

$$S_n = 1 + \frac{1}{2} + \dots + \frac{1}{n}$$

$$\lim_{n \rightarrow \infty} P\left(\bigcup_{i=1}^n A_i\right) \\ = P\left(\bigcup_{i=1}^{\infty} A_i\right)$$

Lecture 3: Random Variables

Ziyu Shao

School of Information Science and Technology
ShanghaiTech University

October 15, 2024

Outline

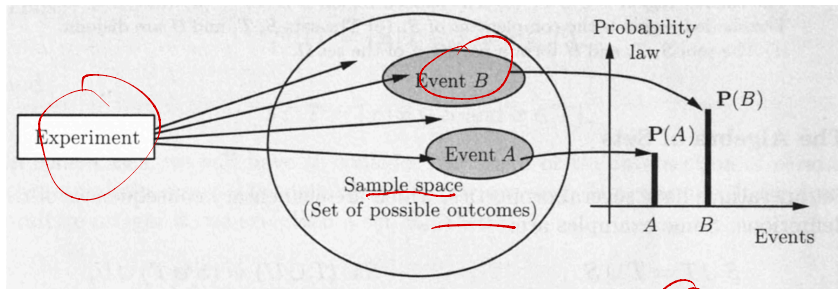
- 1 Random Variables
- 2 Bernoulli and Binomial
- 3 Hypergeometric
- 4 Discrete Uniform & Zipf Distribution
- 5 Cumulative Distribution Functions
- 6 Functions of Random Variables: Random Variables
- 7 Independence of R.V.s
- 8 Binomial & Hypergeometric
- 9 Information Theory & Entropy

Outline

- 1 Random Variables
- 2 Bernoulli and Binomial
- 3 Hypergeometric
- 4 Discrete Uniform & Zipf Distribution
- 5 Cumulative Distribution Functions
- 6 Functions of Random Variables: Random Variables
- 7 Independence of R.V.s
- 8 Binomial & Hypergeometric
- 9 Information Theory & Entropy

Probabilistic Model ^{Motivation}

U, Ω, \dots
 $+, -, \times, \div$



Motivation 2: event \rightarrow \mathbb{R}
Language \rightarrow Numerical value

Definition of Random Variables

$$X: S \rightarrow R$$

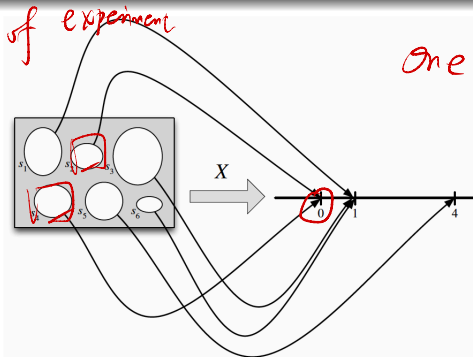
Definition

\otimes . deterministic

Given an experiment with sample space S , a random variable (r.v.) is a function from the sample space S to the real numbers R . It is common, but not required, to denote random variables by capital letters.

X : Summary of experiment

~~one-one mapping~~



Example: Coin Tosses

$$S \rightarrow \mathcal{R}$$

$$\textcircled{1} X(CHH) = 2; X(HT) = 1; X(TH) = 1;$$

$$X(TT) = 0; \quad \{0, 1, 2\}$$

Consider an experiment where we toss a fair coin twice. The sample space consists of four possible outcomes: $S = \{HH, HT, TH, TT\}$. Here are some random variables on this space (for practice, you can think up some of your own). Each r.v. is a numerical summary of some aspect of the experiment

• X : the number of Heads.

• Y : the number of Tails.

• I : equals 1 if the first toss lands Heads and 0 otherwise.

$\textcircled{2}$ $Y(s)$ sample point $s, s \in S$,
 $= 2 - X(s)$ $(Y(CHH) = 0)$
 $Y(CHH) = 2 - X(CHH) = 2 - 2 = 0$

$$\textcircled{3} I(CHH) = 1; I(HT) = 1;$$

$$I(TH) = 0; I(TT) = 0;$$

Discrete Random Variable

$$X = a_j \quad \underline{A}$$

$$P(X = a_j)$$

event

$$\{s : X(s) = a_j\}$$

$$= P(\{s : X(s) = a_j\})$$

set of sample points

$$\subseteq S$$

Definition

A random variable X is said to be *discrete* if there is a finite list of values a_1, a_2, \dots, a_n or an infinite list of values a_1, a_2, \dots such that $P(X = a_j \text{ for some } j) = 1$. If X is a discrete r.v., then the finite or countably infinite set of values x such that $P(X = x) > 0$ is called the support of X .

fair

Support of X is $\{0, 1\}$

coin tossing $\{H, T\}$ X : # of heads.

$$P(X=1) = P(\text{"H"}) = \frac{1}{2} > 0$$

$$P(X=0) = P(\text{"T"}) = \frac{1}{2} > 0$$

$X = 1$ or 0

Probability Mass Function

$$X = x$$

$$\underline{p_X(x)} \stackrel{\Delta}{=} P(X=x) \\ = P(\{s: X(s)=x\})$$

Definition

The probability mass function (PMF) of a discrete r.v. X is the function p_X given by $p_X(x) = P(X = x)$. Note that this is positive if x is in the support of X , and 0 otherwise.

$$X: \underbrace{S} \rightarrow \underline{R_X} \subset \mathbb{R}$$

$$x \in R_X$$

Example: Coin Tosses $\{X=x\} \triangleq \{s: X(s)=x, s \in S\}$

① $X=0 \Leftrightarrow \{s\}=\{TT\}; P_X(0) = P(X=0) = P(\{TT\}) = \frac{1}{4}$.

