

Answer (Ex. 5.2) —

We want $1 - \alpha = 0.95$, and from the standard Normal Table we know that the corresponding $z_{\alpha/2} = 1.96$. Then we can get the right sample size n from the CLT implied Equation (61) in the lecture notes, which is,

$$n = \left(\frac{\sqrt{V(X_1) z_{\alpha/2}}}{\epsilon} \right)^2$$

$$= \left(\frac{(1/4) \times 1.96}{0.98} \right)^2 = \left(\frac{0.49}{0.98} \right)^2 = 0.04$$

Finally, by rounding up to the next integer we need $n = 97$ measurements to meet the specifications of your boss (at least up to the approximation provided by the CLT).

Answer (Ex. 5.3) —

This work is licensed under the Creative Commons Attribution-Noncommercial-Share Alike 4.0 International License. To view a copy of this license, visit <https://creativecommons.org/licenses/by-nc-sa/4.0/>.
By CLT, $\frac{\sqrt{V(X_n) - E(X_1)}}{\sqrt{n}} \rightsquigarrow Z \sim \text{Normal}(0, 1)$. So we need to apply the "standardization" to both sides of the inequality that is defining the event of interest:

$$\{X_n < 5.5\},$$

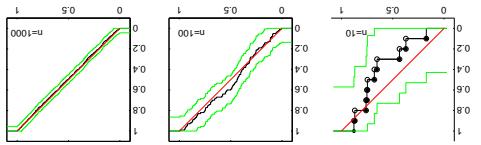
in order to find its probability $P(X_n < 5.5)$.

(source: Wasserman, *All of Statistics*, Springer, p. 78, 2003)

$$\begin{aligned} P(X_n > 5.5) &= P\left(\frac{\sqrt{V(X_n) - E(X_1)}}{\sqrt{n}} > \frac{\sqrt{V(X_1)} - \frac{1}{\sqrt{n}}}{\sqrt{n}}\right) \\ &\approx P\left(Z > \frac{\sqrt{V(X_1)} - \frac{1}{\sqrt{n}}}{\sqrt{n}}\right) \quad [\text{since we know/assume that } E(X_1) = V(X_1) = \bar{X}] \\ &= P\left(Z > \frac{\sqrt{V(X_1)} - \bar{X}}{\sqrt{n}}\right) \quad [\text{Since, } \bar{X} = 5 \text{ and } n = 125 \text{ in this Example}] \\ &= P(Z > 2.5) = \Phi(2.5) = 0.9938. \end{aligned}$$

Answer (Ex. 5.4) — HINT: Use the LLN after finding the population mean and variance of X_i .

Answer (Ex. 5.5) — HINT: Use the CLT after finding the population mean and variance of X_i .



Answer (Ex. 5.6) —

This work is licensed under the Creative Commons Attribution-Noncommercial-Share Alike 4.0 International License. To view a copy of this license, visit <https://creativecommons.org/licenses/by-nc-sa/4.0/>.
By CLT, $\frac{\sqrt{V(X_n) - E(X_1)}}{\sqrt{n}} \rightsquigarrow Z \sim \text{Normal}(0, 1)$. So we need to apply the "standardization" to both sides of the inequality that is defining the event of interest:

Version Date: October 22, 2019

Variation

©2008-2013 Razeeh Samiuddin. ©2008-2013 Dominic Lee.

School of Mathematics and Statistics, University of Canterbury, Christchurch, New Zealand
*Department of Mathematics, Uppsala University, Uppsala, Sweden
Laboratory for Mathematical Statistical Experiments, Uppsala Centre, and Razeeh Samiuddin and Dominic Lee*,
Razeeh Samiuddin* and Dominic Lee*,
Probability Theory I

Version Date: October 22, 2019

Variation

©2007-2019 Razeeh Samiuddin. ©2008-2013 Dominic Lee.

Answer (Ex. 3.53) —

1.

$$\varphi_X(t) = E(e^{itX}) = \sum_{x=0}^{\infty} e^{itx} \frac{\lambda^x}{x!} e^{-\lambda} = e^{-\lambda} \sum_{x=0}^{\infty} \frac{e^{itx} \lambda^x}{x!} = e^{-\lambda} \sum_{x=0}^{\infty} \frac{(\lambda e^{it})^x}{x!} = e^{-\lambda} e^{\lambda e^{it}} = e^{\lambda e^{it} - \lambda} .$$

The second-last equality above is using $e^\alpha = \sum_{x=0}^{\infty} \frac{\alpha^x}{x!}$ with $\alpha = \lambda e^{it}$.

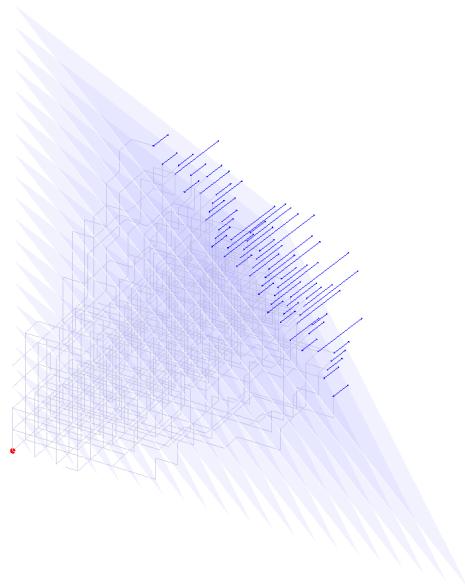
2. To find $V(X)$ using $\varphi_X(t)$ we need $E(X)$ and $E(X^2)$.

$$\begin{aligned} E(X) &= \frac{1}{i} \left[\frac{d}{dt} \varphi_X(t) \right]_{t=0} = \frac{1}{i} \left[\frac{d}{dt} e^{\lambda e^{it} - \lambda} \right]_{t=0} = \frac{1}{i} \left[e^{\lambda e^{it} - \lambda} \frac{d}{dt} (\lambda e^{it} - \lambda) \right]_{t=0} \\ &= \frac{1}{i} \left[e^{\lambda e^{it} - \lambda} \lambda i e^{it} \right]_{t=0} = \frac{1}{i} (e^{\lambda - \lambda} \lambda i) = \frac{1}{i} (\lambda i) = \lambda . \end{aligned}$$

$$\begin{aligned} E(X^2) &= \frac{1}{i^2} \left[\frac{d^2}{dt^2} \varphi_X(t) \right]_{t=0} = \frac{1}{i^2} \left[\frac{d}{dt} \left(e^{\lambda e^{it} - \lambda} \lambda i e^{it} \right) \right]_{t=0} = \frac{1}{i^2} \left[\frac{d}{dt} \left(\lambda i e^{\lambda e^{it} - \lambda + it} \right) \right]_{t=0} \\ &= \frac{1}{i^2} \left[\lambda i e^{\lambda e^{it} - \lambda + it} \frac{d}{dt} (\lambda e^{it} - \lambda + it) \right]_{t=0} = \frac{1}{i^2} \left[\lambda i e^{\lambda e^{it} - \lambda + it} (\lambda i e^{it} - 0 + i) \right]_{t=0} \\ &= \frac{1}{i^2} \left(\lambda i e^{\lambda e^0 - \lambda + 0} (\lambda i e^0 + i) \right) = \frac{1}{i^2} (\lambda i (\lambda i + i)) = \frac{1}{i^2} (\lambda i^2 (\lambda + 1)) = \lambda^2 + \lambda . \end{aligned}$$

Finally,

$$V(X) = E(X^2) - (E(X))^2 = \lambda^2 + \lambda - \lambda^2 = \lambda .$$



Answer (Ex. 3.54) —

$$\varphi_W(t) = \varphi_{X+Y}(t) = \varphi_X(t) \times \varphi_Y(t) = e^{\lambda e^{it} - \lambda} \times e^{\mu e^{it} - \mu} = e^{\lambda e^{it} - \lambda + \mu e^{it} - \mu} = e^{(\lambda + \mu) e^{it} - (\lambda + \mu)} .$$

So, W is a Poisson($\lambda + \mu$) RV. Thus the sum of two independent Poisson RVs is also a Poisson RV with parameter given by the sum of the parameters of the two RVs being added. The same idea generalizes to the sum of more than two Poisson RVs.

Answer (Ex. 3.55) —

We can use the facts by noting $Y = -Z = 0 + (-1) \times Z$, with $a = 0$ and $b = -1$ in $Y = a + bX$ and get

$$\varphi_Y(t) = e^{i \times 0 \times t} \varphi_Z(-1 \times t) = \varphi_Z(-t) = e^{-((-1 \times t))^2 / 2} = e^{-t^2 / 2} = \varphi_Z(t)$$

Thus, $\varphi_{-Z}(t) = \varphi_Z(t)$ and therefore the distributions of Z and $-Z$ are the same. This should make sense because by switching signs of a symmetric (about 0) RV you have not changed its distribution! Note: we are not saying $Z = -Z$ but just that their distributions are the same, i.e., $F_Z(z) = F_{-Z}(z)$ for every $z \in \mathbb{R}$.

Answer (Ex. 3.56) — Apply the formulas for the sample mean, \bar{X}_7 , and sample variance.

Answer (Ex. 4.1) — This is nothing but the inversion sampler for the standard Cauchy RV X .

COMBINATIONS Combimbitability concept. Calculation of probabilities. Probability distributions. Independence and conditioned distributions. Expectation and variance. Continuous distributions. Independent and conditioned distributions. Law of large numbers. Practical examples of design of probability models.

ASSESSMENT Written examination at the end of the course combined with written assignments during the course according to instructions delivered at course start.

CONTENT

- account for the axiomatic basis of the probability theory;
 - carry out probability calculations by means of combinatorial principles and be able to use methods for independent events;
 - account for the concepts of stochastic variable and expectation and be able to calculate probabilities, expectations and variance for given distributions;
 - account for the most common probability distributions and how to do simulations with them;
 - handle conditioned probabilities, distributions and expectations as well as moment generating functions;
 - apply the law of large numbers and the central limit theorem;
 - account for probability models within different application fields.

On completion of the course, the student should be able to

LEARNING OUTCOMES

SEE <https://www.uu.se/en/administrations/master/se/lma/kursplan/?kp1d=3891&typ=1>.

Official Course Syllabus

Course Code: MNSO34	Semester: Autumn 2019	Report Code: 10504	33%, DAG, NNL	$\overleftarrow{\text{work}}$	12 = $\frac{33}{100} \cdot 37.5$ hours / week	2019-09-02 - 2019-11-03	Uppsala University	Probability Theory I, 5.0 e	Course Syllabus: Friday, 4th of October 2019	Course Coordinator: Razae Shaimuddin
---------------------	-----------------------	--------------------	---------------	-------------------------------	---	-------------------------	--------------------	-----------------------------	--	--------------------------------------

Course Syllabus and Overview

505

Time Table & Course Overview – In Progress

(KEY,VALUE): (EX, Exercise), (RD, Read), (RW, Review), (UD, Understand), (PM, Program)

Table 1: Time Table for Virtual Student of Probability Theory I

Lec.	Lab.	Week	Topics	Comprehension \mapsto Action \times Content
01	*	36	Preliminaries: Set Theory, Numbers, Functions ,...	RW Sec. 1.1,1.3,1.4 EX 1.2; RW Table. 1.1 RD 1.6; UD 6,7; EX 1.5; RD 1.9 PM 5,9,10,11,21
			Elementary Combinatorics & Number Theory [optional] Introduction to MATLAB	RD 2.1,2.2; EXs 2.3
02			Probability	RD 2.2,2.4; UD 31
03			Probability and Conditional Probability	
04	37		Conditional Probability & Bayes Theorem	RD 2.4.1; UD 32,33,34; RD 2.4.2
05			Conditional Probability & Independence	RD 2.4.2; UD 35,36,37; EXs 2.5
06			Random variables	RD 3.1; EX 3.1; UD 1,4ii, EX 3.2
			Discrete Random variables and IID Bernoulli Trials	RD 3.1, 3.2.1, 3.2.3
			Common Discrete Random variables	RD 3.2.3; UD 45; UD 4 ;EX 3.3;
07	38		Common Discrete Random variables	UD 46, 5, 47, 48, 6, 49, 50, 51
08			Common Discrete Random variables	EX 3.4 & EXs 3.3;
			Continuous RVs	RD 3.4; UD 52, 53, 54, RD 3.4.1
			Common Continuous RVs	RD 3.4.2; UD 55, 56; EX 3.17, 3.18
09			Common Continuous RVs	EX 57, 58; UD 59; EXs 3.5
			Transformations of RVs – Discrete	RD 3.6; UD 60, 61; RW 3.6.1
10	39		Transformations of RVs – Discrete	RD 3.6.2, UD 62, 63, 64
			Transformations of RVs – Continuous (1-to-1 & monotone)	RD 3.6.3; UD 65, 66, 67, 68
11			Transformations of RVs – Continuous (Direct method)	RD 3.6.3; UD 69, 70; EXs 3.7
			Expectations	RD 3.8, 3.8.1; UD 71, 72, 74, 75, 76; EX 3.29; UD 77, 78, 79, 80; EXs 3.9
12	40		Multivariate Random Variables	RD 3.10; UD 83, 84, 85, 86, UD 87, 88, 89, 90, 91, 92
13			Common \mathbb{R}^m -valued RVs	RD 3.10.2, 3.10.3; UD 95, 96
14			Characteristic Functions	RD 3.10.4; UD 97, 98; EX 3.36, 3.37; EXs 3.11
15	41		Statistics & Random Number Generation	RD 3.12; UD 101, 102, 103, 104, 105; EXs 3.13; [non-examinable] UD 169, 170
16			Statistics	RD 3.14, 4.2, 4.3
17			Simulation – Inversion & Rejection Samplers	EXs 3.15
18	42		Convergence of RVs & Limit Laws	EXs 4.4
19			Basic Inequalities & Law of Large Numbers	RD 5.1, UD 157, 158, 159, 160
20			Law of Large Numbers & Central Limit Theorem via CFs	UD 161, 163, 164; EX 5.1
			Model exam with live-scribed solutions	UD 165, 166, 167, 168; EXs 5.4

Table 2: Time Table for Inference Theory I

Lec.	Lab.	Week	Topics	Section/Labworks/Simulations
------	------	------	--------	------------------------------

Answer (Ex. 3.48) —

Let X_1, X_2, \dots, X_{10} denote the fill volumes of 10 cans. The average fill volume is the sample mean

$$\bar{X}_{10} = \frac{1}{10} \sum_{i=1}^{10} X_i$$

By property of Expectations and Variances for linear combinations

$$E(\bar{X}_{10}) = E\left(\frac{1}{10} \sum_{i=1}^{10} X_i\right) = \frac{1}{10} \sum_{i=1}^n E(X_i) = \frac{1}{10} \sum_{i=1}^n E(X_1) = \frac{1}{10} \times 10 \times E(X_1) = E(X_1) = 12.1$$

Or by directly using the “formula” $E(\bar{X}_{10}) = E(X_1) = 12.1$ for these 10 identically distributed RVs. Similarly,

$$V(\bar{X}_{10}) = V\left(\frac{1}{10} \sum_{i=1}^{10} X_i\right) = 10 \times \frac{1}{10^2} V(X_1) = \frac{1}{10} \times 0.01 = 0.001$$

Or by directly using the “formula” $V(\bar{X}_{10}) = V(X_1)/10$ for these 10 independently and identically distributed RVs.

By the special property of Normal RVs – a linear combination of independent normal RVs is also normal – we know that \bar{X}_{10} is a Normal(12.1, 0.001) RV. Consequently, the probability of interest is

$$\begin{aligned} P(\bar{X}_{10} < 12.01) &= P\left(\frac{\bar{X}_{10} - E(\bar{X}_{10})}{\sqrt{0.001}} < \frac{12.01 - E(\bar{X}_{10})}{\sqrt{0.001}}\right) = P\left(Z < \frac{12.01 - 12.1}{0.0316}\right) \\ &\cong P(Z < -2.85) = 1 - P(Z < 2.85) = 1 - \Phi(2.85) = 1 - 0.9978 = 0.0022 \end{aligned}$$

Answer (Ex. 3.49) —

Using the Multinomial($n = 10, \theta_1 = 0.6, \theta_2 = 0.3, \theta_3 = 0.08, \theta_4 = 0.02$) RV as our model

$$\begin{aligned} P((X_1, X_2, X_3, X_4) = (6, 2, 2, 0); n = 10, \theta_1 = 0.6, \theta_2 = 0.3, \theta_3 = 0.08, \theta_4 = 0.02) \\ = \frac{10!}{6! \times 2! \times 2! \times 0!} \times 0.6^6 \times 0.3^2 \times 0.08^2 \times 0.02^0 \\ = \frac{10 \times 9 \times 8 \times 7 \times 6 \times 5 \times 4 \times 3 \times 2 \times 1}{(6 \times 5 \times 4 \times 3 \times 2 \times 1) \times (2 \times 1) \times (2 \times 1)} \times 0.6^6 \times 0.3^2 \times 0.08^2 \times 1 \approx 0.03386 \end{aligned}$$

Answer (Ex. 3.50) —

1. The CF of the discrete RV X is

$$\begin{aligned} \varphi_X(t) &= E(e^{itX}) = \sum_{x \in \{0,1,2\}} e^{itx} f_X(x) = e^{it \times 0} \times \frac{1}{3} + e^{it \times 1} \times \frac{1}{3} + e^{it \times 2} \times \frac{1}{3} \\ &= \frac{1}{3} (1 + e^{it} + e^{2it}). \end{aligned}$$

Contents

3 Random Variables	59
3.1 Basic Definitions	61
3.2 Discrete Random Variables	64
3.2.1 An Elementary Family of Bernoulli Random Variables	68
3.2.2 Independent Bernoulli Trials	69
3.2.3 Some Common Discrete Random Variables	70
3.3 Exercises in Discrete Random Variables	79
3.4 Continuous Random Variables	81
3.4.1 An Elementary Continuous Random Variable	84
3.4.2 Some Common Continuous Random Variables	85
3.5 Exercises in Continuous Random Variables	91
3.6 Transformations of random variables	91
3.6.1 A Review of Inverse Images	92
3.6.2 Transformations of discrete random variables	94
3.6.3 Transformations of continuous random variables	95
3.7 Exercises in Transformations of Random Variables	102
3.8 Expectations	102
3.8.1 Expectations of functions of random variables	103
3.8.2 Properties of expectations	106
3.8.3 Expectation of Common Random Variables	107
3.9 Exercises in Expectations of Random Variables	112
3.10 Multivariate Random Variables	113
3.10.1 \mathbb{R}^2 -valued Random Variables	114
3.10.2 Conditional Random Variables	124
3.10.3 \mathbb{R}^m -valued Random Variables	126
3.10.4 Some Common \mathbb{R}^m -valued RVs	130
3.10.5 Dependent Random Variables	135
3.11 Exercises in Multivariate Random Variables	136
3.12 Characteristic Functions	139
3.12.1 Obtaining Moments from Characteristic Function	139
3.12.2 Moment Generating Function	144
3.13 Exercises in Characteristic Functions	144
3.14 Statistics	145
3.14.1 Data and Statistics	145
3.14.2 Univariate Data	152
3.14.3 Bivariate Data	154
3.14.4 Trivariate Data	155
3.14.5 Multivariate Data	156
3.14.6 Loading and Exploring Real-world Data	157
3.14.7 Geological Data	157
3.14.8 Meteorological Data	160
3.14.9 Textual Data	164
3.14.10 Machine Sensor Data	165
3.15 Exercises in Statistics	166

$$\begin{aligned}
 E((X, Y)) &= \sum_{(x,y) \in S_{X,Y}} (x, y) \times f_{X,Y}(x, y) \\
 &= (0, 0) \times 0.2 + (1, 1) \times 0.1 + (1, 2) \times 0.1 + (2, 1) \times 0.1 + (2, 2) \times 0.1 + (3, 3) \times 0.4 \\
 &= (0, 0) + (0.1, 0.1) + (0.1, 0.2) + (0.2, 0.1) + (0.2, 0.2) + (1.2, 1.2) \\
 &= (1.8, 1.8)
 \end{aligned}$$

Since addition is component-wise $E((X, Y)) = (E(X), E(Y))$ and therefore $E(X) = E(Y) = 1.8$. Alternatively, you can first find the marginal PMFs f_X and f_Y for X and Y and then take the expectations $E(X) = \sum_x x \times f_X(x)$ and $E(Y) = \sum_y y \times f_Y(y)$.

Finally,

$$\text{Cov}(X, Y) = E(XY) - E(X)E(Y) = 4.5 - 1.8^2 = 1.26.$$

Answer (Ex. 3.44) —

Since X and Y are independent, $F_{X,Y}(x, y) = F_X(x)F_Y(y)$ for all $(x, y) \in \mathbb{R}^2$, and we get:

$$F_{X,Y}(x, y) = \begin{cases} 0 & \text{if } x < 1 \text{ or } y < 0 \\ \frac{1}{2} \left(1 - \frac{1}{2}e^{-y} - \frac{1}{2}e^{-2y}\right) & \text{if } 1 \leq x < 2 \text{ and } y \geq 0 \\ 1 - \frac{1}{2}e^{-y} - \frac{1}{2}e^{-2y} & \text{if } x \geq 2 \text{ and } y \geq 0 \end{cases}$$

You can arrive at the answer by partitioning x -axis into $(-\infty, 1)$, $[1, 2)$ and $[2, \infty)$ where $F_X(x)$ takes distinct values. Similarly, partition the y -axis into $(-\infty, 0)$ and $[0, \infty)$ where $F_Y(y)$ takes distinct values. Now (x, y) can take values in one of these $3 \times 2 = 6$ partitions of the $x \times y$ plane as follows (make a picture!):

$$(-\infty, 1) \times (-\infty, 0), [1, 2) \times (-\infty, 0), [2, \infty) \times (-\infty, 0), (-\infty, 1) \times [0, \infty), [1, 2) \times [0, \infty), [2, \infty) \times [0, \infty).$$

Now work out what $F_{X,Y}(x, y) = F_X(x)F_Y(y)$ is for (x, y) in each of the above six partitions of the plane and you will get the expression for $F_{X,Y}(x, y)$ given above.

Answer (Ex. 3.45) —

First obtain marginal PDF of Y . If $y \in [2, 3]$ then

$$f_Y(y) = \int_{-\infty}^{\infty} f_{X,Y}(x, y) dx = \int_0^{\infty} e^{-x} dx = [-e^{-x}]_0^{\infty} = 0 - (-1) = 1.$$

Therefore,

$$f_Y(y) = \begin{cases} 1 & \text{if } y \in [2, 3] \\ 0 & \text{otherwise.} \end{cases}$$

Now, obtain the marginal PDF of X . If $x \in [0, \infty)$ then

$$f_X(x) = \int_{-\infty}^{\infty} f_{X,Y}(x, y) dy = \int_2^3 e^{-x} dy = e^{-x} \int_2^3 1 dy = e^{-x} [y]_2^3 = e^{-x} (3 - 2) = e^{-x}.$$

Therefore,

$$f_X(x) = \begin{cases} e^{-x} & \text{if } x \in [0, \infty) \\ 0 & \text{otherwise.} \end{cases}$$

Finally, verifying that $f_{X,Y}(x, y) = f_X(x)f_Y(y)$ for any $(x, y) \in \mathbb{R}^2$ is done case by case. Draw a picture on the plane to work out the cases from the distinct expressions taken by $f_{X,Y}(x, y)$. There are only two cases to consider (when $f_{X,Y}(x, y)$ takes zero values and when $f_{X,Y}(x, y)$ takes non-zero values):

$$E(XY) = \sum_{x,y} x \times y \times f_{XY}(x,y)$$

Find $E(XY)$, $E(X)$ and $E(Y)$ to get $\text{Cov}(X, Y) = E(XY) - E(X)E(Y)$ as follows:

Answer (Ex. 3.43) —

$$= 0 \times 0 \times 0.2 + 1 \times 1 \times 0.1 + 1 \times 2 \times 0.1 + 2 \times 1 \times 0.1 + 2 \times 2 \times 0.1 + 3 \times 3 \times 0.4 = 4.5$$

$$\text{Var}(X) = I_2 V(X_1) + I_2 V(X_2) + I_2 V(X_3) = 25 + 40 + 30 = 95 \text{nm}^2$$

By the property that $V(\sum_{i=1}^n a_i X_i) = \sum_{i=1}^n a_i^2 V(X_i)$, Variance of X is

$$X = X_1 + X_2 + X_3$$

respectively. Let X denote the thickness of the final product. Then let X_1, X_2, X_3 be independent RVs that denote the thicknesses of the first, second and third layer,

Answer (Ex. 3.42) —

which in turn is equal to the PMF $f_{XY}(x, y)$ in the question. Therefore we have shown that the component RVs X and Y in the RV (X, Y) are indeed independent.

$$f_{XY}(x, y) = \begin{cases} \frac{1}{2} \times \frac{1}{2} = \frac{1}{4} & \text{if } (x, y) = (1, 1) \\ \frac{1}{2} \times \frac{1}{2} = \frac{1}{4} & \text{if } (x, y) = (0, 1) \\ \frac{1}{2} \times \frac{1}{2} = \frac{1}{4} & \text{if } (x, y) = (0, 0) \\ 0 & \text{otherwise} \end{cases}$$

Finally, the product of $f_X(x)$ and $f_Y(y)$ is

$$f_Y(y) = \begin{cases} \sum_{x \in S^{X,Y}} f_{XY}(x, y) = f_{XY}(0, 1) + f_{XY}(1, 1) = \frac{1}{4} + \frac{1}{4} = \frac{1}{2} & \text{if } y = 1 \\ 0 & \text{otherwise} \end{cases}$$

Similarly,

$$f_X(x) = \begin{cases} \frac{1}{2} & \text{if } x = 1 \\ 0 & \text{otherwise} \end{cases}$$

Thus,

$$f_X(1) = \sum_{y \in S^{X,Y}} f_{XY}(1, y) = f_{XY}(1, 0) + f_{XY}(1, 1) = \frac{1}{1} + \frac{1}{4} = \frac{5}{4}$$

and

First derive the marginal PMF of X and Y and then check if the PMF is the product of the marginal PMFs.

$f_X(0) = \sum_{y \in S^{X,Y}} f_{XY}(0, y) = f_{XY}(0, 0) + f_{XY}(0, 1) = \frac{1}{1} + \frac{1}{4} = \frac{5}{4}$

CHAPTER 6. FINITE MARKOV CHAINS

Answer (Ex. 3.41) —

4 Simulation 4.1 Physical Random Number Generators 167 4.2 Pseudo-Random Number Generators 167 4.2.1 Linear Congruential Generators 168 4.2.2 Generalized Feedback Shift Register and the "Mersenne Twister" PRNG 171 4.3 Simulation of non-Uniform(0, 1) Random Variables 173 4.3.1 Inversion Sampler for Continuous Random Variables 173 4.3.2 Inversion Sampler for Discrete Random Variables 181 4.3.3 von Neumann Rejection Sampler (RS) 190 4.4 Exercises in Simulation 195 5.1 Convergence of Random Variables 196 5.2 Law of Large Numbers 202 5.2.1 Application: Point Estimation of $E(X_1)$ 205 5.3 Central Limit Theorem 207 5.3.1 Application: Tolerating Errors in our estimate of $E(X_1)$ 208 5.3.2 Application: Set Estimation of $E(X_1)$ 210 5.4 Exercises in Limit Laws of Statistics 211 6.1 Stochastic Processes 212 6.2 Introduction 213 6.3 irreducibility and Aperiodicity 221 6.4 Stationarity 223 6.5 reversibility 225 6.6 Standard normal distribution function table 229	Answers to Selected Exercises 235 230 229 228 227 226 225 224 223 222 221 220 219 218 217 216 215 214 213 212 211 210 209 208 207 206 205 204 203 202 201 200 199 198 197 196 195 194 193 192 191 190 189 188 187 186 185 184 183 182 181 180 179 178 177 176 175 174 173 172 171 170 169 168 167 166 165 164 163 162 161 160 159 158 157 156 155 154 153 152 151 150 149 148 147 146 145 144 143 142 141 140 139 138 137 136 135 134 133 132 131 130 129 128 127 126 125 124 123 122 121 120 119 118 117 116 115 114 113 112 111 110 109 108 107 106 105 104 103 102 101 100 99 98 97 96 95 94 93 92 91 90 89 88 87 86 85 84 83 82 81 80 79 78 77 76 75 74 73 72 71 70 69 68 67 66 65 64 63 62 61 60 59 58 57 56 55 54 53 52 51 50 49 48 47 46 45 44 43 42 41 40 39 38 37 36 35 34 33 32 31 30 29 28 27 26 25 24 23 22 21 20 19 18 17 16 15 14 13 12 11 10 9 8 7 6 5 4 3 2 1 0
---	--

List of Tables

1	Time Table for Virtual Student of Probability Theory I	4
2	Time Table for Inference Theory I	4
1.1	Symbol Table: Sets and Numbers	22
3.1	The 8 ω 's in the sample space Ω of the experiment \mathcal{E}_θ^3 are given in the first row above. The RV Y is the number of 'Heads' in the 3 tosses and the RV Z is the number of 'Tails' in the 3 tosses. Finally, the RVs Y' and Z' are the indicator functions of the event that 'all three tosses were Heads' and the event that 'all three tosses were Tails', respectively.	125
6.1	Symbol Table: Probability and Statistics	232
6.2	Random Variables with PDF and PMF (using indicator function), Mean and Variance	233
6.3	Symbol Table: Sets and Numbers	233
6.4	Symbol Table: Probability and Statistics	234

8.

$$\begin{aligned}
 E(Y) &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} y f_{X,Y}(x,y) dx dy \\
 &= \int_0^1 \int_0^1 y \frac{6}{5} (x^2 + y) dx dy = \frac{6}{5} \int_0^1 \int_0^1 (x^2 y + y^2) dx dy = \frac{6}{5} \int_0^1 \left[\frac{x^3 y}{3} + y^2 x \right]_{x=0}^1 dy \\
 &= \frac{6}{5} \int_0^1 \left(\frac{y}{3} + y^2 - 0 - 0 \right) dy = \frac{6}{5} \left[\frac{y^2}{6} + \frac{y^3}{3} \right]_{y=0}^1 = \frac{6}{5} \left(\frac{1}{6} + \frac{1}{3} + -0 - 0 \right) = \frac{6}{5} \times \frac{3}{6} = \frac{3}{5}
 \end{aligned}$$

9.

$$\begin{aligned}
 E(XY) &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} xy f_{X,Y}(x,y) dx dy \\
 &= \frac{6}{5} \int_0^1 \int_0^1 xy (x^2 + y) dx dy = \frac{6}{5} \int_0^1 \int_0^1 x^3 y + xy^2 dx dy = \frac{6}{5} \int_0^1 \left[\frac{x^4 y}{4} + \frac{x^2 y^2}{2} \right]_{x=0}^1 dy \\
 &= \frac{6}{5} \int_0^1 \left(\frac{y}{4} + \frac{y^2}{2} - 0 - 0 \right) dy = \frac{6}{5} \left[\frac{y^2}{8} + \frac{y^3}{6} \right]_{y=0}^1 = \frac{6}{5} \left(\frac{1}{8} + \frac{1}{6} - 0 - 0 \right) \\
 &= \frac{6}{5} \left(\frac{3}{24} + \frac{4}{24} \right) = \frac{6}{5} \times \frac{7}{24} = \frac{7}{20}
 \end{aligned}$$

10.

$$\text{Cov}(X, Y) = E(XY) - E(X)E(Y) = \frac{7}{20} - \left(\frac{3}{5} \times \frac{3}{5} \right) = \frac{7}{20} - \frac{9}{25} = \frac{35}{100} - \frac{36}{100} = -\frac{1}{100}$$

Answer (Ex. 3.40) — Note that $Y = (X - \mu)^2$ is not on-to-one so it is better to use the direct method by differentiating the distribution function of Y , $F_Y(y)$, to obtain $f_Y(y)$. If $y \geq 0$,

$$\begin{aligned}
 F_Y(y) &= \mathbf{P}(Y \leq y) \\
 &= \mathbf{P}((X - \mu)^2 \leq y) \\
 &= \mathbf{P}(-\sqrt{y} \leq X - \mu \leq \sqrt{y}) \\
 &= \mathbf{P}(\mu - \sqrt{y} \leq X \leq \mu + \sqrt{y}) \\
 &= F_X(\mu + \sqrt{y}) - F_X(\mu - \sqrt{y})
 \end{aligned}$$

Differentiating this expression gives

$$\begin{aligned}
 f_Y(y) &= \frac{d}{dy} (F_X(\mu + \sqrt{y}) - F_X(\mu - \sqrt{y})) \\
 &= \frac{1}{2} y^{-\frac{1}{2}} f_X(\mu + \sqrt{y}) - \left(-\frac{1}{2} y^{-\frac{1}{2}} \right) f_X(\mu - \sqrt{y}) \\
 &= \frac{1}{2\sqrt{y}} (f_X(\mu + \sqrt{y}) + f_X(\mu - \sqrt{y}))
 \end{aligned}$$

Note: If $y < 0$ then $f_Y(y) = 0$ since $F_Y(y) = 0$ in this case.

- Finally, the marginal PDF of the RV X in the second component of the RV (X, Y) is
- $$f_X(x) = \begin{cases} \frac{6}{9}(x^2 + \frac{1}{3}) & \text{if } 0 > x < 1 \\ 0 & \text{otherwise} \end{cases}$$
- Similarly, the marginal PDF $f_Y(y)$ for any $y \in (0, 1)$ by integrating over x is
- $$f_Y(y) = \int_{-\infty}^{\infty} f_{X,Y}(x, y) dx = \int_1^0 \frac{6}{9}(x^2 + \frac{1}{3})[y^3/3 + y]/(y+1/3).$$
4. The product of marginal PDFs of X and Y does not equal the joint PDF of (X, Y) for values of $(x, y) \in (0, 1)^2$
5. The joint distribution function $F_{X,Y}(x, y)$ for any $(x, y) \in (0, 1)^2$ is
- Therefore X and Y are not independent random variables (they are dependent).
- $$F_{X,Y}(x, y) = \int_y^x \int_a^{\infty} f_{X,Y}(u, v) du dv = \int_y^x \int_a^{\infty} \frac{6}{9}(u^2 + \frac{1}{3}u)(v^3/3 + v)/(v+1/3) du dv$$
- of the total probability theorem in Proposition 14 for the four event case.
- 2.3 Reference to the Venn diagram will help you understand this idea behind the proof of the indicator function of event $A \in \mathcal{F}$ is a RV \mathbb{I}_A with DF
- 3.1 The indicator function of event $A \in \mathcal{F}$ is a RV \mathbb{I}_A with DF
- 3.2 A RV X from a sample space Ω with 8 elements to \mathbb{R} and its DF
- 3.3 $f(x)$ and $F(x)$ of the fair coin toss random variable X , a discrete uniform RV on $\{0, 1\}$.
- 3.4 $f(x)$ and $F(x)$ of the fair die toss random variable X , a discrete uniform RV on $\{1, 2, 3, 4, 5, 6\}$.
- 3.5 $f(x)$ and $F(x)$ of sumised astrogal toss random variable X , a discrete (non-uniform) RV on $\{1, 2, 3, 4\}$.
- 3.6 PMF $f(x; \theta)$ and DF $F(x; \theta)$ with $\theta = 0.33$. You should see how PMF and DF change as θ goes from 0 to 1.
- 3.7 PMF of $X \sim \text{Geometric}(\theta = 0.5)$ and the relative frequency histogram based on 100 and 1000 samples from X according to Simulation 144 and Labwork 145 you will see in the sequel.
- 3.8 PDF of $X \sim \text{Binomial}(n = 10, \theta = 0.5)$ and the relative frequency histogram based on 100,000 samples from X obtained according to Simulation 148.
- 3.9 Figures from Sir Francis Galton, F.R.S., *Natural Inheritance*, Macmillan, 1889.
- 3.10 $f(x)$ and $F(x)$ of the Uniform($0, 1$) random variable X .

List of Figures

3.11 A convenient but mathematically imprecise Matlab plot from a polyline interpolation for the PDF and DF or CDF of the Uniform(0, 1) continuous RV X	85
3.12 $f(x; \lambda)$ and $F(x; \lambda)$ of an exponential random variable where $\lambda = 1$ (dotted) and $\lambda = 2$ (dashed).	86
3.13 Density and distribution functions of Exponential(λ) RVs, for $\lambda = 1, 10, 10^{-1}$, in four different axes scales.	87
3.14 $f(x)$ and $F(x)$ of the Uniform(θ_1, θ_2) random variable X	88
3.15 PDF and DF of a Normal(μ, σ^2) RV for different values of μ and σ^2	99
3.16 Mean ($\mathbf{E}_\theta(X)$), variance ($\mathbf{V}_\theta(X)$) and the rate of change of variance ($\frac{d}{d\theta}\mathbf{V}_\theta(X)$) of a Bernoulli(θ) RV X as a function of the parameter θ	107
3.17 Mean and variance of a Geometric(θ) RV X as a function of the parameter θ	109
3.18 PDF of $X \sim \text{Poisson}(\lambda = 10)$ and the relative frequency histogram based on 1000 samples from X according to Simulation 149.	110
3.19 Diagrams done on the board!	124
3.20 Visual Cognitive Tool GUI: Quincunx & Septcunx.	132
3.21 Quincunx on the Cartesian plane. Simulations of Binomial($n = 10, \theta = 0.5$) RV as the x-coordinate of the ordered pair resulting from the culmination of sample trajectories formed by the accumulating sum of $n = 10$ IID Bernoulli($\theta = 0.5$) random vectors over $\{(1, 0), (0, 1)\}$ with probabilities $\{\theta, 1 - \theta\}$, respectively. The blue lines and black asterisks perpendicular to and above the diagonal line, i.e. the line connecting $(0, 10)$ and $(10, 0)$, are the density histogram of the samples and the PMF of our Binomial($n = 10, \theta = 0.5$) RV, respectively.	133
3.22 JPDF, Marginal PDFs and Frequency Histogram of Bivariate Standard Normal R \vec{V} . .	135
3.23 JPDF, Marginal PDFs and Frequency Histogram of a Bivariate Normal R \vec{V} for lengths of girths of cylindrical shafts in a manufacturing process (in cm).	136
3.24 Sample Space, Random Variable, Realisation, Data, and Data Space.	146
3.25 Data Spaces $\mathbb{X} = \{0, 1\}^2$ and $\mathbb{X} = \{0, 1\}^3$ for two and three Bernoulli trials, respectively. .	146
3.26 Plot of the DF of Uniform(0, 1), five IID samples from it, and the ECDF \hat{F}_5 for these five data points $x = (x_1, x_2, x_3, x_4, x_5) = (0.5164, 0.5707, 0.0285, 0.1715, 0.6853)$ that jumps by $1/5 = 0.20$ at each of the five samples.	151
3.27 Frequency, Relative Frequency and Density Histograms	153
3.28 Frequency, Relative Frequency and Density Histograms	154
3.29 2D Scatter Plot	155
3.30 3D Scatter Plot	156
3.31 Plot Matrix of uniformly generated data in $[0, 1]^5$	156
3.32 Matrix of Scatter Plots of the latitude, longitude, magnitude and depth of the 22-02-2011 earth quakes in Christchurch, New Zealand.	159
3.33 Google Earth Visualisation of the earth quakes	161
3.34 Daily rainfalls in Christchurch since March 27 2010	162
3.35 Daily temperatures in Christchurch for one year since March 27 2010	163
3.36 Wordle of JOE 2010	164
3.37 Double Pendulum	165

Answer (Ex. 3.38) — The probability that *one* light bulb doesn't need to be replaced in 1200 hours is:

$$\begin{aligned}\mathbf{P}(X > 1.2) &= 1 - \mathbf{P}(X < 1.2) \\ &= 1 - \int_1^{1.2} 6(0.25 - (x - 1.5)^2) dx \\ &= 1 - \int_1^{1.2} 6(0.25 - x^2 + 3x - 2.25) dx \\ &= 1 - \int_1^{1.2} (-6x^2 + 18x - 12) dx \\ &= 1 - [-2x^3 + 9x^2 - 12x]_1^{1.2} \\ &= 1 - 0.1040 \\ &= 0.8960\end{aligned}$$

Assuming that the three light bulbs function independently of each other, the probability that none of them need to be replaced in the first 1200 hours is

$$\mathbf{P}(\{X_1 > 1.2\} \cap \{X_2 > 1.2\} \cap \{X_3 > 1.2\}) = 0.8960^3 = 0.7193$$

where X_i is the length of time that bulb i lasts.

Answer (Ex. 3.39) —

1. To find a we simply set $1 = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f_{X,Y}(x, y) dx dy$ and solve for a as follows:

$$\begin{aligned}1 &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f_{X,Y}(x, y) dx dy = \int_0^1 \int_0^1 a(x^2 + y) dx dy \\ &= a \int_0^1 \int_0^1 (x^2 + y) dx dy = a \int_0^1 \left[\frac{1}{3}x^3 + xy \right]_{x=0}^1 dy \\ &= a \int_0^1 \left(\frac{1}{3} + y - 0 \right) dy = a \left[\frac{y}{3} + \frac{1}{2}y^2 \right]_{y=0}^1 \\ &= a \left(0 - \left(\frac{1}{3} + \frac{1}{2}1^2 \right) \right) = a \left(\frac{1}{3} + \frac{1}{2} \right) \\ &= a \left(\frac{5}{6} \right)\end{aligned}$$

Therefore $a = 6/5$ and the joint PDF is

$$f_{X,Y}(x, y) = \begin{cases} \frac{6}{5}(x^2 + y) & \text{if } 0 < x < 1 \text{ and } 0 < y < 1 \\ 0 & \text{otherwise} \end{cases}$$

2. First compute the marginal PDF $f_X(x)$ for any $x \in (0, 1)$ by integrating over y

$$\begin{aligned}f_X(x) &= \int_{-\infty}^{\infty} f_{X,Y}(x, y) dy = \int_0^1 \frac{6}{5}(x^2 + y) dy = \left[\frac{6}{5}(yx^2 + y^2/2) \right]_{y=0}^1 \\ &= \frac{6}{5}((1 \times x^2 + 1^2/2) - 0) = \frac{6}{5} \left(x^2 + \frac{1}{2} \right)\end{aligned}$$

- 6.2 Transition Diagram of Flippant Freddy's Jumps 214
 6.1 Transition Diagram of Flippant Freddy 214
 6.2 The probability of being back in roulette in three steps after having started there under transition matrix P with (i) $p = b = 0.5$, (ii) $p = 0.85$, (iii) $p = 0.35$ [black line with dots] and (iv) $p = 0.15$, $b = 0.95$ (red line with pluses) 216

205

25.3	PDF $f_{X_1}(x) = \mathbb{I}_{(0,1)}(x)(1 - \cos(2\pi x))$ of the RV X_1 , [the left sub-figure] and its DF $F_{X_1}(x) = \int_x^{\infty} \mathbb{I}_{(0,1)}(u)(1 - \cos(2\pi u))du$ [the right sub-figure]. One can see that convergence of the DFs F_n to $F_{X_1}(x)$, the PDF of the Uniform(0,1) RV, while the corresponding PDFs $f_n(x)$ keep oscillating wildly across [0,1], RV, which corresponds to the claim that uniform(0,1) RV X , thus giving a counterexample to the claim that convergence in DFs does not imply convergence in PDFs.
25.2	Distribution functions of several Normal(μ, σ^2) RVs for $\sigma^2 = 1, \frac{1}{10}, \frac{1}{100}, \frac{1}{1000}, \dots$. 197
25.1	Sequence of Point Mass($\sum_{j=1}^J \delta_j$) RVs (left panel) and {Point Mass($1/\sum_{j=1}^J \delta_j$) RVs (only the first seven are shown on right panel)} and their limiting RVs in red. 197

4.10 The PDF $F(x; 0.3, 0.7)$ of the de Moivre($0.3, 0.7$) RV and its inverse $F_{[1]}^{-1}(u; 0.3, 0.7)$.	183
4.11 Visual Configuration Tool GUI: Rejection Sampling from $X \sim \text{Normal}(0, 1)$ with PDF f based on proposals from $Y \sim \text{Normal}(0, 1)$ with PDF g .	188
4.12 Rejection Sampling from $Y \sim \text{Normal}(0, 1)$ with PDF f based on 100 proposals from $Y \sim \text{Laplace}(1)$ with PDF g .	192
4.13	191

14.1	A plot of the PDF, D_F or CDF, and inverse DF of the Uniform(-1, 1) RV X .	174
14.5	Visual Cognitive Tool GUI: Inverse Sampling from $X \sim \text{Uniform}(-5, 5)$.	175
14.6	The PDF f , CDF F , and inverse DF F^{-1} of the Exponential($\lambda = 1.0$) RV.	176
14.7	Visual Cognitive Tool GUI: Inverse Sampling from $X \sim \text{Exponential}(0.5)$.	177
14.8	Visual Cognitive Tool GUI: Inverse Sampling from $X \sim \text{Laplace}(5)$.	179
14.9	Visual Cognitive Tool GUI: Inverse Sampling from $X \sim \text{Cauchy}$.	180

¹ $\alpha = \alpha_1 + \alpha_2 + \dots + \alpha_n$, $\beta = \beta_1 + \beta_2 + \dots + \beta_n$, $\gamma = \gamma_1 + \gamma_2 + \dots + \gamma_n$, $\delta = \delta_1 + \delta_2 + \dots + \delta_n$.

- The number of paths that lead to a (x_1, x_2) with $x_1 + x_2 = n$ is equal to $\binom{x_1}{n}$. We have already seen this as random walks in Maharashtra.
 - $\binom{lx}{n} \theta^{x_1} (1-\theta)^{x_2}$

Answer (Exercise 3.36) — (You should also be able to do this by inspection.)

we can use earlier terms to get

3. Since $f(x)$ is the density function of an Exponentiaal(λ) random variable with parameter $\lambda = 2$, we can use our results for $\text{Exponentiaal}(\lambda)$ to get

Therefore, $\mathbb{E}(X) = \frac{2}{8+0} = \frac{1}{4}$, and

$$g(x) = \begin{cases} x & \text{otherwise} \\ 8 & x > 0 \end{cases} = (x)f$$

9 The density function of the uniform distribution X on $[0, 8]$ is

$$V(X) = \mathbb{E}(X^2) - (\mathbb{E}(X))^2 = 15.1667 - 3.5^2 = 2.9167.$$

$$\mathbb{E}(X_{\bar{c}}) = \frac{1}{1} + \frac{9}{4} + \frac{9}{16} + \frac{25}{64} + \frac{36}{36} = \frac{9}{16} = 15.1667$$

and so the variance is

$$\mathcal{Z} \mathcal{E} = \frac{9}{9} + \frac{9}{\mathcal{Z}} + \frac{9}{\frac{9}{4}} + \frac{9}{\frac{9}{3}} + \frac{9}{\frac{9}{2}} + \frac{9}{\frac{9}{1}} = (^t x) f ^t x \sum_{9}^{1=t} = (X) \mathbf{E}$$

$$\begin{aligned} g &= x \quad \frac{9}{1} \\ g &= x \quad \frac{9}{1} \\ f &= x \quad \frac{9}{1} \end{aligned} \quad \text{os parte}$$

$$\begin{matrix} \xi = x \\ \zeta = x \\ \eta = x \end{matrix} \left. \begin{array}{c} \scriptstyle 9 \\ \scriptstyle 1 \end{array} \right\} = (x) f$$

Answer (Ex. 3.35) — 1. The probability mass function of X is

Answer (Ex. 3.35) — 1. The probability mass function of X is

2. The profit per conditioner is \$55, and so the expected daily profit given by

$$E(55X) = 55E(X) = 55 \times 11.7 = 643.50,$$

is \$643.50.

Answer (Ex. 3.31) — The expected value $\mathbf{E}(X)$ is

$$\begin{aligned}\mathbf{E}(X) &= \int_0^1 6x(1-x)x \, dx \\ &= \int_0^1 (6x^2 - 6x^3) \, dx \\ &= 2x^3 - \frac{6}{4}x^4 \Big|_0^1 \\ &= 2(1^3 - 0) - \frac{6}{4}(1^4 - 0) \\ &= 0.5\end{aligned}$$

Chapter 1

Preliminaries

1.1 Elementary Set Theory

A **set** is a collection of distinct objects. We write a set by enclosing its elements with curly braces. For example, we denote a set of the two objects \circ and \bullet by:

$$\{\circ, \bullet\}.$$

Sometimes, we give names to sets. For instance, we might call the first example set A and write:

$$A = \{\circ, \bullet\}.$$

We do not care about the order of elements within a set, i.e. $A = \{\circ, \bullet\} = \{\bullet, \circ\}$. We do not allow a set to contain multiple copies of any of its elements unless the copies are distinguishable, say by labels. So, $B = \{\circ, \bullet, \bullet\}$ is not a set unless the two copies of \bullet in B are labelled or marked to make them distinct, e.g. $B = \{\circ, \tilde{\bullet}, \bullet'\}$. Names for sets that arise in a mathematical discourse are given upper-case letters (A, B, C, D, \dots). Special symbols are reserved for commonly encountered sets.

Here is the set \mathfrak{G} of twenty two Greek lower-case alphabets that we may encounter later:

$$\mathfrak{G} = \{ \alpha, \beta, \gamma, \delta, \epsilon, \zeta, \eta, \theta, \kappa, \lambda, \mu, \nu, \xi, \pi, \rho, \sigma, \tau, \upsilon, \phi, \chi, \psi, \omega \}.$$

They are respectively named alpha, beta, gamma, delta, epsilon, zeta, eta, theta, kappa, lambda, mu, nu, xi, pi, rho, sigma, tau, upsilon, phi, chi, psi and omega. *LHS* and *RHS* are abbreviations for objects on the Left and Right Hand Sides, respectively, of some binary relation. By the notation:

$$LHS := RHS,$$

we mean that *LHS is equal, by definition, to RHS*.

The set which does not contain any element (the collection of nothing) is called the **empty set**:

$$\emptyset := \{ \}.$$

We say an element b **belongs to** a set B , or simply that b belongs to B or that b is an element of B , if b is one of the elements that make up the set B , and write:

$$b \in B.$$

When b **does not belong to** B , we write:

$$b \notin B.$$

$$\begin{aligned}\mathbf{E}(X^2) &= \int_0^1 6x(1-x)x^2 \, dx \\ &= \int_0^1 (6x^3 - 6x^4) \, dx \\ &= \frac{6}{4}x^4 - \frac{6}{5}x^5 \Big|_0^1 \\ &= \frac{6}{4}(1^4 - 0) - \frac{6}{5}(1^5 - 0) \\ &= 0.3\end{aligned}$$

the variance is $\mathbf{V}(X) = \mathbf{E}(X^2) - (\mathbf{E}(X))^2 = 0.3 - 0.5^2 = 0.05$

Answer (Ex. 3.32) — This was already done in Sec. 3.8.2 on Properties of Expectation. Make sure you understand each step.

Answer (Ex. 3.33) — Using the definition of variance, expectations and by completing the square, we get:

$$\begin{aligned}V(aX + b) &= E((aX + b)^2) - (E(aX + b))^2 = E((aX)^2 + 2aXb + b^2) - (aE(X) + b)^2 \\ &= a^2E(X^2) + 2abE(X) + b^2 - a^2(E(X))^2 - 2abE(X) - b^2 = a^2(E(X^2) - (E(X))^2) = a^2V(X).\end{aligned}$$

Answer (Ex. 3.27) — Since $y = g(x) = \sqrt{x}$ is a monotone increasing function for $x \geq 0$, we can apply the change of variable formula. Now $x = g^{-1}(y) = y^2$ is a monotone increasing function for $y \geq 0$ so on this interval

$$f_Y(y) = f_X(g^{-1}(y)) \times |2y| = Ae^{-\lambda y^2} \times 2y = 2Aye^{-\lambda y^2}.$$

Therefore

So the probability density function of Y is given by

$$f_Y(y) = \begin{cases} 0 & y < 0 \\ 2Aye^{-\lambda y^2} & y \geq 0 \end{cases}$$

So the probability density function of Y is given by

Answer (Ex. 3.28) — First note that $y = g(x) = \log_e(x)$ is a monotone increasing function over $a \leq x \leq b$, so we can apply the change of variable formula.

We can formally express our definition of set union as:

The union of two sets C and D , written as $C \cup D$, is the set of elements that

$$C \cup D := \{x : x \in C \text{ or } x \in D\}.$$

When a colon ($:$) appears inside a set, it stands for such that. Thus, the above expression is read as, C union D is equal by definition to the set of all elements x , such that x belongs to C or D .

$$C \neq D$$

When two sets C and D are not equal by the above definition, we say that C is not equal to D and write:

$$C = D \iff C \subset D, D \subset C.$$

We say that two sets C and D are equal (as sets) and write $C = D$ if and only if (\iff) every element of C is also an element of D . By this definition, any set is a subset of itself.

$$C \subset D$$

We say that a set C is a subset of another set D and write:

For our example set $A = \{\circ, \bullet\}$, $\star \notin A$ but $\bullet \in A$.

Classwork 1 (Fruits and colours) Consider a set of fruits $F = \{\text{orange, banana, apple}\}$ and a set of colours $C = \{\text{red, green, blue, orange}\}$. Then,

$$(A \cup B)^c = A^c \cup B^c \text{ and } (A \cap B)^c = A^c \cap B^c$$

We say two sets C and D are disjoint if they have no elements in common, i.e. $C \cap D = \emptyset$.

$$B^c := U \setminus B$$

the set of all elements of U that don't belong to B , i.e.:

When a universal set, e.g. U is well-defined, the complement of a given set B denoted by B^c is

$$C \setminus D := \{x : x \in C \text{ and } x \notin D\}.$$

The set-difference or difference of two sets C and D , written as $C \setminus D$, is the set of elements in C that do not belong to D . Formally:

Figure 1.1: Union and intersection of sets shown by Venn diagrams

Venn diagrams are visual aids for set operations as in the diagrams below.

$$C \cup D := \{x : x \in C \text{ or } x \in D\}.$$

Similarly, the intersection of two sets C and D , written as $C \cap D$, is the set of elements that belong to both C and D . Formally:

$$C \cap D := \{x : x \in C \text{ and } x \in D\}.$$

When a colon ($:$) appears inside a set, it stands for such that. Thus, the above expression is read as, C union D is equal by definition to the set of all elements x , such that x belongs to C or D .

$$C \cup D := \{x : x \in C \text{ or } x \in D\}.$$

We can formally express our definition of set union as:

The union of two sets C and D , written as $C \cup D$, is the set of elements that

$$C \neq D$$

and write:

$$C = D \iff C \subset D, D \subset C.$$

of set equality is notationally summarised as follows:

We say that two sets C and D are equal (as sets) and write $C = D$ if and only if (\iff) every element of C is also an element of D , and every element of D is also an element of C . This definition

$$C \subset D$$

it every element of C is also an element of D . By this definition, any set is a subset of itself.

$$C = D \iff C \subset D, D \subset C.$$

$$C \subset D$$

Classwork 1 (Fruits and colours) Consider a set of fruits $F = \{\text{orange, banana, apple}\}$ and a set of colours $C = \{\text{red, green, blue, orange}\}$. Then,

$$(A \cup B)^c = A^c \cup B^c \text{ and } (A \cap B)^c = A^c \cap B^c$$

By drawing Venn diagrams, let us check De Morgan's Laws:

We say two sets C and D are disjoint if they have no elements in common, i.e. $C \cap D = \emptyset$.

$$B^c := U \setminus B$$

the set of all elements of U that don't belong to B , i.e.:

When a universal set, e.g. U is well-defined, the complement of a given set B denoted by B^c is

$$C \setminus D := \{x : x \in C \text{ and } x \notin D\}.$$

The set-difference or difference of two sets C and D , written as $C \setminus D$, is the set of elements in C that do not belong to D . Formally:

Figure 1.1: Union and intersection of sets shown by Venn diagrams

Venn diagrams are visual aids for set operations as in the diagrams below.

$$C \cup D := \{x : x \in C \text{ or } x \in D\}.$$

Similarly, the intersection of two sets C and D , written as $C \cap D$, is the set of elements that belong to both C and D . Formally:

$$C \cap D := \{x : x \in C \text{ and } x \in D\}.$$

When a colon ($:$) appears inside a set, it stands for such that. Thus, the above expression is read as, C union D is equal by definition to the set of all elements x , such that x belongs to C or D .

$$C \cup D := \{x : x \in C \text{ or } x \in D\}.$$

We can formally express our definition of set union as:

The union of two sets C and D , written as $C \cup D$, is the set of elements that

$$C \neq D$$

and write:

$$C = D \iff C \subset D, D \subset C.$$

$$C \subset D$$

of set equality is notationally summarised as follows:

We say that two sets C and D are equal (as sets) and write $C = D$ if and only if (\iff) every element of C is also an element of D , and every element of D is also an element of C . This definition

$$C = D \iff C \subset D, D \subset C.$$

$$C \subset D$$

Classwork 1 (Fruits and colours) Consider a set of fruits $F = \{\text{orange, banana, apple}\}$ and a set of colours $C = \{\text{red, green, blue, orange}\}$. Then,

$$(A \cup B)^c = A^c \cup B^c \text{ and } (A \cap B)^c = A^c \cap B^c$$

By drawing Venn diagrams, let us check De Morgan's Laws:

We say two sets C and D are disjoint if they have no elements in common, i.e. $C \cap D = \emptyset$.

$$B^c := U \setminus B$$

the set of all elements of U that don't belong to B , i.e.:

When a universal set, e.g. U is well-defined, the complement of a given set B denoted by B^c is

$$C \setminus D := \{x : x \in C \text{ and } x \notin D\}.$$

The set-difference or difference of two sets C and D , written as $C \setminus D$, is the set of elements in C that do not belong to D . Formally:

Figure 1.1: Union and intersection of sets shown by Venn diagrams

Venn diagrams are visual aids for set operations as in the diagrams below.

$$C \cup D := \{x : x \in C \text{ or } x \in D\}.$$

Similarly, the intersection of two sets C and D , written as $C \cap D$, is the set of elements that belong to both C and D . Formally:

$$C \cap D := \{x : x \in C \text{ and } x \in D\}.$$

When a colon ($:$) appears inside a set, it stands for such that. Thus, the above expression is read as, C union D is equal by definition to the set of all elements x , such that x belongs to C or D .

$$C \cup D := \{x : x \in C \text{ or } x \in D\}.$$

We can formally express our definition of set union as:

The union of two sets C and D , written as $C \cup D$, is the set of elements that

$$C \neq D$$

and write:

$$C = D \iff C \subset D, D \subset C.$$

$$C \subset D$$

of set equality is notationally summarised as follows:

We say that two sets C and D are equal (as sets) and write $C = D$ if and only if (\iff) every element of C is also an element of D , and every element of D is also an element of C . This definition

$$C = D \iff C \subset D, D \subset C.$$

$$C \subset D$$

Classwork 1 (Fruits and colours) Consider a set of fruits $F = \{\text{orange, banana, apple}\}$ and a set of colours $C = \{\text{red, green, blue, orange}\}$. Then,

$$(A \cup B)^c = A^c \cup B^c \text{ and } (A \cap B)^c = A^c \cap B^c$$

By drawing Venn diagrams, let us check De Morgan's Laws:

We say two sets C and D are disjoint if they have no elements in common, i.e. $C \cap D = \emptyset$.

$$B^c := U \setminus B$$

the set of all elements of U that don't belong to B , i.e.:

When a universal set, e.g. U is well-defined, the complement of a given set B denoted by B^c is

$$C \setminus D := \{x : x \in C \text{ and } x \notin D\}.$$

The set-difference or difference of two sets C and D , written as $C \setminus D$, is the set of elements in C that do not belong to D . Formally:

Figure 1.1: Union and intersection of sets shown by Venn diagrams

Venn diagrams are visual aids for set operations as in the diagrams below.

$$C \cup D := \{x : x \in C \text{ or } x \in D\}.$$

Similarly, the intersection of two sets C and D , written as $C \cap D$, is the set of elements that belong to both C and D . Formally:

$$C \cap D := \{x : x \in C \text{ and } x \in D\}.$$

When a colon ($:$) appears inside a set, it stands for such that. Thus, the above expression is read as, C union D is equal by definition to the set of all elements x , such that x belongs to C or D .

$$C \cup D := \{x : x \in C \text{ or } x \in D\}.$$

We can formally express our definition of set union as:

The union of two sets C and D , written as $C \cup D$, is the set of elements that

$$C \neq D$$

and write:

$$C = D \iff C \subset D, D \subset C.$$

$$C \subset D$$

of set equality is notationally summarised as follows:

We say that two sets C and D are equal (as sets) and write $C = D$ if and only if (\iff) every element of C is also an element of D , and every element of D is also an element of C . This definition

$$C = D \iff C \subset D, D \subset C.$$

$$C \subset D$$

Classwork 1 (Fruits and colours) Consider a set of fruits $F = \{\text{orange, banana, apple}\}$ and a set of colours $C = \{\text{red, green, blue, orange}\}$. Then,

$$(A \cup B)^c = A^c \cup B^c \text{ and } (A \cap B)^c = A^c \cap B^c$$

By drawing Venn diagrams, let us check De Morgan's Laws:

We say two sets C and D are disjoint if they have no elements in common, i.e. $C \cap D = \emptyset$.

$$B^c := U \setminus B$$

the set of all elements of U that don't belong to B , i.e.:

When a universal set, e.g. U is well-defined, the complement of a given set B denoted by B^c is

$$C \setminus D := \{x : x \in C \text{ and } x \notin D\}.$$

The set-difference or difference of two sets C and D , written as $C \setminus D$, is the set of elements in C that do not belong to D . Formally:

Figure 1.1: Union and intersection of sets shown by Venn diagrams

Venn diagrams are visual aids for set operations as in the diagrams below.

$$C \cup D := \{x : x \in C \text{ or } x \in D\}.$$

Similarly, the intersection of two sets C and D , written as $C \cap D$, is the set of elements that belong to both C and D . Formally:

$$C \cap D := \{x : x \in C \text{ and } x \in D\}.$$

When a colon ($:$) appears inside a set, it stands for such that. Thus, the above expression is read as, C union D is equal by definition to the set of all elements x , such that x belongs to C or D .

$$C \cup D := \{x : x \in C \text{ or } x \in D\}.$$

We can formally express our definition of set union as:

The union of two sets C and D , written as $C \cup D$, is the set of elements that

$$C \neq D$$

and write:

$$C = D \iff C \subset D, D \subset C.$$

$$C \subset D$$

of set equality is notationally summarised as follows:

We say that two sets C and D are equal (as sets) and write $C = D$ if and only if (\iff) every element of C is also an element of D , and every element of D is also an element of C . This definition

$$C = D \iff C \subset D, D \subset C.$$

$$C \subset D$$

Classwork 1 (Fruits and colours) Consider a set of fruits $F = \{\text{orange, banana, apple}\}$ and a set of colours $C = \{\text{red, green, blue, orange}\}$. Then,

$$(A \cup B)^c = A^c \cup B^c \text{ and } (A \cap B)^c = A^c \cap B^c$$

By drawing Venn diagrams, let us check De Morgan's Laws:

We say two sets C and D are disjoint if they have no elements in common, i.e. $C \cap D = \emptyset$.

$$B^c := U \setminus B$$

the set of all elements of U that don't belong to B , i.e.:

When a universal set, e.g. U is well-defined, the complement of a given set B denoted by B^c is

$$C \setminus D := \{x : x \in C \text{ and } x \notin D\}.$$

The set-difference or difference of two sets C and D , written as $C \setminus D$, is the set of elements in C that do not belong to D . Formally:

Figure 1.1: Union and intersection of sets shown by Venn diagrams

Venn diagrams are visual aids for set operations as in the diagrams below.

$$C \cup D := \{x : x \in C \text{ or } x \in D\}.$$

Similarly, the intersection of two sets C and D , written as $C \cap D$, is the set of elements that belong to both C and D . Formally:

$$C \cap D := \{x : x \in C \text{ and } x \in D\}.$$

When a colon ($:$) appears inside a set, it stands for such that. Thus, the above expression is read as, C union D is equal by definition to the set of all elements x , such that x belongs to C or D .

$$C \cup D := \{x : x \in C \text{ or } x \in D\}.$$

We can formally express our definition of set union as:

The union of two sets C and D , written as $C \cup D$, is the set of elements that

$$C \neq D$$

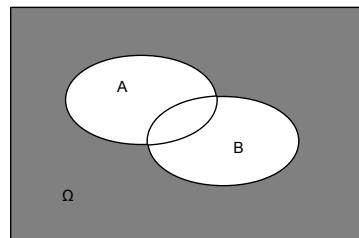
and write:

$$C = D \iff C \subset D, D \subset C.$$

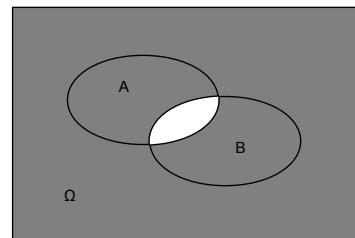
$$C \subset D$$

of set equality is notationally summarised as follows:

We say that two sets C and D are equal (as sets) and write $C = D$ if and only if (\iff) every element of C is also an element of D , and every element of D is also an element of C .



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



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

Figure 1.2: These Venn diagram illustrate De Morgan's Laws.

$$1. F \cap C =$$

$$2. F \cup C =$$

$$3. F \setminus C =$$

$$4. C \setminus F =$$

Classwork 2 (Subsets of a universal set) Suppose we are given a universal set U , and three of its subsets, A , B and C . Also suppose that $A \subset B \subset C$. Find the circumstances, if any, under which each of the following statements is true (T) and justify your answer:

- | | | | |
|---------------------------|--------------------------------|---------------------------|------------------------|
| (1) $C \subset B$ | T when $B = C$ | (2) $A \subset C$ | T by assumption |
| (3) $C \subset \emptyset$ | T when $A = B = C = \emptyset$ | (4) $\emptyset \subset A$ | T always |
| (5) $C \subset U$ | T by assumption | (6) $U \subset A$ | T when $A = B = C = U$ |

1.2 Exercises

Ex. 1.1 — Let Ω be the universal set of students, lecturers and tutors involved in a course. Now consider the following subsets:

- The set of 50 students, $S = \{S_1, S_2, S_3, \dots, S_{50}\}$.
- The set of 3 lecturers, $L = \{L_1, L_2, L_3\}$.
- The set of 4 tutors, $T = \{T_1, T_2, T_3, T_4\}$.

Note that one of the lecturers also tutors in the course. Find the following sets:

- | | |
|-----------------------|------------------|
| (a) $T \cap L$ | (f) $S \cap L$ |
| (b) $T \cap S$ | (g) $S^c \cap L$ |
| (c) $T \cup L$ | (h) T^c |
| (d) $T \cup L \cup S$ | (i) $T^c \cap L$ |
| (e) S^c | (j) $T^c \cap T$ |

Ex. 1.2 — Using Venn diagram, sketch and check the rule:

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

Note that the second equality above is emphasizing that the inverse image $g[-1](y)$ is indeed the inverse function $g^{-1}(y)$ for this $g : \{1, 2, 3, \dots\} \rightarrow \{2^{-1}, 2^{-2}, 2^{-3}, \dots\}$. Therefore,

$$f_Y(y) = \begin{cases} f_X(-\log_2(y)) = \theta(1 - \theta)^{-\log_2(y)-1} & \text{if } y \in \{2^{-1}, 2^{-2}, 2^{-3}, \dots\} \\ 0 & \text{otherwise.} \end{cases}$$

Answer (Ex. 3.25) — Since X is a Poisson(λ) random variable (by suppressing the ' λ ' in the argument to $f_X(\cdot)$ for notational ease), we get

$$f_X(x) = \mathbf{P}(X = x) = \begin{cases} \frac{\lambda^x e^{-\lambda}}{x!} & \text{for } x = 0, 1, 2, \dots \\ 0 & \text{otherwise.} \end{cases}$$

If $Y = (X + 1)^{-2} = 1/(X + 1)^2$ then

$$\{\dots, 2, 1, 0\} \ni x \xrightarrow{(x+1)^{-2}} y \in \left\{1, \frac{1}{4}, \frac{1}{9}, \dots\right\}$$

and since $y = g(x) = (x + 1)^{-2}$ as it maps or associates each $y \in \{1, \frac{1}{4}, \frac{1}{9}, \dots\}$ to exactly one $x \in \{0, 1, 2, \dots\}$ given by $g^{-1}(y) = y^{-1/2} - 1 = \frac{1}{\sqrt{y}} - 1 = x$, its inverse function, this is because g is *injective* or *one-to-one* as explained here if you want to recall quickly https://en.wikipedia.org/wiki/Injective_function again, SO WE GET:

$$f_Y(y) = \mathbf{P}(Y = y) = \sum_{\{x: g(x)=y\}} f_X(x) = f_X\left(\frac{1}{\sqrt{y}} - 1\right) = \frac{\lambda^{(\frac{1}{\sqrt{y}}-1)} e^{-\lambda}}{(\frac{1}{\sqrt{y}}-1)!}$$

for $y = 1, \frac{1}{4}, \frac{1}{9}, \dots$, and 0 otherwise.

CAUTION: This is a discrete RV and so don't just blindly apply the change of variable formula that only applies to continuous RV with a monotone and one-to-one function g with inverse g^{-1} ; it's just that in this discrete RV setting the inverse image also happens to satisfy these properties. But Poisson is discrete and 'change of variable formula' is hence inapplicable.

Answer (Ex. 3.26) — Since $y = g(x) = e^x$ is a monotone increasing function for $x \geq 0$, we can apply the change of variable formula.

Now $x = g^{-1}(y) = \log_e(y)$ is a monotone increasing function for y in $[1, \infty)$

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

Therefore

$$f_Y(y) = f_X(\log_e(y)) \times \left| \frac{1}{y} \right| = \log_e(y) e^{-\log_e(y)} \times \frac{1}{y} = \log_e(y) \frac{1}{y^2}$$

since $e^{-\log_e(y)} = e^{\log_e(y^{-1})} = y^{-1}$.

So the probability density function of Y is given by

$$f_Y(y) = \begin{cases} \frac{\log_e(y)}{y^2} & \text{if } x \geq 1 \\ 0 & \text{otherwise} \end{cases}.$$

$$\int_2^0 h e^{-x} dx = 1$$

$$[-h e^{-x}]_2^0 = 1$$

$$h(-e^{-2} + 1) = 1$$

$$h = \frac{1 - e^{-2}}{1} \quad (\approx 1.1565)$$

$$\mathbf{P}(X \geq 1) = 1 - \mathbf{P}(X < 1)$$

$$\begin{aligned} &= 1 - \int_1^0 h e^{-x} dx \\ &= 1 - h \left[e^{-x} \right]_1^0 \\ &= 1 + \frac{h}{e-1} - \frac{1}{e-1} \\ &\approx 0.2689 \end{aligned}$$

y	$f_X(3) = \frac{1}{6}$	$f_X(2) + f_X(4) = \frac{4}{6}$	$f_X(1) + f_X(5) = \frac{9}{6}$	$f_X(6) = \frac{9}{6}$
0	1	4	9	9

Answer (Ex. 3.23) — The probability mass function $f_Y(y; n)$ for $Y = |X|$, the absolute value of

$$f_Y(y) = \sum_{x \in g(x)} f_X(x; n) \quad \text{as follows:}$$

$$\left\{ \begin{array}{ll} \text{if } y = 0 & \sum_{x \in g(x)} f_X(x; n) = 0 \\ \text{if } y \in \{1, 2, \dots, n\} & \sum_{x \in g(x)} f_X(x; n) = \frac{2n+1}{2} \\ \text{otherwise} & \sum_{x \in g(x)} f_X(x; n) = \frac{2n+1}{2} \end{array} \right. \quad \text{if } y = 0$$

Answer (Ex. 3.24) — We are given that $Y = 2-X$. Define the function

$$g : \{1, 2, 3, \dots\} \rightarrow \{2-1, 2-2, 2-3, \dots\}$$

Answer (Ex. 3.25) — Then y is one-to-one and so by Equation (3.37),

$$\sum_{x \in g^{-1}(y)} f_X(x; n) = f_X(y; n) \quad \text{by } y = g(x) = 2-x.$$

A product set is the **Cartesian product** (\times) of two or more possibly distinct sets:

$$A \times B := \{(a, b) : a \in A \text{ and } b \in B\}$$

For example, if $A = \{\circ, \bullet\}$ and $B = \{\ast\}$, then $A \times B = \{(\circ, \ast), (\bullet, \ast)\}$. Elements of $A \times B$ are called **ordered pairs**.

$$\mathbb{Z}^+ := \mathbb{N} \setminus \{0\} = \{1, 2, 3, \dots\}.$$

The set of **non-negative integers** is:

$$0 = \#\emptyset = \#\{\}$$

For our example sets, $A = \{\circ, \bullet\}$ and the set of Greek alphabets \mathcal{G} , $\#A = 2$ and $\#\mathcal{G} = 22$. The number zero may be defined as the size of an empty set:

⋮

$$\begin{aligned} 2 &= \#\{\ast, \bullet\} = \#\{\bullet, \circ\} = \#\{\circ, \ast\} = \#\{\circ, \circ\} = \#\{\ast, \ast\} = \cdots, \\ 1 &= \#\{\bullet\} = \#\{\circ\} = \#\{\{\bullet\}\} = \{\{\circ\}\} = \cdots, \end{aligned}$$

$$\mathbb{N} := \{1, 2, 3, 4, \dots\}, \text{ may be defined using } \# \text{ as follows:}$$

In fact, the Hindu-Arab numerals we have inherited are based on this intuition of the size of a collection. The elements of the set of **natural numbers**:

$$\#\mathbb{N} = \text{Number of elements in the set } B.$$

We denote the number of elements in a set named B by:

1.3 Natural Numbers, Integers and Rational Numbers

SET SUMMARY

$\{a_1, a_2, \dots, a_n\}$	a set containing the elements, a_1, a_2, \dots, a_n .
$a \in A$	a is an element of the set A .
$A \subseteq B$	the set A is a subset of B .
$A \cup B$	“union”, meaning the set of all elements which are in A or B , or both.
$A \cap B$	“intersection”, meaning the set of all elements which are not in A .
$\{\}$ or \emptyset	empty set.
\mathcal{U}	universal set.
A^c	the complement of A , meaning the set of all elements in \mathcal{U} which are not in A .
\mathcal{Q}	universal set, which are not in A .
\mathcal{O}	empty set.
$\#\mathcal{O}$	“intersection”, meaning the set of all elements which are in both A and B .
$\#\mathcal{B}$	“union”, meaning the set of all elements which are in A or B , or both.
$\#\mathcal{A}$	a is an element of the set A .
$\#\mathcal{U}$	the set A is a subset of B .
$\#\mathcal{Q}$	“union”, meaning the set of all elements which are not in A .
$\#\mathcal{O}^c$	“intersection”, meaning the set of all elements which are in A and only if $A \cup B = B$.

Ex. 1.4 — Using a Venn diagram, illustrate the idea that $A \subseteq B$ if and only if $A \cup B = B$.

Ex. 1.3 — Using Venn diagram, sketch and check the rule:

The binary arithmetic operation of **addition** (+) between a pair of non-negative integers $c, d \in \mathbb{Z}_+$ can be defined via sizes of disjoint sets. Suppose, $c = \#C$, $d = \#D$ and $C \cap D = \emptyset$, then:

$$c + d = \#C + \#D := \#(C \cup D).$$

For example, if $A = \{\circ, \bullet\}$ and $B = \{\star\}$, then $A \cap B = \emptyset$ and $\#A + \#B = \#(A \cup B) \iff 2 + 1 = 3$.

The binary arithmetic operation of **multiplication** (\cdot) between a pair of non-negative integers $c, d \in \mathbb{Z}_+$ can be defined via sizes of product sets. Suppose, $c = \#C$, $d = \#D$, then:

$$c \cdot d = \#C \cdot \#D := \#(C \times D).$$

For example, if $A = \{\circ, \bullet\}$ and $B = \{\star\}$, then $\#A \cdot \#B = \#(A \times B) \iff 2 \cdot 1 = 2$.

More generally, a product set of A_1, A_2, \dots, A_m is:

$$A_1 \times A_2 \times \cdots \times A_m := \{(a_1, a_2, \dots, a_m) : a_1 \in A_1, a_2 \in A_2, \dots, a_m \in A_m\}$$

Elements of an m -product set are called **ordered m -tuples**. When we take the product of the same set we abbreviate as follows:

$$A^m := \underbrace{A \times A \times \cdots \times A}_{m \text{ times}} := \{(a_1, a_2, \dots, a_m) : a_1 \in A, a_2 \in A, \dots, a_m \in A\}$$

Classwork 3 (Cartesian product of sets) 1. Let $A = \{\circ, \bullet\}$. What are the elements of A^2 ? 2. Suppose $\#A = 2$ and $\#B = 3$. What is $\#(A \times B)$? 3. Suppose $\#A_1 = s_1, \#A_2 = s_2, \dots, \#A_m = s_m$. What is $\#(A_1 \times A_2 \times \cdots \times A_m)$?

Now, let's recall the definition of a function. A **function** is a “mapping” that associates each element in some set \mathbb{X} (the domain) to exactly one element in some set \mathbb{Y} (the range). Two different elements in \mathbb{X} can be mapped to or associated with the same element in \mathbb{Y} , and not every element in \mathbb{Y} needs to be mapped. Suppose $x \in \mathbb{X}$. Then we say $f(x) = y \in \mathbb{Y}$ is the **image** of x . To emphasise that f is a **function** from $\mathbb{X} \ni x$ to $\mathbb{Y} \ni y$, we write:

$$f(x) = y : \mathbb{X} \rightarrow \mathbb{Y}.$$

And for some $y \in \mathbb{Y}$, we call the set:

$$f^{[-1]}(y) := \{x \in \mathbb{X} : f(x) = y\} \subset \mathbb{X},$$

the **pre-image** or **inverse image** of y , and

$$f^{[-1]} := f^{[-1]}(y \in \mathbb{Y}) = X \subset \mathbb{X},$$

That is,

$$k x \Big|_{-4}^4 = 1$$

$$k(4 - (-4)) = 1$$

$$8k = 1$$

$$k = \frac{1}{8}$$

2. First note that if $x < -4$, then

$$F(x) = \int_{-\infty}^x 0 \, dv = 0.$$

If $-4 \leq x \leq 4$, then

$$\begin{aligned} F(x) &= \int_{-\infty}^{-4} 0 \, dv + \int_{-4}^x \frac{1}{8} \, dv \\ &= 0 + \left[\frac{1}{8} v \right]_{-4}^x \\ &= \frac{1}{8}(x + 4) \end{aligned}$$

If $x \geq 4$, then

$$\begin{aligned} F(x) &= \int_{-\infty}^{-4} 0 \, dv + \int_{-4}^4 \frac{1}{8} \, dv + \int_4^x 0 \, dv \\ &= 0 + \left[\frac{1}{8} v \right]_{-4}^4 + 0 \\ &= 1 \end{aligned}$$

Hence

$$F(x) = \begin{cases} 0 & x < -4 \\ \frac{1}{8}(x + 4) & -4 \leq x \leq 4 \\ 1 & x \geq 4 \end{cases}$$

3. The graphs of $f(x)$ and $F(x)$ for random variable X are as follows:

Answer (Ex. 3.20) — 1. Since the distribution function is $F(t; \lambda) = 1 - \exp(-\lambda t)$,

$$\mathbf{P}(t > \tau) = 1 - \mathbf{P}(t < \tau) = 1 - F(\tau; \lambda = 0.01) = 1 - (1 - e^{-0.01\tau}) = e^{-0.01\tau}.$$

2. Set

$$\mathbf{P}(t > \tau) = e^{-0.01\tau} = \frac{1}{2}$$

and solve for τ to get then $\tau = -100 \times \log(0.5) = 69.3$ (3 sig. fig.).

whose irreducible unique expression is $1/2$.
For example, $1/2$, $2/4$, $3/6$, and $1001/2002$ are different expressions for the same rational number $1/2$, where y is positive and as small as possible. Rational number has a unique irreducible expression p/y , where y is positive and as small as possible. The expressions p/y and $p \cdot q/q$ denote the same rational number if and only if $p \cdot q = p \cdot q$. Every rational number has a unique irreducible expression p/q .

$$\mathbb{Q} := \{p/q : p \in \mathbb{Z}, q \in \mathbb{Z} \setminus \{0\}\}$$

If the magnitude of the entity's position is measured in units (e.g. metres) that can be rationaly divided into q pieces with $q \in \mathbb{N}$, then we have the set of rational numbers:

$$+ : \mathbb{Z} \times \mathbb{Z} \hookrightarrow$$

What is its range?

$$\mathbb{Z}^2 := \mathbb{Z} \times \mathbb{Z} = \{(a, b) : a \in \mathbb{Z}, b \in \mathbb{Z}\}$$

Cartesian product of \mathbb{Z} :

Try to set up the arithmetic operation of addition as a function. The domain for addition is the

Classwork 4 (Addition over integers) Consider the set of integers $\mathbb{Z} = \{\dots, -3, -2, -1, 0, 1, 2, \dots\}$.

Every integer is either positive, negative, or zero. In terms of this we define the notion of integers. The reader is assumed to be familiar with such arithmetic operations with pairs of integers. We say an integer a is **less than or equal to** an integer b and write $a \leq b$ if $b - a$ is positive. We say an integer a is **less than** an integer b and write $a < b$ if $b - a$ is positive. We say an integer a is **greater than** an integer b and write $a > b$ if $b - a$ is negative. Finally, we say that a is greater than b and write $a > b$ if $a - b$ is positive or zero. In the set of integers are **well-ordered**, i.e., for every integer a there is a next largest integer $a + 1$.

$$\mathbb{Z} := \{\dots, -3, -2, -1, 0, +1, +2, +3, \dots\}.$$

We motivated the non-negative integers \mathbb{Z}^+ via the size of a set. With the notion of two directions as the **inverse** of f .

we can motivate the set of integers:

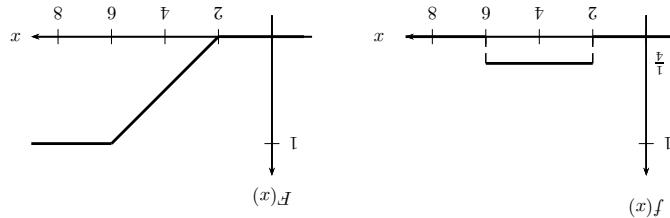
with a **plus or positive sign** (+) before them are called positive integers. Conventionally, + signs are dropped. Some examples of fractions you may have encountered are **arithmetic operations** such as **addition** (+), **subtraction** (-), **multiplication** (-) and **division** (/) of ordered pairs of integers. The reader is assumed to be familiar with such arithmetic operations with pairs of integers. The reader is assumed to be familiar with such arithmetic operations with pairs of integers. The integers with a **minus or negative sign** (-) before them are called negative integers and those

Figure 1.3: A function f ("father of") from \mathbb{X} (a set of children) to \mathbb{Y} (their fathers) and its inverse ("children of").

Answer (Ex. 3.19) — 1. Since $f(x)$ is a (continuous) probability density function which inte-

grates to one,

$$\int_{-4}^{-4} k dx = 1.$$



Graphs of $f(x)$ and $F(x)$.

so the graphs are:

$$F(x) = \begin{cases} 0 & x < 2 \\ \frac{1}{4}(x-2) & 2 \leq x < 6 \\ 1 & x \geq 6 \end{cases}$$

(b) Now

$$\begin{aligned} h &= \frac{1}{4} \\ 1 &= 4h \\ 1 &= 6h - 2h \\ 1 &= kx^6 \\ 1 &= \int_6^2 f(x) dx = \int_6^2 k dx \end{aligned}$$

(a) Since $f(x)$ is a density function which integrates to one,

Answer (Exercise 3.18) —

Answer (Exercise 3.17) — Let $X \sim \text{Exponential}(\lambda = 0.1)$ denote the time taken to serve any given customer in an IID manner. Let y denote the unknown time that the current customer being served has already been served before your arrival. By memorylessness of Exponential($\lambda = 0.1$) RV X , we know $\mathbf{P}(X < 2 + y | X < y) = \mathbf{P}(X < 2) = e^{-\lambda 2} = e^{-2/10}$.

4. Occurrence of lacunae may not always be independent. For example, a machine malfunction may cause them to be clumped.

Figure 1.4: A pictorial depiction of addition and its inverse. The domain is plotted in orthogonal **Cartesian coordinates**.

Addition and multiplication are defined for rational numbers by:

$$\frac{p}{q} + \frac{p'}{q'} = \frac{p \cdot q' + p' \cdot q}{q \cdot q'} \quad \text{and} \quad \frac{p}{q} \cdot \frac{p'}{q'} = \frac{p \cdot p'}{q \cdot q'}.$$

The rational numbers form a **field** under the operations of addition and multiplication defined above in terms of addition and multiplication over integers. This means that the following properties are satisfied:

1. Addition and multiplication are each **commutative**

$$a + b = b + a, \quad a \cdot b = b \cdot a,$$

and associative

$$a + (b + c) = (a + b) + c, \quad a \cdot (b \cdot c) = (a \cdot b) \cdot c.$$

2. Multiplication **distributes** over addition

$$a \cdot (b + c) = (a \cdot b) + (a \cdot c).$$

3. 0 is the **additive identity** and 1 is the multiplicative identity

$$0 + a = a \quad \text{and} \quad 1 \cdot a = a.$$

4. Every rational number a has a negative, $a + (-a) = 0$ and every non-zero rational number a has a reciprocal, $a \cdot 1/a = 1$.

The field axioms imply the usual laws of arithmetic and allow subtraction and division to be defined in terms of addition and multiplication as follows:

$$\frac{p}{q} - \frac{p'}{q'} := \frac{p}{q} + \frac{-p'}{q'} \quad \text{and} \quad \frac{p}{q} / \frac{p'}{q'} := \frac{p}{q} \cdot \frac{q'}{p'}, \quad \text{provided } p' \neq 0.$$

We will see later that the theory of finite fields is necessary for the study of pseudo-random number generators (PRNGs) and PRNGs are the heart-beat of randomness and statistics with computers.

Answer (Ex. 3.13) — This is a Binomial experiment with parameters $\theta = 0.1$ and $n = 10$, and so

$$\mathbf{P}(X \geq 1) = 1 - \mathbf{P}(X < 1) = 1 - \mathbf{P}(X = 0),$$

where

$$\mathbf{P}(X = 0) = \binom{10}{0} 0.1^0 0.9^{10} \approx 0.3487.$$

Therefore, the probability that the target will be hit at least once is

$$1 - 0.3487 \approx 0.6513.$$

Answer (Ex. 3.14) — Since 2 defects exist on every 100 meters, we would expect 6 defects on a 300 meter tape. If X is the number of defects on a 300 meter tape, then X is Poisson with $\lambda = 6$ and so the probability of zero defects is

$$\mathbf{P}(X = 0; 6) = \frac{6^0}{0!} e^{-6} = 0.0025.$$

Answer (Ex. 3.15) — Since X is Poisson(λ) random variable with $\lambda = 0.5$, $\mathbf{P}(X \geq 2)$ is the probability of observing two or more particles during any given second.

$$\mathbf{P}(X \geq 2) = 1 - \mathbf{P}(X < 2) = 1 - \mathbf{P}(X = 1) - \mathbf{P}(X = 0),$$

where $\mathbf{P}(X = 1)$ and $\mathbf{P}(X = 0)$ can be carried out by the Poisson probability mass function

$$\mathbf{P}(X = x) = f(x) = \frac{\lambda^x}{x!} e^{-\lambda}.$$

Now

$$\mathbf{P}(X = 0) = \frac{0.5^0}{0!} \times e^{-0.5} = 0.6065$$

and

$$\mathbf{P}(X = 1) = \frac{0.5^1}{1!} \times e^{-0.5} = 0.3033$$

and so

$$\mathbf{P}(X \geq 2) = 1 - 0.9098 = 0.0902.$$

Answer (Ex. 3.16) — 1.The Probability mass function for Poisson(λ) random variable X is

$$\mathbf{P}(X = x) = f(x; \lambda) = \frac{e^{-\lambda} \lambda^x}{x!}$$

where λ is the mean number of lacunae per specimen and X is the random variable “number of lacunae on a specimen”.

2.If $x = 0$ then $x! = 0! = 1$ and $\lambda^x = \lambda^0 = 1$, and the formula becomes $\mathbf{P}(X = 0) = e^{-\lambda}$.

3.Since $\mathbf{P}(X \geq 1) = 0.1$,

$$\mathbf{P}(X = 0) = 1 - \mathbf{P}(X \geq 1) = 0.9.$$

Using (b) and solving for λ gives:

$$e^{-\lambda} = 0.9 \quad \text{that is, } \lambda = -\ln(0.9) = 0.1 \text{ (approximately.)}$$

Hence

$$\mathbf{P}(X = 2) = \frac{e^{-0.1} (0.1)^2}{2!} = 0.45\% \text{ (approximately.)}$$

Figure 1.5: A depiction of the real line segment $[-10, 10]$.

