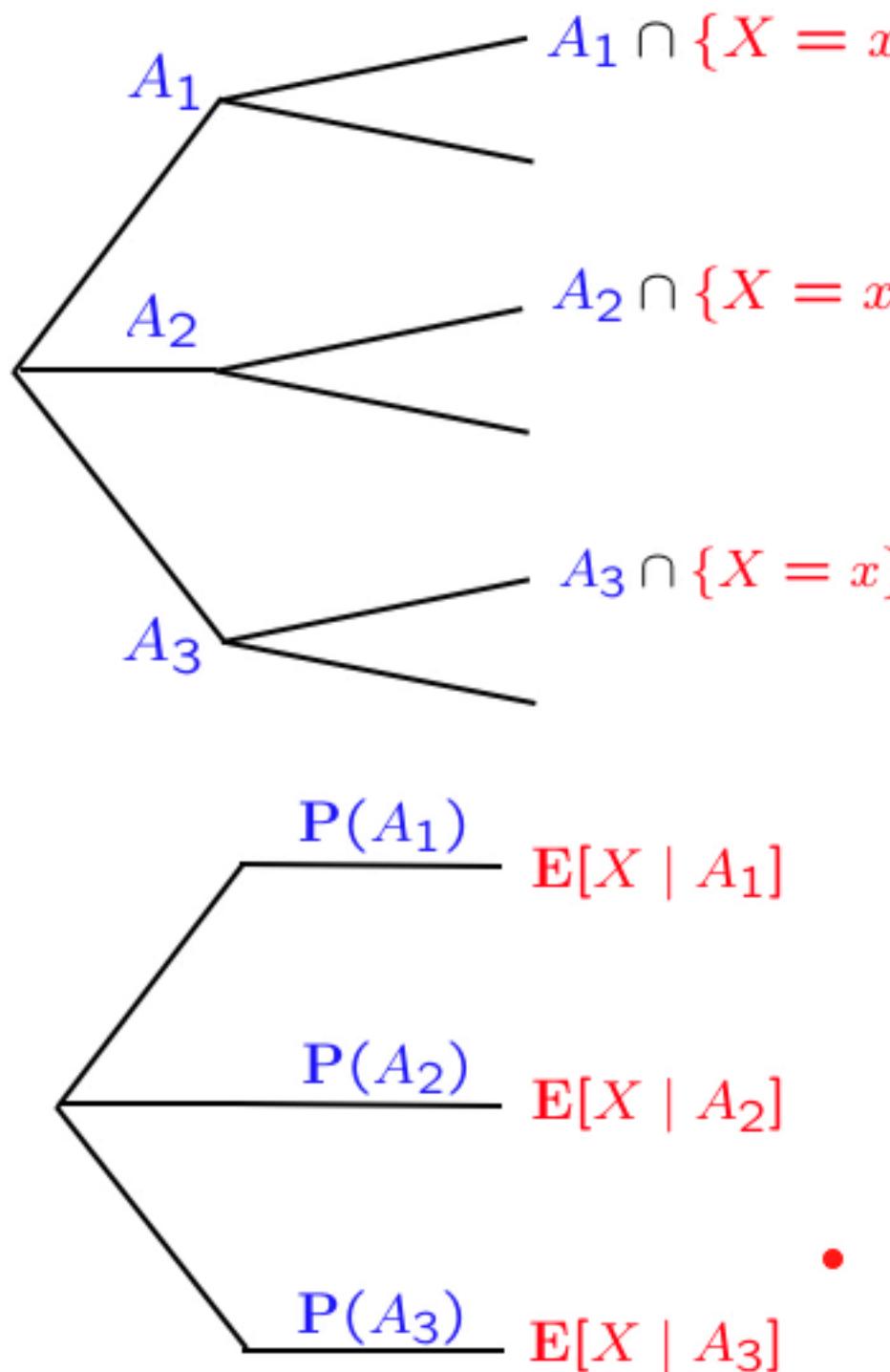
$$\begin{aligned}\text{var}(X | A) &= \frac{1}{3} (4-3)^2 + \frac{1}{3} (3-3)^2 \\ &\quad + \frac{1}{3} (2-3)^2 = \frac{2}{3}\end{aligned}$$

Total expectation theorem



$$P(B) = P(A_1)P(B | A_1) + \dots + P(A_n)P(B | A_n)$$
$$\mathcal{B} = \{x = \alpha\}$$

Total expectation theorem



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

$$B = \{x = x\}$$

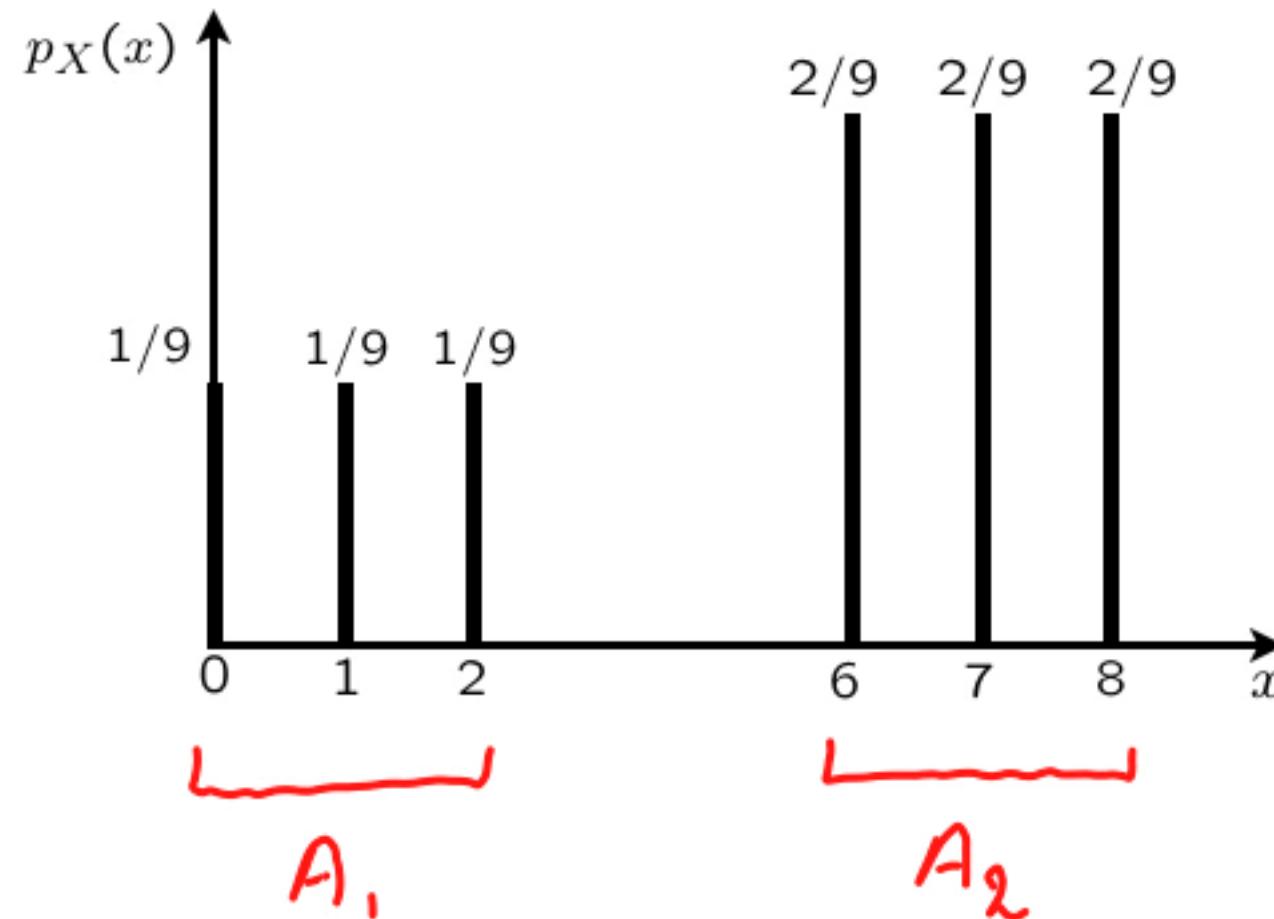
$$p_X(x) = P(A_1)p_{X|A_1}(x) + \cdots + P(A_n)p_{X|A_n}(x)$$

for all x

$$\sum_x x p_X(x) = P(A_1) \underbrace{\sum_x x p_{X|A_1}(x)}_{E[X | A_1]} + \cdots$$

$$E[X] = P(A_1)E[X | A_1] + \cdots + P(A_n)E[X | A_n]$$

Total expectation example



$$P(A_1) = \frac{1}{3}$$

$$P(A_2) = \frac{2}{3}$$

$$E[x|A_1] = 1$$

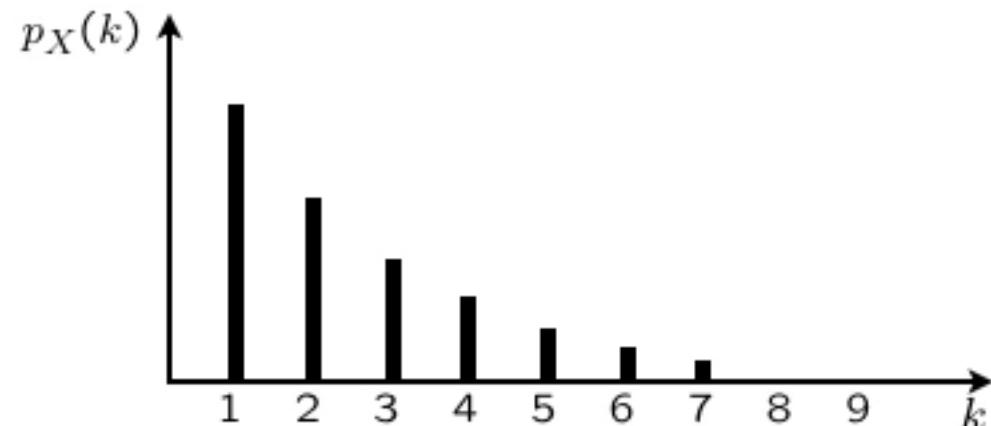
$$E[x|A_2] = 7$$

$$E[x] = \frac{1}{3} \cdot 1 + \frac{2}{3} \cdot 7 .$$

Conditioning a geometric random variable

- X : number of independent coin tosses until first head; $P(H) = p$

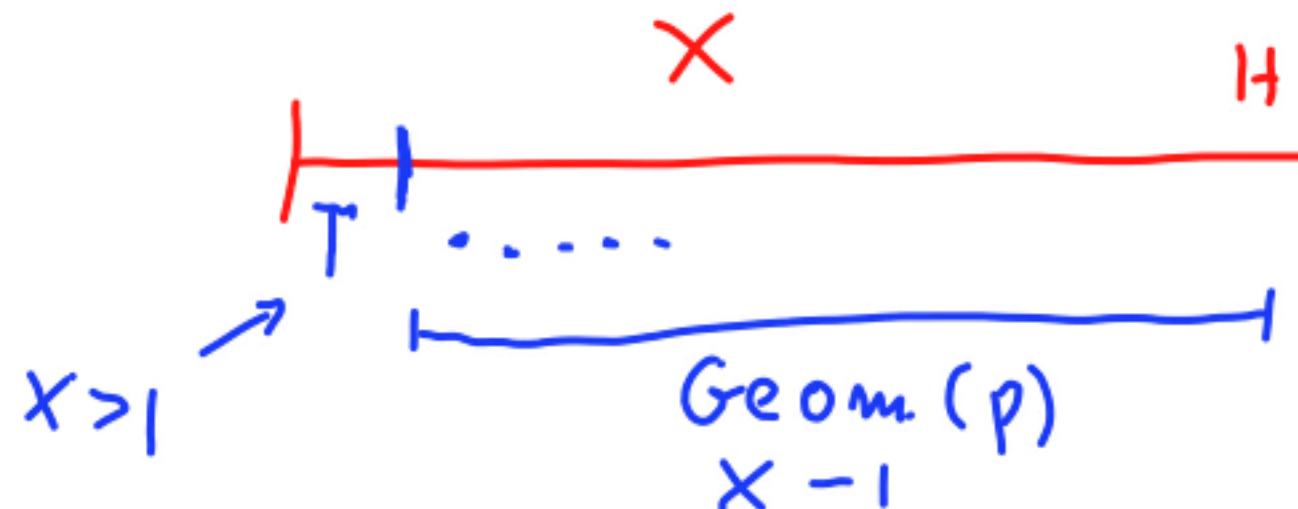
$$p_X(k) = (1 - p)^{k-1} p, \quad k = 1, 2, \dots$$



Memorylessness:

Number of **remaining** coin tosses, conditioned on Tails in the first toss, is **Geometric**, with parameter p

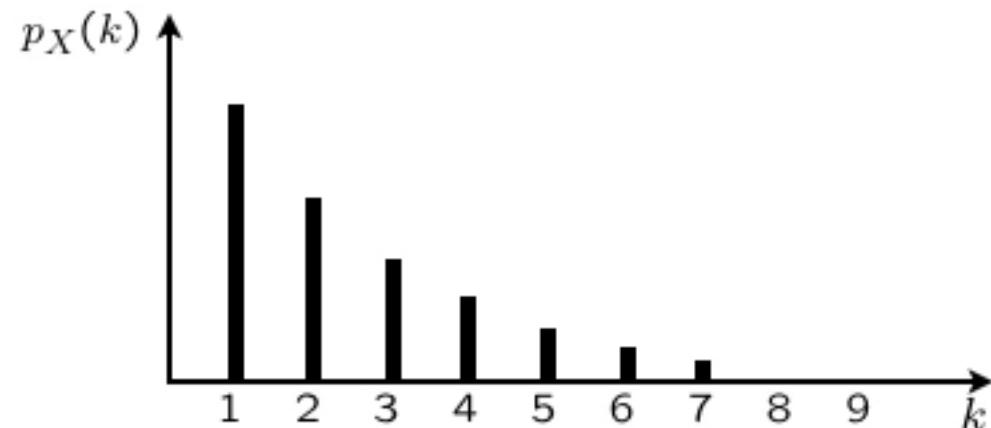
Conditioned on $X > 1$, $X - 1$ is geometric with parameter p



Conditioning a geometric random variable

- X : number of independent coin tosses until first head; $P(H) = p$

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



Memorylessness:

Number of **remaining** coin tosses, conditioned on Tails in the first toss, is **Geometric**, with parameter p

Conditioned on $X > 1$, $X - 1$ is geometric with parameter p

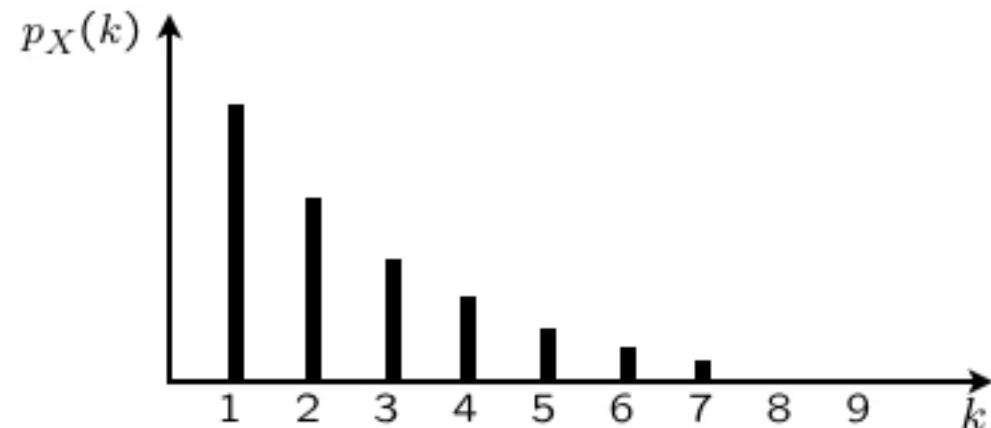
$$P_{X-1|X>1}(3) = P(X-1=3 | X>1) = P(T_2 T_3 H_4 | \bar{T}_1) = P(T_2 T_3 H_4)$$

$$= (1-p)^2 p = p_X(3)$$
$$P_{X-1|X>1}(k) = p_X(k)$$

Conditioning a geometric random variable

- X : number of independent coin tosses until first head; $P(H) = p$

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



Memorylessness:

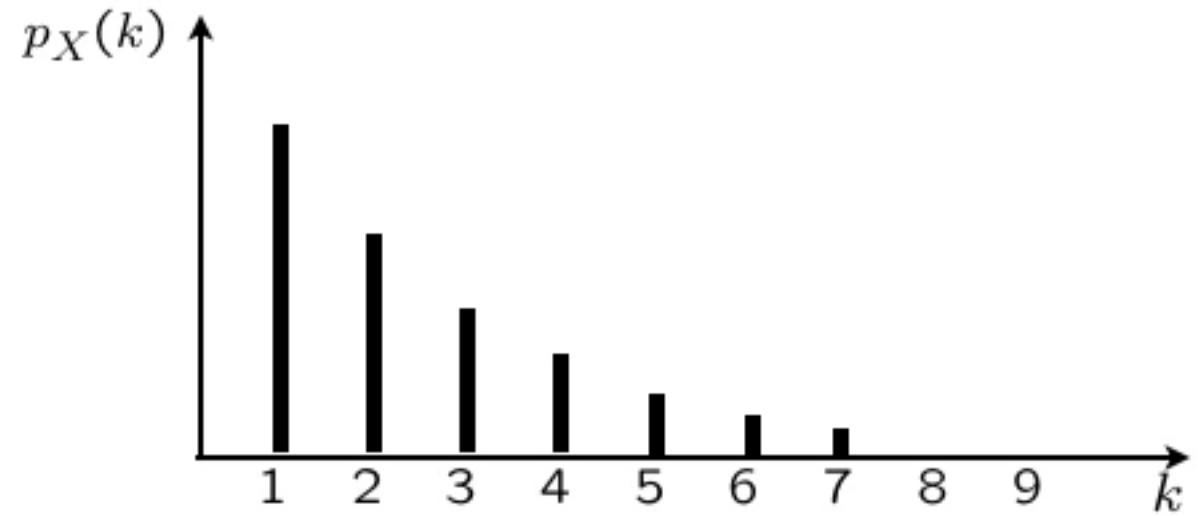
Number of **remaining** coin tosses, conditioned on Tails in the first toss, is **Geometric**, with parameter p

Conditioned on $X > n$, $X - n$ is geometric with parameter p

$$P_{X-1|X>1}(3) = P(X-1=3 | X>1) = P(T_2 T_3 H_4 | \bar{T}_1) = P(T_2 T_3 H_4)$$

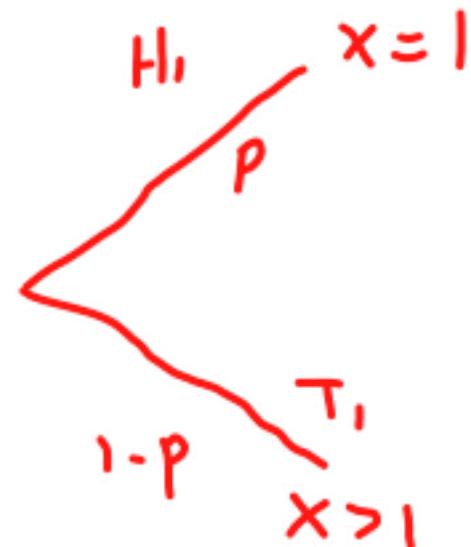
$$P_{X-1|X>1}(k) = P_X(k) = P_{X-n|X>n}(k) \cdot \\ = (1-p)^{k-n} p = p_X(n)$$

The mean of the geometric



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

$$E[X] = \frac{1}{p}$$



$$\begin{aligned} E[X] &= 1 + E[X-1] \\ &= 1 + p \cdot E[X-1 | X=1] + (1-p) E[X-1 | X>1] \\ &= 1 + 0 + (1-p) E[X] \end{aligned}$$

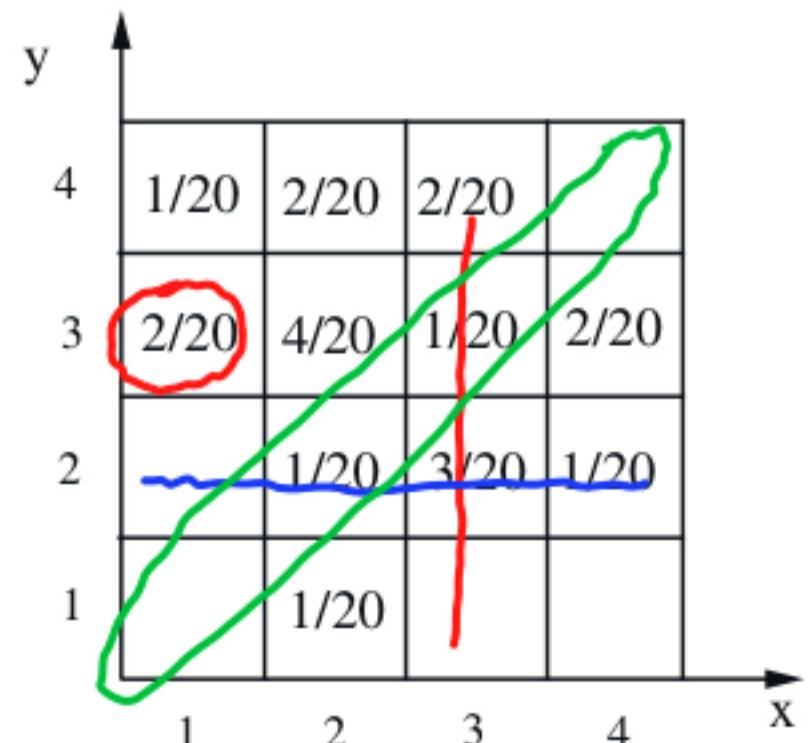
Multiple random variables and joint PMFs

$X : p_X$
 $Y : p_Y$

$\leftarrow \text{marginal pmfs}$

$P(X = Y) = \frac{2}{20}$

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



$$P_x(3)$$

$$P_{x,y}(1,3) = \frac{2}{20}$$

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

$$P_x(4) = \frac{1}{20} + \frac{2}{20}$$

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

$$P_y(2) = \frac{1}{20} + \frac{3}{20} + \frac{1}{20}$$

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

More than two random variables

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

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

$$p_X(\textcolor{red}{x}) = \sum_{\textcolor{blue}{y}} \sum_z p_{X,Y,Z}(\textcolor{red}{x}, \textcolor{red}{y}, \textcolor{blue}{z})$$

$$p_{X,Y}(x,y) = \sum_z p_{X,Y,Z}(x,y,\textcolor{blue}{z})$$

Functions of multiple random variables

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

PMF: $p_Z(z) = \mathbf{P}(Z = z) = \mathbf{P}(g(X, Y) = z) = \sum_{(x, y) : g(x, y) = z} p_{X,Y}(x, y)$

Expected value rule: $\mathbf{E}[g(X, Y)] = \sum_x \sum_y g(x, y) p_{X,Y}(x, y)$

$$\mathbf{E}[g(x)]$$

Linearity of expectations

$$E[aX + b] = aE[X] + b$$

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

$$E[x + y] = E[g(x, y)]$$

$$(g(x, y) = x + y)$$

$$= \sum_x \sum_y (x + y) P_{x,y}(x, y)$$

$$= \underbrace{\sum_x \sum_y x P_{x,y}(x, y)} + \underbrace{\sum_x \sum_y y P_{x,y}(x, y)}$$

$$= \underbrace{\sum_x x \sum_y P_{x,y}(x, y)} + \underbrace{\dots}$$

$$= \sum_x x P_x(x) + \sum_y y P_y(y) = E[X] + E[Y]$$

Linearity of expectations

$$\mathbb{E}[aX + b] = a\mathbb{E}[X] + b$$

$$\mathbb{E}[X + Y] = \mathbb{E}[X] + \mathbb{E}[Y]$$

$$\mathbb{E}[X_1 + \dots + X_n] = \mathbb{E}[X_1] + \dots + \mathbb{E}[X_n]$$

$$\mathbb{E}[2X + 3Y - Z] = E[2x] + E[3y] - E[z] = 2E[x] + 3E[y] - E[z]$$

The mean of the binomial

- X : binomial with parameters n, p
 - number of successes in n independent trials

$X_i = 1$ if i th trial is a success;
 $X_i = 0$ otherwise

\xrightarrow{P} (indicator variable)
 $\xrightarrow{1-p}$

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

$P_X(k)$

$$E[X] = np$$

$$X = X_1 + \cdots + X_n$$

$$E[X] = \underbrace{E[X_1]}_p + \cdots + \underbrace{E[X_n]}_p = np$$

LECTURE 7: Conditioning on a random variable; Independence of r.v.'s

- Conditional PMFs
 - Conditional expectations
 - Total expectation theorem
- Independence of r.v.'s
 - Expectation properties
 - Variance properties
- The variance of the binomial
- The hat problem: mean and variance

Conditional PMFs

$$A = \{Y = y\}$$

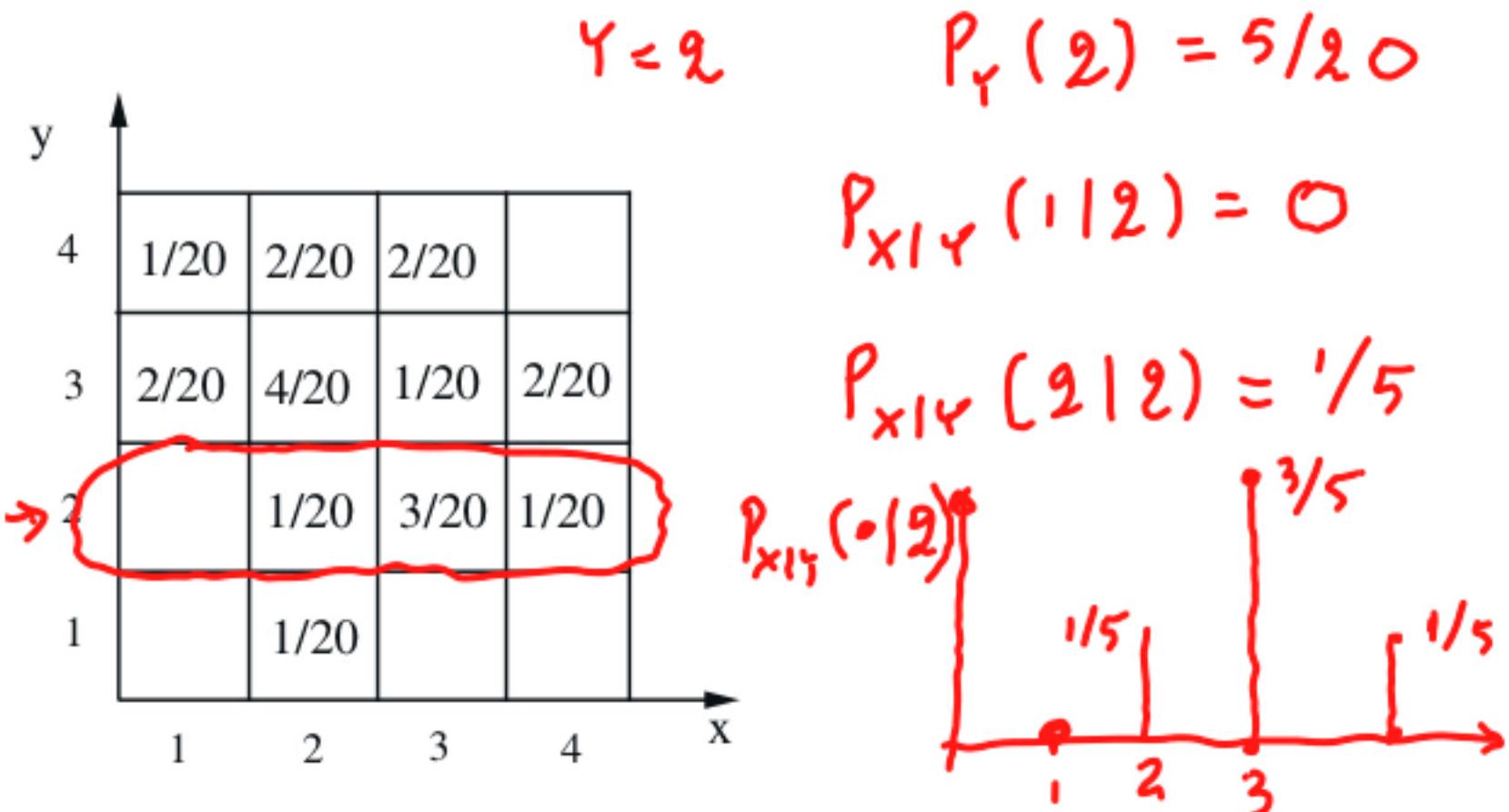
$$p_{X|A}(x | A) = P(X = x | A)$$

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

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

defined for y such that $p_Y(y) > 0$

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



Conditional PMFs involving more than two r.v.'s

- Self-explanatory notation

$$p_{X|Y,Z}(x | y, z) = \underline{P(X=x | Y=y, Z=z)} = \frac{\underline{P(X=x, Y=y, Z=z)}}{\underline{P(Y=y, Z=z)}} = \frac{P_{x,y,z}(x, y, z)}{P_{y,z}(y, z)}$$

$$p_{X,Y|Z}(x, y | z) = \underline{P(X=x, Y=y | Z=z)}$$

- Multiplication rule

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

$$A = \{X=x\} \quad B = \{Y=y\} \quad C = \{Z=z\}$$

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

Conditional expectation

$$A = \{Y = y\}$$

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

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

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

- Expected value rule

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

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

Total probability and expectation theorems

- A_1, \dots, A_n : partition of Ω $Y = \{y_1, \dots, y_n\}$ $A_i = \{Y = y_i\}$
- $p_X(x) = P(A_1)p_{X|A_1}(x) + \dots + P(A_n)p_{X|A_n}(x)$

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

- $E[X] = P(A_1)E[X | A_1] + \dots + P(A_n)E[X | A_n]$

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

- Fine print:
Also valid when Y is a discrete r.v. that ranges over an infinite set,
as long as $E[|X|] < \infty$

Independence

- of two events:

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

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

- of a r.v. and an event:

$$P(\underline{X=x} \text{ and } \underline{A}) = P(X=x) \cdot P(A), \quad \text{for all } \underline{\underline{x}}$$

$$p_{x|A}(x) = p_x(x), \text{ for all } x$$

$$P(A | X=x) = P(A), \text{ for all } x$$

- of two r.v.'s:

$$P(\underline{X=x} \text{ and } \underline{Y=y}) = P(X=x) \cdot P(Y=y), \quad \text{for all } \underline{\underline{x, y}}$$

$$p_{x|y}(x|y) = p_x(x)$$

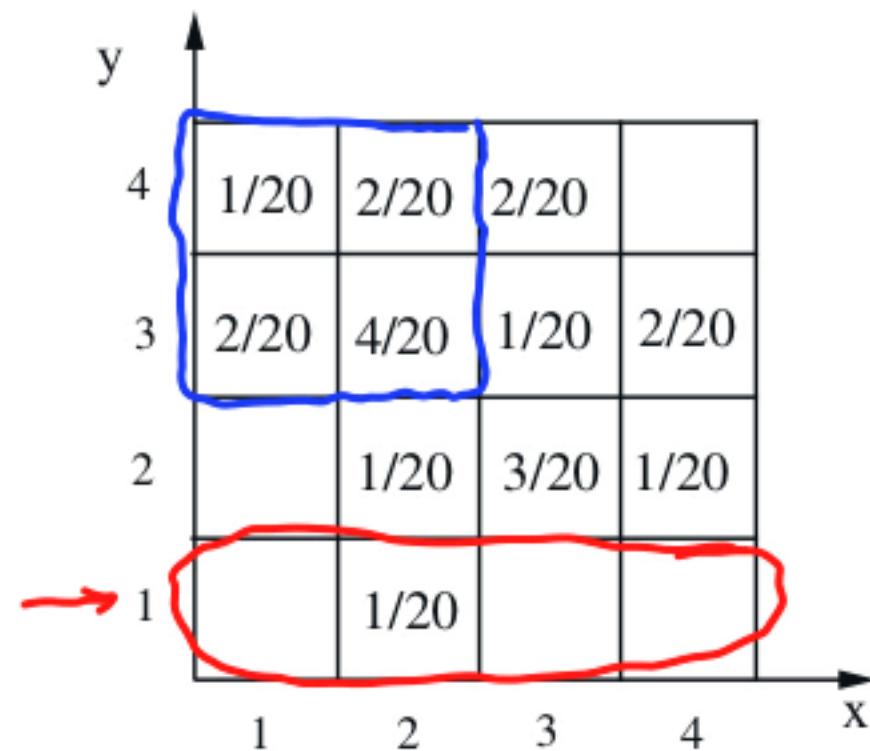
$$p_{X,Y}(x,y) = p_X(x)p_Y(y), \quad \text{for all } x, y$$

$$p_{y|x}(y|x) = p_y(y)$$

X, Y, Z are **independent** if:

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

Example: independence and conditional independence



- Independent? **No**
- $P_X(1) = 3/20$
- $P_{X|Y}(1|1) = 0$
- What if we condition on $X \leq 2$ and $Y \geq 3$?

1/9	2/9
2/9	4/9

Yes .

Independence and expectations

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

always true

- Exceptions: $E[aX + b] = aE[X] + b$

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

If X, Y are **independent**: $E[XY] = E[X]E[Y]$

$g(X)$ and $h(Y)$ are also independent: $E[g(X)h(Y)] = E[g(X)] \cdot E[h(Y)]$

$$E[g(x, y)] \quad g(x, y) = xy$$

$$= \sum_x \sum_y xy P_{x,y}(x, y) = \sum_x \sum_y \underbrace{xy}_{\text{xy}} P_x(x) P_y(y)$$

$$= \sum_x x P_x(x) \underbrace{\sum_y y P_y(y)}_{\text{E[Y]}} = E[x] E[y]$$

Independence and variances

- Always true: $\text{var}(aX) = a^2\text{var}(X)$ $\text{var}(X + a) = \text{var}(X)$
- In general: $\text{var}(X + Y) \neq \text{var}(X) + \text{var}(Y)$

If X, Y are independent: $\text{var}(X + Y) = \text{var}(X) + \text{var}(Y)$

assume
 $E[X] = E[Y] = 0$
 $E[XY] = E[X]E[Y] = 0$

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

- Examples:

- If $X = Y$: $\text{var}(X + Y) = \text{var}(2X) = 4\text{var}(X)$

- If $X = -Y$: $\text{var}(X + Y) = \text{var}(0) = 0$

- If X, Y independent: $\text{var}(X - 3Y) = \text{var}(X) + \text{var}(-3Y) = \text{var}(X) + 9\text{var}(Y)$

Variance of the binomial

- X : binomial with parameters n, p
 - number of successes in n independent trials

$$\begin{aligned} X_i &= 1 \text{ if } i\text{th trial is a success;} \\ X_i &= 0 \text{ otherwise} \end{aligned} \quad (\text{indicator variable})$$

$$X = X_1 + \cdots + X_n$$

$$\begin{aligned}\text{var}(x) &= \text{var}(X_1) + \dots + \text{var}(X_n) \\ &= n \cdot \text{var}(X_1) = \boxed{n p(1-p)}\end{aligned}$$

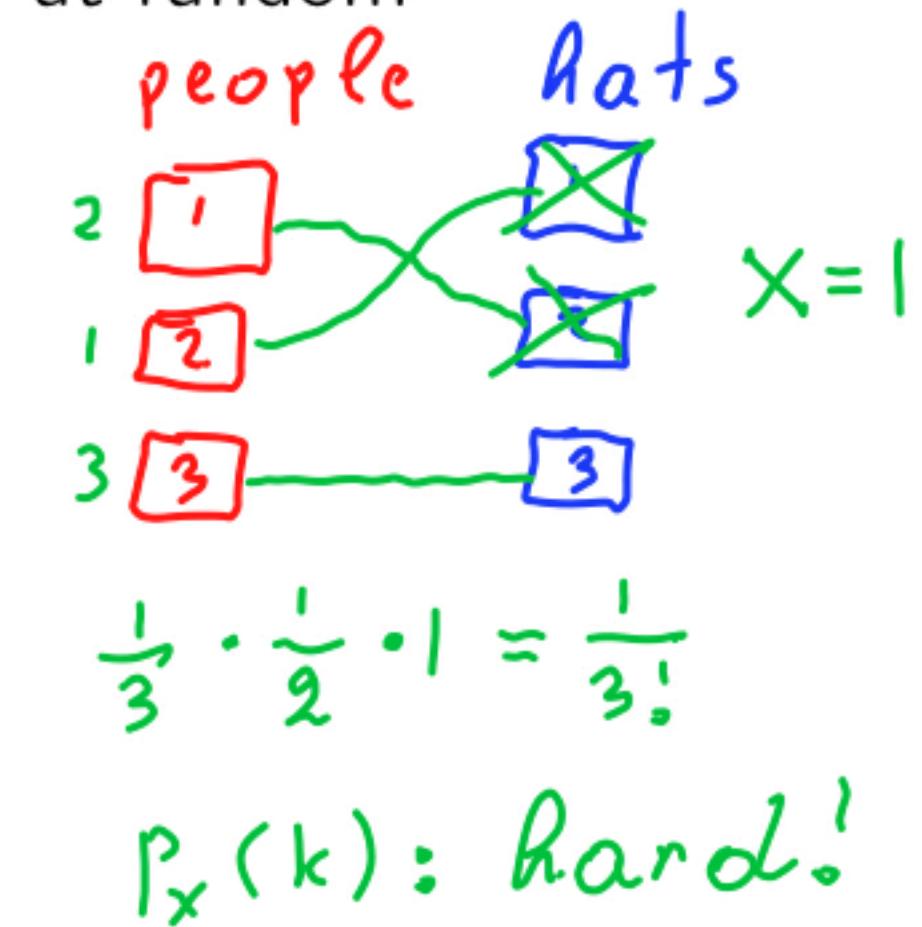
The hat problem

- n people throw their hats in a box and then pick one at random
 - All permutations equally likely $1/n!$
 - Equivalent to picking one hat at a time
- X : number of people who get their own hat
 - Find $E[X] = E[X_1] + \dots + E[X_n] = n \cdot \frac{1}{n} = 1$

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

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

$$\bullet E[X_i] = E[X_1] = P(X_1 = 1) = \frac{1}{n}$$



$$\sum_k k p_X(k)$$

The variance in the hat problem

- X : number of people who get their own hat
 - Find $\text{var}(X)$

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

$$\begin{aligned} n &= 2 \\ X_1 = 1 &\Rightarrow X_2 = 1 \\ X_1 = 0 &\Rightarrow X_2 = 0 \end{aligned}$$

$$X = X_1 + X_2 + \cdots + X_n$$

$$n(n-1)$$

$$n^2 - n$$

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

$$X^2 = \underbrace{\sum_i X_i^2}_{n} + \underbrace{\sum_{i,j: i \neq j} X_i X_j}_{n^2 - n}.$$

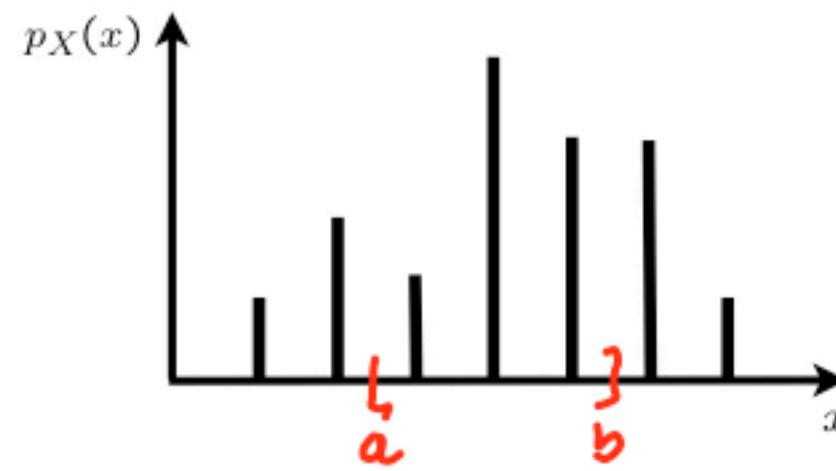
$$\bullet E[X_i^2] = E[X_1^2] = E[X_1] = 1/n \quad E[X^2] = n \cdot \frac{1}{n} + n(n-1) \cdot \frac{1}{n} \cdot \frac{1}{n-1},$$

$$\bullet \text{For } i \neq j: E[X_i X_j] = E[X_1 X_2] = P(X_1, X_2 = 1) = P(X_1 = 1, X_2 = 1) \\ = P(X_1 = 1) P(X_2 = 1 | X_1 = 1) = \frac{1}{n} \cdot \frac{1}{n-1},$$

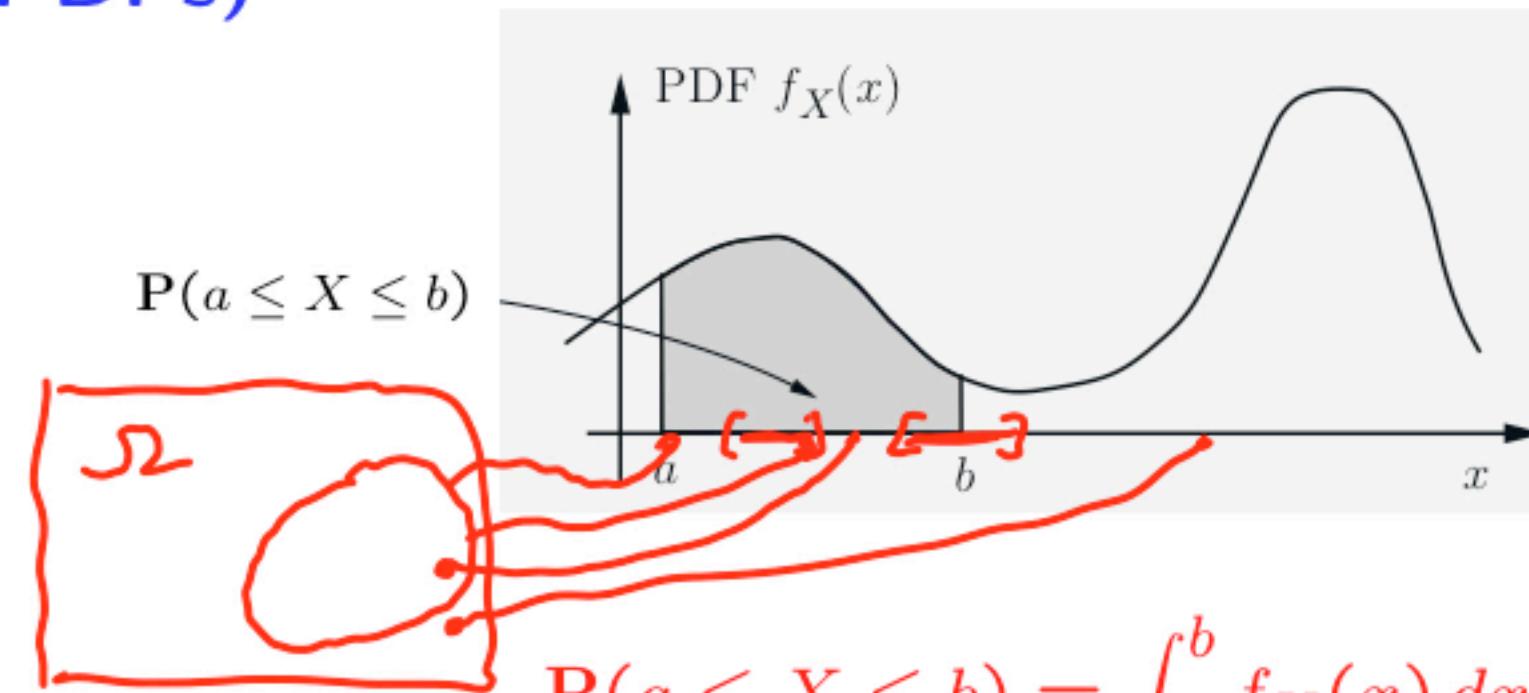
LECTURE 8: Continuous random variables and probability density functions

- Probability density functions
 - Properties
 - Examples
- Expectation and its properties
 - The expected value rule
 - Linearity
- Variance and its properties
- Uniform and exponential random variables
- Cumulative distribution functions
- Normal random variables
 - Expectation and variance
 - Linearity properties
 - Using tables to calculate probabilities

Probability density functions (PDFs)



$$P(a \leq X \leq b) = \sum_{x: a \leq x \leq b} p_X(x)$$



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

$$p_X(x) \geq 0 \quad \sum_x p_X(x) = 1$$

- $f_X(x) \geq 0$
- $\int_{-\infty}^{\infty} f_X(x) dx = 1$

Definition: A random variable is **continuous** if it can be described by a **PDF**

$$P(1 \leq X \leq 3 \text{ or } 4 \leq X \leq 5) = P(1 \leq X \leq 3) + P(4 \leq X \leq 5)$$

Probability density functions (PDFs)

$\delta > 0$, small

$$P(a \leq X \leq a + \delta)$$

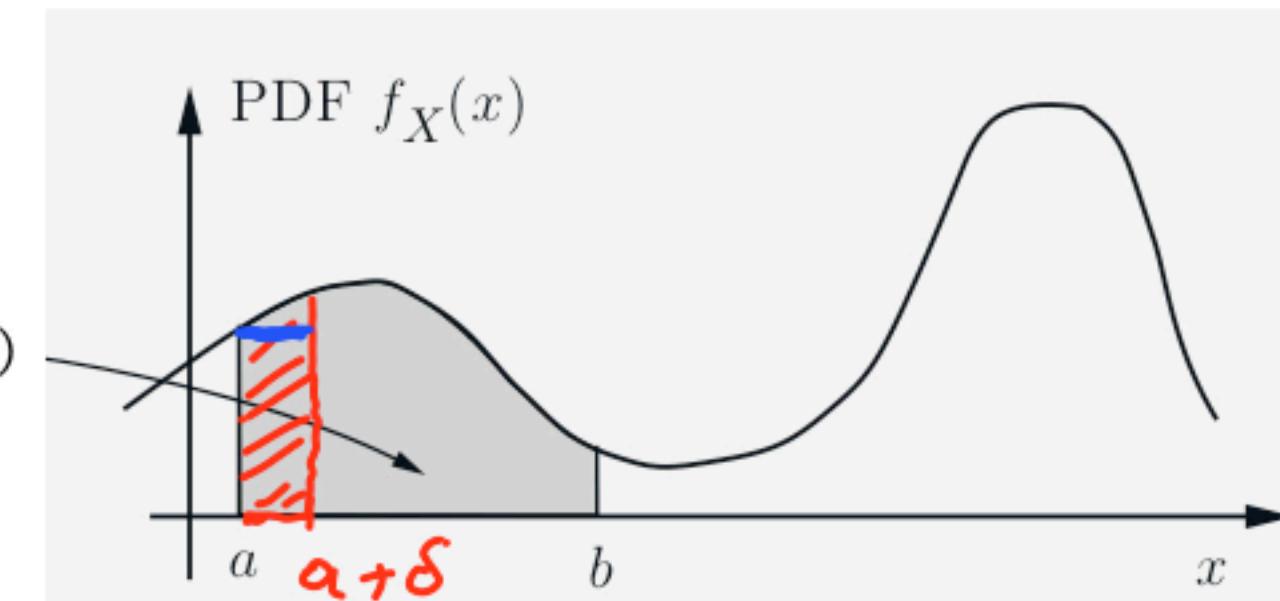
$$\approx f_X(a) \cdot \delta$$

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

$$P(X = a) = 0$$

~~$$P(a \leq X \leq b) = P(x=a) + P(x=b) + P(a < X < b)$$~~

$$P(a \leq X \leq b)$$

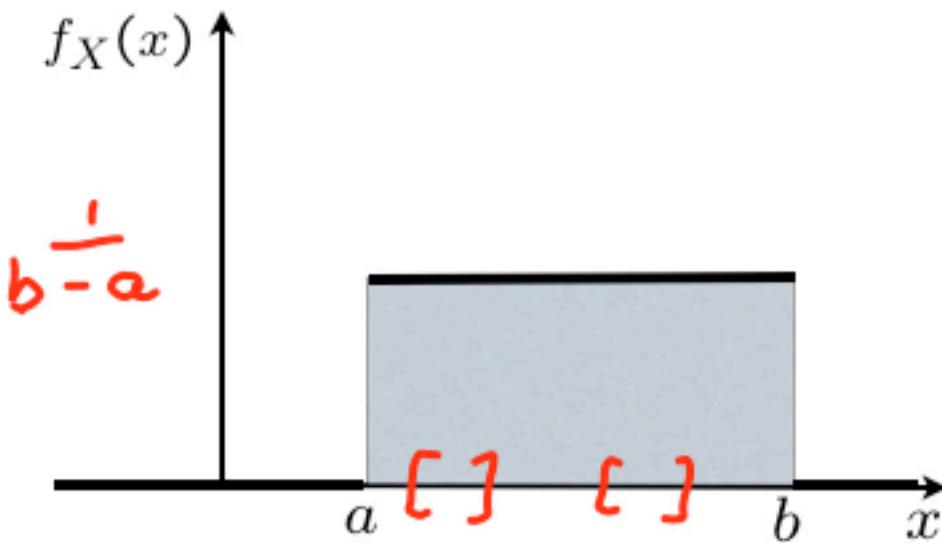
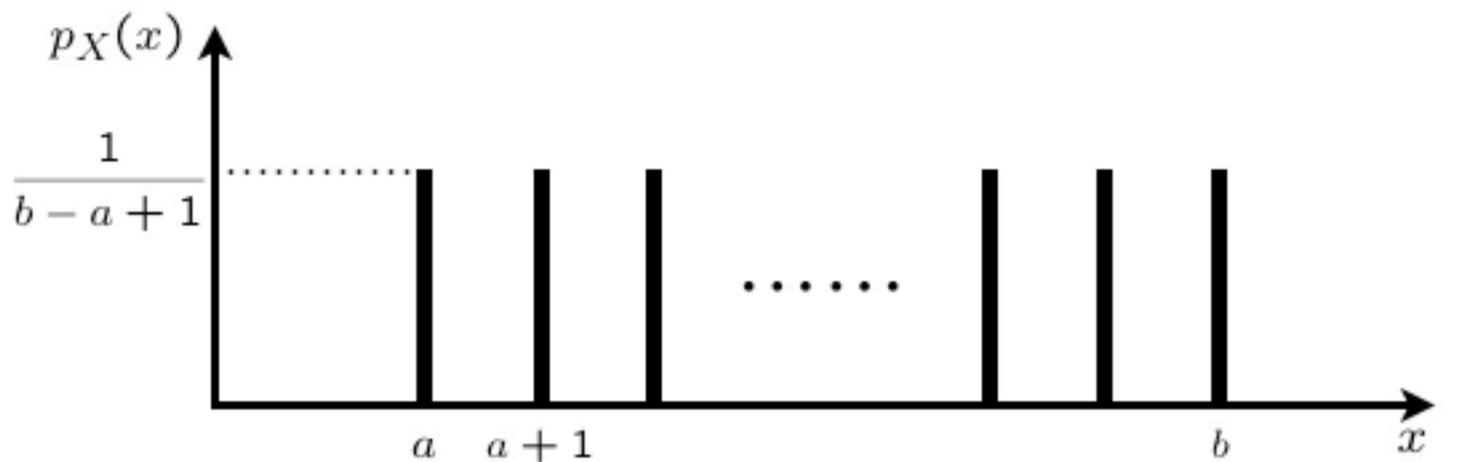


$$P(a \leq X \leq b) = \int_a^b f_X(x) dx$$

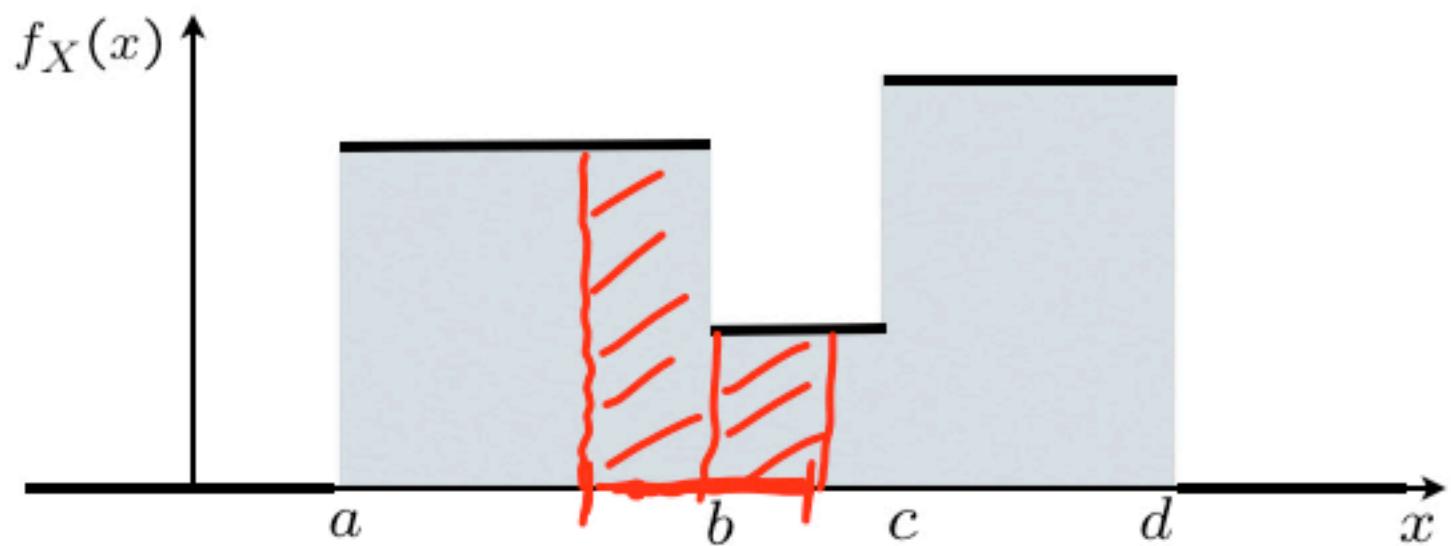
$$f_X(x) \geq 0$$

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

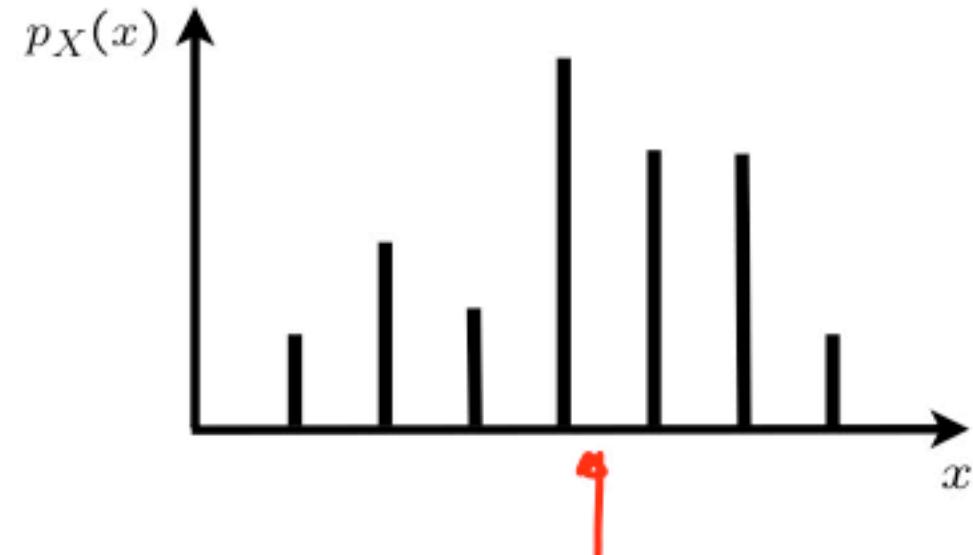
Example: continuous uniform PDF



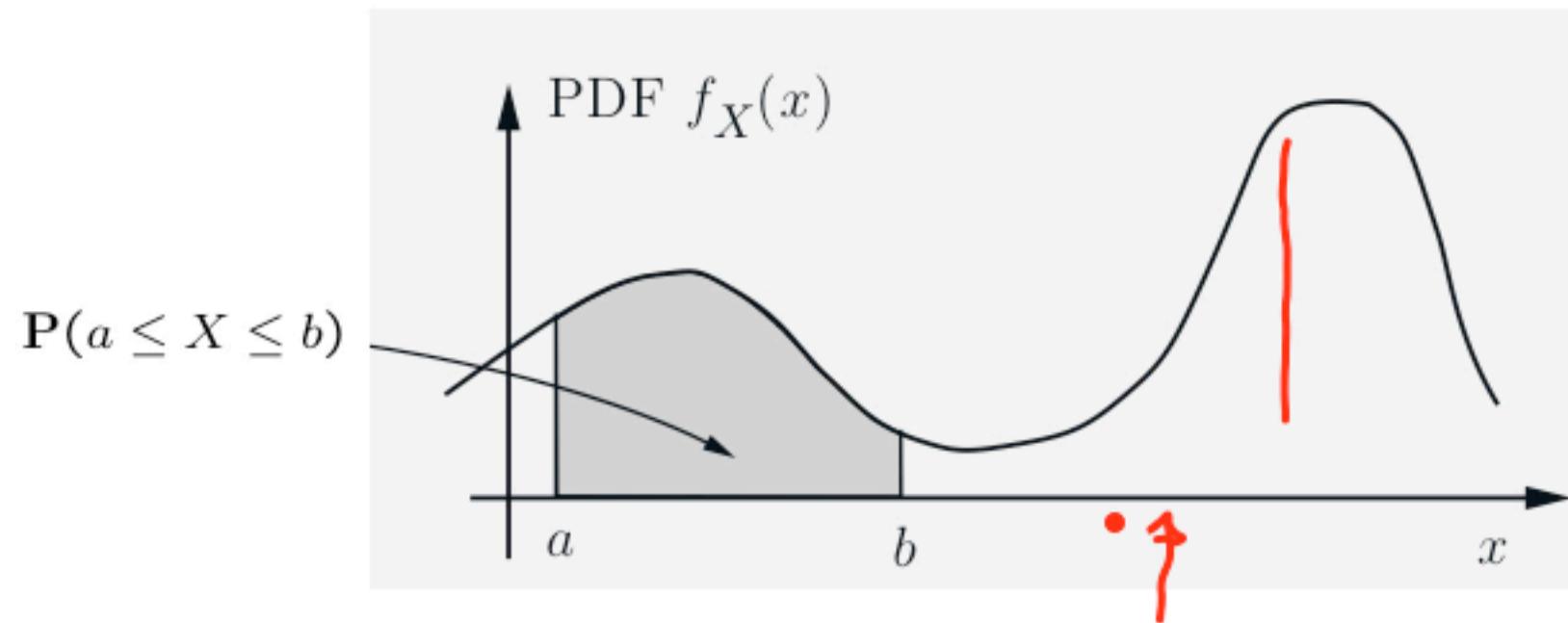
- Generalization: piecewise constant PDF