$X=1 \Leftrightarrow s = HT, TH; P_X(1) = P(X=1) = P(\{HT, TH\}) = \frac{1}{2};$

$X=2 \Leftrightarrow s = HH; P_X(2) = P(X=2) = P(\{HH\}) = \frac{1}{4}$

Consider an experiment where we toss a fair coin twice. The sample space consists of four possible outcomes: $S = \{HH, HT, TH, TT\}$.

Here are some random variables on this space (for practice, you can think up some of your own). Each r.v. is a numerical summary of some aspect of the experiment

$X(s)$ • X: the number of Heads.

$Y(s)$ • Y: the number of Tails.

$Z(s)$ • Z: equals 1 if the first toss lands Heads and 0 otherwise.

PMF. $R_X = \{0, 1, 2\}$ $P_X(0) = \frac{1}{4}$
 $X \in R_X = \{0, 1, 2\}$ $P_X(1) = \frac{1}{2}$
 $P_X(2) = \frac{1}{4}$ }
 Support of $X = (R_X)$.

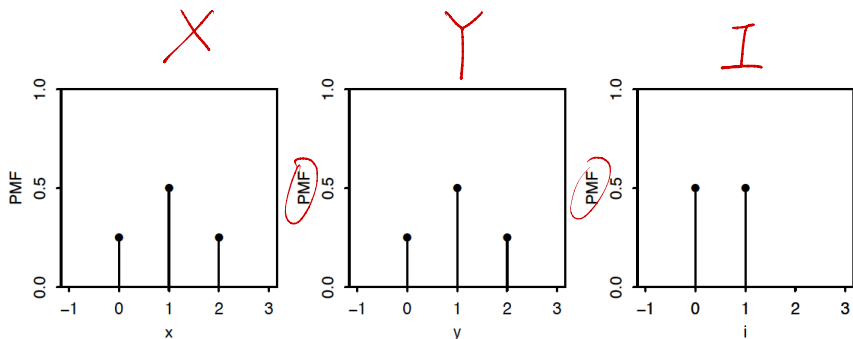
② $P_Y(x) = P_X(x), \forall x \in \{0, 1, 2\}$

③ $P_Z(0) = P(Z=0) = P(\{TH, TT\}) = \frac{1}{2}$

$P_Z(1) = P(Z=1) = P(\{HH, HT\}) = \frac{1}{2}$

PMF $\Rightarrow \left\{ \begin{array}{l} P_Z(x) = \frac{1}{2} \\ x=0 \text{ or } 1. \end{array} \right.$

Example: Coin Tosses



distribution of X

Valid PMFs

Theorem

Let X be a discrete r.v. with support x_1, x_2, \dots (assume these values are distinct and, for notational simplicity, that the support is countably infinite; the analogous results hold if the support is finite). The PMF p_X of X must satisfy the following two criteria:

- Nonnegative: $p_X(x) > 0$ if $x = x_j$ for some j , and $p_X(x) = 0$ otherwise;
- Sums to 1: $\sum_{j=1}^{\infty} p_X(x_j) = 1$.

Outline

- 1 Random Variables
- 2 Bernoulli and Binomial
- 3 Hypergeometric
- 4 Discrete Uniform & Zipf Distribution
- 5 Cumulative Distribution Functions
- 6 Functions of Random Variables: Random Variables
- 7 Independence of R.V.s
- 8 Binomial & Hypergeometric
- 9 Information Theory & Entropy

Math:

$S \xrightarrow{\text{r.v.}} R$ r.v. mapping

Modeling:

Random Variable.

Bernoulli Distribution

(PMF)

$$P(X=x) = p^x (1-p)^{1-x}, \quad x=0 \text{ or } 1$$

$$P(X; \theta) = \theta^x (1-\theta)^{1-x}$$

$0 < \theta < 1$

Family of Bernoulli
distributions

Definition

An r.v. X is said to have the *Bernoulli distribution* with parameter p if $P(X=1) = p$ and $P(X=0) = 1-p$, where $0 < p < 1$. We write this as $X \sim \text{Bern}(p)$. The symbol \sim is read "is distributed as".

Indicator Random Variable

$$I_A = \begin{cases} 1 & \text{if } A \text{ occurs.} \\ 0 & \text{otherwise.} \end{cases}$$

$P(A) = p$

Definition

The indicator random variable of an event A is the r.v. which equals 1 if A occurs and 0 otherwise. We will denote the indicator r.v. of A by I_A or $I(A)$. Note that $I_A \sim \text{Bern}(p)$ with $p = P(A)$.

indicator function \neq r.v.

$$I_A(x) = \begin{cases} 1 & \text{if } x \in A. \\ 0 & \text{otherwise.} \end{cases}$$

Story: Bernoulli Trial

An experiment that can result in either a “success” or a “failure” (but not both) is called a *Bernoulli trial*. A Bernoulli random variable can be thought of as the *indicator of success* in a Bernoulli trial: it equals 1 if success occurs and 0 if failure occurs in the trial.

Story: Binomial Distribution

$$X \in \{0, 1, \dots, n\}$$

Suppose that n *independent* Bernoulli trials are performed, each with the same success probability p . Let X be the number of successes. The distribution of X is called the Binomial distribution with parameters n and p . We write $X \sim \text{Bin}(n, p)$ to mean that X has the Binomial distribution with parameters n and p , where n is a positive integer and $0 < p < 1$.