The **half-open interval** $(y, z]$ or $[y, z)$ and the **open interval** (y, z) are defined analogously:

$$\begin{aligned}(y, z] &:= \{x : y < x \leq z\}, \\ [y, z) &:= \{x : y \leq x < z\}, \\ (y, z) &:= \{x : y < x < z\}.\end{aligned}$$

We also allow y to be **minus infinity** (denoted $-\infty$) or z to be **infinity** (denoted ∞) at an open endpoint of an interval, meaning that there is no lower or upper bound. With this allowance we get the set of **real numbers** $\mathbb{R} := (-\infty, \infty)$, the **non-negative real numbers** $\mathbb{R}_+ := [0, \infty)$ and the **positive real numbers** $\mathbb{R}_{>0} := (0, \infty)$ as follows:

$$\begin{aligned}\mathbb{R} &:= (-\infty, \infty) = \{x : -\infty < x < \infty\}, \\ \mathbb{R}_+ &:= [0, \infty) = \{x : 0 \leq x < \infty\}, \\ \mathbb{R}_{>0} &:= (0, \infty) = \{x : 0 < x < \infty\}.\end{aligned}$$

For a positive real number $b \in \mathbb{R}_{>0}$ and an integer $n \in \mathbb{Z}$, the n -th **power** or **exponent** of b is:

$$b^0 = 1, \quad b^n = b^{n-1} \cdot b \quad \text{if } n > 0, \quad b^n = b^{n+1}/b \quad \text{if } n < 0.$$

The following **laws of exponents** hold by mathematical induction when $m, n \in \mathbb{Z}$:

$$b^{m+n} = b^m \cdot b^n, \quad (b^m)^n = b^{m \cdot n}.$$

If $y \in \mathbb{R}$ and $m \in \mathbb{N}$, the unique positive real number $z \in \mathbb{R}_{>0}$ such that $z^m = y$ is called the m -th **root** of y and denoted by $\sqrt[m]{y}$, i.e.,

$$z^m = y \implies z = \sqrt[m]{y}.$$

For a rational number $r = p/q \in \mathbb{Q}$, we define the r -th power of $b \in \mathbb{R}$ as follows:

$$b^r = b^{p/q} := \sqrt[q]{b^p}$$

The laws of exponents hold for this definition and different expressions for the same rational number $r = ap/aq$ yield the same power, i.e., $b^{p/q} = b^{ap/aq}$. Recall that a real number $x = n + 0.d_1d_2d_3\dots \in \mathbb{R}$ can be arbitrarily precisely enclosed by the rational numbers $\underline{x}_k := n + \frac{d_1}{10} + \frac{d_2}{100} + \dots + \frac{d_k}{10^k}$ and $\bar{x}_k := n + \frac{d_1}{10} + \frac{d_2}{100} + \dots + \frac{d_k}{10^k} + \frac{1}{10^k}$ by increasing k . Suppose first that $b > 1$. Then, using rational powers, we can enclose b^x ,

$$b^{n+\frac{d_1}{10}+\frac{d_2}{100}+\dots+\frac{d_k}{10^k}} = b^{\underline{x}_k} \leq b^x < b^{\bar{x}_k} := b^{n+\frac{d_1}{10}+\frac{d_2}{100}+\dots+\frac{d_k}{10^k}+\frac{1}{10^k}},$$

within an interval of width $b^{n+\frac{d_1}{10}+\frac{d_2}{100}+\dots+\frac{d_k}{10^k}} (b^{\frac{1}{10^k}} - 1) < b^{n+1} (b - 1)/10^k$. By taking a large enough k we can evaluate b^x to any accuracy. Finally, when $b < 1$ we define $b^x := (1/b)^{-x}$ and when $b = 0$, $b^x := 1$.

Now $\sum_{x=0}^{\infty} \frac{1}{2^x}$ is a geometric series with common ratio $r = \frac{1}{2}$ and first term $a = 1$, and so has sum

$$S = \frac{a}{1-r} = \frac{1}{1-\frac{1}{2}} = 2$$

Therefore,

$$2k = 1, \quad \text{that is, } k = \frac{1}{2}.$$

2. From (a), the probability mass function of f is

$$f(x) = \frac{\frac{1}{2}}{2^x} = \frac{1}{2^{x+1}}. \quad (x = 0, 1, 2, \dots)$$

Now

$$\mathbf{P}(X \geq 4) = 1 - \mathbf{P}(X < 4) = 1 - \mathbf{P}(X \leq 3)$$

where

$$\begin{aligned}\mathbf{P}(X \leq 3) &= \sum_{x=0}^3 \frac{1}{2^{x+1}} \\ &= \frac{1}{2} + \frac{1}{4} + \frac{1}{8} + \frac{1}{16} \\ &= \frac{8}{16} + \frac{4}{16} + \frac{2}{16} + \frac{1}{16} \\ &= \frac{15}{16}.\end{aligned}$$

$$\text{That is, } \mathbf{P}(X \geq 4) = \frac{1}{16}.$$

Answer (Ex. 3.11) — Note that $\theta = \frac{1}{2}$ here.

1. X has probability mass function

$$f(x) = \begin{cases} \binom{4}{0} \frac{1}{2}^0 \frac{1}{2}^4 = \frac{1}{16} & x = 0 \\ \binom{4}{1} \frac{1}{2}^1 \frac{1}{2}^3 = \frac{4}{16} & x = 1 \\ \binom{4}{2} \frac{1}{2}^2 \frac{1}{2}^2 = \frac{6}{16} & x = 2 \\ \binom{4}{3} \frac{1}{2}^3 \frac{1}{2}^1 = \frac{4}{16} & x = 3 \\ \binom{4}{4} \frac{1}{2}^4 \frac{1}{2}^0 = \frac{1}{16} & x = 4 \end{cases}$$

2. The required probabilities are:

$$\begin{aligned}\mathbf{P}(X = 0) &= f(0) = \frac{1}{16} \\ \mathbf{P}(X = 1) &= f(1) = \frac{4}{16}\end{aligned}$$

$$\sum_{k=0}^{\infty} \frac{z^k}{k!} = 1, \quad \text{that is,} \quad k \sum_{k=0}^{\infty} \frac{z^k}{k!} = 1$$

Answer (Ex. 3.10) — I. Since f is a probability mass function,

$$\mathbf{P}(X < 1) = \mathbf{P}(X = 2) = \frac{3}{1}$$

$$\mathbf{P}(X \geq 1) = \mathbf{P}(X = 1) + \mathbf{P}(X = 2) = \frac{15}{8} + \frac{3}{1} = \frac{15}{8} + \frac{3}{1} = \frac{18}{8} = \frac{9}{4}$$

$$\mathbf{P}(X \leq 1) = \mathbf{P}(X = 0) + \mathbf{P}(X = 1) = \frac{2}{15} + \frac{8}{15} = \frac{2}{3}$$

The required probabilities are:

$$\left\{ \begin{array}{ll} \frac{3}{1} & \text{if } x = 2 \\ \frac{15}{8} & \text{if } x = 1 \\ \frac{2}{15} & \text{if } x = 0 \end{array} \right\} = f(x) \mathbf{P}(x) =$$

So the probability mass function of X is:

$$\mathbf{P}(X = 2) = \frac{10}{6} \cdot \frac{9}{5} = \frac{3}{1} \quad (\text{one way of drawing two left screws}).$$

$$\mathbf{P}(X = 1) = \frac{10}{6} \cdot \frac{4}{4} + \frac{10}{6} \cdot \frac{9}{8} = \frac{15}{8} \quad (\text{two ways of drawing one left and one right screw}),$$

$$\mathbf{P}(X = 0) = \frac{10}{6} \cdot \frac{9}{2} = \frac{15}{2} \quad (\text{one way of drawing two right screws}),$$

Answer (Ex. 3.9) — Since we are sampling without replacement,

$$(vi) \mathbf{P}(4 < X \leq 11) = \mathbf{P}(4 < X \leq 10) = F(10) - F(4) = 0.82 - 0.08 = 0.74$$

$$(v) \mathbf{P}(4 < X \leq 9) = F(9) - F(4) = 0.79 - 0.08 = 0.71$$

$$(iv) \mathbf{P}(X \geq 9) = 1 - \mathbf{P}(X < 9) = 1 - \mathbf{P}(X \leq 8) = 1 - 0.59 = 0.41$$

$$(iii) \mathbf{P}(X \geq 9) = 1 - \mathbf{P}(X \leq 9) = 1 - F(9) = 1 - 0.79 = 0.21$$

$$(ii) \mathbf{P}(X < 12) = \mathbf{P}(X \leq 11) = F(11) = 0.84$$

$$(i) \mathbf{P}(X \leq 5) = F(5) = 0.17$$

x	$\mathbf{P}(X \leq x)$	0.07	0.08	0.17	0.18	0.34	0.59	0.79	0.82	0.84	0.95	1.00
3												
4												
5												
6												
7												
8												
9												
10												
11												
12												
13												

Answer (Ex. 3.8) — (a)

Suppose $y \in \mathbb{R}^>0$ and $b \in \mathbb{R} \setminus \{1\}$ then the real number x such that $y = b^x$ is called the **logarithm** of y to the base b and we write this as:

$$y = b^x \iff x = \log_b y$$

The definition implies:

$$\log_b(xy) = \log_b x + \log_b y, \quad \text{if } x < 0, y < 0 \text{ and}$$

and the laws of exponents imply:

$$\log_b(b^x) = x \log_b b,$$

upper bound of a non-empty set of real numbers A to be the **supremum** of A and denote it as: For example, $\inf(0, 1) = 0$ and $\inf\{10.333 \cup [-99, 1001.33]\} = -99$. We similarly define the **least upper bound** above then $\sup A = \infty$. Finally, if a set A is not bounded below then $\inf A = -\infty$ and if a set A is not bounded above then $\sup A = \infty$.

$$\sup A = \text{least upper bound of } A$$

For example, $\inf(0, 1) = 0$ and $\inf\{10.333 \cup [-99, 1001.33]\} = -99$. By convention, we define the **greatest lower bound** of a set of real numbers A is called the **minimum** of A and is denoted by: For example, $\min(0, 1) = 0$ and $\min\{10.333 \cup [-99, 1001.33]\} = -99$. We say that $a \in A$ is a **lower bound** if it is at least as large as any other lower bound. A **greatest lower bound** is the **greatest lower bound** if it is at least as large as any other lower bound. The **greatest lower bound** is the **minimum** of a set of real numbers A .

For example, $\min\{1, 4, -9, 345\} = -9$, $\min\{-93.8889, 1002.786\} = -93.8889$. We need a slightly more sophisticated notion for the extremal elements of a set A that may not belong to A . We say that a real number x is a **lower bound** for a non-empty set of real numbers A , provided $x \leq a$ for every $a \in A$. We say that the set A is **bounded below** if it has at least one lower bound. A **greatest lower bound** is the **greatest lower bound** if it is at least as large as any other lower bound. The **greatest lower bound** is the **minimum** of a set of real numbers A .

$$\min A = \text{least element in } A$$

For example, $\max\{1, 4, -9, 345\} = 345$, $\max\{-93.8889, 1002.786\} = 1002.786$.

$$\max A = \text{greatest element in } A$$

Familiar extremal elements of a set of real numbers, say A , are the following:

...). We sometimes denote the special power function e^y by $\exp(y)$. You are assumed to be familiar with trigonometric functions ($\sin(x)$, $\cos(x)$, $\tan(x)$) to mean $\log_e(y)$, where e is the Euler's constant. Since we will mostly work with $\log_e(y)$ we use $\log_e(y)$ to is $\log_e(y)$, where e is the Euler's constant. The common logarithm is $\log_{10}(y)$, the binary logarithm is $\log_2(y)$ and the natural logarithm

$$\log_b(c^y) = y \log_b c, \quad \text{if } c < 0.$$

$$\log_b(xy) = \log_b x + \log_b y, \quad \text{if } x < 0, y < 0 \text{ and}$$

and the laws of exponents imply:

$$\log_b(b^x) = x \log_b b,$$

Symbol	Meaning
$A = \{\star, \circ, \bullet\}$	A is a set containing the elements \star, \circ and \bullet
$\circ \in A$	\circ belongs to A or \circ is an element of A
$A \ni \circ$	\circ belongs to A or \circ is an element of A
$\circ \notin A$	\circ does not belong to A
$\#A$	Size of the set A , for e.g. $\#\{\star, \circ, \bullet, \odot\} = 4$
\mathbb{N}	The set of natural numbers $\{1, 2, 3, \dots\}$
\mathbb{Z}	The set of integers $\{\dots, -3, -2, -1, 0, 1, 2, 3, \dots\}$
\mathbb{Z}_+	The set of non-negative integers $\{0, 1, 2, 3, \dots\}$
\emptyset	Empty set or the collection of nothing or {}
$A \subset B$	A is a subset of B or A is contained by B , e.g. $A = \{\circ\}, B = \{\bullet\}$
$A \supset B$	A is a superset of B or A contains B e.g. $A = \{\circ, \star, \bullet\}, B = \{\circ, \bullet\}$
$A = B$	A equals B , i.e. $A \subset B$ and $B \subset A$
$Q \implies R$	Statement Q implies statement R or If Q then R
$Q \iff R$	$Q \implies R$ and $R \implies Q$
$\{x : x \text{ satisfies property } R\}$	The set of all x such that x satisfies property R
$A \cup B$	A union B , i.e. $\{x : x \in A \text{ or } x \in B\}$
$A \cap B$	A intersection B , i.e. $\{x : x \in A \text{ and } x \in B\}$
$A \setminus B$	A minus B , i.e. $\{x : x \in A \text{ and } x \notin B\}$
$A := B$	A is equal to B by definition
$A :=: B$	B is equal to A by definition
A^c	A complement, i.e. $\{x : x \in U, \text{ the universal set, but } x \notin A\}$
$A_1 \times A_2 \times \dots \times A_m$	The m -product set $\{(a_1, a_2, \dots, a_m) : a_1 \in A_1, a_2 \in A_2, \dots, a_m \in A_m\}$
A^m	The m -product set $\{(a_1, a_2, \dots, a_m) : a_1 \in A, a_2 \in A, \dots, a_m \in A\}$
$f := f(x) = y : \mathbb{X} \rightarrow \mathbb{Y}$	A function f from domain \mathbb{X} to range \mathbb{Y}
$f^{[-1]}(y)$	Inverse image of y
$f^{[-1]} := f^{[-1]}(y \in \mathbb{Y}) = X \subset \mathbb{X}$	Inverse of f
$a < b$ or $a \leq b$	a is less than b or a is less than or equal to b
$a > b$ or $a \geq b$	a is greater than b or a is greater than or equal to b
\mathbb{Q}	Rational numbers
(x, y)	the open interval (x, y) , i.e. $\{r : x < r < y\}$
$[x, y]$	the closed interval (x, y) , i.e. $\{r : x \leq r \leq y\}$
$(x, y]$	the half-open interval $(x, y]$, i.e. $\{r : x < r \leq y\}$
$[x, y)$	the half-open interval $[x, y)$, i.e. $\{r : x \leq r < y\}$
$\mathbb{R} := (-\infty, \infty)$	Real numbers, i.e. $\{r : -\infty < r < \infty\}$
$\mathbb{R}_+ := [0, \infty)$	Real numbers, i.e. $\{r : 0 \leq r < \infty\}$
$\mathbb{R}_{>0} := (0, \infty)$	Real numbers, i.e. $\{r : 0 < r < \infty\}$

Table 1.1: Symbol Table: Sets and Numbers

x	observed frequency	observed relative frequency	Prob of x hits
0	229	$229/576 = 0.398$	$f(0; 0.933) = 0.394$
1	211	$211/576 = 0.366$	$f(1; 0.933) = 0.367$
2	93	$93/576 = 0.161$	$f(2; 0.933) = 0.171$
3	35	$35/576 = 0.0608$	$f(3; 0.933) = 0.0532$
4	7	$7/576 = 0.0122$	$f(4; 0.933) = 0.0124$
≥ 5	1	$1/576 = 0.00174$	$1 - \sum_{x=0}^4 f(x; 0.933) = 0.00275$

Answer (Ex. 3.5) — $\mathbf{P}(X = 3)$ does not satisfy the condition that $0 \leq \mathbf{P}(A) \leq 1$ for any event A . If Ω is the sample space, then $\mathbf{P}(\Omega) = 1$ and so the correct probability is

$$\mathbf{P}(X = 3) = 1 - 0.07 - 0.10 - 0.32 - 0.40 = 0.11 .$$

Answer (Ex. 3.6) — 1. Tabulate the values for the probability mass function as follows:

x	1	2	3	4	5
$\mathbf{P}(X = x)$	0.1	0.2	0.2	0.2	0.3

so the distribution function is:

$$F(x) = \mathbf{P}(X \leq x) = \begin{cases} 0 & \text{if } 0 \leq x < 1 \\ 0.1 & \text{if } 1 \leq x < 2 \\ 0.3 & \text{if } 2 \leq x < 3 \\ 0.5 & \text{if } 3 \leq x < 4 \\ 0.7 & \text{if } 4 \leq x < 5 \\ 1 & \text{if } x \geq 5 \end{cases}$$

2. The probability that the machine needs to be replaced during the first 3 years is:

$$\mathbf{P}(X \leq 3) = \mathbf{P}(X = 1) + \mathbf{P}(X = 2) + \mathbf{P}(X = 3) = 0.1 + 0.2 + 0.2 = 0.5 .$$

(This answer is easily seen from the distribution function of X .)

3. The probability that the machine needs no replacement during the first three years is

$$\mathbf{P}(X > 3) = 1 - \mathbf{P}(X \leq 3) = 0.5 .$$

Answer (Ex. 3.7) — Assuming that the probability model is being built from the observed relative frequencies, the probability mass function is:

$$f(x) = \begin{cases} \frac{176}{200} & x = 1 \\ \frac{22}{200} & x = 2 \\ \frac{2}{200} & x = 3 \end{cases}$$

1.5 Introduction to MATLAB

23

CHAPTER I. PRELIMINARIES

Hardware 3 (Basics of MATLAB) Let us familiarize ourselves with MATLAB in this session. This help. The command window within the MATLAB widow is where you need to type commands. You need to launch MATLAB from your terminal. Since this is system dependent, ask your tutor for help. Here is a minimal set of commands you need to familiarize yourself with in this session.

The summand 37 of 13 and 24 is stored in the default variable called `ans` which is short for answer.

2. We can write comments in MATLAB following the % character. All the characters in a given line that follow the percent character % are ignored by MATLAB. It is very helpful to comment what is being done in the code. You won't get full credit without sufficient comments in your coding assignments. For example we could have added the comment to the previous addition.

To save space in these notes, we suppress the blank lines and excessive line breaks present in MATLAB's command window.

```
>>> diary off % turn off the current diary file blah.txt  
>>> diary blah.txt % start a diary file named blah.txt  
>>> ans = 59  
>>> 3+96
```

Answer (Exercise 3A) — We are given that 537 flying bombs made up of $24 \times 24 = 576$ small equal-sized areas, say A_1, A_2, \dots, A_{576} . Let X denote the number of hits that a particular bomb randomly lands over area A . The probability that a particular bomb falls over area A_i is $\frac{1}{576}$. Finally, we can model X as Binomial($n = 537$, $p = \frac{1}{576}$). By Poisson(λ) trials, finally, we can approximate this Binomial distribution with Poisson($\lambda = 0.933$) random variable X where the formula for Poisson(λ) is $\mu = n\theta = \frac{537}{576} \approx 0.933$.

Answer (Exercise 33) — This was done in week 1. Get notes from your mates or wait until Razza scribbles for virtual conference.

$$\left. \begin{array}{l} \infty > x \geq 1, \\ 1 > x \geq 0, \\ 0 > x > -1, \end{array} \right\} = (x \geq X)D = (x)H$$

$$\left. \begin{array}{l} \infty > x \geqslant 1 \quad \text{if} \quad 1 = (\mathfrak{U})_D = (\{\perp, \top\})_D \\ 1 > x \geqslant 0 \quad \text{if} \quad \frac{\mathfrak{U}}{1} = (\{\perp\})_D \\ 0 > x \geqslant -1 \quad \text{if} \quad 0 = (\emptyset)_D \end{array} \right\} = (\{x \geqslant 0\})_D = (x \geqslant 0)_D = (x)_D$$

The distribution function for X is:

$$P(x = X) = \begin{cases} 0 & \text{otherwise} \\ \frac{1}{2} & \text{if } x = 1 \\ \frac{1}{2} & \text{if } x = 0 \end{cases}$$

$$\left\{ \begin{array}{l} I = x \cdot j \\ 0 = x \cdot j \\ \{1, 0\} \not\ni x \end{array} \right. \left. \begin{array}{l} \frac{\zeta}{j} = (\{H\})_D \\ \frac{\zeta}{j} = (\{\perp\})_D \\ 0 = (\emptyset)_D \end{array} \right\} = (\{x = (\sigma)X : \sigma\})_D = (x = X)_D$$

Answer (Exercise 3.2) — The probability that X takes on a specific value x is:

The second mistake is failing to emphasize that the ratio taken by the function at 0 and 1 is $P(\text{not } A)$, respectively. So it is best to introduce an empty circle like \circ at 0 (0) and $P(\text{not } A)$ to indicate the points of discontinuity. The same mistakes should be fixed in the part [Exercises 3-9](#).

Answer (Exercise 3.1) — The first mistake is the solid vertical lines (blue) from 0 in the domain or x -axis to $P(\text{not } A)$, in the range of y -axis and from 1 in the domain 1 to 1 in the range. This is ill-defined for any function if we are to interpret that the elements in the domain, namely 0 and 1, are to be associated with the inconclusability many image values in the range of the function, namely [0, $P(\text{not } A)$] and $[P(\text{not } A), 1]$, respectively. So we should first replace them by dotted lines which merely help us track where the function jumped to at 0 and 1.

242

```
>> type blah.txt % this allows you to see the contents of blah.txt
3+56
ans = 59
diary off
>> diary blah.txt % reopen the existing diary file blah.txt
>> 45-54
ans = -9
>> diary off % turn off the current diary file blah.txt again
>> type blah.txt % see its contents
3+56
ans = 59
diary off
45-54
ans = -9
diary off
```

4. Let's learn to store values in variables of our choice. Type the following at the command prompt :

```
>> VariableCalledX = 12
```

Upon hitting enter you should see that the number 12 has been assigned to the variable named `VariableCalledX` :

```
VariableCalledX = 12
```

5. MATLAB stores default value for some variables, such as `pi` (π), `i` and `j` (complex numbers).

```
>> pi
ans = 3.1416
>> i
ans = 0 + 1.0000i
>> j
ans = 0 + 1.0000i
```

All predefined symbols (variables, constants, function names, operators, etc) in MATLAB are written in lower-case. Therefore, it is a good practice to name the variable you define using upper and mixed case letters in order to prevent an unintended overwrite of some predefined MATLAB symbol.

6. We could have stored the sum of 13 and 24 in the variable `X`, by entering:

```
>> X = 13 + 24
X = 37
```

7. Similarly, you can store the outcome of multiplication (via operation `*`), subtraction (via operation `-`), division (via `/`) and exponentiation (via `^`) of any two numbers of your choice in a variable name of your choice. Evaluate the following expressions in MATLAB :

$$\begin{aligned} p &= 45.89 * 1.00009 & d &= 89.0 / 23.3454 \\ m &= 5376.0 - 6.00 & p &= 2^{0.5} \end{aligned}$$

8. You may compose the elementary operations to obtain rational expressions by using parenthesis to specify the order of the operations. To obtain $\sqrt{2}$, you can type the following into MATLAB's command window.

Answer (Ex. 2.11) — (a) Let $F1$ be the event a gale of force 1 occurs, let $F2$ be the event a gale of force 2 occurs and $F3$ be the event a gale of force 3 occurs. Now we know that

$$P(F1) = \frac{2}{3}, \quad P(F2) = \frac{1}{4}, \quad P(F3) = \frac{1}{12}.$$

If D is the event that a gale causes damage, then we also know the following conditional probabilities:

$$P(D|F1) = \frac{1}{4}, \quad P(D|F2) = \frac{2}{3}, \quad P(D|F3) = \frac{5}{6}.$$

The probability that a reported gale causes damage is

$$P(D) = P(D \cap F1) + P(D \cap F2) + P(D \cap F3)$$

where

$$P(D \cap F1) = P(D|F1)P(F1) = \frac{1}{4} \times \frac{2}{3} = \frac{1}{6},$$

$$P(D \cap F2) = P(D|F2)P(F2) = \frac{2}{3} \times \frac{1}{4} = \frac{1}{6},$$

and

$$P(D \cap F3) = P(D|F3)P(F3) = \frac{5}{6} \times \frac{1}{12} = \frac{5}{72}.$$

Hence

$$P(D) = \frac{1}{6} + \frac{1}{6} + \frac{5}{72} = \frac{29}{72}$$

(b) Knowing that the gale did cause damage we can calculate the probabilities that it was of the various forces using the probabilities in (a) as follows (Note: $P(D \cap F1) = P(F1 \cap D)$ etc.):

$$P(F1|D) = \frac{P(F1 \cap D)}{P(D)} = \frac{1/6}{29/72} = \frac{12}{29}$$

$$P(F2|D) = \frac{P(F2 \cap D)}{P(D)} = \frac{1/6}{29/72} = \frac{12}{29}$$

$$P(F3|D) = \frac{P(F3 \cap D)}{P(D)} = \frac{5/72}{29/72} = \frac{5}{29}$$

(c) First note that the probability that a reported gale does NOT cause damage is:

$$P(D^c) = 1 - P(D) = 1 - \frac{29}{72} = \frac{43}{72}.$$

Now we need to find probabilities like $P(F1 \cap D^c)$. The best way to do this is to use the partitioning idea of the "Total Probability Theorem", and write:

$$P(F1) = P(F1 \cap D^c) + P(F1 \cap D),$$

Rearranging this gives

$$P(F1 \cap D^c) = P(F1) - P(F1 \cap D)$$

and so

$$P(F1|D^c) = \frac{P(F1 \cap D^c)}{P(D^c)} = \frac{P(F1) - P(F1 \cap D)}{P(D^c)} = \frac{2/3 - 1/6}{43/72} = \frac{36}{43}.$$

Similarly,

$$P(F2|D^c) = \frac{P(F2 \cap D^c)}{P(D^c)} = \frac{P(F2) - P(F2 \cap D)}{P(D^c)} = \frac{1/4 - 1/6}{43/72} = \frac{6}{43},$$

and

$$P(F3|D^c) = \frac{P(F3 \cap D^c)}{P(D^c)} = \frac{P(F3) - P(F3 \cap D)}{P(D^c)} = \frac{1/12 - 5/72}{43/72} = \frac{1}{43}.$$

Answer (Ex. 2.10) — Let the event that a micro-chip is defective be D , and the event that the test is correct be C . So the probability that the micro-chip is defective is $P(D) = 0.05$, and the test is correct be C . The probability that the test correctly detects a defective micro-chip is the conditional probability $P(C|D) = 0.8$, and the probability that the test fails to detect a good micro-chip is the conditional probability $P(C|D^c) = 0.2$, and has been declared as effective is $P(C|D^c) = 0.9$.

Moreover, the probability that a micro-chip is defective, and has been declared as effective is $P(D) = 0.2$, and $P(C|D) = 0.9$. Therefore, we also have the probabilities $P(C|D^c) = 0.1$. Therefore, we also have the probabilities $P(C|D^c) = 0.1$. The probability that the test correctly detects a defective micro-chip is the conditional probability $P(C|D) = 0.8$, and the probability that the test fails to detect a good micro-chip is the conditional probability $P(C|D^c) = 0.2$, and has been declared as effective is $P(C|D^c) = 0.9$.

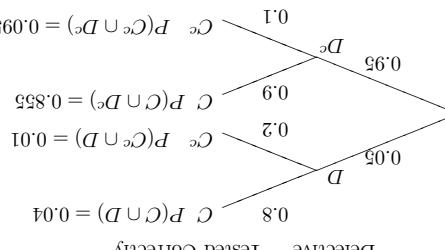
The probability that a micro-chip is effective, and has been declared as effective is $P(C \cup D) = P(C|D)P(D) = 0.2 \times 0.05 = 0.01$.

The probability that a micro-chip is effective, and has been declared as effective is $P(C \cup D^c) = P(C|D^c)P(D^c) = 0.9 \times 0.95 = 0.855$.

The probability that a micro-chip is effective, and has been declared as effective is $P(C^c \cup D) = P(C^c|D)P(D) = 0.8 \times 0.05 = 0.04$.

The probability that a micro-chip is effective, and has been declared as effective is $P(C^c \cup D^c) = P(C^c|D^c)P(D^c) = 0.1 \times 0.95 = 0.095$.

The tree diagram for these events and probabilities is:



10. We can clear the value we have assigned to a particular variable and reuse it. We demonstrate it by the following commands and their output:

```

ans = Inf
>>> 10/0

```

9. When you try to divide by 0, MATLAB returns Inf for infinity.

operations.

MATLAB first takes the last power of 2 and then divides it by 2 using its default precedence rule for binary operators in the absence of parentheses. The order of operations and roots); 3. division and multiplication; 4. addition and subtraction. The innermost bedmas can be handy. When in doubt, use parentheses to force the intended order of operations.

11. We can suppress the output on the screen by ending the command with a semi-colon. Take a look at the simple command that sets X to $\sin(3.145678)$ with and without the ; at the end:

```

>> X = sin(3.145678);
X = -0.0041
>> X = sin(3.145678);
X = 0.0041

```

12. If you do not understand a MATLAB function or command then type help or doc followed by the function or command. For example:

```

>> help sin
SIN SIN of argument in radians.
>> doc sin
SIN(X) is the sine of the elements of X.
See also asin, sind,
darray/sin
Detailed methods:
Referenced page in Help browser
doc sin
>> doc sin

```

(b) Similarly, the probability that a micro-chip is tested to be defective, but it was good is

$$P(C^c \cup D) = \frac{0.01}{0.01 + 0.855} \approx 0.012$$

(a) If a micro-chip is tested to be good, it could be defective but tested incorrectly, or it could be effective and tested correctly. Therefore, the probability that the micro-chip is tested good, but it is actually defective is

$$P(C \cup D) = \frac{0.095}{0.095 + 0.04} \approx 0.704$$

(c) The probability that both the micro-chips are effective, and have been tested and determined to be good, is

$$\left(\frac{P(C \cup D)}{P(C \cup D^c)} \right)^2$$

It is a good idea to use the help files before you ask your tutor.

$$1 - \left(\frac{P(C \cup D)}{P(C \cup D^c)} \right)^2 = 1 - \left(\frac{0.01 + 0.855}{0.855} \right)^2 \approx 0.023$$

and so the probability that at least one is defective is:

13. Set the variable `x` to equal 17.13 and evaluate $\cos(x)$, $\log(x)$, $\exp(x)$, $\arccos(x)$, $\text{abs}(x)$, $\text{sign}(x)$ using the MATLAB commands `cos`, `log`, `exp`, `acos`, `abs`, `sign`, respectively. Read the help files to understand what each function does.
14. When we work with real numbers (floating-point numbers) or really large numbers, we might want the output to be displayed in concise notation. This can be controlled in MATLAB using the `format` command with the `short` or `long` options with/without `e` for scientific notation. `format compact` is used for getting compacted output and `format` returns the default format. For example:

```
>> format compact
>> Y=15;
>> Y = Y + acos(-1)
Y = 18.1416
>> format short
>> Y
Y = 18.1416
>> format short e
>> Y
Y = 1.8142e+001
>> format long
>> Y
Y = 18.141592653589793
>> format long e
>> Y
Y = 1.814159265358979e+001
>> format
>> Y
Y = 18.1416
```

15. Finally, to quit from MATLAB just type `quit` or `exit` at the prompt.

```
>> quit
```

16. An **M-file** is a special text file with a `.m` extension that contains a set of code or instructions in MATLAB . In this course we will be using two types of M-files: **script** and **function** files. A script file is simply a list of commands that we want executed and saves us from retyping code modules we are pleased with. A function file allows us to write specific tasks as functions with input and output. These functions can be called from other script files, function files or command window. We will see such examples shortly.

By now, you are expected to be familiar with arithmetic operations, simple function evaluations, format control, starting and stopping a diary file and launching and quitting MATLAB .

1.6 Elementary Combinatorics

Combinatorics is the branch of mathematics that specialises in counting. We will give a more intuitive treatment with examples and then formally define the most primitive ideas called permutations and combinations. We also use several commonly encountered notations.

The most basic counting rule we use enables us to determine the number of distinct elements in a set that is constructed from taking two or more steps, where each step uses elements of another set. This is a lot easier than it sounds. Let's understand this through the analogy of performing several tasks.

To do this with formulae, partitioning according to the midterm test result and using the multiplication rule, we get:

$$\begin{aligned}\mathbf{P}(A) &= \mathbf{P}(A \cap B) + \mathbf{P}(A \cap B^c) \\ &= \mathbf{P}(A|B)\mathbf{P}(B) + \mathbf{P}(A|B^c)\mathbf{P}(B^c) \\ &= (0.8)(0.7) + (0.4)(0.3) = 0.68\end{aligned}$$

Answer (Ex. 2.9) — Let A be the event that bottles are produced by machine 1; and A^c is the event that bottles are produced by machine 2. R denotes the event that the bottles are rejected; and R^c denotes the event that the bottles are accepted. We know the following probabilities:

$$\begin{aligned}\mathbf{P}(A) &= 0.75 \quad \text{and} \quad \mathbf{P}(A^c) = 0.25 \\ \mathbf{P}(R|A) &= \frac{1}{20} \quad \text{and} \quad \mathbf{P}(R^c|A) = \frac{19}{20} \\ \mathbf{P}(R|A^c) &= \frac{1}{30} \quad \text{and} \quad \mathbf{P}(R^c|A^c) = \frac{29}{30}\end{aligned}$$

We want $\mathbf{P}(A|R^c)$ which is give by

$$\mathbf{P}(A|R^c) = \frac{\mathbf{P}(R^c \cap A)}{\mathbf{P}(R^c)} = \frac{\mathbf{P}(R^c \cap A)}{\mathbf{P}(R^c \cap A) + \mathbf{P}(R^c \cap A^c)}$$

where,

$$\mathbf{P}(R^c \cap A) = \mathbf{P}(R^c|A)\mathbf{P}(A) = \frac{19}{20} \times 0.75$$

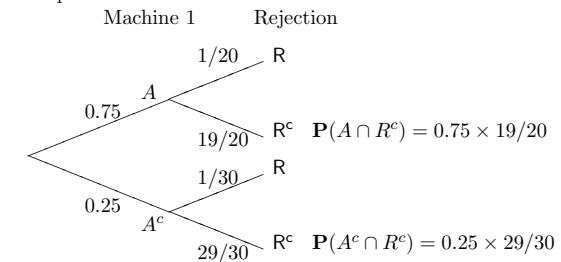
and,

$$\mathbf{P}(R^c \cap A^c) = \mathbf{P}(R^c|A^c)\mathbf{P}(A^c) = \frac{29}{30} \times 0.25$$

Therefore,

$$\mathbf{P}(A|R^c) = \frac{\frac{19}{20} \times 0.75}{\frac{19}{20} \times 0.75 + \frac{29}{30} \times 0.25} \approx 0.747$$

The tree diagram for this problem is:

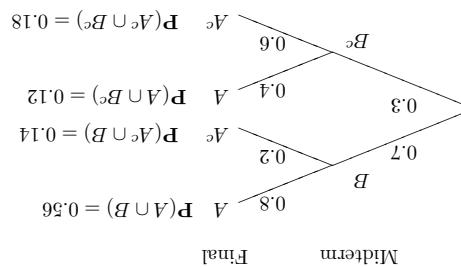


So the required probability is

$$\mathbf{P}(A|R^c) = \frac{\mathbf{P}(R^c \cap A)}{\mathbf{P}(R^c)} = \frac{\frac{19}{20} \times 0.75}{\frac{19}{20} \times 0.75 + \frac{29}{30} \times 0.25} \approx 0.747$$

$$P(A) = 0.56 + 0.12 = 0.68.$$

Then the probability of passing the final exam is:



1. Repetition is allowed, as in choosing the letters (unrestricted choice) in the PIN of Example 6.
More generally, when you have n objects to choose from, you have n choices each time, so when choosing r of them, the number of permutations are n^r .

Permutations: There are basically two types of permutations:

- A selection of objects in which the order is *not* important is called a **combination**.
- A selection of objects in which the order is *important* is called a **permutation**.

So in mathematics, we use more precise language:
I do care about order. A different order gives a different PIN.

"The combination of my PIN is mathagg!"

then I don't care (usually) about what order they are in, but in the statement
"I have 17 probability texts on my bottom shelf"

When does order matter? In English we use the word "combination" loosely. If I say

$$26 \times 25 \times 24 \times 23 \times 9 \times 10 = 32,292,000.$$

Example 7 Suppose we now put restrictions on the letters and digits we use. For example, we might say that the first digit cannot be zero, and letters cannot be repeated. This time the total number of possible PINs is:

$$26 \times 26 \times 26 \times 26 \times 10 \times 10 = 26^4 \times 10^2 = 45,697,600.$$

So in total, the total number of possible PINs is:
First letter: 26 possibilities
Second letter: 26 possibilities
Third letter: 26 possibilities
Fourth letter: 26 possibilities
Second digit: 10 possibilities
First digit: 10 possibilities

PINs are there? There are six selections to be made:
which the first four entries are letters (lowercase) and the last two entries are digits. How many PINs are there?

	1	2	3	4	5	6
1	2	3	4	5	6	7
2	3	4	5	6	7	8
3	4	5	6	7	8	9
4	5	6	7	8	9	10
5	6	7	8	9	10	11
6	7	8	9	10	11	12

The multiplication principle: If a task can be performed in n_1 ways, a second task in n_2 ways, a third task in n_3 ways, etc., then the total number of distinct ways of performing all tasks together is $n_1 \times n_2 \times n_3 \times \dots$

Answer (Ex. 2.7) — 1. The sample space is

be the event that the student passes the final exam and let B be the event that the student passes the second branch of the mid-term test. Note that the probabilities involved in this test and the second branch of the final exam. Note that the outcome of the mid-term test is the outcome of the final exam that depends on the outcome of the mid-term test. Let A be the event that the student passes the final exam and let B be the event that the student passes the mid-term test.

$$P(A \cup B) = P(A) + P(B) = \frac{36}{4} + \frac{5}{1} = \frac{36}{4}$$

Therefore,

$$P(A) = \frac{36}{4} \quad \text{and} \quad P(B) = \frac{36}{5}.$$

Let A be the event the sum is 5 and B be the event the sum is 6, then A and B are mutually exclusive events with probabilities

6	7	8	9	10	11	12
5	6	7	8	9	10	11
4	5	6	7	8	9	10
3	4	5	6	7	8	9
2	3	4	5	6	7	8
1	2	3	4	5	6	7

2. First tabulate all possible sums as follows:

Note: Order matters here. For example, the outcome "16" refers to a "1" on the first die and a "6" on the second, whereas the outcome "61" refers to a "6" on the first die and a "1" on a "6" on the second.

$$\{(5, 1), (5, 2), (5, 3), (5, 4), (5, 5), (5, 6), (6, 1), (6, 2), (6, 3), (6, 4), (6, 5), (6, 6), (3, 1), (3, 2), (3, 3), (3, 4), (3, 5), (3, 6), (4, 1), (4, 2), (4, 3), (4, 4), (4, 5), (4, 6), (1, 1), (1, 2), (1, 3), (1, 4), (1, 5), (1, 6), (2, 1), (2, 2), (2, 3), (2, 4), (2, 5), (2, 6)\}$$

2. No repetition is allowed, as in the restricted PIN Example 7. Here you have to reduce the number of choices. If we had a 26 letter PIN then the total permutations would be

$$26 \times 25 \times 24 \times 23 \times \dots \times 3 \times 2 \times 1 = 26!$$

but since we want four letters only here, we have

$$\frac{26!}{22!} = 26 \times 25 \times 24 \times 23$$

choices.

The number of distinct **permutations** of n objects taking r at a time is given by

$${}^n P_r = \frac{n!}{(n-r)!}$$

Combinations: There are also two types of combinations:

1. Repetition is allowed such as the coins in your pocket, say, (10c, 50c, 50c, \$1, \$2, \$2).
2. No repetition is allowed as in the lottery numbers (2, 9, 11, 26, 29, 31). The numbers are drawn one at a time, and if you have the lucky numbers (no matter what order) you win!

The number of distinct **combinations** of n objects taking r at a time is given by

$${}^n C_r = \binom{n}{r} = \frac{n!}{(n-r)!r!}$$

Example 8 Let us imagine being in the lower Manhattan in New York city with its perpendicular grid of streets and avenues. If you start at a given intersection and are asked to only proceed in a north-easterly direction then how many ways are there to reach another intersection by walking exactly two blocks or exactly three blocks?

Solution:

Let us answer this question of combinations by drawing Fig. 1.6. Let us denote the number of easterly turns you take by r and the total number of blocks you are allowed to walk either easterly or northerly by n . From Fig. 1.6(a) it is clear that the number of ways to reach each of the three intersections labeled by r is given by $\binom{n}{r}$, with $n = 2$ and $r \in \{0, 1, 2\}$. Similarly, from Fig. 1.6(b) it is clear that the number of ways to reach each of the four intersections labeled by r is given by $\binom{n}{r}$, with $n = 3$ and $r \in \{0, 1, 2, 3\}$.

Exercise 1.5 (Choosing Volunteers) Suppose we need three students to be the class representatives in this course. Assume that everyone wants to be selected to keep it simple. In how many ways can we choose these three people from the class of 50 students?

Now, we give more formal definitions and notations that will help us make precise arguments faster when we study sampling schemes in Inference Theory.

2. Since there are eleven letters in WAIMAKARIRI the probabilities are:

$$P(\{W\}) = \frac{1}{11}, P(\{A\}) = \frac{3}{11}, P(\{I\}) = \frac{3}{11}, P(\{M\}) = \frac{1}{11}, P(\{K\}) = \frac{1}{11}, P(\{R\}) = \frac{2}{11}.$$

3. By the complementation rule, the probability of not choosing the letter R is:

$$1 - P(\text{choosing the letter R}) = 1 - \frac{2}{11} = \frac{9}{11}.$$

Answer (Ex. 2.5) — 1. First, the sample space is: $\Omega = \{B, I, N, G, O\}$.

2. The probabilities of simple events are:

$$P(B) = P(I) = P(N) = P(G) = P(O) = \frac{15}{75} = \frac{1}{5}.$$

3. Using the addition rule for mutually exclusive events,

$$\begin{aligned} P(\Omega) &= P(\{B, I, N, G, O\}) \\ &= P(\{B\} \cup \{I\} \cup \{N\} \cup \{G\} \cup \{O\}) \\ &= P(B) + P(I) + P(N) + P(G) + P(O) \quad \text{simplifying notation} \\ &= \frac{1}{5} + \frac{1}{5} + \frac{1}{5} + \frac{1}{5} + \frac{1}{5} \\ &= 1 \end{aligned}$$

4. Since the events $\{B\}$ and $\{I\}$ are disjoint,

$$P(\{B\} \cup \{I\}) = P(B) + P(I) = \frac{1}{5} + \frac{1}{5} = \frac{2}{5}.$$

5. Using the addition rule for two arbitrary events we get,

$$\begin{aligned} P(C \cup D) &= P(C) + P(D) - P(C \cap D) \\ &= P(\{B, I, G\}) + P(\{G, I, N\}) - P(\{G, I\}) \\ &= \frac{3}{5} + \frac{3}{5} - \frac{2}{5} \\ &= \frac{4}{5}. \end{aligned}$$

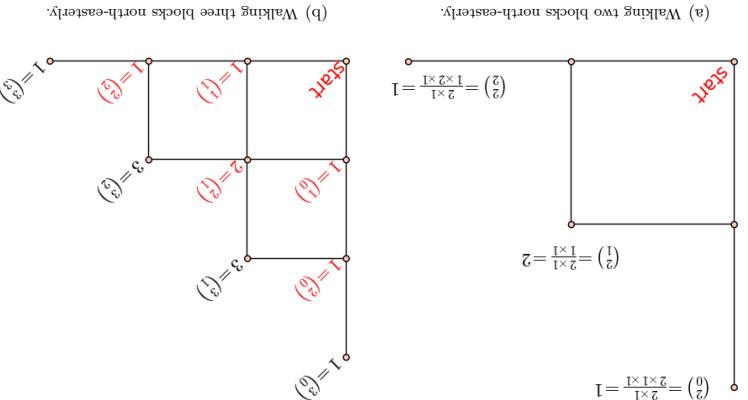
Answer (Ex. 2.6) — We can assume that the first shot is independent of the second shot so we can multiply the probabilities here.

For case A, there is only one shot so the probability of hitting at least once is $\frac{1}{2}$.

For case B, the probability of missing both shots is $\frac{2}{3} \cdot \frac{2}{3} = \frac{4}{9}$, so the probability hitting some target at least once is

$$1 - P(\text{missing the target both times}) = 1 - \frac{4}{9} = \frac{5}{9}.$$

Therefore, case B has the greater probability of hitting the target at least once.



Definition 1 (Permutations and Factorials) A **permutation** of n objects is an arrangement of n distinct objects in a row. For example, there are 2 permutations of the two objects {1, 2}:

Let the number of ways to choose k objects out of n and to arrange them in a row be denoted by $P_{n,k}$. For example, we can choose two ($k = 2$) objects out of three ($n = 3$) objects, {a, b, c}, and arrange them in a row in six ways (p.2):

$$abc, \quad acb, \quad bac, \quad bca, \quad cab, \quad cba.$$

and 6 permutations of the three objects {a, b, c}:

$$12, \quad 21,$$

Given n objects, there are n ways to choose the left-most object, and once this choice has been

$$ab, \quad ac, \quad ba, \quad bc, \quad ca, \quad cb.$$

made there are $n - 1$ ways to select a different object to place next to the left-most one. Thus, there are $n(n - 1)$ possible choices for the first two positions. Similarly, when $n > 2$, there are $n - 2$ choices for the third object that is distinct from the first two. Thus, there are $n(n - 1)(n - 2)$ possible ways to choose three distinct objects from a set of n objects and arrange them in a row. In general,

and the total number of permutations called **$n!$ factorial** and denoted by $n!$ is

$$p_{n,k} = n(n - 1)(n - 2) \cdots (n - k + 1)$$

In addition rule for two arbitrary events, rule (3).

Some factorials to bear in mind

$$n! := p_{n,n} = n(n - 1)(n - 2) \cdots (n - 2)(n - 1)(2)(1) = \prod_{i=1}^n i.$$

When n is large we can get a good idea of $n!$ without laboriously carrying out the $n - 1$ multiplications via Stirling's approximation (*Methods of Approximation* (1930), p. 133):

$$0! = 1, \quad 1! = 1, \quad 2! = 2, \quad 3! = 6, \quad 4! = 24, \quad 5! = 120, \quad 10! = 3,628,800.$$

Answer (Ex. 2.4) — 1. The sample space $\Omega = \{W, A, I, M, K, R\}$.

3. {6, 26, 36, 46, 56, 116, 126, 136, 146, 156, 216, 226, 236, 246, 256, ...}

2. LRLR, RLRL, RRLR, LRRL, RLRL, RRLR, RRRL, RRLR, RLRL, RLLR,

Answer (Ex. 2.3) — 1. {BB, BW, WB, WW}

rule for mutually exclusive events, rule (2).

(e) $P(\text{picking any letter in the word WAZZZUP}) = P(W, A, Z, U, P) = 14.4\%$, by the addition rule for two arbitrary events, rule (3).

(f) $P(\text{picking any letter in the word WAZZZUP or a vowel}) = P(W, A, Z, U, P) + P(A, E, I, O, U) = 14.4\% + 37.8\% - 10\% = 42.2\%$, by the addition rule for two arbitrary events, rule (3).

(g) $P(\text{picking any letter in the word WAZZZUP}) = P(W, A, Z, U, P) = 14.4\%$, by the addition rule for mutually exclusive events, rule (2).

(h) $P(\text{picking any letter in the word WAZZZUP or a vowel}) = P(W, A, Z, U, P) + P(A, E, I, O, U) = (7.3\% + 13.0\%) + 7.4\% + 2.7\% = 37.8\%$, by the addition rule for two arbitrary events, rule (3).

(i) $P(E, Z) = P(\{E\} \cap \{Z\}) = P(\{E\}) + P(\{Z\}) = 0.13 + 0.001 = 0.131$, by Axiom (3)

(j) $P(\text{picking any letter}) = P(\Omega) = 1$

Answer (Ex. 2.2) — (a) $P(Z) = 0.1\% = \frac{1}{100} = 0.001$

Answer (Exercise 2.1) — This is an optional exercise. You will understand this as you progress through your mathematics programme. The explanation in the said item was (or will be explained in person again) in the lectures. This exercise was created to answer natural questions that were asked by students who wanted to know.

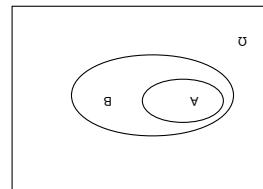
Placed in order $3!$ ways, so the required number of ways of choosing the class representatives is:

factor representing the number of ways the objects could be in order. Here, three students can be placed in order $3!$ ways, so the required number of ways of choosing the class representatives is:

But, because order doesn't matter, all we have to do is to adjust our permutation formula by a factor representing the number of ways the objects could be in order. Here, three students can be placed in order $3!$ ways, so the required number of ways of choosing the class representatives is:

$${}_{50}P_3 = \frac{(50-3)!}{50!} = \frac{47!}{50!} = \frac{3 \cdot 2 \cdot 1}{50 \cdot 49 \cdot 48} = 19,600$$

permutation, so that the number of ways we can select the three class representatives is



Definition 2 (Combinations) The combinations of n objects taken k at a time are the possible choices of k different elements from a collection of n objects, disregarding order. They are called the k -combinations of the collection. The combinations of the three objects $\{a, b, c\}$ taken two at a time, called the 2-combinations of $\{a, b, c\}$, are

$$ab, \quad ac, \quad bc,$$

and the combinations of the five objects $\{1, 2, 3, 4, 5\}$ taken three at a time, called the 3-combinations of $\{1, 2, 3, 4, 5\}$ are

$$123, \quad 124, \quad 125, \quad 134, \quad 135, \quad 145, \quad 234, \quad 235, \quad 245, \quad 345.$$

The total number of k -combination of n objects, called a **binomial coefficient**, denoted $\binom{n}{k}$ and read “ n choose k ,” can be obtained from $p_{n,k} = n(n-1)(n-2)\dots(n-k+1)$ and $k! := p_{k,k}$. Recall that $p_{n,k}$ is the number of ways to choose the first k objects from the set of n objects and arrange them in a row with regard to order. Since we want to disregard order and each k -combination appears exactly $p_{k,k}$ or $k!$ times among the $p_{n,k}$ many permutations, we perform a division:

$$\binom{n}{k} := \frac{p_{n,k}}{p_{k,k}} = \frac{n(n-1)(n-2)\dots(n-k+1)}{k(k-1)(k-2)\dots2\ 1}.$$

Binomial coefficients are often called “Pascal’s Triangle” and attributed to Blaise Pascal’s *Traité du Triangle Arithmétique* from 1653, but they have many “fathers”. There are earlier treatises of the binomial coefficients including Szu-yüan Yü-chien (“The Precious Mirror of the Four Elements”) by the Chinese mathematician Chu Shih-Chieh in 1303, and in an ancient Hindu classic, *Piṅgala’s Chandadhāśṭra*, due to Halāyudha (10-th century AD).

1.7 Array, Sequence, Limit, ...

In this section we will study a basic data structure in MATLAB called an **array** of numbers. Arrays are finite sequences and they can be processed easily in MATLAB. The notion of infinite sequences lead to **limits**, one of the most fundamental concepts in mathematics.

For any natural number n , we write

$$\langle x_{1:n} \rangle := x_1, x_2, \dots, x_{n-1}, x_n$$

to represent the **finite sequence** of real numbers $x_1, x_2, \dots, x_{n-1}, x_n$. For two integers m and n such that $m \leq n$, we write

$$\langle x_{m:n} \rangle := x_m, x_{m+1}, \dots, x_{n-1}, x_n$$

to represent the **finite sequence** of real numbers $x_m, x_{m+1}, \dots, x_{n-1}, x_n$. In mathematical analysis, finite sequences and their countably infinite counterparts play a fundamental role in limiting processes. Given an integer m , we denote an **infinite sequence** or simply a sequence as:

$$\langle x_{m:\infty} \rangle := x_m, x_{m+1}, x_{m+2}, x_{m+3}, \dots$$

Given index set \mathcal{I} which may be finite or infinite in size, a sequence can either be seen as a set of ordered pairs:

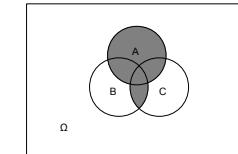
$$\{(i, x_i) : i \in \mathcal{I}\},$$

Answers to Selected Exercises

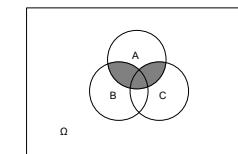
Answer (Ex. 1.1) — By operating with Ω , T , L and S we can obtain the answers as follows:

- | | |
|---|--|
| (a) $T \cap L = \{L_3\}$ | (f) $S \cap L = \emptyset$ |
| (b) $T \cap S = \emptyset$ | (g) $S^c \cap L = \{L_1, L_2, L_3\} = L$ |
| (c) $T \cup L = \{T_1, T_2, T_3, L_3, L_1, L_2\}$ | (h) $T^c = \{L_1, L_2, S_1, S_2, S_3, \dots, S_{50}\}$ |
| (d) $T \cup L \cup S = \Omega$ | (i) $T^c \cap L = \{L_1, L_2\}$ |
| (e) $S^c = \{T_1, T_2, T_3, L_3, L_1, L_2\}$ | (j) $T^c \cap T = \emptyset$ |

Answer (Ex. 1.3) — We can check $A \cup (B \cap C) = (A \cup B) \cap (A \cup C)$ from the following sketch:



Answer (Ex. 1.3) — We can check $A \cap (B \cup C) = (A \cap B) \cup (A \cap C)$ from the following sketch:
We can check $A \cup (B \cap C) = (A \cup B) \cap (A \cup C)$ from the following sketch:



Answer (Ex. 1.4) — To illustrate the idea that $A \subseteq B$ if and only if $A \cup B = B$, we need to illustrate two implications:

- 1.if $A \subseteq B$ then $A \cup B = B$ and
- 2.if $A \cup B = B$ then $A \subseteq B$.

The following Venn diagram illustrates the two implications clearly.

or as a function that maps the index set to the set of real numbers:

$$x(i) = x_i : i \in \mathcal{I},$$

$$(x_j:k) = x_j, x_{j+1}, \dots, x_{k-1}, x_k \quad \text{where,} \quad m \leq j \leq k < \infty.$$

A rectangular arrangement of $m \times n$ real numbers in m rows and n columns is called a **matrix**. The, $m \times n$, represents the **size** of the matrix. We use bold upper-case letters to denote matrices, for e.g.:

$$\mathbf{X} = \begin{bmatrix} x_{1,1} & x_{1,2} & \dots & x_{1,n-1} & x_{1,n} \\ x_{2,1} & x_{2,2} & \dots & x_{2,n-1} & x_{2,n} \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ x_{m-1,1} & x_{m-1,2} & \dots & x_{m-1,n-1} & x_{m-1,n} \\ x_{m,n} \end{bmatrix}$$

Matrices with only one row or only one column are called **vectors**. An $1 \times n$ matrix is called a **row vector** since there is only one row and an $m \times 1$ matrix is called a **column vector** since there is only one column. We use bold-face lowercase letters to denote row and column vectors.

$$\text{A row vector } \mathbf{x} = [x_1 \ x_2 \ \dots \ x_n] = (x_1, x_2, \dots, x_n)$$

$$\text{and a column vector } \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix} = [y_1 \ y_2 \ \dots \ y_m].$$

The superscripting by $'$ is the transpose operation and simply means that the rows and columns are exchanged. Thus the transpose of the matrix \mathbf{X} is:

$$\mathbf{X}' = \begin{bmatrix} x_{1,1} & x_{2,1} & \dots & x_{m-1,1} & x_{m,n} \\ x_{1,2} & x_{2,2} & \dots & x_{m-1,2} & x_{m,2} \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ x_{1,m} & x_{2,m} & \dots & x_{m-1,m} & x_{m,m} \end{bmatrix}$$

In linear algebra and calculus, it is natural to think of vectors and matrices as points (ordered m -tuples) and ordered collection of points in Cartesian co-ordinates. We assume that the reader has heard of operations with matrices and vectors such as matrix multiplication, determinants, transposes, etc. Such concepts will be introduced as they are needed in the sequel.

Labwork 9 (Sequences as arrays) Let us learn to represent, visualize and operate finite sequences as MATLAB arrays. Try out the commands and read the comments for clarification.

Table 6.4: Symbol Table: Probability and Statistics

\mathbb{R}^d	$\mathbb{R}^d := (-\infty, \infty)^d$	$\mathbb{A}(x)$	$\mathbb{I}(A)$	$\mathbb{S}_{X,Y}$	$f_{X,Y}(x,y)$	$E(X,Y)$	$E(X)$	$E(X_s)$	$E(X,Y_s)$	$E(g(X,Y))$	$E(g(X))$	$F_{X,Y}(x,y)$	$Cov(X,Y) = E(XY) - E(X)E(Y)$
\mathbb{R}^d	\mathbb{R}^d	Measuring	$\text{Indicates or set membership function that returns 1 if } x \in A \text{ and 0 otherwise}$	random vector	$\text{Joint distribution function (JDF) of the RV } (X, Y)$	$\text{Joint cumulative distribution function (CDF) of the RV } (X, Y)$	$\text{Joint probability mass function (PMF) of the discrete RV } (X, Y)$	$\text{Marginal probability density/mass function (MPDF/MPDF) of } X$	$\text{Expectation of a function } g(x, y) \text{ for continuous RV } Y$	$\text{Expectation of a function } g(x, y) \text{ for discrete RV } Y$	Joint moment	$\text{Covariance of } X \text{ and } Y, \text{ provided } E(X^2) > \infty \text{ and } E(Y^2) > \infty$	
\mathbb{R}^d	\mathbb{R}^d	Meaning	$\text{for every } (x_1, x_2, \dots, x_n)$	$\text{for every } (x, y)$	$f_{X,Y}(x,y)$	$F_{X,Y}(x,y)$	$E(X,Y)$	$E(g(X,Y))$	$E(g(X,Y))$	$E(g(X))$	$E(g(X))$	$F_{X,Y}(x,y)$	$\text{for every } (x, y)$
\mathbb{R}^d	\mathbb{R}^d	$\text{d-dimensional Real Space}$	$= \prod_{i=1}^n f_{X_i}(x_i)$	$f_{X,Y}(x,y)$	$f_{X,Y}(x,y)$	$F_{X,Y}(x,y)$	$E(X)$	$E(g(X))$	$E(g(X))$	$E(g(X))$	$E(g(X))$	$F_{X,Y}(x,y)$	$\text{for every } (x, y)$

The finite sequence $(x_{m:n})$ has $\mathcal{I} = \{m, m+1, m+2, m+3, \dots, n\}$ as its index set while an infinite sequence $(x_{m:\infty})$ has $\mathcal{I} = \{m, m+1, m+2, m+3, \dots\}$ as its index set. A **sub-sequence** $(x_{j:k})$ is:

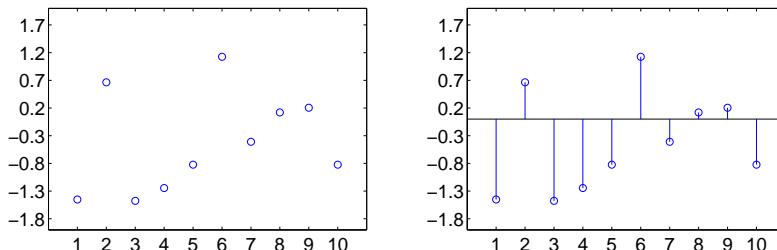
```

>> a = [17] % Declare the sequence of one element 17 in array a
a = 17
>> % Declare the sequence of 10 numbers in array b
>> b=[-1.4508 0.6636 -1.4768 -1.2455 -0.8235 1.1254 -0.4093 0.1199 0.2043 -0.8236]
b =
-1.4508 0.6636 -1.4768 -1.2455 -0.8235 1.1254 -0.4093 0.1199 0.2043 -0.8236
>> c = [1 2 3] % Declare the sequence of 3 consecutive numbers 1,2,3
c =
1 2 3
>> % linspace(x1, x2, n) generates n points linearly spaced between x1 and x2
>> r = linspace(1, 3, 3) % Declare sequence r = c using linspace
r =
1 2 3
>> s1 = 1:10 % declare an array s1 starting at 1, ending by 10, in increments of 1
s =
1 2 3 4 5 6 7 8 9 10
>> s2 = 1:2:10 % declare an array s2 starting at 1, ending by 10, in increments of 2
s =
1 3 5 7 9
>> s2(3) % obtain the third element of the finite sequence s2
ans =
5
>> s2(2:4) % obtain the subsequence from second to fourth elements of the finite sequence s2
ans =
3 5 7

```

We may visualise (as per Figure 1.6) the finite sequences $\langle b_{1:n} \rangle$ stored in the array b as the set of ordered pairs $\{(1, b_1), (2, b_2), \dots, (10, b_{10})\}$ representing the function $b(i) = b_i : \{1, 2, \dots, n\} \rightarrow \{b_1, b_2, \dots, b_n\}$ via **point plot** and **stem plot** using Matlab's `plot` and `stem` commands, respectively.

Figure 1.6: Point plot and stem plot of the finite sequence $\langle b_{1:10} \rangle$ declared as an array.



```

>> display(b) % display the array b in memory
b =
-1.4508 0.6636 -1.4768 -1.2455 -0.8235 1.1254 -0.4093 0.1199 0.2043 -0.8236
>> plot(b,'o') % point plot of ordered pairs (1,b(1)), (2,b(2)), ..., (10,b(10))
>> stem(b) % stem plot of ordered pairs (1,b(1)), (2,b(2)), ..., (10,b(10))
>> plot(b,'-o') % point plot of ordered pairs (1,b(1)), (2,b(2)), ..., (10,b(10)) connected by lines

```

Labwork 10 (Vectors and matrices as arrays) Let us learn to represent, visualise and operate vectors as MATLAB arrays. Syntactically, a vector is stored in an array exactly in the same way we stored a finite sequence. However, mathematically, we think of a vector as an ordered m -tuple that can be visualised as a point in Cartesian co-ordinates. Try out the commands and read the comments for clarification.

```

>> a = [1 2] % an 1 X 2 row vector
>> z = [1 2 3] % Declare an 1 X 3 row vector z with three numbers

```

Model	PDF or PMF	Mean	Variance
Bernoulli(θ)	$\theta^x(1-\theta)^{1-x}\mathbb{1}_{\{0,1\}}(x)$	θ	$\theta(1-\theta)$
Binomial(n, θ)	$\binom{n}{\theta}\theta^x(1-\theta)^{n-x}\mathbb{1}_{\{0,1,\dots,n\}}(x)$	$n\theta$	$n\theta(1-\theta)$
Geometric(θ)	$\theta(1-\theta)^{x-1}\mathbb{1}_{\mathbb{Z}_+}(x)$	$\frac{1}{\theta}-1$	$\frac{1-\theta}{\theta^2}$
Poisson(λ)	$\frac{\lambda^x e^{-\lambda}}{x!}\mathbb{1}_{\mathbb{Z}_+}(x)$	λ	λ
Uniform(θ_1, θ_2)	$\mathbb{1}_{[\theta_1,\theta_2]}(x)/(\theta_2 - \theta_1)$	$\frac{\theta_1+\theta_2}{2}$	$\frac{(\theta_2-\theta_1)^2}{12}$
Exponential(λ)	$\lambda e^{-\lambda x}$	λ^{-1}	λ^{-2}
Normal(μ, σ^2)	$\frac{1}{\sigma\sqrt{2\pi}}e^{-(x-\mu)^2/(2\sigma^2)}$	μ	σ^2

Table 6.2: Random Variables with PDF and PMF (using indicator function), Mean and Variance

Symbol	Meaning
$A = \{\star, \circ, \bullet\}$	A is a set containing the elements \star, \circ and \bullet
$\circ \in A$	\circ belongs to A or \circ is an element of A
$A \ni \circ$	\circ belongs to A or \circ is an element of A
$\circ \notin A$	\circ does not belong to A
$\#A$	Size of the set A , for e.g. $\#\{\star, \circ, \bullet, \circ\} = 4$
\mathbb{N}	The set of natural numbers $\{1, 2, 3, \dots\}$
\mathbb{Z}	The set of integers $\{\dots, -3, -2, -1, 0, 1, 2, 3, \dots\}$
\mathbb{Z}_+	The set of non-negative integers $\{0, 1, 2, 3, \dots\}$
\emptyset	Empty set or the collection of nothing or $\{\}$
$A \subset B$	A is a subset of B or A is contained by B , e.g. $A = \{\circ\}, B = \{\bullet\}$
$A \supset B$	A is a superset of B or A contains B e.g. $A = \{\circ, \star, \bullet\}, B = \{\circ, \bullet\}$
$A = B$	A equals B , i.e. $A \subset B$ and $B \subset A$
$Q \implies R$	Statement Q implies statement R or If Q then R
$Q \iff R$	$Q \implies R$ and $R \implies Q$
$\{x : x \text{ satisfies property } R\}$	The set of all x such that x satisfies property R
$A \cup B$	A union B , i.e. $\{x : x \in A \text{ or } x \in B\}$
$A \cap B$	A intersection B , i.e. $\{x : x \in A \text{ and } x \in B\}$
$A \setminus B$	A minus B , i.e. $\{x : x \in A \text{ and } x \notin B\}$
$A := B$	A is equal to B by definition
$A :=: B$	B is equal to A by definition
A^c	A complement, i.e. $\{x : x \in U, \text{ the universal set, but } x \notin A\}$
$A_1 \times A_2 \times \dots \times A_m$	The m -product set $\{(a_1, a_2, \dots, a_m) : a_1 \in A_1, a_2 \in A_2, \dots, a_m \in A_m\}$
A^m	The m -product set $\{(a_1, a_2, \dots, a_m) : a_1 \in A, a_2 \in A, \dots, a_m \in A\}$
$f := f(x) = y : \mathbb{X} \rightarrow \mathbb{Y}$	A function f from domain \mathbb{X} to range \mathbb{Y}
$f^{[-1]}(y)$	Inverse image of y
$f^{[-1]} := f^{[-1]}(y) \in \mathbb{Y} = X \subset \mathbb{X}$	Inverse of f
$a < b$ or $a \leq b$	a is less than b or a is less than or equal to b
$a > b$ or $a \geq b$	a is greater than b or a is greater than or equal to b
\mathbb{Q}	Rational numbers
(x, y)	the open interval (x, y) , i.e. $\{r : x < r < y\}$
$[x, y]$	the closed interval $[x, y]$, i.e. $\{r : x \leq r \leq y\}$
$(x, y]$	the half-open interval $(x, y]$, i.e. $\{r : x < r \leq y\}$
$[x, y)$	the half-open interval $[x, y)$, i.e. $\{r : x \leq r < y\}$
$\mathbb{R} := (-\infty, \infty)$	Real numbers, i.e. $\{r : -\infty < r < \infty\}$
$\mathbb{R}_+ := [0, \infty)$	Real numbers, i.e. $\{r : 0 \leq r < \infty\}$
$\mathbb{R}_{>0} := (0, \infty)$	Real numbers, i.e. $\{r : 0 < r < \infty\}$

Table 6.3: Symbol Table: Sets and Numbers

```

z = 1 2 3
<>> x = linspace(x1, x2, n) % generates n points linearly spaced between x1 and x2
<>> c = [1; 2; 3] % Declare a 3 X 1 column vector c with three numbers. Semicolons delineate columns
c =
    1
    2
    3
x =
    1 2 3
<>> z = linspace(1, 3, 3) % Declares an 1 X 3 row vector x = z using linspace
z =
    1 2 3
<>> C = [1; 2; 3] % The column vector (1,2,3), by taking the transpose of x via x'
C =
    1
    2
    3
xT =
    1
    2
    3
<>> xT = x', % The column vector (1,2,3), by taking the transpose of x via x',
y =
    1 1 1
<>> ans = ones(1,10) % ones(m,n) is an m X n matrix of ones. Useful when m or n is Large.
ans =
    1 1 1 1 1 1 1 1 1 1
<>> Z = zeros(2,10) % the 2 X 10 matrix of zeros
Z =
    0 0 0 0 0 0 0 0 0 0
<>> D=ones(4,5) % the 4 X 5 matrix of ones
D =
    1 1 1 1 1
    1 1 1 1 1
    1 1 1 1 1
    1 1 1 1 1
<>> E=eye(4) % the 4 X 4 identity matrix
E =
    0 0 0 0
    0 0 1 0
    0 1 0 0
    1 0 0 0
<>> x = 1 2 3
<>> y = x + z
x =
    2 3 4
y =
    2 3 4
<>> p = z * y
p =
    2 4 6
<>> d = z ./ y
d =
    2 4 6
<>> a = 0.5000 1.0000 1.5000
a =
    0.5000 1.0000 1.5000
<>> t = -10.0000 -5.3333 3.3333 10.0000
t = -10.0000 -5.3333 3.3333 10.0000
<>> e=linspace(-10,10,4)
e=linspace(-10,10,4)
<>> s = sin(t)
s =
    0.5440 0.1906 -0.1906 -0.5440
<>> s2 = sin(t)
s2 =
    0.2960 0.0363 -0.0363 -0.2960
<>> s3 = cos(t)
s3 =
    0.7040 0.9637 0.9637 0.7040
<>> c3 = cos(t)
c3 =
    0.2 % c3 is an array obtained from term-wise squaring ( .^ 2 ) of the sin(t) array
<>> c3 = cos(t)
c3 =
    0.2960 0.0363 0.0363 -0.2960
% s3 is a vector obtained from the term-wise sin of the vector t
% s2 is a vector obtained by term-by-term division of z and y
% d is the vector obtained by term-by-term product of z and y
% p is the vector obtained by term-by-term product of z and y
% y is updated to 2 * y (each term of y is multiplied by 2)
% x is the sum of vectors y and z (with same size 1 X 3)
% x is also known as

```

We can also perform operations with arrays representing vectors, finite sequences, or matrices.

```

<>> y % the array y is
<>> z % the array z is
<>> x = y + z
x =
    1 1 1
    1 1 1
    1 1 1
<>> E=eye(4) % the 4 X 4 identity matrix
E =
    0 0 0 0
    0 0 1 0
    0 1 0 0
    1 0 0 0
<>> D=zeros(4,5) % the 4 X 5 matrix of ones
D =
    1 1 1 1 1
    1 1 1 1 1
    1 1 1 1 1
    1 1 1 1 1
<>> Z=zeros(2,10) % the 2 X 10 matrix of zeros
Z =
    0 0 0 0 0 0 0 0 0 0
<>> D=ones(4,5) % the 4 X 5 matrix of ones
D =
    1 1 1 1 1
    1 1 1 1 1
    1 1 1 1 1
    1 1 1 1 1
<>> E=eye(4) % the 4 X 4 identity matrix
E =
    0 0 0 0
    0 0 1 0
    0 1 0 0
    1 0 0 0
<>> x = 1 2 3
<>> y = x + z
x =
    2 3 4
y =
    2 3 4
<>> p = z * y
p =
    2 4 6
<>> d = z ./ y
d =
    2 4 6
<>> a = 0.5000 1.0000 1.5000
a =
    0.5000 1.0000 1.5000
<>> t = -10.0000 -5.3333 3.3333 10.0000
t = -10.0000 -5.3333 3.3333 10.0000
<>> e=linspace(-10,10,4)
e=linspace(-10,10,4)
<>> s = sin(t)
s =
    0.5440 0.1906 -0.1906 -0.5440
<>> s2 = sin(t)
s2 =
    0.2960 0.0363 -0.0363 -0.2960
<>> s3 = cos(t)
s3 =
    0.7040 0.9637 0.9637 0.7040
<>> c3 = cos(t)
c3 =
    0.2 % c3 is an array obtained from term-wise squaring ( .^ 2 ) of the sin(t) array
<>> c3 = cos(t)
c3 =
    0.2960 0.0363 0.0363 -0.2960
% s3 is a vector obtained from the term-wise sin of the vector t
% s2 is a vector obtained by term-by-term division of z and y
% d is the vector obtained by term-by-term product of z and y
% p is the vector obtained by term-by-term product of z and y
% y is updated to 2 * y (each term of y is multiplied by 2)
% x is the sum of vectors y and z (with same size 1 X 3)
% x is also known as

```

We can use two dimensional arrays to represent matrices. Some useful built-in commands to generate standard matrices are:

```

<>> ans = ones(1,10) % ones(m,n) is an m X n matrix of ones. Useful when m or n is Large.
ans =
    1 1 1 1 1 1 1 1 1 1
<>> y = [1 1 1]
y =
    1 1 1
<>> xT = x', % The column vector (1,2,3), by taking the transpose of x via x',
xT =
    1
    2
    3
<>> x = 1 2 3
x =
    1 2 3
<>> c = [1; 2; 3] % Declare a 3 X 1 column vector c with three numbers. Semicolons delineate columns
c =
    1
    2
    3
<>> z = linspace(1, 3, 3) % Declares an 1 X 3 row vector x = z using linspace
z =
    1 2 3
<>> D=ones(4,5) % the 4 X 5 matrix of ones
D =
    1 1 1 1 1
    1 1 1 1 1
    1 1 1 1 1
    1 1 1 1 1
<>> E=eye(4) % the 4 X 4 identity matrix
E =
    0 0 0 0
    0 0 1 0
    0 1 0 0
    1 0 0 0
<>> x = 1 2 3
<>> y = x + z
x =
    2 3 4
y =
    2 3 4
<>> p = z * y
p =
    2 4 6
<>> d = z ./ y
d =
    2 4 6
<>> a = 0.5000 1.0000 1.5000
a =
    0.5000 1.0000 1.5000
<>> t = -10.0000 -5.3333 3.3333 10.0000
t = -10.0000 -5.3333 3.3333 10.0000
<>> e=linspace(-10,10,4)
e=linspace(-10,10,4)
<>> s = sin(t)
s =
    0.5440 0.1906 -0.1906 -0.5440
<>> s2 = sin(t)
s2 =
    0.2960 0.0363 -0.0363 -0.2960
<>> s3 = cos(t)
s3 =
    0.7040 0.9637 0.9637 0.7040
<>> c3 = cos(t)
c3 =
    0.2 % c3 is an array obtained from term-wise squaring ( .^ 2 ) of the sin(t) array
<>> c3 = cos(t)
c3 =
    0.2960 0.0363 0.0363 -0.2960
% s3 is a vector obtained from the term-wise sin of the vector t
% s2 is a vector obtained by term-by-term division of z and y
% d is the vector obtained by term-by-term product of z and y
% p is the vector obtained by term-by-term product of z and y
% y is updated to 2 * y (each term of y is multiplied by 2)
% x is the sum of vectors y and z (with same size 1 X 3)
% x is also known as

```

$\mathbb{M}_A(x)$	Membership Function that returns 1 if $x \in A$ and 0 otherwise
\mathbb{M}_d : $(-\infty, \infty)^d$	d-dimensional Real Space

Table 6.1: Symbol Table: Probability and Statistics

x^k	$E(X^k) = \frac{1}{P_X} \int_{-\infty}^{\infty} x^k f(x) dx$
e_{nX}	$\phi(X)(t) := E(e^{nX})$
x	$E(X) = \frac{1}{P_X} \int_{-\infty}^{\infty} x f(x) dx$
$g(x)$	$E(g(X)) = \begin{cases} \int_x^{\infty} g(x) f(x) dx & \text{if } X \text{ is a discrete RV} \\ \int_{-\infty}^{\infty} g(x) f(x) dx & \text{if } X \text{ is a continuous RV} \end{cases}$
definition	Some Common Expectations
also known as	

$$\text{Expectation of a function } g(X) \text{ of a random variable } X \text{ is defined as:}$$

$P(a > X \geq b) = F(b) - F(a)$	$P(a > X) = \int_a^{\infty} f(x) dx$
$P(x) \leq 1$	$P(x) \leq 1$
$P(x) = P(X \leq x)$ is a probability	$P(x) = \int_{-\infty}^x f(x) dx$
$P(x) \geq 0$	$P(x) \geq 0$
$f(x)$: Distribution function (DF)	
$f(x)$: Probability density function (PDF)	
$\int_a^b f(x) dx$	Areas underneath $f(x)$ measure probabilities.
$\int_{-\infty}^{\infty} f(x) dx = 1$	• Areas underneath $f(x)$ measure probabilities.
$P(a < X \leq b) = F(b) - F(a)$	• $P(a < X \leq b) = \int_a^b f(x) dx$
$P(a < X) = \int_a^{\infty} f(x) dx$	• $P(a < X) = \int_a^{\infty} f(x) dx$ for every x where $f(x)$ is continuous

CHAPTER 6. FINITE MARKOV CHAINS

```
>> sSq + cSq % we can add the two arrays sSq and cSq to get the array of 1's
ans = 1 1 1 1
>> n = sin(t) .^2 + cos(t) .^2           % we can directly do term-wise operation sin^2(t) + cos^2(t) of t as well
n = 1 1 1 1
>> t2 = (-10:6.666665:10)             % t2 is similar to t above but with ':' syntax of (start:increment:stop)
t2 = -10.0000 -3.3333 3.3333 10.0000
```

Similarly, operations can be performed with matrices.

```
>> (0+0) .^ (1/2) % term-by-term square root of the matrix obtained by adding 0=ones(4,5) to itself
ans =
1.4142 1.4142 1.4142 1.4142 1.4142
1.4142 1.4142 1.4142 1.4142 1.4142
1.4142 1.4142 1.4142 1.4142 1.4142
1.4142 1.4142 1.4142 1.4142 1.4142
```

We can access specific rows or columns of a matrix as follows:

```
>> % declare a 3 X 3 array A of row vectors
>> A = [0.2760 0.4984 0.7513; 0.6797 0.9597 0.2551; 0.1626 0.5853 0.6991]
A =
0.2760 0.4984 0.7513
0.6797 0.9597 0.2551
0.1626 0.5853 0.6991
>> A(2,:) % access the second row of A
ans =
0.6797 0.9597 0.2551
>> B = A(2:3,:); % store the second and third rows of A in matrix B
B =
0.6797 0.9597 0.2551
0.1626 0.5853 0.6991
>> C = A(:,[1 3]) % store the first and third columns of A in matrix C
C =
0.2760 0.7513
0.6797 0.2551
```

Labwork 11 (Plotting a function as points of ordered pairs in two arrays) Next we plot the function $\sin(x)$ from several ordered pairs $(x_i, \sin(x_i))$. Here x_i 's are from the domain $[-2\pi, 2\pi]$. We use the `plot` function in MATLAB. Create an M-file called `MySineWave.m` and copy the following commands in it. By entering `MySineWave` in the command window you should be able to run the script and see the figure in the figure window.

SineWave.m

```
x = linspace(-2*pi,2*pi,100); % x has 100 points equally spaced in [-2*pi, 2*pi]
y = sin(x); % y is the term-wise sin of x, ie sin of every number in x is in y, resp.
plot(x,y,'.');
xlabel('x'); % label x-axis with the single quote enclosed string x
ylabel('sin(x)',FontSize',16); % label y-axis with the single quote enclosed string
title('Sine Wave in [-2 pi, 2 pi]',FontSize',16); % give a title; click Figure window to see changes
set(gca,'XTick',-8:1:8,FontSize',16) % change the range and size of X-axis ticks
% you can go to the Figure window's File menu to print/save the plot
```

The plot was saved as an encapsulated postscript file from the File menu of the Figure window and is displayed below.

CONDITIONAL PROBABILITY SUMMARY

$P(A|B)$ means the probability that A occurs given that B has occurred.

$$P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{P(A)P(B|A)}{P(B)} \quad \text{if } P(B) \neq 0$$

$$P(B|A) = \frac{P(A \cap B)}{P(A)} = \frac{P(B)P(A|B)}{P(A)} \quad \text{if } P(A) \neq 0$$

Conditional probabilities obey the 4 axioms of probability.

DISCRETE RANDOM VARIABLE SUMMARY

Probability mass function

$$f(x) = P(X = x_i)$$

Distribution function

$$F(x) = \sum_{x_i \leq x} f(x_i)$$

Random Variable	Possible Values	Probabilities	Modeled situations
Discrete uniform	$\{x_1, x_2, \dots, x_k\}$	$P(X = x_i) = \frac{1}{k}$	Situations with k equally likely values. Parameter: k .
Bernoulli(θ)	$\{0, 1\}$	$P(X = 0) = 1 - \theta$ $P(X = 1) = \theta$	Situations with only 2 outcomes, coded 1 for success and 0 for failure. Parameter: $\theta = P(\text{success}) \in (0, 1)$.
Geometric(θ)	$\{1, 2, 3, \dots\}$	$P(X = x) = (1 - \theta)^{x-1}\theta$	Situations where you count the number of trials until the first success in a sequence of independent trials with a constant probability of success. Parameter: $\theta = P(\text{success}) \in (0, 1)$.
Binomial(n, θ)	$\{0, 1, 2, \dots, n\}$	$P(X = x) = \binom{n}{x} \theta^x (1 - \theta)^{n-x}$	Situations where you count the number of success in n trials where each trial is independent and there is a constant probability of success. Parameters: $n \in \{1, 2, \dots\}$; $\theta = P(\text{success}) \in (0, 1)$.
Poisson(λ)	$\{0, 1, 2, \dots\}$	$P(X = x) = \frac{\lambda^x e^{-\lambda}}{x!}$	Situations where you count the number of events in a continuum where the events occur one at a time and are independent of one another. Parameter: $\lambda = \text{rate} \in (0, \infty)$.

$$\langle x_{1:\infty} \rangle = -\frac{1}{1} + \frac{2}{1} - \frac{3}{1} + \dots \quad \text{and} \quad \langle x_{1:\infty}^2 \rangle = -\frac{1}{1} + \frac{4}{1} - \frac{9}{1} + \dots$$

and finally some that approach 0 from either side are:

$$\langle x_{1:\infty} \rangle = \frac{1}{1} - \frac{2}{1} - \frac{3}{1} \dots \quad \text{and} \quad \langle x_{1:\infty}^2 \rangle = \frac{1}{1} - \frac{4}{1} - \frac{9}{1} \dots$$

and some that approach the limit 0 from the left are:

$$\langle x_{1:\infty} \rangle = \frac{1}{1} \frac{1}{1} \frac{1}{1} \dots \quad \text{and} \quad \langle x_{1:\infty}^2 \rangle = \frac{1}{1} \frac{8}{1} \frac{27}{1} \dots$$

the limit 0 from the right are:

However, several other sequences also approach the limit 0. Some such sequences that approach

$$\text{for every } j \geq N_m = m, |x_j - 0| = \left| \frac{j}{1} - 0 \right| = \frac{j}{1} \leq \frac{m}{1}$$

because for every $m \in \mathbb{N}$, we can take $N_m = m$ and satisfy the definition of the limit, i.e.:

Example 13 (Limit of $1/i$) Let $\langle x_i \rangle_{i=1}^{\infty} = \frac{1}{1}, \frac{1}{2}, \frac{1}{3}, \dots$, i.e. $x_i = \frac{1}{i}$, then $\lim_{i \rightarrow \infty} x_i = 0$. This is

$$\text{for every } j \geq N_m = 1, |x_j - 17| = |17 - 17| = 0 \leq \frac{m}{1}$$

This is because for every $m \in \mathbb{N}$, we can take $N_m = 1$ and satisfy the definition of the limit, i.e.:

Example 12 (Limit of a sequence of 17s) Let $\langle x_i \rangle_{i=1}^{\infty} = 17, 17, 17, \dots$. Then $\lim_{i \rightarrow \infty} x_i = 17$.

sequence beyond the N_m -th element is within distance $\frac{1}{m}$ of the limit a .
In words, $\lim_{i \rightarrow \infty} x_i = a$ means the following: no matter how small you make $\frac{1}{m}$ by picking a

$$|x_j - a| \leq \frac{1}{m}.$$

if for every natural number $m \in \mathbb{N}$, a natural number $N_m \in \mathbb{N}$ exists such that for every $j \geq N_m$,

$$\lim_{i \rightarrow \infty} x_i = a,$$

x_1, x_2, \dots is said to converge to a limit $a \in \mathbb{R}$ and denoted by:

Definition 3 (Convergent sequence of real numbers) A sequence of real numbers $\langle x_i \rangle_{i=1}^{\infty} :=$

Let us first recall some elementary ideas from real analysis.

1.8.1 Limits of Real Numbers – A Review

1.8 Elementary Real Analysis

Figure 1.7: A plot of the sine wave over $[-2\pi, 2\pi]$.

Summary of Probability Theory I

SET SUMMARY	
$\{a_1, a_2, \dots, a_n\}$	a set containing the elements a_1, a_2, \dots, a_n .
$a \in A$	a is an element of the set A .
$A \subseteq B$	the set A is a subset of B .
$A \cup B$	union, meaning the set of all elements which are in A or B , or both.
$\{\} \text{ or } \emptyset$	empty set, "intersection", meaning the set of all elements in both A and B .
A^c	the complement of A , meaning the set of all elements in Ω which are not in A .
Ω	universal set, a set of all outcomes of the experiment.
$A \bar{\in} Q$	A is a subset of Q , an individual outcome in Q , called a simple event.
ω	one particular outcome of an experiment resulting in 1 outcome.
EXPERIMENT SUMMARY	
Ω	an activity producing distinct outcomes.
ω	an activity producing distinct outcomes.
$A \subseteq \Omega$	A is an individual outcome in Ω , called a simple event.
Ω	an activity producing distinct outcomes.
PROBABILITY SUMMARY	
$P(A \cup B)$	$P(A) + P(B) - P(A \cap B)$ [always true]
$P(A^c)$	$1 - P(A)$
Rules:	
1. If A, A_2, \dots are disjoint, then $P(A_1 \cup A_2 \cup \dots) = P(A_1) + P(A_2) + \dots$	[This is true only when A and B are disjoint.]
2. If A, B are disjoint events, then $P(A \cup B) = P(A) + P(B)$.	
3. If A_1, A_2, \dots are disjoint then $P(A_1 \cap A_2 \cap \dots) = P(A_1) \cdot P(A_2) \cdot \dots$	

When we do not particularly care about the specifics of a sequence of real numbers $\langle x_{1:\infty} \rangle$, in terms of the exact values it takes for each i , but we are only interested that it converges to a limit a we write:

$$x \rightarrow a$$

and say that x approaches a . If we are only interested in those sequences that converge to the limit a from the right or left, we write:

$$x \rightarrow a^+ \quad \text{or} \quad x \rightarrow a^-$$

and say x approaches a from the right or left, respectively.

Definition 4 (Limits of Functions) We say a function $f(x) : \mathbb{R} \rightarrow \mathbb{R}$ has a **limit** $L \in \mathbb{R}$ as x approaches a and write:

$$\lim_{x \rightarrow a} f(x) = L ,$$

provided $f(x)$ is arbitrarily close to L for all values of x that are sufficiently close to, but not equal to, a . We say that f has a **right limit** L_R or **left limit** L_L as x approaches a from the left or right, and write:

$$\lim_{x \rightarrow a^+} f(x) = L_R \quad \text{or} \quad \lim_{x \rightarrow a^-} f(x) = L_L ,$$

provided $f(x)$ is arbitrarily close to L_R or L_L for all values of x that are sufficiently close to, but not equal to, a from the right of a or the left of a , respectively. When the limit is not an element of \mathbb{R} or when the left and right limits are distinct, we say that the limit does not exist.

Example 14 (Limit of $1/x^2$) Consider the function $f(x) = \frac{1}{x^2}$. Then

$$\lim_{x \rightarrow 1} f(x) = \lim_{x \rightarrow 1} \frac{1}{x^2} = 1$$

exists since the limit $1 \in \mathbb{R}$, and the right and left limits are the same:

$$\lim_{x \rightarrow 1^+} f(x) = \lim_{x \rightarrow 1^+} \frac{1}{x^2} = 1 \quad \text{and} \quad \lim_{x \rightarrow 1^-} f(x) = \lim_{x \rightarrow 1^-} \frac{1}{x^2} = 1 .$$

However, the following limit does not exist:

$$\lim_{x \rightarrow 0} f(x) = \lim_{x \rightarrow 0} \frac{1}{x^2} = \infty$$

since $\infty \notin \mathbb{R}$.

Let us next look at some limits of functions that exist despite the function itself being undefined at the limit point.