Expectation/mean of a continuous random variable



$$\mathbb{E}[X] = \sum_x x \underline{p_X(x)}$$



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

- **Interpretation:** Average in large number of independent repetitions of the experiment

Fine print:
Assume $\int_{-\infty}^{\infty} |x| f_X(x) dx < \infty$

Properties of expectations

- If $X \geq 0$, then $E[X] \geq 0$
- If $a \leq X \leq b$, then $a \leq E[X] \leq b$
- Expected value rule:

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

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

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

- Linearity

$$E[aX + b] = aE[X] + b$$

Variance and its properties

- **Definition of variance:** $\text{var}(X) = E[(X - \mu)^2]$

$$\mu = E[X]$$

- Calculation using the expected value rule, $E[g(X)] = \int_{-\infty}^{\infty} g(x)f_X(x) dx$

$$\text{var}(X) = \int_{-\infty}^{\infty} (x - \mu)^2 f_X(x) dx$$

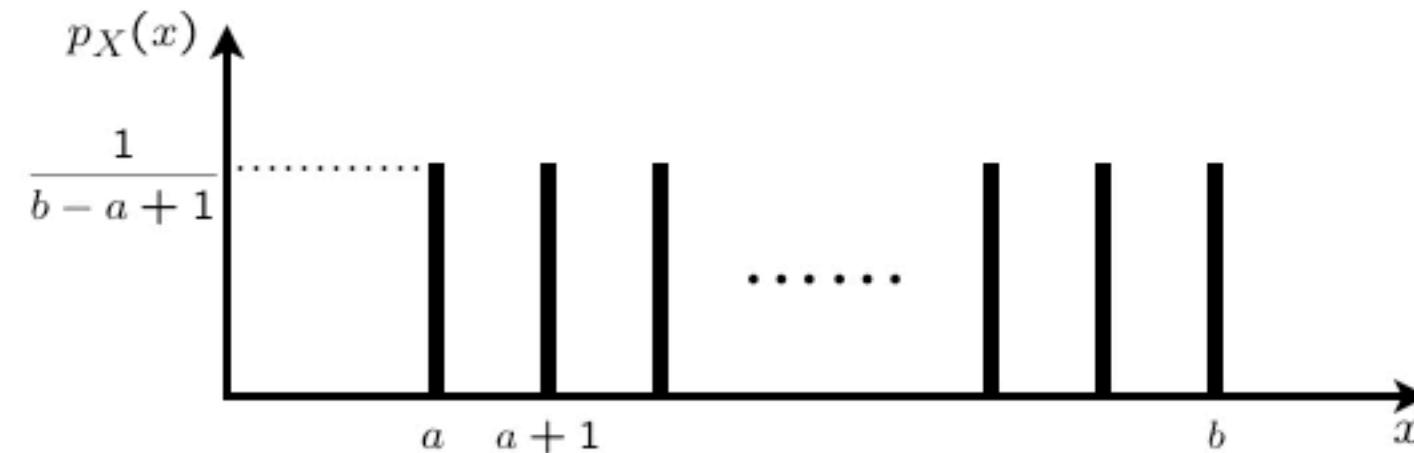
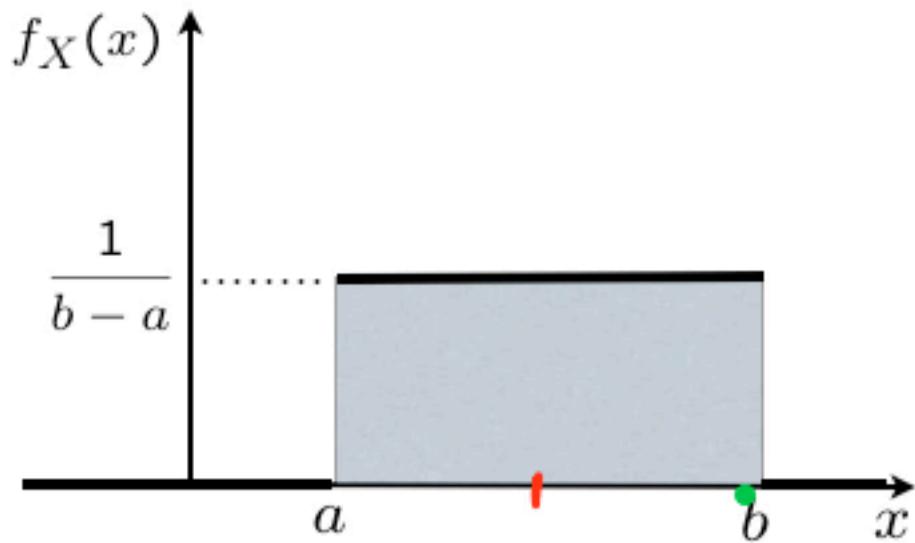
$$g(x) = (x - \mu)^2$$

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

✓ $\text{var}(aX + b) = a^2\text{var}(X)$

✓ A useful formula: $\text{var}(X) = E[X^2] - (E[X])^2$

Continuous uniform random variable; parameters a, b



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

$$= \int_a^b x \cdot \frac{1}{b-a} dx = \frac{a+b}{2}$$

$$E[X^2] = \int_a^b x^2 \cdot \frac{1}{b-a} dx = \frac{1}{b-a} \left(\frac{b^3}{3} - \frac{a^3}{3} \right)$$

$$E[X] = \frac{a+b}{2}$$

$$\text{var}(X) = \frac{1}{12}(b-a)(b-a+2)$$

$$\text{var}(X) = E[X^2] - (E[X])^2 = \boxed{\frac{(b-a)^2}{12}}$$

$$\sigma = \frac{b-a}{\sqrt{12}}$$

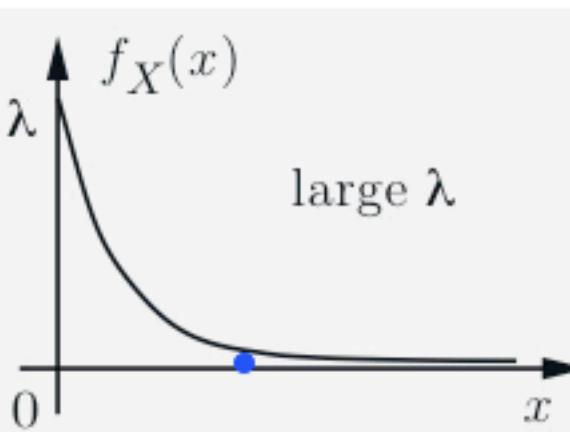
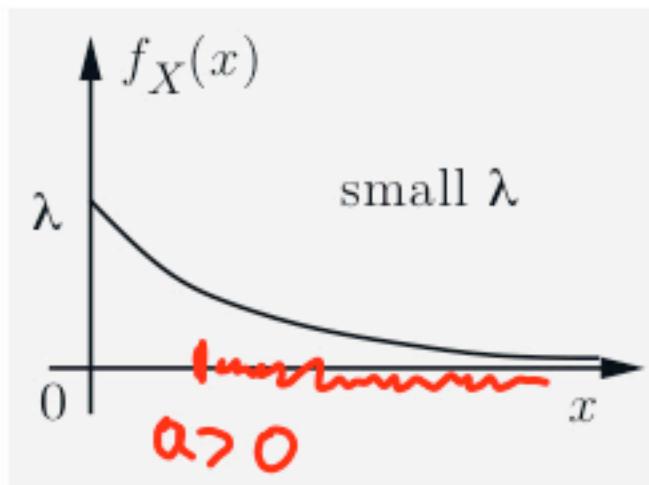
Exponential random variable; parameter $\lambda > 0$

$$f_X(x) = \begin{cases} \lambda e^{-\lambda x}, & x \geq 0 \\ 0, & x < 0 \end{cases}$$

$$\int f_X(x) dx = 1$$

$$E[X] = 1/p$$

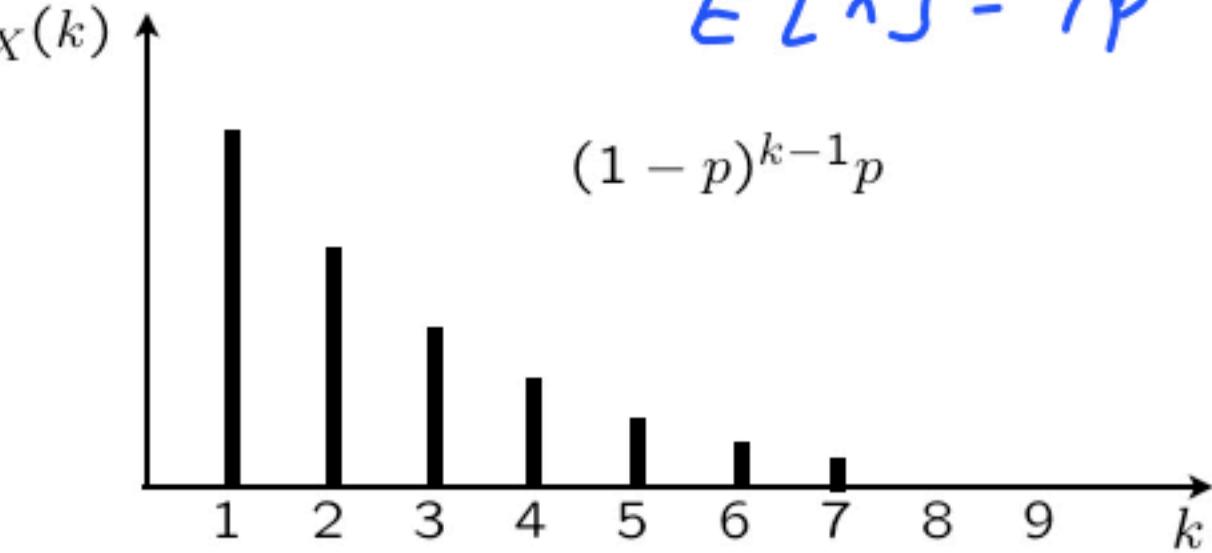
$$p_X(k) = (1 - p)^{k-1} p$$



$$E[X] = \int_0^{\infty} x \cdot \lambda e^{-\lambda x} dx = 1/\lambda$$

$$E[X^2] = \int_0^{\infty} x^2 \lambda e^{-\lambda x} dx = 2/\lambda^2$$

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



$$\boxed{P(X \geq a)} = \int_a^{\infty} \lambda e^{-\lambda x} dx$$

$$\left[\int e^{ax} dx = \frac{1}{a} e^{ax} \quad a \leftrightarrow -\lambda \right]$$

$$= \lambda \cdot \left(-\frac{1}{\lambda} \right) e^{-\lambda x} \Big|_a^{\infty}$$

$$= -e^{-\lambda \cdot 0} + e^{-\lambda a} = \boxed{e^{-\lambda a}}$$

Cumulative distribution function (CDF)

CDF definition: $F_X(x) = P(X \leq x)$

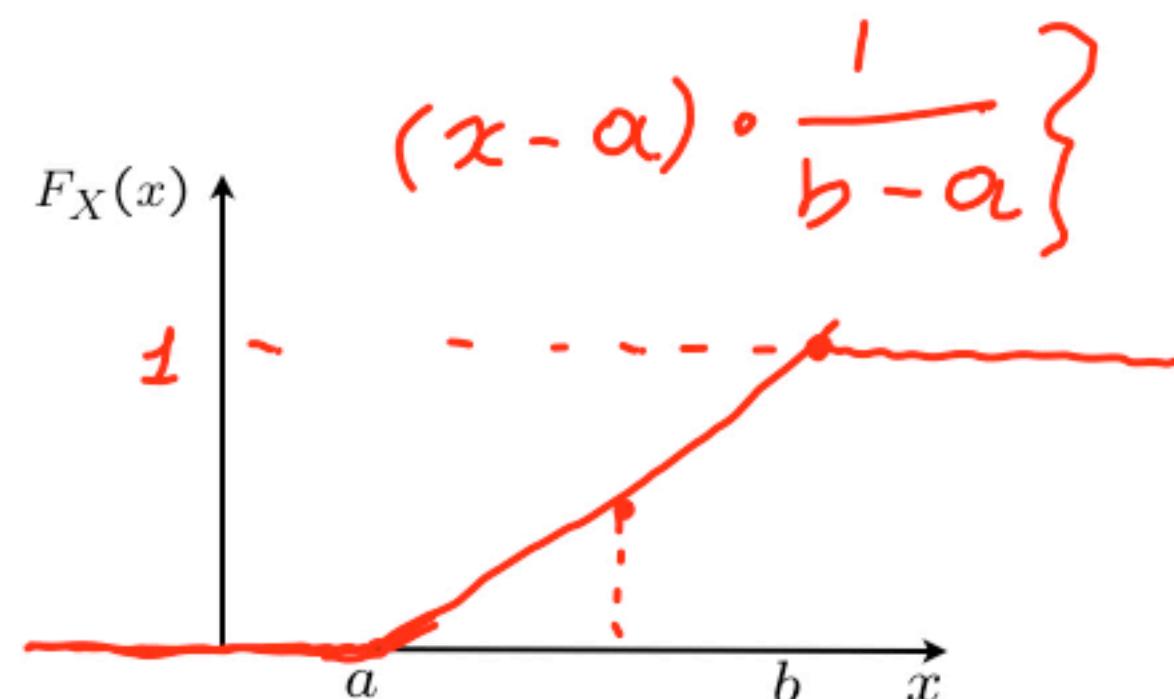
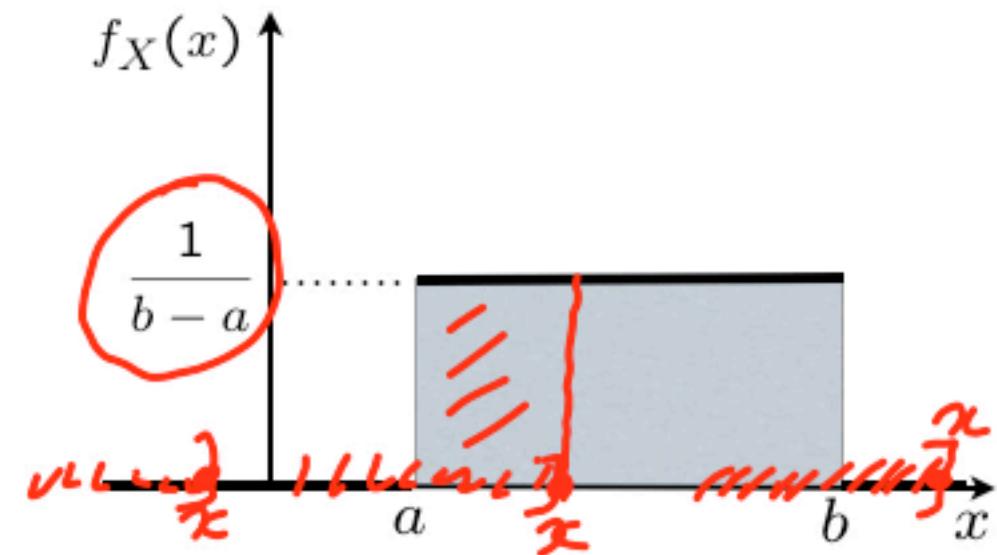
- Continuous random variables:

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



$$P(X \leq 4) = P(X \leq 3) + P(3 < X \leq 4)$$

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

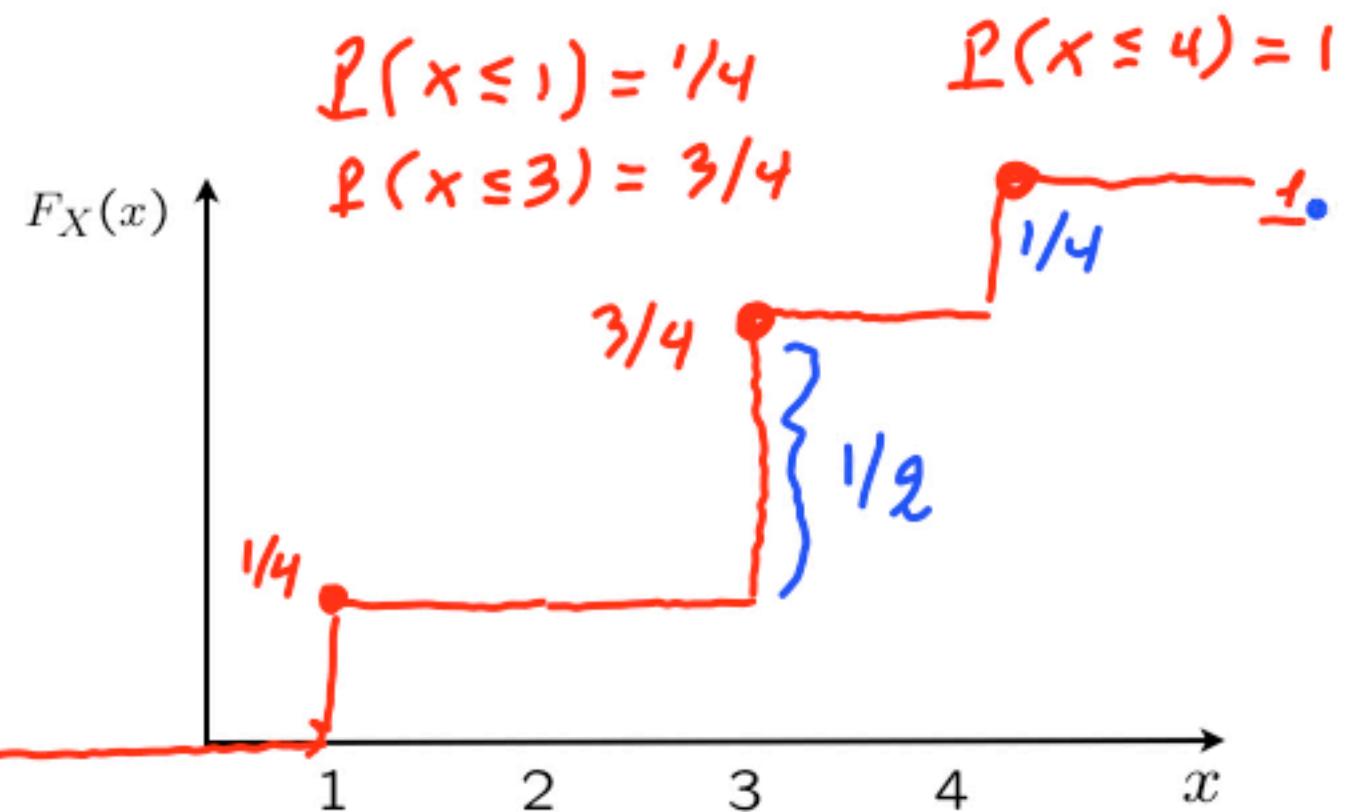
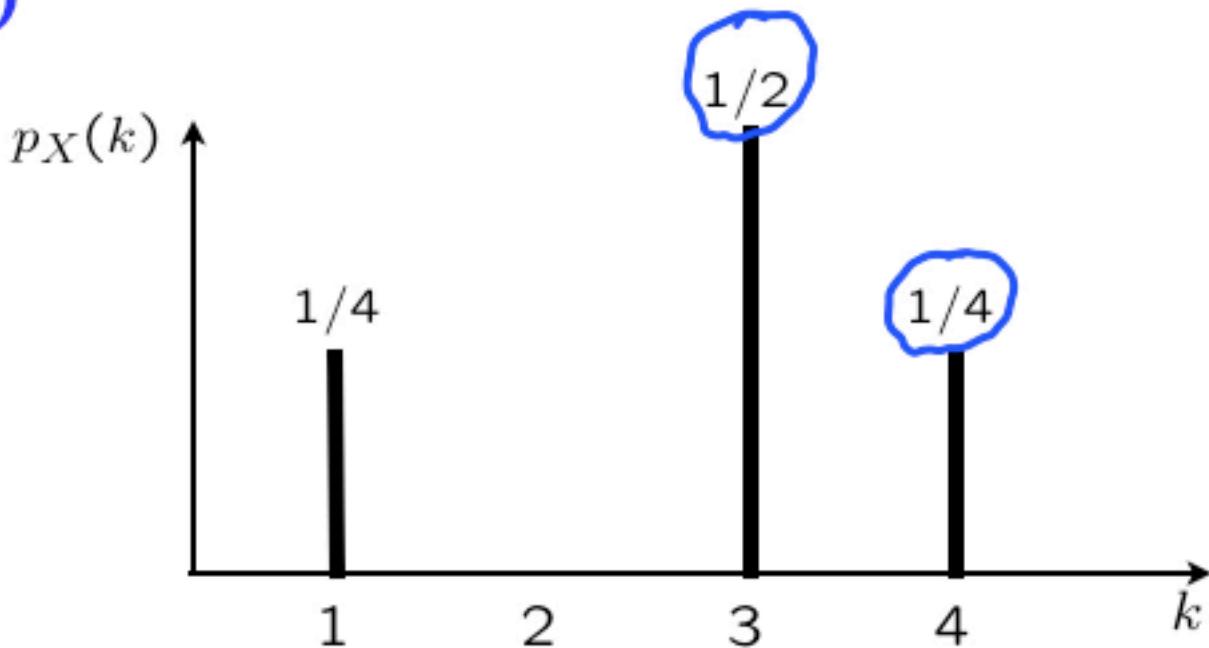


Cumulative distribution function (CDF)

CDF definition: $F_X(x) = P(X \leq x)$

- Discrete random variables:

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



General CDF properties

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



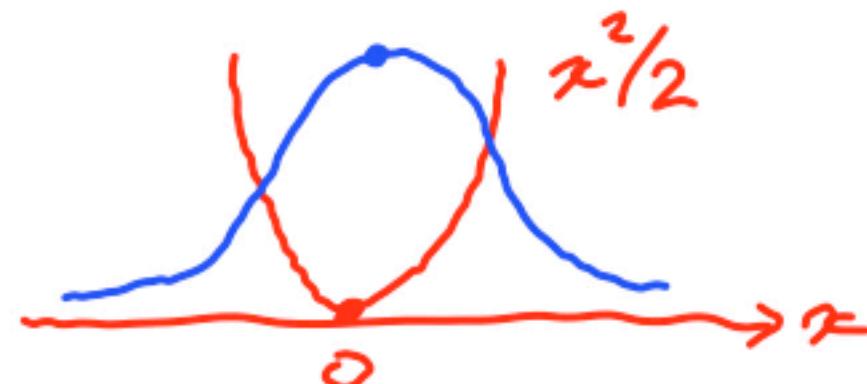
- Non-decreasing If $y \geq x \Rightarrow F_X(y) \geq F_X(x)$
- $F_X(x)$ tends to 1, as $x \rightarrow \infty$
- $F_X(x)$ tends to 0, as $x \rightarrow -\infty$

Normal (Gaussian) random variables

- Important in the theory of probability
 - Central limit theorem
- Prevalent in applications
 - Convenient analytical properties
 - Model of noise consisting of many, small independent noise terms

Standard normal (Gaussian) random variables

- Standard normal $N(0, 1)$: $f_X(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}$



calculus:

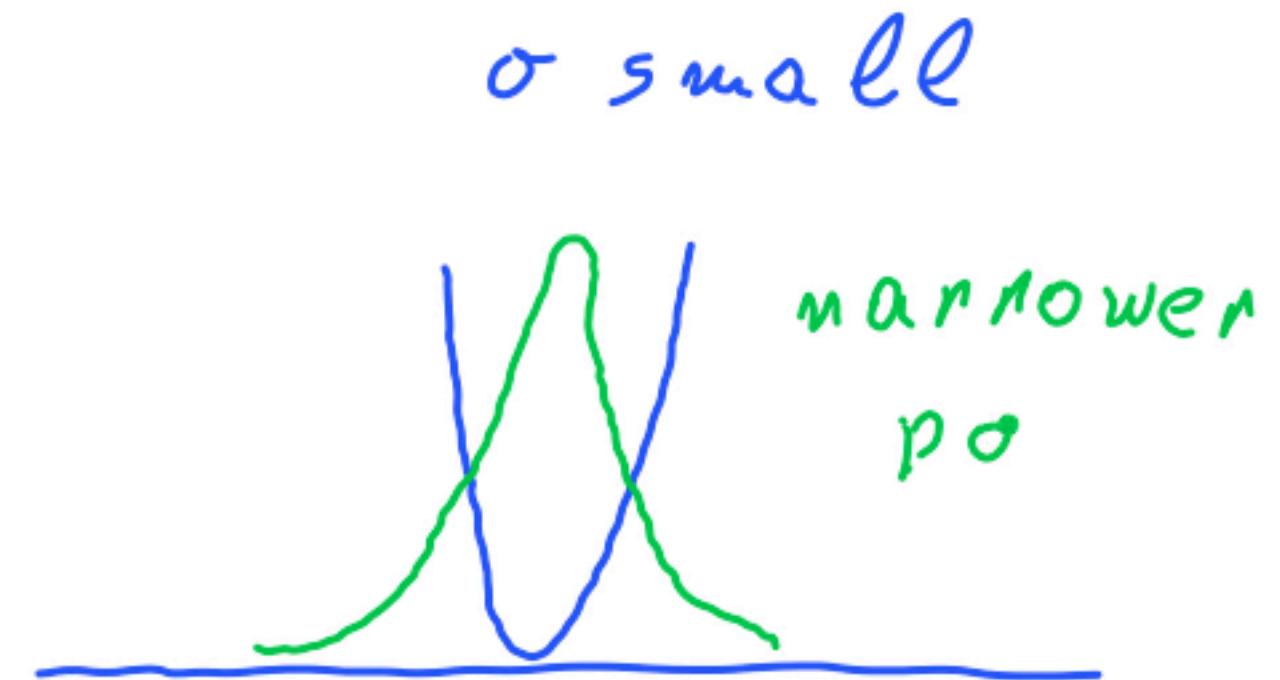
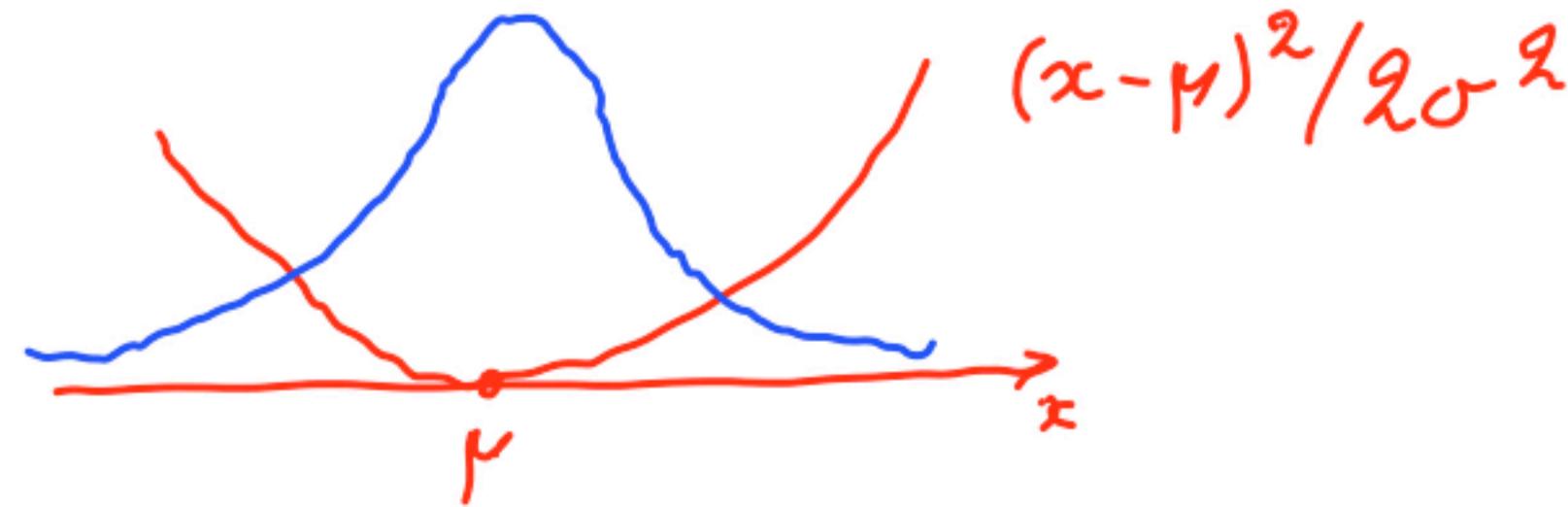
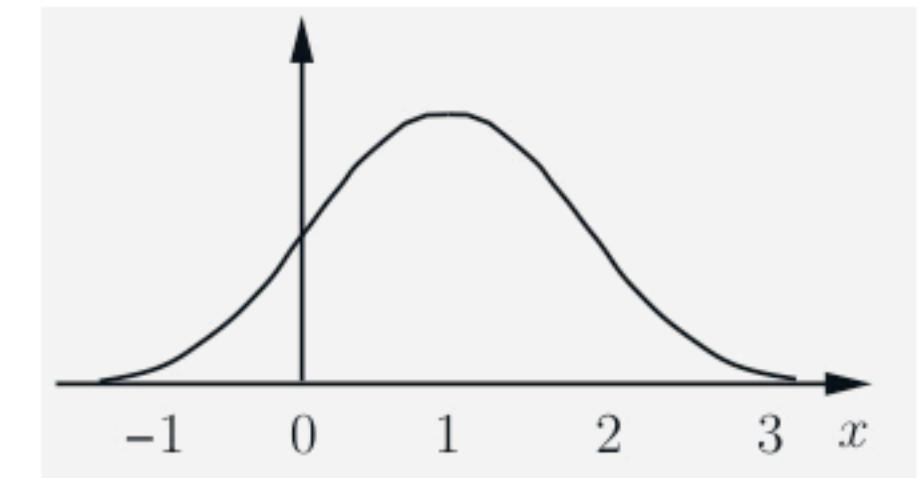
$$\int_{-\infty}^{\infty} e^{-x^2/2} dx = \sqrt{2\pi}$$

- $E[X] = 0$

- $\text{var}(X) = 1$ integrate by parts

General normal (Gaussian) random variables

- General normal $N(\mu, \sigma^2)$: $f_X(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-(x-\mu)^2/2\sigma^2}$
 $\sigma > 0$



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

Linear functions of a normal random variable

- Let $Y = aX + b \quad X \sim N(\mu, \sigma^2)$

$$E[Y] = a\mu + b$$

$$\text{Var}(Y) = a^2 \sigma^2$$

- Fact (will prove later in this course):

$$Y \sim N(a\mu + b, a^2\sigma^2)$$

- Special case: $a = 0$?

$$Y = b \quad \text{discrete}$$

\nearrow

$$N(b, 0)$$

Standard normal tables

- No closed form available for CDF

but have tables, for the standard normal

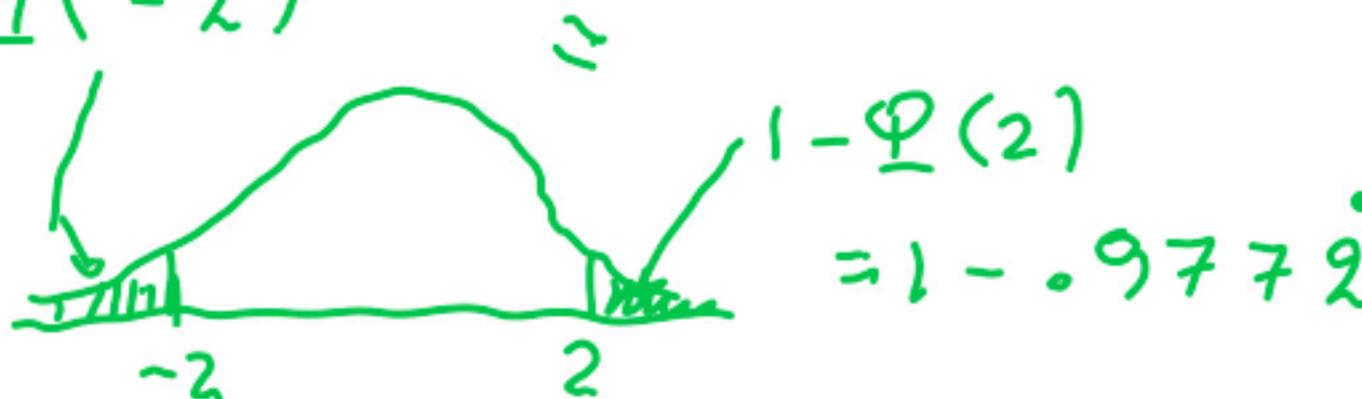
$$Y \sim N(0, 1)$$

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


$$\Phi(0) = P(Y \leq 0) = 0.5$$

$$\Phi(1.16) = 0.8770 \quad \Phi(2.9) = 0.9981$$

$$\Phi(-2)$$



	.00	.01	.02	.03	.04	.05	.06	.07	.08	.09
0.0	5000	5040	5080	5120	5160	5199	5239	5279	5319	5359
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	8186	8212	8238	8264	8289	8315	8340	8365	8389
1.0	8413	8438	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
2.1	9821	9826	9830	9834	9838	9842	9846	9850	9854	9857
2.2	9861	9864	9868	9871	9875	9878	9881	9884	9887	9890
2.3	9893	9896	9898	9901	9904	9906	9909	9911	9913	9916
2.4	9918	9920	9922	9925	9927	9929	9931	9932	9934	9936
2.5	9938	9940	9941	9943	9945	9946	9948	9949	9951	9952
2.6	9953	9955	9956	9957	9959	9960	9961	9962	9963	9964
2.7	9965	9966	9967	9968	9969	9970	9971	9972	9973	9974
2.8	9974	9975	9976	9977	9977	9978	9979	9979	9980	9981
2.9	9981	9982	9982	9983	9984	9984	9985	9985	9986	9986

Standardizing a random variable

- Let X have mean μ and variance $\sigma^2 > 0$

- Let $Y = \frac{X - \mu}{\sigma}$ $E[Y] = 0$ $\text{Var}(Y) = \frac{1}{\sigma^2} \text{Var}(x) = 1$

$$x = \mu + \sigma Y$$

- If also X is normal, then: $Y \sim N(0, 1)$

Calculating normal probabilities

- Express an event of interest in terms of standard normal

$$X \sim N(6, 4) \quad \sigma = 2$$

st. normalized

$$\frac{2 - 6}{2} \leq \frac{X - 6}{2} \leq \frac{8 - 6}{2}$$

$$P(2 \leq X \leq 8) = P(-2 \leq Y \leq 1)$$

$$\begin{aligned}
 &= P(Y \leq 1) - P(Y \leq -2) \\
 &\quad \text{Graph: A bell curve with vertical grid lines. The mean is at 0. The region from -2 to 1 is shaded green. The x-axis is labeled -2, 0, 1. The y-axis is unlabeled.} \\
 &= P(Y \leq 1) - (1 - P(Y \leq 2))
 \end{aligned}$$

	.00	.01	.02	.03	.04	.05	.06	.07	.08	.09
0.0	.5000	.5040	.5080	.5120	.5160	.5199	.5239	.5279	.5319	.5359
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	.8186	.8212	.8238	.8264	.8289	.8315	.8340	.8365	.8389
1.0	.8413	.8438	.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
2.1	.9821	.9826	.9830	.9834	.9838	.9842	.9846	.9850	.9854	.9857
2.2	.9861	.9864	.9868	.9871	.9875	.9878	.9881	.9884	.9887	.9890
2.3	.9893	.9896	.9898	.9901	.9904	.9906	.9909	.9911	.9913	.9916
2.4	.9918	.9920	.9922	.9925	.9927	.9929	.9931	.9932	.9934	.9936
2.5	.9938	.9940	.9941	.9943	.9945	.9946	.9948	.9949	.9951	.9952
2.6	.9953	.9955	.9956	.9957	.9959	.9960	.9961	.9962	.9963	.9964
2.7	.9965	.9966	.9967	.9968	.9969	.9970	.9971	.9972	.9973	.9974
2.8	.9974	.9975	.9976	.9977	.9977	.9978	.9979	.9979	.9980	.9981
2.9	.9981	.9982	.9982	.9983	.9984	.9984	.9985	.9985	.9986	.9986

LECTURE 9: Conditioning on an event; Multiple continuous r.v.'s

- Conditioning a r.v. on an event
 - Conditional PDF
 - Conditional expectation and the expected value rule
 - Exponential PDF: memorylessness
 - Total probability and expectation theorems
 - Mixed distributions
- Jointly continuous r.v.'s and joint PDFs
 - From the joints to the marginals
 - Uniform joint PDF example
 - The expected value rule and linearity of expectations
 - The joint CDF

Conditional PDF, given an event

$$P(A) > 0$$

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

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

$$p_{X|A}(x) = P(X = x | A)$$

$$\underline{f_{X|A}(x)} \cdot \delta \approx P(x \leq X \leq x + \delta | A)$$

$$P(X \in B) = \sum_{x \in B} p_X(x)$$

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

$$P(X \in B | A) = \sum_{x \in B} p_{X|A}(x)$$

$$P(X \in B | A) = \int_B f_{X|A}(x) dx$$

Def

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

$$\int \limits_{\bullet} f_{X|A}(x) dx = 1$$

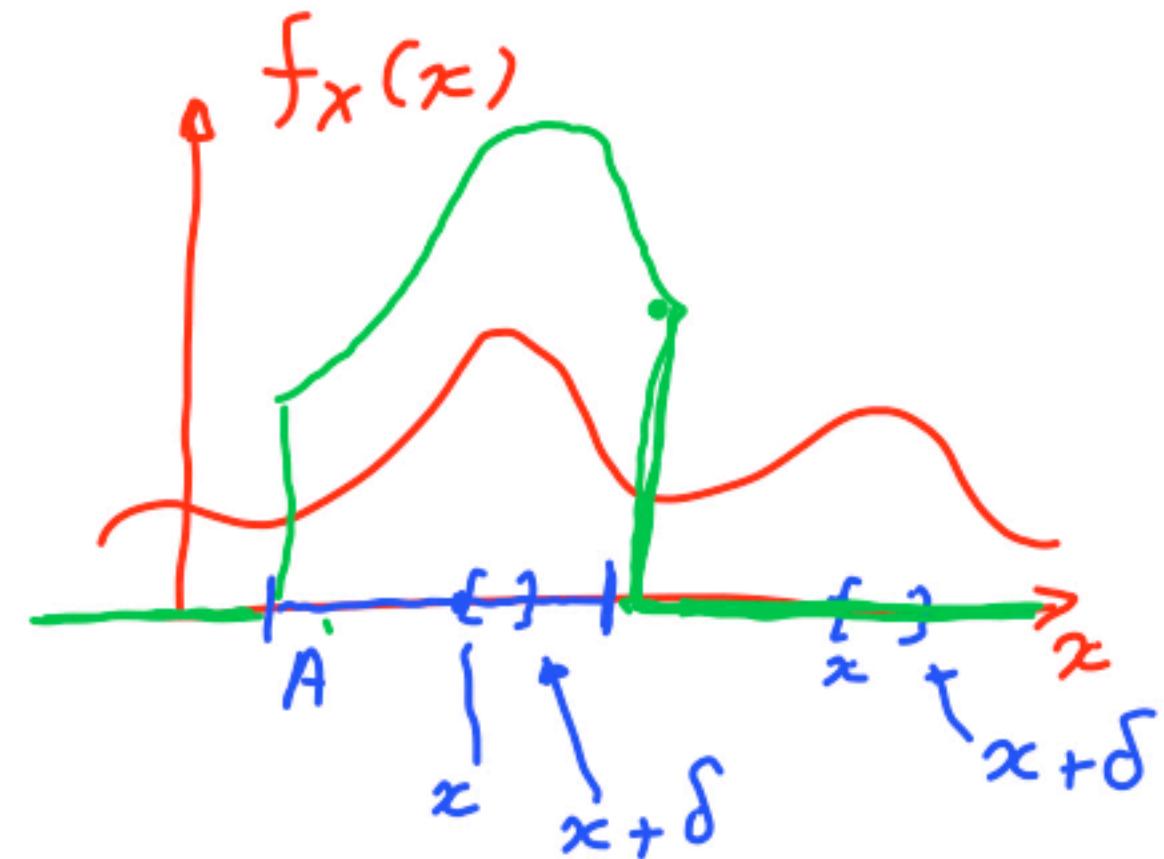
Conditional PDF of X , given that $\underline{X \in A}$

$$\mathbf{P}(x \leq X \leq x + \delta \mid X \in A) \approx f_{X|X \in A}(x) \cdot \cancel{\delta}$$

$$= \frac{\mathbf{P}(x \leq X \leq x + \delta, X \in A)}{\mathbf{P}(A)}$$

$$= \frac{\mathbf{P}(x \leq X \leq x + \delta)}{\mathbf{P}(A)} \approx \frac{f_X(x) \cancel{\delta}}{\mathbf{P}(A)}$$

$$f_{X|X \in A}(x) = \begin{cases} 0, & \text{if } x \notin A \\ \frac{f_X(x)}{\mathbf{P}(A)}, & \text{if } x \in A \end{cases}$$



Conditional expectation of X , given an event

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

$$\mathbb{E}[X] = \int x f_X(x) dx$$

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

$$\mathbb{E}[X | A] = \int x f_{X|A}(x) dx \quad \text{Def}$$

Expected value rule:

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

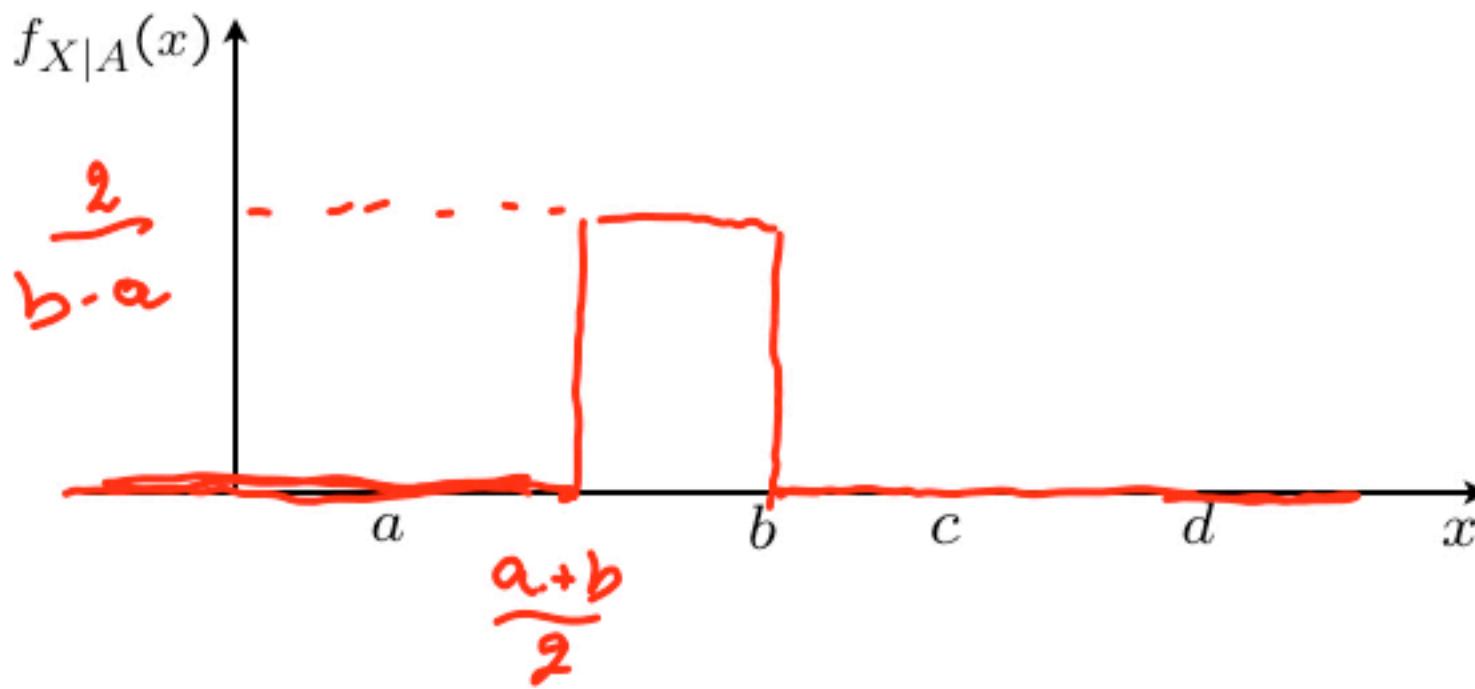
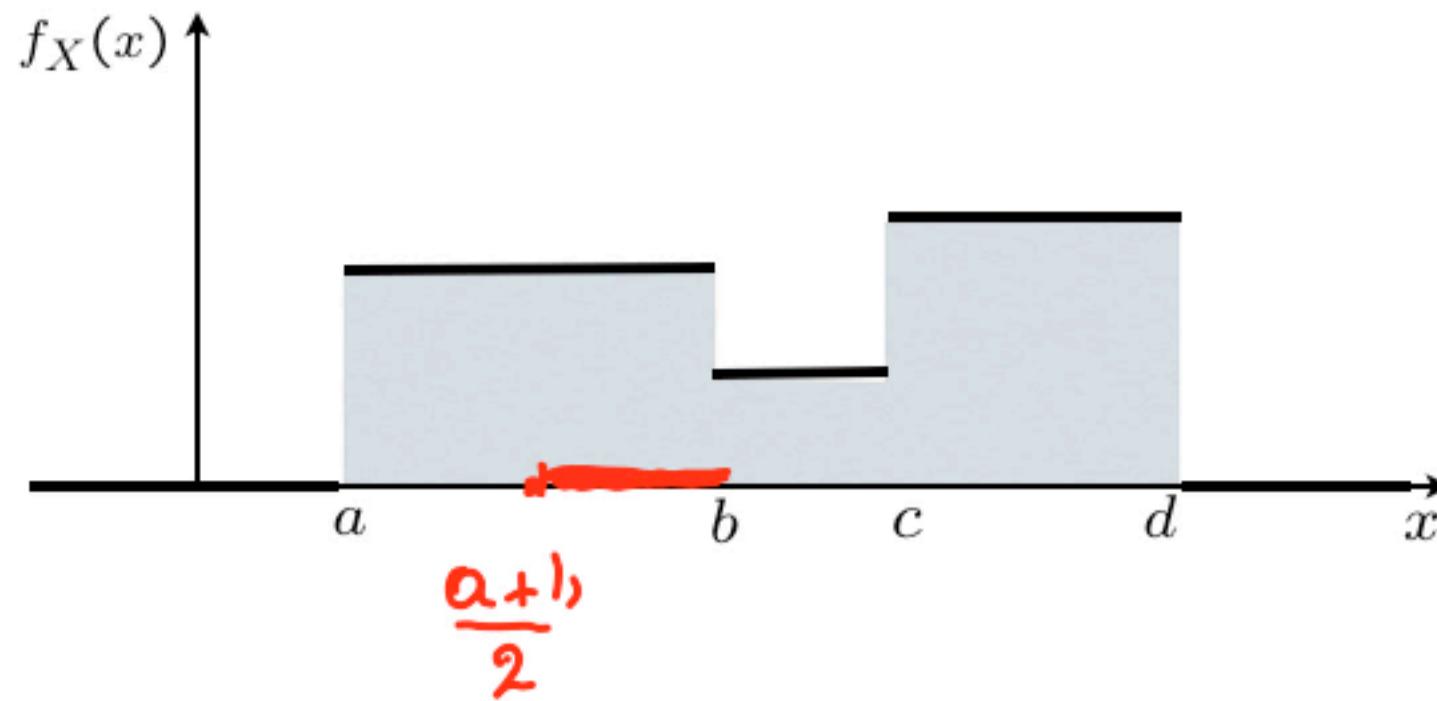
$$\mathbb{E}[g(X)] = \int g(x) f_X(x) dx$$

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

$$\mathbb{E}[g(X) | A] = \int g(x) f_{X|A}(x) dx$$

Example

$$A : \frac{a+b}{2} \leq X \leq b$$



$$\mathbb{E}[X | A] = \frac{1}{2} \cdot \frac{a+b}{2} + \frac{1}{2} b$$

$$= \frac{1}{4} a + \frac{3}{4} b$$

$$\mathbb{E}[X^2 | A] =$$

$$\left. \frac{2}{b-a} \cdot x^2 dx \right|_{\frac{a+b}{2}}^b$$

Memorylessness of the exponential PDF

- Do you prefer a used or a new “exponential” light bulb? **Probabilistically identical!**

- Bulb lifetime T : $\text{exponential}(\lambda)$

$$P(T > x) = e^{-\lambda x}, \text{ for } x \geq 0$$

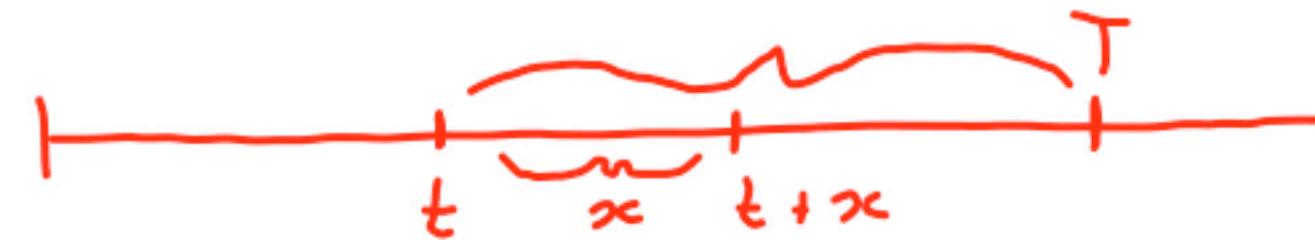
— we are told that $T > t$

— r.v. X : remaining lifetime $= T - t$

$$P(X > x | T > t) = e^{-\lambda x}, \text{ for } x \geq 0$$

$$= \frac{P(T-t > x, T > t)}{P(T > t)} = \frac{P(T > t+x, T > t)}{P(T > t)} = \frac{P(T > t+x)}{P(T > t)}$$

$$= \frac{e^{-\lambda(t+x)}}{e^{-\lambda t}} = e^{-\lambda x}$$



Memorylessness of the exponential PDF

$$f_T(x) = \lambda e^{-\lambda x}, \quad \text{for } x \geq 0$$

$$\mathbf{P}(0 \leq T \leq \delta) \approx f_T(0) \cdot \delta = \lambda \delta$$

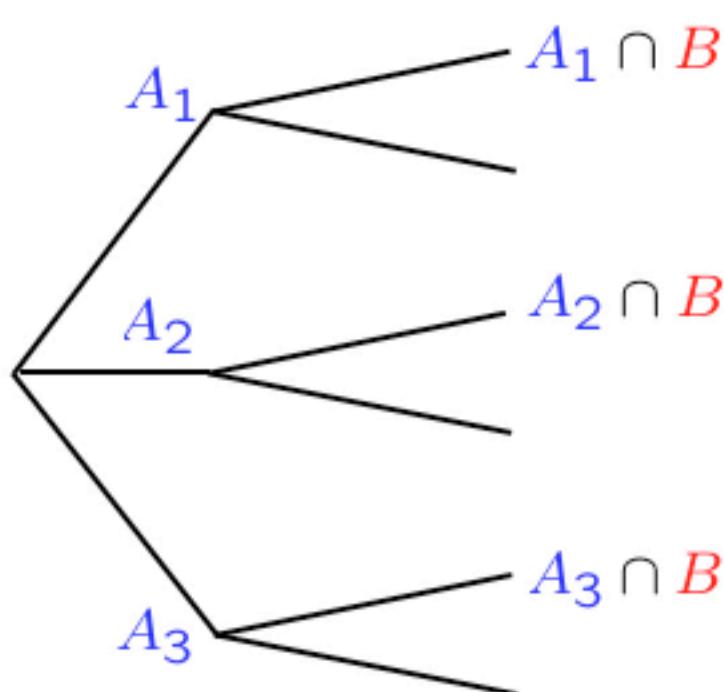
(approximation)

$$\mathbf{P}(t \leq T \leq t + \delta \mid T > t) = \approx \lambda \delta$$



similar to an independent coin flip,
every δ time steps,
with $\mathbf{P}(\text{success}) \approx \lambda \delta$

Total probability and expectation theorems

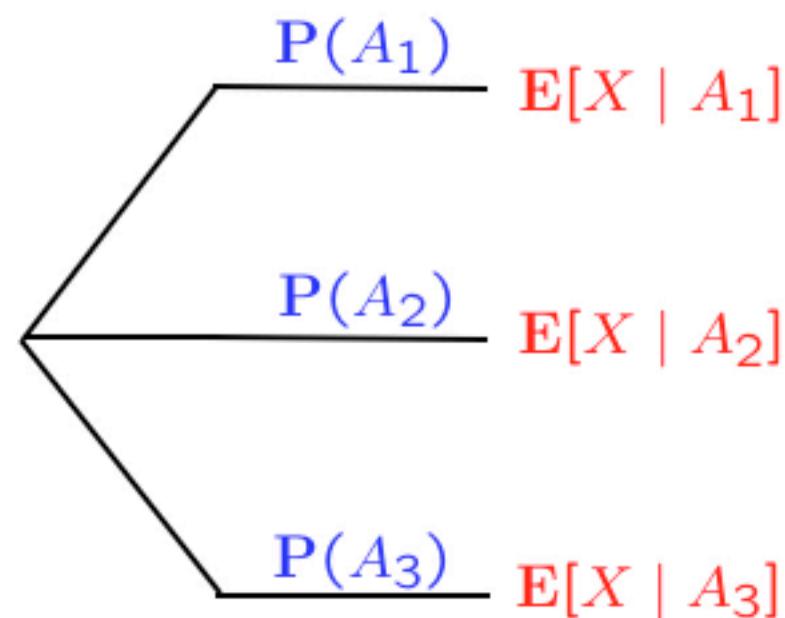


$$P(B) = P(A_1)P(B|A_1) + \cdots + P(A_n)P(B|A_n)$$

$$p_X(x) = P(A_1)p_{X|A_1}(x) + \cdots + P(A_n)p_{X|A_n}(x)$$

$$\begin{aligned} F_X(x) &= P(X \leq x) = P(A_1)P(X \leq x | A_1) + \cdots \\ &= P(A_1)F_{X|A_1}(x) + \cdots \end{aligned}$$

$$f_X(x) = P(A_1)f_{X|A_1}(x) + \cdots + P(A_n)f_{X|A_n}(x)$$



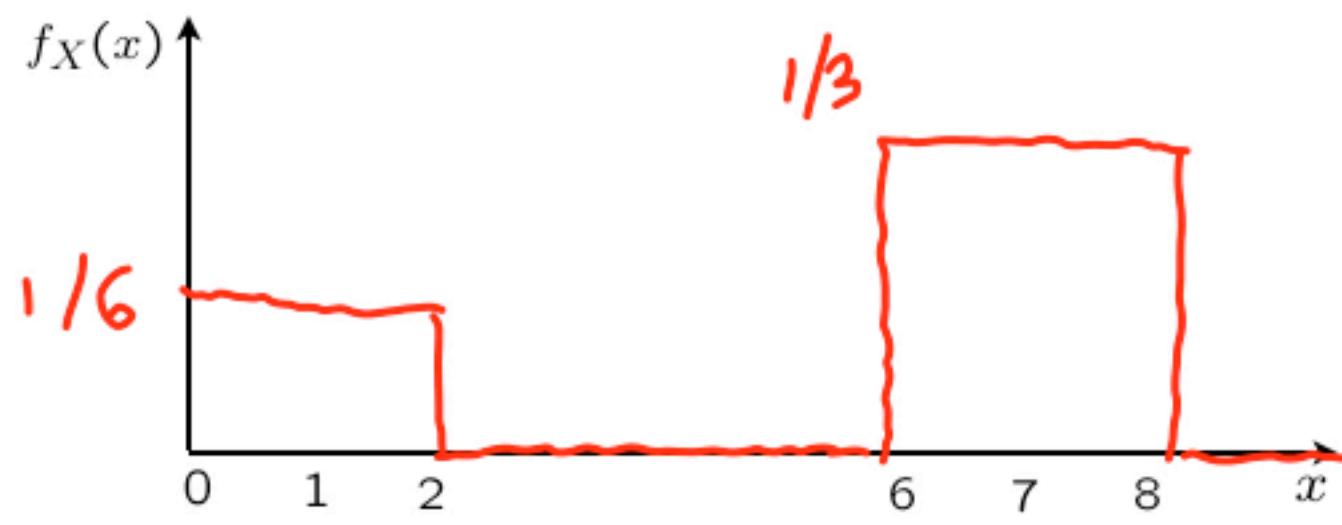
$$\int x f_X(x) dx = P(A_1) \int x f_{X|A_1}(x) dx + \cdots$$

$$E[X] = P(A_1)E[X | A_1] + \cdots + P(A_n)E[X | A_n]$$

Example

- Bill goes to the supermarket shortly, with probability $1/3$, at a time uniformly distributed between 0 and 2 hours from now; or with probability $2/3$, later in the day at a time uniformly distributed between 6 and 8 hours from now

$$\Pr(A_1) = \frac{1}{3} \quad f_{X|A_1} \sim \text{unif}[0, 2] \quad \Pr(A_2) = \frac{2}{3} \quad f_{X|A_2} \sim U[6, 8]$$



$$f_X(x) = \Pr(A_1)f_{X|A_1}(x) + \dots + \Pr(A_n)f_{X|A_n}(x)$$

$$\bullet \quad E[X] = \Pr(A_1)E[X | A_1] + \dots + \Pr(A_n)E[X | A_n]$$

$$\frac{1}{3} \cdot 1 + \frac{2}{3} \cdot 7$$

Mixed distributions

$X = \begin{cases} \text{uniform on } [0, 2], & \text{with probability } 1/2 \\ 1, & \text{with probability } 1/2 \end{cases}$	Is X discrete? No
Y discrete Z continuous	$X = \begin{cases} Y, & \text{with probability } p \\ Z, & \text{with probability } 1 - p \end{cases}$ Is X continuous? No
	$P(X=1) = 1/2$ X is mixed

$$\begin{aligned} F_X(x) &= P(Y \leq x) + (1-p) P(Z \leq x) \\ &= p F_Y(x) + (1-p) F_Z(x) \end{aligned}$$

