

Introduction to Data Science 1MS041

Benny Avelin, Raazesh Sainudiin*, Dominic Lee[†] and Michael Nussbaum[•],

*Laboratory for Mathematical Statistical Experiments, Uppsala Centre, and

[†]Department of Mathematics, Uppsala University, Uppsala, Sweden

[•]School of Mathematics and Statistics, University of Canterbury, Christchurch, New Zealand

[•]Department of Mathematics, Cornell University, Ithaca, New York, USA

Version Date: July 5, 2021

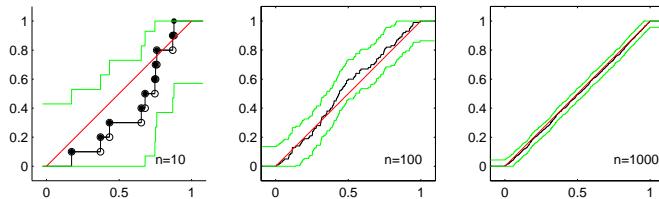
©2021–2021 Benny Avelin ©2007–2021 Raazesh Sainudiin. ©2008–2021 Dominic Lee. ©2010–2021

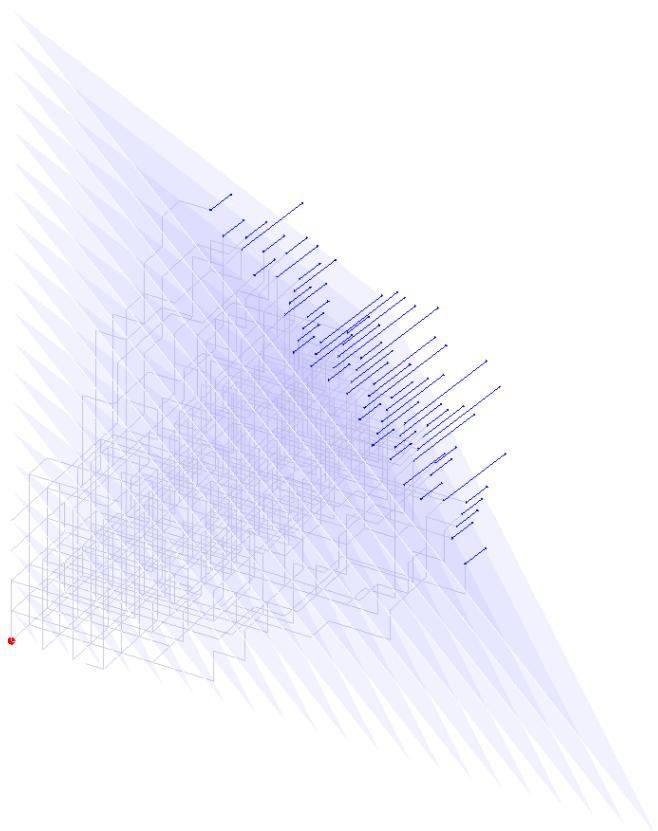
This work is licensed under the Creative Commons Attribution-Noncommercial-Share Alike 4.0

International License. To view a copy of this license, visit

<https://creativecommons.org/licenses/by-nc-sa/4.0/>.

This work was partially supported by NSF grant DMS-03-06497, NSF/NIGMS grant DMS-02-01037 and Erskine Fellowship at Department of Statistics, University of Oxford..





Contents

1 Preliminaries	11
1.1 Elementary Set Theory	11
1.2 Exercises	13
1.3 Natural Numbers, Integers and Rational Numbers	14
1.4 Real Numbers	18
1.5 Introduction to MATLAB	22
1.6 Elementary Combinatorics	26
1.7 Array, Sequence, Limit,	29
1.8 Elementary Real Analysis	34
1.8.1 Limits of Real Numbers – A Review	34
1.9 Elementary Number Theory	37
2 Probability Model	38
2.1 Experiments	38
2.2 Probability	40
2.2.1 Consequences of our Definition of Probability	42
2.2.2 Sigma Algebras of Typical Experiments*	44
2.3 Exercises in Probability	46
2.4 Conditional Probability	47
2.4.1 Bayes' Theorem	49
2.4.2 Independence and Dependence	53
2.5 Exercises in Conditional Probability	55
3 Random Variables	58
3.1 Basic Definitions	60
3.2 Discrete Random Variables	63
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	86
3.5 Exercises in Continuous Random Variables	91
3.6 Transformations of random variables	92
3.6.1 A Review of Inverse Images	92
3.6.2 Transformations of discrete random variables	95
3.6.3 Transformations of continuous random variables	96

3.7	Exercises in Transformations of Random Variables	103
3.8	Expectations	103
3.8.1	Expectations of functions of random variables	104
3.8.2	Properties of expectations	107
3.8.3	Expectation of Common Random Variables	107
3.9	Exercises in Expectations of Random Variables	113
3.10	Multivariate Random Variables	114
3.10.1	\mathbb{R}^2 -valued Random Variables	115
3.10.2	Conditional Random Variables	126
3.10.3	\mathbb{R}^m -valued Random Variables	127
3.10.4	Some Common \mathbb{R}^m -valued RVs	131
3.10.5	Dependent Random Variables	137
3.11	Exercises in Multivariate Random Variables	137
3.12	Characteristic Functions	141
3.12.1	Obtaining Moments from Characteristic Function	141
3.12.2	Moment Generating Function	146
3.13	Exercises in Characteristic Functions	146
4	Statistics	148
4.1	Data and Statistics	148
4.1.1	Univariate Data	155
4.1.2	Bivariate Data	157
4.1.3	Trivariate Data	158
4.1.4	Multivariate Data	159
4.1.5	Loading and Exploring Real-world Data	160
4.1.6	Geological Data	160
4.1.7	Metereological Data	164
4.1.8	Textual Data	166
4.1.9	Machine Sensor Data	168
4.2	Exercises in Statistics	169
4.3	Fundamentals of Estimation	169
4.3.1	Introduction	169
4.3.2	Point Estimation	169
4.3.3	Some Properties of Point Estimators	171
4.3.4	Confidence Set Estimation	174
4.4	Parameter Estimation and Likelihood	177
4.4.1	Point and Set Estimation – A General Likelihood Approach	178
4.4.2	Likelihood	178
4.4.3	Moment Estimator (MME)	189
4.5	Practical Excursion in One-dimensional Optimisation	190
4.6	More Properties of the Maximum Likelihood Estimator	194
4.7	Fisher Information	194
4.8	Delta Method	201
5	Maximum Likelihood Estimation for Multiparameter Models	205
5.1	Introduction	205
5.2	Practical Excursion in Multi-dimensional Optimisation	205
5.3	Confidence Sets for Multiparameter Models	209

6 Simulation	213
6.1 Physical Random Number Generators	213
6.2 Pseudo-Random Number Generators	213
6.2.1 Linear Congruential Generators	214
6.2.2 Generalized Feedback Shift Register and the “Mersenne Twister” PRNG .	217
6.3 Simulation of non-Uniform(0, 1) Random Variables	219
6.3.1 Inversion Sampler for Continuous Random Variables	219
6.3.2 Inversion Sampler for Discrete Random Variables	227
6.3.3 von Neumann Rejection Sampler (RS)	236
6.4 Exercises in Simulation	241
7 Limit Laws of Statistics	242
7.1 Convergence of Random Variables	242
7.1.1 Properties of Convergence of RVs**	248
7.2 Law of Large Numbers	248
7.2.1 Application: Point Estimation of $E(X_1)$	251
7.3 Central Limit Theorem	253
7.3.1 Application: Tolerating Errors in our estimate of $E(X_1)$	254
7.3.2 Application: Set Estimation of $E(X_1)$	256
7.4 Exercises in Limit Laws of Statistics	258
8 Finite Markov Chains	259
8.1 Introduction	259
8.2 Random Mapping Representation and Simulation	267
8.3 Irreducibility and Aperiodicity	273
8.4 Stationarity	276
8.5 Reversibility	278
8.6 Metropolis-Hastings Markov chain	282
8.7 Glauber Dynamics	286
8.7.1 Random Walks on \mathbb{Z} and the reflection principle	291
8.8 Coupling from the past	291
8.8.1 <i>Algorithm – Coupling from the past.</i>	292
8.9 Non-parametric DF Estimation	295
8.9.1 Estimating DF	296
8.10 Plug-in Estimators of Statistical Functionals	301
8.11 Bootstrap	303
8.11.1 Non-parametric Bootstrap for Confidence Sets	303
8.11.2 Parametric Bootstrap for Confidence Sets	306

List of Tables

1.1	Symbol Table: Sets and Numbers	21
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.	126
5.1	Summary of the Method of Moment Estimator (MME) and the Maximum Likelihood Estimator (MLE) for some IID Experiments.	209

List of Figures

1.1	Union and intersection of sets shown by Venn diagrams	12
1.2	These Venn diagram illustrate De Morgan's Laws.	13
1.3	A function f ("father of") from \mathbb{X} (a set of children) to \mathbb{Y} (their fathers) and its inverse ("children of").	16
1.4	A pictorial depiction of addition and its inverse. The domain is plotted in orthogonal Cartesian coordinates	17
1.5	A depiction of the real line segment $[-10, 10]$	19
1.6	Point plot and stem plot of the finite sequence $\langle b_{1:10} \rangle$ declared as an array.	31
1.7	A plot of the sine wave over $[-2\pi, 2\pi]$	34
2.1	A binary tree whose leaves are all possible outcomes.	40
2.2	First ball number in 1114 NZ Lotto draws from 1987 to 2008.	43
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.	51
3.1	The Indicator function of event $A \in \mathcal{F}$ is a RV $\mathbf{1}_A$ with DF F	62
3.2	A RV X from a sample space Ω with 8 elements to \mathbb{R} and its DF F	63
3.3	$f(x)$ and $F(x)$ of the <i>fair coin toss</i> random variable X , a discrete uniform RV on $\{0, 1\}$	65
3.4	$f(x)$ and $F(x)$ of the <i>fair die toss</i> random variable X , a discrete uniform RV on $\{1, 2, 3, 4, 5, 6\}$	67
3.5	$f(x)$ and $F(x)$ of surmised <i>astragali toss</i> random variable X , a discrete (non-uniform) RV on $\{1, 2, 3, 4\}$	67
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	68
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 170 and Labwork 171 you will see in the sequel.	71
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 174.	73
3.9	Figures from Sir Francis Galton, F.R.S., <i>Natural Inheritance</i> , , Macmillan, 1889.	75
3.10	$f(x)$ and $F(x)$ of the $\text{Uniform}(0, 1)$ random variable X	85
3.11	A convenient but mathematically imprecise Matlab plot from a polyline interpolation for the PDF and DF or CDF of the $\text{Uniform}(0, 1)$ 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).	87
3.13	Density and distribution functions of $\text{Exponential}(\lambda)$ RVs, for $\lambda = 1, 10, 10^{-1}$, in four different axes scales.	87
3.14	$f(x)$ and $F(x)$ of the $\text{Uniform}(\theta_1, \theta_2)$ random variable X	89
3.15	PDF and DF of a $\text{Normal}(\mu, \sigma^2)$ RV for different values of μ and σ^2	100

3.16	Mean ($E_\theta(X)$), variance ($V_\theta(X)$) and the rate of change of variance ($\frac{d}{d\theta} V_\theta(X)$) of a Bernoulli(θ) RV X as a function of the parameter θ	108
3.17	Mean and variance of a Geometric(θ) RV X as a function of the parameter θ	110
3.18	PDF of $X \sim \text{Poisson}(\lambda = 10)$ and the relative frequency histogram based on 1000 samples from X according to Simulation 175.	111
3.19	Diagrams done on the board!	125
3.20	Visual Cognitive Tool GUI: Quincunx & Septcunx.	134
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.	134
3.22	JPDF, Marginal PDFs and Frequency Histogram of Bivariate Standard Normal \vec{R}	137
3.23	JPDF, Marginal PDFs and Frequency Histogram of a Bivariate Normal \vec{R} for lengths of girths of cylindrical shafts in a manufacturing process (in cm).	138
4.1	Sample Space, Random Variable, Realisation, Data, and Data Space.	148
4.2	Data Spaces $\mathbb{X} = \{0, 1\}^2$ and $\mathbb{X} = \{0, 1\}^3$ for two and three Bernoulli trials, respectively.	149
4.3	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.	154
4.4	Frequency, Relative Frequency and Density Histograms	156
4.5	Frequency, Relative Frequency and Density Histograms	157
4.6	2D Scatter Plot	158
4.7	3D Scatter Plot	159
4.8	Plot Matrix of uniformly generated data in $[0, 1]^5$	160
4.9	Matrix of Scatter Plots of the latitude, longitude, magnitude and depth of the 22-02-2011 earth quakes in Christchurch, New Zealand.	162
4.10	Google Earth Visualisation of the earth quakes	164
4.11	Daily rainfalls in Christchurch since March 27 2010	165
4.12	Daily temperatures in Christchurch for one year since March 27 2010	167
4.13	Wordle of JOE 2010	168
4.14	Double Pendulum	169
4.15	Density and Confidence Interval of the Asymptotically Normal Point Estimator .	176
4.16	Data Spaces $\mathbb{X}_1 = \{0, 1\}$, $\mathbb{X}_2 = \{0, 1\}^2$ and $\mathbb{X}_3 = \{0, 1\}^3$ for one, two and three IID Bernoulli trials, respectively and the corresponding likelihood functions.	179
4.17	Plot of $\log(L(\lambda))$ as a function of the parameter λ and the MLE $\hat{\lambda}_{132}$ of 0.1102 for Fenemore-Wang Orbiter Waiting Times Experiment from STAT 218 S2 2007. The density or PDF and the DF at the MLE of 0.1102 are compared with a histogram and the empirical DF.	184
4.18	Comparing the Exponential($\hat{\lambda}_{6128} = 28.6694$) PDF and DF with a histogram and empirical DF of the times (in units of days) between earth quakes in NZ. The epicenters of 6128 earth quakes are shown in left panel.	184

4.19	Plots of the log likelihood $\ell_n(\theta) = \log(L(1, 0, 0, 0, 1, 1, 0, 0, 1, 0; \theta))$ as a function of the parameter θ over the parameter space $\Theta = [0, 1]$ and the MLE $\hat{\theta}_{10}$ of 0.4 for the coin-tossing experiment shown in standard scale (left panel) and log scale for x -axis (right panel).	188
4.20	100 realizations of 95% confidence intervals based on samples of size $n = 10, 100$ and 1000 simulated from IID Bernoulli($\theta^* = 0.5$) RVs. The MLE $\hat{\theta}_n$ (cyan dot) and the log-likelihood function (magenta curve) for each of the 100 replications of the experiment for each sample size n are depicted. The approximate normal-based 95% confidence intervals with blue boundaries are based on the exact $\text{se}_n = \sqrt{\theta^*(1 - \theta^*)/n} = \sqrt{1/4}$, while those with red boundaries are based on the estimated $\widehat{\text{se}}_n = \sqrt{\hat{\theta}_n(1 - \hat{\theta}_n)/n} = \sqrt{\frac{\bar{x}_n(1 - \bar{x}_n)}{n}}$. The fraction of times the true parameter $\theta^* = 0.5$ was contained by the exact and approximate confidence interval (known as <i>empirical coverage</i>) over the 100 replications of the simulation experiment for each of the three sample sizes are given by the numbers after <code>Cvrg.=</code> and <code>~</code> , above each sub-plot, respectively.	188
4.21	The ML fitted Rayleigh($\hat{\alpha}_{10} = 2$) PDF and a histogram of the ocean wave heights.	192
4.22	Plot of $\log(L(\lambda))$ as a function of the parameter λ , the MLE $\hat{\lambda}_{132} = 0.1102$ and 95% confidence interval $C_n = [0.0914, 0.1290]$ for Fenemore-Wang Orbiter Waiting Times Experiment from STAT 218 S2 2007. The PDF and the DF at (1) the MLE 0.1102 (black), (2) lower 95% confidence bound 0.0914 (red) and (3) upper 95% confidence bound 0.1290 (blue) are compared with a histogram and the empirical DF.	200
5.1	Plot of Levy density as a function of the parameter $(x, y) \in [-10, 10]^2$ scripted in Labwork ???.	206
5.2	Plot of the “well-behaved” (uni-modal and non-spiky) $\log(L((x_1, x_2, \dots, x_{100}); \lambda, \zeta))$, based on 100 samples $(x_1, x_2, \dots, x_{100})$ drawn from the Lognormal($\lambda^* = 10.36, \zeta^* = 0.26$) as per Labwork ???.	208
6.1	The linear congruential sequence of <code>LinConGen(256, 137, 0, 123, 257)</code> 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.	215
6.2	The LCG called <code>RANDU</code> 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.	217
6.3	Triplet point clouds from the “Mersenne Twister” with two different seeds (see Labwork 156).	218
6.4	A plot of the PDF, DF or CDF and inverse DF of the Uniform($-1, 1$) RV X	220
6.5	Visual Cognitive Tool GUI: Inversion Sampling from $X \sim \text{Uniform}(-5, 5)$	221
6.6	The PDF f , DF F , and inverse DF $F^{[-1]}$ of the the Exponential($\lambda = 1.0$) RV.	222
6.7	Visual Cognitive Tool GUI: Inversion Sampling from $X \sim \text{Exponential}(0.5)$	223
6.8	Visual Cognitive Tool GUI: Inversion Sampling from $X \sim \text{Laplace}(5)$	225
6.9	Visual Cognitive Tool GUI: Inversion Sampling from $X \sim \text{Cauchy}$	226
6.10	The DF $F(x; 0.3, 0.7)$ of the de Moivre($0.3, 0.7$) RV and its inverse $F^{[-1]}(u; 0.3, 0.7)$	229
6.11	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	237
6.12	Rejection Sampling from $X \sim \text{Normal}(0, 1)$ with PDF f based on 100 proposals from $Y \sim \text{Laplace}(1)$ with PDF g	238

7.1	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.	242
7.2	Distribution functions of several $\text{Normal}(\mu, \sigma^2)$ RVs for $\sigma^2 = 1, \frac{1}{10}, \frac{1}{100}, \frac{1}{1000}$	243
7.3	PDF $f_{X_n}(x) := \mathbb{1}_{(0,1)}(x)(1 - \cos(2\pi nx))$ of the RV X_n [the left sub-figure] and its DF $F_n(x) := \int_{-\infty}^x \mathbb{1}_{(0,1)}(v)(1 - \cos(2\pi nv))dv$ [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 $\mathbb{1}_{(0,1)}(x)x$, the DF of the $\text{Uniform}(0, 1)$ RV, while the corresponding PDFs $f_n(x)$ keep oscillating wildly with n across $[0, 2]$ about $\mathbb{1}_{(0,1)}(x)$, the PDF of the $\text{Uniform}(0, 1)$ RV X . Thus giving a counter-example to the claim that convergence in DFs does not imply convergence in PDFs.	245
7.4	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), $\text{Uniform}(0, 30)$ RV (green) and $\text{Exponential}(0.1)$ 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.	250
7.5	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 $\text{Uniform}(0, 10)$ RV, for each of five replicate simulations (magenta lines).	251
8.1	Transition Diagram of Flippant Freddy's Jumps.	260
8.2	The probability of being back in rollopia 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).	262
8.3	Transition Diagram of Dry and Wet Days in Christchurch.	264
8.4	Transition diagram over six lounges (without edge probabilities).	273
8.5	Stochastic Optimization with Metropolis chain.	285
8.6	The sample at time step 10^6 from the Glauber dynamics for the hardcore model on 100×100 regular torus grid. A red site is occupied while a blue site is vacant.	289
8.7	Plots of ten distinct ECDFs \hat{F}_n based on 10 sets of n IID samples from $\text{Uniform}(0, 1)$ RV X , as n increases from 10 to 100 to 1000. The DF $F(x) = x$ over $[0, 1]$ is shown in red. The script of Labwork ?? was used to generate this plot.	296
8.8	The empirical DFs $\hat{F}_n^{(1)}$ from sample size $n = 10, 100, 1000$ (black), is the point estimate of the fixed and known DF $F(x) = x, x \in [0, 1]$ of $\text{Uniform}(0, 1)$ RV (red). The 95% confidence band for each \hat{F}_n are depicted by green lines.	298
8.9	The empirical DF \hat{F}_{6128} for the inter earth quake times and the 95% confidence bands for the non-parametric experiment.	298
8.10	The empirical DF \hat{F}_{132} for the Orbiter waiting times and the 95% confidence bands for the non-parametric experiment.	299
8.11	The empirical DFs $\hat{F}_{n_1}^{(1)}$ with $n_1 = 56485$, for the web log times starting October 1, and $\hat{F}_{n_2}^{(2)}$ with $n_2 = 53966$, for the web log times starting October 2. Their 95% confidence bands are indicated by the green.	300
8.12	Data from Bradley Efrons LSAT and GPA scores for fifteen individuals (left). The confidence interval of the sample correlation, the plug-in estimate of the population correlation, is obtained from the sample correlation of one thousand bootstrapped data sets (right).	306

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 $\text{E}\mathbb{E}G$ of twenty two Greek lower-case alphabets that we may encounter later:

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

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:

$$b \notin B.$$

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

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

$$C \subset D$$

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

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 of set equality is notationally summarised as follows:

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

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

$$C \neq D .$$

The **union** of two sets C and D , written as $C \cup D$, is the set of elements that belong to C or D . We can formally express our definition of set union as:

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

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 x belongs to 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:

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

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

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

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:

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

When a universal set, e.g. U is well-defined, the **complement** of a given set B denoted by B^c is the set of all elements of U that don’t belong to B , i.e.:

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

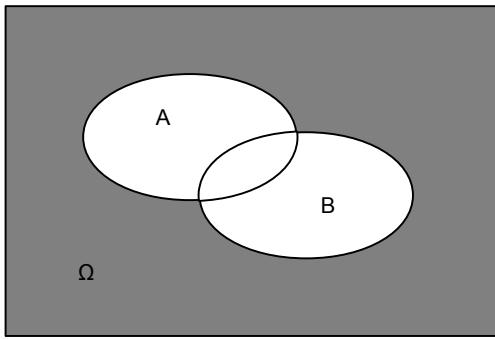
We say two sets C and D are **disjoint** if they have no elements in common, i.e. $C \cap D = \emptyset$.

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

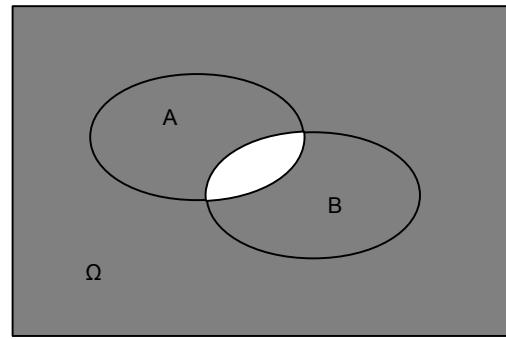
$$(A \cup B)^c = A^c \cap B^c \text{ and } (A \cap B)^c = A^c \cup B^c$$

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,

$$1. F \cap C =$$



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



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

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

$$2. F \cup C =$$

$$3. F \setminus C =$$

$$4. C \setminus F =$$

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, $T = \{T_1, T_2, T_3, L_3\}$.

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:

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

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

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

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

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 in both A and B .
$\{\} \text{ or } \emptyset$	— empty set.
Ω	— universal set.
A^c	— the complement of A , meaning the set of all elements in Ω , the universal set, which are not in A .

1.3 Natural Numbers, Integers and Rational Numbers

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

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

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**:

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

$$\begin{aligned} 1 &:= \#\{\star\} = \#\{\bullet\} = \#\{\alpha\} = \#\{\{\bullet\}\} = \#\{\{\bullet, \bullet'\}\} = \dots, \\ 2 &:= \#\{\star', \star\} = \#\{\bullet, \circ\} = \#\{\alpha, \omega\} = \#\{\{\circ\}, \{\alpha, \star, \bullet\}\} = \dots, \\ &\vdots \end{aligned}$$

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

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

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

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

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 $B = \{\star\}$, then $A \times B = \{(\circ, \star), (\bullet, \star)\}$. Elements of $A \times B$ are called **ordered pairs**.

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** (·) 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

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

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”).

as the **inverse** of f .

We motivated the non-negative integers \mathbb{Z}_+ via the size of a set. With the notion of two directions (+ and -) and the magnitude of the current position from the origin zero (0) of a dynamic entity, we can motivate the set of **integers**:

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

The integers with a **minus or negative sign** (-) before them are called negative integers and those with a **plus or positive sign** (+) before them are called positive integers. Conventionally, + signs are dropped. Some examples of functions you may have encountered are **arithmetic operations** such as **addition** (+), **subtraction** (-), **multiplication** (\cdot) and **division** ($/$) of ordered pairs of integers. The reader is assumed to be familiar with such arithmetic operations with pairs of integers. Every integer is either positive, negative, or zero. In terms of this we define the notion of **order**. 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 **less than or equal to** an integer b and write $a \leq b$ if $b - a$ is positive or zero. Finally, we say that a is greater than b and write $a > b$ if $b < a$. Similarly, a is greater than equal to b , i.e. $a \geq b$, if $b \leq a$. The set of integers are **well-ordered**, i.e., for every integer a there is a next largest integer $a + 1$.

Classwork 4 (Addition over integers) Consider the set of integers $\mathbb{Z} = \{\dots, -3, -2, -1, 0, 1, 2, 3, \dots\}$. Try to set up the arithmetic operation of addition as a function. The domain for addition is the Cartesian product of \mathbb{Z} :

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

What is its range ?

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

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

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

The expressions p/q and p'/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 , where q is positive and as small as possible. For example, $1/2$, $2/4$, $3/6$, and $1001/2002$ are different expressions for the same rational number whose irreducible unique expression is $1/2$.

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.

1.4 Real Numbers

Unlike rational numbers which are expressible in their reduced forms by p/q , it is fairly tricky to define or express real numbers. It is possible to define real numbers formally and constructively via equivalence classes of Cauchy sequence of rational numbers. For this all we need are notions of (1) infinity, (2) sequence of rational numbers and (3) distance between any two rational numbers in an infinite sequence of them. These are topics usually covered in an introductory course in real analysis and are necessary for a firm foundation in computational statistics. Instead of a formal constructive definition of real numbers, we give a more concrete one via decimal expansions. See Donald E. Knuth's treatment [*Art of Computer Programming, Vol. I, Fundamental Algorithms*, 3rd Ed., 1997, pp. 21-25] for a fuller story. A **real number** is a numerical quantity x that has a decimal expansion:

$$x = n + 0.d_1d_2d_3\dots, \text{ where, each } d_i \in \{0, 1, \dots, 9\}, n \in \mathbb{Z},$$

and the sequence $0.d_1d_2d_3\dots$ does not terminate with infinitely many consecutive 9s. By the above decimal representation, the following arbitrarily accurate enclosure of the real number x by rational numbers is implied:

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

for every $k \in \mathbb{N}$. Thus, rational arithmetic $(+, -, \cdot, /)$ can be extended with arbitrary precision to any ordered pair of real numbers x and y by operations on their rational enclosures \underline{x}, \bar{x} and \underline{y}, \bar{y} .

Some examples of real numbers that are not rational (**irrational numbers**) are:

$\sqrt{2} = 1.41421356237309\dots$ the side length of a square with area of 2 units

$\pi = 3.14159265358979\dots$ the ratio of the circumference to diameter of a circle

$e = 2.71828182845904\dots$ Euler's constant

We can think of π as being enclosed by the following pairs of rational numbers:

$$\begin{aligned} 3 + \frac{1}{10} &=: \underline{\pi}_1 \leq \pi < \bar{\pi}_1 := 3 + \frac{1}{10} + \frac{1}{10^1} \\ 3 + \frac{1}{10} + \frac{4}{100} &=: \underline{\pi}_2 \leq \pi < \bar{\pi}_2 := 3 + \frac{1}{10} + \frac{4}{100} + \frac{1}{100} \\ 3 + \frac{1}{10} + \frac{4}{100} + \frac{1}{10^3} &=: \underline{\pi}_3 \leq \pi < \bar{\pi}_3 := 3 + \frac{1}{10} + \frac{4}{100} + \frac{1}{10^3} + \frac{1}{10^3} \\ &\vdots \\ 3.14159265358979 &=: \underline{\pi}_{14} \leq \pi < \bar{\pi}_{14} := 3.14159265358979 + \frac{1}{10^{14}} \\ &\vdots \end{aligned}$$

Think of the real number system as the continuum of points that make up a line, as shown in Figure 1.5.