Example 15 (Limit of $(1+x)^{\frac{1}{x}}$) The limit of $f(x) = (1+x)^{\frac{1}{x}}$ as x approaches 0 exists and it is

6.6 Standard normal distribution function table

For any given value z , its cumulative probability $\Phi(z)$.

z	$\Phi(z)$										
0.01	0.5040	0.51	0.6950	1.01	0.8438	1.51	0.9345	2.01	0.9778	2.51	0.9940
0.02	0.5080	0.52	0.6985	1.02	0.8461	1.52	0.9357	2.02	0.9783	2.52	0.9941
0.03	0.5120	0.53	0.7019	1.03	0.8485	1.53	0.9370	2.03	0.9788	2.53	0.9943
0.04	0.5160	0.54	0.7054	1.04	0.8508	1.54	0.9382	2.04	0.9793	2.54	0.9945
0.05	0.5199	0.55	0.7088	1.05	0.8531	1.55	0.9394	2.05	0.9798	2.55	0.9946
0.06	0.5239	0.56	0.7123	1.06	0.8554	1.56	0.9406	2.06	0.9803	2.56	0.9948
0.07	0.5279	0.57	0.7157	1.07	0.8577	1.57	0.9418	2.07	0.9808	2.57	0.9949
0.08	0.5319	0.58	0.7190	1.08	0.8599	1.58	0.9429	2.08	0.9812	2.58	0.9951
0.09	0.5359	0.59	0.7224	1.09	0.8621	1.59	0.9441	2.09	0.9817	2.59	0.9952
0.10	0.5398	0.60	0.7257	1.10	0.8643	1.60	0.9452	2.10	0.9821	2.60	0.9953
0.11	0.5438	0.61	0.7291	1.11	0.8665	1.61	0.9463	2.11	0.9826	2.61	0.9955
0.12	0.5478	0.62	0.7324	1.12	0.8686	1.62	0.9474	2.12	0.9830	2.62	0.9956
0.13	0.5517	0.63	0.7357	1.13	0.8708	1.63	0.9484	2.13	0.9834	2.63	0.9957
0.14	0.5557	0.64	0.7389	1.14	0.8729	1.64	0.9495	2.14	0.9838	2.64	0.9959
0.15	0.5596	0.65	0.7422	1.15	0.8749	1.65	0.9505	2.15	0.9842	2.65	0.9960
0.16	0.5636	0.66	0.7454	1.16	0.8770	1.66	0.9515	2.16	0.9846	2.66	0.9961
0.17	0.5675	0.67	0.7486	1.17	0.8790	1.67	0.9525	2.17	0.9850	2.67	0.9962
0.18	0.5714	0.68	0.7517	1.18	0.8810	1.68	0.9535	2.18	0.9854	2.68	0.9963
0.19	0.5753	0.69	0.7549	1.19	0.8830	1.69	0.9545	2.19	0.9857	2.69	0.9964
0.20	0.5793	0.70	0.7580	1.20	0.8849	1.70	0.9554	2.20	0.9861	2.70	0.9965
0.21	0.5832	0.71	0.7611	1.21	0.8869	1.71	0.9564	2.21	0.9864	2.71	0.9966
0.22	0.5871	0.72	0.7642	1.22	0.8888	1.72	0.9573	2.22	0.9868	2.72	0.9967
0.23	0.5910	0.73	0.7673	1.23	0.8907	1.73	0.9582	2.23	0.9871	2.73	0.9968
0.24	0.5948	0.74	0.7704	1.24	0.8925	1.74	0.9591	2.24	0.9875	2.74	0.9969
0.25	0.5987	0.75	0.7734	1.25	0.8944	1.75	0.9599	2.25	0.9878	2.75	0.9970
0.26	0.6026	0.76	0.7764	1.26	0.8962	1.76	0.9608	2.26	0.9881	2.76	0.9971
0.27	0.6064	0.77	0.7794	1.27	0.8980	1.77	0.9616	2.27	0.9884	2.77	0.9972
0.28	0.6103	0.78	0.7823	1.28	0.8997	1.78	0.9625	2.28	0.9887	2.78	0.9973
0.29	0.6141	0.79	0.7852	1.29	0.9015	1.79	0.9633	2.29	0.9890	2.79	0.9974
0.30	0.6179	0.80	0.7881	1.30	0.9032	1.80	0.9641	2.30	0.9893	2.80	0.9974
0.31	0.6217	0.81	0.7910	1.31	0.9049	1.81	0.9649	2.31	0.9896	2.81	0.9975
0.32	0.6255	0.82	0.7939	1.32	0.9066	1.82	0.9656	2.32	0.9898	2.82	0.9976
0.33	0.6293	0.83	0.7967	1.33	0.9082	1.83	0.9664	2.33	0.9901	2.83	0.9977
0.34	0.6331	0.84	0.7995	1.34	0.9099	1.84	0.9671	2.34	0.9904	2.84	0.9977
0.35	0.6368	0.85	0.8023	1.35	0.9115	1.85	0.9678	2.35	0.9906	2.85	0.9978
0.36	0.6406	0.86	0.8051	1.36	0.9131	1.86	0.9686	2.36	0.9909	2.86	0.9979
0.37	0.6443	0.87	0.8078	1.37	0.9147	1.87	0.9693	2.37	0.9911	2.87	0.9979
0.38	0.6480	0.88	0.8106	1.38	0.9162	1.88	0.9699	2.38	0.9913	2.88	0.9980
0.39	0.6517	0.89	0.8133	1.39	0.9177	1.89	0.9706	2.39	0.9916	2.89	0.9981
0.40	0.6554	0.90	0.8159	1.40	0.9192	1.90	0.9713	2.40	0.9918	2.90	0.9981
0.41	0.6591	0.91	0.8186	1.41	0.9207	1.91	0.9719	2.41	0.9920	2.91	0.9982
0.42	0.6628	0.92	0.8212	1.42	0.9222	1.92	0.9726	2.42	0.9922	2.92	0.9982
0.43	0.6664	0.93	0.8238	1.43	0.9236	1.93	0.9732	2.43	0.9925	2.93	0.9983
0.44	0.6700	0.94	0.8264	1.44	0.9251	1.94	0.9738	2.44	0.9927	2.94	0.9984
0.45	0.6736	0.95	0.8289	1.45	0.9265	1.95	0.9744	2.45	0.9929	2.95	0.9984
0.46	0.6772	0.96	0.8315	1.46	0.9279	1.96	0.9750	2.46	0.9931	2.96	0.9985
0.47	0.6808	0.97	0.8340	1.47	0.9292	1.97	0.9756	2.47	0.9932	2.97	0.9985
0.48	0.6844	0.98	0.8365	1.48	0.9306	1.98	0.9761	2.48	0.9934	2.98	0.9986
0.49	0.6879	0.99	0.8389	1.49	0.9319	1.99	0.9767	2.49	0.9936	2.99	0.9986
0.50	0.6915	1.00	0.8413	1.50	0.9332	2.00	0.9772	2.50	0.9938	3.00	0.9987

Example 187 (Drunkard's biased walk around the block) Consider the Markov chain $(X_t)_{t \in \mathbb{Z}^+}$ on $\mathbb{X} = \{0, 1, 2, 3\}$ with initial distribution $\mathbb{I}_{\{3\}}(x)$ and transition matrix

$$P = \begin{pmatrix} 0 & 1/3 & 0 & 2/3 \\ 1/3 & 0 & 2/3 & 0 \\ 0 & 1/3 & 0 & 2/3 \\ 2 & 0 & 1/3 & 0 & 2/3 \end{pmatrix}.$$

Draw the transition diagram for this Markov chain that corresponds to a drunkard who flips a biased coin to make his next move at each corner. The stationary distribution is $\pi = (1/4, 1/4, 1/4, 1/4)$. We will show that $(X_t)_{t \in \mathbb{Z}^+}$ is not a reversible Markov chain. Since $(X_t)_{t \in \mathbb{Z}^+}$ is irreducible (apply Proposition 97), it has to be a reversible distribution in order for $(X_t)_{t \in \mathbb{Z}^+}$ to be a reversible Markov chain. But reversibility fails for π since,

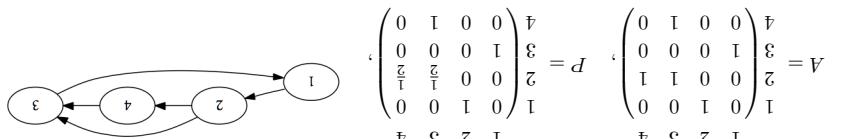
$$\pi(0)P(0, 1) = \frac{1}{4} \times \frac{1}{3} = \frac{1}{12} < \frac{1}{1} \times \frac{2}{3} = \pi(1)P(1, 0).$$

Exercise 188 Find the stationary distribution of the Markov chain in Exercise 180.

Model 22 (Random Walk on a Directed Graph) A random walk on a directed graph $G = (V, E)$ is a Markov chain with state space $V = \{v_1, v_2, \dots, v_k\}$ and transition matrix given by:

$$P(v_i, v_j) = \begin{cases} \frac{1}{\deg(v_i)} & \text{if } (v_i, v_j) \in E \\ 0 & \text{otherwise} \end{cases}$$

Example 189 (Directed Triangulated Quadrangle) The random walk on the directed graph



depicted below with adjacency matrix A is a Markov chain on $\{1, 2, 3, 4\}$ with transition matrix P :

$$G = \{(1, 2, 3, 4), (1, 2), (3, 1), (2, 3), (2, 4), (4, 3)\}$$

Exercise 190 Show that there is no reversible distribution for the Markov chain in Example 189.

Example 191 (Random surf on the word wide web) Consider the huge graph with vertices as webpages and hyper-links as undirected edges. Then Model 21 gives a random walk on this graph. However if a page has no links to other pages, it becomes a sink and therefore terminates the random walk. Let us modify this random walk into a **random surf** to avoid getting stuck. If the random surfer arrives at a sink page, she picks another page at random and continues surfing the random walk. Finally the random walk ends at a point $a \in D$, provided that the random surfer never gets stuck. The random walk on this random walk on this graph is no different than the random walk on the word wide web is a very successful model for ranking pages.

Finally, f is said to be continuous if f is continuous at every $a \in D$.

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

respectively. We say f is continuous at $a \in D$, provided:

$$\lim_{x \rightarrow a^+} f(x) = f(a) \quad \text{or} \quad \lim_{x \rightarrow a^-} f(x) = f(a).$$

Definition 5 (Continuity of a function) We say a real-valued function $f(x) : D \rightarrow \mathbb{R}$ with the domain $D \subset \mathbb{R}$ is **right continuous or left continuous** at a point $a \in D$, provided:

$$\lim_{n \rightarrow \infty} f(n) = \lim_{n \leftarrow \infty} \left(1 - \frac{n}{a} \right) = 1.$$

proaches ∞ exists and it is 1 :

$$\lim_{n \leftarrow \infty} f(n) = \lim_{n \leftarrow \infty} \left(1 - \frac{n}{a} \right) = e^{-a}.$$

Example 17 (Limit of $(1 - \frac{n}{a})^n$) The limit of $f(n) = (1 - \frac{n}{a})^n$ as n approaches ∞ exists and it is e^{-a} :

Next we look at some examples of limits at infinity.

despite the fact that $f(1) = \frac{1}{1-1} = \frac{1}{0}$ itself is undefined and does not exist.

$$\lim_{x \rightarrow 1^+} f(x) = \lim_{x \rightarrow 1^+} x^3 - 1 = \lim_{x \rightarrow 1^+} (x-1)(x^2 + x + 1) = \lim_{x \rightarrow 1^+} x^2 + x + 1 = 3$$

Example 16 (Limit of $\frac{x^3-1}{x-1}$) For $f(x) = \frac{x^3-1}{x-1}$, this limit exists:

not exist.

Notice that the above limit exists despite the fact that $f(0) = (1+0)^{\frac{1}{0}}$ itself is undefined and does not exist.

$$\lim_{x \leftarrow 0^0} f(x) = \lim_{x \leftarrow 0^0} (1+x)^{\frac{1}{x}} = e \approx 2.71828.$$

$$= \exp(1) = e \approx 2.71828.$$

$$= \exp(1/\lim_{x \rightarrow 0^0}(x+1)) \quad \text{The limit of } x+1 \text{ as } x \text{ approaches 0 is 1}$$

$$= \exp(\lim_{x \rightarrow 0^0} 1/(x+1)) \quad \text{limit of a quotient is the quotient of the limits}$$

$$= \exp\left(\lim_{x \rightarrow 0^0} d \log(x+1)/dx\right) \quad \text{Applying L'Hospital's rule}$$

$$= \exp\left(\lim_{x \rightarrow 0^0} (\log(x+1))/x\right) \quad \text{Indeterminate form of type } 0/0.$$

$$= \exp\left(\lim_{x \rightarrow 0^0} \log((x+1)^{1/x})\right) \quad \text{Transformed using } \exp(\lim_{x \rightarrow 0^0} \log((x+1)^{1/x}))$$

$$= \lim_{x \rightarrow 0^0} (x+1)^{1/x} \quad \text{Indeterminate form of type } 1^\infty.$$

$$\lim_{x \rightarrow 0^0} f(x) = \lim_{x \rightarrow 0^0} (1+x)^{\frac{1}{x}}$$

the Euler's constant e :

Example 19 (Discontinuity of $f(x) = (1+x)^{\frac{1}{x}}$ at 0) Let us reconsider the function $f(x) = (1+x)^{\frac{1}{x}} : \mathbb{R} \rightarrow \mathbb{R}$. Clearly, $f(x)$ is continuous at 1, since:

$$\lim_{x \rightarrow 1} f(x) = \lim_{x \rightarrow 1} (1+x)^{\frac{1}{x}} = 2 = f(1) = (1+1)^{\frac{1}{1}},$$

but it is not continuous at 0, since:

$$\lim_{x \rightarrow 0} f(x) = \lim_{x \rightarrow 0} (1+x)^{\frac{1}{x}} = e \approx 2.71828 \neq f(0) = (1+0)^{\frac{1}{0}}.$$

Thus, $f(x)$ is not a continuous function over \mathbb{R} .

1.9 Elementary Number Theory

We introduce basic notions that we need from elementary number theory here. These notions include integer functions and modular arithmetic as they will be needed later on.

For any real number x :

$\lfloor x \rfloor := \max\{y : y \in \mathbb{Z} \text{ and } y \leq x\}$, i.e., the greatest integer less than or equal to x (the **floor** of x),
 $\lceil x \rceil := \min\{y : y \in \mathbb{Z} \text{ and } y \geq x\}$, i.e., the least integer greater than or equal to x (the **ceiling** of x).

Example 20 (Floors and ceilings)

$$\lfloor 1 \rfloor = 1, \quad \lceil 1 \rceil = 1, \quad \lfloor 17.8 \rfloor = 17, \quad \lceil -17.8 \rceil = -18, \quad \lceil \sqrt{2} \rceil = 2, \quad \lfloor \pi \rfloor = 3, \quad \lceil \frac{1}{10^{100}} \rceil = 1.$$

Labwork 21 (Floors and ceilings in MATLAB) We can use MATLAB functions `floor` and `ceil` to compute $\lfloor x \rfloor$ and $\lceil x \rceil$, respectively. Also, the argument x to these functions can be an array.

```
>> sqrt(2) % the square root of 2 is
ans = 1.4142
>> ceil(sqrt(2)) % ceiling of square root of 2
ans = 2
>> floor(-17.8) % floor of -17.8
ans = -18
>> ceil([1 sqrt(2) pi -17.8 1/(10^100)]) %the ceiling of each element of an array
ans = 1 2 4 -17 1
>> floor([1 sqrt(2) pi -17.8 1/(10^100)]) % the floor of each element of an array
ans = 1 1 3 -18 0
```

Classwork 22 (Relations between floors and ceilings) Convince yourself of the following formulae. Use examples, plots and/or formal arguments.

$$\begin{aligned} \lceil x \rceil &= \lfloor x \rfloor \iff x \in \mathbb{Z} \\ \lceil x \rceil &= \lfloor x \rfloor + 1 \iff x \notin \mathbb{Z} \\ \lfloor -x \rfloor &= -\lceil x \rceil \\ x - 1 < \lfloor x \rfloor &\leq x \leq \lceil x \rceil < x + 1 \end{aligned}$$

Let us define modular arithmetic next. Suppose x and y are any real numbers, i.e. $x, y \in \mathbb{R}$, we define the binary operation called “ $x \bmod y$ ” as:

$$x \bmod y := \begin{cases} x - y\lfloor x/y \rfloor & \text{if } y \neq 0 \\ x & \text{if } y = 0 \end{cases}$$

Proposition 99 The random walk on an undirected graph $\mathbb{G} = (\mathbb{V}, \mathbb{E})$, with vertex set $\mathbb{V} := \{v_1, v_2, \dots, v_k\}$ and degree sum $d = \sum_{i=1}^k \deg(v_i)$ is a reversible Markov chain with the reversible distribution π given by:

$$\pi = \left(\frac{\deg(v_1)}{d}, \frac{\deg(v_2)}{d}, \dots, \frac{\deg(v_k)}{d} \right).$$

Proof: First note that π is a probability distribution provided that $d > 0$. To show that π is reversible we need to verify Equation 6.12 for each $(v_i, v_j) \in \mathbb{V}^2$. Fix a pair of states $(v_i, v_j) \in \mathbb{V}^2$, then

$$\pi(v_i)P(v_i, v_j) = \begin{cases} \frac{\deg(v_i)}{d} \frac{1}{\deg(v_j)} = \frac{1}{d} = \frac{\deg(v_j)}{d} \frac{1}{\deg(v_i)} = \pi(v_j)P(v_j, v_i) & \text{if } \langle v_i, v_j \rangle \in \mathbb{E} \\ 0 = \pi(v_j)P(v_j, v_i) & \text{otherwise.} \end{cases}$$

By Proposition 97 π is also the stationary distribution.

Exercise 184 Prove Proposition 99 by directly showing that $\pi P = \pi$, i.e., for each $v_i \in \mathbb{V}$, $\sum_{i=1}^k \pi(v_i)P(v_i, v_j) = \pi(v_j)$.

Example 185 (Random Walk on a regular graph) A graph $\mathbb{G} = (\mathbb{V}, \mathbb{E})$ is called regular if every vertex in $\mathbb{V} = \{v_1, v_2, \dots, v_k\}$ has the same degree δ , i.e., $\deg(v_i) = \delta$ for every $v_i \in \mathbb{V}$. Consider the random walk on a regular graph with symmetric transition matrix

$$Q(v_i, v_j) = \begin{cases} \frac{1}{\delta} & \text{if } \langle v_i, v_j \rangle \in \mathbb{E} \\ 0 & \text{otherwise} \end{cases}.$$

By Proposition 99, the stationary distribution of the random walk on \mathbb{G} is the uniform distribution on \mathbb{V} given by

$$\pi = \left(\frac{\delta}{\delta \# \mathbb{V}}, \dots, \frac{\delta}{\delta \# \mathbb{V}} \right) = \left(\frac{1}{\# \mathbb{V}}, \dots, \frac{1}{\# \mathbb{V}} \right).$$

Example 186 (Triangulated Quadrangle) The random walk on the undirected graph

$$\mathbb{G} = (\{1, 2, 3, 4\}, \{(1, 2), (3, 1), (2, 3), (2, 4), (4, 3)\})$$

depicted below with adjacency matrix A is a Markov chain on $\{1, 2, 3, 4\}$ with transition matrix P :

$$A = \begin{pmatrix} 1 & 2 & 3 & 4 \\ 0 & 1 & 1 & 0 \\ 2 & 1 & 0 & 1 & 1 \\ 3 & 1 & 1 & 0 & 1 \\ 4 & 0 & 1 & 1 & 0 \end{pmatrix}, \quad P = \begin{pmatrix} 1 & 2 & 3 & 4 \\ 0 & \frac{1}{2} & \frac{1}{2} & 0 \\ 2 & \frac{1}{3} & 0 & \frac{1}{3} & \frac{1}{3} \\ 3 & \frac{1}{3} & \frac{1}{3} & 0 & \frac{1}{3} \\ 4 & 0 & \frac{1}{2} & \frac{1}{2} & 0 \end{pmatrix},$$

By Proposition 99, the stationary distribution of the random walk on \mathbb{G} is

$$\pi = \left(\frac{\deg(v_1)}{d}, \frac{\deg(v_2)}{d}, \frac{\deg(v_3)}{d}, \frac{\deg(v_4)}{d} \right) = \left(\frac{2}{10}, \frac{3}{10}, \frac{3}{10}, \frac{2}{10} \right).$$

The outcome of a random experiment is **uncertain** until it is performed and those outcomes have to be *discreteable in some well-specified sense* to the experiment. An experiment's sample space is merely a collection of distinct elements called outcomes and these sample spaces need to reflect the problem in hand. The example below is to convince you that

The simple events of Ω are $\{1\}, \{2\}, \{3\}, \{4\}, \{5\}$, and $\{6\}$.

Some examples of events are the set of odd numbered outcomes $A = \{1, 3, 5\}$, and the set of even numbered outcomes $B = \{2, 4, 6\}$.

- If our experiment is to roll a die then there are six outcomes corresponding to the number that shows on the top. For this experiment, $\Omega = \{1, 2, 3, 4, 5, 6\}$.

This time, $\Omega = \{\omega_1, \omega_2\}$ where $\omega_1 = \text{Heads}$ and $\omega_2 = \text{Tails}$.

- $\Omega = \{\text{Heads}, \text{Tails}\}$ if our experiment is to note the outcome of a coin toss.
- $\Omega = \{\text{Defective}, \text{Non-defective}\}$ if our experiment is to inspect a light bulb.

Example 23 Some standard examples of experiments are the following:

- Ω are only two outcomes here, so $\Omega = \{\omega_1, \omega_2\}$ where $\omega_1 = \text{Defective and } \omega_2 = \text{Non-defective}$.
- Ω is a **sample space**, ω is an **outcome**, ω_i is a **simple event**.
- The subsets of Ω are called **events**. A single outcome, ω , when seen as a subset of Ω , as in $\{\omega\}$, is called a **simple event**.
- The subsets of Ω are called **events**. A set of all outcomes is called the **sample space**, and is denoted by Ω .
- Possibilities called **outcomes**. The set of all outcomes is called the **sample space**, and is denoted by Ω .
- Definition 6 An **experiment** is an activity or procedure that produces distinct, well-defined outcomes based on a proper theoretical basis by Fermat and Pascal in the early 17th century.

Ideas about chance events and random behaviour arose out of thousands of years of game playing, long before any attempt was made to use mathematical reasoning about them. Board and dice games were well known in Egyptian times, and Augustus Caesar gambled with dice. Calculations of odds for gamblers were put on a proper theoretical basis by Fermat and Pascal in the early 17th century.

2.1 Experiments

Probability Model

Chapter 2

In words, $\pi(x, y) = \pi(y, x)$ says that if you start the chain at the reversible distribution π , i.e., $y_0 = \pi$, then the probability of going from x to y is the same as that of going from y to x .

Proof: Suppose π is a reversible distribution for $(X^i)_{i \in \mathbb{Z}^+}$. Then π is a probability distribution for $(X^i)_{i \in \mathbb{Z}^+}$ and $\pi(x, y) = \pi(y, x)$ for each $(x, y) \in \mathbb{X}^2$. We need to show that for any $y \in \mathbb{X}$ we have

Proposition 97 (A reversible distribution π is a stationary π) Let $(X^i)_{i \in \mathbb{Z}^+}$ be a Markov chain on $\mathbb{X} = \{s_1, s_2, \dots, s_N\}$ with transition matrix P . If π is a reversible distribution for $(X^i)_{i \in \mathbb{Z}^+}$, then π is a stationary distribution for $(X^i)_{i \in \mathbb{Z}^+}$.

In words, $\pi(x, y) = \pi(y, x)$ says that if you start the chain at the reversible distribution π , i.e., $y_0 = \pi$, then the probability of going from x to y is the same as that of going from y to x .

$$\begin{aligned} LHS &= \pi(y) = \pi(y) \sum_{x \in \mathbb{X}} \pi(y, x), \text{ since } P \text{ is a stochastic matrix} \\ &= \sum_{x \in \mathbb{X}} \pi(y)P(y, x) = \sum_{x \in \mathbb{X}} \pi(x)P(x, y), \text{ by reversibility} \\ &= RHS. \end{aligned}$$

Fix a $y \in \mathbb{X}$,

$$\begin{aligned} \pi(y) &= \sum_{x \in \mathbb{X}} \pi(x)P(x, y). \end{aligned}$$

Proof: Suppose π is a reversible distribution for $(X^i)_{i \in \mathbb{Z}^+}$. Then π is a probability distribution for $(X^i)_{i \in \mathbb{Z}^+}$ and $\pi(x, y) = \pi(y, x)$ for each $(x, y) \in \mathbb{X}^2$. We need to show that for any $y \in \mathbb{X}$ we have

Proposition 97 (A reversible distribution π is a stationary π) Let $(X^i)_{i \in \mathbb{Z}^+}$ be a Markov chain on $\mathbb{X} = \{s_1, s_2, \dots, s_N\}$ with transition matrix P . If π is a reversible distribution for $(X^i)_{i \in \mathbb{Z}^+}$, then π is a stationary distribution for $(X^i)_{i \in \mathbb{Z}^+}$.

In words, $\pi(x, y) = \pi(y, x)$ says that if you start the chain at the reversible distribution π , i.e., $y_0 = \pi$, then the probability of going from x to y is the same as that of going from y to x .

Model 21 (Random Walk on an Undirected Graph) A random walk on an undirected graph $G = (\mathbb{V}, \mathbb{E})$ is a Markov chain with state space $\mathbb{V} = \{v_1, v_2, \dots, v_n\}$ and the following transition rules: if the chain is at vertex v_i at time t then it moves uniformly at random to one of the neighbors of v_i at vertex v_i is denoted by $\deg(v_i)$ respectively. Note that a transition diagram of a graph G has **in-edges** that come into it and **out-edges** that go out of it. The number of in-edges and out-edges of v_i is denoted by $\deg(v_i)$ and $\deg(v_i)$ respectively. Thus the adjacency matrix of an undirected graph is symmetric. In a directed graph, each vertex

Markov chain is a weighted directed graph and is represented by the transition probability matrix. Thus the adjacency matrix of an undirected graph is symmetric. In a directed graph, each vertex v_i has **in-edges** that come into it and **out-edges** that go out of it. The number of in-edges and out-edges of v_i is denoted by $\deg(v_i)$ and $\deg(v_i)$ respectively. Thus the adjacency matrix of an undirected graph is given by:

$$A := (A(v_i, v_j))_{(v_i, v_j) \in \mathbb{V} \times \mathbb{V}}, \quad A(v_i, v_j) = \begin{cases} 1 & \text{if } (v_i, v_j) \in \mathbb{E} \\ 0 & \text{otherwise} \end{cases}.$$

We can represent a directed graph by its **adjacency matrix** given by: or they can be **undirected**. An edge can be **weighted** by being associated with a real number or they can be **directed**. An edge can be **directed** to preserve the order of the pair of vertices they connect we only allow one edge per pair of vertices but in a **multigraph** we allow more than one edge per pair of vertices. An edge can be **directed** to prevent the edge v_i from being part of a cycle. In a graph graph is called its **degree** and is denoted by $\deg(v_i)$. Note that $\deg(v_i) = \# \text{nbhd}(v_i)$. In a directed graph the set of neighbours of a vertex v_i in an undirected graph is called its **neighbourhood** of v_i . The number of neighbours of a vertex v_i in a directed graph is denoted by $\deg(v_i)$. Two vertices are **neighbours** if they share an edge and is denoted by $(v_i, v_j) \in \mathbb{E}$. The set of neighbours of a vertex v_i in the set of neighbours of v_i is denoted by $\text{nbhd}(v_i) = \{(v_j : (v_i, v_j) \in \mathbb{E})\}$ is called its **closed neighbourhood** of v_i . The number of edges per vertex v_i is denoted by $\deg(v_i)$. Each edge connects two of the vertices in \mathbb{V} . All together with an edge set $\mathbb{E} = \{e_1, e_2, \dots, e_l\}$. Each edge connects two of the vertices in \mathbb{V} . A vertex v_i with an edge set $\{v_i\}$ consists of a vertex v_i and $\# \text{nbhd}(v_i) = l$. The number of edges in a graph is denoted by $\# \text{edges}$.

$$\begin{aligned} &= \sum_{x \in \mathbb{X}} \pi(x)P(x, y) = \sum_{x \in \mathbb{X}} \pi(x)P(y, x), \text{ by reversibility} \\ &= RHS. \end{aligned}$$

$$\begin{aligned} LHS &= \pi(y) = \pi(y) \sum_{x \in \mathbb{X}} P(y, x), \text{ since } P \text{ is a stochastic matrix} \\ &= \sum_{x \in \mathbb{X}} \pi(y)P(y, x) = \sum_{x \in \mathbb{X}} \pi(x)P(x, y), \text{ by reversibility} \\ &= RHS. \end{aligned}$$

$$\begin{aligned} &= \sum_{x \in \mathbb{X}} \pi(x)P(x, y) = \sum_{x \in \mathbb{X}} \pi(x)P(y, x), \text{ by reversibility} \\ &= RHS. \end{aligned}$$

$$\begin{aligned} &= \sum_{x \in \mathbb{X}} \pi(x)P(x, y) = \sum_{x \in \mathbb{X}} \pi(x)P(y, x), \text{ by reversibility} \\ &= RHS. \end{aligned}$$

$$\begin{aligned} &= \sum_{x \in \mathbb{X}} \pi(x)P(x, y) = \sum_{x \in \mathbb{X}} \pi(x)P(y, x), \text{ by reversibility} \\ &= RHS. \end{aligned}$$

$$\begin{aligned} &= \sum_{x \in \mathbb{X}} \pi(x)P(x, y) = \sum_{x \in \mathbb{X}} \pi(x)P(y, x), \text{ by reversibility} \\ &= RHS. \end{aligned}$$

$$\begin{aligned} &= \sum_{x \in \mathbb{X}} \pi(x)P(x, y) = \sum_{x \in \mathbb{X}} \pi(x)P(y, x), \text{ by reversibility} \\ &= RHS. \end{aligned}$$

$$\begin{aligned} &= \sum_{x \in \mathbb{X}} \pi(x)P(x, y) = \sum_{x \in \mathbb{X}} \pi(x)P(y, x), \text{ by reversibility} \\ &= RHS. \end{aligned}$$

$$\begin{aligned} &= \sum_{x \in \mathbb{X}} \pi(x)P(x, y) = \sum_{x \in \mathbb{X}} \pi(x)P(y, x), \text{ by reversibility} \\ &= RHS. \end{aligned}$$

$$\begin{aligned} &= \sum_{x \in \mathbb{X}} \pi(x)P(x, y) = \sum_{x \in \mathbb{X}} \pi(x)P(y, x), \text{ by reversibility} \\ &= RHS. \end{aligned}$$

$$\begin{aligned} &= \sum_{x \in \mathbb{X}} \pi(x)P(x, y) = \sum_{x \in \mathbb{X}} \pi(x)P(y, x), \text{ by reversibility} \\ &= RHS. \end{aligned}$$

$$\begin{aligned} &= \sum_{x \in \mathbb{X}} \pi(x)P(x, y) = \sum_{x \in \mathbb{X}} \pi(x)P(y, x), \text{ by reversibility} \\ &= RHS. \end{aligned}$$

$$\begin{aligned} &= \sum_{x \in \mathbb{X}} \pi(x)P(x, y) = \sum_{x \in \mathbb{X}} \pi(x)P(y, x), \text{ by reversibility} \\ &= RHS. \end{aligned}$$

$$\begin{aligned} &= \sum_{x \in \mathbb{X}} \pi(x)P(x, y) = \sum_{x \in \mathbb{X}} \pi(x)P(y, x), \text{ by reversibility} \\ &= RHS. \end{aligned}$$

$$\begin{aligned} &= \sum_{x \in \mathbb{X}} \pi(x)P(x, y) = \sum_{x \in \mathbb{X}} \pi(x)P(y, x), \text{ by reversibility} \\ &= RHS. \end{aligned}$$

$$\begin{aligned} &= \sum_{x \in \mathbb{X}} \pi(x)P(x, y) = \sum_{x \in \mathbb{X}} \pi(x)P(y, x), \text{ by reversibility} \\ &= RHS. \end{aligned}$$

$$\begin{aligned} &= \sum_{x \in \mathbb{X}} \pi(x)P(x, y) = \sum_{x \in \mathbb{X}} \pi(x)P(y, x), \text{ by reversibility} \\ &= RHS. \end{aligned}$$

$$\begin{aligned} &= \sum_{x \in \mathbb{X}} \pi(x)P(x, y) = \sum_{x \in \mathbb{X}} \pi(x)P(y, x), \text{ by reversibility} \\ &= RHS. \end{aligned}$$

$$\begin{aligned} &= \sum_{x \in \mathbb{X}} \pi(x)P(x, y) = \sum_{x \in \mathbb{X}} \pi(x)P(y, x), \text{ by reversibility} \\ &= RHS. \end{aligned}$$

$$\begin{aligned} &= \sum_{x \in \mathbb{X}} \pi(x)P(x, y) = \sum_{x \in \mathbb{X}} \pi(x)P(y, x), \text{ by reversibility} \\ &= RHS. \end{aligned}$$

$$\begin{aligned} &= \sum_{x \in \mathbb{X}} \pi(x)P(x, y) = \sum_{x \in \mathbb{X}} \pi(x)P(y, x), \text{ by reversibility} \\ &= RHS. \end{aligned}$$

$$\begin{aligned} &= \sum_{x \in \mathbb{X}} \pi(x)P(x, y) = \sum_{x \in \mathbb{X}} \pi(x)P(y, x), \text{ by reversibility} \\ &= RHS. \end{aligned}$$

$$\begin{aligned} &= \sum_{x \in \mathbb{X}} \pi(x)P(x, y) = \sum_{x \in \mathbb{X}} \pi(x)P(y, x), \text{ by reversibility} \\ &= RHS. \end{aligned}$$

$$\begin{aligned} &= \sum_{x \in \mathbb{X}} \pi(x)P(x, y) = \sum_{x \in \mathbb{X}} \pi(x)P(y, x), \text{ by reversibility} \\ &= RHS. \end{aligned}$$

$$\begin{aligned} &= \sum_{x \in \mathbb{X}} \pi(x)P(x, y) = \sum_{x \in \mathbb{X}} \pi(x)P(y, x), \text{ by reversibility} \\ &= RHS. \end{aligned}$$

$$\begin{aligned} &= \sum_{x \in \mathbb{X}} \pi(x)P(x, y) = \sum_{x \in \mathbb{X}} \pi(x)P(y, x), \text{ by reversibility} \\ &= RHS. \end{aligned}$$

$$\begin{aligned} &= \sum_{x \in \mathbb{X}} \pi(x)P(x, y) = \sum_{x \in \mathbb{X}} \pi(x)P(y, x), \text{ by reversibility} \\ &= RHS. \end{aligned}$$

$$\begin{aligned} &= \sum_{x \in \mathbb{X}} \pi(x)P(x, y) = \sum_{x \in \mathbb{X}} \pi(x)P(y, x), \text{ by reversibility} \\ &= RHS. \end{aligned}$$

$$\begin{aligned} &= \sum_{x \in \mathbb{X}} \pi(x)P(x, y) = \sum_{x \in \mathbb{X}} \pi(x)P(y, x), \text{ by reversibility} \\ &= RHS. \end{aligned}$$

$$\begin{aligned} &= \sum_{x \in \mathbb{X}} \pi(x)P(x, y) = \sum_{x \in \mathbb{X}} \pi(x)P(y, x), \text{ by reversibility} \\ &= RHS. \end{aligned}$$

$$\begin{aligned} &= \sum_{x \in \mathbb{X}} \pi(x)P(x, y) = \sum_{x \in \mathbb{X}} \pi(x)P(y, x), \text{ by reversibility} \\ &= RHS. \end{aligned}$$

$$\begin{aligned} &= \sum_{x \in \mathbb{X}} \pi(x)P(x, y) = \sum_{x \in \mathbb{X}} \pi(x)P(y, x), \text{ by reversibility} \\ &= RHS. \end{aligned}$$

$$\begin{aligned} &= \sum_{x \in \mathbb{X}} \pi(x)P(x, y) = \sum_{x \in \mathbb{X}} \pi(x)P(y, x), \text{ by reversibility} \\ &= RHS. \end{aligned}$$

$$\begin{aligned} &= \sum_{x \in \mathbb{X}} \pi(x)P(x, y) = \sum_{x \in \mathbb{X}} \pi(x)P(y, x), \text{ by reversibility} \\ &= RHS. \end{aligned}$$

$$\begin{aligned} &= \sum_{x \in \mathbb{X}} \pi(x)P(x, y) = \sum_{x \in \mathbb{X}} \pi(x)P(y, x), \text{ by reversibility} \\ &= RHS. \end{aligned}$$

$$\begin{aligned} &= \sum_{x \in \mathbb{X}} \pi(x)P(x, y) = \sum_{x \in \mathbb{X}} \pi(x)P(y, x), \text{ by reversibility} \\ &= RHS. \end{aligned}$$

$$\begin{aligned} &= \sum_{x \in \mathbb{X}} \pi(x)P(x, y) = \sum_{x \in \mathbb{X}} \pi(x)P(y, x), \text{ by reversibility} \\ &= RHS. \end{aligned}$$

$$\begin{aligned} &= \sum_{x \in \mathbb{X}} \pi(x)P(x, y) = \sum_{x \in \mathbb{X}} \pi(x)P(y, x), \text{ by reversibility} \\ &= RHS. \end{aligned}$$

$$\begin{aligned} &= \sum_{x \in \mathbb{X}} \pi(x)P(x, y) = \sum_{x \in \mathbb{X}} \pi(x)P(y, x), \text{ by reversibility} \\ &= RHS. \end{aligned}$$

$$\begin{aligned} &= \sum_{x \in \mathbb{X}} \pi(x)P(x, y) = \sum_{x \in \mathbb{X}} \pi(x)P(y, x), \text{ by reversibility} \\ &= RHS. \end{aligned}$$

$$\begin{aligned} &= \sum_{x \in \mathbb{X}} \pi(x)P(x, y) = \sum_{x \in \mathbb{X}} \pi(x)P(y, x), \text{ by reversibility} \\ &= RHS. \end{aligned}$$

$$\begin{aligned} &= \sum_{x \in \mathbb{X}} \pi(x)P(x, y) = \sum_{x \in \mathbb{X}} \pi(x)P(y, x), \text{ by reversibility} \\ &= RHS. \end{aligned}$$

$$\begin{aligned} &= \sum_{x \in \mathbb{X}} \pi(x)P(x, y) = \sum_{x \in \mathbb{X}} \pi(x)P(y, x), \text{ by reversibility} \\ &= RHS. \end{aligned}$$

$$\begin{aligned} &= \sum_{x \in \mathbb{X}} \pi(x)P(x, y) = \sum_{x \in \mathbb{X}} \pi(x)P(y, x), \text{ by reversibility} \\ &= RHS. \end{aligned}$$

$$\begin{aligned} &= \sum_{x \in \mathbb{X}} \pi(x)P(x, y) = \sum_{x \in \mathbb{X}} \pi(x)P(y, x), \text{ by reversibility} \\ &= RHS. \end{aligned}$$

$$\begin{aligned} &= \sum_{x \in \mathbb{X}} \pi(x)P(x, y) = \sum_{x \in \mathbb{X}} \pi(x)P(y, x), \text{ by reversibility} \\ &= RHS. \end{aligned}$$

$$\begin{aligned} &= \sum_{x \in \mathbb{X}} \pi(x)P(x, y) = \sum_{x \in \mathbb{X}} \pi(x)P(y, x), \text{ by reversibility} \\ &= RHS. \end{aligned}$$

$$\begin{aligned} &= \sum_{x \in \mathbb{X}} \pi(x)P(x, y) = \sum_{x \in \mathbb{X}} \pi(x)P(y, x), \text{ by reversibility} \\ &= RHS. \end{aligned}$$

$$\begin{aligned} &= \sum_{x \in \mathbb{X}} \pi(x)P(x, y) = \sum_{x \in \mathbb{X}} \pi(x)P(y, x), \text{ by reversibility} \\ &= RHS. \end{aligned}$$

$$\begin{aligned} &= \sum_{x \in \mathbb{X}} \pi(x)P(x, y) = \sum_{x \in \mathbb{X}} \pi(x)P(y, x), \text{ by reversibility} \\ &= RHS. \end{aligned}$$

$$\begin{aligned} &= \sum_{x \in \mathbb{X}} \pi(x)P(x, y) = \sum_{x \in \mathbb{X}} \pi(x)P(y, x), \text{ by reversibility} \\ &= RHS. \end{aligned}$$

$$\begin{aligned} &= \sum_{x \in \mathbb{X}} \pi(x)P(x, y) = \sum_{x \in \mathbb{X}} \pi(x)P(y, x), \text{ by reversibility} \\ &= RHS. \end{aligned}$$

$$\begin{aligned} &= \sum_{x \in \mathbb{X}} \pi(x)P(x, y) = \sum_{x \in \mathbb{X}} \pi(x)P(y, x), \text{ by reversibility} \\ &= RHS. \end{aligned}$$

$$\begin{aligned} &= \sum_{x \in \mathbb{X}} \pi(x)P(x, y) = \sum_{x \in \mathbb{X}} \pi(x)P(y, x), \text{ by reversibility} \\ &= RHS. \end{aligned}$$

$$\begin{aligned} &= \sum_{x \in \mathbb{X}} \pi(x)P(x, y) = \sum_{x \in \mathbb{X}} \pi(x)P(y, x), \text{ by reversibility} \\ &= RHS. \end{aligned}$$

$$\begin{aligned} &= \sum_{x \in \mathbb{X}} \pi(x)P(x, y) = \sum_{x \in \mathbb{X}} \pi(x)P(y, x), \text{ by reversibility} \\ &= RHS. \end{aligned}$$

$$\begin{aligned} &= \sum_{x \in \mathbb{X}} \pi(x)P(x, y) = \sum_{x \in \mathbb{X}} \pi(x)P(y, x), \text{ by reversibility} \\ &= RHS. \end{aligned}$$

$$\begin{aligned} &= \sum_{x \in \mathbb{X}} \pi(x)P(x, y) = \sum_{x \in \mathbb{X}} \pi(x)P(y, x), \text{ by reversibility} \\ &= RHS. \end{aligned}$$

$$\begin{aligned} &= \sum_{x \in \mathbb{X}} \pi(x)P(x, y) = \sum_{x \in \mathbb{X}} \pi(x)P(y, x), \text{ by reversibility} \\ &= RHS. \end{aligned}$$

$$\begin{aligned} &= \sum_{x \in \mathbb{X}} \pi(x)P(x, y) = \sum_{x \in \mathbb{X}} \pi(x)P(y, x), \text{ by reversibility} \\ &= RHS. \end{aligned}$$

$$\begin{aligned} &= \sum_{x \in \mathbb{X}} \pi(x)P(x, y) = \sum_{x \in \mathbb{X}} \pi(x)P(y, x), \text{ by reversibility} \\ &= RHS. \end{aligned}$$

$$\begin{aligned} &= \sum_{x \in \mathbb{X}} \pi(x)P(x, y) = \sum_{x \in \mathbb{X}} \pi(x)P(y, x), \text{ by reversibility} \\ &= RHS. \end{aligned}$$

$$\begin{aligned} &= \sum_{x \in \mathbb{X}} \pi(x)P(x, y) = \sum_{x \in \mathbb{X}} \pi(x)P(y, x), \text{ by reversibility} \\ &= RHS. \end{aligned}$$

$$\begin{aligned} &= \sum_{x \in \mathbb{X}} \pi$$

Example 24 Consider a generic die-tossing experiment by a human experimenter. Here $\Omega = \{\omega_1, \omega_2, \omega_3, \dots, \omega_6\}$, but the experiment might correspond to rolling a die whose faces are:

1. sprayed with six different scents (nose!), or
2. studded with six distinctly flavoured candies (tongue!), or
3. contoured with six distinct bumps and pits (touch!), or
4. acoustically discernible at six different frequencies (ears!), or
5. painted with six different colours (eyes!), or
6. marked with six different numbers 1, 2, 3, 4, 5, 6 (eyes!), or , ...

These six experiments are equivalent as far as probability goes.

Definition 7 A **trial** is a single performance of an experiment and it results in an outcome.

Example 25 Some standard examples of a trial are:

- A roll of a die.
- A toss of a coin.
- A release of a chaotic double pendulum.

An experimenter often performs more than one trial. Repeated trials of an experiment forms the basis of science and engineering as the experimenter learns about the phenomenon by repeatedly performing the same mother experiment with possibly different outcomes. This repetition of trials in fact provides the very motivation for the definition of probability.

Definition 8 An **n-product experiment** is obtained by repeatedly performing n trials of some experiment. The experiment that is repeated is called the “mother” experiment.

Example 26 (Toss a coin n times) Suppose our experiment entails tossing a coin n times and recording H for Heads and T for Tails. When $n = 3$, one possible outcome of this experiment is HHT, ie. a Head followed by another Head and then a Tail. Seven other outcomes are possible.

The sample space for “toss a coin three times” experiment is:

$$\Omega = \{H, T\}^3 = \{HHH, HHT, HTH, HTT, THH, THT, TTH, TTT\},$$

with a particular sample point or outcome $\omega = HTH$, and another distinct outcome $\omega' = HHH$. An event, say A , that ‘at least two Heads occur’ is the following subset of Ω :

$$A = \{HHH, HHT, HTH, THH\}.$$

Another event, say B , that ‘no Heads occur’ is:

$$B = \{TTT\}$$

Note that the event B is also an outcome or sample point. Another interesting event is the empty set $\emptyset \subset \Omega$. The event that ‘nothing in the sample space occurs’ is \emptyset .

Proposition 94 (Markov chain convergence theorem) Let $(X_t)_{t \in \mathbb{Z}_+}$ be an irreducible aperiodic Markov chain with state space $\mathbb{X} = \{s_1, s_2, \dots, s_k\}$, transition matrix $P = (P(x, y))_{(x,y) \in \mathbb{X}^2}$ and initial distribution μ_0 . Then for any distribution π which is stationary for the transition matrix P , we have

$$\mu_t \xrightarrow{\text{TV}} \pi. \quad (6.11)$$

Proof: TBD

Proposition 95 (Uniqueness of stationary distribution) Any irreducible aperiodic Markov chain has a unique stationary distribution.

Proof: TBD

Exercise 181 Consider the Markov chain on $\{1, 2, 3, 4, 5, 6\}$ with the following transition matrix:

$$P = \begin{pmatrix} 1 & 2 & 3 & 4 & 5 & 6 \\ 1 & \frac{1}{2} & \frac{1}{2} & 0 & 0 & 0 \\ 2 & \frac{1}{2} & \frac{1}{2} & 0 & 0 & 0 \\ 3 & 0 & 0 & \frac{1}{4} & \frac{3}{4} & 0 \\ 4 & 0 & 0 & \frac{3}{4} & \frac{1}{4} & 0 \\ 5 & 0 & 0 & 0 & 0 & \frac{3}{4} \\ 6 & 0 & 0 & 0 & 0 & \frac{1}{4} \end{pmatrix}.$$

Show that this chain is reducible and it has three stationary distributions:

$$(1/2, 1/2, 0, 0, 0, 0), \quad (0, 0, 1/2, 1/2, 0, 0), \quad (0, 0, 0, 0, 1/2, 1/2).$$

Exercise 182 If there are two stationary distributions π and π' then show that there is a infinite family of stationary distributions $\{\pi_p : p \in [0, 1]\}$, called the convex combinations of π and π' .

Exercise 183 Show that for a drunkard’s walk chain started at state 0 around a polygonal block with k corners labelled $\{0, 1, 2, \dots, k-1\}$, the state probability vector at time step t

$$\mu_t \xrightarrow{\text{TV}} \pi$$

if and only if k is odd. Explain what happens to μ_t when k is even.

6.5 Reversibility

We introduce another specific property called reversibility. This property will assist in conjuring Markov chains with a desired stationary distribution.

Definition 96 (Reversible) A probability distribution π on $\mathbb{X} = \{s_1, s_2, \dots, s_k\}$ is said to be a **reversible distribution** for a Markov chain $(X_t)_{t \in \mathbb{Z}}$ on \mathbb{X} with transition matrix P if for every pair of states $(x, y) \in \mathbb{X}^2$:

$$\pi(x)P(x, y) = \pi(y)P(y, x). \quad (6.12)$$

A Markov chain that has a reversible distribution is said to be a reversible Markov chain.

$$N(H \cap T, n) = \frac{n}{H} = 1.$$

denoted by $H \cap T$. The probability that "something happens" is 1. More formally:

1. **Something Happens:** Each time we toss a coin, we are certain to observe Heads or Tails,

Other crucial assumptions that we have made here are:

of 0.1 of landing Heads. We might think that it is fair to have observed $N(H, n) \rightarrow 0.5$ as $n \rightarrow \infty$. We might, at least intuitively, think that the coin is unfair and has a lower "probability" $N(H, n) \rightarrow 0.1$ as million and found that this number approached closer to 0.1, or, more generally, $N(H, n) \rightarrow 0.1$ as 1000 tosses, then $N(H, 1000) = 9/1000 = 0.009$. Suppose we continued the number of tosses to a after conducting the tossing experiment 1000 times, we rarely observed Heads. Suppose that it n times and call $N(H, n)$ the fraction of times we observed Heads out of n tosses. We can toss fairness of a coin, i.e., if landing Heads has the same "probability" as landing Tails. We can toss the range of $dTV(v_1, v_2) = 0$ in [0, 1]. If $dTV(v_1, v_2) = 1$ then v_1 and v_2 have disjoint supports, i.e., we observe that if $dTV(v_1, v_2) = 0$ then $v_1 = v_2$. The constant $1/2$ in Equation 6.10 ensures that

fruitful record. In fact, you are here for exactly this reason.

The mathematical model for probability or the probability model is an axiomatic system that may the application of probability models to real-world problems through statistical experiments has a these axioms and definitions. No attempt to define probability in the real world is made. However, are intuitively motivated, the probability model simply follows from the application of logic to be motivated by the intuitive idea of long-term relative frequency. If the axioms and definitions be ranges of dTV is in $[0, 1]$. If $dTV(v_1, v_2) = 1$ then v_1 and v_2 have disjoint supports, i.e., we interpret this as $n \rightarrow \infty$ and write $v_1 \xrightarrow{dTV} v_2$, i.e.,

In words, Proposition 94 says that if you run the chain for a sufficient long enough time t , then, as $t \rightarrow \infty$.

This is referred to as the Markov chain **approaching equilibrium or stationarity** regardless of the initial distribution v_0 , the distribution at time t will be close to the stationary interpretation means that the total variation distance between v_1 and v_2 is the maximal difference in probabilities that the two distributions assign to any one event $A \in \sigma(\mathbb{X}) = 2^{\mathbb{X}}$.

This interpretation means that the total variation distance between v_1 and v_2 is the maximal interpretation of dTV is in $[0, 1]$. The total variation distance gets its name from the following natural

and $\sum_{x \in \mathbb{X}} v_2(x) = 1$. The total variation distance gets its name from the following natural interpretation:

$$dTV(v_1, v_2) = \max_{A \in \sigma(\mathbb{X})} |v_1(A) - v_2(A)|.$$

2.2 Probability

EXPERIMENT SUMMARY	
Trial	one performance of an experiment resulting in 1 outcome.
$A \subseteq \Omega$	a subset of Ω is an event.
ω	an individual outcome in Ω , called a simple event.
Event	set of all outcomes of the experiment.

Figure 2.1, using the caption as a hint (in other words, draw your own Figure 2.1). Enumerate the outcomes of the Experiment 26? Draw a diagram of this under the caption of Classwork 27 (A three-bifurcating tree of outcomes) Can you think of a graphical way to

Definition 93 (**Total variation distance**) If $v_1 := (v_1(x))_{x \in \mathbb{X}}$ and $v_2 := (v_2(x))_{x \in \mathbb{X}}$ are elements of $P(\mathbb{X})$, the set of all probability distributions on $\mathbb{X} := \{s_1, s_2, \dots, s_k\}$, then we define the total variation distance between v_1 and v_2 as

$$dTV(v_1, v_2) := \frac{1}{2} \sum_{x \in \mathbb{X}} |v_1(x) - v_2(x)|, \quad dTV : P(\mathbb{X}) \times P(\mathbb{X}) \rightarrow [0, 1]. \quad (6.10)$$

Proof: TBD

Proposition 92 (**Existence of Stationary distribution**) For any irreducible and aperiodic Markov chain there exists at least one stationary distribution.

$$\tau(x, y) < \infty.$$

and the mean hitting time is finite, i.e.,

$$\mathbf{P}(\tau(x, y) < \infty) = 1.$$

Proposition 91 (**Hitting times of irreducible aperiodic Markov chains**) If $(X_t)_{t \in \mathbb{Z}_+}$ is an irreducible aperiodic Markov chain with state space $\mathbb{X} = \{s_1, s_2, \dots, s_k\}$, transition matrix $P = (P(x, y))_{(x, y) \in \mathbb{X}^2}$ then for any pair of states $(x, y) \in \mathbb{X}^2$,

return time to state x .

be the expected time taken to reach y after having started at x . Note that $\tau(x, x)$ is the mean

$$\tau(x, y) := \mathbf{E}(\tau(x, y)),$$

the mean hitting time

and let $T(x, y) = \min\{t \mid Y_t = y\}$ be the Markov chain never visits y after having started from x . Let

Figure 2.1: A binary tree whose leaves are all possible outcomes.

This is an intuitively reasonable assumption that simply says that one of the possible outcomes is certain to occur, provided the coin is not so thick that it can land on or even roll along its circumference.

2. **Addition Rule:** Heads and Tails are mutually exclusive events in any given toss of a coin, i.e. they cannot occur simultaneously. The intersection of mutually exclusive events is the empty set and is denoted by $H \cap T = \emptyset$. The event $H \cup T$, namely that the event that “coin lands Heads **or** coin lands Tails” satisfies:

$$N(H \cup T, n) = N(H, n) + N(T, n).$$

3. The coin-tossing experiment is repeatedly performed in an **independent** manner, i.e. the outcome of any individual coin-toss does not affect that of another. This is an intuitively reasonable assumption since the coin has no memory and the coin is tossed identically each time.

We will use the LTRF idea more generally to motivate a mathematical model of probability called probability model. Suppose A is an event associated with some experiment \mathcal{E} , so that A either does or does not occur when the experiment is performed. We want the probability that event A occurs in a specific performance of \mathcal{E} , denoted by $\mathbf{P}(A)$, to intuitively mean the following: if one were to perform a super-experiment \mathcal{E}^∞ by independently repeating the experiment \mathcal{E} and recording $N(A, n)$, the fraction of times A occurs in the first n performances of \mathcal{E} within the super-experiment \mathcal{E}^∞ . Then the LTRF idea suggests:

$$N(A, n) := \frac{\text{Number of times } A \text{ occurs}}{n = \text{Number of performances of } \mathcal{E}} \rightarrow \mathbf{P}(A), \text{ as } n \rightarrow \infty \quad (2.1)$$

Now, we are finally ready to define probability.

Definition 10 (Probability) Let \mathcal{E} be an experiment with sample space Ω . Let \mathcal{F} denote a suitable collection of events in Ω that satisfy the following conditions:

1. It (the collection) contains the sample space: $[\Omega \in \mathcal{F}]$.
2. It is closed under complementation: $A \in \mathcal{F} \implies A^c \in \mathcal{F}$.
3. It is closed under countable unions: $A_1, A_2, \dots \in \mathcal{F} \implies \bigcup_i A_i := A_1 \cup A_2 \cup \dots \in \mathcal{F}$.

Formally, this collection of events is called a **sigma field** or a **sigma algebra**. Our experiment \mathcal{E} has a sample space Ω and a collection of events \mathcal{F} that satisfy the three condition.

Given a double, e.g. (Ω, \mathcal{F}) , **probability** is just a function \mathbf{P} which assigns each event $A \in \mathcal{F}$ a number $\mathbf{P}(A)$ in the real interval $[0, 1]$, i.e. $[\mathbf{P} : \mathcal{F} \rightarrow [0, 1]]$, such that:

1. The ‘Something Happens’ axiom holds, i.e. $[\mathbf{P}(\Omega) = 1]$.
2. The ‘Addition Rule’ axiom holds, i.e. for events A and B :

$$[A \cap B = \emptyset \implies \mathbf{P}(A \cup B) = \mathbf{P}(A) + \mathbf{P}(B)].$$

Proof: TBD

Exercise 179 (King’s random walk on a chessboard) Consider the squares in the chessboard as the state space $\mathbb{X} = \{0, 1, 2, \dots, 7\}^2$ with a randomly walking black king, i.e., for each move from current state $(u, v) \in \mathbb{X}$ the king chooses one of his $k(u, v)$ possible moves uniformly at random. Is the Markov chain corresponding to the randomly walking black king on the chessboard irreducible and/or aperiodic?

Exercise 180 (King’s random walk on a chesstorus) We can obtain a chesstorus from a pliable chessboard by identifying the eastern edge with the western edge (roll the chessboard into a cylinder) and then identifying the northern edge with the southern edge (gluing the top and bottom end of the cylinder together by turning into a doughnut or torus). Consider the squares in the chesstorus as the state space $\mathbb{X} = \{0, 1, 2, \dots, 7\}^2$ with a randomly walking black king, i.e., for each move from current state $(x, y) \in \mathbb{X}$ the king chooses one of his 8 possible moves uniformly at random according to the scheme: $X_t \leftarrow X_{t-1} + W_t$, where W_t is independent and identically distributed as follows:

$$\mathbf{P}(W_t = w) = \begin{cases} \frac{1}{8} & \text{if } w \in \{(1, 1), (1, 0), (1, -1), (0, -1), (-1, -1), (-1, 0), (-1, 1), (0, 1)\}, \\ 0 & \text{otherwise.} \end{cases}$$

Is the Markov chain corresponding to the randomly walking black king on the chesstorus irreducible and/or aperiodic? Write a MATLAB script to simulate a sequence of n states visited by the king if he started from $(0, 0)$ on the chesstorus.

6.4 Stationarity

We are interested in statements about a Markov chain that has been running for a long time. For any nontrivial Markov chain (X_0, X_1, \dots) the value of X_t will keep fluctuating in the state space \mathbb{X} as $t \rightarrow \infty$ and we cannot hope for convergence to a fixed point state $x^* \in \mathbb{X}$ or to a k -cycle of states $\{x_1, x_2, \dots, x_k\} \subset \mathbb{X}$. However, we can look one level up into the space of probability distributions over \mathbb{X} that give the probability of the Markov chain visiting each state $x \in \mathbb{X}$ at time t , and hope that the distribution of X_t over \mathbb{X} settles down as $t \rightarrow \infty$. The Markov chain convergence theorem indeed sates that the distribution of X_t over \mathbb{X} settles down as $t \rightarrow \infty$, provided the Markov chain is irreducible and aperiodic.

Definition 89 (Stationary distribution) Let $(X_t)_{t \in \mathbb{Z}_+}$ be a Markov chain with state space $\mathbb{X} = \{s_1, s_2, \dots, s_k\}$ and transition matrix $P = (P(x, y))_{(x,y) \in \mathbb{X}^2}$. A row vector

$$\pi = (\pi(s_1), \pi(s_2), \dots, \pi(s_k)) \in \mathbb{R}^{1 \times k}$$

is said to be a **stationary distribution** for the Markov chain, if it satisfies the conditions of being:

1. a **probability distribution**: $\pi(x) \geq 0$ for each $x \in \mathbb{X}$ and $\sum_{x \in \mathbb{X}} \pi(x) = 1$, and
2. a **fixed point**: $\pi P = \pi$, i.e., $\sum_{x \in \mathbb{X}} \pi(x)P(x, y) = \pi(y)$ for each $y \in \mathbb{X}$.

Definition 90 (Hitting times) If a Markov chain $(X_t)_{t \in \mathbb{Z}_+}$ with state space $\mathbb{X} = \{s_1, s_2, \dots, s_k\}$ and transition matrix $P = (P(x, y))_{(x,y) \in \mathbb{X}^2}$ starts at state x , then we can define the **hitting time**

$$T(x, y) = \min\{t \geq 1 : X_t = y\}.$$

6. For a sequence of mutually disjoint events $A_1, A_2, A_3, \dots, A_n$:

$$A_i \cap A_j = \emptyset \text{ for any } i \neq j \implies P(A_1 \cup A_2 \cup \dots \cup A_n) = P(A_1) + P(A_2) + \dots + P(A_n).$$

Proof: If A_1, A_2, A_3 are mutually disjoint events, then $A_1 \cup A_2$ is disjoint from A_3 . Thus, two applications of the addition rule for disjoint events yields:

$$P(A_1 \cup A_2 \cup A_3) = P((A_1 \cup A_2) \cup A_3) \stackrel{+ \text{ rule}}{=} P(A_1 \cup A_2) + P(A_3) \stackrel{+ \text{ rule}}{=} P(A_1) + P(A_2) + P(A_3)$$

The n -event case follows by mathematical induction.

We have formally defined the **probability model** specified by the **probability triple** (Ω, \mathcal{F}, P) that can be used to model an **experiment** \mathcal{E} .

Example 28 (First Ball out of NZ Lotto) Let us observe the number on *the first ball that pops out in a New Zealand Lotto trial*. There are forty balls labelled 1 through 40 for this experiment and so the sample space is

$$\Omega = \{1, 2, 3, \dots, 39, 40\}.$$

Because the balls are vigorously whirled around inside the Lotto machine, modelled as a well-stirred urn, before the first one pops out, we can model each ball to pop out first with the same probability. So, we assign each outcome $\omega \in \Omega$ the same probability of $\frac{1}{40}$, i.e., our probability model for this experiment is:

$$P(\omega) = \frac{1}{40}, \text{ for each } \omega \in \Omega = \{1, 2, 3, \dots, 39, 40\}.$$

Note: We sometimes abuse notation and write $P(\omega)$ instead of the more accurate but cumbersome $P(\{\omega\})$ when writing down probabilities of simple events.

Crucially, by $\omega = 17$ for example, we mean all the detailed dynamics inside the Lotto machine that lead to the event that the ball labelled by the number 17 ends up popping out. So, Ω here is indeed a more complicated set although it only leads to 40 possible outcomes.

Figure 2.2 (a) shows the frequency of the first ball number in 1114 NZ Lotto draws. Figure 2.2 (b) shows the relative frequency, i.e., the frequency divided by 1114, the number of draws. Figure 2.2 (b) also shows the equal probabilities under our model.

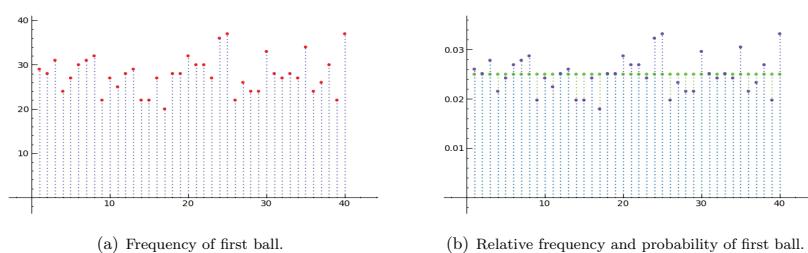


Figure 2.2: First ball number in 1114 NZ Lotto draws from 1987 to 2008.

Next, let us take a detour into how one might interpret it in the real world. The following is an adaptation from Williams D, *Weighing the Odds: A Course in Probability and Statistics*, Cambridge University Press, 2001, which henceforth is abbreviated as WD2001.

```
>> mu1 = mu0 * Phot % distribution for tomorrow since today is hot
mu1 =
    0.9500    0.0500
>> mu2 = mu1 * Pcold % distribution for day after tomorrow since tomorrow is supposed to be cold
mu2 =
    0.6400    0.3600
>> mu2 = mu0 * Phot * Pcold % we can also get the distribution for day after tomorrow directly
mu2 =
    0.6400    0.3600
```

Exercise 176 For the Markov chain in Example 175 compute the probability that the day after tomorrow is wet if today is dry and hot but tomorrow is supposed to be cold.

6.3 Irreducibility and Aperiodicity

The utility of our mathematical constructions with Markov chains depends on a delicate balance between generality and specificity. We introduce two specific conditions called irreducibility and aperiodicity that make Markov chains more useful to model real-word phenomena.

Definition 81 (Communication between states) Let $(X_t)_{t \in \mathbb{Z}_+}$ be a homogeneous Markov chain with transition matrix P on state space $\mathbb{X} := \{s_1, s_2, \dots, s_k\}$. We say that a state s_i **communicates** with a state s_j and write $s_i \rightarrow s_j$ or $s_j \leftarrow s_i$ if there exists an $\eta(s_i, s_j) \in \mathbb{N}$ such that:

$$P(X_{t+\eta(s_i, s_j)} = s_j | X_t = s_i) = P^{\eta(s_i, s_j)}(s_i, s_j) > 0.$$

In words, s_i communicates with s_j if you can eventually reach s_j from s_i . If $P^\eta(s_i, s_j) = 0$ for every $\eta \in \mathbb{N}$ then we say that s_i **does not communicate** with s_j and write $s_i \not\rightarrow s_j$ or $s_j \not\leftarrow s_i$.

We say that two states s_i and s_j **intercommunicate** and write $s_i \leftrightarrow s_j$ if $s_i \rightarrow s_j$ and $s_j \rightarrow s_i$. In words, two states intercommunicate if you can eventually reach one from another and vice versa. When s_i and s_j do not intercommunicate we write $s_i \not\leftrightarrow s_j$.

Definition 82 (Irreducible) A homogeneous Markov chain $(X_t)_{t \in \mathbb{Z}_+}$ with transition matrix P on state space $\mathbb{X} := \{s_1, s_2, \dots, s_k\}$ is said to be **irreducible** if $s_i \leftrightarrow s_j$ for each $(s_i, s_j) \in \mathbb{X}^2$. Otherwise the chain is said to be **reducible**.

We have already seen examples of reducible and irreducible Markov chains. For example, Flippant Freddy's family of Markov chains with the (p, q) -parametric family of transition matrices, $\{P_{(p,q)} : (p, q) \in [0, 1]^2\}$, where each $P_{(p,q)}$ is given by Equation 6.3. If $(p, q) \in (0, 1)^2$, then the corresponding Markov chain is irreducible because we can go from rollopia to flippopia or vice versa in just one step with a positive probability. Thus, the Markov chains with transition matrices in $\{P_{(p,q)} : (p, q) \in (0, 1)^2\}$ are irreducible. But if p or q take probability values at the boundary of $[0, 1]$, i.e., $p \in \{0, 1\}$ or $q \in \{0, 1\}$ then we have to be more careful because we may never get from at least one state to the other and the corresponding Markov chains may be reducible. For instance, if $p = 0$ or $q = 0$ then we will be stuck in either rollopia or flippopia, respectively. However, if $p = 1$ and $q \neq 0$ or $q = 1$ and $p \neq 0$ then we can get from each state to the other. Therefore, only the transition matrices in $\{P_{(p,q)} : p \in \{0\} \text{ or } q \in \{0\}\}$ are reducible.

The simplest way to verify whether a Markov chain is irreducible is by looking at its transition diagram (without the positive edge probabilities) and checking that from each state there is a sequence of arrows leading to any other state.

Exercise 177 Revisit all the Markov chains we have considered up to now and determine whether they are reducible or irreducible by checking that from each state there is a sequence of arrows leading to any other state in their transition graphs.

Classwork 30 (The trivial sigma algebra) Note that $\mathcal{F} = \{\emptyset, \Omega\}$ is also a sigma algebra of the sample space $\Omega = \{\text{H}, \text{T}\}$. Can you think of a probability for the collection \mathcal{F} ?

$\Omega = \{\text{H}, \text{T}\}$	\bullet	\emptyset
$\text{H} \bullet$	\leftarrow	$\frac{1}{2}$
$\text{T} \bullet$	\leftarrow	$1 - \frac{1}{2}$
$\text{Event } A \in \mathcal{F}$	\leftarrow	1
$\mathbf{P}: \mathcal{F} \rightarrow [0, 1]$		$\mathbf{P}(A) \in [0, 1]$

A function that will satisfy the definition of probability for this collection of events \mathcal{F} and assign for \mathbf{P} with arrows that map elements in the domain \mathcal{F} given above to elements in its range.

$$\Omega = \{\text{H}, \text{T}\}, \quad \mathcal{F} = \{\text{H}, \text{T}, \emptyset\},$$

is this sample space Ω and a reasonable collection of events \mathcal{F} that underpin this experiment? Consider the Toss a fair coin once experiment. What

2.2.2 Sigma Algebras of Typical Experiments*

$$\mathbf{P}(E) = \mathbf{P}(\{\omega_1, \omega_2, \dots, \omega_k\}) = \mathbf{P}\left(\bigcup_{i=1}^k \{\omega_i\}\right) = \sum_{i=1}^k \frac{1}{k} = \frac{k}{40}.$$

Let $E = \{\omega_1, \omega_2, \dots, \omega_k\}$ be an event with k outcomes (simple events). Then by the addition rule for mutually exclusive events we get:

$$\mathbf{P}(E) = \frac{k}{40} \times \text{number of elements in } E.$$

In the probability model of Example 28, show that for any event $E \subseteq \Omega$,

Events in Probability Model	Real-world Interpretation	Sample space Ω	The certain even, something happens	The impossible event, nothing happens	At least one of A and B occurs	All of the events A_1, A_2, \dots, A_n occurs	At least one of the events A_1, A_2, \dots, A_n occurs	A does not occur	A occurs, but B does not occur	If A occurs, then B must occur	A occurs, then B must occur	$A \cap B$	$A \setminus B$	$A \cup B$	$A_1 \cup A_2 \cup \dots \cup A_m$	$A_1 \cap A_2 \cap \dots \cap A_m$	$A_1 \cup A_2 \cup \dots \cup A_m$	$A_1 \cap A_2 \cap \dots \cap A_m$	$A \setminus B$	$A \cap B$
Possible outcome ω	Set of all outcomes of an experiment	Sample point ω	Possible outcome of an experiment	No outcome	Actual outcome ω^* of an experiment	Actual outcome of an experiment	Event A , a (stable) subset of Ω	The real-world event A	Both A and B occur	The impossible event, nothing happens	The certain even, something happens	The intersection $A \cap B$	The union $A \cup B$	The intersection $A_1 \cap A_2 \cap \dots \cap A_m$	The union $A_1 \cup A_2 \cup \dots \cup A_m$	The intersection $A_1 \cap A_2 \cap \dots \cap A_m$	The union $A_1 \cup A_2 \cup \dots \cup A_m$	The intersection $A \cap B$	$A \setminus B$	
Possible outcome ω	Set of all outcomes of an experiment	Sample space Ω	Possible outcome ω	No outcome	Actual outcome ω^* of an experiment	Actual outcome ω^* of an experiment	Event A , a (stable) subset of Ω	The real-world event A	Both A and B occur	The impossible event, nothing happens	The certain even, something happens	The intersection $A \cap B$	The union $A \cup B$	The intersection $A_1 \cap A_2 \cap \dots \cap A_m$	The union $A_1 \cup A_2 \cup \dots \cup A_m$	The intersection $A_1 \cap A_2 \cap \dots \cap A_m$	The union $A_1 \cup A_2 \cup \dots \cup A_m$	The intersection $A \cap B$	$A \setminus B$	
Possible outcome ω	Set of all outcomes of an experiment	Sample space Ω	Possible outcome ω	No outcome	Actual outcome ω^* of an experiment	Actual outcome ω^* of an experiment	Event A , a (stable) subset of Ω	The real-world event A	Both A and B occur	The impossible event, nothing happens	The certain even, something happens	The intersection $A \cap B$	The union $A \cup B$	The intersection $A_1 \cap A_2 \cap \dots \cap A_m$	The union $A_1 \cup A_2 \cup \dots \cup A_m$	The intersection $A_1 \cap A_2 \cap \dots \cap A_m$	The union $A_1 \cup A_2 \cup \dots \cup A_m$	The intersection $A \cap B$	$A \setminus B$	
Possible outcome ω	Set of all outcomes of an experiment	Sample space Ω	Possible outcome ω	No outcome	Actual outcome ω^* of an experiment	Actual outcome ω^* of an experiment	Event A , a (stable) subset of Ω	The real-world event A	Both A and B occur	The impossible event, nothing happens	The certain even, something happens	The intersection $A \cap B$	The union $A \cup B$	The intersection $A_1 \cap A_2 \cap \dots \cap A_m$	The union $A_1 \cup A_2 \cup \dots \cup A_m$	The intersection $A_1 \cap A_2 \cap \dots \cap A_m$	The union $A_1 \cup A_2 \cup \dots \cup A_m$	The intersection $A \cap B$	$A \setminus B$	

Probability $P(A)$, a number between 0 and 1 Probability that A will occur for an experiment yet to be performed

Probability $P(A)$, a number between 0 and 1 Probability that A will occur for an experiment yet to be performed

Stochastic matrices satisfy the conditions in Equation 6.2. Then, the stochastic sequence $(X_t)_{t \in \mathbb{Z}}^+ = (X_0, X_1, \dots)$ with finite state space $\mathbb{X} = \{s_1, s_2, \dots, s_k\}$ is called an homogeneous Markov chain with transition matrices P_1, P_2, \dots ; it for all pairs of states $(x, y) \in \mathbb{X} \times \mathbb{X}$, all integers $t \geq 1$, and all probable historical events $H_{t-1} = \bigcup_{n=0}^{t-1} \{X_n = x_n\}$ with $\mathbf{P}(H_{t-1} \cup \{X_t = x\}) < 0$,

the following **Markov Property** is satisfied:

```

m0 = 1
0
>>> today_is_dry
0.4500 0.5500
0.9500 0.0500
0.7500 0.2500
0.6500 0.3500
Pcold = [0.65 0.35; 0.45 0.55] % Transition Probability Matrix for cold day
Pheat = [0.95 0.05; 0.75 0.25] % Transition Probability Matrix for hot day
>>> Pheat = [0.95 0.05; 0.75 0.25] % Transition Probability Matrix for hot day
>>> Pcold = [0.65 0.35; 0.45 0.55] % Transition Probability Matrix for cold day
We say that a day is hot if its maximum temperature is more than 20°C. Otherwise it is cold. We use the transition matrix for today to obtain the state probabilities for tomorrow. If today is dry and hot tomorrow is supposed to be cold then what is the probability that the day after tomorrow will be wet? We can use (6.9) to obtain the answer as 0.36:

```

```

Pheat = d w
d   w
Pcold = d w
d   w
>>> Pheat = [0.95 0.05; 0.75 0.25]
>>> Pcold = [0.65 0.35; 0.45 0.55]

```

We say that a day is hot if its maximum temperature is more than 20°C. Otherwise it is cold. We use the transition matrix for today to obtain the state probabilities for tomorrow. If today is dry and hot tomorrow is supposed to be cold then what is the probability that the day after tomorrow will be wet? We can use (6.9) to obtain the answer as 0.36:

Exercise 174 Prove Proposition 80 using induction as done for Proposition 78.

Proof: Left as Exercise 174.

$$u_t = u_0 P_1 P_2 \dots P_t. \quad (6.9)$$

where $u_t(s_i) = \mathbf{P}(X_t = s_i)$, satisfies:

$$u_t = (u_t(s_1), u_t(s_2), \dots, u_t(s_k)),$$

we have for any $t \in \mathbb{Z}^+$ that the distribution at time t given by:

$$(P_1, P_2, \dots), \quad P_t = (P(s_i, s_j))_{(s_i, s_j) \in \mathbb{X} \times \mathbb{X}}, \quad t \in \{1, 2, \dots\}$$

where $u_0(s_i) = \mathbf{P}(X_0 = s_i)$, and transition matrices

$$u_0 = (u_0(s_1), u_0(s_2), \dots, u_0(s_k)),$$

Proposition 80 For a finite homogeneous Markov chain $(X_t)_{t \in \mathbb{Z}}^+$ with state space $\mathbb{X} = \{s_1, s_2, \dots, s_k\}$,

$$\mathbf{P}(X_{t+1} = y | H_{t-1} \cup \{X_t = x\}) = \mathbf{P}(X_{t+1} = y | X_t = x) =: P_{t+1}(x, y). \quad (6.8)$$

the following **Markov Property** is satisfied:

Definition 79 (Inhomogeneous Markov chain) Let P_1, P_2, \dots be a sequence of $k \times k$ stochastic matrices satisfying the conditions in Equation 6.2. Then, the stochastic sequence $(X_t)_{t \in \mathbb{Z}}^+ = (X_0, X_1, \dots)$ with finite state space $\mathbb{X} = \{s_1, s_2, \dots, s_k\}$ is called an inhomogeneous Markov chain with transition matrices P_1, P_2, \dots ; it for all pairs of states $(x, y) \in \mathbb{X} \times \mathbb{X}$, all integers $t \geq 1$, and all probable historical events $H_{t-1} = \bigcup_{n=0}^{t-1} \{X_n = x_n\}$ with $\mathbf{P}(H_{t-1} \cup \{X_t = x\}) < 0$,

Event $A \in \mathcal{F}'$	$\mathbf{P} : \mathcal{F}' \rightarrow [0, 1]$	$\mathbf{P}(A) \in [0, 1]$
$\Omega = \{\text{H, T}\}$ •	→	
\emptyset •	→	

Thus, \mathcal{F} and \mathcal{F}' are two distinct sigma algebras over our $\Omega = \{\text{H, T}\}$. Moreover, $\mathcal{F}' \subset \mathcal{F}$ and is called a sub sigma algebra. Try to show that $\{\Omega, \emptyset\}$ is the smallest possible sigma algebra over all possible sigma algebras over any given sample space Ω (think of intersecting an arbitrary family of sigma algebras)?

Generally one encounters four types of sigma algebras (you will understand the last two types after taking more advanced courses in mathematics, so it is fine to understand the ideas intuitively for now!) and they are:

- When the sample space $\Omega = \{\omega_1, \omega_2, \dots, \omega_k\}$ is a finite set with k outcomes and $\mathbf{P}(\omega_i)$, the probability for each outcome $\omega_i \in \Omega$ is known, then one typically takes the sigma-algebra \mathcal{F} to be the set of all subsets of Ω called the **power set** and denoted by 2^Ω . The probability of each event $A \in 2^\Omega$ can be obtained by adding the probabilities of the outcomes in A , i.e., $\mathbf{P}(A) = \sum_{\omega_i \in A} \mathbf{P}(\omega_i)$. Clearly, 2^Ω is indeed a sigma-algebra and it contains $2^{\#\Omega}$ events in it.
- When the sample space $\Omega = \{\omega_1, \omega_2, \dots\}$ is a countable set then one typically takes the sigma-algebra \mathcal{F} to be the set of all subsets of Ω . Note that this is very similar to the case with finite Ω except now $\mathcal{F} = 2^\Omega$ could have uncountably many events in it.
- If $\Omega = \mathbb{R}^d$ for finite $d \in \{1, 2, 3, \dots\}$ then the **Borel sigma-algebra** is the smallest sigma-algebra containing all **half-spaces**, i.e., sets of the form

$$\{x = (x_1, x_2, \dots, x_d) \in \mathbb{R}^d : x_1 \leq c_1, x_2 \leq c_2, \dots, x_d \leq c_d\}, \quad \text{for any } c = (c_1, c_2, \dots, c_d) \in \mathbb{R}^d,$$

When $d = 1$ the half-spaces are the half-lines $\{(-\infty, c] : c \in \mathbb{R}\}$ and when $d = 2$ the half-spaces are the south-west quadrants $\{(-\infty, c_1] \times (-\infty, c_2] : (c_1, c_2) \in \mathbb{R}^2\}$, etc. (Equivalently, the Borel sigma-algebra is the smallest sigma-algebra containing all open sets in \mathbb{R}^d).

- Given a finite set $\mathbb{S} = \{s_1, s_2, \dots, s_k\}$, let Ω be the sequence space $\mathbb{S}^\infty := \mathbb{S} \times \mathbb{S} \times \mathbb{S} \times \dots$, i.e., the set of sequences of infinite length that are made up of elements from \mathbb{S} . A set of the form

$$A_1 \times A_2 \times \dots \times A_n \times \mathbb{S} \times \mathbb{S} \times \dots, \quad A_k \subset \mathbb{S} \text{ for all } k \in \{1, 2, \dots, n\},$$

is called a **cylinder set**. The set of events in \mathbb{S}^∞ is the smallest sigma-algebra containing the cylinder sets.

- **A most primitive sigma-algebra for probability theory:** For example if $\mathbb{S} = \{0, 1\}$, then $\Omega = \{0, 1\}^\infty$ is the set of all infinite sequences made of 0's and 1's. To take advantage of arithmetic and analysis, Ω can be seen as the binary representation of all real numbers in the unit interval $[0, 1]$. We can take advantage of combinatorics and algebra if we further represent the dyadic partition of $[0, 1]$ by a binary tree (as drawn in lectures). Then, a cylinder set such as $1 \times 1 \times 0 \times \{0, 1\} \times \{0, 1\} \times \dots$, an event here, can be interpreted as the finite binary sequence $(1, 1, 0)$ — corresponding to the third leaf of a finite binary tree with four leaves obtained by splitting the right-most leaf twice. This cylindrical event $(1, 1, 0)$ contains all real numbers in the interval $[\frac{3}{4}, \frac{7}{8}] \subset [0, 1] =: \Omega$.

Exercise 2.1 (Intuiting a most primitive sigma-algebra – this is optional) Try to carefully recollect and understand the most primitive sigma-algebra in the last item above as it was explained in lectures.

Suppose you sell \$100 of lemonade at a road-side stand on a hot day but only \$50 on a cold day. Then we can compute your expected sales tomorrow if today is dry as follows:

```
>> P=[0.75 0.25; 0.5 0.5] % Transition Probability Matrix
P =
    0.7500    0.2500
    0.5000    0.5000
>> f = [100; 50] % sales of lemonade in dollars on a dry and wet day
f =
    100
    50
>> P*f % expected sales tomorrow
ans =
    87.5000
    75.0000
>> mu0 = [1 0] % today is dry
mu0 =
    1
    0
>> mu0*P*f % expected sales tomorrow if today is dry
ans =
    87.5000
```

Exercise 171 (Freddy discovers a gold coin) Flippant Freddy of Example 169 found a gold coin at the bottom of the pond. Since this discovery he jumps around differently in the enchanted pond. He can be found now in one of three states: flipopia, rollopia and hydropia (when he dives into the pond). His state space is $\mathbb{X} = \{r, f, h\}$ now and his transition mechanism is as follows: If he rolls an odd number with his fair die in rollopia he will jump to flipopia but if he rolls an even number then he will stay in rollopia only if the outcome is 2 otherwise he will dive into hydropia. If the fair gold coin toss at the bottom of hydropia is Heads then Freddy will swim to flipopia otherwise he will remain in hydropia. Finally, if he is in flipopia he will remain there if the silver coin lands Heads otherwise he will jump to rollopia.

Make a Markov chain model of the new jumping mechanism adopted by Freddy. Draw the transition diagram, produce the transition matrix P and compute using MATLAB the probability that Freddy will be in hydropia after one, two, three, four and five jumps given that he starts in hydropia.

Exercise 172 Let $(X_t)_{t \in \mathbb{Z}_+}$ be a Markov chain with state space $\{a, b, c\}$, initial distribution $\mu_0 = (1/3, 1/3, 1/3)$ and transition matrix

$$P = \begin{pmatrix} a & b & c \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ c & 1 & 0 \end{pmatrix}.$$

For each t , define $Y_t = \mathbb{1}_{\{b,c\}}(X_t)$. Show that $(Y_t)_{t \in \mathbb{Z}_+}$ is not a Markov chain.

Exercise 173 Let $(X_t)_{t \in \mathbb{Z}_+}$ be a (homogeneous) Markov chain on $\mathbb{X} = \{s_1, s_2, \dots, s_k\}$ with transition matrix P and initial distribution μ_0 . For a given $m \in \mathbb{N}$, let $(Y_t)_{t \in \mathbb{Z}_+}$ be a stochastic sequence with $Y_t = X_{mt}$. Show that $(Y_t)_{t \in \mathbb{Z}_+}$ is a Markov chain with transition matrix P^m . This establishes that Markov chains that are sampled at regular time steps are also Markov chains.

Until now our Markov chains have been **homogeneous** in time according to Definition 77, i.e., the transition matrix P does not change with time. We define inhomogeneous Markov chains that allow their transition matrices to possibly change with time. Such Markov chains are more realistic as models in some situations and more flexible as algorithms in the sequel.

Ex. 2.5 — There are seventy five balls in total inside the Bingo Machine. Each ball is labelled by one of the following five letters: B, I, N, G, and O. There are fifteen balls labelled by each letter. The letter on the first ball that comes out of a BINGO machine after it has been well-mixed is the outcome of our experiment.

- (a) Write down the sample space of this experiment.
- (b) Find the probabilities of each simple event.
- (c) Show that $\mathbf{P}(\Omega)$ is indeed 1.
- (d) Check that the addition rule for mutually exclusive events holds for the simple events $\{B\}$ and $\{I\}$.
- (e) Consider the following events: $C = \{B, I, G\}$ and $D = \{G, I, N\}$. Using the addition rule for two arbitrary events, find $\mathbf{P}(C \cup D)$.

2.4 Conditional Probability

Conditional probabilities arise when we have partial information about the result of an experiment which restricts the sample space to a range of outcomes. For example, if there has been a lot of recent seismic activity in Christchurch, then the probability that an already damaged building will collapse tomorrow is clearly higher than if there had been no recent seismic activity.

Conditional probabilities are often expressed in English by phrases such as:

“If A happens, what is the probability that B happens?”

or

“What is the probability that A happens if B happens?”

or

“What is the probability that A occurs given that B occurs?”

Next, we define conditional probability and the notion of independence of events. We use the LTRF idea to motivate the definition.

Idea 11 (LTRF intuition for conditional probability) Let A and B be any two events associated with our experiment \mathcal{E} with $\mathbf{P}(A) \neq 0$. The ‘conditional probability that B occurs given that A occurs’ denoted by $\mathbf{P}(B|A)$ is again intuitively underpinned by the super-experiment \mathcal{E}^∞ which is the ‘independent’ repetition of our original experiment \mathcal{E} ‘infinitely’ often. The LTRF idea is that $\mathbf{P}(B|A)$ is the long-term proportion of those experiments on which A occurs that B also occurs.

Recall that $N(A, n)$ as defined in (2.1) is the fraction of times A occurs out of n independent repetitions of our experiment \mathcal{E} (ie. the experiment \mathcal{E}^n). If $A \cap B$ is the event that ‘ A and B occur simultaneously’, then we intuitively want

$$\mathbf{P}(B|A) \quad “\rightarrow” \quad \frac{N(A \cap B, n)}{N(A, n)} = \frac{N(A \cap B, n)/n}{N(A, n)/n} = \frac{\mathbf{P}(A \cap B)}{\mathbf{P}(A)}$$

as our $\mathcal{E}^n \rightarrow \mathcal{E}^\infty$. So, we **define** conditional probability as we want.

Now let us generalise the lessons learned from Example 169.

Proposition 78 For a finite Markov chain $(X_t)_{t \in \mathbb{Z}_+}$ with state space $\mathbb{X} = \{s_1, s_2, \dots, s_k\}$, initial distribution

$$\mu_0 := (\mu_0(s_1), \mu_0(s_2), \dots, \mu_0(s_k)),$$

where $\mu_0(s_i) = \mathbf{P}(X_0 = s_i)$, and transition matrix

$$P := (P(s_i, s_j))_{(s_i, s_j) \in \mathbb{X}^2},$$

we have for any $t \in \mathbb{Z}_+$ that the distribution at time t given by:

$$\mu_t := (\mu_t(s_1), \mu_t(s_2), \dots, \mu_t(s_k)),$$

where $\mu_t(s_i) = \mathbf{P}(X_t = s_i)$, satisfies:

$$\mu_t = \mu_0 P^t. \quad (6.7)$$

Proof: We will prove this by induction on $\mathbb{Z}_+ := \{0, 1, 2, \dots\}$. First consider the case when $t = 0$. Since P^0 is the identity matrix I , we get the desired equality:

$$\mu_0 P^0 = \mu_0 I = \mu_0.$$

Next consider the case when $t = 1$. We get for each $j \in \{1, 2, \dots, k\}$, that

$$\begin{aligned} \mu_1(s_j) &= \mathbf{P}(X_1 = s_j) = \sum_{i=1}^k \mathbf{P}(X_1 = s_j, X_0 = s_i) \\ &= \sum_{i=1}^k \mathbf{P}(X_1 = s_j | X_0 = s_i) \mathbf{P}(X_0 = s_i) \\ &= \sum_{i=1}^k P(s_i, s_j) \mu_0(s_i) \\ &= (\mu_0 P)(s_j), \quad \text{the } j\text{-th entry of the row vector } (\mu_0 P). \end{aligned}$$

Hence, $\mu_1 = \mu_0 P$. Now, we will fix m and suppose that (6.7) holds for $t = m$ and prove that (6.7) also holds for $t = m + 1$. For each $j \in \{1, 2, \dots, k\}$, we get

$$\begin{aligned} \mu_{m+1}(s_j) &= \mathbf{P}(X_{m+1} = s_j) = \sum_{i=1}^k \mathbf{P}(X_{m+1} = s_j, X_m = s_i) \\ &= \sum_{i=1}^k \mathbf{P}(X_{m+1} = s_j | X_m = s_i) \mathbf{P}(X_m = s_i) \\ &= \sum_{i=1}^k P(s_i, s_j) \mu_m(s_i) \\ &= (\mu_m P)(s_j), \quad \text{the } j\text{-th entry of the row vector } (\mu_m P). \end{aligned}$$

Hence, $\mu_{m+1} = \mu_m P$. But $\mu_m = \mu_0 P^m$ by the induction hypothesis, and therefore:

$$\mu_{m+1} = \mu_m P = \mu_0 P^m P = \mu_0 P^{m+1}.$$

Thus by the principle of mathematical induction we have proved the proposition.