Binomial PMF

1: success p

0: failure $(1-p)$

0001010101...

length n ... binary sequence

Theorem

If $X \sim \text{Bin}(n, p)$, then the PMF of X is

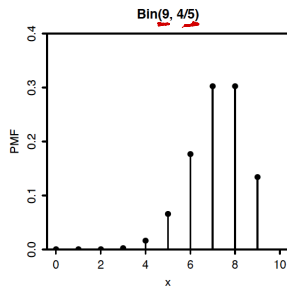
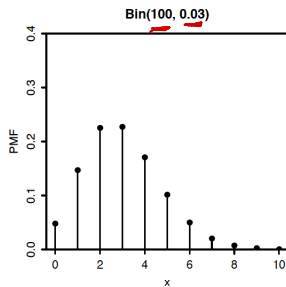
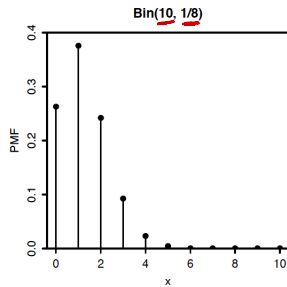
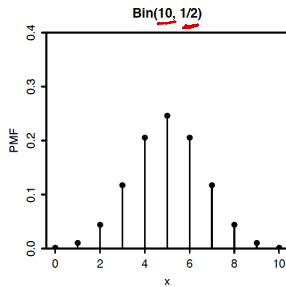
$$P(X = k) = \binom{n}{k} p^k (1-p)^{n-k}$$

k "1"

$n-k$ "0"

for $k = 0, 1, \dots, n$ (and $P(X = k) = 0$ otherwise).

Binomial PMF



Binomial PMF

Theorem

Let $X \sim \text{Bin}(n, p)$, and $q = 1 - p$ (we often use q to denote the failure probability of a Bernoulli trial). Then $n - X \sim \text{Bin}(n, q)$.

Example: Statistical Multiplexing

X : # of active users at the same time.

User (sender) in or out Bernoulli trial.

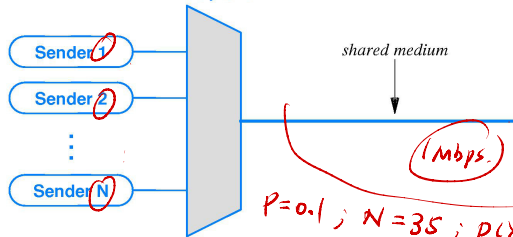
$$X \sim \text{Bin}(N, p)$$

each sender (user) request μ kbps.
active prob. (p) .

NO statistical multiplexing

$$P(X=k) = \binom{N}{k} p^k (1-p)^{n-k}$$

$k \in \{0, 1, \dots, N\}$

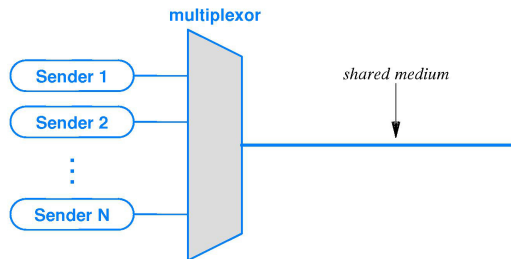


$$p=0.1; N=35; P(X>10) < 0.0004$$

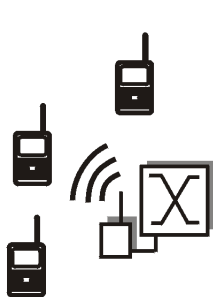
System overload: $P(X > 10) = 1 - P(X \leq 10)$

$$= 1 - \sum_{k=0}^{10} \binom{N}{k} p^k (1-p)^{n-k} = 1 - \sum_{k=0}^{10} P(X=k)$$

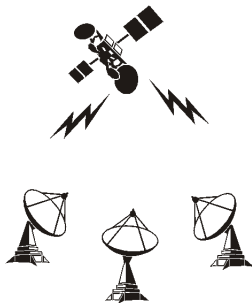
Example: Statistical Multiplexing



Example: Multiple Access (Aloha Protocol)



shared wireless



satellite

Example: Multiple Access (Aloha Protocol)



X = # of devices transmit simultaneously.

time-slotted Aloha

$$X \sim \text{Bin}(N, p)$$

- N smart devices sharing a WiFi access point (e.g., in starbucks)
- ≥ 2 devices transmit simultaneously lead to collision
- Aloha Protocol: proposed by Norman Abramson in the later 1960s
- Each device transmits with probability p independently
- What is the transmission rate (the number of successful transmissions per unit time)?

$$f(p) = \frac{P(X=1) \cdot 1}{1} = P(X=1)$$

Maximize $f(p)$, $0 \leq p \leq 1$.

$$f'(p) = 0; \quad f'(p) < 0 \Rightarrow p^* = \frac{1}{N} \quad = \binom{N}{1} p^1 (1-p)^{N-1} = N p (1-p)^{N-1}$$

$$\Rightarrow f^*(p) = \left(1 - \frac{1}{N}\right)^{N-1} \xrightarrow{N \rightarrow \infty} e^{-1} \approx 0.36$$

Example: Multiple Access (Aloha Protocol)

Outline

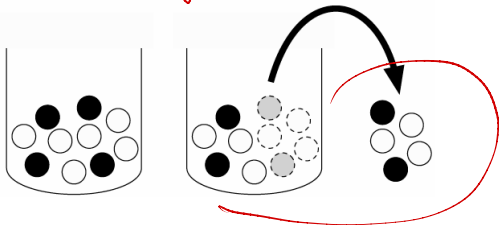
- 1 Random Variables
- 2 Bernoulli and Binomial
- 3 Hypergeometric**
- 4 Discrete Uniform & Zipf Distribution
- 5 Cumulative Distribution Functions
- 6 Functions of Random Variables: Random Variables
- 7 Independence of R.V.s
- 8 Binomial & Hypergeometric
- 9 Information Theory & Entropy

Urn Model

$$\text{Bern}\left(\frac{w}{w+b}\right)$$

An urn is filled with w white and b black balls, then drawing n balls out of the urn

- with replacement: $\text{Bin}(n, w/(w+b))$ distribution for the number of white balls obtained
- without replacement: Hypergeometric distribution



Story: Hypergeometric Distribution

$$P(X=k) = \frac{\binom{w}{k} \cdot \binom{b}{n-k}}{\binom{w+b}{n}} \quad \left\{ \begin{array}{l} 0 \leq k \leq w \\ 0 \leq n-k \leq b \end{array} \right.$$

$n \leq w+b$

Consider an urn with w white balls and b black balls. We draw n balls out of the urn at random without replacement, such that all $\binom{w+b}{n}$ samples are equally likely. Let X be the number of white balls in the sample. Then X is said to have the *Hypergeometric distribution* with parameters w , b , and n ; we denote this by $X \sim \text{HGeom}(w, b, n)$.