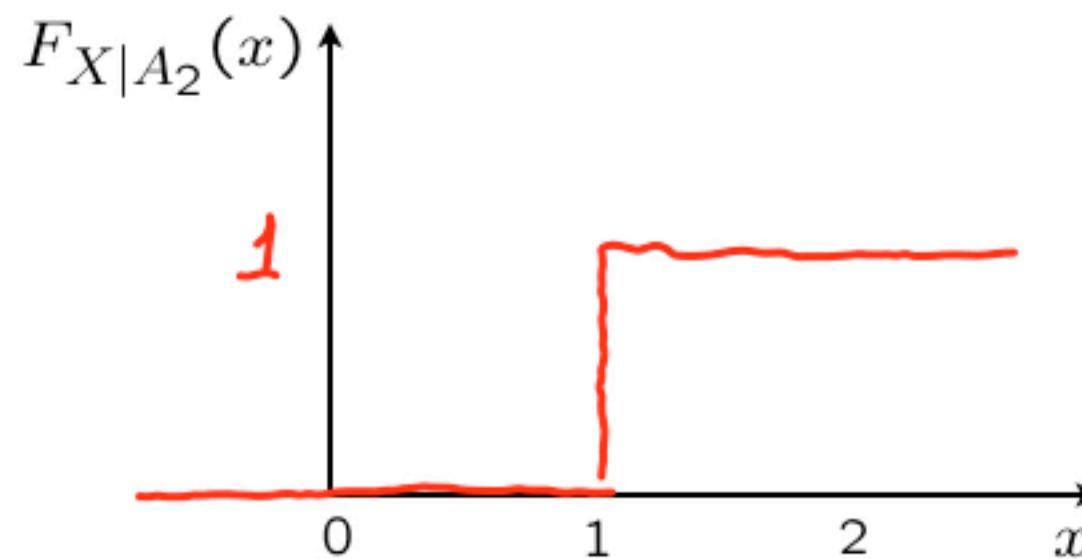
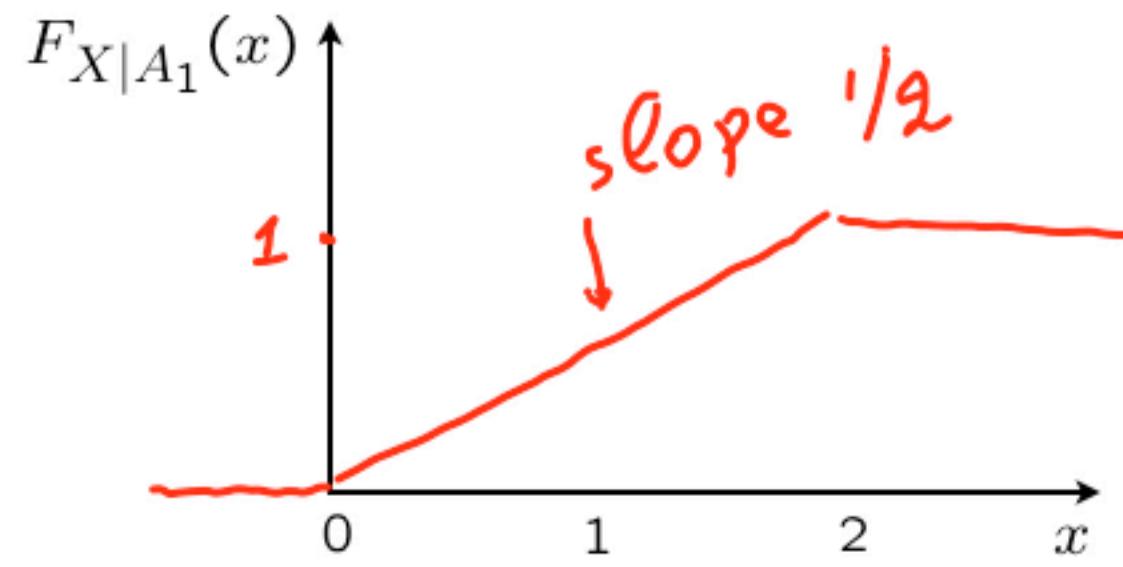
$$E[X] = p E[Y] + (1-p) E[Z]$$

Mixed distributions

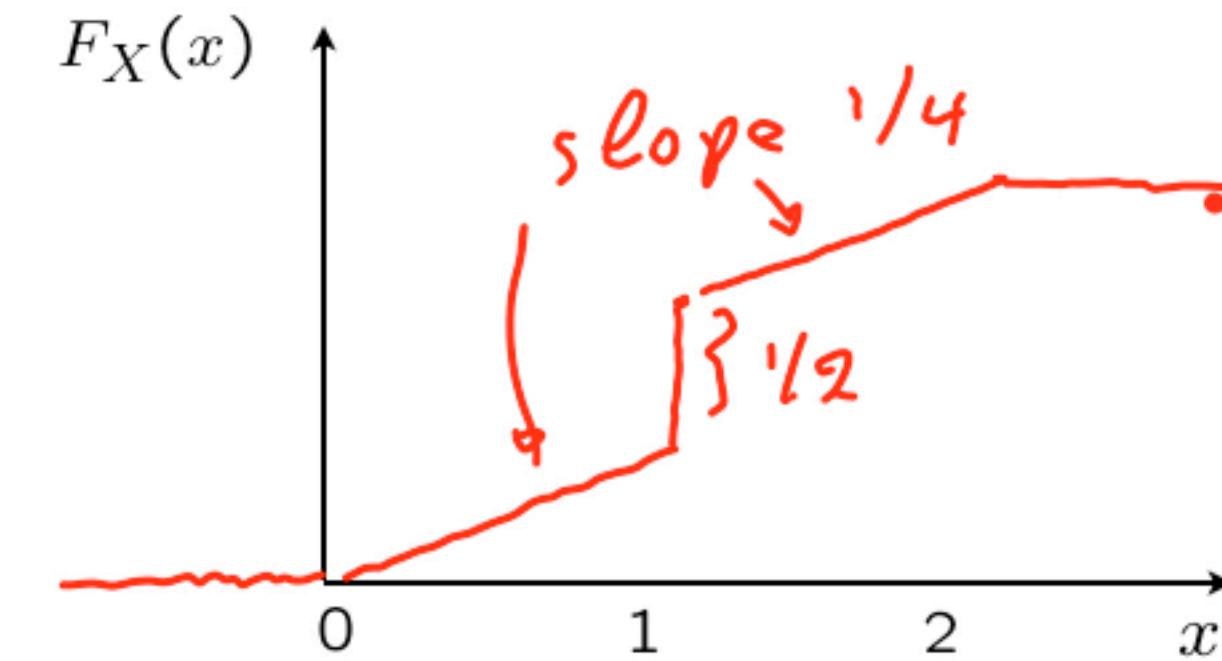
$X = \begin{cases} \text{uniform on } [0, 2], & \text{with probability } 1/2 \\ 1, & \text{with probability } 1/2 \end{cases}$

A_1

A_2



$$F_X(x) = P(A_1)F_{X|A_1}(x) + P(A_2)F_{X|A_2}(x)$$



Jointly continuous r.v.'s and joint PDFs

$p_X(x)$	$f_X(x)$
$p_{X,Y}(x,y)$	$f_{X,Y}(x,y)$

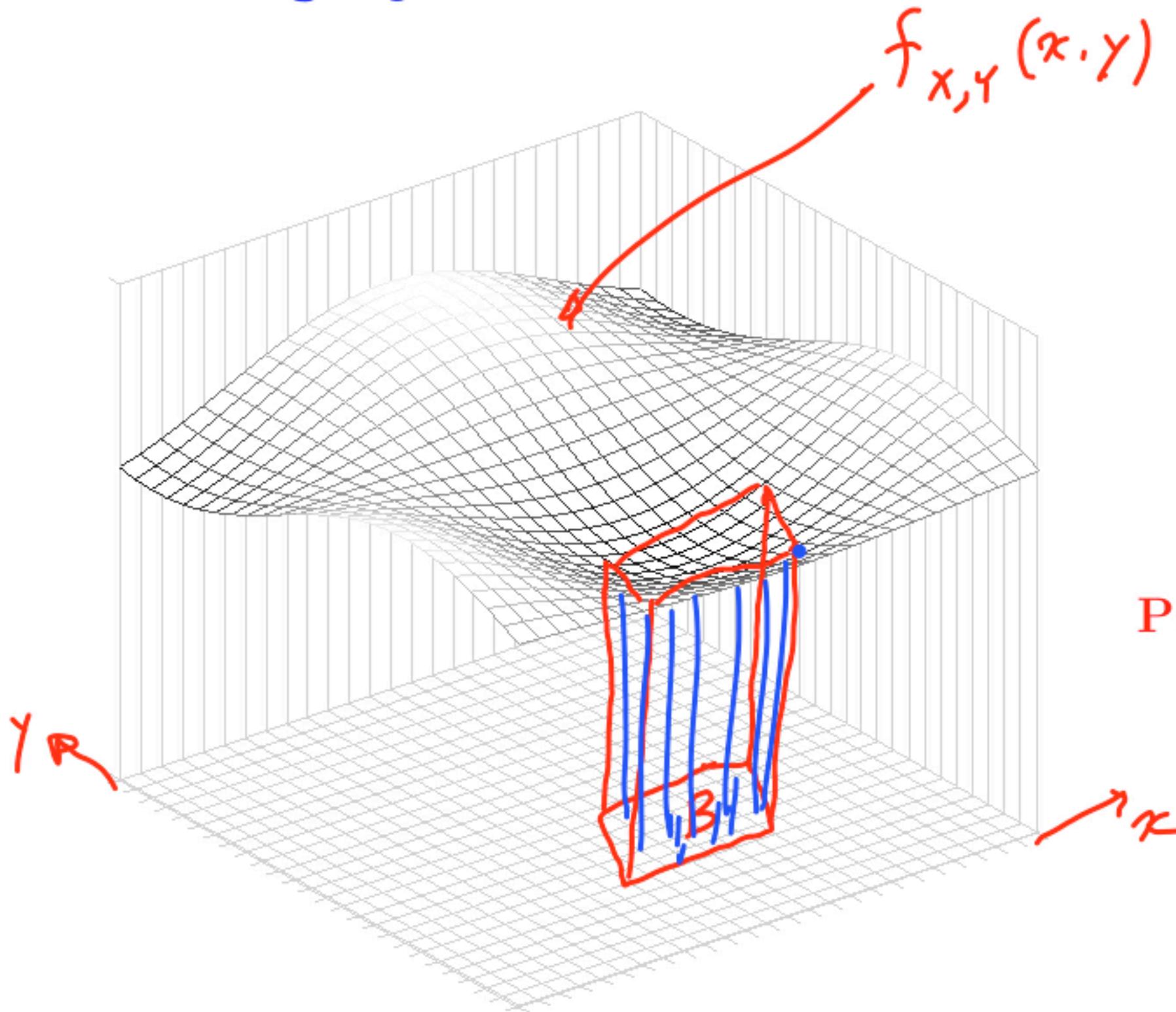
$$p_{X,Y}(x,y) = \mathbf{P}(X = x \text{ and } Y = y) \geq 0 \quad f_{X,Y}(x,y) \geq 0$$

$$\mathbf{P}((X,Y) \in B) = \sum_{(x,y) \in B} p_{X,Y}(x,y) \quad \mathbf{P}((X,Y) \in B) = \iint_{(x,y) \in B} f_{X,Y}(x,y) dx dy \quad •$$

$$\sum_x \sum_y p_{X,Y}(x,y) = 1 \quad \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f_{X,Y}(x,y) dx dy = 1$$

Definition: Two random variables are **jointly continuous** if they can be described by a joint PDF

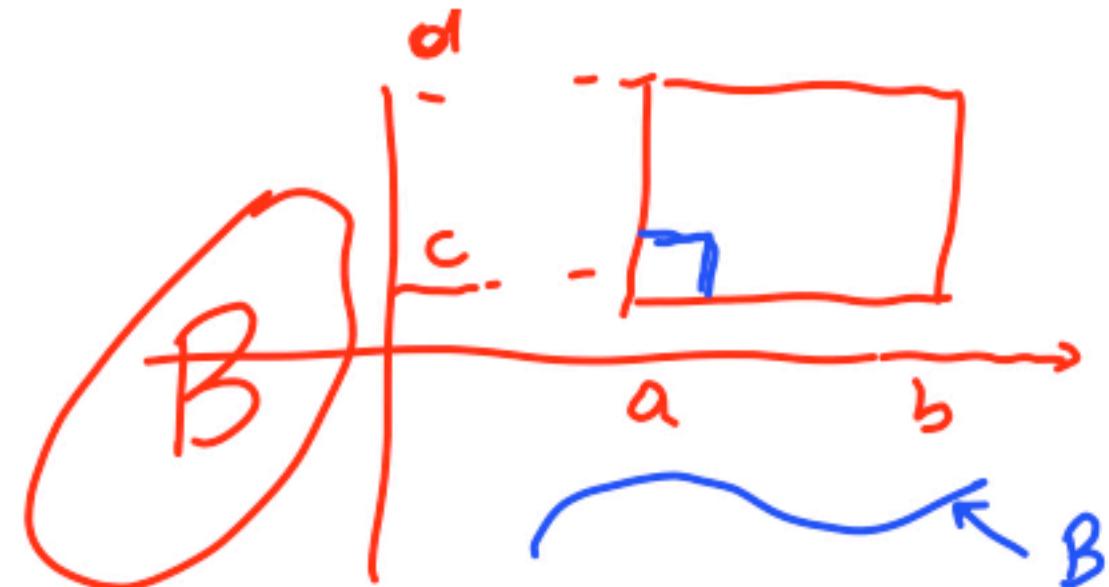
Visualizing a joint PDF



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

On joint PDFs

$$\mathbf{P}((X, Y) \in B) = \iint_{(x,y) \in B} f_{X,Y}(x, y) dx dy$$



$$\mathbf{P}(a \leq X \leq b, c \leq Y \leq d) = \int_c^d \int_a^b f_{X,Y}(x, y) dx dy$$

$$\mathbf{P}(a \leq X \leq a + \delta, c \leq Y \leq c + \delta) \approx f_{X,Y}(a, c) \cdot \delta^2$$

$$Y = X$$

$f_{X,Y}(x, y)$: probability per unit area

$$\text{area}(B) = 0 \Rightarrow \mathbf{P}((X, Y) \in B) = 0$$



From the joint to the marginals

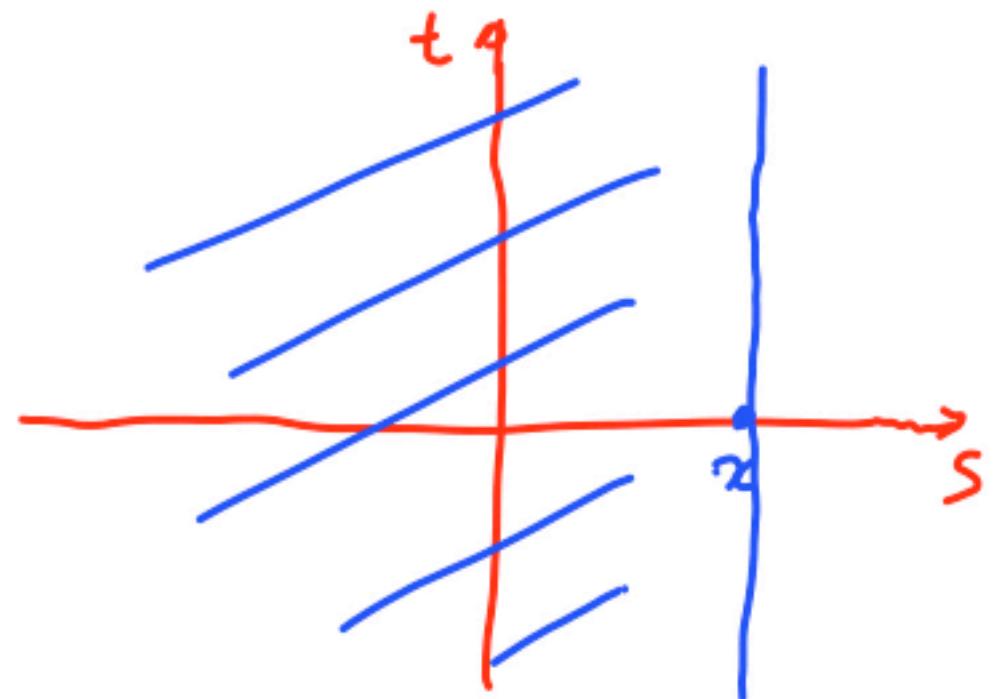
$$p_X(x) = \sum_y p_{X,Y}(x,y)$$

$$f_X(x) = \int f_{X,Y}(x,y) dy$$

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

$$f_Y(y) = \int f_{X,Y}(x,y) dx$$

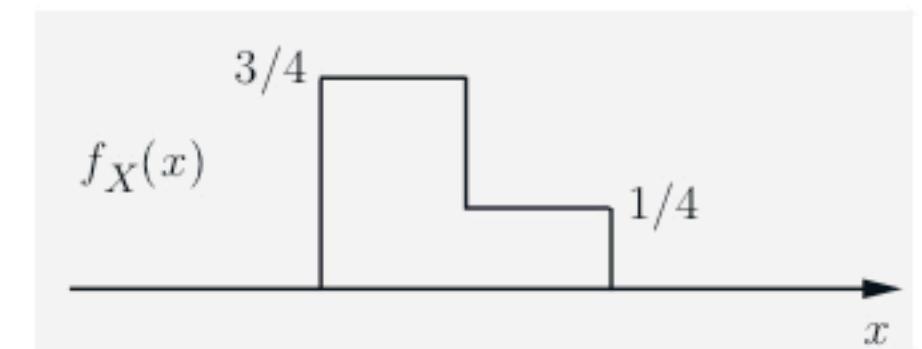
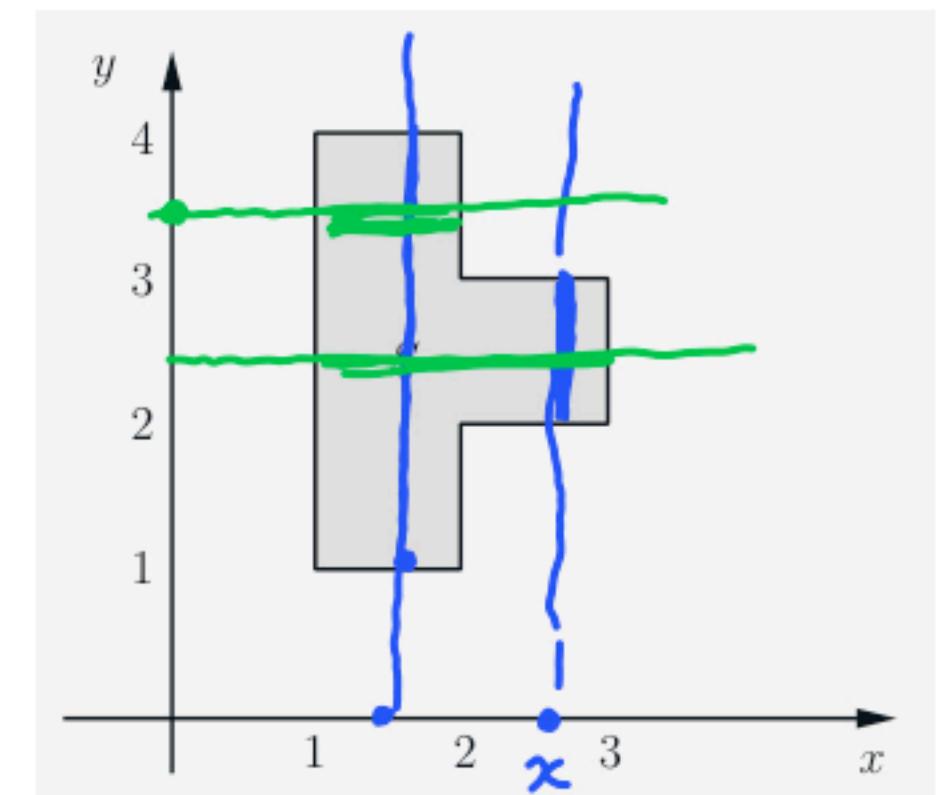
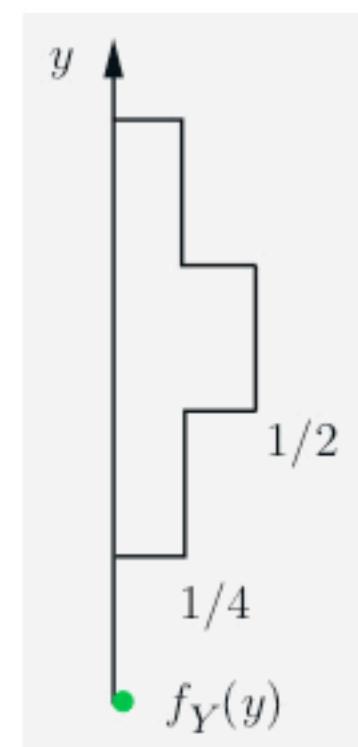
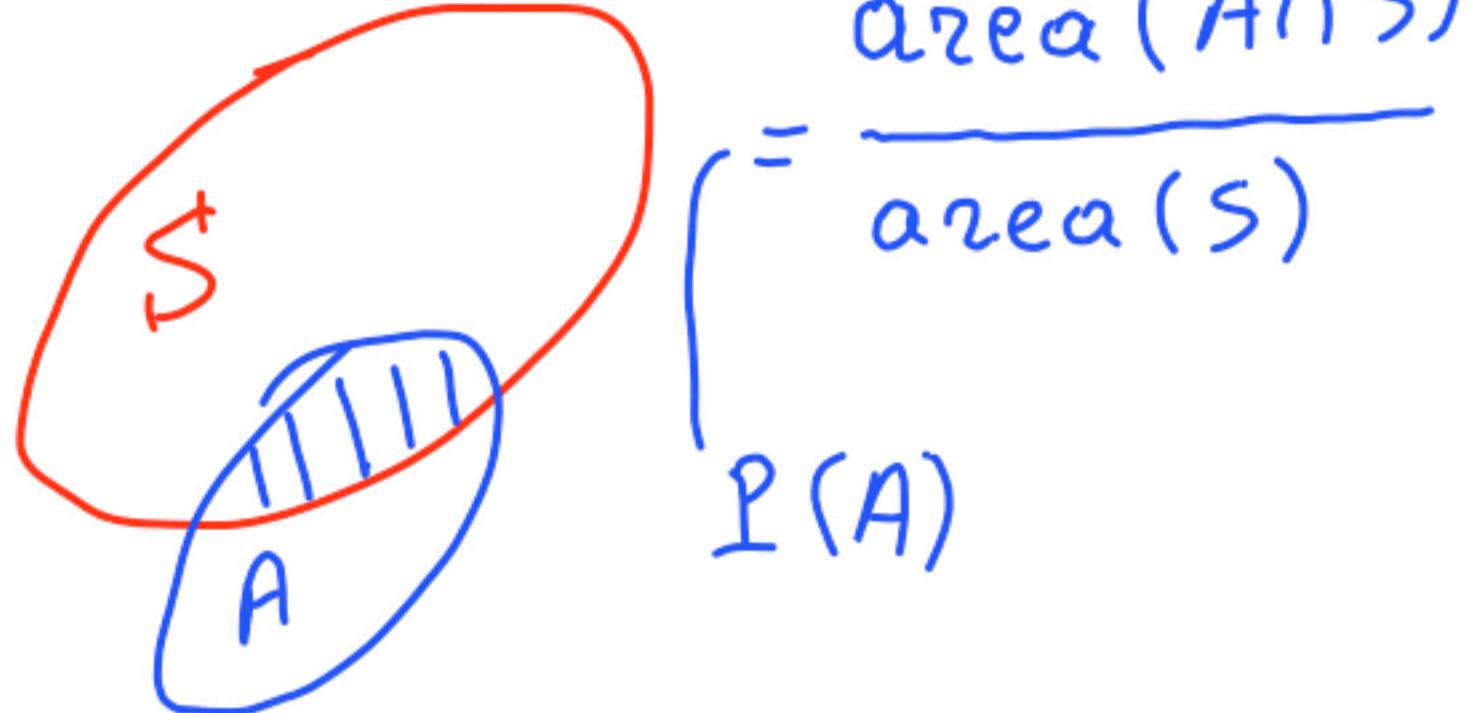
$$F_x(x) = P(X \leq x) = \int_{-\infty}^x \left[\int_{-\infty}^{\infty} f_{X,Y}(s,t) dt \right] ds$$



$$f_x(x) = \frac{d}{dx} F_x(x) = []$$

Uniform joint PDF on a set S

$$f_{X,Y}(x,y) = \begin{cases} \frac{1}{\text{area of } S}, & \text{if } (x,y) \in S, \\ 0, & \text{otherwise.} \end{cases}$$



$$f_{X,Y} = \frac{1}{4}$$

More than two random variables

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

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

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

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

$$p_{X,Y}(x, y) = \sum_z p_{X,Y,Z}(x, y, z)$$

Functions of multiple random variables

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

Expected value rule:

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

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

Linearity of expectations

$$\mathbb{E}[aX + b] = a\mathbb{E}[X] + b$$

$$\mathbb{E}[X + Y] = \mathbb{E}[X] + \mathbb{E}[Y]$$

$$\mathbb{E}[X_1 + \dots + X_n] = \mathbb{E}[X_1] + \dots + \mathbb{E}[X_n]$$

The joint CDF

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

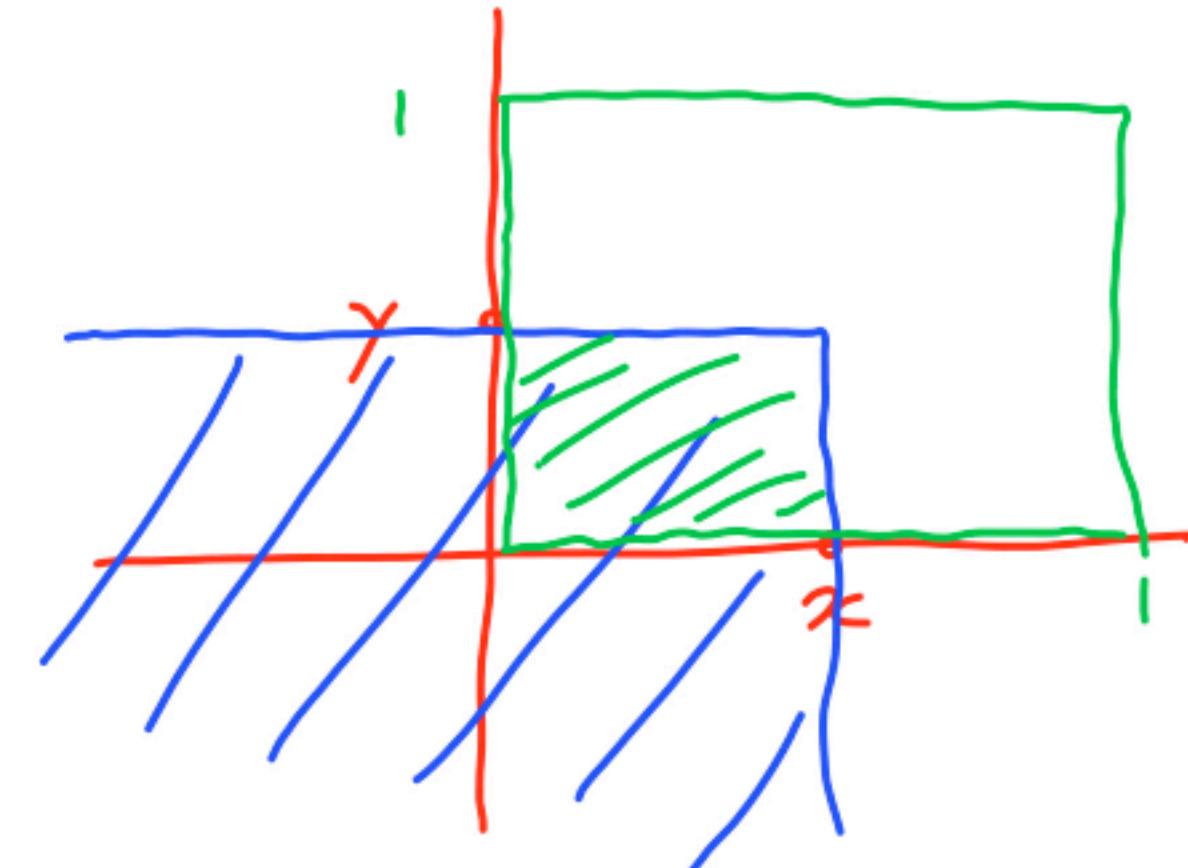
$$f_X(x) = \frac{dF_X}{dx}(x)$$

$$F_{X,Y}(x, y) = \mathbf{P}(X \leq x, Y \leq y) = \int_{-\infty}^y \left[\int_{-\infty}^x f_{X,Y}(s, t) ds \right] dt$$

$$f_{X,Y}(x, y) = \frac{\partial^2 F_{X,Y}}{\partial x \partial y}(x, y)$$

$$F_{X,Y}(x, y) = xy$$

$$f_{X,Y}(x, y) = 1$$



LECTURE 10: Conditioning on a random variable; Independence; Bayes' rule

- Conditioning X on Y
 - Total probability theorem
 - Total expectation theorem
- Independence
 - independent normals
- A comprehensive example
- Four variants of the Bayes rule

Conditional PDFs, given another r.v.

$$p_{X|Y}(x | y) = \mathbf{P}(X = x | Y = y) = \frac{p_{X,Y}(x, y)}{p_Y(y)}, \quad \text{if } p_Y(y) > 0$$

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

$$\mathbf{P}(x \leq X \leq x + \delta | A) \approx f_{X|A}(x) \cdot \delta, \quad \text{where } \mathbf{P}(A) > 0$$

$$\overbrace{Y=y}^{\gamma=y} \quad \downarrow \quad Y \approx y$$

$$\mathbf{P}(x \leq X \leq x + \delta | y \leq Y \leq y + \epsilon) \approx \frac{f_{x,y}(x, y) \delta}{f_y(y) \delta} = f_{x|y}(x | y) \delta$$

Definition: $\mathbf{P}(X \in A | Y = y) = \int_A f_{X|Y}(x | y) dx$

$p_{X,Y}(x, y)$	$f_{X,Y}(x, y)$
$p_{X A}(x)$	$f_{X A}(x)$
$p_{X Y}(x y)$	$f_{X Y}(x y)$

Comments on conditional PDFs

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

- $f_{X|Y}(x | y) \geq 0$

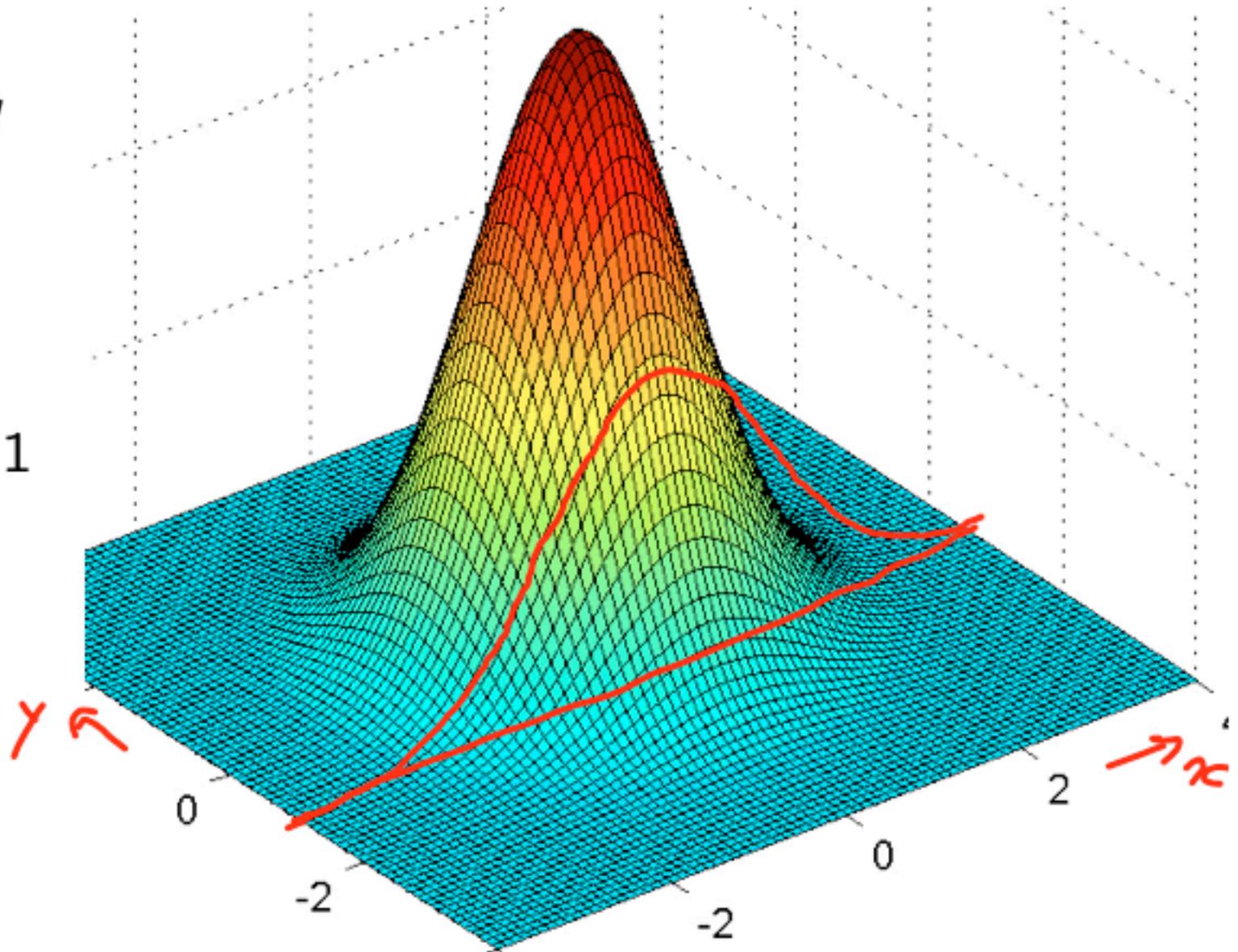
- Think of value of Y as fixed at some y
shape of $f_{X|Y}(\cdot | y)$: slice of the joint

- $\int_{-\infty}^{\infty} f_{X|Y}(x | y) dx = \frac{\int_{-\infty}^{\infty} f_{X,Y}(x, y) dx}{f_Y(y)} = 1$

- Multiplication rule:

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

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



Total probability and expectation theorems

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

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

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

- Expected value rule...

$$\mathbb{E}[g(X) | Y = y]$$

$$= \int_{-\infty}^{\infty} g(x) f_{X|Y}(x|y) dx$$

$$f_X(x) = \underbrace{\int_{-\infty}^{\infty} f_Y(y) f_{X|Y}(x|y) dy}_{f_{X,Y}(x,y)}$$

Thm.

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

Def.

$$\begin{aligned} \mathbb{E}[X] &= \int_{-\infty}^{\infty} f_Y(y) \mathbb{E}[X | Y = y] dy \\ &= \int_{-\infty}^{\infty} f_Y(y) \int_{-\infty}^{\infty} x f_{X|Y}(x|y) dx dy \end{aligned}$$

$$= \int_{-\infty}^{\infty} x \int_{-\infty}^{\infty} \cancel{f_Y(y)} f_{X|Y}(x|y) dy dx$$

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

Independence

$$p_{X,Y}(x,y) = p_X(x)p_Y(y), \quad \text{for all } x, y$$

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

$$f_{Y|X} = f_Y$$

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

- equivalent to: $f_{X|Y}(x|y) = f_X(x)$, for all y with $f_Y(y) > 0$ and all x

If X, Y are **independent**: $\mathbf{E}[XY] = \mathbf{E}[X]\mathbf{E}[Y]$

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

$g(X)$ and $h(Y)$ are also independent: $\mathbf{E}[g(X)h(Y)] = \mathbf{E}[g(X)] \cdot \mathbf{E}[h(Y)]$

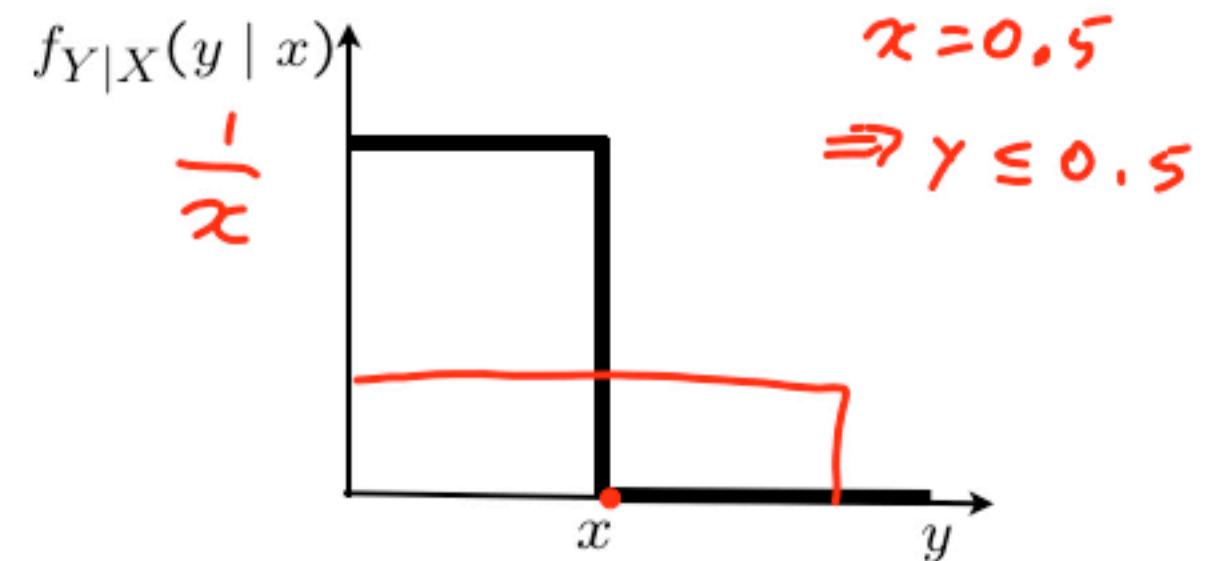
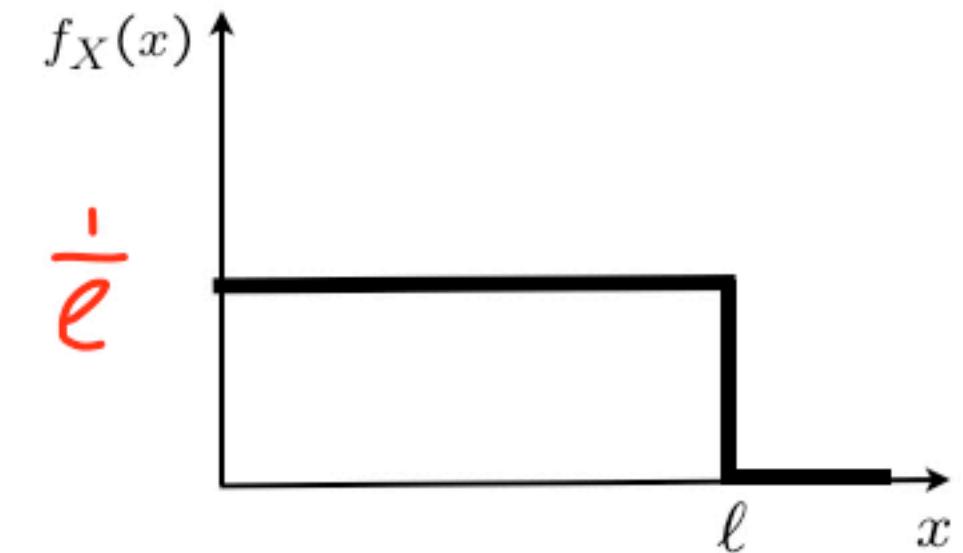
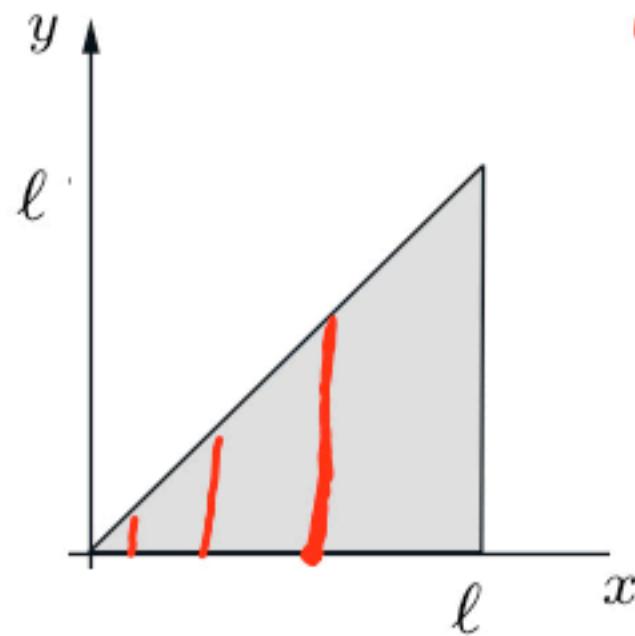
Stick-breaking example



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

$$f_{X,Y}(x,y) = f_X(x)f_{Y|X}(y|x) = \frac{1}{\ell x}$$

$$0 \leq y \leq x \leq \ell$$

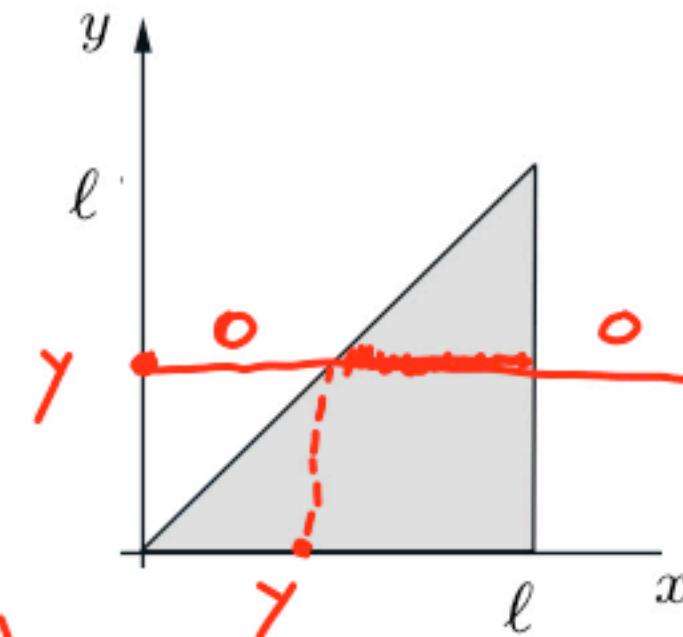


Stick-breaking example

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

$$f_Y(y) = \int_{y}^{\ell} f_{x,y}(x,y) dx = \int_y^{\ell} \frac{1}{\ell x} dx = \frac{1}{\ell} \log\left(\frac{\ell}{y}\right)$$

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



- Using total expectation theorem:

$$E[Y] = \int_0^{\ell} \frac{1}{\ell} E[Y|x=x] dx = \int_0^{\ell} \left(\frac{1}{\ell} \right) \frac{x}{2} dx = \frac{1}{2} E[X] = \frac{1}{2} \cdot \frac{\ell}{2} = \frac{\ell}{4}$$

$f_X(x)$



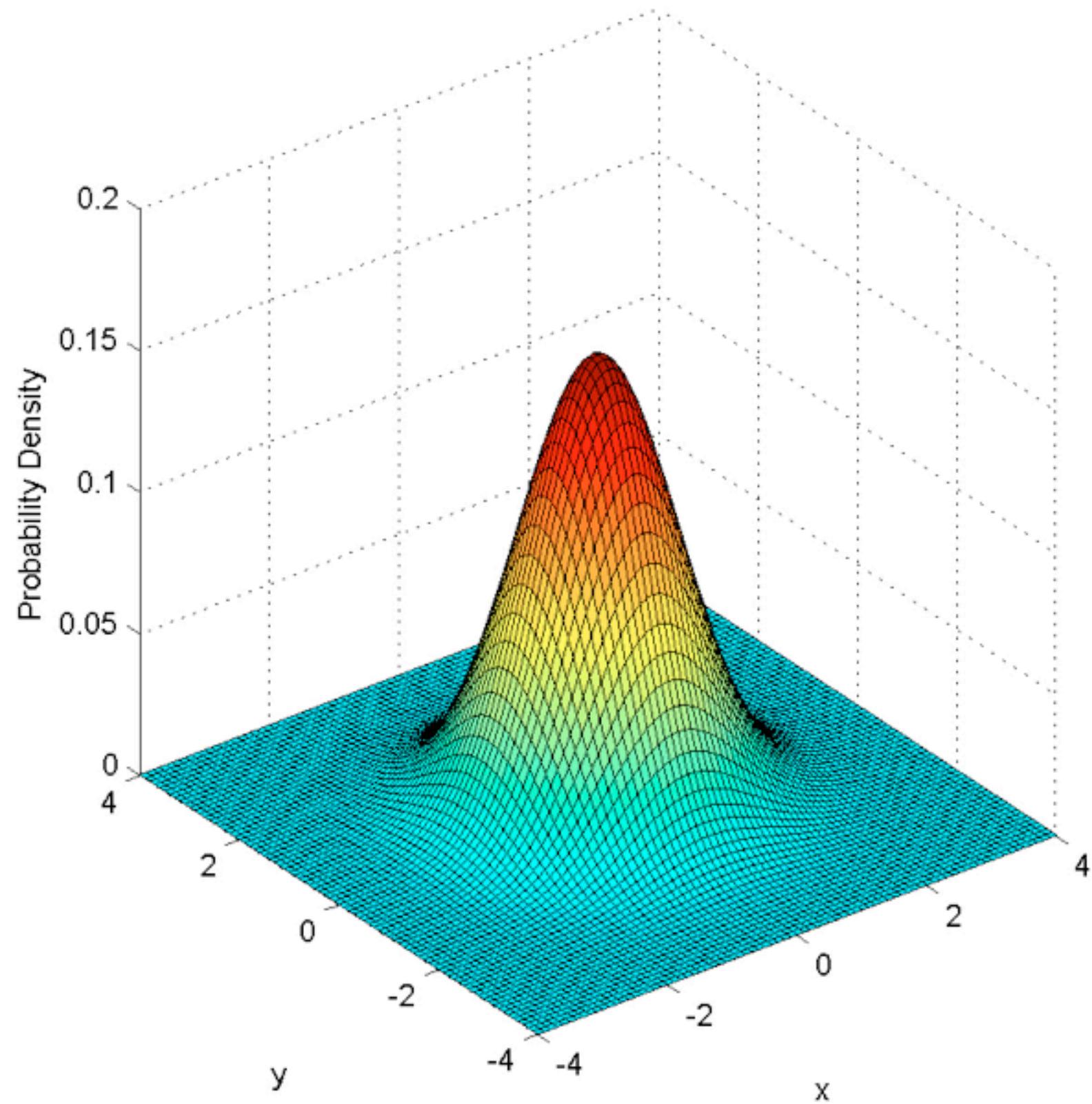
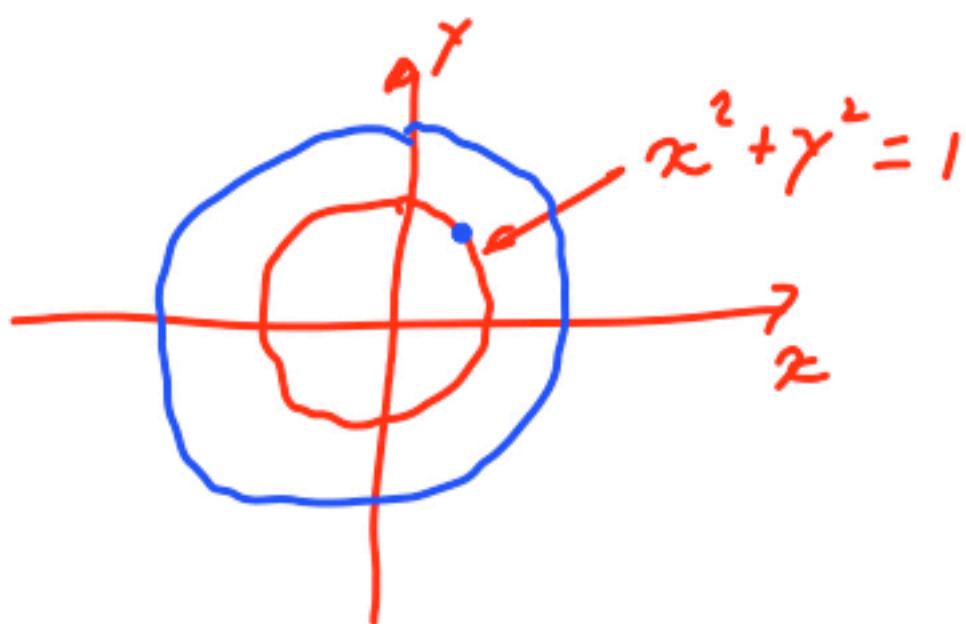
Independent standard normals

$\mu_X = \mu_Y = 0; \sigma_X^2 = \sigma_Y^2 = 1$

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

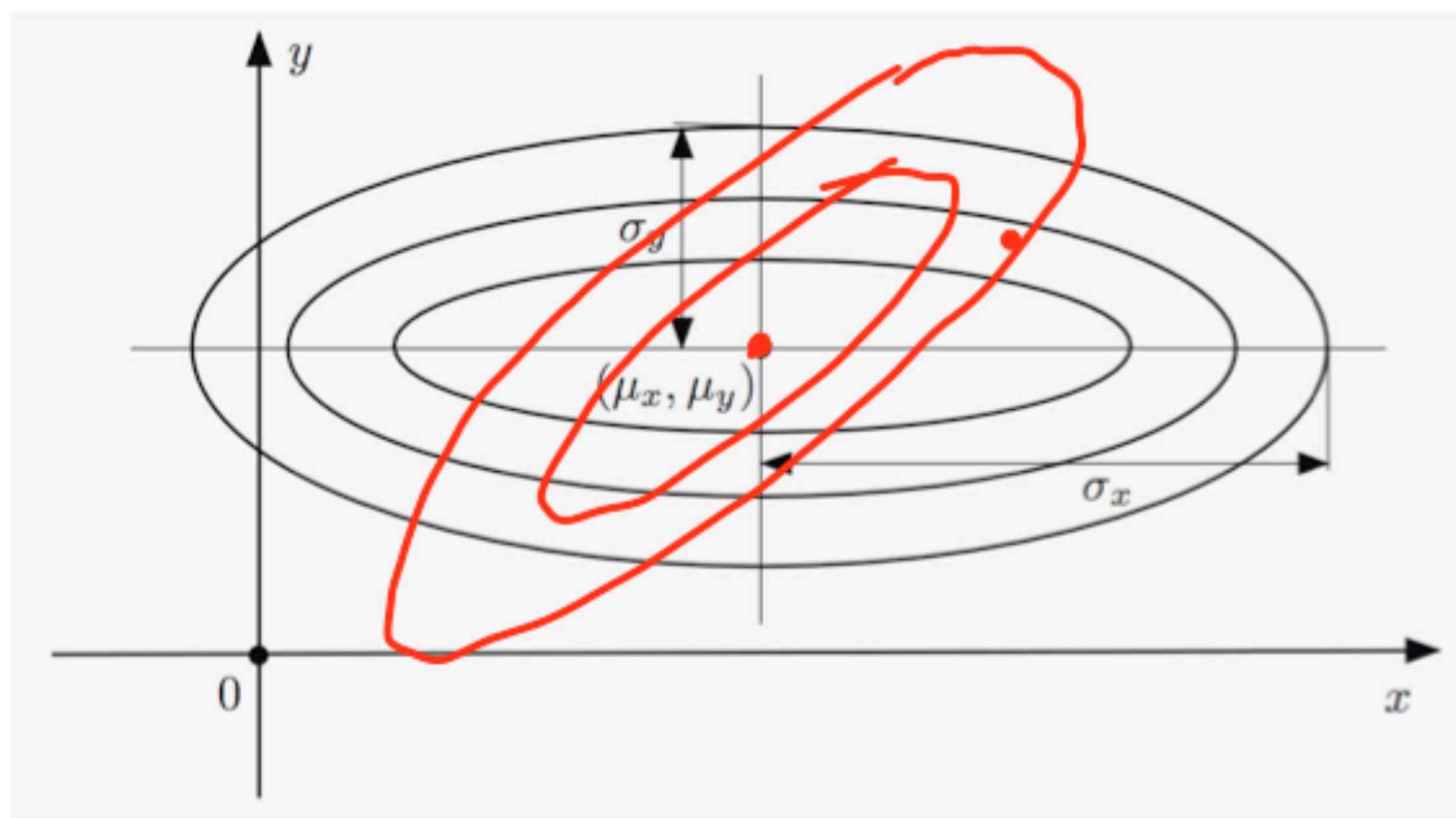
$$= \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{x^2}{2}\right\} \cdot \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{y^2}{2}\right\}$$

$$= \frac{1}{2\pi} \exp\left\{-\frac{1}{2}(x^2 + y^2)\right\}$$

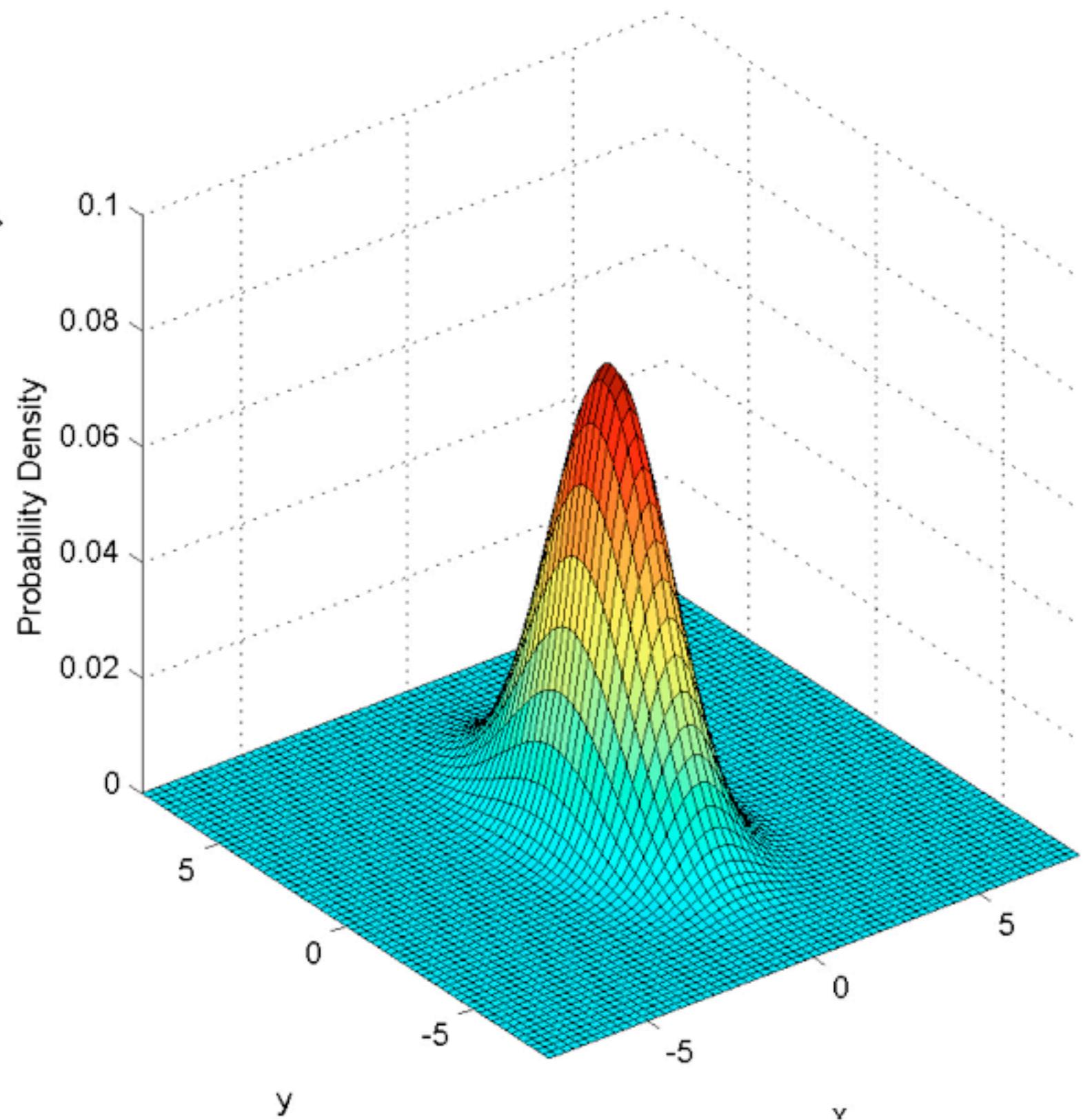


Independent normals

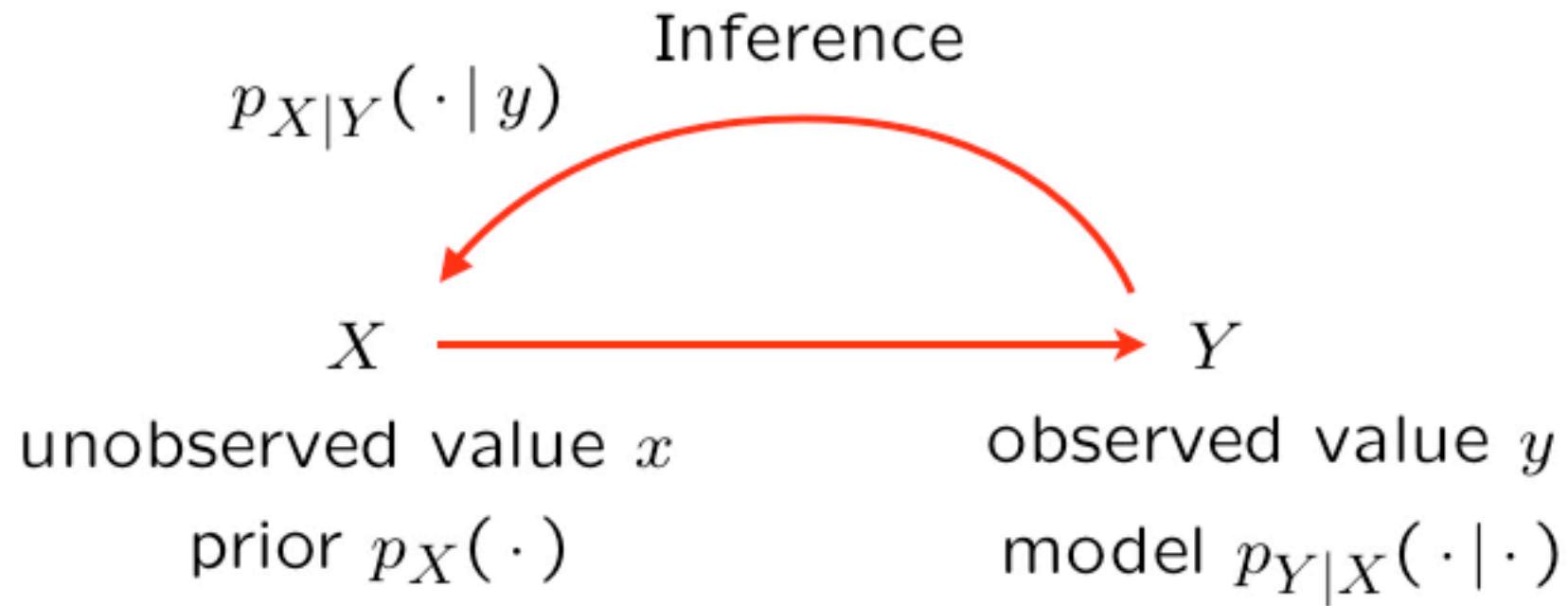
$$f_{X,Y}(x,y) = f_X(x)f_Y(y)$$
$$= \frac{1}{2\pi\sigma_x\sigma_y} \exp \left\{ -\frac{(x-\mu_x)^2}{2\sigma_x^2} - \frac{(y-\mu_y)^2}{2\sigma_y^2} \right\}$$