Let y and z be two real numbers such that $y \leq z$. Then, the **closed interval** $[y, z]$ is the set of real numbers x such that $y \leq x \leq z$:

$$[y, z] := \{x : y \leq x \leq z\}.$$

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 $b = 0$, $b^x := 1$.

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:

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

and the laws of exponents imply:

$$\begin{aligned}\log_b(xy) &= \log_b x + \log_b y, & \text{if } x > 0, y > 0 \text{ and} \\ \log_b(c^y) &= y \log_b c, & \text{if } c > 0.\end{aligned}$$

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

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

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

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

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

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 lower bound is the **greatest lower bound** if it is at least as large as any other lower bound. The greatest lower bound of a set of real numbers A is called the **infimum** of A and is denoted by:

$$\inf A := \text{greatest lower bound of } A$$

For example, $\inf(0, 1) = 0$ and $\inf\{10.333 \cup [-99, 1001.33]\} = -99$. We similarly define the **least upper bound** of a non-empty set of real numbers A to be the **supremum** of A and denote it as:

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

For example, $\sup(0, 1) = 1$ and $\sup\{10.333 \cup [-99, 1001.33]\} = 1001.33$. By convention, we define $\inf \emptyset := \infty$, $\sup \emptyset := -\infty$. Finally, if a set A is not bounded below then $\inf A := -\infty$ and if a set A is not bounded above then $\sup A := \infty$.

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

1.5 Introduction to MATLAB

We use MATLAB to perform computations and visualisations. MATLAB is a numerical computing environment and programming language that is optimised for vector and matrix processing. STAT 218/313 students will have access to Maths & Stats Department's computers that are licensed to run MATLAB . You can remotely connect to these machines from home by following instructions at <http://www.math.canterbury.ac.nz/php/resources/compdocs/remote>.

Labwork 5 (Basics of MATLAB) Let us familiarize ourselves with MATLAB in this session. First, you need to launch MATLAB from your terminal. Since this is system dependent, ask your tutor for help. The command window within the MATLAB window is where you need to type commands. Here is a minimal set of commands you need to familiarize yourself with in this session.

1. Type the following command to add 2 numbers in the command window right after the command prompt `>>` .

```
>> 13+24
```

Upon hitting **Enter** or **Return** on your keyboard, you should see:

```
ans =
```

```
37
```

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.

```
>> 13+24 % adding 13 to 24 using the binary arithmetic operator +
ans =      37
```

3. You can **create or reopen a diary file** in MATLAB to record your work. Everything you typed or input and the corresponding output in the command window will be recorded in the diary file. You can create or reopen a diary file by typing `diary filename.txt` in the command window. When you have finished recording, simply type `diary off` in the command window **to turn off the diary file**. The diary file with .txt extension is simply a text-file. It can be edited in different editors after the diary is turned off in MATLAB . You need to type `diary LabWeek1.txt` to start recording your work for electronic submission if needed.

```

>> diary blah.txt % start a diary file named blah.txt
>> 3+56
ans =      59
>> diary off % turn off the current diary file blah.txt
>> 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 :

$$p = 45.89 * 1.00009$$
$$m = 5376.0 - 6.00$$

$$d = 89.0 / 23.3454$$
$$p = 2^{0.5}$$

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.

```
>> 2^(1/2)
ans =      1.4142
```

The omission of parenthesis about $1/2$ means something else and you get the following output:

```
>> 2^1/2
ans =      1
```

MATLAB first takes the 1st power of 2 and then divides it by 2 using its default precedence rules for binary operators in the absence of parenthesis. The order of operations or default precedence rule for arithmetic operations is 1. brackets or parentheses; 2. exponents (powers and roots); 3. division and multiplication; 4. addition and subtraction. The mnemonic **bedmas** can be handy. When in doubt, use parenthesis to force the intended order of operations.

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

```
>> 10/0
ans =      Inf
```

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:

```
>> X
X =      37
>> clear X
>> X
??? Undefined function or variable 'X'.
```

Entering **X** after **clearing** it gives the above self-explanatory error message preceded by **???**.

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);
```

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

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

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.

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

$$n_1 \times n_2 \times n_3 \times \dots$$

Example 6 Suppose that a Personal Identification Number (PIN) is a six-symbol code word in which the first four entries are letters (lowercase) and the last two entries are digits. How many PINS are there? There are six selections to be made:

First letter: 26 possibilities

Fourth letter: 26 possibilities

Second letter: 26 possibilities

First digit: 10 possibilities

Third letter: 26 possibilities

Second digit: 10 possibilities

So in total, the total number of possible PINS is:

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

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 the total number of possible PINS is:

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

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

“I have 17 probability texts on my bottom shelf”

then I don't care (usually) about what order they are in, but in the statement

“The combination of my PIN is math99”

I do care about order. A different order gives a different PIN.

So in mathematics, we use more precise language:

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

Permutations: There are basically two types of permutations:

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 .
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!

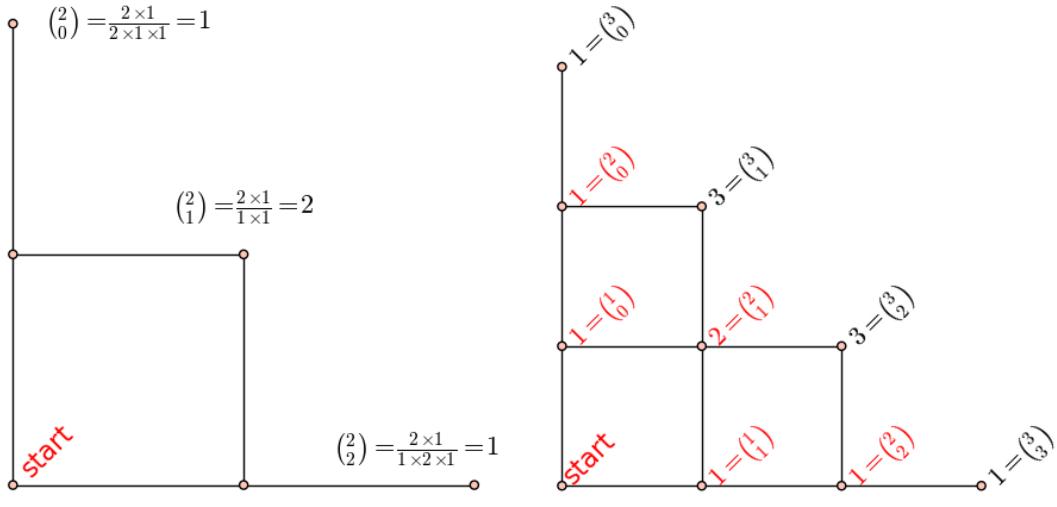
The number of distinct **combinations** of n objects taking r at a time is given by

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

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 $r \in \{0, 1, 2, 3\}$.



(a) Walking two blocks north-easterly.

(b) Walking three blocks north-easterly.

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.

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 $\{1, 2\}$:

$$12, \quad 21,$$

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

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

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 ($p_{3,2}$):

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

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

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

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

$$n! := p_{n,n} = n(n - 1)(n - 2) \dots (n - n + 1) = n(n - 1)(n - 2) \dots (3)(2)(1) =: \prod_{i=1}^n i.$$

Some factorials to bear in mind

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

When n is large we can get a good idea of $n!$ without laboriously carrying out the $n - 1$ multiplications via Stirling's approximation (*Methodus Differentialis* (1730), p. 137) :

$$n! \approx \sqrt{2\pi n} \left(\frac{n}{e}\right)^n .$$

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:

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

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, *Pingala’s Chandadhśāstra*, 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:

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

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

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

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

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

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

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

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.

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

$$\text{and a column vector } \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_{m-1} \\ y_m \end{bmatrix} = [y_1 \ y_2 \ \dots \ y_m]' = (y_1, y_2, \dots, y_m)'.$$

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:

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

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.

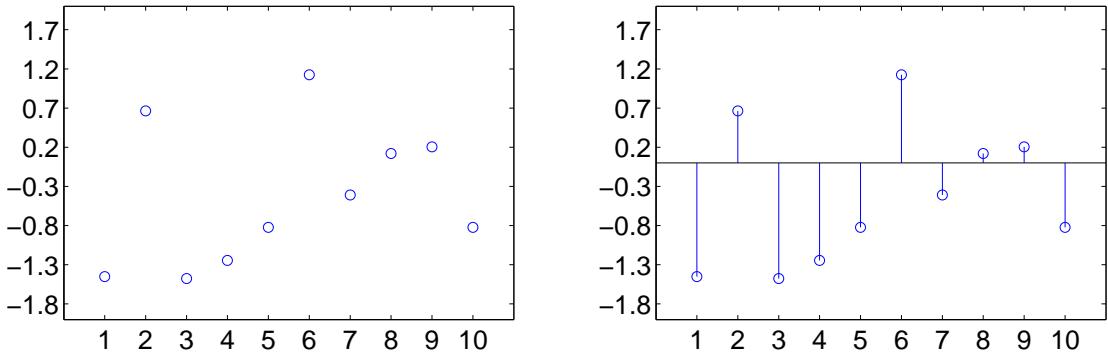
Finite sequences, vectors and matrices can be represented in a computer by an elementary data structure called an **array**.

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

```
>> 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
z =
    1      2      3
>> % linspace(x1, x2, n) generates n points linearly spaced between x1 and x2
>> r = linspace(1, 3, 3)          % Declare an 1 X 3 row vector r = z using linspace
r =
    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
>> rT = r'               % The column vector (1,2,3)' by taking the transpose of r via r'
rT =
    1
    2
    3
>> y = [1 1 1]           % y is a sequence or row vector of 3 1's
y =
    1      1      1
>> 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
```

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

```
>> Z=zeros(2,10) % the 2 X 10 matrix of zeros
Z =
    0      0      0      0      0      0      0      0      0      0
    0      0      0      0      0      0      0      0      0      0
>> O=ones(4,5) % the 4 X 5 matrix of ones
O =
    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 =
    1      0      0      0
    0      1      0      0
    0      0      1      0
    0      0      0      1
```