Hypergeometric PMF

Theorem

If $X \sim \text{HGeom}(w, b, n)$, then the PMF of X is

$$P(X = k) = \frac{\binom{w}{k} \binom{b}{n-k}}{\binom{w+b}{n}},$$

for integers k satisfying $0 \leq k \leq w$ and $0 \leq n - k \leq b$, and $P(X = k) = 0$ otherwise.

Identical Distribution

$$X \sim \underline{HGeom(w, b, n)}$$

first sampled tag;

second color tag;

X : # of white balls in sampled balls.

Color tag: white/black

Sampled tag: Yes or No

Theorem

The $\underline{HGeom(w, b, n)}$ and $\underline{HGeom(n, w + b - n, w)}$ distributions are identical. That is, if $X \sim \underline{HGeom(w, b, n)}$ and $Y \sim \underline{HGeom(n, w + b - n, w)}$, then X and Y have the same distribution.

Y : first color tag.

second sampled tag;

Outline

- 1 Random Variables
- 2 Bernoulli and Binomial
- 3 Hypergeometric
- 4 Discrete Uniform & Zipf Distribution**
- 5 Cumulative Distribution Functions
- 6 Functions of Random Variables: Random Variables
- 7 Independence of R.V.s
- 8 Binomial & Hypergeometric
- 9 Information Theory & Entropy

Story: Discrete Uniform Distribution

Let C be a finite, nonempty set of numbers. Choose one of these numbers uniformly at random (i.e., all values in C are equally likely). Call the chosen number X . Then X is said to have the *Discrete Uniform distribution* with parameter C ; we denote this by $X \sim \text{DUnif}(C)$.

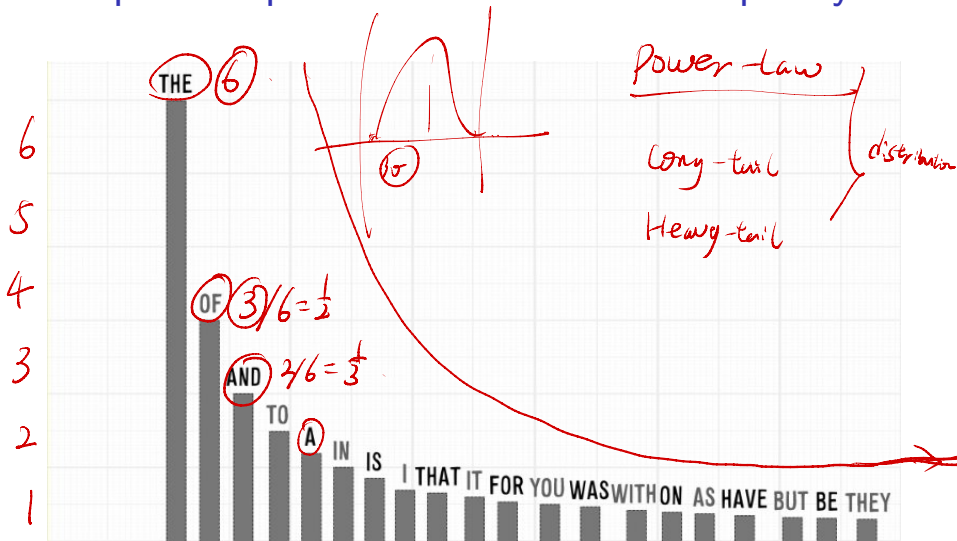
$$P(X=x) = \frac{1}{|C|} \quad (x \in C)$$

Zipf Distribution

- Zipf's Law & Zipf distribution: American linguist George Kingsley Zipf (1902-1950)
- Popularity distribution: popularity of the i^{th} most popular term is proportional to $1/i$.
- If $X \sim \text{Zipf}(\alpha > 0)$, then PMF of X is:

$$P(X = k) = \frac{1}{k^{\alpha+1} \sum_{j=1}^{\infty} (\frac{1}{j})^{\alpha+1}}, k = 1, 2, \dots$$

Example of Zipf Distribution: Word Frequency



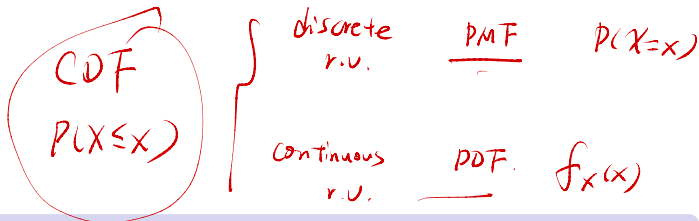
Examples of Zipf Distribution

- The world population lives in several large cities, a greater number of medium-sized cities, and a vast number of small towns.
- There are a few websites that get lots of hits, a greater number of websites that get a moderate number of hits, and a vast number of websites that hardly get any hits at all.
- A library has a few books that everyone wants to borrow (best sellers), a greater number of books that get borrowed occasionally (classics), and a vast number of books that hardly ever get borrowed.