$\mu_X=\mu_Y=0; \sigma_X^2=1, \sigma_Y^2=4$



The Bayes rule — a theme with variations



$$\begin{aligned} p_{X,Y}(x,y) &= p_X(x) p_{Y|X}(y | x) \\ &= p_Y(y) p_{X|Y}(x | y) \end{aligned}$$

$$\begin{aligned} f_{X,Y}(x,y) &= f_X(x) f_{Y|X}(y | x) \\ &= f_Y(y) f_{X|Y}(x | y) \end{aligned}$$

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

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

posterior $p_Y(y) = \sum_{x'} p_X(x') p_{Y|X}(y | x')$

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

The Bayes rule — one discrete and one continuous random variable

K : discrete

Y : continuous

$$\begin{aligned}
 & P(K=k, y \leq Y \leq y+\delta) \quad \delta > 0, \delta \approx 0 \\
 & = P(K=k) P(y \leq Y \leq y+\delta | K=k) \quad \approx \quad p_K(k) f_{Y|K}(y|k) \\
 & = P(y \leq Y \leq y+\delta) P(K=k | y \leq Y \leq y+\delta) \approx f_Y(y) p_{K|Y}(k|y)
 \end{aligned}$$

$$p_{K|Y}(k|y) = \frac{p_K(k) f_{Y|K}(y|k)}{f_Y(y)}$$

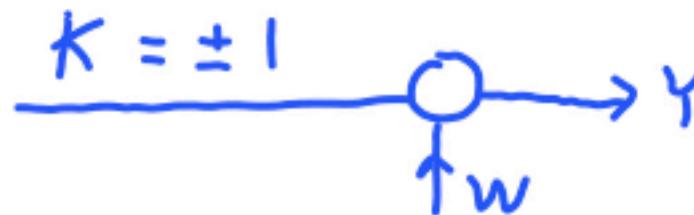
$$f_{Y|K}(y|k) = \frac{f_Y(y) p_{K|Y}(k|y)}{\bullet p_K(k)}$$

$$f_Y(y) = \sum_{k'} p_K(k') f_{Y|K}(y|k')$$

$$p_K(k) = \int f_Y(y') p_{K|Y}(k|y') dy'$$

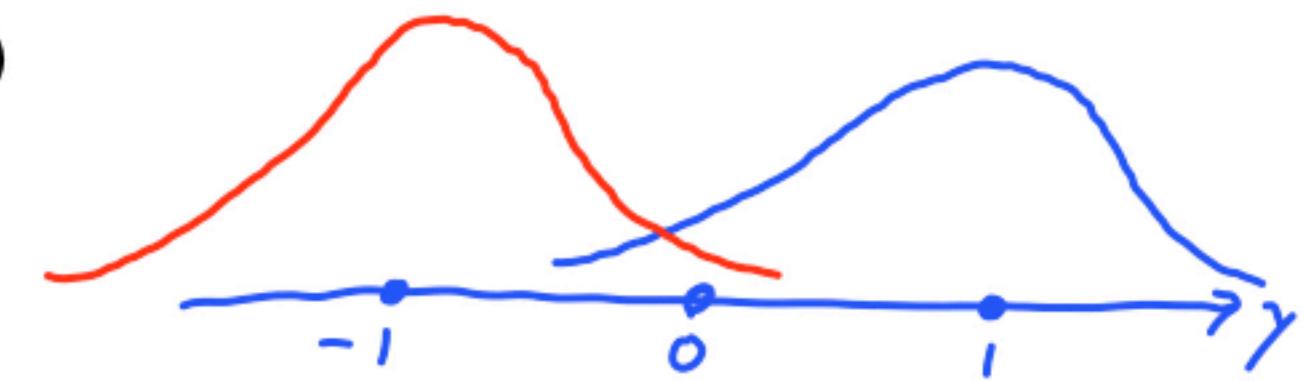
The Bayes rule — discrete unknown, continuous measurement

- unkown K : equally likely to be -1 or $+1$
- measurement Y : $Y = K + W$; $W \sim \mathcal{N}(0, 1)$



$$Y|K=1 \sim \mathcal{N}(1, 1)$$

$$Y|K=-1 \sim \mathcal{N}(-1, 1)$$



- Probability that $K = 1$, given that $Y = y$? $P_{K|Y}(1|y)$

$$p_K(k) = 1/2 \quad f_{Y|K}(y|k) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(y-k)^2}$$

$k = -1, +1$

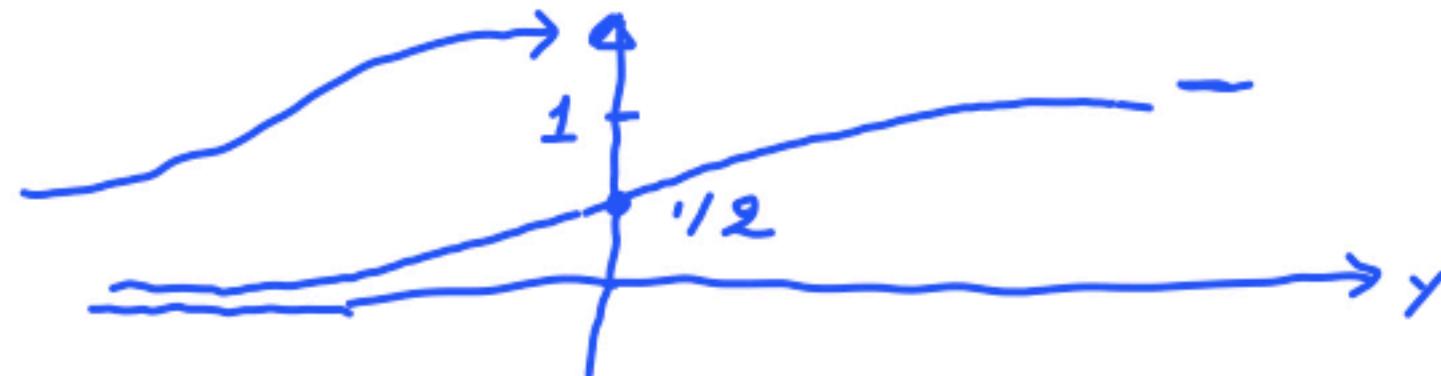
$$f_Y(y) = \frac{1}{2} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(y+1)^2} + \frac{1}{2} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(y-1)^2}$$

$$p_{K|Y}(k|y) = \frac{p_K(k) f_{Y|K}(y|k)}{f_Y(y)}$$

$$f_Y(y) = \sum_{k'} p_K(k') f_{Y|K}(y|k')$$

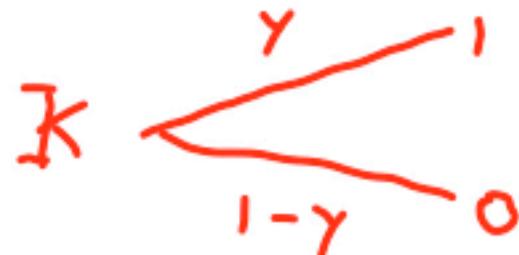
$$p_{K|Y}(1|y) = \frac{1}{1 + e^{-2y}}$$

algebra

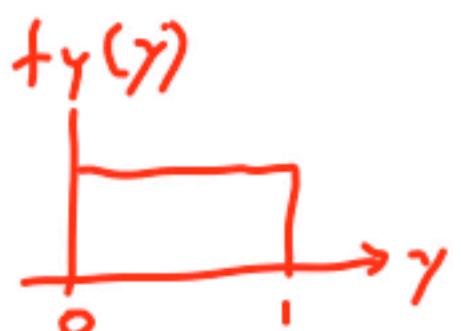


The Bayes rule — continuous unknown, discrete measurement

- measurement K : Bernoulli with parameter Y



- unkown Y : uniform on $[0, 1]$



- Distribution of Y given that $K = 1$?

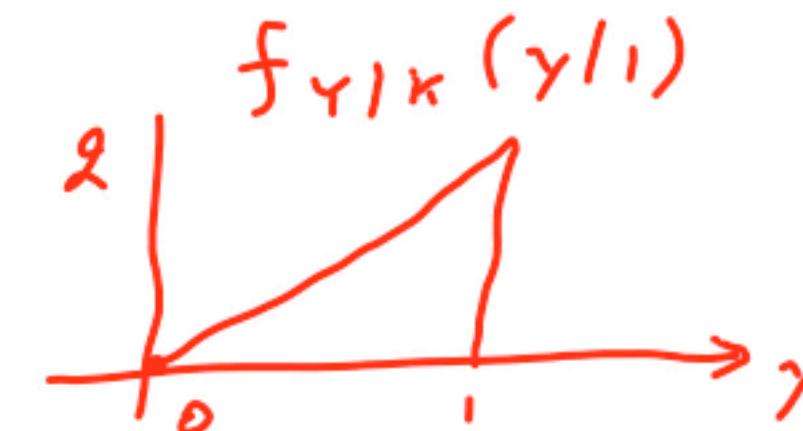
$$f_{Y|K}(y|1)$$

$$f_Y(y) = \begin{cases} 1 & y \in [0, 1] \\ 0 & \text{otherwise} \end{cases}$$

$$p_{K|Y}(1|y) =$$

$$p_K(1) = \int_0^1 1 \cdot y \, dy = \frac{y^2}{2} \Big|_0^1 = \frac{1}{2}$$

$$f_{Y|K}(y|1) = \frac{1 \cdot y}{1/2} = 2y, \quad y \in [0, 1]$$



$$f_{Y|K}(y|k) = \frac{f_Y(y) p_{K|Y}(k|y)}{p_K(k)}$$

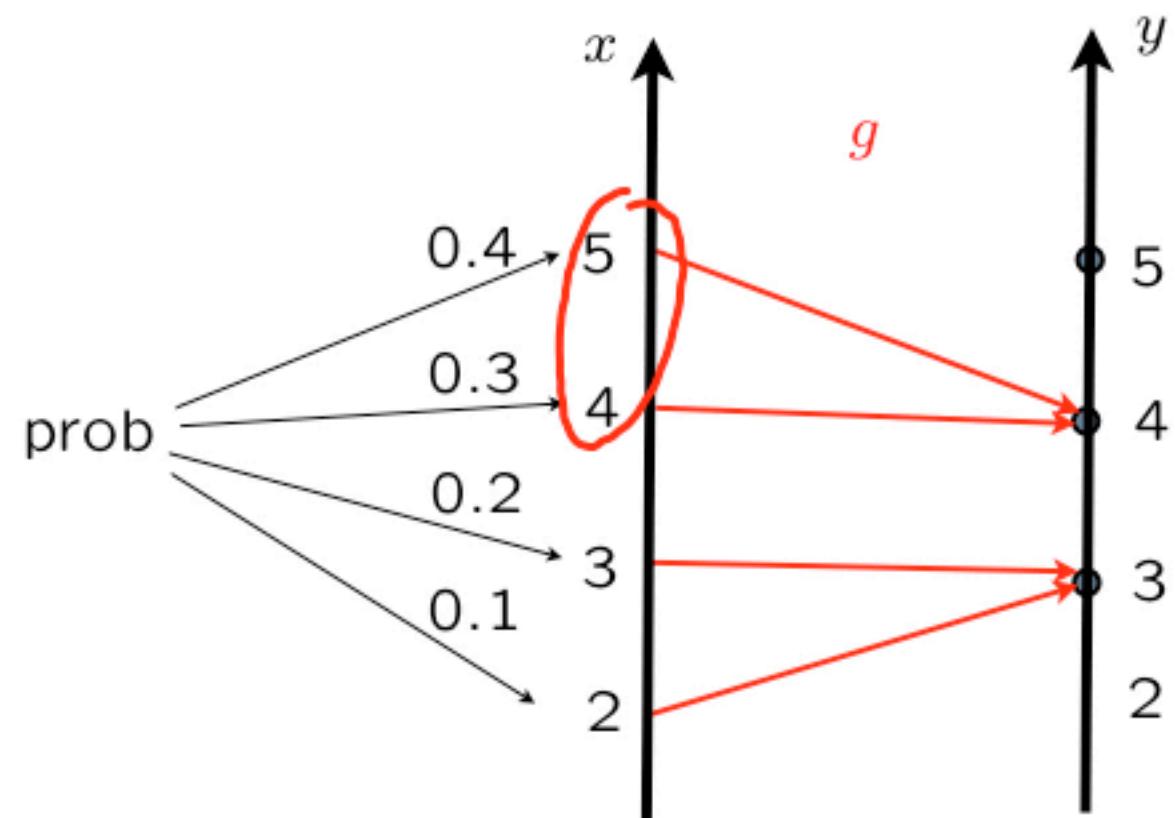
$$p_K(k) = \int f_Y(y') p_{K|Y}(k|y') dy'$$

LECTURE 11: Derived distributions

- Given the distribution of X ,
find the distribution of $Y = g(X)$
 - the discrete case
 - the continuous case
 - general approach, using CDFs
 - the linear case: $Y = aX + b$
 - general formula when g is monotonic
- Given the (joint) distribution of X and Y ,
find the distribution of $Z = g(X, Y)$

Derived distributions — the discrete case

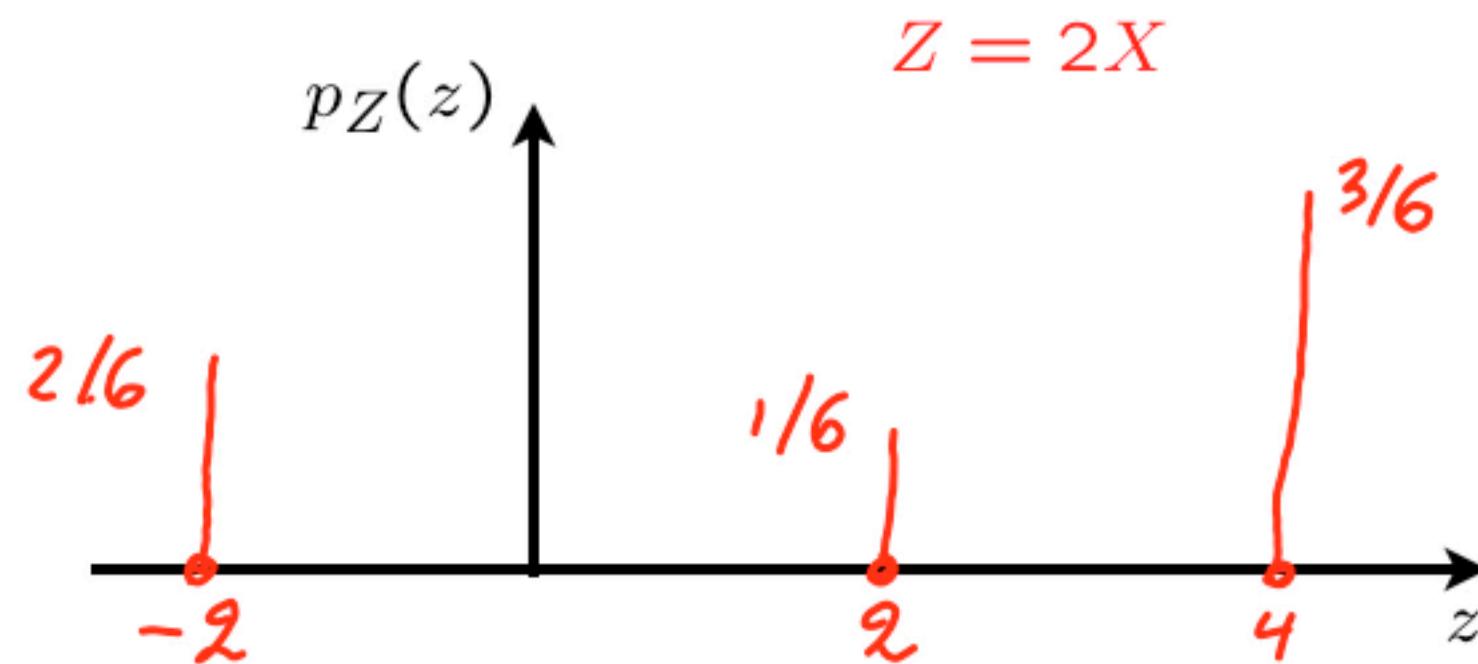
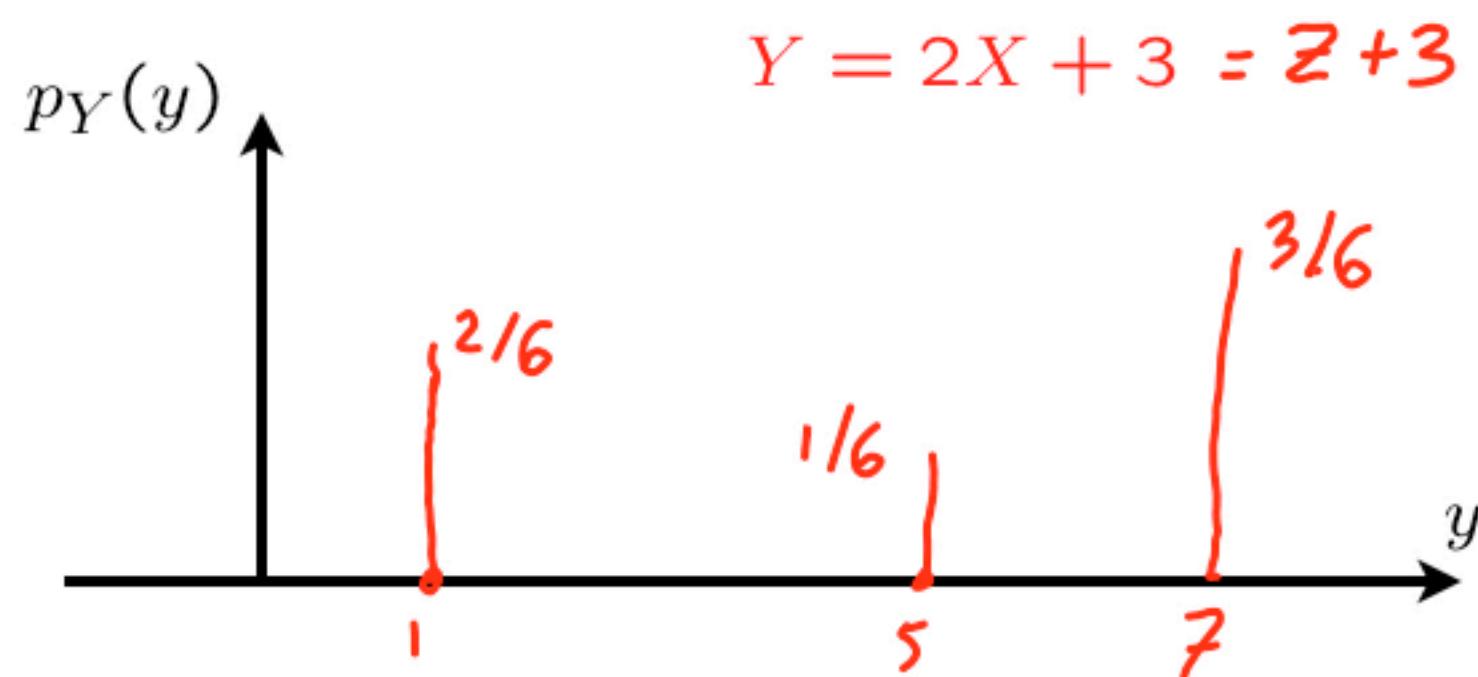
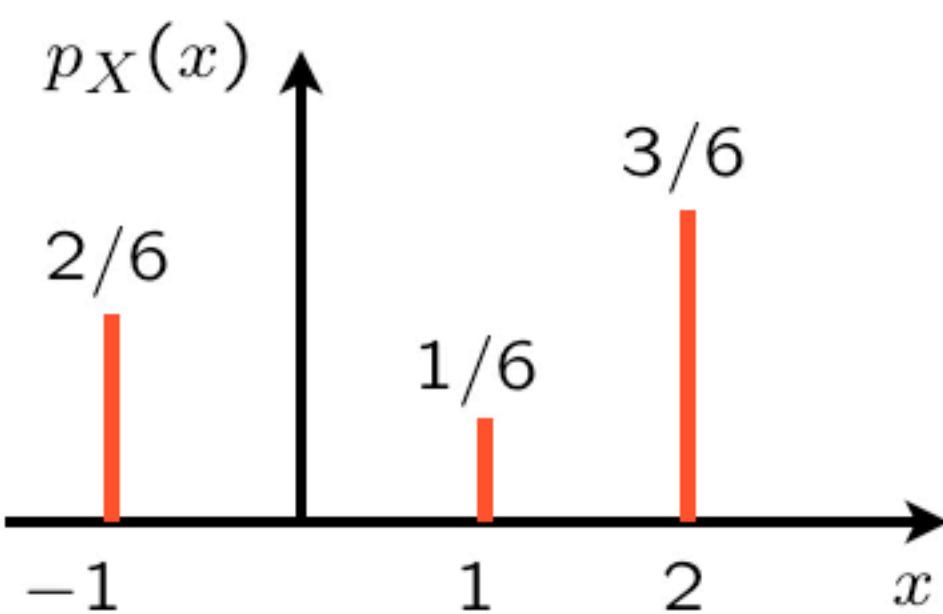
$$Y = g(X)$$



$$\begin{aligned} p_Y(4) &= P(Y=4) \\ &= P(X=4) + P(X=5) \\ &= p_X(4) + p_X(5) = 0.3 + 0.4 \end{aligned}$$

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

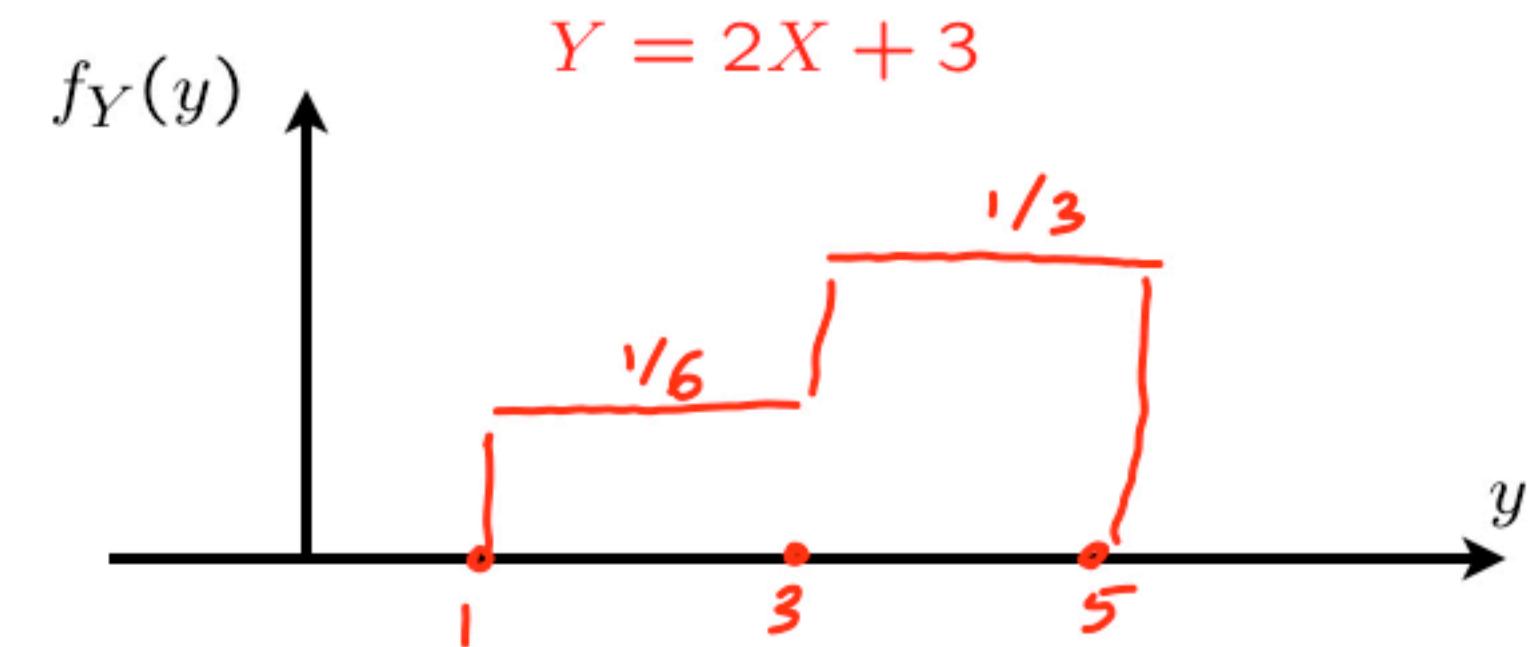
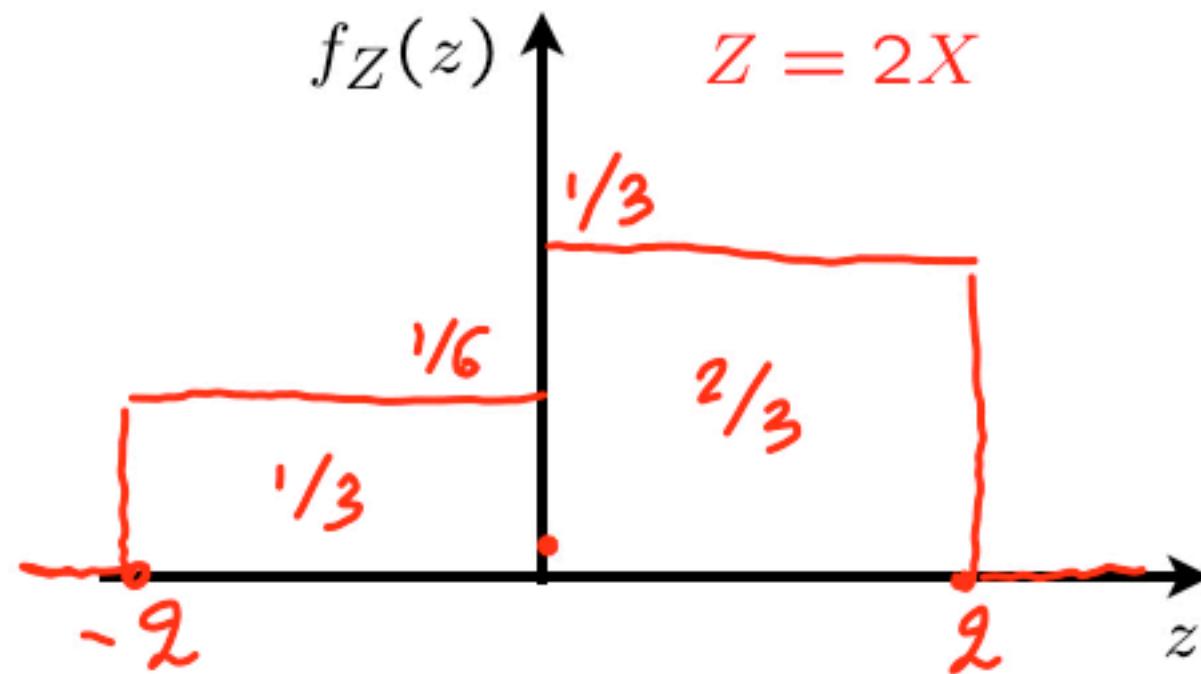
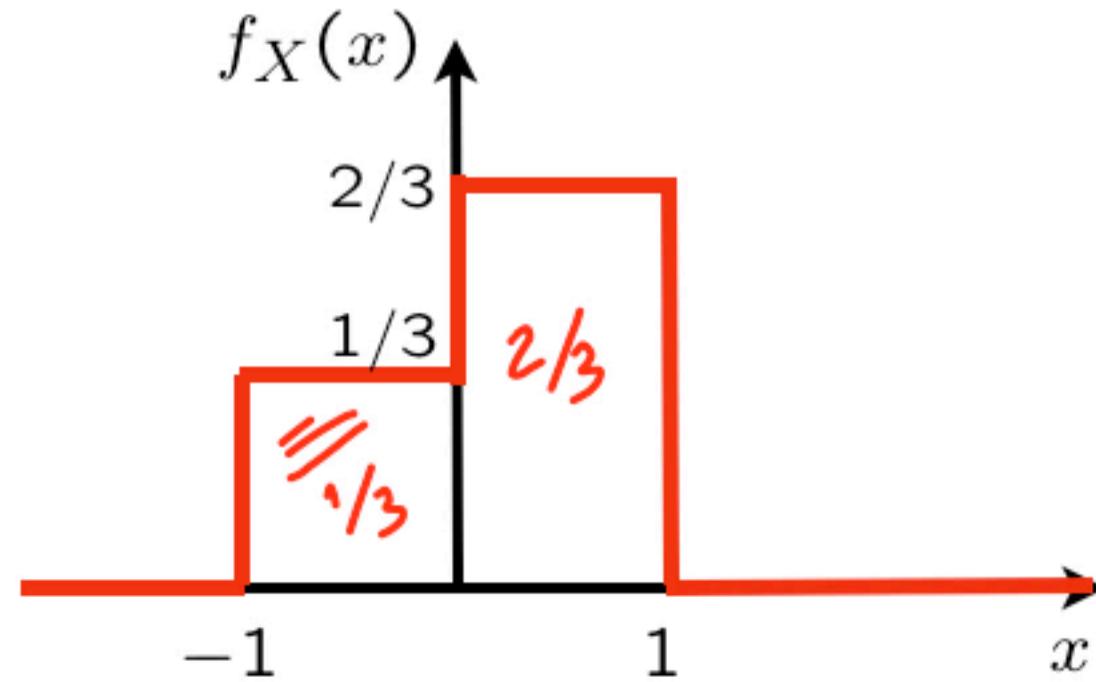
A linear function of a discrete r.v.



$$\begin{aligned}
 P_Y(y) &= P(Y=y) = P(2X+3=y) \\
 &= P\left(X = \frac{y-3}{2}\right) = P_X\left(\frac{y-3}{2}\right)
 \end{aligned}$$

$$Y = aX + b : \quad p_Y(y) = p_X\left(\frac{y-b}{a}\right)$$

A linear function of a continuous r.v.



A linear function of a continuous r.v.

$$Y = aX + b$$

$$a > 0$$

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

$$F_Y(y) = P(Y \leq y) = P(aX+b \leq y)$$

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

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

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

$$a < 0$$

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

$$= 1 - P\left(X \leq \frac{y-b}{a}\right)$$

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

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

$$p_Y(y) = p_X\left(\frac{y-b}{a}\right) \cdot$$

A linear function of a normal r.v. is normal

$$X \sim \mathcal{N}(\mu, \sigma^2)$$

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

$$Y = aX + b, \quad a \neq 0$$

$$\begin{aligned} f_Y(y) &= \frac{1}{|a|} \frac{1}{\sqrt{2\pi}\sigma} e^{-\left(\frac{y-b}{a}-\mu\right)^2/2\sigma^2} \\ &= \frac{1}{\sqrt{2\pi}\sigma|a|} e^{-\frac{(y-b-a\mu)^2}{2\sigma^2 a^2}} \end{aligned}$$

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

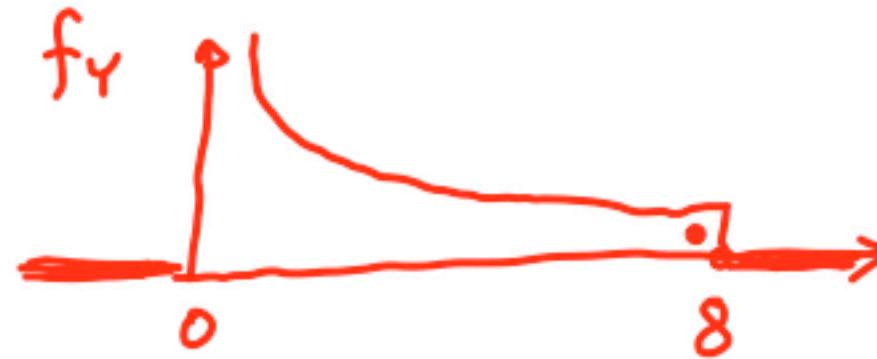
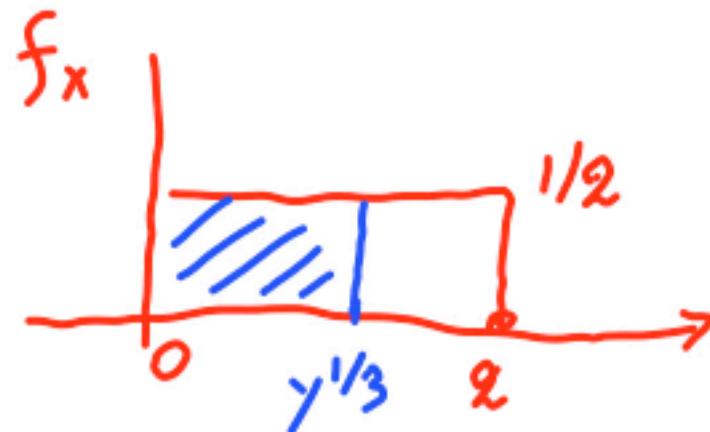
If $X \sim \mathcal{N}(\mu, \sigma^2)$, then $aX + b \sim \mathcal{N}(a\mu + b, a^2\sigma^2)$

A general function $g(X)$ of a continuous r.v.

- Two-step procedure:

- Find the CDF of Y : $F_Y(y) = \mathbf{P}(Y \leq y) = \mathbf{P}(g(X) \leq y)$
- Differentiate: $f_Y(y) = \frac{dF_Y}{dy}(y)$

Example: $Y = X^3$; X uniform on $[0, 2]$



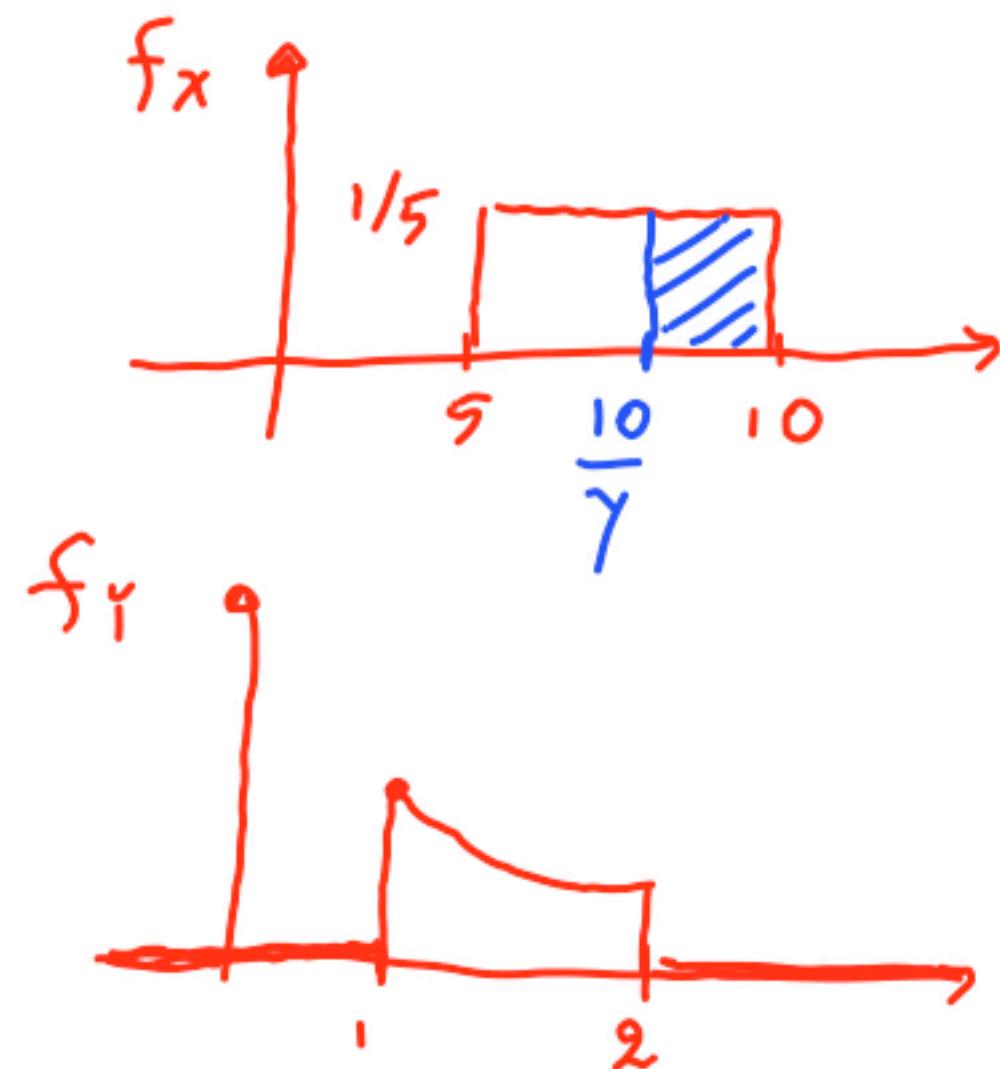
$$0 \leq y \leq 8$$

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

$$f_Y(y) = \frac{d}{dy} F_Y(y) = \frac{1}{2} \cdot \frac{1}{3} y^{-2/3} = \frac{1}{6} \cdot \frac{1}{y^{2/3}}$$

Example: $Y = a/X$

- You go to the gym and set the speed X of the treadmill to a number between 5 and 10 km/hr (with a uniform distribution). Find the PDF of the time it takes to run 10km.



$$\text{time} = Y = \frac{10}{X} \quad 1 \leq Y \leq 2$$

$$F_Y(y) = P(Y \leq y) = P\left(\frac{10}{X} \leq y\right)$$

$$= P\left(X \geq \frac{10}{y}\right) = \frac{1}{5} \left(10 - \frac{10}{y}\right)$$

$$f_Y(y) = \frac{1}{5} \frac{(-10)}{-y^2} = \frac{2}{y^2}, \quad 1 \leq y \leq 2$$

$$= 0, \quad \text{otherwise}$$

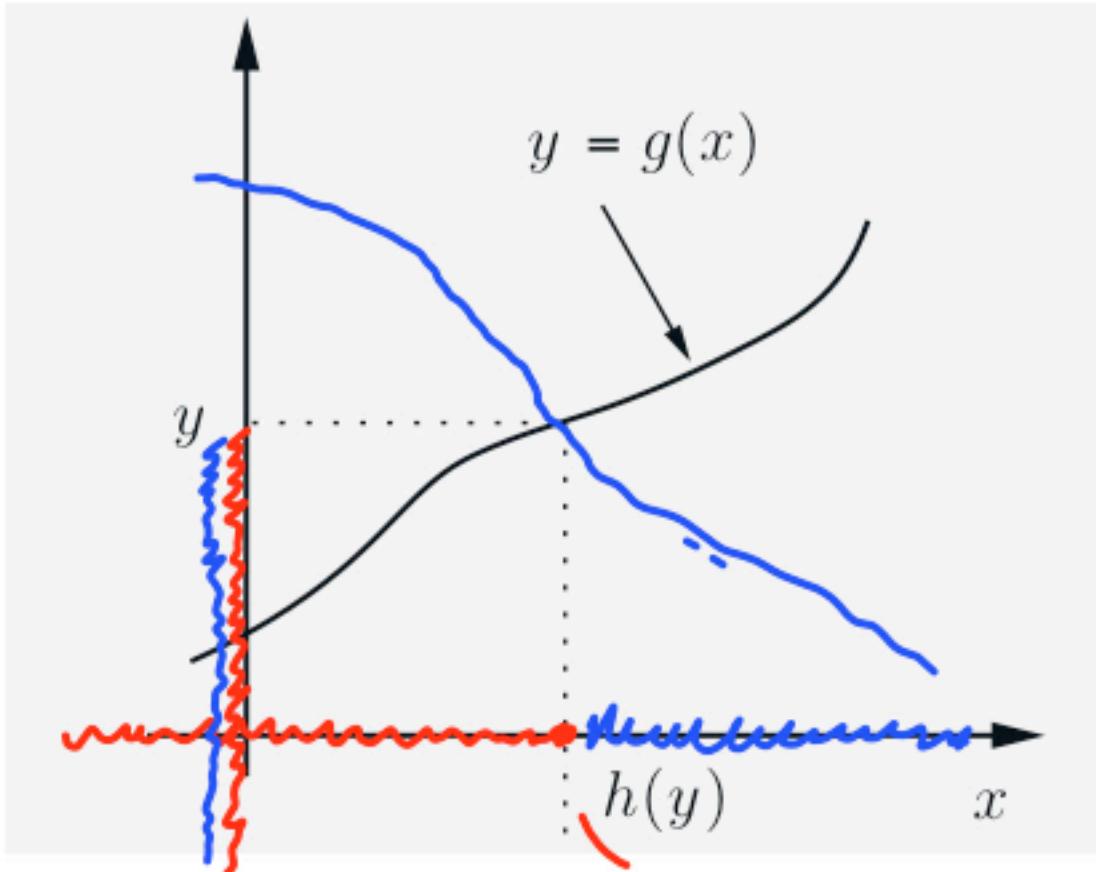
A general formula for the PDF of $Y = g(X)$ when g is monotonic

$$x^3 \frac{\alpha}{x}$$

~~decreasing~~ $x < x' \Rightarrow g(x) < g(x')$

Assume g strictly increasing

and differentiable



$$F_Y(y) = P(Y \leq y) = P(X \leq h(y)) = F_X(h(y))$$

$$f_Y(y) = f_X(h(y)) \left| \frac{dh}{dy}(y) \right|$$

$$\begin{aligned} F_Y(y) &= P(Y \leq y) = P(X \geq h(y)) \\ &= 1 - P(X \leq h(y)) = 1 - F_X(h(y)) \end{aligned}$$

$$f_Y(y) = f_X(h(y)) \left| \frac{dh}{dy}(y) \right|$$

inverse function $h \rightarrow$ decreasing

$$f_Y(y) = f_X(h(y)) \left| \frac{dh}{dy}(y) \right|$$

Example: $Y = X^2$; X uniform on $[0, 1]$

$$f_Y(y) = f_X(h(y)) \left| \frac{dh}{dy}(y) \right|$$

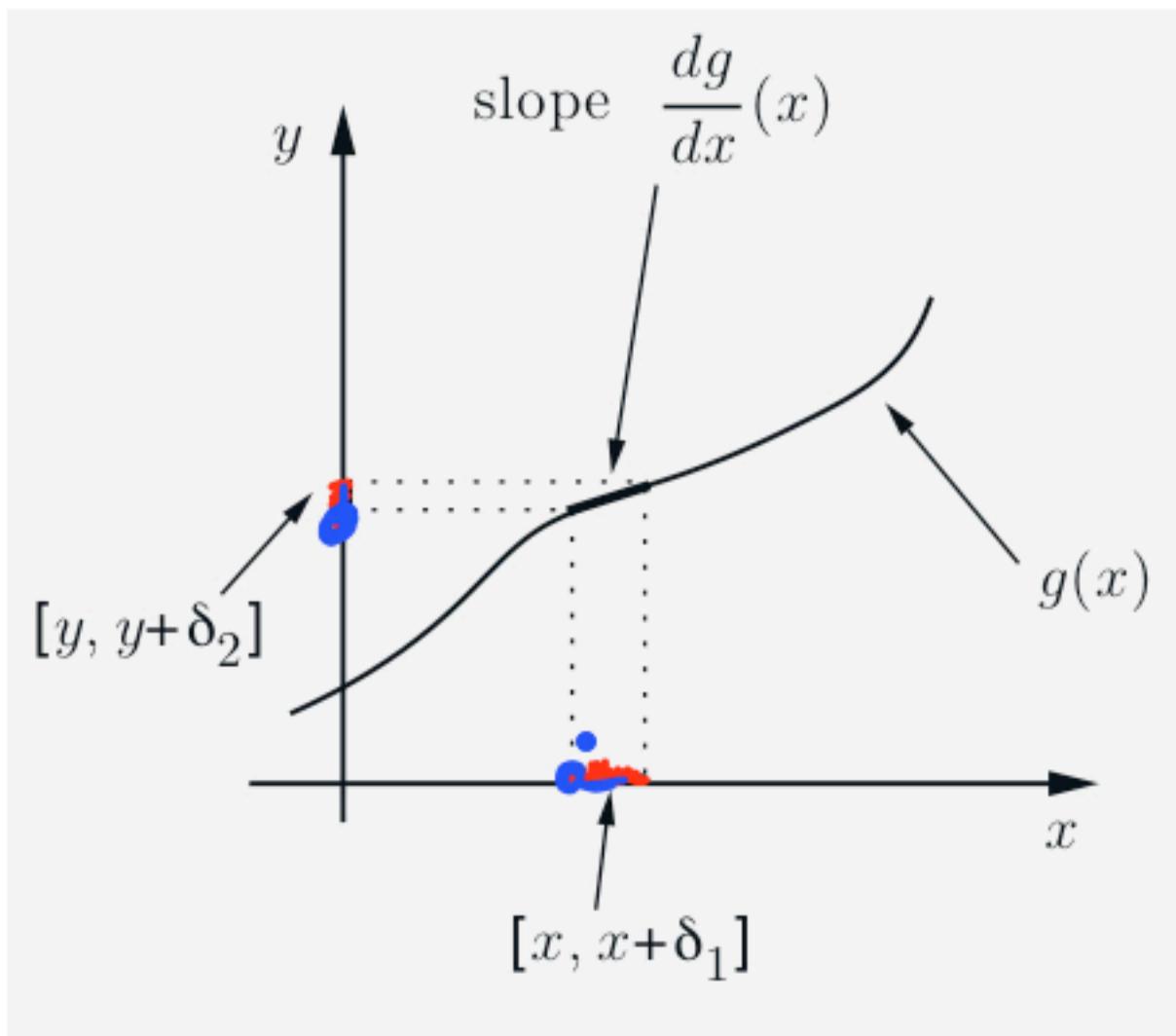


$$y = x^2 \Leftrightarrow x = \sqrt{y} \quad h(y) = \sqrt{y}$$

$$f_Y(y) = \frac{1}{2\sqrt{y}}$$

$$0 \leq y \leq 1$$

An intuitive explanation for the monotonic case



$$y = g(x)$$

$$x = h(y)$$

$$\delta_2 \approx \delta_1 \cdot \frac{\partial g}{\partial x}(x)$$

$$\delta_1 \approx \delta_2 \cdot \frac{\partial h}{\partial y}(y) \quad \textcircled{R}$$

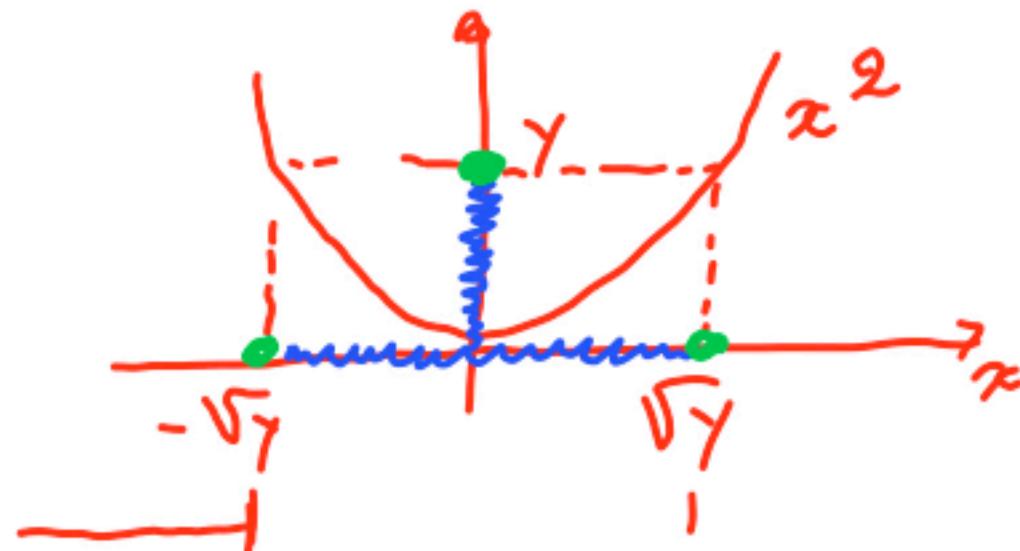
$$f_y(y) \cancel{\delta_2} \stackrel{P}{=} P(y \leq Y \leq y + \delta_2) = P(x \leq X \leq x + \delta_1)$$

$$\approx f_x(x) \delta_1 \approx f_x(x) \cancel{\delta_2} \frac{\partial h}{\partial y}(y)$$

$$f_y(y) = f_x(x) \frac{\partial h}{\partial y}(y)$$

$$= f_x(h(y)) \frac{\partial h}{\partial y}(y)$$

A nonmonotonic example: $Y = X^2$



- The discrete case:

$$p_Y(9) = P(X=3) + P(X=-3)$$

$$p_Y(y) = P_x(\sqrt{y}) + P_x(-\sqrt{y})$$

- The continuous case:

$$Y \geq 0$$

$$F_Y(y) = P(Y \leq y) = P(X^2 \leq y) = P(|X| \leq \sqrt{y}) = P(-\sqrt{y} \leq X \leq \sqrt{y})$$

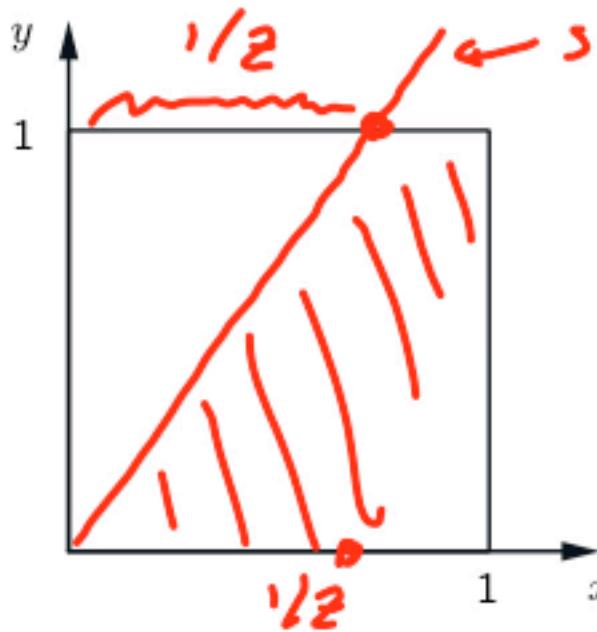
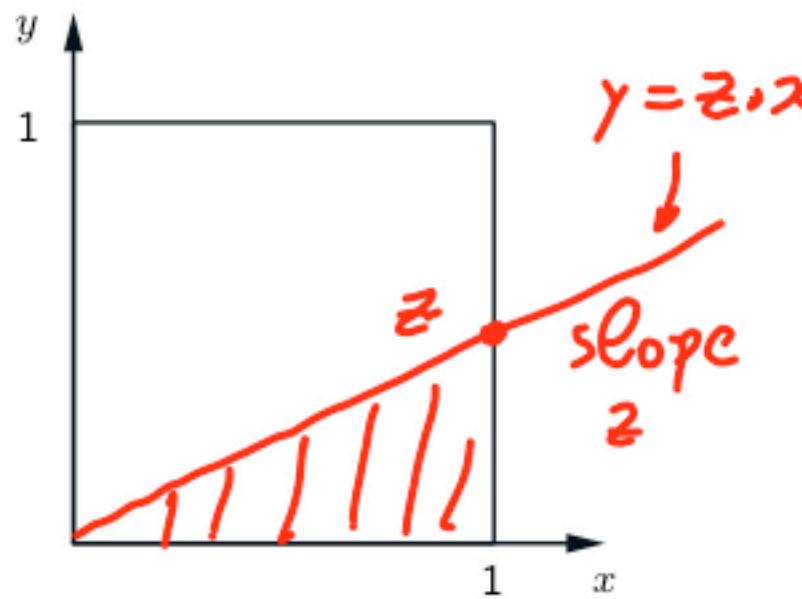
$$= F_X(\sqrt{y}) - F_X(-\sqrt{y})$$

$$f_Y(y) = f_X(\sqrt{y}) \frac{1}{2\sqrt{y}} + f_X(-\sqrt{y}) \frac{\cancel{-1}}{2\sqrt{y}}$$

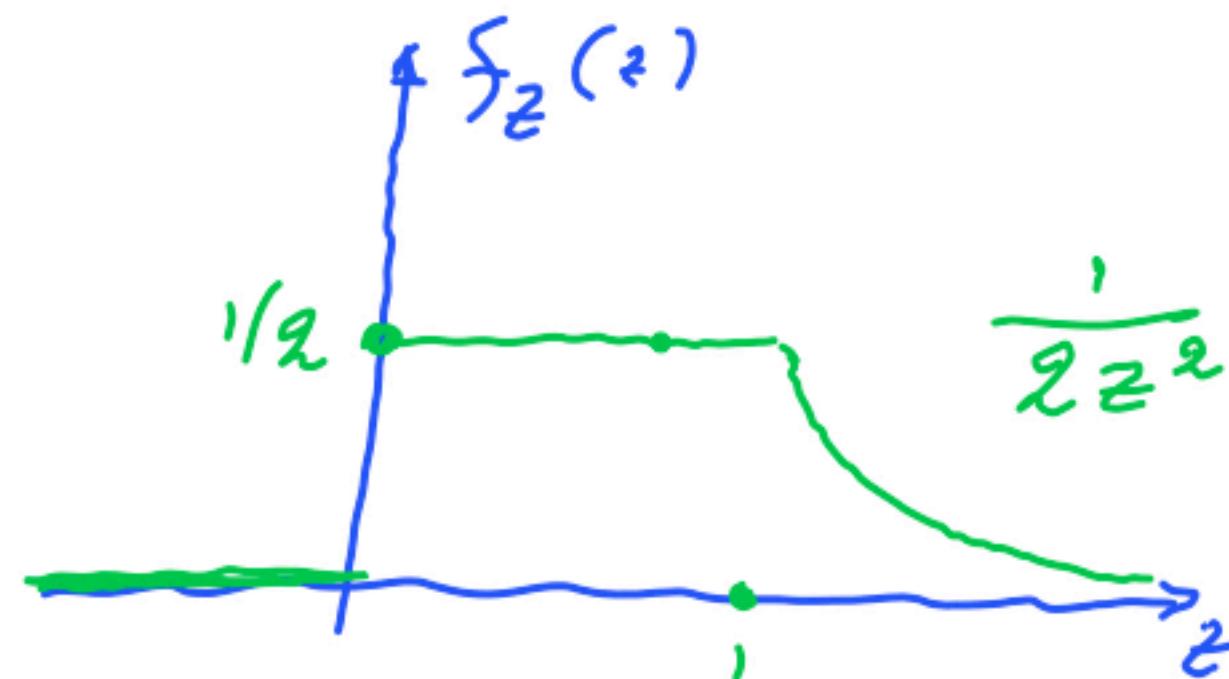
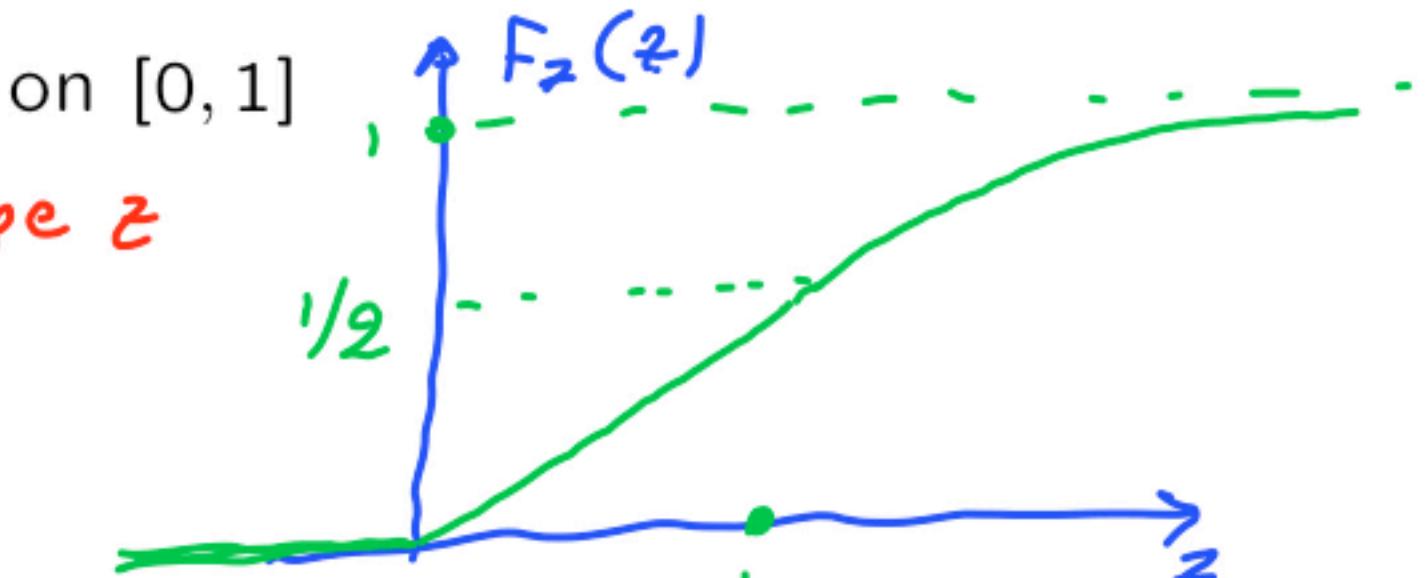
A function of multiple r.v.'s: $Z = g(X, Y)$