	Have Disease (D)	Don't have disease (D^c)
Test positive (+)	0.009	0.99
Test negative (-)	0.001	0.991

Example 31 (Wasserman3, p. 11) A medical test for a disease D has outcomes + and - . the probabilities are:

$$P(A \cup B) = P(A)P(B|A) = P(B)P(A|B).$$

If A and B are events, and if $P(A) \neq 0$ and $P(B) \neq 0$, then

Multiplication rule for two likely events:

Solving for $P(A \cup B)$ with these definitions of conditional probability gives another rule:

$$P(B_1 \cup B_2 | A) = P(B_1 | A) + P(B_2 | A) - P(B_1 \cap B_2 | A).$$

Addition rule for two arbitrary events B_1 and B_2 :

$$\text{Complement rule: } P(B|A) = 1 - P(B^c|A).$$

From the definition of conditional probability we get the following properties or rules:

$$P(B_1 \cup B_2 \cup \dots | A) = P(B_1 | A) + P(B_2 | A) + \dots.$$

Axiom (4): For mutually exclusive events, B_1, B_2, \dots ,

$$B_1 \cup B_2 = \emptyset \text{ implies } P(B_1 \cup B_2 | A) = P(B_1 | A) + P(B_2 | A).$$

Axiom (3): The Addition Rule axiom holds, ie. for events $B_1, B_2 \in \mathcal{F}$,

Axiom (2): $P(Q|A) = 1$ Meaning, Something Happens given the event A happens

Axiom (1): For any event B , $0 \leq P(B|A) \leq 1$.

satisfied:

$$P(B|A) : \mathcal{F} \rightarrow [0, 1]$$

Note that A serves as the new reduced sample space so that conditional probabilities given A are indeed probabilities. Thus, for a fixed event $A \in \mathcal{F}$ with $P(A) > 0$ and any event $B \in \mathcal{F}$, the conditional probability $P(B|A)$ is a probability as in Definition 10, ie. a function:

that assigns to each $B \in \mathcal{F}$ a number in the interval $[0, 1]$, such that, the axioms of probability are satisfied:

$$P(B|A) = \frac{P(A \cap B)}{P(A)}. \quad (2.2)$$

conditional probability of B given A by,

Definition 12 (Conditional Probability) Suppose we are given an experiment \mathcal{E} with a triple (Ω, \mathcal{F}, P) . Let A and B be events, ie. $A, B \in \mathcal{F}$, such that $P(A) \neq 0$. Then, we define the

three values of p and q according to the following script:

$$\begin{aligned} \text{(iii) } p = 0.15, q = 0.95, \quad P(X_i = r) \leftarrow \pi(r) = \frac{d}{q} = 0.15 + 0.95 = 0.8636. \\ \text{(ii) } p = 0.85, q = 0.35, \quad P(X_i = r) \leftarrow \pi(r) = \frac{d}{q} = 0.85 + 0.35 = 0.2917, \\ \text{(i) } p = 0.50, q = 0.50, \quad P(X_i = r) \leftarrow \pi(r) = \frac{d}{q} = 0.50 + 0.50 = 0.5000, \end{aligned}$$

In Figure 6.2 we see that $P(X_i = r)$ approaches $\pi(r) = \frac{d}{p+q}$ for the three cases of p and q :

$$\pi(r) = \frac{d}{p+q}, \quad \pi(f) = \frac{d}{p+q}.$$

that gives the solution:

$$\pi P = \pi,$$

fixed point condition:

It is evident from Figure 6.2 that as $t \rightarrow \infty$, π_t approaches a stationary distribution, say π , that depends on p and q in P . Such a limit distribution is called the **stationary distribution** and must satisfy the

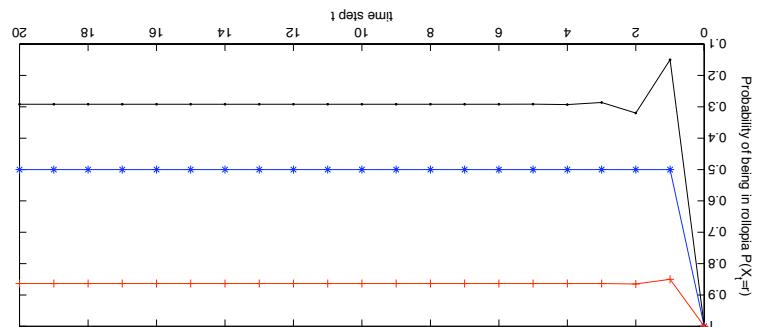


Figure 6.2: The probability of being back in roulette in t time steps after having started there under transition matrix P with (i) $p = q = 0.5$ (blue line with asterisks), (ii) $p = 0.85, q = 0.35$ (black line with dots) and (iii) $p = 0.15, q = 0.95$ (red line with pluses).

```

p=0.5; q=0.95; P = [1-p; p; q-1-q]; % assume an unfair coin and another unfair die
t=0:1:20; % vector of time steps t
mu0 = [1, 0]; % initial state vector since Freddy started in roulette
mu0 = [1, 0]; % assume a fair coin and a fair die
for t = 1:1:21, mu(t,:)= mu0*p; end
plot(t,mu(:,1),'r','*-');
plot(t,mu(:,2),'b','*--');
plot(t,mu(:,3),'k','.-');
for t = 1:1:21, mu(t,:)= mu0*p; % assume another unfair coin and another unfair die
mu0 = [1, 0]; % initial state vector since Freddy started in roulette
mu0 = [1, 0]; % assume a fair coin and a fair die
for t = 1:1:21, mu(t,:)= mu0*p; end
plot(t,mu(:,1),'r','*-');
plot(t,mu(:,2),'b','*--');
plot(t,mu(:,3),'k','.-');
for t = 1:1:21, mu(t,:)= mu0*p; % assume an unfair coin and an unfair die
mu0 = [1, 0]; % initial state vector since Freddy started in roulette
mu0 = [1, 0]; % assume a fair coin and a fair die
for t = 1:1:21, mu(t,:)= mu0*p; end
plot(t,mu(:,1),'r','*-');
plot(t,mu(:,2),'b','*--');
plot(t,mu(:,3),'k','.-');
xlabel('time step t'), ylabel('Probability of being in roulette B(X_{t+1})')

```

Using the definition of conditional probability, we can compute the conditional probability that you test positive given that you have the disease:

$$\mathbf{P}(+|D) = \frac{\mathbf{P}(+ \cap D)}{\mathbf{P}(D)} = \frac{0.009}{0.009 + 0.001} = 0.9 ,$$

and the conditional probability that you test negative given that you don't have the disease:

$$\mathbf{P}(-|D^c) = \frac{\mathbf{P}(- \cap D^c)}{\mathbf{P}(D^c)} = \frac{0.891}{0.099 + 0.891} \approx 0.9 .$$

Thus, the test is quite accurate since sick people test positive 90% of the time and healthy people test negative 90% of the time.

Now, suppose you go for a test and test positive. What is the probability that you have the disease?

$$\mathbf{P}(D|+) = \frac{\mathbf{P}(D \cap +)}{\mathbf{P}(+)} = \frac{0.009}{0.009 + 0.099} \approx 0.08$$

Most people who are not used to the definition of conditional probability would intuitively associate a number much bigger than 0.08 for the answer. Interpret conditional probability in terms of the meaning of the numbers that appear in the numerator and denominator of the above calculations.

2.4.1 Bayes' Theorem

Next we look at one of the most elegant applications of the definition of conditional probability along with the addition rule for a partition of Ω called *Bayes' Theorem*. We will present a two event case first called *Bayes' Rule* and then present the more general case of the Theorem.

This is useful because many problems involve reversing the order of conditional probabilities. Suppose we want to investigate some phenomenon A and have an observation B that is evidence about A : for example, A may be breast cancer and B may be a positive mammogram. Then Bayes' Theorem tells us how we should update our probability of A , given the new evidence B .

Or, put more simply, Bayes' Rule is useful when you know $P(B|A)$ but want $P(A|B)$!

Proposition 13 (Bayes' Rule)

$$P(A|B) = \frac{P(A)P(B|A)}{P(B)} . \quad (2.3)$$

Proof: From the definition of conditional probability and the multiplication rule for two likely events A and B we get:

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

Example 32 (Mammogram) Approximately 1% of women aged 40–50 have breast cancer. A woman with breast cancer has a 90% chance of a positive test from a mammogram, while a woman without breast cancer has a 10% chance of a false positive result from the test. What is the probability that a woman indeed has breast cancer given that she just had a positive test?

Definition of conditional probability and the Markov property, we see that,

$$\begin{aligned} \mathbf{P}(X_2 = f|X_0 = r) &= \mathbf{P}(X_2 = f, X_1 = f|X_0 = r) + \mathbf{P}(X_2 = f, X_1 = r|X_0 = r) \\ &= \frac{\mathbf{P}(X_2 = f, X_1 = f, X_0 = r)}{\mathbf{P}(X_0 = r)} + \frac{\mathbf{P}(X_2 = f, X_1 = r, X_0 = r)}{\mathbf{P}(X_0 = r)} \\ &= \mathbf{P}(X_2 = f|X_1 = f, X_0 = r) \frac{\mathbf{P}(X_1 = f, X_0 = r)}{\mathbf{P}(X_0 = r)} \\ &\quad + \mathbf{P}(X_2 = f|X_1 = r, X_0 = r) \frac{\mathbf{P}(X_1 = r, X_0 = r)}{\mathbf{P}(X_0 = r)} \\ &= \mathbf{P}(X_2 = f|X_1 = f, X_0 = r) \mathbf{P}(X_1 = f|X_0 = r) \\ &\quad + \mathbf{P}(X_2 = f|X_1 = r, X_0 = r) \mathbf{P}(X_1 = r|X_0 = r) \\ &= \mathbf{P}(X_2 = f|X_1 = f) \mathbf{P}(X_1 = f|X_0 = r) \\ &\quad + \mathbf{P}(X_2 = f|X_1 = r) \mathbf{P}(X_1 = r|X_0 = r) \\ &= P(f, f)P(r, f) + P(r, f)P(r, r) \\ &= (1 - q)p + p(1 - p) \end{aligned} \quad (6.5)$$

Similarly,

$$\mathbf{P}(X_2 = r|X_0 = r) = P(f, r)P(r, f) + P(r, r)P(r, r) = qp + (1 - p)(1 - p) \quad (6.6)$$

Instead of elaborate computations of the probabilities of being in a given state after Freddy's t -th restless moment, we can store the state probabilities at time t in a row vector:

$$\mu_t := (\mathbf{P}(X_t = r|X_0 = r), \mathbf{P}(X_t = f|X_0 = r)) ,$$

Now, we can conveniently represent Freddy starting in rollovia by the **initial distribution** $\mu_0 = (1, 0)$ and obtain the 1-step **state probability vector** in (6.4) from $\mu_1 = \mu_0 P$ and the 2-step state probabilities in (6.5) and (6.6) by $\mu_2 = \mu_1 P = \mu_0 P P = \mu_0 P^2$. In general, multiplying μ_t , the state probability vector at time t , by the transition matrix P on the right updates the state probabilities by another step:

$$\mu_t = \mu_{t-1} P \quad \text{for all } t \geq 1 .$$

And for any initial distribution μ_0 ,

$$\mu_t = \mu_0 P^t \quad \text{for all } t \geq 0 .$$

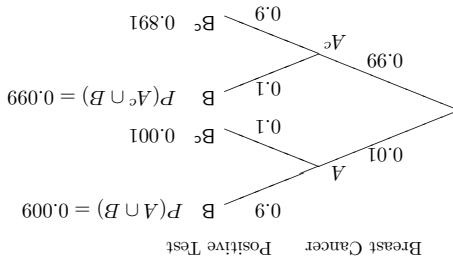
This can be easily implemented in MATLAB as follows:

```
>> p=0.85; q=0.35; P = [1-p p; q 1-q] % assume an unfair coin and an unfair die
P =
    0.1500    0.8500
    0.3500    0.6500
>> mu0 = [1, 0] % initial state vector since Freddy started in rollovia
mu0 =
    1         0
>> mu0*P^0 % initial state distribution at t=0 is just mu0
ans =
    1         0
>> mu0*P^1 % state distribution at t=1
ans =
    0.1500    0.8500
>> mu0*P^2 % state distribution at t=2
ans =
    0.3200    0.6800
>> mu0*P^3 % state distribution at t=3
ans =
    0.2860    0.7140
```

* In the exam, there won't be any need for electronic calculators and you may leave the steps in your reasoning.

$$P(A|B) = \frac{P(B)}{P(A \cup B)} = \frac{0.009 + 0.099}{0.009} = \frac{9}{108}$$

So the probability that a woman has breast cancer given that she has just had a positive test is



Alternative solution using a tree diagram:

$$\text{aritive the answer } 9/(9+99).$$

This answer is somewhat surprising. Indeed when ninety-five physicians were asked this question they average answer was 75%. The two statisticians who carried out this survey indicated that physicians were better able to see the answer when the data was presented in frequency format. 10 out of 1000 women have breast cancer. Of these 9 will have a positive mammogram. However, of the remaining 990 women without breast cancer 99 will have a positive reaction, and again we arrive at the answer 9/(9+99).

for a healthy woman, which has probability 0.009.

$$P(A|B) = \frac{P(B)}{P(A \cup B)} = \frac{0.009 + 0.099}{0.009} = \frac{9}{108}$$

Now $P(B) = P(A \cup B) + P(A^c \cup B)$ so

$$P(A^c \cup B) = P(A^c)P(B|A^c) = 0.99 \times 0.1 = 0.099$$

Similarly,

$$P(A \cup B) = P(A)P(B|A) = 0.01 \times 0.9 = 0.009$$

To evaluate the numerator we use the multiplication rule

$$P(A|B) = P(A \cup B)/P(B)$$

By the definition of conditional probability,

$$\text{We want } P(A|B) \text{ but what we are given is } P(B|A) = 0.9.$$

Let A = "the woman has breast cancer", and B = "a positive test."

Solution:

What happens when he is restless for the second time? By considering the two possibilities for X_1 , we know from the first row of P that he will leave to Hippopota with probability p and stay with Hippopota in rollopia, i.e., $X_0 = r$. When he gets restless for the first time suppose we first see Freddy in rollopia, i.e., $X_0 = r$.

$$\mathbf{P}(X_1 = f | X_0 = r) = p, \quad \mathbf{P}(X_1 = r | X_0 = r) = 1 - p. \quad (6.4)$$

Probability $1 - p$, i.e., we know from the first row of P that he will leave to Hippopota with probability p and stay with Hippopota in rollopia, i.e., $X_0 = r$. When he gets restless for the first time suppose we first see Freddy in rollopia, i.e., $X_0 = r$.

$$\begin{aligned} & \cdot \begin{pmatrix} f & b \\ r & f \end{pmatrix} = \begin{pmatrix} (f, f, r) & P(f, f, r) \\ (r, f, f) & P(r, f, f) \end{pmatrix} \begin{pmatrix} f \\ r \end{pmatrix} = P \\ & \cdot \begin{pmatrix} f & b \\ r & f \end{pmatrix} = \begin{pmatrix} f & b \\ r & f \end{pmatrix} \begin{pmatrix} f \\ r \end{pmatrix} = P \end{aligned} \quad (6.3)$$

Then Freddy's sequence of jumps (X_0, X_1, \dots) is a Markov chain on \mathbb{X} with transition matrix:

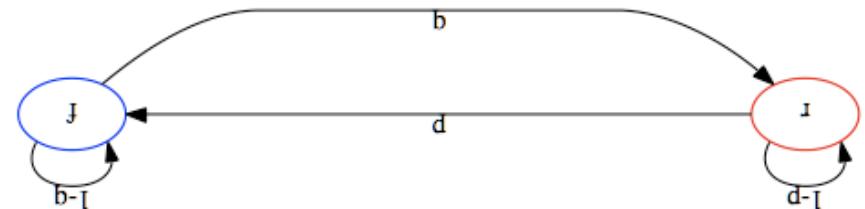


Figure 6.1: Transition Diagram of Hippoat Freddy's jumps.

Let the state space $\mathbb{X} = \{r, f\}$, and let (X_0, X_1, \dots) be the sequence ofify pads occupied by Freddy jumps by the following **transition diagram**:

Let the die on Hippopota f has probability q of turning up heads. We can visualise the rules of Freddy's restless moments. Say the die on rollopia r has probability p of turning up odd and the die on Hippopota f has probability q of turning up even. We can visualise the rules of Freddy's restless moments by the following transition diagram:

After his restless moments, Freddy left the die and if the die landed odd he would leave the die behind and jump to Hippopota. When Freddy got restless in Hippopota he would roll the die and if the die landed odd he would leave the die behind and jump to Hippopota, otherwise flip the coin and if it landed Heads he would leave the coin behind and jump to rollopia, otherwise flip the coin and if it landed Tails he would stay put. When Freddy got restless in Hippopota he would leave the coin behind and jump to rollopia, otherwise flip the coin to Hippopota, otherwise he would stay put. We can visualise the rules of Freddy's restless moments by the following transition diagram:

Two silly pads, rollopia and Hippopota. A wizard gave a die and a silver coin to help Hippoat Freddy decide where to jump next. Freddy left the die on rollopia and the coin on Hippopota. When Freddy got restless in Hippopota he would roll the die and if the die landed odd he would leave the die behind and jump to Hippopota, otherwise he would stay put. When Freddy got restless in Hippopota he would leave the coin behind and jump to rollopia, otherwise flip the coin and if it landed Heads he would leave the coin behind and jump to rollopia, otherwise flip the coin and if it landed Tails he would stay put. We can visualise the rules of Freddy's restless moments by the following transition diagram:

Thus, for a Markov chain $(X_n)_{n \in \mathbb{Z}_+}$, the distribution of X_{t+1} given X_0, \dots, X_t depends on X_t alone. Because of this dependence on the previous state, the stochastic sequence, (X_0, X_1, \dots) , are not independent. We introduce the most important concepts using a simple example.

Example 169 (Freddy the Hippoat frog lives in an enchanted pond with only two silly pads, rollopia and Hippopota. A wizard gave a die and a silver coin to help Hippoat Freddy decide where to jump next. Freddy left the die on rollopia and the coin on Hippopota. When Freddy got restless in Hippopota he would roll the die and if the die landed odd he would leave the die behind and jump to Hippopota, otherwise he would stay put. When Freddy got restless in Hippopota he would leave the coin behind and jump to rollopia, otherwise flip the coin and if it landed Heads he would leave the coin behind and jump to rollopia, otherwise flip the coin and if it landed Tails he would stay put. We can visualise the rules of Freddy's restless moments by the following transition diagram:

disttribution $P(x, \cdot) := (P(x, y))_{y \in \mathbb{X}}$. For this reason P is called a **stochastic matrix**, i.e.,

$$P(x, y) \geq 0 \quad \text{for all } (x, y) \in \mathbb{X}^2 \quad \text{and} \quad \sum_{y \in \mathbb{X}} P(x, y) = 1 \quad \text{for all } x \in \mathbb{X}. \quad (6.2)$$

$|\mathbb{X}| \times |\mathbb{X}|$ matrix P is enough to obtain the state transitions since the x -th row of P is the probability distribution $P(x, \cdot) = (P(x, y))_{y \in \mathbb{X}}$.

Before we see the more general form of Bayes' Rule, let us make a simple observation called the *total probability theorem*.

Proposition 14 (Total probability theorem) Suppose $A_1 \cup A_2 \dots \cup A_k$ is a sequence of events with positive probability that partition the sample space, that is, $A_1 \cup A_2 \dots \cup A_k = \Omega$ and $A_i \cap A_j = \emptyset$ for any $i \neq j$, then for some arbitrary event B .

$$P(B) = \sum_{h=1}^k P(B \cap A_h) = \sum_{h=1}^k P(B|A_h)P(A_h) \quad (2.4)$$

Proof: The first equality is due to the addition rule for mutually exclusive events,

$$B \cap A_1, B \cap A_2, \dots, B \cap A_k$$

and the second equality is due to the multiplication rule for two likely events.

Reference to the Venn diagram (you should draw below as done in lectures) will help you understand this idea for the four event case.

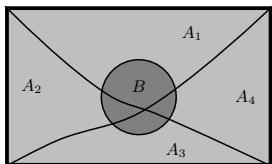


Figure 2.3: Reference to the Venn diagram will help you understand this idea behind the proof of the total probability theorem in Proposition 14 for the four event case.

Example 33 (Urn with red and black balls) A well-mixed urn contains five red and ten black balls. We draw two balls from the urn without replacement. What is the probability that the second ball drawn is red?

This is easy to see if we draw a probability tree diagram. The first split in the tree is based on the outcome of the first draw and the second on the outcome of the last draw. The outcome of the first draw dictates the probabilities for the second one since we are sampling without replacement. We multiply the probabilities on the edges to get probabilities of the four endpoints, and then sum the ones that correspond to red in the second draw, that is

$$P(\text{second ball is red}) = 4/42 + 10/42 = 1/3 .$$

Definition 75 (Independent and Identically Distributed (IID) Process) The finite or infinite sequence of RVs or the stochastic process X_1, X_2, \dots is said to be independent and identically distributed or IID if :

- they are independently distributed according to Definition 43, and
- $F(X_1) = F(X_2) = \dots$, ie. all the X_i 's have the same DF $F(X_1)$.

This is perhaps the most elementary class of stochastic processes and we succinctly denote it by

$$(X_i)_{i \in [n]} := X_1, X_2, \dots, X_n \stackrel{\text{IID}}{\sim} F, \quad \text{or} \quad (X_i)_{i \in \mathbb{N}} := X_1, X_2, \dots \stackrel{\text{IID}}{\sim} F .$$

We sometimes replace the DF F above by the name of the RV.

Definition 76 (Independently Distributed) The sequence of RVs or the stochastic process $(X_i)_{i \in \mathbb{N}} := X_1, X_2, \dots$ is said to be independently distributed if :

- X_1, X_2, \dots is independently distributed according to Definition 43.

This is a class of stochastic processes that is more general than the IID class.

As an example of such a class consider the sequence of RVs that are independent but non-identically distributed with each $X_i \sim \text{Bernoulli}(\theta_i)$.

When a stochastic process $(X_\alpha)_{\alpha \in \mathbb{A}}$ is not independent it is said to be dependent. So far we have mostly concerned ourselves with independent processes. In this chapter we introduce finite Markov chains and their simulation methods. Finite Markov chains are among the simplest stochastic processes with a 'first-order' dependence called Markov dependence.

6.2 Introduction

A finite Markov chain is a stochastic process that moves among elements in a finite set \mathbb{X} as follows: when at $x \in \mathbb{X}$ the next position is chosen at random according to a fixed probability distribution $P(\cdot|x)$. We define such a process more formally below.

Definition 77 (Finite Markov Chain) A stochastic sequence,

$$(X_n)_{n \in \mathbb{Z}_+} := (X_0, X_1, \dots),$$

is a homogeneous **Markov chain** with **state space** \mathbb{X} and **transition matrix** $P := (P(x,y))_{(x,y) \in \mathbb{X}^2}$ if for all pair of **states** $(x,y) \in \mathbb{X}^2 := \mathbb{X} \times \mathbb{X}$, all integers $t \geq 1$, and all probable historical events $H_{t-1} := \bigcap_{n=0}^{t-1} \{X_n = x_n\}$ with $\mathbf{P}(H_{t-1} \cap \{X_t = x\}) > 0$, the following **Markov property** is satisfied:

$$\mathbf{P}(X_{t+1} = y | H_{t-1} \cap \{X_t = x\}) = \mathbf{P}(X_{t+1} = y | X_t = x) =: P(x,y) . \quad (6.1)$$

The Markov property means that the conditional probability of going to state y at time $t+1$ from state x at current time t is always given by the (x,y) -th entry $P(x,y)$ of the transition matrix P , no matter what sequence of states $(x_0, x_1, \dots, x_{t-1})$ preceded the current state x . Thus, the

This is exactly the sequence of RVs associated with our product experiment $\mathcal{G}_{\otimes \infty} := (\Omega, \mathcal{F}_{\infty}, P_{\infty})_{\otimes \infty}$.

Distributed or IID Process or merely **IID Sequence of RVs** when the index set is a subset of \mathbb{N} .
absolutely simplest but extremely useful assumption is that of the **Independent and Identically Distributed** RVs. Generally, we cannot produce useful models without making simple assumptions. The Of course the above process is quite general and can allow for arbitrary dependence among the RVs. Of course the above process is quite general and can allow for arbitrary dependence among the RVs. Generally, we cannot produce useful models without making simple assumptions. The etc.

If $A \subset \mathbb{R}$ then we have a **continuous time stochastic process**, typically denoted by $\{X_t\}_{t \in \mathbb{R}}$,

$$(X_i)_{i \in [n]} := X_1, X_2, \dots, X_n, \text{ where, } [n] := \{1, 2, \dots, n\}.$$

$$(X_i)_{i \in \mathbb{Z}} := X_1, X_2, \dots, \text{ or}$$

$$(X_i)_{i \in \mathbb{Z}} := \dots, X_{-2}, X_{-1}, X_0, X_1, X_2, \dots, \text{ or}$$

denoted by

X_a is a RV. If the index set $A \subset \mathbb{A}$, the index set of the stochastic process, typically is called a **stochastic process**. Thus, for every $a \in A$, the collection of RVs

$$(X_a)_{a \in A} := (X_a : a \in A)$$

Definition 74 (Stochastic Process) A collection of RVs

6.1 Stochastic Processes

as done in lectures.

- Example 170.

- Example 169 and

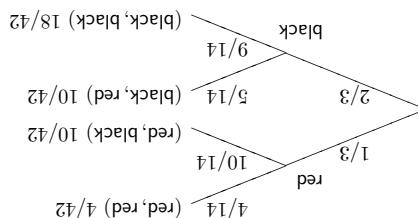
the ideas behind:

This topic is only introduced briefly in Probability Theory I to give concrete instances of dependent sequence of random variables. We will revisit these ideas in the sequel. You only need to understand

NOTE: No materials from this Chapter will be in Probability Theory I exam!

Finite Markov Chains

Chapter 6



We call $\mathbf{P}(A_h)$ the **prior probability** of A_h , i.e., before observing B or *a priori*, and $\mathbf{P}(A_h|B)$ the **posterior probability** of A_h , i.e., after observing B or *a posteriori*.

This theorem is at the heart of solving Bayesian *Decision Problems* which fall into several sub-problems called *inference*, *learning* and *control* problems. Let's see one of the simplest such *learning problems* called *prediction*, more specifically *classification*, where we need to choose between finitely many possible choices based on past information next.

Example 34 (Wasserman2003 p.12) Suppose Larry divides his email into three categories: $A_1 = \text{"spam"}$, $A_2 = \text{"low priority"}$, and $A_3 = \text{"high priority"}$. From previous experience, he finds that $\mathbf{P}(A_1) = 0.7$, $\mathbf{P}(A_2) = 0.2$ and $\mathbf{P}(A_3) = 0.1$. Note that $\mathbf{P}(A_1 \cup A_2 \cup A_3) = \mathbf{P}(\Omega) = 0.7 + 0.2 + 0.1 = 1$. Let B be the event that the email contains the word "free." From previous experience, $\mathbf{P}(B|A_1) = 0.9$, $\mathbf{P}(B|A_2) = 0.01$ and $\mathbf{P}(B|A_3) = 0.01$. Note that $\mathbf{P}(B|A_1) + \mathbf{P}(B|A_2) + \mathbf{P}(B|A_3) = 0.9 + 0.01 + 0.01 \neq 1$. Now, suppose Larry receives an email with the word "free." What is the probability that it is "spam," "low priority," and "high priority"?

Solution:

$$\begin{aligned}\mathbf{P}(A_1|B) &= \frac{\mathbf{P}(B|A_1)\mathbf{P}(A_1)}{\mathbf{P}(B|A_1)\mathbf{P}(A_1) + \mathbf{P}(B|A_2)\mathbf{P}(A_2) + \mathbf{P}(B|A_3)\mathbf{P}(A_3)} = \frac{0.9 \times 0.7}{(0.9 \times 0.7) + (0.01 \times 0.2) + (0.01 \times 0.1)} = \frac{0.63}{0.633} \approx 0.995 \\ \mathbf{P}(A_2|B) &= \frac{\mathbf{P}(B|A_2)\mathbf{P}(A_2)}{\mathbf{P}(B|A_1)\mathbf{P}(A_1) + \mathbf{P}(B|A_2)\mathbf{P}(A_2) + \mathbf{P}(B|A_3)\mathbf{P}(A_3)} = \frac{0.01 \times 0.2}{(0.9 \times 0.7) + (0.01 \times 0.2) + (0.01 \times 0.1)} = \frac{0.002}{0.633} \approx 0.003 \\ \mathbf{P}(A_3|B) &= \frac{\mathbf{P}(B|A_3)\mathbf{P}(A_3)}{\mathbf{P}(B|A_1)\mathbf{P}(A_1) + \mathbf{P}(B|A_2)\mathbf{P}(A_2) + \mathbf{P}(B|A_3)\mathbf{P}(A_3)} = \frac{0.01 \times 0.1}{(0.9 \times 0.7) + (0.01 \times 0.2) + (0.01 \times 0.1)} = \frac{0.001}{0.633} \approx 0.002\end{aligned}$$

Note that $\mathbf{P}(A_1|B) + \mathbf{P}(A_2|B) + \mathbf{P}(A_3|B) = 0.995 + 0.003 + 0.002 = 1$.

This is essentially the idea behind *Bayes classifiers*, that are used to solve such *prediction* problems across different problem domains in *statistical machine learning*, where solutions are given from computer programs.

2.4.2 Independence and Dependence

In general, $P(A|B)$ and $P(A)$ are different, but sometimes the occurrence of B makes no difference, and gives no new information about the chances of A occurring. This is the idea behind independence. Events like "having blue eyes" and "having blond hair" are associated due to common genetic ancestry, but events like "my neighbour wins Lotto" and "I win Lotto" are not due to the Lotto machine being chaotically whirled around before ejection (as modelled by a well-stirred urn).

Definition 16 (Independence of two events) Any two events A and B are said to be **independent** if and only if

$$\mathbf{P}(A \cap B) = \mathbf{P}(A)\mathbf{P}(B). \quad (2.6)$$

Let us make sense of this definition in terms of our previous definitions. When $\mathbf{P}(A) = 0$ or $\mathbf{P}(B) = 0$, both sides of the above equality are 0. If $\mathbf{P}(A) \neq 0$, then rearranging the above equation we get:

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

But, the LHS is $\mathbf{P}(B|A)$ by definition 2.2, and thus for independent events A and B , we get:

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

This says that information about the occurrence of A does not affect the occurrence of B . If $\mathbf{P}(B) \neq 0$, then an analogous argument:

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

Example 168 We model the tosses of a coin with unknown $E(X_1) = \theta^*$ as

$$X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \text{Bernoulli}(\theta^*)$$

and observed the following data, sample mean, sample variance and sample standard deviation:

$$(x_1, x_2, \dots, x_7) = (0, 1, 1, 0, 0, 1, 0), \bar{x}_7 = 0.4286, s_7^2 = 0.2857, s_7 = 0.5345,$$

respectively. Our point estimate and $1 - \alpha = 95\%$ confidence interval for $E(X_1)$ are:

$$\bar{x}_7 = 0.4286 \quad \text{and} \quad (\bar{x}_7 \pm z_{\alpha/2}s_7/\sqrt{7}) = (0.4286 \pm 1.96 \times 0.5345/\sqrt{7}) = (0.0326, 0.8246),$$

respectively. So with 95% probability the true population mean $E(X_1) = \theta^*$ is contained in $(0.0326, 0.8246)$ and since $1/2$ is contained in this interval of width 0.792 we cannot rule out that the flipped coin is not fair with $\theta^* = 1/2$.

Remark 73 The normal-based confidence interval for θ^* (as well as λ^* in the previous example) may not be a valid approximation here with just $n = 7$ samples. After all, the CLT only tells us that the point estimator $\hat{\Theta}_n$ can be approximated by a normal distribution for large sample sizes. When the sample size n was increased from 7 to 100 by tossing the same coin another 93 times, a total of 57 trials landed as Heads. Thus the point estimate and confidence interval for $E(X_1) = \theta^*$ based on the sample mean and sample standard deviations are:

$$\hat{\theta}_{100} = \frac{57}{100} = 0.57 \quad \text{and} \quad (0.57 \pm 1.96 \times 0.4975/\sqrt{100}) = (0.4725, 0.6675).$$

Thus our confidence interval shrank considerably from a width of 0.792 to 0.195 after an additional 93 Bernoulli trials. Thus, we can make the width of the confidence interval as small as we want by making the number of observations or sample size n as large as we can.

5.4 Exercises in Limit Laws of Statistics

Ex. 5.2 — Suppose you plan to obtain a simple random sequence (SRS) — also known as independent and identically distributed (IID) sequence — of n measurements from an instrument. This instrument has been calibrated so that the distribution of measurements made with it have population variance of $1/4$. Your boss wants you to make a point estimate of the unknown population mean from a SRS of sample size n . He also insists that the tolerance for error has to be $1/10$ and the probability of meeting this tolerance should be just above 95%. Use CLT to find how large should n be to meet the specifications of your boss.

Ex. 5.3 — Suppose the collection of RVs X_1, X_2, \dots, X_n model the number of errors in n computer programs named $1, 2, \dots, n$, respectively. Suppose that the RV X_i modeling the number of errors in the i -th program is the Poisson($\lambda = 5$) for any $i = 1, 2, \dots, n$. Further suppose that they are independently distributed. Succinctly, we suppose that

$$X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \text{Poisson}(\lambda = 5).$$

Suppose we have $n = 125$ programs and want to make a probability statement about \bar{X}_{125} which is the average error per program out of these 125 programs. Since $E(X_i) = \lambda = 5$ and $V(X_i) = \lambda = 5$, we want to know how often our sample mean \bar{X}_{125} differs from the expectation of 5 errors per program. Using the CLT find the $P(\bar{X}_{125} < 5.5)$.

Ex. 5.4 — What is the distribution of $\sum_{i=1}^n X_i/n$ as $n \rightarrow \infty$ when $X_i \stackrel{iid}{\sim} \text{Uniform}(-10, 10)$?

Ex. 5.5 — What is the distribution of $\sum_{i=1}^n X_i / \sqrt{\mathbf{V}(X_i)}$ as $n \rightarrow \infty$ when $X_i \stackrel{iid}{\sim} \text{Uniform}(-10, 10)$?

Now, let C be the event that the sum of the two dice equals seven. Then

$$\mathbf{P}(C \cap B) = \mathbf{P}(\{(4, 3)\}) = \frac{1}{36},$$

while

$$\begin{aligned}\mathbf{P}(C \cap B) &= \mathbf{P}(\{(1, 6), (2, 5), (3, 4), (4, 3), (5, 2), (6, 1)\})\mathbf{P}(\{(4, 1), (4, 2), (4, 3), (4, 4), (4, 5), (4, 6)\}) \\ &= \frac{6}{36} \times \frac{6}{36} = \frac{1}{36},\end{aligned}$$

and therefore C and B are independent events. Once again this is clear because the chance of getting a total of seven does not depend any more on the outcome of the first die (it is allowed to be any one of the six possible outcomes).

Example 37 (Pairwise independent events that are not jointly independent) Let a ball be drawn from an well-stirred urn containing four balls labelled 1,2,3,4. Consider the events $A = \{1, 2\}$, $B = \{1, 3\}$ and $C = \{1, 4\}$. Then,

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

but,

$$\frac{1}{4} = \mathbf{P}(\{1\}) = \mathbf{P}(A \cap B \cap C) \neq \mathbf{P}(A)\mathbf{P}(B)\mathbf{P}(C) = \frac{2}{4} \times \frac{2}{4} \times \frac{2}{4} = \frac{1}{8}.$$

Therefore, inspite of being pairwise independent, the events A , B and C are not jointly independent.

CONDITIONAL PROBABILITY SUMMARY

$\mathbf{P}(A|B)$ means the probability that A occurs given that B has occurred.

$$\mathbf{P}(A|B) = \frac{\mathbf{P}(A \cap B)}{\mathbf{P}(B)} = \frac{\mathbf{P}(A)\mathbf{P}(B|A)}{\mathbf{P}(B)} \quad \text{if } \mathbf{P}(B) \neq 0$$

$$\mathbf{P}(B|A) = \frac{\mathbf{P}(A \cap B)}{\mathbf{P}(A)} = \frac{\mathbf{P}(B)\mathbf{P}(A|B)}{\mathbf{P}(A)} \quad \text{if } \mathbf{P}(A) \neq 0$$

Conditional probabilities obey the axioms and rules of probability.

2.5 Exercises in Conditional Probability

Ex. 2.6 — What gives the greater probability of hitting some target at least once:

- 1.hitting in a shot with probability $\frac{1}{2}$ and firing 1 shot, or
- 2.hitting in a shot with probability $\frac{1}{3}$ and firing 2 shots?

First guess. Then calculate.

Example 165 Suppose an IID sequence of observations $(x_1, x_2, \dots, x_{80})$ was drawn from a distribution with variance $V(X_1) = 4$. What is the probability that the error in \bar{x}_n used to estimate $E(X_1)$ is less than 0.1?

By CLT,

$$P(\text{error} < 0.1) \cong P\left(-\frac{0.1}{\sqrt{4/80}} < Z < \frac{0.1}{\sqrt{4/80}}\right) = P(-0.447 < Z < 0.447) = 0.345.$$

Suppose you want the error to be less than tolerance = ϵ with a certain probability $1 - \alpha$. Then we can use CLT to do such **sample size calculations**. Recall the DF $\Phi(z) = P(Z < z)$ is tabulated in the standard normal table and now we want

$$P\left(-\frac{\epsilon}{\sqrt{V(X_1)/n}} < Z < \frac{\epsilon}{\sqrt{V(X_1)/n}}\right) = 1 - \alpha.$$

We know,

$$P(-z_{\alpha/2} < Z < z_{\alpha/2}) = 1 - \alpha,$$

make the picture here of $f_Z(z) = \Phi'(z)$ to recall what $z_{\alpha/2}$, $z_{-\alpha/2}$, and the various areas below $f_Z(\cdot)$ in terms of $\Phi(\cdot)$ from the table really mean... (See Example 59).

where, $\Phi(z_{\alpha/2}) = 1 - \alpha/2$ and $\Phi(z_{-\alpha/2}) = 1 - \Phi(z_{\alpha/2}) = \alpha/2$. So, we set

$$\frac{\epsilon}{\sqrt{V(X_1)/n}} = z_{\alpha/2}$$

and rearrange to get

$$n = \left(\frac{\sqrt{V(X_1)} z_{\alpha/2}}{\epsilon} \right)^2 \quad (5.9)$$

for the needed sample size that will ensure that our error is less than our tolerance = ϵ with probability $1 - \alpha$. Of course, if n given by Equation (5.9) is not a natural number then we naturally round up to make it one!

A useful $z_{\alpha/2}$ value to remember: If $\alpha = 0.05$ when the probability of interest $1 - \alpha = 0.95$ then $z_{\alpha/2} = z_{0.025} = 1.96$.

Example 166 How large a sample size is needed to make the error in our estimate of the population mean $E(X_1)$ to be less than 0.1 with probability $1 - \alpha = 0.95$ if we are observing IID samples from a distribution with a population variance $V(X_1)$ of 4?

Using Equation (5.9) we see that the needed sample size is

$$n = \left(\frac{\sqrt{4} \times 1.96}{0.1} \right)^2 \cong 1537$$

Thus, it pays to check the sample size needed in advance of experimentation, provided you already know the population variance of the distribution whose population mean you are interested in estimating within a given tolerance and with a high probability.

where $Z \sim \text{Normal}(0, 1)$.

$$P(-e < \underline{X}^n - E(X^1) < e) \approx P\left(\frac{\underline{X}^n - E(X^1)}{\sqrt{n}} < \frac{Z}{\sqrt{n}} < \frac{e}{\sqrt{n}}\right) = P\left(\frac{Z}{\sqrt{n}} > \frac{e}{\sqrt{n}}\right) =$$

Due to the Central Limit Theorem (CLT) we now know that (assuming n is large)

$$P(\text{error} < \text{tolerance}) = P(|\underline{X}^n - E(X^1)| < e) = P(-e < \underline{X}^n - E(X^1) < e) = 1 - \alpha.$$

Required tolerance = e and make the following probability statement:

Recall that we wanted to ensure the error = $|\underline{X}^n - E(X^1)|$ in our estimate of $E(X^1)$ is within a

5.3.1 Application: Tolerating Errors in our estimate of $E(X^1)$

then $Z = \frac{\underline{X}^n - E(X^1)}{\sigma_{\underline{X}^n}} \sim \text{Normal}(0, 1)$ through the linear transformation $W = \sigma_Z Z + \mu$ of Example 67.

which is equivalent to Equation (5.7) by a standardization argument that if $W \sim \text{Normal}(\mu, \sigma^2)$

For the last limit we have used $(1 + \frac{y}{x})^x \rightarrow e^y$ as $x \rightarrow \infty$. Thus, we have proved Equation (5.8).

$$\phi_{U^n}(t) = \left(\frac{\sqrt{n}}{t} \right)^n = \left(1 + \frac{\sqrt{n}}{t} \times 0 + \frac{2n}{t^2} + o\left(\frac{n}{t^2}\right) \right)^n = \left(1 - \frac{2n}{t^2} + o\left(\frac{n}{t^2}\right) \right)^n \rightarrow e^{-2t^2/n} = \phi_Z(t).$$

we finally get

$$\phi_Y(t) = \left(\frac{\sqrt{n}}{t} \right)^n = 1 + \frac{\sqrt{n}}{t} E(Y) + \frac{2n}{t^2} E(Y^2) + o\left(\frac{n}{t^2}\right),$$

which implies

$$\phi_Y(t) = 1 + tE(Y) + \frac{t^2}{2} E(Y^2) + o(t^2),$$

and since we can Taylor expand $\phi_Y(t)$ as follows:

$$\phi_{U^n}(t) = \left(\frac{\sqrt{n}}{t} \right)^n,$$

So, the CF of U_n is

$$Y = \frac{\sqrt{A(X^1)}}{\underline{X}^n - E(X^1)}$$

Now, if we let

$$\begin{aligned} \phi_{U^n}(t) &= \left(\frac{\sqrt{n}}{t} \right)^n \frac{\sqrt{A(X^1)}}{\underline{X}^n - E(X^1)} \\ &= \left(\frac{\exp\left(\frac{t}{\sqrt{n}}\sqrt{A(X^1)}\right)}{\underline{X}^n - E(X^1)} \right)^n = \left(\prod_{k=1}^n \exp\left(\frac{t}{\sqrt{k}}\sqrt{\frac{A(X^1)}{A(X^k)}}\right) \right)^n = \left(\prod_{k=1}^n \left(\frac{\sqrt{A(X^1)}}{\underline{X}^k - E(X^1)} \right)^{\sqrt{\frac{A(X^1)}{A(X^k)}}} \right)^n = \left(\prod_{k=1}^n \left(\frac{\sqrt{A(X^1)}}{\underline{X}^k - E(X^1)} \right)^{\sqrt{\frac{A(X^1)}{A(X^k)}}} \right)^n = \left(\frac{\sqrt{A(X^1)}}{\underline{X}^n - E(X^1)} \right)^n = \end{aligned}$$

Therefore, the CF of U_n is

$$U_n := \underline{X}^n - E(X^1) = \frac{\sqrt{n} A(X^1)/n}{\sum_{k=1}^n X^k - nE(X^1)} = \frac{\sqrt{n} A(X^1)/n}{1 - \frac{1}{n} \sum_{k=1}^{n-1} X^k} =$$

Second,

If the detection rate is 0.99 and the false alarm rate is 0.001, and the probability of an intrusion occurring is 0.01, find

$$\text{false alarm rate} = P(\text{detection declared} | \text{no intrusion}),$$

,

$$\text{detection rate} = P(\text{detection declared} | \text{intrusion}),$$

,

Ex. 2.13 — **The detection rate and false alarm rate of an intrusion sensor are defined as

(b) the probability that a patient who tests negative is free from the disease.

(a) the probability that a patient who tests positive actually has the disease,

Suppose that a medical test has a sensitivity of 0.7 and a specificity of 0.95. If the prevalence of the disease in the general population is 1%, and

$$\text{specificity} = P(\text{test is negative} | \text{patient does not have the disease}),$$

$$\text{sensitivity} = P(\text{test is positive} | \text{patient has the disease}),$$

as follows:

Ex. 2.12 — **The sensitivity and specificity of a medical diagnostic test for a disease are defined

3.

(c) If the gale did NOT cause damage, find the probabilities that it was of: force 1; force 2; force 3.

,

(b) If the gale caused damage, find the probabilities that it caused damage?

,

(a) If a gale is reported, what is the probability that force 3 gales cause damage?

,

is $\frac{3}{5}$ and the probability that force 1 gales cause damage $\frac{1}{5}$, the probability that force 2 gales cause damage

the probability that force 2 and force 3 gales cause damage $\frac{1}{5}$ are force 3. Furthermore,

Ex. 2.11 — Suppose that $\frac{3}{5}$ of all galaxies are force 1, $\frac{1}{5}$ are force 2 and $\frac{1}{5}$ are force 3. Furthermore,

one is in fact defective?

,

(c) If 2 micro-chips are tested and determined to be good, what is the probability that at Least

it was good anyway?

(b) If a micro-chip is chosen at random, and tested to be defective, what was the probability that it was defective anyway?

,

(a) If a micro-chip is chosen at random, and tested to be good, what was the probability that it

was defective?

micro-chip is tested, and the test will correctly detect a defective one $\frac{4}{5}$ of the time, and if a good

chip is tested, and the test will correctly detect a defective with probability $\frac{1}{10}$.

Ex. 2.10 — A process producing micro-chips, produces 5% defective, at random. Each micro-

chip is randomly selected bottle comes from machine 1 given that it is accepted?

,

reason, while one out of every 30 bottles filled by machine 2 is rejected. What is the probability

and machine 2 produces 25%. One out of every 20 bottles filled by machine 1 is rejected for some

and machine 2 produces 25%. One out of every 30 bottles filled by machine 2 is rejected. What is the probability

that a random sample of 30 bottles filled by machine 1 is accepted?

,

The final exam is passed by 80% of those who passed the mid-term test, but only by 40% of those

who fail the mid-term test. What fraction of students pass the final exam?

,

Ex. 2.8 — Based on past experience, 70% of students in a certain course pass the mid-term test.

,

The final exam is passed by 80% of those who passed the mid-term test, but only by 40% of those

who fail the mid-term test. What is the probability of obtaining a sum greater than 7?

,

Ex. 2.7 — Suppose we independently roll two fair dice each of whose faces are marked by numbers

1, 2, 3, 4, 5 and 6.

CHAPTE R 2. PROBABILITY MODEL

57

- (a) the probability that there is an intrusion when a detection is declared,
- (b) the probability that there is no intrusion when no detection is declared.

Ex. 2.14 — **Let A and B be events such that $\mathbf{P}(A) \neq 0$ and $\mathbf{P}(B) \neq 0$. When A and B are disjoint, are they also independent? Explain clearly why or why not.

trajectories for simulated tosses of a fair coin from IID Bernoulli($\theta^* = 1/2$) RVs and the twenty red sample mean trajectories for simulated waiting times from IID Exponential($\lambda^* = 1/10$) RVs in Figure 5.4 with $n = 7$. Clearly, the point estimates for such a small sample size are fluctuating wildly! However, the fluctuations in the point estimates settles down for larger sample sizes.

The next natural question is how large should the sample size be in order to have a small interval of width, say 2ϵ , “contain” $E(X_1)$, the quantity of interest, with a high probability, say $1 - \alpha$? If we can answer this then we can make probability statements like the following:

$$P(\text{error} < \text{tolerance}) = P(|\bar{X}_n - E(X_1)| < \epsilon) = P(-\epsilon < \bar{X}_n - E(X_1) < \epsilon) = 1 - \alpha .$$

In order to ensure the error = $|\bar{X}_n - E(X_1)|$ in our estimate of $E(X_1)$ is within a required tolerance = ϵ we need to know the full distribution of $\bar{X}_n - E(X_1)$ itself. The Central Limit Theorem (CLT) helps us here.

5.3 Central Limit Theorem

What if we scale the sum of X_i 's by \sqrt{n} instead of n ?

Exercise 5.1 (What if we scale by \sqrt{n}) After reading Sec. 5.1 up to now, think carefully about what you need to be able to show that $Z_n := 1/\sqrt{n} \sum_{i=1}^n X_i$ converges in distribution to the Normal(0, 1/3) RV, where $X_i \stackrel{\text{IID}}{\sim} \text{Uniform}(-1, 1)$. Hint: Characteristic functions

Proposition 71 (Central Limit Theorem (CLT)) If we are given a sequence of independently and identically distributed (IID) RVs, $X_1, X_2, \dots \stackrel{\text{IID}}{\sim} X_1$ and if $E(X) < \infty$ and $V(X_1) < \infty$, then the sample mean \bar{X}_n converges in distribution to the Normal RV with mean given by any one of the IID RVs, say $\mathbf{E}(X_1)$ by convention, and variance given by $\frac{1}{n}$ times the variance of any one of the IID RVs, say $V(X_1)$ by convention. More formally, we write:

$$\begin{aligned} \text{If } X_1, X_2, \dots \stackrel{\text{IID}}{\sim} X_1 \text{ and if } \mathbf{E}(X_1) < \infty, V(X_1) < \infty \\ \text{then } \bar{X}_n \rightsquigarrow \text{Normal}\left(E(X_1), \frac{V(X_1)}{n}\right) \text{ as } n \rightarrow \infty , \end{aligned} \quad (5.7)$$

or equivalently after standardization:

$$\begin{aligned} \text{If } X_1, X_2, \dots \stackrel{\text{IID}}{\sim} X_1 \text{ and if } \mathbf{E}(X_1) < \infty, V(X_1) < \infty \\ \text{then } \frac{\bar{X}_n - E(X_1)}{\sqrt{V(X_1)/n}} \rightsquigarrow Z \sim \text{Normal}(0, 1) \text{ as } n \rightarrow \infty . \end{aligned} \quad (5.8)$$

Proof: Our proof is based on the convergence of characteristic functions (CFs). We will prove the standardized form of the CLT in Equation (5.8) by showing that the CF of

$$U_n := \frac{\bar{X}_n - E(X_1)}{\sqrt{V(X_1)/n}}$$

converges to the CF of Z , the Normal(0, 1) RV. First, note from Equation (3.72) that the CF of $Z \sim \text{Normal}(0, 1)$ is:

$$\varphi_Z(t) = E(e^{itZ}) = e^{-t^2/2} .$$

$$X(\omega) = \begin{cases} 0, & \text{if } \omega = \text{not rain} \\ 1, & \text{if } \omega = \text{rain} \end{cases}$$

create a random variable X with this experiment as follows:

Example 38 (Rain or Shine) Suppose our experiment is to observe whether it will rain or not rain tomorrow. The sample space of this experiment is $\Omega = \{\text{rain}, \text{not rain}\}$. We can associate a random variable X with this experiment as follows:

Thus, we want a **random variable** to be a function from the sample space Ω to the set of real numbers \mathbb{R} , that is, $X : \Omega \rightarrow \mathbb{R}$. Let us go through some examples before giving the formal definition of such a real-valued random variable.

Experiment	Possible measured outcomes
Counting the number of typos up to now	$\mathbb{Z}_+ = \{0, 1, 2, \dots\} \subset \mathbb{R}$
Length in centimetres of some shells on New Brighton beach	$(0, +\infty) \subset \mathbb{R}$
Waiting time in minutes for the next Orbiter bus to arrive	$\mathbb{R}_+ = [0, \infty) \subset \mathbb{R}$
Vertical displacement from current position of a pollen on water	$\mathbb{R} \subset \mathbb{R}$

Crucially, it can become inconvenient to work with a set of outcomes Ω upon which arithmetic is not possible. We are often measuring our outcomes with subsets of real numbers. Some examples include:

Random variables, unlike classical deterministic variables, can take a bunch of different values. We may take different values in a non-deterministic manner. **Random variables** do this job for us. We need a different kind of variable to deal with real-world situations where the same variable

What these *classical variables* have in common is that they take a fixed or deterministic value when we can solve for them.

$$\{a_n\}_{n=1}^{\infty} = \{1, 2, 3, \dots\}$$

Yet another example is the use of variables to represent sequences such as:

$$\text{over the real line } \mathbb{R} = (-\infty, \infty).$$

where the variable y for the y -axis is determined by the value taken by the variable x , as x varies over the real line.

$$y = 3x - 2,$$

We also use classical variables to represent geometric objects such as a line:

$$\text{We are used to classical variables such as } x \text{ as an "unknown" in the equation: } x + 3 = 7.$$

Random Variables

Chapter 3

Of course, if we tossed the same coin in the same IID manner another seven times or if we observed another seven waiting times of orbits on a different bus-stop or on a different day we may get a different point estimate for $E(X_1)$. See the interpretation of the twenty magenta sample mean which is the same as the probability of Heads is $\frac{x_7}{7} = 3/7$.

Then you can use the observed sample mean $\bar{x}_7 = (0 + 1 + 0 + 0 + 1 + 0)/7 = 3/7 \approx 0.4286$ as a **point estimate** of the population mean $E(X_1) = \theta$. Thus, our "single best guess" for $E(X_1)$ and therefore, we can use the same coin in the same IID manner another seven times or if we observed

$$(x_1, x_2, \dots, x_7) = (0, 1, 0, 1, 0).$$

and have the following realization as your observed data:

$$X_1, X_2, \dots, X_7 \stackrel{\text{iid}}{\sim} \text{Bernoulli}(\theta^*)$$

Now, suppose you model seven coin tosses (encoding Heads as 1 with probability θ^* and Tails as 0 with probability $1 - \theta^*$) as follows:

and therefore, we can use the sample mean \bar{X}_n as a point estimator of $E(X_1) = \theta^*$.

$$X_n \sim \text{Point Mass}(E(X_1))$$

Typically, we do not know the "true" parameter $\theta \in \Theta = [0, 1]$, which is the same as the population mean $E(X_1) = \theta^*$. But by LNN, we know that

$$X_1, X_2, \dots, X_n \stackrel{\text{iid}}{\sim} \text{Bernoulli}(\theta^*)$$

Example 164 Let $X_1, X_2, \dots, X_n \stackrel{\text{iid}}{\sim} X_1$, where X_1 is an $\text{Bernoulli}(\theta^*)$ RV, i.e., let

typerally distinctly especially for small n as shown in Figure 5.4. The sample means from n replicates of the experiment are $\bar{x}_1, \bar{x}_2, \dots, \bar{x}_n$, but the point estimate \bar{x}_n for 20 replicates of $E(X_1)$ is still X_1 , that is different from our first data vector (x_1, x_2, \dots, x_n) , our point estimator of $E(X_1)$ is still X_1 , but the point estimator from the first random variable X_1 , called the point estimator of $E(X_1)$. Therefore, when we observe a new data vector (x_1, x_2, \dots, x_n) and its corresponding realization of the data RV , i.e., our observed data vector (x_1, x_2, \dots, x_n) and its corresponding realization of the data RV , i.e., our observed sample mean \bar{x}_n is a point estimate of $E(X_1)$. In other words, the point estimate \bar{x}_n is a realization of the random variable X_1 , which is a random variable that depends on the data RV (X_1, X_2, \dots, X_n) , is a point estimator of $E(X_1)$. But once we have a realization of the data RV , i.e., our observed data vector (x_1, x_2, \dots, x_n) and its corresponding realization of the data RV , i.e., our observed sample mean \bar{x}_n is a point estimate of $E(X_1)$. We say the statistics are realized as sample mean \bar{x}_n is a point estimate of the "true" parameter θ^* from $1/\bar{x}_7 = 7/7 \approx 0.986$.

Then you can use the observed sample mean $\bar{x}_7 = (2 + 12 + 8 + 9 + 14 + 15 + 11)/7 = 71/7 \approx 10.14$ as a **point estimate** of the population mean $E(X_1) = 1/\bar{x}_7$. By rearranging $\bar{x}_7 = 1/E(X_1)$, we can also obtain a point estimate of the "true" parameter θ^* from $1/\bar{x}_7 = 7/7 \approx 0.986$.

$$(x_1, x_2, \dots, x_7) = (2, 12, 8, 9, 14, 15, 11)$$

and have the following realization as your observed data:

$$X_1, X_2, \dots, X_7 \stackrel{\text{iid}}{\sim} \text{Exponential}(\lambda^*)$$

Now, suppose you model seven waiting times in nearest minutes between Orbiter buses at Bangalore street as follows:

and therefore, we can use the sample mean \bar{X}_n as a point estimator of $E(X_1) = 1/\lambda^*$.

Thus, X will take the value 1 if it will rain tomorrow and 0 otherwise. Note that another equally valid (though possibly not so useful) random variable, say Y , for this experiment is:

$$Y(\omega) = \begin{cases} \pi, & \text{if } \omega = \text{rain} \\ \sqrt{2}, & \text{if } \omega = \text{not rain} \end{cases}$$

Example 39 (Rain Fall on Angstrom) Suppose our experiment instead is to measure the volume of rain that falls into a large funnel stuck on top of a graduated cylinder that is placed on top of the middle of House 1 of Angstrom Laboratory. Suppose the cylinder is graduated in millimeters then our random variable $X(\omega)$ can report a non-negative real number given by the lower miniscus of the water column, if any, in the cylinder tomorrow. Thus, $X(\omega)$ will measure the volume of rain in millilitres that will fall into our funnel tomorrow.

Example 40 (Counting Seedlings) Suppose ten seeds are planted. Perhaps fewer than ten will actually germinate. The number which do germinate, say X , must be one of the integer numbers in \mathbb{R} given by the set:

$$\mathbb{X} := \{0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10\} .$$

But until the seeds are actually planted and allowed to germinate it is impossible to say which number $X(\omega) : \Omega \rightarrow \mathbb{X}$ will take. The number of seeds which germinate is a variable, but it is not necessarily the same for each group of ten seeds planted, but takes values from the same set \mathbb{X} . As X is not known in advance it is called a **random variable**. Its value cannot be known until we actually perform the experiment, i.e., plant the seeds.

Certain things can be said about the value a random variable might take. In the case of these ten seeds we can be sure the number that germinate is less than eleven, and not less than zero! It may also be known that the probability of seven seeds germinating is greater than the probability of one seed; or perhaps that the number of seeds germinating averages eight. These statements are based on probabilities unlike the sort of statements made about deterministic variables.

Discrete versus continuous random variables.

A **discrete** random variable is one in which the set of possible values of the random variable is finite or at most countably infinite, whereas a **continuous** random variable may take on any value in some range, and its value may be any real value in that range (Think: uncountably infinite). Examples 38 and 40 are about discrete random variables and Example 39 is about a continuous random variable.

Discrete random variables are usually generated from experiments where things are “counted” rather than “measured” such as the seed planting experiment in Example 40. Continuous random variables appear in experiments in which we measure, such as the amount of rain, in millilitres in Example 39.

Random variables as functions.

In fact, random variables are actually functions, more formally measurable maps from \mathcal{F} to certain subsets of \mathbb{R} that you will learn carefully in more advanced courses. They take you from the “world of random processes and phenomena” to the world of real numbers. In other words, a random variable is a numerical value determined by the outcome of the experiment.

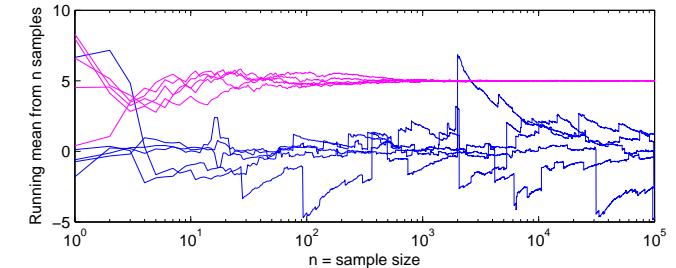
We said that a random variable can take one of many values, but we cannot be certain of which value it will take. However, *we can make probabilistic statements about the value x the random variable X will take*. A question like,

```

u=rand(1,N);           % draw N IID samples from Uniform(0,1)
x=tan(pi * u);        % draw N IID samples from Standard cauchy RV using inverse CDF
n=1:N;                 % make a vector n of current sample size [1 2 3 ... N-1 N]
CSx=cumsum(x); % CSx is the cumulative sum of the array x (type 'help cumsum')
% Runnign Means <- vector division of cumulative sum of samples by n
RunningMeanStdCauchy = CSx ./ n; % Running Mean for Standard Cauchy samples
RunningMeanUnif010 = cumsum(u*10.0) ./ n; % Running Mean for Uniform(0,10) samples
semilogx(n, RunningMeanStdCauchy) %
hold on;
semilogx(n, RunningMeanUnif010, 'm')
end
xlabel('n = sample size');
ylabel('Running mean from n samples')

```

Figure 5.5: Unending fluctuations of the running means based on n IID samples from the Standard Cauchy RV X in each of five replicate simulations (blue lines). The running means, based on n IID samples from the Uniform(0, 10) RV, for each of five replicate simulations (magenta lines).



The resulting plot is shown in Figure 5.5. Notice that the running means or the sample mean of n samples as a function of n , for each of the five replicate simulations, never settles down to a particular value. This is because of the “thick tails” of the density function for this RV which produces extreme observations. Compare them with the running means, based on n IID samples from the Uniform(0, 10) RV, for each of five replicate simulations (magenta lines). The latter sample means have settled down stably to the mean value of 5 after about 700 samples.

5.2.1 Application: Point Estimation of $E(X_1)$

LLN gives us a method to obtain a **point estimator** that gives “the single best guess” for the possibly unknown population mean $E(X_1)$ based on \bar{X}_n , the sample mean, of a simple random sequence (SRS) or independent and identically distributed (IID) sequence of n RVs $X_1, X_2, \dots, X_n \stackrel{\text{IID}}{\sim} X_1$.

Example 163 Let $X_1, X_2, \dots, X_n \stackrel{\text{IID}}{\sim} X_1$, where X_1 is an $\text{Exponential}(\lambda^*)$ RV, i.e., let

$$X_1, X_2, \dots, X_n \stackrel{\text{IID}}{\sim} \text{Exponential}(\lambda^*) .$$

Typically, we do not know the “true” parameter $\lambda^* \in \Lambda = (0, \infty)$ or the population mean $E(X_1) = 1/\lambda^*$. But by LLN, we know that

$$\bar{X}_n \rightsquigarrow \text{Point Mass}(E(X_1)) ,$$

($c \geq X \geq 7$)

Remark 20 (Notation) It is enough to understand the idea of random variables as explained above, and work with random variables using simplified notation like

Thus, $F(x)$ or simply F is a non-decreasing, right continuous, $[0, 1]$ -valued function over \mathbb{R} . When a RV X has DF F we write $X \sim F$.

$$(3.2) \quad \cdot \quad \text{for any } x \in \mathbb{A}, \quad \{x > (w)X : w\} = P = (x > X)P =: (x)F(x).$$

Definition 19 (Distribution Function) The Distribution Function (DF) or Cumulative Distribution Function (CDF) of any RV X , over a probability triple (Ω, \mathcal{F}, P) , denoted by F is:

$$P(X \leq x) = P(\{\omega : X(\omega) \leq x\}) \quad (3.1)$$

This definition can be summarized by the statement that a RV is an F -measurable map. We assign probability to the RV X as follows:

such that for every $x \in E$, the inverse image of the half-open real interval $(-\infty, x]$ is an element of the collection of events \mathcal{F} , i.e.:

To take advantage of our measurements over the real numbers, in terms of its metric structure and arithmetic, we need to formally define this measurement Process using the notion of a random variable.

3.1 Basic Definitions

With this motivation we are ready to formally define such a random variable.

65 (6) 22

or, more simply,

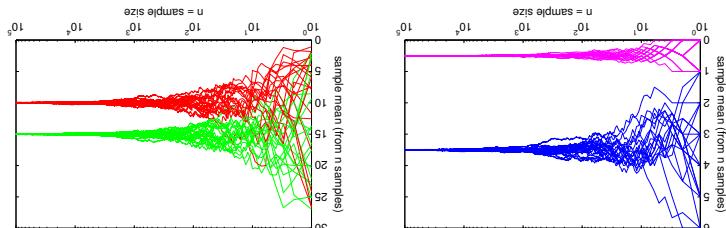
"*i*($\{\Gamma = (\omega)X : \omega\}$)". What is *P*?

"What is the probability of it raining tomorrow?"

Finally, we have shown that $E(e^{n\bar{X}_n}) = e^{nE(\bar{X}_1)}$, the CF of the n -sample mean RV \bar{X}_n , converges to infinity.

Heuristic Interpretation of LBN

Figure 5A: Sample mean \bar{X}_n as a function of sample size n for 20 replications from independent realizations of a fair die (blue), fair coin (magenta), Uniform(0, 30) RV (green) and Exponential(0.1) RV (red) with population means $(1+2+3+4+5+6)/6 = 21/6 = 3.5$, $(0+1)/2 = 0.5$, $(30-0)/2 = 15$ and $1/0.1 = 10$, respectively.



Labwork 162 (Running mean of the Standard Cauchy RV) Let us see what happens when we plot the running sample mean for an increasing sequence of IID samples from the Standard Cauchy RV X by implementing the following script file:

Recall that the mean of the Cauchy RV X does not exist since $\int |x| dF(x) = \infty$ (3.55). We will

Cauchy whose expectations does not exist has no Law of Large Numbers

Example 161 (Bernoulli WLLN and Galton's Quincunx) We can appreciate the WLLN for Bernoulli's Quincunx. We drop n balls into a binomial distribution $\text{Bin}(n, p)$. Using the paths of balls dropped into the binomial distribution, we can approximate the WLLN for Bernoulli's Quincunx.

rather than

$$P(\{\omega : 2 \leq X(\omega) \leq 3\})$$

but note that when learning or doing more advanced work this sample space notation is usually needed to clarify the true meaning of the simplified notation at least to yourself! But in the exam you can use the simpler notation as done in the solutions to exercises.

From the idea of a distribution function, we get:

Proposition 21 The probability that the random variable X takes a value x in the half-open interval $(a, b]$, i.e., $a < x \leq b$, is:

$$P(a < X \leq b) = F(b) - F(a) . \quad (3.3)$$

Proof: Since $(X \leq a)$ and $(a < X \leq b)$ are disjoint events whose union is the event $(X \leq b)$,

$$F(b) = P(X \leq b) = P(X \leq a) + P(a < X \leq b) = F(a) + P(a < X \leq b) .$$

Subtraction of $F(a)$ from both sides of the above equation yields Equation 3.3.

A special RV that often plays the role of ‘building-block’ in Probability and Statistics is the indicator function of an event A that tells us whether the event A has occurred or not. Recall that an event belongs to the collection of possible events \mathcal{F} for our experiment.

Definition 22 (Indicator Function) Given a probability triple $(\Omega, \mathcal{F}, \mathbf{P})$, the **Indicator Function** of an event $A \in \mathcal{F}$ which is denoted $\mathbb{1}_A$ is defined as follows:

$$\mathbb{1}_A(\omega) := \begin{cases} 1 & \text{if } \omega \in A \\ 0 & \text{if } \omega \notin A \end{cases} \quad (3.4)$$

Model 1 (Indicator of an event as Bernoulli RV) This is the most primitive RV from which all others are obtained. Let us convince ourselves that $\mathbb{1}_A$ is really a RV. For $\mathbb{1}_A$ to be a RV, we need to verify that for any real number $x \in \mathbb{R}$, the inverse image $\mathbb{1}_A^{[-1]}((-\infty, x])$ is an event, ie :

$$\mathbb{1}_A^{[-1]}((-\infty, x]) := \{\omega : \mathbb{1}_A(\omega) \leq x\} \in \mathcal{F} .$$

All we can assume about the collection of events \mathcal{F} is that it contains the event A and that it is a sigma algebra. A careful look at the Figure 3.1 yields:

$$\mathbb{1}_A^{[-1]}((-\infty, x]) := \{\omega : \mathbb{1}_A(\omega) \leq x\} = \begin{cases} \emptyset & \text{if } x < 0 \\ A^c & \text{if } 0 \leq x < 1 \\ A \cup A^c = \Omega & \text{if } 1 \leq x \end{cases}$$

Thus, $\mathbb{1}_A^{[-1]}((-\infty, x])$ is one of the following three sets that belong to \mathcal{F} ; (1) \emptyset , (2) A^c and (3) Ω depending on the value taken by x relative to the interval $[0, 1]$. We have proved that $\mathbb{1}_A$ is indeed a RV.

Model 1 is called the Bernoulli RV for event A with a known probability $\mathbf{P}(A)$. We will define as our next model the Bernoulli(θ) RV by introducing a parameter $\theta \in [0, 1]$ for the typically unknown probability $\mathbf{P}(A)$.

Therefore, for any given $\epsilon > 0$,

$$\begin{aligned} \mathbf{P}(|\bar{X}_n - \mathbf{E}(X_1)| \geq \epsilon) &= \mathbf{P}(|\bar{X}_n - \mathbf{E}(\bar{X}_n)| \geq \epsilon) \quad [\text{by the IID assumption of } X_1, X_2, \dots, \mathbf{E}(\bar{X}_n) = \mathbf{E}(X_1), \text{ as per (3.76)}] \\ &= \frac{\frac{1}{n}\mathbf{V}(X_1)}{\epsilon^2} \rightarrow 0, \quad \text{as } n \rightarrow \infty , \end{aligned}$$

or equivalently, $\lim_{n \rightarrow \infty} \mathbf{P}(|\bar{X}_n - \mathbf{E}(X_1)| \geq \epsilon) = 0$. And the last statement is the definition of the claim made by the law of large numbers (LLN), namely that $\bar{X}_n \xrightarrow{\mathbf{P}} \mathbf{E}(X_1)$.

Proposition 69 (Weak Law of Large Numbers (WLLN): $\bar{X}_n \rightsquigarrow \text{Point Mass}(\mathbf{E}(X_1))$) If we are given a sequence of independently and identically distributed (IID) RVs, $X_1, X_2, \dots \stackrel{\text{IID}}{\sim} X_1$ and if $\mathbf{E}(X_1)$ exists, i.e. $\mathbf{E}(\text{abs}(X)) < \infty$, then the sample mean \bar{X}_n converges in distribution to the expectation of any one of the IID RVs, say $\text{Point Mass}(\mathbf{E}(X_1))$ by convention. More formally, we write:

If $X_1, X_2, \dots \stackrel{\text{IID}}{\sim} X_1$ and if $\mathbf{E}(X_1)$ exists, then $\bar{X}_n \rightsquigarrow \text{Point Mass}(\mathbf{E}(X_1))$ as $n \rightarrow \infty$.

Proof: Our proof now is based on the convergence of characteristic functions (CFs) pointwise to the CF of the limiting RV, as this implies, by Lévy’s Continuity Theorem on CFs⁷, the convergence of the corresponding distribution functions (DFs).

First, the CF of $\text{Point Mass}(\mathbf{E}(X_1))$ is

$$\mathbf{E}(e^{it\mathbf{E}(X_1)}) = e^{it\mathbf{E}(X_1)} ,$$

since $\mathbf{E}(X_1)$ is just a constant, i.e., a Point Mass RV that puts all of its probability mass at $\mathbf{E}(X_1)$.

Second, the CF of \bar{X}_n is

$$\begin{aligned} \mathbf{E}(e^{it\bar{X}_n}) &= \mathbf{E}\left(e^{it\frac{1}{n}\sum_{k=1}^n X_k}\right) = \mathbf{E}\left(\prod_{k=1}^n e^{itX_k/n}\right) = \prod_{k=1}^n \mathbf{E}\left(e^{itX_k/n}\right) = \prod_{k=1}^n \varphi_{X_1}(t/n) \\ &= \prod_{k=1}^n \varphi_{X_1}(t/n) = (\varphi_{X_1}(t/n))^n . \end{aligned}$$

Let us recall Landau’s “small o” notation for the relation between two functions. We say, $f(x)$ is small o of $g(x)$ if f is dominated by g as $x \rightarrow \infty$, i.e., $\frac{|f(x)|}{|g(x)|} \rightarrow 0$ as $x \rightarrow \infty$. More formally, for every $\epsilon > 0$, there exists an x_ϵ such that for all $x > x_\epsilon$ $|f(x)| < \epsilon |g(x)|$. For example, $\log(x)$ is $o(x)$, x^2 is $o(x^3)$ and x^m is $o(x^{m+1})$ for $m \geq 1$.

Third, we can expand any CF whose expectation exists as a Taylor series with a remainder term that is $o(t)$ as follows:

$$\varphi_X(t) = 1 + it\mathbf{E}(X) + o(t) .$$

Hence,

$$\varphi_{X_1}(t/n) = 1 + it\frac{1}{n}\mathbf{E}(X_1) + o\left(\frac{t}{n}\right)$$

and

$$E\left(e^{it\bar{X}_n}\right) = \left(1 + it\frac{1}{n}\mathbf{E}(X_1) + o\left(\frac{t}{n}\right)\right)^n \rightarrow e^{it\mathbf{E}(X_1)} \text{ as } n \rightarrow \infty .$$

⁷https://en.wikipedia.org/wiki/L%C3%A9vy_continuity_theorem

probabilistic
the derivative case
https://en.wikipedia.org/wiki/Proofs_of_convergence_in_probability
https://en.wikipedia.org/wiki/Random_variables#Convergence_in_distribution
https://en.wikipedia.org/wiki/Convergence_in_probability
https://en.wikipedia.org/wiki/Convergence_in_mean
https://en.wikipedia.org/wiki/Convergence_in_probability#Implications_of_convergence_in_probability

$$\mathbf{P}(\underline{X}_n - \mathbf{E}(X_n) \geq \epsilon) = \frac{\mathbf{E}[e^{\underline{X}_n}]}{\mathbf{E}(e^{\underline{X}_n})}$$

[by applying Chebyshev's inequality (5.6) to the RV \underline{X}_n]

[by the IID assumption of X_1, X_2, \dots , we can apply (3.77)]

If $X_1, X_2, \dots \stackrel{\text{iid}}{\sim} X_1$ and if $\mathbf{E}(X_1)$ exists, then $\underline{X}_n \xrightarrow{\text{P}} \mathbf{E}(X_1)$.

More formally, we write:

X_n converges in probability to the expectation of any one of the IID RVs, say $\mathbf{E}(X_1)$ by convention. (X_n , i.e., $\mathbf{E}(\text{abs}(X_1)) < \infty$, and the variance is finite, i.e., $\mathbf{V}(X_1) < \infty$, then the sample mean (\bar{X}_n) converges in probability distributed RVs, $X_1, X_2, \dots, \bar{X}_n$, and if $\mathbf{E}(X_1)$ exists, as per it independent and identically distributed RVs, $X_1, X_2, \dots, \bar{X}_n$, say $\mathbf{E}(X_1)$ by convention.

Proposition 68 (Law of Large Numbers (LLN): $\underline{X}_n \xrightarrow{\text{P}} \mathbf{E}(X_1)$) If we are given a sequence

5.2 Law of Large Numbers

- In general, convergence in distribution does not imply convergence in probability.

$$X_n \rightsquigarrow \text{Point Mass}(\theta) \iff X_n \xrightarrow{\text{P}} \text{Point Mass}(\theta).$$

- Convergence in distribution to a constant θ implies convergence in probability to θ .⁶

$$\underline{X}_n \xrightarrow{\text{P}} X \iff X_n \rightsquigarrow X.$$

- Convergence in probability implies convergence in distribution.⁵

By the Borel-Cantelli Lemma³, convergence in probability does not imply almost sure convergence in the discrete case.⁴

$$X_n \xrightarrow{\text{a.s.}} X \iff X_n \xrightarrow{\text{P}} X.$$

- Convergence almost surely implies convergence in probability.²

We will merely state some properties (without proofs) that are hyper-linked for the curious student as they are advanced for this course) and relations between the three notions of convergence with some examples to better appreciate the subtleties among them. You will study the proofs of these statements in Probability Theory II. Just remember that subtle implications exist between the three notions.

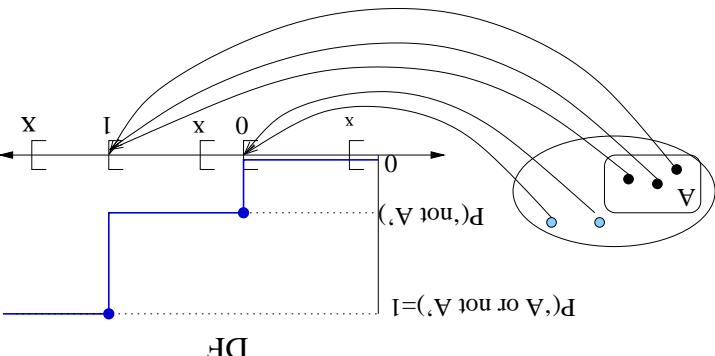


Figure 3.1: The indicator function of event A is a RV \mathbb{I}_A with DF

We slightly abuse notation when A is a single element set by ignoring the curly braces.

$$\mathbb{I}_A^c = 1 - \mathbb{I}_A, \quad \mathbb{I}_{A \cup B} = \mathbb{I}_A + \mathbb{I}_B - \mathbb{I}_A \mathbb{I}_B$$

Some useful properties of the indicator function are:

5.1.1 Properties of Convergence of RVs**

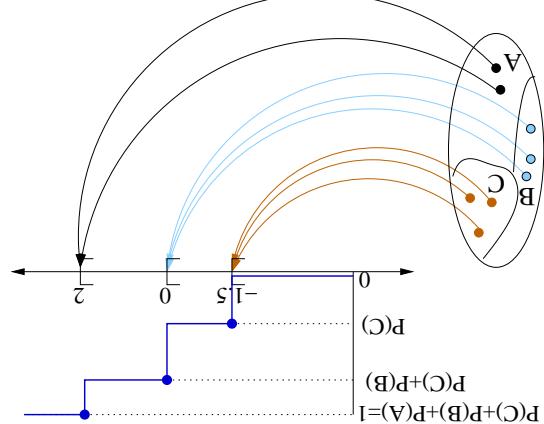


Figure 3.2: A RV X from a sample space Ω with 8 elements to \mathbb{R} and its DF F .

Classwork 41 (A random variable with three values and eight sample points) Consider the RV X of Figure 3.2. First draw this property as done in Ex. 3.1. Let the events $A = \{w_1, w_2\}$, $B = \{w_3, w_4, w_5\}$ and $C = \{w_6, w_7, w_8\}$. Define the RV X formally. What sets should F minimally include? What do you need to do to make sure that F is a sigma algebra?

Exercise 3.1 (Drawing discontinuous functions) Identify the mistakes in how the LA is drawn as a discontinuous function in Figure 3.1.

Exercise 3.1 (Drawing discontinuous functions) Identify the mistakes in how the DF is drawn as a discontinuous function in Figure 3.1.

Exercise 3.2 (Fair coin toss RV) Consider the *fair coin toss experiment* with $\Omega = \{\text{H}, \text{T}\}$ and $P(\text{H}) = P(\text{T}) = 1/2$.

We can associate a Bernoulli random variable X (in Model 1) for the event that the coin lands as H, with this experiment as follows:

$$X(\omega) = \begin{cases} 1, & \text{if } \omega = \text{H} \\ 0, & \text{if } \omega = \text{T} \end{cases}$$

Find the distribution function for X .

3.2 Discrete Random Variables

When a RV takes at most countably many values from a discrete set \mathbb{X} , we call it a **discrete** RV. Recall that a set \mathbb{X} is said to be discrete if we can enumerate its elements, i.e., find an enumerating or counting function $\mathbb{X} \ni x \mapsto i \in \mathbb{N}$ that associates each element $x \in \mathbb{X}$ to a natural number $i \in \mathbb{N}$. So, \mathbb{X} is either finite with k elements in $\mathbb{X} = \{x_1, x_2, \dots, x_k\}$ or countably infinite with the same cardinality as \mathbb{N} with $\mathbb{X} = \{x_1, x_2, \dots\}$. When $\mathbb{X} \subset \mathbb{R}$, we have a real-valued or \mathbb{R} -valued discrete random variable.

Definition 23 (probability mass function (PMF)) Let X be a \mathbb{R} -valued discrete RV over a probability triple $(\Omega, \mathcal{F}, \mathbf{P})$. We define the **probability mass function** (PMF) f of X to be the function $f : \mathbb{R} \rightarrow [0, 1]$ defined as follows:

$$f(x) := \mathbf{P}(X = x) = \mathbf{P}(\{\omega : X(\omega) = x\}) = \begin{cases} \theta_i & \text{if } x = x_i \in \mathbb{X}. \\ 0 & \text{otherwise.} \end{cases} \quad (3.5)$$

The DF F and PMF f for a discrete RV X satisfy the following:

1. For any $x \in \mathbb{R}$,

$$\mathbf{P}(X \leq x) = F(x) = \sum_{x_i \leq x} f(x_i) = \sum_{x_i \leq x} \theta_i. \quad (3.6)$$

2. For any $a, b \in \mathbb{R}$ with $a < b$,

$$\mathbf{P}(a < X \leq b) = F(b) - F(a) = \sum_{a < x_i \leq b} \theta_i. \quad (3.7)$$

This is just the sum of all probabilities θ_i for which x_i satisfies $a < x_i \leq b$.

3. From the fact that $\mathbf{P}(\Omega) = 1$, we get that the sum of all the probabilities is 1:

$$\sum_i \theta_i = 1. \quad (3.8)$$

Proposition 66 (Chebychev's Inequality) For any RV X and any $\epsilon > 0$,

$$\mathbf{P}(|X| > \epsilon) \leq \frac{\mathbf{E}(|X|)}{\epsilon} \quad (5.4)$$

$$\mathbf{P}(|X| > \epsilon) = \mathbf{P}(X^2 \geq \epsilon^2) \leq \frac{\mathbf{E}(X^2)}{\epsilon^2} \quad (5.5)$$

$$\mathbf{P}(|X - \mathbf{E}(X)| \geq \epsilon) = \mathbf{P}((X - \mathbf{E}(X))^2 \geq \epsilon^2) \leq \frac{\mathbf{E}(X - \mathbf{E}(X))^2}{\epsilon^2} = \frac{\mathbf{V}(X)}{\epsilon^2} \quad (5.6)$$

Proof: All three forms of Chebychev's inequality are mere corollaries (careful reapplications) of Markov's inequality.

Armed with Markov's inequality we next enquire the convergence in probability for the sequence of RVs in Classwork 157 and Example 158.

Example 160 (Convergence in probability) Does the the sequence of RVs $\{X_n\}_{n=1}^{\infty}$, where $X_n \sim \text{Normal}(0, 1/n)$, converge in probability to $X \sim \text{Point Mass}(0)$, i.e. does $X_n \xrightarrow{\mathbf{P}} X$?

To find out if $X_n \xrightarrow{\mathbf{P}} X$, we need to show that for any $\epsilon > 0$, $\lim_{n \rightarrow \infty} \mathbf{P}(|X_n - X| > \epsilon) = 0$.

Let ϵ be any real number greater than 0, then

$$\begin{aligned} \mathbf{P}(|X_n| > \epsilon) &= \mathbf{P}(|X_n|^2 > \epsilon^2) \\ &\leq \frac{\mathbf{E}(X_n^2)}{\epsilon^2} \quad [\text{by Markov's Inequality (5.2)}] \\ &= \frac{1}{\epsilon^2} \rightarrow 0, \quad \text{as } n \rightarrow \infty \quad [\text{in the sense of Definition 4}]. \end{aligned}$$

Hence, we have shown that for any $\epsilon > 0$, $\lim_{n \rightarrow \infty} \mathbf{P}(|X_n - X| > \epsilon) = 0$ and therefore by Definition 64, $X_n \xrightarrow{\mathbf{P}} X$ or $X_n \xrightarrow{\mathbf{P}} 0$.

Convention: When X has a Point Mass(θ) distribution and $X_n \xrightarrow{\mathbf{P}} X$, we simply write $X_n \xrightarrow{\mathbf{P}} \theta$.

Definition 67 (Convergence Almost Surely (or with Probability 1)) To say that the sequence of RVs $\{X_n\}_{n=1}^{\infty}$ converges almost surely (or with probability 1 or strongly) towards another RV X on the same probability space $(\Omega, \mathcal{F}, \mathbf{P})$, as denoted by

$$X_n \xrightarrow{a.s.} X$$

means that

$$\mathbf{P}\left(\lim_{n \rightarrow \infty} X_n = X\right) = 1 \iff \mathbf{P}\left(\{\omega \in \Omega : \lim_{n \rightarrow \infty} X_n(\omega) = X(\omega)\}\right) = 1.$$

This means that the values of X_n approach the value of X , in the sense that events for which X_n does not converge to X have probability 0.

Other notions of convergence are termed sure convergence or pointwise convergence, such as convergence in mean. But the above three types of convergence are elementary.

<p>Model 2 (Discrete Uniform) We say that a discrete random variable X is uniformly distributed over k possible values in $\mathbb{X} = \{x_1, x_2, \dots, x_k\}$ if its probability mass function is:</p> $f(x) = \begin{cases} 0 & \text{otherwise} \\ \frac{1}{k} & \text{if } x = x_i, \text{ where } i = 1, 2, \dots, k \end{cases} \quad (3.11)$
--

- Discrete non-uniform random variables with countably infinite many possibilities
- Discrete non-uniform random variables with finitely many possibilities
- Discrete uniform random variables with finitely many possibilities

Out of the class of discrete random variables we will define specific kinds as they arise often in applications. We classify discrete random variables into three types for convenience as follows:

Note that this table hides the more complex notation but it is still there, under the surface. In Probability Theory I, you should be able to work with and manipulate discrete random variables using the simplified notation given above. The same comment applies to the continuous random variables discussed later. But you are free students of mathematics and should know more about what is "under the hood".

It is customary to use P , instead of θ , for the probabilities. But we try to avoid it as it will hurt us when we start doing Measure Theory soon!

Probability: $P(X = x_i) = \theta_i$	θ_1	θ_2	θ_3	\dots
Possible values: x_i	x_1	x_2	x_3	\dots

Hence, we can describe a discrete random variable by the table:
 Examples it is convenient to associate the possible values x_1, x_2, \dots with the outcomes $\omega_1, \omega_2, \dots$.
 Variables are defined as functions, is much reduced. The reason is that in straightforward
 Notice that in equations (3.5), (3.6) and (3.7) the use of the " $\omega \in \Omega$ " notation, where random
 It is customary to use P , instead of θ , for the probabilities. But we try to avoid it as it will
 hurt us when we start doing Measure Theory soon!

DISCRETE RANDOM VARIABLES - SIMPLIFIED NOTATION

See bits: <https://en.wikipedia.org/wiki/Simplex> for the images scribed on the board.

$$\Delta_1 := \{(\theta, 1-\theta) \in \mathbb{R}^2 : 0 \leq \theta \leq 1\}. \quad (3.10)$$

In particular when X has only two possible values with $\mathbb{X} = \{x_1, x_2\}$ then $\theta_2 = 1 - \theta_1$, so we can avoid subscripts and take $\theta := \theta_1$ and realize that the probability P is now specified by the point $(\theta, 1 - \theta)$ in the unit 1 simplex:

$$\Delta_{k-1} := \{(\theta_1, \theta_2, \dots, \theta_k) \in \mathbb{R}^k : \sum_i \theta_i = 1 \text{ and } \theta_i \geq 0, \text{ for all } i\} \quad (3.9)$$

think of the probability P specified by $(\theta_1, \theta_2, \dots, \theta_k)$ as a point in the unit $(k-1)$ simplex:

4. When X only has finitely many possibilities, say k with $\mathbb{X} = \{x_1, x_2, \dots, x_k\}$, then we may

As $n \rightarrow \infty$, the expression below the first overbrace $\rightarrow 1$, while the second overbrace

being independent of n remains the same. By the elementary examples of limits 17 and 18, as $n \rightarrow \infty$, the expression over the first underbrace approaches $e^{-\lambda}$ while that over the second underbrace

approaches 1. Finally, we get the desired limit:

Let us look at some immediate consequences of Markov's inequality.

$$E(X) \leq E(\mathbb{I}_{\{\theta_i \geq \epsilon\}}(x)) = \mathbb{P}(X \geq \epsilon).$$

get the desired result:

$$\begin{aligned} & \geq \mathbb{P}(\mathbb{I}_{\{\theta_i \geq \epsilon\}}(x)) \\ & \geq X \mathbb{I}_{\{\theta_i \geq \epsilon\}}(x) + X \mathbb{I}_{\{\theta_i < \epsilon\}}(x) \\ X & = X \mathbb{I}_{\{\theta_i \geq \epsilon\}}(x) + X \mathbb{I}_{\{\theta_i < \epsilon\}}(x) \end{aligned}$$

Proof:

$$P(X \geq \epsilon) \leq \frac{\epsilon}{E(X)}, \quad \text{for any } \epsilon > 0. \quad (5.2)$$

be a non-negative RV. Then,

Proposition 65 (Markov's Inequality) Let (Ω, \mathcal{F}, P) be a probability triple and let $X = X(\omega)$ need some elementary inequalities in Probability to help us answer this question. We visit these inequalities next.