Outline

- 1 Random Variables
- 2 Bernoulli and Binomial
- 3 Hypergeometric
- 4 Discrete Uniform & Zipf Distribution
- 5 Cumulative Distribution Functions**
- 6 Functions of Random Variables: Random Variables
- 7 Independence of R.V.s
- 8 Binomial & Hypergeometric
- 9 Information Theory & Entropy

Definition



Theorem

The cumulative distribution function (CDF) of an r.v. X is the function F_X given by $F_X(x) = P(X \leq x)$. When there is no risk of ambiguity, we sometimes drop the subscript and just write F (or some other letter) for a CDF.

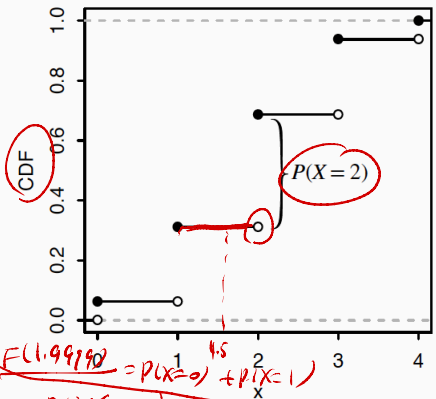
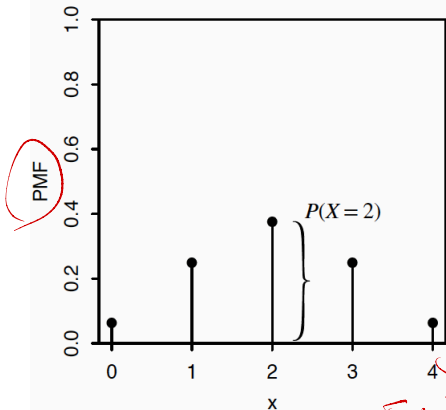
Example

$$F(1.5) = P(X \leq 1.5)$$

$$= P(X=0) + P(X=1)$$

Let $X \sim \text{Bin}(4, 1/2)$, the PMF and CDF of X :

$$F(1.6) = F(1.5)$$



$$F(1.9999) = P(X=0) + P(X=1)$$

$$F(2) = P(X \leq 2)$$

$$= P(X=0) + P(X=1) + P(X=2)$$

Example

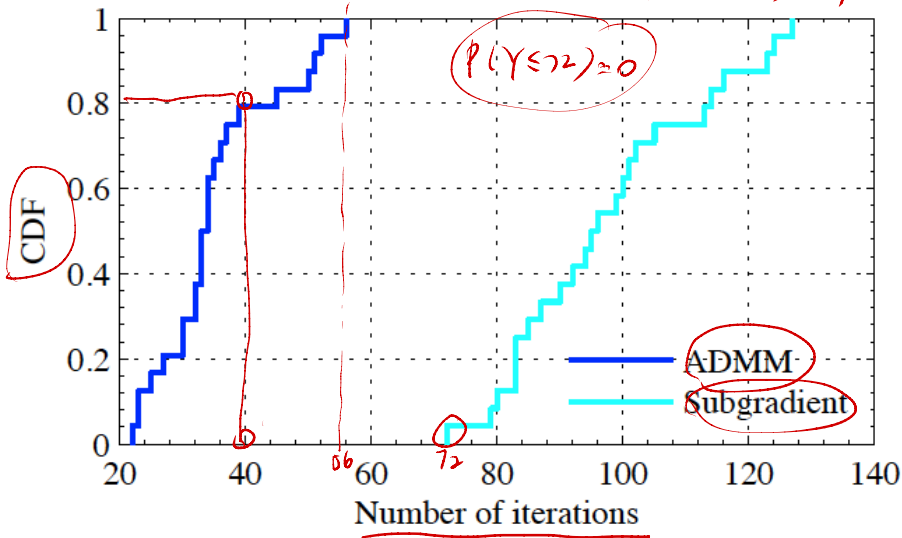
$X = \#$ of iterations of ADMM

$$P(X \leq 40) = 0.8$$

Y_2 Subgradient

$$P(X \leq 56) = 1$$

$$P(Y \leq 72) = 0$$



Valid CDFs

Any CDF F has the following properties.

- Increasing: If $x_1 \leq x_2$, then $F(x_1) \leq F(x_2)$.
- Right-continuous: the CDF is continuous except possibly for having some jumps. Wherever there is a jump, the CDF is continuous from the right. That is, for any a , we have

$$F(a) = \lim_{x \rightarrow a^+} F(x).$$

- Convergence to 0 and 1 in the limits:

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

$$\lim_{x \rightarrow -\infty} F(x) = 0 \quad \text{and} \quad \lim_{x \rightarrow \infty} F(x) = 1$$

Outline

- 1 Random Variables
- 2 Bernoulli and Binomial
- 3 Hypergeometric
- 4 Discrete Uniform & Zipf Distribution
- 5 Cumulative Distribution Functions
- 6 Functions of Random Variables: Random Variables**
- 7 Independence of R.V.s
- 8 Binomial & Hypergeometric
- 9 Information Theory & Entropy

Definition

$$\begin{array}{ccccc} S & \xrightarrow{X} & \mathbb{R} & \xrightarrow{g} & \mathbb{R} \\ s & \longrightarrow & X(s) & \longrightarrow & g(X(s)) \end{array}$$

Theorem

For an experiment with sample space S , an r.v. X , and a function $g : \mathbb{R} \rightarrow \mathbb{R}$, $g(X)$ is the r.v. that maps s to $g(X(s))$ for all $s \in S$.

PMF of $g(X)$

$$g(X) = X^2$$

Theorem

Let X be a discrete r.v. and $g : \mathbb{R} \rightarrow \mathbb{R}$. Then the support of $g(X)$ is the set of all y such that $g(x) = y$ for at least one x in the support of X , and the PMF of $g(X)$ is

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

for all y in the support of $g(X)$.