- Same methodology: find CDF of Z

- Let $Z = Y/X$; X, Y independent, uniform on $[0, 1]$



$$\begin{aligned}
 F_Z(z) &= P\left(\frac{Y}{X} \leq z\right) = 0, \quad z < 0 \\
 &= \frac{1}{2} \cdot z, \quad 0 \leq z \leq 1 \\
 &= 1 - \frac{1}{2z}, \quad z > 1
 \end{aligned}$$



LECTURE 12: Sums of independent random variables; Covariance and correlation

- The PMF/PDF of $X + Y$ (X and Y independent)
 - the discrete case
 - the continuous case
 - the mechanics
 - the sum of independent normals
- Covariance and correlation
 - definitions
 - mathematical properties
 - interpretation

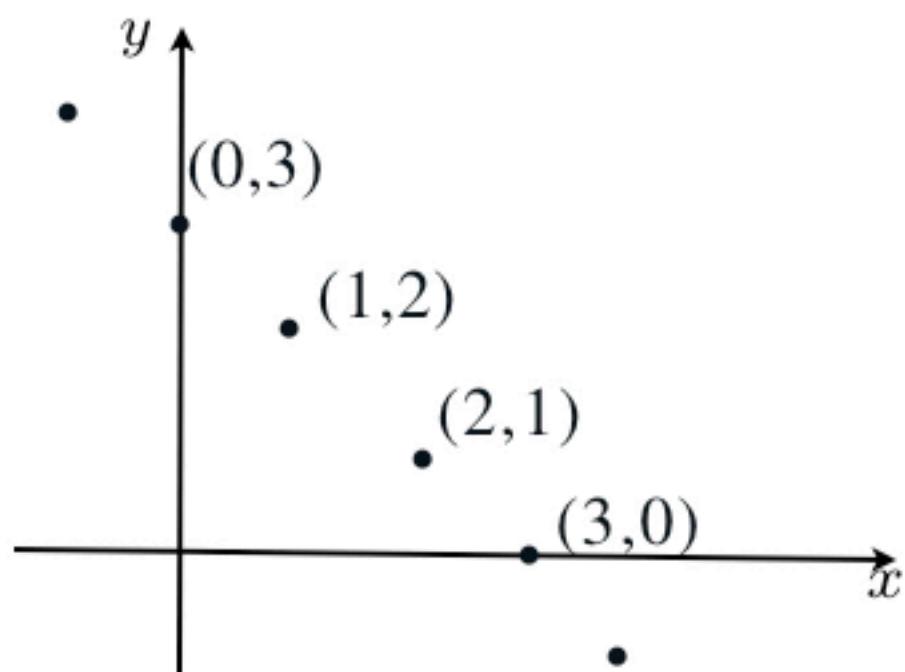
The distribution of $X + Y$: the discrete case

- $Z = X + Y$; X, Y independent, discrete

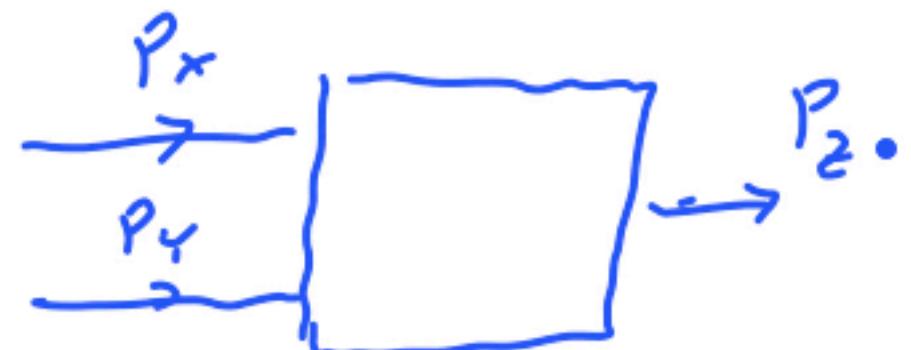
$p(x,y)$

known PMFs

$$\begin{aligned} p_Z(3) &= \dots + P(X=0, Y=3) + P(X=1, Y=2) + \dots \\ &= \dots + P_X(0) P_Y(3) + P_X(1) P_Y(2) + \dots \end{aligned}$$



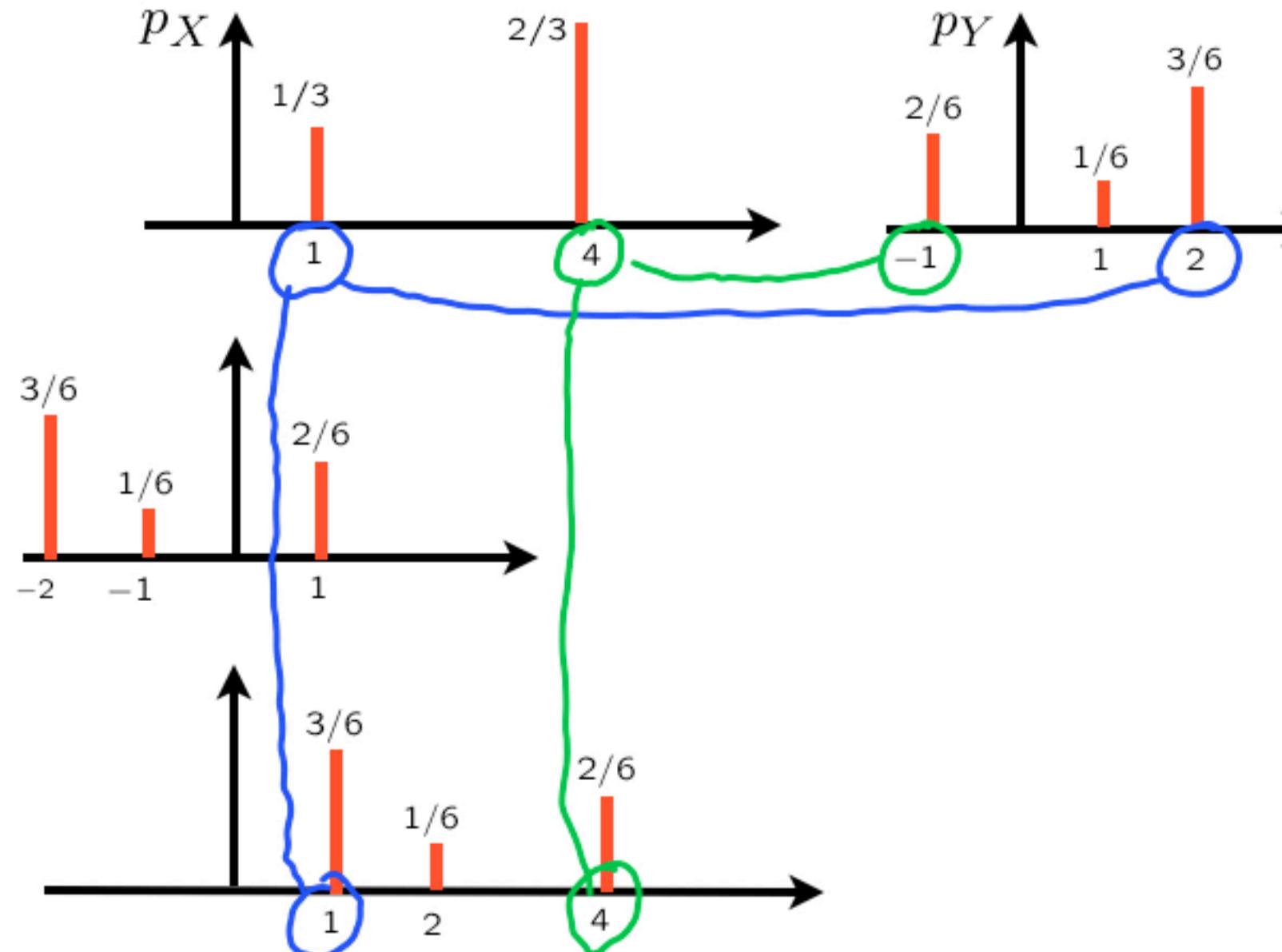
$$p_Z(z) = \sum_x p_X(x) p_Y(z-x)$$



$$\begin{aligned} P_Z(z) &= \sum_x P(X=x, Y=z-x) \\ &= \sum_x P_X(x) P_Y(z-x) \end{aligned}$$

Discrete convolution mechanics

$$p_Z(z) = \sum_x p_X(x) p_Y(z - x)$$



- To find $p_Z(3)$:

- Flip (horizontally) the PMF of Y
- Put it underneath the PMF of X
- Right-shift the flipped PMF by 3
- Cross-multiply and add
- Repeat for other values of z

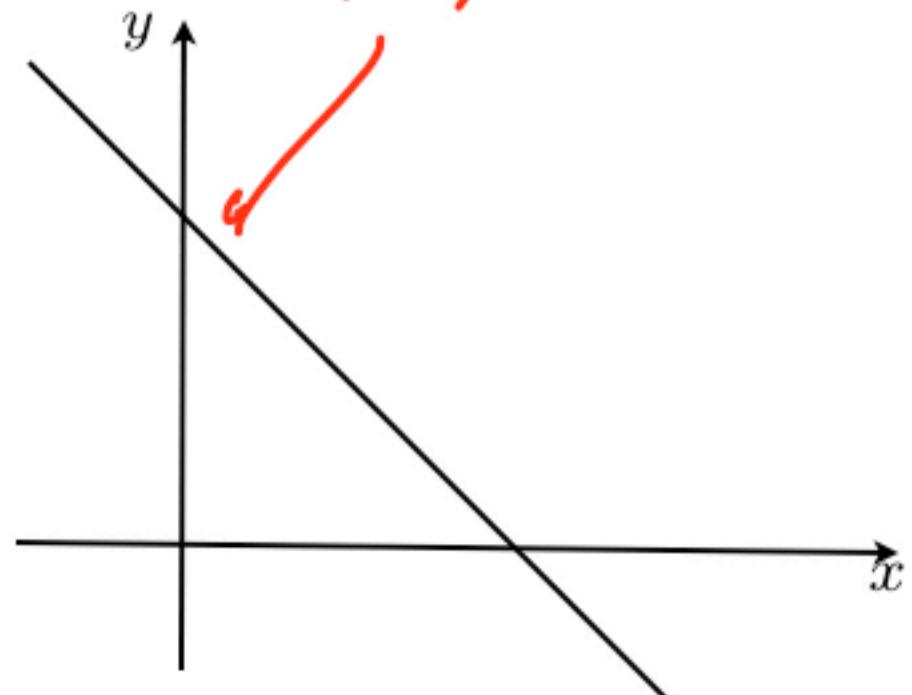
The distribution of $X + Y$: the continuous case

- $Z = X + Y$; X, Y independent, continuous known PDFs

$$p_Z(z) = \sum_x p_X(x) p_Y(z - x)$$

$$f_Z(z) = \int_{-\infty}^{\infty} f_X(x) f_Y(z - x) dx$$

$x+y=\text{constant}$



Conditional on $X = x$: $Z = x + Y$ $x = 3$ $Z = Y + 3$

$$f_{Z|X}(z|x) = f_{Y+3|X}(z|x) = f_{Y+3}(z) = f_Y(z-3)$$

$$f_{Z|X}(z|x) = f_Y(z-x) \quad f_{X+b}(x) = f_X(x-b)$$

Joint PDF of Z and X :

$$f_{X,Z}(x,z) = f_X(x) f_Y(z-x)$$

From joint to the marginal: $f_Z(z) = \int_{-\infty}^{\infty} f_{X,Z}(x,z) dx$

- Same mechanics as in discrete case (flip, shift, etc.)

The sum of independent normal r.v.'s

$$f_Z(z) = \int_{-\infty}^{\infty} f_X(x)f_Y(z-x) dx$$

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

$$f_X(x) = \frac{1}{\sqrt{2\pi}\sigma_x} e^{-(x-\mu_x)^2/2\sigma_x^2}$$

$$f_Y(y) = \frac{1}{\sqrt{2\pi}\sigma_y} e^{-(y-\mu_y)^2/2\sigma_y^2}$$

$$f_Z(z) = \int_{-\infty}^{\infty} f_X(x)f_Y(z-x) dx$$

$X + Y + W$

$$= \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}\sigma_x} \exp\left\{-\frac{(x-\mu_x)^2}{2\sigma_x^2}\right\} \frac{1}{\sqrt{2\pi}\sigma_y} \exp\left\{-\frac{(z-x-\mu_y)^2}{2\sigma_y^2}\right\} dx$$

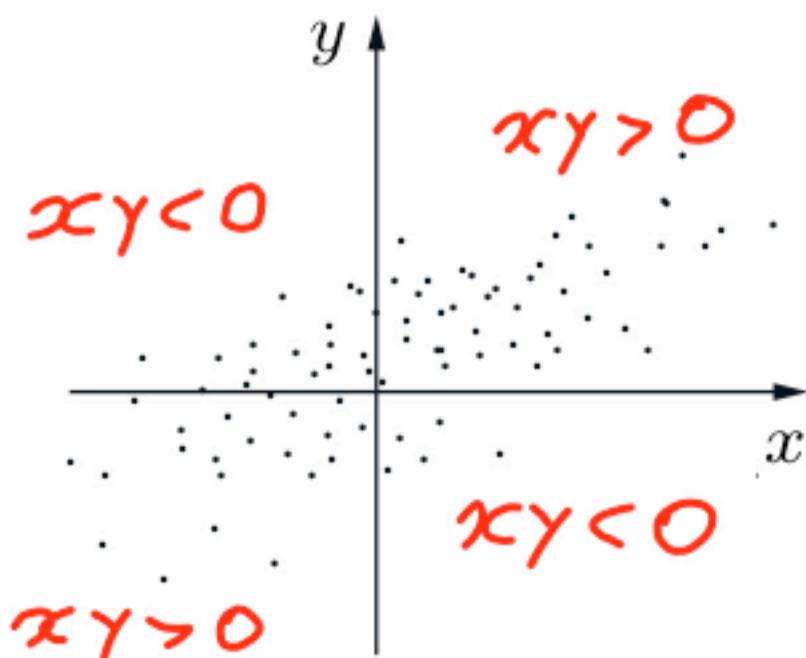
$$\text{(algebra)} = \frac{1}{\sqrt{2\pi(\sigma_x^2 + \sigma_y^2)}} \exp\left\{-\frac{(z-\mu_x-\mu_y)^2}{2(\sigma_x^2 + \sigma_y^2)}\right\} \quad N(\mu_x + \mu_y, \sigma_x^2 + \sigma_y^2)$$

The sum of finitely many independent normals is normal

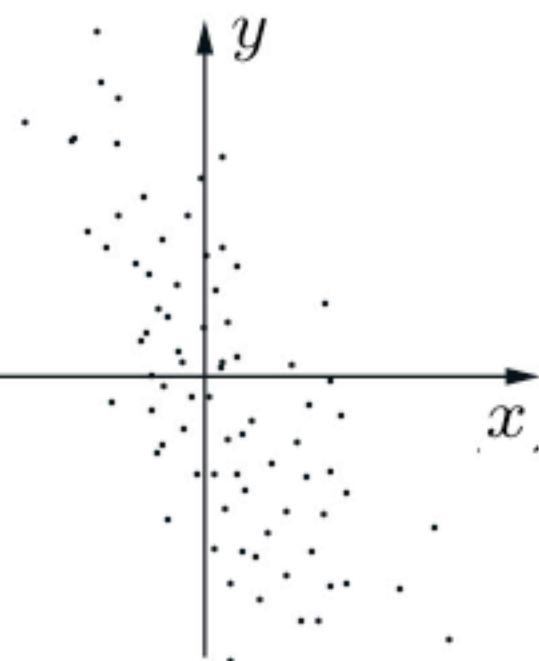
Covariance

- Zero-mean, discrete X and Y
 - if independent: $E[XY] =$

$$= E[X] E[Y] = 0$$



$$E[XY] > 0$$



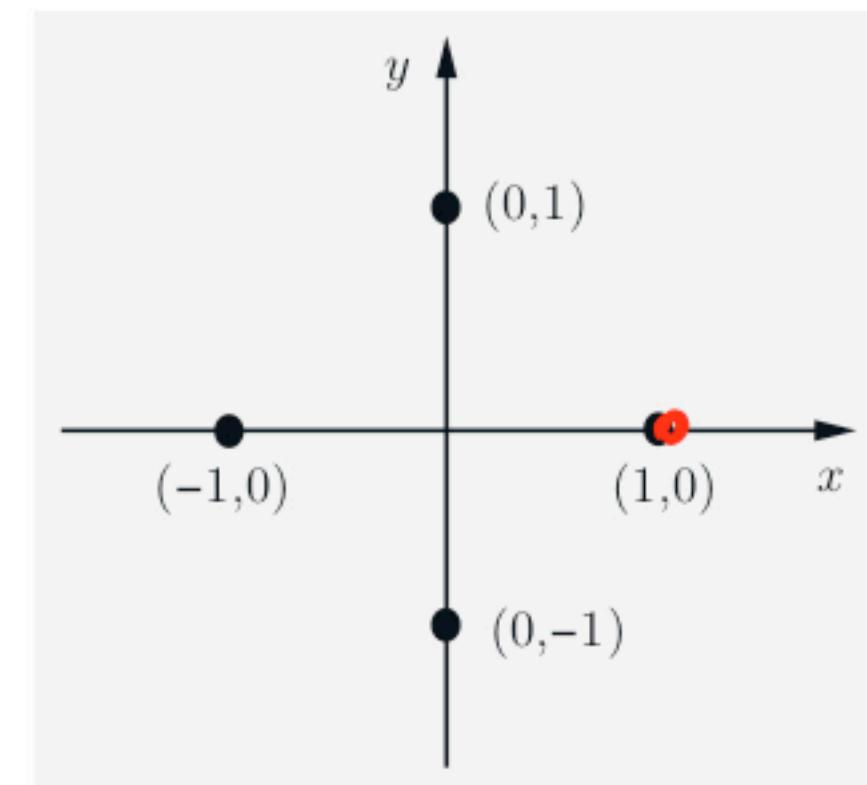
$$E[XY] < 0$$

Definition for general case:

$$\text{cov}(X, Y) = E[(\underline{X - E[X]}) \cdot (\underline{Y - E[Y]})]$$

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

- independent $\Rightarrow \text{cov}(X, Y) = 0$
(converse is not true)



$$XY = 0$$

$$\text{Cov} = 0$$

$$X = 1 \Rightarrow Y = 0$$

Covariance properties

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

$$\text{cov}(aX + b, Y) =$$

(assume 0 means)

$$= E[(aX+b)Y] = aE[XY] + bE[Y]$$
$$= a \cdot \text{cov}(X, Y)$$

$$\text{cov}(X, Y + Z) = E[X(Y + Z)]$$
$$= E[XY] + E[XZ] = \text{cov}(X, Y) + \text{cov}(X, Z)$$

$$\text{cov}(X, Y) = E[(X - E[X]) \cdot (Y - E[Y])]$$

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

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

The variance of a sum of random variables

$$\begin{aligned}\text{var}(X_1 + X_2) &= E[(X_1 + X_2 - E[X_1 + X_2])^2] \\&= E\left[\left(\underline{(X_1 - E[X_1])} + \underline{(X_2 - E[X_2])}\right)^2\right] \\&= E\left[\left(X_1 - E[X_1]\right)^2 + \left(X_2 - E[X_2]\right)^2\right. \\&\quad \left.+ 2(X_1 - E[X_1])(X_2 - E[X_2])\right] \\&= \text{Var}(X_1) + \text{Var}(X_2) + 2 \text{cov}(X_1, X_2).\end{aligned}$$

The variance of a sum of random variables

$$\text{var}(X_1 + X_2) = \text{var}(X_1) + \text{var}(X_2) + 2 \text{cov}(X_1, X_2)$$

$$\begin{aligned}\text{var}(X_1 + \dots + X_n) &= E[(X_1 + \dots + X_n)^2] \\ (\text{assume 0 means}) &= E\left[\sum_{i=1}^n X_i^2 + \sum_{\substack{i=1, \dots, n \\ j=1, \dots, n \\ i \neq j}} X_i X_j\right] \\ &= \sum_i \text{Var}(X_i) + \sum_{i \neq j} \text{Cov}(X_i, X_j)\end{aligned}$$

$\left. \begin{array}{l} i=1, \dots, n \\ j=1, \dots, n \\ i \neq j \end{array} \right\} n^2 - n \text{ terms}$

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

The Correlation coefficient

- Dimensionless version of covariance:

$$-1 \leq \rho \leq 1$$

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

- Measure of the degree of “association” between X and Y
- Independent $\Rightarrow \rho = 0$, “uncorrelated”
(converse is not true)
- $|\rho| = 1 \Leftrightarrow (X - E[X]) = c(Y - E[Y])$ (linearly related)

$$\bullet \quad \rho(X, X) = \frac{\text{var}(X)}{\sigma_X^2} = 1$$

$$\bullet \quad \text{cov}(aX + b, Y) = a \cdot \text{cov}(X, Y) \Rightarrow \rho(aX + b, Y) = \frac{a \text{cov}(X, Y)}{|a| \sigma_X \sigma_Y} = \frac{\text{sign}(a)}{\cdot \rho(X, Y)}$$

Proof of key properties of the correlation coefficient

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

$$-1 \leq \rho \leq 1$$

- Assume, for simplicity, zero means and unit variances, so that $\rho(X, Y) = E[XY]$

$$\begin{aligned} E[(X - \rho Y)^2] &= E[X^2] - 2\rho E[XY] + \rho^2 E[Y^2] \\ 0 \leq &= 1 - 2\rho^2 + \rho^2 = \underline{\underline{1 - \rho^2}} \quad 1 - \rho^2 \geq 0 \Rightarrow \rho^2 \leq 1 \end{aligned}$$

If $|\rho| = 1$, then $X = \rho Y \Rightarrow X = Y$ or $X = -Y$

Interpreting the correlation coefficient

- Association does not imply causation or influence

X : math aptitude

Y : musical ability

- Correlation often reflects underlying, common, hidden factor

- Assume, Z, V, W are independent

$$X = \underline{\underline{Z}} + V$$

$$Y = \underline{\underline{Z}} + W$$

Assume, for simplicity, that Z, V, W have zero means, unit variances

$$\text{var}(x) = \text{var}(Z) + \text{var}(v) = 2 \Rightarrow \sigma_x = \sqrt{2} \quad \sigma_y = \sqrt{2}$$

$$\begin{aligned}\text{cov}(x, y) &= E[(Z+v)(Z+w)] = E[Z^2] + E[vz] + E[zw] + E[vw] \\ &= 1 + 0 + 0 + 0\end{aligned}$$

$$\rho(X, Y) = \frac{\text{cov}(X, Y)}{\sigma_X \sigma_Y}$$

Correlations matter...

- A real-estate investment company invests \$10M in each of 10 states. At each state i , the return on its investment is a random variable X_i , with mean 1 and standard deviation 1.3 (in millions).

$$\text{var}(X_1 + \dots + X_{10}) = \sum_{i=1}^{10} \text{var}(X_i) + \sum_{\{(i,j): i \neq j\}} \text{cov}(X_i, X_j)$$

$$E[X_1 + \dots + X_{10}] = 10$$

- If the X_i are uncorrelated, then:

$$\text{var}(X_1 + \dots + X_{10}) = 10 \cdot (1.3)^2 = 16.9 \quad \sigma(X_1 + \dots + X_{10}) = 4.1$$

- If for $i \neq j$, $\rho(X_i, X_j) = 0.9$: $\text{cov}(X_i, X_j) = \rho \sigma_{X_i} \sigma_{X_j} = 0.9 \times 1.3 \times 1.3$

$$\text{var}(X_1 + \dots + X_{10}) = 10 \cdot (1.3)^2 + 90 \cdot 1.52 = 154 \quad = 1.52$$

$$\sigma(X_1 + \dots + X_{10}) = 12.4$$

LECTURE 13: Conditional expectation and variance revisited;

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

- A more abstract version of the conditional expectation
 - view it as a random variable
 - the law of iterated expectations
- A more abstract version of the conditional variance
 - view it as a random variable
 - the law of total variance
- Sum of a random number
of independent r.v.'s
 - mean
 - variance

Conditional expectation as a random variable

- Function h
e.g., $h(x) = x^2$, for all x
 - Random variable X ; what is $h(X)$?
 $\cancel{= X^2}$
 - $h(X)$ is the r.v. that takes the value x^2 , if X happens to take the value x
- $g(y) = \mathbb{E}[X | Y = y] = \sum_x x p_{X|Y}(x | y)$
 $\cancel{=} \quad$ (integral in continuous case)
 - $g(Y)$: is the r.v. that takes the value $\underline{\mathbb{E}[X | Y = y]}$, if Y happens to take the value y
- Remarks:
 - It is a function of Y
 - It is a random variable
 - Has a distribution, mean, variance, etc.

Definition: $\underline{\mathbb{E}[X|Y]} = g(Y)$

The mean of $E[X | Y]$: Law of iterated expectations

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

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

$$E[X | Y] \stackrel{\Delta}{=} g(Y)$$

$$E[\underbrace{E[X | Y]}_{\text{}}] = E[g(Y)]$$

$$= \sum_y g(y) P_Y(y)$$

exp. value rule

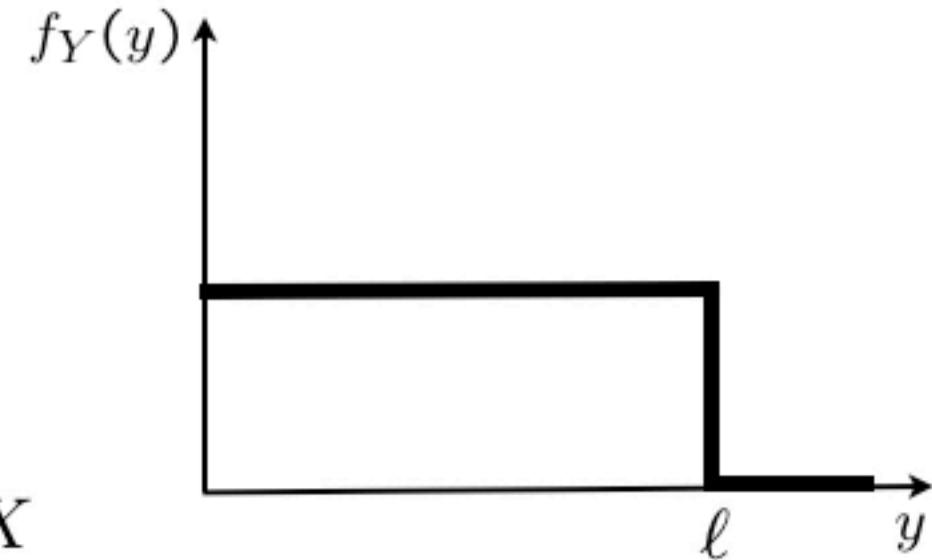
$$= \sum_y E[X | Y=y] P_Y(y)$$

• total exp thm

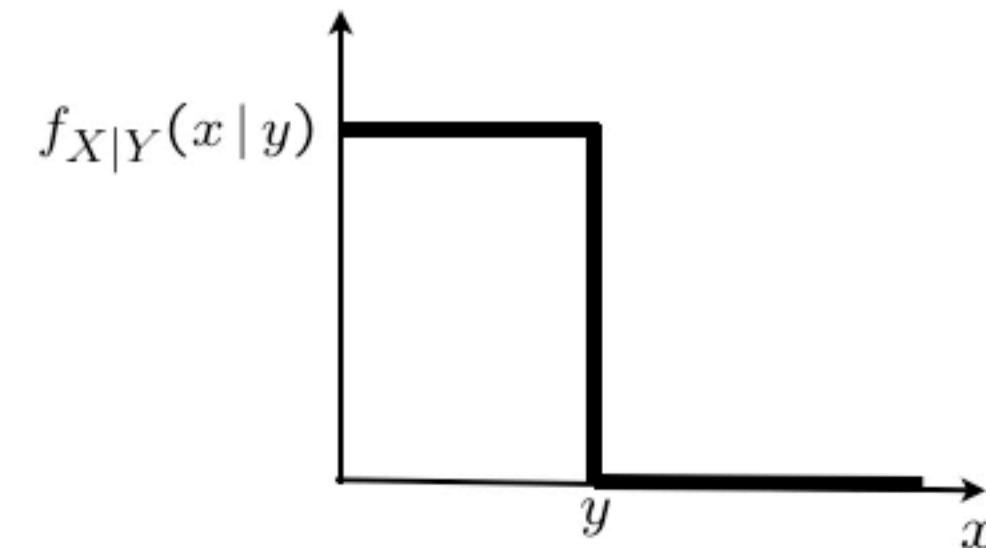
$$= E[X]$$

Stick-breaking example

- Stick example: stick of length ℓ
break at uniformly chosen point Y
break what is left at uniformly chosen point X



- $E[X | Y = y] = \frac{y}{2}$



- $E[X | Y] = \frac{\gamma}{2}$

$$E[X] = E[E[X | Y]] = E\left[\frac{Y}{2}\right] = \frac{1}{2} E[Y] = \frac{1}{2} \cdot \frac{\ell}{2} = \frac{\ell}{4}$$

Forecast revisions

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

- Suppose forecasts are made by calculating expected value, given any available information

- X : February sales

- Forecast in the beginning of the year: $E[X]$

- End of January: will get new information, value y of Y

Revised forecast:



$$E[X|Y=y]$$

$$E[X|Y]$$

- Law of iterated expectations:

$$E[\text{revised forecast}] = E[X] = \text{original forecast}$$

The conditional variance as a random variable

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

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

\uparrow

$\text{var}(X | Y)$ is the r.v. that takes the value $\text{var}(X | Y = y)$, when $Y = y$

- Example: X uniform on $[0, Y]$

$$\text{var}(X | Y = y) = \frac{y^2}{12}$$

$$\text{var}(X | Y) = \frac{Y^2}{12}$$

Law of total variance: $\text{var}(X) = \mathbb{E}[\text{var}(X | Y)] + \text{var}(\mathbb{E}[X | Y])$

Derivation of the law of total variance

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

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

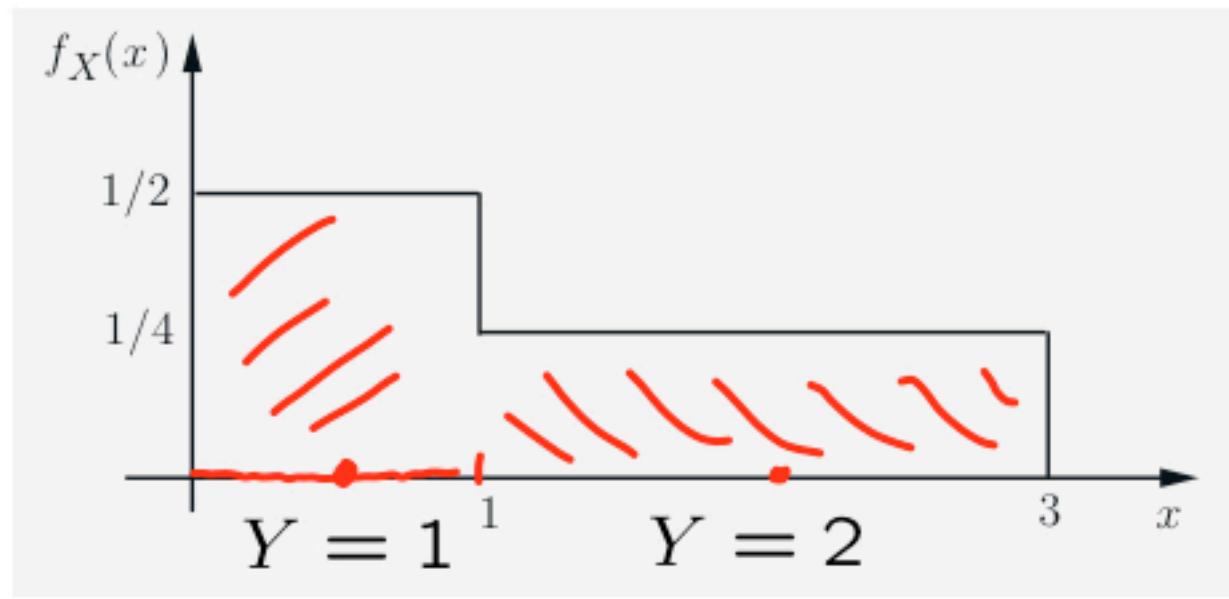
$$\text{var}(X | Y = y) = E[X^2 | Y = y] - (E[X | Y = y])^2 \text{ for all } y$$

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

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

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

A simple example



$$\begin{aligned} \text{var}(X) &= \mathbb{E}[\text{var}(X | Y)] + \text{var}(\mathbb{E}[X | Y]) = \frac{37}{48} \\ &= 5/24 + 9/16 \end{aligned}$$

$$\begin{aligned} \text{var}(X | Y) &= \frac{1/2}{1/2} \text{ var}(X | Y = 1) = \frac{1/12}{1/2} \\ &\quad \frac{1/2}{1/2} \text{ var}(X | Y = 2) = \frac{2^2/12}{1/2} = \frac{4}{12} \\ \mathbb{E}[\text{var}(X | Y)] &= \frac{1}{2} \cdot \frac{1}{12} + \frac{1}{2} \cdot \frac{4}{12} = \frac{5}{24} \end{aligned}$$

$$\begin{aligned} \mathbb{E}[X | Y] &= \frac{1/2}{1/2} \mathbb{E}[X | Y = 1] = \frac{1}{2} \\ &\quad \frac{1/2}{1/2} \mathbb{E}[X | Y = 2] = 2 \end{aligned}$$

$$\mathbb{E}[\mathbb{E}[X | Y]] = \frac{1}{2} \cdot \frac{1}{2} + \frac{1}{2} \cdot 2 = \frac{5}{4} = \mathbb{E}[X]$$

$$\begin{aligned} \text{var}(\mathbb{E}[X | Y]) &= \frac{1}{2} \left(\frac{1}{2} - \frac{5}{4} \right)^2 \\ &\quad + \frac{1}{2} \left(2 - \frac{5}{4} \right)^2 = \frac{9}{16} \end{aligned}$$

Section means and variances

- Two sections of a class: $y = 1$ (10 students); $y = 2$ (20 students)
 x_i : score of student i
- Experiment: pick a student at random (uniformly)
random variables: X and Y
- Data: $y = 1 : \frac{1}{10} \sum_{i=1}^{10} x_i = 90$ $y = 2 : \frac{1}{20} \sum_{i=11}^{30} x_i = 60$

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

$$E[X | Y = 1] = 90$$

$$E[X | Y = 2] = 60$$

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

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

Section means and variances (ctd.)

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

$$E[E[X | Y]] = 70 = E[X]$$

$$\text{var}(E[X | Y]) = \frac{1}{3}(90 - 70)^2 + \frac{2}{3}(60 - 70)^2 = 200$$

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

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

$$\text{var}(X | Y) = \frac{\cancel{1/3} \cdot 10}{\cancel{2/3} \cdot 20}$$

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

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

$$\text{var}(X) = E[\text{var}(X | Y)] + \text{var}(E[X | Y]) = \frac{50}{3} + 200$$

$\text{var}(X)$ = (average variability **within** sections) + (variability **between** sections)

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

$$E[Y] = E[N] \cdot E[X]$$

- N : number of stores visited (N is a nonnegative integer r.v.)
- Let $Y = X_1 + \dots + X_N$
- X_i : money spent in store i
 - X_i independent, identically distributed
 - independent of N

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

- Total expectation theorem:

$$E[Y] = \sum_n p_N(n) E[Y | N = n] = \underbrace{\sum_n p_N(n) n}_{\text{Total expectation}} E[X] = E[N] E[X]$$

- Law of iterated expectations:

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

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

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

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

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

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

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

$$\bullet \quad \text{var}(Y | N = n) = \text{var}(X_1 + \dots + X_n | N = n) = \text{var}(X_1 + \dots + X_n) \\ = n \text{var}(X)$$

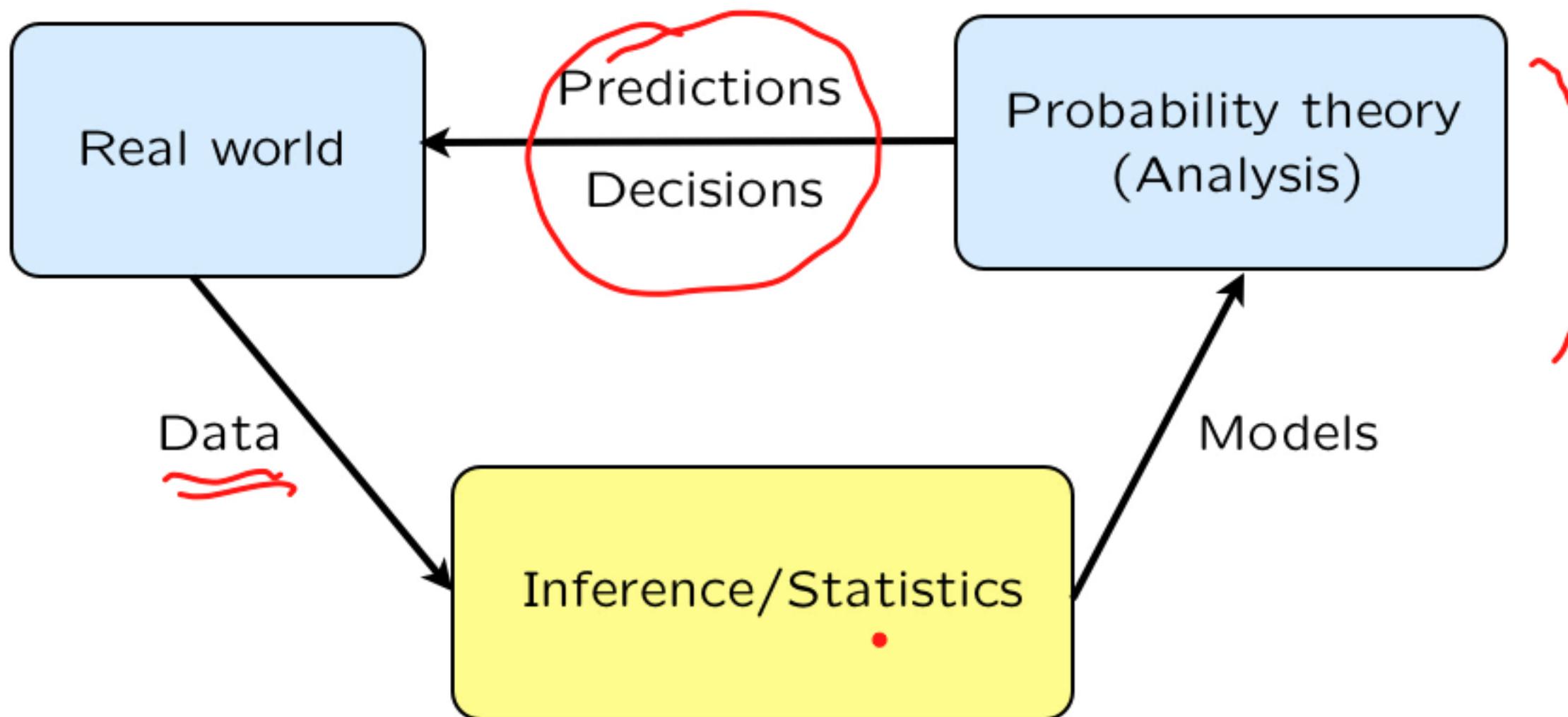
($\text{var}(Y | N) = N \text{var}(X)$)

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

LECTURE 14: Introduction to Bayesian inference

- The big picture
 - motivation, applications
 - problem types (hypothesis testing, estimation, etc.)
- The general framework
 - Bayes' rule → posterior
(4 versions)
 - point estimates (MAP, LMS)
 - performance measures)
(prob. of error; mean squared error)
 - examples

Inference: the big picture



Inference then and now

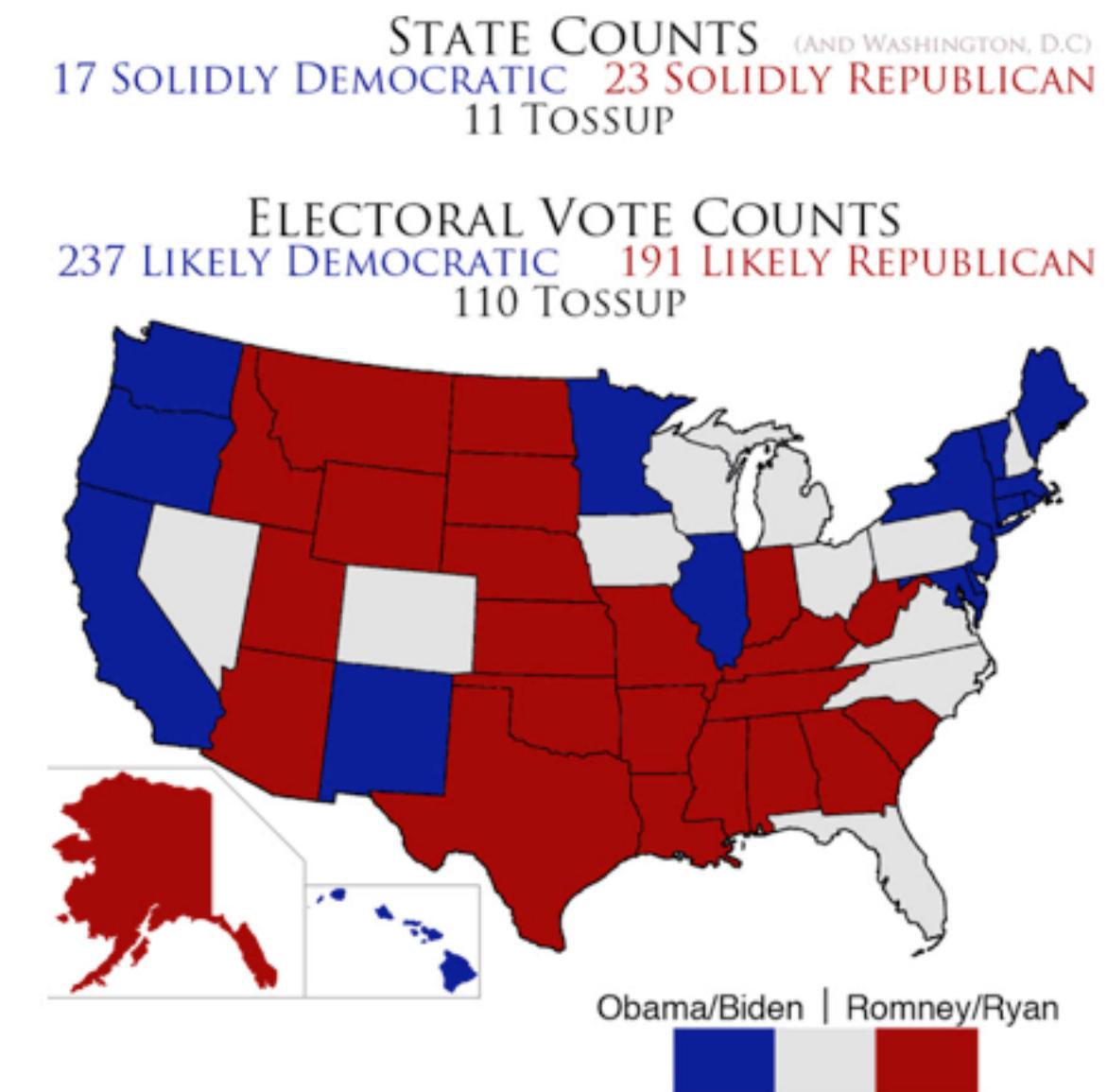
- Then:
 - 10 patients were treated: 3 died
 - 10 patients were not treated: 5 died
- Therefore ...

Now:

- Big data
- Big models
- Big computers

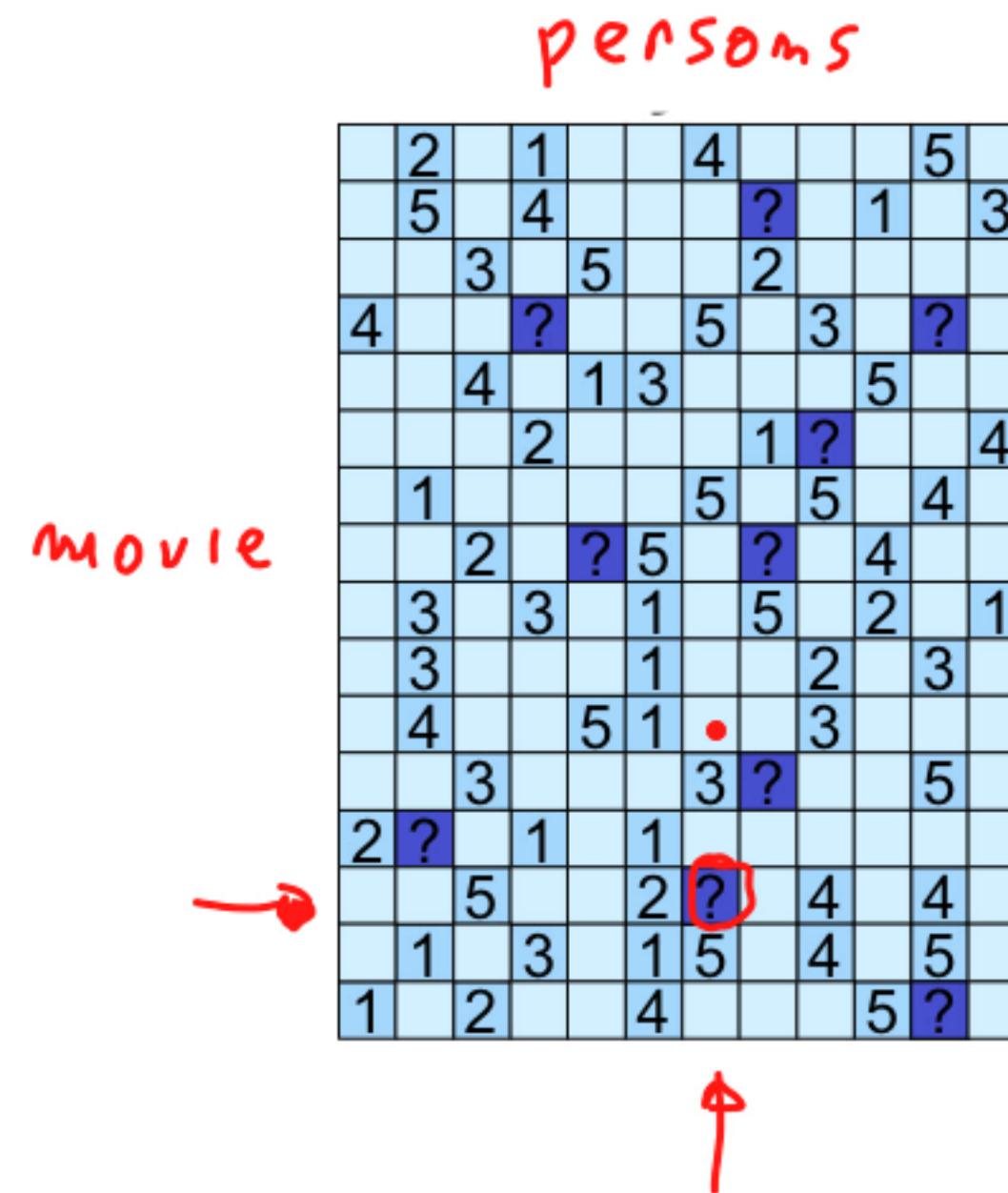
A sample of application domains

- Design and interpretation of experiments
 - polling •



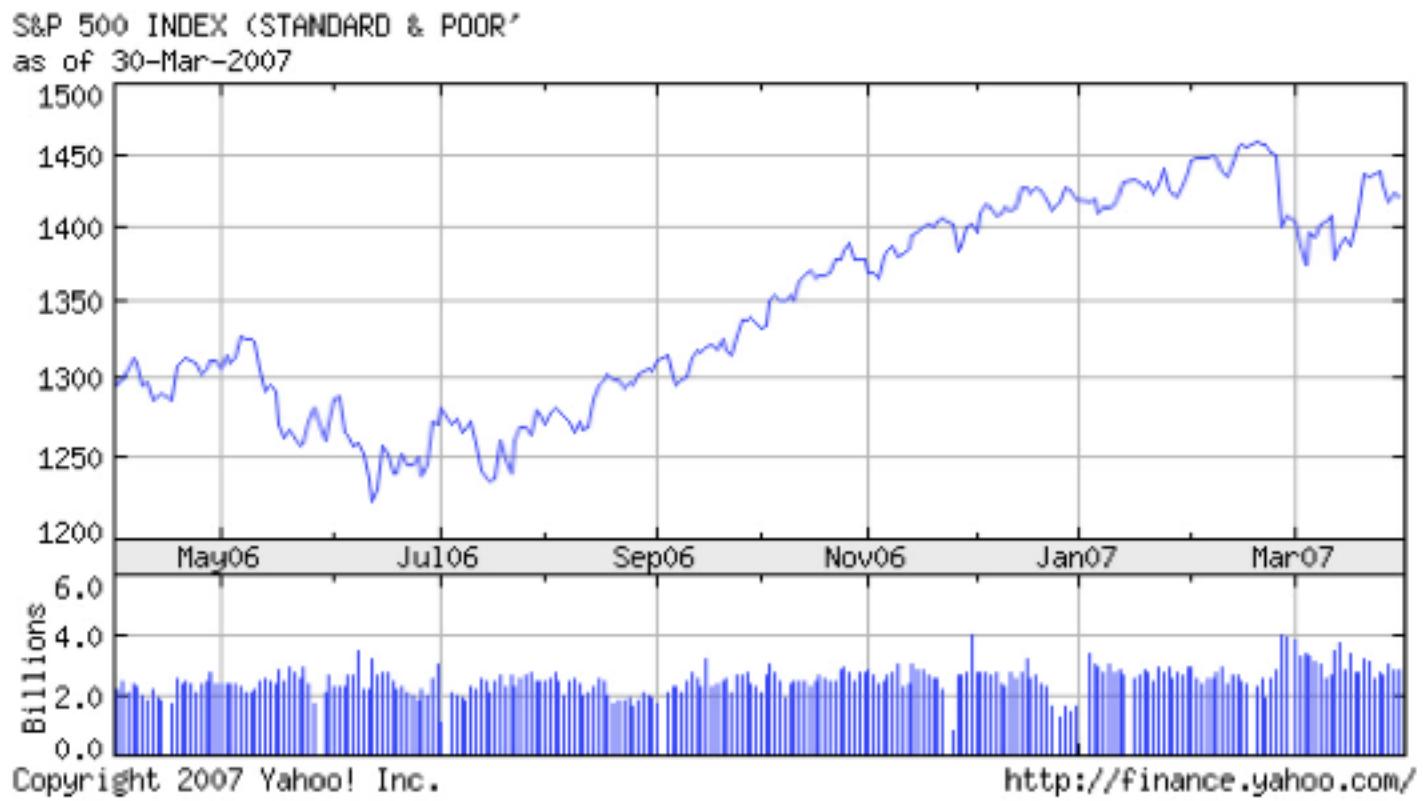
A sample of application domains

- marketing, advertising
- recommendation systems
 - Netflix competition



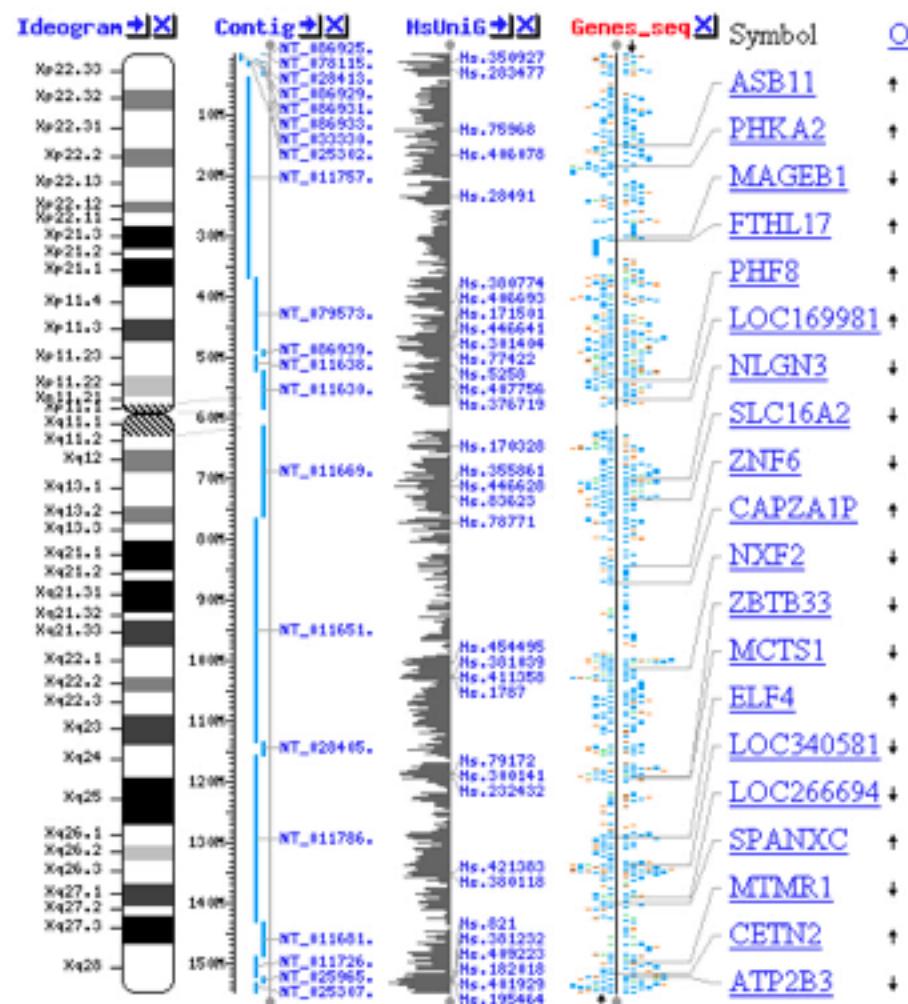
A sample of application domains

- Finance

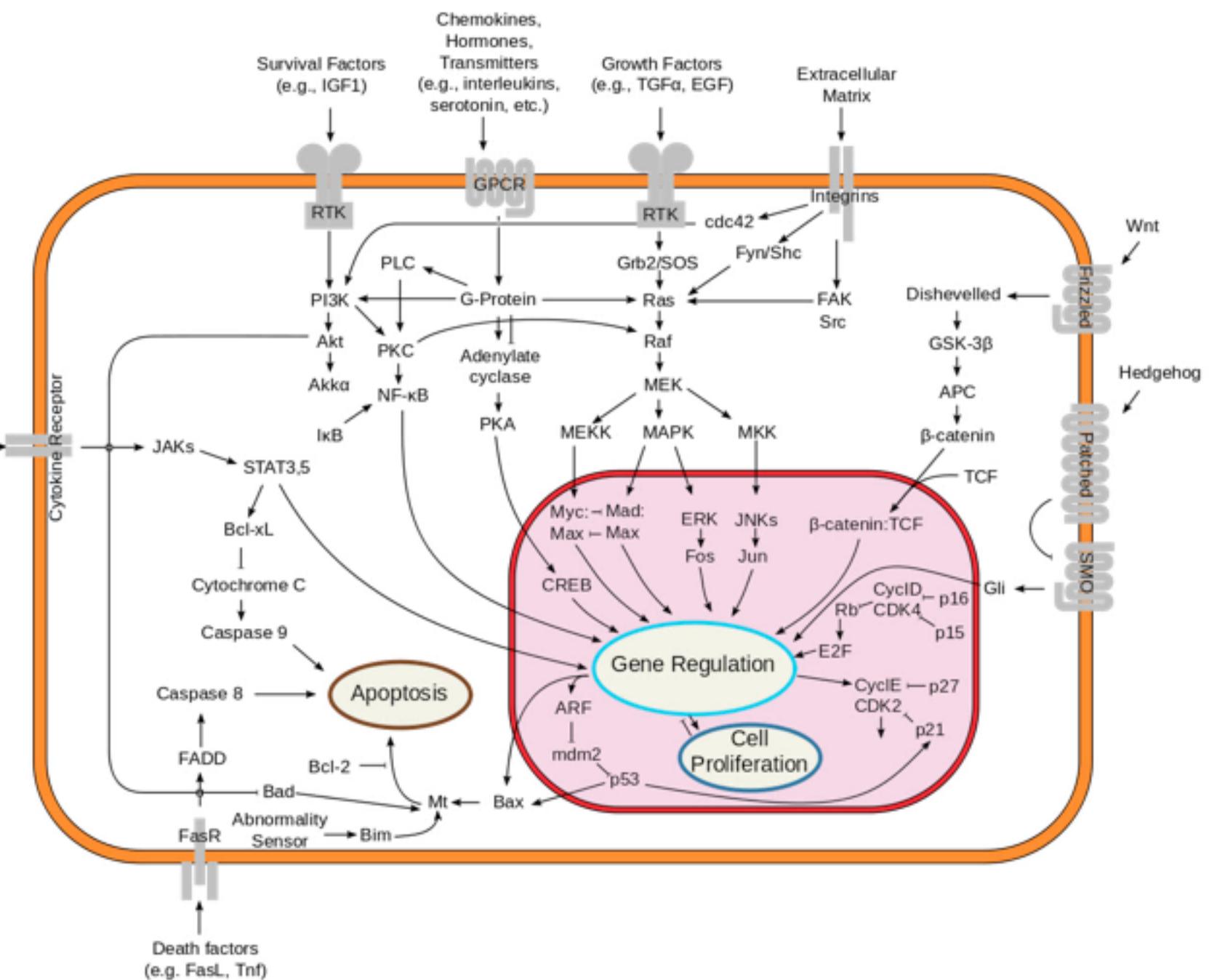


A sample of application domains

- Life sciences
 - genomics



— systems biology



- neuroscience, etc., etc.