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

```
>> y % the array y is
y =
    1      1      1
>> z % the array z is
z =
    1      2      3
>> x = y + z           % x is the sum of vectors y and z (with same size 1 X 3)
x =
    2      3      4
>> y = y * 2            % y is updated to 2 * y (each term of y is multiplied by 2)
y =
    2      2      2
>> p = z .* y           % p is the vector obtained by term-by-term product of z and y
p =
    2      4      6
>> d = z ./ y           % d is the vector obtained by term-by-term division of z and y
```

```

d =      0.5000    1.0000    1.5000
>> t=linspace(-10,10,4)          % t has 4 numbers equally-spaced between -10 and 10
t =   -10.0000   -3.3333    3.3333    10.0000
>> s = sin(t)                  % s is a vector obtained from the term-wise sin of the vector t
s =    0.5440    0.1906   -0.1906   -0.5440
>> sSq = sin(t) .^ 2           % sSq is an array obtained from term-wise squaring ( .^ 2 ) of the sin(t) array
sSq =   0.2960    0.0363    0.0363    0.2960
>> cSq = cos(t) .^ 2           % cSq is an array obtained from term-wise squaring ( .^ 2 ) of the cos(t) array
cSq =   0.7040    0.9637    0.9637    0.7040
>> 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,'.');                      % plot x versus y as dots should appear in the Figure window
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.

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

1.8 Elementary Real Analysis

1.8.1 Limits of Real Numbers – A Review

Let us first recall some elementary ideas from real analysis.

Definition 3 (Convergent sequence of real numbers) A sequence of real numbers $\langle x_i \rangle_{i=1}^{\infty} := x_1, x_2, \dots$ is said to converge to a limit $a \in \mathbb{R}$ and denoted by:

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

if for every natural number $m \in \mathbb{N}$, a natural number $N_m \in \mathbb{N}$ exists such that for every $j \geq N_m$, $|x_j - a| \leq \frac{1}{m}$.

In words, $\lim_{i \rightarrow \infty} x_i = a$ means the following: no matter how small you make $\frac{1}{m}$ by picking as large an m as you wish, I can find an N_m , that may depend on m , such that every number in the sequence beyond the N_m -th element is within distance $\frac{1}{m}$ of the limit a .

Example 12 (Limit of a sequence of 17s) Let $\langle x_i \rangle_{i=1}^{\infty} = 17, 17, 17, \dots$. Then $\lim_{i \rightarrow \infty} x_i = 17$. This is because for every $m \in \mathbb{N}$, we can take $N_m = 1$ and satisfy the definition of the limit, i.e.:

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

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

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

However, several other sequences also approach the limit 0. Some such sequences that approach the limit 0 from the right are:

$$\langle x_{1:\infty} \rangle = \frac{1}{1}, \frac{1}{4}, \frac{1}{9}, \dots \quad \text{and} \quad \langle x_{1:\infty} \rangle = \frac{1}{1}, \frac{1}{8}, \frac{1}{27}, \dots ,$$

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

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

and finally some that approach 0 from either side are:

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

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 the Euler's constant e :

$$\begin{aligned}
\lim_{x \rightarrow 0} f(x) &= \lim_{x \rightarrow 0} (1+x)^{\frac{1}{x}} \\
&= \lim_{x \rightarrow 0} (x+1)^{(1/x)} \quad \text{Indeterminate form of type } 1^\infty. \\
&= \exp \left(\lim_{x \rightarrow 0} \log((x+1)^{(1/x)}) \right) \quad \text{Transformed using } \exp(\lim_{x \rightarrow 0} \log((x+1)^{(1/x)})) \\
&= \exp \left(\lim_{x \rightarrow 0} (\log(x+1))/x \right) \quad \text{Indeterminate form of type } 0/0. \\
&= \exp \left(\lim_{x \rightarrow 0} \frac{d \log(x+1)/dx}{dx/dx} \right) \quad \text{Applying L'Hospital's rule} \\
&= \exp \left(\lim_{x \rightarrow 0} 1/(x+1) \right) \quad \text{limit of a quotient is the quotient of the limits} \\
&= \exp \left(1/(\lim_{x \rightarrow 0} (x+1)) \right) \quad \text{The limit of } x+1 \text{ as } x \text{ approaches 0 is 1} \\
&= \exp(1) = e \approx 2.71828 .
\end{aligned}$$

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

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

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

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

despite the fact that $f(1) = \frac{1^3-1}{1-1} = \frac{0}{0}$ itself is undefined and does not exist.

Next we look at some examples of limits at infinity.

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

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

Example 18 (Limit of $(1 - \frac{\lambda}{n})^{-\alpha}$) The limit of $f(n) = (1 - \frac{\lambda}{n})^{-\alpha}$, for some $\alpha > 0$, as n approaches ∞ exists and it is 1 :

$$\lim_{n \rightarrow \infty} f(n) = \lim_{n \rightarrow \infty} \left(1 - \frac{\lambda}{n}\right)^{-\alpha} = 1 .$$

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:

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

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

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

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

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)

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

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:

$$x \bmod y := \begin{cases} x - y \lfloor x/y \rfloor & \text{if } y \neq 0 \\ x & \text{if } y = 0 \end{cases}$$

Chapter 2

Probability Model

2.1 Experiments

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.

Definition 6 An **experiment** is an activity or procedure that produces distinct, well-defined possibilities called **outcomes**. The set of all outcomes is called the **sample space**, and is denoted by Ω .

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**.

Events, $E_1, E_2 \dots E_n$, that cannot occur at the same time are called **mutually exclusive** events, or **pair-wise disjoint** events. This means that $E_i \cap E_j = \emptyset$ where $i \neq j$.

Example 23 Some standard examples of experiments are the following:

- $\Omega = \{\text{Defective, Non-defective}\}$ if our experiment is to inspect a light bulb.

There are only two outcomes here, so $\Omega = \{\omega_1, \omega_2\}$ where $\omega_1 = \text{Defective}$ and $\omega_2 = \text{Non-defective}$.

- $\Omega = \{\text{Heads, Tails}\}$ if our experiment is to note the outcome of a coin toss.

This time, $\Omega = \{\omega_1, \omega_2\}$ where $\omega_1 = \text{Heads}$ and $\omega_2 = \text{Tails}$.

- 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, $\Omega = \{1, 2, 3, 4, 5, 6\}$.

Some examples of events are the set of odd numbered outcomes $A = \{1, 3, 5\}$, and the set of even numbered outcomes $B = \{2, 4, 6\}$.

The simple events of Ω are $\{1\}, \{2\}, \{3\}, \{4\}, \{5\}$, and $\{6\}$.

The outcome of a random experiment is uncertain until it is performed and observed. Note that sample spaces need to reflect the problem in hand. The example below is to convince you that an experiment's sample space is merely a collection of distinct elements called outcomes and these outcomes have to be *discernible in some well-specified sense* to the experimenter!

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 .

Figure 2.1: A binary tree whose leaves are all possible outcomes.

Classwork 27 (A thrice-bifurcating tree of outcomes) Can you think of a graphical way to enumerate the outcomes of the Experiment 26? Draw a diagram of this under the caption of Figure 2.1, using the caption as a hint (in other words, draw your own Figure 2.1).

EXPERIMENT SUMMARY

Experiment	—	an activity producing distinct outcomes.
Ω	—	set of all outcomes of the experiment.
ω	—	an individual outcome in Ω , called a simple event.
$A \subseteq \Omega$	—	a subset A of Ω is an event.
Trial	—	one performance of an experiment resulting in 1 outcome.

2.2 Probability

The mathematical model for probability or the probability model is an axiomatic system that may be motivated by the intuitive idea of ‘long-term relative frequency’. If the axioms and definitions are intuitively motivated, the probability model simply follows from the application of logic to these axioms and definitions. No attempt to define probability in the real world is made. However, the application of probability models to real-world problems through statistical experiments has a fruitful track record. In fact, you are here for exactly this reason.

Idea 9 (The long-term relative frequency (LTRF) idea) Suppose we are interested in the fairness of a coin, i.e. if landing Heads has the same “probability” as landing Tails. We can toss it n times and call $N(H, n)$ the fraction of times we observed Heads out of n tosses. Suppose that after conducting the tossing experiment 1000 times, we rarely observed Heads, e.g. 9 out of the 1000 tosses, then $N(H, 1000) = 9/1000 = 0.009$. Suppose we continued the number of tosses to a million and found that this number approached closer to 0.1, or, more generally, $N(H, n) \rightarrow 0.1$ as $n \rightarrow \infty$. We might, at least intuitively, think that the coin is unfair and has a lower “probability” of 0.1 of landing Heads. We might think that it is fair had we observed $N(H, n) \rightarrow 0.5$ as $n \rightarrow \infty$. Other crucial assumptions that we have made here are:

1. **Something Happens:** Each time we toss a coin, we are certain to observe Heads **or** Tails, denoted by $H \cup T$. The probability that “something happens” is 1. More formally:

$$N(H \cup T, n) = \frac{n}{n} = 1.$$

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 $P(A)$, 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:

$$N(A, n) := \frac{\text{Number of times } A \text{ occurs}}{n = \text{Number of performances of } \mathcal{E}} \rightarrow P(A), \text{ as } n \rightarrow \infty \quad (2.1)$$

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: $\boxed{\Omega \in \mathcal{F}}$.
2. It is closed under complementation: $\boxed{A \in \mathcal{F} \implies A^c \in \mathcal{F}}$.
3. It is closed under countable unions: $\boxed{A_1, A_2, \dots \in \mathcal{F} \implies \bigcup_i A_i := A_1 \cup A_2 \cup \dots \in \mathcal{F}}$.

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 P which assigns each event $A \in \mathcal{F}$ a number $P(A)$ in the real interval $[0, 1]$, i.e. $\boxed{P : \mathcal{F} \rightarrow [0, 1]}$, such that:

1. The ‘Something Happens’ axiom holds, i.e. $\boxed{P(\Omega) = 1}$.
2. The ‘Addition Rule’ axiom holds, i.e. for events A and B :

$$\boxed{A \cap B = \emptyset \implies P(A \cup B) = P(A) + P(B)}.$$

2.2.1 Consequences of our Definition of Probability

It is important to realize that we accept the ‘addition rule’ as an axiom in our mathematical definition of probability (or our probability model) and we do **not** prove this rule. However, the facts which are stated (with proofs) below, are logical consequences of our definition of probability:

1. For any event A , $\boxed{P(A^c) = 1 - P(A)}.$

Proof: One line proof.

$$\overbrace{P(A) + P(A^c)}^{LHS} \underset{\substack{+ \text{ rule } A \cap A^c = \emptyset}}{=} P(A \cup A^c) \underset{A \cup A^c = \Omega}{=} P(\Omega) \underset{P(\Omega) = 1}{=} \overbrace{1}^{RHS} \Rightarrow_{LHS - P(A) \text{ & } RHS - P(A)} P(A^c) = 1 - P(A)$$

- If $A = \Omega$ then $A^c = \Omega^c = \emptyset$ and $\boxed{P(\emptyset) = 1 - P(\Omega) = 1 - 1 = 0}.$

2. For any two events A and B , we have the **inclusion-exclusion principle**:

$$\boxed{P(A \cup B) = P(A) + P(B) - P(A \cap B)}.$$

Proof: Since:

$$\begin{aligned} A &= (A \setminus B) \cup (A \cap B) && \text{and} && (A \setminus B) \cap (A \cap B) = \emptyset, \\ A \cup B &= (A \setminus B) \cup B && \text{and} && (A \setminus B) \cap B = \emptyset \end{aligned}$$

the addition rule implies that:

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

Substituting the first equality above into the second, we get:

$$P(A \cup B) = P(A \setminus B) + P(B) = P(A) - P(A \cap B) + P(B)$$

3. From inclusion-exclusion principle we get **Boole's inequality**: for any two events A, B

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

4. The inclusion-exclusion principle extends similarly to any three events A_1, A_2, A_3 as follows:

$$P(A_1 \cup A_2 \cup A_3) = P(A_1) + P(A_2) + P(A_3) - P(A_1 \cap A_2) - P(A_1 \cap A_3) - P(A_2 \cap A_3) + P(A_1 \cap A_2 \cap A_3)$$

and generalises to any n events A_1, A_2, \dots, A_n as follows:

$$P\left(\bigcup_{i=1}^n A_i\right) = \sum_{i=1}^n P(A_i) - \sum_{i < j} P(A_i \cap A_j) + \sum_{i < j < k} P(A_i \cap A_j \cap A_k) + \dots + (-1)^{n-1} \sum_{i < \dots < n} P\left(\bigcap_{i=1}^n A_i\right)$$

Proof: See the counting argument in https://en.wikipedia.org/wiki/Inclusion%20%23exclusion_principle if you are curious.

5. Once again by the inclusion-exclusion principle, the Boole's inequality generalises to any n events A_1, A_2, \dots, A_n as follows:

$$P\left(\bigcup_{i=1}^n A_i\right) \leq \sum_{i=1}^n P(A_i)$$

6. For a sequence of mutually disjoint events $A_1, A_2, A_3, \dots, A_n$:

$$A_i \cap A_j = \emptyset \quad \text{for any } i \neq j \quad \Rightarrow \quad P(A_1 \cup A_2 \cup \dots \cup A_n) = P(A_1) + P(A_2) + \dots + P(A_n).$$

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:

$$P(A_1 \cup A_2 \cup A_3) = P((A_1 \cup A_2) \cup A_3) \underset{+ \text{ rule}}{\underset{\curvearrowleft}{=}} P(A_1 \cup A_2) + P(A_3) \underset{+ \text{ rule}}{\underset{\curvearrowleft}{=}} P(A_1) + P(A_2) + P(A_3)$$

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

$$\Omega = \{1, 2, 3, \dots, 39, 40\}.$$

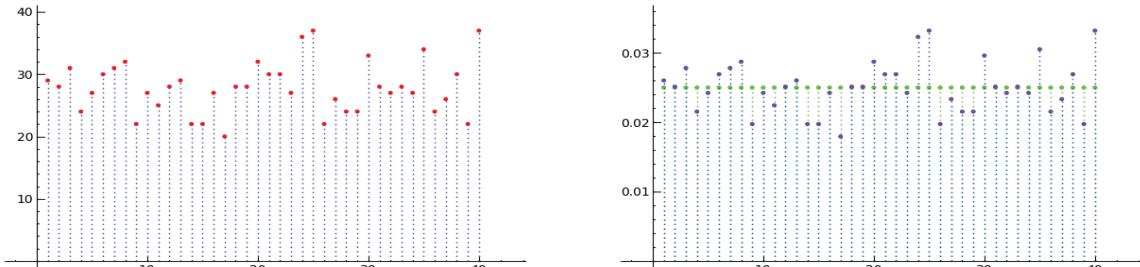
Because the balls are vigorously whirled around inside the Lotto machine, modelled as a well-stirred urn, before the first one pops out, we can model each ball to pop out first with the same probability. So, we assign each outcome $\omega \in \Omega$ the same probability of $\frac{1}{40}$, i.e., our probability model for this experiment is:

$$P(\omega) = \frac{1}{40}, \quad \text{for each } \omega \in \Omega = \{1, 2, 3, \dots, 39, 40\}.$$

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.



(a) Frequency of first ball.

(b) Relative frequency and probability of first ball.

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.

Probability Model

Sample space Ω	Set of all outcomes of an experiment
Sample point ω	Possible outcome of an experiment
(No counterpart)	Actual outcome ω^* of an experiment
Event A, a (suitable) subset of Ω	The real-world event corresponding to A occurs if and only if $\omega^* \in A$
$P(A)$, a number between 0 and 1	Probability that A will occur for an experiment yet to be performed

Real-world Interpretation

Sample space Ω	Set of all outcomes of an experiment
Sample point ω	Possible outcome of an experiment
(No counterpart)	Actual outcome ω^* of an experiment
Event A, a (suitable) subset of Ω	The real-world event corresponding to A occurs if and only if $\omega^* \in A$
$P(A)$, a number between 0 and 1	Probability that A will occur for an experiment yet to be performed

Events in Probability Model

Sample space Ω	The certain even ‘something happens’
The \emptyset of Ω	The impossible event ‘nothing happens’
The intersection $A \cap B$	‘Both A and B occur’
$A_1 \cap A_2 \cap \dots \cap A_n$	‘All of the events A_1, A_2, \dots, A_n occur simultaneously’
The union $A \cup B$	‘At least one of A and B occurs’
$A_1 \cup A_2 \cup \dots \cup A_n$	‘At least one of the events A_1, A_2, \dots, A_n occurs’
A^c , the complement of A	‘A does not occur’
$A \setminus B$	‘A occurs, but B does not occur’
$A \subset B$	‘If A occurs, then B must occur’

In the probability model of Example 28, show that for any event $E \subset \Omega$,

$$P(E) = \frac{1}{40} \times \text{number of elements in } E .$$

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:

$$P(E) = P(\{\omega_1, \omega_2, \dots, \omega_k\}) = P\left(\bigcup_{i=1}^k \{\omega_i\}\right) = \sum_{i=1}^k P(\{\omega_i\}) = \sum_{i=1}^k \frac{1}{40} = \frac{k}{40} .$$

2.2.2 Sigma Algebras of Typical Experiments*

Example 29 (‘Toss a fair coin once’) Consider the ‘Toss a fair coin once’ experiment. What is its sample space Ω and a reasonable collection of events \mathcal{F} that underpin this experiment?

$$\Omega = \{H, T\}, \quad \mathcal{F} = \{H, T, \Omega, \emptyset\} ,$$

A function that will satisfy the definition of probability for this collection of events \mathcal{F} and assign $P(H) = \frac{1}{2}$ is summarized below. First check that the above \mathcal{F} is a sigma-algebra. Draw a picture for P with arrows that map elements in the domain \mathcal{F} given above to elements in its range.

Event $A \in \mathcal{F}$	$P : \mathcal{F} \rightarrow [0, 1]$	$P(A) \in [0, 1]$
$\Omega = \{H, T\} \bullet$	\longrightarrow	1
T \bullet	\longrightarrow	$1 - \frac{1}{2}$
H \bullet	\longrightarrow	$\frac{1}{2}$
$\emptyset \bullet$	\longrightarrow	0

Classwork 30 (The trivial sigma algebra) Note that $\mathcal{F}' = \{\Omega, \emptyset\}$ is also a sigma algebra of the sample space $\Omega = \{H, T\}$. Can you think of a probability for the collection \mathcal{F}' ?

Event $A \in \mathcal{F}'$	$P : \mathcal{F}' \rightarrow [0, 1]$	$P(A) \in [0, 1]$
$\Omega = \{H, T\} \bullet$	→	
$\emptyset \bullet$	→	

Thus, \mathcal{F} and \mathcal{F}' are two distinct sigma algebras over our $\Omega = \{H, T\}$. 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:

1. When the sample space $\Omega = \{\omega_1, \omega_2, \dots, \omega_k\}$ is a finite set with k outcomes and $P(\omega_i)$, 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., $P(A) = \sum_{\omega_i \in A} P(\omega_i)$. Clearly, 2^Ω is indeed a sigma-algebra and it contains $2^{\#\Omega}$ events in it.
2. 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.
3. 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).

4. 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.

PROBABILITY SUMMARY

Axioms:

1. If $A \subseteq \Omega$ then $0 \leq P(A) \leq 1$ and $P(\Omega) = 1$.
2. If A, B are disjoint events, then $P(A \cup B) = P(A) + P(B)$.
[This is true only when A and B are disjoint.]
3. If A_1, A_2, \dots are disjoint then $P(A_1 \cup A_2 \cup \dots) = P(A_1) + P(A_2) + \dots$

Rules:

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

$$P(A \cup B) = P(A) + P(B) - P(A \cap B) \quad [\text{always true}]$$

2.3 Exercises in Probability

Ex. 2.2 — In English language text, the twenty six letters in the alphabet occur with the following frequencies:

E	13%	R	7.7%	A	7.3%	H	3.5%	F	2.8%	M	2.5%	W	1.6%	X	0.5%	J	0.2%
T	9.3%	O	7.4%	S	6.3%	L	3.5%	P	2.7%	Y	1.9%	V	1.3%	K	0.3%	Z	0.1%
N	7.8%	I	7.4%	D	4.4%	C	3%	U	2.7%	G	1.6%	B	0.9%	Q	0.3%		

Suppose you pick one letter at random from a randomly chosen English book from our central library with $\Omega = \{A, B, C, \dots, Z\}$ (ignoring upper/lower cases), then what is the probability of these events?

- (a) $P(\{Z\})$
- (b) $P(\text{'picking any letter'})$
- (c) $P(\{E, Z\})$
- (d) $P(\text{'picking a vowel'})$
- (e) $P(\text{'picking any letter in the word WAZZZUP'})$
- (f) $P(\text{'picking any letter in the word WAZZZUP or a vowel'})$.

Ex. 2.3 — Find the sample spaces for the following experiments:

1. Tossing 2 coins whose faces are sprayed with black paint denoted by B and white paint denoted by W .
2. Drawing 4 screws from a bucket of left-handed and right-handed screws denoted by L and R , respectively.
3. Rolling a die and recording the number on the upturned face until the first 6 appears.

Ex. 2.4 — Suppose we pick a letter at random from the word WAIMAKARIRI.

1. What is the sample space Ω ?

- 2.What probabilities should be assigned to the outcomes?
- 3.What is the probability of *not* choosing the letter R?

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 $P(\Omega)$ 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 $P(C \cup D)$.

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 $P(A) \neq 0$. The ‘conditional probability that B occurs given that A occurs’ denoted by $P(B|A)$ 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 $P(B|A)$ 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

$$P(B|A) \quad “\rightarrow” \quad \frac{N(A \cap B, n)}{N(A, n)} = \frac{N(A \cap B, n)/n}{N(A, n)/n} = \frac{P(A \cap B)}{P(A)}$$

as our $\mathcal{E}^n \rightarrow \mathcal{E}^\infty$. So, we **define** conditional probability as we want.

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 **conditional probability** of B given A by,

$$P(B|A) := \frac{P(A \cap B)}{P(A)} . \quad (2.2)$$

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:

$$P(B|A) : \mathcal{F} \rightarrow [0, 1]$$

that assigns to each $B \in \mathcal{F}$ a number in the interval $[0, 1]$, such that, the axioms of probability are satisfied:

Axiom (1): For any event B , $0 \leq P(B|A) \leq 1$.

Axiom (2): $P(\Omega|A) = 1$ Meaning ‘Something Happens given the event A happens’

Axiom (3): The ‘Addition Rule’ axiom holds, ie. for events $B_1, B_2 \in \mathcal{F}$,

$$B_1 \cap B_2 = \emptyset \text{ implies } P(B_1 \cup B_2|A) = P(B_1|A) + P(B_2|A) .$$

Axiom (4): For mutually exclusive events, B_1, B_2, \dots ,

$$P(B_1 \cup B_2 \cup \dots | A) = P(B_1|A) + P(B_2|A) + \dots .$$

From the definition of conditional probability we get the following properties or rules:

Complementation rule: $P(B|A) = 1 - P(B^c|A)$.

Addition rule for two arbitrary events B_1 and B_2 :

$$P(B_1 \cup B_2|A) = P(B_1|A) + P(B_2|A) - P(B_1 \cap B_2|A) .$$

Solving for $P(A \cap B)$ with these definitions of conditional probability gives another rule:

Multiplication rule for two likely events:

If A and B are events, and if $P(A) \neq 0$ and $P(B) \neq 0$, then

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

Example 31 (Wasserman03, p. 11) A medical test for a disease D has outcomes + and -. the probabilities are:

	Have Disease (D)	Don't have disease (D^c)
Test positive (+)	0.009	0.099
Test negative (-)	0.001	0.891

Using the definition of conditional probability, we can compute the conditional probability that you test positive given that you have the disease:

$$P(+|D) = \frac{P(+ \cap D)}{P(D)} = \frac{0.009}{0.009 + 0.001} = 0.9 ,$$

and the conditional probability that you test negative given that you don't have the disease:

$$P(-|D^c) = \frac{P(- \cap D^c)}{P(D^c)} = \frac{0.891}{0.099 + 0.891} \approx 0.9 .$$

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?

$$P(D|+) = \frac{P(D \cap +)}{P(+)} = \frac{0.009}{0.009 + 0.099} \approx 0.08$$

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?

Solution:

Let A = “the woman has breast cancer”, and B = “a positive test.”

We want $P(A|B)$ but what we are given is $P(B|A) = 0.9$.

By the definition of conditional probability,

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

To evaluate the numerator we use the multiplication rule

$$P(A \cap B) = P(A)P(B|A) = 0.01 \times 0.9 = 0.009$$

Similarly,

$$P(A^c \cap B) = P(A^c)P(B|A^c) = 0.99 \times 0.1 = 0.099$$

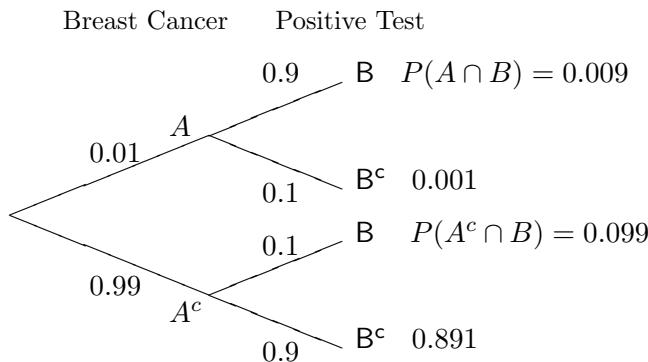
Now $P(B) = P(A \cap B) + P(A^c \cap B)$ so

$$P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{0.009}{0.009 + 0.099} = \frac{9}{108}$$