Example: Maximum of Two Die Rolls

$$\max\{X, Y\} \text{ r.v. } \in \{1, 2, 3, 4, 5, 6\}$$

$$\begin{aligned} 1^\circ. P(\max(X, Y) = 1) &= P(X=1, Y=1) = P(X=1)P(Y=1) \\ &= \frac{1}{6} \times \frac{1}{6} = \frac{1}{36}. \end{aligned}$$

We roll two fair 6-sided dice. Let X be the number on the first die and Y the number on the second die. What is the PMF of $\max(X, Y)$.

$$\begin{aligned} 2^\circ. P(\max(X, Y) = 2) &= P(X=1, Y=2) + P(X=2, Y=1) \\ &\quad + P(X=2, Y=2) = \frac{3}{36} = \frac{1}{12}; \end{aligned}$$

$$3^\circ. P(\max(X, Y) = c) = \begin{cases} \frac{5}{36} & c=3 \\ \frac{4}{36} & c=5 \\ \frac{3}{36} & c=4 \end{cases} \quad \frac{11}{36} \quad c=6$$

Example: Sympathetic Magic

$$\textcircled{1} \quad Y = 2X \quad P(Y=y) \neq 2P(X=y)$$

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

- Given an r.v. X , trying to get the PMF of $2X$ by multiplying the PMF of X by 2.
- Claiming that because X and Y have the same distribution, X must always equal Y , i.e., $P(X=Y) = 1$.

$\textcircled{2}$ Toss coin ^(H, T)
fair coin
 $X =$ indicator of Head. event.
 $Y =$ Tail.
 $P(X=1) = P(X=0) = \frac{1}{2}$
 $P(Y=1) = P(Y=0) = \frac{1}{2}$

$X, Y \sim \text{Bern}(1/2)$. But $X \neq Y$.

$X+Y=1$ \Rightarrow if $X=Y \Rightarrow X=Y=1/2$

Outline

- 1 Random Variables
- 2 Bernoulli and Binomial
- 3 Hypergeometric
- 4 Discrete Uniform & Zipf Distribution
- 5 Cumulative Distribution Functions
- 6 Functions of Random Variables: Random Variables
- 7 Independence of R.V.s**
- 8 Binomial & Hypergeometric
- 9 Information Theory & Entropy

Independence of Two R.V.s

Definition

Random variables X and Y are said to be *independent* if

$$P(X \leq x, Y \leq y) = P(X \leq x) P(Y \leq y),$$

for all $x, y \in \mathbb{R}$. In the discrete case, this is equivalent to the condition

$$P(X = x, Y = y) = P(X = x) P(Y = y)$$

for all x, y with x in the support of X and y in the support of Y .

Independence of Many R.V.s

Events A, B, C independent

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

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

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

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

Definition

Random variables X_1, \dots, X_n are independent if $\prod_{i=1}^n$ individual CDF;

$$P(X_1 \leq x_1, \dots, X_n \leq x_n) = P(X_1 \leq x_1) \cdots P(X_n \leq x_n)$$

for all $x_1, \dots, x_n \in \mathbb{R}$. For infinitely many r.v.s, we say that they are independent if every finite subset of the r.v.s is independent.

I.I.D.

We will often work with random variables that are independent and have the same distribution. We call such r.v.s independent and identically distributed, or i.i.d. for short.

- Independent & Identically Distributed
- Independent & NOT Identically Distributed
- Dependent & Identically Distributed
- Dependent & NOT Identically Distributed

I.I.D.

Binomial Distribution

Theorem

If $X \sim \text{Bin}(n, p)$, viewed as the number of successes in n independent Bernoulli trials with success probability p , then we can write $X = X_1 + \dots + X_n$ where the X_i are i.i.d. $\text{Bern}(p)$.

$$X_i = \mathbb{I}_{\{\text{trial } i \text{ is successful}\}}$$

Binomial Distribution

Theorem

If $X \sim \text{Bin}(n, p)$, $Y \sim \text{Bin}(m, p)$, and X is independent of Y , then $X + Y \sim \text{Bin}(n + m, p)$

Proof 1: LOTP

$$X \sim \text{Bin}(n, p) ; Y \sim \text{Bin}(m, p).$$

$$\underline{X \perp\!\!\!\perp Y}$$

$$0 \leq k \leq n+m$$

$$P(X+Y=k)$$

$$\stackrel{\text{LOTP}}{=} \sum_{j=0}^k P(X+Y=k | X=j) \cdot P(X=j)$$

$$= \sum_{j=0}^k P(Y=k-j | X=j) \cdot P(X=j)$$

$$= \sum_{j=0}^k P(Y=k-j) \cdot P(X=j)$$

$X \perp\!\!\!\perp Y.$

$$= \sum_{j=0}^k \binom{m}{k-j} \cdot p^{k-j} (1-p)^{m-k+j} \cdot \binom{n}{j} \cdot p^j (1-p)^{n-j}$$

$$= \left(\sum_{j=0}^k \binom{m}{k-j} \cdot \binom{n}{j} \right) p^k (1-p)^{n+m-k}$$

$$= \binom{n+m}{k} \cdot p^k (1-p)^{n+m-k}$$

$X+Y \sim \text{Bin}(n+m, p)$

$$\binom{n+m}{k}$$

$$= \sum_{j=0}^k \binom{m}{k-j} \binom{n}{j}$$

Strong proof.