A sample of application domains

- Modeling and monitoring the oceans
- Modeling and monitoring global climate
- Modeling and monitoring pollution
- Interpreting data from physics experiments •
- Interpreting astronomy data

A sample of application domains

- Signal processing
 - communication systems (noisy ...)
 - speech processing and understanding
 - image processing and understanding
 - tracking of objects
 - positioning systems (e.g., GPS)
 - detection of abnormal events

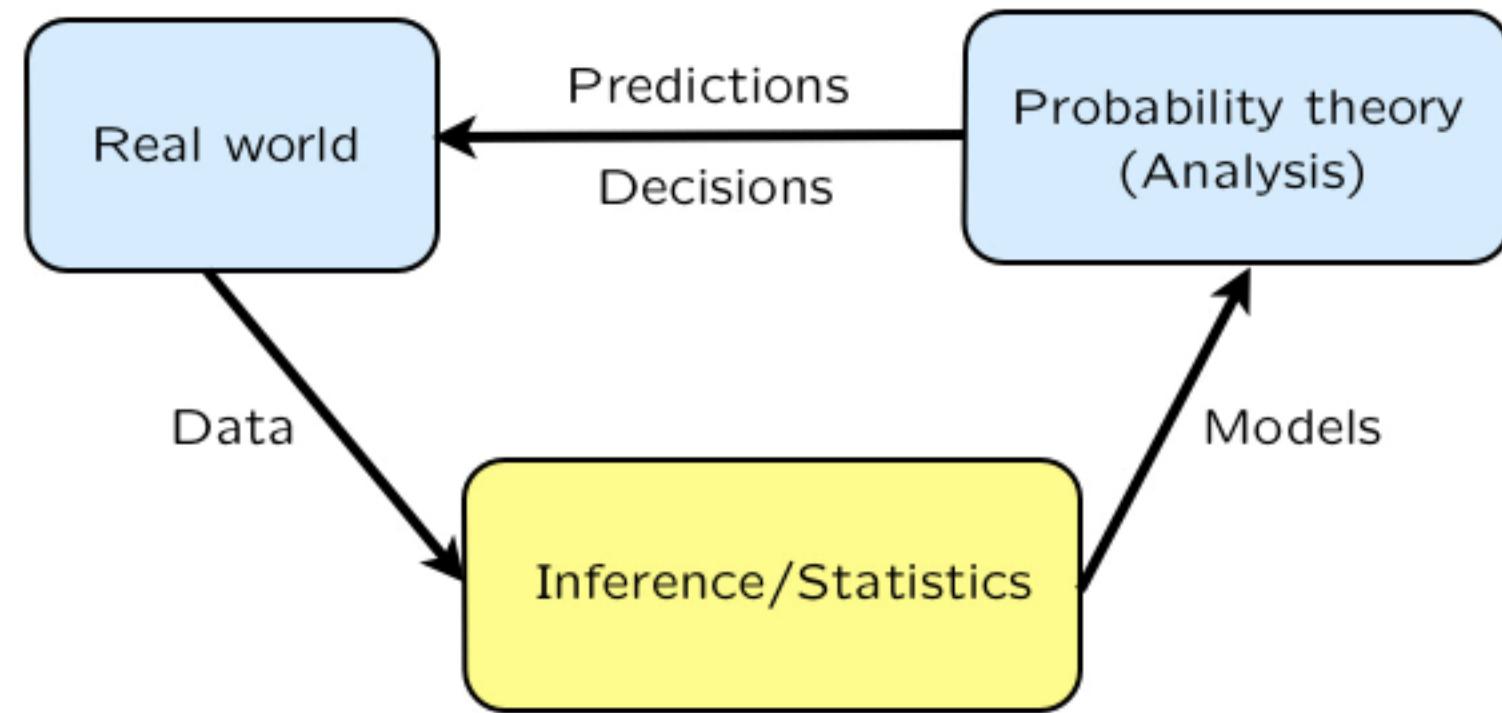
•

Model building versus inferring unobserved variables



$$X = aS + W$$

- Model building:
 - know “signal” S , observe X
 - infer a
- Variable estimation:
 - know a , observe X
 - infer S



Hypothesis testing versus estimation

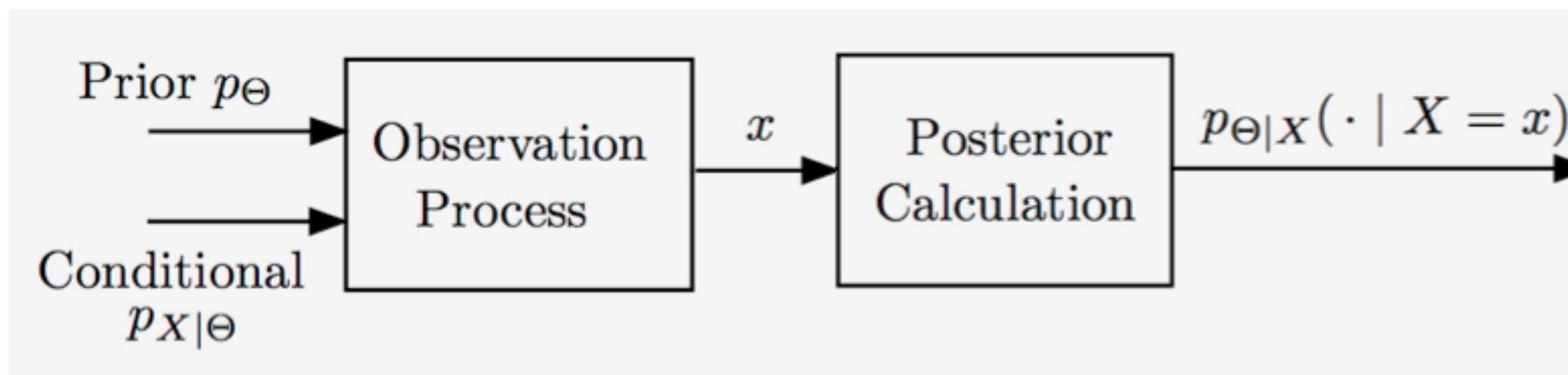
- Hypothesis testing:
 - unknown takes one of few possible values
 - aim at small probability of incorrect decision

Is it an airplane or a bird?

- Estimation:
 - numerical unknown(s)
 - aim at an estimate that is “close” to the true but unknown value

The Bayesian inference framework

- Unknown Θ
 - treated as a random variable
 - prior distribution p_Θ or f_Θ
- Observation X
 - observation model $p_{X|\Theta}$ or $f_{X|\Theta}$
- Use appropriate version of the Bayes rule to find $p_{\Theta|X}(\cdot | X = x)$ or $f_{\Theta|X}(\cdot | X = x)$
- Where does the prior come from?
 - symmetry
 - known range
 - earlier studies
 - subjective or arbitrary
-



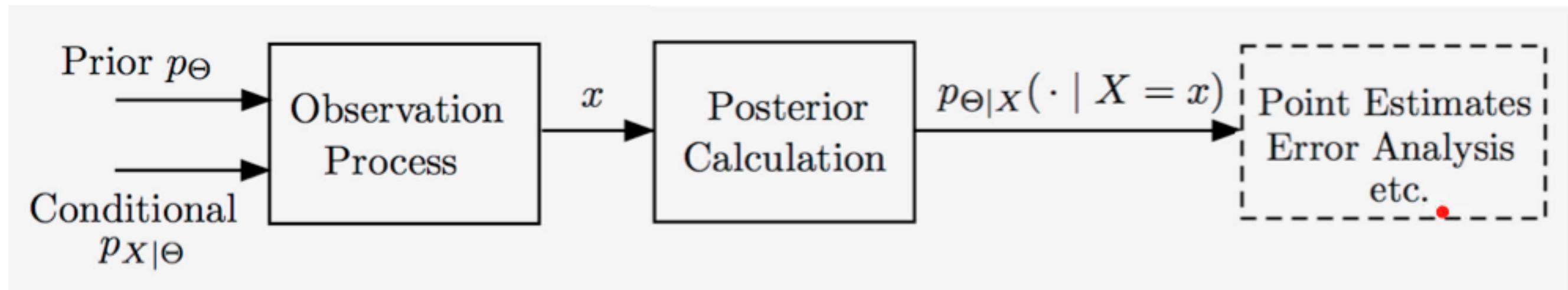
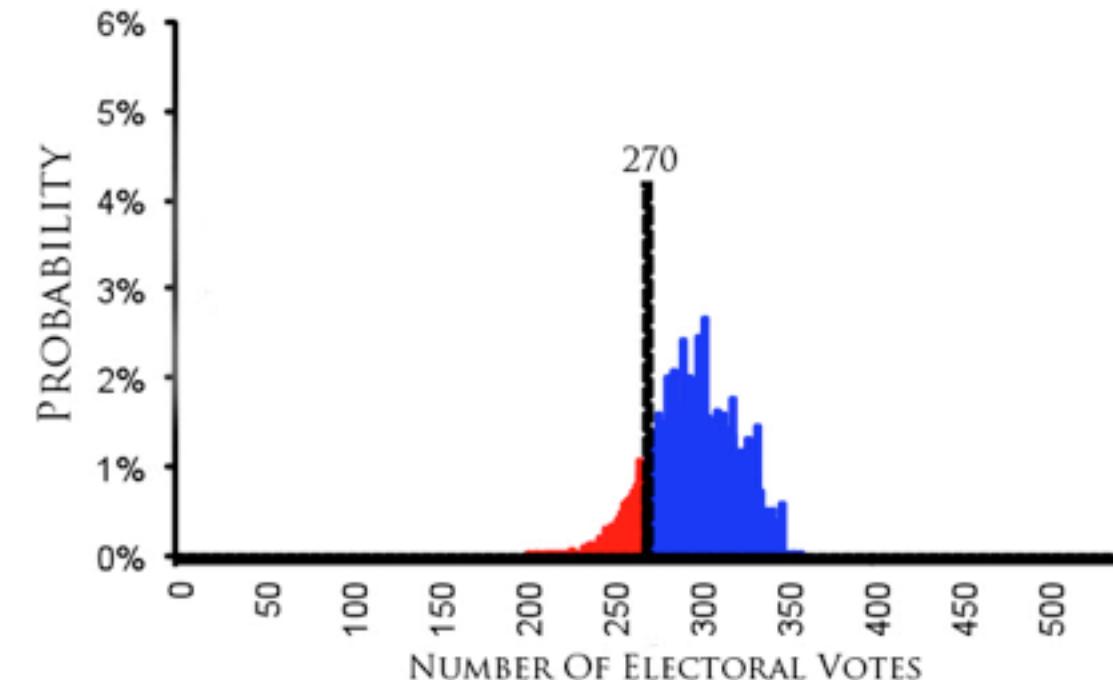
The output of Bayesian inference

The complete answer is a posterior distribution:

PMF $p_{\Theta|X}(\cdot | x)$ or PDF $f_{\Theta|X}(\cdot | x)$



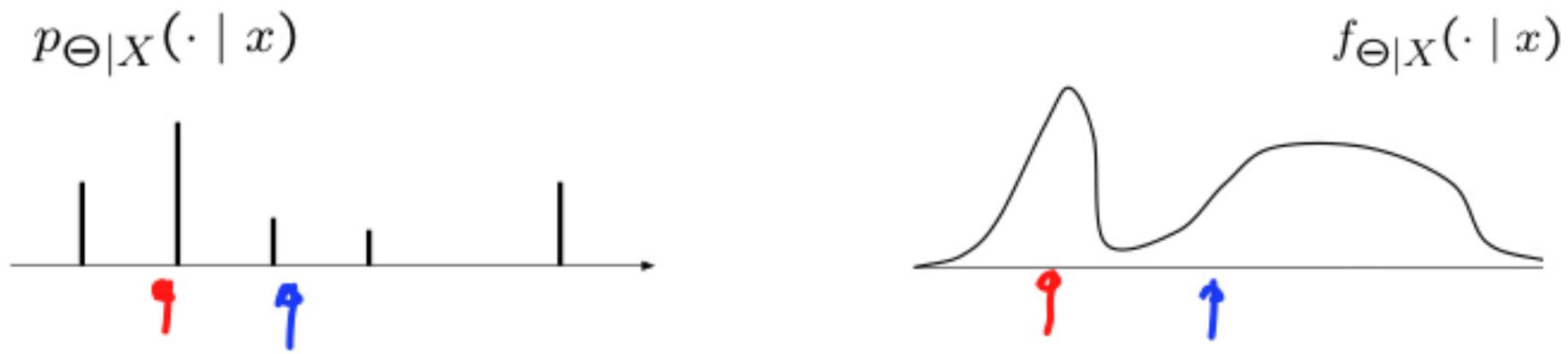
ELECTORAL VOTE DISTRIBUTION FOR OBAMA
ROMNEY 14.62% 84.59% OBAMA
0.79% TIE



Point estimates in Bayesian inference

The complete answer is a posterior distribution:

PMF $p_{\Theta|X}(\cdot | x)$ or PDF $f_{\Theta|X}(\cdot | x)$



- Maximum a posteriori probability (MAP):

$$p_{\Theta|X}(\theta^* | x) = \max_{\theta} p_{\Theta|X}(\theta | x),$$

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

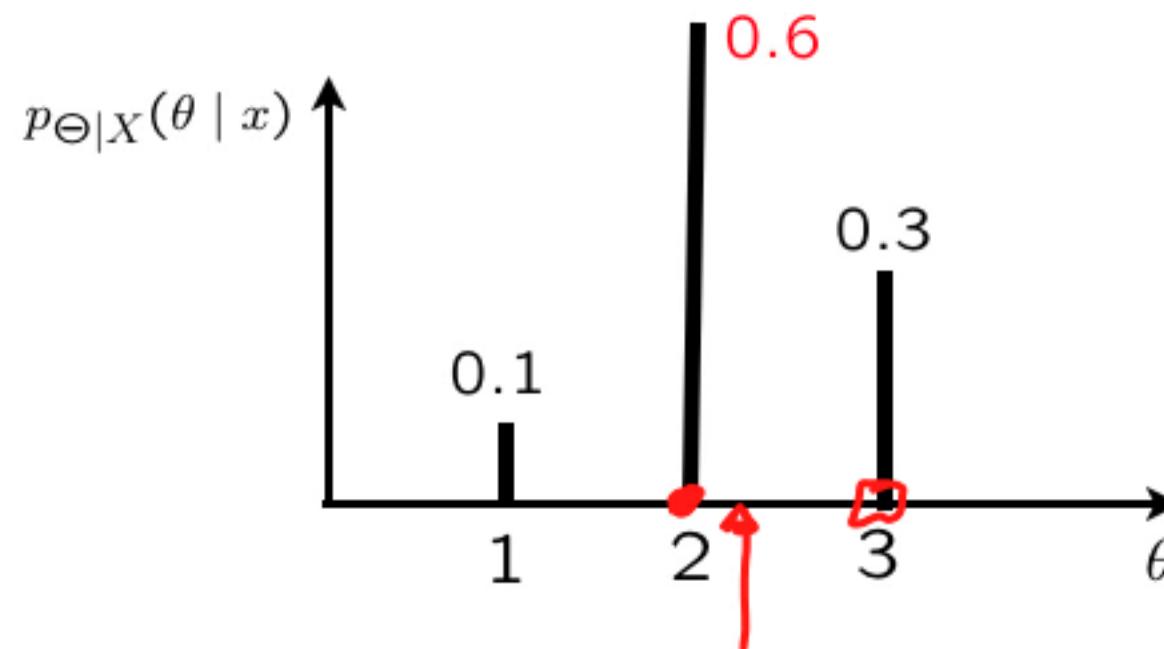
- Conditional expectation: $\mathbf{E}[\Theta | X = x]$ (LMS: Least Mean Squares)

estimate: $\hat{\theta} = g(x)$
(number)

estimator: $\widehat{\Theta} = g(X)$
(random variable)

Discrete Θ , discrete X

- values of Θ : alternative hypotheses



- MAP rule: $\hat{\theta} = 2$

$$LMS: \hat{\theta} = E[\Theta | X=x] = 2.2$$

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

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

- conditional prob of error:

$$P(\hat{\theta} \neq \Theta | X=x) = 0.4$$

smallest under the MAP rule

- overall probability of error:

$$P(\hat{\Theta} \neq \Theta) = \sum_x P(\hat{\Theta} \neq \Theta | X=x) p_X(x)$$

$$= \sum_{\theta} P(\hat{\Theta} \neq \Theta | \Theta = \theta) p_{\Theta}(\theta)$$

Discrete Θ , continuous X

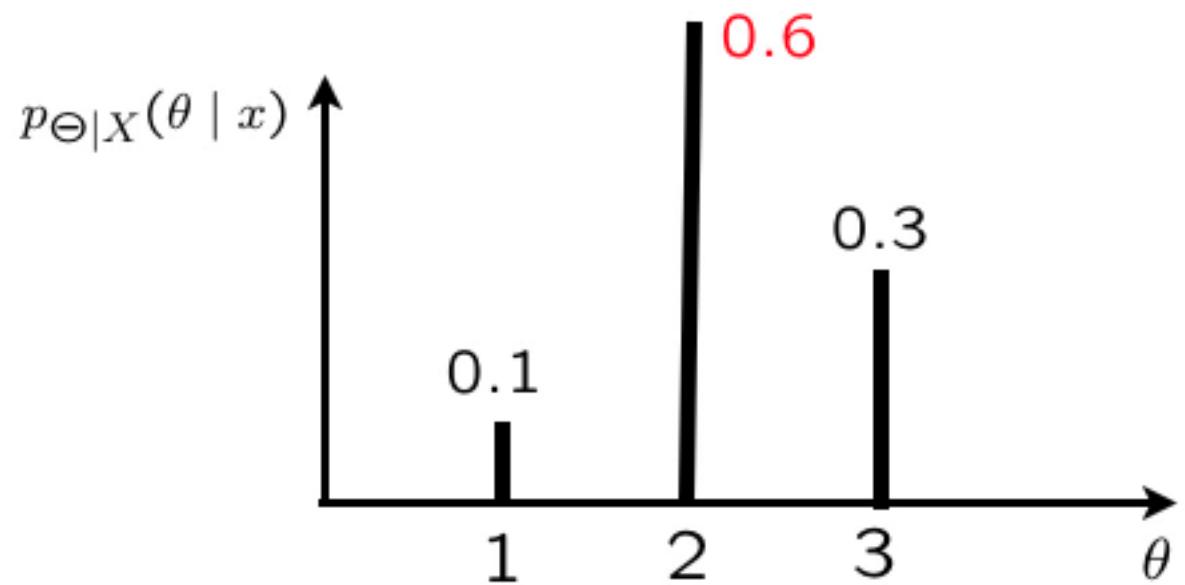
- Standard example:

- send signal $\Theta \in \{1, 2, 3\}$

$$X = \Theta + W$$

$W \sim N(0, \sigma^2)$, indep. of Θ

$$f_{X|\Theta}(x | \theta) = f_W(x - \theta)$$



- MAP rule: $\hat{\theta} = 2$

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

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

- conditional prob of error:

$$\mathbf{P}(\hat{\theta} \neq \Theta | X = x)$$

→ **smallest under the MAP rule** •

- overall probability of error:

$$\begin{aligned}\mathbf{P}(\hat{\Theta} \neq \Theta) &= \int \underbrace{\mathbf{P}(\hat{\Theta} \neq \Theta | X = x)}_{\text{overall probability}} f_X(x) dx \\ &= \sum_{\theta} \mathbf{P}(\hat{\Theta} \neq \theta | \Theta = \theta) p_{\Theta}(\theta)\end{aligned}$$

Continuous Θ , continuous X

- linear normal models
estimation of a noisy signal

$$X = \Theta + W$$

Θ and W : independent normals

multi-dimensional versions (many normal parameters, many observations)

- estimating the parameter of a uniform

X : uniform[0, Θ]

Θ : uniform [0, 1]

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

$$f_X(x) = \int f_{\Theta}(\theta') f_{X|\Theta}(x | \theta') d\theta'$$

- $\widehat{\Theta} = g(X)$ MAP
 LMS
- interested in:

$$\left\{ \begin{array}{l} E[(\widehat{\Theta} - \Theta)^2 | X = x] \\ E[(\widehat{\Theta} - \Theta)^2] \end{array} \right.$$

Inferring the unknown bias of a coin and the Beta distribution

- Standard example:
 - coin with bias Θ ; prior $f_\Theta(\cdot)$
 - fix n ; K = number of heads
- Assume $f_\Theta(\cdot)$ is uniform in $[0, 1]$

$$f_{\Theta|K}(\theta | k) = \frac{f_\Theta(\theta) p_{K|\Theta}(k | \theta)}{p_K(k)}$$

$$p_K(k) = \int f_\Theta(\theta') p_{K|\Theta}(k | \theta') d\theta'$$

$$f_{\Theta|K}(\theta | k) = \frac{1 \cdot \binom{n}{k} \theta^k (1-\theta)^{n-k}}{p_K(k)}$$

$$= \frac{1}{d(n, k)} \theta^k (1-\theta)^{n-k}$$

“Beta distribution, with parameters $(k+1, n-k+1)$ ”

$$\theta \in [0, 1]$$

\Rightarrow

- If prior is Beta: $f_\Theta(\theta) = \frac{1}{c} \theta^\alpha (1-\theta)^\beta$ $\alpha, \beta > 0$

$$f_{\Theta|K}(\theta | k) = \frac{\frac{1}{c} \theta^\alpha (1-\theta)^\beta \binom{n}{k} \theta^k (1-\theta)^{n-k}}{p_K(k)} = d \theta^{\alpha+k} (1-\theta)^{\beta+n-k}$$

Inferring the unknown bias of a coin: point estimates

- Standard example:
 - coin with bias Θ ; prior $f_\Theta(\cdot)$
 - fix n ; K = number of heads
- Assume $f_\Theta(\cdot)$ is uniform in $[0, 1]$

$$f_{\Theta|K}(\theta | k) = \frac{1}{d(n, k)} \underline{\theta^k (1-\theta)^{n-k}}$$

- MAP estimate:

$$\hat{\theta}_{\text{MAP}} = \boxed{k/n}$$

$$\max_{\theta} [k \log \theta + (n-k) \log(1-\theta)]$$

$$\frac{\partial}{\partial \theta} \frac{k/\theta + (n-k)/(1-\theta)}{=} 0$$

$$\widehat{\Theta}_{\text{MAP}} = \boxed{k/n}$$

$$\int_0^1 \theta^\alpha (1-\theta)^\beta d\theta = \frac{\alpha! \beta!}{(\alpha+\beta+1)!} \quad \begin{matrix} \alpha \geq 0 \\ \beta \geq 0 \end{matrix}$$

$$\begin{aligned} E[\Theta | K = k] &= \int_0^1 \theta f_{\Theta|K}(\theta | k) d\theta \\ &= \frac{1}{d(n, k)} \int_0^1 \theta^{k+1} (1-\theta)^{n-k} d\theta \\ &= \frac{1}{\cancel{k!} \cancel{(n-k)!} (n+1)!} \cdot \frac{(k+1)! (n-k)!}{(n+2)!} \\ &= \boxed{\frac{k+1}{n+2}} \approx \frac{k}{n} \quad (\text{if } n \text{ large}) \end{aligned}$$

Summary

- Problem data: $p_{\Theta}(\cdot)$, $p_{X|\Theta}(\cdot | \cdot)$
- Given the value x of X : find, e.g., $p_{\Theta|X}(\cdot | x)$
 - using appropriate version of the Bayes rule (4 choices)
- Estimator $\widehat{\Theta} = g(X)$ Estimate $\widehat{\theta} = g(x)$
 - MAP: $\widehat{\theta}_{\text{MAP}} = g_{\text{MAP}}(x)$ maximizes $p_{\Theta|X}(\theta | x)$
 - LMS: $\widehat{\theta}_{\text{LMS}} = g_{\text{LMS}}(x) = \mathbb{E}[\Theta | X = x]$
- Performance evaluation of an estimator $\widehat{\Theta}$
 - $P(\widehat{\Theta} \neq \Theta | X = x)$
 - $P(\widehat{\Theta} \neq \Theta)$
 - $E[(\widehat{\Theta} - \Theta)^2 | X = x]$
 - $E[(\widehat{\Theta} - \Theta)^2]$ • total \sum_{exp} prob } thru.

LECTURE 15: Linear models with normal noise

$$X_i = \sum_{j=1}^m a_{ij} \Theta_j + W_i$$

W_i, Θ_j : independent, normal

- Very common and convenient model
- Bayes' rule: normal posteriors
- MAP and LMS estimates coincide
 - simple formulas
(linear in the observations)
- Many nice properties
- Trajectory estimation example

Recognizing normal PDFs

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

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

$$2\alpha x + \beta = 0$$

$$c \cdot e^{-8(x-3)^2}$$

$$\mu = 3$$

$$\frac{1}{2\sigma^2} = 8 \Rightarrow \sigma^2 = \frac{1}{16}$$

$$c = \frac{1}{\frac{1}{4}\sqrt{2\pi}}$$

$$f_X(x) = c \cdot e^{-(\alpha x^2 + \beta x + \gamma)}$$

$$\alpha > 0$$

Normal with mean $-\beta/2\alpha$ and variance $1/2\alpha$

$$\alpha x^2 + \beta x + \gamma = \alpha \left(x^2 + \frac{\beta}{\alpha} x + \frac{\gamma}{\alpha} \right) = \alpha \left(\left(x + \frac{\beta}{2\alpha} \right)^2 - \frac{\beta^2}{4\alpha^2} + \frac{\gamma}{\alpha} \right)$$

$$f_X(x) = c e^{-\alpha \left(x + \frac{\beta}{2\alpha} \right)^2} e^{-\alpha \left(-\frac{\beta^2}{4\alpha^2} + \frac{\gamma}{\alpha} \right)}$$

$$\mu = -\frac{\beta}{2\alpha}$$

$$\frac{1}{2\sigma^2} = \alpha \Rightarrow \sigma^2 = 1/2\alpha$$

Estimating a normal random variable in the presence of additive normal noise

$$X = \Theta + W \quad \Theta, W : N(0, 1), \text{ independent}$$

$$f_{X|\Theta}(x|\theta) : X = \theta + W \quad N(\theta, 1)$$

$$f_{\Theta|X}(\theta|x) = \frac{1}{f_X(x)} c e^{-\frac{1}{2}\theta^2} c e^{-\frac{1}{2}(x-\theta)^2} = \underline{\underline{c(x)}} e^{-\text{quadratic}(\theta)}$$

$$\text{Fix } x \quad \min_{\theta} \left[\frac{1}{2}\theta^2 + \frac{1}{2}(x-\theta)^2 \right] \quad \begin{matrix} \text{Normals!} \\ \theta + (\theta - x) = 0 \end{matrix}$$

$$\hat{\theta}_{\text{MAP}} = \hat{\theta}_{\text{LMS}} = E[\Theta | X = x] = x/2$$

$$\widehat{\Theta}_{\text{MAP}} = E[\Theta | X] = x/2$$

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

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

Estimating a normal parameter in the presence of additive normal noise

$$X = \Theta + W \quad \Theta, W : N(0, 1), \text{ independent}$$

$$\widehat{\Theta}_{\text{MAP}} = \widehat{\Theta}_{\text{LMS}} = \mathbf{E}[\Theta | X] = \frac{X}{2}$$

- Even with general means and variances:
 - posterior is normal
 - LMS and MAP estimators coincide
 - these estimators are “linear,” of the form $\widehat{\Theta} = aX + b$ •

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

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

The case of multiple observations

$$X_1 = \Theta + W_1 \quad \Theta \sim N(x_0, \sigma_0^2) \quad W_i \sim N(0, \sigma_i^2)$$

$$\begin{matrix} \vdots \\ X_n = \Theta + W_n \end{matrix} \quad \Theta, W_1, \dots, W_n \text{ independent}$$

$$f_{X_i|\Theta}(x_i|\theta) = c_i e^{-\frac{(x_i - \theta)^2}{2\sigma_i^2}}$$

given $\Theta = \theta$: $X_i = \theta + W_i \sim N(\theta, \sigma_i^2)$

$$f_{X|\Theta}(x|\theta) = f_{X_1, \dots, X_n|\Theta}(x_1, \dots, x_n|\theta) = \prod_{i=1}^n f_{X_i|\Theta}(x_i|\theta)$$

given $\Theta = \theta$: W_i independent $\Rightarrow X_i$ independent

$$f_{\Theta|X}(\theta|x) = \frac{1}{f_X(x)} \cdot c_0 e^{-\frac{(\theta-x_0)^2}{2\sigma_0^2}} \prod_{i=1}^n c_i e^{-\frac{(x_i - \theta)^2}{2\sigma_i^2}}$$

Normal!

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

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

The case of multiple observations

$$f_{\Theta|X}(\theta|x) = c \cdot \exp \left\{ -\text{quad}(\theta) \right\} \quad \text{quad}(\theta) = \frac{(\theta - x_0)^2}{2\sigma_0^2} + \frac{(\theta - x_1)^2}{2\sigma_1^2} + \cdots + \frac{(\theta - x_n)^2}{2\sigma_n^2}$$

find peak

$$\frac{d}{d\theta} \text{quad}(\theta) = 0: \sum_{i=0}^n \frac{(\theta - x_i)}{\sigma_i^2} = 0 \Rightarrow \theta \sum_{i=0}^n \frac{1}{\sigma_i^2} = \sum_{i=0}^n \frac{x_i}{\sigma_i^2}$$

$$\hat{\theta}_{\text{MAP}} = \hat{\theta}_{\text{LMS}} = \mathbb{E}[\Theta | X = \underline{x}] = \frac{\sum_{i=0}^n \frac{x_i}{\sigma_i^2}}{\sum_{i=0}^n \frac{1}{\sigma_i^2}}$$

The case of multiple observations

- Key conclusions:
 - posterior is normal
 - LMS and MAP estimates coincide
 - these estimates are “linear,” of the form $\hat{\theta} = a_0 + a_1x_1 + \cdots + a_nx_n$
- Interpretations:
 - estimate $\hat{\theta}$: weighted average of x_0 (prior mean) and x_i (observations)
 - weights determined by variances

$$\hat{\theta}_{\text{MAP}} = \hat{\theta}_{\text{LMS}} = E[\Theta | X = x] = \frac{\sum_{i=0}^n \frac{x_i}{\sigma_i^2}}{\sum_{i=0}^n \frac{1}{\sigma_i^2}}$$

σ_i^2 large
 x_i very noisy
 \Rightarrow small weight

The mean squared error

$$f_{\Theta|X}(\theta | x) = c \cdot \exp \left\{ -\text{quad}(\theta) \right\}$$

$$X_i = \tilde{\theta} + w_i$$

$$\text{quad}(\theta) = \frac{(\theta - x_0)^2}{2\sigma_0^2} + \frac{(\theta - x_1)^2}{2\sigma_1^2} + \dots + \frac{(\theta - x_n)^2}{2\sigma_n^2}$$

$$\hat{\theta} = \frac{\sum_{i=0}^n \frac{x_i}{\sigma_i^2}}{\sum_{i=0}^n \frac{1}{\sigma_i^2}}$$

- Performance measures:

$$\mathbb{E}[(\Theta - \hat{\Theta})^2 | X = x] = \mathbb{E}[(\Theta - \hat{\theta})^2 | X = x] = \text{var}(\Theta | X = x) = \boxed{1 / \sum_{i=0}^n \frac{1}{\sigma_i^2}}$$

$$\mathbb{E}[(\Theta - \hat{\Theta})^2] = \boxed{\int \mathbb{E}[(\Theta - \hat{\theta})^2 | X = x] f_x(x) dx}$$

$$f_X(x) = c \cdot e^{-(\alpha x^2 + \beta x + \gamma)} \quad \alpha > 0 \quad \text{Normal with mean } -\beta/2\alpha \text{ and variance } 1/2\alpha$$

$$\alpha = \frac{1}{2\sigma_0^2} + \dots + \frac{1}{2\sigma_n^2}$$

some σ_i^2 small \rightarrow MSE small
all σ_i^2 large \rightarrow MSE large

The mean squared error

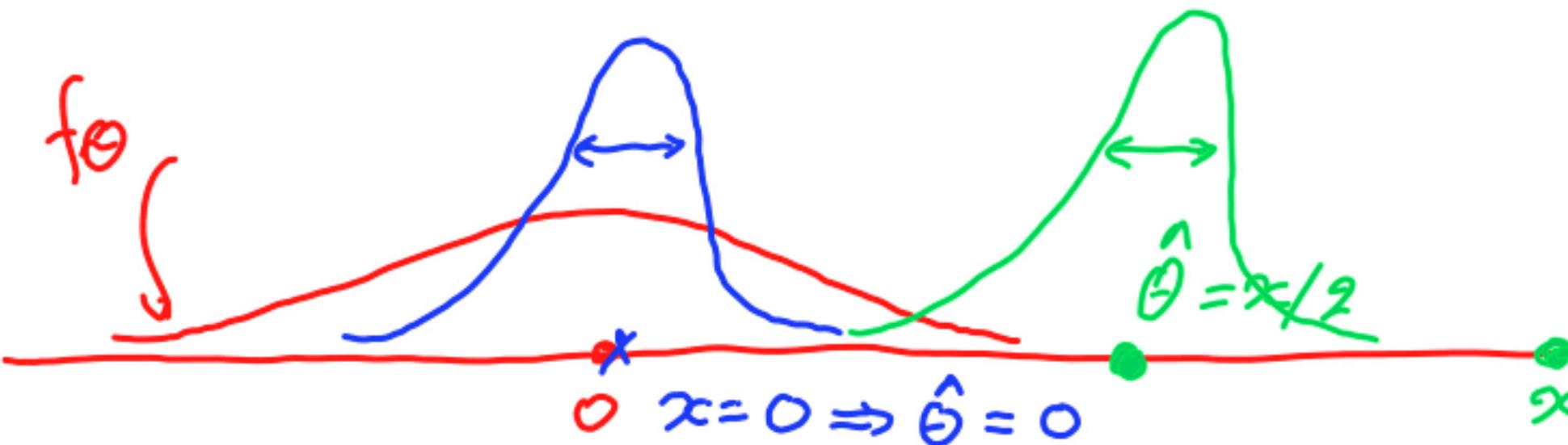
$$E[(\Theta - \hat{\Theta})^2 | X = \underline{x}] \underset{?}{=} E[(\Theta - \hat{\Theta})^2] = 1 / \sum_{i=0}^n \frac{1}{\sigma_i^2}$$

$$\hat{\theta} = \frac{\sum_{i=0}^n \frac{x_i}{\sigma_i^2}}{\sum_{i=0}^n \frac{1}{\sigma_i^2}}$$

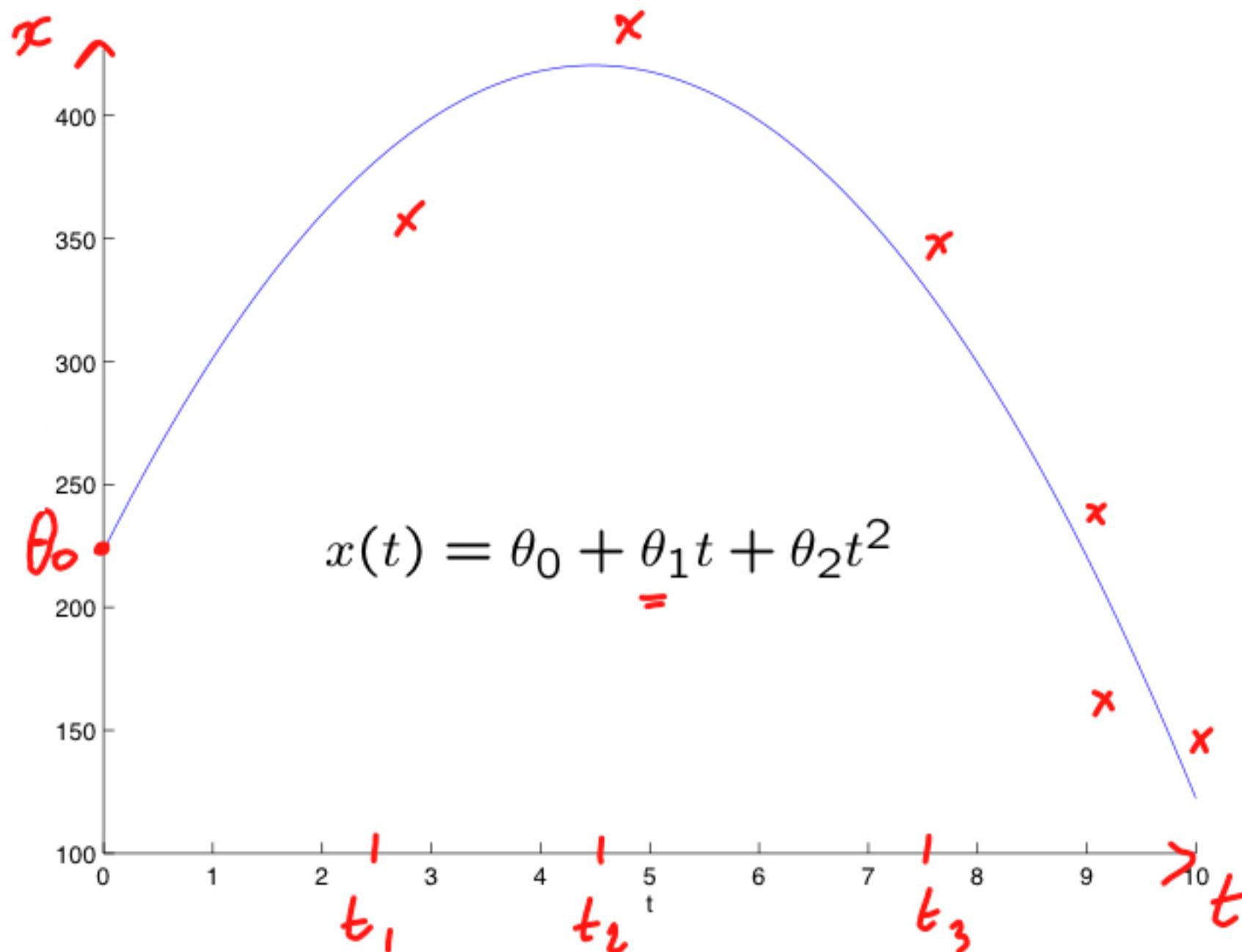
- Example: $\sigma_0^2 = \sigma_1^2 = \dots = \sigma_n^2 = \sigma^2$ $\frac{1}{(n+1) \frac{1}{\sigma^2}} = \frac{\sigma^2}{n+1}$
- conditional mean squared error same for all x
- Example: $X = \Theta + W$ $\Theta \sim N(0, \frac{1}{\sigma^2})$, $W \sim N(0, 1)$
independent Θ, W

$$\hat{\Theta} = X/2$$

$$E[(\Theta - \hat{\Theta})^2 | X = \underline{x}] = \underline{1/2}$$



The case of multiple parameters: trajectory estimation



- Random variables $\Theta_0, \Theta_1, \Theta_2$ independent; priors f_{Θ_j}
- Measurements at times t_1, \dots, t_n
$$X_i = \underline{\Theta_0 + \Theta_1 t_i + \Theta_2 t_i^2} + W_i$$
noise model: f_{W_i} independent W_i ; independent from Θ_j

A model with normality assumptions

$$X_i = \underline{\Theta_0 + \Theta_1 t_i + \Theta_2 t_i^2} + W_i \quad i = 1, \dots, n$$

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

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

- assume $\Theta_j \sim N(0, \sigma_j^2)$, $W_i \sim N(0, \sigma^2)$; independent

- Given $\Theta = \theta = (\theta_0, \theta_1, \theta_2)$, X_i is: $N(\theta_0 + \theta_1 t_i + \theta_2 t_i^2, \sigma^2)$

$$f_{X_i|\Theta}(x_i | \theta) = c \cdot \exp \left\{ - \underbrace{(x_i - \theta_0 - \theta_1 t_i - \theta_2 t_i^2)^2 / 2\sigma^2}_{\text{red arrow}} \right\}$$

- posterior: $f_{\Theta|X}(\theta | x) = \frac{1}{f_X(x)} \prod_{j=0}^2 f_{\Theta_j}(\theta_j) \prod_{i=1}^n f_{X_i|\Theta}(x_i | \theta)$

$$c(x) \exp \left\{ - \frac{1}{2} \left(\frac{\theta_0^2}{\sigma_0^2} + \frac{\theta_1^2}{\sigma_1^2} + \frac{\theta_2^2}{\sigma_2^2} \right) - \frac{1}{2\sigma^2} \sum_{i=1}^n \underbrace{(x_i - \theta_0 - \theta_1 t_i - \theta_2 t_i^2)^2}_{\text{red arrow}} \right\}$$

A model with normality assumptions

$$\underline{f_{\Theta|X}(\theta|x) = c(x) \exp \left\{ -\frac{1}{2} \left(\frac{\theta_0^2}{\sigma_0^2} + \frac{\theta_1^2}{\sigma_1^2} + \frac{\theta_2^2}{\sigma_2^2} \right) - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \theta_0 - \theta_1 t_i - \theta_2 t_i^2)^2 \right\}}$$

- MAP estimate: maximize over $(\theta_0, \theta_1, \theta_2)$;
(minimize quadratic function)

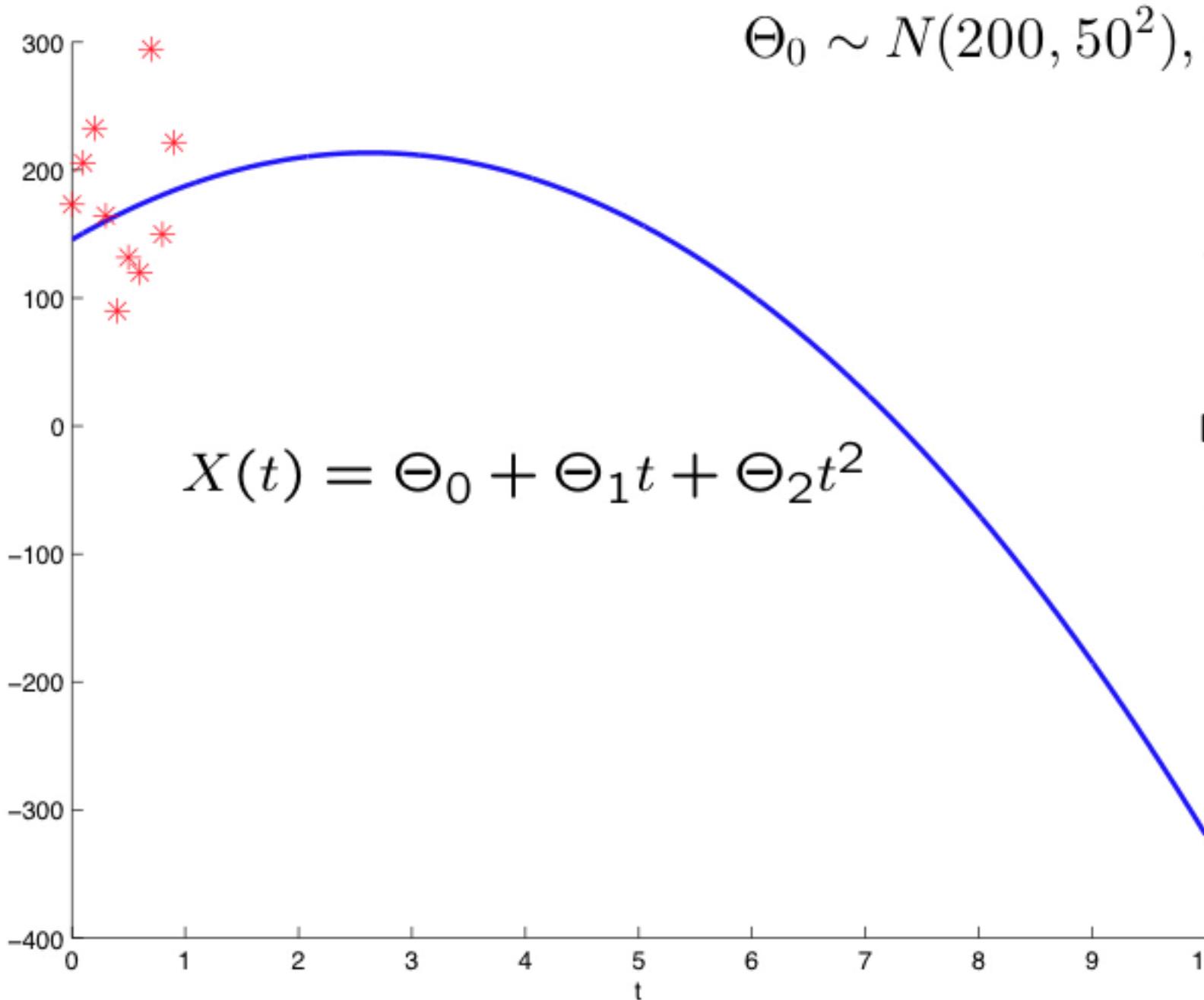
$$\frac{\partial}{\partial \theta_j} (\text{quad}(\theta)) = 0 \quad \begin{matrix} 3 \text{ equations, 3 unknowns} \\ \uparrow \text{linear} \end{matrix} .$$

Linear normal models •

- Θ_j and X_i are linear functions of independent normal random variables
- $f_{\Theta|X}(\theta | x) = c(x) \exp \left\{ -\text{quadratic}(\theta_1, \dots, \theta_m) \right\}$ *linear regression*
- MAP estimate: maximize over $(\theta_1, \dots, \theta_m)$; *linear equations*
(minimize quadratic function)
- $\widehat{\Theta}_{\text{MAP},j}$: linear function of $X = (X_1, \dots, X_n)$
- Facts:
 - $\widehat{\Theta}_{\text{MAP},j} = \mathbf{E}[\Theta_j | X]$
 - marginal posterior PDF of Θ_j : $f_{\Theta_j|X}(\theta_j | x)$, is normal
 - MAP estimate based on the joint posterior PDF:
same as MAP estimate based on the marginal posterior PDF
 - $\mathbf{E}[(\widehat{\Theta}_{i,\text{MAP}} - \Theta_i)^2 | X = x]$: same for all x

An illustration

Estimating the trajectory of a free-falling object



$$\Theta_0 \sim N(200, 50^2), \Theta_1 \sim N(50, 50^2), \Theta_2 = -9.81, W_i \sim N(0, 50^2)$$

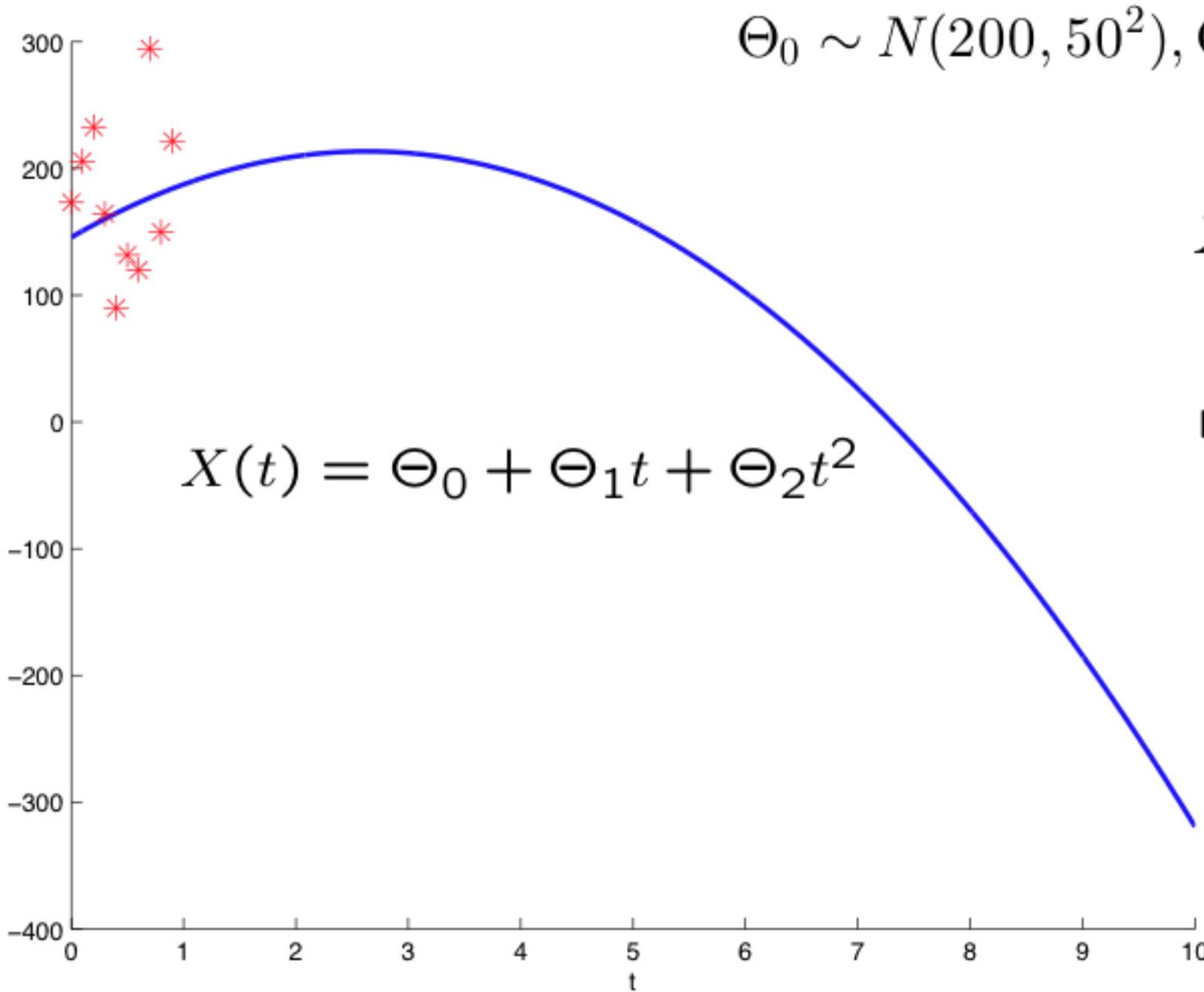
$$X_i = \Theta_0 + \Theta_1 t_i + \Theta_2 t_i^2 + W_i$$

$$\begin{aligned} & \text{minimize} \\ & \theta_0, \theta_1, \theta_2 \end{aligned}$$

$$\begin{aligned} & \frac{1}{2} \left(\frac{\theta_0^2}{\sigma_0^2} + \frac{\theta_1^2}{\sigma_1^2} + \frac{\theta_2^2}{\sigma_2^2} \right) \\ & + \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \theta_0 - \theta_1 t_i - \theta_2 t_i^2)^2 \end{aligned}$$

An illustration

Estimating the trajectory of a free-falling object



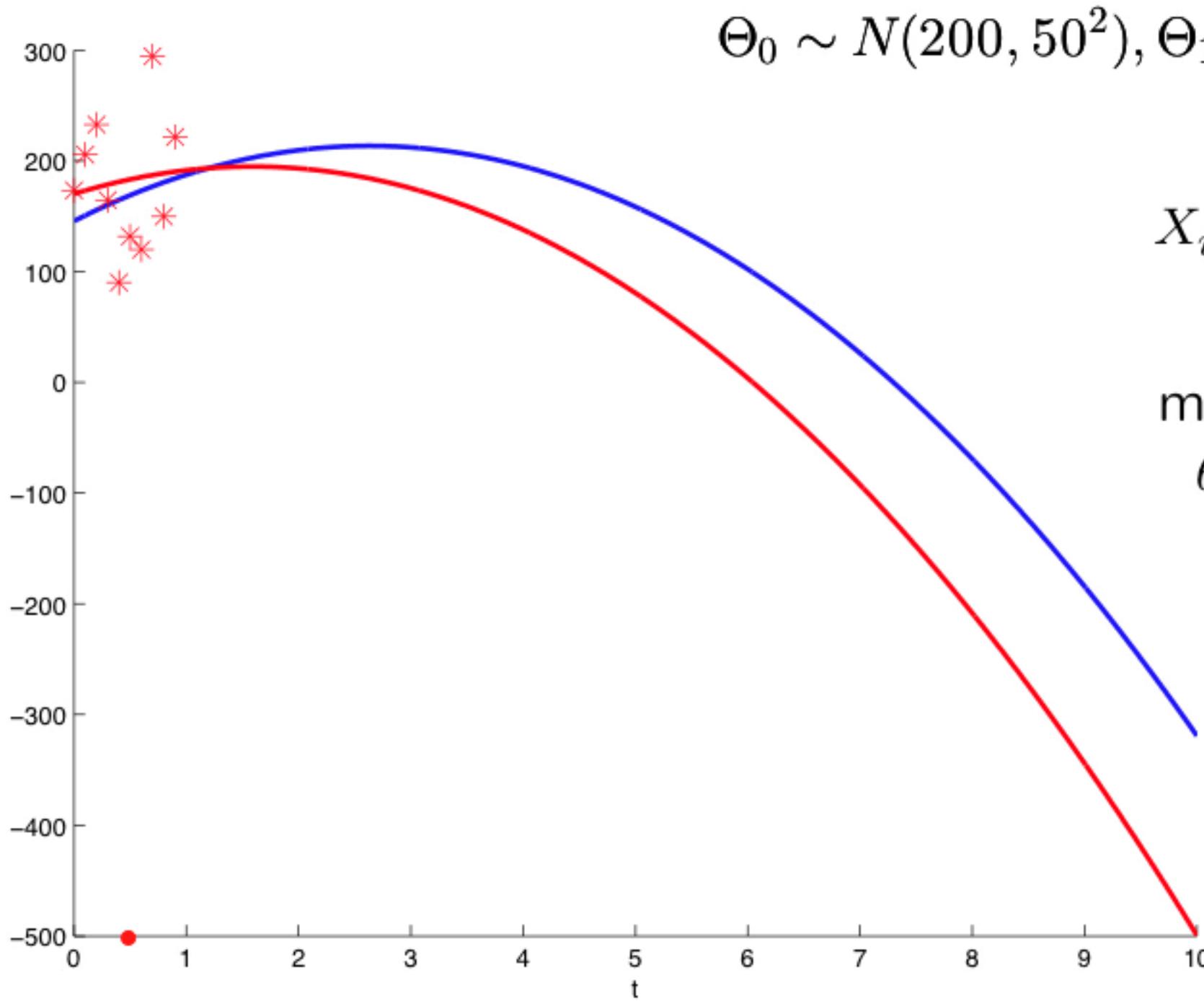
$$\Theta_0 \sim N(200, 50^2), \Theta_1 \sim N(50, 50^2), \Theta_2 = -9.81, W_i \sim N(0, 50^2)$$

$$X_i = \Theta_0 + \Theta_1 t_i + \Theta_2 t_i^2 + W_i$$

$$\begin{aligned} & \text{minimize}_{\theta_0, \theta_1} && (\theta_0 - 200)^2 + (\theta_1 - 50)^2 \\ & && + \sum_{i=1}^n (x_i - \theta_0 - \theta_1 t_i + 9.81 t_i^2)^2 \end{aligned}$$

An illustration

Estimating the trajectory of a free-falling object



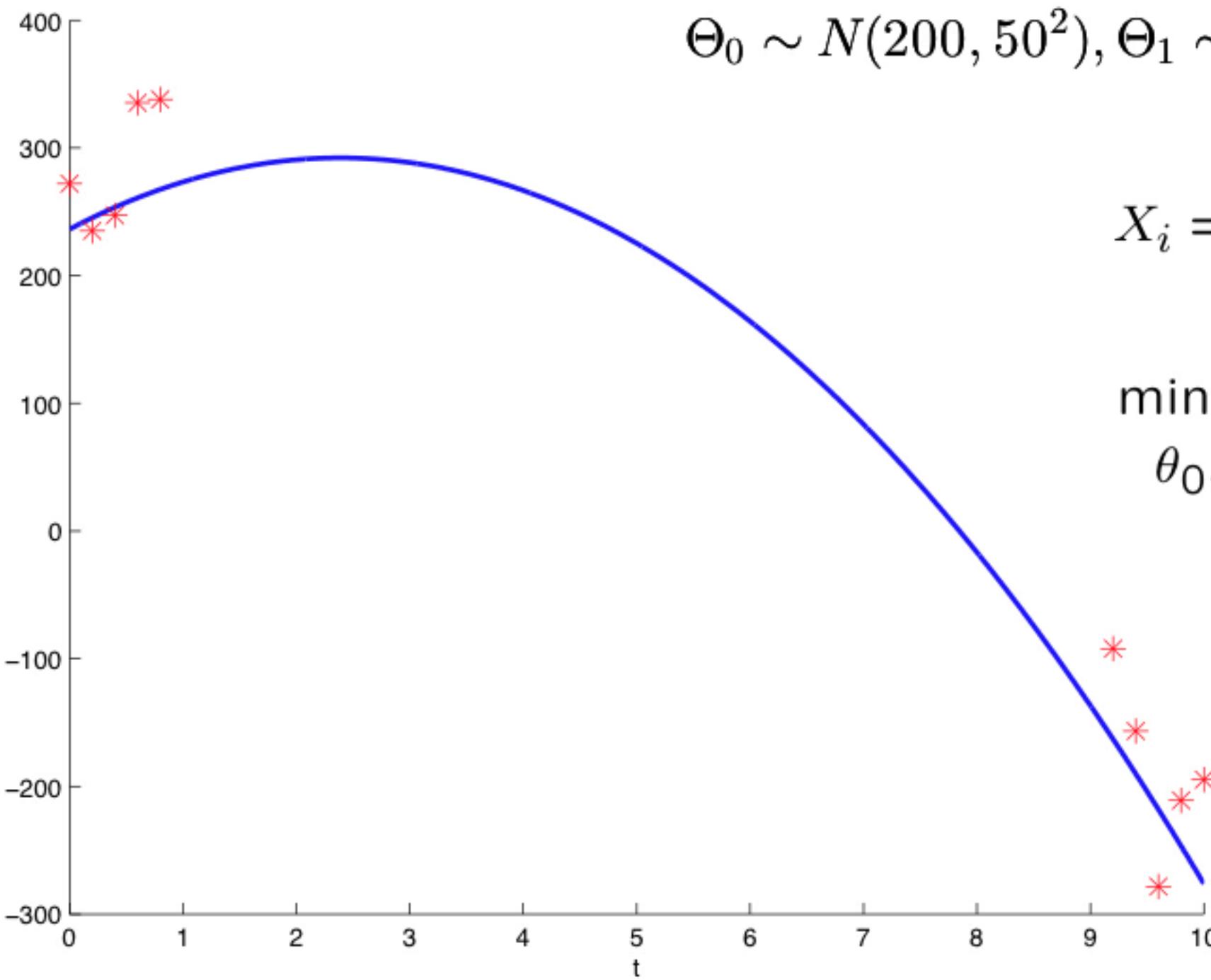
$$\Theta_0 \sim N(200, 50^2), \Theta_1 \sim N(50, 50^2), \Theta_2 = -9.81, W_i \sim N(0, 50^2)$$

$$X_i = \Theta_0 + \Theta_1 t_i + \Theta_2 t_i^2 + W_i$$

$$\begin{aligned} & \text{minimize}_{\theta_0, \theta_1} && (\theta_0 - 200)^2 + (\theta_1 - 50)^2 \\ & && + \sum_{i=1}^n (x_i - \theta_0 - \theta_1 t_i + 9.81 t_i^2)^2 \end{aligned}$$