or a little less than 9%. This situation comes about because it is much easier to have a false positive for a healthy woman, which has probability 0.099, than to find a woman with breast cancer having a positive test, which has probability 0.009.

This answer is somewhat surprising. Indeed when ninety-five physicians were asked this question their 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 the answer $9/(9 + 99)$.

Alternative solution using a tree diagram:



So the probability that a woman has breast cancer given that she has just had a positive test is

$$P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{0.009}{0.009 + 0.099} = \frac{9}{108}$$

**In the exam, there won't be any need for electronic calculators and you may leave the answer in either of the last two numerical forms for full credit, provided you show the steps in your reasoning.*

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.

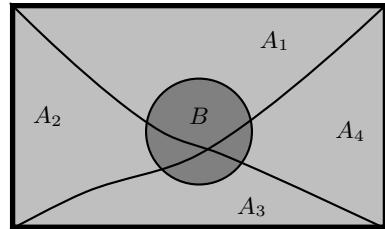
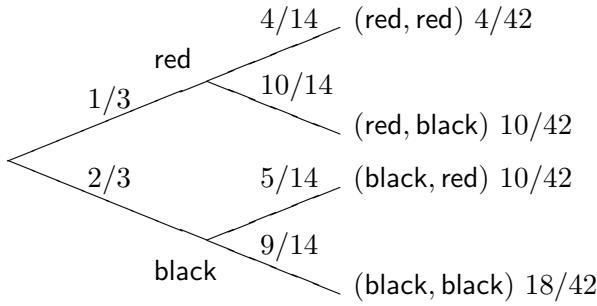


Figure 2.3: Reference to the Venn diagram will help you understand this idea behind the proof of the total probability theorem in Proposition 14 for the four event case.

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

$$P(\text{second ball is red}) = 4/42 + 10/42 = 1/3 .$$



Alternatively, use the total probability theorem to break the problem down into manageable pieces. Let $R_1 = \{(\text{red}, \text{red}), (\text{red}, \text{black})\}$ and $R_2 = \{(\text{red}, \text{red}), (\text{black}, \text{red})\}$ be the events corresponding to a red ball in the 1st and 2nd draws, respectively, and let $B_1 = \{(\text{black}, \text{red}), (\text{black}, \text{black})\}$ be the event of a black ball on the first draw.

Now R_1 and B_1 partition Ω so we can write:

$$\begin{aligned} P(R_2) &= P(R_2 \cap R_1) + P(R_2 \cap B_1) \\ &= P(R_2|R_1)P(R_1) + P(R_2|B_1)P(B_1) \\ &= (4/14)(1/3) + (5/14)(2/3) = 1/3 . \end{aligned}$$

Proposition 15 (Bayes' Theorem, 1763) Suppose the events $A_1, A_2, \dots, A_k \in \mathcal{F}$, with $P(A_h) > 0$ for each $h \in \{1, 2, \dots, k\}$, partition the sample space Ω , ie. they are mutually exclusive (disjoint) and exhaustive events with positive probability:

$$A_i \cap A_j = \emptyset, \text{ for any distinct } i, j \in \{1, 2, \dots, k\}, \quad \bigcup_{h=1}^k A_h = \Omega, \quad P(A_h) > 0$$

Thus, precisely one of the A_h 's will occur on any performance of our experiment \mathcal{E} .

Let $B \in \mathcal{F}$ be some event with $P(B) > 0$, then

$$P(A_h|B) = \frac{P(B|A_h)P(A_h)}{\sum_{h=1}^k P(B|A_h)P(A_h)} \quad (2.5)$$

Proof: We apply elementary set theory, the definition of conditional probability $k+2$ times and the addition rule once:

$$\begin{aligned} P(A_h|B) &= \frac{P(A_h \cap B)}{P(B)} = \frac{P(B \cap A_h)}{P(B)} = \frac{P(B|A_h)P(A_h)}{P(B)} \\ &= \frac{P(B|A_h)P(A_h)}{P\left(\bigcup_{h=1}^k (B \cap A_h)\right)} = \frac{P(B|A_h)P(A_h)}{\sum_{h=1}^k P(B \cap A_h)} \\ &= \frac{P(B|A_h)P(A_h)}{\sum_{h=1}^k P(B|A_h)P(A_h)} \end{aligned}$$

The operations done to the denominator in the proof above is merely the total probability theorem:

$$P(B) = \sum_{h=1}^k P(B|A_h)P(A_h)$$

We call $P(A_h)$ the **prior probability** of A_h , i.e., before observing B or *a priori*, and $P(A_h|B)$ 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 = “spam”, A_2 = “low priority”, and A_3 = “high priority”. From previous experience, he finds that $P(A_1) = 0.7$, $P(A_2) = 0.2$ and $P(A_3) = 0.1$. Note that $P(A_1 \cup A_2 \cup A_3) = P(\Omega) = 0.7 + 0.2 + 0.1 = 1$. Let B be the event that the email contains the word “free.” From previous experience, $P(B|A_1) = 0.9$, $P(B|A_2) = 0.01$ and $P(B|A_3) = 0.01$. Note that $P(B|A_1) + P(B|A_2) + P(B|A_3) = 0.9 + 0.01 + 0.01 \neq 1$. Now, suppose Larry receives an email with the word “free.” What is the probability that it is “spam,” “low priority,” and “high priority”?

Solution:

$$\begin{aligned} P(A_1|B) &= \frac{P(B|A_1) P(A_1)}{P(B|A_1) P(A_1) + P(B|A_2) P(A_2) + P(B|A_3) P(A_3)} = \frac{0.9 \times 0.7}{(0.9 \times 0.7) + (0.01 \times 0.2) + (0.01 \times 0.1)} = \frac{0.63}{0.633} \approx 0.995 \\ P(A_2|B) &= \frac{P(B|A_2) P(A_2)}{P(B|A_1) P(A_1) + P(B|A_2) P(A_2) + P(B|A_3) P(A_3)} = \frac{0.01 \times 0.2}{(0.9 \times 0.7) + (0.01 \times 0.2) + (0.01 \times 0.1)} = \frac{0.002}{0.633} \approx 0.003 \\ P(A_3|B) &= \frac{P(B|A_3) P(A_3)}{P(B|A_1) P(A_1) + P(B|A_2) P(A_2) + P(B|A_3) P(A_3)} = \frac{0.01 \times 0.1}{(0.9 \times 0.7) + (0.01 \times 0.2) + (0.01 \times 0.1)} = \frac{0.001}{0.633} \approx 0.002 \end{aligned}$$

Note that $P(A_1|B) + P(A_2|B) + P(A_3|B) = 0.995 + 0.003 + 0.002 = 1$.

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

$$P(A \cap B) = P(A) P(B) . \quad (2.6)$$

Let us make sense of this definition in terms of our previous definitions. When $P(A) = 0$ or $P(B) = 0$, both sides of the above equality are 0. If $P(A) \neq 0$, then rearranging the above equation we get:

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

But, the LHS is $P(B|A)$ by definition 2.2, and thus for independent events A and B , we get:

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

This says that information about the occurrence of A does not affect the occurrence of B . If $P(B) \neq 0$, then an analogous argument:

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

says that information about the occurrence of B does not affect the occurrence of A . Therefore, the probability of their joint occurrence $P(A \cap B)$ is simply the product of their individual probabilities $P(A)P(B)$.

Definition 17 (Independence of a sequence of events) We say that a finite or infinite sequence of events A_1, A_2, \dots are independent if whenever i_1, i_2, \dots, i_k are distinct elements from the set of indices \mathbb{N} , such that $A_{i_1}, A_{i_2}, \dots, A_{i_k}$ are defined (elements of \mathcal{F}), then

$$P(A_{i_1} \cap A_{i_2} \dots \cap A_{i_k}) = P(A_{i_1})P(A_{i_2}) \dots P(A_{i_k})$$

Example 35 (Some Standard Examples) A sequence of events in a sequence of independent trials is independent.

- (a) Suppose you toss a fair coin twice such that the first toss is independent of the second. Then,

$$P(\text{Heads on the first toss} \cap \text{Tails on the second toss}) = P(H)P(T) = \frac{1}{2} \times \frac{1}{2} = \frac{1}{4} .$$

- (b) Suppose you independently toss a fair die three times. Let E_i be the event that the outcome is an even number on the i -th trial. The probability of getting an even number in all three trials is:

$$\begin{aligned} P(E_1 \cap E_2 \cap E_3) &= P(E_1)P(E_2)P(E_3) \\ &= (P(\{2, 4, 6\}))^3 \\ &= (P(\{2\}) \cup \{4\} \cup \{6\}))^3 \\ &= (P(\{2\}) + P(\{4\}) + P(\{6\}))^3 \\ &= \left(\frac{1}{6} + \frac{1}{6} + \frac{1}{6}\right)^3 = \left(\frac{1}{2}\right)^3 = \frac{1}{8} . \end{aligned}$$

- (c) Suppose you toss a fair coin independently m times. Then each of the 2^m possible outcomes in the sample space Ω has equal probability of $\frac{1}{2^m}$ due to independence.

Example 36 (dependence and independence) Suppose we toss two fair dice. Let A denote the event that the sum of the dice is six and B denote the event that the first die equals four. The sample space encoding the thirty six ordered pairs of outcomes for the two dice is $\Omega = \{(1, 1), (1, 2), \dots, (1, 6), (2, 1), \dots, (2, 6), \dots, (5, 6), (6, 6)\}$ and due to independence $P(\omega) = 1/36$ for each $\omega \in \Omega$. Then

$$P(A \cap B) = P(\{(4, 2)\}) = \frac{1}{36} ,$$

but

$$\begin{aligned} P(A)P(B) &= P(\{(1, 5), (2, 4), (3, 3), (4, 2), (5, 1)\})P(\{(4, 1), (4, 2), (4, 3), (4, 4), (4, 5), (4, 6)\}) \\ &= \frac{5}{36} \times \frac{6}{36} = \frac{5}{36} \times \frac{1}{6} = \frac{5}{216} , \end{aligned}$$

and therefore A and B are not independent. The reason for the events A and B being dependent is clear because the chance of getting a total of six depends on the outcome of the first die (not being six).

Now, let C be the event that the sum of the two dice equals seven. Then

$$P(C \cap B) = P(\{(4, 3)\}) = \frac{1}{36},$$

while

$$\begin{aligned} P(C \cap B) &= P(\{(1, 6), (2, 5), (3, 4), (4, 3), (5, 2), (6, 1)\}) P(\{(4, 1), (4, 2), (4, 3), (4, 4), (4, 5), (4, 6)\}) \\ &= \frac{6}{36} \times \frac{6}{36} = \frac{1}{36}, \end{aligned}$$

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,

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

but,

$$\frac{1}{4} = P(\{1\}) = P(A \cap B \cap C) \neq P(A) P(B) P(C) = \frac{2}{4} \times \frac{2}{4} \times \frac{2}{4} = \frac{1}{8}.$$

Therefore, inspite of being pairwise independent, the events A , B and C are not jointly independent.

CONDITIONAL PROBABILITY SUMMARY

$P(A|B)$ 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 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.

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.

- 1.List the sample space for the experiment if we note the numbers on the 2 upturned faces.
- 2.What is the probability of obtaining a sum greater than 4 but less than 7?

Ex. 2.8 — Based on past experience, 70% of students in a certain course pass the midterm test. The final exam is passed by 80% of those who passed the midterm test, but only by 40% of those who fail the midterm test. What fraction of students pass the final exam?

Ex. 2.9 — A small brewery has two bottling machines. Machine 1 produces 75% of the bottles and machine 2 produces 25%. One out of every 20 bottles filled by machine 1 is rejected for some reason, while one out of every 30 bottles filled by machine 2 is rejected. What is the probability that a randomly selected bottle comes from machine 1 given that it is accepted?

Ex. 2.10 — A process producing micro-chips, produces 5% defective, at random. Each micro-chip is tested, and the test will correctly detect a defective one $4/5$ of the time, and if a good micro-chip is tested the test will declare it is defective with probability $1/10$.

- (a)If a micro-chip is chosen at random, and tested to be good, what was the probability that it was defective anyway?
- (b)If a micro-chip is chosen at random, and tested to be defective, what was the probability that it was good anyway?
- (c)If 2 micro-chips are tested and determined to be good, what is the probability that at least one is in fact defective?

Ex. 2.11 — Suppose that $\frac{2}{3}$ of all gales are force 1, $\frac{1}{4}$ are force 2 and $\frac{1}{12}$ are force 3. Furthermore, the probability that force 1 gales cause damage is $\frac{1}{4}$, the probability that force 2 gales cause damage is $\frac{2}{3}$ and the probability that force 3 gales cause damage is $\frac{5}{6}$.

- (a)If a gale is reported, what is the probability of it causing damage?
- (b)If the gale caused damage, find the probabilities that it was of: force 1; force 2; force 3.
- (c)If the gale did NOT cause damage, find the probabilities that it was of: force 1; force 2; force 3.

Ex. 2.12 — **The sensitivity and specificity of a medical diagnostic test for a disease are defined as follows:

$$\begin{aligned}\text{sensitivity} &= P(\text{test is positive} \mid \text{patient has the disease}) , \\ \text{specificity} &= P(\text{test is negative} \mid \text{patient does not have the disease}) .\end{aligned}$$

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%, find

- (a)the probability that a patient who tests positive actually has the disease,
- (b)the probability that a patient who tests negative is free from the disease.

Ex. 2.13 — **The detection rate and false alarm rate of an intrusion sensor are defined as

$$\begin{aligned}\text{detection rate} &= P(\text{detection declared} \mid \text{intrusion}) , \\ \text{false alarm rate} &= P(\text{detection declared} \mid \text{no intrusion}) .\end{aligned}$$

If the detection rate is 0.999 and the false alarm rate is 0.001, and the probability of an intrusion occurring is 0.01, find

- (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 $P(A) \neq 0$ and $P(B) \neq 0$. When A and B are disjoint, are they also independent? Explain clearly why or why not.

Chapter 3

Random Variables

We are used to classical variables such as x as an “unknown” in the equation: $x + 3 = 7$.

We also use classical variables to represent geometric objects such as a line:

$$y = 3x - 2,$$

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 $\mathbb{R} = (-\infty, \infty)$.

Yet another example is the use of variables to represent sequences such as:

$$\{a_n\}_{n=1}^{\infty} = a_1, a_2, a_3, \dots .$$

What these *classical variables* have in common is that they *take a fixed or deterministic value* when we can solve for them.

We need a different kind of variable to deal with real-world situations where the same variable may take different values in a non-deterministic manner. **Random variables** do this job for us. Random variables, unlike classical deterministic variables, can take a bunch of different values.

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:

Experiment	Possible measured outcomes
Counting the number of typos up to now	$\mathbb{Z}_+ := \{0, 1, 2, \dots\} \subset \mathbb{R}$
Length in centi-meters 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}

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, $X : \Omega \rightarrow \mathbb{R}$ that should satisfy certain conditions to keep the meaning of the underlying probability space (Ω, \mathcal{F}, P) . Let us go through some examples before giving the formal definition of such a real-valued or \mathbb{R} -valued random variable.

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 $\Omega = \{\text{rain, not rain}\}$. We can associate a random variable X with this experiment as follows:

$$X(\omega) = \begin{cases} 1, & \text{if } \omega = \text{rain} \\ 0, & \text{if } \omega = \text{not rain} \end{cases}$$

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 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,

“What is the probability of it raining tomorrow?”

in the rain/not experiment of Example 38 becomes

“What is $P(\{\omega : X(\omega) = 1\})$?”

or, more simply,

“What is $P(X = 1)$?”

With this motivation we are ready to formally define such a random variable.

3.1 Basic Definitions

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.

Definition 18 (Random Variable) Let (Ω, \mathcal{F}, P) be some probability triple. Then, a **Random Variable (RV)**, say X , is a function from the sample space Ω to the set of real numbers \mathbb{R}

$$X : \Omega \rightarrow \mathbb{R}$$

such that for every $x \in \mathbb{R}$, the inverse image of the half-open real interval $(-\infty, x]$ is an element of the collection of events \mathcal{F} , i.e.:

$$\text{for every } x \in \mathbb{R}, \quad X^{[-1]}((-\infty, x]) := \{\omega : X(\omega) \leq x\} \in \mathcal{F}.$$

This definition can be summarised by the statement that a RV is an \mathcal{F} -measurable map. We assign probability to the RV X as follows:

$$P(X \leq x) = P(X^{[-1]}((-\infty, x])) := P(\{\omega : X(\omega) \leq x\}). \quad (3.1)$$

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:

$$F(x) := P(X \leq x) = P(\{\omega : X(\omega) \leq x\}), \quad \text{for any } x \in \mathbb{R}. \quad (3.2)$$

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 $X \sim F$.

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

$$P(2 \leq X \leq 3)$$

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 (Ω, \mathcal{F}, 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.

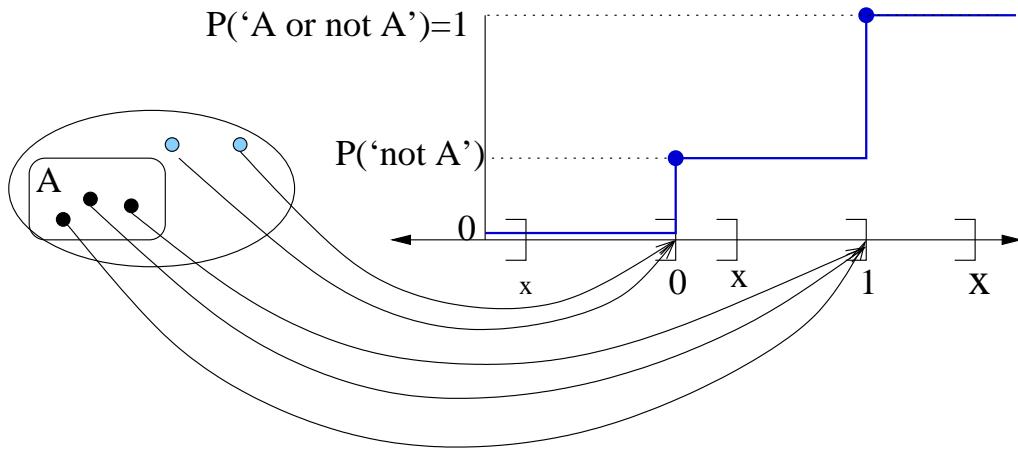
Model 1 is called the Bernoulli RV for event A with a known probability $P(A)$. We will define as our next model the Bernoulli(θ) RV by introducing a parameter $\theta \in [0, 1]$ for the typically unknown probability $P(A)$.

Some useful properties of the Indicator Function are:

$$\mathbb{1}_{A^c} = 1 - \mathbb{1}_A, \quad \mathbb{1}_{A \cap B} = \mathbb{1}_A \mathbb{1}_B, \quad \mathbb{1}_{A \cup B} = \mathbb{1}_A + \mathbb{1}_B - \mathbb{1}_A \mathbb{1}_B$$

We slightly abuse notation when A is a single element set by ignoring the curly braces.

Figure 3.1: The Indicator function of event $A \in \mathcal{F}$ is a RV $\mathbb{1}_A$ with DF F
DF



Exercise 3.1 (Drawing discontinuous functions) Identify the mistakes in how the $\mathbb{1}_A$ is drawn as a discontinuous function in Figure 3.1.

Classwork 41 (A random variable with three values and eight sample points) Consider the RV X of Figure 3.2. First draw this properly as done in Ex. 3.1. Let the events $A = \{\omega_1, \omega_2\}$, $B = \{\omega_3, \omega_4, \omega_5\}$ and $C = \{\omega_6, \omega_7, \omega_8\}$. Define the RV X formally. What sets should \mathcal{F} minimally include? What do you need to do to make sure that \mathcal{F} is a sigma algebra?

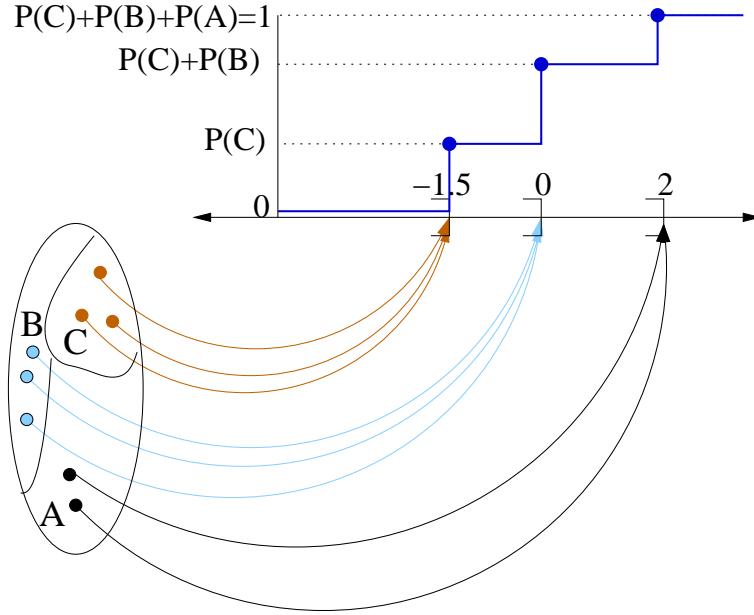
Exercise 3.2 (Fair coin toss RV) Consider the *fair coin toss experiment* with $\Omega = \{H, T\}$ and $P(H) = P(T) = 1/2$.

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:

$$X(\omega) = \begin{cases} 1, & \text{if } \omega = H \\ 0, & \text{if } \omega = T \end{cases}$$

Find the distribution function for X .

Figure 3.2: A RV X from a sample space Ω with 8 elements to \mathbb{R} and its DF F .



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 (Ω, \mathcal{F}, P) . We define the **probability mass function** (PMF) f of X to be the function $f : \mathbb{R} \rightarrow [0, 1]$ defined as follows:

$$f(x) := P(X = x) = P(\{\omega : X(\omega) = x\}) = \begin{cases} \theta_i & \text{if } x = x_i \in \mathbb{X}. \\ 0 & \text{otherwise.} \end{cases} \quad (3.5)$$

The DF F and PMF f for a discrete RV X satisfy the following:

1. For any $x \in \mathbb{R}$,

$$P(X \leq x) = F(x) = \sum_{x_i \leq x} f(x_i) = \sum_{x_i \leq x} \theta_i. \quad (3.6)$$

2. For any $a, b \in \mathbb{R}$ with $a < b$,

$$P(a < X \leq b) = F(b) - F(a) = \sum_{a < x_i \leq b} \theta_i. \quad (3.7)$$

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 $P(\Omega) = 1$, we get that the sum of all the probabilities is 1:

$$\sum_i \theta_i = 1 . \quad (3.8)$$

4. When X only has finitely many possibilities, say k with $\mathbb{X} = \{x_1, x_2, \dots, x_k\}$, then we may think of the probability P specified by $(\theta_1, \theta_2, \dots, \theta_k)$ as a point in the **unit** $(k-1)$ **simplex**:

$$\Delta^{k-1} := \{(\theta_1, \theta_2, \dots, \theta_k) \in \mathbb{R}^k : \sum_i \theta_i = 1 \text{ and } \theta_i \geq 0, \text{ for all } i\} \quad (3.9)$$

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**:

$$\Delta^1 := \{(\theta, 1 - \theta) \in \mathbb{R}^2 : 0 \leq \theta \leq 1\} . \quad (3.10)$$

See <https://en.wikipedia.org/wiki/Simplex> for the images scribed on the board.

DISCRETE RANDOM VARIABLES - SIMPLIFIED NOTATION

Notice that in equations (3.5), (3.6) and (3.7) the use of the “ $\omega \in \Omega$ ” notation, where random variables are defined as functions, is much reduced. The reason is that in straightforward examples it is convenient to associate the possible values x_1, x_2, \dots with the outcomes $\omega_1, \omega_2, \dots$. Hence, we can describe a discrete random variable by the table:

Possible values: x_i	x_1	x_2	x_3	\dots
Probability: $P(X = x_i) = \theta_i$	θ_1	θ_2	θ_3	\dots

It is customary to use p_i instead of θ_i for the probabilities. But we try to avoid it as it will hurt us when we start doing Inference Theory soon!

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 students of mathematics and should know more about what is “under the hood”.

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:

- Discrete uniform random variables with finitely many possibilities
- Discrete non-uniform random variables with finitely many possibilities
- Discrete non-uniform random variables with (countably) infinitely many possibilities

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:

$$f(x) = \begin{cases} \theta_i = \frac{1}{k} & \text{if } x = x_i, \text{ where } i = 1, 2, \dots, k, \\ 0 & \text{otherwise.} \end{cases} \quad (3.11)$$

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

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

and its distribution function is given by

$$F(x) = P(X \leq x) = \begin{cases} 0, & \text{if } -\infty < x < 0 \\ \frac{1}{2}, & \text{if } 0 \leq x < 1 \\ 1, & \text{if } 1 \leq x < \infty \end{cases}$$

Let us sketch the probability mass function and distribution function for X below.

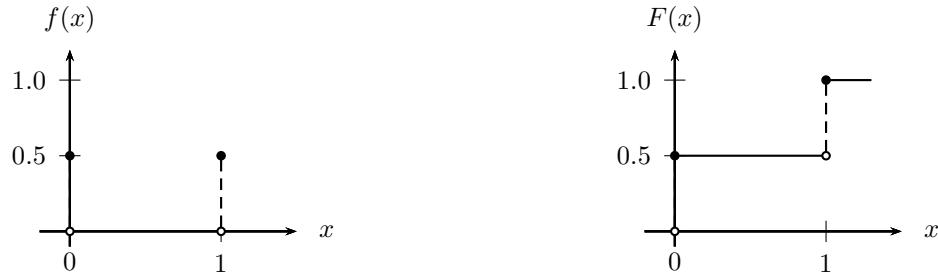


Figure 3.3: $f(x)$ and $F(x)$ of the *fair coin toss* random variable X , a discrete uniform RV on $\{0, 1\}$.

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	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$

Find the probability mass function and distribution function for this random variable, and sketch their graphs.

Solution:

The probability mass function of this random variable is:

$$f(x) = P(X = x) = \begin{cases} \frac{1}{6} & \text{if } x = 1 \\ \frac{1}{6} & \text{if } x = 2 \\ \frac{1}{6} & \text{if } x = 3 \\ \frac{1}{6} & \text{if } x = 4 \\ \frac{1}{6} & \text{if } x = 5 \\ \frac{1}{6} & \text{if } x = 6 \\ 0 & \text{otherwise} \end{cases}$$

and the distribution function is:

$$F(x) = P(X \leq x) = \begin{cases} 0, & \text{if } -\infty < x < 1 \\ \frac{1}{6}, & \text{if } 1 \leq x < 2 \\ \frac{1}{3}, & \text{if } 2 \leq x < 3 \\ \frac{1}{2}, & \text{if } 3 \leq x < 4 \\ \frac{2}{3}, & \text{if } 4 \leq x < 5 \\ \frac{5}{6}, & \text{if } 5 \leq x < 6 \\ 1, & \text{if } 6 \leq x < \infty \end{cases}$$

Example 44 (Astragali with a Kiwi sheep ankle bone) Astragali. Board games involving chance were known in Egypt, 3000 years before Christ. The element of chance needed for these games was at first provided by tossing astragali, the ankle bones of sheep. These bones could come to rest on only four sides, the other two sides being rounded. The upper side of the bone, broad 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.

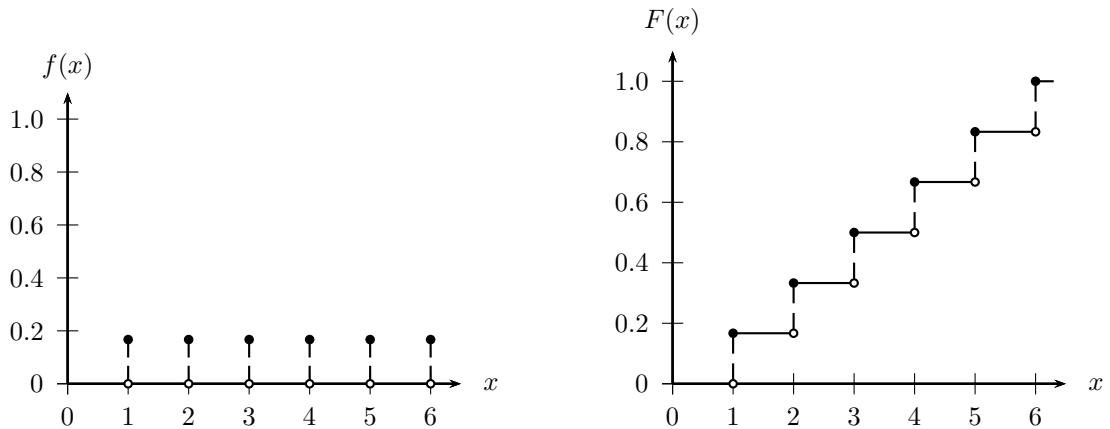


Figure 3.4: $f(x)$ and $F(x)$ of the fair die toss random variable X , a discrete uniform RV on $\{1, 2, 3, 4, 5, 6\}$.

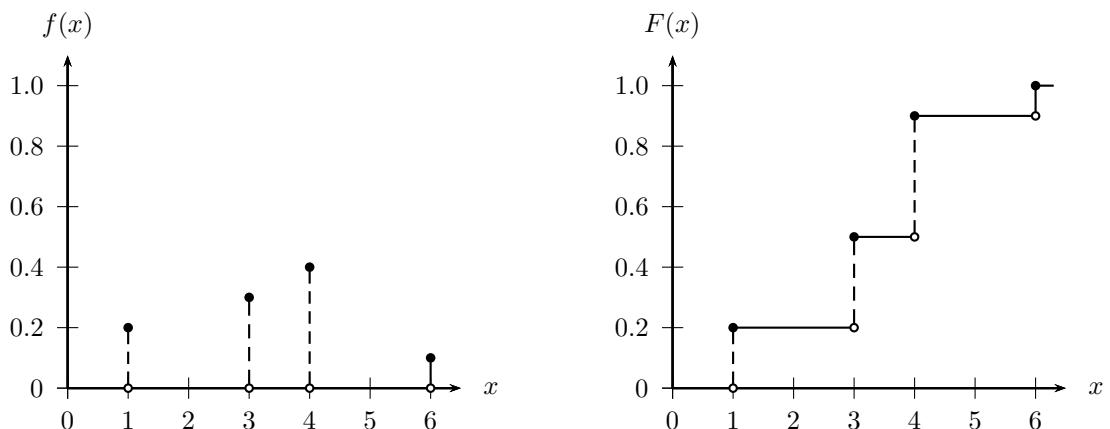


Figure 3.5: $f(x)$ and $F(x)$ of surmised astragali toss random variable X , a discrete (non-uniform) RV on $\{1, 2, 3, 4\}$.

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 (Ω, \mathcal{F}, 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 $\text{Bernoulli}(\theta)$ 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 $\text{Bernoulli}(\theta)$ RV X is:

$$f(x; \theta) = \theta^x(1 - \theta)^{1-x}\mathbf{1}_{\{0,1\}}(x) = \begin{cases} \theta & \text{if } x = 1, \\ 1 - \theta & \text{if } x = 0, \\ 0 & \text{otherwise} \end{cases} \quad (3.13)$$

and its DF is:

$$F(x; \theta) = \begin{cases} 1 & \text{if } 1 \leq x, \\ 1 - \theta & \text{if } 0 \leq x < 1, \\ 0 & \text{otherwise} \end{cases} \quad (3.14)$$

We emphasise 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. $P_\theta(X = 1) = \theta$ and $P_\theta(X = 0) = 1 - \theta$.



Figure 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

3.2.2 Independent Bernoulli Trials

Random variables make sense for a series of trials as well as just a single trial of an experiment. We now look at what happens when we perform a sequence of independent Bernoulli trials. For instance:

- Flip a coin 10 times; count the number of heads
 - by possibly allowing for the coin's $P(H)$ to change each time because each of them are manufactured in a terrible mint.
- Test 50 randomly selected circuits from an assembly line; count the number of defective circuits.
- Roll a die 100 times; count the number of sixes you throw.
- 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.

Since the $\text{Bernoulli}(\theta)$ 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 deterministic variable or parameter θ . Now consider doing more than one trial so we have sequence of $\text{Bernoulli}(\theta_i)$ trials, say,

$$X_i \sim \text{Bernoulli}(\theta_i) \text{ with } i \in \mathbb{N},$$

with each $\theta_i \in [0, 1]$ being possibly unknown but fixed as a parameter. Now, if we assume independence across trials, so one trial's outcome does not affect the outcome of any of the other trials, in the sense of Definition 17 about *independence of a sequence of events*, then we can obtain the probability of the entire sequence of outcomes for this sequence of **independently distributed** $\text{Bernoulli}(\theta_i)$ **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:

$$\begin{aligned} P(x; (\theta_1, \theta_2, \dots)) &= \prod_i f(x_i; \theta_i) = \prod_i \theta_i^{x_i} (1 - \theta_i)^{1-x_i} \mathbf{1}_{\{0,1\}}(x_i), \\ &\text{where } x := (x_1, x_2, \dots) \in \mathbb{X}_\infty = \{0, 1\}^\infty := \{0, 1\} \times \{0, 1\} \times \cdots \end{aligned} \quad (3.15)$$

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)** $\text{Bernoulli}(\theta)$ **trails** as follows:

$$\begin{aligned} P(x; \theta) &= \prod_i f(x_i; \theta) = \prod_i \theta^{x_i} (1 - \theta)^{1-x_i} \mathbf{1}_{\{0,1\}}(x_i) = \mathbf{1}_{\{0,1\}^\infty}(x) \theta^{\sum_i x_i} (1 - \theta)^{\sum_i (1-x_i)} \\ &= \begin{cases} \theta^{\sum_i x_i} (1 - \theta)^{n - \sum_i x_i} & \text{if } x := (x_1, x_2, \dots) \in \mathbb{X}_\infty = \{0, 1\}^\infty \\ 0 & \text{otherwise} \end{cases} \end{aligned} \quad (3.16)$$

Remembering that all other RVs can be derived from such IID $\text{Bernoulli}(\theta)$ 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

$$\begin{aligned} f(1) &= P(X = 1) = P(H) = \frac{1}{2}, \\ f(2) &= P(X = 2) = P(TH) = \frac{1}{2} \cdot \frac{1}{2} = \left(\frac{1}{2}\right)^2, \\ f(3) &= P(X = 3) = P(TTH) = \frac{1}{2} \cdot \frac{1}{2} \cdot \frac{1}{2} = \left(\frac{1}{2}\right)^3, \quad \text{etc.} \end{aligned}$$

the probability mass function of X is:

$$f(x) = P(X = x) = \left(\frac{1}{2}\right)^x, \quad x = 1, 2, \dots.$$

In the previous Example, noting that we have independent trials, we get:

$$f(x) = P(X = x) = P(\underbrace{TT\dots T}_n H) = P(T)^{x-1} P(H) = \left(\frac{1}{2}\right)^{x-1} \frac{1}{2}.$$

More generally, let there be two possibilities, success (S) or failure (F), with $P(S) = \theta$ and $P(F) = 1 - \theta$ so that:

$$P(X = x) = P(\underbrace{FF\dots F}_x S) = (1 - \theta)^{x-1} \theta.$$

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.

Model 4 (Geometric(θ) RV) Given a parameter $\theta \in (0, 1)$, the PMF of the Geometric(θ) RV X is

$$f(x; \theta) = \begin{cases} \theta(1 - \theta)^x & \text{if } x \in \mathbb{Z}_+ := \{0, 1, 2, \dots\} \\ 0 & \text{otherwise} \end{cases} \quad (3.17)$$

It is straightforward to verify that $f(x; \theta)$ is indeed a PMF :

$$\sum_{x=0}^{\infty} f(x; \theta) = \sum_{x=0}^{\infty} \theta(1 - \theta)^x = \theta \left(\frac{1}{1 - (1 - \theta)} \right) = \theta \left(\frac{1}{\theta} \right) = 1$$

The above equality is a consequence of the geometric series identity (3.18) with $a = \theta$ and $\vartheta := 1 - \theta$:

$$\sum_{x=0}^{\infty} a\vartheta^x = a \left(\frac{1}{1 - \vartheta} \right), \quad \text{provided, } 0 < \vartheta < 1. \quad (3.18)$$

Proof:

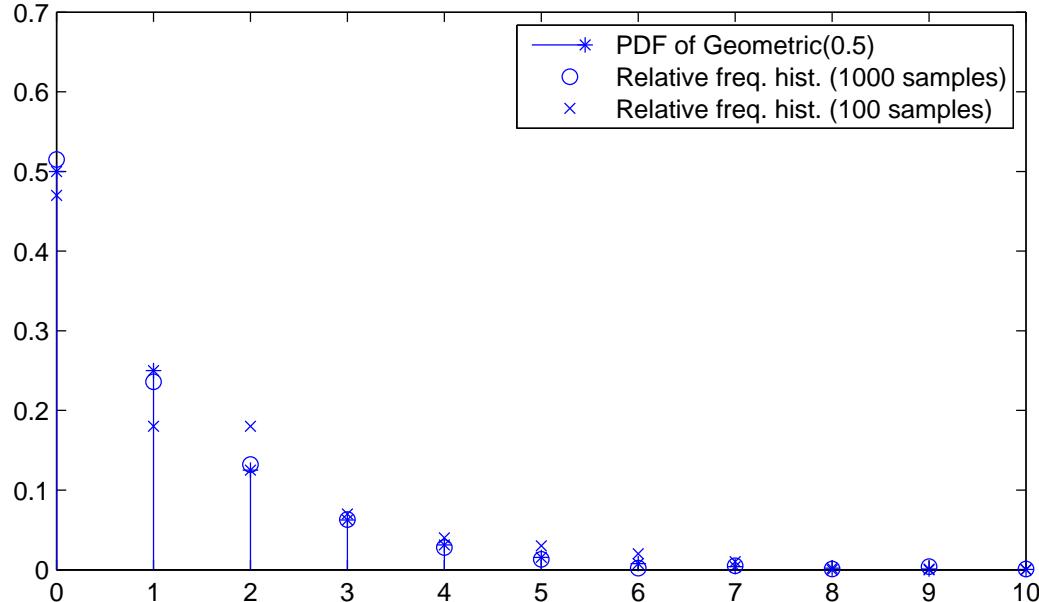
$$a + a\vartheta + a\vartheta^2 + \cdots + a\vartheta^n = \sum_{0 \leq x \leq n} a\vartheta^x = a + \sum_{1 \leq x \leq n} a\vartheta^x = a + \vartheta \sum_{1 \leq x \leq n} a\vartheta^{x-1} = a + \vartheta \sum_{0 \leq x \leq n-1} a\vartheta^x = a + \vartheta \sum_{0 \leq x \leq n} a\vartheta^x - a\vartheta^{n+1}$$

Therefore,

$$\begin{aligned} \sum_{0 \leq x \leq n} a\vartheta^x &= a + \vartheta \sum_{0 \leq x \leq n} a\vartheta^x - a\vartheta^{n+1} \\ \left(\sum_{0 \leq x \leq n} a\vartheta^x \right) - \left(\vartheta \sum_{0 \leq x \leq n} a\vartheta^x \right) &= a - a\vartheta^{n+1} \\ \left(\sum_{0 \leq x \leq n} a\vartheta^x \right) (1 - \vartheta) &= a(1 - \vartheta^{n+1}) \\ \sum_{0 \leq x \leq n} a\vartheta^x &= a \left(\frac{1 - \vartheta^{n+1}}{1 - \vartheta} \right) \\ \sum_{x=0}^{\infty} a\vartheta^x := \lim_{n \rightarrow \infty} \sum_{0 \leq x \leq n} a\vartheta^x &= a \left(\frac{1}{1 - \vartheta} \right), \text{ provided, } 0 < \vartheta < 1 \end{aligned}$$

The outcome of a Geometric(θ) RV can be thought of as “the number of tosses needed before the appearance of the first ‘Head’ when tossing a coin with probability of ‘Heads’ equal to θ in a independent and identical manner.”

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 170 and Labwork 171 you will see in the sequel.



Exercise 3.3 (Coupon Collector’s Problem) Recall the Coupon Collector’s Problem from lectures.

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

$$\underbrace{\frac{1}{2} \frac{1}{2} \frac{1}{2}}_{3 \text{ successes}} \quad \underbrace{\frac{1}{2} \frac{1}{2} \cdots \frac{1}{2}}_{7 \text{ failures}}.$$

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

$$P(\text{'3 heads'}) = \binom{10}{3} \left(\frac{1}{2}\right)^3 \left(\frac{1}{2}\right)^7 \approx 0.117$$

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:

$$f(x; n, \theta) = \begin{cases} \binom{n}{x} \theta^x (1 - \theta)^{n-x} & \text{if } x \in \{0, 1, 2, 3, \dots, n\}, \\ 0 & \text{otherwise} \end{cases}, \quad (3.19)$$

where, $\binom{n}{x}$ is:

$$\binom{n}{x} = \frac{n(n-1)(n-2)\dots(n-x+1)}{x(x-1)(x-2)\cdots(2)(1)} = \frac{n!}{x!(n-x)!}.$$

$\binom{n}{x}$ is read as “ n choose x .”

A Quick Justification: The argument from Example 46 generalises as follows. Since the trials are independent and identical, the probability of x successes followed by $n - x$ failures, *in order*, is given by

$$\underbrace{\text{S}\text{S}\dots\text{S}}_x \underbrace{\text{F}\text{F}\dots\text{F}}_{n-x} = \theta^x (1 - \theta)^{n-x}.$$

Since the n symbols $\text{S}\text{S}\dots\text{S}\text{F}\text{F}\dots\text{F}$ may be arranged in

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

ways, the probability of x successes and $n - x$ failures, *in any order*, is given by

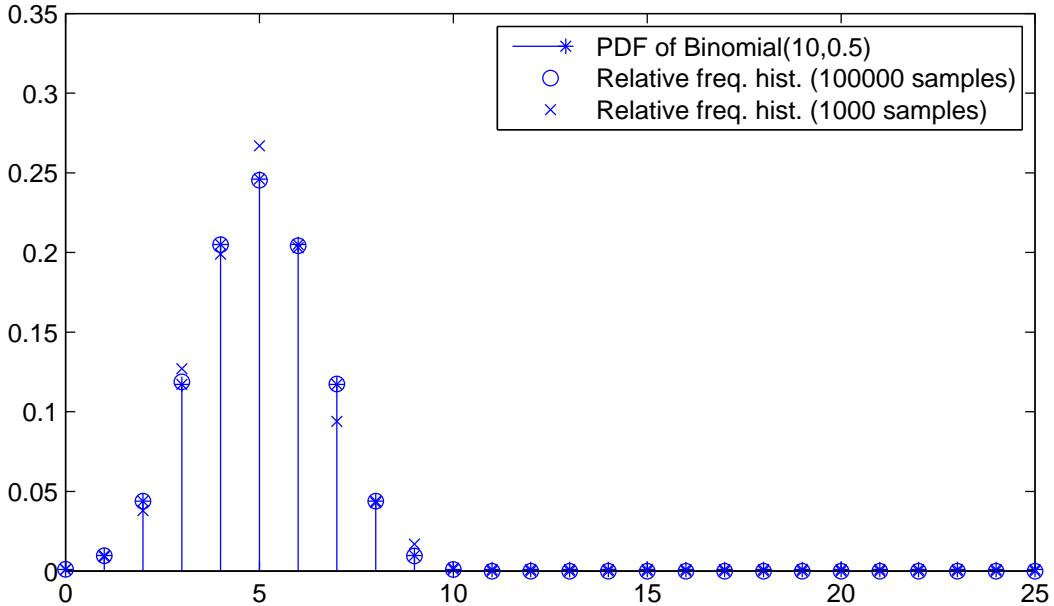
$$\binom{n}{x} \theta^x (1 - \theta)^{n-x}.$$

Proof: This is only a sketch. A formal proof should start with the mathematical induction for the very formula for the binomial coefficient.

Observe that for the Binomial(n, θ) RV X , $P(X = x) = f(x; n, \theta)$ is the probability that x of the n Bernoulli(θ) trials result in an outcome of 1's. Next note that if all $n X_i$'s are 0's, then $X = 0$, and if all $n X_i$'s are 1's, then $X = n$. In general, if some of the $n X_i$'s are 1's and the others are 0, then X can only take values in $\{0, 1, 2, \dots, n\}$ and therefore $f(x; n, \theta) = 0$ if $x \notin \{0, 1, 2, \dots, n\}$.

Now, let us compute $f(x; n, \theta)$ when $x \in \{0, 1, 2, \dots, n\}$. Consider the set of indices $\{1, 2, 3, \dots, n\}$ for the n IID Bernoulli(θ) RVs $\{X_1, X_2, \dots, X_n\}$. Now choose x indices from $\{1, 2, \dots, n\}$ to mark those trials in a particular realization of $\{x_1, x_2, \dots, x_n\} \in \{0, 1\}^n := \{\text{all binary } (0-1) \text{ strings of length } n\}$, specified by a choice of x trial indices with Bernoulli outcome 1, the binomial RV $X = \sum_{i=1}^n X_i$ takes the value x . Since there are exactly $\binom{n}{x}$ many ways in which we can choose x trial indices (with outcome 1) from the set of n trial indices $\{1, 2, \dots, n\}$, we get the desired product for $f(x; n, \theta) = \binom{n}{x} \theta^x (1 - \theta)^{n-x}$ when $x \in \{0, 1, \dots, n\}$.

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 174.



Example 47 Find the probability that seven of ten persons will recover from a tropical disease where the probability is identically 0.80 that any one of them will recover from the disease.

Solution:

We can assume independence here, so we have a binomial situation with $x = 7$, $n = 10$, and $\theta = 0.8$. Substituting these into the formula for the probability mass function for Binomial(10, 0.8)

random variable, we get:

$$\begin{aligned}
f(7; 10, 0.8) &= \binom{10}{7} \times (0.8)^7 \times (1 - 0.8)^{10-7} \\
&= \frac{10!}{(10-7)!7!} \times (0.8)^7 \times (1 - 0.8)^{10-7} \\
&= 120 \times (0.8)^7 \times (1 - 0.8)^{10-7} \quad * \text{In the exam you can give your answer as such an expression.} \\
&\approx 0.20
\end{aligned}$$

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.

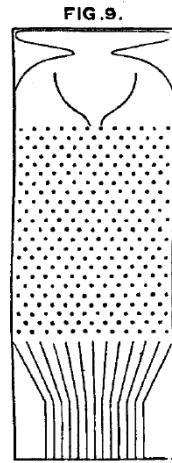
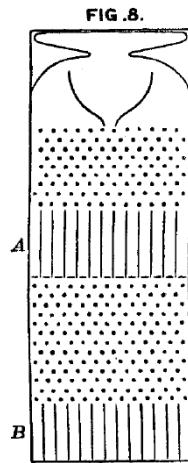
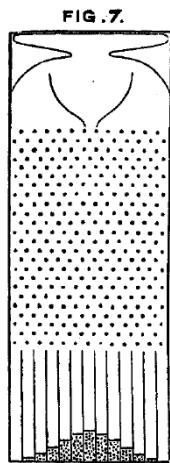
Sir Francis Galton's Quincunx

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 in the sequel as a way of giving precise meanings to Galton's observations with his Quincunx. “The Charms of Statistics.—*It is difficult to understand why statisticians commonly limit their inquiries to Averages, and do not revel in more comprehensive views. Their souls seem as dull to the charm of variety as that of the native of one of our flat English counties, whose retrospect of Switzerland was that, it its mountains could be thrown into its lakes, two nuances would be got rid of at once. An Average is but a solitary fact, whereas if a single other fact be added to it, an entire Normal Scheme, which nearly corresponds to the observed one, starts potentially into existence.*

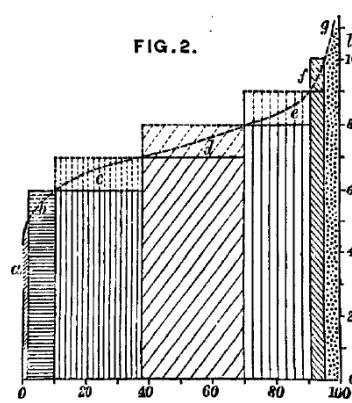
Some people hate the very name of statistics, but I find them full of beauty and interest. Whenever they are not brutalised, but delicately handled by the higher methods, and are warily interpreted, their power of dealing with complicated phenomenon is extraordinary. They are the

only tools by which an opening can be cut through the formidable thicket of difficulties that bars the path of those who pursue the Science of man.

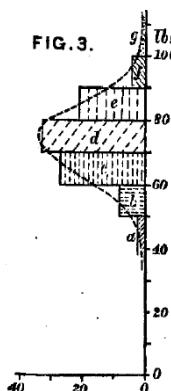
Figure 3.9: Figures from Sir Francis Galton, F.R.S., *Natural Inheritance*, Macmillan, 1889.



(a) FIG. 7, FIG. 8, and FIG. 9 (p. 63)



(b) FIG. 2 and FIG. 3 (p. 38)



Mechanical Illustration of the Cause of the Curve of Frequency.—*The Curve of Frequency, and that of Distribution, are convertible: therefore if the genesis of either of them can be made clear, that of the other also becomes intelligible. I shall now illustrate the origin of the Curve of Frequency, by means of an apparatus shown in Fig. 7, that mimics in a very pretty way the conditions on which Deviation depends.* It is a frame glazed in front, leaving a depth of about a quarter of an inch behind the glass. Strips are placed in the upper part to act as a funnel. Below the outlet of the funnel stand a succession of rows of pins stuck squarely into the backboard, and below these again are a series of vertical compartments. A charge of small shot is inclosed. When the frame is held topsy-turvy, all the shot runs to the upper end; then, when it is turned back into its working position, the desired action commences. Lateral strips, shown in the diagram, have the effect of directing all the shot that had collected at the upper end of the frame to run into the wide mouth of the funnel. The shot passes through the funnel and issuing from its narrow end, scampers deviously down through the pins in a curious and interesting way; each of them darting a step to the right or left, as the case may be, every time it strikes a pin. The pins are disposed in a quincunx fashion, so that every descending shot strikes against a pin in each successive row. The cascade issuing from the funnel broadens as it descends, and, at length, every shot finds itself caught in a compartment immediately after freeing itself from the last row of pins. The outline of the columns of shot that accumulate in the successive compartments approximates to the Curve of Frequency (Fig. 3, p. 38), and is closely of the same shape however often the experiment is repeated. The outline of the columns would 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 .

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:

$$\begin{aligned}
 P(X \leq 3) &= f(0; 2) + f(1; 2) + f(2; 2) + f(3; 2) \\
 &= e^{-2} \left(\frac{2^0}{0!} + \frac{2^1}{1!} + \frac{2^2}{2!} + \frac{2^3}{3!} \right) \quad * \text{This is a perfectly fine answer in the exam.} \\
 &= 0.857 \quad (3 \text{ sig. fig.})
 \end{aligned}$$

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?

Solution:

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

$$\lambda = \frac{8}{4} = 2.$$

(a)

$$P(X = 0) = f(0; 2) = \frac{e^{-2} 2^0}{0!} = 0.135 \quad (3 \text{ sig. fig})$$

(b)

$$\begin{aligned}
 P(X \geq 3) &= 1 - P(X < 3) \\
 &= 1 - P(X = 0) - P(X = 1) - P(X = 2) \\
 &= 1 - f(0; 2) - f(1; 2) - f(2; 2) \\
 &= 1 - 0.1353 - 0.2707 - 0.2707 \\
 &= 0.323 \quad (3 \text{ sig. fig.})
 \end{aligned}$$

Remark 24 In the binomial case where θ is small and n is large, it can be shown that the binomial distribution with parameters n and θ is closely approximated by the Poisson distribution having $\lambda = n\theta$. The smaller the value of θ and larger the value of n , the better the approximation.

In the sequel we will see more formally, after understanding notions of convergence of RVs, that the sum of a sequence of n IID Bernoulli(θ) RVs with $\lambda = n\theta$ converges to the Poisson(λ) RV as $n \rightarrow \infty$ and $\theta \rightarrow 0$ in a specific sense.

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

$$P(X \leq 1) = P(X = 0) + P(X = 1) = f(0; 0.5) + f(1; 0.5) = 0.910 \quad (\text{3 sig. fig.})$$

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

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

Distribution function

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

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: 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: $\lambda = \text{rate} \in (0, \infty)$.

3.3 Exercises in Discrete Random Variables

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?

x	1	2	3	4	5
$P(X = x)$	0.07	0.10	1.10	0.32	0.40

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 $f(1) = 0.1$, $f(2) = 0.2$, $f(3) = 0.2$, $f(4) = 0.2$, $f(5) = 0.3$.

1. Find the distribution function, F , for X , and graph both f and F .
2. Find the probability that the machine needs to be replaced during the first 3 years.
3. Find the probability that the machine needs no replacement during the first 3 years.

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

x	3	4	5	6	7	8	9	10	11	12	13
$P(X = x)$	0.07	0.01	0.09	0.01	0.16	0.25	0.20	0.03	0.02	0.11	0.05

(a) Construct a row of cumulative probabilities for this table, that is, find the distribution function of X .

(b) Find the following probabilities.

$$(i) P(X \leq 5)$$

$$(iii) P(X > 9)$$

$$(v) P(4 < X \leq 9)$$

$$(ii) P(X < 12)$$

$$(iv) P(X \geq 9)$$

$$(vi) P(4 < X < 11)$$

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:

$$1. P(X \leq 1)$$

$$2. P(X \geq 1)$$

$$3. P(X > 1)$$

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 $P(X \geq 4)$?

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?

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 *emphasised sentence* 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.]

3.4 Continuous Random Variables

If X is a measurement of a continuous quantity, such as,

- the maximum diameter in millimeters of a venus shell I picked up at New Brighton beach,
- the distance you transported yourself to lectures today in meters,
- the volume of rain that fell on the roof of this building over the past 365 days in litres,
- the vertical position (in micro meters above sea-level) since the release of a pollen grain at a location in Lake Rogen in Härjedalen, as it traces through Göta älvKlarälven, the longest river of Sweden before discharging in a delta into Vänern at Karlstad.
- the volume of water (in cubic meters) that fell on the southern Alps of the South Island of New Zealand throughout last year.
- etc.,

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.

Example 52 Suppose that X is the time, in minutes, before the next student leaves the lecture room. This is an example of a continuous random variable that takes one of (uncountably) infinitely 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 2.9999999 minutes, etc., depending the measurement precision of the clock being used to measure time.

Finding $P(X = 2)$, for example, doesn't make sense because how can it ever be *exactly* 2.00000... minutes? It is more sensible to consider probabilities like $P(X > x)$ or $P(X < x)$ or $P(a < X < b)$ with $a < b$, up to measurement precision of the time-measuring clock rather than the discrete approach of trying to compute $P(X = x)$.

The characteristics of continuous random variables are:

- 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

$$F(x) = P(X \leq x) = \int_{-\infty}^x f(v) dv \quad (3.21)$$

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 $x \in \mathbb{R}$, $P(X = x) = P(X \in [x, x]) = \int_x^x f(v)dv = 0$.
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 $a < b$,

$$P(a < X < b) = P(a < X \leq b) = P(a \leq X \leq b) = P(a \leq X < b) \quad (3.23)$$

$$= F(b) - F(a) = \int_a^b f(v)dv . \quad (3.24)$$

4. And $P(\Omega) = 1$ implies that:

$$\int_{-\infty}^{\infty} f(x) dx = P(-\infty < X < \infty) = 1 .$$

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

(a) First note that if $x \leq 0$, then

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

If $0 < x < 1$, then

$$\begin{aligned} F(x) &= \int_{-\infty}^0 0 dv + \int_0^x 3v^2 dv \\ &= 0 + [v^3]_0^x \\ &= x^3 \end{aligned}$$

If $x \geq 1$, then

$$\begin{aligned} F(x) &= \int_{-\infty}^0 0 dv + \int_0^1 3v^2 dv + \int_1^x 0 dv \\ &= 0 + [v^3]_0^1 + 0 \\ &= 1 \end{aligned}$$

Hence

$$F(x) = \begin{cases} 0 & x \leq 0 \\ x^3 & 0 < x < 1 \\ 1 & x \geq 1 \end{cases}$$

(b)

$$\begin{aligned} P\left(\frac{1}{3} \leq X \leq \frac{2}{3}\right) &= F\left(\frac{2}{3}\right) - F\left(\frac{1}{3}\right) \\ &= \left(\frac{2}{3}\right)^3 - \left(\frac{1}{3}\right)^3 \\ &= \frac{7}{27} \end{aligned}$$

Example 54 Consider the continuous random variable, X , whose distribution function is:

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

(a) Find the probability density function, $f(x)$.

(b) Find $P(X > \frac{\pi}{4})$

Solution

- (a) The probability density function, $f(x)$ is given by

$$f(x) = F'(x) = \begin{cases} 0 & x < 0 \\ \cos x & 0 < x < \frac{\pi}{2} \\ 0 & x \geq \frac{\pi}{2} \end{cases}$$

- (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 6. 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.

****tossing a fair coin infinitely often, i.e., IID sequence of Bernoulli(1/2) trials, and the fundamental model**

— 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 intuiting a most primitive sigma-algebra)

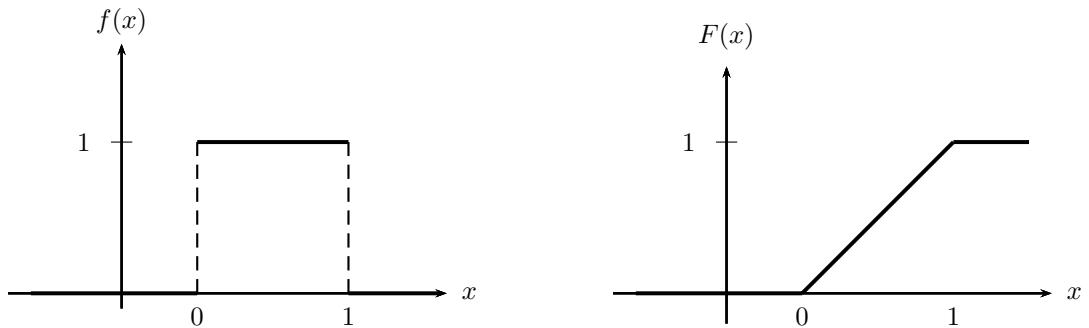
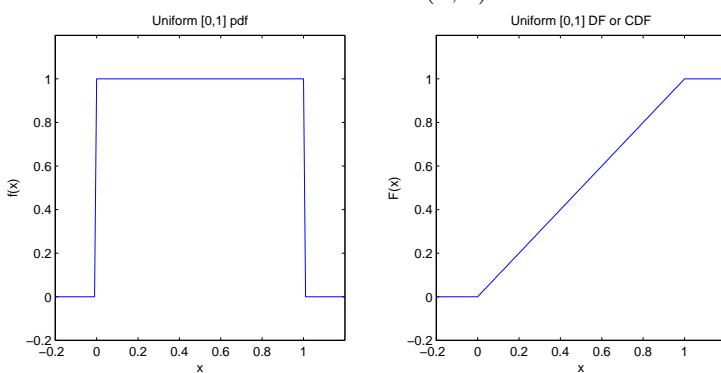


Figure 3.10: $f(x)$ and $F(x)$ of the Uniform(0, 1) random variable X .

Figure 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 .



— 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.

**universality of the fundamental model

— one can obtain any other random variable from the fundamental model whose unique DF is its own inverse, i.e., $F(x) = F^{[-1]}(x)$, as you will See from von Neumann's Fundamental Theorem of Simulation in Chapter 6.

3.4.2 Some Common Continuous Random Variables

Let us warm-up with an example.

Example 55 Let X have density function $f(x) = e^{-x}$, if $x \geq 0$, and zero otherwise.

- (a) Find the distribution function.
- (b) Find the probabilities, $P(\frac{1}{4} \leq X \leq 2)$ and $P(-\frac{1}{2} \leq X \leq \frac{1}{2})$.
- (c) Find x such that $P(X \leq x) = 0.95$.

Solution:

(a)

$$F(x) = \int_0^x e^{-v} dv = -e^{-v}]_0^x = -e^{-x} + 1 = 1 - e^{-x} \quad \text{if } x \geq 0$$

Therefore,

$$F(x) = \begin{cases} 1 - e^{-x} & \text{if } x \geq 0 \\ 0 & \text{otherwise} \end{cases} .$$

(b)

$$\begin{aligned} P\left(\frac{1}{4} \leq X \leq 2\right) &= F(2) - F\left(\frac{1}{4}\right) = 0.634 \text{ (3 sig. fig.)} \\ P\left(-\frac{1}{2} \leq X \leq \frac{1}{2}\right) &= F\left(\frac{1}{2}\right) - F\left(-\frac{1}{2}\right) = 0.394 \text{ (3 sig. fig.)} \end{aligned}$$

(c)

$$P(X \leq x) = F(x) = 1 - e^{-x} = 0.95$$

Therefore,

$$x = -\log(1 - 0.95) = 3.00 \text{ (3 sig. fig.)} .$$

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 $\bar{F}(x; \lambda) := P(X > x) = 1 - F(x; \lambda)$:

$$f(x; \lambda) = \mathbf{1}_{(0, \infty)} \lambda e^{-\lambda x} = \begin{cases} \lambda \exp(-\lambda x) & x > 0 \\ 0 & \text{otherwise} \end{cases} , \quad (3.27)$$

$$F(x; \lambda) = 1 - e^{-\lambda x} , \quad (3.28)$$

$$\bar{F}(x; \lambda) = e^{-\lambda x} . \quad (3.29)$$

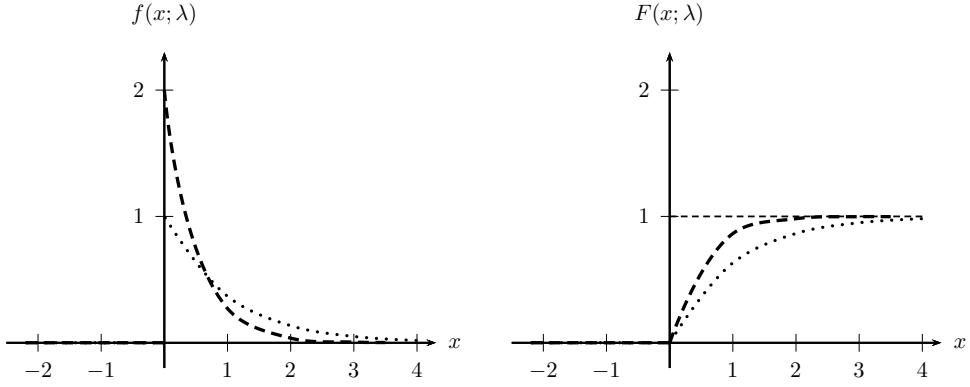


Figure 3.12: $f(x; \lambda)$ and $F(x; \lambda)$ of an exponential random variable where $\lambda = 1$ (dotted) and $\lambda = 2$ (dashed).

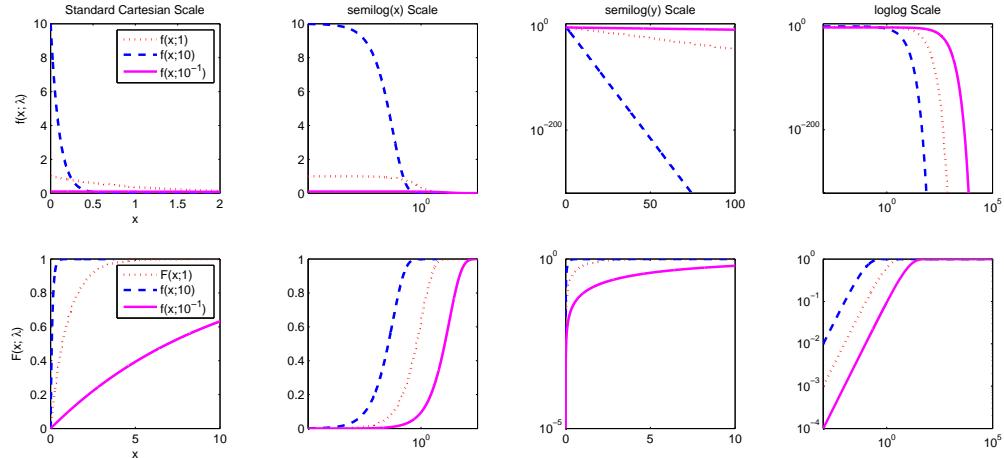
The last two equations are derived from definitions as follows:

$$\begin{aligned} F(x; \lambda) &= \int_{-\infty}^x \mathbb{1}_{(0, \infty)} \lambda e^{-\lambda v} dv = \lambda \int_0^x e^{-\lambda v} dv = \lambda \left(-\frac{1}{\lambda} e^{-\lambda v} \right]_0^x = \left(-e^{-\lambda v} \right]_0^x \\ &= -e^{-\lambda x} - (-e^{-0}) = -e^{-\lambda x} - (-1/e^0) = -e^{-\lambda x} - (-1/1) = -e^{-\lambda x} - (-1) = -e^{-\lambda x} + 1 \end{aligned}$$

$$P(X > x) = 1 - P(X \leq x) = 1 - F(x; \lambda) = 1 - \left(1 - e^{-\lambda x} \right) = e^{-\lambda x}$$

This distribution is unique because of its property of **memorylessness**, i.e., $P(X > x + y | X > y) = e^{-\lambda x}$, 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.

Figure 3.13: Density and distribution functions of $\text{Exponential}(\lambda)$ RVs, for $\lambda = 1, 10, 10^{-1}$, in four different axes scales.



Example 56 (On a dark desert highway) At a certain location on a dark desert highway, the time in minutes between arrival of cars that exceed the speed limit is an $\text{Exponential}(\lambda = 1/60)$ random variable. If you just saw a car that exceeded the speed limit then what is the probability of waiting less than 5 minutes before seeing another car that will exceed the speed limit?

Solution:

The waiting time in minutes is simply given by the $\text{Exponential}(\lambda = 1/60)$ random variable. Thus, the desired probability is

$$P(0 \leq X < 5) = \int_0^5 \frac{1}{60} e^{-\frac{1}{60}x} dx = -e^{-\frac{1}{60}x} \Big|_0^5 = -e^{-\frac{1}{12}} + 1 \approx 0.07996.$$

In exam you can stop at the expression $-e^{-\frac{1}{12}} + 1$ for full credit. You may need a calculator for the last step (with answer 0.07996).

Note: We could use the distribution function directly:

$$P(0 \leq X < 5) = F\left(5; \frac{1}{60}\right) - F\left(0; \frac{1}{60}\right) = F\left(5; \frac{1}{60}\right) = 1 - e^{-\frac{1}{60}5} = 1 - e^{-\frac{1}{12}} \approx 0.07996$$

Proposition 26 (Memorylessness of $\text{Exponential}(\lambda)$ RV) If $X \sim \text{Exponential}(\lambda)$, then X has the property of **memorylessness**, i.e.,

$$\boxed{P(X > x + y | X > y) = P(X > x)} . \quad (3.30)$$

Proof: By the definition of conditional probability,

$$P(X > x + y | X > y) = P(\{X > x + y\} | \{X > y\}) = \frac{P(\{X > x + y\} \cap \{X > y\})}{P(\{X > y\})}$$

Due to redundancy, i.e., $\{X > x + y\} \subset \{X > y\} \implies \{X > x + y\} \cap \{X > y\} = \{X > x + y\}$, so

$$\begin{aligned} P(X > x + y | X > y) &= \frac{P(\{X > x + y\})}{P(\{X > y\})} = \frac{\bar{F}(x + y; \lambda)}{\bar{F}(y; \lambda)} = \frac{e^{-\lambda(x+y)}}{e^{-\lambda y}} = \frac{e^{-\lambda x} e^{-\lambda y}}{e^{-\lambda y}} = e^{-\lambda x} = \bar{F}(x; \lambda) \\ &= P(X > x) \end{aligned}$$

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 $\text{Exponential}(\lambda)$ 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 ($\text{Uniform}(\theta_1, \theta_2)$) Given two real parameters $\theta_1, \theta_2 \in \mathbb{R}$, such that $\theta_1 < \theta_2$, the PDF of the $\text{Uniform}(\theta_1, \theta_2)$ 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 .

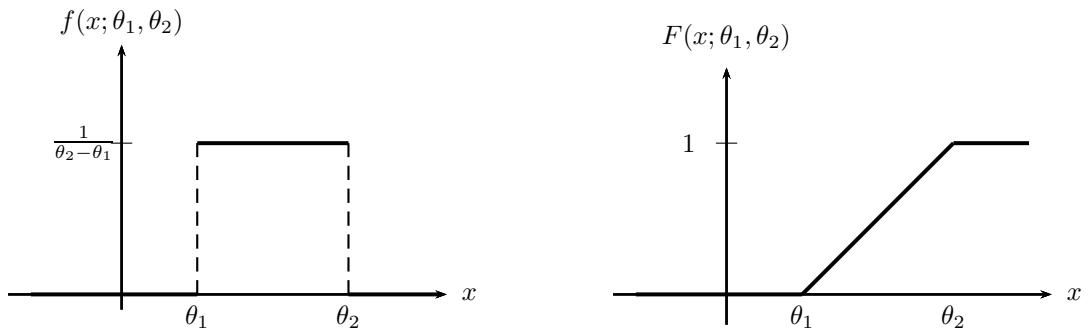


Figure 3.14: $f(x)$ and $F(x)$ of the $\text{Uniform}(\theta_1, \theta_2)$ random variable X .

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 $F(x)$.

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 - 1885). 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.

Model 10 ($\text{Normal}(0, 1)$ or **standard normal** or **Gaussian RV**) A continuous random variable Z is called **standard normal** or **standard Gaussian** if its probability density function is

$$\phi(z) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{z^2}{2}\right). \quad (3.33)$$

An exercise in calculus yields the first two derivatives of ϕ as follows:

$$\frac{d\phi}{dz} = -\frac{1}{\sqrt{2\pi}}z \exp\left(-\frac{z^2}{2}\right) = -z\phi(z), \quad \frac{d^2\phi}{dz^2} = \frac{1}{\sqrt{2\pi}}(z^2 - 1) \exp\left(-\frac{z^2}{2}\right) = (z^2 - 1)\phi(z).$$

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 centred at the origin over the real line.

Classwork 57 From the above exercise in calculus let us draw the graph of ϕ by hand now!

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!

The distribution function of Z is given by

$$\Phi(z) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^z e^{-v^2/2} dv. \quad (3.34)$$

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. ??.

We can express $\Phi(z)$ in terms of the error function (erf) as follows:

$$\Phi(z) = \frac{1}{2} \operatorname{erf}\left(\frac{z}{\sqrt{2}}\right) + \frac{1}{2} \quad (3.35)$$

And use MATLAB's `erf` function to get $\Phi(z)$ numerically instead of looking up the Table.

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.
just do it!

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
- $F(x) = \int_{-\infty}^x f(v)dv$
- $P(a < X \leq b) = F(b) - F(a) = \int_a^b f(v)dv$

3.5 Exercises in Continuous Random Variables

Ex. 3.19 — Consider the probability density function

$$f(x) = \begin{cases} k & -4 \leq x \leq 4 \\ 0 & \text{otherwise} \end{cases}.$$

1. Find the value of k .

2. Find the distribution function, F .

3. Graph f and F .

Ex. 3.20 — Assume that a new light bulb will burn out at time t hours according to the probability density function given by

$$f(t) = \begin{cases} \lambda e^{-\lambda t} & \text{if } t > 0 , \\ 0 & \text{otherwise .} \end{cases}$$

In this context, λ is often called the failure rate of the bulb.

(a) Assume that $\lambda = 0.01$, and find the probability that the bulb will not burn out before τ hours. This τ -specific probability is often called the reliability of the bulb.

Hint: Use the distribution function for an Exponential(λ) random variable (recall, $F(\tau; \lambda) = \int_{-\infty}^{\tau} f(t)dt$!)

(b) For what value of τ is the reliability of the bulb exactly $\frac{1}{2}$?

Ex. 3.21 — Let the random variable X be the time after which certain ball bearings wear out, with density

$$f(x) = \begin{cases} ke^{-x} & 0 \leq x \leq 2 \\ 0 & \text{otherwise} \end{cases} .$$

Note: X is measured in years.

1. Find k .

2. Find the probability that a bearing will last at least 1 year.

3.6 Transformations of random variables

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?

Example 60 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

$$Y = 5X - 500 .$$

Example 61 In a cell-phone system a mobile signal may have a signal-to-noise-ratio of X , but engineers prefer to express such ratios in decibels, i.e.,

$$Y = 10 \log_{10}(X) .$$

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;** $g^{-1}(y) = \frac{y}{2}$
- if $g(x) = 2x + 1$ then $g : \mathbb{R} \rightarrow \mathbb{R}$ is injective and surjective and therefore its inverse function is:
Step 1; $y = 2x + 1$, **Step 2;** $x = \frac{y-1}{2}$, **Step 3;** $g^{-1}(y) = \frac{y-1}{2}$

- if $g(x) = x^3$ then $g : \mathbb{R} \rightarrow \mathbb{R}$ is injective and surjective and therefore its inverse function is:

$$\text{Step 1; } y = x^3, \text{ Step 2; } x = y^{\frac{1}{3}}, \text{ Step 3; } g^{-1}(y) = y^{\frac{1}{3}}$$

However, you need to be careful by limiting the domain to obtain the inverse function for the following examples:

- if $g(x) = x^2$ and domain of g is $[0, +\infty)$ then its inverse function is $g^{-1}(y) = \sqrt{y}$, i.e., if $g(x) = x^2 : [0, +\infty) \rightarrow [0, +\infty)$ then the inverse image $g^{[-1]}(y)$ for $y \in [0, +\infty)$ is given by the inverse function $g^{-1}(y) = \sqrt{y} : [0, +\infty) \rightarrow [0, +\infty)$.
- if $g(x) = x^2$ and domain of g is $(-\infty, 0]$ then its inverse function is $g^{-1}(y) = -\sqrt{y}$, i.e., if $g(x) = x^2 : (-\infty, 0] \rightarrow [0, +\infty)$ then the inverse image $g^{[-1]}(y)$ for $y \in [0, +\infty)$ is given by the inverse function $g^{-1}(y) = -\sqrt{y} : [0, +\infty) \rightarrow (-\infty, 0]$.
- if $g(x) = \sin(x)$ and domain of g is $[0, \frac{\pi}{2}]$ then its inverse function $g^{-1}(y) = \arcsin(y)$, i.e., if $g(x) = \sin(x) : [0, \frac{\pi}{2}] \rightarrow [0, 1]$ then the inverse image $g^{[-1]}(y)$ for $y \in [0, 1]$ is given by the inverse function $g^{-1}(y) = \arcsin(y) : [0, 1] \rightarrow [0, \frac{\pi}{2}]$.
- if $g(x) = \sin(x)$ and domain of g is $[-\frac{\pi}{2}, \frac{\pi}{2}]$ then its inverse function $g^{-1}(y) = \arcsin(y)$, i.e., if $g(x) = \sin(x) : [-\frac{\pi}{2}, \frac{\pi}{2}] \rightarrow [-1, 1]$ then the inverse image $g^{[-1]}(y)$ for $y \in [-1, 1]$ is given by the inverse function $g^{-1}(y) = \arcsin(y) : [-1, 1] \rightarrow [-\frac{\pi}{2}, \frac{\pi}{2}]$.
- if ...

Now, let us return to our question of determining the distribution of the transformation $g(X)$. To answer this question we must first observe that the inverse image $g^{[-1]}$ satisfies the following properties:

- $g^{[-1]}(\mathbb{Y}) = \mathbb{X}$
- For any set A , $g^{[-1]}(A^c) = (g^{[-1]}(A))^c$
- For any collection of sets $\{A_1, A_2, \dots\}$,

$$g^{[-1]}(A_1 \cup A_2 \cup \dots) = g^{[-1]}(A_1) \cup g^{[-1]}(A_2) \cup \dots .$$

Consequentially,

$$P(g(X) \in A) = P\left(X \in g^{[-1]}(A)\right)$$

(3.36)

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 event given by the inverse image of A underpinned by the random variable X . It is crucial to understand this from the sample space Ω of the underlying experiment in the sense that Equation (3.36) is just short-hand for its actual meaning:

$$P(\{\omega \in \Omega : g(X(\omega)) \in A\}) = P\left(\left\{\omega \in \Omega : X(\omega) \in g^{[-1]}(A)\right\}\right) .$$

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 that of Y by $F_Y(y)$.

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) &= P(Y = y) = 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:

$f_Y(y) = P(Y = y) = \sum_{x \in g^{[-1]}(y)} f_X(x) = \sum_{x \in \{x: g(x)=y\}} f_X(x) .$		(3.37)
---	--	--------

Example 62 Let X be the discrete random variable with probability mass function f_X as tabulated below:

x	-1	0	1
$f_X(x) = P(X = x)$	$\frac{1}{4}$	$\frac{1}{2}$	$\frac{1}{4}$

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:

$$f_Y(y) = P(Y = y) = \sum_{x \in g^{[-1]}(y)} f_X(x) = \sum_{x \in \{y/2\}} f_X(x) = f_X(y/2) ,$$

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,

$$f_Y(y) = P(Y = y) = \sum_{x \in g^{[-1]}(y)} f_X(x) = \sum_{x \in \{(y-1)/2\}} f_X(x) = f_X((y-1)/2) ,$$

and tabulated below:

y	-1	1	3
$f_Y(y)$	$\frac{1}{4}$	$\frac{1}{2}$	$\frac{1}{4}$

In fact, obtaining the probability of a one-to-one transformation of a discrete random variable as in Examples 62 and 63 is merely a matter of looking up the probability at the image of the inverse function. This is because there is only one term in the sum that appears in Equation (3.37). When the transformation is not one-to-one the number of terms in the sum can be more than one as shown in the next Example.

Example 64 Reconsider the random variable X of the last two Examples and let $Y = X^2$. Recall that $g(x) = x^2$ does not have an inverse function unless the domain is restricted to the positive or the negative parts of the real line. Since our random variable X takes values on both sides of the real line, namely $\{-1, 0, 1\}$, let us note that the transformation $g(X) = X^2$ is no longer a one-to-one function. Then, by Equation (3.37) the probability mass function of Y is expressed in terms of the known probabilities of X as:

$$f_Y(y) = P(Y = y) = \sum_{x \in g^{-1}(y)} f_X(x) = \sum_{\{x: g(x)=y\}} f_X(x) = \sum_{\{x: x^2=y\}} f_X(x) ,$$

computed for each $y \in \{0, 1\}$ as follows:

$$\begin{aligned} f_Y(0) &= \sum_{\{x: x^2=0\}} f_X(x) = f_X(0) = \frac{1}{2} , \\ f_Y(1) &= \sum_{\{x: x^2=1\}} f_X(x) = f_X(-1) + f_X(1) = \frac{1}{4} + \frac{1}{4} = \frac{1}{2} , \end{aligned}$$

and finally tabulated below:

y	0	1
$f_Y(y)$	$\frac{1}{2}$	$\frac{1}{2}$

3.6.3 Transformations of continuous random variables

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 .

One-to-one transformations

The easiest case for transformations of continuous random variables is when g is **one-to-one and monotone**.

- 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

$$F_Y(y) = P(Y \leq y) = P(g(X) \leq y) = P(X \leq g^{-1}(y)) = F_X(g^{-1}(y)) .$$

Now, let us use a form of chainrule to compute the density of Y as follows:

$$f_Y(y) = \frac{d}{dy} F_Y(y) = \frac{d}{dy} F_X(g^{-1}(y)) = f_X(g^{-1}(y)) \frac{d}{dy} (g^{-1}(y)) .$$

- 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 .

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

(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 $0 \leq y \leq 1$

Thus, we have shown that if X is a Uniform(0, 1) random variable then $Y = 1 - X$ is also a Uniform(0, 1) random variable.

Example 66 Let X be a Uniform(0, 1) random variable and let $Y = g(X) = -\log(X)$. We are interested in the density of the tranformed random variable Y . Once again, since g is a one-to-one monotone function let us follow the four steps and use the change of variable formula to obtain f_Y from f_X and g .

1. $y = g(x) = -\log(x)$ is a monotone decreasing function over $0 < x < 1$, the range of X .
So, we can apply the change of variable formula.
2. $x = g^{-1}(y) = \exp(-y)$ is a monotone decreasing function over $0 < y < \infty$.
3. For $0 < y < \infty$,

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

4. We can use Equation (3.38) to find the density of Y as follows:

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

Note that $0 < \exp(-y) < 1$ for $0 < y < \infty$.

Thus, we have shown that if X is a Uniform(0, 1) random variable then $Y = -\log(X)$ is an random variable with PDF $f_Y(y) = \mathbb{1}_{(0,\infty)}(y) \exp(-y)$. We can similarly show that for a parameter $\lambda > 0$, if $X \sim \text{Uniform}(0, 1)$ then $Y = -\lambda^{-1} \log(X)$ yields a probability model of RVs that are parameterized by λ and extremely useful in applications. This is noting but our Exponential(λ) RV.

The next example yields the *location-scale* family of normal random variables via a family of linear transformations of the standard normal random variable.

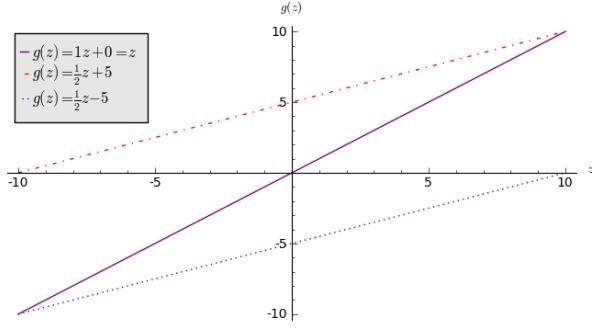
Example 67 Let Z be the standard Gaussian or standard normal random variable with probability density function $\phi(z)$ given by Equation (3.33). For real numbers $\sigma > 0$ and μ consider the linear transformation of Z given by

$$Y = g(Z) = \sigma Z + \mu .$$

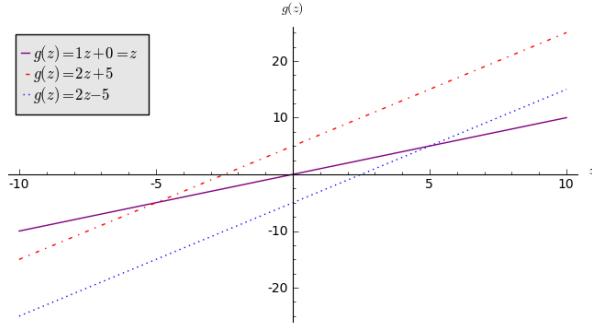
Some graphs of such linear transformations of Z are shown in Figures (a) and (b).

We are interested in the density of the tranformed random variable $Y = g(Z) = \sigma Z + \mu$. Once again, since g is a one-to-one monotone function let us follow the four steps and use the change of variable formula to obtain f_Y from $f_Z = \phi$ and g .

1. $y = g(z) = \sigma z + \mu$ is a monotone increasing function over $-\infty < z < \infty$, the range of Z .
So, we can apply the change of variable formula.
2. $z = g^{-1}(y) = (y - \mu)/\sigma$ is a monotone increasing function over the range of y given by,
 $-\infty < y < \infty$.



(a) $g(z) = z$, $g(z) = \frac{1}{2}z + 5$ and $g(z) = \frac{1}{2}z - 5$.



(b) $g(z) = z$, $g(z) = (z + 5)/0.5$ and $g(z) = (z - 5)/0.5$.

3. For $-\infty < y < \infty$,

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

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 .

The normal distribution has the **distribution function**

$$F(x; \mu, \sigma^2) = \frac{1}{\sigma\sqrt{2\pi}} \int_{-\infty}^x \exp\left[-\frac{1}{2}\left(\frac{v-\mu}{\sigma}\right)^2\right] dv . \quad (3.40)$$

Here we need x as the upper limit of integration and so we write v in the integrand.

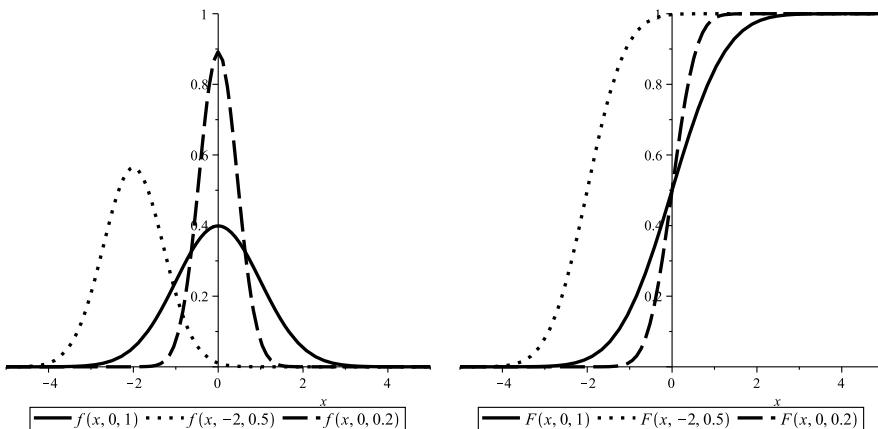


Figure 3.15: PDF and DF of a $\text{Normal}(\mu, \sigma^2)$ RV for different values of μ and σ^2

Using the direct method's Equation 3.41, we can obtain the distribution function of the $\text{Normal}(\mu, \sigma^2)$ random variable from that of the tabulated distribution function of the $\text{Normal}(0, 1)$ in the Standard normal distribution function table in Sec. ??.

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:

$$F_X(x; \mu, \sigma^2) = F_Z\left(\frac{x-\mu}{\sigma}\right) = \Phi\left(\frac{x-\mu}{\sigma}\right) .$$

Proof: Let Z be a $\text{Normal}(0, 1)$ random variable with distribution function $\Phi(z) = P(Z \leq z)$. We know that if $X = g(Z) = \sigma Z + \mu$ then X is the $\text{Normal}(\mu, \sigma^2)$ random variable. Therefore,

$$\begin{aligned} F_X(x; \mu, \sigma^2) &= P(X \leq x) = P(g(Z) \leq x) = P(\sigma Z + \mu \leq x) = P\left(Z \leq \frac{x - \mu}{\sigma}\right) \\ &= F_Z\left(\frac{x - \mu}{\sigma}\right) = \Phi\left(\frac{x - \mu}{\sigma}\right). \end{aligned}$$

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:

$$Z = \frac{X - \mu}{\sigma}.$$

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:

$$\begin{aligned} F_Y(y) &= P(Y \leq y) = P(Y \in (-\infty, y]) \\ &= P(g(X) \in (-\infty, y]) = P\left(X \in g^{[-1]}((-\infty, y])\right) = P(X \in \{x : g(x) \in (-\infty, y]\}) \end{aligned}$$

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:

- If $y < 0$ then $f_Y(y) = \frac{d}{dy}(F_Y(y)) = \frac{d}{dy}0 = 0$.
- If $y \geq 0$ then

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

Therefore, the distribution function of $Y = X^2$ is:

$$F_Y(y) = \begin{cases} 0 & \text{if } y < 0 \\ F_X(\sqrt{y}) - F_X(-\sqrt{y}) & \text{if } y \geq 0 \end{cases} . \quad (3.42)$$

and the probability density function of $Y = X^2$ is:

$$f_Y(y) = \begin{cases} 0 & \text{if } y < 0 \\ \frac{1}{2\sqrt{y}}(f_X(\sqrt{y}) + f_X(-\sqrt{y})) & \text{if } y \geq 0 \end{cases} . \quad (3.43)$$

Example 70 If X is the standard normal random variable with density

$$f_X(x) = \phi(x) = \frac{1}{\sqrt{2\pi}} \exp(-x^2/2)$$

then by Equation (3.43) the density of $Y = X^2$ is:

$$f_Y(y) = \begin{cases} 0 & \text{if } y < 0 \\ \frac{1}{2\sqrt{y}}(f_X(\sqrt{y}) + f_X(-\sqrt{y})) = \frac{1}{\sqrt{2\pi y}} \exp\left(-\frac{y}{2}\right) & \text{if } y \geq 0 \end{cases} .$$

Y 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 Inference Theory and was derived at the beginning of last century to settle “supposedly evidence-based disputes” among scientists using mathematics.

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.

Definition 29 (Expectation of a RV) The **expectation**, or **expected value**, or **mean**, or **first moment**, of a random variable X , with distribution function F and density f , is defined to be

$$E(X) := \int x dF(x) = \begin{cases} \sum_x x f(x) & \text{if } X \text{ is discrete} \\ \int x f(x) dx & \text{if } X \text{ is continuous,} \end{cases} \quad (3.44)$$

provided the sum or integral is well-defined. We say the expectation exists if

$$\int |x| dF(x) < \infty. \quad (3.45)$$

Sometimes, we denote $E(X)$ by $E X$ for brevity. Thus, the expectation is a single-number summary of the RV X and may be thought of as the average. We subscript E to specify the parameter $\theta \in \Theta$ with respect to which the integration is undertaken.

$$E_\theta(X) := \int x dF(x; \theta)$$

Definition 30 (Variance of a RV) Let X be a RV with mean or expectation $E(X)$. Variance of X denoted by $V(X)$ or simply $V X$ is

$$V(X) := E((X - E(X))^2) = \int (x - E(X))^2 dF(x),$$

provided this expectation exists. The **standard deviation** denoted by $sd(X) := \sqrt{V(X)}$. Thus variance is a measure of “spread” of a distribution.

Definition 31 (k -th moment of a RV) We call

$$E(X^k) = \int x^k dF(x)$$

as the k -th moment of the RV X and say that the k -th moment exists when $E(|X|^k) < \infty$. We call the following expectation as the k -th central moment:

$$E((X - E(X))^k).$$

3.8.1 Expectations of functions of random variables

More generally, by taking the expected value of various functions of a random variable, we can measure many interesting features of its distribution, including spread and correlation.

Definition 32 (Expectation of a function of a RV) The **Expectation** of a function $g(X)$ of a random variable X is defined as:

$$E(g(X)) := \int g(x) dF(x) = \begin{cases} \sum_x g(x) f(x) & \text{if } X \text{ is a discrete RV} \\ \int_{-\infty}^{\infty} g(x) f(x) dx & \text{if } X \text{ is a continuous RV} \end{cases}$$

provided $E(g(X))$ exists, i.e., $\int |g(x)| dF(x) < \infty$.

The **mean** which characterises the central location of the random variable X is merely the expectation of the identity function $g(x) = x$:

$$E(X) = \begin{cases} \sum_x xf(x) & \text{if } X \text{ is a discrete RV} \\ \int_{-\infty}^{\infty} xf(x)dx & \text{if } X \text{ is a continuous RV} \end{cases}$$

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 $g(x) = (x - E(X))^2$:

$$V(X) = E((X - E(X))^2) = \begin{cases} \sum_x (x - E(X))^2 f(x) & \text{if } X \text{ is a discrete RV} \\ \int_{-\infty}^{\infty} (x - E(X))^2 f(x) dx & \text{if } X \text{ is a continuous RV} \end{cases}$$

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

$$E(X) = -10 \times \frac{1}{2} + 10 \times \frac{1}{2} = 0.$$

So what you expect (nothing) and the weighted average ($E(X) = 0$) 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 $E(X)$ 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

$$E(Y) = E(r(X)) = \int r(x) dF(x).$$

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 $r(X) = \mathbf{1}_A(x)$. Recall $\mathbf{1}_A(x)$ is 1 if $x \in A$ and $\mathbf{1}_A(x) = 0$ if $x \notin A$. Then

$$E(\mathbf{1}_A(X)) = \int \mathbf{1}_A(x) dF(x) = \int_A dF(x) = P(X \in A) = P(A) \quad (3.46)$$

Thus, probability is a special case of expectation. Recall our LTRF motivation for the definition of probability and make the connection.

Expectations of functions of \mathbb{R}^2 -valued random variables

In the case of a single random variable we saw that its expectation gives the population mean, a measure of the center of the distribution of the variable in some sense. Similarly, by taking the expected value of various functions of a \mathbb{R}^2 -valued random variable, we can measure many interesting features of its joint distribution.

Definition 33 The **Expectation** of a function $g(X, Y)$ of the \mathbb{R}^2 -valued RV (X, Y) is defined as:

$$E(g(X, Y)) = \begin{cases} \sum_{(x,y)} g(x, y) f_{X,Y}(x, y) & \text{if } (X, Y) \text{ is a discrete RV} \\ \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} g(x, y) f_{X,Y}(x, y) dx dy & \text{if } (X, Y) \text{ is a continuous RV} \end{cases}$$

Some typical expectations for \mathbb{R}^2 -valued random variables are:

1. Joint Moments

$$E(X^r Y^s)$$

When $r = s = 1$, we have $E(XY)$, the expectation of the product of two RVs.

2. We need a new notion for the variance of two RVs.

If $E(X^2) < \infty$ and $E(Y^2) < \infty$ then $E(|XY|) < \infty$ and $E(|(X - E(X))(Y - E(Y))|) < \infty$. This allows the definition of **covariance** of X and Y as

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

The same ideas naturally extend, via multiple sums and integrals, to define the expectation of functions of \mathbb{R}^k -valued random variables with $k > 2$.

Viewing a deterministic real variable as a random variable

Consider the class of discrete RVs with distributions that place all probability mass on a single real number. This is the probability model for the deterministic real variable, which is often thought of as an unknown constant $\theta \in \mathbb{R}$.

Model 12 (Point Mass(θ)) Given a specific point $\theta \in \mathbb{R}$, we say an RV X has point mass at θ or is Point Mass(θ) distributed if the DF is:

$$F(x; \theta) = \begin{cases} 0 & \text{if } x < \theta \\ 1 & \text{if } x \geq \theta \end{cases} \quad (3.47)$$

and the PMF is:

$$f(x; \theta) = \begin{cases} 0 & \text{if } x \neq \theta \\ 1 & \text{if } x = \theta \end{cases} \quad (3.48)$$

Thus, Point Mass(θ) RV X is deterministic in the sense that every realisation of X is exactly equal to $\theta \in \mathbb{R}$. We will see that this distribution plays a central limiting role in asymptotic statistics.

Example 73 (Mean and variance of Point Mass(θ) RV) Let $X \sim \text{Point Mass}(\theta)$. Then:

$$E(X) = \sum_x x f(x) = \theta \times 1 = \theta , \quad V(X) = E(X^2) - (E(X))^2 = \theta^2 - \theta^2 = 0 .$$

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:

$$\boxed{\mathbb{E}(a) = a}$$

$$\boxed{\mathbb{E}(a g(X)) = a \mathbb{E}(g(X))}$$

$$\boxed{\mathbb{E}(g(X) + h(X)) = \mathbb{E}(g(X)) + \mathbb{E}(h(X))}$$

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:

$$\begin{aligned} \mathbb{V}(X) &= \mathbb{E}((X - \mathbb{E}(X))^2) \\ &= \mathbb{E}(X^2 - 2X\mathbb{E}(X) + (\mathbb{E}(X))^2) \\ &= \mathbb{E}(X^2) - \mathbb{E}(2X\mathbb{E}(X)) + \mathbb{E}((\mathbb{E}(X))^2) \\ &= \mathbb{E}(X^2) - 2\mathbb{E}(X)\mathbb{E}(X) + (\mathbb{E}(X))^2 \\ &= \mathbb{E}(X^2) - 2(\mathbb{E}(X))^2 + (\mathbb{E}(X))^2 \\ &= \mathbb{E}(X^2) - (\mathbb{E}(X))^2 . \end{aligned}$$

That is,

$$\boxed{\mathbb{V}(X) = \mathbb{E}(X^2) - (\mathbb{E}(X))^2}.$$