Proof 2: Representation

$$\underline{X \sim \text{Bin}(n, p)}. \quad X = X_1 + \dots + X_n, \quad \begin{array}{l} X_i \sim \text{i.i.d.} \\ \text{Bern}(p). \end{array}$$

$$Y \sim \text{Bin}(m, p), \quad Y = Y_1 + \dots + Y_m, \quad \begin{array}{l} Y_j \sim \text{i.i.d.} \\ \text{Bern}(p). \end{array}$$

$$X + Y = \underbrace{(X_1 + \dots + X_n) + (Y_1 + \dots + Y_m)}_{n+m \text{ i.i.d. Bern}(p)}$$

$$\Rightarrow X + Y \sim \text{Bin}(n+m, p)$$

Proof 3: Story

X : # of success in n independent Bernoulli trials.

each trial.
Success with prob. p

Y :

m

$X+Y$: # of total success in $n+m$

$$X+Y \sim \text{Bin}(n+m, p)$$

Conditional Independence of R.V.s

$$\hat{P}(\cdot) = P(\cdot | Z = z)$$

Definition

Random variables X and Y are *conditionally independent* given an r.v. Z if for all $x, y \in \mathbb{R}$ and all z in the support of Z ,

$$P(X \leq x, Y \leq y | Z = z) = P(X \leq x | Z = z) P(Y \leq y | Z = z).$$

For discrete r.v.s, an equivalent definition is to require

$$P(X = x, Y = y | Z = z) = P(X = x | Z = z) P(Y = y | Z = z).$$

Conditional PMF

PMF

$$\{ p(X=x), x \in \text{support} \}$$

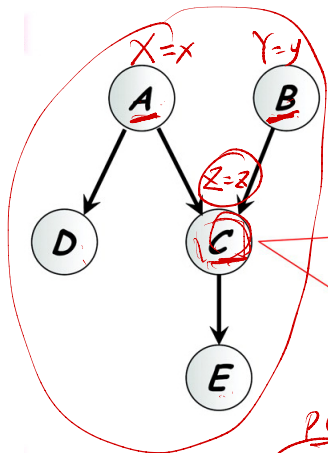
$$\{ p(X=x|Z=z), x \in \text{support} \}$$

Conditional PMF

Definition

For any discrete r.v.s X and Z , the function $P(X = x|Z = z)$, when considered as a function of x for fixed z , is called the conditional PMF of X given $Z = z$.

Example: Bayesian Network



C

$$P(X=x | Y=y, Z=z)$$

| | | $P(C A, B)$ | |
|---|---|---------------|------|
| A | B | 0 | 1 |
| 0 | 0 | 0.9 | 0.1 |
| 0 | 1 | 0.2 | 0.8 |
| 1 | 0 | 0.9 | 0.1 |
| 1 | 1 | 0.01 | 0.99 |

$$P(A, B, C, D, E) = P(A) P(B) P(C | A, B) P(D | A) P(E | C)$$

Example: Bayesian Network

- A probabilistic graphical model proposed by Judea Pearl in 1985
- Represents a set of random variables and their conditional dependencies
- Node: random variables
- Edge: conditional dependency
- Topology: a directed acyclic graph (DAG)
- Each node has a conditional probability table (CPT) with input from its parent nodes.
- Popular models for inference and learning

Outline

- 1 Random Variables
- 2 Bernoulli and Binomial
- 3 Hypergeometric
- 4 Discrete Uniform & Zipf Distribution
- 5 Cumulative Distribution Functions
- 6 Functions of Random Variables: Random Variables
- 7 Independence of R.V.s
- 8 Binomial & Hypergeometric**
- 9 Information Theory & Entropy

Connection

- Binomial \implies Hypergeometric: **conditioning**
- Hypergeometric \implies Binomial: **taking a limit**

Connection

$$\textcircled{1} X+Y \sim \text{Bin}(n+m, p)$$

$$\begin{aligned} \textcircled{2} P(X=x | X+Y=r) &= \frac{P(X=x, X+Y=r)}{P(X+Y=r)} \\ &= \frac{P(X=x, Y=r-x)}{P(X+Y=r)} \stackrel{X \perp Y}{=} \frac{P(X=x) P(Y=r-x)}{P(X+Y=r)} \end{aligned}$$

Theorem

If $X \sim \text{Bin}(n, p)$, $Y \sim \text{Bin}(m, p)$, and X is independent of Y , then the conditional distribution of X given $X+Y=r$ is $\text{HGeom}(n, m, r)$.

$$\begin{aligned} &= \frac{\binom{n}{x} \cancel{p^x (1-p)^{n-x}} \cdot \binom{m}{r-x} \cancel{p^{r-x} (1-p)^{m-r+x}}}{\binom{n+m}{r} \cancel{p^r (1-p)^{n+m-r}}} \\ &= \frac{\binom{n}{x} \binom{m}{r-x}}{\binom{n+m}{r}} \end{aligned}$$

Connection

$$N = w + b \rightarrow \infty.$$

Sampling with/without replacement

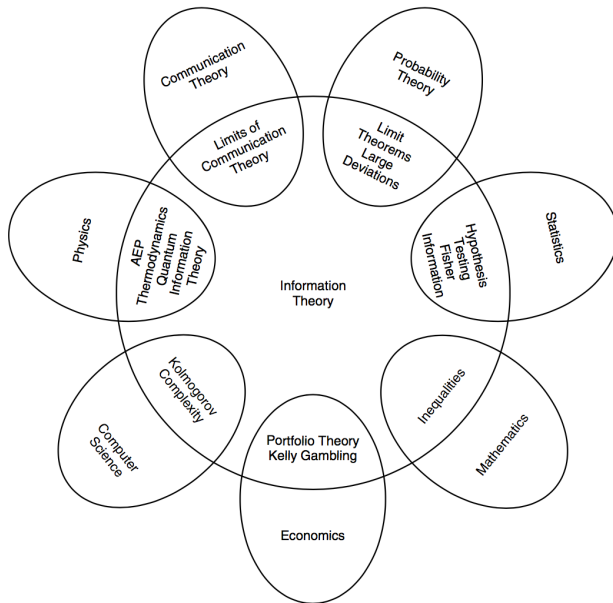
Theorem

If $X \sim \text{HGeom}(w, b, n)$ and $N = w + b \rightarrow \infty$ such that $p = w / (w + b)$ remains fixed, then the PMF of X converges to the $\text{Bin}(n, p)$ PMF.