For the same sequence of RVs in Classwork 157 and Example 158 we are tempted to ask whether $X_n \sim \text{Normal}(0, 1/n)$ converges in probability to $X \sim \text{Point Mass}(0)$, i.e. whether $X_n \xrightarrow{P} X$. We

$$\lim_{n \rightarrow \infty} P(\{\omega : |X_n(\omega) - X(\omega)| < \epsilon\}) = 0, \quad \text{ie, } P(\{\omega : |X_n(\omega) - X(\omega)| < \epsilon\}) \rightarrow 0, \quad \text{as } n \rightarrow \infty.$$

Once again, the above limit, by (3.1) in our Definition 18 of a RV, can be equivalently expressed as follows:

$$\lim_{n \rightarrow \infty} P(|X_n - X| < \epsilon) = 0 \quad [\text{in the sense of Definition 4}].$$

If for every real number $\epsilon > 0$,

$$X_n \xrightarrow{P} X$$

converges to X in probability, and write:

be another RV. Let F_n denote the DF of X_n and F denote the DF of X . The we say that X_n Definition 64 (Convergence in Probability) Let X_1, X_2, \dots , be a sequence of RVs and let X_n

The second notion of convergence of RVs is convergence in probability.

$$\lim_{n \rightarrow \infty} P(X = x) = e^{-\lambda x}.$$

Finally, we get the desired limit:

As $n \rightarrow \infty$, the expression below the first overbrace approaches $e^{-\lambda}$ while that over the second underbrace approaches 1. By the elementary examples of limits 17 and 18, as $n \rightarrow \infty$, the expression over the first underbrace approaches $e^{-\lambda}$ while that over the second underbrace

The distribution function for the discrete uniform random variable X is:

$$F(x) = \sum_{x_i \leq x} f(x_i) = \sum_{x_i \leq x} \theta_i = \begin{cases} 0 & \text{if } -\infty < x < x_1 , \\ \frac{1}{k} & \text{if } x_1 \leq x < x_2 , \\ \frac{2}{k} & \text{if } x_2 \leq x < x_3 , \\ \vdots & \\ \frac{k-1}{k} & \text{if } x_{k-1} \leq x < x_k , \\ 1 & \text{if } x_k \leq x < \infty . \end{cases} \quad (3.12)$$

The discrete uniform RV with values in $\mathbb{X} = \{1, 2, \dots, k\}$ is called the equi-probable de Moivre(k) RV as we will see in the sequel.

Example 42 The *fair coin toss experiment* of Exercise 3.2 is an example of a discrete uniform random variable with finitely many possibilities. Its probability mass function is given by

$$f(x) = \mathbf{P}(X = x) = \begin{cases} \frac{1}{2} & \text{if } x = 0 \\ \frac{1}{2} & \text{if } x = 1 \\ 0 & \text{otherwise} \end{cases}$$

and its distribution function is given by

$$F(x) = \mathbf{P}(X \leq x) = \begin{cases} 0, & \text{if } -\infty < x < 0 \\ \frac{1}{2}, & \text{if } 0 \leq x < 1 \\ 1, & \text{if } 1 \leq x < \infty \end{cases}$$

Let us sketch the probability mass function and distribution function for X below.

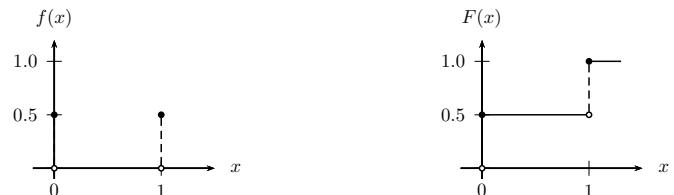


Figure 3.3: $f(x)$ and $F(x)$ of the *fair coin toss* random variable X , a discrete uniform RV on $\{0, 1\}$.

Example 43 (Fair dice RV) Now consider the *toss a fair die* experiment and define X to be the number that shows up on the top face. Note that here Ω is the set of numerical symbols $\{1, 2, 3, 4, 5, 6\}$ that label each face while each of these symbols are associated with the real number $x \in \{1, 2, 3, 4, 5, 6\}$. We can describe this random variable by the table

Possible values, x_i	1	2	3	4	5	6
Probability, θ_i	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$

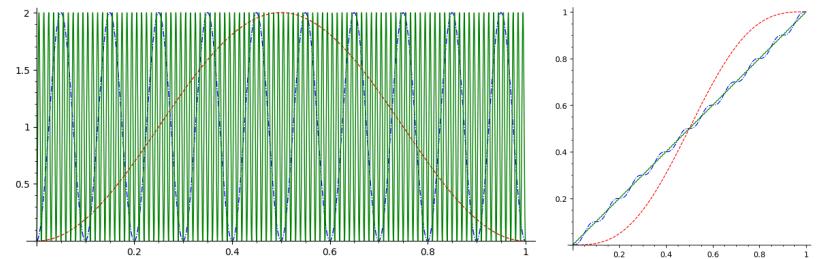


Figure 5.3: $\text{PDF } f_{X_n}(x) := \mathbf{1}_{(0,1)}(x)(1 - \cos(2\pi nx))$ of the RV X_n [the left sub-figure] and its $F_n(x) := \int_{-\infty}^x \mathbf{1}_{(0,1)}(v)(1 - \cos(2\pi nv))dv$ [the right sub-figure], for $n = 1$ [red '---'], $n = 10$ [blue '-.-'], and $n = 100$ [green '-'], respectively. One can see clear convergence of the DFs F_n to $\mathbf{1}_{(0,1)}(x)x$, the DF of the Uniform($0, 1$) RV, while the corresponding PDFs $f_n(x)$ keep oscillating wildly with n across $[0, 2]$ about $\mathbf{1}_{(0,1)}(x)$, the PDF of the Uniform($0, 1$) RV X . Thus giving a counter-example to the claim that convergence in DFs does not imply convergence in PDFs.

Since $F(x) = \mathbf{P}(X \leq x)$, convergence in distribution means that the probability for X_n to be in a given range is approximately equal to the probability that the value of the limiting RV X is in that range, provided n is sufficiently large.

Thus, for a discrete sequence of RVs X_n 'n to converge in distribution to another discrete RV X taking values in $\mathbb{Z}_+ = \{0, 1, 2, \dots\}$, it is sufficient to show that $\lim_{n \rightarrow \infty} \mathbf{P}(X_n = x) = \mathbf{P}(X = x)$ for each $x \in \mathbb{Z}_+$. We will use this fact to prove why we can approximate Binomial RVs by a Poisson under some limiting conditions.

Example 159 ($\text{Binomial}(n, \lambda/n) \rightsquigarrow \text{Poisson}(\lambda)$) In several situations, as we saw already, it becomes cumbersome to model the events using the $\text{Binomial}(n, \theta)$ RV, especially when the parameter $\theta \propto 1/n$ and the events become rare.

$\text{Binomial}(n, \lambda/n)$ converges in distribution to $\text{Poisson}(\lambda)$ as $n \rightarrow \infty$, $\theta = \lambda/n \rightarrow 0$

However, for some real parameter $\lambda > 0$, the $\text{Binomial}(n, \lambda/n)$ RV with probability of the number of successes in n trials, with per-trial success probability λ/n , approaches the Poisson distribution with expectation λ , as n approaches ∞ (actually, it converges in distribution). The $\text{Poisson}(\lambda)$ RV is much simpler to work with than the combinatorially laden $\text{Binomial}(n, \theta = \lambda/n)$ RV. We sketch the details of this next.

Let $X_n \sim \text{Binomial}(n, \theta = \lambda/n)$ and $Y \sim \text{Poisson}(\lambda)$ and let $\lambda = n\theta$ remain constant as $n \rightarrow \infty$, $\theta \rightarrow 0$. We need to show that $\lim_{n \rightarrow \infty} \mathbf{P}(X_n = x) = \mathbf{P}(Y = x) = e^{-\lambda} \lambda^x / x!$ for any $x \in \{0, 1, 2, 3, \dots, n\}$.

$$\begin{aligned} \mathbf{P}(X = x) &= \binom{n}{x} \left(\frac{\lambda}{n}\right)^x \left(1 - \frac{\lambda}{n}\right)^{n-x} \\ &= \frac{n(n-1)(n-2)\cdots(n-x+1)}{x(x-1)(x-2)\cdots(2)(1)} \left(\frac{\lambda^x}{n^x}\right) \left(1 - \frac{\lambda}{n}\right)^n \left(1 - \frac{\lambda}{n}\right)^{-x} \\ &= \overbrace{\binom{n}{x} \left(\frac{n-1}{n}\right) \left(\frac{n-2}{n}\right) \cdots \left(\frac{n-x+1}{n}\right)}^{\left(\frac{\lambda^x}{x!}\right)} \underbrace{\left(1 - \frac{\lambda}{n}\right)^n}_{\left(\frac{1-\lambda}{n}\right)^n} \underbrace{\left(1 - \frac{\lambda}{n}\right)^{-x}}_{\left(\frac{1-\lambda}{n}\right)^{-x}} \end{aligned} \quad (5.1)$$

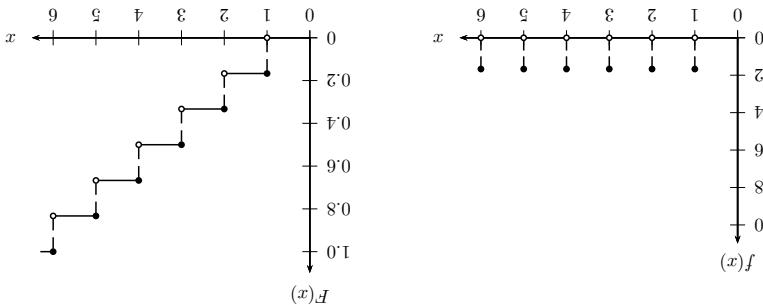
to rest on only four sides, the other two sides being rounded. The upper side of the bone, broad games was at first provided by tossing astreaghi, the ankle bones of sheep. These bones could come chance were known in Egypt, 3000 years before Christ. The location of chance needed for these RVs X_n , to $f(x)$, the PMF of another discrete RV X , implies convergence in their corresponding DFs, i.e., $F_n(x) \rightarrow F(x)$ for each x as $n \rightarrow \infty$.

Example 44 (Astreaghi with a Kiwi sheep ankle bone) Astreaghi. Board games involving

[See \[https://en.wikipedia.org/wiki/Scheff%C3%A9%27s_lemma\]\(https://en.wikipedia.org/wiki/Scheff%C3%A9%27s_lemma\)](https://en.wikipedia.org/wiki/Scheff%C3%A9%27s_lemma).

in generality¹. However, you should be able to see why convergence of PMFs $f_n(x)$ for discrete RVs X_n , to $f(x)$, the PMF of another discrete RV X , implies convergence in their corresponding DFs, i.e., $F_n(x) \rightarrow F(x)$ for each x as $n \rightarrow \infty$.

Figure 3.4: $f(x)$ and $F(x)$ of the fair die toss random variable X , a discrete uniform RV on $\{1, 2, 3, 4, 5, 6\}$.



$$F(x) = \begin{cases} 0, & \text{if } -\infty < x < 1 \\ \frac{1}{6}, & \text{if } 1 \leq x < 2 \\ \frac{2}{6}, & \text{if } 2 \leq x < 3 \\ \frac{3}{6}, & \text{if } 3 \leq x < 4 \\ \frac{4}{6}, & \text{if } 4 \leq x < 5 \\ \frac{5}{6}, & \text{if } 5 \leq x < 6 \\ 1, & \text{if } 6 \leq x < \infty \end{cases}$$

and the distribution function is:

$$f(x) = \begin{cases} 0 & \text{otherwise} \\ \frac{1}{6} & \text{if } x = 1 \\ \frac{1}{6} & \text{if } x = 2 \\ \frac{1}{6} & \text{if } x = 3 \\ \frac{1}{6} & \text{if } x = 4 \\ \frac{1}{6} & \text{if } x = 5 \\ \frac{1}{6} & \text{if } x = 6 \\ 0 & \text{if } x = 7 \end{cases}$$

The probability mass function of this random variable is:

Solution:

Find the probability mass function and distribution function for this random variable, and sketch their graphs.

We can formalize our observation in Classwork 157 that X_n is concentrating about 0 as $n \rightarrow \infty$ by

Proposition 63 (Scheffé's Theorem) According to Scheffé's Theorem convergence of the probability density function (for a continuous RV) or probability mass function (for a discrete RV) implies convergence in distribution.

Convergence in distribution does not in general imply that the sequence of corresponding probability density functions will also converge. Consider for example RV X_n with density $\mathbb{I}_{(0,1)}(x)(1 - \cos(2\pi nx))$. These RVs converge in distribution to $X \sim \text{Uniform}(0, 1)$, but their densities (PDFs) do not converge at all as evident in Figure 3.3.

However, note that

$$\lim_{n \rightarrow \infty} F_n(t) = \lim_{n \rightarrow \infty} \mathbf{P}(Z > \sqrt{n}t) = 1 = F(t).$$

and so convergence fails at 0, i.e. $\lim_{n \rightarrow \infty} F_n(t) \neq F(t)$ at $t = 0$. But, $t = 0$ is not a continuity point of F and the definition of convergence in distribution only requires the convergence to hold at continuity points of F .

Thus, we have proved that $X_n \rightsquigarrow X$ by verifying that for any t at which the Point Mass(0) DF F is continuous, we also have the desired equality: $\lim_{n \rightarrow \infty} F_n(t) = F(t)$.

And, when $t < 0$, $F(t)$, being the constant 1 function over the interval $(0, \infty)$, is again continuous at t . Since $\sqrt{n}t \rightarrow -\infty$, as $n \rightarrow \infty$,

$$\lim_{n \rightarrow \infty} F_n(t) = \lim_{n \rightarrow \infty} \mathbf{P}(Z > \sqrt{n}t) = 0 = F(t).$$

When $t < 0$, $F(t)$, being the constant 0 function over the interval $(-\infty, 0)$, is continuous at t . Since the only discontinuous point of F is 0 where F jumps from 0 to 1.

$$F_n(t) = \mathbf{P}(X_n > t) = \mathbf{P}(\sqrt{n}X_n > \sqrt{n}t) = \mathbf{P}(Z > \sqrt{n}t).$$

$$X_n \sim \text{Normal}(0, 1/n) \iff Z = \sqrt{n}X_n \sim \text{Normal}(0, 1),$$

First note that we need to verify that for any continuity point t of the Point Mass(0) DF F , $\lim_{n \rightarrow \infty} F_n(t) = F(t)$. Hence in distribution is satisfied for our sequence of RVs X_1, X_2, \dots and the limiting RV X . Thus, we need to verify that for any continuity point t of the Point Mass(0) DF F , $\lim_{n \rightarrow \infty} F_n(t) = F(t)$.

and slightly convex counted four; the opposite side broad and slightly concave counted three; the lateral side flat and narrow, one, and the opposite narrow lateral side, which is slightly hollow, six. You may examine an astragali of a kiwi sheep.

This is an example of a discrete non-uniform random variable with finitely many possibilities. A surmised probability mass function with $f(4) = \frac{4}{10}$, $f(3) = \frac{3}{10}$, $f(1) = \frac{2}{10}$, $f(6) = \frac{1}{10}$ and distribution function are shown below.

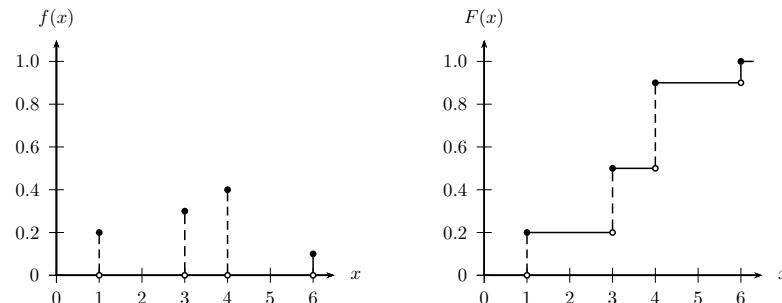


Figure 3.5: $f(x)$ and $F(x)$ of surmised *astragali toss* random variable X , a discrete (non-uniform) RV on $\{1, 2, 3, 4\}$.

3.2.1 An Elementary Family of Bernoulli Random Variables

In many experiments there are only two outcomes. For instance:

- Flip a coin to see whether it is defective.
- Roll a die and determine whether it is a 6 or not.
- Determine whether it will be below 0 degrees Celsius at 0600 hours in Uppsala tomorrow or not.

Performing such an experiment \mathcal{E} once to see if an event of interest A occurs is called a **Bernoulli trial** and its probability model over a triple $(\Omega, \mathcal{F}, \mathbf{P})$, with $A \in \mathcal{F}$, given by the Indicator Function $\mathbf{1}_A$ in Model 1 is called the Bernoulli RV.

If we do not know the probability θ that ‘ A occurs’, i.e., the Bernoulli RV will equal 1, then we can define a whole family of Bernoulli RVs for each $\theta \in [0, 1]$ or more precisely for each $(\theta, 1 - \theta) \in \Delta^1$, the unit 1-Simplex. Note that this family includes the fair Bernoulli trial of Example 42 when $\theta = 0.5$. Let us formalise this as the Bernoulli(θ) RV for each $\theta \in [0, 1]$ next.

Model 3 (Bernoulli(θ) RV) Given a parameter $\theta \in [0, 1]$, the probability mass function (PMF) for the Bernoulli(θ) RV X is:

$$f(x; \theta) = \theta^x (1 - \theta)^{1-x} \mathbf{1}_{\{0,1\}}(x) = \begin{cases} \theta & \text{if } x = 1, \\ 1 - \theta & \text{if } x = 0, \\ 0 & \text{otherwise} \end{cases} \quad (3.13)$$

Figure 5.1: Sequence of $\{\text{Point Mass}(17)\}_{i=1}^{\infty}$ RVs (left panel) and $\{\text{Point Mass}(1/i)\}_{i=1}^{\infty}$ RVs (only the first seven are shown on right panel) and their limiting RVs in red.

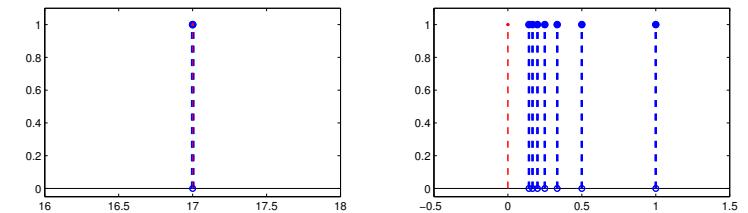
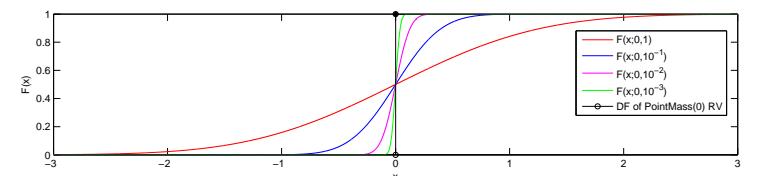


Figure 5.2: Distribution functions of several $\text{Normal}(\mu, \sigma^2)$ RVs for $\sigma^2 = 1, \frac{1}{10}, \frac{1}{100}, \frac{1}{1000}$.



Thus, we need more sophisticated notions of convergence for sequences of RVs. Two such notions are formalized next as they are minimal prerequisites for a clear understanding of two basic propositions in Statistics :

1. Law of Large Numbers,
2. Central Limit Theorem,

Definition 62 (Convergence in Distribution (or Weakly, or in Law)) Let X_1, X_2, \dots be a sequence of RVs and let X be another RV. Let F_n denote the DF of X_n and F denote the DF of X . Then we say that X_n converges to X in distribution, and write:

$$X_n \rightsquigarrow X$$

if for any real number t at which F is continuous,

$$\lim_{n \rightarrow \infty} F_n(t) = F(t) \quad [\text{in the sense of Definition 4}].$$

The above limit, by (3.2) in our Definition 19 of a DF, can be equivalently expressed as follows:

$$\lim_{n \rightarrow \infty} \mathbf{P}(\{\omega : X_n(\omega) \leq t\}) = \mathbf{P}(\{\omega : X(\omega) \leq t\}),$$

i.e. $\mathbf{P}(\{\omega : X_n(\omega) \leq t\}) \rightarrow \mathbf{P}(\{\omega : X(\omega) \leq t\}), \text{ as } n \rightarrow \infty$.

Let us revisit the problem of convergence in Classwork 157 armed with our new notions of convergence.

Example 158 (Convergence in distribution) Suppose you are given an independent sequence of RVs $\{X_i\}_{i=1}^n$, where $X_i \sim \text{Normal}(0, 1/i)$ with DF F_n and let $X \sim \text{Point Mass}(0)$ with DF F .

in the sense of Definition 1.7 about independence of a sequence of events, then we can obtain the dependence across trials, so one trial's outcome does not affect the outcome of any of the other trials, with each $\theta \in [0, 1]$ being possibly unknown but fixed as a parameter. Now, if we assume independence across trials, then we can obtain the number of Bernoulli trials, say,

$$X_i \sim \text{Bernoulli}(\theta), \quad i \in \mathbb{N},$$

Since the Bernoulli(θ) RV has only two outcomes, i.e., simple events, we know how to obtain the probability of each of the two outcomes in a given Bernoulli trial with the probability given by the probability of Bernoulli(θ) trials, say,

probabilistic variable or parameter θ . Now consider doing more than one trial so we have sequence of Bernoulli(θ) trials, say.

probabilistic variable or parameter θ . Note: we assume that flooding is independent from year to year, and that the probability of flooding is the same each year.

- Provide a property near a particular bridge in our archipelago with flood insurance for 20 years; count the number of years, during the 20-year period, during which the property is flooded. Note: we assume that flooding is independent from year to year, and that the probability of flooding is the same each year.
- Roll a die 100 times; count the number of sixes you throw.
- Test 50 randomly selected circuits from an assembly line; count the number of defective circuits.
- Manufactured in a terrible mint.
- Possibility allowing for the coins $P(H)$ to change each time because each of them are manufactured in a terrible mint.

We now look at what happens when we perform a sequence of independent Bernoulli trials. For instance:

Random variables make sense for a series of trials as well as just a single trial of an experiment.

3.2.2 Independent Bernoulli Trials



$$\text{and } P_\theta(X = 0) = 1 - \theta.$$

We emphasize the dependence of the probabilities on the parameter θ by specifying it following the semicolon in the argument for f and F and by subscripting the probabilities, i.e., $P_\theta(X = 1) = \theta$

$$(3.14) \quad F(x; \theta) = \begin{cases} 0 & \text{otherwise} \\ 1 - \theta & \text{if } 0 \leq x < 1, \\ 1 & \text{if } 1 \leq x, \end{cases}$$

and its DF is:

number in its support, such as 0. In other words, a continuous RV, such as X_n , has 0 probability of realizing any single real large. The answer is no. This is because $P(X_n = X) = 0$ for any n , since $X \sim \text{Point Mass}(0)$ is a discrete RV with exactly one outcome 0 and $X_n \sim \text{Normal}(0, 1/n)$ is a continuous RV for every n , however

limiting RV $X \sim \text{Point Mass}(0)$?

0, as depicted in Figure 3.2. Based on this observation, can we expect $\lim_{n \rightarrow \infty} X_n = X$, where the masses of X_n increasingly concentrates about 0 as n approaches ∞ ? Take a look at Figure 3.2 for a broad idea of $X_n \sim \text{Normal}(0, 1/n)$ as an approximation to X . The probability sequence of RVs $\{X_i\}_{i=1}^{\infty}$, where $X_i \sim \text{Normal}(0, 1/i)$, How would you talk about the convergence sequence of RVs $\{X_i\}_{i=1}^{\infty}$? Suppose you are given an independent

Y-axis not – just move to space of distributions over the reals! See Figure 3.1.

Examples 12 and 13?

Can the sequences of $\{\text{Point Mass}_i(\theta) = 17\}_{i=1}^{\infty}$ and $\{\text{Point Mass}_i(\theta) = 1/\iota\}_{i=1}^{\infty}$ RVs be the same as the two sequences of real numbers $\{x_i\}_{i=1}^{\infty} = 17, 17, 17, \dots$ and $\{x_i\}_{i=1}^{\infty} = \frac{1}{1}, \frac{2}{2}, \frac{3}{3}, \dots$ we saw in Section 1.8.1 before proceeding further.

Let us first refresh ourselves with notions of convergence, limits and continuity in the real line convergence_of_random_variables.

We need different notions of convergence to characterize such a behavior: two simplest behaviors are that values in the sequence eventually takes a constant value θ , i.e., X_n approaches $X \sim \text{Point Mass}(\theta)$ are that the sequence converges to a constant value θ , i.e., X_n approaches $X \sim F(x)$. See https://en.wikipedia.org/wiki/Probability_distributions.

From a statistical or decision-making viewpoint, as you will see in Inferential Theory I course, $n \rightarrow \infty$ is associated with the amount of data or information $\rightarrow \infty$. More abstractly, we are interested in what happens to the limiting RV $X := \lim_{n \rightarrow \infty} X_n$ when given the DFs $F_n(x)$ for each X_n .

$$\{X_i\}_{i=1}^{\infty} := X_1, X_2, X_3, \dots, X_{n-1}, X_n \quad \text{as } n \rightarrow \infty.$$

This important topic is concerned with the limiting behavior of sequences of RVs. We want to understand what it means for a sequence of random variables $\{X_n\}_{n=1}^{\infty} := X_1, X_2, \dots$ to converge to another random variable X , when all RVs are defined on the same probability space $(\Omega, \mathcal{F}, \mathbf{P})$. This means that the limit of the sequence of random variables $\{X_n\}_{n=1}^{\infty}$ is the same as the limit of the sequence of random variables $\{F_n\}_{n=1}^{\infty}$ to another random variable F , as you will see in Inferential Theory I course, $n \rightarrow \infty$.

5.1 Convergence of Random Variables

Limit Laws of Statistics

Chapter 5

probability of the entire sequence of outcomes for this sequence of **independently distributed Bernoulli(θ_i) trails** which can be any infinite sequence of 0's and 1's, i.e., any element of $\{0, 1\}^\infty$, by simply multiplying the corresponding probabilities given by θ_i 's in $(\theta_1, \theta_2, \dots) \in [0, 1]^\infty$, an infinite dimensional parameter space, as follows:

$$\mathbf{P}(x; (\theta_1, \theta_2, \dots)) = \prod_i f(x_i; \theta_i) = \prod_i \theta_i^{x_i} (1 - \theta_i)^{1-x_i} \mathbf{1}_{\{0,1\}}(x_i), \quad (3.15)$$

where $x := (x_1, x_2, \dots) \in \mathbb{X}_\infty = \{0, 1\}^\infty := \{0, 1\} \times \{0, 1\} \times \dots$

By further assuming that all the θ_i 's are identical, say $\theta = \theta_1 = \theta_2 = \dots$, with $\theta \in [0, 1]$, a one-dimensional parameter space, we get the much simpler expression for the **independent and identically distributed (IID) Bernoulli(θ) trails** as follows:

$$\begin{aligned} \mathbf{P}(x; \theta) &= \prod_i f(x_i; \theta) = \prod_i \theta^{x_i} (1 - \theta)^{1-x_i} \mathbf{1}_{\{0,1\}}(x_i) = \mathbf{1}_{\{0,1\}^\infty}(x) \theta^{\sum_i x_i} (1 - \theta)^{\sum_i (1-x_i)} \\ &= \begin{cases} \theta^{\sum_i x_i} (1 - \theta)^{n - \sum_i x_i} & \text{if } x := (x_1, x_2, \dots) \in \mathbb{X}_\infty = \{0, 1\}^\infty \\ 0 & \text{otherwise.} \end{cases} \end{aligned} \quad (3.16)$$

Remembering that all other RVs can be derived from such IID Bernoulli(θ) trials using $\theta = 1/2$, as we will see in the sequel, we are ready to take a tour through some common discrete and continuous random variables that are useful in many applications.

3.2.3 Some Common Discrete Random Variables

Let us start with the simplest example to fix ideas carefully.

Example 45 (Waiting For the First Heads) Suppose our experiment is to toss a fair coin independently and identically (that is, the same coin is tossed in essentially the same manner independent of the other tosses in each trial) as often as necessary until we have a head, H. Let the random variable X denote the *Number of trials until the first H appears*.

Let's first find the probability mass function of X .

Now X can take on the values $\{1, 2, 3, \dots\}$, so we have a non-uniform random variable with infinitely many possibilities. Since

$$\begin{aligned} f(1) &= \mathbf{P}(X = 1) = P(H) = \frac{1}{2}, \\ f(2) &= \mathbf{P}(X = 2) = P(TH) = \frac{1}{2} \cdot \frac{1}{2} = \left(\frac{1}{2}\right)^2, \\ f(3) &= \mathbf{P}(X = 3) = P(TTH) = \frac{1}{2} \cdot \frac{1}{2} \cdot \frac{1}{2} = \left(\frac{1}{2}\right)^3, \quad \text{etc.} \end{aligned}$$

the probability mass function of X is:

$$f(x) = \mathbf{P}(X = x) = \left(\frac{1}{2}\right)^x, \quad x = 1, 2, \dots.$$

In the previous Example, noting that we have independent trials, we get:

4. A related distribution, denoted by $N_-(\mu, \tau, \sigma^2)$, is the right-truncated normal distribution truncated on the right at τ . Describe how samples from $N_-(\mu, \tau, \sigma^2)$ can be obtained by simulating from an appropriate left-truncated normal distribution.
5. Write a MATLAB function that provides samples from a truncated normal distribution. The function should have the following inputs: number of samples required, left or right truncation, μ , σ^2 and τ .

4.4 Exercises in Simulation

Ex. 4.1 — Suppose the continuous RV X has PDF:

$$f_X(x) = (\pi(1 + x^2))^{-1}$$

Devise an algorithm to transform samples from Uniform(0, 1) RV to those from X . Present your answer as pseudo-code.

The appearance of a Geometric(θ) RV can be thought of as "the number of tosses needed before the outcome of the first 'Head', when tossing a coin with probability of 'Heads', equal to θ in a independent and identical manner".

$$\begin{aligned}
 & \lim_{x \rightarrow 0^+} \sum_{n=0}^{\infty} x^{n\theta} = \sum_{n=0}^{\infty} \lim_{x \rightarrow 0^+} x^{n\theta} = \sum_{n=0}^{\infty} 1^n = \sum_{n=0}^{\infty} 1 = \infty
 \end{aligned}$$

The above equality is a consequence of the geometric series identity (3.18) with $a = \theta$ and $\theta = 1 - \theta$:

$$f(x; \theta) = \begin{cases} \theta(1 - \theta)^x & \text{if } x \in \mathbb{Z}_+ \\ 0 & \text{otherwise} \end{cases} = (\theta : x) f$$

(3.17)

Model 4 (Geometric(θ) RV) Given a parameter $\theta \in (0, 1)$, the PDF of the Geometric(θ) RV X

This is called a **geometric random variable** with success probability parameter θ . We can spot a geometric distribution because there will be a sequence of independent trials with a constant probability of success. We are counting the number of trials until the first success appears. Let us define this random variable formally next.

More generally, let there be two possibilities, success (S) or failure (F), with $\mathbf{P}(S) = \theta$. $\mathbf{P}(F) = 1 - \theta$ so that:

$$\mathbf{P}(x) = \mathbf{P}(\overbrace{\mathbf{E}\cdots\mathbf{E}}^{x-1} S) = (\mathbf{I} - \theta)^{1-x} \mathbf{I} \theta.$$

¹ See also the discussion of the relationship between the two in the section on "Theoretical Implications" below.

2. Find the maximum expected acceptance probabilities for the following recruitment points,
 $t = 0, 0.5, 1, 1.5, 2, 2.5$ and 3 . What can you conclude about efficiency as t gets further out

$$\alpha = \frac{2}{\tau + \sqrt{\tau^2 + 4}}$$

¹¹. Shows that for sufficiently small $\tau_0 + \tau_1 = 0$, $\tau_0 = 1$ when $\tau_1 \geq 0$, the best choice of τ_0 maximizes the expected acceptance probability for the rejection sampler given by

$$\cdot \cdot \bar{\zeta}^{\bar{f}} \mathbb{I}((\underline{\zeta} - \bar{f}) \vee -) dx \partial \vee = (\underline{\zeta} \cdot \vee | \bar{f})$$

When $\tau < \mu$, the rejection sampler can readily be used to simulate from $N_{\tau}^{+}(\mu, \tau, \sigma^2)$ by simulating from $N_{\tau}(\mu, \sigma^2)$ until a number larger than τ is obtained. When $\tau > \mu$, however, this can be inefficient and increasesimply so as τ gets further out into the right tail. In this case, a more efficient approach is to use the rejection sampler with the following translated exponential distribution:

$$\int_{\mathbb{R}^n} \frac{|(\varphi/(n-\cdot))(\Phi - 1)|}{(\varphi(\zeta)/\zeta(n-x)-)dx} = (\varphi|_H, \tau_x)_J$$

Labwork 15c (Sampling from truncated normal distributions) [Histograms; Boxplots; Simulation of truncated normal variables; Statistics and Computing (1995), 5, 121-125] Let $N(\mu, \tau^2)$ denote the left-truncated normal distribution with truncation point τ and density given by

```

<>> mu=pi % set the desired mean parameter
<>> sigma=sqrt(2) % set the desired standard deviation parameter sigma
<>> n=1000 % number of samples
<>> x=randn(n,1) % generate n random numbers from a standard normal distribution
<>> y=x.*sigma+mu % add the mean and multiply by sigma
<>> hist(y,20) % plot a histogram of the generated data

```

Suppose we want samples from $X \sim \text{Normal}(\mu = 1, \sigma^2 = 2)$, then we can do the following:

: (1,0)IMPLICTION $\approx z$ ‘ $z \theta + \eta \Rightarrow z$ ’

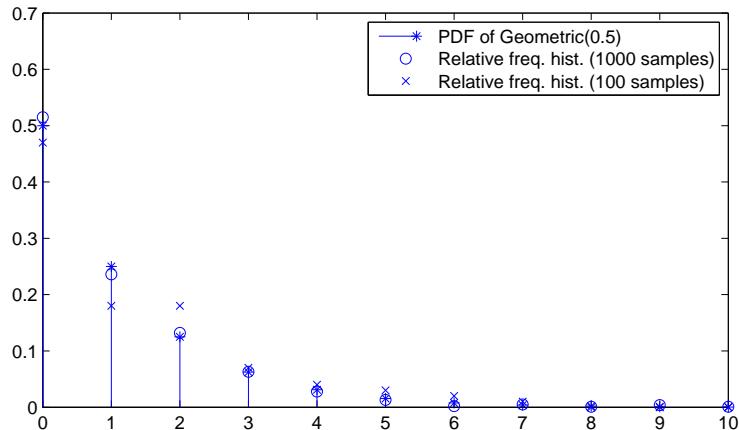
we write as products of samples from $X \sim \text{Normal}(0, 1)$, with some secret parameters θ_1 and θ_2 , then we can use the following relationship between X and $Z \sim \text{Normal}(0, 1)$:

```

ans =
ans = 1.587
>>> random % produce 1 sample from Normal(0,1) RV
>>> random % produce an 2 x 8 array of samples from Normal(0,1)
>>> random % generate 1 sample from Uniform(0,1) RV
>>> random % initialize the seed at 67678 and method as ZIGGURAT -- TYPE HELP random
>>> random % generate 1 sample from Uniform(0,1) RV

```

Figure 3.7: PMF of $X \sim \text{Geometric}(\theta = 0.5)$ and the relative frequency histogram based on 100 and 1000 samples from X according to Simulation 144 and Labwork 145 you will see in the sequel.



Example 46 Suppose we flip a coin 10 times and count the number of heads. Let's consider the probability of getting three heads, say. The probability that the first three flips are heads and the last seven flips are tails, *in order*, is

$$\underbrace{\frac{1}{2} \frac{1}{2} \frac{1}{2}}_{3 \text{ successes}} \underbrace{\frac{1}{2} \frac{1}{2} \cdots \frac{1}{2}}_{7 \text{ failures}}.$$

But there are

$$\binom{10}{3} = \frac{10!}{7!3!} = 120$$

ways of ordering three heads and seven tails, so the probability of getting three heads and seven tails *in any order*, is

$$\mathbf{P}(\text{'3 heads'}) = \binom{10}{3} \left(\frac{1}{2}\right)^3 \left(\frac{1}{2}\right)^7 \approx 0.117$$

We can describe this sort of situation by considering a random variable X which counts the number of successes, as follows:

Model 5 (Binomial(n, θ) RV) Let the RV $X = \sum_{i=1}^n X_i$ be the sum of n independent and identically distributed Bernoulli(θ) RVs, i.e.:

$$X = \sum_{i=1}^n X_i, \quad X_1, X_2, \dots, X_n \stackrel{\text{IID}}{\sim} \text{Bernoulli}(\theta).$$

Given two parameters n and θ , the PMF of the Binomial(n, θ) RV X is:

$$f(x; n, \theta) = \begin{cases} \binom{n}{x} \theta^x (1-\theta)^{n-x} & \text{if } x \in \{0, 1, 2, 3, \dots, n\}, \\ 0 & \text{otherwise} \end{cases} \quad (3.19)$$

Proof: For the continuous case:

$$\mathbf{P}(\text{'accept } y) = \mathbf{P}\left(u \leq \frac{f(y)}{ag(y)}\right) = \int_{-\infty}^{\infty} \left(\int_0^{f(y)/ag(y)} du\right) g(y) dy = \int_{-\infty}^{\infty} \frac{f(y)}{ag(y)} g(y) dy = \frac{1}{a}.$$

And the number of proposals before acceptance is the Geometric($1/a$) RV with expectation $\frac{1}{1/a} = a$.

The closer $ag(x)$ is to $f(x)$, especially in the tails, the closer a will be to 1, and hence the more efficient the rejection method will be.

The rejection method can still be used only if the un-normalised form of f or g (or both) is known. In other words, if we use:

$$f(x) = \frac{\tilde{f}(x)}{\int \tilde{f}(x) dx} \text{ and } g(x) = \frac{\tilde{g}(x)}{\int \tilde{g}(x) dx}$$

we know only $\tilde{f}(x)$ and/or $\tilde{g}(x)$ in closed-form. Suppose the following are satisfied:

- (a) we can generate random variables from g ;
- (b) the support of g contains the support of f , i.e. $\mathbb{Y} \supset \mathbb{X}$;
- (c) a constant $\tilde{a} > 0$ exists, such that:

$$\tilde{f}(x) \leq \tilde{a}\tilde{g}(x), \quad (4.7)$$

for any $x \in \mathbb{X}$, the support of X . Then x can be generated from Algorithm 8.

Algorithm 8 Rejection Sampler (RS) of von Neumann – target shape

1: *input:*

- (1) shape of a target density $\tilde{f}(x) = \left(\int \tilde{f}(x) dx\right) f(x)$,
- (2) a proposal density $g(x)$ satisfying (a), (b) and (c) above.

2: *output:* a sample x from RV X with density f

3: **repeat**

4: Generate $y \sim g$ and $u \sim \text{Uniform}(0, 1)$

5: **until** $u \leq \frac{\tilde{f}(y)}{\tilde{a}\tilde{g}(y)}$

6: **return:** $x \leftarrow y$

Now, the expected number of iterations to get an x is no longer \tilde{a} but rather the integral ratio:

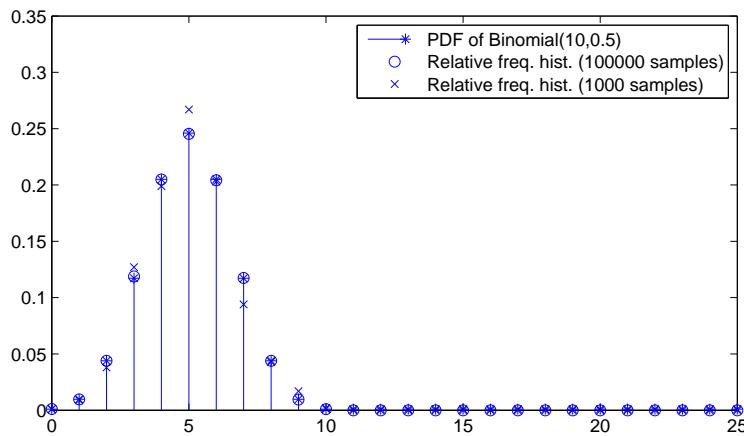
$$\left(\frac{\int_{\mathbb{X}} \tilde{f}(x) dx}{\int_{\mathbb{Y}} \tilde{a}\tilde{g}(y) dy} \right)^{-1}.$$

The **Ziggurat Method** [G. Marsaglia and W. W. Tsang, SIAM Journal of Scientific and Statistical Programming, volume 5, 1984] is a rejection sampler that can efficiently draw samples from the $Z \sim \text{Normal}(0, 1)$ RV. The MATLAB function `randn` uses this method to produce samples from Z .¹

Labwork 155 (Gaussian Sampling with `randn`) We can use MATLAB function `randn` that implements the Ziggurat method to draw samples from an RV $Z \sim \text{Normal}(0, 1)$ as follows:

¹ See http://en.wikipedia.org/wiki/Ziggurat_algorithm for more details.

Figure 3.8: PDF of $X \sim \text{Binomial}(n = 10, \theta = 0.5)$ and the relative frequency histogram based on 100,000 samples from X obtained according to Simulation 148.



Example 48 Compute the probability of obtaining *at least two 6's* in rolling a fair die independently and identically four times.

Solution:

In any given toss let $\theta = P(\{6\}) = 1/6$, $1 - \theta = 5/6$, $n = 4$.

The event *at least two 6's* occurs if we obtain two or three or four 6's. Hence the answer is:

$$\begin{aligned} P(\text{at least two 6's}) &= f\left(2; 4, \frac{1}{6}\right) + f\left(3; 4, \frac{1}{6}\right) + f\left(4; 4, \frac{1}{6}\right) \\ &= \binom{4}{2} \left(\frac{1}{6}\right)^2 \left(\frac{5}{6}\right)^{4-2} + \binom{4}{3} \left(\frac{1}{6}\right)^3 \left(\frac{5}{6}\right)^{4-3} + \binom{4}{4} \left(\frac{1}{6}\right)^4 \left(\frac{5}{6}\right)^{4-4} \\ &= \frac{1}{6^4} (6 \cdot 25 + 4 \cdot 5 + 1) \\ &\approx 0.132 \end{aligned}$$

To make concrete sense of the $\text{Binomial}(n, \theta)$ and other more sophisticated concepts in the sequel, let us take a historical detour into some origins of statistical thinking in 19th century England.

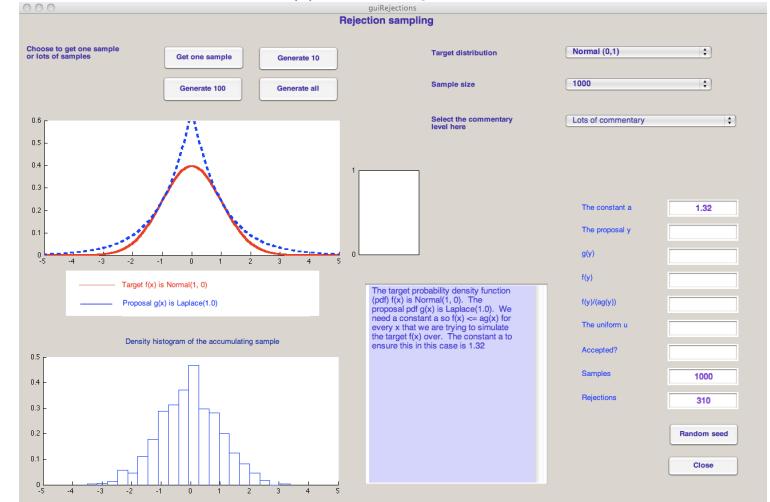
Sir Francis Galton's Quincunx

This section is introduced to provide some forms for a kinesthetic (hands-on) and visual understanding of some elementary statistical distributions and laws. The following words are from Sir Francis Galton, F.R.S., *Natural Inheritance*, pp. 62-65, Macmillan, 1889. In here you will already find the kernels behind the construction of $\text{Binomial}(\theta)$ RV as sum of IID $\text{Bernoulli}(\theta)$ RVs, Weak Law of Large Numbers, Central Limit Theorem, and more. We will mathematically present these concepts

```
>> guiRejections
```

The M-file `guiRejections.m` will bring a graphical user interface (GUI) as shown in Figure 4.11. Try various buttons and see how the output changes with explanations. Try switching the “Target distribution” to “Mywavy4” and generate several rejection samples and see the density histogram of the accumulating samples.

Figure 4.11: Visual Cognitive Tool GUI: Rejection Sampling from $X \sim \text{Normal}(0, 1)$ with PDF f based on proposals from $Y \sim \text{Laplace}(1)$ with PDF g .



Simulation 153 (Rejection Sampling Normal(0, 1) with Laplace(1) proposals) Suppose we wish to generate from $X \sim \text{Normal}(0, 1)$. Consider using the rejection sampler with proposals from $Y \sim \text{Laplace}(1)$ (using inversion sampler of Simulation 137). The support of both RVs is $(-\infty, \infty)$. Next:

$$f(x) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{x^2}{2}\right) \text{ and } g(x) = \frac{1}{2} \exp(-|x|)$$

and therefore:

$$\frac{f(x)}{g(x)} = \sqrt{\frac{2}{\pi}} \exp\left(|x| - \frac{x^2}{2}\right) \leq \sqrt{\frac{2}{\pi}} \exp\left(\frac{1}{2}\right) = a \approx 1.3155 .$$

Hence, we can use the rejection method with:

$$\frac{f(y)}{ag(y)} = \frac{f(y)}{g(y)a} = \sqrt{\frac{2}{\pi}} \exp\left(|y| - \frac{y^2}{2}\right) \frac{1}{\sqrt{\frac{2}{\pi}} \exp\left(\frac{1}{2}\right)} = \exp\left(|y| - \frac{y^2}{2} - \frac{1}{2}\right)$$

Let us implement a rejection sampler as a function in the M-file `RejectionNormalLaplace.m` by reusing the function in `LaplaceInvCDF.m`.

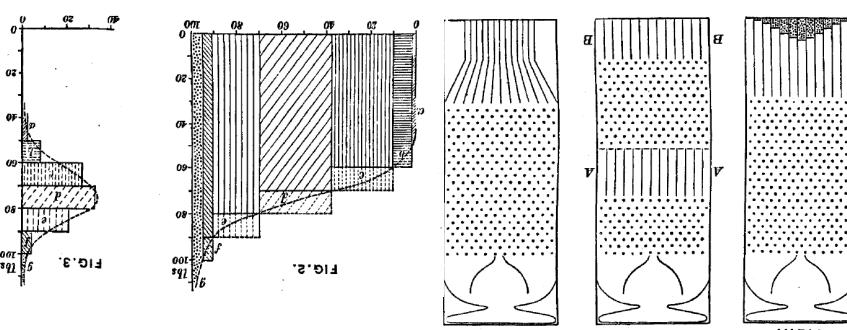


Figure 3-9: Figures from Sir Francis Galton, F.R.S., *Natural Inheritance*, Macmillan, 1889.

Charm of Statistics.—It is difficult to understand why statisticians commonly limit their inquiries to a way of getting precise meanings to Galtor's observations with his Qumcunx. The Qumcunx of Statistics.—It is difficult to understand why statisticians commonly limit their inquiries to a way of getting precise meanings to Galtor's observations with his Qumcunx. The Qumcunx of Statistics.—It is difficult to understand why statisticians commonly limit their inquiries to a way of getting precise meanings to Galtor's observations with his Qumcunx. The Qumcunx of Statistics.—It is difficult to understand why statisticians commonly limit their inquiries to a way of getting precise meanings to Galtor's observations with his Qumcunx.

CHAPTER 3. RANDOM VARIABLES

GL

Labwork 152 (Rejection Sampler Demo) Let us understand the rejection sampler by calling the interactive visual cognitive tool:

Proof: We shall prove the result for the continuous case. For any real number t :

Proposition 60 (Fundamental Theorem of Simulation) The von Neumann rejection sam-
pler of Algorithm 7 produces a sample x from the random variable X with density $f(x)$.

1: *input:* f
 2: *output:* a sample x from RV X with density f
 3: **repeat**
 4: Generate $y \sim g$ and $u \sim \text{Uniform}(0, 1)$
 5: until $u < \frac{f(y)}{f(x)}$
 6: **return:** $x \rightarrow y$

for any $x \in \mathbb{X}$, the support of X . Then x can be generated from Algorithm 7.

(4.6)
$$f(x) \leq ag(x).$$

- (a) we can generate random variables from g ;
- (b) the support of g contains the support of f , i.e. $\mathbb{Y} \subseteq \mathbb{X}$;
- (c) a constant $a > 1$ exists, such that:

Suppose we have another density or mass function g for which the following are true:
 Typically, the target density f is only known up to a constant and therefore the (normalized) density f itself may be unknown and it is difficult to generate samples directly from X .
 draw independent samples from a target RV X with probability density $f(x)$, where $x \in \mathbb{C}^N$.
 Los Alamos Scientific Laboratory, 1987, p. 135-136] is a Monte Carlo method to
 Neglecton sampling [John von Neumann, 1941, in *Statistical Ulam 1909-1984*, a special issue of

4.3.3 von Neumann Rejection Sampler (Rs)

(c) draw 100 samples from the geometric(θ) RV and report the sample mean. Note: the inputs θ_0 and $C(i)$ for the Geometric(θ) RV should be derived and the workings

(b) Set the variable Mytheta=pi/4.

become more nearly identical with the Normal Curve of Frequency, if the rows of pins were much more numerous, the shot smaller, and the compartments narrower; also if a larger quantity of shot was used.

The principle on which the action of the apparatus depends is, that a number of small and independent accidents befall each shot in its career. In rare cases, a long run of luck continues to favour the course of a particular shot towards either outside place, but in the large majority of instances the number of accidents that cause Deviation to the right, balance in a greater or less degree those that cause Deviation to the left. Therefore most of the shot finds its way into the compartments that are situated near to a perpendicular line drawn from the outlet of the funnel, and the Frequency with which shots stray to different distances to the right or left of that line diminishes in a much faster ratio than those distances increase. This illustrates and explains the reason why mediocrity is so common."

We now consider the last of our common discrete random variables for now, the **Poisson** case. A Poisson random variable counts the number of times an event occurs.

We might, for example, ask:

- How many customers visit Cafe Angstrom each day?
- How many sixes are scored in a cricket season? Cricket is a game played in the English-speaking worlds.
- How many bombs hit a city block in south London during World War II?

A Poisson experiment has the following characteristics:

- The average rate of an event occurring is known. This rate is constant.
- The probability that an event will occur during a short continuum is proportional to the size of the continuum.
- Events occur independently.

The number of events occurring in a Poisson experiment is referred to as a **Poisson random variable**.

Model 6 (Poisson(λ) RV) Given a real parameter $\lambda > 0$, the discrete RV X is said to be Poisson(λ) distributed if X has PDF:

$$f(x; \lambda) = \begin{cases} \frac{e^{-\lambda} \lambda^x}{x!} & \text{if } x \in \mathbb{Z}_+ := \{0, 1, 2, \dots\}, \\ 0 & \text{otherwise.} \end{cases} \quad (3.20)$$

Note that the PDF integrates to 1:

$$\sum_{x=0}^{\infty} f(x; \lambda) = \sum_{x=0}^{\infty} \frac{e^{-\lambda} \lambda^x}{x!} = e^{-\lambda} \sum_{x=0}^{\infty} \frac{\lambda^x}{x!} = e^{-\lambda} e^{\lambda} = 1,$$

where we exploit the Taylor series of e^λ to obtain the second-last equality above.

We interpret X as the number of times an event occurs during a specified continuum given that the average value in the continuum is λ .

Algorithm 6 allows us to simulate from any member of the class of non-negative discrete RVs as specified by the probabilities $(\theta_0, \theta_1, \dots)$. When an RV X takes values in another countable set $\mathbb{X} \neq \mathbb{Z}_+$, then we can still use the above algorithm provided we have a one-to-one and onto mapping $D(i) = x : \mathbb{Z}_+ \rightarrow \mathbb{X}$ that allows us to think of $(0, 1, 2, \dots)$ as indices of an array D giving $\mathbb{X} = (D(0), D(1), \dots)$.

Algorithm 6 Inversion Sampler for $GD(\theta_0, \theta_1, \dots)$ RV X

- 1: *input:*
 2. θ_0 and $\{C(i) = \theta_i / \theta_{i-1}\}$ for any $i \in \{1, 2, 3, \dots\}$.
 3. $u \sim \text{Uniform}(0, 1)$
 - 2: *output:* a sample from X
 - 3: *initialise:* $p \leftarrow \theta_0$, $q \leftarrow \theta_0$, $i \leftarrow 0$
 - 4: **while** $u > q$ **do**
 - 5: $i \leftarrow i + 1$, $p \leftarrow p C(i)$, $q \leftarrow q + p$
 - 6: **end while**
 - 7: *return:* $x = i$
-

Simulation 150 (Binomial(n, θ)) To simulate from a Binomial(n, θ) RV X , we can use Algorithm 6 with:

$$\theta_0 = (1 - \theta)^n, \quad C(x+1) = \frac{\theta(n-x)}{(1-\theta)(x+1)}, \quad \text{Mean Efficiency: } O(1+n\theta).$$

Similarly, with the appropriate θ_0 and $C(x+1)$, we can also simulate from the Geometric(θ) and Poisson(λ) RVs.

Labwork 151 This is a challenging exercise for the student who is finding the other Labworks too easy. So those who are novice to MATLAB may skip this Labwork.

1. Implement Algorithm 6 via a function named **MyGenDiscInvSampler** in MATLAB. Hand in the **M-file** named **MyGenDiscInvSampler.m** giving detailed comments explaining your understanding of each step of the code. [Hint: $C(i)$ should be implemented as a function (use function handles via \emptyset) that can be passed as a parameter to the function **MyGenDiscInvSampler**].
2. Show that your code works for drawing samples from a Binomial(n, p) RV by doing the following:
 - (a) Seed the fundamental sampler by your Student ID (if your ID is 11424620 then type `rand('twister', 11424620);`)
 - (b) Draw 100 samples from the Binomial($n = 20, p = 0.5$) RV and report the results in an 2×2 table with column headings x and No. of observations. [Hint: the inputs θ_0 and $C(i)$ for the Binomial(n, p) RV is given above].
3. Show that your code works for drawing samples from a Geometric(p) RV by doing the following:
 - (a) Seed the fundamental sampler by your Student ID.

Example 49 If on the average, 2 cars enter a certain parking lot per minute, what is the probability that during any given minute three cars or fewer will enter the lot?

Think: Why are the assumptions for a Poisson random variable likely to be correct here?

Note: Use calculators, or Excel or Maple, etc. In an exam you may be given needed values from Poisson tables.

Let the random variable X denote the number of cars arriving per minute. Note that the continuum is 1 minute here. Then X can be considered to have a Poisson distribution with $\lambda = 2$ because 2 cars enter on average.

The probability that three cars or fewer enter the lot is:

$$\begin{aligned} P(X \leq 3) &= f(0; 2) + f(1; 2) + f(2; 2) + f(3; 2) \\ &= e^{-2} \left(\frac{2^0}{0!} + \frac{2^1}{1!} + \frac{2^2}{2!} + \frac{2^3}{3!} \right) \\ &= 0.857 \quad (\text{3 sig. fig.}) \end{aligned}$$

Example 50 (Arrivals at a Service Station) The proprietor of a service station finds that, on average, 8 cars arrive per hour on Saturdays. What is the probability that during a randomly chosen 15 minute period on a Saturday:

- (a) No cars arrive?
- (b) At least three cars arrive?

Let the random variable X denote the number of cars arriving in a 15 minute interval. The continuum is 15 minutes here so we need the average number of cars that arrive in a 15 minute period, or $\frac{1}{4}$ of an hour. We know that 8 cars arrive per hour, so X has a Poisson distribution with

$$\lambda = \frac{8}{4} = 2.$$

Here is a call to simulate 10 samples from $\text{Poisson}(\lambda = 10.0)$ and $\text{Poisson}(\lambda = 0.1)$ RVs:

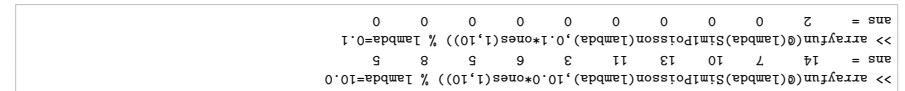
$$f(x; \theta_0, \theta_1, \dots) = \begin{cases} \theta_1, & \text{if } x = 1 \\ \theta_0, & \text{if } x = 0 \\ 0, & \text{if } x \notin \{0, 1, 2, \dots\} \end{cases}$$

is defined as follows:

Model 20 ($GD(\theta_0, \theta_1, \dots)$) We say X is a General Discrete($\theta_0, \theta_1, \dots$) or $GD(\theta_0, \theta_1, \dots)$ RV over the countable discrete state space $\mathbb{Z}_+ = \{0, 1, 2, \dots\}$ with parameters $(\theta_0, \theta_1, \dots)$ if the PMF of X

Simulating from a Poisson(λ) RV is also a special case of simulating from the following more general RV.

Figure 3.18 depicts a comparison of the PDF of Poisson($\lambda = 10$) RV and a relative frequency histogram based on 1000 simulations from it.



Here is a call to simulate 10 samples from $\text{Poisson}(\lambda = 10.0)$ and $\text{Poisson}(\lambda = 0.1)$ RVs:

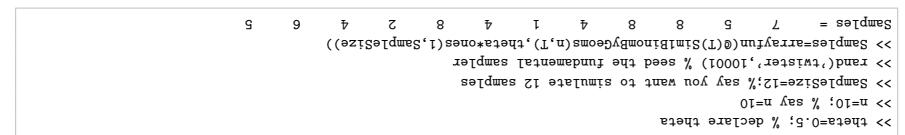
```
function x = SimPoisson(lamdba)
    % Simulate one sample Poisson(lamdba) via Exponentials
    YSum=0; k=0; % initialize
    while (YSum < 1),
        k=k+1;
        x=k-1; % return x
    end
    YSum = YSum + -(1/lamdba) * Log(rand);
```

We implement the above algorithm via the following M-file:

Simulation 149 (Poisson(λ) from IID Exponential(λ) RVs) By this principle, we can simulate from the Poisson(λ) X by Step 1: generating IID Exponential(λ) RVs X_1, X_2, \dots , Step 2: stopping as soon as $\sum_{i=1}^k X_i \geq 1$ and Step 3: setting $x = k - 1$.

Let us simulate from the Poisson(λ) RV of Model 6 as shown in Figure 3.18.

Figure 3.8 depicts a comparison of the PDF of Binomial($n = 10, \theta = 0.5$) RV and a relative frequency histogram based on 100,000 simulations from it.



Here is a call to simulate 12 samples from Binomial($n = 10, \theta = 0.5$) RV:

CHAPTER 4. SIMULATION

Example 51 (Still-born Babies) About 0.01% of babies are stillborn in a certain hospital. We find the probability that of the next 5000 babies born, there will be no more than 1 stillborn baby.

Let the random variable X denote the number of stillborn babies. Then X has a binomial distribution with parameters $n = 5000$ and $\theta = 0.0001$. Since θ is so small and n is large, this binomial distribution may be approximated by a Poisson distribution with parameter

$$\lambda = n\theta = 5000 \times 0.0001 = 0.5.$$

Hence

$$\mathbf{P}(X \leq 1) = \mathbf{P}(X = 0) + \mathbf{P}(X = 1) = f(0; 0.5) + f(1; 0.5) = 0.910 \quad (\text{3 sig. fig.})$$

Exercise 3.4 (Nazi Bombs on London) Feller discusses the probability and statistics of flying bomb hits in an area of southern London during II world war. The area in question was partitioned into $24 \times 24 = 576$ small squares. The total number of hits was 537. There were 229 squares with 0 hits, 211 with 1 hit, 93 with 2 hits, 35 with 3 hits, 7 with 4 hits and 1 with 5 or more hits. Assuming the hits were purely random, use the Poisson approximation to find the probability that a particular square would have exactly k hits. Compute the expected number of squares that would have 0, 1, 2, 3, 4, and 5 or more hits and compare this with the observed results (Snell 9.2.14).

THINKING POISSON

The Poisson distribution has been described as a limiting version of the Binomial. In particular, Exercise 49 thinks of a Poisson distribution as a model for the number of events (cars) that occur in a period of time (1 minute) when in each little chunk of time one car arrives with constant probability, independently of the other time intervals. This leads to the general view of the Poisson distribution as a good model when:

You count the number of events in a continuum when the events occur at constant rate, one at a time and independent of each other.

DISCRETE RANDOM VARIABLE SUMMARY

Probability mass function

$$f(x) = \mathbf{P}(X = x_i)$$

Distribution function

$$F(x) = \sum_{x_i \leq x} f(x_i)$$

```
>> x=0:1:8
x =      0      1      2      3      4      5      6      7      8
>> BinomialPdf(x,8,0.5)
ans =    0.0039    0.0312    0.1094    0.2188    0.2734    0.2188    0.1094    0.0312    0.0039
```

Simulation 147 (Binomial(n, θ) as $\sum_{i=1}^n$ Bernoulli(θ)) Since the Binomial(n, θ) RV X is the sum of n IID Bernoulli(θ) RVs we can also simulate from X by first simulating n IID Bernoulli(θ) RVs and then adding them up as follows:

```
>> rand('twister',17678);
>> theta=0.5; % give some desired theta value, say 0.5
>> n=5; % give the parameter n for Binomial(n,theta) RV X, say n=5
>> xis=floor(rand(1,n)*theta) % produce n IID samples from Bernoulli(theta=0.5) RVs X1,X2,...Xn
xis =
     1     1     0     0
>> x=sum(xis) % sum up the xis to get a sample from Binomial(n=5,theta=0.5) RV X
x =
     2
```

It is straightforward to produce more than one sample from X by exploiting the column-wise summing property of MATLAB's `sum` function when applied to a two-dimensional array:

```
>> rand('twister',17);
>> theta=0.25; % give some desired theta value, say 0.25 this time
>> n=3; % give the parameter n for Binomial(n,theta) RV X, say n=3 this time
>> xis10 = floor(rand(n,10)*theta) % produce an n by 10 array of IID samples from Bernoulli(theta=0.25) RVs
xis10 =
     0     0     0     0     1     0     0     0     0     0
     0     1     0     1     1     0     0     0     0     0
     0     0     0     0     0     0     0     1     0     0
>> x=sum(xis10) % sum up the array column-wise to get 10 samples from Binomial(n=3,theta=0.25) RV X
x =
     0     1     0     1     2     0     0     1     0     0
```

In Simulation 147, the number of IID Bernoulli(θ) RVs needed to simulate one sample from the Binomial(n, θ) RV is exactly n . Thus, as n increases, the amount of time needed to simulate from Binomial(n, θ) is $O(n)$, i.e. linear in n . We can simulate more efficiently by exploiting a simple relationship between the Geometric(θ) RV and the Binomial(n, θ) RV.

The Binomial(n, θ) RV X is related to the IID Geometric(θ) RV Y_1, Y_2, \dots : X is the number of successful Bernoulli(θ) outcomes (outcome is 1) that occur in a total of n Bernoulli(θ) trials, with the number of trials between consecutive successes distributed according to IID Geometric(θ) RV.

Simulation 148 (Binomial(θ) from IID Geometric(θ) RVs) By this principle, we can simulate from the Binomial(θ) X by Step 1: generating IID Geometric(θ) RVs Y_1, Y_2, \dots , Step 2: stopping as soon as $\sum_{i=1}^k (Y_i + 1) > n$ and Step 3: setting $x \leftarrow k - 1$.

We implement the above algorithm via the following M-file:

```
function x = Sim1BinomByGeoms(n,theta)
% Simulate one sample from Binomial(n,theta) via Geometric(theta) RVs
YSum=0; k=0; % initialise
while (YSum <= n),
    TrialsToSuccess=floor(log(rand)/log (1-theta)) + 1; % sample from Geometric(theta)+1
    YSum = YSum + TrialsToSuccess; % total number of trials
    k=k+1; % number of Bernoulli successes
end
x=k-1; % return x
```

Random Variable	Possible Values	Probabilities	Modelled situations
Discrete uniform	$\{x_1, x_2, \dots, x_k\}$	$P(X = x_i) = \frac{1}{k}$	Situations with k equally likely values. Parameter: $\theta = k$.
Bernoulli(θ)	{0, 1}	$P(X = 0) = 1 - \theta$ $P(X = 1) = \theta$	Situations with only 2 outcomes, coded 1 for success and 0 for failure. Parameter: $\theta = P(\text{success}) \in (0, 1)$.
Geometric(θ)	{1, 2, 3, ...}	$P(X = x) = (1 - \theta)^{x-1} \theta$	Situations where you count the number of trials until the first success in a sequence of independent trials with a constant probability of success.
Binomial(n, θ)	{0, 1, 2, ..., n}	$P(X = x) = \binom{n}{x} \theta^x (1-\theta)^{n-x}$	Situations where you count the number of successes in n trials where each trial is independent and there is a constant probability of success.
Poisson(λ)	{0, 1, 2, ...}	$P(X = x) = \frac{\lambda^x e^{-\lambda}}{x!}$	Situations where you count the number of events in a continuum where the events occur one at a time and are independent of one another. Parameter: $\lambda = \text{rate } \in (0, \infty)$.

Ex. 3.5 — One number in the following table for the probability function of a random variable X is incorrect. Which is it, and what should the correct value be?

Ex. 3.6 — Let X be the number of years before a particular type of machine will need replacement. Assume that X has the probability function $f(1) = 0.1, f(2) = 0.2, f(3) = 0.2, f(4) = 0.2$, $f(5) = 0.3$.

Ex. 3.7 — Of 200 adults, 176 own one TV set, 22 own two TV sets, and 2 own three TV sets. A person is chosen at random. What is the probability mass function of X , the number of TV sets owned by that person?

Ex. 3.8 — Suppose a discrete random variable X has probability function give by

$$P(X = x) \begin{array}{|c|ccccccccc|} \hline x & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 & 13 \\ \hline 0.07 & 0.01 & 0.09 & 0.01 & 0.16 & 0.25 & 0.20 & 0.03 & 0.02 & 0.11 & 0.05 & \\ \hline \end{array}$$

It is a good idea to make a relative frequency histogram of simulated Geometric(θ) RV that is similar to the PDF of the discrete RV we are simulating from. We use the following script to create

Figure 3.7: PMF versus relative frequency histogram of simulated Geometric(θ) RV

Let us simulate from the Binomial(n, θ) RV of Model 5.

With the following M-L file:

```
function BC = BinomialCoefficient(n,x)
    % returns the binomial coefficient of n choose x
    % i.e. the combination of n objects taken x at a time
    % x and n are scalar integers and 0 <= x <= n
    % returns the binomial coefficient of n choose x
    % i.e. the combination of n objects taken x at a time
    % x and n are scalar integers and 0 <= x <= n
    % NumeratorPostCancel = prod(n:-1:max([NumInsx,X]+1));
    % DenominatorPostCancel = prod(2:min([NumInsx,X]));;
    BC = NumeratorPostCancel/DenominatorPostCancel;
```

Labwork 146 (Binomial coefficient) The MATLAB function BinomialCoefficient can be used to compute:

```
function BC = BinomialCoefficient(n,x)
    % returns the binomial coefficient of n choose x
    % i.e. the combination of n objects taken x at a time
    % x and n are scalar integers and 0 <= x <= n
    % NumeratorPostCancel = prod(n:-1:max([NumInsx,X]+1));
    % DenominatorPostCancel = prod(2:min([NumInsx,X]));;
    BC = NumeratorPostCancel/DenominatorPostCancel;
```

and call BinomialCoefficient in the function BinomialPDF to compute the PDF $f(x; n, \theta)$ of the Binomial(n, θ) RV X as follows:

```
function fx = BinomialPDF(x,n,theta)
    % Binomial probability mass function. Needs BinomialCoefficient(n,x)
    % f is the prob mass function for the Binomial(x;n,theta)
    % and x is array of samples.
    fx = zeros(size(x));
    for m=0:n
        fx(m+1)=BinomialCoefficient(n,m,theta);
    end
    fx = fx.* (theta.^x .* (1-theta).^(n-x));
```

- (a) Construct a row of cumulative probabilities for this table, that is, find the distribution function of X .
 (b) Find the following probabilities.

$$\begin{array}{lll} \text{(i)} \mathbf{P}(X \leq 5) & \text{(iii)} \mathbf{P}(X > 9) & \text{(v)} \mathbf{P}(4 < X \leq 9) \\ \text{(ii)} \mathbf{P}(X < 12) & \text{(iv)} \mathbf{P}(X \geq 9) & \text{(vi)} \mathbf{P}(4 < X < 11) \end{array}$$

Ex. 3.9 — A box contains 4 right-handed and 6 left-handed screws. Two screws are drawn at random without replacement. Let X be the number of left-handed screws drawn. Find the probability mass function for X , and then calculate the following probabilities:

1. $\mathbf{P}(X \leq 1)$
2. $\mathbf{P}(X \geq 1)$
3. $\mathbf{P}(X > 1)$

Ex. 3.10 — Suppose that a random variable X has geometric probability mass function,

$$f(x) = \frac{k}{2^x} \quad (x = 0, 1, 2, \dots).$$

1. Find the value of k .
2. What is $\mathbf{P}(X \geq 4)$?

Ex. 3.11 — Four fair coins are tossed simultaneously. If we count the number of heads that appear then we have a binomial random variable, $X = \text{the number of heads}$.

1. Find the probability mass function of X .
2. Compute the probabilities of obtaining no heads, precisely 1 head, at least 1 head, not more than 3 heads.

Ex. 3.12 — The distribution of blood types in a certain population is as follows:

Blood type	Type O	Type A	Type B	Type AB
Proportion	0.45	0.40	0.10	0.05

A random sample of 15 blood donors is observed from this population. Find the probabilities of the following events.

1. Only one type AB donor is included.
2. At least three of the donors are type B .
3. More than ten of the donors are either type O or type A .
4. Fewer than five of the donors are not type A .

Ex. 3.13 — If the probability of hitting a target in a single shot is 10% and 10 shots are fired independently, what is the probability that the target will be hit at least once?

Ex. 3.14 — Suppose that a certain type of magnetic tape contains, on the average, 2 defects per 100 meters. What is the probability that a roll of tape 300 meters long will contain no defects?

Ex. 3.15 — In 1910, E. Rutherford and H. Geiger showed experimentally that the number of alpha particles emitted per second in a radioactive process is a random variable X having a Poisson distribution. If the average number of particles emitted per second is 0.5, what is the probability of observing two or more particles during any given second?

```
>> rand('twister',1777); % initialise the fundamental sampler
>> f2=rand(1,10); % create an arbitrary array
>> f2=f2/sum(f2); % normalize to make a probability mass function
>> disp(f2); % display the weights of our 10-faced die
    0.0073   0.0188   0.1515   0.1311   0.1760   0.1121   ...
    0.1718   0.1213   0.0377   0.0723
>> disp(sum(f2)); % the weights sum to 1
    1.0000
>> disp(arrayfun(@(u)(SimdeMoivreOnce(u,f2)),rand(5,5))) % the samples from f2 are
    4     3     4     7     3
    6     7     4     5     3
    5     8     7    10     6
    2     3     5     7     7
    6     5     9     5     7
```

Note that the principal work here is the sequential search, in which the mean number of comparisons until success is:

$$1\theta_1 + 2\theta_2 + 3\theta_3 + \dots + k\theta_k = \sum_{i=1}^k i\theta_i$$

For the de Moivre($1/k, 1/k, \dots, 1/k$) RV, the right-hand side of the above expression is:

$$\sum_{i=1}^k i \frac{1}{k} = \frac{1}{k} \sum_{i=1}^k i = \frac{1}{k} \frac{k(k+1)}{2} = \frac{k+1}{2},$$

indicating that the average-case efficiency is linear in k . This linear dependence on k is denoted by $O(k)$. In other words, as the number of faces k increases, one has to work linearly harder to get samples from de Moivre($1/k, 1/k, \dots, 1/k$) RV using Algorithm 5. Using the simpler Algorithm 4, which exploits the fact that all values of θ_i are equal, we generated samples in constant time, which is denoted by $O(1)$.

Simulation 144 (Geometric(θ)) We can simulate a sample x from a Geometric(θ) RV X using the following simple algorithm:

$$x \leftarrow \lfloor \log(u)/\log(1-\theta) \rfloor, \quad \text{where, } u \sim \text{Uniform}(0, 1).$$

To verify that the above procedure is valid, note that:

$$\begin{aligned} \lfloor \log(U)/\log(1-\theta) \rfloor = x &\iff x \leq \log(U)/\log(1-\theta) < x+1 \\ &\iff x \leq \log_{1-\theta}(U) < x+1 \\ &\iff (1-\theta)^x \geq U > (1-\theta)^{x+1} \end{aligned}$$

The inequalities are reversed since the base being exponentiated is $1-\theta \leq 1$. The uniform event $(1-\theta)^x \geq U > (1-\theta)^{x+1}$ happens with the desired probability:

$$(1-\theta)^x - (1-\theta)^{x+1} = (1-\theta)^x(1-(1-\theta)) = \theta(1-\theta)^x =: f(x; \theta), \quad X \sim \text{Geometric}(\theta).$$

We implement the sampler to generate samples from Geometric(θ) RV with $\theta = 0.5$, for instance:

```
>> theta=0.5; u=rand(); % choose some theta and uniform(0,1) variate
>> % Simulate from a Geometric(theta) RV
>> floor(log(u) / log(1-theta))
ans =
    0
>> floor(log(rand(1,10)) / log(1-0.5)) % theta=0.5, 10 samples
ans =
    0     0     1     0     2     1     0     0     0     0
```

Ex. 3.16 — The number of lacunae (surface pits) on specimens of steel, polished and examined in a metallurgical laboratory, is thought to have a Poisson distribution.

1. Write down the formula for the probability that a specimen has x defects, explaining the meanings of the symbols you use.

2. Simplify the formula in the case $x = 0$.

3. In a large homogeneous collection of specimens, 10% have one or more lacunae. Find (approximately) the percentage having exactly two.

4. Why might the Poisson distribution not apply in this situation?

[HINT: Recall the *marginalised entropy* in THINKING POISSON and what the continuum on which the number of events occur is for the problem, and what could possibly go wrong in your imagination of the manufacturing process of the steel specimens (normally you need to melt and manipulate iron with other elements and cast them in moulds and this needs energy and raw materials of possibly varying quality and the machines used in the process could break down, etc), to violate the Poisson assumption about the occurrence of pits on the surface of the specimens].

If X is a measurement of a continuous quantity, such as,

3.4 Continuous Random Variables

- the volume of water (in cubic metres) that fell on the southern Alps of the South Island of New Zealand throughout last year.
- the volume of water (in cubic metres) that fell on the next student leaves the lecture room. This is an example of a continuous random variable that takes one of (uncountably) infinitly many values. When a student leaves, X will take on the value x and this x could be 2.1 minutes, or 2.1000000001 minutes, or 2.99999999 minutes, etc., depending on the measurement precision of the time-measuring clock rather than the discrete approach of trying to compute $P(X = x)$.
- the volume of rain that fell on the roof of this building over the past 365 days in litres.
- the distance you transported yourself to lectures today in metres.
- the maximum diameter in millimetres of a venus shell I picked up at New Brighton beach, a location in Lake Rotoiti in Hokianga, as it traces through Gata alVikaravenu, the longest river of Sweden before discharging in a delta into Vänern at Karlstad.
- the vertical position (in metres above sea-level) since the release of a pollen grain in a river of Sweden before discharging in a delta into Vänern at Karlstad.
- the volume of water (in cubic metres) that fell on the next student leaves the lecture room. This is an example of a continuous random variable that takes one of (uncountably) infinitly many values. When a student leaves, X will take on the value x and this x could be 2.1 minutes, or 2.1000000001 minutes, or 2.99999999 minutes, etc., depending on the measurement precision of the time-measuring clock rather than the discrete approach of trying to compute $P(X = x)$.

then X is a continuous random variable. Continuous random variables are based on measurements in a continuous scale of a given precision as opposed to discrete random variables that are based on counting.

The characteristics of continuous random variables are:

Clearly, Algorithm 5 may be used to sample any de Moivre($\theta_1, \dots, \theta_k$) RV X . We demonstrate this by producing five samples from a randomly generated PMF f_Z .

```

Algorithm 5 Inversion Sampler for de Moivre( $\theta_1, \theta_2, \dots, \theta_k$ ) RV  $X$ 

Input:  $f_Z$ : discrete distribution
Output:  $x$ : realization of  $X$ 

1:  $i \leftarrow 1$ 
2:  $U \sim \text{Uniform}(0, 1)$ 
3: initialize:  $F \rightarrow \theta_i$ ,  $i \rightarrow 1$ 
4: while  $U < F$  do
5:    $i \rightarrow i + 1$ 
6:    $F \rightarrow F + \theta_i$ 
7: end while
8: return:  $x \rightarrow i$ 

Algorithm 5 Inversion Sampler for de Moivre( $\theta_1, \theta_2, \dots, \theta_k$ ) RV  $X$ 

Input:  $f_Z$ : discrete distribution
Output:  $x$ : realization of  $X$ 

1:  $i \leftarrow 1$ 
2:  $U \sim \text{Uniform}(0, 1)$ 
3: initialize:  $F \rightarrow \theta_i$ ,  $i \rightarrow 1$ 
4: while  $U < F$  do
5:    $i \rightarrow i + 1$ 
6:    $F \rightarrow F + \theta_i$ 
7: end while
8: return:  $x \rightarrow i$ 

```

Let us use the function `deMoivreEqn` to draw five samples from a fair seven-faced dice.

```

The M-File implementing Algorithm 5 is:

Algorithm 5 Inversion Sampler for de Moivre( $\theta_1, \dots, \theta_k$ ) RV  $X$ 

Input:  $f_Z$ : discrete distribution
Output:  $x$ : realization of  $X$ 

1:  $i \leftarrow 1$ 
2:  $U \sim \text{Uniform}(0, 1)$ 
3: initialize:  $F \rightarrow \theta_i$ ,  $i \rightarrow 1$ 
4: while  $U < F$  do
5:    $i \rightarrow i + 1$ 
6:    $F \rightarrow F + \theta_i$ 
7: end while
8: return:  $x \rightarrow i$ 

Algorithm 5 Inversion Sampler for de Moivre( $\theta_1, \dots, \theta_k$ ) RV  $X$ 

Input:  $f_Z$ : discrete distribution
Output:  $x$ : realization of  $X$ 

1:  $i \leftarrow 1$ 
2:  $U \sim \text{Uniform}(0, 1)$ 
3: initialize:  $F \rightarrow \theta_i$ ,  $i \rightarrow 1$ 
4: while  $U < F$  do
5:    $i \rightarrow i + 1$ 
6:    $F \rightarrow F + \theta_i$ 
7: end while
8: return:  $x \rightarrow i$ 

```

Let us use the function `deMoivreEqn` to draw five samples from a fair seven-faced dice.

```

Algorithm 5 Inversion Sampler for de Moivre( $\theta_1, \dots, \theta_k$ ) RV  $X$ 

Input:  $f_Z$ : discrete distribution
Output:  $x$ : realization of  $X$ 

1:  $i \leftarrow 1$ 
2:  $U \sim \text{Uniform}(0, 1)$ 
3: initialize:  $F \rightarrow \theta_i$ ,  $i \rightarrow 1$ 
4: while  $U < F$  do
5:    $i \rightarrow i + 1$ 
6:    $F \rightarrow F + \theta_i$ 
7: end while
8: return:  $x \rightarrow i$ 

Algorithm 5 Inversion Sampler for de Moivre( $\theta_1, \dots, \theta_k$ ) RV  $X$ 

Input:  $f_Z$ : discrete distribution
Output:  $x$ : realization of  $X$ 

1:  $i \leftarrow 1$ 
2:  $U \sim \text{Uniform}(0, 1)$ 
3: initialize:  $F \rightarrow \theta_i$ ,  $i \rightarrow 1$ 
4: while  $U < F$  do
5:    $i \rightarrow i + 1$ 
6:    $F \rightarrow F + \theta_i$ 
7: end while
8: return:  $x \rightarrow i$ 

```

Algorithm 5 Inversion Sampler for de Moivre($\theta_1, \dots, \theta_k$) RV X

RV X when $(\theta_1, \theta_2, \dots, \theta_k)$ are specified as an input vector via the following algorithm.

Simulations 143 (de Moivre($\theta_1, \theta_2, \dots, \theta_k$)) We can generate samples from a de Moivre($\theta_1, \theta_2, \dots, \theta_k$) RV X which are specified as an input vector via the following algorithm.

- The outcomes are measured, not counted.
- Geometrically, the probability of an outcome is equal to an area under a mathematical curve.
- Each individual value has zero probability of occurring. So we find the probability that the value is between two endpoints of an interval, or a set of intervals, including half-lines in \mathbb{R} .

Definition 25 (probability density function (PDF)) A RV X with distribution function (DF) given by F is said to be **continuous** if there exists a piecewise-continuous function f , called the **probability density function (PDF)** of X , such that

$$F(x) = \mathbf{P}(X \leq x) = \int_{-\infty}^x f(v) dv \quad (3.21)$$

where $f : \mathbb{R} \rightarrow \mathbb{R}$ is a non-negative function, i.e., $f(x) \geq 0$. We write v because x is needed as the upper limit of the integral. Piecewise-continuity of f means f is continuous, perhaps possibly at the x -values where f is discontinuous between the continuous pieces (see <https://en.wikipedia.org/wiki/Piecewise>).

The following hold for a continuous RV X with PDF f :

1. For any $x \in \mathbb{R}$, $\mathbf{P}(X = x) = \mathbf{P}(X \in [x, x]) = \int_x^x f(v)dv = 0$.
2. By the fundamental theorem of calculus:

$$f(x) = \frac{d}{dx} F(x) =: F'(x), \quad (3.22)$$

for every x at which $f(x)$ is continuous.

3. Consequentially, for any $a, b \in \mathbb{R}$ with $a < b$,

$$\mathbf{P}(a < X < b) = \mathbf{P}(a < X \leq b) = \mathbf{P}(a \leq X \leq b) = \mathbf{P}(a \leq X < b) \quad (3.23)$$

$$= F(b) - F(a) = \int_a^b f(v)dv. \quad (3.24)$$

4. And $P(\Omega) = 1$ implies that:

$$\int_{-\infty}^{\infty} f(x) dx = \mathbf{P}(-\infty < X < \infty) = 1.$$

The next set of examples illustrate notation and typical applications of the formulae above.

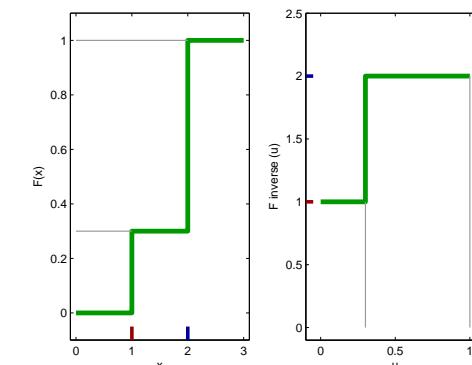
Example 53 Consider the continuous random variable, X , whose probability density function is:

$$f(x) = \begin{cases} 3x^2 & 0 < x < 1 \\ 0 & \text{otherwise} \end{cases}$$

- (a) Find the distribution function, $F(x)$.
- (b) Find $P(\frac{1}{3} \leq X \leq \frac{2}{3})$.

Solution

Figure 4.10: The DF $F(x; 0.3, 0.7)$ of the de Moivre(0.3, 0.7) RV and its inverse $F^{-1}(u; 0.3, 0.7)$.



Algorithm 4 Inversion Sampler for de Moivre($1/k, 1/k, \dots, 1/k$) RV

- 1: *input:*
1. k in de Moivre($1/k, 1/k, \dots, 1/k$) RV X
 2. $u \sim \text{Uniform}(0, 1)$
- 2: *output:* a sample from X
- 3: *return:* $x \leftarrow \lceil ku \rceil$
-

```
function x = SimdeMoivreEqui(u,k);
% return samples from de Moivre(1/k,1/k,...,1/k) RV X
% Call Syntax: x = SimdeMoivreEqui(u,k);
% Input      : u = array of uniform random numbers eg. rand
%               k = number of equi-probabble outcomes of X
% Output     : x = samples from X
x = ceil(k * u); % ceil(y) is the smallest integer larger than y
%x = floor(k * u); if outcomes are in {0,1,...,k-1}
```

Let us use the function `SimdeMoivreEqui` to draw five samples from a fair seven-faced cylindrical dice.

```
>> k=7; % number of faces of the fair dice
>> n=5; % number of trials
>> rand('twister',78657); % initialise the fundamental sampler
>> u=rand(1,n); % draw n samples from Uniform(0,1)
>> % inverse transform samples from Uniform(0,1) to samples
>> % from de Moivre(1/7,1/7,1/7,1/7,1/7,1/7,1/7)
>> outcomes=SimdeMoivreEqui(u,k); % save the outcomes in an array
>> disp(outcomes);
6 5 5 5 2
```

Now, let us consider the more general problem of implementing a sampler for an arbitrary but specified de Moivre($\theta_1, \theta_2, \dots, \theta_k$) RV. That is, the values of θ_i need not be equal to $1/k$.

$$f(x) = \begin{cases} 0 & x < 0 \\ \cos x & 0 < x < \frac{\pi}{2} \\ 0 & x > \frac{\pi}{2} \end{cases}$$

(a) The probability density function, $f(x)$, is given by

Solution

$$(b) \text{ Find } P(X < \frac{\pi}{4})$$

(a) Find the probability density function, $f(x)$.

$$F(x) = \begin{cases} 0 & x \leq 0 \\ \sin(x) & 0 < x < \frac{\pi}{2} \\ 1 & x \geq \frac{\pi}{2} \end{cases}$$

Example 54 Consider the continuous random variable, X , whose distribution function is:

$$\begin{aligned} p &= F\left(\frac{3}{2}\right) - F\left(\frac{1}{2}\right) \\ &= \left(\frac{2}{3}\right)^3 - \left(\frac{1}{2}\right)^3 \\ &= \frac{7}{8} \end{aligned}$$

(b)

$$F(x) = \begin{cases} 0 & x < 0 \\ x^3 & 0 < x < 1 \\ 1 & x \geq 1 \end{cases}$$

Hence

$$F(x) = \int_x^1 3a^2 da + \int_0^0 0 da = 1 - x^3$$

If $x \geq 1$, then

$$\begin{aligned} F(x) &= \int_x^1 3a^2 da + 0 \\ &= x^3 \end{aligned}$$

If $0 < x < 1$, then

$$F(x) = \int_x^0 0 da = 0$$

(a) First note that if $x \leq 0$, then

The M-File implementing Algorithm 4 is:

Algorithm 4 produces samples from the de Moivre($1/k, 1/k, \dots, 1/k$) RV X .
 RV X with a discrete uniform distribution over $[k] = \{1, 2, \dots, k\}$ can be efficiently sampled using the ceiling function. Recall that $[y]$ is the smallest integer larger than or equal to y , e.g. $[13.1] = 14$.
 First we simulate from an equal-probable special case of the de Moivre($\theta_1, \theta_2, \dots, \theta_k$) RV, with $\theta_1 = \theta_2 = \dots = \theta_k = 1/k$.