An illustration

Estimating the trajectory of a free-falling object



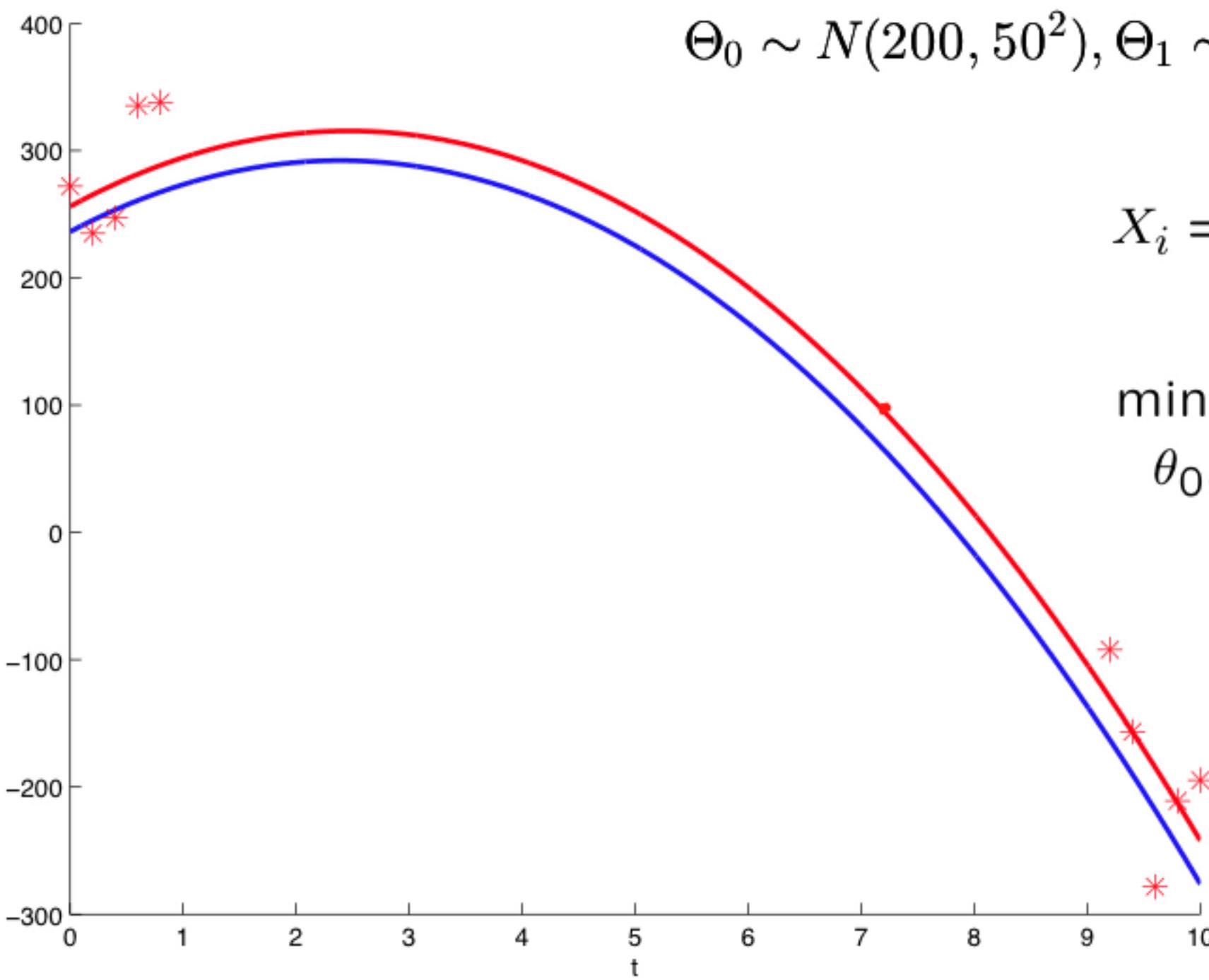
$$\Theta_0 \sim N(200, 50^2), \Theta_1 \sim N(50, 50^2), \Theta_2 = -9.81, W_i \sim N(0, 50^2)$$

$$X_i = \Theta_0 + \Theta_1 t_i + \Theta_2 t_i^2 + W_i$$

$$\begin{aligned} & \text{minimize}_{\theta_0, \theta_1} && (\theta_0 - 200)^2 + (\theta_1 - 50)^2 \\ & && + \sum_{i=1}^n (x_i - \theta_0 - \theta_1 t_i + 9.81 t_i^2)^2 \end{aligned}$$

An illustration

Estimating the trajectory of a free-falling object



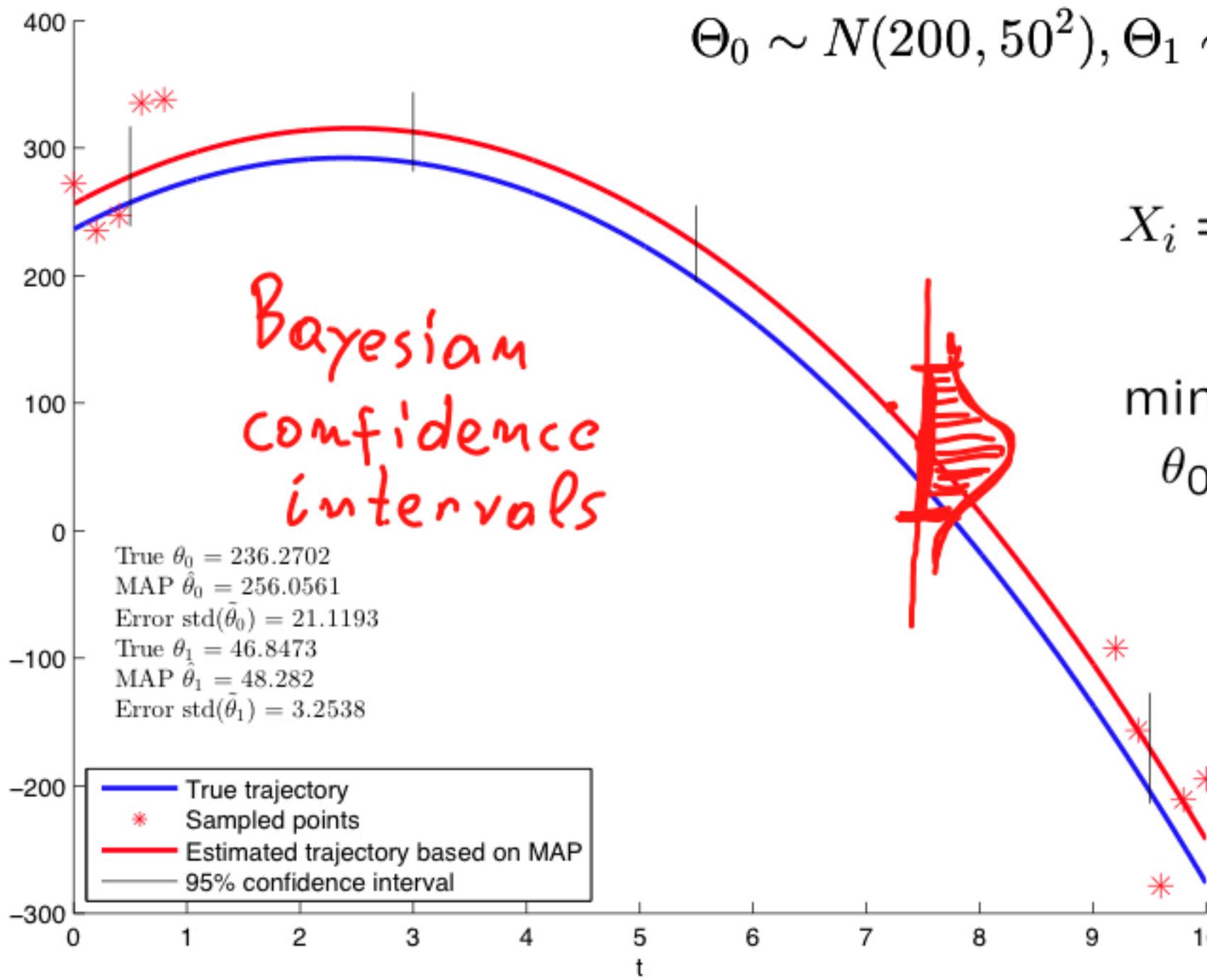
$$\Theta_0 \sim N(200, 50^2), \Theta_1 \sim N(50, 50^2), \Theta_2 = -9.81, W_i \sim N(0, 50^2)$$

$$X_i = \Theta_0 + \Theta_1 t_i + \Theta_2 t_i^2 + W_i$$

$$\begin{aligned} & \text{minimize}_{\theta_0, \theta_1} && (\theta_0 - 200)^2 + (\theta_1 - 50)^2 \\ & && + \sum_{i=1}^n (x_i - \theta_0 - \theta_1 t_i + 9.81 t_i^2)^2 \end{aligned}$$

An illustration

Estimating the trajectory of a free-falling object



$$\Theta_0 \sim N(200, 50^2), \Theta_1 \sim N(50, 50^2), \Theta_2 = -9.81, W_i \sim N(0, 50^2)$$

$$x(t) \\ X_i = \Theta_0 + \Theta_1 t_i + \Theta_2 t_i^2 + W_i$$

$$\begin{aligned} & \text{minimize}_{\theta_0, \theta_1} (\theta_0 - 200)^2 + (\theta_1 - 50)^2 \\ & + \sum_{i=1}^n (x_i - \theta_0 - \theta_1 t_i + 9.81 t_i^2)^2 \end{aligned}$$

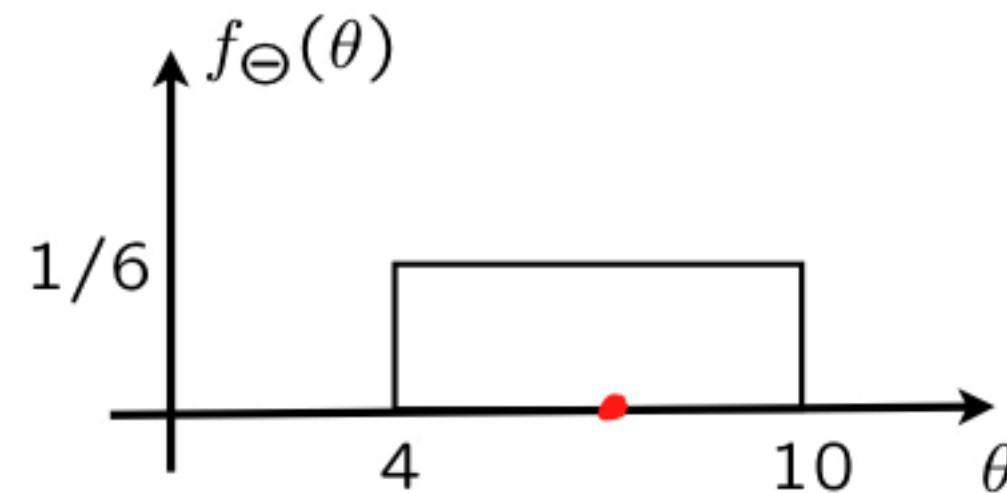
$$\begin{aligned} & P(x(t) \in \text{interval} | \text{data}) \\ & = 0.95 \end{aligned}$$

LECTURE 16: Least mean squares (LMS) estimation

- minimize (conditional) mean squared error $E[(\Theta - \hat{\theta})^2 | X = x]$
 - solution: $\hat{\theta} = E[\Theta | X = x]$
 - general estimation method
- Mathematical properties
- Example

LMS estimation in the absence of observations

- unknown Θ ; prior $p_\Theta(\theta)$
 - interested in a point estimate $\hat{\theta}$
 - no observations available
 - MAP rule: $\text{any } \hat{\theta} \in [4, 10]$
 - (Conditional) expectation: $\hat{\theta} = 7$
- Criterion: Mean Squared Error (MSE): $E[(\Theta - \hat{\theta})^2]$
minimize mean squared error

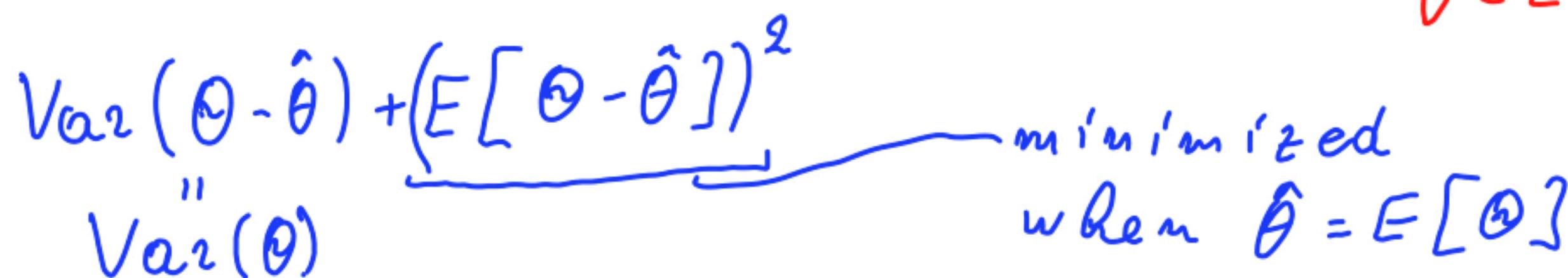


LMS estimation in the absence of observations

- Least mean squares formulation:

minimize mean squared error (MSE), $E[(\Theta - \hat{\theta})^2]$: $\hat{\theta} = E[\Theta]$.

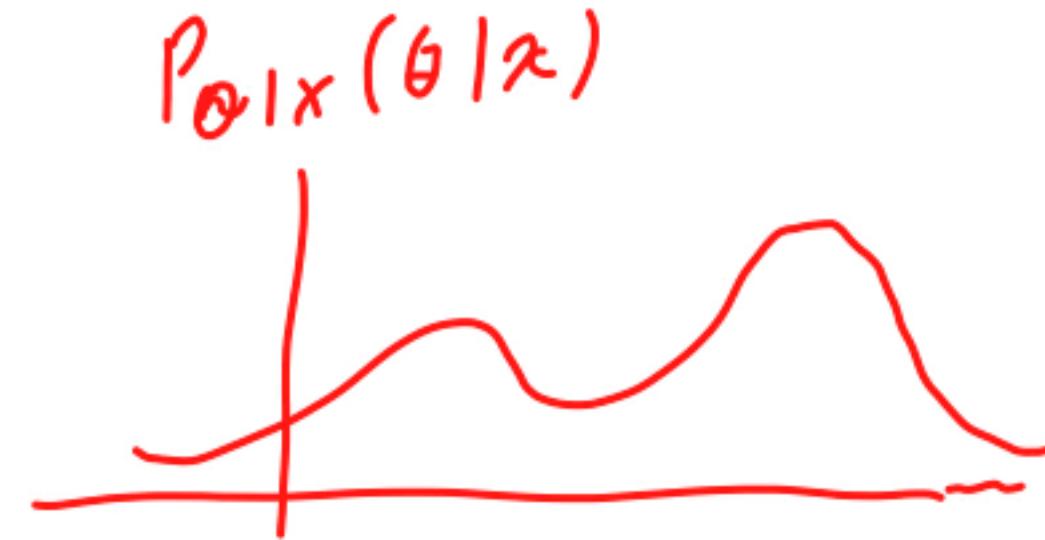
$$E[\Theta^2] - 2E[\Theta]\hat{\theta} + \hat{\theta}^2 \quad \frac{d}{d\hat{\theta}} = 0 : -2E[\Theta] + 2\hat{\theta} = 0$$
$$\hat{\theta} = E[\Theta]$$



- Optimal mean squared error: $E[(\Theta - E[\Theta])^2] = \text{var}(\Theta)$

LMS estimation of Θ based on X

- unknown Θ ; prior $p_\Theta(\theta)$
 - interested in a point estimate $\hat{\theta}$
- observation X ; model $p_{X|\Theta}(x|\theta)$
 - observe that $X = x$



minimize mean squared error (MSE), $E[(\Theta - \hat{\theta})^2]$: $\hat{\theta} = E[\Theta]$

minimize conditional mean squared error, $E[(\Theta - \hat{\theta})^2 | X = x]$: $\hat{\theta} = E[\Theta | X = x]$

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

estimator: $\widehat{\Theta} = E[\Theta | X]$

LMS estimation of Θ based on X

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

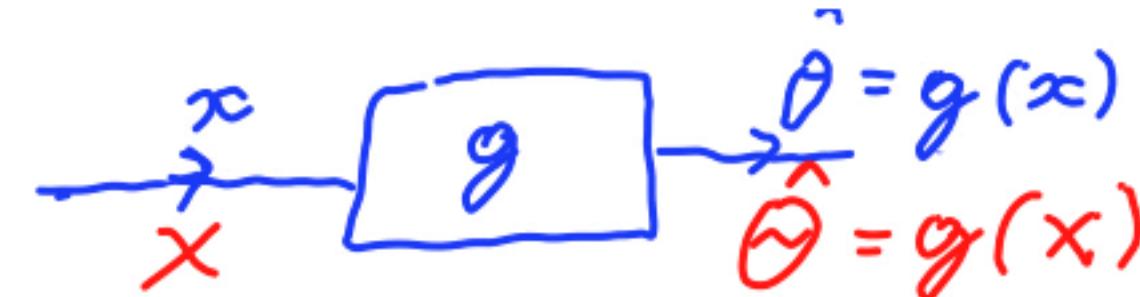
$$E[(\Theta - E[\Theta])^2] \leq E[(\Theta - c)^2], \text{ for all } c$$

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

$$E[(\Theta - E[\Theta | X = x])^2 | X = x] \leq E[(\Theta - g(x))^2 | X = x] \text{ for all } x$$

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

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



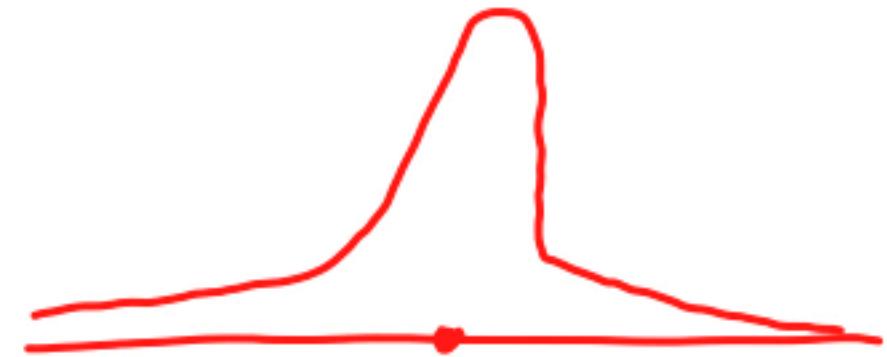
$\widehat{\Theta}_{\text{LMS}} = E[\Theta | X]$ minimizes $E[(\Theta - g(X))^2]$, over all estimators $\widehat{\Theta} = g(X)$

LMS performance evaluation

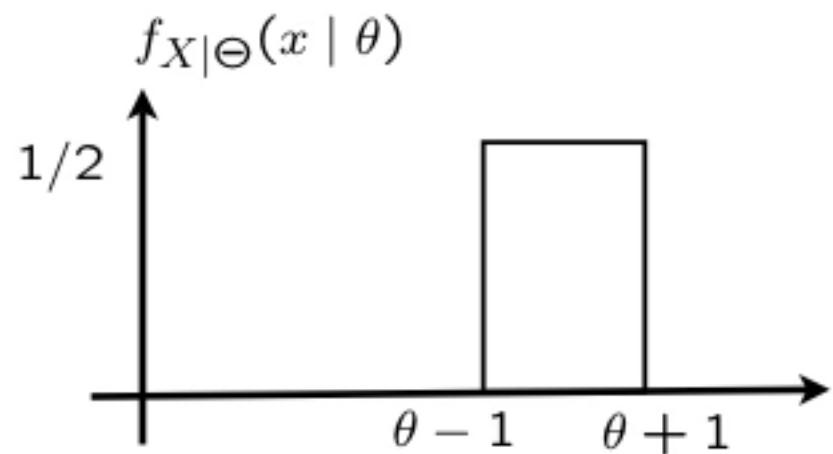
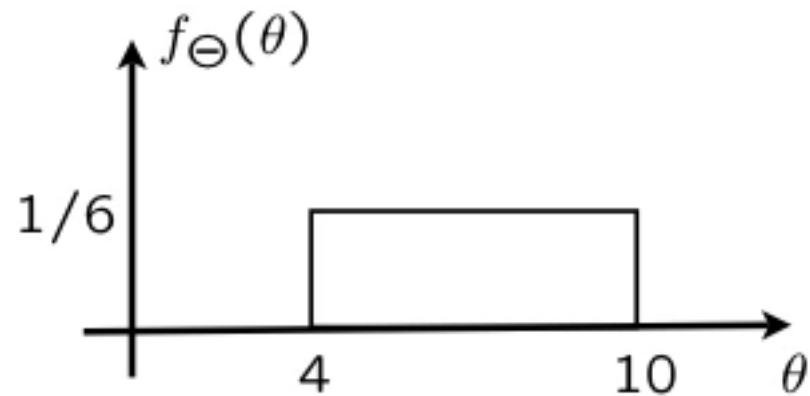
- LMS estimate: $\hat{\theta} = \mathbb{E}[\Theta | X = x]$
estimator: $\widehat{\Theta} = \mathbb{E}[\Theta | X]$
- Expected performance, once we have a measurement:
$$\text{MSE} = \mathbb{E}\left[\left(\underline{\Theta} - \mathbb{E}[\Theta | X = x]\right)^2 | X = x\right] = \underline{\text{var}(\Theta | X = x)}$$
- Expected performance of the design:
$$\text{MSE} = \mathbb{E}\left[\left(\Theta - \mathbb{E}[\Theta | X]\right)^2\right] = \mathbb{E}\left[\underline{\text{var}(\Theta | X)}\right]$$

LMS estimation of Θ based on X

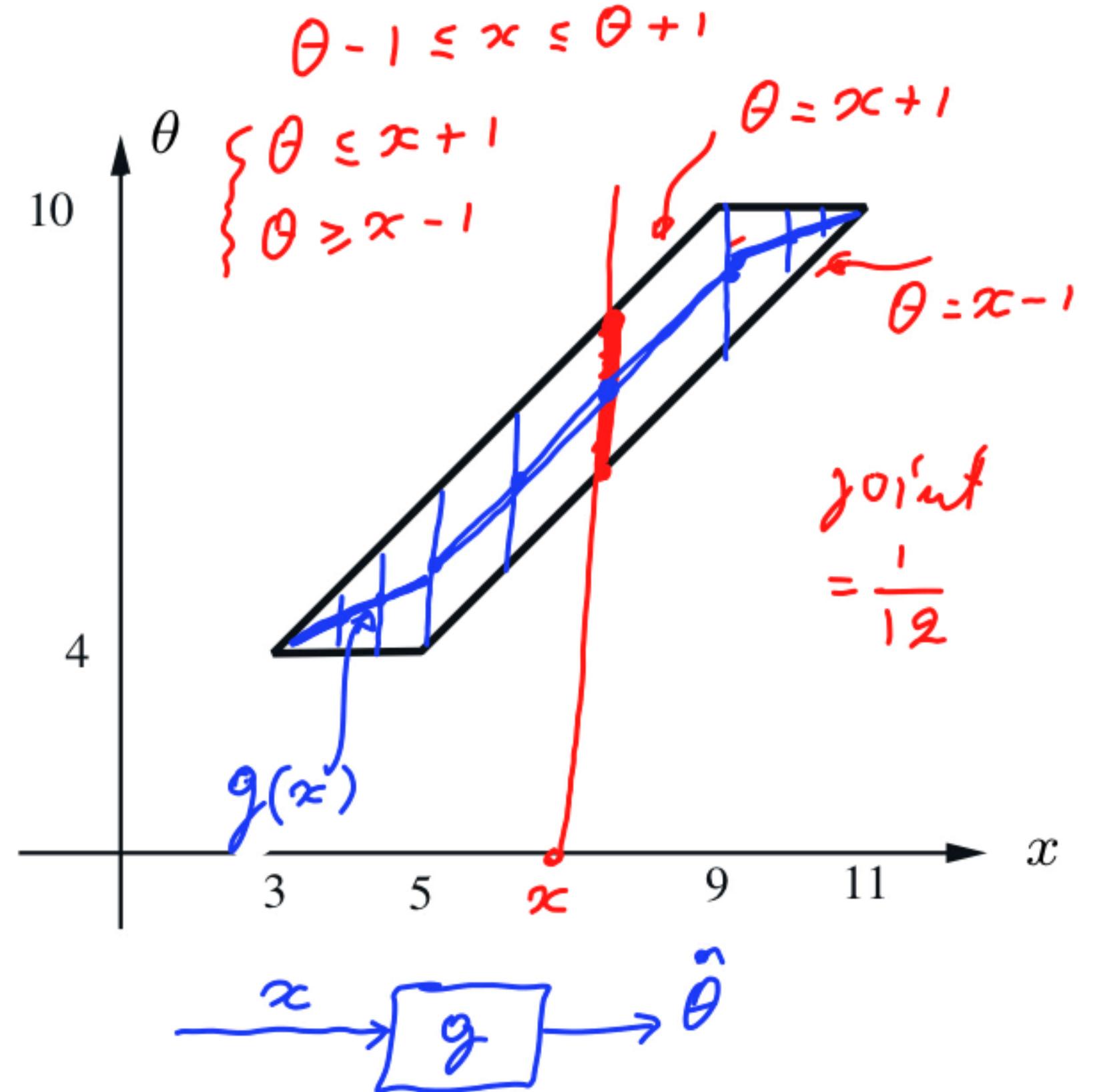
- LMS relevant to estimation (not hypothesis testing)
- Same as MAP if the posterior is unimodal and symmetric around the mean
 - e.g., when posterior is normal (the case in “linear–normal” models)



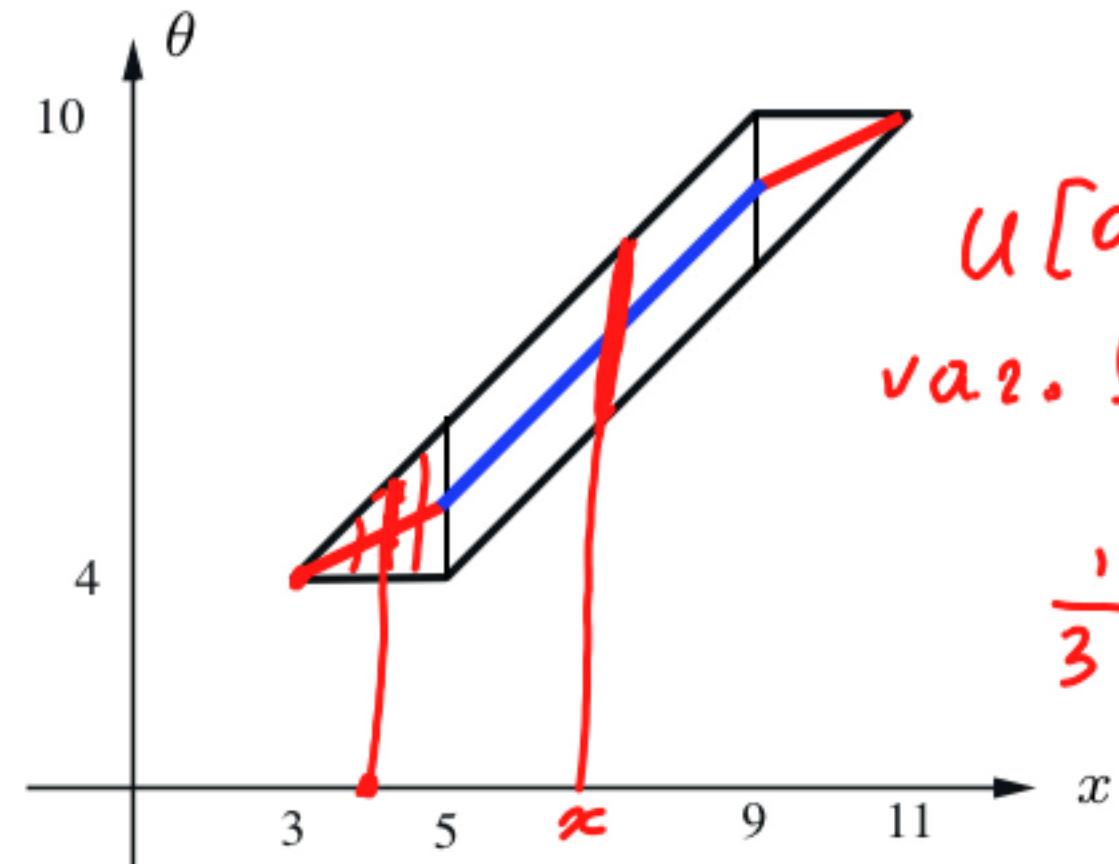
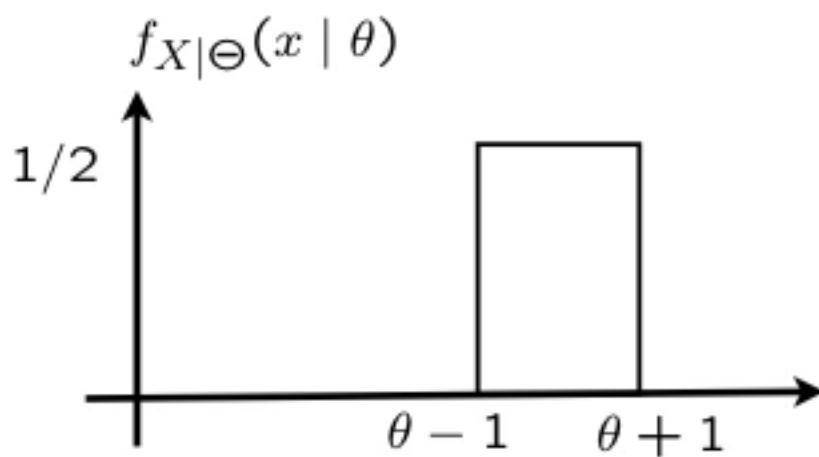
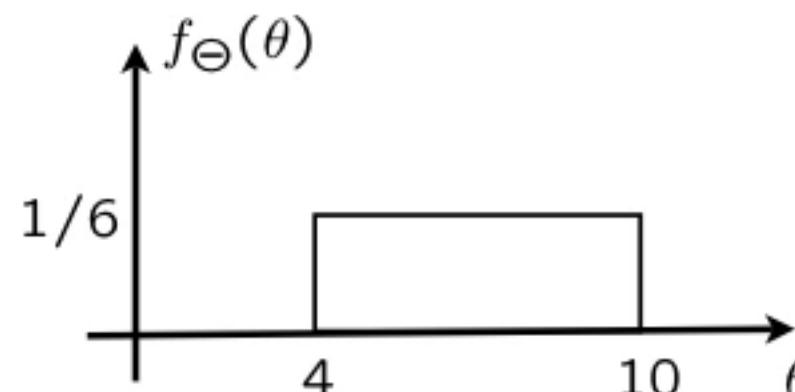
Example



$$x = \theta + u \quad u \sim \text{unif}(-1, 1)$$

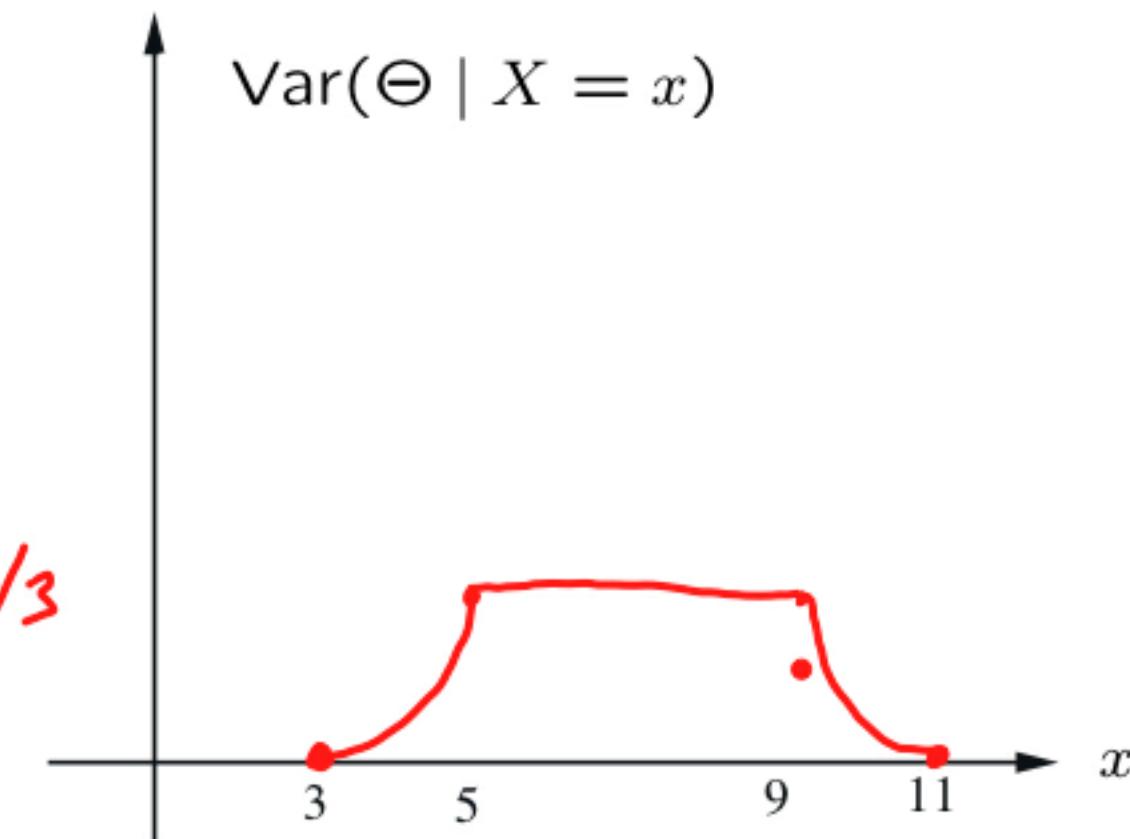


Conditional mean squared error



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

$$E[\text{Var}(\Theta | x)] = \int f_x(x) \text{Var}_\Theta(\Theta | x=x) dx$$



LMS estimation with multiple observations or unknowns

- unknown Θ ; prior $p_\Theta(\theta)$
 - interested in a point estimate $\hat{\theta}$
- observations $X = (X_1, X_2, \dots, X_n)$; model $p_{X|\Theta}(x | \theta)$
 - observe that $X = x$
 - new universe: condition on $X = x$
- LMS estimate: $E[\Theta | X_1 = x_1, \dots, X_n = x_n]$
- If Θ is a vector, apply to each component separately

$$\Theta = (\theta_1, \dots, \theta_m) \quad \hat{\theta}_j = E[\theta_j | X_1 = x_1, \dots, X_n = x_n]$$

Some challenges in LMS estimation

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

$$f_X(x) = \int f_{\Theta}(\theta') f_{X|\Theta}(x | \theta') d\theta'$$

- Full correct model, $f_{X|\Theta}(x | \theta)$, may not be available •
- Can be hard to compute/implement/analyze

$$E[\theta_j | x=x] = \iiint \theta_j \cdot f_{\Theta|X}(\theta | x) d\theta_1 \dots d\theta_m$$

Properties of the estimation error in LMS estimation

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

$$E[\tilde{\theta} | X = x] = 0$$

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

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

$$E[\tilde{\theta}] = 0$$

$$E[\hat{\theta} - \theta | X = x] = \hat{\theta} - E[\theta | X = x] = 0$$

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

$$\underline{E[\tilde{\theta} \hat{\theta}]} - \cancel{E[\tilde{\theta}] E[\hat{\theta}]} \circ$$

$$E[\tilde{\theta} \hat{\theta} | X = x] = \hat{\theta} E[\tilde{\theta} | X = x] = 0$$

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

•

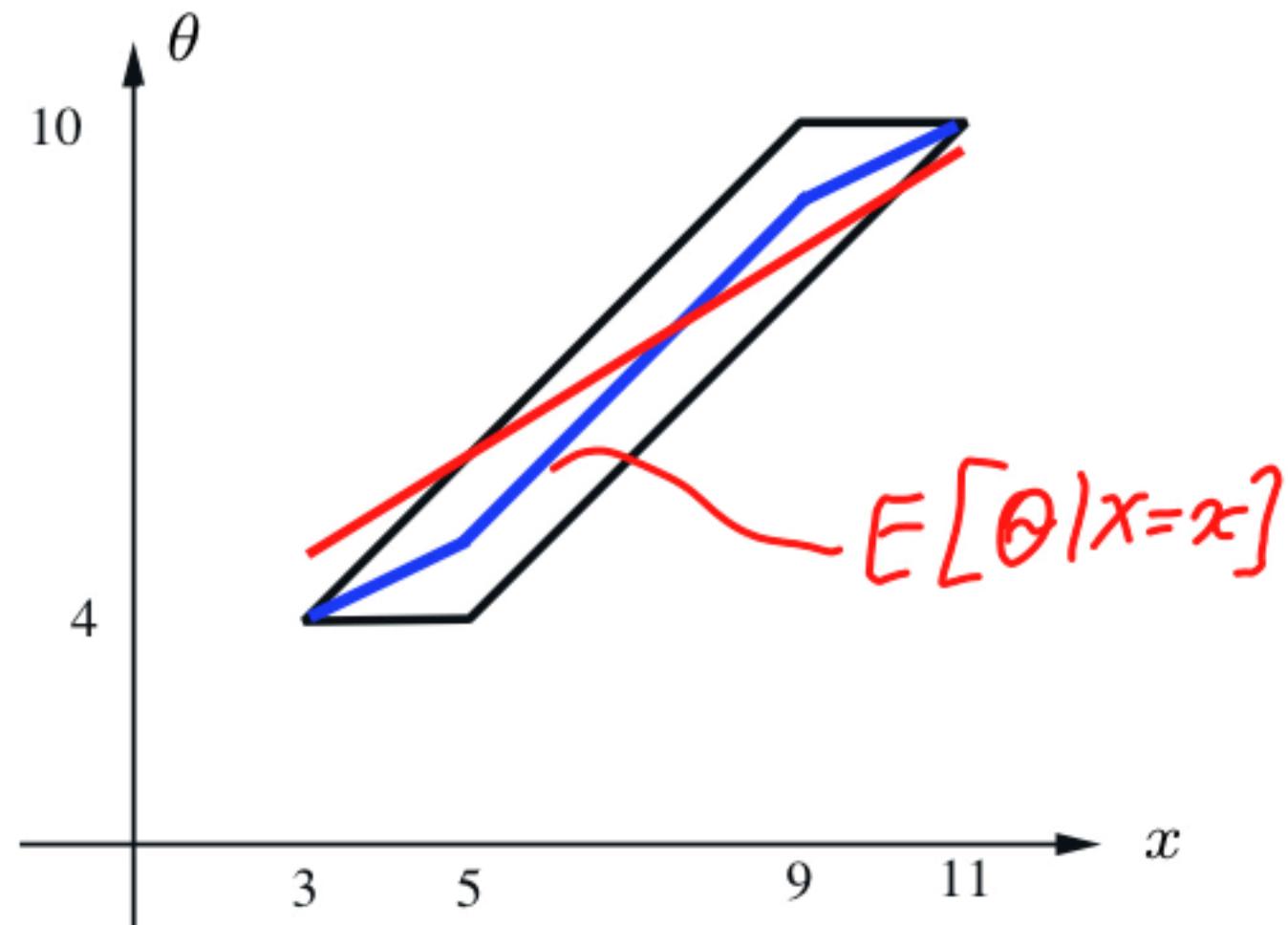
$$\theta = \hat{\theta} - \tilde{\theta}$$

LECTURE 17: Linear least mean squares (LLMS) estimation

- Conditional expectation $E[\Theta | X]$ may be hard to compute/implement
- Restrict to estimators $\widehat{\Theta} = aX + b$
 - minimize mean squared error
- Simple solution
- Mathematical properties
- Example

LLMS formulation

- Unknown Θ ; observation X



- Minimize $E[(\hat{\Theta} - \Theta)^2]$
- Estimators $\hat{\Theta} = g(X) \rightarrow \hat{\Theta}_{LLMS} = E[\Theta | X]$
- Consider estimators of Θ , of the form $\hat{\Theta} = aX + b$
- Minimize $E[(\Theta - aX - b)^2]$, w.r.t. a, b
- If $E[\Theta | X]$ is linear in X , then $\hat{\Theta}_{LLMS} = \hat{\Theta}_{LMS}$

Solution to the LLMS problem

- Minimize $E[(\Theta - aX - b)^2]$, w.r.t. a, b

– suppose a has already been found:

$$\begin{aligned} \min \quad & E[(\Theta - aX - E[\Theta - aX])^2] = \text{var}(\Theta - aX) \\ & = \text{var}(\Theta) + a^2 \text{var}(X) - 2a \text{cov}(\Theta, X) \\ \frac{d}{da} &= 0 : 2a \text{var}(X) - 2 \text{cov}(\Theta, X) = 0 \quad \left| \begin{array}{l} \rho = \frac{\text{cov}(\Theta, X)}{\sigma_\Theta \sigma_X} \\ a = \frac{\rho \sigma_\Theta \sigma_X}{\sigma_X^2} \end{array} \right. \\ da & \quad a = \frac{\text{cov}(\Theta, X) / \text{var}(X)}{2} \end{aligned}$$

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

Remarks on the solution and on the error variance

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

- Only means, variances, covariances matter

- $\rho > 0: X > E[X] \Rightarrow \widehat{\Theta}_L > E[\Theta]$

$$|\rho| = 1$$

$$\widehat{\Theta}_L = \Theta$$

- $\rho = 0: \widehat{\Theta}_L = E[\Theta]$

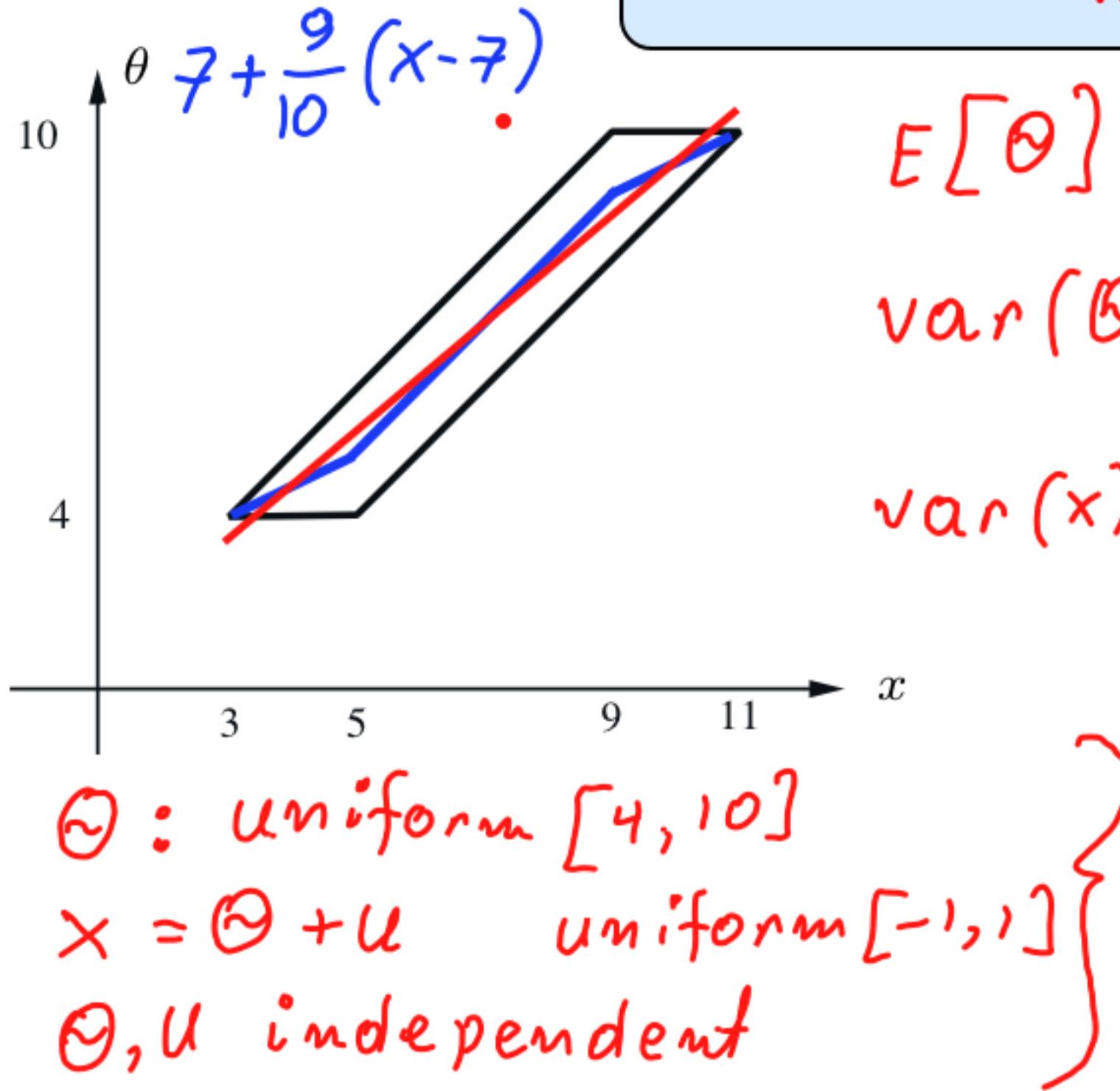
assume $E[\Theta] = E[X] = 0$

$$E[(\widehat{\Theta}_L - \Theta)^2] = (1 - \rho^2) \text{var}(\Theta)$$

$$E\left[(\Theta - \rho \frac{\sigma_\Theta}{\sigma_X} X)^2\right] = \sigma_\Theta^2 - 2\rho \frac{\sigma_\Theta}{\sigma_X} \cancel{\rho \sigma_\Theta \sigma_X} + \rho^2 \frac{\sigma_\Theta^2}{\cancel{\sigma_X^2}} \cancel{\sigma_X^2}$$

Example

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



$$E[\Theta] = 7 \quad E[U] = 0 \quad E[X] = 7$$

$$\text{var}(\Theta) = \frac{6^2}{12} = 3 \quad \text{var}(U) = \frac{2^2}{12} = \frac{1}{3}$$

$$\text{var}(X) = 3 + \frac{1}{3} = \frac{10}{3}$$

$$\begin{aligned} \text{cov}(\Theta, \Theta + U) &= \\ &= \text{cov}(\Theta, \Theta) + \cancel{\text{cov}(\Theta, U)} = 3 \end{aligned}$$

LLMS for inferring the parameter of a coin

- Standard example:
 - coin with bias Θ ; prior $f_\Theta(\cdot)$
 - fix n ; $X =$ number of heads
- Assume $f_\Theta(\cdot)$ is uniform in $[0, 1]$

$$\widehat{\Theta}_{\text{LMS}} = \frac{X + 1}{n + 2} = \widehat{\Theta}_{\text{LLMS}}$$

$$\widehat{\Theta}_{\text{LLMS}} = \mathbf{E}[\Theta] + \frac{\text{Cov}(\Theta, X)}{\text{var}(X)}(X - \mathbf{E}[X])$$

LLMS for inferring the parameter of a coin

- Θ : uniform on $[0, 1]$ $E[\Theta] = \frac{1}{2}$ $\text{var}(\Theta) = \frac{1}{12}$ $E[\Theta^2] = \frac{1}{12} + \frac{1}{2^2} = \frac{1}{3}$

- $p_{X|\Theta}$: $\text{Bin}(n, \Theta)$ $E[X | \Theta] = n\Theta$ $\text{var}(X | \Theta) = n\Theta(1 - \Theta)$

$$E[X] = E[n\Theta] = n/2 \quad E[X^2 | \Theta] = n\Theta(1 - \Theta) + n^2\Theta^2$$

- $E[X^2] = E[E[X^2 | \Theta]] = E[n\Theta + (n^2 - n)\Theta^2] = \frac{n}{2} + \frac{n^2 - n}{3} = \frac{n}{6} + \frac{n^2}{3}$

$$\text{var}(X) = E[X^2] - (E[X])^2 = \frac{n}{6} + \frac{n^2}{3} - \frac{n^2}{4} = \frac{n}{6} + \frac{n^2}{12} = \frac{n(n+2)}{12}$$

$$E[\Theta X | \Theta] = \Theta E[X | \Theta] = n\Theta^2$$

$$E[\Theta X] = E[E[\Theta X | \Theta]] = E[n\Theta^2] = n/3$$

$$\text{cov}(\Theta, X) = E[\Theta X] - E[\Theta]E[X] = \frac{n}{3} - \frac{n}{4} = \frac{n}{12}$$

LLMS for inferring the parameter of a coin

$$\widehat{\Theta}_{LLMS} = \mathbf{E}[\Theta] + \frac{\text{Cov}(\Theta, X)}{\text{var}(X)}(X - \mathbf{E}[X])$$

$$\text{cov}(\Theta, X) = \frac{n}{12} \quad \text{var}(X) = \frac{n(n+2)}{12} \quad \mathbf{E}[X] = \frac{n}{2}$$

$$\widehat{\Theta}_{LLMS} = \frac{X+1}{n+2} = \widehat{\Theta}_{LMS}$$

LLMS with multiple observations

- Unknown Θ ; observations $X = (X_1, \dots, X_n)$
- Consider estimators of the form: $\widehat{\Theta} = a_1X_1 + \dots + a_nX_n + b$
- Find best choices of a_1, \dots, a_n, b
minimize: $E[(a_1X_1 + \dots + a_nX_n + b - \Theta)^2] = a_1^2 E[X^2] + 2a_1 a_2 E[X, X_2] + \dots + a_n^2 E[X, \Theta] + \dots$
- If $E[\Theta | X]$ is linear in X , then $\widehat{\Theta}_{\text{LMS}} = \widehat{\Theta}_{\text{LLMS}}$
- Solve linear system in b and the a_i •
- Only means, variances, covariances matter
- If multiple unknown Θ_j , apply to each one, separately

The simplest LLMS example with multiple observations

$$\begin{aligned} X_1 &= \Theta + W_1 & \Theta &\sim x_0, \sigma_0^2 & W_i &\sim 0, \sigma_i^2 \\ &\vdots \\ X_n &= \Theta + W_n & \Theta, W_1, \dots, W_n &\text{ uncorrelated} \end{aligned}$$

- Suppose Θ, W_1, \dots, W_n are independent normal

$$\hat{\theta}_{\text{LMS}} = \mathbb{E}[\Theta | X = x] = \frac{\sum_{i=0}^n \frac{x_i}{\sigma_i^2}}{\sum_{i=0}^n \frac{1}{\sigma_i^2}} \quad \widehat{\Theta}_{\text{LMS}} = \mathbb{E}[\Theta | X] = \frac{\frac{x_0}{\sigma_0^2} + \sum_{i=1}^n \frac{X_i}{\sigma_i^2}}{\sum_{i=0}^n \frac{1}{\sigma_i^2}} = \widehat{\Theta}_{\text{LLMS}}$$

- Suppose general (not normal) distributions, but same means, variances, as in normal example
 - all covariances also the same
 - solution must be the same

The representation of the data matters in LLMS

- Estimation based on X versus X^3
 - LMS: $\underline{E[\Theta | X]}$ is the same as $\underline{E[\Theta | X^3]}$
 - LLMS is different: estimator $\widehat{\Theta} = \underline{aX + b}$ versus $\widehat{\Theta} = \underline{aX^3 + b}$
 $\text{cov}(\Theta, X^3) \quad \text{var}(x^3)$
 - can also consider $\widehat{\Theta} = \underline{a_1 X} + \underline{a_2 X^2} + \underline{a_3 X^3} + b$
 - can also consider $\widehat{\Theta} = \underline{a_1 X} + \underline{a_2 e^X} + \underline{a_3 \log X} + b$

LECTURE 18: Inequalities, convergence, and the Weak Law of Large Numbers

- Inequalities
 - bound $\mathbf{P}(X \geq a)$ based on limited information about a distribution
 - Markov inequality (based on the mean)
 - Chebyshev inequality (based on the mean and variance)
- WLLN: X, X_1, \dots, X_n i.i.d.

$$\frac{X_1 + \dots + X_n}{n} \longrightarrow \mathbf{E}[X]$$

- application to polling
- Precise defn. of convergence
 - convergence “in probability”

The Markov inequality

- Use a bit of information about a distribution to learn something about probabilities of “extreme events”
- “If $X \geq 0$ and $E[X]$ is small, then X is unlikely to be very large”

Markov inequality: If $X \geq 0$ and $a > 0$, then $P(X \geq a) \leq \frac{E[X]}{a}$

$$Y = \begin{cases} 0, & \text{if } X < a \\ a, & \text{if } X \geq a \end{cases} \quad \text{and } P(X \geq a) = E[Y] \leq E[X]$$

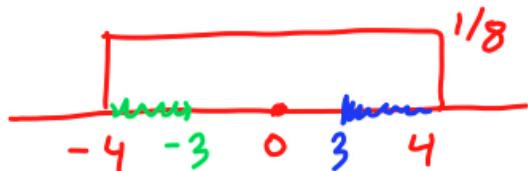
The Markov inequality

Markov inequality: If $X \geq 0$ and $a > 0$, then $P(X \geq a) \leq \frac{E[X]}{a}$

- **Example:** X is Exponential($\lambda = 1$): $P(X \geq a) \leq \frac{1}{a}$



- **Example:** X is Uniform[-4, 4]: $P(|X| \geq 3) \leq \frac{E[|X|]}{3} = \frac{2}{3}$



$$= \frac{1}{2} P(|X| \geq 3) \leq \frac{1}{3}$$



The Chebyshev inequality

- Random variable X , with finite mean μ and variance σ^2
- “If the variance is small, then X is unlikely to be too far from the mean”

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

Markov inequality: If $X \geq 0$ and $a > 0$, then $P(X \geq a) \leq \frac{E[X]}{a}$

$$P(|X - \mu| \geq c) = P(\underbrace{(X - \mu)^2}_{\geq c^2} \geq c^2) \leq \frac{E[(X - \mu)^2]}{c^2} = \frac{\sigma^2}{c^2}$$

The Chebyshev inequality

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

$$P(|X - \mu| \geq k\sigma) \leq \frac{\sigma^2}{k^2\sigma^2} = \frac{1}{k^2} \quad k=3 \leq \frac{1}{9}$$

- **Example:** X is Exponential($\lambda = 1$): $P(X \geq a) \leq \frac{1}{a}$ (Markov)



$$P(X \geq a) = P(X-1 \geq a-1) \leq P(|X-1| \geq a-1) \leq \frac{1}{(a-1)^2} \sim \frac{1}{a^2}$$

The Weak Law of Large Numbers (WLLN)

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

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

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

- $\text{Var}(M_n) = \frac{\text{Var}(X_1 + \dots + X_n)}{n^2} = \frac{n\sigma^2}{n^2} = \frac{\sigma^2}{n}$

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

WLLN: For $\epsilon > 0$, $P(|M_n - \mu| \geq \epsilon) = P\left(\left|\frac{X_1 + \dots + X_n}{n} - \mu\right| \geq \epsilon\right) \rightarrow 0$, as $n \rightarrow \infty$

Interpreting the WLLN

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

WLLN: For $\epsilon > 0$, $P(|M_n - \mu| \geq \epsilon) = P\left(\left|\frac{X_1 + \dots + X_n}{n} - \mu\right| \geq \epsilon\right) \rightarrow 0, \text{ as } n \rightarrow \infty$

- One experiment
 - many measurements $X_i = \mu + W_i$
 - W_i : measurement noise; $E[W_i] = 0$; independent W_i
 - **sample mean** M_n is unlikely to be far off from **true mean** μ
- Many independent repetitions of the same experiment
 - event A , with $p = P(A)$
 - X_i : indicator of event A
 - the sample mean M_n is the **empirical frequency** of event A

$$x_i = \begin{cases} 1, & \text{if } A \text{ occurs} \\ 0, & \text{o.w.} \end{cases} \quad E[\overset{\bullet}{x}_i] = p$$

The pollster's problem

- p : fraction of population that will vote "yes" in a referendum

- i th (randomly selected) person polled: $X_i = \begin{cases} 1, & \text{if yes,} \\ 0, & \text{if no.} \end{cases}$
- uniformly, independently*

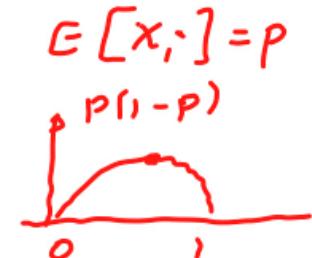
- $M_n = (X_1 + \dots + X_n)/n$: fraction of "yes" in our sample

- Would like "small error," e.g.: $|M_n - p| < 0.01$

- Try $n = 10,000$

$$\bullet P(|M_{10,000} - p| \geq 0.01) \leq \frac{\sigma^2}{n \varepsilon^2} = \frac{p(1-p)}{10^4 \cdot 10^{-4}} \leq \frac{1}{4} \quad \leftarrow \text{want } \leq 5\% \quad \overbrace{p}^{=}$$

$$\frac{1/4}{n \cdot 10^{-4}} \leq \frac{5}{10^2} \Leftrightarrow n \geq \frac{10^6}{20} = 50,000 \quad \leftarrow \text{will suffice}$$



Convergence “in probability”

WLLN: For any $\epsilon > 0$, $P(|M_n - \mu| \geq \epsilon) \rightarrow 0$, as $n \rightarrow \infty$