Outline

- 1 Random Variables
- 2 Bernoulli and Binomial
- 3 Hypergeometric
- 4 Discrete Uniform & Zipf Distribution
- 5 Cumulative Distribution Functions
- 6 Functions of Random Variables: Random Variables
- 7 Independence of R.V.s
- 8 Binomial & Hypergeometric
- 9 Information Theory & Entropy**

Information Theory & Other Fields



Entropy

$$\underline{X \sim p(x) = \begin{cases} 1 & x = x_0 \\ 0 & \text{otherwise.} \end{cases}}$$

$$X = \{x_0\}$$

$$H(X) = p(x_0) \log_2 \frac{1}{p(x_0)}$$

$$X = x_0$$

as.

$$= 1 \cdot \log_2 \frac{1}{1} = \log_2 1 = 0$$

Definition

Given a random variable X with a probability mass function $p(x)$ and a support \mathcal{X} . The entropy of X is defined by

bit

$$\underline{H(X) = - \sum_{x \in \mathcal{X}} p(x) \log_2 p(x)}$$

$$= \sum_{x \in \mathcal{X}} p(x) \log_2 \frac{1}{p(x)}$$

Entropy of Discrete Uniform Distribution

- X has a uniform distribution over k outcomes.
- $p(x) = 1/k$
- Then the entropy of X is

$$H(X) = - \sum_{x=1}^k p(x) \log_2 p(x) = - \sum_{x=1}^k \frac{1}{k} \log_2 \frac{1}{k} = \log_2 k$$

$$\log_2 2 = 1 \text{ (bit)}$$

X discrete r.v.

$$H(X) \leq \log_2 k$$

Balance Puzzle

①

one weight

\log_2^3 (bits)

$\left. \begin{matrix} = \\ \leq \\ > \end{matrix} \right\}$

$$k \leq \frac{3^n - 1}{2}$$

$$k=13$$

$$\Rightarrow 3^n \geq 27$$

$$\Rightarrow n \geq 3$$

② k

n weights

$n \log_2^3$ (bits)

You have 13 apparently identical gold coins. One of them is false but is virtually indistinguishable from the others. You also have a balance with two pans, but without weights. Accordingly, any measurement will tell you if the loaded pans weight the same or, if not, which weighs more. How many measurements are needed to find the false coin?

③

k coins - one of them is false.

$$\log_2 k + \log_2^2 = \log_2^{2k} \text{ (bits)}$$

④

$$n \log_2^3 \geq \log_2^{2k} \Rightarrow \underline{3^n} \geq \underline{2k} \Rightarrow \underline{3^n} \geq \underline{2k+1}$$

Solution

Entropy Bounds in General

$$n \log_2 3 \geq \log_2 k \Rightarrow 3^n \geq k$$

$$n \geq \log_3 k$$

| Known | Goal | Maximum Coins for n weighings | Number of Weighings for c coins |
|-----------------------------------------------------------------|----------------------------------------------------------------------|---------------------------------|-----------------------------------|
| Whether target coin is lighter or heavier than others | Identify coin | 3^n | $\lceil \log_3(c) \rceil$ |
| Target coin is different from others | Identify coin | $\frac{3^n - 1}{2}$ [1] | $\lceil \log_3(2c + 1) \rceil$ |
| Target coin is different from others, or all coins are the same | Identify if unique coin exists, and whether it is lighter or heavier | $\frac{3^n - 1}{2} - 1$ | $\lceil \log_3(2c + 3) \rceil$ |

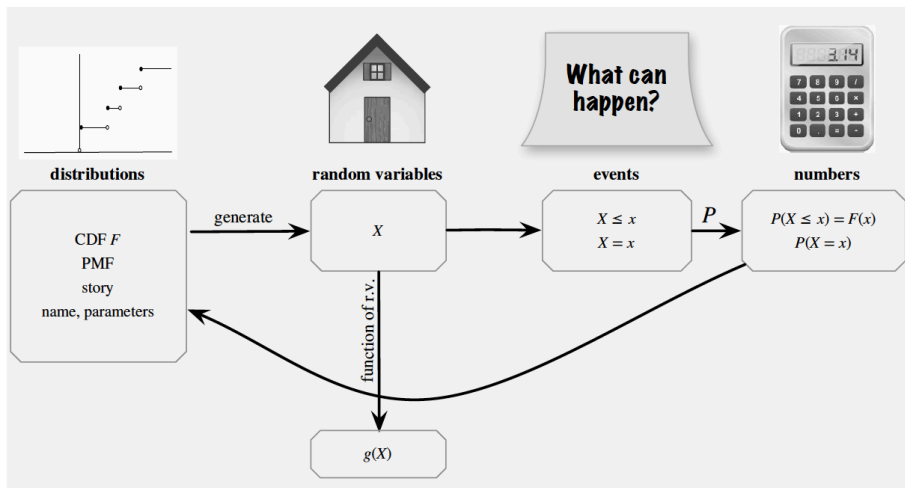
$$k < \frac{3^n - 3}{2}$$

$$n \log_2 3 \geq \log_2(2k+1) + \log_2 2 = \log_2(2k+2)$$

$$\Rightarrow 3^n \geq 2k+2$$

$$\Rightarrow 3^n \geq 2k+3$$

Summary 1



References

- Chapter 3 of **BH**
- Chapter 2 of **BT**