When $k = 2$ in the de Moivre(θ_1, θ_2) model, we have an RV that is similar to the Bernoulli($p = \theta_1$) RV. The DF F and its inverse $F_{[-1]}$ for a specific $\theta_1 = 0.3$ are depicted in Figure 4.10.

$$F_{[-1]}(u; \theta_1, \theta_2, \dots, \theta_k) = \begin{cases} 1 & \text{if } \theta_1 + \theta_2 + \dots + \theta_{k-1} \leq u < 1 \\ \vdots & \vdots \\ 3 & \text{if } \theta_1 + \theta_2 > u > \theta_1 + \theta_2 + \theta_3 \\ 2 & \text{if } \theta_1 \leq u > \theta_1 + \theta_2 \\ 1 & \text{if } 0 \leq u < \theta_1 \end{cases} \quad (4.5)$$

Given by:

$$F_{[-1]} : [0, 1] \rightarrow [k] := \{1, 2, \dots, k\},$$

Next we simulate from the de Moivre($\theta_1, \theta_2, \dots, \theta_k$) RV X of Model 14 via its inverse DF

```
ans = >> arrayfun(@(u)(simpotimass(u,17)),zeros(2,8))
ans = 17 17 17 17 17 17 17 17
```

Note that it is not necessary to have input IID samples from Uniform(0, 1) RV rand in order to draw samples from the Point Mass(θ) RV. For instance, an input matrix of zeros can do the job:

```
ans = >> arrayfun(@(u)(simpotimass(u,17)),rand(2,10))
ans = 17 17 17 17 17 17 17 17 17 17
```

Here is call to the function.

```
function x = simpotimass(u,theta)
    % Returns one sample from the Point Mass(theta) RV X
    % Input: x = simpotimass(u,theta);
    % Call Syntax: x = simpotimass(u,theta);
    % Output: : u = one uniform random number e.g. rand();
    % Output : x = sample from X
    % Input : theta = a real number (scalar)
    % Note: % we can use arrayfun to apply simpotimass to any array of uniform(0,1) samples
```

this RV produces the same realization θ we can implement it via the following M-File:

(b)

$$P\left(X > \frac{\pi}{4}\right) = 1 - P\left(X \leq \frac{\pi}{4}\right) = 1 - F\left(\frac{\pi}{4}\right) = 1 - \sin\left(\frac{\pi}{4}\right) = 0.293 \text{ (3 sig. fig.)}$$

* You may stop at $1 - \sin\left(\frac{\pi}{4}\right)$ for full credit in the exam.

Note: $f(x)$ is not defined at $x = 0$ as $F(x)$ is not differentiable at $x = 0$. There is a “kink” in the distribution function at $x = 0$ causing this problem. It is standard to define $f(0) = 0$ in such situations, as $f(x) = 0$ for $x < 0$. This choice is arbitrary but it simplifies things and makes no difference to the calculated probability.

Now that we have warmed-up with two examples of continuous RVs, let us define the most elementary continuous RV next.

3.4.1 An Elementary Continuous Random Variable

An elementary and fundamental example of a continuous RV is the Uniform(0, 1) RV of Model 7. It forms the foundation for all non-uniform random variate generation and simulation as we will see in Chapter 4. In fact, it is appropriate to call this the fundamental model since every other probability model can be obtained from this one!

Model 7 (The Fundamental Model) The probability density function (PDF) of the fundamental model or the Uniform(0, 1) RV is

$$f(x) = \mathbb{1}_{[0,1]}(x) = \begin{cases} 1 & \text{if } 0 \leq x \leq 1, \\ 0 & \text{otherwise} \end{cases} \quad (3.25)$$

and its distribution function (DF) or cumulative distribution function (CDF) is:

$$F(x) := \int_{-\infty}^x f(y) dy = \begin{cases} 0 & \text{if } x < 0, \\ x & \text{if } 0 \leq x \leq 1, \\ 1 & \text{if } x > 1 \end{cases} \quad (3.26)$$

Note that the DF is the identity map in [0, 1]. The PDF and DF are depicted in Figure 3.11.

Let us draw the PDF and DF for Uniform(0, 1) RV next by hand.

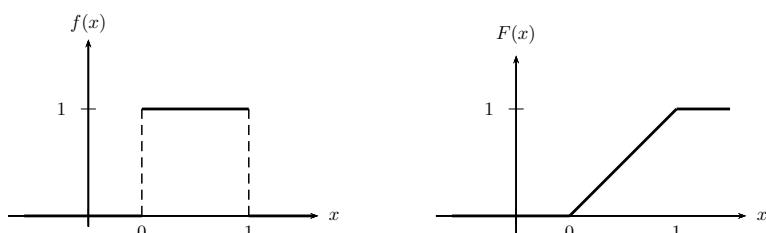


Figure 3.10: $f(x)$ and $F(x)$ of the Uniform(0, 1) random variable X .

Simulation 139 (Cauchy) We can draw n IID samples from the Cauchy RV X by transforming n IID samples from Uniform(0, 1) RV U using the inverse DF as follows:

```
>> rand('twister',2435567); % initialise the fundamental sampler
>> u=rand(1,5); % draw 5 IID samples from Uniform(0,1) RV
>> disp(u);
0.7176 0.6655 0.9405 0.9198 0.2598
>> x=tan(pi * u); % draw 5 samples from Standard cauchy RV using inverse CDF
>> disp(x); % display the samples in x
-1.2272 -1.7470 -0.1892 -0.2575 1.0634
```

4.3.2 Inversion Sampler for Discrete Random Variables

Next, consider the problem of **sampling from a random variable X with a discontinuous or discrete DF** using the inversion sampler. We need to define the inverse more carefully here.

Proposition 59 (Inversion sampler with compact support) Let the support of the RV X be over some real interval $[a, b]$ and let its inverse DF be defined as follows:

$$F^{[-1]}(u) := \inf\{x \in [a, b] : F(x) \geq u, 0 \leq u \leq 1\}.$$

If $U \sim \text{Uniform}(0, 1)$ then $F^{[-1]}(U)$ has the DF F , i.e. $F^{[-1]}(U) \sim F \sim X$.

Proof: The proof is a consequence of the following equalities:

$$\mathbf{P}(F^{[-1]}(U) \leq x) = \mathbf{P}(U \leq F(x)) = F(x) := \mathbf{P}(X \leq x)$$

Simulation 140 (Bernoulli(θ)) Consider the problem of simulating from a Bernoulli(θ) RV based on an input from a Uniform(0, 1) RV. Recall that $\lfloor x \rfloor$ (called the ‘floor of x ’) is the largest integer that is smaller than or equal to x , e.g. $\lfloor 3.8 \rfloor = 3$. Using the floor function, we can simulate a Bernoulli(θ) RV X as follows:

```
>> theta = 0.3; % set theta = Prob(X=1)
% return x -- floor(y) is the largest integer less than or equal to y
>> x = floor(rand + theta); % rand is the Fundamental Sampler
>> disp(x) % display the outcome of the simulation
0
>> n=10; % set the number of IID Bernoulli(theta=0.3) trials you want to simulate
>> x = floor(rand(1,n)+theta); % vectorize the operation
>> disp(x) % display the outcomes of the simulation
0 0 1 0 0 0 0 0 1 1
```

Again, it is straightforward to do replicate experiments, e.g. to demonstrate the Central Limit Theorem for a sequence of n IID Bernoulli(θ) trials.

```
>> % a demonstration of Central Limit Theorem --
>> % the sample mean of a sequence of n IID Bernoulli(theta) RVs is Gaussian(theta,theta*(1-theta)/n)
>> theta=0.5; Reps=10000; n=10; hist(sum(floor(rand(n,Reps)+theta))/n)
>> theta=0.5; Reps=10000; n=100; hist(sum(floor(rand(n,Reps)+theta))/n,20)
>> theta=0.5; Reps=10000; n=1000; hist(sum(floor(rand(n,Reps)+theta))/n,30)
```

Recall the Point Mass(θ) RV. Formally, we can simulate from it trivially as follows.

$$F(x) = \begin{cases} 0 & \text{otherwise} \\ 1 - e^{-x} & \text{if } x \geq 0 \end{cases}$$

Therefore,

$$F(x) = \int_x^0 e^{-a} da = -e^{-a}\Big|_0^x = -e^{-x} + 1 = 1 - e^{-x} \quad \text{if } x \geq 0$$

(a)

Solution:

$$(c) \text{ Find } x \text{ such that } P(X \leq x) = 0.95.$$

$$(b) \text{ Find the probabilities, } P\left(\frac{1}{4} \leq X \leq 2\right) \text{ and } P\left(-\frac{2}{3} \leq X \leq \frac{2}{3}\right).$$

$$(a) \text{ Find the distribution function.}$$

Example 3.5 Let X have density function $f(x) = e^{-x}$, if $x \geq 0$, and zero otherwise.

Let us warm-up with an example.

3.4.2 Some Common Continuous Random Variables

Simulation in Chapter 4.

— one can obtain any other random variable from the fundamental model whose unique DF is its own inverse, i.e., $F(x) = F_{-\text{I}}(x)$, as you will see from von Neumann's Fundamental Theorem of

**universality of the fundamental model

— The fundamental model has infinitely many copies of itself within it! You can see this since its DF F is the identity function on $[0, 1]$ or equivalently how the dyadic binary tree is identical below a given node in the tree no matter which node in the tree you choose.

— The fundamental model is equivalent to infinite tosses of a fair coin (see using binary expansion of any $x \in (0, 1)$ if you want as suggested in optional Exercise 2.1 on infinite primitive sigma-algebra).

— The fundamental model is equivalent to infinite tosses of a fair coin (see using binomial expansion of Bernoulli($1/2$) trials, and the

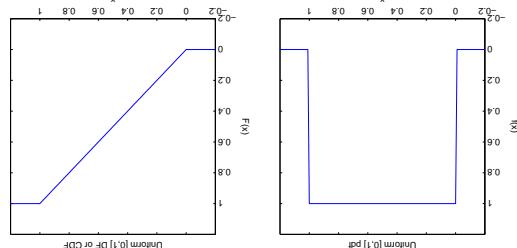


Figure 3.11: A convenient but mathematically imprecise Matlab plot from a polyline interpolation for the PDF and DF of the Uniform($0, 1$) continuous RV X .

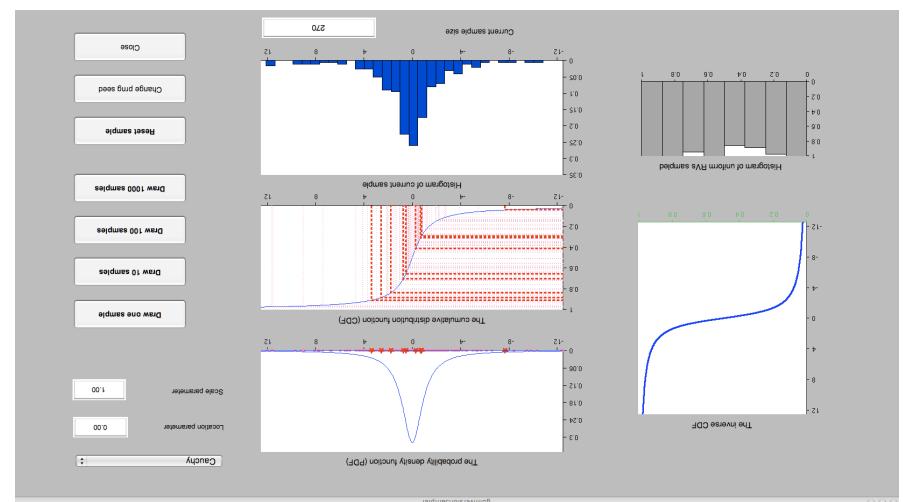


Figure 4.9: Visual Cognitive Tool GUI: Inversion Sampling from $X \sim \text{Cauchy}$.

The M-file `guiInversionSampler.m` will bring a graphical user interface (GUI) as shown in Figure 4.9. Using the drop-down menu change from the default target distribution Uniform(-5, 5) to Cauchy RV of Model 13. Now repeatedly push the "Draw one sample" button several times and compare with the simulation process. You can also press "Draw 10 samples" several times and compare with the simulation process. Next try changing the numbers in the Scale parameter and Location parameter boxes from the default values of 1.00 and 0.00, respectively. Although our formulation of Cauchy RV is also called Standard Cauchy as it implicitly had a location parameter of 0.00 and scale parameter of 1. With a pencil and paper (in conjunction with a wikipedia search if you have to) try to rewrite the PDF in (3.53) with an additional location parameter μ and scale parameter σ .

by calling the interactive visual cognitive tool:

```
>> guiInversionSampler
```

We can simply call the function to draw a sample from, say the Laplace($\lambda = 1.0$) RV by

(b)

$$P\left(\frac{1}{4} \leq X \leq 2\right) = F(2) - F\left(\frac{1}{4}\right) = 0.634 \text{ (3 sig. fig.)}$$

$$P\left(-\frac{1}{2} \leq X \leq \frac{1}{2}\right) = F\left(\frac{1}{2}\right) - F\left(-\frac{1}{2}\right) = 0.394 \text{ (3 sig. fig.)}$$

(c)

$$P(X \leq x) = F(x) = 1 - e^{-x} = 0.95$$

Therefore,

$$x = -\log(1 - 0.95) = 3.00 \text{ (3 sig. fig.)}.$$

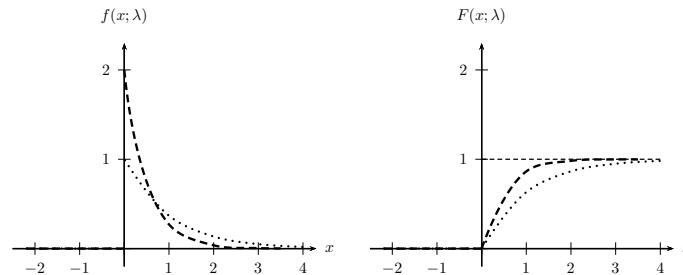


Figure 3.12: $f(x; \lambda)$ and $F(x; \lambda)$ of an exponential random variable where $\lambda = 1$ (dotted) and $\lambda = 2$ (dashed).

The previous example is a special case of the following parametric family of random variables.

Model 8 (Exponential(λ)) For a given $\lambda > 0$, an Exponential(λ) RV has the following PDF f and DF F and its complementary distribution function denoted by $\bar{F}(x; \lambda) := \mathbf{P}(X > x) = 1 - F(x; \lambda)$:

$$f(x; \lambda) = \mathbb{1}_{(0, \infty)} \lambda e^{-\lambda x} = \begin{cases} \lambda \exp(-\lambda x) & x > 0 \\ 0 & \text{otherwise} \end{cases}, \quad (3.27)$$

$$F(x; \lambda) = 1 - e^{-\lambda x}, \quad (3.28)$$

$$\bar{F}(x; \lambda) = e^{-\lambda x}. \quad (3.29)$$

The last two equations are derived from definitions as follows:

$$\begin{aligned} F(x; \lambda) &= \int_{-\infty}^x \mathbb{1}_{(0, \infty)} \lambda e^{-\lambda v} dv = \lambda \int_0^x e^{-\lambda v} dv = \lambda \left(-\frac{1}{\lambda} e^{-\lambda v} \right)_0^x = \left(-e^{-\lambda v} \right)_0^x \\ &= -e^{-\lambda x} - (-e^0) = -e^{-\lambda x} - (-1/e^0) = -e^{-\lambda x} - (-1/1) = -e^{-\lambda x} - (-1) = -e^{-\lambda x} + 1 \end{aligned}$$

$$\mathbf{P}(X > x) = 1 - \mathbf{P}(X \leq x) = 1 - F(x; \lambda) = 1 - \left(1 - e^{-\lambda x} \right) = e^{-\lambda x}$$

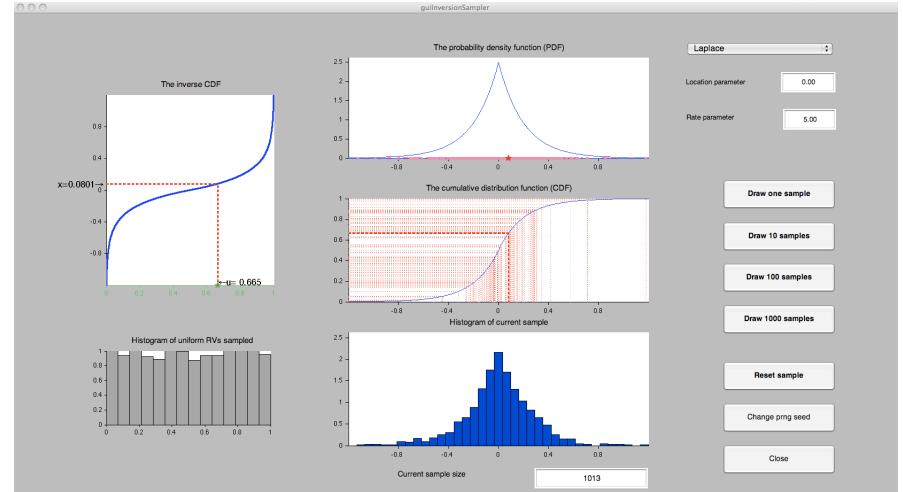
This distribution is unique because of its property of **memorylessness**, i.e., $\mathbf{P}(X > x+y | X > y) = e^{-\lambda x}$, and plays a fundamental role in modeling continuous time processes, such as time between occurrence of events of interest, as we will see in the sequel.

Labwork 136 (Rejection Sampler Demo – Laplace(5)) Let us comprehend the rejection sampler by calling the interactive visual cognitive tool:

```
>> guiInversionSampler
```

The M-file `guiInversionSampler.m` will bring a graphical user interface (GUI) as shown in Figure 4.8. Using the drop-down menu change from the default target distribution Uniform($-5, 5$) to Laplace(5). Now repeatedly push the “Draw one sample” button several times and comprehend the simulation process. You can also press “Draw 1000 samples” and see the density histogram of the generated samples. Next try changing the numbers in the “Rate parameter” box from 5.00 to 1.00 in order to alter the parameter λ of Laplace(λ) RV. If you are more adventurous then try to alter the number in the “Location parameter” box from 0.00 to some thing else, say 10.00. Although our formulation of Laplace(λ) implicitly had a location parameter of 0.00, we can easily introduce a location parameter μ into the PDF. With a pencil and paper try to rewrite the PDF in (4.2) with an additional location parameter μ .

Figure 4.8: Visual Cognitive Tool GUI: Inversion Sampling from $X \sim \text{Laplace}(5)$.



Simulation 137 (Laplace(λ)) Here is an implementation of an inversion sampler to draw IID samples from a Laplace(λ) RV X by transforming IID samples from the Uniform($0, 1$) RV U :

```
LaplaceInvCDF.m
function x = LaplaceInvCDF(u,lambda);
% Call Syntax: x = LaplaceInvCDF(u,lambda);
% or LaplaceInvCDF(u,lambda);
%
% Input      : lambda = rate parameter > 0,
%               u = an 1 X n array of IID samples from Uniform[0,1] RV
% Output     : an 1Xn array x of IID samples from Laplace(lambda) RV
%               or Inverse CDF of Laplace(lambda) RV
x=-(1/lambda)*sign(u-0.5).* log(1-2*abs(u-0.5));
```


Exercise 3.17 (Memoryless Server Times) Suppose customers in a Queue are served one at a time by a server whose service time is an independent and identical Exponential(λ) RV, with $\lambda = 1/10$. The server is immediately free to serve the next customer once the current customer being served is done. Suppose you just arrive and are the first in the queue and know that the server is busy serving another customer. You do not know how long the customer has already been in service. What is the probability that the server will be free after 2 units of time?

Let us introduce parameters for the lower and upper bounds of the interval upon which a continuous RV is uniformly distributed using the following probability model.

Model 9 (Uniform(θ_1, θ_2)) Given two real parameters $\theta_1, \theta_2 \in \mathbb{R}$, such that $\theta_1 < \theta_2$, the PDF of the Uniform(θ_1, θ_2) RV X is:

$$f(x; \theta_1, \theta_2) = \begin{cases} \frac{1}{\theta_2 - \theta_1} & \text{if } \theta_1 \leq x \leq \theta_2, \\ 0 & \text{otherwise} \end{cases} \quad (3.31)$$

and its DF given by $F(x; \theta_1, \theta_2) = \int_{-\infty}^x f(y; \theta_1, \theta_2) dy$ is:

$$F(x; \theta_1, \theta_2) = \begin{cases} 0 & \text{if } x < \theta_1 \\ \frac{x - \theta_1}{\theta_2 - \theta_1} & \text{if } \theta_1 \leq x \leq \theta_2, \\ 1 & \text{if } x > \theta_2 \end{cases} \quad (3.32)$$

Recall that we emphasise the dependence of the probabilities on the two parameters θ_1 and θ_2 by specifying them following the semicolon in the argument for f and F .

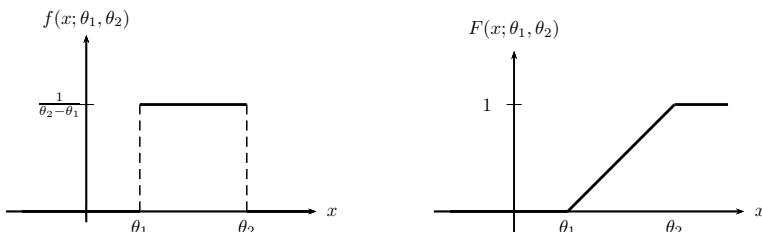


Figure 3.14: $f(x)$ and $F(x)$ of the Uniform(θ_1, θ_2) random variable X .

Exercise 3.18 Consider a random variable with a probability density function

$$f(x) = \begin{cases} k & \text{if } 2 \leq x \leq 6, \\ 0 & \text{otherwise} \end{cases}$$

- (a) Find the value of k .
- (b) Sketch the graphs of $f(x)$ and $F(x)$.

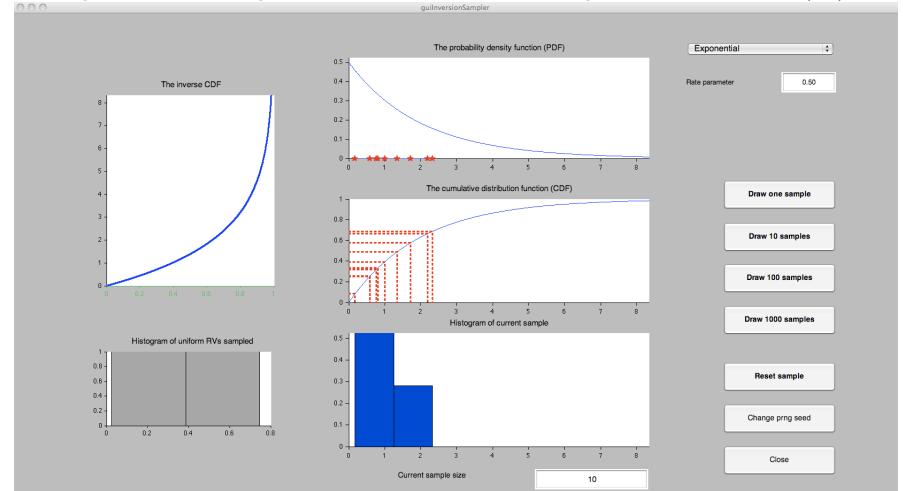
```
>> rand('twister',46678); % initialise the fundamental sampler
>> Lambda=1.0; % declare Lambda=1.0
>> x=ExpInvSam(rand(1,5),Lambda); % pass an array of 5 Uniform(0,1) samples from rand
>> disp(x); % display the Exponential(1.0) distributed samples
0.5945 2.5956 0.9441 1.9015 1.3973
```

Labwork 134 (Inversion Sampler Demo – Exponential(0.5)) Let us understand the inversion sampler by calling the interactive visual cognitive tool:

```
>> guiInversionSampler
```

The M-file `guiInversionSampler.m` will bring a graphical user interface (GUI) as shown in Figure 4.7. First change the target distribution from the default Uniform($-5, 5$) to Exponential(0.5) from the drop-down menu. Now push the “Draw 10 samples” button and comprehend the simulation process. Next try changing the “Rate Parameter” from 0.5 to 10.0 for example and generate several inversion samples and see the density histogram of the accumulating samples. You can press “Draw one sample” to really comprehend the inversion sampler in action one step at a time.

Figure 4.7: Visual Cognitive Tool GUI: Inversion Sampling from $X \sim \text{Exponential}(0.5)$.



It is straightforward to do replicate experiments. Consider the experiment of drawing five independent samples from the Exponential($\lambda = 1.0$) RV. Suppose we want to repeat or replicate this experiment seven times and find the sum of the five outcomes in each of these replicates. Then we may do the following:

```
>> rand('twister',1973); % initialise the fundamental sampler
>> % store 7 replications of 5 IID draws from Exponential(1.0) RV in array a
>> lambda=1.0; a= -1/lambda * log(rand(5,7)); disp(a);
0.7267 0.3226 1.2649 0.4786 0.3774 0.0394 1.8210
1.2698 0.4401 1.6745 1.4571 0.1786 0.4738 3.3690
```

Classwork 58 Note that the curve of $\phi(z)$ is S -shaped, increasing in a strictly monotone way from 0 at $-\infty$ to 1 at ∞ , and intersects the vertical axis at $1/2$. Draw this by hand too.

And use MATLAB's `erf` function to get $\phi(z)$ numerically instead of looking up the Table.

$$\Phi(z) = \frac{1}{2} \operatorname{erf}\left(\frac{z}{\sqrt{2}}\right) + \frac{1}{2} \quad (3.35)$$

We can express $\phi(z)$ in terms of the error function (`erf`) as follows:

Remark 27 The integral for $\phi(z)$ has no closed form expression and cannot be evaluated exactly by standard methods of calculus, but its values can be obtained numerically and tabulated. Values of $\phi(z)$ are tabulated in the "Standard Normal Distribution Function Table" in Sec. 6.6.

$$\Phi(z) = \frac{1}{\sqrt{2\pi}} \int_z^{\infty} e^{-\frac{a^2}{2}} da \quad (3.34)$$

The distribution function of Z is given by

$$\Phi(z) = \frac{1}{\sqrt{2\pi}} \int_z^{\infty} e^{-\frac{a^2}{2}} da$$

Do it step by step: $z^2, -z^2, -z^2/2, \exp(-z^2/2), \phi(z) = \frac{1}{\sqrt{2\pi}} \exp(-z^2/2)$ now!

Classwork 57 From the above exercise in calculus let us draw the graph of ϕ by hand now!

Thus, ϕ has a global maximum at 0, it is concave down if $z \in (-1, 1)$ and concave up if $z \in (-\infty, -1) \cup (1, \infty)$. This shows that the graph of ϕ is shaped like a smooth symmetric bell curve over the real line.

$$\frac{d\phi}{dz} = -\frac{\sqrt{2\pi}}{1} z \exp\left(-\frac{z^2}{2}\right) = -z\phi(z), \quad \frac{d^2\phi}{dz^2} = \frac{\sqrt{2\pi}}{1} (z^2 - 1) \exp\left(-\frac{z^2}{2}\right) = (z^2 - 1)\phi(z).$$

An exercise in calculus yields the first two derivatives of ϕ as follows:

$$\phi(z) = \frac{\sqrt{2\pi}}{1} \exp\left(-\frac{z^2}{2}\right). \quad (3.33)$$

Model 10 ($\text{Normal}(0, 1)$ or standard normal or Gaussian RV) A continuous random variable Z is called **standard normal** or **Gaussian** if its probability density function is

The standard normal distribution is the most important continuous probability distribution. It was first described by De Moivre in 1733 and subsequently by C. F. Gauss (1777 - 1855). Many random variables have a normal distribution, or they are approximately normal, or can be transformed into normal random variables in a relatively simple fashion. Furthermore, the normal distribution is a useful approximation of more complicated distributions.

For a given $\lambda > 0$, an $\text{Exponential}(\lambda)$ RV has the following

```
% Output : x = array of numbers in [0,1] from uniform[0,1] RV
% Input  : Lambda = rate parameter,
%          or ExpInvSam(u,Lambda);
% Call Syntex: x = ExpInvSam(u,Lambda);
% Return the inverse CDF based Sample from Exponential(Lambda) RV x
function x = ExpInvSam(u,Lambda);
    x=-((1/Lambda)* log(u));

```

is exactly how we defined X as the $\text{Exponential}(\lambda)$ RV in Model 8. This is implemented as the following function.

we could save a subtraction operation in the above algorithm by replacing $-(1/\lambda)\log(U)$

$$U \sim \text{Uniform}(0,1) \iff -U \sim \text{Uniform}(-1,0) \iff 1-U \sim \text{Uniform}(0,1),$$

Because of the following (recall Example 65):

```
% Lambda=1.0;
% same value for Lambda
% rand is the Fundamenta1(1) RV via function in ExpInvCDF.m
function x = ExpInvCDF(u,Lambda);
    x=-((1/Lambda)* log(1-u));

```

We can simply call the function to draw a sample from, say the $\text{Exponential}(\lambda = 1.0)$ RV by:

```
% Lambda=1.0;
% same value for Lambda
% rand is the Fundamenta1(1) RV via function in ExpInvCDF.m
function x = ExpInvCDF(u,Lambda);
    x=-((1/Lambda)* log(1-u));

```

Inversion Sampler for $\text{Exponential}(\lambda)$ as a function in the M-file:

$$f(x; \lambda) = \lambda e^{-\lambda x}, \quad F(x; \lambda) = 1 - e^{-\lambda x}, \quad F_{[-1]}(u; \lambda) = \frac{-1}{\lambda} \log_e(1-u) \quad (4.1)$$

Simulation 133 ($\text{Exponential}(\lambda)$) For a given $\lambda > 0$, an $\text{Exponential}(\lambda)$ RV has the following

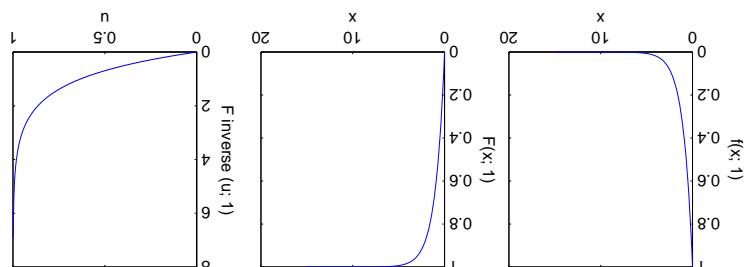


Figure 4.6: The PDF f , CDF F , and inverse CDF $F_{[-1]}$ of the $\text{Exponential}(\lambda = 1.0)$ RV.

Example 59 Find the probabilities, using normal tables, that a random variable having the standard normal distribution will take on a value:

- | | |
|---------------------|----------------------------|
| (a) less than 1.72 | (c) between 1.30 and 1.75 |
| (b) less than -0.88 | (d) between -0.25 and 0.45 |
| (a) | |

$$P(Z < 1.72) = \Phi(1.72) = 0.9573$$

- (b) First note that $P(Z < 0.88) = 0.8106$, so that

$$\begin{aligned} P(Z < -0.88) &= P(Z > 0.88) \\ &= 1 - P(Z < 0.88) \\ &= 1 - \Phi(0.88) \\ &= 1 - 0.8106 = 0.1894 \end{aligned}$$

(c) $P(1.30 < Z < 1.75) = \Phi(1.75) - \Phi(1.30) = 0.9599 - 0.9032 = 0.0567$

(d)

$$\begin{aligned} P(-0.25 < Z < 0.45) &= P(Z < 0.45) - P(Z < -0.25) \\ &= P(Z < 0.45) - (1 - P(Z < 0.25)) \\ &= \Phi(0.45) - (1 - \Phi(0.25)) \\ &= (0.6736) - (1 - 0.5987) \\ &= 0.2723 \end{aligned}$$

CONTINUOUS RANDOM VARIABLES: NOTATION

$f(x)$: Probability density function (PDF)

- $f(x) \geq 0$
- Areas underneath $f(x)$ measure probabilities.

$F(x)$: Distribution function (DF)

- $0 \leq F(x) \leq 1$
- $F(x) = P(X \leq x)$ is a probability
- $F'(x) = f(x)$ for every x where $f(x)$ is continuous
- $F(x) = \int_{-\infty}^x f(v)dv$
- $P(a < X \leq b) = F(b) - F(a) = \int_a^b f(v)dv$

It is just as easy to draw n IID samples from $\text{Uniform}(\theta_1, \theta_2)$ RV X by transforming n IID samples from the $\text{Uniform}(0, 1)$ RV as follows:

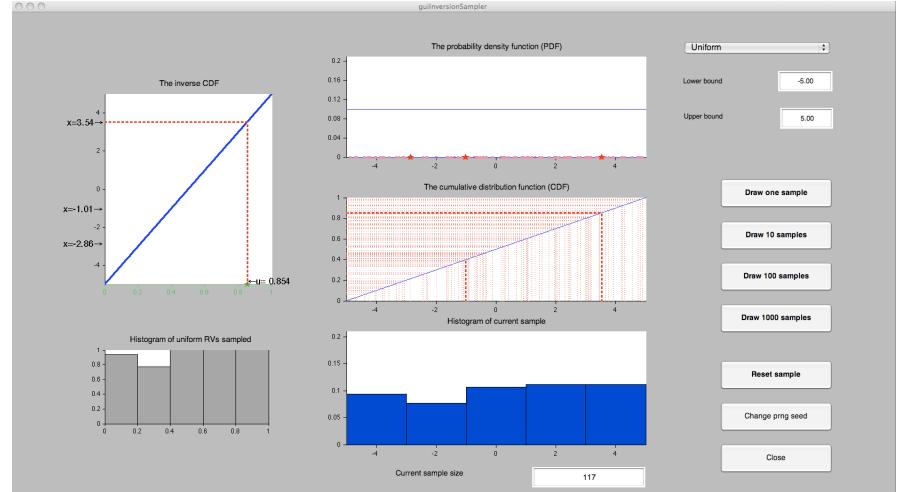
```
>> rand('twister',786543); % initialise the fundamental sampler
>> theta1=-83; theta2=1004; % declare values for parameters a and b
>> u=rand(1,5); % now u is an array of 5 samples from Uniform(0,1)
>> x=theta1+(theta2 - theta1)*u; % x is an array of 5 samples from Uniform(-83,1004) RV
>> disp(x); % display the 5 samples just drawn from Uniform(-83,1004) RV
465.3065 111.4994 14.3535 724.8881 254.0168
```

Labwork 132 (Inversion Sampler Demo – Uniform($-5, 5$)) Let us comprehend the inversion sampler by calling the interactive visual cognitive tool built by Jennifer Harlow under a grant from University of Canterbury's Centre for Teaching and Learning (UCTL):

```
>> guiInversionSampler
```

The M-file `guiInversionSampler.m` will bring a graphical user interface (GUI) as shown in Figure 4.5. The default target distribution is $\text{Uniform}(-5, 5)$. Now repeatedly push the “Draw one sample” button several times and comprehend the simulation process. You can press “Draw 100 samples” to really comprehend the inversion sampler in action after 100 samples are drawn and depicted in the density histogram of the accumulating samples. Next try changing the numbers in the “Lower bound” and “Upper bound” boxes in order to alter the parameters θ_1 and θ_2 of $\text{Uniform}(\theta_1, \theta_2)$ RV.

Figure 4.5: Visual Cognitive Tool GUI: Inversion Sampling from $X \sim \text{Uniform}(-5, 5)$.



Recall the $\text{Exponential}(\lambda)$ RV of Model 8. Let us simulate from it using the inversion sampler.

Let us consider the problem of simulating from an $\text{Exponential}(\lambda)$ RV with realisations in $\mathbb{R}_+ := [0, \infty) := \{x : x \geq 0, x \in \mathbb{R}\}$ to model the waiting time for a bus at a bus stop.

$$Y = 10^{\log_{10}(X)}$$

Example 2 If we can produce S_1 with a yield of 50% , how many moles of S_1 are needed to produce 100 g of S_2 ?

$$\cdot 009 - X^9 = \lambda$$

Example 6.6 Consider a simple financial example where an individual sells X items per day; the profit per item is \$5 and the overhead costs are \$500 per day. The original random variable is X , but the random variable Y which gives the daily profit is of more interest, where

Suppose we know the distribution of a random variable X . How do we find the distribution of a transformation of X , say $g(X)$? Before we answer this question let us ask a motivational question. Why are we interested in functions of random variables?

3.6 Transformations of random variables

1. Find k .
2. Find the probability that a bearing will last at least 1 year.

$$\begin{cases} 0 & \text{otherwise} \\ x \geq y & \end{cases} = (x)f$$

Ex. 3.21 — Let the random variable X be the time after which certain ball bearings wear out,

(b) For what value of τ is the reliability of the bulb exactly $\frac{1}{2}$?

(a) Assume that $\lambda = 0.01$, and that the probability that the bird will fly before t hours, $F(t)$, is given by the distribution function for an Exponential(λ) random variable (recall, $F(t; \lambda) = \int_{-\infty}^t f(t; \lambda) dt$).
 Hint: Use the distribution function for an Exponential(λ) random variable (recall, $F(t; \lambda) = 1 - e^{-\lambda t}$).
 This t -specific probability is often called the reliability of the bulb.

$$\left\{ \begin{array}{ll} 0 & \text{otherwise} \\ 0 \leq y & 0 \end{array} \right\} = (t)f$$

Ex. 2-27 — Assume that a new light bulb will burn out at time t hours according to the proba-

2. Find the distribution function, F_F .

$$\begin{cases} \frac{x}{y} & \text{otherwise} \\ 0 \end{cases} = (x)f$$

Ex. 3.19 — Consider the probability density function

3.5 Exercises in Continuous Random Variables

Proof: The “one-line proof” of the proposition is due to the following equalities:

$\exists x \forall y \forall z ((x)_{\bar{H}} = ((x)_{\bar{H}} \supseteq z) \bar{H} = (x \supseteq \{(z = (y)_{\bar{H}} : y\}_{y \in \bar{H}}) \bar{H} = (x \supseteq (\cup_{y \in \bar{H}} y) \bar{H})$

This yields the inversion sampler or the inverse (CI) sampler, where we (i) generate $u \sim$ Uniform(0, 1) and (ii) return $x = F_{\{I\}}^{-1}(u)$, as formalised by the following algorithm.

Algorithm	3 Inversion Sampler or inverse (CDF) Sampler
1: input:	(1) $F^{-1}(x)$, inverse of the DF of the target RV X , (2) the fundamental sampler
2: initialize:	set the seed, if any, for the fundamental sampler
3: output:	a sample X distributed according to F
4: draw $u \sim \text{Uniform}(0, 1)$	

This algorithm emphasizes the fundamental sampler's availability in an *input* step, and its set-up needs in an *initialise* step. In the following sections, we will not mention these initial steps; they will be taken for granted. In the following sections, we will not mention these initial steps; they will be taken for granted. The direct applicability of Algorithm 3 is limited to univariate densities for which the inverse of the cumulative distribution function is explicitly known. The next section will consider some examples.

Recall the Uniform(0,1) RV of Model 9 with the following PDF, DF and inverse DF. Let us simulate from it using the inversion sampler.

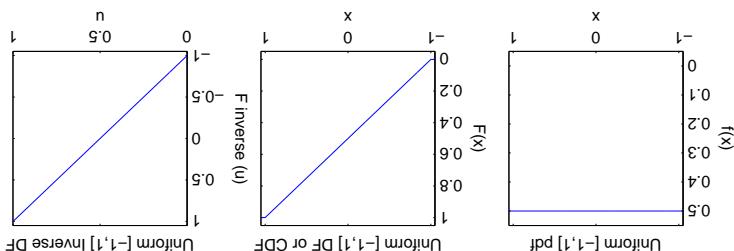


Figure 4.4: A plot of the PDF, DF or CDF and inverse DF of the Uniform(-1, 1) RV X .

Here is a simple implementation of the bivariate Smoother for the Unifitorm(θ_1, θ_2) RV in MATLAB:

$$n(\iota_\theta - \iota_\theta)(\iota_\theta - \iota_\theta) = n(\iota_\theta - \iota_\theta) \iff \iota_\theta + n(\iota_\theta - \iota_\theta) = x \iff \frac{\iota_\theta - \iota_\theta}{\iota_\theta - x} = n$$

Simulation 131 (Uniform(θ_1, θ_2)) To simulate from Uniform(θ_1, θ_2) RV X using the Inversion method, we first need to find $F_{[-1]}(u)$ by solving for x in terms of $u = F(x; \theta_1, \theta_2)$:

```

    >>> rand(0,1)
    >>> read('testset',786); % initialize the fundamental sample for Unifrom(0,1)
    >>> theta1=1; theta2=1; % declare values for parameters theta1 and theta2
    >>> xtheta1=(theta1+(theta2- theta1)*rand(1,1)); % draw 1 sample from Unifrom(-1,1)
    >>> xtheta2=(theta2-(theta1+theta2)*rand(1,1)); % draw 1 sample from Unifrom(-1,1)
    >>> xtheta1+(xtheta2-theta2)*sample; % calculate the sample from Unifrom(-1,1)
    >>> disp(x); % display the sample from Unifrom(-1,1) RV

```

```
>>> rand('twister',786); % initialize the fundamental sample for uniform(0,1)
>>> theta1=-pi;% declare values for parameters theta1 and theta2
```

3.6.1 A Review of Inverse Images

Hence in a great many situations we are more interested in functions of random variables. Let us return to our original question of determining the distribution of a transformation or function of X . First note that this transformation of X is itself another random variable, say $Y = g(X)$, where g is a function from a subset \mathbb{X} of \mathbb{R} to a subset \mathbb{Y} of \mathbb{R} , i.e., $g : \mathbb{X} \rightarrow \mathbb{Y}$, $\mathbb{X} \subset \mathbb{R}$ and $\mathbb{Y} \subset \mathbb{R}$.

The **inverse image** of a set A is the set of all real numbers in \mathbb{X} whose image is in A , i.e.,

$$g^{[-1]}(A) = \{x \in \mathbb{X} : g(x) \in A\}.$$

In other words,

$$x \in g^{[-1]}(A) \text{ if and only if } g(x) \in A.$$

For example,

- if $g(x) = 2x$ then $g^{[-1]}([4, 6]) = [2, 3]$
- if $g(x) = 2x + 1$ then $g^{[-1]}([5, 7]) = [2, 3]$
- if $g(x) = x^3$ then $g^{[-1]}([1, 8]) = [1, 2]$
- if $g(x) = x^2$ then $g^{[-1]}([1, 4]) = [-2, -1] \cup [1, 2]$
- if $g(x) = \sin(x)$ then $g^{[-1]}([-1, 1]) = \mathbb{R}$
- if ...

For the singleton set $A = \{y\}$, we write $g^{[-1]}(y)$ instead of $g^{[-1]}(\{y\})$. For example,

- if $g(x) = 2x$ then $g^{[-1]}(4) = \{2\}$
- if $g(x) = 2x + 1$ then $g^{[-1]}(7) = \{3\}$
- if $g(x) = x^3$ then $g^{[-1]}(8) = \{2\}$
- if $g(x) = x^2$ then $g^{[-1]}(4) = \{-2, 2\}$
- if $g(x) = \sin(x)$ then $g^{[-1]}(0) = \{k\pi : k \in \mathbb{Z}\} = \{\dots, -3\pi, -2\pi, -\pi, 0, \pi, 2\pi, 3\pi, \dots\}$
- if ...

If $g : \mathbb{X} \rightarrow \mathbb{Y}$ is one-to-one (injective) and onto (surjective), then the inverse image of a singleton set is itself a singleton set. Thus, the inverse image of such a function g becomes itself a function and is called the **inverse function**. One can find the inverse function, if it exists by the following steps:

Step 1; write $y = g(x)$

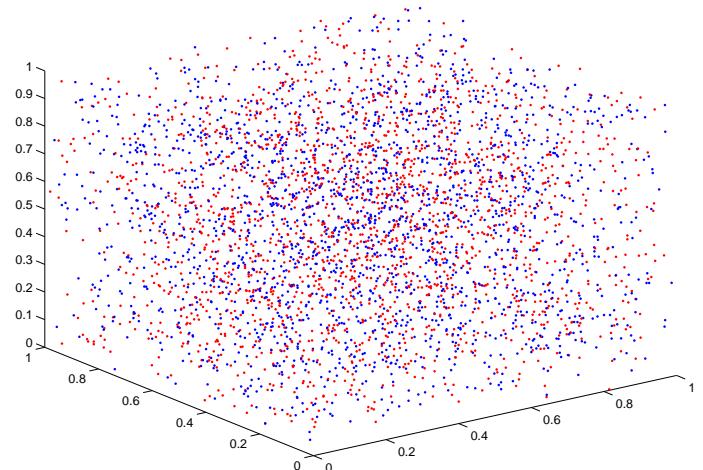
Step 2; solve for x in terms of y

Step 3; set $g^{-1}(y)$ to be this solution

We write g^{-1} whenever the inverse image $g^{[-1]}$ exists as an inverse function of g . Thus, the inverse function g^{-1} is a specific type of inverse image $g^{[-1]}$. For example,

- if $g(x) = 2x$ then $g : \mathbb{R} \rightarrow \mathbb{R}$ is injective and surjective and therefore its inverse function is:
Step 1; $y = 2x$, Step 2; $x = \frac{y}{2}$, Step 3; $g^{-1}(y) = \frac{y}{2}$

Figure 4.3: Triplet point clouds from the “Mersenne Twister” with two different seeds (see Lab-work 130). .



4.3 Simulation of non-Uniform(0, 1) Random Variables

The Uniform(0, 1) RV of Model 7 forms the foundation for random variate generation and simulation. This is appropriately called the fundamental model or experiment, since every other experiment can be obtained from this one.

Next, we simulate or generate samples from other RVs by making the following two assumptions:

1. independent samples from the Uniform(0, 1) RV can be generated, and
2. real arithmetic can be performed exactly in a computer.

Both these assumptions are, in fact, not true and require a more careful treatment of the subject. We may return to these careful treatments later on.

4.3.1 Inversion Sampler for Continuous Random Variables

Proposition 58 (Inversion sampler) Let $F(x) := \int_{-\infty}^x f(y) dy : \mathbb{R} \rightarrow [0, 1]$ be a continuous DF with density f , and let its inverse $F^{[-1]} : [0, 1] \rightarrow \mathbb{R}$ be:

$$F^{[-1]}(u) := \inf\{x : F(x) = u\}.$$

Then, $F^{[-1]}(U)$ has the distribution function F , provided U is a Uniform(0, 1) RV. Recall $\inf(A)$ or infimum of a set A of real numbers is the greatest lower bound of every element of A .

To answer this question we must first observe that the inverse image g_{-1}^{-1} satisfies the following properties.

- If $y = \sin(x)$ and domain of g is $[0, \frac{\pi}{2}] \leftarrow [0, 1]$ then its inverse function $g_{-1}(y) = \arcsin(y)$, i.e., if $g(x) = \sin(x) : [0, \frac{\pi}{2}] \leftarrow [0, 1]$ then the inverse image $g_{-1}(y)$ for $y \in [0, 1]$ is given by the inverse function $g_{-1}(y) = \arcsin(y) : [0, 1] \leftarrow [0, \frac{\pi}{2}]$.
- If $y = \sin(x)$ and domain of g is $[-\frac{\pi}{2}, \frac{\pi}{2}] \leftarrow [-1, 1]$ then its inverse function $g_{-1}(y) = \arcsin(y)$, i.e., if $g(x) = \sin(x) : [-\frac{\pi}{2}, \frac{\pi}{2}] \leftarrow [-1, 1]$ then the inverse image $g_{-1}(y)$ for $y \in [-1, 1]$ is given by the inverse function $g_{-1}(y) = \arcsin(y) : [-1, 1] \leftarrow [-\frac{\pi}{2}, \frac{\pi}{2}]$.
- If $y = x^2$ and domain of g is $(-\infty, 0] \leftarrow (-\infty, 0]$ then its inverse function $g_{-1}(y) = -\sqrt{y}$, i.e., if $g(x) = x^2 : (-\infty, 0] \leftarrow [0, +\infty)$ then the inverse image $g_{-1}(y)$ for $y \in [0, +\infty)$ is given by the inverse function $g_{-1}(y) = -\sqrt{y} : [0, +\infty) \leftarrow (-\infty, 0]$.
- If $y = x^2$ and domain of g is $[0, +\infty) \leftarrow [0, +\infty)$ then its inverse function is $g_{-1}(y) = \sqrt{y}$, i.e., if $g(x) = x^2 : [0, +\infty) \leftarrow [0, +\infty)$ then the inverse image $g_{-1}(y)$ for $y \in [0, +\infty)$ is given by the inverse function $g_{-1}(y) = \sqrt{y} : [0, +\infty) \leftarrow [0, +\infty)$.

• if $g(x) = 2x + 1$ then $g^{-1}(y) = \frac{y-1}{2}$ is injective and subjective and therefore its inverse function is:

$$\text{Step 1; } y = 2x + 1, \text{ Step 2; } x = \frac{y-1}{2}, \text{ Step 3; } g^{-1}(y) = \frac{y-1}{2}$$

• if $g(x) = x^3$ then $g^{-1}(y) = \sqrt[3]{y}$ is injective and subjective and therefore its inverse function is:

$$\text{Step 1; } y = x^3, \text{ Step 2; } x = \sqrt[3]{y}, \text{ Step 3; } g^{-1}(y) = \sqrt[3]{y}$$

• if $g(x) = 2x - 1$ then $g^{-1}(y) = \frac{y+1}{2}$ is injective and subjective and therefore its inverse function is:

$$\text{Step 1; } y = 2x - 1, \text{ Step 2; } x = \frac{y+1}{2}, \text{ Step 3; } g^{-1}(y) = \frac{y+1}{2}$$

• if $g(x) = 2x^2 + 1$ then $g^{-1}(y) = \pm\sqrt{\frac{y-1}{2}}$ is injective and subjective and therefore its inverse function is:

$$\text{Step 1; } y = 2x^2 + 1, \text{ Step 2; } x^2 = \frac{y-1}{2}, \text{ Step 3; } x = \pm\sqrt{\frac{y-1}{2}}$$

• if $g(x) = x^3 - 1$ then $g^{-1}(y) = \sqrt[3]{y+1}$ is injective and subjective and therefore its inverse function is:

$$\text{Step 1; } y = x^3 - 1, \text{ Step 2; } x^3 = y+1, \text{ Step 3; } g^{-1}(y) = \sqrt[3]{y+1}$$

• if $g(x) = x^3 + 1$ then $g^{-1}(y) = \sqrt[3]{y-1}$ is injective and subjective and therefore its inverse function is:

$$\text{Step 1; } y = x^3 + 1, \text{ Step 2; } x^3 = y-1, \text{ Step 3; } g^{-1}(y) = \sqrt[3]{y-1}$$

• if $g(x) = x^3 - 3x$ then $g^{-1}(y) = \sqrt[3]{\frac{y+3}{2}} - \sqrt[3]{\frac{y-3}{2}}$ is injective and subjective and therefore its inverse function is:

$$\text{Step 1; } y = x^3 - 3x, \text{ Step 2; } x^3 - 3x = y, \text{ Step 3; } x = \sqrt[3]{\frac{y+3}{2}} - \sqrt[3]{\frac{y-3}{2}}$$

• if $g(x) = x^3 + 3x$ then $g^{-1}(y) = \sqrt[3]{\frac{y-3}{2}} + \sqrt[3]{\frac{y+3}{2}}$ is injective and subjective and therefore its inverse function is:

$$\text{Step 1; } y = x^3 + 3x, \text{ Step 2; } x^3 + 3x = y, \text{ Step 3; } x = \sqrt[3]{\frac{y-3}{2}} + \sqrt[3]{\frac{y+3}{2}}$$

- if $g(x) = 2x + 1$ then $g : \mathbb{R} \rightarrow \mathbb{R}$ is injective and subjective and therefore its inverse function

```

<> rand(1,10) % generate a 1 x 10 array of PRNs
ans =
0.9649 0.9058 0.1270 0.9134 0.6324 0.0975 0.2785 0.5469 0.9575 0.9595
<> rand(1,10) % generate another 1 x 10 array of PRNs
ans =
0.9576 0.9706 0.9572 0.4864 0.8003 0.1419 0.4218 0.9157 0.7922 0.9595
<> rand(1,10) % generate a 1 x 10 array of PRNs
ans =
0.9649 0.9058 0.1270 0.9134 0.6324 0.0975 0.2785 0.5469 0.9575 0.9595
<> rand(1,10) % reproduce the first array
ans =
0.9649 0.9058 0.1270 0.9134 0.6324 0.0975 0.2785 0.5469 0.9575 0.9595
<> rand(1,10) % reproduce the second array
ans =
0.9576 0.9706 0.9572 0.4864 0.8003 0.1419 0.4218 0.9157 0.7922 0.9595
<> set the seed.

```

In general, you can use any seed value to initialize your PRNG. You may use the `clock` command to set the seed:

```

>>> x=rand(2,2000); % store PRNs in a 3x2000 matrix named x
>>> plot3(x(1,:),x(2,:),x(3,:));
>>> xlabel('twister',1234)
>>> title('Mersenne Twister',1234)

```

by the following code:

Labwork 130 (3D plots of triples generated by the "Mersenne Twister" by rotating the 3D plot generated correlation between triplets generated by the "Mersenne Twister") Try to find an

Compare this with the 3D plot of triplets from RANDU of Labwork 127. Which of these two PRNGs do you think is "more random" looking and why?

satisfies the axioms of probability and gives the desired probability of the event A from the transformation $Y = g(X)$ in terms of the probability of the random variable X . It is crucial to understand this from the sample space underpinning by the random variable X . The event given by the inverse image of A of the underlying experiment in the sense that Equation (3.36) is just short-hand for its actual meaning:

$$\left((A)_{[\mathbb{I}-]} \delta \ni X \right) d = (A \ni (X) \delta) d$$

Consequently,

$$g_{[-l]}(A_1 \cup A_2 \cup \dots) = g_{[-l]}(A_1) \cup g_{[-l]}(A_2) \cup \dots$$

- For any collection of sets $\{A_1, A_2, \dots\}$,
 - For any set A , $g_{-[I]}(A^c) = (g_{-[I]}(A))^c$
 - $g_{-[I]}(\mathbb{X}) = \mathbb{X}$

Now, let us return to our question of determining the distribution of the transformation $y(X)$. To answer this question we must first observe that the inverse image $g^{-1}[y]$ satisfies the following properties:

- If $y = \sin(x)$ and domain of y is $(-\infty, \infty)$, then its inverse function is $y = g^{-1}(x) = \arcsin(y)$, i.e., if $y = \sin(x)$ for $y \in [-1, 1]$, then $x = \arcsin(y)$.
- If $y = x^2$ and domain of y is $[0, \infty)$, then its inverse function is $y = g^{-1}(x) = \sqrt{x}$, i.e., if $y = x^2$ for $y \in [0, \infty)$, then $x = \sqrt{y}$.
- If $y = x^2$ and domain of y is $(-\infty, 0]$, then its inverse function is $y = g^{-1}(x) = -\sqrt{-x}$, i.e., if $y = x^2$ for $y \in [0, \infty)$, then $x = -\sqrt{-y}$.
- If $y = \sin(x)$ and domain of y is $[0, \frac{\pi}{2}]$, then its inverse function is $y = g^{-1}(x) = \arcsin(x)$, i.e., if $y = \sin(x)$ for $y \in [0, \frac{\pi}{2}]$, then $x = \arcsin(y)$.
- If $y = \sin(x)$ and domain of y is $[-\frac{\pi}{2}, 0]$, then its inverse function is $y = g^{-1}(x) = -\arcsin(-x)$, i.e., if $y = \sin(x)$ for $y \in [-\frac{\pi}{2}, 0]$, then $x = -\arcsin(-y)$.
- If $y = \sin(x)$ and domain of y is $[-1, 1]$, then its inverse function is $y = g^{-1}(x) = \arcsin(x)$, i.e., if $y = \sin(x)$ for $y \in [-1, 1]$, then $x = \arcsin(y)$.
- If $y = \sin(x)$ and domain of y is $[-\frac{\pi}{2}, \frac{\pi}{2}]$, then its inverse function is $y = g^{-1}(x) = \arcsin(x)$, i.e., if $y = \sin(x)$ for $y \in [-\frac{\pi}{2}, \frac{\pi}{2}]$, then $x = \arcsin(y)$.

However, you need to be careful by naming the domain to obtain the inverse function for the following examples:

Step 1: $y = x^3$, Step 2: $x = y^{\frac{1}{3}}$, Step 3: $y^{-\frac{1}{3}}(y) = y^{\frac{2}{3}}$

- if $g(x) = x^3$ then $g^{-1} : y \leftarrow \sqrt[3]{x}$ is injective and subjective and therefore its inverse function is:

- if $g(x) = 2x + 1$ then $g : \mathbb{R} \rightarrow \mathbb{R}$ is injective and surjective and therefore its inverse function is:

CHAPTER 3. RANDOM VARIABLES

Because we have more than one random variable to consider, namely, X and its transformation $Y = g(X)$, we will subscript the probability density or mass function and the distribution function by the random variable itself. For example we denote the distribution function of X by $F_X(x)$ and by the random variable itself. For example we denote the distribution function of Y by $F_Y(y)$.

$$\cdot \left(\left\{ (V)_{[I]} b \ni (\omega)X : \mathcal{U} \ni \omega \right\} \right)_d = (\{A \ni ((\omega)X)b : \mathcal{U} \ni \omega\})_d$$

3.6.2 Transformations of discrete random variables

For a discrete random variable X with probability mass function f_X we can obtain the probability mass function f_Y of $Y = g(X)$ using Equation (3.36) as follows:

$$\begin{aligned} f_Y(y) &= \mathbf{P}(Y = y) = \mathbf{P}(Y \in \{y\}) \\ &= P(g(X) \in \{y\}) = P\left(X \in g^{[-1]}\(\{y\}\)\right) \\ &= P\left(X \in g^{[-1]}(y)\right) = \sum_{x \in g^{[-1]}(y)} f_X(x) = \sum_{x \in \{x: g(x)=y\}} f_X(x) . \end{aligned}$$

This gives the formula:

$$f_Y(y) = \mathbf{P}(Y = y) = \sum_{x \in g^{[-1]}(y)} f_X(x) = \sum_{x \in \{x: g(x)=y\}} f_X(x) . \quad (3.37)$$

Example 62 Let X be the discrete random variable with probability mass function f_X as tabulated below:

x	-1	0	1
$f_X(x) = \mathbf{P}(X = x)$	$\frac{1}{4}$	$\frac{1}{2}$	$\frac{1}{4}$

If $Y = 2X$ then the transformation $g(X) = 2X$ has inverse image $g^{[-1]}(y) = \{y/2\}$. Then, by Equation (3.37) the probability mass function of Y is expressed in terms of the known probabilities of X as:

$$f_Y(y) = \mathbf{P}(Y = y) = \sum_{x \in g^{[-1]}(y)} f_X(x) = \sum_{x \in \{y/2\}} f_X(x) = f_X(y/2) ,$$

and tabulated below:

y	-2	0	2
$f_Y(y)$	$\frac{1}{4}$	$\frac{1}{2}$	$\frac{1}{4}$

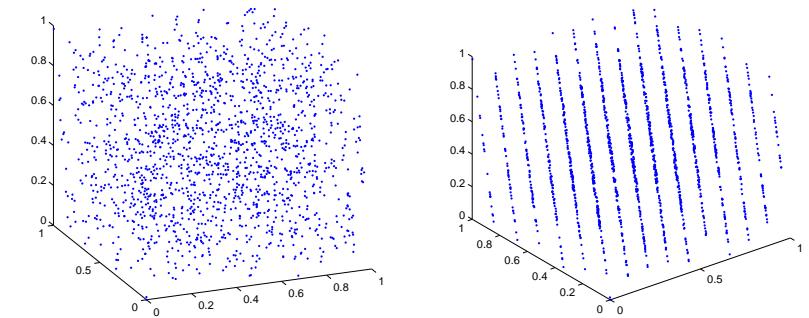
Example 63 If X is the random variable in the previous Example then what is the probability mass function of $Y = 2X + 1$? Once again,

$$f_Y(y) = \mathbf{P}(Y = y) = \sum_{x \in g^{[-1]}(y)} f_X(x) = \sum_{x \in \{(y-1)/2\}} f_X(x) = f_X((y-1)/2) ,$$

and tabulated below:

y	-1	1	3
$f_Y(y)$	$\frac{1}{4}$	$\frac{1}{2}$	$\frac{1}{4}$

Figure 4.2: The LCG called `RANDU` with $(m, a, c) = (2147483648, 65539, 0)$ has strong correlation between three consecutive points as: $x_{i+2} = 6x_{k+1} - 9x_k$. The two plots are showing (x_i, x_{i+1}, x_{i+2}) from two different view points. .



The number of random numbers n should at most be about $m/1000$ in order to avoid the future sequence from behaving like the past. Thus, if $m = 2^{32}$ then a new generator, with a new suitable set of (m, a, c, x_0, n) should be adopted after the consumption of every few million pseudo-random numbers.

The LCGs are the least sophisticated type of PRNGs. They are easier to understand but are not recommended for intensive simulation purposes. The next section briefly introduces a more sophisticated PRNG we will be using in this course. Moreover our implementation of LCGs using the variable precision integer package is extremely slow in MATLAB and is only of pedagogical interest.

4.2.2 Generalized Feedback Shift Register and the “Mersenne Twister” PRNG

The following generator termed `twister` in MATLAB is recommended for use in simulation. It has extremely long periods, low correlation and passes most statistical tests (the DIEHARD statistical tests). The `twister` random number generator of Makoto Matsumoto and Takuji Nishimura is a variant of the twisted generalized feedback shift-register algorithm, and is known as the “Mersenne Twister” generator [Makoto Matsumoto and Takuji Nishimura, *Mersenne Twister: A 623-dimensionally equidistributed uniform pseudorandom number generator*, ACM Transactions on Modeling and Computer Simulation, Vol. 8, No. 1 (Jan. 1998), Pages 3–30]. It has a Mersenne prime period of $2^{19937} - 1$ (about 10^{6000}) and is **equi-distributed** in 623 dimensions. It uses 624 words of state per generator and is comparable in speed to the other generators. The recommended default seed is 5489. See <http://www.math.sci.hiroshima-u.ac.jp/~m-mat/MT/emt.html> and http://en.wikipedia.org/wiki/Mersenne_twister for details.

Let us learn to implement the MATLAB function that generates PRNs. In MATLAB the function `rand` produces a deterministic PRN sequence. First, read `help rand`. We can generate PRNs as follows.

Labwork 129 (Calling PRNG in MATLAB) In MATLAB `rand` is basic PRNG command.

$$\cdot ((\kappa)_{1-\delta}) \frac{\kappa p}{p} ((\kappa)_{1-\delta}) x f = ((\kappa)_{1-\delta}) x \mathcal{A} \frac{\kappa p}{p} = (\kappa) x \mathcal{A} \frac{\kappa p}{p} = (\kappa) x f$$

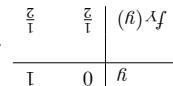
Now, let us use a form of chainrule to compute the density of Y as follows:

- First, let us consider the case when g is **monotone and increasing** on the range of the random variable X . In this case g^{-1} is also an increasing function and we can obtain the distribution function of $Y = g(X)$ in terms of the distribution function of X as

The easiest case for transformations of continuous random variables is when y is one-to-one and

Suppose we know F_x and/or f_x of a continuous random variable X . Let $Y = g(X)$ be a transformation of X . Our objective is to obtain F_Y and/or f_Y of Y from F_X and/or f_X .

3.6.3 Transformations of continuous random variables



and finally tabulated below:

$$\begin{aligned} & \frac{\zeta}{\zeta} = (0)Xf = (x)Xf \sum_{\{0=x^x\}} = (0)\lambda f \\ & + \frac{\frac{4}{4}}{\frac{4}{4}} = (1)Xf + (-1)Xf = (x)Xf \sum_{\{1=x^x\}} = (1)\lambda f \end{aligned}$$

computed for each $y \in \{0, 1\}$ as follows:

$$(x)Xf \sum_{\{\tilde{n}=\varepsilon x:x\}} = (x)Xf \sum_{\{\tilde{n}=(x)\delta x\}} = (x)Xf \sum_{\{(\tilde{n})_{[1]}=\delta x\}} = (\tilde{n}=\lambda) \mathbf{d} = (\tilde{n})\lambda f$$

terms of the known probabilities of X as:

In fact, obtaining the probability of a one-to-one transformation of a discrete random variable as shown in the next Example.

96

```

ans = 0.0041 0.5239 0.0755 0.7624 0.6496 0.0769 0.9030 0.4259 0.9948 0.8868
ans = 0.0006 0.3934 0.4117 0.1234567,10) / . 214743347
ans = LincogenGen(214743367,4271,10),087874458,10) / . 214743347
ans = LincogenGen(214743399,40692,0,01234567,10) / . 214743399
ans = 0.0006 0.3934 0.4117 0.1234567,10) / . 214743399
ans = 0.0006 0.3934 0.4117 0.1234567,10) / . 214743399

```

In Kurnitski's Art of Computer Programming, vol. 2, for generating pseudo-random numbers for simple laboratory work (§28 (Fishman20 and Leucymer21 LCGs)), the following two LCGs are recommended

Labwork 126 (TCG with maximum period length of 256) Consider the linear congruential sequence with $(m, a, c, x_0, n) = (256, 137, 123, 13, 256)$. First check that these parameters do indeed satisfy the three condition above and therefore can produce the maximal period length of only $m = 256$. Modify the input parameter to `lincmgen` and repeat Labwork 125 in order to first produce a sequence of length 257. Do you see that the period is of maximal length of 256 as opposed to the generator of Labwork 125? Next produce a figure to visualise the sequence as done

3. $a - 1$ is a multiple of 4 if m is a multiple of 4
4. $a + 1$ is a multiple of 4 if m is a multiple of 4

1. c and m are relatively prime, if m as shown in the examples considered earlier. The LCG will have a full period if and only

Since f maps a three minute set $\{1, 2, \dots, m-1\}$ into itself, such systems are bound to have a repeating cycle of numbers called the **period**. In Labwork 12, the generator $L_{10, 12}$ has period (10, 1) of length 2, the generator $L_{11, 12}$ has period (8, 12), the generator $L_{12, 13}$ has period (9, 13), the generator $L_{13, 14}$ has period (8, 12), the generator $L_{14, 15}$ has period (9, 13), the generator $L_{15, 16}$ has period (8, 12), the generator $L_{16, 17}$ has period (9, 13), the generator $L_{17, 18}$ has period (8, 12), the generator $L_{18, 19}$ has period (9, 13), the generator $L_{19, 20}$ has period (8, 12), the generator $L_{20, 21}$ has period (9, 13), the generator $L_{21, 22}$ has period (8, 12), the generator $L_{22, 23}$ has period (9, 13), the generator $L_{23, 24}$ has period (8, 12), the generator $L_{24, 25}$ has period (9, 13), the generator $L_{25, 26}$ has period (8, 12), the generator $L_{26, 27}$ has period (9, 13), the generator $L_{27, 28}$ has period (8, 12), the generator $L_{28, 29}$ has period (9, 13), the generator $L_{29, 30}$ has period (8, 12), the generator $L_{30, 31}$ has period (9, 13), the generator $L_{31, 32}$ has period (8, 12), the generator $L_{32, 33}$ has period (9, 13), the generator $L_{33, 34}$ has period (8, 12), the generator $L_{34, 35}$ has period (9, 13), the generator $L_{35, 36}$ has period (8, 12), the generator $L_{36, 37}$ has period (9, 13), the generator $L_{37, 38}$ has period (8, 12), the generator $L_{38, 39}$ has period (9, 13), the generator $L_{39, 40}$ has period (8, 12), the generator $L_{40, 41}$ has period (9, 13), the generator $L_{41, 42}$ has period (8, 12), the generator $L_{42, 43}$ has period (9, 13), the generator $L_{43, 44}$ has period (8, 12), the generator $L_{44, 45}$ has period (9, 13), the generator $L_{45, 46}$ has period (8, 12), the generator $L_{46, 47}$ has period (9, 13), the generator $L_{47, 48}$ has period (8, 12), the generator $L_{48, 49}$ has period (9, 13), the generator $L_{49, 50}$ has period (8, 12), the generator $L_{50, 51}$ has period (9, 13), the generator $L_{51, 52}$ has period (8, 12), the generator $L_{52, 53}$ has period (9, 13), the generator $L_{53, 54}$ has period (8, 12), the generator $L_{54, 55}$ has period (9, 13), the generator $L_{55, 56}$ has period (8, 12), the generator $L_{56, 57}$ has period (9, 13), the generator $L_{57, 58}$ has period (8, 12), the generator $L_{58, 59}$ has period (9, 13), the generator $L_{59, 60}$ has period (8, 12), the generator $L_{60, 61}$ has period (9, 13), the generator $L_{61, 62}$ has period (8, 12), the generator $L_{62, 63}$ has period (9, 13), the generator $L_{63, 64}$ has period (8, 12), the generator $L_{64, 65}$ has period (9, 13), the generator $L_{65, 66}$ has period (8, 12), the generator $L_{66, 67}$ has period (9, 13), the generator $L_{67, 68}$ has period (8, 12), the generator $L_{68, 69}$ has period (9, 13), the generator $L_{69, 70}$ has period (8, 12), the generator $L_{70, 71}$ has period (9, 13), the generator $L_{71, 72}$ has period (8, 12), the generator $L_{72, 73}$ has period (9, 13), the generator $L_{73, 74}$ has period (8, 12), the generator $L_{74, 75}$ has period (9, 13), the generator $L_{75, 76}$ has period (8, 12), the generator $L_{76, 77}$ has period (9, 13), the generator $L_{77, 78}$ has period (8, 12), the generator $L_{78, 79}$ has period (9, 13), the generator $L_{79, 80}$ has period (8, 12), the generator $L_{80, 81}$ has period (9, 13), the generator $L_{81, 82}$ has period (8, 12), the generator $L_{82, 83}$ has period (9, 13), the generator $L_{83, 84}$ has period (8, 12), the generator $L_{84, 85}$ has period (9, 13), the generator $L_{85, 86}$ has period (8, 12), the generator $L_{86, 87}$ has period (9, 13), the generator $L_{87, 88}$ has period (8, 12), the generator $L_{88, 89}$ has period (9, 13), the generator $L_{89, 90}$ has period (8, 12), the generator $L_{90, 91}$ has period (9, 13), the generator $L_{91, 92}$ has period (8, 12), the generator $L_{92, 93}$ has period (9, 13), the generator $L_{93, 94}$ has period (8, 12), the generator $L_{94, 95}$ has period (9, 13), the generator $L_{95, 96}$ has period (8, 12), the generator $L_{96, 97}$ has period (9, 13), the generator $L_{97, 98}$ has period (8, 12), the generator $L_{98, 99}$ has period (9, 13), the generator $L_{99, 100}$ has period (8, 12), the generator $L_{100, 101}$ has period (9, 13).

- Second, let us consider the case when g is **monotone and decreasing** on the range of the random variable X . In this case g^{-1} is also a decreasing function and we can obtain the distribution function of $Y = g(X)$ in terms of the distribution function of X as

$$F_Y(y) = P(Y \leq y) = P(g(X) \leq y) = P(X \geq g^{-1}(y)) = 1 - F_X(g^{-1}(y)) ,$$

and the density of Y as

$$f_Y(y) = \frac{d}{dy} F_Y(y) = \frac{d}{dy} (1 - F_X(g^{-1}(y))) = -f_X(g^{-1}(y)) \frac{d}{dy} (g^{-1}(y)) .$$

For a monotonic and decreasing g , its inverse function g^{-1} is also decreasing and consequently the density f_Y is indeed positive because $\frac{d}{dy} (g^{-1}(y))$ is negative.

We can combine the above two cases and obtain the following **change of variable formula** for the probability density of $Y = g(X)$ when g is one-to-one and monotone on the range of X .

$$f_Y(y) = f_X(g^{-1}(y)) \left| \frac{d}{dy} g^{-1}(y) \right| . \quad (3.38)$$

The steps involved in finding the density of $Y = g(X)$ for a one-to-one and monotone g are:

1. Write $y = g(x)$ for x in range of x and check that $g(x)$ is monotone over the required range to apply the change of variable formula.
2. Write $x = g^{-1}(y)$ for y in range of y .
3. Obtain $\left| \frac{d}{dy} g^{-1}(y) \right|$ for y in range of y .
4. Finally, from Equation (3.38) get $f_Y(y) = f_X(g^{-1}(y)) \left| \frac{d}{dy} g^{-1}(y) \right|$ for y in range of y .

Let us use these four steps to obtain the density of monotone transformations of continuous random variables.

Example 65 Let X be Uniform(0,1) random variable and let $Y = g(X) = 1 - X$. We are interested in the density of the tranformed random variable Y . Let us follow the four steps and use the change of variable formula to obtain f_Y from f_X and g .

1. $y = g(x) = 1 - x$ is a monotone decreasing function over $0 \leq x \leq 1$, the range of X . So, we can apply the change of variable formula.
2. $x = g^{-1}(y) = 1 - y$ is a monotone decreasing function over $1 - 0 \geq 1 - x \geq 1 - 1$, i.e., $0 \leq y \leq 1$.
3. For $0 \leq y \leq 1$,

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

4. we can use Equation (3.38) to find the density of Y as follows:

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

for $0 \leq y \leq 1$

```
x(i) = double(X); % convert to double
end
```

We can call it for some arbitrary input arguments as follows:

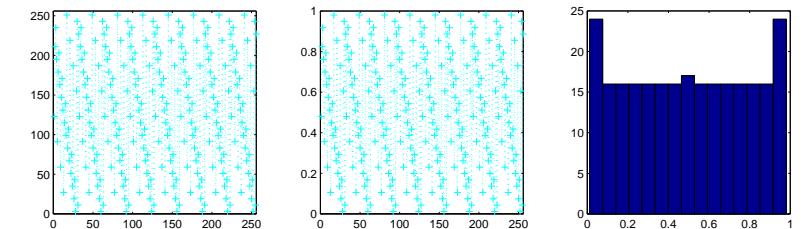
```
>> LinConGen(13,12,11,10,12)
ans =    10      1      10      1      10      1      10      1      10      1      10      1
>> LinConGen(13,10,9,8,12)
ans =     8     11      2      3      0      9      8     11      2      3      0      9
```

and observe that the generated sequences are not “random” for input values of (m, a, c, x_0, n) equalling $(13, 12, 11, 10, 12)$ or $(13, 10, 9, 8, 12)$. Thus, we need to do some work to determine the *suitable* input integers (m, a, c, x_0, n) .

Labwork 125 (LCG with period length of 32) Consider the linear congruential sequence with $(m, a, c, x_0, n) = (256, 137, 0, 123, 257)$ with period length of only $32 < m = 256$. We can visualise the sequence as plots in Figure 4.1 after calling the following M-file.

```
LCGSeq=LinConGen(256,137,0,123,257)
subplot(1,3,1)
plot(LCGSeq,'c+:')
axis([0 256 0 256]); axis square
LCGSeqIn01=LCGSeq ./ 256
subplot(1,3,2)
plot(LCGSeqIn01,'c+:')
axis([0 256 0 1]); axis square
subplot(1,3,3)
hist(LCGSeqIn01,15)
axis square
```

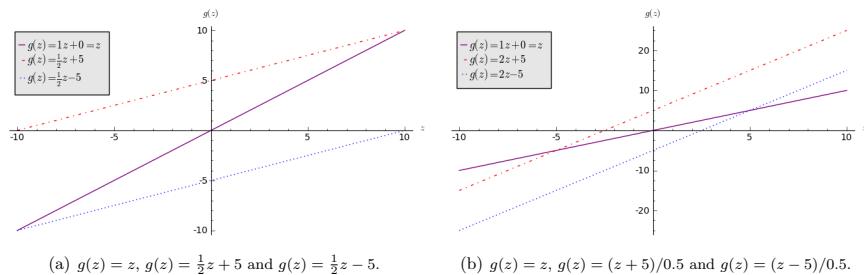
Figure 4.1: The linear congruential sequence of $\text{LinConGen}(256, 137, 0, 123, 257)$ with non-maximal period length of 32 as a line plot over $\{0, 1, \dots, 256\}$, scaled over $[0, 1]$ by a division by 256 and a histogram of the 256 points in $[0, 1]$ with 15 bins.



Choosing the *suitable* magic input (m, a, c, x_0, n)

The linear congruential generator is a special case of a *discrete dynamical system*:

$$x_i = f(x_{i-1}), \quad f : \{0, 1, 2, \dots, m-1\} \rightarrow \{0, 1, 2, \dots, m-1\} \text{ and } f(x_{i-1}) = (ax_{i-1} + c) \pmod{m} .$$



4. we can use Equation (3.38) and Equation (3.33) which gives

$$f_Z(z) = \phi(z) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{z^2}{2}\right),$$

to find the density of Y as follows:

$$f_Y(y) = f_Z(g^{-1}(y)) \left| \frac{d}{dy} g^{-1}(y) \right| = \phi\left(\frac{y-\mu}{\sigma}\right) \frac{1}{\sigma} = \frac{1}{\sigma\sqrt{2\pi}} \exp\left[-\frac{1}{2}\left(\frac{y-\mu}{\sigma}\right)^2\right],$$

for $-\infty < y < \infty$.

Thus, we have obtained the expression for the probability density function of the linear transformation $\sigma Z + \mu$ of the standard normal random variable Z . This analysis leads to the following definition.

Model 11 (Normal(μ, σ^2) RV) Given a location parameter $\mu \in (-\infty, +\infty)$ and a scale parameter $\sigma^2 > 0$, the Normal(μ, σ^2) or Gaussian(μ, σ^2) random variable X has probability density function:

$$f(x; \mu, \sigma^2) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left[-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2\right] \quad (\sigma > 0). \quad (3.39)$$

This is simpler than it may at first look. $f(x; \mu, \sigma^2)$ has the following features.

- μ is the expected value or mean parameter and σ^2 is the variance parameter. These concepts, mean and variance, are described in more detail in the next section on expectations.
- $1/(\sigma\sqrt{2\pi})$ is a constant factor that makes the area under the curve of $f(x)$ from $-\infty$ to ∞ equal to 1, as it must be.
- The curve of $f(x)$ is symmetric with respect to $x = \mu$ because the exponent is quadratic. Hence for $\mu = 0$ it is symmetric with respect to the y -axis $x = 0$.
- The exponential function decays to zero very fast — the faster the decay, the smaller the value of σ .

Chapter 4

Simulation

“Anyone who considers arithmetical methods of producing random digits is, of course, in a state of sin.” — John von Neumann (1951)

4.1 Physical Random Number Generators

Physical devices such as the BINGO machine demonstrated in class can be used to produce an integer uniformly at random from a finite set of possibilities. Such “ball bouncing machines” used in the British national lottery as well as the New Zealand LOTTO are complex nonlinear systems that are extremely sensitive to initial conditions (“chaotic” systems) and are physical approximations of the probability model called a “well-stirred urn” or an equi-probable de Moivre($1/k, \dots, 1/k$) random variable.

Let us look at the New Zealand LOTTO draws at <http://lotto.nzpages.co.nz/statistics.html> and convince ourselves that all forty numbers $\{1, 2, \dots, 39, 40\}$ seem to be drawn uniformly at random. The British lottery animation at <http://understandinguncertainty.org/node/39> shows how often each of the 49 numbers came up in the first 1240 draws. Are these draws really random? We will answer these questions in the sequel (see <http://understandinguncertainty.org/node/40> if you can’t wait).

4.2 Pseudo-Random Number Generators

Our probability model and the elementary continuous Uniform(0, 1) RV are built from the abstract concept of a random variable over a probability triple. A direct implementation of these ideas on a computing machine is not possible. In practice, random variables are typically simulated by **deterministic** methods or algorithms. Such algorithms generate sequences of numbers whose behavior is virtually indistinguishable from that of truly random sequences. In computational statistics, simulating realisations from a given RV is usually done in two distinct steps. First, sequences of numbers that imitate independent and identically distributed (IID) Uniform(0, 1) RVs are generated. Second, appropriate transformations are made to these imitations of IID Uniform(0, 1) random variates in order to imitate IID random variates from other random variables or other random structures. These two steps are essentially independent and are studied by two non-overlapping groups of researchers. The formal term **pseudo-random number generator** (PRNG) or simply **random number generator** (RNG) usually refers to some deterministic algorithm used in the first step to produce pseudo-random numbers (PRNs) that imitate IID Uniform(0, 1) random variates.

$$\cdot \frac{\sigma}{\mu - X} = Z$$

Hence we often transform a general $\text{Normal}(\mu, \sigma^2)$ random variable, X , to a standardised $\text{Normal}(0, 1)$ random variable, Z , by the substitution:

$$\begin{aligned} \left(\frac{\sigma}{\mu - X} \right) \Phi &= \left(\frac{\sigma}{\mu - x} \right) F_Z \\ P(X(x; \mu, \sigma^2) > Z) &= P(x > \bar{x}) = P(Z > \bar{Z}) \end{aligned}$$

We know that if $X = g(Z) = \sigma Z + \mu$ then X is the $\text{Normal}(\mu, \sigma^2)$ random variable. Therefore, **Proof:** Let Z be a $\text{Normal}(0, 1)$ random variable with distribution function $\Phi(z) = P(Z \leq z)$.

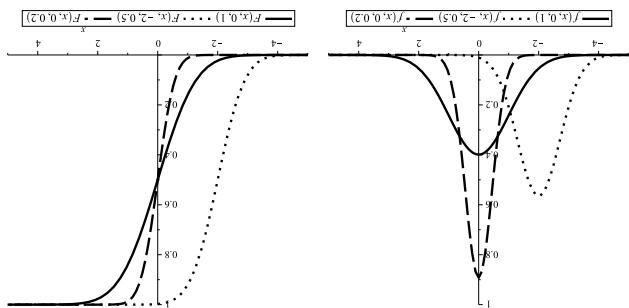
$$F_X(x; \mu, \sigma^2) = \Phi\left(\frac{x - \mu}{\sigma}\right)$$

normal random variable Z are related by:

Proposition 28 (One Table to Rule Them All Gaussians) The distribution function $F_X(x; \mu, \sigma^2)$ of the $\text{Normal}(\mu, \sigma^2)$ random variable X and the distribution function $F_Z(z) = \Phi(z)$ of the standard normal random variable Z are related by:

Using the direct method's Equation 3.41, we can obtain the distribution function of the $\text{Normal}(0, 1)$ in the Standard normal distribution function table in Sec. 6.6.

Figure 3.15: PDF and DF of a $\text{Normal}(\mu, \sigma^2)$ RV for different values of μ and σ^2



Here we need x as the upper limit of integration and so we write v in the integrand.

$$F(x; \mu, \sigma^2) = \frac{1}{\sigma\sqrt{2\pi}} \int_x^\infty \exp\left[-\frac{1}{2} \left(\frac{v - \mu}{\sigma}\right)^2\right] dv \quad (3.40)$$

The normal distribution has the distribution function

1, 3, 2, 1, 2, 3, 3

Ex. 3.56 — What is the sample mean and sample variance of the following dataset:

3.15 Exercises in Statistics

Example 68 Suppose that the amount of cosmic radiation to which a person is exposed when flying by jet across the United States is a random variable, X , having a normal distribution with a mean of 4.35 mrem and a standard deviation of 0.59 mrem. What is the probability that a person will be exposed to more than 5.20 mrem of cosmic radiation on such a flight?

Solution:

$$\begin{aligned} P(X > 5.20) &= 1 - P(X \leq 5.20) \\ &= 1 - F(5.20) \\ &= 1 - \Phi\left(\frac{5.20 - 4.35}{0.59}\right) \\ &= 1 - \Phi(1.44) \\ &= 1 - 0.9251 \\ &= 0.0749 \end{aligned}$$

After some more notions you will see that $\text{Normal}(0, 1)$ RV can actually be obtained from an IID process of $\text{Bernoulli}(\theta)$ RVs. This is an instance of the central limit theorem. To appreciate this we first need to understand what we mean by statistics and then familiarise ourselves with notions of convergence of random variables.

Direct method

If the transformation g in $Y = g(X)$ is not necessarily one-to-one then special care is needed to obtain the distribution function or density of Y . For a continuous random variable X with a known distribution function F_X we can obtain the distribution function F_Y of $Y = g(X)$ using Equation (3.36) as follows:

$$\begin{aligned} F_Y(y) &= P(Y \leq y) = P(Y \in (-\infty, y]) \\ &= P(g(X) \in (-\infty, y]) = P\left(X \in g^{[-1]}((-\infty, y])\right) = P(X \in \{x : g(x) \in (-\infty, y]\}) \quad (3.41) \end{aligned}$$

In words, the above equalities just mean that the probability that $Y \leq y$ is the probability that X takes a value x that satisfies $g(x) \leq y$. We can use this approach if it is reasonably easy to find the set $g^{[-1]}((-\infty, y]) = \{x : g(x) = (-\infty, y]\}$.

Example 69 Let X be any random variable with distribution function F_X . Let $Y = g(X) = X^2$. Then we can find F_Y , the distribution function of Y from F_X as follows:

- Since $Y = X^2 \geq 0$, if $y < 0$ then $F_Y(y) = P(X \in \{x : x^2 < y\}) = P(X \in \emptyset) = 0$.
- If $y \geq 0$ then

$$\begin{aligned} F_Y(y) = P(Y \leq y) &= P(X^2 \leq y) \\ &= P(-\sqrt{y} \leq X \leq \sqrt{y}) \\ &= F_X(\sqrt{y}) - F_X(-\sqrt{y}) . \end{aligned}$$

By differentiation we get:

Labwork 122 (favourite word cloud) This is just for fun. Produce a “word cloud” of your honours thesis or summer project or any other document that fancies your interest by using *wordle* from <http://www.wordle.net/>. Play with the aesthetic features to change colour, shapes, etc.

3.14.10 Machine Sensor Data

Instrumentation of modern machines, such as planes, rockets and cars allow the sensors in the machines to collect live data and dynamically take *decisions* and subsequent *actions* by executing algorithms to drive their devices in response to the data that is streaming into their sensors. For example, a rocket may have to adjust its boosters to compensate for the prevailing directional changes in wind in order to keep going up and launch a satellite. These types of decisions and actions, theorised by *controlled Markov processes*, typically arise in various fields of engineering such as, aerospace, civil, electrical, mechanical, robotics, etc.

In an observational setting, without an associated control problem, one can use machine sensor data to get information about some state of the system or phenomenon, i.e., what is it doing? or where is it?, etc. Sometimes sensors are attached to a sample of individuals from a wild population, say Emperor Penguins in Antarctica where the phenomenon of interest may be the diving habits of this species after the eggs hatch. As an other example we can attach sensors to a double pendulum and find what it is doing when we give it a spin.

Based on such observational data the experimenter typically tries to learn about the behaviour of the system from the sensor data to estimate parameters, test hypotheses, etc. Such types of experiments are typically performed by scientists in various fields of science, such as, astronomy, biology, chemistry, geology, physics, etc.

Chaotic Time Series of a Double Pendulum

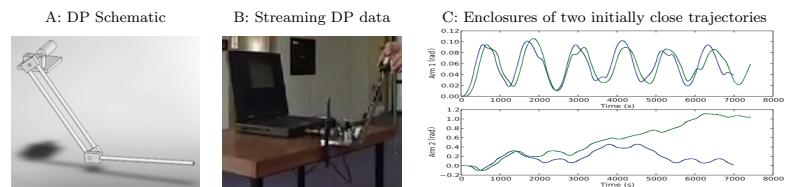


Figure 3.37: Double Pendulum

Sensors called *optical encoders* have been attached to the top end of each arm of a chaotic double pendulum in order to obtain the angular position of each arm through time as shown in Figure 3.37. Time series of the angular position of each arm for two trajectories that were initialized very similarly, say the angles of each arm of the double pendulum are almost the same at the initial time of release. Note how quickly the two trajectories diverge! System with such a sensitivity to initial conditions are said to be *chaotic*.

Labwork 123 (A Challenging Task) Try this if you are interested. Read any of the needed details about the design and fabrication of the double pendulum at <http://www.math.canterbury.ac.nz/~r.sainudiin/lmse/double-pendulum/>. Then use MATLAB to generate a plot similar to Figure 3.37(C) using time series data of trajectory 1 and trajectory 2 linked from the bottom of the above URL.

3.14.9 Textual Data

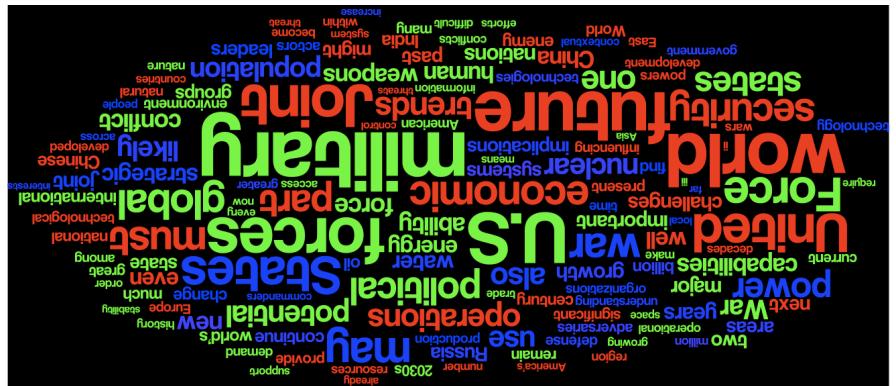
Processing and analyzing textual data to make a decision is another important computational task in this field. An obvious example is machine translation and a less obvious one is exploratory data analysis of the text content.