The above properties of expectations imply that for constants a and b ,

$$\boxed{\mathbb{V}(aX + b) = a^2 \mathbb{V}(X)} . \quad (3.49)$$

More generally, for random variables X_1, X_2, \dots, X_n and constants a_1, a_2, \dots, a_n

- $\mathbb{E}\left(\sum_{i=1}^n a_i X_i\right) = \sum_{i=1}^n a_i \mathbb{E}(X_i) . \quad (3.50)$

- $\mathbb{V}\left(\sum_{i=1}^n a_i X_i\right) = \sum_{i=1}^n a_i^2 \mathbb{V}(X_i)$, provided X_1, X_2, \dots, X_n are independent . $\quad (3.51)$

- Let X_1, X_2, \dots, X_n be independent RVs, then

$$\mathbb{E}\left(\prod_{i=1}^n X_i\right) = \prod_{i=1}^n \mathbb{E}(X_i), \text{ provided } X_1, X_2, \dots, X_n \text{ are independent} . \quad (3.52)$$

3.8.3 Expectation of Common Random Variables

Let us compute the mean and variance of our familiar RVs.

Example 74 (Mean and variance of Bernoulli(θ) RV) Let $X \sim \text{Bernoulli}(\theta)$. Then,

$$E(X) = \sum_{x=0}^1 xf(x) = (0 \times (1 - \theta)) + (1 \times \theta) = 0 + \theta = \theta ,$$

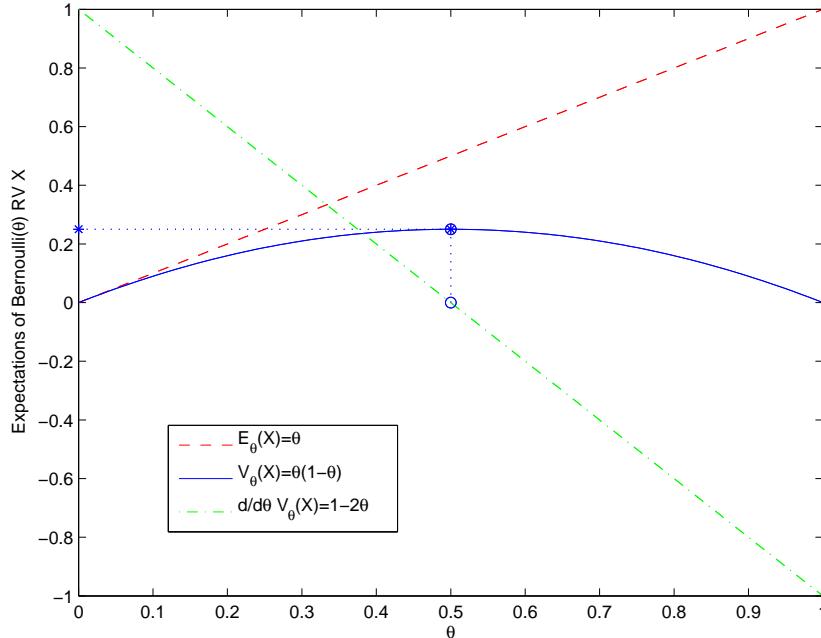
$$E(X^2) = \sum_{x=0}^1 x^2 f(x) = (0^2 \times (1 - \theta)) + (1^2 \times \theta) = 0 + \theta = \theta ,$$

$$V(X) = E(X^2) - (E(X))^2 = \theta - \theta^2 = \theta(1 - \theta) .$$

Parameter specifically,

$$E_\theta(X) = \theta \quad \text{and} \quad V_\theta(X) = \theta(1 - \theta) .$$

Figure 3.16: Mean ($E_\theta(X)$), variance ($V_\theta(X)$) and the rate of change of variance ($\frac{d}{d\theta} V_\theta(X)$) of a Bernoulli(θ) RV X as a function of the parameter θ .



Maximum of the variance $V_\theta(X)$ is found by setting the derivative to zero, solving for θ and showing the second derivative is locally negative, i.e. $V_\theta(X)$ is concave down:

$$V'_\theta(X) := \frac{d}{d\theta} V_\theta(X) = 1 - 2\theta = 0 \iff \theta = \frac{1}{2} , \quad V''_\theta(X) := \frac{d}{d\theta} \left(\frac{d}{d\theta} V_\theta(X) \right) = -2 < 0 ,$$

$$\max_{\theta \in [0,1]} V_\theta(X) = \frac{1}{2} \left(1 - \frac{1}{2} \right) = \frac{1}{4} , \text{ since } V_\theta(X) \text{ is maximized at } \theta = \frac{1}{2}$$

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 $V_\theta(X)$ attains its maximum value of $1/4$ at $\theta = 0.5$

where $\frac{d}{d\theta} V_\theta(X) = 0$. Furthermore, we know that we don't have a minimum at $\theta = 0.5$ since the second derivative $V''_\theta(X) = -2$ is negative for any $\theta \in [0, 1]$. This confirms that $V_\theta(X)$ is concave down and therefore we have a maximum of $V_\theta(X)$ 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,

$$\begin{aligned} E(X) &= \int_{x=0}^1 x f(x) dx = \int_{x=0}^1 x \cdot 1 dx = \frac{1}{2} (x^2) \Big|_{x=0}^{x=1} = \frac{1}{2} (1 - 0) = \frac{1}{2}, \\ E(X^2) &= \int_{x=0}^1 x^2 f(x) dx = \int_{x=0}^1 x^2 \cdot 1 dx = \frac{1}{3} (x^3) \Big|_{x=0}^{x=1} = \frac{1}{3} (1 - 0) = \frac{1}{3}, \\ V(X) &= E(X^2) - (E(X))^2 = \frac{1}{3} - \left(\frac{1}{2}\right)^2 = \frac{1}{3} - \frac{1}{4} = \frac{1}{12}. \end{aligned}$$

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 $E(X)$ and $V(X)$ 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 $E(Y)$.

We can simply apply the definition of $E(r(X))$, since $Y = r(X)$, is just a function of X , as follows:

$$E(Y) = \int_0^1 e^x f(x) dx = \int_0^1 e^x \cdot 1 dx = e - 1.$$

Example 77 (Mean and variance of Exponential(λ)) Show that the mean of an Exponential(λ) RV X is:

$$E_\lambda(X) = \int_0^\infty x f(x; \lambda) dx = \int_0^\infty x \lambda e^{-\lambda x} dx = \frac{1}{\lambda},$$

and the variance is:

$$V_\lambda(X) = \left(\frac{1}{\lambda}\right)^2.$$

Example 78 (Mean and variance of Geometric(θ) RV) Let $X \sim \text{Geometric}(\theta)$ RV. Then,

$$E(X) = \sum_{x=0}^{\infty} x \theta (1 - \theta)^x = \theta \sum_{x=0}^{\infty} x (1 - \theta)^x$$

In order to simplify the RHS above, let us employ differentiation with respect to θ :

$$\frac{-1}{\theta^2} = \frac{d}{d\theta} \left(\frac{1}{\theta}\right) = \frac{d}{d\theta} \sum_{x=0}^{\infty} (1 - \theta)^x = \sum_{x=0}^{\infty} -x (1 - \theta)^{x-1}$$

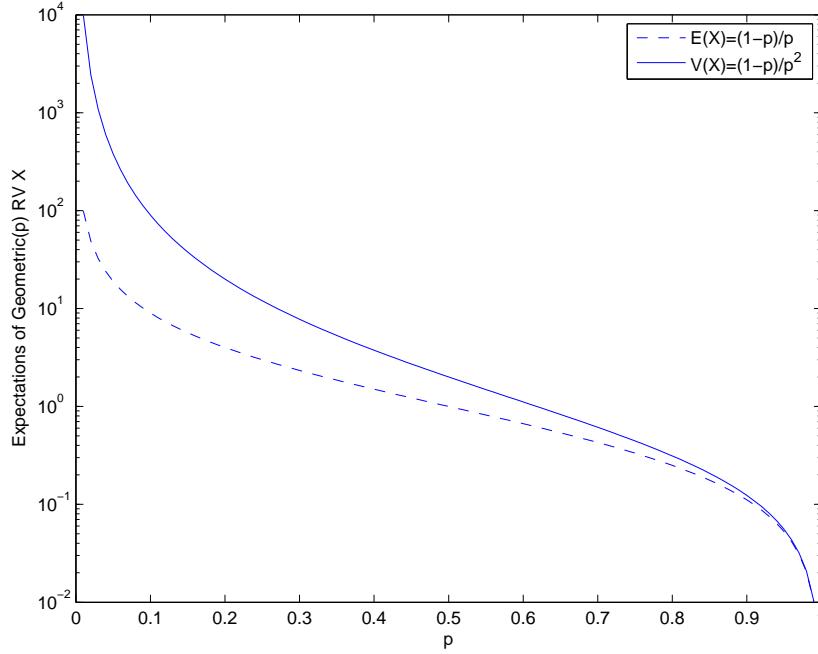
Multiplying the LHS and RHS above by $-(1 - \theta)$ and substituting in $E(X) = \theta \sum_{x=0}^{\infty} x (1 - \theta)^x$, we get a much simpler expression for $E(X)$:

$$\frac{1 - \theta}{\theta^2} = \sum_{x=0}^{\infty} x (1 - \theta)^x \implies E(X) = \theta \left(\frac{1 - \theta}{\theta^2}\right) = \frac{1 - \theta}{\theta}.$$

Similarly, it can be shown that

$$V(X) = \frac{1-\theta}{\theta^2} .$$

Figure 3.17: Mean and variance of a Geometric(θ) RV X as a function of the parameter θ .



Example 79 (Mean and variance of Binomial(n, θ) RV) Let $X \sim \text{Binomial}(n, \theta)$. Based on the definition of expectation:

$$E(X) = \int x dF(x; n, \theta) = \sum_x x f(x; n, \theta) = \sum_{x=0}^n x \binom{n}{x} \theta^x (1-\theta)^{n-x} .$$

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 $\{X_1, X_2, \dots, X_n\} \stackrel{\text{IID}}{\sim} \text{Bernoulli}(\theta)$, $E(X_i) = \theta$ and $V(X_i) = \theta(1-\theta)$:

$$\begin{aligned} E(X) &= E(X_1 + X_2 + \dots + X_n) = E\left(\sum_{i=1}^n X_i\right) = \sum_{i=1}^n E(X_i) = n\theta , \\ V(X) &= V\left(\sum_{i=1}^n X_i\right) = \sum_{i=1}^n V(X_i) = \sum_{i=1}^n \theta(1-\theta) = n\theta(1-\theta) . \end{aligned}$$

Example 80 (Mean and variance of Poisson(λ) RV) Let $X \sim \text{Poisson}(\lambda)$. Then:

$$E(X) = \sum_{x=0}^{\infty} x f(x; \lambda) = \sum_{x=0}^{\infty} x \frac{e^{-\lambda} \lambda^x}{x!} = e^{-\lambda} \sum_{x=0}^{\infty} x \frac{\lambda^x}{x!} = e^{-\lambda} \sum_{x-1=0}^{\infty} \frac{\lambda \lambda^{x-1}}{(x-1)!} = e^{-\lambda} \lambda e^{\lambda} = \lambda .$$

Similarly,

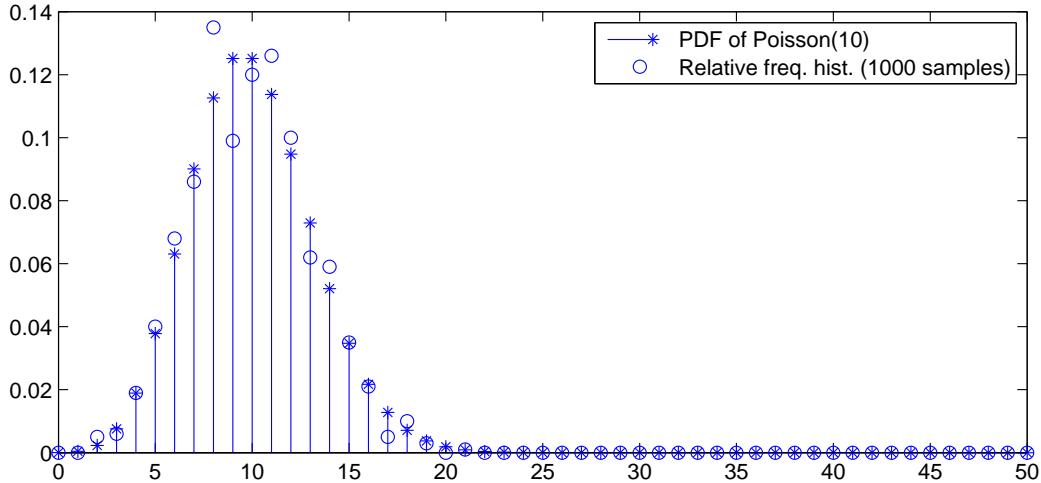
$$V(X) = E(X^2) - (E(X))^2 = \lambda + \lambda^2 - \lambda^2 = \lambda .$$

since

$$\begin{aligned}
E(X^2) &= \sum_{x=0}^{\infty} x^2 \frac{e^{-\lambda} \lambda^x}{x!} = \lambda e^{-\lambda} \sum_{x=1}^{\infty} \frac{x \lambda^{x-1}}{(x-1)!} = \lambda e^{-\lambda} \left(1 + \frac{2\lambda}{1} + \frac{3\lambda^2}{2!} + \frac{4\lambda^3}{3!} + \dots \right) \\
&= \lambda e^{-\lambda} \left(\left(1 + \frac{\lambda}{1} + \frac{\lambda^2}{2!} + \frac{\lambda^3}{3!} + \dots \right) + \left[\frac{\lambda}{1} + \frac{2\lambda^2}{2!} + \frac{3\lambda^3}{3!} + \dots \right] \right) \\
&= \lambda e^{-\lambda} \left((e^\lambda) + \lambda \left(1 + \frac{2\lambda}{2!} + \frac{3\lambda^2}{3!} + \dots \right) \right) = \lambda e^{-\lambda} \left(e^\lambda + \lambda \left(1 + \lambda + \frac{\lambda^2}{2!} + \dots \right) \right) \\
&= \lambda e^{-\lambda} (e^\lambda + \lambda e^\lambda) = \lambda e^{-\lambda} (e^\lambda + \lambda e^\lambda) = \lambda(1 + \lambda) = \lambda + \lambda^2
\end{aligned}$$

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 175.



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 $E(X) = \mu$ and $V(X) = \sigma^2$ follows directly from the properties of Expectations, provided $E(Z) = 0$ and $V(Z) = E(Z^2) - (E(Z))^2 = E(Z^2) = 1$.

The mean of a $\text{Normal}(0, 1)$ RV Z is:

$$E(Z) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} z \exp\left(-\frac{1}{2}z^2\right) dz = \frac{1}{\sqrt{2\pi}} \left[-\exp\left(-\frac{1}{2}z^2\right) \right]_{-\infty}^{\infty} = 0,$$

and the variance is:

$$V(Z) = E(Z^2) - (E(Z))^2 = E(Z^2) - 0 = E(Z^2) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} z^2 e^{-z^2/2} dz.$$

Using integration by parts with $u = z, dv = ze^{-z^2/2} \implies du = 1, v = -e^{-z^2/2}$, $\int uv \, dv = uv - \int v \, du$

$$\frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} z^2 e^{-z^2/2} dz = \frac{1}{\sqrt{2\pi}} \left(-ze^{-z^2/2} \right)_{-\infty}^{\infty} + \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{-z^2/2} dz = 0 + 1 = 1$$

The first term after the first equality above equals 0 because the exponential goes to 0 much faster than z grows to $\pm\infty$. The second term equals 1 because it is exactly the total probability integral of the PDF of the $\text{Normal}(0, 1)$ RV.

Next, let us become familiar with an RV for which the expectation does not exist.

Model 13 (Cauchy) The density of the Cauchy RV Y is:

$$f(y) = \frac{1}{\pi(1+y^2)}, \quad -\infty < y < \infty , \quad (3.53)$$

and its DF is:

$$F(y) = \frac{1}{\pi} \tan^{-1}(y) + \frac{1}{2} . \quad (3.54)$$

Randomly spinning a LASER emitting improvisation of “Darth Maul’s double edged lightsaber” that is centered 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.

The Cauchy RV Y can be derived from a RV $X \sim \text{Uniform}(-\pi/2, \pi/2)$ by the simple transformation $Y = \tan(X)$ for the above construction. Since $\tan(x)$ is one-to-one and monotone on the range of X given by $(-\pi/2, \pi/2)$, we can use the change of variable formula in Equation 3.38 to obtain the PDF $f_Y(y)$ from the PDF $f_X(x) = \frac{1}{\pi} \mathbf{1}_{(-\pi/2, \pi/2)}(x)$ as follows:

$$f_Y(y) = f_X(g^{-1}(y)) \left| \frac{d}{dy} g^{-1}(y) \right| = f_X(\tan^{-1}(y)) \left| \frac{d}{dy} \tan^{-1}(y) \right| = \frac{1}{\pi} \left| \frac{1}{1+y^2} \right|$$

Note that the construction is valid even if we sample X uniformly from $(0, \pi)$ and take its $\tan(X)$.

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:

$$\int |x| \, dF(x) = \frac{2}{\pi} \int_0^\infty \frac{x}{1+x^2} dx = (x \tan^{-1}(x)]_0^\infty - \int_0^\infty \tan^{-1}(x) dx = \infty . \quad (3.55)$$

Note that we consider symmetry of integral about the origin and take twice the integral over $(0, \infty)$ above. Variance and higher moments cannot be defined when the expectation itself is undefined.

Next let us consider a natural generalization of the $\text{Bernoulli}(\theta)$ RV with more than two outcomes but in the set $\{1, 2, \dots, k\}$.

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:

$$\Delta^{k-1} := \{ (\theta_1, \theta_2, \dots, \theta_k) : \theta_1 \geq 0, \theta_2 \geq 0, \dots, \theta_k \geq 0, \sum_{i=1}^k \theta_i = 1 \} ,$$

we say that an RV X is de Moivre($\theta_1, \theta_2, \dots, \theta_k$) distributed if its PMF is:

$$f(x; \theta_1, \theta_2, \dots, \theta_k) = \begin{cases} 0 & \text{if } x \notin [k] := \{1, 2, \dots, k\}, \\ \theta_x & \text{if } x \in [k]. \end{cases}$$

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,

$$\mathbb{E}(X) = \sum_{x=1}^k x\theta(x) = \theta_1 + 2\theta_2 + \dots + k\theta_k, \text{ and}$$

$$\mathbb{V}(X) = \mathbb{E}(X^2) - (\mathbb{E}(X))^2 = (\theta_1 + 2^2\theta_2 + \dots + k^2\theta_k) - (\theta_1 + 2\theta_2 + \dots + k\theta_k)^2.$$

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,

$$\mathbb{E}(X) = \frac{1}{k} (1 + 2 + \dots + k) = \frac{1}{k} \frac{k(k+1)}{2} = \frac{k+1}{2},$$

$$\mathbb{E}(X^2) = \frac{1}{k} (1^2 + 2^2 + \dots + k^2) = \frac{1}{k} \frac{k(k+1)(2k+1)}{6} = \frac{2k^2 + 3k + 1}{6},$$

$$\begin{aligned} \mathbb{V}(X) &= \mathbb{E}(X^2) - (\mathbb{E}(X))^2 = \frac{2k^2 + 3k + 1}{6} - \left(\frac{k+1}{2}\right)^2 = \frac{2k^2 + 3k + 1}{6} - \left(\frac{k^2 + 2k + 1}{4}\right) \\ &= \frac{8k^2 + 12k + 4 - 6k^2 - 12k - 6}{24} = \frac{2k^2 - 2}{24} = \frac{k^2 - 1}{12}. \end{aligned}$$

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 .

Ex. 3.32 — Starting from the definition of the variance of a random variable (Definition 30) show that

$$\text{V}(X) = \text{E}(X^2) - (\text{E}(X))^2 .$$

Ex. 3.33 — Show that $\text{V}(aX + b) = a^2\text{V}(X)$ for constants a and b and a random variable X .

Ex. 3.34 — **Let X be a discrete random variable with PMF given by

$$f(x) = \begin{cases} \frac{x}{10} & \text{if } x \in \{1, 2, 3, 4\}, \\ 0 & \text{otherwise.} \end{cases}$$

(a)Find:

- (i) $\text{P}(X = 0)$
- (ii) $\text{P}(2.5 < X < 5)$
- (iii) $\text{E}(X)$
- (iv) $\text{V}(X)$

(b)Write down the DF (or CDF) of X .

(c)Plot the PMF and CDF of X .

Ex. 3.35 — Find the mean and the variance of the following random variables.

1. X a discrete uniform random variable on $\{1, 2, 3, 4, 5, 6\}$, i.e., *the number a fair die turns up.*
2. X is a Uniform($0, 8$) random variable, i.e., *a continuous uniform random variable from the interval $[0, 8]$.*
3. X has a density function

$$f(x) = \begin{cases} 2e^{-2x} & \text{if } x \geq 0 \\ 0 & \text{otherwise.} \end{cases}$$

3.10 Multivariate Random Variables

Often, in 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 $\omega \in \Omega$ 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 $\mathbb{R}^3 = (-\infty, \infty) \times (-\infty, \infty) \times (-\infty, \infty)$:

$$\omega \mapsto (X(\omega), Y(\omega)) : \Omega \rightarrow \mathbb{R}^2 \quad \omega \mapsto (X(\omega), Y(\omega), Z(\omega)) : \Omega \rightarrow \mathbb{R}^3$$

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 thus need to make m measurements of the outcome in \mathbb{R}^m using a “measurable mapping” from $\Omega \rightarrow \mathbb{R}^m$. To deal with such multivariate measurements we need the notion of **random vectors** ($\vec{\text{RVs}}$), i.e. ordered pairs of random variables (X, Y) , ordered triples of random variables (X, Y, Z) , or more generally ordered m -tuples of random variables (X_1, X_2, \dots, X_m) .

3.10.1 \mathbb{R}^2 -valued Random Variables

We first focus on understanding (X, Y) , a bivariate $\vec{\text{RV}}$ 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

$$\begin{aligned} F_{X,Y}(x, y) &= P(X \leq x \cap Y \leq y) = P(X \leq x, Y \leq y) \\ &= P(\{\omega : X(\omega) \leq x, Y(\omega) \leq y\}), \text{ for any } (x, y) \in \mathbb{R}^2 , \end{aligned} \quad (3.57)$$

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 $p_{i,j} = P(X = x_i, Y = y_j) > 0$, then its **joint probability mass function** (or JPMF) is:

$$f_{X,Y}(x_i, y_j) = P(X = x_i, Y = y_j) = \begin{cases} p_{i,j} & \text{if } (x_i, y_j) \in \mathcal{S}_{X,Y} \\ 0 & \text{otherwise} \end{cases} . \quad (3.58)$$

Since $P(\Omega) = 1$, $\sum_{(x_i, y_j) \in \mathcal{S}_{X,Y}} f_{X,Y}(x_i, y_j) = 1$.

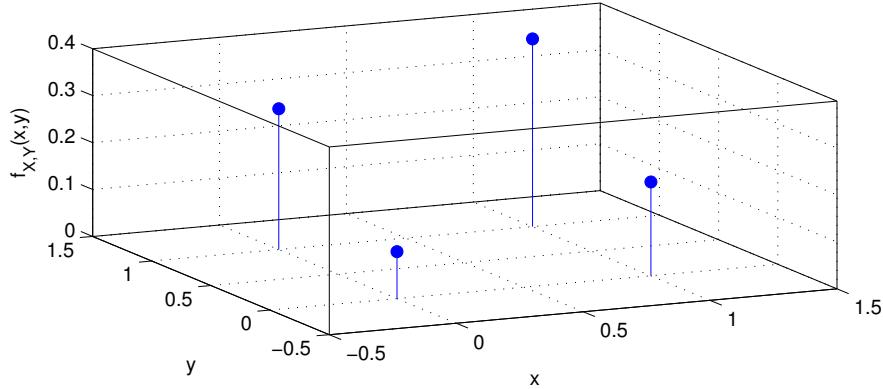
From JPMF $f_{X,Y}$ we can get the values of the JDF $F_{X,Y}(x, y)$ and the probability of any event B by simply taking sums,

$$F_{X,Y}(x, y) = \sum_{x_i \leq x, y_j \leq y} f_{X,Y}(x_i, y_j) ,$$

$$P(B) = \sum_{(x_i, y_j) \in B \cap \mathcal{S}_{X,Y}} f_{X,Y}(x_i, y_j) , \quad (3.59)$$

Example 83 Let (X, Y) be a discrete bivariate R.V with the following joint probability mass function (JPMF):

$$f_{X,Y}(x,y) := P(X = x, Y = y) = \begin{cases} 0.1 & \text{if } (x,y) = (0,0) \\ 0.3 & \text{if } (x,y) = (0,1) \\ 0.2 & \text{if } (x,y) = (1,0) \\ 0.4 & \text{if } (x,y) = (1,1) \\ 0.0 & \text{otherwise.} \end{cases}$$



It is helpful to write down the JPMF $f_{X,Y}(x,y)$ in a tabular form:

	$Y = 0$	$Y = 1$
$X = 0$	0.1	0.3
$X = 1$	0.2	0.4

From the above Table we can read for instance that the joint probability $f_{X,Y}(0,0) = 0.1$.

Find $P(B)$ for the event $B = \{(0,0), (1,1)\}$, $F_{X,Y}(1/2, 1/2)$, $F_{X,Y}(3/2, 1/2)$, $F_{X,Y}(4, 5)$ and $F_{X,Y}(-4, -1)$.

1. $P(B) = \sum_{(x,y) \in \{(0,0), (1,1)\}} f_{X,Y}(x,y) = f_{X,Y}(0,0) + f_{X,Y}(1,1) = 0.1 + 0.4$
2. $F_{X,Y}(1/2, 1/2) = \sum_{\{(x,y): x \leq 1/2, y \leq 1/2\}} f_{X,Y}(x,y) = f_{X,Y}(0,0) = 0.1$
3. $F_{X,Y}(3/2, 1/2) = \sum_{\{(x,y): x \leq 3/2, y \leq 1/2\}} f_{X,Y}(x,y) = f_{X,Y}(0,0) + f_{X,Y}(1,0) = 0.1 + 0.2 = 0.3$
4. $F_{X,Y}(4, 5) = \sum_{\{(x,y): x \leq 4, y \leq 5\}} f_{X,Y}(x,y) = f_{X,Y}(0,0) + f_{X,Y}(0,1) + f_{X,Y}(1,0) + f_{X,Y}(1,1) = 1$
5. $F_{X,Y}(-4, -1) = \sum_{\{(x,y): x \leq -4, y \leq -1\}} f_{X,Y}(x,y) = 0$

Definition 36 (JPDF) We say (X, Y) is a **continuous \mathbb{R}^2 -valued random variable** if its JDF $F_{X,Y}(x, y)$ is differentiable and its **joint probability density function (JPDF)** is given by:

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

For notational convenience, we sometimes suppress the subscripting when the random variables are clear from the context and write $f(x, y)$ and $F(x, y)$ instead of $f_{X,Y}(x, y)$ and $F_{X,Y}(x, y)$, respectively.

From JPDF $f_{X,Y}$ we can compute the JDF $F_{X,Y}$ at any point $(x, y) \in \mathbb{R}^2$ and more generally we can compute the probability of any event B , that can be cast as a region in \mathbb{R}^2 , by simply taking two-dimensional integrals:

$$F_{X,Y}(x, y) = \int_{-\infty}^y \int_{-\infty}^x f_{X,Y}(u, v) dudv , \quad (3.60)$$

and

$$\text{P}(B) = \int \int_B f_{X,Y}(x, y) dx dy . \quad (3.61)$$

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 $\delta \rightarrow 0$:

$$\text{P}((X, Y) \in \mathbb{B}_\delta(x, y)) \approx \delta f_{X,Y}(x, y) . \quad (3.62)$$

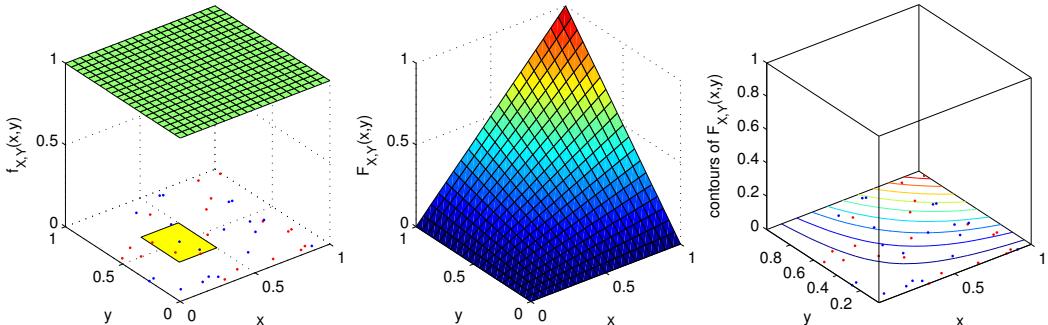
The JPDF satisfies the following two properties:

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 R.V 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) $\text{P}(X \leq 1/3, Y \leq 1/2)$, (3) $\text{P}((X, Y) \in [1/4, 1/2] \times [1/3, 2/3])$.



Let us begin to find the needed expressions.

1. Let $(x, y) \in [0, 1]^2$ then by Equation (3.60):

$$F_{X,Y}(x, y) = \int_{-\infty}^y \int_{-\infty}^x f_{X,Y}(u, v) dudv = \int_0^y \int_0^x 1 dudv = \int_0^y [u]_{u=0}^x dv = \int_0^y x dv = [xv]_{v=0}^y = xy$$

2. We can obtain $P(X \leq 1/3, Y \leq 1/2)$ by evaluating $F_{X,Y}$ at $(1/3, 1/2)$:

$$P(X \leq 1/3, Y \leq 1/2) = F_{X,Y}(1/3, 1/2) = \frac{1}{3} \frac{1}{2} = \frac{1}{6}$$

We can also find $P(X \leq 1/3, Y \leq 1/2)$ 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 $P(A) = (1/3)(1/2) = 1/6$.

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):

$$\begin{aligned} P((X, Y) \in [1/4, 1/2] \times [1/3, 2/3]) &= \int \int_B f_{X,Y}(x, y) dx dy = \int_{1/3}^{2/3} \int_{1/4}^{1/2} 1 dx dy \\ &= \int_{1/3}^{2/3} [x]_{1/4}^{1/2} dy = \int_{1/3}^{2/3} \left[\frac{1}{2} - \frac{1}{4} \right] dy = \left(\frac{1}{2} - \frac{1}{4} \right) [y]_{1/3}^{2/3} \\ &= \left(\frac{1}{2} - \frac{1}{4} \right) \left(\frac{2}{3} - \frac{1}{3} \right) = \frac{1}{4} \left(\frac{1}{3} \right) = \frac{1}{12} \end{aligned}$$

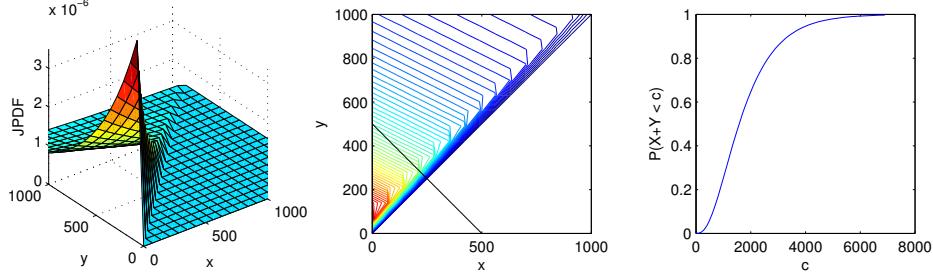
In general, for a bivariate uniform R \vec{V} on the unit square the $P([a, b] \times [c, d]) = (b-a)(d-c)$ for any event given by the rectangular region $[a, b] \times [c, d]$ inside the unit square $[0, 1] \times [0, 1]$. Thus any two events with the same rectangular area have the same probability (imagine sliding a small rectangle inside the unit square... no matter where you slide this rectangle to while remaining in the unit square, the probability of $\omega \mapsto (X(\omega), Y(\omega)) = (x, y)$ falling inside this “slidable” rectangle is the same...).

Example 85 Let the RV X denote the time until a web server connects to your computer, and let the RV Y denote the time until the server authorizes you as a valid user. Each of these RVs measures the waiting time from a common starting time (in milliseconds) and $X < Y$. From past response times of the web server we know that a good approximation for the JPDF of the R \vec{V} (X, Y) is

$$f_{X,Y}(x, y) = \begin{cases} \frac{6}{10^6} \exp\left(-\frac{1}{1000}x - \frac{2}{1000}y\right) & \text{if } x > 0, y > 0, x < y \\ 0 & \text{otherwise.} \end{cases}$$

Answer the following:

1. identify the support of (X, Y) , i.e., the region in the plane where $f_{X,Y}$ takes positive values
2. check that $f_{X,Y}$ indeed integrates to 1 as it should
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 $P(X + Y < 10^2)$?



Let us answer the questions.

1. The support is the intersection of the positive quadrant with the $y > x$ half-plane.
- 2.

$$\begin{aligned}
 \int_{y=-\infty}^{\infty} \int_{x=-\infty}^{\infty} f_{X,Y}(x,y) dx dy &= \int_{x=0}^{\infty} \int_{y=x}^{\infty} f_{X,Y}(x,y) dy dx \\
 &= \int_{x=0}^{\infty} \int_{y=x}^{\infty} \frac{6}{10^6} \exp\left(-\frac{1}{1000}x - \frac{2}{1000}y\right) dy dx \\
 &= \frac{6}{10^6} \int_{x=0}^{\infty} \left(\int_{y=x}^{\infty} \exp\left(-\frac{2}{1000}y\right) dy \right) \exp\left(-\frac{1}{1000}x\right) dx \\
 &= \frac{6}{10^6} \int_{x=0}^{\infty} \left[-\frac{1000}{2} \exp\left(-\frac{2}{1000}y\right) \right]_{y=x}^{\infty} \exp\left(-\frac{1}{1000}x\right) dx \\
 &= \frac{6}{10^6} \int_{x=0}^{\infty} \left[0 + \frac{1000}{2} \exp\left(-\frac{2}{1000}x\right) \right] \exp\left(-\frac{1}{1000}x\right) dx \\
 &= \frac{6}{10^6} \int_{x=0}^{\infty} \frac{1000}{2} \exp\left(-\frac{2}{1000}x - \frac{1}{1000}x\right) dx \\
 &= \frac{6}{10^6} \frac{1000}{2} \left[-\frac{1000}{3} \exp\left(-\frac{3}{1000}x\right) \right]_{x=0}^{\infty} \\
 &= \frac{6}{10^6} \frac{1000}{2} \left[0 + \frac{1000}{3} \right] \\
 &= 1
 \end{aligned}$$

3. First, identify the region with positive JPDF for the event $(X \leq 400, Y \leq 800)$

$$\begin{aligned}
 P(X \leq 400, Y \leq 800) &= \int_{x=0}^{400} \int_{y=x}^{800} f_{X,Y}(x,y) dy dx \\
 &= \int_{x=0}^{400} \int_{y=x}^{800} \frac{6}{10^6} \exp\left(-\frac{1}{1000}x - \frac{2}{1000}y\right) dy dx \\
 &= \frac{6}{10^6} \int_{x=0}^{400} \left[-\frac{1000}{2} \exp\left(-\frac{2}{1000}y\right) \right]_{y=x}^{800} \exp\left(-\frac{1}{1000}x\right) dx \\
 &= \frac{6}{10^6} \frac{1000}{2} \int_{x=0}^{400} \left(-\exp\left(-\frac{1600}{1000}\right) + \exp\left(-\frac{2}{1000}x\right) \right) \exp\left(-\frac{1}{1000}x\right) dx \\
 &= \frac{6}{10^6} \frac{1000}{2} \int_{x=0}^{400} \left(\exp\left(-\frac{3}{1000}x\right) - e^{-8/5} \exp\left(-\frac{1}{1000}x\right) \right) dx \\
 &= \frac{6}{10^6} \frac{1000}{2} \left(\left(-\frac{1000}{3} \exp\left(-\frac{3}{1000}x\right) \right)_{x=0}^{400} - e^{-8/5} \left(-1000 \exp\left(-\frac{1}{1000}x\right) \right)_{x=0}^{400} \right) \\
 &= \frac{6}{10^6} \frac{1000}{2} 1000 \left(\frac{1}{3} \left(1 - e^{-6/5} \right) - e^{-8/5} \left(1 - e^{-2/5} \right) \right) \\
 &= 3 \left(\frac{1}{3} \left(1 - e^{-6/5} \right) - e^{-8/5} \left(1 - e^{-2/5} \right) \right) \\
 &\approx 0.499
 \end{aligned}$$

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 the JPDF over our triangular event as follows:

$$\begin{aligned}
P(X + Y \leq c) &= \int_{x=0}^{c/2} \int_{y=x}^{c-x} f_{X,Y}(x,y) dy dx \\
&= \int_{x=0}^{c/2} \int_{y=x}^{c-x} \frac{6}{10^6} \exp\left(-\frac{1}{1000}x - \frac{2}{1000}y\right) dy dx \\
&= \frac{6}{10^6} \int_{x=0}^{c/2} \int_{y=x}^{c-x} \exp\left(-\frac{1}{1000}x - \frac{2}{1000}y\right) dy dx \\
&= \frac{6}{10^6} \frac{1000}{2} \int_{x=0}^{c/2} \left[-\exp\left(-\frac{2}{1000}y\right) \right]_{y=x}^{c-x} \exp\left(-\frac{1}{1000}x\right) dx \\
&= \frac{3}{10^3} \int_{x=0}^{c/2} \left[-\exp\left(-\frac{2c-2x}{1000}\right) + \exp\left(-\frac{2x}{1000}\right) \right] \exp\left(-\frac{x}{1000}\right) dx \\
&= \frac{3}{10^3} \int_{x=0}^{c/2} \left(\exp\left(-\frac{3x}{1000}\right) - \exp\left(\frac{x-2c}{1000}\right) \right) dx \\
&= 3 \left(\left[-\frac{1}{3} \exp\left(-\frac{3x}{1000}\right) \right]_{x=0}^{c/2} - \left[e^{-2c/1000} \exp\left(\frac{x}{1000}\right) \right]_{x=0}^{c/2} \right) \\
&= 3 \left(\frac{1}{3} (1 - e^{-3c/2000}) - e^{-2c/1000} (e^{c/2000} - 1) \right) \\
&= 1 - e^{-3c/2000} + 3e^{-2c/1000} - 3e^{-3c/2000} \\
&= 1 - 4e^{-3c/2000} + 3e^{-c/500}
\end{aligned}$$

5. $P(X + Y < 100) = 1 - 4e^{-300/2000} + 3e^{-100/500} \approx 0.134$. This means only about one in one hundred requests to this server will be processed within 100 milliseconds.

We can obtain $P(X + Y < c)$ for several values of c using MATLAB and note that about 96% of requests are processed in less than 3000 milliseconds or 3 seconds.