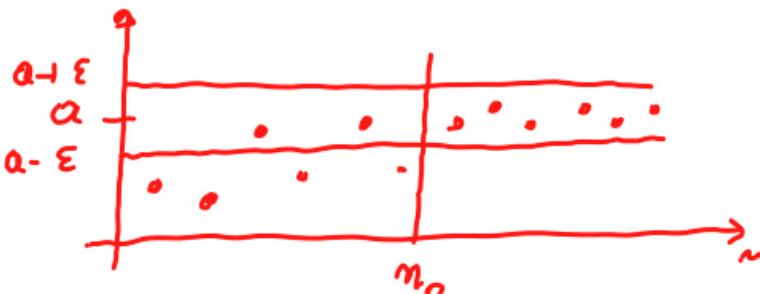
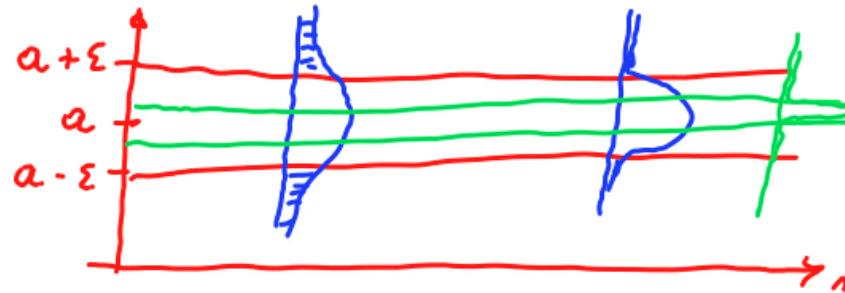
- Would like to say that “ M_n converges to μ ”
- Need to define the word “converges”
- Sequence of random variables Y_n ; not necessarily independent

$$M_n \xrightarrow[n \rightarrow \infty]{i.p} \mu$$

Definition: A sequence Y_n converges in probability to a number a if:

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

Understanding convergence “in probability”

- Ordinary convergence
 - Sequence a_n ; number a
 $a_n \rightarrow a$
 “ a_n eventually gets and stays (arbitrarily) close to a ”
 - Convergence in probability
 - Sequence Y_n ; number a
 $Y_n \rightarrow a$
 - for any $\epsilon > 0$, $P(|Y_n - a| \geq \epsilon) \rightarrow 0$
- For every $\epsilon > 0$, there exists n_0 , such that for every $n \geq n_0$, we have $|a_n - a| \leq \epsilon$
- “(almost all) of the PMF/PDF of Y_n eventually gets concentrated (arbitrarily) close to a ”

Some properties

- Suppose that $X_n \rightarrow a$, $Y_n \rightarrow b$, in probability

- If g is continuous, then $g(X_n) \rightarrow g(a)$

$$X_n^2 \rightarrow a^2$$

- $X_n + Y_n \rightarrow a + b$

- **But:** $E[X_n]$ need not converge to a

Convergence in probability examples



$Y_n \xrightarrow[n \rightarrow \infty]{i.p.} 0.$

$$\varepsilon > 0 \quad P(|Y_n - 0| \geq \varepsilon) = 1/n \xrightarrow{n \rightarrow \infty} 0$$

$$E[Y_n] = n^2 \cdot \frac{1}{n} = n \xrightarrow{n \rightarrow \infty}$$

- convergence in probability does **not** imply convergence of expectations

Convergence in probability examples

- X_i : i.i.d., uniform on $[0, 1]$
- $Y_n = \min\{X_1, \dots, X_n\}$



$$Y_{n+1} \leq Y_n$$

$$\mathbb{P}(|Y_n - 0| \geq \epsilon) = \mathbb{P}(Y_n \geq \epsilon).$$

$$\epsilon > 0$$

$$= \mathbb{P}(X_1 \geq \epsilon, \dots, X_n \geq \epsilon)$$

$$Y_n \xrightarrow[n \rightarrow \infty]{i.p.} 0$$

$$\epsilon > 1$$

$$= \mathbb{P}(X_1 \geq \epsilon) \cdots \mathbb{P}(X_n \geq \epsilon)$$

$$\epsilon \leq 1$$

$$= (1 - \epsilon)^n \xrightarrow{n \rightarrow \infty} 0$$

Related topics

- Better bounds/approximations on tail probabilities

- Markov and Chebyshev inequalities

$$\Pr(|M_n - \mu| \geq a) \leq e^{-n \frac{\lambda(a)}{a^2}} \sim \frac{1}{n}$$

- Chernoff bound

$$\Pr(|M_n - \mu| \geq a) \leq e^{-n \frac{\lambda(a)}{a^2}} \sim \frac{1}{n}$$

- Central limit theorem " $M_n \sim N(\mu, \sigma^2/n)$ "

- Different types of convergence

- Convergence in probability

- Convergence "with probability 1" $\Pr\left(\left\{\omega : Y_n(\omega) \rightarrow Y(\omega)\right\}\right) = 1$

- Strong law of large numbers $M_n \xrightarrow[n \rightarrow \infty]{\text{a.s.}} \mu$

- Convergence of a sequence of distributions (CDFs) to a limiting CDF

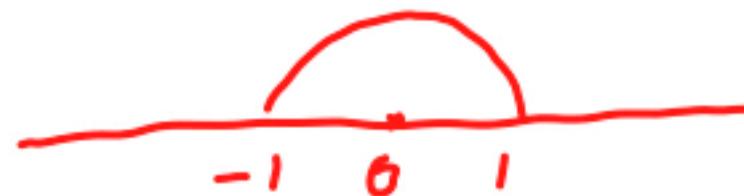
LECTURE 19: The Central Limit Theorem (CLT)

- WLLN: $\frac{X_1 + \dots + X_n}{n} \rightarrow E[X]$

- CLT: $X_1 + \dots + X_n \approx \text{normal}$
 - precise statement
 - universality, usefulness
 - many examples
 - refinement for discrete r.v.s
 - application to polling

Different scalings of the sum of i.i.d. random variables

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



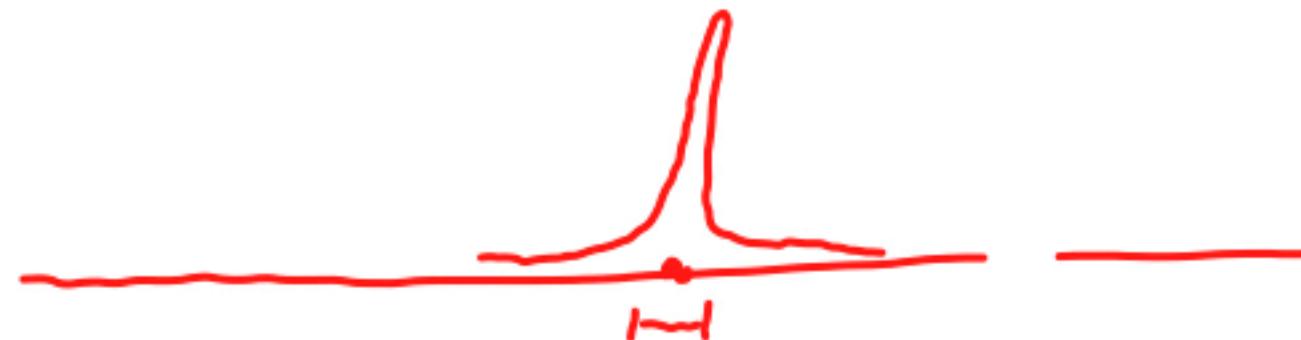
- $S_n = X_1 + \dots + X_n$

variance: $n\sigma^2$



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

variance: $\frac{\sigma^2}{n} \xrightarrow{n \rightarrow \infty} 0$



- $\frac{S_n}{\sqrt{n}} = \frac{X_1 + \dots + X_n}{\sqrt{n}}$

variance: $\sigma^2 = \frac{n\sigma^2}{n}$



The Central Limit Theorem (CLT)

- X_1, \dots, X_n i.i.d., finite mean μ and variance σ^2
- $S_n = X_1 + \dots + X_n$ variance: $n\sigma^2$
- $\frac{S_n}{\sqrt{n}} = \frac{X_1 + \dots + X_n}{\sqrt{n}}$ variance: σ^2

$$Z_n = \frac{S_n - n\mu}{\sqrt{n}\sigma} \quad E[Z_n] = 0 \quad \text{var}(Z_n) = 1$$

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

Central Limit Theorem: For every z : $\lim_{n \rightarrow \infty} P(Z_n \leq z) = P(Z \leq z)$

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

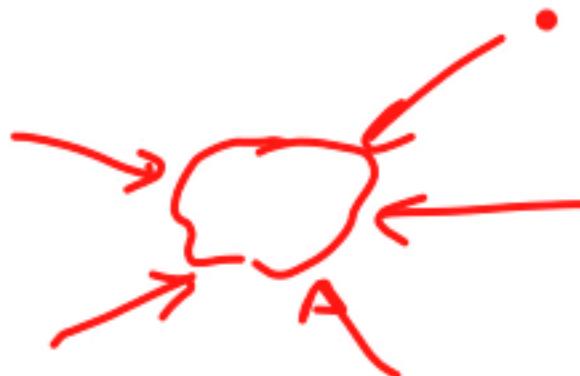
Usefulness of the CLT

$$S_n = X_1 + \dots + X_n$$

$$Z_n = \frac{S_n - n\mu}{\sqrt{n} \sigma} \quad Z \sim N(0, 1)$$

Central Limit Theorem: For every z : $\lim_{n \rightarrow \infty} P(Z_n \leq z) = P(Z \leq z)$

- universal and easy to apply; only means, variances matter
- fairly accurate computational shortcut
- justification of normal models



What exactly does the CLT say? — Theory

$$S_n = X_1 + \dots + X_n \quad Z_n = \frac{S_n - n\mu}{\sqrt{n}\sigma} \quad Z \sim N(0, 1)$$

Central Limit Theorem: For every z : $\lim_{n \rightarrow \infty} P(Z_n \leq z) = P(Z \leq z)$

- CDF of Z_n converges to normal CDF
- results for convergence of PDFs or PMFs (with more assumptions)
- results without assuming that the X_i are identically distributed
- results under “weak dependence”
- proof: uses “transforms”: $E[e^{sZ_n}] \rightarrow E[e^{sZ_\bullet}]$, for all s

What exactly does the CLT say? — Practice

$$S_n = X_1 + \dots + X_n$$

$$Z_n = \frac{S_n - n\mu}{\sqrt{n}\sigma} \quad Z \sim N(0, 1)$$

Central Limit Theorem: For every z : $\lim_{n \rightarrow \infty} P(Z_n \leq z) = P(Z \leq z)$

- The **practice** of normal approximations:

- treat Z_n as if it were normal

$$S_n = \sqrt{n}\sigma Z_n + n\mu$$

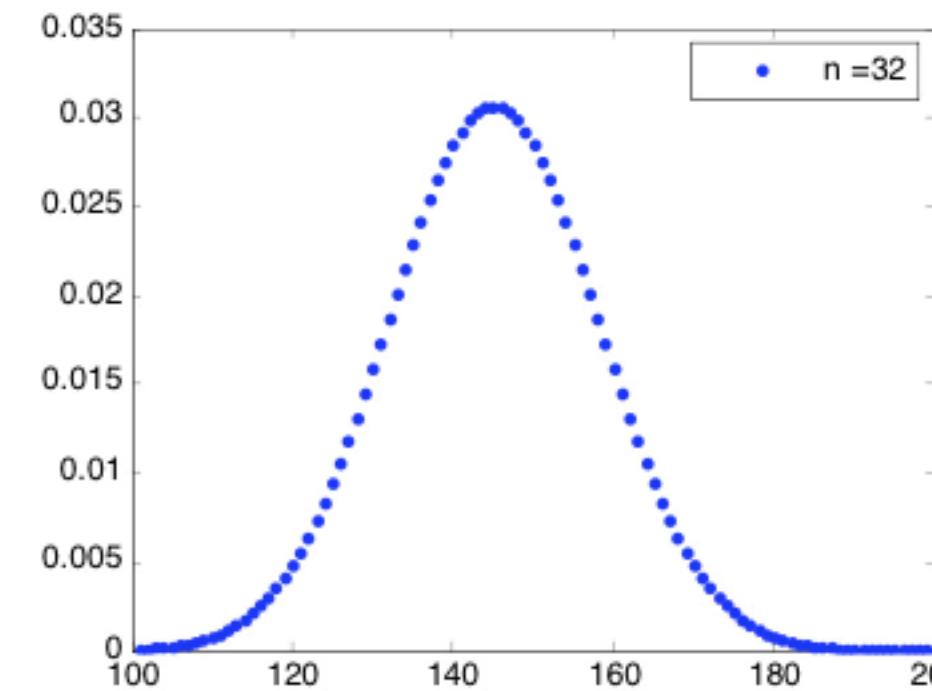
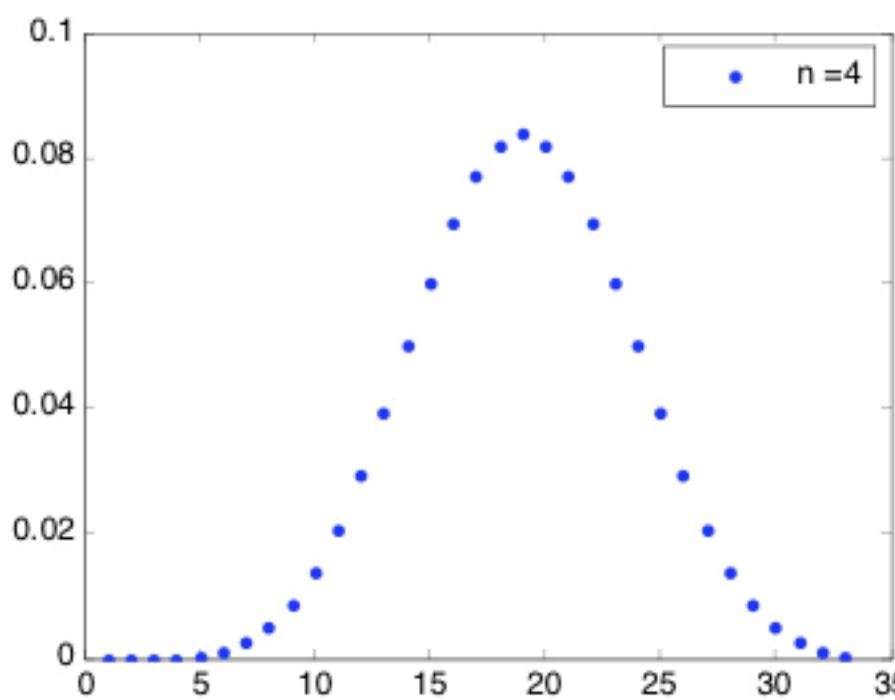
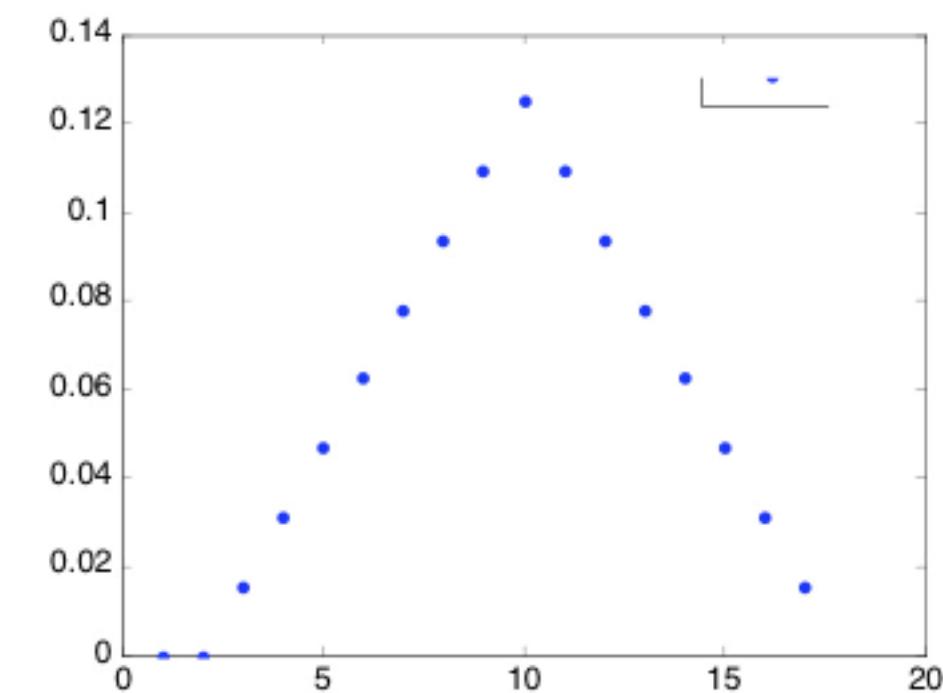
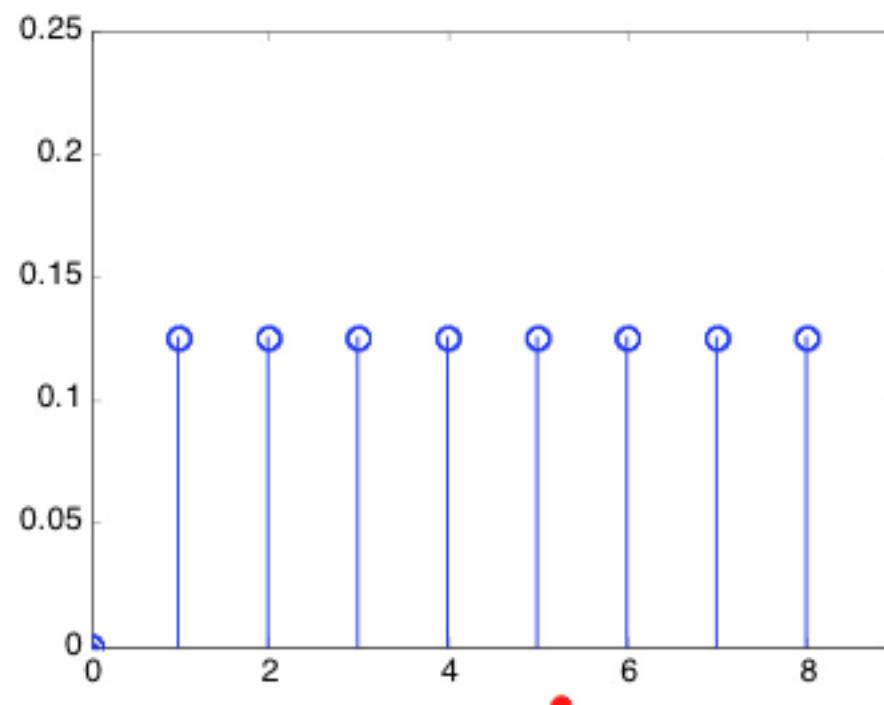
- hence treat S_n as if normal: $N(n\mu, n\sigma^2)$

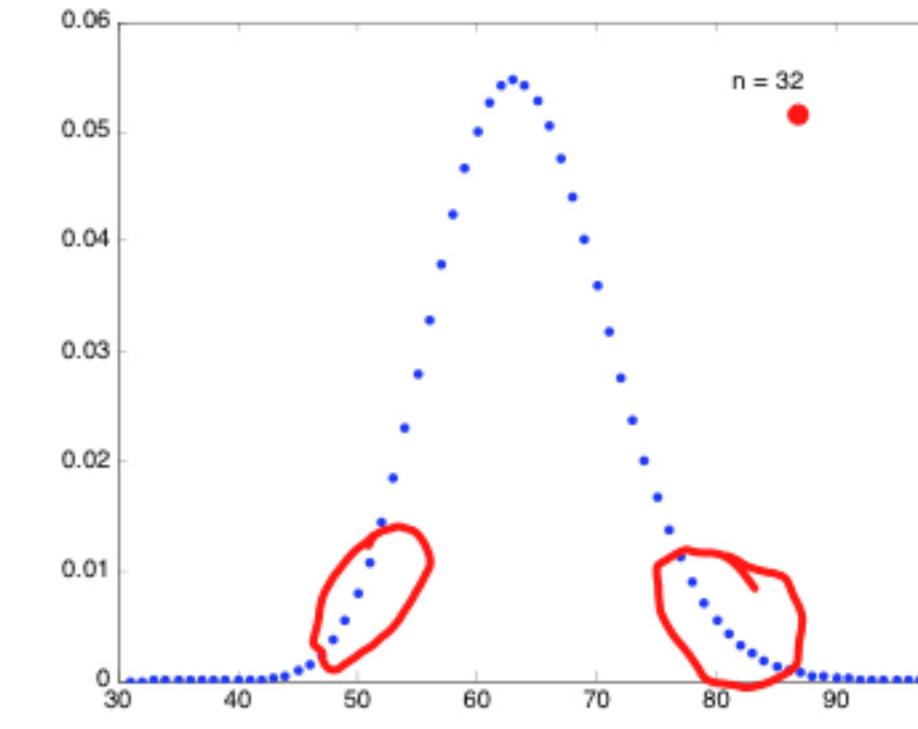
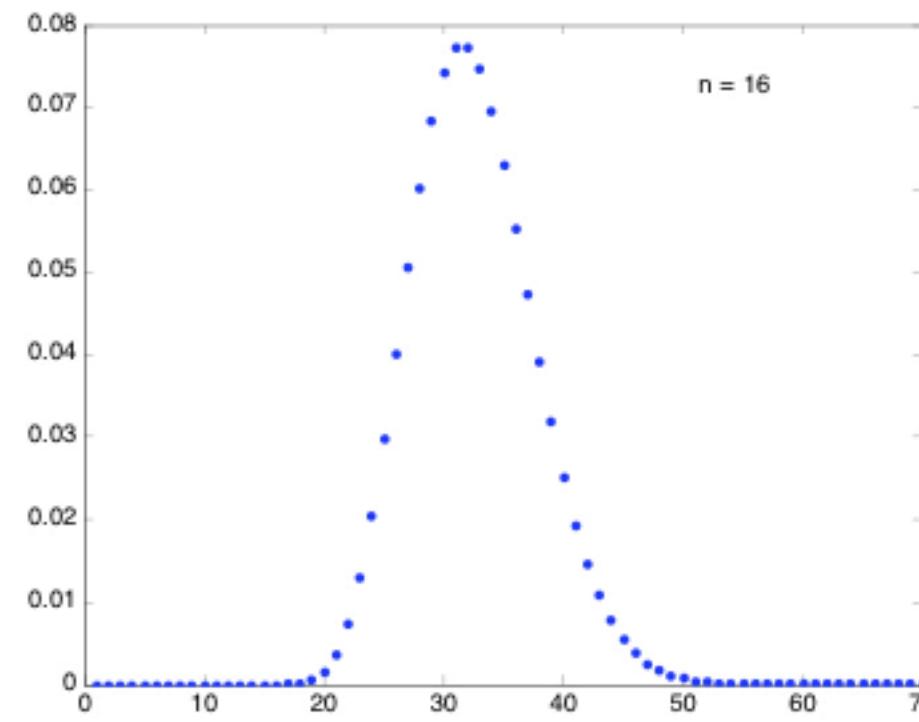
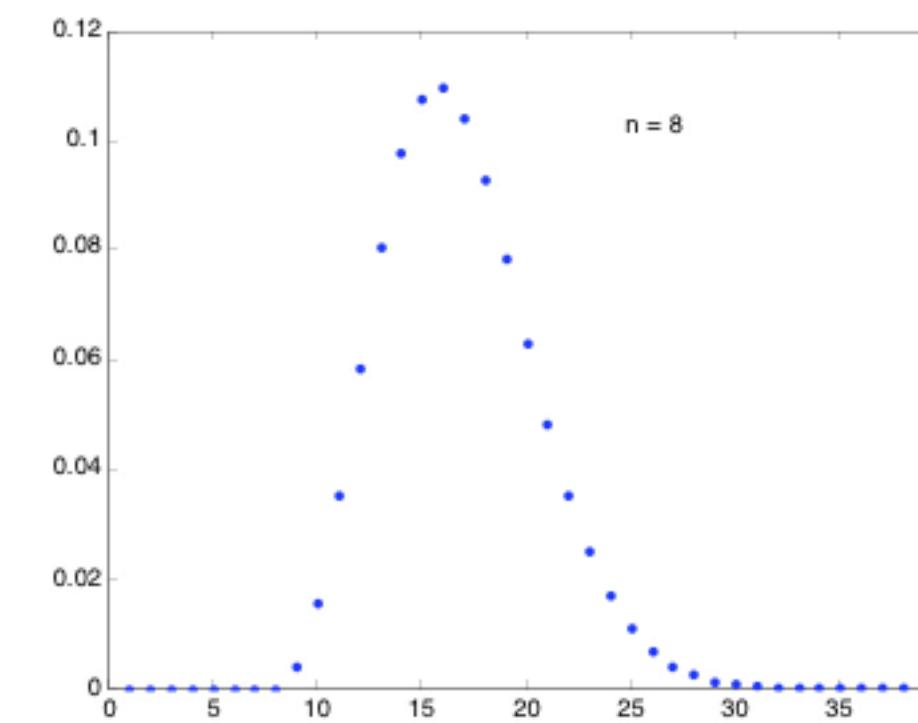
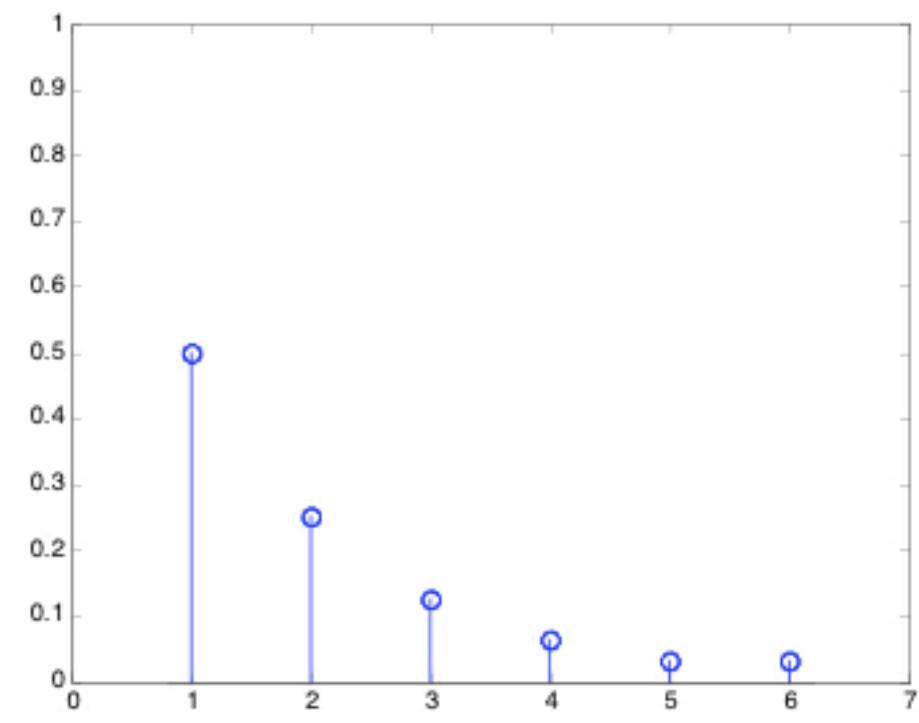
- Can we use the CLT when n is “moderate”?

$n = 30$?

- usually, yes

- symmetry and unimodality help





Example 1

- $P(S_n \leq a) \approx b$ given two parameters, find the third
- Package weights X_i , i.i.d. exponential, $\lambda = 1/2$;
- Load container with $n = 100$ packages

$$Z_n = \frac{S_n - n\mu}{\sqrt{n}\sigma}$$

$$\mu = \sigma = 2$$

$$P(S_n \geq 210)$$

$$= P\left(\frac{S_n - 200}{20} > \frac{210 - 200}{20}\right)$$

$$= P(Z_n > 0.5) \approx P(Z > 0.5)$$

$$= 1 - P(Z < 0.5) = 1 - \Phi(0.5)$$

$$= 1 - 0.6915 = 0.3085$$

	.00	.01	.02	.03	.04	.05	.06	.07	.08	.09
0.0	.5000	.5040	.5080	.5120	.5160	.5199	.5239	.5279	.5319	.5359
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	0.6915	.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	.8186	.8212	.8238	.8264	.8289	.8315	.8340	.8365	.8389
1.0	.8413	.8438	.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

Example 2

- $P(S_n \leq a) \approx b$ given two parameters, find the third
- Package weights X_i , i.i.d. exponential, $\lambda = 1/2$;
- Let $n = 100$. Choose the “capacity” a , so that $P(S_n \geq a) \approx 0.05$.

$$0.05 \approx P\left(\frac{S_n - 200}{20} > \frac{a - 200}{20}\right)$$

$$\approx 1 - \Phi\left(\frac{a - 200}{20}\right)$$

0.95

$$\frac{a - 200}{20} = 1.645 \quad a = 232.9$$

$$Z_n = \frac{S_n - n\mu}{\sqrt{n}\sigma}$$

$$\mu = \sigma = 2$$

	.00	.01	.02	.03	.04	.05	.06	.07	.08	.09
0.0	.5000	.5040	.5080	.5120	.5160	.5199	.5239	.5279	.5319	.5359
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	.8186	.8212	.8238	.8264	.8289	.8315	.8340	.8365	.8389
1.0	.8413	.8438	.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

Example 3

- $P(S_n \leq a) \approx b$ given two parameters, find the third
- Package weights X_i , i.i.d. exponential, $\lambda = 1/2$;
- How large can n be,
so that $P(S_n \geq 210) \approx 0.05$?

$$Z_n = \frac{S_n - n\mu}{\sqrt{n}\sigma}$$

$$\mu = \sigma = 2$$

$$P\left(\frac{S_n - 2n}{2\sqrt{n}} \geq \frac{210 - 2n}{2\sqrt{n}}\right)$$

$$\approx 1 - \Phi\left(\frac{210 - 2n}{2\sqrt{n}}\right) \approx 0.05$$

0.95

$$\frac{210 - 2n}{2\sqrt{n}} = 1.645$$

$n = 89$

	.00	.01	.02	.03	.04	.05	.06	.07	.08	.09
0.0	.5000	.5040	.5080	.5120	.5160	.5199	.5239	.5279	.5319	.5359
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	.8186	.8212	.8238	.8264	.8289	.8315	.8340	.8365	.8389
1.0	.8413	.8438	.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

Example 4

- $P(S_n \leq a) \approx b$ given two parameters, find the third
- Package weights X_i , i.i.d. exponential, $\lambda = 1/2$;

$$Z_n = \frac{S_n - n\mu}{\sqrt{n}\sigma}$$

$$\mu = \sigma = 2$$

- Load container until weight exceeds 210

N : number of packages loaded

- $P(N > 100)$

$$= P\left(\sum_{i=1}^{100} X_i \leq 210\right)$$

$$\approx \Phi\left(\frac{210 - 200}{20}\right) = \Phi(0.5)$$

$$= 0.6915$$

	.00	.01	.02	.03	.04	.05	.06	.07	.08	.09
0.0	.5000	.5040	.5080	.5120	.5160	.5199	.5239	.5279	.5319	.5359
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	.8186	.8212	.8238	.8264	.8289	.8315	.8340	.8365	.8389
1.0	.8413	.8438	.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

Normal approximation to the binomial

- X_i : independent, Bernoulli(p); $0 < p < 1$
- $S_n = X_1 + \dots + X_n$: Binomial(n, p)
 - mean np , variance $np(1 - p)$
- $n = 36, p = 0.5$; find $P(S_n \leq 21)$

$$np = 18 \quad \sqrt{np(1 - p)} = 3$$

$$P\left(\frac{S_n - 18}{3} \leq \frac{21 - 18}{3}\right)$$

$$= P(Z_n \leq 1) \approx \Phi(1) = .8413$$

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

$$\sum_{k=0}^{21} \binom{36}{k} \left(\frac{1}{2}\right)^{36} = 0.8785$$

	.00	.01	.02	.03	.04	.05	.06	.07	.08	.09
0.0	.5000	.5040	.5080	.5120	.5160	.5199	.5239	.5279	.5319	.5359
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	.8186	.8212	.8238	.8264	.8289	.8315	.8340	.8365	.8389
1.0	.8413	.8438	.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

The $1/2$ correction for integer random variables

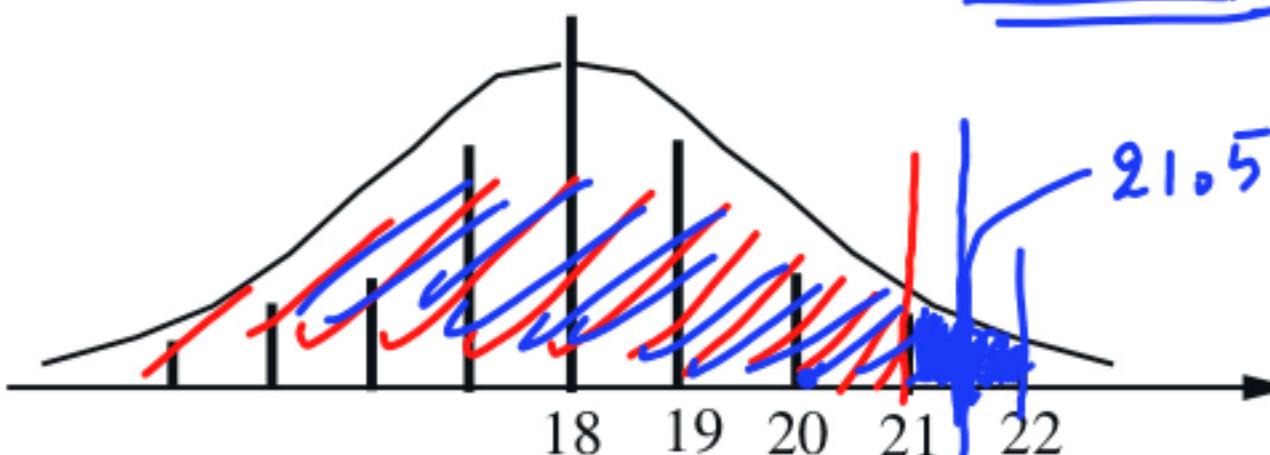
- $0.8413 \approx P(S_n \leq 21) = P(S_n < 22)$, because S_n is integer

$$= P\left(\frac{S_n - 18}{3} < \frac{22 - 18}{3}\right)$$

true value 0.8785

$$= P(Z_n < 1.33) \approx \Phi(1.33) = 0.9082$$

$$P(S_n \leq 21.5) = P(Z_n \leq \frac{21.5 - 18}{3}) \\ \approx \Phi(1.17) = 0.8790$$



	.00	.01	.02	.03	.04	.05	.06	.07	.08	.09
0.0	.5000	.5040	.5080	.5120	.5160	.5199	.5239	.5279	.5319	.5359
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	.8186	.8212	.8238	.8264	.8289	.8315	.8340	.8365	.8389
1.0	.8413	.8438	.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

De Moivre–Laplace CLT to the binomial

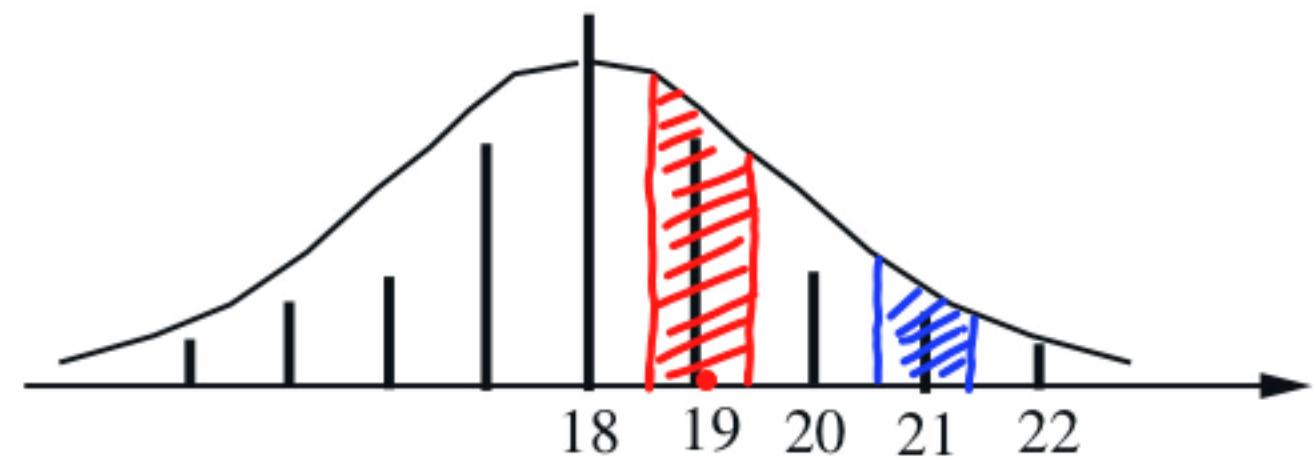
$$P(S_n = 19) = P(18.5 \leq S_n \leq 19.5)$$

$$= P\left(\frac{18.5 - 18}{3} \leq Z_n \leq \frac{19.5 - 18}{3}\right)$$

$$= P(0.17 \leq Z_n \leq 0.5)$$

$$\approx \Phi(0.5) - \Phi(0.17)$$

$$= 0.6915 - 0.5675 = 0.124$$



- Exact answer:

$$\binom{36}{19} \left(\frac{1}{2}\right)^{36} = 0.1251$$

- When the 1/2 correction is used, the CLT can also approximate the binomial PMF (not just the binomial CDF)

The pollster's problem revisited

- p : fraction of population that will vote “yes” in a referendum

- i th (randomly selected) person polled:

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

$$E[X_i] = p = \mu$$

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

- $M_n = (X_1 + \dots + X_n)/n$: fraction of “yes” in our sample

- Would like “small error,” e.g.: $|M_n - p| < 0.01$

$$P(|M_n - p| \geq .01) = P\left(|Z_n| \geq \frac{.01\sqrt{n}}{\sigma}\right) \approx P\left(|Z| \geq \frac{.01\sqrt{n}}{\sigma}\right)$$

$$Z_n = \frac{S_n - n\mu}{\sqrt{n}\sigma}$$

$$\left| \frac{S_n - np}{n} \right| \geq .01$$

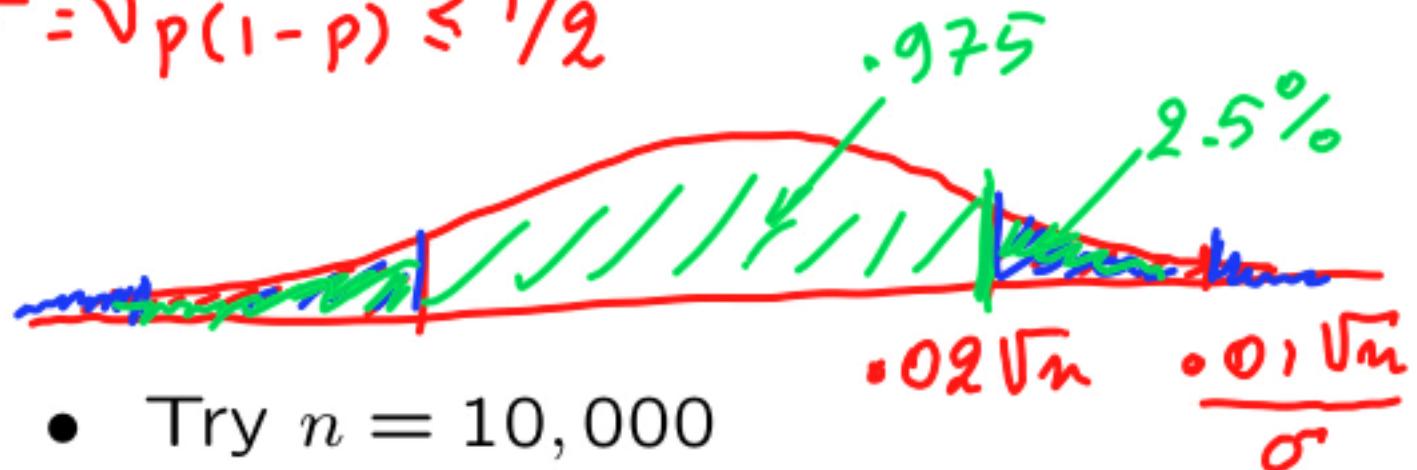
$N(0, 1)$

$$\underbrace{Z_n}_{\frac{S_n - np}{\sqrt{n}\sigma}} \geq \frac{.01\sqrt{n}}{\sigma}$$

The pollster's problem revisited

$$P(|M_n - p| \geq .01) \approx P\left(\left|Z\right| \geq \frac{.01\sqrt{n}}{\sigma}\right) \leq P\left(\left|Z\right| \geq .02\sqrt{n}\right) = 2\left(1 - \Phi\left(\frac{.02\sqrt{n}}{\sigma}\right)\right) = 0.05$$

$$\sigma = \sqrt{p(1-p)} \leq 1/2$$



- Try $n = 10,000$

$$\text{prob} \leq 2(1 - \Phi(2)) =$$

$$= 2(1 - 0.9772) = 0.046$$

- Specs: $P(|M_n - p| \geq .01) \leq .05$

$$\Phi(.02\sqrt{n}) = 0.975$$

$$.02\sqrt{n} = 1.96 \Rightarrow n = 9604$$

	.00	.01	.02	.03	.04	.05	.06	.07	.08	.09
0.0	.5000	.5040	.5080	.5120	.5160	.5199	.5239	.5279	.5319	.5359
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	.8186	.8212	.8238	.8264	.8289	.8315	.8340	.8365	.8389
1.0	.8413	.8438	.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

LECTURE 20: An introduction to classical statistics

- Unknown constant θ (not a r.v.)
- if $\theta = \mathbb{E}[X]$: estimate using the sample mean $(X_1 + \dots + X_n)/n$
 - terminology and properties
- Confidence intervals (CIs)
 - CIs using the CLT
 - CIs when the variance is unknown
- Other uses of sample means
- Maximum Likelihood estimation

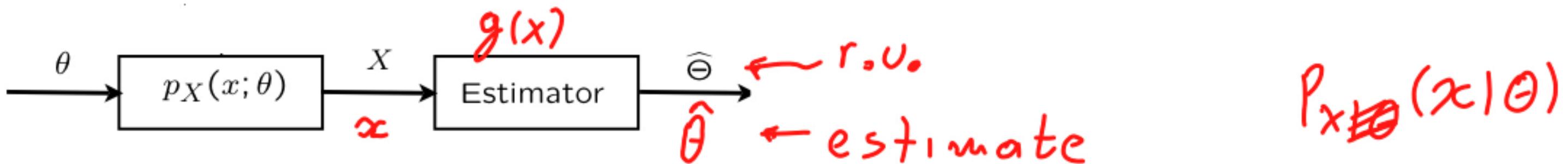
Classical statistics

- Inference using the Bayes rule:

unknown Θ and observation X are both random variables

– Find $p_{\Theta|X}$

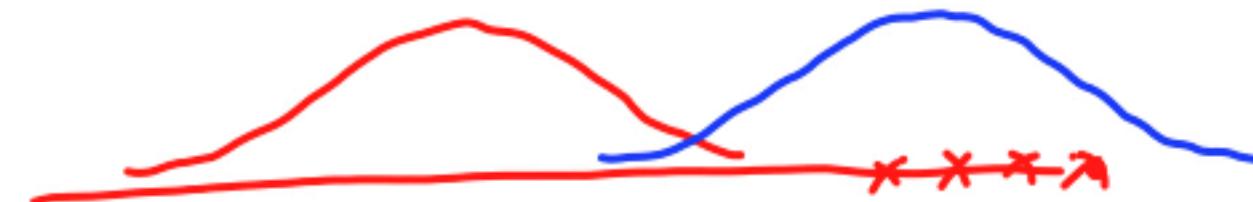
- Classical statistics: unknown constant θ



– also for vectors X and θ : $p_{X_1, \dots, X_n}(x_1, \dots, x_n; \theta_1, \dots, \theta_m)$

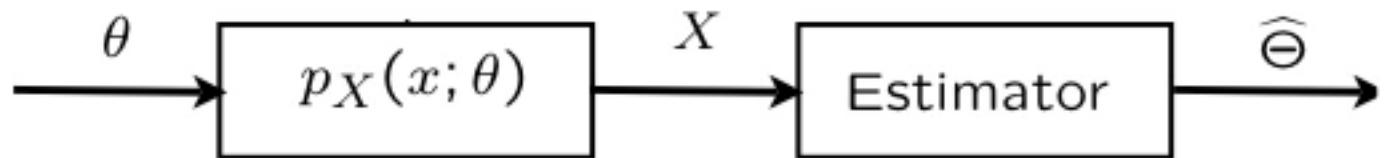
– $p_X(x; \theta)$ are NOT conditional probabilities; θ is NOT random

– mathematically: many models, one for each possible value of θ



Problem types in classical statistics

- Classical statistics: unknown constant θ



- 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** $\widehat{\Theta}$, to “keep estimation **error** $\widehat{\Theta} - \theta$ small”

Art! .

Estimating a mean

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

$$\widehat{\Theta}_n = \text{sample mean} = M_n = \frac{X_1 + \dots + X_n}{n}$$

$\widehat{\Theta}_n$: estimator (a random variable)

Properties and terminology:

- $E[\widehat{\Theta}_n] = \theta$ (unbiased)
for all θ
- WLLN: $\widehat{\Theta}_n \xrightarrow{i.p.} \theta$ (consistency)
for all θ
- mean squared error (MSE): $E[(\widehat{\Theta}_n - \theta)^2] = \text{var}(\widehat{\Theta}_n) = \frac{\sigma^2}{n}$

$$\widehat{\Theta} = g(x)$$
$$E[\widehat{\Theta}] = \sum_x g(x) P_x(x; \theta)$$

On the mean squared error of an estimator

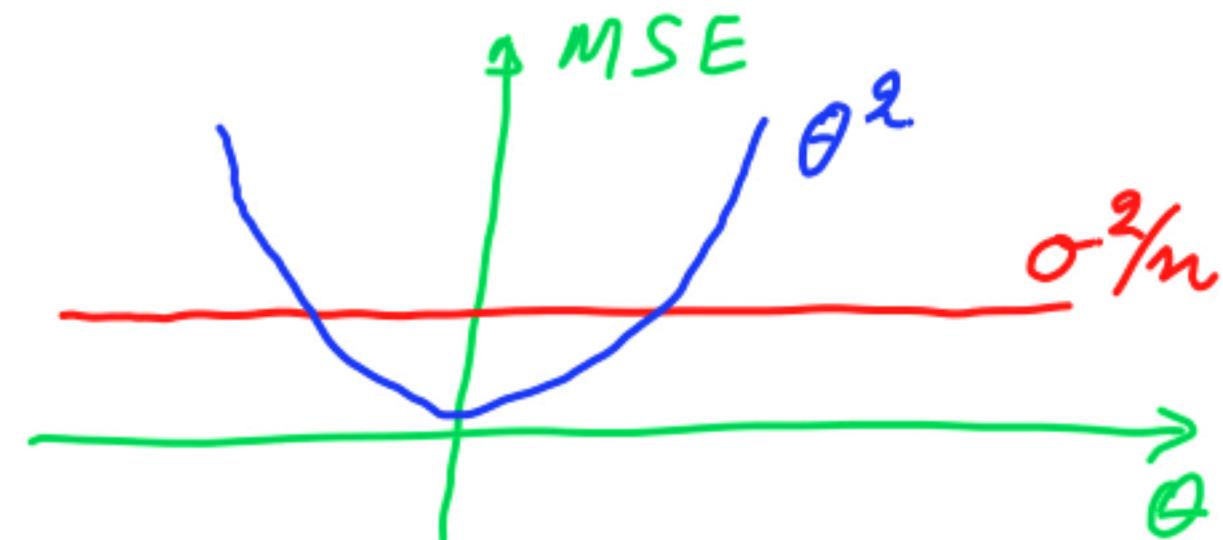
- For any estimator, using $E[Z^2] = \text{var}(Z) + (E[Z])^2$:

$$Z = \hat{\Theta} - \theta$$

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

$$\hat{\Theta}_n = M_n : \text{MSE} = \sigma^2/n + 0$$

$$\hat{\Theta} = 0 : \text{MSE} = 0 + \theta^2$$

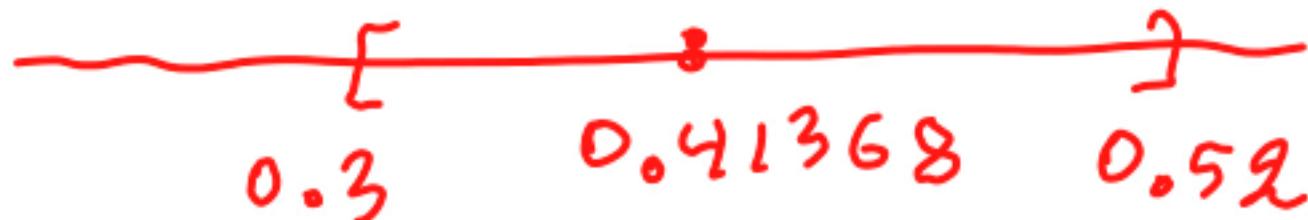
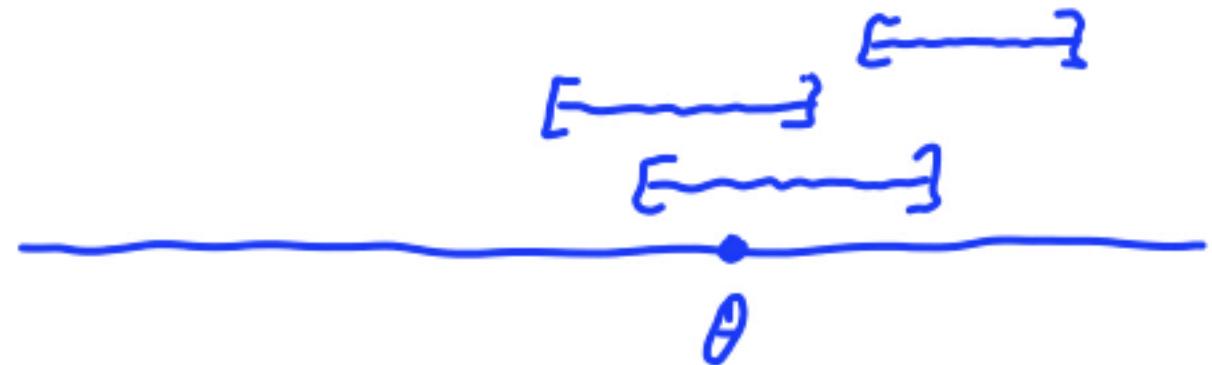


- $\sqrt{\text{var}(\hat{\Theta})}$ is called the standard error



Confidence intervals (CIs)

- The value of an estimator $\widehat{\Theta}$ may not be informative enough
95%
- An $1 - \alpha$ **confidence interval** is an interval $[\widehat{\Theta}^-, \widehat{\Theta}^+]$,
s.t. $P(\widehat{\Theta}^- \leq \theta \leq \widehat{\Theta}^+) \geq 1 - \alpha$, for all θ
 - often $\alpha = 0.05$, or 0.025 , or 0.01
 - interpretation is subtle



$$P(0.3 < \theta < 0.52) \geq 0.95$$

CI for the estimation of the mean

$$\widehat{\Theta}_n = \text{sample mean} = M_n = \frac{X_1 + \cdots + X_n}{n}$$

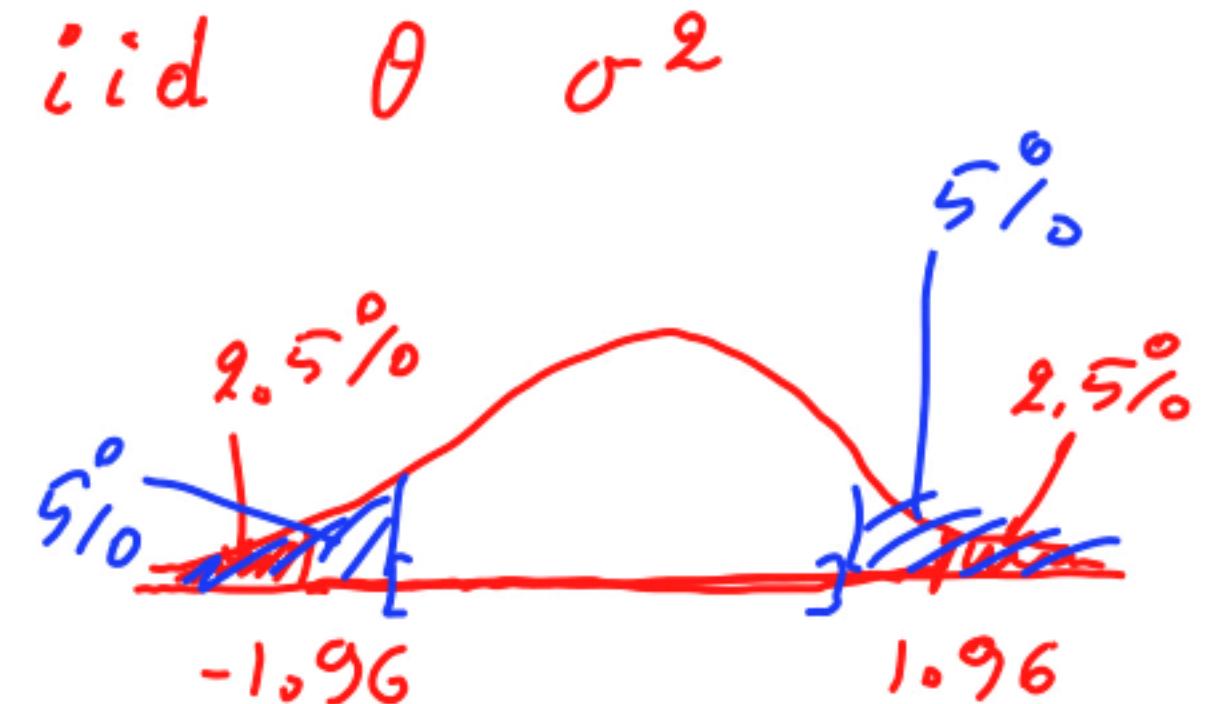
95%

normal tables: $\Phi(1.96) = 0.975 = 1 - 0.025$

90%

$$\Phi(1.645) = 0.90$$

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



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

$\widehat{\Theta}^-$

$\widehat{\Theta}^+$

Confidence intervals for the mean when σ is unknown

$$\widehat{\Theta}_n = \text{sample mean} = M_n = \frac{X_1 + \cdots + X_n}{n}$$

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

- **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: $\hat{\sigma} = \sqrt{\widehat{\Theta}_n(1 - \widehat{\Theta}_n)}$

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

Confidence intervals for the mean when σ is unknown

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

Start from $\sigma^2 = E[(X_i - \theta)^2]$

- **Option 3:** Use sample mean estimate of the variance

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

- Two approximations involved here:
 - CLT: approximately normal
 - using estimate of σ
- correction for second approximation (t -tables)
used when n is small

(but do not know θ)

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

Other natural estimators

- $\theta_X = \mathbb{E}[X]$ $\widehat{\Theta}_X = \frac{1}{n} \sum_{i=1}^n X_i$ • $\theta = \mathbb{E}[g(X)]$ $\widehat{\Theta} = \frac{1}{n} \sum_{i=1}^n g(X_i)$
- $v_X = \text{var}(X) = \mathbb{E}[(X - \theta_X)^2]$ $\widehat{v}_X = \frac{1}{n} \sum_{i=1}^n (X_i - \widehat{\Theta}_X)^2$
- $\text{cov}(X, Y) = \mathbb{E}[(X - \theta_X)(Y - \theta_Y)]$ $\widehat{\text{cov}}(X, Y) = \frac{1}{n} \sum_{i=1}^n (X_i - \widehat{\Theta}_X)(Y_i - \widehat{\Theta}_Y)$
 (x_i, y_i)
- $\rho = \frac{\text{cov}(X, Y)}{\sqrt{v_X} \cdot \sqrt{v_Y}}$ $\widehat{\rho} = \frac{\widehat{\text{cov}}(X, Y)}{\sqrt{\widehat{v}_X} \cdot \sqrt{\widehat{v}_Y}}$
- next steps: find the distribution of $\widehat{\Theta}$, MSE, confidence intervals,...

Maximum Likelihood (ML) estimation

- Pick θ that “makes data most likely”

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

- also applies when x, θ are vectors or x is continuous

- compare to Bayesian posterior: $p_{\Theta|X}(\theta | x) = \frac{p_{X|\Theta}(x | \theta)p_{\Theta}(\theta)}{p_X(x)}$

constant

- interpretation is very different

Comments on ML

- maximize $p_X(x; \theta)$
- maximization is usually done numerically
- if have n i.i.d. data drawn from model $p_X(x; \theta)$, then, under mild assumptions:
 - consistent: $\widehat{\Theta}_n \rightarrow \theta$
 - asymptotically normal: $\frac{\widehat{\Theta}_n - \theta}{\sigma(\widehat{\Theta}_n)} \rightarrow N(0, 1)$ (CDF convergence)
- analytical and simulation methods for calculating $\widehat{\sigma} \approx \sigma(\widehat{\Theta}_n)$
 - hence confidence intervals $P\left(\widehat{\Theta}_n - 1.96 \widehat{\sigma} \leq \theta \leq \widehat{\Theta}_n + 1.96 \widehat{\sigma}\right) \approx 0.95$
 - asymptotically “efficient” (“best”)

ML estimation example: parameter of binomial

- K : binomial with parameters n (known), and θ (unknown)

k

$$p_K(k; \theta) = \binom{n}{k} \theta^k (1 - \theta)^{n-k}$$

$$\log \left[\binom{n}{k} \right] + k \log \theta + (n - k) \log (1 - \theta)$$

$$0 + \frac{k}{\theta} - \frac{n - k}{1 - \theta} = 0 \Rightarrow k - k\theta = n\theta - k\theta$$

$$\hat{\theta}_{\text{ML}} = \frac{k}{n} \quad \widehat{\Theta}_{\text{ML}} = \frac{K}{n}$$

- same as MAP estimator with uniform prior on θ

ML estimation example — normal mean and variance

- X_1, \dots, X_n : i.i.d., $N(\mu, v)$ $f_X(x; \mu, v) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi v}} \exp\left\{-\frac{(x_i - \mu)^2}{2v}\right\}$

minimize $\frac{n}{2} \log v + \sum_{i=1}^n \frac{(x_i - \mu)^2}{2v}$

– minimize w.r.t. μ : $\hat{\mu} = \frac{x_1 + \dots + x_n}{n}$

$$\frac{1}{n} \sum_{i=1}^n (x_i - \mu) = 0 \Rightarrow \sum x_i = n\mu$$

– minimize w.r.t. v : $\hat{v} = \frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu})^2$

$$\cancel{\frac{n}{2} \cdot \frac{1}{v}} \rightarrow \sum_{i=1}^n \frac{(x_i - \mu)^2}{\cancel{2v^2}} = 0$$