- a large document
 - etc.

2010 Report by the US Department of Defense. This document was downloaded from <http://www.jfcom.mil/newslink/storyarchive/2010/JOE-2010-o.pdf>. The first paragraph of this 2010 document reads:

AROOT THIS STUDY THE joint operating environment is intended to inform joint concepts development and experimentation throughout the Department of Defense. It provides a perspective on future trends, shocks, contexts, and implications for joint force commanders and other leaders and professionals in the national security field. This document is speculative in nature and does not suppose to predict what will happen in the next twenty-five years. Rather, it is intended to serve as a starting point for discussions about the future security environment at the operational level of war. inquiries about the Joint Operating Environment should be directed to USJFCOM Public Affairs, 1562 Miltcher Avenue, Suite 200, Norfolk, VA 23551-2488, (757) 386-6555.

Figure 3.36: Wordle of JOE 2010



Wordle is a toy for generating word clouds from text that you provide. The clouds give greater prominence to words that appear more frequently in the source text. You can tweak your clouds with different fonts, layouts, and color schemes. The images you create with Wordle are free to use however you like. You can print them out, or save them to the Wordle gallery to share with your friends.

We can try to produce a statistic of this document by recording textual content. Then we can produce a "word histogram" or "word visula

χ^2 is called the **chi-square** random variable with one degree of freedom. This distribution plays a fundamental role in hypothesis testing as we will see in Inferential Theory and was derived at the beginning of last century to settle “supposedly evidence-based disputes” among scientists using mathematics.

$$\left. \begin{aligned} 0 &\leqslant f_j \\ 0 &> f_j \end{aligned} \right\} = (f_j)_{\lambda f}$$

then by Equation (3.43) the density of $Y = X^2$ is:

$$(\partial/\partial x -) dx \otimes \frac{\sqrt{-g}}{1} = (x)\phi = (x)Xf$$

Example 70 If X is the standard normal random variable with density

$$\left. \begin{array}{l} 0 \leqslant \underline{\kappa} \text{ 且 } ((\underline{\kappa}^\wedge -)Xf + (\underline{\kappa}^\wedge)Xf) \frac{\underline{\kappa}^\wedge \zeta}{1} \\ 0 > \underline{\kappa} \text{ 且 } 0 \end{array} \right\} = (\underline{\kappa})^X f$$

and the probability density function of $Y = X^z$ is:

$$\left. \begin{array}{l} 0 < \mu \in \\ 0 > \mu \end{array} \right\} = (\mu)_{\mathcal{A}}$$

Therefore, the distribution function of $Y = X^2$ is:

$$\begin{aligned} & \cdot ((\underline{\mathcal{R}} \wedge -)Xf + (\underline{\mathcal{R}} \wedge)Xf) \frac{\underline{\mathcal{R}} \wedge \zeta}{1} = \\ & \left((\underline{\mathcal{R}} \wedge -)Xf \frac{\zeta}{1} - \underline{\mathcal{R}} \wedge \frac{\zeta}{1} \right) - (\underline{\mathcal{R}} \wedge)Xf \frac{\zeta}{1} - \underline{\mathcal{R}} \wedge \frac{\zeta}{1} = \\ & ((\underline{\mathcal{R}} \wedge -)XJ) \frac{\underline{\mathcal{R}} p}{p} - ((\underline{\mathcal{R}} \wedge)XJ) \frac{\underline{\mathcal{R}} p}{p} = \\ & ((\underline{\mathcal{R}} \wedge -)XJ - (\underline{\mathcal{R}} \wedge)XJ) \frac{\underline{\mathcal{R}} p}{p} = ((\underline{\mathcal{R}} \wedge)XJ) \frac{\underline{\mathcal{R}} p}{p} = (\underline{\mathcal{R}} \wedge)XJ \end{aligned}$$

- If $y < 0$ then •

• If $y > 0$ then $f(y) = ((\hbar)(\lambda)H)\frac{y}{p} = (\hbar y)(\lambda H)\frac{1}{p}$

3.7 Exercises in Transformations of Random Variables

Ex. 3.22 — Let X be the outcome of a fair die roll with probability mass function given by

$$f_X(x) = \begin{cases} \frac{1}{6} & \text{if } x \in \{1, 2, 3, 4, 5, 6\} \\ 0 & \text{otherwise.} \end{cases}$$

If $Y = (X - 3)^2$ then find the probability mass function of Y , $f_Y(y)$.

Ex. 3.23 — Given a natural number n as a parameter, i.e., given a parameter $n \in \{1, 2, 3, \dots\}$, let X be a discrete uniform random variable on the finite set

$$\mathbb{X} = \{-n, -n+1, \dots, -1, 0, 1, \dots, n-1, n\}$$

i.e. the probability mass function of X is:

$$f_X(x; n) = \begin{cases} \frac{1}{2n+1} & \text{if } x \in \mathbb{X} \\ 0 & \text{otherwise.} \end{cases}$$

Find the probability mass function $f_Y(y; n)$ for $Y = |X|$, the absolute value of X .

Ex. 3.24 — If X is a Geometric(θ) random variable and $Y = (\frac{1}{2})^X$ then find an expression for $f_Y(y)$.

Ex. 3.25 — If X is a Poisson(λ) random variable find the probability mass function, $f_Y(y)$, of

$$Y = \frac{1}{(X+1)^2}.$$

Ex. 3.26 — If X is a continuous random variable with probability density function

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

find the probability density function of $Y = e^X$.

Ex. 3.27 — If X , the received power at an antenna is an Exponential(λ) random variable then find the probability density function of the amplitude $Y = \sqrt{X}$.

Ex. 3.28 — If X is a Uniform(a, b) random variable where $0 < a < b$, find the probability density function, $f_Y(y)$, of

$$Y = \log_e(X).$$

3.8 Expectations

Expectation is perhaps the most fundamental concept in probability theory. In fact, probability is itself an expectation as you will soon see!

Expectation is one of the fundamental concepts in probability. The expected value of a real-valued random variable gives the population mean, a measure of the centre of the distribution of the variable in some sense. Its variance measures its spread and so on.

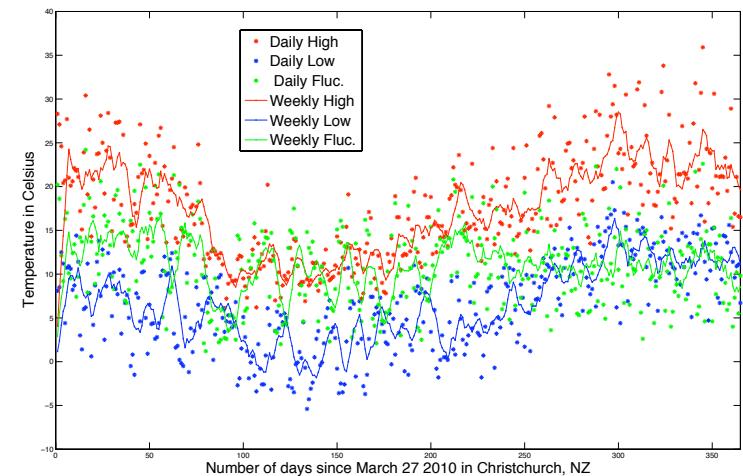
```

clf % clears all current figures

% Daily max and min temperature in the 100 days with good data
% before last date in this data, i.e., March 27 2011 in Christchurch NZ
H365Days = T(end-365:end,2);
L365Days = T(end-365:end,3);
F365Days = H365Days-L365Days; % assign the maximal fluctuation, i.e. max-min
plot(H365Days,'r*') % plot daily high or maximum temperature = Tmax
hold on; % hold the Figure so that we can overlay more plots on it
plot(L365Days,'b*') % plot daily low or minimum temperature = Tmin
plot(F365Days, 'g*') % plot daily Fluctuation = Tmax - Tmin
% filter for running means
windowSize = 7;
WeeklyHighs = filter(ones(1,windowSize)/windowSize,1,H365Days);
plot(WeeklyHighs,'r.-')
WeeklyLows = filter(ones(1,windowSize)/windowSize,1,L365Days);
plot(WeeklyLows,'b.-')
WeeklyFlucs = filter(ones(1,windowSize)/windowSize,1,F365Days);
plot(WeeklyFlucs,'g.-')
xlabel('Number of days since March 27 2010 in Christchurch, NZ','FontSize',20);
ylabel('Temperature in Celsius','FontSize',20);
MyLeg = legend('Daily High','Daily Low',' Daily Fluc. ','Weekly High','Weekly Low',...
    'Weekly Fluc. ','Location','NorthEast')
% Create legend
% legend1 = legend(axes1,'show');
set(MyLeg,'FontSize',20);
xlim([0 365]); % set the limits or boundary on the x-axis of the plots
hold off % turn off holding so we stop overlaying new plots on this Figure

```

Figure 3.35: Daily temperatures in Christchurch for one year since March 27 2010



The **mean** which characterises the central location of the random variable X is merely the expectation of the identity function $g(x) = x$:

$$\mathbf{E}(X) = \begin{cases} \sum_x xf(x) & \text{if } X \text{ is a discrete RV} \\ \int_{-\infty}^{\infty} xf(x)dx & \text{if } X \text{ is a continuous RV} \end{cases}$$

Often, mean is denoted by μ .

The **variance** which characterises the spread or the variability of the random variable X is also the expectation of the function $g(x) = (x - \mathbf{E}(X))^2$:

$$\mathbf{V}(X) = \mathbf{E}((X - \mathbf{E}(X))^2) = \begin{cases} \sum_x (x - \mathbf{E}(X))^2 f(x) & \text{if } X \text{ is a discrete RV} \\ \int_{-\infty}^{\infty} (x - \mathbf{E}(X))^2 f(x) dx & \text{if } X \text{ is a continuous RV} \end{cases}$$

Often, variance is denoted by σ^2 .

INTUITIVELY, WHAT IS EXPECTATION?

Definition 32 gives expectation as a “weighted average” of the possible values. This is true but some intuitive idea of expectation is also helpful.

- Expectation is what you expect.

Consider tossing a fair coin. If it is heads you lose \$10. If it is tails you win \$10. What do you expect to win? Nothing. If X is the amount you win then

$$\mathbf{E}(X) = -10 \times \frac{1}{2} + 10 \times \frac{1}{2} = 0.$$

So what you expect (nothing) and the weighted average ($\mathbf{E}(X) = 0$) agree.

- Expectation is a long run average.

Suppose you are able to repeat an experiment independently, over and over again. Each experiment produces one value x of a random variable X . If you take the average of the x values for a large number of trials, then this average converges to $\mathbf{E}(X)$ as the number of trials grows. In fact, this is called the **law of large numbers**.

We can concretize the above two intuitive insights by the following two examples.

Example 71 (Winnings on Average) Let $Y = r(X)$. Then

$$\mathbf{E}(Y) = \mathbf{E}(r(X)) = \int r(x) dF(x).$$

Think of playing a game where we draw $x \sim X$ and then I pay you $y = r(x)$. Then your average income is $r(x)$ times the chance that $X = x$, summed (or integrated) over all values of x .

Example 72 (Probability is an Expectation) Let A be an event and let $r(X) = \mathbb{1}_A(x)$. Recall $\mathbb{1}_A(x)$ is 1 if $x \in A$ and $\mathbb{1}_A(x) = 0$ if $x \notin A$. Then

$$\mathbf{E}(\mathbb{1}_A(X)) = \int \mathbb{1}_A(x) dF(x) = \int_A dF(x) = \mathbf{P}(X \in A) = \mathbf{P}(A) \quad (3.46)$$

Thus, probability is a special case of expectation. Recall our LTRF motivation for the definition of probability and make the connection.

CHAPTER 3. RANDOM VARIABLES

161

Figure 3.33: Google Earth Visualisation of the earth quakes

We will explore some data of rainfall and temperatures from NIWA.

Daily Rainfalls in Christchurch

Automagic downloading of the data by Method B can be done if the data provider allows automated queries. It can be accomplished by `urlread` for instance.

Paul Brouwers has a basic CliFlo datafeed on <http://www.math.canterbury.ac.nz/php/lib/cliflo/rainfall.php>. This returns the date and rainfall in milli meters as measured from the CHCH aeroclub station. It is assumed that days without readings would not be listed. The data doesn't go back much before 1944.

Labwork 120 Understand how Figure 3.34 is obtained by the script file `RainFallsInChch.m` by typing and following the comments:

```
>> RainFallsInChch
RainFallsChch = [24312x1 int32] [24312x1 double]
ans =
24312
FirstDayOfData =
19430802
LastDayOfData =
20100721
```

```
RainFallsInChch.m
%% How to download data from an URL directly without having to manually
%% fill out forms
% first make a string of the data using urlread (read help urlread if you want details)
StringData = urlread('http://www.math.canterbury.ac.nz/php/lib/cliflo/rainfall.php');
RainFallsChch = textscan(StringData, '%d %f', 'delimiter', ',')
RC = [RainFallsChch{1} RainFallsChch{2}]; % assign Matlab cells as a matrix
size(RC) % find the size of the matrix

FirstDayOfData = min(RC(:,1))
LastDayOfData = max(RC(:,1))

plot(RC(:,2),':.')
xlabel('Days in Christchurch, NZ since August 2nd of 1943','FontSize',20);
ylabel('Rainfall in millimeters','FontSize',20)
```

3.8.2 Properties of expectations

The following results, where a is a constant, may easily be proved using the properties of summations and integrals:

$$\mathbf{E}(a) = a$$

$$\mathbf{E}(a g(X)) = a \mathbf{E}(g(X))$$

$$\mathbf{E}(g(X) + h(X)) = \mathbf{E}(g(X)) + \mathbf{E}(h(X))$$

Note that here $g(X)$ and $h(X)$ are functions of the random variable X : e.g. $g(X) = X^2$.

Using these results we can obtain the following useful formula for variance:

$$\begin{aligned} V(X) &= E((X - \mathbf{E}(X))^2) \\ &= E(X^2 - 2X\mathbf{E}(X) + (\mathbf{E}(X))^2) \\ &= E(X^2) - E(2X\mathbf{E}(X)) + E((\mathbf{E}(X))^2) \\ &= E(X^2) - 2\mathbf{E}(X)\mathbf{E}(X) + (\mathbf{E}(X))^2 \\ &= E(X^2) - 2(\mathbf{E}(X))^2 + (\mathbf{E}(X))^2 \\ &= \mathbf{E}(X^2) - (\mathbf{E}(X))^2. \end{aligned}$$

That is,

$$V(X) = \mathbf{E}(X^2) - (\mathbf{E}(X))^2.$$

The above properties of expectations imply that for constants a and b ,

$$\mathbf{V}(aX + b) = a^2\mathbf{V}(X). \quad (3.49)$$

More generally, for random variables X_1, X_2, \dots, X_n and constants a_1, a_2, \dots, a_n

- $\mathbf{E}\left(\sum_{i=1}^n a_i X_i\right) = \sum_{i=1}^n a_i \mathbf{E}(X_i).$ (3.50)
- $\mathbf{V}\left(\sum_{i=1}^n a_i X_i\right) = \sum_{i=1}^n a_i^2 \mathbf{V}(X_i),$ provided X_1, X_2, \dots, X_n are independent. (3.51)

- Let X_1, X_2, \dots, X_n be independent RVs, then

$$\mathbf{E}\left(\prod_{i=1}^n X_i\right) = \prod_{i=1}^n \mathbf{E}(X_i), \text{ provided } X_1, X_2, \dots, X_n \text{ are independent.} \quad (3.52)$$

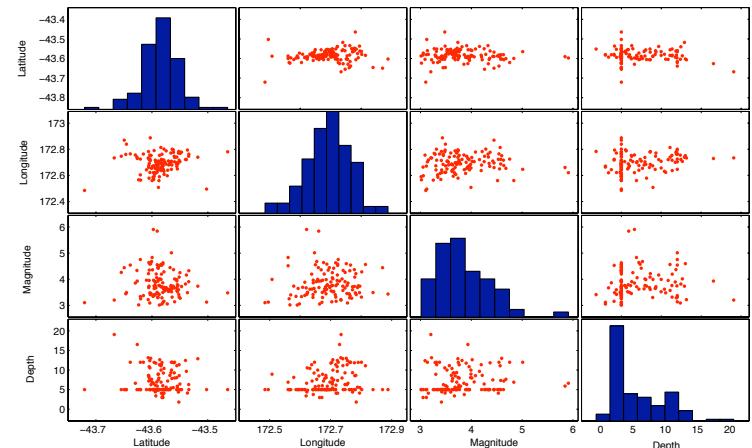
```
>> LatData=EQ(:,2); LonData=EQ(:,3); MagData=EQ(:,12); DepData=EQ(:,13);
>> % finally make a plot matrix of these 124 4-tuples as red points
>> plotmatrix([LatData,LonData,MagData,DepData], 'r.');
```

All of these commands have been put in a script M-file `NZEQChCch20110222.m` and you can simply call it from the command window to automatically load the data and assign it to the variables EQAll EQ, LatData, LonData, MagData and DepData, instead of retyping each command above every time you need these matrices in MATLAB, as follows:

```
>> NZEQChCch20110222
ans = 145 14
ans = 145 14
ans = 124 14
```

In fact, we will do exactly this to conduct more exploratory data analysis with these earth quake measurements in Labwork 119.

Figure 3.32: Matrix of Scatter Plots of the latitude, longitude, magnitude and depth of the 22-02-2011 earth quakes in Christchurch, New Zealand.



Labwork 119 Try to understand how to manipulate time stamps of events in MATLAB and the Figures being output by following the comments in the script file `NZEQChCch20110222EDA.m`.

```
>> NZEQChCch20110222
ans = 145 14
ans = 145 14
ans = 124 14
ans = 145 14
ans = 145 14
ans = 124 14
ans = 22-Feb-2011 00:00:31
ans = 22-Feb-2011 23:50:01
```

• $0 > \zeta - = \left((X)^\theta \Lambda \frac{\partial p}{\partial \theta} \right) \frac{\partial p}{p} =: (X)_\#^\theta \Lambda$ • $\frac{\zeta}{1} = \theta \iff 0 = 1 - \zeta = (X)^\theta \Lambda \frac{\partial p}{p} =: (X)^\theta \Lambda$

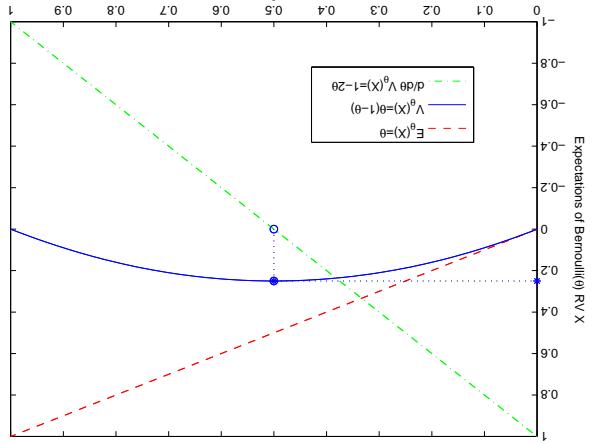


Figure 3.16: Mean ($E_\theta(X)$), variance ($V_\theta(X)$) and the rate of change of variance ($\frac{d\theta}{d\theta} V_\theta(X)$) of a Bernoulli(θ) RV X as a function of the parameter θ .

For a unital \mathbb{A} , $\theta = (X)^\theta \mathbb{A}$ and $\theta = (X)^\theta$.

$$\theta = \theta + 0 = (\theta \times 1) + ((\theta - 1) \times \varepsilon 0) = (x)f_\varepsilon x \sum_{i=0}^{0=x} = (\varepsilon X)\mathbf{E}$$

$$\theta = \theta + 0 = (\theta \times 1) + ((\theta - 1) \times 0) = (x)fx \sum_{i=0}^{0=x} = (X)\mathbf{E}$$

Example 7.4 (Mean and variance of Bernoulli(θ) RV) Let $X \sim \text{Bernoulli}(\theta)$. Then,

Let us compute the mean and variance of our familiar RVs.

3.8.3 Expectation of Common Random Variables

401

```

>>> Eqa11 = [145 14
>>> size(Eqa11)
>>> % report the size of Eqa11 and see if it is different from matrix Eqa11
>>> Eqa11(1,any(isnan(Eqa11),2,:)) = []
>>> % remove any rows containing Nans from the matrix Eqa11
>>> % report the size of Eqa11 and see if it is different from matrix Eqa11
>>> Eqa11 = [145 14
>>> size(Eqa11)
>>> % remove locations outside chch and assign it to a new variable called Eqa
>>> Eqa = Eqa11(:,1:3)*Eqa11(:,2) * Eqa11(:,2)'*Eqa11(:,1:3);
>>> % remove locations outside chch and assign it to a new variable called Eqa
>>> Eqa = Eqa11(:,1:3)*Eqa11(:,2) * Eqa11(:,2)'*Eqa11(:,1:3);
>>> size(Eqa)
>>> % now report the size of the earthquakes in Christchurch in variable Eq
>>> size(Eqa)
>>> % assign the four variables of interest

```

When three are units in the CSV file that can't be converted to floating-point numbers, it is customary to load them as a NaN or Not-a-Number value in MATLAB. So, let's check it is customary to load them as a NaN or Not-a-Number value in MATLAB, it three are any rows with NaN values and remove them from our analysis. Note that this is not the only way to deal with missing data! After that let's remove any locations outside Christchurch and its suburbs (we can find the latitude and longitude bounds from online resources easily) and finally view the 4-tuples of (latitude, longitude, magnitude, depth) for each measured earth quake in Christchurch on February 22 of 2011 as a scatter plot shown in Figure 3-32 (the axes labels were subsequently added from clicking <Edit> and <Figure Properties...> tabs of the output Figure window).

```
DLIMREAD Read ASCII delimited file.  
">>>> help dlimread
```

. In order to understand the syntax in detail get help from MATLAB !

```
% Load the data from the comma delimited text file N20210122earthquakes.csv' with
%% each row representing an earthquake
%% Using MATLAB's dlmread command we can assign the data as a matrix to Eq1.
%% Note that the option 1,0 to dlmread skips first row of column descriptors
%% the variable Eq1 is about to be assigned the data as a matrix
% size(Eq1) % report the dimensions or size of the matrix Eq1
Eq1 = dlmread('N20210122earthquakes.csv', ',', 1, 0);
% the variable Eq11 is about to be assigned the data as a matrix
Eq11 = dlmread('N20210122earthquakes.csv', ',', 1, 0);
% size(Eq11) % report the dimensions or size of the matrix Eq11
ans =
145 14
```

The thirteen columns correspond to nearly self-descriptive features of each measured earth quake given in the first line or row. They will become clear in the sequel. Note that the comma character (‘,’) separates each unit or measurement or description in any CSV file.

The plot depicting these expectations as well as the rate of change of the variance are depicted in Figure 3.16. Note from this Figure that $\mathbf{V}_\theta(X)$ attains its maximum value of $1/4$ at $\theta = 0.5$ where $\frac{d}{d\theta}\mathbf{V}_\theta(X) = 0$. Furthermore, we know that we don't have a minimum at $\theta = 0.5$ since the second derivative $\mathbf{V}''_\theta(X) = -2$ is negative for any $\theta \in [0, 1]$. This confirms that $\mathbf{V}_\theta(X)$ is concave down and therefore we have a maximum of $\mathbf{V}_\theta(X)$ at $\theta = 0.5$. We will revisit this example when we employ a numerical approach called Newton-Raphson method to solve for the maximum of a differentiable function by setting its derivative equal to zero.

Example 75 (Mean and variance of Uniform(0, 1) RV) Let $X \sim \text{Uniform}(0, 1)$. Then,

$$\begin{aligned}\mathbf{E}(X) &= \int_{x=0}^1 x f(x) dx = \int_{x=0}^1 x \cdot 1 dx = \frac{1}{2} (x^2) \Big|_{x=0}^{x=1} = \frac{1}{2} (1 - 0) = \frac{1}{2}, \\ \mathbf{E}(X^2) &= \int_{x=0}^1 x^2 f(x) dx = \int_{x=0}^1 x^2 \cdot 1 dx = \frac{1}{3} (x^3) \Big|_{x=0}^{x=1} = \frac{1}{3} (1 - 0) = \frac{1}{3}, \\ \mathbf{V}(X) &= \mathbf{E}(X^2) - (\mathbf{E}(X))^2 = \frac{1}{3} - \left(\frac{1}{2}\right)^2 = \frac{1}{3} - \frac{1}{4} = \frac{1}{12}.\end{aligned}$$

Exercise 3.29 (Mean and variance of Uniform(θ_1, θ_2) RV) Let $X \sim \text{Uniform}(\theta_1, \theta_2)$ of Model 9. Derive expressions for $\mathbf{E}(X)$ and $\mathbf{V}(X)$ in terms of the parameters θ_1 and θ_2 . Make sure that when $\theta_2 = 1$ and $\theta_1 = 0$ you recover the expectation and variance of the Uniform(0, 1) RV in Example 75.

Example 76 (Expected Exponential of the Uniform(0, 1) RV) Let $X \sim \text{Uniform}(0, 1)$ and $Y = r(X) = e^X$. Compute $\mathbf{E}(Y)$.

We can simply apply the definition of $\mathbf{E}(r(X))$, since $Y = r(X)$, is just a function of X , as follows:

$$\mathbf{E}(Y) = \int_0^1 e^x f(x) dx = \int_0^1 e^x \cdot 1 dx = e - 1.$$

Example 77 (Mean and variance of Exponential(λ) RV) Show that the mean of an Exponential(λ) RV X is:

$$\mathbf{E}_\lambda(X) = \int_0^\infty x f(x; \lambda) dx = \int_0^\infty x \lambda e^{-\lambda x} dx = \frac{1}{\lambda},$$

and the variance is:

$$\mathbf{V}_\lambda(X) = \left(\frac{1}{\lambda}\right)^2.$$

Example 78 (Mean and variance of Geometric(θ) RV) Let $X \sim \text{Geometric}(\theta)$ RV. Then,

$$\mathbf{E}(X) = \sum_{x=0}^{\infty} x \theta (1 - \theta)^x = \theta \sum_{x=0}^{\infty} x (1 - \theta)^x$$

In order to simplify the RHS above, let us employ differentiation with respect to θ :

$$\frac{-1}{\theta^2} = \frac{d}{d\theta} \left(\frac{1}{\theta} \right) = \frac{d}{d\theta} \sum_{x=0}^{\infty} (1 - \theta)^x = \sum_{x=0}^{\infty} -x (1 - \theta)^{x-1}$$

Multiplying the LHS and RHS above by $-(1 - \theta)$ and substituting in $\mathbf{E}(X) = \theta \sum_{x=0}^{\infty} x (1 - \theta)^x$, we get a much simpler expression for $\mathbf{E}(X)$:

$$\frac{1 - \theta}{\theta^2} = \sum_{x=0}^{\infty} x (1 - \theta)^x \implies \mathbf{E}(X) = \theta \left(\frac{1 - \theta}{\theta^2} \right) = \frac{1 - \theta}{\theta}.$$

Labwork 117 Let us make matrix plots from a uniformly generated sequence of 100 points in 5D unit cube $[0, 1]^5$ as shown in Figure 3.31.

```
>> rand('twister',5489);
>> % generate a sequence of 1000 points uniformly distributed in 5D unit cube [0,1]x[0,1]x[0,1]x[0,1]x[0,1]
>> x=rand(1000,5);
>> x(1:6,:) % first six points in our 5D unit cube, i.e., the first six rows of x
ans =
    0.8147    0.6312    0.7449    0.3796    0.4271
    0.9058    0.3551    0.8923    0.3191    0.9554
    0.1270    0.9970    0.2426    0.9861    0.7242
    0.9134    0.2242    0.1296    0.7182    0.5809
    0.6324    0.6525    0.2251    0.4132    0.5403
    0.0975    0.6050    0.3500    0.0986    0.7054
>> plotmatrix(x(1:5,:),'r*') % make a plot matrix
>> plotmatrix(x) % make a plot matrix of all 1000 points
```

3.14.6 Loading and Exploring Real-world Data

All of the data we have played with so far were computer-generated. It is time to get our hands dirty with real-world data. The first step is to obtain the data. Often, publicly-funded institutions allow the public to access their databases. Such data can be fetched from appropriate URLs in one of the two following ways:

Method A: Manually download by filling the appropriate fields in an online request form.

Method B: Automagically download directly from your MATLAB session.

Then we want to inspect it for inconsistencies, missing values and replace them with NaN values in MATLAB that stand for not-any-number. Finally, we can visually explore, transform and interact with the data to discover interesting patterns that are hidden in the data. This process is called *exploratory data analysis* and is the foundational first step towards subsequent computational statistical experiments [John W. Tukey, *Exploratory Data Analysis*, Addison-Wesley, New York, 1977].

3.14.7 Geological Data

Let us focus on the data of earth quakes that heavily damaged Christchurch on February 22 2011. This data can be fetched from the URL <http://magma.geonet.org.nz/resources/quakesearch/> by Method A and loaded into MATLAB for exploratory data analysis as done in Labwork 118.

Labwork 118 Let us go through the process one step at a time using Method A.

- Download the data as a CSV or *comma separated variable* file in plain ASCII text (this has been done for this data already for you and saved as `NZ20110222earthquakes.csv` in the `CSEMatlabScripts` directory).
- Open the file in a simple text editor such as `Note Pad` in Windows or one of the following editors in OS X, Unix, Solaris, Linux/GNU variants such as Ubuntu, SUSE, etc: `vi`, `vim`, `emacs`, `geany`, etc. The first three and last two lines of this file look as follows:

$$\Lambda(X) = \mathbb{E}(X^2) - (\mathbb{E}(X))^2 = \mathbb{E}(X^2) - \mathbb{E}^2(X) = \mathbb{E}(X^2) - \mathbb{E}^2(X) = \mathbb{V}(X).$$

Similarly,

$$\mathbb{E}(X) = \sum_{x=0}^{\infty} x f(x; \lambda) = e^{-\lambda} \lambda \sum_{x=0}^{\infty} x \frac{x!}{\lambda^x x!} = e^{-\lambda} \lambda \sum_{x=0}^{\infty} x \frac{x!}{\lambda^x x!} = e^{-\lambda} \lambda = \lambda.$$

Example 80 (Mean and variance of Poisson(λ) RV) Let $X \sim \text{Poisson}(\lambda)$. Then:

$$\mathbb{E}(X) = \mathbb{E}\left(\sum_{i=1}^n X_i\right) = n \mathbb{E}(X_i) = n \theta.$$

$$\mathbb{E}(X) = \mathbb{E}\left(\sum_{i=1}^n X_i\right) = \mathbb{E}\left(\sum_{i=1}^n \mathbb{E}(X_i)\right) = n \theta.$$

However, this is a nontrivial sum to evaluate. Instead, we may use (3.50) and (3.51) by noting that $X = \sum_{i=1}^n X_i$, where the $\{X_1, X_2, \dots, X_n\} \stackrel{\text{iid}}{\sim} \text{Bernoulli}(\theta)$, $\mathbb{E}(X_i) = \theta$ and $\mathbb{V}(X_i) = \theta(1-\theta)$:

$$\mathbb{E}(X) = \int x dF(x; n, \theta) = \sum_{x=0}^n x \theta^x (1-\theta)^{n-x}.$$

the definition of expectation:

Example 79 (Mean and variance of Binomial(n, θ) RV) Let $X \sim \text{Binomial}(n, \theta)$. Based on

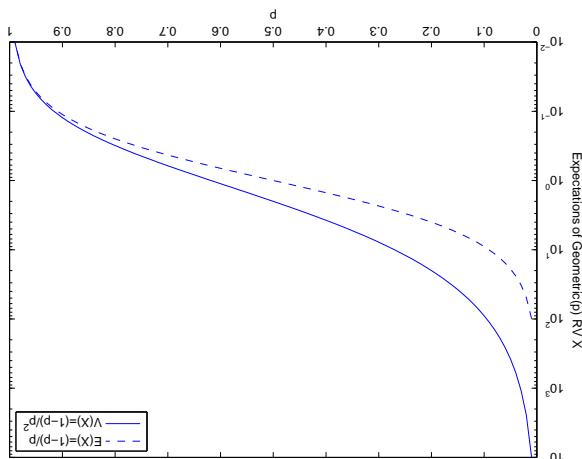


Figure 3.17: Mean and variance of a Geometric(θ) RV X as a function of the parameter θ .

$$\Lambda(X) = \frac{\theta^2}{1-\theta}.$$

Similarly, it can be shown that

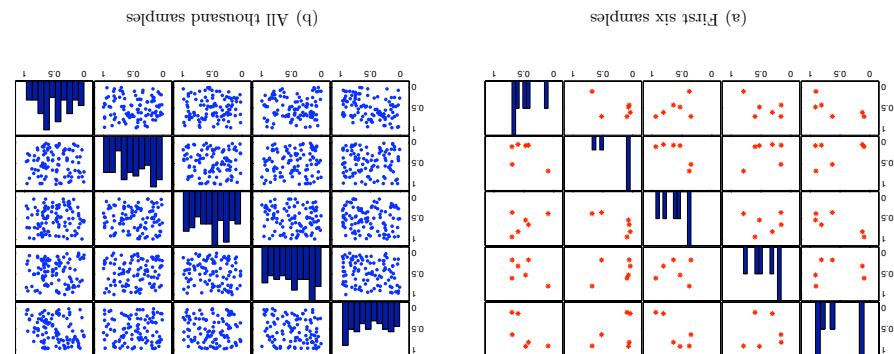


Figure 3.31: Plot Matrix of uniformly generated data in $[0, 1]^5$

For high-dimensional data in d -dimensional space \mathbb{R}^d with $d \geq 3$ you have to look at several lower dimensional projections of the data. We can simultaneously look at 2D scatter plots for every pair of co-ordinates $(i, j) \in \{1, 2, \dots, d\} : i \neq j\}$ and at histograms for every co-ordinate $i \in \{1, 2, \dots, d\}$ of the n data points in \mathbb{R}^d . Such a set of low-dimensional projections can be conveniently represented in a $d \times d$ matrix of plots called a **matrix plot**.

3.14.5 Multivariate Data

There are several other techniques for visualizing trivariate data, including, iso-surface plots, moving surface or heat plots, and you will encounter some of them in the future.

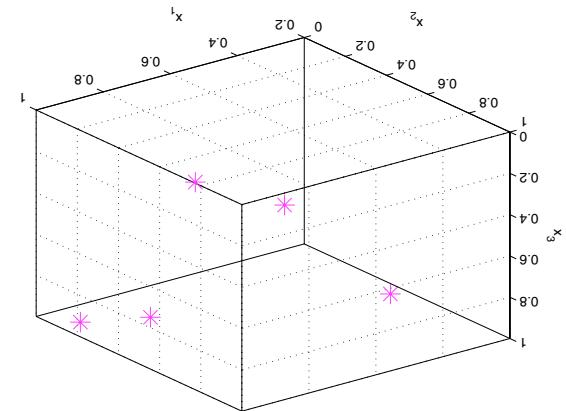


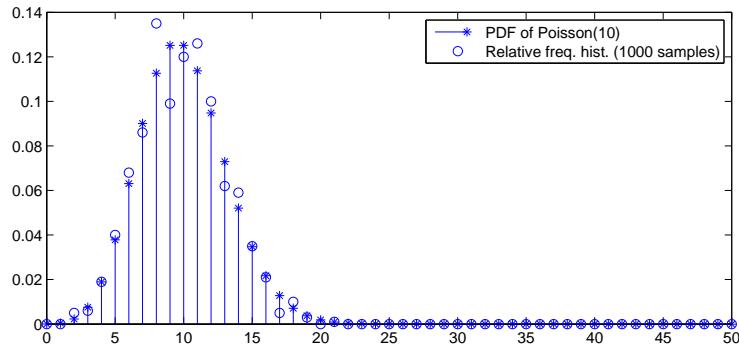
Figure 3.30: 3D Scatter Plot

since

$$\begin{aligned}\mathbf{E}(X^2) &= \sum_{x=0}^{\infty} x^2 \frac{e^{-\lambda} \lambda^x}{x!} = \lambda e^{-\lambda} \sum_{x=1}^{\infty} \frac{x \lambda^{x-1}}{(x-1)!} = \lambda e^{-\lambda} \left(1 + \frac{2\lambda}{1} + \frac{3\lambda^2}{2!} + \frac{4\lambda^3}{3!} + \dots \right) \\ &= \lambda e^{-\lambda} \left(\left(1 + \frac{\lambda}{1} + \frac{\lambda^2}{2!} + \frac{\lambda^3}{3!} + \dots \right) + \left[\frac{\lambda}{1} + \frac{2\lambda^2}{2!} + \frac{3\lambda^3}{3!} + \dots \right] \right) \\ &= \lambda e^{-\lambda} \left((e^\lambda) + \lambda \left(1 + \frac{2\lambda}{2!} + \frac{3\lambda^2}{3!} + \dots \right) \right) = \lambda e^{-\lambda} \left(e^\lambda + \lambda \left(1 + \lambda + \frac{\lambda^2}{2!} + \dots \right) \right) \\ &= \lambda e^{-\lambda} (e^\lambda + \lambda e^\lambda) = \lambda e^{-\lambda} (e^\lambda + \lambda e^\lambda) = \lambda(1 + \lambda) = \lambda + \lambda^2\end{aligned}$$

Note that Poisson(λ) distribution is one whose mean and variance are the same, namely λ .

Figure 3.18: PDF of $X \sim \text{Poisson}(\lambda = 10)$ and the relative frequency histogram based on 1000 samples from X according to Simulation 149.



The Poisson(λ) RV X is also related to the IID Exponential(λ) RV Y_1, Y_2, \dots : X is the number of occurrences, per unit time, of an instantaneous event whose inter-occurrence time is the IID Exponential(λ) RV. For example, the number of buses arriving at our bus-stop in the next minute, with exponentially distributed inter-arrival times, has a Poisson distribution.

Example 81 (Mean and variance of Normal(μ, σ^2) RV) The location-scale family of RVs is indeed parameterised by its mean and variance, i.e., if $X \sim \text{Normal}(\mu, \sigma^2)$ where $X = g(Z) = \sigma Z + \mu$ and $Z \sim \text{Normal}(0, 1)$ then $\mathbf{E}(X) = \mu$ and $\mathbf{V}(X) = \sigma^2$ follows directly from the properties of Expectations, provided $\mathbf{E}(Z) = 0$ and $\mathbf{V}(Z) = \mathbf{E}(Z^2) - (\mathbf{E}(Z))^2 = \mathbf{E}(Z^2) = 1$.

The mean of a Normal(0, 1) RV Z is:

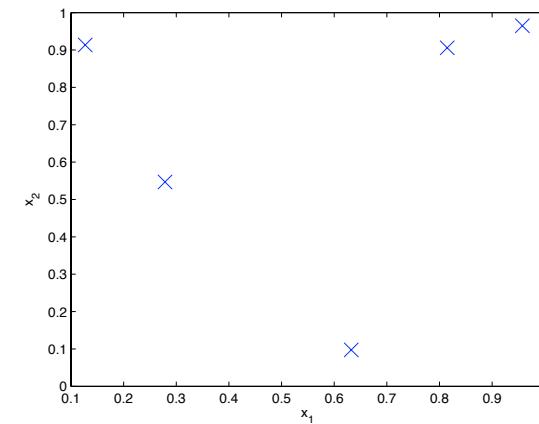
$$\mathbf{E}(Z) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} z \exp\left(-\frac{1}{2}z^2\right) dz = \frac{1}{\sqrt{2\pi}} \left[-\exp\left(-\frac{1}{2}z^2\right) \right]_{-\infty}^{\infty} = 0,$$

and the variance is:

$$\mathbf{V}(Z) = \mathbf{E}(Z^2) - (\mathbf{E}(Z))^2 = \mathbf{E}(Z^2) - 0 = \mathbf{E}(Z^2) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} z^2 e^{-z^2/2} dz.$$

```
0.9058 0.9134 0.0975 0.5469 0.9649
>> plot(x(1,:),x(2,:),'x') % a 2D scatter plot with marker cross or 'x'
>> plot(x(1,:),x(2,:),'x', 'MarkerSize',15) % a 2D scatter plot with marker cross or 'x' and larger Marker size
>> xlabel('x_1'); ylabel('x_2'); % label the axes
```

Figure 3.29: 2D Scatter Plot



There are several other techniques for visualising bivariate data, including, 2D histograms, surface plots, heat plots, and we will encounter some of them in the sequel.

3.14.4 Trivariate Data

Trivariate data is more difficult to visualise on paper but playing around with the rotate 3D feature in MATLAB's Figure window can help bring a lot more perspective.

Labwork 116 (Visualising trivariate data) We can make **3D scatter plots** as shown in Figure 3.30 as follows:

```
>> rand('twister',5489);
>> x=rand(3,5)% create a sequence of 5 ordered triples uniformly from unit cube [0,1]X[0,1]X[0,1]
x =
    0.8147 0.9134 0.2785 0.9649 0.9572
    0.9058 0.6324 0.5469 0.1576 0.4854
    0.1270 0.0975 0.9575 0.9706 0.8003
>> plot3(x(1,:),x(2,:),x(3,:),'x') % a simple 3D scatter plot with marker 'x'
>>% a more interesting one with options that control marker type, line-style,
>>% colour in [Red Green Blue] values and marker size - read help plot3 for more options
>> plot3(x(1,:),x(2,:),x(3,:),'Marker','*', 'LineStyle','none','Color',[1 0 1],'MarkerSize',15)
>> plot3(x(1,:),x(2,:),x(3,:),'m*','MarkerSize',15) % makes same figure as before but shorter to write
>> box on % turn on the box and see the effect on the Figure
>> grid on % turn on the grid and see the effect on the Figure
>> xlabel('x_1'); ylabel('x_2'); zlabel('x_3'); % assign labels to x,y and z axes
```

Repeat the visualisation below with a larger array, say $x=\text{rand}(3,1000)$, and use the rotate 3D feature in the Figure window to visually explore the samples in the unit cube. Do they seem to be uniformly distributed inside the unit cube?

$$f(x; \theta_1, \theta_2, \dots, \theta_k) = \begin{cases} \theta & \text{if } x \in [k], \\ 0 & \text{if } x \notin [k], \end{cases}$$

we say that an RV X is de Moivre($\theta_1, \theta_2, \dots, \theta_k$) distributed if its PMF is:

$$\nabla_{k-1} := \{(\theta_1, \theta_2, \dots, \theta_k) : \theta_1 \geq 0, \theta_2 \geq 0, \dots, \theta_k \geq 0, \sum_{i=1}^k \theta_i = 1\},$$

Simplex:

Model 14 (de Moivre($\theta_1, \theta_2, \dots, \theta_k$)) Given a specific point $(\theta_1, \theta_2, \dots, \theta_k)$ in the unit $k-1$ -

simplex. Variable and higher moments cannot be defined when the expectation is defined over $(0, \infty)$

Note that we consider symmetry of integrals about the origin and take twice the integral over $(0, \infty)$ above. Variable and higher moments cannot be defined when the expectation is defined over $(0, \infty)$

Example 82 (Mean of Cauchy RV) The expectation of the Cauchy RV X , obtained via integration by parts (set $u = x$ and $v = \tan^{-1}(x)$) does not exist, since:

Note that the construction is valid even if we sample X uniformly from $(0, \infty)$ and take its $\tan(X)$.

$$f_Y(y) = f_X(g_{-1}(y)) \left| \frac{dy}{dx} \right| = f_X(\tan^{-1}(y)) \left| \frac{dy}{d \tan^{-1}(y)} \right| = \frac{1}{1 + y^2}$$

PDF $f_Y(y)$ from the PDF $f_X(x) = \frac{1}{\pi(1+x^2)}$ as follows: X given by $(-\pi/2, \pi/2)$, we can use the change of variable formula in Equation 3.38 to obtain the $Y = \tan(X)$ for the above construction. Since $\tan(x)$ is one-to-one and monotone on the range of X randomly spinning a LASER emitting impurity of Darth Mall's double edged lightsaber

The Cauchy RV Y can be derived from a RV $X \sim \text{Uniform}(-\pi/2, \pi/2)$ by the simple transformation that is centred at $(1, 0)$ in the plane \mathbb{R}^2 and recording its intersection with the y -axis, in terms of the y coordinates of the point $(0, y)$, gives rise to the Standard Cauchy RV.

Randomly spinning a LASER emitting impurity of Darth Mall's double edged lightsaber

$$F(y) = \frac{1}{\pi} \tan^{-1}(y) + \frac{1}{2}.$$

and its DF is:

$$f(y) = \frac{\pi(1+y^2)^{-1}}{1}, \quad -\infty < y < \infty,$$

Model 13 (Cauchy) The density of the Cauchy RV Y is:

Next, let us become familiar with an RV for which the expectation does not exist.

The first term after the first equality above equals 0 because the exponential goes to 0 much faster than z grows to ∞ . The second term equals 1 because it is exactly the total probability integrated

of the PDF of the Normal(0, 1) RV.

$$\int_{-\infty}^{\infty} z^2 e^{-z^2/2} dz = \frac{\sqrt{2\pi}}{1} \left(-ze^{-z^2/2} \Big|_{-\infty}^{\infty} + \int_{-\infty}^{\infty} e^{-z^2/2} dz \right) = 0 + 1 = 1$$

Using integration by parts with $u = z$, $du = ze^{-z^2/2} \iff du = 1$, $v = -e^{-z^2/2}$, $\int v du = ue^{-z^2/2}$

CHAPTER 3. RANDOM VARIABLES

```
0.847 0.1270 0.6324 0.2785 0.9575
x =
>>> x=rand(2,5) % create a sequence of 5 ordered pairs uniformly from unit square [0,1]x[0,1]
>>> rand('twister',5489);
```

2D scatter plot as shown in Figure 3.29 as follows:

of 5 ordered pairs sampled uniformly at random over the unit square $[0, 1] \times [0, 1]$. We can make 2D scatter plots of real numbers or equivalently n ordered pairs in statistics.

By bivariate data array x we mean a $2 \times n$ matrix of real numbers of equivalent n ordered pairs in the form $(x_{1,i}, x_{2,i})$ as $i = 1, 2, \dots, n$. The most elementary visualization of these n ordered pairs is in orthogonal Cartesian co-ordinates. Such plots are termed 2D scatter plots in statistics.

Labwork 115 (Visualising bivariate data) Let us generate a 2×5 array representing samples in the Stats Toolbox of MATLAB. These family of plots display a set of sample quantiles, typically in the Statistics Toolbox of MATLAB. They are median, the first and third quartiles and the minimum and maximum values they are in include, the median, the first and third quartiles and the minimum and maximum values in the Statistics Toolbox of MATLAB. We can also visually summarise univariate data using the box-whisker plot available in the Statistics Toolbox of MATLAB. The box plot displays the median, the first and third quartiles, typically in the Statistics Toolbox of MATLAB. We can also visually summarise univariate data using the box plot or box-whisker plot available in the Statistics Toolbox of MATLAB. These family of plots display a set of sample quantiles, typically in the Statistics Toolbox of MATLAB. They are median, the first and third quartiles and the minimum and maximum values in the Statistics Toolbox of MATLAB. By bivariate data array x we mean a $2 \times n$ matrix of real numbers or equivalently n ordered pairs in the Statistics Toolbox of MATLAB. Such plots are termed 2D scatter plots in statistics.

3.14.3 Bivariate Data

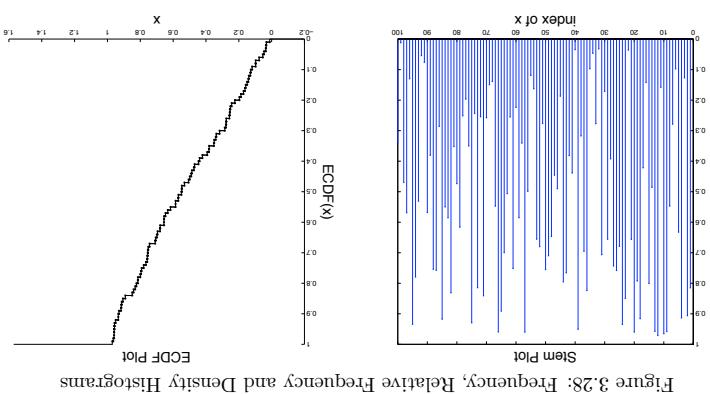


Figure 3.28: Frequency, Relative Frequency and Density Histograms

```
>>> x=rand(1,100); % produce 100 samples with random
>>> stem(x,6); % make a stem plot of the 100 data points in x (the option ' ', gives solid circles for x)
>>> ecdf(x,6); % make a step ECDF plot in x (the option ' ', gives solid circles for x)
>>> % (second parameter 6 makes the dots in the plot smaller).
>>> % (step ECDF) plot is extended to left and right by .2 and .6, respectively
>>> % (stem(x,6), ECDF(x,6), hist(x,6))
```

100 data points in the array x using stem plot and ECDF plot as shown in Figure 3.28 as follows:

We can also visualise the 100 data points in the array x using stem plot and ECDF plot as shown in Figure 3.28 as follows:

Try making a density histogram with 1000 samples from with 15 bins. You can specify the number of bins by adding an extra argument to hist, for e.g., $[fs, cs] = \text{hist}(x, 15)$ will produce 15 bins of equal width over the data range $R(x)$.

CHAPTER 3. RANDOM VARIABLES

The DF for de Moivre($\theta_1, \theta_2, \dots, \theta_k$) RV X is:

$$F(x; \theta_1, \theta_2, \dots, \theta_k) = \begin{cases} 0 & \text{if } -\infty < x < 1 \\ \theta_1 & \text{if } 1 \leq x < 2 \\ \theta_1 + \theta_2 & \text{if } 2 \leq x < 3 \\ \vdots & \\ \theta_1 + \theta_2 + \dots + \theta_{k-1} & \text{if } k-1 \leq x < k \\ \theta_1 + \theta_2 + \dots + \theta_{k-1} + \theta_k = 1 & \text{if } k \leq x < \infty \end{cases} \quad (3.56)$$

The de Moivre($\theta_1, \theta_2, \dots, \theta_k$) RV can be thought of as a probability model for “the outcome of rolling a polygonal cylindrical die with k rectangular faces that are marked with $1, 2, \dots, k$ ”. The parameters $\theta_1, \theta_2, \dots, \theta_k$ specify how the die is loaded and may be idealised as specifying the cylinder’s centre of mass with respect to the respective faces. Thus, when $\theta_1 = \theta_2 = \dots = \theta_k = 1/k$, we have a probability model for the outcomes of a fair die.

Mean and variance of de Moivre($\theta_1, \theta_2, \dots, \theta_k$) RV: The not too useful expressions for the first two moments of $X \sim \text{de Moivre}(\theta_1, \theta_2, \dots, \theta_k)$ are,

$$\mathbf{E}(X) = \sum_{x=1}^k x\theta(x) = \theta_1 + 2\theta_2 + \dots + k\theta_k, \text{ and}$$

$$\mathbf{V}(X) = \mathbf{E}(X^2) - (\mathbf{E}(X))^2 = (\theta_1 + 2^2\theta_2 + \dots + k^2\theta_k) - (\theta_1 + 2\theta_2 + \dots + k\theta_k)^2.$$

However, if $X \sim \text{de Moivre}(1/k, 1/k, \dots, 1/k)$, then the mean and variance for the fair k -faced die based on Faulhaber’s formula for $\sum_{i=1}^k i^m$, with $m \in \{1, 2\}$, are,

$$\mathbf{E}(X) = \frac{1}{k} (1 + 2 + \dots + k) = \frac{1}{k} \frac{k(k+1)}{2} = \frac{k+1}{2},$$

$$\mathbf{E}(X^2) = \frac{1}{k} (1^2 + 2^2 + \dots + k^2) = \frac{1}{k} \frac{k(k+1)(2k+1)}{6} = \frac{2k^2 + 3k + 1}{6},$$

$$\begin{aligned} \mathbf{V}(X) &= \mathbf{E}(X^2) - (\mathbf{E}(X))^2 = \frac{2k^2 + 3k + 1}{6} - \left(\frac{k+1}{2}\right)^2 = \frac{2k^2 + 3k + 1}{6} - \left(\frac{k^2 + 2k + 1}{4}\right) \\ &= \frac{8k^2 + 12k + 4 - 6k^2 - 12k - 6}{24} = \frac{2k^2 - 2}{24} = \frac{k^2 - 1}{12}. \end{aligned}$$

3.9 Exercises in Expectations of Random Variables

Ex. 3.30 — Let X be the number of air conditioners a store sells each day, and assume that X has probability mass function $f(10) = 0.1$, $f(11) = 0.3$, $f(12) = 0.4$, $f(13) = 0.2$.

1. Find the expected number of conditioners that the store sells each day.
2. If the profit per conditioner is \$55, what is the expected daily profit?

Ex. 3.31 — A small petrol station is supplied with fuel every Saturday afternoon. Assume that its volume of sales X , in ten thousands of litres, has density

$$f(x) = \begin{cases} 6x(1-x) & 0 \leq x \leq 1 \\ 0 & \text{otherwise} \end{cases}.$$

Determine the mean and variance of X .

For a given partition of the data range $\mathcal{R}(x)$ or some superset of $\mathcal{R}(x)$, three types of histograms are possible: frequency histogram, relative frequency histogram and density histogram. Typically, the partition b is assumed to be composed of m overlapping intervals of the same width $w = \bar{b}_i - \underline{b}_i$ for all $i = 1, 2, \dots, m$. Thus, a histogram can be obtained by a set of bins along with their corresponding heights:

$$h = (h_1, h_2, \dots, h_m), \text{ where } h_k := g(\#\{x_i : x_i \in b_k\})$$

Thus, h_k , the height of the k -th bin, is some function g of the number of data points that fall in the bin b_k . Formally, a histogram is a sequence of ordered pairs:

$$((b_1, h_1), (b_2, h_2), \dots, (b_m, h_m)).$$

Given a partition b , a **frequency histogram** is the histogram:

$$((b_1, h_1), (b_2, h_2), \dots, (b_m, h_m)), \text{ where } h_k := \#\{x_i : x_i \in b_k\},$$

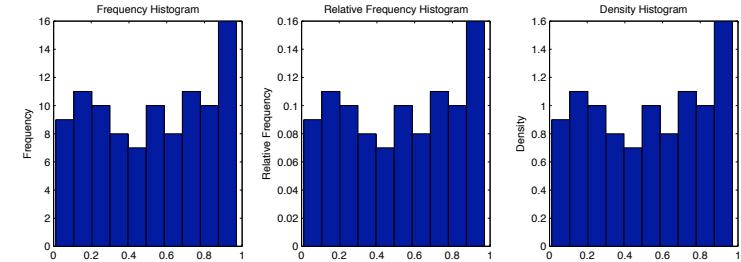
a **relative frequency histogram** is the histogram:

$$((b_1, h_1), (b_2, h_2), \dots, (b_m, h_m)), \text{ where } h_k := n^{-1} \#\{x_i : x_i \in b_k\},$$

and a **density histogram** is the histogram:

$$((b_1, h_1), (b_2, h_2), \dots, (b_m, h_m)), \text{ where } h_k := (w_k n)^{-1} \#\{x_i : x_i \in b_k\}, w_k := \bar{b}_k - \underline{b}_k.$$

Figure 3.27: Frequency, Relative Frequency and Density Histograms



Labwork 113 (Histograms with specified number of bins for univariate data) Let us use samples from the `rand('twister', 5489)` as our data set x and plot various histograms. Let us use `hist` function (read `help hist`) to make a default histogram with ten bins. Then we can make three types of histograms as shown in Figure 3.27 as follows:

```
>> rand('twister', 5489);
>> x=rand(1,100); % generate 100 PRNs
>> hist(x) % see what default hist does in Figure Window
>> % Now let us look deeper into the last hist call
>> [Fs, Cs] = hist(x) % Cs is the bin centers and Fs is the frequencies of data set x
Fs =
    9     11     10     8      7     10     8     11     10     16
Cs =
    0.0598    0.1557    0.2516    0.3474    0.4433    0.5392    0.6351    0.7309    0.8268    0.9227
>> % produce a histogram plot the last argument 1 is the width value for immediately adjacent bars -- help bar
>> bar(Cs,Fs,1) % create a frequency histogram
>> bar(Cs,Fs/100,1) % create a relative frequency histogram
>> bar(Cs,Fs/(0.1*100),1) % create a density histogram (area of bars sum to 1)
>> sum(Fs/(0.1*100)) .* ones(1,10)*0.1 % checking if area does sum to 1
>> ans = 1
```

Ex. 3.32 — Starting from the definition of the variance of a random variable (Definition 30) show that

$$\mathbf{A}(X) = \mathbf{E}(X^2) - (\mathbf{E}(X))^2.$$

Ex. 3.33 — Show that $V(aX + b) = a^2V(X)$ for constants a and b and a random variable X .

Ex. 3.34 — Let X be a discrete random variable with PMF given by

$$f(x) = \begin{cases} \frac{1}{10} & \text{if } x \in \{1, 2, 3, 4\}, \\ 0 & \text{otherwise.} \end{cases}$$

Ex. 3.35 — Find the mean and the variance of the following random variables.

(a) Write down the DF (or CDF) of X .
 (b) Plot the PMF and CDF of X .
 (c) Plot the PMF and CDF of X .

$$\begin{aligned} & \text{(i) } \mathbf{P}(X = 0) \\ & \text{(ii) } \mathbf{P}(2.5 > X > 5) \\ & \text{(iii) } \mathbf{E}(X) \\ & \text{(iv) } \mathbf{V}(X) \end{aligned}$$

2. X is a discrete uniform random variable on $\{1, 2, 3, 4, 5, 6\}$, i.e., the number a fair die turns up. Of course, the experiments we are measuring two or more aspects simultaneously. For example, we may be measuring the diameters and lengths of cylindrical shafts manufactured in a plant or heights, weights and blood-sugar levels of individuals in a clinical trial. Thus, the underlying outcome Ω needs to be mapped to measurements as realizations of random vectors in the real plane $\mathbb{R}^2 = (-\infty, \infty) \times (-\infty, \infty)$ or the real space $\mathbb{R}^3 = (-\infty, \infty) \times (-\infty, \infty) \times (-\infty, \infty)$:

More generally, we may be interested in heights, weights, blood-sugar levels, family medical history, known allergies, etc. of individuals in the clinical trial and this need to make m measurements of variables (X_1, X_2, \dots, X_m) , ordered triples of random variables (X, Y, Z) , or more generally ordered m -tuples of random variables (X, Y, Z) , ordered triples of random variables (X, Y, Z) , or more generally ordered m -tuples of random variables (X_1, X_2, \dots, X_m) . These measurements we need the notion of **random vectors** (RVs), i.e., ordered pairs of random variables the outcome in \mathbb{R}^m using a “measurable mapping” from $\Omega \rightarrow \mathbb{R}^m$. To deal with such multivariate known allergies, etc. of individuals in the clinical trial and this need to make m measurements of More generally, we may be interested in heights, weights, blood-sugar levels, family medical history,

$$\omega \mapsto (X(\omega), Y(\omega)) : \Omega \rightarrow \mathbb{R}^2 \quad \omega \mapsto (X(\omega), Y(\omega), Z(\omega)) : \Omega \rightarrow \mathbb{R}^3$$

Elements of this partition \mathfrak{b} are called bins, their mid-points are called **bin centres**:

of the **data range** of x given by the closed interval:

$R(x) = [\min\{x_1, x_2, \dots, x_n\}, \max\{x_1, x_2, \dots, x_n\}]$

and their overlapping boundaries, i.e., $\underline{b}_i = \overline{b}_{i+1}$ for $1 \leq i < m$, are called **bin edges**:

$c = (c_1, c_2, \dots, c_m) = ((\overline{b}_1 + \underline{b}_2)/2, (\overline{b}_2 + \underline{b}_3)/2, \dots, (\overline{b}_m + \underline{b}_{m+1})/2)$

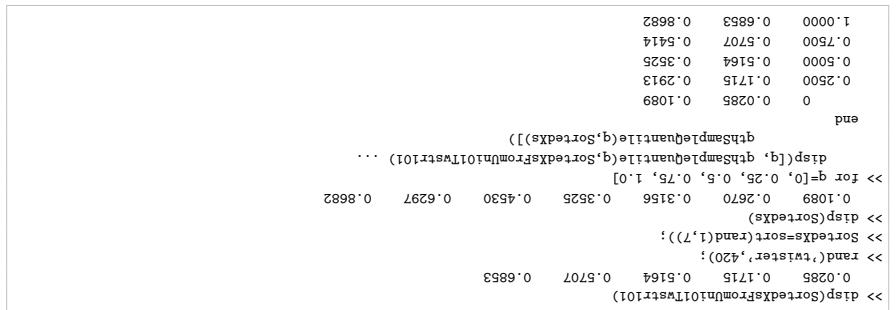
$$d = (d_1, d_2, \dots, d_{m+1}) = (\underline{b}_1, \underline{b}_2, \dots, \underline{b}_{m-1}, \underline{b}_m, \underline{b}_m).$$

of real numbers fall within each of the m intervals or **bins** of some **interval partition**:

$$x = (x_1, x_2, \dots, x_n),$$

A **histogram** is a graphical representation of the frequency with which elements of a data array:

3.14.2 Univariate Data



1: *input:* $F_{[-1]}^n(q)$, the q th sample quantile
 2. *order statistic* $(x(1), x(2), \dots, x(n))$, i.e. the sorted (x_1, x_2, \dots, x_n) , where $n > 0$.
 1. q in the q th sample quantile, i.e. the argument q of $F_{[-1]}^n(q)$,

Algorithm 1 q th Sample Quantile of Order Statistics

Algorithm 1 q th Sample Quantile of Order Statistics

3.10.1 \mathbb{R}^2 -valued Random Variables

We first focus on understanding (X, Y) , a bivariate R \vec{V} or \mathbb{R}^2 -valued RV that is obtained from a pair of discrete or continuous RVs. We then generalize to \mathbb{R}^m -valued RVs with $m > 2$ in the next section.

Definition 34 (JDF) The joint distribution function (JDF) or joint cumulative distribution function (JCDF), $F_{X,Y}(x, y) : \mathbb{R}^2 \rightarrow [0, 1]$, of the bivariate random vector (X, Y) is

$$\begin{aligned} F_{X,Y}(x, y) &= \mathbf{P}(X \leq x \cap Y \leq y) = \mathbf{P}(X \leq x, Y \leq y) \\ &= P(\{\omega : X(\omega) \leq x, Y(\omega) \leq y\}), \text{ for any } (x, y) \in \mathbb{R}^2, \end{aligned} \quad (3.57)$$

where the right-hand side represents the probability that the random vector (X, Y) takes on a value in $\{(x', y') : x' \leq x, y' \leq y\}$, the set of points in the plane that are south-west of the point (x, y) .

The JDF $F_{X,Y}(x, y) : \mathbb{R}^2 \rightarrow \mathbb{R}$ satisfies the following conditions to remain a probability:

1. $0 \leq F_{X,Y}(x, y) \leq 1$
2. $F_{X,Y}(x, y)$ is a non-decreasing function of both x and y
3. $F_{X,Y}(x, y) \rightarrow 1$ as $x \rightarrow \infty$ and $y \rightarrow \infty$
4. $F_{X,Y}(x, y) \rightarrow 0$ as $x \rightarrow -\infty$ and $y \rightarrow -\infty$

Definition 35 (JPMF) If (X, Y) is a discrete random vector that takes values in a discrete support set $\mathcal{S}_{X,Y} = \{(x_i, y_j) : i = 1, 2, \dots, j = 1, 2, \dots\} \subset \mathbb{R}^2$ with probabilities $p_{i,j} = \mathbf{P}(X = x_i, Y = y_j) > 0$, then its joint probability mass function (or JPMF) is:

$$f_{X,Y}(x_i, y_j) = \mathbf{P}(X = x_i, Y = y_j) = \begin{cases} p_{i,j} & \text{if } (x_i, y_j) \in \mathcal{S}_{X,Y} \\ 0 & \text{otherwise} \end{cases}. \quad (3.58)$$

Since $\mathbf{P}(\Omega) = 1$, $\sum_{(x_i, y_j) \in \mathcal{S}_{X,Y}} f_{X,Y}(x_i, y_j) = 1$.

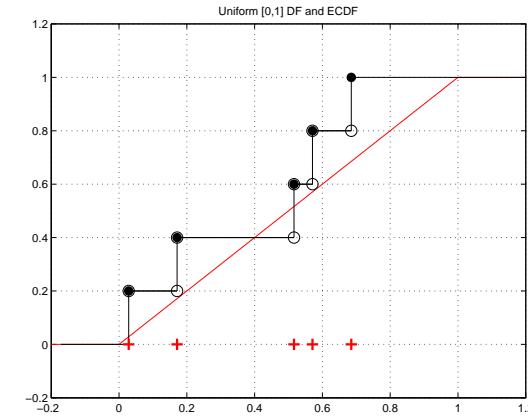
From JPMF $f_{X,Y}$ we can get the values of the JDF $F_{X,Y}(x, y)$ and the probability of any event B by simply taking sums,

$$F_{X,Y}(x, y) = \sum_{x_i \leq x, y_j \leq y} f_{X,Y}(x_i, y_j), \quad \mathbf{P}(B) = \sum_{(x_i, y_j) \in B \cap \mathcal{S}_{X,Y}} f_{X,Y}(x_i, y_j), \quad (3.59)$$

Example 83 Let (X, Y) be a discrete bivariate R \vec{V} with the following joint probability mass function (JPMF):

$$f_{X,Y}(x, y) := P(X = x, Y = y) = \begin{cases} 0.1 & \text{if } (x, y) = (0, 0) \\ 0.3 & \text{if } (x, y) = (0, 1) \\ 0.2 & \text{if } (x, y) = (1, 0) \\ 0.4 & \text{if } (x, y) = (1, 1) \\ 0.0 & \text{otherwise.} \end{cases}$$

Figure 3.26: Plot of the DF of Uniform(0, 1), five IID samples from it, and the ECDF \hat{F}_5 for these five data points $x = (x_1, x_2, x_3, x_4, x_5) = (0.5164, 0.5707, 0.0285, 0.1715, 0.6853)$ that jumps by $1/5 = 0.20$ at each of the five samples.



Definition 57 (q^{th} Sample Quantile) For some $q \in [0, 1]$ and n IID RVs $X_1, X_2, \dots, X_n \stackrel{\text{IID}}{\sim} F$, we can obtain the ECDF \hat{F}_n using (3.82). The q^{th} sample quantile is defined as the statistic (statistical functional):

$$T(\hat{F}_n) = \hat{F}_n^{[-1]}(q) := \inf \{x : \hat{F}_n^{[-1]}(x) \geq q\}. \quad (3.83)$$

By replacing q in this definition of the q^{th} sample quantile by 0.25, 0.5 or 0.75, we obtain the first, second (**sample median**) or third (**sample quartile**), respectively.

The following algorithm can be used to obtain the q^{th} sample quantile of n IID samples (x_1, x_2, \dots, x_n) on the basis of their order statistics $(x_{(1)}, x_{(2)}, \dots, x_{(n)})$.

The q^{th} sample quantile, $\hat{F}_n^{[-1]}(q)$, is found by interpolation from the order statistics $(x_{(1)}, x_{(2)}, \dots, x_{(n)})$ of the n data points (x_1, x_2, \dots, x_n) , using the formula:

$$\hat{F}_n^{[-1]}(q) = (1 - \delta)x_{(i+1)} + \delta x_{(i+2)}, \quad \text{where,} \quad i = \lfloor (n-1)q \rfloor \quad \text{and} \quad \delta = (n-1)q - \lfloor (n-1)q \rfloor.$$

Thus, the **sample minimum** of the data points (x_1, x_2, \dots, x_n) is given by $\hat{F}_n^{[-1]}(0)$, the **sample maximum** is given by $\hat{F}_n^{[-1]}(1)$ and the **sample median** is given by $\hat{F}_n^{[-1]}(0.5)$, etc.

Labwork 112 (The q^{th} sample quantile) Use the implementation of Algorithm 1 as the MATLAB function `qthSampleQuantile` to find the q^{th} sample quantile of two simulated data arrays:

1. `SortedXsFromUni01Twstr101`, the order statistics that was constructed in Labwork 110 and
2. Another sorted array of 7 samples called `SortedXs`

In particular, if $\mathbb{B}_\delta(x, y)$ denotes a square of a small area $\delta > 0$ that is centered at (x, y) , then the following approximate equality holds and improves as $\delta \rightarrow 0$:

$$\mathbf{P}((X, Y) \in \mathbb{B}_\delta(x, y)) \cong \delta f_{X,Y}(x, y). \quad (3.62)$$

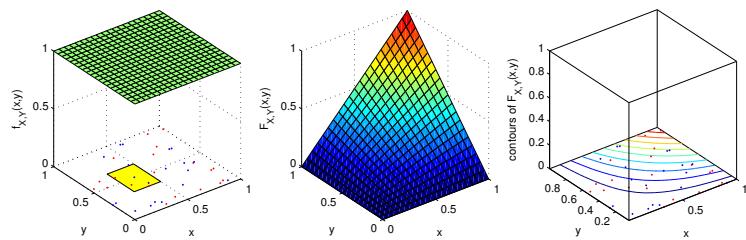
The JPDF satisfies the following two properties:

1. integrates to 1, i.e., $\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f_{X,Y}(x, y) dx dy = 1$
2. is a non-negative function, i.e., $f_{X,Y}(x, y) \geq 0$ for every $(x, y) \in \mathbb{R}^2$.

Example 84 Let (X, Y) be a continuous RV that is uniformly distributed on the unit square $[0, 1]^2 := [0, 1] \times [0, 1]$ with following JPDF:

$$f(x, y) = \mathbb{1}_{[0,1]^2}(x) \begin{cases} 1 & \text{if } (x, y) \in [0, 1]^2 \\ 0 & \text{otherwise.} \end{cases}$$

Find explicit expressions for the following: (1) DF $F(x, y)$ for any $(x, y) \in [0, 1]^2$, (2) $\mathbf{P}(X \leq 1/3, Y \leq 1/2)$, (3) $P((X, Y) \in [1/4, 1/2] \times [1/3, 2/3])$.



Let us begin to find the needed expressions.

1. Let $(x, y) \in [0, 1]^2$ then by Equation (3.60):

$$F_{X,Y}(x, y) = \int_{-\infty}^y \int_{-\infty}^x f_{X,Y}(u, v) du dv = \int_0^y \int_0^x 1 du dv = \int_0^y [u]_{u=0}^x dv = \int_0^y x dv = [xv]_{v=0}^y = xy$$

2. We can obtain $\mathbf{P}(X \leq 1/3, Y \leq 1/2)$ by evaluating $F_{X,Y}$ at $(1/3, 1/2)$:

$$\mathbf{P}(X \leq 1/3, Y \leq 1/2) = F_{X,Y}(1/3, 1/2) = \frac{1}{3} \frac{1}{2} = \frac{1}{6}$$

We can also find $\mathbf{P}(X \leq 1/3, Y \leq 1/2)$ by integrating the JPDF over the rectangular event $A = \{X < 1/3, Y < 1/2\} \subset [0, 1]^2$ according to Equation (3.61). This amounts here to finding the area of A , we compute $\mathbf{P}(A) = (1/3)(1/2) = 1/6$.

3. We can find $P((X, Y) \in [1/4, 1/2] \times [1/3, 2/3])$ by integrating the JPDF over the rectangular event $B = [1/4, 1/2] \times [1/3, 2/3]$ according to Equation (3.61):

$$\begin{aligned} P((X, Y) \in [1/4, 1/2] \times [1/3, 2/3]) &= \int \int_B f_{X,Y}(x, y) dx dy = \int_{1/3}^{2/3} \int_{1/4}^{1/2} 1 dx dy \\ &= \int_{1/3}^{2/3} [x]_{1/4}^{1/2} dy = \int_{1/3}^{2/3} \left[\frac{1}{2} - \frac{1}{4} \right] dy = \left(\frac{1}{2} - \frac{1}{4} \right) [y]_{1/3}^{2/3} \\ &= \left(\frac{1}{2} - \frac{1}{4} \right) \left(\frac{2}{3} - \frac{1}{3} \right) = \frac{1}{4} \left(\frac{1}{3} \right) = \frac{1}{12} \end{aligned}$$

Labwork 109 (Sample variance and sample standard deviation) We can compute the sample variance and sample standard deviation for the five samples stored in the array `XsFromUni01Twstr101` from Labwork 108 using MATLAB's functions `var` and `std`, respectively.

```
>> disp(XsFromUni01Twstr101); % The data-points x_1,x_2,x_3,x_4,x_5 are :
    0.5164 0.5707 0.0285 0.1715 0.6853
>> SampleVar=var(XsFromUni01Twstr101);% find sample variance
>> SampleStd=std(XsFromUni01Twstr101);% find sample standard deviation
>> disp(SampleVar) % The sample variance is:
    0.0785
>> disp(SampleStd) % The sample standard deviation is:
    0.2802
```

It is important to bear in mind that the statistics such as sample mean and sample variance are random variables and have an underlying distribution.

Definition 54 (Order Statistics) Suppose $X_1, X_2, \dots, X_n \stackrel{\text{IID}}{\sim} F$, where F is the DF from the set of all DFs over the real line. Then, the n -sample **order statistics** $X_{(n)}$ is:

$$X_{(n)}((X_1, X_2, \dots, X_n)) := (X_{(1)}, X_{(2)}, \dots, X_{(n)}) , \text{ such that, } X_{(1)} \leq X_{(2)} \leq \dots \leq X_{(n)} . \quad (3.80)$$

For brevity, we write $X_{(n)}((X_1, X_2, \dots, X_n))$ as $X_{(n)}$ and its realisation $X_{(n)}((x_1, x_2, \dots, x_n))$ as $x_{(n)} = (x_{(1)}, x_{(2)}, \dots, x_{(n)})$.

Without going into the details of how to sort the data in ascending order to obtain the order statistics (an elementary topic of an Introductory Computer Science course), we simply use MATLAB's function `sort` to obtain the order statistics, as illustrated in the following example.

Labwork 110 (Order statistics and sorting) The order statistics for the five samples stored in `XsFromUni01Twstr101` from Labwork 108 can be computed using `sort` as follows:

```
>> disp(XsFromUni01Twstr101); % display the sample points
    0.5164 0.5707 0.0285 0.1715 0.6853
>> SortedXsFromUni01Twstr101=sort(XsFromUni01Twstr101); % sort data
>> disp(SortedXsFromUni01Twstr101); % display the order statistics
    0.0285 0.1715 0.5164 0.5707 0.6853
```

Therefore, we can use `sort` to obtain our order statistics $x_{(1)}, x_{(2)}, \dots, x_{(n)}$ from n sample points x_1, x_2, \dots, x_n .

Next, we will introduce a family of common statistics, called the q^{th} quantile, by first defining the function:

Definition 55 (Inverse DF or Inverse CDF or Quantile Function) Let X be an RV with DF F . The **inverse DF** or **inverse CDF** or **quantile function** is:

$$F^{[-1]}(q) := \inf \{x : F(x) > q\}, \quad \text{for some } q \in [0, 1] . \quad (3.81)$$

If F is strictly increasing and continuous then $F^{[-1]}(q)$ is the unique $x \in \mathbb{R}$ such that $F(x) = q$.

A **functional** is merely a function of another function. Thus, $T(F) : \{\text{All DFs}\} \rightarrow \mathbb{T}$, being a map or function from the space of DFs to its range \mathbb{T} , is a functional. Some specific examples of functionals we have already seen include:

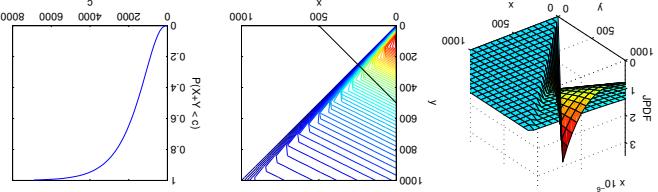
In general, for a bivariate uniform RV on the unit square the $H(a, b | c, d) = (b-a)(d-c)$ for any event given by the rectangular region $[a, b] \times [c, d]$ inside the unit square the $P([a, b] \times [c, d]) = (b-a)(d-c)$. This is true for any two events with the same rectangular area have the same probability (imagine sliding a small rectangle inside the unit square... no matter where you slide this rectangle to while remaining inside the unit square, the probability of $w \hookrightarrow (X(w), Y(w)) = (x, y)$ falling inside this "slidable" rectangle is the same...).

1. Identify the support of (X, Y) , i.e., the region in the plane where $f_{X,Y}$ takes positive values. Answer the following:

3. Find $P(X \leq 400, Y \leq 800)$

4. It is known that humans prefer a response time of under $1/10$ seconds (10^2 milliseconds) from the web server before they get impatient. What is $P(X + Y < 10^2)$?

^{1.} The support is the intersection of the positive quadrant with the $y > x$ half-plane.



Let us answer the questions.

1. The support is the intersection of the positive quadrant with the $y > x$ half-plane.

Once again, if X_1, X_2, \dots, X_n are $\sim X_1$, the expectation of the sample variance is:

For brevity, we write $S_n((X_1, X_2, \dots, X_n))$ as S_n and its realization $S_n((x_1, x_2, \dots, x_n))$ as s_n .

$$S^n((X^1, X^2, \dots, X^n)) = \bigwedge S^2_n((X^1, X^2, \dots, X^n))$$

$$S^u(X_1, X_2, \dots, X_n) = \bigwedge^u S$$

Sample standard deviation is simply the square root of sample variance:

For brevity, we write $S^2_n((X_1, X_2, \dots, X_n))$ as S^2_n and its realisation $S^2_n((x_1, x_2, \dots, x_n))$ as s^2_n .

$$L^u((X_1, X_2, \dots, X_u)) = S_{\epsilon}^u(X_1, X_2, \dots, X_u) \quad (3.78)$$

simply the sample variance:

Definition 53 (Sample Variance & Standard Deviation) From a given a sequence of random variables X_1, X_2, \dots, X_n , we may obtain another statistic called the n -samples variance or

```

>>> size(XStarmandjitske101) % size(SomeMatrix) gives the size or dimensions of the array SomeMatrix
ans = 5 5

>>> ones(5,1) % here ones(5,1) is an array of 1's with size or dimension 5 X 1
ans = 1 1 1 1 1

>>> size(XStarmandjitske101 * ones(5,1)) % multiplying an 1 X 5 matrix with a 5 X 1 matrix of 1's
ans = 5 5

>>> XStarmandjitske101 * (ones(5,1) * 1/5) % multiplying an 1 X 5 matrix with a 5 X 1 matrix of 1/5
ans = 1.9723

```

We can also obtain the sample mean via matrix product or multiplication as follows:

```

>>> sum(xPerformance101) / 5 % take the sum of the elements of the xPerformance101 array
ans =
1.923
>>> sum(xPerformance101) / 5 % divide the sum by the sample size 5
ans =
0.3945

```

We may also obtain the sample mean using the `sum` function and a division by sample size:

Labwork 108 (Sample mean) After initializing the fundamental samples, we draw five samples and then obtain the sample mean using the MATLAB function `mean`. In the following, we will reuse the samples stored in the array `X$Performance108$`.

CHAPTER 3. RANDOM VARIABLES

2.

$$\begin{aligned}
\int_{y=-\infty}^{\infty} \int_{x=-\infty}^{\infty} f_{X,Y}(x,y) dy dx &= \int_{x=0}^{\infty} \int_{y=x}^{\infty} f_{X,Y}(x,y) dy dx \\
&= \int_{x=0}^{\infty} \int_{y=x}^{\infty} \frac{6}{10^6} \exp\left(-\frac{1}{1000}x - \frac{2}{1000}y\right) dy dx \\
&= \frac{6}{10^6} \int_{x=0}^{\infty} \left(\int_{y=x}^{\infty} \exp\left(-\frac{2}{1000}y\right) dy \right) \exp\left(-\frac{1}{1000}x\right) dx \\
&= \frac{6}{10^6} \int_{x=0}^{\infty} \left[-\frac{1000}{2} \exp\left(-\frac{2}{1000}y\right) \right]_{y=x}^{\infty} \exp\left(-\frac{1}{1000}x\right) dx \\
&= \frac{6}{10^6} \int_{x=0}^{\infty} \left[0 + \frac{1000}{2} \exp\left(-\frac{2}{1000}x\right) \right] \exp\left(-\frac{1}{1000}x\right) dx \\
&= \frac{6}{10^6} \int_{x=0}^{\infty} \frac{1000}{2} \exp\left(-\frac{2}{1000}x - \frac{1}{1000}x\right) dx \\
&= \frac{6}{10^6} \frac{1000}{2} \left[-\frac{1000}{3} \exp\left(-\frac{3}{1000}x\right) \right]_{x=0}^{\infty} \\
&= \frac{6}{10^6} \frac{1000}{2} \left[0 + \frac{1000}{3} \right] \\
&= 1
\end{aligned}$$

3. First, identify the region with positive JPDF for the event $(X \leq 400, Y \leq 800)$

$$\begin{aligned}
\mathbf{P}(X \leq 400, Y \leq 800) &= \int_{x=0}^{400} \int_{y=x}^{800} f_{X,Y}(x,y) dy dx \\
&= \int_{x=0}^{400} \int_{y=x}^{800} \frac{6}{10^6} \exp\left(-\frac{1}{1000}x - \frac{2}{1000}y\right) dy dx \\
&= \frac{6}{10^6} \int_{x=0}^{400} \left[-\frac{1000}{2} \exp\left(-\frac{2}{1000}y\right) \right]_{y=x}^{800} \exp\left(-\frac{1}{1000}x\right) dx \\
&= \frac{6}{10^6} \frac{1000}{2} \int_{x=0}^{400} \left(-\exp\left(-\frac{1600}{1000}x\right) + \exp\left(-\frac{2}{1000}x\right) \right) \exp\left(-\frac{1}{1000}x\right) dx \\
&= \frac{6}{10^6} \frac{1000}{2} \int_{x=0}^{400} \left(\exp\left(-\frac{3}{1000}x\right) - e^{-8/5} \exp\left(-\frac{1}{1000}x\right) \right) dx \\
&= \frac{6}{10^6} \frac{1000}{2} \left(\left(-\frac{1000}{3} \exp\left(-\frac{3}{1000}x\right) \right)_{x=0}^{400} - e^{-8/5} \left(-1000 \exp\left(-\frac{1}{1000}x\right) \right)_{x=0}^{400} \right) \\
&= \frac{6}{10^6} \frac{1000}{2} 1000 \left(\frac{1}{3} \left(1 - e^{-6/5} \right) - e^{-8/5} \left(1 - e^{-2/5} \right) \right) \\
&= 3 \left(\frac{1}{3} \left(1 - e^{-6/5} \right) - e^{-8/5} \left(1 - e^{-2/5} \right) \right) \\
&\approx 0.499 .
\end{aligned}$$

4. First, identify the region with positive JPDF for the event $(X + Y \leq c)$, say $c = 500$ (but generally c can be any positive number). This is the triangular region at the intersection of the four half-planes: $x > 0$, $x < c$, $y > x$ and $y < c - x$. (Draw picture here) Let's integrate

data $X(\omega) = x$, a statistic? In other words, is the data a statistic? [Hint: consider the identity map $T(x) = x : \mathbb{X} \rightarrow \mathbb{T} = \mathbb{X}$.]

Next, we define two important statistics called the **sample mean** and **sample variance**. Since they are obtained from the sample data, they are called **sample moments**, as opposed to the **population moments**. The corresponding population moments are $\mathbf{E}(X_1)$ and $\mathbf{V}(X_1)$, respectively.

Definition 52 (Sample Mean) From a given a sequence of RVs X_1, X_2, \dots, X_n , we may obtain another RV called the n -samples mean or simply the sample mean:

$$T_n((X_1, X_2, \dots, X_n)) = \bar{X}_n((X_1, X_2, \dots, X_n)) := \frac{1}{n} \sum_{i=1}^n X_i . \quad (3.75)$$

For brevity, we write

$$\bar{X}_n((X_1, X_2, \dots, X_n)) \text{ as } \bar{X}_n ,$$

and its realisation

$$\bar{X}_n((x_1, x_2, \dots, x_n)) \text{ as } \bar{x}_n .$$

Note that the expectation and variance of \bar{X}_n are:

$$\begin{aligned}
\mathbf{E}(\bar{X}_n) &= \mathbf{E}\left(\frac{1}{n} \sum_{i=1}^n X_i\right) && [\text{by definition (3.75)}] \\
&= \frac{1}{n} \sum_{i=1}^n \mathbf{E}(X_i) && [\text{by property (3.50)}]
\end{aligned}$$

Furthermore, if every X_i in the original sequence of RVs X_1, X_2, \dots is **identically distributed** with the same expectation, by convention $\mathbf{E}(X_1)$, then:

$$\mathbf{E}(\bar{X}_n) = \frac{1}{n} \sum_{i=1}^n \mathbf{E}(X_i) = \frac{1}{n} \sum_{i=1}^n \mathbf{E}(X_1) = \frac{1}{n} n \mathbf{E}(X_1) = \mathbf{E}(X_1) . \quad (3.76)$$

Similarly, we can show that:

$$\begin{aligned}
\mathbf{V}(\bar{X}_n) &= \mathbf{V}\left(\frac{1}{n} \sum_{i=1}^n X_i\right) && [\text{by definition (3.75)}] \\
&= \left(\frac{1}{n}\right)^2 \mathbf{V}\left(\sum_{i=1}^n X_i\right) && [\text{by property (3.49)}]
\end{aligned}$$

Furthermore, if the original sequence of RVs X_1, X_2, \dots is **independently distributed** then:

$$\mathbf{V}(\bar{X}_n) = \left(\frac{1}{n}\right)^2 \mathbf{V}\left(\sum_{i=1}^n X_i\right) = \frac{1}{n^2} \sum_{i=1}^n \mathbf{V}(X_i) \quad [\text{by property (3.51)}]$$

Finally, if the original sequence of RVs X_1, X_2, \dots is **independently and identically distributed** with the same variance ($\mathbf{V}(X_1)$) by convention then:

$$\mathbf{V}(\bar{X}_n) = \frac{1}{n^2} \sum_{i=1}^n \mathbf{V}(X_i) = \frac{1}{n^2} \sum_{i=1}^n \mathbf{V}(X_1) = \frac{1}{n^2} n \mathbf{V}(X_1) = \frac{1}{n} \mathbf{V}(X_1) . \quad (3.77)$$

$X = 0$	$Y = 0$	$X = 1$
0.1	0.3	0.2
0.3	0.1	0.4

Example 86 Obtain the marginal PMFs $f_Y(y)$ and $f_X(x)$ from the joint PMF $f_{X,Y}(x,y)$ of the discrete RV in Example 83. Just sum $f_{X,Y}(x,y)$ over x 's and y 's (reported in a tabular form):

$$\left\{ \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f_X(x, y) dx dy \quad \text{if } (X, Y) \text{ is a continuous RV} \right. \\ \left. \sum_x f_X(x, y) \quad \text{if } (X, Y) \text{ is a discrete RV} \right\} = f_Y(y)$$

and the marginal PDF or PMF of Y is defined by:

$$\left\{ \int_{-\infty}^y f(X, x, y) dy \quad \text{if } (X, Y) \text{ is a continuous RV} \right. \\ \left. \sum_y f(X, x, y) \quad \text{if } (X, Y) \text{ is a discrete RV} \right\} = (x)^X f$$

Definition 37 (Marginal PDF or PMF) If the RV (X, Y) has joint PDF $f_{X,Y}(x,y)$, as its marginal PDF or joint PMF, then the marginal PDF or PMF of a random vector (X, Y) is defined by :

```

>>> c = [100 1000 2000 3000 4000]
>>> p = 0.0134
      0.5135
      0.8558
      0.9630
      0.9911

>>> p = 1 - 4 * exp(-3*c/2000) + 3 * exp(-c/500)

```

We can obtain $\mathbf{P}(X + Y < c)$ for several values of c by using MATLAB and note that about 96% of requests are processed in less than 3000 milliseconds or 3 seconds.

$$\begin{aligned}
&= 1 - \frac{4e^{-\zeta c}/2000 + 3e^{-\zeta c}/500}{1 - e^{-\zeta c}/2000 + 3e^{-\zeta c}/1000 - 3e^{-\zeta c}/2000} \\
&= 1 - \frac{3}{1} \left(1 - e^{-\zeta c}/2000 \right) - e^{-\zeta c}/1000 \left(e^{\zeta c}/2000 - 1 \right) \\
&= 3 \left(1 - e^{-\zeta c}/2000 \right) - e^{-\zeta c}/1000 \left(e^{\zeta c}/2000 - 1 \right) \\
&= \left(\frac{0=x}{\zeta/c} \left[\left(\frac{1000}{x} \right) dx \right] - \frac{0=x}{\zeta/c} \left[\left(\frac{0001}{x\zeta} - \frac{x}{\zeta} \right) dx \right] \right) \varepsilon = \\
&\quad xp \left(\left(\frac{0001}{x\zeta - x} \right) dx - \left(\frac{0001}{x\zeta} \right) dx \right) \varepsilon = \\
&xp \left(\frac{0001}{x} - \frac{0001}{x\zeta} \right) dx \left[\left(\frac{0001}{x\zeta} - \frac{0001}{x\zeta - x} \right) dx + \left(\frac{0001}{x\zeta - x} - \frac{0001}{x} \right) dx \right] \varepsilon = \\
&\quad xp \left(\frac{x-0001}{1} - \frac{0001}{1} \right) dx \left[\frac{x-0001}{2} \left(\frac{0001}{x-x} - \frac{0001}{x-\zeta} \right) dx \right] \varepsilon = \\
&\quad xp \#ip \left(\frac{0001}{2} - \frac{0001}{1} - x \frac{0001}{1} - x \frac{0001}{2} \right) dx \varepsilon = \\
&\quad xp \#ip \left(\frac{0001}{2} - \frac{0001}{1} - x \frac{0001}{1} \right) dx \varepsilon = \frac{901}{9} \int_{x-\zeta/c}^{0=x} \frac{901}{z/c} dz = \\
&\quad xp \#ip \left(\zeta \cdot x \right) X' X f \int_{x-\zeta/c}^{0=x} \frac{901}{z/c} dz = (\omega \gtrdot X + X) \omega
\end{aligned}$$

the JPDF over our triangular event as follows:

Classwork 107 (Is data a statistic?) Is the RV X , for which the realization is the observed

Thus, a statistic T is also an RV that takes values in the space \mathbb{I} . When $x \in \mathcal{E}$ is the realisation of X , we use $T_n(x)$ and \mathbb{I}_n to emphasise that X is an n -dimensional random vector, i.e. $\mathbb{X} \in \mathbb{R}^n$. Some times an experiment, we let $T(x) = t$ denote the corresponding realisation of the statistic T . Sometimes we use $T_n(x)$ and \mathbb{I}_n to emphasise that X is an n -dimensional random vector, i.e. $\mathbb{X} \in \mathbb{R}^n$.

$\cdot \mathbb{L} \leftarrow \mathbb{X} : (x)_{\mathbb{L}}$

Definition 51 (Statistic) A statistic T is any function of the data:

Example 106 (Lossing a coin n times) For some given parameter $\theta \in \Theta = [0, 1]$, consider n IID Bernoulli(θ) trials, i.e., X_1, X_2, \dots, X_n . Then the random vector $X = (X_1, X_2, \dots, X_n)$, which takes values in the data space $\mathbb{X} = \{0, 1\}^n = \{(x_1, x_2, \dots, x_n) : x_i \in \{0, 1\}, i = 1, 2, \dots, n\}$, made up of vertices of the n -dimensional hypercube, measures the outcomes of this experiment. A particular realization of X , upon permuting hyper-cube, measures the observations, data or data vector (x_1, x_2, \dots, x_n) . For instance, if we observed $n - 1$ tails and 1 heads, in that order, then our data vector $(x_1, x_2, \dots, x_{n-1}, x_n) = (0, \dots, 0, 1)$.

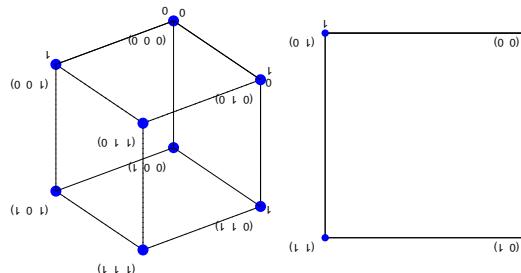


Figure 3.25: Data Spaces $X = \{0, 1\}^2$ and $\tilde{X} = \{0, 1\}^2$ for two and three Bernoulli trials, respectively.

Figure 3.24: Sample Space, Random Variable, Realization, Data, and Data Space.

From the above Table we can find:

$$\begin{aligned} f_X(x) &= \mathbf{P}(X = x) = \sum_y f_{X,Y}(x,y) \\ &= f_{X,Y}(x,0) + f_{X,Y}(x,1) = \begin{cases} f_{X,Y}(0,0) + f_{X,Y}(0,1) = 0.1 + 0.3 = 0.4 & \text{if } x = 0 \\ f_{X,Y}(1,0) + f_{X,Y}(1,1) = 0.2 + 0.4 = 0.6 & \text{if } x = 1 \end{cases} \end{aligned}$$

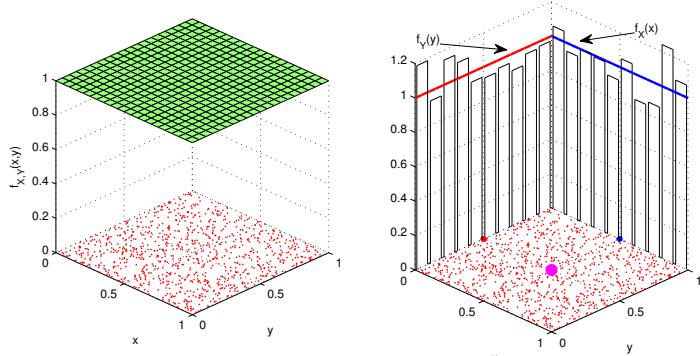
Similarly,

$$\begin{aligned} f_Y(y) &= \mathbf{P}(Y = y) = \sum_x f_{X,Y}(x,y) \\ &= f_{X,Y}(0,y) + f_{X,Y}(1,y) = \begin{cases} f_{X,Y}(0,0) + f_{X,Y}(1,0) = 0.1 + 0.2 = 0.3 & \text{if } y = 0 \\ f_{X,Y}(0,1) + f_{X,Y}(1,1) = 0.3 + 0.4 = 0.7 & \text{if } y = 1 \end{cases} \end{aligned}$$

Just report the marginal probabilities as row and column sums of the JPDF table.

Thus marginal PMF gives us the probability of a specific RV, within a R \vec{V} , taking a value irrespective of the value taken by the other RV in this R \vec{V} .

Example 87 Obtain the marginal PDFs $f_Y(y)$ and $f_X(x)$ from the joint PDF $f_{X,Y}(x,y)$ of the continuous R \vec{V} in Example 84 (the bivariate uniform R \vec{V} on $[0, 1]^2$).



Let us suppose $(x, y) \in [0, 1]^2$ and note that $f_{X,Y} = 0$ if $(x, y) \notin [0, 1]^2$. We can obtain marginal PMFs $f_X(x)$ and $f_Y(y)$ by integrating the JPDF $f_{X,Y} = 1$ along y and x , respectively.

$$f_X(x) = \int_{-\infty}^{\infty} f_{X,Y}(x,y) dy = \int_0^1 f_{X,Y}(x,y) dy = \int_0^1 1 dy = [y]_0^1 = 1 - 0 = 1$$

Similarly,

$$f_Y(y) = \int_{-\infty}^{\infty} f_{X,Y}(x,y) dx = \int_0^1 f_{X,Y}(x,y) dx = \int_0^1 1 dx = [x]_0^1 = 1 - 0 = 1$$

We are seeing a histogram of the **marginal samples** and their marginal PDFs in the Figure.

1. Find the characteristic function (CF) of X

2. Using the CF find $V(X)$, the variance of X . Hint: $V(X) = E(X^2) - (E(X))^2$

Ex. 3.51 — Recall that the Geometric(θ) RV X has the following PMF

$$f_X(x; \theta) = \begin{cases} \theta(1-\theta)^x & \text{if } x \in \{0, 1, 2, 3, \dots\} \\ 0 & \text{otherwise} \end{cases}$$

- Find the CF of X . (Hint: the sum of the infinite geometric series $\sum_{x=0}^{\infty} ar^x = \frac{a}{1-r}$.)
- Using the CF find $E(X)$.

Ex. 3.52 — Let X be the Uniform(a, b) RV with the following probability density function (PDF)

$$f_X(x; a, b) = \begin{cases} \frac{1}{b-a} & \text{if } x \in [a, b] \\ 0 & \text{otherwise} \end{cases}$$

Find the CF of X .

Ex. 3.53 — Recall that the Poisson(λ) RV has the following PMF

$$f_X(x; \lambda) = \begin{cases} \frac{\lambda^x}{x!} e^{-\lambda} & \text{if } x \in \{0, 1, 2, 3, \dots\} \\ 0 & \text{otherwise} \end{cases}$$

Hint: the power series of $e^\alpha = \sum_{x=0}^{\infty} \frac{\alpha^x}{x!}$.

- Find the CF of X .
- Find the variance of X using its CF.

Ex. 3.54 — Let X be a Poisson(λ) RV and Y be another Poisson(μ) RV. Suppose X and Y are independent. Use Eqn. (3.74) to first find the CF of the RV $W = X + Y$. From the CF of W try to identify what RV it is.

Ex. 3.55 — Recall from lecture that if $Y = a + bX$ for some constants a and b with $b \neq 0$ then $\varphi_Y(t) = e^{iat}\varphi_X(bt)$ and that $\varphi_Z(t) = e^{-t^2/2}$ if Z is the Normal($0, 1$) RV. Using these facts find the CF of $-Z$, the RV obtained from Z by simply switching its sign. From the CF of $-Z$ identify what RV it is.

3.14 Statistics

3.14.1 Data and Statistics

Definition 50 (Data) The function X measures the outcome ω of an experiment with sample space Ω [Often, the sample space is also denoted by S]. Formally, X is a random variable [or a random vector $X = (X_1, X_2, \dots, X_n)$, i.e. a vector of random variables] taking values in the **data space** \mathbb{X} :

$$X(\omega) : \Omega \rightarrow \mathbb{X}$$

The realisation of the RV X when an experiment is performed is the observation or data $x \in \mathbb{X}$. That is, when the experiment is performed once and it yields a specific $\omega \in \Omega$, the data $X(\omega) = x \in \mathbb{X}$ is the corresponding realisation of the RV X .

We have seen the notion of independence of two events in Definition 16 or of a sequence of events in Definition 17. Recall that independence amounts to having the probability of the joint occurrence of the events to be given by the product of the probabilities of each of the events.

We can use the definition of independence of two events to define the independence of two random variables using their distribution functions.

Definition 38 (Independence of Two RVs) Consider an \mathbb{R}^2 -valued RV $X := (X_1, X_2)$. Then the \mathbb{R} -valued RVs X_1 and X_2 are said to be independent or independently distributed if and only if

$$\mathbf{P}(X_1 \leq x_1, X_2 \leq x_2) = \mathbf{P}(X_1 \leq x_1)\mathbf{P}(X_2 \leq x_2)$$

or equivalently,

$$F_{X_1, X_2}(x_1, x_2) = F_{X_1}(x_1)F_{X_2}(x_2),$$

for any pair of real numbers $(x_1, x_2) \in \mathbb{R}^2$.

By the above definition, for **discrete** RVs X_1, X_2 that are independent, the following equality is satisfied between the joint and marginal PMFs:

$$f_{X_1, X_2}(x_1, x_2) = \mathbf{P}(X_1 = x_1, X_2 = x_2) = \mathbf{P}(X_1 = x_1)\mathbf{P}(X_2 = x_2) = f_{X_1}(x_1)f_{X_2}(x_2) \text{ for any } (x_1, x_2) \in \mathbb{R}^2,$$

and for **continuous** RVs X_1, X_2 that are independent, the following equality is satisfied between the joint and marginal PDFs:

$$f_{X_1, X_2}(x_1, x_2) = f_{X_1}(x_1)f_{X_2}(x_2) \text{ for any } (x_1, x_2) \in \mathbb{R}^2.$$

In summary, two RVs X and Y are said to be **independent** if and only if for every (x, y)

$$F_{X, Y}(x, y) = F_X(x) \times F_Y(y) \quad \text{or} \quad f_{X, Y}(x, y) = f_X(x) \times f_Y(y)$$

Let us confirm that our familiar experiment of tossing a fair coin twice independently when encoded by a pair of independent Bernoulli(1/2) RVs satisfies the above definition.

Example 89 (Pair of independent Bernoulli(1/2) RVs) Let X_1 and X_2 be a pair of independent Bernoulli(1/2) RVs each taking values in the set $\{0, 1\}$ with the following tabulated probabilities. Verify that the JPMF $f_{X_1, X_2}(x_1, x_2) = 1/4$ for each $(x_1, x_2) \in \{0, 1\}^2$ is indeed given by the marginal PMF $f_{X_i}(x_i) = 1/2$ for each $i \in \{1, 2\}$ and each $x_i \in \{0, 1\}$.

	$X_2 = 0$	$X_2 = 1$	
$X_1 = 0$	1/4	1/4	1/2
$X_1 = 1$	1/4	1/4	1/2
	1/2	1/2	1

From the above Table we can read for instance that the *joint probability* that \mathbb{R}^2 -valued RV (X_1, X_2) takes the value or realization $(0, 0)$ is 1/4 from the first entry of the inner-most tabulated rectangle, i.e., $\mathbf{P}((X_1, X_2) = (0, 0)) = 1/4$, and that the *marginal probability* that the RV X_1 takes the value or realization 0 is 1/2, i.e., $\mathbf{P}(X_1 = 0) = 1/2$. Clearly, $1/4 = 1/2 \times 1/2$, and so our familiar experiment when seen as an \mathbb{R}^2 -valued RV is indeed composed of two independent \mathbb{R} -valued Bernoulli(1/2) RVs.

Thus the CF of the standard normal RV Z is

$$\boxed{\varphi_Z(t) = e^{-t^2/2}} \quad (3.72)$$

Let X be a RV with CF $\varphi_X(t)$. Let Y be a linear transformation of X

$$Y = a + bX$$

where a and b are two constant real numbers and $b \neq 0$. Then the CF of Y is

$$\boxed{\varphi_Y(t) = \exp(itb)\varphi_X(bt)} \quad (3.73)$$

Proof: This is easy to prove using the definition of CF as follows:

$$\begin{aligned} \varphi_Y(t) &= E(\exp(itY)) = E(\exp(it(a + bX))) = E(\exp(itbX + ita)) \\ &= E(\exp(itbX))\exp(itbX) = \exp(itbX)E(\exp(itbX)) = \exp(itbX)\varphi_X(bt) \end{aligned}$$

Example 103 Let Y be a Normal(μ, σ^2) RV. Recall that Y is a linear transformation of Z , i.e., $Y = \mu + \sigma Z$ where Z is a Normal(0, 1) RV. Using Equations (3.72) and (3.73) find the CF of Y .

Solution:

$$\begin{aligned} \varphi_Y(t) &= \exp(i\mu t)\varphi_Z(\sigma t), \quad \text{since } Y = \mu + \sigma Z \\ &= e^{i\mu t}e^{(-\sigma^2 t^2)/2}, \quad \text{since } \varphi_Z(t) = e^{-t^2/2} \\ &= e^{i\mu t - (\sigma^2 t^2)/2} \end{aligned}$$

A generalization of (3.73) is the following. If X_1, X_2, \dots, X_n are independent RVs and a_1, a_2, \dots, a_n are some constants, then the CF of the linear combination $Y = \sum_{i=1}^n a_i X_i$ is

$$\boxed{\varphi_Y(t) = \varphi_{X_1}(a_1 t) \times \varphi_{X_2}(a_2 t) \times \cdots \times \varphi_{X_n}(a_n t) = \prod_{i=1}^n \varphi_{X_i}(a_i t).} \quad (3.74)$$

Example 104 Using the following three facts:

- Eqn. (3.74)
- the Binomial(n, θ) RV Y is the sum of n independent Bernoulli(θ) RVs (from Probability Course)
- the CF of Bernoulli(θ) RV (from lecture notes for Inference Course)

find the CF of the Binomial(n, θ) RV Y .

Solution:

Let X_1, X_2, \dots, X_n be independent Bernoulli(θ) RVs with CF $(1 - \theta + \theta e^{it})$ then $Y = \sum_{i=1}^n X_i$ is the Binomial(n, θ) RV and by Eqn. (3.74) with $a_1 = a_2 = \dots = 1$, we get

$$\varphi_Y(t) = \varphi_{X_1}(t) \times \varphi_{X_2}(t) \times \cdots \times \varphi_{X_n}(t) = \prod_{i=1}^n \varphi_{X_i}(t) = \prod_{i=1}^n (1 - \theta + \theta e^{it}) = (1 - \theta + \theta e^{it})^n.$$

$$f_{X_1, X_2}(x_1, x_2) = \frac{1}{\pi} \int_{[0, \pi]^2} \int_{[0, \pi]^2} \frac{1}{\pi} \int_{[0, \pi]^2} f_{X_1, X_2}(x_1, x_2) d\theta_1 d\theta_2 d\phi_1 d\phi_2$$

Let X_1 be the angle between the needle and the direction of the rulings, and let X_2 be the distance between the bottom point of the needle and the nearest line above this point (see left sub-figure). Hence assuming that the RVs X_1 and X_2 are independent, we find that their joint probability density function (JPDF) is:

Solution:

What is the probability that the needle intersects one of the parallel lines? Can you use repeated trials of this experiment to find an approximation to π ?

Example 92 (bottom s Needie experiment to physically estimate L) Suppose a needle is tossed at random onto a plane ruled with parallel lines a distance L apart. By a "needle" we mean a line segment of length $l \leq L$.

done in lectures...

Example 19 (distance between random numbers in a manufactured line) Suppose two points are tossed independently and uniformly at random onto a line segment of unit length. What is the probability that the distance between the two points does not exceed a given length l ?

Now, let us take advantage of independent random variables and solve some problems.

zero density when $x > y$, but the product of the marginal densities won't.

autocorrelation of the marginal PDFs. But intuitively, we know that these RVs (connection in time and space) are dependent – one is strictly greater than the other. Also the PDF has

the product of the marginal PDFs. But intuitively, we know that these RVs (connection in time and space) are dependent – one is strictly greater than the other. Also the PDF has

we can compute $J_X(x)$ and use the already computed $J_Y(y)$ to interchangeably change the PDFs.

dependent in the server times RV from Example 85?

This can be shown by checking that the joint PDF is indeed equal to the product of the marginal PDFs of $U_{11}(0, 1)$ RVs as follows:

Example 90 Recall the \mathbb{R}^n -valued continuous RV (x, γ) of Example 84 that is uniformly distributed on the unit square $[0, 1]^2$. First show that X and Y independent. Then show both X and Y are identically distributed according to the Uniform(0, 1) RV.

$$\begin{aligned}
& = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} e^{-\frac{(z-n)^2}{2}} dz dp(z) \\
& = \int_{-\infty}^{\infty} e^{-\frac{(z-n)^2}{2}} dz = \sqrt{2\pi} \\
& = \int_{-\infty}^{\infty} e^{-\frac{(z-n)^2}{2}} dz = \sqrt{2\pi} \int_{-\infty}^{\infty} e^{-\frac{(y-n)^2}{2}} dy = \sqrt{2\pi} \\
& = e^{-\frac{n^2}{2}} \int_{-\infty}^{\infty} e^{-\frac{(y-n)^2}{2}} dy = e^{-\frac{n^2}{2}} \text{ using the normalizing constant in PDF of } \text{Normal}(0, 1) \text{ RV}
\end{aligned}$$

(1999) A Probability Path, Birkhäuser). Thus, if we can show that two RVs have the same CDF then we know they are the same. This can be much more challenging or impossible to do directly with their DFs. Let Z be $\text{Normal}(0, 1)$, the standard normal RV. We can find the CDF for Z using couple of tricks as follows.

Let's check that this is what we had as variance for the Exponentia(l λ) RV when we first introduced it and directly computed using integrals for definition of expectation.

Finally, from the first and second moments we can get the variance as follows:

$$\begin{aligned} & \cdot \frac{\chi}{\zeta} = \left(\frac{\chi}{\zeta} \right)^0 = \left[\frac{\chi - \mu}{2\chi^2} \right]^0 = 0 = \varepsilon_{-(\mu - \chi)} \varepsilon_{-\chi} \left(\frac{\chi - \mu}{\zeta} \right) \frac{\zeta}{1} = \\ & = \left[\left((\mu - \chi) \frac{\mu p}{p} - \varepsilon_{-\chi} \varepsilon_{-\chi} \right) \varepsilon \right] \frac{\zeta}{1} = 0 = \left[\varepsilon_{-(\mu - \chi)} \frac{\mu p}{p} \times \varepsilon_{-\chi} \right] \frac{\zeta}{1} = \\ & = \left[\varepsilon_{-\chi} \frac{\mu p}{p} \right] \frac{\zeta}{1} = 0 = \left[(\chi X) \cancel{\varepsilon} \frac{\mu p}{p} \right] \frac{\zeta}{1} = 0 = \left[(\chi X) \cancel{\varepsilon} \frac{\mu p}{p} \right] \frac{\zeta}{1} = (\chi X) E \end{aligned}$$

Similarly from Equation (3.71) we can get $E(X^2)$ as follows:

$$\frac{\chi}{l} = \left(\frac{\chi}{i}\right) \frac{i}{l} = \left(\frac{\bar{z}\chi}{i\chi}\right) \frac{i}{l} = \begin{matrix} 0=i \\ \left[\frac{\bar{z}(i\chi - \chi)}{i\chi}\right] \end{matrix} \frac{i}{l} = \begin{matrix} 0=i \\ \left[(i)X \cancel{\sigma} \frac{ip}{p}\right] \end{matrix} \frac{i}{l} = (X)\cancel{\sigma}$$

We get $E(X)$ by evaluating $\frac{d}{dt}\phi X(t)$ at $t = 0$ and dividing by i according to Equation (3.71) as follows.

$$\frac{\varepsilon^{(\mu - \chi)}}{i\chi} = \left((i-) \times \frac{\varepsilon^{(\mu - \chi)}}{I-} \right) \chi = \\ \left((\mu - \chi) \frac{ip}{p} \times \varepsilon_{-} (\mu - \chi) \times I- \right) \chi = \left(\frac{\mu - \chi}{\chi} \right) \frac{ip}{p} = (i) x \not{=} \frac{ip}{p}$$

Let us differentiate the Cf to get moments using Equation (3.1) (Cf has to be once and twice differentiable at $t = 0$ to get the first and second moments).

Figure 3.19: Diagrams done on the board!

The event A that the needle intersects one of the parallel ruled lines occurs if and only if

$$X_2 \leq l \sin(X_1) ,$$

i.e., if and only if the corresponding point $X := (X_1, X_2)$ falls in the region B , where B is part of the rectangle $[0, \pi] \times [0, L]$ lying between the x_1 -axis and the curve $x_2 = \sin(x_1)$ (area under the curve in right-subfigure of Figure 3.19). Hence, we can integrate the JPDF to get the probability of the event A of interest:

$$\mathbf{P}(A) = \mathbf{P}((X_1, X_2) \in B) = \int_B \frac{dx_1 dx_2}{\pi L} = \frac{2l}{\pi L}$$

where,

$$l \int_0^\pi \pi \sin(x_1) dx_1 = l(-\cos(x_1))|_0^\pi = l(1 - (-1)) = l(1 + 1) = 2l ,$$

is the area of B .

Thus, if the needle is repeatedly tossed onto the ruled plane and $n(A)$ is the number of times A occurs out of n trials, then the relative frequency of the event A should approach $\mathbf{P}(A)$ as $n \rightarrow \infty$ (we will see this as the Law of Large Numbers in the sequel, but recall that this is also how we motivated the LTRF or long-term relative frequency idea of probability):

$$\frac{n(A)}{n} \rightarrow \frac{2l}{\pi L}$$

Hence, for large n ,

$$\frac{2l}{L} \frac{n}{n(A)}$$

should be a good approximation to $\pi = 3.14\dots$. This is indeed the case.

3.10.2 Conditional Random Variables

Often we will have a condition where one of the two random variables that make up a random vector (X_1, X_2) already occurs and takes a value. And we might want to compute the probability of the occurrence of the other random variable given this conditional information. For this all we need to do is extend the idea of conditional probabilities to \mathbb{R}^2 -valued random variables as defined below.

Definition 39 (Conditional PDF or PMF) Let (X_1, X_2) be a discrete bivariate R.V. The conditional PMF of $X_1 | X_2 = x_2$, where $f_{X_2}(x_2) := \mathbf{P}(X_2 = x_2) > 0$ is:

$$f_{X_1|X_2}(x_1|x_2) := \mathbf{P}(X_1 = x_1 | X_2 = x_2) = \frac{\mathbf{P}(X_1 = x_1, X_2 = x_2)}{\mathbf{P}(X_2 = x_2)} = \frac{f_{X_1, X_2}(x_1, x_2)}{f_{X_2}(x_2)} .$$

Part 2:

Let's differentiate CF

$$\frac{d}{dt} \varphi_X(t) = \frac{d}{dt} (1 - \theta + \theta e^{it}) = \theta i \exp(it)$$

We get $E(X)$ by evaluating $\frac{d}{dt} \varphi_X(t)$ at $t = 0$ and dividing by i according to Equation (3.71) as follows:

$$E(X) = \frac{1}{i} \left[\frac{d}{dt} \varphi_X(t) \right]_{t=0} = \frac{1}{i} [\theta i \exp(it)]_{t=0} = \frac{1}{i} (\theta i \exp(i0)) = \theta .$$

Similarly from Equation (3.71) we can get $E(X^2)$ as follows:

$$\begin{aligned} E(X^2) &= \frac{1}{i^2} \left[\frac{d^2}{dt^2} \varphi_X(t) \right]_{t=0} = \frac{1}{i^2} \left[\frac{d}{dt} \frac{d}{dt} \varphi_X(t) \right]_{t=0} = \frac{1}{i^2} \left[\frac{d}{dt} \theta i \exp(it) \right]_{t=0} \\ &= \frac{1}{i^2} [\theta i^2 \exp(it)]_{t=0} = \frac{1}{i^2} (\theta i^2 \exp(i0)) = \theta . \end{aligned}$$

Finally, from the first and second moments we can get the variance as follows:

$$V(X) = E(X^2) - (E(X))^2 = \theta - \theta^2 = \theta(1 - \theta) .$$

Let's check that this is what we have as variance for the Bernoulli(θ) RV if we directly computed it using weighted sums in the definition of expectations: $E(X) = 1 \times \theta + 0 \times (1 - \theta) = \theta$, $E(X^2) = 1^2 \times \theta + 0^2 \times (1 - \theta) = \theta$ and thus giving the same $V(X) = E(X^2) - (E(X))^2 = \theta - \theta^2 = \theta(1 - \theta)$.

Example 102 Let X be an Exponential(λ) RV. First show that its CF is $\lambda/(\lambda - it)$. Then use CF to find $E(X)$, $E(X^2)$ and from this obtain the variance $V(X) = E(X^2) - (E(X))^2$.

Solution:

Recall that the PDF of an Exponential(λ) RV for a given parameter $\lambda \in (0, \infty)$ is $\lambda e^{-\lambda x}$ if $x \in [0, \infty)$ and 0 if $x \notin [0, \infty)$.

Part 1: Find the CF.

We will use the fact that

$$\int_0^\infty \alpha e^{-\alpha x} dx = [-e^{-\alpha x}]_0^\infty = 1$$

$$\begin{aligned} \varphi_X(t) &= E(\exp(itX)) = E(e^{itX}) = \int_{-\infty}^\infty e^{itx} \lambda e^{-\lambda x} dx = \lambda \int_0^\infty e^{-(\lambda - it)x} dx \\ &= \frac{\lambda}{\lambda - it} \int_0^\infty (\lambda - it)e^{-(\lambda - it)x} dx = \frac{\lambda}{\lambda - it} \int_0^\infty \alpha e^{-\alpha x} dx = \frac{\lambda}{\lambda - it} , \end{aligned}$$

where $\alpha = \lambda - it$ with $\lambda > 0$.

Alternatively, you can use $e^{itx} = \cos(tx) + i \sin(tx)$ and do integration by parts to arrive at the same answer starting from:

$$\varphi_X(t) = \int_{-\infty}^\infty e^{itx} \lambda e^{-\lambda x} dx = \int_{-\infty}^\infty \cos(tx) e^{-\lambda x} dx + i \int_{-\infty}^\infty \sin(tx) e^{-\lambda x} dx = \frac{\lambda}{\lambda - it} .$$

Part 2:

Classwork 94 (The number of Heads, given there is at least one Tails) Consider the following two questions.

$$\left\{ \begin{array}{l} \text{P}(Y = y) \\ \text{P}(Y' = y') \end{array} \right\} = \left\{ \begin{array}{l} \text{P}(Y = y) \\ \text{P}(Y' = y') \end{array} \right\}$$

your spare time.

Classwork 93 (Two random variables of toss a coin three, experiment) Describe the probability of the RV Y and Y' of Table 3.1 in terms of its PMF. Repeat the process for the RV Z in Part 1 Solution.

	Z(ω):							
	0	0	0	0	0	0	0	1
$Y'(\omega)$:	1	0	0	0	0	0	0	$Y' = X_1 X_2 X_3$
$Z(\omega)$:	0	1	1	2	1	2	2	$Z = (1 - X_1) + (1 - X_2) + (1 - X_3)$
$Y(\omega)$:	3	2	2	1	2	1	1	$Y = X_1 + X_2 + X_3$
$\text{P}(\omega)$:	$\frac{1}{8}$	$\frac{8}{1}$	$\frac{1}{8}$	$\frac{1}{8}$	$\frac{1}{8}$	$\frac{1}{8}$	$\frac{1}{8}$	$X, \text{Ber}_k(\frac{1}{2})$

Table 3.1: The 8 ω 's in the sample space of the experiment \mathcal{E}_3 are given in the first row above. The RV Y is the number of Heads, the RVs Y' and Z are the indicator functions of the event that all three tosses were Heads, and the event that all three tosses were Tails, respectively.

Let us consider a few discrete RVs for the simple coin tossing experiment \mathcal{E}_3 that build on the Bernoulli(θ) RV X_i for the i -th toss in an independent and identically distributed (IID) manner.

$$f_{X^2|X_1}(x_2|x_1) = \frac{f_{X^2}(x_2)}{\int_A f_{X^2|X_1}(x_2|x_1) dx_2}, \quad \text{P}(X^2 \in A|X_1 = x_1) = \int_A f_{X^2|X_1}(x_2|x_1) dx_2.$$

Similarly, if $f_{X^1}(x_1) > 0$, then the conditional PDF of $X^2|X_1 = x_1$ is:

$$f_{X^1|X^2}(x_1|x_2) = \frac{f_{X^1}(x_1)}{\int_A f_{X^1|X^2}(x_1|x_2) dx_1}, \quad \text{P}(X_1 \in A|X^2 = x_2) = \int_A f_{X^1|X^2}(x_1|x_2) dx_1.$$

of $X_1|X^2 = x_2$ is:

$$f_{X^1|X^2}(x_1|x_2) = \text{P}(X_1 = x_1|X^2 = x_2) = \frac{\text{P}(X_1 = x_1, X^2 = x_2)}{\text{P}(X_1 = x_1, X^2)}$$

Similarly, if $f_{X^1}(x_1) := \text{P}(X_1 = x_1) > 0$, then the conditional PMF of $X^2|X_1 = x_1$ is:

$$= \exp(\mu \times 0)(1 - \theta) + \exp(\mu \times 1)\theta = \exp(0)(1 - \theta) + \exp(\mu)\theta = 1 - \theta + \theta \exp(\mu)$$

$$\phi_X(t) = \mathbb{E}(\exp(tX)) = \sum_x \exp(tx)f(x; \theta) \quad \text{By Defn. in Equation (3.70)}$$

$$f_X(x; \theta) = \begin{cases} 0 & \text{otherwise,} \\ \theta & \text{if } x = 1 \\ 1 - \theta & \text{if } x = 0 \end{cases}$$

Recall the PMF for this discrete RV with parameter $\theta \in (0, 1)$ is

Part 1 Solution. Let X be the Bernoulli(θ) RV. Find the CF of X . Then use CF to find $E(X)$.

Example 101 Let X be the Bernoulli(θ) RV. Find the CF of X . Then use CF to find $E(X)$, $E(X^2)$ and from this obtain the variance $V(X) = E(X^2) - (E(X))^2$.

Proof: For proof see e.g., Ushakov, N. G. (1999) Selected topics in characteristic functions, VSP p. 39).

$$\text{where } \left[\frac{d^k \phi_X(t)}{dt^k} \right]_{t=0} \text{ is the } k\text{-th derivative of } \phi_X(t) \text{ with respect to } t, \text{ evaluated at the point } t = 0. \quad (3.71)$$

In both cases,

2. if k is odd, the n -th moment of X exists and is finite for any $0 \leq n \leq k - 1$.

1. if k is even, the n -th moment of X exists and is finite for any $0 \leq n \leq k$:

If $\phi_X(t)$ is k times differentiable at the point $t = 0$, then

Proposition 49 (Moments from CF). Let X be a random variable and $\phi_X(t)$ be its CF.

The above Theorem gives us the relationship between the moments and the derivatives of the CF if we already know that the moment exists. When one wants to compute a moment of a random variable, what we need is the following Theorem.

This completes the sketch of the proof.

$$\left[\frac{d^k \phi_X(t)}{dt^k} \right]_{t=0} = \int_0^{\infty} \mathbb{E}(X^k) \exp(itX) dt = \left[\left(\frac{d^k \phi_X(t)}{dt^k} \right) \mathbb{E}(X^k) \right]_{t=0}.$$

The RHS evaluated at $t = 0$ is

$$\left(\frac{d^k \phi_X(t)}{dt^k} \right) \mathbb{E}(X^k) = \left((X^k) \frac{d^k \phi_X(t)}{dt^k} \right) \mathbb{E}(X^k) = \left((X^k) \frac{d^k \exp(itX)}{dt^k} \right) \mathbb{E}(X^k) = \left((X^k) \frac{d^k}{dt^k} \exp(itX) \right) \mathbb{E}(X^k) = \left(\frac{d^k}{dt^k} \exp(itX) \right) \mathbb{E}(X^k) = \frac{d^k}{dt^k} (\exp(itX)) \mathbb{E}(X^k) = \frac{d^k}{dt^k} t^k \mathbb{E}(X^k) = k! \mathbb{E}(X^k).$$

the linearity of the expectation (integral) and the derivative operators, we can change the order of operations:

The proper proof is very messy so we just give a sketch of the ideas in the proof. Due to

the linearity of the expectation (integral) and the derivative operators, we can change the order of

operations:

the proper proof is very messy so we just give a sketch of the ideas in the proof. Due to

the linearity of the expectation (integral) and the derivative operators, we can change the order of

operations:

the proper proof is very messy so we just give a sketch of the ideas in the proof. Due to

the linearity of the expectation (integral) and the derivative operators, we can change the order of

operations:

the proper proof is very messy so we just give a sketch of the ideas in the proof. Due to

the linearity of the expectation (integral) and the derivative operators, we can change the order of

operations:

the proper proof is very messy so we just give a sketch of the ideas in the proof. Due to

the linearity of the expectation (integral) and the derivative operators, we can change the order of

operations:

the proper proof is very messy so we just give a sketch of the ideas in the proof. Due to

the linearity of the expectation (integral) and the derivative operators, we can change the order of

operations:

the proper proof is very messy so we just give a sketch of the ideas in the proof. Due to

the linearity of the expectation (integral) and the derivative operators, we can change the order of

operations:

the proper proof is very messy so we just give a sketch of the ideas in the proof. Due to

the linearity of the expectation (integral) and the derivative operators, we can change the order of

operations:

the proper proof is very messy so we just give a sketch of the ideas in the proof. Due to

the linearity of the expectation (integral) and the derivative operators, we can change the order of

operations:

the proper proof is very messy so we just give a sketch of the ideas in the proof. Due to

the linearity of the expectation (integral) and the derivative operators, we can change the order of

operations:

the proper proof is very messy so we just give a sketch of the ideas in the proof. Due to

the linearity of the expectation (integral) and the derivative operators, we can change the order of

operations:

the proper proof is very messy so we just give a sketch of the ideas in the proof. Due to

the linearity of the expectation (integral) and the derivative operators, we can change the order of

operations:

the proper proof is very messy so we just give a sketch of the ideas in the proof. Due to

the linearity of the expectation (integral) and the derivative operators, we can change the order of

operations:

the proper proof is very messy so we just give a sketch of the ideas in the proof. Due to

the linearity of the expectation (integral) and the derivative operators, we can change the order of

operations:

the proper proof is very messy so we just give a sketch of the ideas in the proof. Due to

the linearity of the expectation (integral) and the derivative operators, we can change the order of

operations:

the proper proof is very messy so we just give a sketch of the ideas in the proof. Due to

the linearity of the expectation (integral) and the derivative operators, we can change the order of

operations:

the proper proof is very messy so we just give a sketch of the ideas in the proof. Due to

the linearity of the expectation (integral) and the derivative operators, we can change the order of

operations:

the proper proof is very messy so we just give a sketch of the ideas in the proof. Due to

the linearity of the expectation (integral) and the derivative operators, we can change the order of

operations:

the proper proof is very messy so we just give a sketch of the ideas in the proof. Due to

the linearity of the expectation (integral) and the derivative operators, we can change the order of

operations:

the proper proof is very messy so we just give a sketch of the ideas in the proof. Due to

the linearity of the expectation (integral) and the derivative operators, we can change the order of

operations:

the proper proof is very messy so we just give a sketch of the ideas in the proof. Due to

the linearity of the expectation (integral) and the derivative operators, we can change the order of

operations:

the proper proof is very messy so we just give a sketch of the ideas in the proof. Due to

the linearity of the expectation (integral) and the derivative operators, we can change the order of

operations:

the proper proof is very messy so we just give a sketch of the ideas in the proof. Due to

the linearity of the expectation (integral) and the derivative operators, we can change the order of

operations:

the proper proof is very messy so we just give a sketch of the ideas in the proof. Due to

the linearity of the expectation (integral) and the derivative operators, we can change the order of

operations:

the proper proof is very messy so we just give a sketch of the ideas in the proof. Due to

the linearity of the expectation (integral) and the derivative operators, we can change the order of

operations:

the proper proof is very messy so we just give a sketch of the ideas in the proof. Due to

the linearity of the expectation (integral) and the derivative operators, we can change the order of

operations:

the proper proof is very messy so we just give a sketch of the ideas in the proof. Due to

the linearity of the expectation (integral) and the derivative operators, we can change the order of

operations:

the proper proof is very messy so we just give a sketch of the ideas in the proof. Due to

the linearity of the expectation (integral) and the derivative operators, we can change the order of

operations:

the proper proof is very messy so we just give a sketch of the ideas in the proof. Due to

the linearity of the expectation (integral) and the derivative operators, we can change the order of

operations:

the proper proof is very messy so we just give a sketch of the ideas in the proof. Due to

the linearity of the expectation (integral) and the derivative operators, we can change the order of

operations:

the proper proof is very messy so we just give a sketch of the ideas in the proof. Due to

the linearity of the expectation (integral) and the derivative operators, we can change the order of

operations:

the proper proof is very messy so we just give a sketch of the ideas in the proof. Due to

the linearity of the expectation (integral) and the derivative operators, we can change the order of

operations:

the proper proof is very messy so we just give a sketch of the ideas in the proof. Due to

the linearity of the expectation (integral) and the derivative operators, we can change the order of

operations:

the proper proof is very messy so we just give a sketch of the ideas in the proof. Due to

the linearity of the expectation (integral) and the derivative operators, we can change the order of

operations:

the proper proof is very messy so we just give a sketch of the ideas in the proof. Due to

the linearity of the expectation (integral) and the derivative operators, we can change the order of

operations:

the proper proof is very messy so we just give a sketch of the ideas in the proof. Due to

the linearity of the expectation (integral) and the derivative operators, we can change the order of

operations:

the proper proof is very messy so we just give a sketch of the ideas in the proof. Due to

the linearity of the expectation (integral) and the derivative operators, we can change the order of

operations:

the proper proof is very messy so we just give a sketch of the ideas in the proof. Due to

the linearity of the expectation (integral) and the derivative operators, we can change the order of

operations:

the proper proof is very messy so we just give a sketch of the ideas in the proof. Due to

the linearity of the expectation (integral) and the derivative operators, we can change the order of

operations:

the proper proof is very messy so we just give a sketch of the ideas in the proof. Due to

the linearity of the expectation (integral) and the derivative operators, we can change the order of

operations:

the proper proof is very messy so we just give a sketch of the ideas in the proof. Due to

the linearity of the expectation (integral) and the derivative operators, we can change the order of

operations:

the proper proof is very messy so we just give a sketch of the ideas in the proof. Due to

the linearity of the expectation (integral) and the derivative operators, we can change the order of

operations:

the proper proof is very messy so we just give a sketch of the ideas in the proof. Due to

the linearity of the expectation (integral) and the derivative operators, we can change the order of

operations:

the proper proof is very messy so we just give a sketch of the ideas in the proof. Due to

the linearity of the expectation (integral) and the derivative operators, we can change the order of

operations:

the proper proof is very messy so we just give a sketch of the ideas in the proof. Due to

the linearity of the expectation (integral) and the derivative operators, we can change the order of

operations:

the proper proof is very messy so we just give a sketch of the ideas in the proof. Due to

the linearity of the expectation (integral) and the derivative operators, we can change the order of

operations:

the proper proof is very messy so we just give a sketch of the ideas in the proof. Due to

the linearity of the expectation (integral) and the derivative operators, we can change the order of

operations:

the proper proof is very messy so we just give a sketch of the ideas in the proof. Due to

the linearity of the expectation (integral) and the derivative operators, we can change the order of

operations:

the proper proof is very messy so we just give a sketch of the ideas in the proof. Due to

the linearity of the expectation (integral) and the derivative operators, we can change the order of

operations:

the proper proof is very messy so we just give a sketch of the ideas in the proof. Due to

the linearity of the expectation (integral) and the derivative operators, we can change the order of

operations:

1. What is conditional probability $\mathbf{P}(Y|Y' = 0)$?

$\mathbf{P}(Y = y Y' = 0)$	$= \frac{\mathbf{P}(Y=y, Y'=0)}{\mathbf{P}(Y'=0)}$	$= \frac{\mathbf{P}(\{\omega: Y(\omega)=y \cap Y'(\omega)=0\})}{\mathbf{P}(\{\omega: Y'(\omega)=0\})}$	$= ?$
$\mathbf{P}(Y = 0 Y' = 0)$	$\frac{\mathbf{P}(Y=0, Y'=0)}{\mathbf{P}(Y'=0)}$	$\frac{\frac{1}{8}}{\frac{1}{8} + \frac{1}{8} + \frac{1}{8} + \frac{1}{8} + \frac{1}{8} + \frac{1}{8}}$	$\frac{1}{6}$
$\mathbf{P}(Y = 1 Y' = 0)$	$\frac{\mathbf{P}(Y=1, Y'=0)}{\mathbf{P}(Y'=0)}$	$\frac{\frac{1}{8} + \frac{1}{8} + \frac{1}{8}}{\frac{1}{8} + \frac{1}{8} + \frac{1}{8} + \frac{1}{8} + \frac{1}{8} + \frac{1}{8}}$	$\frac{3}{6}$
$\mathbf{P}(Y = 2 Y' = 0)$	$\frac{\mathbf{P}(Y=2, Y'=0)}{\mathbf{P}(Y'=0)}$	$\frac{\frac{1}{8} + \frac{1}{8} + \frac{1}{8}}{\frac{1}{8} + \frac{1}{8} + \frac{1}{8} + \frac{1}{8} + \frac{1}{8} + \frac{1}{8}}$	$\frac{3}{6}$
$\mathbf{P}(Y = 3 Y' = 0)$	$\frac{\mathbf{P}(Y=3, Y'=0)}{\mathbf{P}(Y'=0)}$	$\frac{\mathbf{P}(\emptyset)}{\frac{1}{8} + \frac{1}{8} + \frac{1}{8} + \frac{1}{8} + \frac{1}{8} + \frac{1}{8}}$	0
$\mathbf{P}(Y \in \{0, 1, 2, 3\} Y' = 0)$	$\frac{\sum_{y=0}^3 \mathbf{P}(Y=y, Y'=0)}{\mathbf{P}(Y'=0)}$	$\frac{\frac{1}{8} + \frac{1}{8} + \frac{1}{8} + \frac{1}{8}}{\frac{1}{8} + \frac{1}{8} + \frac{1}{8} + \frac{1}{8} + \frac{1}{8} + \frac{1}{8}}$	1

2. What is $\mathbf{P}(Y|Y' = 1)$?

$$\mathbf{P}(Y = y|Y' = 1) = \begin{cases} 1 & \text{if } y = 3 \\ 0 & \text{otherwise} \end{cases}$$

3.10.3 \mathbb{R}^m -valued Random Variables

Consider the RV \vec{X} whose components are the RVs X_1, X_2, \dots, X_m , i.e., $X := (X_1, X_2, \dots, X_m)$, where $m \geq 2$. A particular realization of this RV is a point (x_1, x_2, \dots, x_m) in \mathbb{R}^m . Now, let us extend the notions of JCDF, JPMF and JPDF to \mathbb{R}^m .

Definition 40 (multivariate JDF) The **joint distribution function (JDF)** or **joint cumulative distribution function (JCDF)**, $F_{X_1, X_2, \dots, X_m}(x_1, x_2, \dots, x_m) : \mathbb{R}^m \rightarrow [0, 1]$, of the multivariate random vector (X_1, X_2, \dots, X_m) is

$$\begin{aligned} F_{X_1, X_2, \dots, X_m}(x_1, x_2, \dots, x_m) &= \mathbf{P}(X \leq x_1 \cap X_2 \leq x_2 \cap \dots \cap X_m \leq x_m) \\ &= \mathbf{P}(X_1 \leq x_1, X_2 \leq x_2, \dots, X_m \leq x_m) \\ &= P(\{\omega: X_1(\omega) \leq x_1, X_2(\omega) \leq x_2, \dots, X_m(\omega) \leq x_m\}), \end{aligned} \quad (3.63)$$

for any $(x_1, x_2, \dots, x_m) \in \mathbb{R}^m$, where the right-hand side represents the probability that the random vector (X_1, X_2, \dots, X_m) takes on a value in $\{(x'_1, x'_2, \dots, x'_m) : x'_1 \leq x_1, x'_2 \leq x_2, \dots, x'_m \leq x_m\}$, the set of points in \mathbb{R}^m that are less than the point (x_1, x_2, \dots, x_m) in each coordinate $1, 2, \dots, m$.

The JDF $F_{X_1, X_2, \dots, X_m}(x_1, x_2, \dots, x_m) : \mathbb{R}^m \rightarrow \mathbb{R}$ satisfies the following conditions to remain a probability:

1. $0 \leq F_{X_1, X_2, \dots, X_m}(x_1, x_2, \dots, x_m) \leq 1$
2. $F_{X_1, X_2, \dots, X_m}(x_1, x_2, \dots, x_m)$ is an increasing function of x_1, x_2, \dots and x_m
3. $F_{X_1, X_2, \dots, X_m}(x_1, x_2, \dots, x_m) \rightarrow 1$ as $x_1 \rightarrow \infty, x_2 \rightarrow \infty, \dots$ and $x_m \rightarrow \infty$

3.12 Characteristic Functions

The characteristic function (CF) of a random variable gives another way to specify its distribution. Thus CF is a powerful tool for analytical results involving random variables (more).

Definition 47 (Characteristic Function (CF)) Let X be a RV and $i = \sqrt{-1}$. The function $\varphi_X(t) : \mathbb{R} \rightarrow \mathbb{C}$ defined by

$$\boxed{\varphi_X(t) := E(\exp(itX)) = \begin{cases} \sum_x \exp(itx) f_X(x) & \text{if } X \text{ is discrete RV} \\ \int_{-\infty}^{\infty} \exp(itx) f_X(x) dx & \text{if } X \text{ is continuous RV} \end{cases}} \quad (3.70)$$

is called the **characteristic function** of X .

NOTE: $\varphi_X(t)$ exists for any $t \in \mathbb{R}$, because

$$\begin{aligned} \varphi_X(t) &= E(\exp(itX)) \\ &= E(\cos(tx) + i \sin(tx)) \\ &= E(\cos(tx)) + iE(\sin(tx)) \end{aligned}$$

and the last two expected values are well-defined, because the sine and cosine functions are bounded by $[-1, 1]$.

For a continuous RV, $\int_{-\infty}^{\infty} \exp(-itx) f_X(x) dx$ is called the *Fourier transform* of f_X . This is the CF but with t replaced by $-t$. You will also encounter Fourier transforms when solving differential equations.

3.12.1 Obtaining Moments from Characteristic Function

Recall that the k -th moment of X is $E(X^k)$ for any $k \in \mathbb{N} := \{1, 2, 3, \dots\}$ is

$$\boxed{E(X^k) = \begin{cases} \sum_x x^k f_X(x) & \text{if } X \text{ is a discrete RV} \\ \int_{-\infty}^{\infty} x^k f_X(x) dx & \text{if } X \text{ is a continuous RV} \end{cases}}$$

The characteristic function can be used to derive the moments of X due to the following nice relationship between the the k -th moment of X and the k -th derivative of the CF of X .

Proposition 48 (Moment & CF.) Let X be a random variable and $\varphi_X(t)$ be its CF. If $E(X^k)$ exists and is finite, then $\varphi_X(t)$ is k times continuously differentiable and

$$E(X^k) = \frac{1}{i^k} \left[\frac{d^k \varphi_X(t)}{dt^k} \right]_{t=0}.$$

where $\left[\frac{d^k \varphi_X(t)}{dt^k} \right]_{t=0}$ is the k -th derivative of $\varphi_X(t)$ with respect to t , evaluated at the point $t = 0$.

$${}^{(u_i x)} {}^{u_i} X_{\tilde{J}} \cdots {}^{(\tilde{e}_i x)} {}^{\tilde{e}_i} X_{\tilde{J}} {}^{(l_i x)} {}^{l_i} X_{\tilde{J}} = {}^{(u_i x, \dots, \tilde{e}_i x, l_i x)} {}^{u_i} X_{\tilde{J}} \cdots {}^{\tilde{e}_i} X_{\tilde{J}} {}^{l_i} X_{\tilde{J}}$$

or equivalently,

$$(\exists_i x \succ \exists_i X)D \wedge \dots \wedge (\exists_i x \succ \exists_i X)D \wedge (\exists_i x \succ \exists_i X)D = (\exists_i x \succ \exists_i X \wedge \dots \wedge \exists_i x \succ \exists_i X \wedge \exists_i x \succ \exists_i X)D$$

is said to be independent or independently distributed if and only if

Definition 43 (Independent sequence of Sequences of RVs) A finite or infinite sequence of RVs X_1, X_2, \dots

The marginal PDF (marginal PMF) is obtained by integrating (summing) the PDF of X_1 over all other random variables. For example, the marginal PDF of X_1 is

2. is a non-negative function, i.e., $f_{X_1, X_2, \dots, X_m}(x_1, x_2, \dots, x_m) \geq 0$.

1. integrates to 1, i.e., $\int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} f_{X_1, X_2, \dots, X_m}(x_1, x_2, \dots, x_m) dx_1 dx_2 \cdots dx_m = 1$

The JPDF satisfies the following two properties:

$$(99.\mathfrak{E}) \quad \cdot \quad \boxed{\mu_{xp} \cdots \varepsilon_{xp} \iota_{xp} (\mu_{x_1} x) \mu_{X' X} \cdots \varepsilon_{X' X} f^B \int \int \cdots \int} = (B)\mathbf{d}$$

pure

from $\text{JPDF}_{X_1, X_2, \dots, X_m}$ we can compute the $\text{JPDF}_{X_1, X_2, \dots, X_m}$ at any point $(x_1, x_2, \dots, x_m) \in \mathbb{R}^m$ and more generally we can compute the probability of any event B , that can be cast as a region $\Omega_B \subset \mathbb{R}^m$, by "simply" taking m -dimensional integrals (you have done such iterated integrals when

$$Q_{x_1, x_2, \dots, x_m} = Q_{x_1, x_2, \dots, x_m}^m F_{X_1, X_2, \dots, X_m}(x_1, x_2, \dots, x_m)$$

Definition 42 (Multivariate PDF) (X_1, X_2, \dots, X_m) is a continuous random vector if its joint PDF $F_{X_1, X_2, \dots, X_m}(x_1, x_2, \dots, x_m)$ is differentiable and the joint probability density function (PDF) is given by:

³⁹ By simply taking sums as in Equation (3.59) but now over all m coordinates.

$$\text{Summe } \mathbf{P}(\mathcal{U}) = 1, \sum_{x_1, x_2, \dots, x_m \in S} f_{X_1, X_2, \dots, X_m}(x_1, x_2, \dots, x_m) = 1.$$

$$(\mathbf{F}_0, \mathbf{C}) = (\mathbf{u}x, \mathbf{w}_Y, \dots, \mathbf{z}x, \mathbf{z}_Y, \mathbf{l}_Y) \cdot \mathbf{A} = (\mathbf{u}x, \dots, \mathbf{z}x, \mathbf{l}_x) \mathbf{w}_X \cdot (\mathbf{z}_X, \mathbf{l}_X)$$

JPMF) is:

Deinition 41 (Multivariate JPMF): If (X_1, X_2, \dots, X_m) is a discrete random vector that takes values in a discrete support set S_{X_1, X_2, \dots, X_m} , then its joint probability mass function (or

4. $F(x_1, x_2, \dots, x_m) \leftarrow 0$ as $x_1 \leftarrow -\infty, x_2 \leftarrow -\infty, \dots$ and $x^m \rightarrow \infty$

Ex. 3.49 Let X_1, X_2, X_3, X_4 be RVs that denote the number of bits received in a digital channel that are classified as *excellent*, *good*, *fair* and *poor*, respectively. In a transmission of 10 bits, what is the probability that at least 6 of the bits received are *excellent*, *good*, *fair* and *poor*, respectively. The number of bits received in a digital channel of n bits is the sum of the number of bits received in each bit being *excellent*, *good*, *fair* and *poor* respectively. The probabilities of each bit being *excellent*, *good*, *fair* and *poor* are $0.3, \theta_2 = 0.6, \theta_3 = 0.08$, and $\theta_4 = 0.02$ respectively. Think of Multinomial($n = 10, \theta_1 = 0.3, \theta_2 = 0.6, \theta_3 = 0.08, \theta_4 = 0.02$) as a model for bit classification in this digital channel.]

Ex. 3.48 — Soft drink cans are filled by an automated filling machine. Assume the full volumes of the cans are independent Normal(121.0, 0.01) RVs. What is the probability that the average volume of ten cans selected from this process is less than 120.1 fluid ounces?

Ex. 3.47 — Suppose the HVs X_1 , X_2 and X_3 represent the thicknesses in micrometres of a substrate, an active layer, and a coating layer of a chemical product. Assume X_1 , X_2 and X_3 are $N_{Normal}(10000, 250^2)$, $N_{Normal}(1000, 20^2)$ and $N_{Normal}(80, 4^2)$ RVs, respectively. Further suppose that they are independent. The required specifications for the thicknesses of the substrate, active layer and coating layer are $[9500, 10500]$, $[75, 85]$, respectively. What proportion of the semicircular products meets all thickness specifications? Hint: this is just $P(9500 < X_1 < 10500, 75 < X_2 < 85 < X_3 < 85)$. Which one of the three thickmesses has the least probability of meeting specifications?

What is the probability that the device operates for more than 1000 hours without any failures? [Hint: The requested probability is $P(X_1 > 1000, X_2 > 1000, X_3 > 1000, X_4 > 1000)$ since each one of the four components of the device must not fail before 1000 hours.]

$$0 \times 10^{-12}e^{-0.001x_1 - 0.002x_2 - 0.0015x_3 - 0.003x_4} \begin{cases} 0 & \text{otherwise} \\ \frac{1}{x_1 x_2 x_3 x_4} & \text{if } x_1 \geq 0, x_2 \geq 0, x_3 \geq 0, x_4 \geq 0 \end{cases} =$$

Ex. 3.46 — In an electronic assembly, let the RVs X_1, X_2, X_3, X_4 denote the lifetimes of four components in hours. Suppose that the PDF of these variables is

Are X and Y independent?

$$f_{X,Y}(x,y) = \begin{cases} 0 & \text{otherwise} \\ e^{-x} & \text{if } x \in [0, \infty) \text{ and } y \in [2, 3] \end{cases}$$

Ex. 3.45 — Let (X, Y) be a continuous RV with joint probability density function (JPDF):

for any distinct subset of indices $\{i_1, i_2, \dots, i_m\}$ of $\{1, 2, \dots\}$, the index set of the sequence of RVs and any sequence of real numbers $x_{i_1}, x_{i_2}, \dots, x_{i_m}$.

By the above definition, the sequence of **discrete** RVs X_1, X_2, \dots taking values in an at most countable set \mathbb{D} are said to be independently distributed if for any distinct subset of indices $\{i_1, i_2, \dots, i_k\}$ such that the corresponding RVs $X_{i_1}, X_{i_2}, \dots, X_{i_k}$ exists as a distinct subset of our original sequence of RVs X_1, X_2, \dots and for any elements $x_{i_1}, x_{i_2}, \dots, x_{i_k}$ in \mathbb{D} , the following equality is satisfied:

$$\mathbf{P}(X_{i_1} = x_{i_1}, X_{i_2} = x_{i_2}, \dots, X_{i_k} = x_{i_k}) = \mathbf{P}(X_{i_1} = x_{i_1})\mathbf{P}(X_{i_2} = x_{i_2}) \cdots \mathbf{P}(X_{i_k} = x_{i_k}),$$

or equivalently,

$$f_{X_{i_1}, X_{i_2}, \dots, X_{i_k}}(x_{i_1}, x_{i_2}, \dots, x_{i_k}) = f_{X_{i_1}}(x_{i_1})f_{X_{i_2}}(x_{i_2}) \cdots f_{X_{i_k}}(x_{i_k}).$$

From Definition 43, we say m random variables X_1, X_2, \dots, X_m are jointly independent or mutually independent if and only if for every $(x_1, x_2, \dots, x_m) \in \mathbb{R}^m$

$$F_{X_1, X_2, \dots, X_m}(x_1, x_2, \dots, x_m) = F_{X_1}(x_1)F_{X_2}(x_2) \cdots F_{X_m}(x_m), \quad (3.67)$$

$$f_{X_1, X_2, \dots, X_m}(x_1, x_2, \dots, x_m) = f_{X_1}(x_1)f_{X_2}(x_2) \cdots f_{X_m}(x_m). \quad (3.68)$$

Proposition 44 (Conditional probability of independent sequence of RVs) For an independent sequence of RVs $\{X_1, X_2, \dots\}$, we have

$$\mathbf{P}(X_{i+1} \leq x_{i+1} | X_i \leq x_i, X_{i-1} \leq x_{i-1}, \dots, X_1 \leq x_1) = \mathbf{P}(X_{i+1} \leq x_{i+1}) \quad (3.69)$$

Proof:

$$\begin{aligned} & \mathbf{P}(X_{i+1} \leq x_{i+1} | X_i \leq x_i, X_{i-1} \leq x_{i-1}, \dots, X_1 \leq x_1) \\ &= \frac{\mathbf{P}(X_{i+1} \leq x_{i+1}, X_i \leq x_i, X_{i-1} \leq x_{i-1}, \dots, X_1 \leq x_1)}{\mathbf{P}(X_i \leq x_i, X_{i-1} \leq x_{i-1}, \dots, X_1 \leq x_1)} \\ &= \frac{\mathbf{P}(X_{i+1} \leq x_{i+1})\mathbf{P}(X_i \leq x_i)\mathbf{P}(X_{i-1} \leq x_{i-1}) \cdots \mathbf{P}(X_1 \leq x_1)}{\mathbf{P}(X_i \leq x_i)\mathbf{P}(X_{i-1} \leq x_{i-1}) \cdots \mathbf{P}(X_1 \leq x_1)} \\ &= \mathbf{P}(X_{i+1} \leq x_{i+1}) \end{aligned}$$

Equation (3.69) simply says that the conditional distribution of the RV X_{i+1} given all previous RVs X_i, X_{i-1}, \dots, X_1 is simply determined by the distribution of X_{i+1} .

Example 95 If X_1 and X_2 are independent random variables then what is their covariance $\mathbf{Cov}(X_1, X_2)$?

Solution:

We know for independent RVs from the properties of expectations that

$$E(X_1 X_2) = E(X_1)E(X_2)$$

From the formula for covariance

$$\begin{aligned} \mathbf{Cov}(X_1, X_2) &= E(X_1 X_2) - E(X_1)E(X_2) \\ &= E(X_1)E(X_2) - E(X_1)E(X_2) \quad \text{due to independence} \\ &= 0 \end{aligned}$$

Remark 45 The converse is not true: two random variables that have zero covariance are not necessarily independent.

7. $E(X)$, the expectation of X or the first moment of X
8. $E(Y)$, the expectation of Y or the first moment of Y
9. $E(XY)$, the expectation of XY
10. $\mathbf{Cov}(X, Y) = E(XY) - E(X)E(Y)$, the covariance of X and Y .

Ex. 3.40 — Logs are milled to have a width of μ . The actual width of a randomly selected item is X . If X is a $\text{Normal}(\mu, \sigma^2)$ random variable then find the probability density function of the *squared-error* of the milling process,

$$Y = (X - \mu)^2.$$

Ex. 3.41 — Let (X, Y) be a discrete random vector ($\text{R}\vec{V}$) with support:

$$\mathcal{S}_{X,Y} = \{(0,0), (0,1), (1,0), (1,1)\}.$$

Let its joint probability mass function (JPMF) be:

$$f_{X,Y}(x, y) = \begin{cases} \frac{1}{4} & \text{if } (x, y) = (0, 0) \\ \frac{1}{4} & \text{if } (x, y) = (0, 1) \\ \frac{1}{4} & \text{if } (x, y) = (1, 0) \\ \frac{1}{4} & \text{if } (x, y) = (1, 1) \\ 0 & \text{otherwise} \end{cases}$$

Are X and Y independent?

Ex. 3.42 — A semiconductor product consists of three layers that are fabricated independently. If the variances in thickness of the first, second and third layers are 25, 40 and 30 nanometers squared, what is the variance of the thickness of the final product?

Ex. 3.43 — Find the covariance for the discrete $\text{R}\vec{V} (X, Y)$ with joint probability mass function

$$f_{X,Y}(x, y) = \begin{cases} 0.2 & \text{if } (x, y) = (0, 0) \\ 0.1 & \text{if } (x, y) = (1, 1) \\ 0.1 & \text{if } (x, y) = (1, 2) \\ 0.1 & \text{if } (x, y) = (2, 1) \\ 0.1 & \text{if } (x, y) = (2, 2) \\ 0.4 & \text{if } (x, y) = (3, 3) \\ 0 & \text{otherwise} \end{cases}$$

[Hint: Recall that $\mathbf{Cov}(X, Y) = E(XY) - E(X)E(Y)$]

Ex. 3.44 — Consider two random variables (RVs) X and Y having marginal distribution functions

$$F_X(x) = \begin{cases} 0 & \text{if } x < 1 \\ \frac{1}{2} & \text{if } 1 \leq x < 2 \\ 1 & \text{if } x \geq 2 \end{cases}$$

$$F_Y(y) = \begin{cases} 0 & \text{if } y < 0 \\ 1 - \frac{1}{2}e^{-y} - \frac{1}{2}e^{-2y} & \text{if } y \geq 0 \end{cases}$$

If X and Y are independent, what is their joint distribution function $F_{X,Y}(x, y)$? [Hint: you need to express $F_{X,Y}(x, y)$ for any $(x, y) \in \mathbb{R}^2$.]

$$\left(\frac{0.1\lambda}{6} < \frac{0.1\lambda}{6 - \lambda} \right) d = (6 - 0 < 6 - \lambda)d = (0 < \lambda)d =$$

$$\left(\frac{0.1\lambda}{6} > Z \right) d =$$

$$\left(\frac{0.1\lambda}{6 - \lambda} > Z \right) d =$$

4. Let $U = 6 - 2Z + X - Y$ and we know U is $\text{Normal}(9, 10)$ RV.

$$\begin{aligned}
 &= \text{Normal}(6 + 0 + 2 + 1, 4 + 2) \\
 &= \text{Normal}(6 + (-2 \times 0) + (1 \times 2) + (-1 \times -1), ((-2)^2 \times 1) + (1^2 \times \frac{1}{4}) + ((-1)^2 \times 2)) \\
 &= \text{Normal}(9, 10)
 \end{aligned}$$

3. From the special property of normal RVs, the distribution of $6 - 2Z + X - Y$ is

$$A(2Y - 3Z) = 2\zeta V(Y) + (-3\zeta) A(Z) = (4 \times 2) + (9 \times 1) = 8 + 9 = 17$$

$$E(3X - 2Y + 4Z) = 3E(X) - 2E(Y) + 4E(Z) = (3 \times 2) + (-2 \times (-1)) + 4 \times 0 = 6 + 2 + 0 = 8$$

1. $E(3X - 2Y + 4Z)$
2. $V(2Y - 3Z)$
3. the distribution of $6 - 2Z + X - Y$
4. the probability that $6 - 2Z + X - Y < 0$
5. $\text{Cov}(X, W)$, where $W = X - Y$.

Example 96 Let A be $\text{NotUnit}(z, t)$, I be $\text{NotUnit}(-1, z)$ and Z be $\text{NotUnit}(0, 1)$. This illustrates the following:

We can get the following special property of normal RVs using Eqn. (3.14). If X_1, X_2, \dots, X_m be jointly independent RVs, where X_i is Normal(μ_i, σ_i^2), for $i = 1, 2, \dots, m$ then $Y = c + \sum_{i=1}^m a_i X_i$ for some constants c, a_1, a_2, \dots, a_m is the Normal($c + \sum_{i=1}^m a_i \mu_i, \sum_{i=1}^m a_i^2 \sigma_i^2$) RV.

Linerar Combination of Independent Normal RVs is a Normal RV

129

Ex. 3.39 — Let (X, Y) be a continuous RV with joint probability density function (jPDF)

$$\left\{ \begin{array}{ll} 0 & \text{otherwise} \\ a(x^2 + y) & \text{if } 0 > x > 1 \text{ and } 0 > y > 1 \end{array} \right\} = (h, x) X' X f$$

Ex. 3.38 — Find the probability that one of the three bins in a tracing signat., that are assumed to have independent life-times (i.e., the time during which they are operational), need to be replaced during the first 1200 hours of operation if the length of time before a single bulb needs to be replaced is a continuous random variable X with density

$$f(x) = \begin{cases} 6(0.25 - (x - 1.5)^2) & 1 < x < 2 \\ 0 & \text{otherwise} \end{cases}$$

Note: X is measured in multiples of 1000 hours.

3.11 Exercises in Multivariate Random Variables

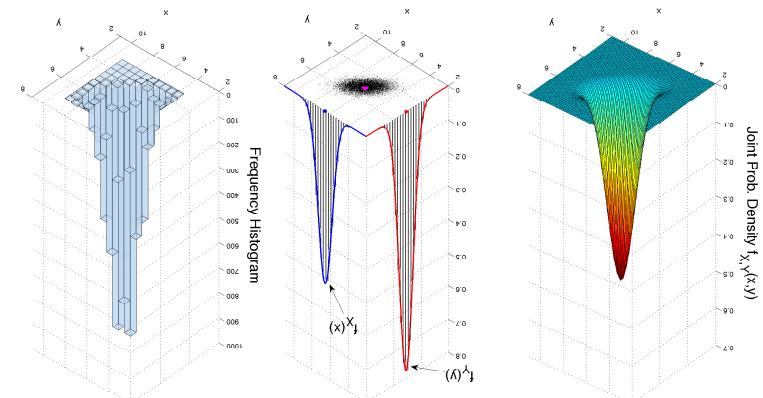


Figure 3.23: PDFs and Marginal PDFs of a Bi-variate Normal RV for lengths of grits of cylindrical shafts in a manufacturing process (in cm).

3.10.4 Some Common \mathbb{R}^m -valued RVs

So far, we have treated our random vectors as random points in \mathbb{R}^m and not been explicit about whether they are row or column vectors. We need to be more explicit now in order to perform arithmetic operations and transformations with them.

Let $X = (X_1, X_2, \dots, X_{m_X})$ be a \vec{X} in $\mathbb{R}^{1 \times m_X}$, i.e., X is a random row vector with 1 row and m_X columns, with JCDF $F_{X_1, X_2, \dots, X_{m_X}}$ and JPDF $f_{X_1, X_2, \dots, X_{m_X}}$. Similarly, let $Y = (Y_1, Y_2, \dots, Y_{m_Y})$ be a \vec{Y} in $\mathbb{R}^{1 \times m_Y}$, i.e., Y is a random row vector with 1 row and m_Y columns, with JCDF $F_{Y_1, Y_2, \dots, Y_{m_Y}}$ and JPDF $f_{Y_1, Y_2, \dots, Y_{m_Y}}$. Let the JCDF of the random vectors X and Y together be $F_{X_1, X_2, \dots, X_{m_X}, Y_1, Y_2, \dots, Y_{m_Y}}$ and JPDF be $f_{X_1, X_2, \dots, X_{m_X}, Y_1, Y_2, \dots, Y_{m_Y}}$.

Independent Random Vectors and their sums

The notion of mutual independence or joint independence of n random vectors is obtained similarly from ensuring the independence of any subset of the n vectors in terms of their JCDFs (JPMFs or JPDFs) being equal to the product of their marginal CDFs (PMFs or PDFs).

Thus, for a given $m_X < \infty$ and $m_Y < \infty$, two random vectors are independent if and only if for any $(x_1, x_2, \dots, x_{m_X}) \in \mathbb{R}^{1 \times m_X}$ and any $(y_1, y_2, \dots, y_{m_Y}) \in \mathbb{R}^{1 \times m_Y}$

$$\begin{aligned} F_{X_1, X_2, \dots, X_{m_X}, Y_1, Y_2, \dots, Y_{m_Y}}(x_1, x_2, \dots, x_{m_X}, y_1, y_2, \dots, y_{m_Y}) \\ = F_{X_1, X_2, \dots, X_{m_X}}(x_1, x_2, \dots, x_{m_X}) \times F_{Y_1, Y_2, \dots, Y_{m_Y}}(y_1, y_2, \dots, y_{m_Y}) \end{aligned}$$

or, equivalently

$$\begin{aligned} f_{X_1, X_2, \dots, X_{m_X}, Y_1, Y_2, \dots, Y_{m_Y}}(x_1, x_2, \dots, x_{m_X}, y_1, y_2, \dots, y_{m_Y}) \\ = f_{X_1, X_2, \dots, X_{m_X}}(x_1, x_2, \dots, x_{m_X}) \times f_{Y_1, Y_2, \dots, Y_{m_Y}}(y_1, y_2, \dots, y_{m_Y}) \end{aligned}$$

Let us consider the natural two-dimensional analogue of the Bernoulli(θ) RV in the real plane $\mathbb{R}^2 := (-\infty, \infty)^2 := (-\infty, \infty) \times (-\infty, \infty)$. A natural possibility is to use the **ortho-normal basis vectors** in \mathbb{R}^2 :

$$e_1 := (1, 0), \quad e_2 := (0, 1).$$

Recall that vector addition and subtraction are done component-wise, i.e. $(x_1, x_2) \pm (y_1, y_2) = (x_1 \pm y_1, x_2 \pm y_2)$. We introduce a useful function called the indicator function of a set, say A .

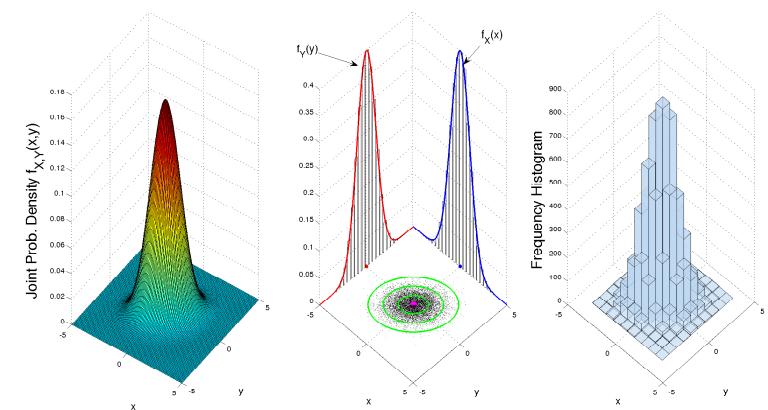
$$\mathbf{1}_A(x) = \begin{cases} 1 & \text{if } x \in A \\ 0 & \text{otherwise.} \end{cases}$$

$\mathbf{1}_A(x)$ returns 1 if x belongs to A and 0 otherwise.

Example 97 Let us recall the geometry and arithmetic of vector addition in the plane.

1. What is $(1, 0) + (1, 0)$, $(1, 0) + (0, 1)$, $(0, 1) + (0, 1)$?
2. What is the relationship between $(1, 0)$, $(0, 1)$ and $(1, 1)$ geometrically?
3. How does the diagonal of the parallelogram relate the its two sides in the geometry of addition in the plane?
4. What is $(1, 0) + (0, 1) + (1, 0)$?

Figure 3.22: JPDF, Marginal PDFs and Frequency Histogram of Bivariate Standard Normal \vec{RV} .



When we have a non-zero mean vector

$$\mu = \begin{pmatrix} \mu_X \\ \mu_Y \end{pmatrix} = \begin{pmatrix} 6.49 \\ 5.07 \end{pmatrix}$$

for the mean lengths and girths of cylindrical shafts from a manufacturing process with variance-covariance matrix

$$\Sigma = \begin{pmatrix} \text{Cov}(X, X) & \text{Cov}(X, Y) \\ \text{Cov}(Y, X) & \text{Cov}(Y, Y) \end{pmatrix} = \begin{pmatrix} V(X) & \text{Cov}(X, Y) \\ \text{Cov}(X, Y) & V(Y) \end{pmatrix} = \begin{pmatrix} 0.59 & 0.24 \\ 0.24 & 0.26 \end{pmatrix}$$

then the $\text{Normal}(\mu, \Sigma)$ \vec{RV} has JPDF, marginal PDFs and samples with frequency histograms as shown in Figure 3.23.

We can use MATLAB to compute for instance the probability that a cylinder has length and girth below 6.0 cms as follows:

```
>> mvncdf([6.0 6.0], [6.49 5.07], [0.59 0.24; 0.24 0.26])
ans =
    0.2615
```

Or find the probability (with numerical error tolerance) that the cylinders are within the rectangular specifications of 6 ± 1.0 along x and y as follows:

```
>> [F err] = mvncdf([5.0 5.0], [7.0 7.0], [6.49 5.07], [0.59 0.24; 0.24 0.26])
F =
    0.3352
err =
    1.0000e-08
```

3.10.5 Dependent Random Variables

When a sequence of RVs are not independent they are said to be **dependent**. The simplest form of dependence is *Markov dependence* that we will briefly see via a couple examples in Chapter 6.

context, this amounts to keeping track of the number of right and left turns made by the ball as it falls through a sequence of n tosses of a coin with probability θ of observing Heads. In the Q -matrix of an IID sequence of n tosses of a coin with probability θ of observing Heads and Tails out the coin-tossing context this can be thought of keeping track of the number of Heads and Tails. In fact, this is the underlying model and the bi in the Binomial(n, θ) does refer to two in Latin. In Remark 46 We can write the Binomial(n, θ) RV X as a Binomial(n, θ) RV $X = (Y, n - Y)$. In

$$\mathbb{E}^\theta(X) = \mathbb{E}^\theta((X_1, X_2)) = \sum_{(x_1, x_2) \in \{e_1, e_2\}} (x_1, x_2) f(x_1, x_2; \theta) = (1, 0)\theta + (0, 1)(1 - \theta) = (0, 1 - \theta).$$

Example 98 Let us find the Expectation of Bernoulli(θ) RV in Model 15.

$$f(x; \theta) = P(X = x) = \begin{cases} 0 & \text{otherwise} \\ \theta & \text{if } x = e_1 = (1, 0) \\ 1 - \theta & \text{if } x = e_2 = (0, 1) \end{cases}$$

realization $x = (x_1, x_2)$ is:

set $\{e_1, e_2\} \subset \mathbb{R}^2$, i.e. $x = (x_1, x_2) \in \{(1, 0), (0, 1)\}$. The PMF of the RV $X = (X_1, X_2)$ with $X = (X_1, X_2)$ is a Bernoulli(θ) random vector (RV) if it has only two possible outcomes in the set $\{e_1, e_2\} \subset \mathbb{R}^2$, i.e. $x = (x_1, x_2) \in \{(1, 0), (0, 1)\}$. The PMF of the RV $X = (X_1, X_2)$ with

Model 15 (Bernoulli(θ) RV) Given a parameter $(\theta, 1 - \theta) \in \Delta_1$, the unit 1-Simplex, we say that

$$(1, 0) + (0, 1) + (1, 0) = (1 + 0 + 1, 0 + 1 + 0) = (2, 1)$$

4.

$$(1, 0) + (0, 1) + (1, 0) = (1 + 0 + 1, 0 + 1 + 0) = (2, 1).$$

its two sides

3. Generally, the diagonal of the parallelogram is the resultant or sum of the vectors representing

2. $(1, 0)$ and $(0, 1)$ are vectors for the two sides of unit square and $(1, 1)$ is its diagonal.

$$(0, 1) + (0, 1) = (0 + 0, 1 + 1) = (0, 2)$$

$$(1, 0) + (0, 1) = (1 + 0, 0 + 1) = (1, 1)$$

$$(1, 0) + (1, 0) = (1 + 1, 0 + 0) = (2, 0)$$

1. addition is component-wise

Solution:

positive definite matrix. Setting $u = 0$ and $\Sigma = I$ gives back the standard multivariate normal RV. where $|\Sigma|$ denotes the determinant of Σ , u is a vector of length m and Σ is a $m \times m$ symmetric

$$f_X(x; \mu, \Sigma) = f_{X_1, X_2, \dots, X_m}(x_1, x_2, \dots, x_m; \mu, \Sigma) = \frac{(2\pi)^{m/2} |\Sigma|^{1/2}}{1} \exp\left(-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu)\right)$$

it has joint probability density function

More generally, a vector X has a multivariate normal distribution denoted by $X \sim \text{Normal}(\mu, \Sigma)$,

1 along the diagonal entries and 0 on all off-diagonal entries.

We say that Z has a standard multivariate normal distribution and write $Z \sim \text{Normal}(0, I)$, where

$$f_Z(z) = f_{Z_1, Z_2, \dots, Z_m}(z_1, z_2, \dots, z_m) = \frac{(2\pi)^{m/2}}{1} \exp\left(-\frac{1}{2} \sum_{j=1}^m z_j^2\right) = \frac{(2\pi)^{m/2}}{1} \exp\left(-\frac{1}{2} z^T z\right)$$

where, Z_1, Z_2, \dots, Z_m are jointly independent $\text{Normal}(0, 1)$ RVs. Then the PDF of Z is

$$\begin{pmatrix} Z_m \\ \vdots \\ Z_2 \\ Z_1 \end{pmatrix} = Z$$

matrix Σ . To begin, let

Model 18 (Normal(μ, Σ) RV) The multivariate Normal(μ, σ^2) RV has two parameters, $\mu \in \mathbb{R}$ and $\sigma^2 \in (0, \infty)$. In the multivariate version $\mu \in \mathbb{R}^{m \times 1}$ is a column vector and σ^2 is replaced by a

Multinomial distributions are at the very foundations of various machine learning algorithms, including, fitting junk email, learning from large knowledge-based resources like www, Wikipedia, word-net, etc.

Multinomial distributions are at the very foundations of various machine learning algorithms, including, fitting junk email, learning from large knowledge-based resources like www, Wikipedia, word-net, etc.

Labwork 100 (Septicum Sampler Demo – Sum of n IID de Molivre(1/3, 1/3, 1/3) RVs) Let us understand the Septicum construction of the Multinomial($n, 1/3, 1/3, 1/3$) RVs as the sum of n independent and identical de Molivre($1/3, 1/3, 1/3$) RVs by calling the interactive visual cognitive tool as follows:

Multinomial($n, \theta_1, \theta_2, \theta_3$) RV.

Once again, we can visualize that the sum of n IID de Molivre($\theta_1, \theta_2, \theta_3$) RVs constitute the n . Once again, we can visualize that the sum of n IID de Molivre($\theta_1, \theta_2, \theta_3$) RVs constitute the n . Of paths that lead to a given triariate sum (y_1, y_2, y_3) with $\sum_{i=1}^3 y_i = n$ as the multinomial coefficient at triplets placed at the integral points in the 3-simplex, $\Sigma^3 = \{(y_1, y_2, y_3) \in \mathbb{Z}_+^3 : \sum_{i=1}^3 y_i = n\}$. In the Septicum, balls choose from one of three paths along e_1 , e_2 and e_3 with probabilities θ_1, θ_2 and θ_3 , respectively, in an IID manner at each of the n levels, before they collect at buckets placed at the integral points in the 3-simplex, $\Sigma^3 = \{(y_1, y_2, y_3) \in \mathbb{Z}_+^3 : \sum_{i=1}^3 y_i = n\}$. We can visualize the Multinomial($n, \theta_1, \theta_2, \theta_3$) process as a sum of n IID de Molivre($\theta_1, \theta_2, \theta_3$) RVs via a three dimensional extension of the Q-matrix called the "Septicum" and relate the number of Paths that lead to a given triariate sum (y_1, y_2, y_3) with $\sum_{i=1}^3 y_i = n$ as the multinomial coefficient $\frac{n!}{y_1! y_2! y_3!}$.

Note that the marginal PMF of X_j is Binomial(n, θ_j) for any $j = 1, 2, \dots, k$.

$$\binom{n}{y_1, y_2, \dots, y_k} := \frac{y_1! y_2! \dots y_k!}{n!}.$$

where, the multinomial coefficients:

drops through n levels of pegs where the probability of a right turn at each peg is independently and identically θ . In other words, the $\text{Binomial}(n, \theta)$ $\text{R}\vec{V} (Y, n - Y)$ is the sum of n IID $\text{Bernoulli}(\theta)$ $\text{R}\vec{V}$ s $X_1 := (X_{1,1}, X_{1,2}), X_2 := (X_{2,1}, X_{2,2}), \dots, X_n := (X_{n,1}, X_{n,2})$:

$$(Y, n - Y) = X_1 + X_2 + \dots + X_n = (X_{1,1}, X_{1,2}) + (X_{2,1}, X_{2,2}) + \dots + (X_{n,1}, X_{n,2})$$

Exercise 3.36 (Random walk in the first Quadrant) Consider an independent and identical random walk starting from $(0,0)$ in the first quadrant where you go east, i.e., add $(1,0)$ to your current position with probability θ , and go north, i.e., add $(0,1)$ to your current position with probability $1 - \theta$. Suppose you take n such IID steps according to the $\text{Bernoulli}(\theta)$ $\text{R}\vec{V}$. Answer the following questions:

1. How does the number of paths that lead to a (x_1, x_2) with $x_1 + x_2 = n$ relate to the binomial coefficient $\binom{n}{x_1}$?
2. What is the probability of taking x_1 steps east and x_2 steps north?

Exercise 3.37 (Random walks in the first Quadrant and Galton's Quincunx) Compare the probability models for the Random walk in the first quadrant and Galton's Quincunx and explain how they are related.

Labwork 99 (Quincunx Sampler Demo – Sum of n IID $\text{Bernoulli}(1/2)$ $\text{R}\vec{V}$ s) Let us understand the Quincunx construction of the $\text{Binomial}(n, 1/2)$ $\text{R}\vec{V} X$ as the sum of n independent and identical $\text{Bernoulli}(1/2)$ $\text{R}\vec{V}$ s by calling the interactive visual cognitive tool as follows:

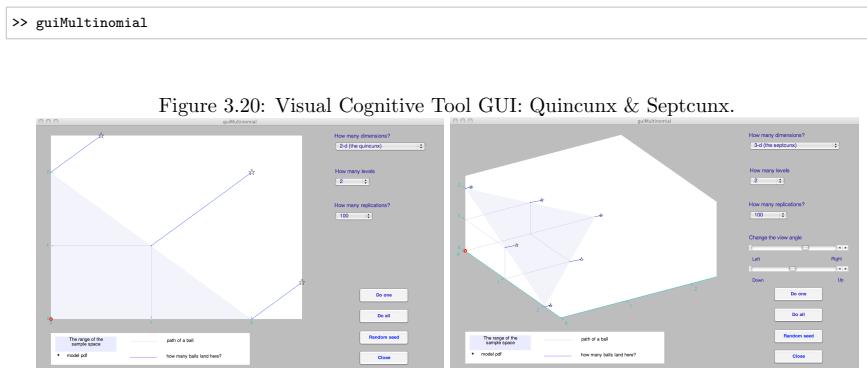
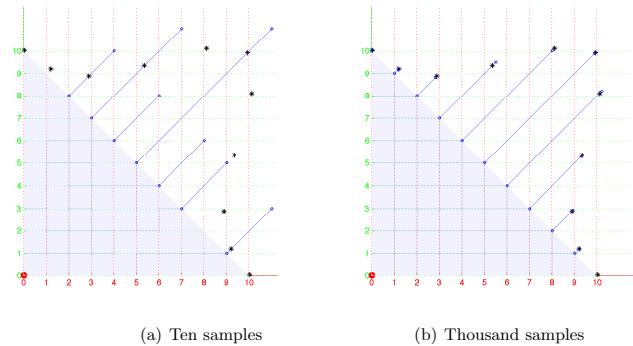


Figure 3.20: Visual Cognitive Tool GUI: Quincunx & Septcunx.

We are now ready to extend the $\text{Binomial}(n, \theta)$ RV or $\text{R}\vec{V}$ to its multivariate version called the $\text{Multinomial}(n, \theta_1, \theta_2, \dots, \theta_k)$ $\text{R}\vec{V}$. We develop this $\text{R}\vec{V}$ as the sum of n IID $\text{de Moivre}(\theta_1, \theta_2, \dots, \theta_k)$ $\text{R}\vec{V}$ s that is defined next by extending $\text{de Moivre}(\theta_1, \theta_2, \dots, \theta_k)$ RV taking values in $\{1, 2, \dots, k\}$ of Model 14 to its vector-valued cousin taking values in $\{e_1, e_2, \dots, e_k\}$, the ortho-normal basis vectors in \mathbb{R}^k .

Figure 3.21: Quincunx on the Cartesian plane. Simulations of $\text{Binomial}(n = 10, \theta = 0.5)$ RV as the x-coordinate of the ordered pair resulting from the culmination of sample trajectories formed by the accumulating sum of $n = 10$ IID $\text{Bernoulli}(\theta = 0.5)$ random vectors over $\{(1, 0), (0, 1)\}$ with probabilities $\{\theta, 1 - \theta\}$, respectively. The blue lines and black asterisks perpendicular to and above the diagonal line, i.e. the line connecting $(0, 10)$ and $(10, 0)$, are the density histogram of the samples and the PMF of our $\text{Binomial}(n = 10, \theta = 0.5)$ RV , respectively.



Model 16 (de Moivre($\theta_1, \theta_2, \dots, \theta_k$) $\text{R}\vec{V}$) The PMF of the $\text{de Moivre}(\theta_1, \theta_2, \dots, \theta_k)$ $\text{R}\vec{V} X := (X_1, X_2, \dots, X_k)$ taking value $x := (x_1, x_2, \dots, x_k) \in \{e_1, e_2, \dots, e_k\}$, where the e_i 's are ortho-normal basis vectors in \mathbb{R}^k is:

$$f(x; \theta_1, \theta_2, \dots, \theta_k) := P(X = x) = \sum_{i=1}^k \theta_i \mathbf{1}_{\{e_i\}}(x) = \begin{cases} \theta_1 & \text{if } x = e_1 := (1, 0, \dots, 0) \in \mathbb{R}^k \\ \theta_2 & \text{if } x = e_2 := (0, 1, \dots, 0) \in \mathbb{R}^k \\ \vdots & \\ \theta_k & \text{if } x = e_k := (0, 0, \dots, 1) \in \mathbb{R}^k \\ 0 & \text{otherwise} \end{cases}$$

Of course, $\sum_{i=1}^k \theta_i = 1$.

When we add n IID $\text{de Moivre}(\theta_1, \theta_2, \dots, \theta_k)$ $\text{R}\vec{V}$ s together, we get the $\text{Multinomial}(n, \theta_1, \theta_2, \dots, \theta_k)$ $\text{R}\vec{V}$ as defined below.

Model 17 (Multinomial($n, \theta_1, \theta_2, \dots, \theta_k$) $\text{R}\vec{V}$) We say that a $\text{R}\vec{V} Y := (Y_1, Y_2, \dots, Y_k)$ obtained from the sum of n IID $\text{de Moivre}(\theta_1, \theta_2, \dots, \theta_k)$ $\text{R}\vec{V}$ s with realizations

$$y := (y_1, y_2, \dots, y_k) \in \mathbb{Y} := \{(y_1, y_2, \dots, y_k) \in \mathbb{Z}_+^k : \sum_{i=1}^k y_i = n\}$$

has the PMF given by:

$$f(y; n, \theta) := f(y; n, \theta_1, \theta_2, \dots, \theta_k) := P(Y = y; n, \theta_1, \theta_2, \dots, \theta_k) = \binom{n}{y_1, y_2, \dots, y_k} \prod_{i=1}^k \theta_i^{y_i},$$