```

>> c = [100 1000 2000 3000 4000]
c = 100      1000      2000      3000      4000

>> p = 1 - 4 * exp(-3*c/2000) + 3 * exp(-c/500)

p = 0.0134    0.5135    0.8558    0.9630    0.9911

```

Definition 37 (Marginal PDF or PMF) If the R \vec{V} (X, Y) has $f_{X,Y}(x, y)$ as its joint PDF or joint PMF, then the **marginal PDF or PMF** of a random vector (X, Y) is defined by :

$$f_X(x) = \begin{cases} \int_{-\infty}^{\infty} f_{X,Y}(x, y) dy & \text{if } (X, Y) \text{ is a continuous R}\vec{V} \\ \sum_y f_{X,Y}(x, y) & \text{if } (X, Y) \text{ is a discrete R}\vec{V} \end{cases}$$

and the **marginal PDF or PMF** of Y is defined by:

$$f_Y(y) = \begin{cases} \int_{-\infty}^{\infty} f_{X,Y}(x, y) dx & \text{if } (X, Y) \text{ is a continuous R}\vec{V} \\ \sum_x f_{X,Y}(x, y) & \text{if } (X, Y) \text{ is a discrete R}\vec{V} \end{cases}$$

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 R \vec{V} in Example 83. Just sum $f_{X,Y}(x, y)$ over x 's and y 's (reported in a tabular form):

	$Y = 0$	$Y = 1$
$X = 0$	0.1	0.3
$X = 1$	0.2	0.4

From the above Table we can find:

$$f_X(x) = P(X = x) = \sum_y f_{X,Y}(x,y)$$

$$= f_{X,Y}(x,0) + f_{X,Y}(x,1) = \begin{cases} f_{X,Y}(0,0) + f_{X,Y}(0,1) = 0.1 + 0.3 = 0.4 & \text{if } x = 0 \\ f_{X,Y}(1,0) + f_{X,Y}(1,1) = 0.2 + 0.4 = 0.6 & \text{if } x = 1 \end{cases}$$

Similarly,

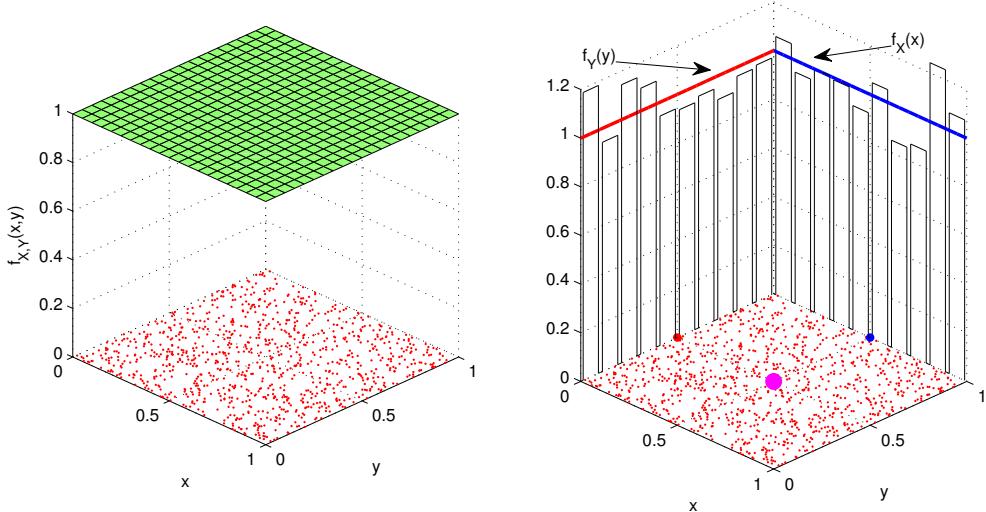
$$f_Y(y) = P(Y = y) = \sum_x f_{X,Y}(x,y)$$

$$= f_{X,Y}(0,y) + f_{X,Y}(1,y) = \begin{cases} f_{X,Y}(0,0) + f_{X,Y}(1,0) = 0.1 + 0.2 = 0.3 & \text{if } y = 0 \\ f_{X,Y}(0,1) + f_{X,Y}(1,1) = 0.3 + 0.4 = 0.7 & \text{if } y = 1 \end{cases}$$

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,

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

We are seeing a histogram of the **marginal samples** and their marginal PDFs in the Figure.

Thus marginal PDF gives us the probability density of a specific RV in a R \vec{V} , irrespective of the value taken by the other RV in this R \vec{V} .

Example 88 Obtain the marginal PDF $f_Y(y)$ from the joint PDF $f_{X,Y}(x,y)$ of the continuous R \vec{V} in Example 85 that gave the response times of a web server.

$$f_{X,Y}(x,y) = \begin{cases} \frac{6}{10^6} \exp\left(-\frac{1}{1000}x - \frac{2}{1000}y\right) & \text{if } x > 0, y > 0, x < y \\ 0 & \text{otherwise.} \end{cases}$$

Use $f_Y(y)$ to compute the probability that Y exceeds 2000 milliseconds.

First we need to obtain an expression for $f_Y(y)$. For $y > 0$,

$$\begin{aligned} f_Y(y) &= \int_{x=-\infty}^{\infty} f_{X,Y}(x,y)dx \\ &= \int_{x=-\infty}^{\infty} 6 \times 10^{-6} e^{-0.001x - 0.002y} dx \\ &= 6 \times 10^{-6} \int_{x=0}^y e^{-0.001x - 0.002y} dx \\ &= 6 \times 10^{-6} e^{-0.002y} \int_{x=0}^y e^{-0.001x} dx \\ &= 6 \times 10^{-6} e^{-0.002y} \left[\frac{e^{-0.001x}}{-0.001} \right]_{x=0}^{x=y} \\ &= 6 \times 10^{-6} e^{-0.002y} \left(\frac{e^{-0.001y}}{-0.001} - \frac{e^{-0.001 \times 0}}{-0.001} \right) \\ &= 6 \times 10^{-6} e^{-0.002y} \left(\frac{1 - e^{-0.001y}}{0.001} \right) \\ &= 6 \times 10^{-3} e^{-0.002y} (1 - e^{-0.001y}) \end{aligned}$$

We have the marginal PDF of Y and from this we can obtain

$$\begin{aligned} P(Y > 2000) &= \int_{2000}^{\infty} f_Y(y)dy \\ &= \int_{2000}^{\infty} 6 \times 10^{-3} e^{-0.002y} (1 - e^{-0.001y}) dy \\ &= 6 \times 10^{-3} \int_{2000}^{\infty} e^{-0.002y} dy - \int_{2000}^{\infty} e^{-0.003y} dy \\ &= 6 \times 10^{-3} \left(\left[\frac{e^{-0.002y}}{-0.002} \right]_{2000}^{\infty} - \left(\left[\frac{e^{-0.003y}}{-0.003} \right]_{2000}^{\infty} \right) \right) \\ &= 6 \times 10^{-3} \left(\frac{e^{-4}}{0.002} - \frac{e^{-6}}{0.003} \right) \\ &= 0.05 \end{aligned}$$

Alternatively, you can obtain $P(Y > 2000)$ by directly integrating the joint PDF $f_{X,Y}(x,y)$ over the appropriate region (but you may now have to integrate two pieces: rectangular infinite strip $(x,y) : 0 < x < 2000, y > 2000$ and a triangular

infinite piece $\{(x, y) : y > x, y > 2000, x > 2000\}$)... more involved but we get the same answer.

$$\begin{aligned} P(Y > 2000) &= \int_{x=0}^{2000} \left(\int_{y=2000}^{\infty} 6 \times 10^{-6} e^{-0.001x-0.002y} dy \right) dx + \\ &\quad \int_{x=2000}^{\infty} \left(\int_{y=x}^{\infty} 6 \times 10^{-6} e^{-0.001x-0.002y} dy \right) dx \\ &\quad \vdots \text{(try as a tutorial problem)} \\ P(Y > 2000) &= 0.0475 + 0.0025 = 0.05 \end{aligned}$$

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

$$P(X_1 \leq x_1, X_2 \leq x_2) = P(X_1 \leq x_1) P(X_2 \leq x_2)$$

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:

$$f_{X_1, X_2}(x_1, x_2) = P(X_1 = x_1, X_2 = x_2) = P(X_1 = x_1) P(X_2 = x_2) = f_{X_1}(x_1) f_{X_2}(x_2) \text{ for any } (x_1, x_2) \in \mathbb{R}^2 ,$$

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., $P((X_1, X_2) = (0, 0)) = 1/4$, and that the *marginal probability* that the RV X_1 takes the value or realization 0 is $1/2$, i.e., $P(X_1 = 0) = 1/2$. 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.

Example 90 Recall the \mathbb{R}^2 -valued continuous RV (X, Y) of Example 84 that is uniformly distributed on the unit square $[0, 1]^2$. First show that X and Y independent. Then show that both X and Y are identically distributed according to the Uniform($0, 1$) RV.

Solution:

This can be shown by checking that the joint PDF is indeed equal to the product of the marginal PDFs of Uniform($0, 1$) RVs as follows:

$$\begin{cases} 1 = f_{X,Y}(x, y) = f_X(x) \times f_Y(y) = 1 \times 1 = 1 & \text{if } (x, y) \in [0, 1]^2 \\ 0 = f_{X,Y}(x, y) = f_X(x) \times f_Y(y) = 0 \times 0 = 0 & \text{if } (x, y) \notin [0, 1]^2 \end{cases}$$

Are X and Y independent in the server times \vec{RV} from Example 85?

We can compute $f_X(x)$ and use the already computed $f_Y(y)$ to mechanically check if the JPDF is the product of the marginal PDFs. But intuitively, we know that these RVs (connection time and authentication time) are dependent – one is strictly greater than the other. Also the JPDF has zero density when $x > y$, but the product of the marginal densities won't.

Now, let us take advantage of independent random variables and solve some problems.

Example 91 (distance between random faults 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 ?

done in lectures...

Example 92 (Buffon's Needle Experiment to Physically Estimate π) 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 $l \leq L$.

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 ?

Solution:

Let X_1 be the angle between the needle and the direction of the rulings, and let X_2 be the

Figure 3.19: Diagrams done on the board!

distance between the bottom point of the needle and the nearest line above this point (see left sub-figure of Figure 3.19). Then the conditions of the “needle tossing at random” experiment are such that the RV X_1 is uniformly distributed in the interval $[0, \pi]$, while the RV X_2 is uniformly distributed in the interval $[0, L]$. Hence *assuming that the RVs X_1 and X_2 are independent*, we find that their joint probability density function (JPDF) is:

$$f_{X_1, X_2}(x_1, x_2) = \frac{1}{\pi} \mathbb{1}_{[0, \pi]}(x_1) \times \frac{1}{L} \mathbb{1}_{[0, L]}(x_2) = \frac{1}{\pi L} \mathbb{1}_{[0, \pi]}(x_1) \mathbb{1}_{[0, L]}(x_2) = \frac{1}{\pi L} \mathbb{1}_{[0, \pi] \times [0, L]}(x_1, x_2) .$$

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:

$$P(A) = P((X_1, X_2) \in B) = \int_B \int \frac{dx_1 dx_2}{\pi L} = \frac{2l}{\pi L}$$

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 $P(A)$ 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 RV. The conditional PMF of $X_1|X_2 = x_2$, where $f_{X_2}(x_2) := P(X_2 = x_2) > 0$ is:

$$f_{X_1|X_2}(x_1|x_2) := P(X_1 = x_1|X_2 = x_2) = \frac{P(X_1 = x_1, X_2 = x_2)}{P(X_2 = x_2)} = \frac{f_{X_1,X_2}(x_1, x_2)}{f_{X_2}(x_2)} .$$

Similarly, if $f_{X_1}(x_1) := P(X_1 = x_1) > 0$, then the conditional PMF of $X_2|X_1 = x_1$ is:

$$f_{X_2|X_1}(x_2|x_1) := P(X_2 = x_2|X_1 = x_1) = \frac{P(X_1 = x_1, X_2 = x_2)}{P(X_1 = x_1)} = \frac{f_{X_1,X_2}(x_1, x_2)}{f_{X_1}(x_1)} .$$

If (X_1, X_2) are continuous RVs such that the marginal PDF $f_{X_2}(x_2) > 0$, then the conditional PDF of $X_1|X_2 = x_2$ is:

$$f_{X_1|X_2}(x_1|x_2) = \frac{f_{X_1,X_2}(x_1, x_2)}{f_{X_2}(x_2)}, \quad P(X_1 \in A|X_2 = x_2) = \int_A f_{X_1|X_2}(x_1|x_2) dx_1 .$$

Similarly, if $f_{X_1}(x_1) > 0$, then the conditional PDF of $X_2|X_1 = x_1$ is:

$$f_{X_2|X_1}(x_2|x_1) = \frac{f_{X_1,X_2}(x_1, x_2)}{f_{X_1}(x_1)}, \quad P(X_2 \in A|X_1 = x_1) = \int_A f_{X_2|X_1}(x_2|x_1) dx_2 .$$

Let us consider a few discrete RVs for the simple coin tossing experiment \mathcal{E}_θ^3 that build on the Bernoulli(θ) RV X_i for the i -th toss in an **independent and identically distributed (IID.)** manner.

Table 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.

$\omega:$	HHH	HHT	HTH	HTT	THH	THT	TTH	TTT	RV Definitions / Model
$P(\omega):$	$\frac{1}{8}$	$X_i \stackrel{\text{IID}}{\sim} \text{Bernoulli}(\frac{1}{2})$							
$Y(\omega):$	3	2	2	1	2	1	1	0	$Y := X_1 + X_2 + X_3$
$Z(\omega):$	0	1	1	2	1	2	2	3	$Z := (1 - X_1) + (1 - X_2) + (1 - X_3)$
$Y'(\omega):$	1	0	0	0	0	0	0	0	$Y' := X_1 X_2 X_3$
$Z'(\omega):$	0	0	0	0	0	0	0	1	$Z' := (1 - X_1)(1 - X_2)(1 - X_3)$

Classwork 93 (Two random variables of ‘toss a coin thrice’ 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 your spare time.

$$P(Y = y) = \left\{ \begin{array}{l} \\ \\ \\ \end{array} \right. \quad P(Y' = y') = \left\{ \begin{array}{l} \\ \\ \\ \end{array} \right.$$

Classwork 94 (The number of ‘Heads’ given there is at least one ‘Tails’) Consider the following two questions.

- What is conditional probability $P(Y|Y' = 0)$?

$P(Y = y Y' = 0)$	$= \frac{P(Y=y, Y'=0)}{P(Y'=0)}$	$= \frac{P(\{\omega : Y(\omega)=y \cap Y'(\omega)=0\})}{P(\{\omega : Y'(\omega)=0\})}$	$= ?$
$P(Y = 0 Y' = 0)$	$\frac{P(Y=0, Y'=0)}{P(Y'=0)}$	$\frac{\frac{1}{8}}{\frac{1}{8} + \frac{1}{8} + \frac{1}{8} + \frac{1}{8} + \frac{1}{8} + \frac{1}{8}}$	$\frac{1}{7}$
$P(Y = 1 Y' = 0)$	$\frac{P(Y=1, Y'=0)}{P(Y'=0)}$	$\frac{\frac{1}{8} + \frac{1}{8} + \frac{1}{8}}{\frac{1}{8} + \frac{1}{8} + \frac{1}{8} + \frac{1}{8} + \frac{1}{8} + \frac{1}{8}}$	$\frac{3}{7}$
$P(Y = 2 Y' = 0)$	$\frac{P(Y=2, Y'=0)}{P(Y'=0)}$	$\frac{\frac{1}{8} + \frac{1}{8} + \frac{1}{8}}{\frac{1}{8} + \frac{1}{8} + \frac{1}{8} + \frac{1}{8} + \frac{1}{8} + \frac{1}{8}}$	$\frac{3}{7}$
$P(Y = 3 Y' = 0)$	$\frac{P(Y=3, Y'=0)}{P(Y'=0)}$	$\frac{P(\emptyset)}{\frac{1}{8} + \frac{1}{8} + \frac{1}{8} + \frac{1}{8} + \frac{1}{8} + \frac{1}{8}}$	0
$P(Y \in \{0, 1, 2, 3\} Y' = 0)$	$\frac{\sum_{y=0}^3 P(Y=y, Y'=0)}{P(Y'=0)}$	$\frac{\frac{1}{8} + \frac{1}{8} + \frac{1}{8} + \frac{1}{8}}{\frac{1}{8} + \frac{1}{8} + \frac{1}{8} + \frac{1}{8} + \frac{1}{8} + \frac{1}{8}}$	1

- What is $P(Y|Y' = 1)$?

$$P(Y = y|Y' = 1) = \begin{cases} 1 & \text{if } y = 3 \\ 0 & \text{otherwise} \end{cases}$$

3.10.3 \mathbb{R}^m -valued Random Variables

Consider the RV 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

$$\begin{aligned} F_{X_1, X_2, \dots, X_m}(x_1, x_2, \dots, x_m) &= P(X \leq x_1 \cap X_2 \leq x_2 \cap \dots \cap X_m \leq x_m) \\ &= P(X_1 \leq x_1, X_2 \leq x_2, \dots, X_m \leq x_m) \\ &= P(\{\omega : X_1(\omega) \leq x_1, X_2(\omega) \leq x_2, \dots, X_m(\omega) \leq x_m\}), \end{aligned} \tag{3.63}$$

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 $x_m \rightarrow \infty$
4. $F_{X_1, X_2, \dots, X_m}(x_1, x_2, \dots, x_m) \rightarrow 0$ as $x_1 \rightarrow -\infty, x_2 \rightarrow -\infty, \dots$ and $x_m \rightarrow -\infty$

Definition 41 (Multivariate JPMF) If (X_1, X_2, \dots, X_m) is a **discrete random vector** that takes values in a discrete support set $\mathcal{S}_{X_1, X_2, \dots, X_m}$, then its **joint probability mass function** (or JPMF) is:

$$f_{X_1, X_2, \dots, X_m}(x_1, x_2, \dots, x_m) = P(X_1 = x_1, X_2 = x_2, \dots, X_m = x_m) . \quad (3.64)$$

Since $P(\Omega) = 1$, $\sum_{(x_1, x_2, \dots, x_m) \in \mathcal{S}_{X_1, X_2, \dots, X_m}} f_{X_1, X_2, \dots, X_m}(x_1, x_2, \dots, x_m) = 1$.

From JPMF f_{X_1, X_2, \dots, X_m} we can get the JCDF $F_{X_1, X_2, \dots, X_m}(x_1, x_2, \dots, x_m)$ and the probability of any event B by simply taking sums as in Equation (3.59) but now over all m coordinates.

Definition 42 (Multivariate JPFD) (X_1, X_2, \dots, X_m) is a **continuous random vector** if its JDF $F_{X_1, X_2, \dots, X_m}(x_1, x_2, \dots, x_m)$ is differentiable and the **joint probability density function (JPDF)** is given by:

$$f_{X_1, X_2, \dots, X_m}(x_1, x_2, \dots, x_m) = \frac{\partial^m}{\partial x_1 \partial x_2 \dots \partial x_m} F_{X_1, X_2, \dots, X_m}(x_1, x_2, \dots, x_m) ,$$

From JPFD f_{X_1, X_2, \dots, X_m} we can compute the JDF F_{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 in \mathbb{R}^m , by “simply” taking m -dimensional integrals (you have done such iterated integrals when $m = 3$):

$$F_{X_1, X_2, \dots, X_m}(x_1, x_2, \dots, x_m) = \int_{-\infty}^{x_m} \dots \int_{-\infty}^{x_2} \int_{-\infty}^{x_1} f_{X_1, X_2, \dots, X_m}(x_1, x_2, \dots, x_m) dx_1 dx_2 \dots dx_m , \quad (3.65)$$

and

$$P(B) = \int \dots \int_B f_{X_1, X_2, \dots, X_m}(x_1, x_2, \dots, x_m) dx_1 dx_2 \dots dx_m . \quad (3.66)$$

The JPFD satisfies the following two properties:

1. integrates to 1, i.e., $\int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} f_{X_1, X_2, \dots, X_m}(x_1, x_2, \dots, x_m) dx_1 dx_2 \dots dx_m = 1$
2. is a non-negative function, i.e., $f_{X_1, X_2, \dots, X_m}(x_1, x_2, \dots, x_m) \geq 0$.

The marginal PDF (marginal PMF) is obtained by integrating (summing) the JPFD (JPMF) over all other random variables. For example, the marginal PDF of X_1 is

$$f_{X_1}(x_1) = \int_{x_2=-\infty}^{\infty} \dots \int_{x_m=-\infty}^{\infty} f_{X_1, X_2, \dots, X_m}(x_1, x_2, \dots, x_m) dx_2 \dots dx_m$$

Definition 43 (Independence of Sequence of RVs) A finite or infinite sequence of RVs X_1, X_2, \dots is said to be independent or independently distributed if and only if

$$P(X_{i_1} \leq x_{i_1}, X_{i_2} \leq x_{i_2}, \dots, X_{i_k} \leq x_{i_k}) = P(X_{i_1} \leq x_{i_1}) P(X_{i_2} \leq x_{i_2}) \cdots P(X_{i_k} \leq x_{i_k})$$

or equivalently,

$$F_{X_{i_1}, X_{i_2}, \dots, X_{i_m}}(x_{i_1}, x_{i_2}, \dots, x_{i_m}) = F_{X_{i_1}}(x_{i_1}) F_{X_{i_2}}(x_{i_2}) \cdots F_{X_{i_m}}(x_{i_m}) ,$$

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:

$$P(X_{i_1} = x_{i_1}, X_{i_2} = x_{i_2}, \dots, X_{i_k} = x_{i_k}) = P(X_{i_1} = x_{i_1}) P(X_{i_2} = x_{i_2}) \cdots P(X_{i_k} = x_{i_k}) ,$$

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

$$P(X_{i+1} \leq x_{i+1} | X_i \leq x_i, X_{i-1} \leq x_{i-1}, \dots, X_1 \leq x_1) = P(X_{i+1} \leq x_{i+1}) \quad (3.69)$$

Proof:

$$\begin{aligned} & P(X_{i+1} \leq x_{i+1} | X_i \leq x_i, X_{i-1} \leq x_{i-1}, \dots, X_1 \leq x_1) \\ &= \frac{P(X_{i+1} \leq x_{i+1}, X_i \leq x_i, X_{i-1} \leq x_{i-1}, \dots, X_1 \leq x_1)}{P(X_i \leq x_i, X_{i-1} \leq x_{i-1}, \dots, X_1 \leq x_1)} \\ &= \frac{P(X_{i+1} \leq x_{i+1}) P(X_i \leq x_i) P(X_{i-1} \leq x_{i-1}) \cdots P(X_1 \leq x_1)}{P(X_i \leq x_i) P(X_{i-1} \leq x_{i-1}) \cdots P(X_1 \leq x_1)} \\ &= P(X_{i+1} \leq x_{i+1}) \end{aligned}$$

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 $\text{Cov}(X_1, X_2)$?

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

$$\begin{aligned}\text{Cov}(X_1, X_2) &= E(X_1 X_2) - E(X_1)E(X_2) \\ &= E(X_1)E(X_2) - E(X_1)E(X_2) \quad \text{due to independence} \\ &= 0\end{aligned}$$

Remark 45 The converse is not true: two random variables that have zero covariance are not necessarily independent.

Linear Combination of Independent Normal RVs is a Normal RV

We can get the following special property of normal RVs using Eqn. (3.74). If X_1, X_2, \dots, X_m be jointly independent RVs, where X_i is $\text{Normal}(\mu_i, \sigma_i^2)$, 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 $\text{Normal}(c + \sum_{i=1}^m a_i \mu_i, \sum_{i=1}^m a_i^2 \sigma_i^2)$ RV.

Example 96 Let X be $\text{Normal}(2, 4)$, Y be $\text{Normal}(-1, 2)$ and Z be $\text{Normal}(0, 1)$ RVs that are jointly independent. Obtain the following:

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 $W = X - Y$.

Solution

1.

$$E(3X - 2Y + 4Z) = 3E(X) - 2E(Y) + 4(Z) = (3 \times 2) + (-2 \times (-1)) + 4 \times 0 = 6 + 2 + 0 = 8$$

2.

$$V(2Y - 3Z) = 2^2 V(Y) + (-3)^2 V(Z) = (4 \times 2) + (9 \times 1) = 8 + 9 = 17$$

3. From the special property of normal RVs, the distribution of $6 - 2Z + X - Y$ is

$$\begin{aligned}&\text{Normal}(6 + (-2 \times 0) + (1 \times 2) + (-1 \times -1), ((-2)^2 \times 1) + (1^2 \times 4) + ((-1)^2 \times 2)) \\ &= \text{Normal}(6 + 0 + 2 + 1, 4 + 4 + 2) \\ &= \text{Normal}(9, 10)\end{aligned}$$

4. Let $U = 6 - 2Z + X - Y$ and we know U is Normal(9, 10) RV.

$$\begin{aligned}
P(6 - 2Z + X - Y > 0) &= P(U > 0) = P(U - 9 > 0 - 9) = P\left(\frac{U - 9}{\sqrt{10}} > \frac{-9}{\sqrt{10}}\right) \\
&= P\left(Z > \frac{-9}{\sqrt{10}}\right) \\
&= P\left(Z < \frac{9}{\sqrt{10}}\right) \\
&\approx P(Z < 2.85) = 0.9978
\end{aligned}$$

5.

$$\begin{aligned}
\text{Cov}(X, W) &= E(XW) - E(X)E(W) = E(X(X - Y)) - E(X)E(X - Y) \\
&= E(X^2 - XY) - E(X)(E(X) - E(Y)) = E(X^2) - E(XY) - 2 \times (2 - (-1)) \\
&= E(X^2) - E(X)E(Y) - 6 = E(X^2) - (2 \times (-1)) - 6 \\
&= (V(X) + (E(X))^2) + 2 - 6 = (4 + 2^2) - 4 = 4
\end{aligned}$$

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 R \vec{V} 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 R \vec{V} 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 .

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

$\mathbb{1}_A(x)$ 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 to its two sides in the geometry of addition in the plane?
4. What is $(1, 0) + (0, 1) + (1, 0)$?

Solution:

1. addition is component-wise

$$(1, 0) + (1, 0) = (1 + 1, 0 + 0) = (2, 0)$$

$$(1, 0) + (0, 1) = (1 + 0, 0 + 1) = (1, 1)$$

$$(0, 1) + (0, 1) = (0 + 0, 1 + 1) = (0, 2)$$

2. $(1, 0)$ and $(0, 1)$ are vectors for the two sides of unit square and $(1, 1)$ is its diagonal.

3. Generally, the diagonal of the parallelogram is the resultant or sum of the vectors representing its two sides

- 4.

$$(1, 0) + (0, 1) + (1, 0) = (1 + 0 + 1, 0 + 1 + 0) = (2, 1)$$

Model 15 (Bernoulli(θ) \vec{RV}) Given a parameter $(\theta, 1 - \theta) \in \Delta^1$, the unit 1-Simplex, we say that $X := (X_1, X_2)$ is a Bernoulli(θ) random vector (\vec{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 \vec{RV} $X := (X_1, X_2)$ with realization $x := (x_1, x_2)$ is:

$$f(x; \theta) := P(X = x) = \theta \mathbf{1}_{\{e_1\}}(x) + (1 - \theta) \mathbf{1}_{\{e_2\}}(x) = \begin{cases} \theta & \text{if } x = e_1 := (1, 0) \\ 1 - \theta & \text{if } x = e_2 := (0, 1) \\ 0 & \text{otherwise} \end{cases}$$

Example 98 Let us find the Expectation of Bernoulli(θ) \vec{RV} in Model 15.

$$\mathbb{E}_\theta(X) = \mathbb{E}_\theta((X_1, X_2)) = \sum_{(x_1, x_2) \in \{e_1, e_2\}} (x_1, x_2) f((x_1, x_2); \theta) = (1, 0)\theta + (0, 1)(1 - \theta) = (\theta, 1 - \theta).$$

Remark 46 We can write the Binomial(n, θ) RV Y as a Binomial(n, θ) \vec{RV} $X := (Y, n - Y)$. In fact, this is the underlying model and the **bi** in the Binomial(n, θ) does refer to two in Latin. In the coin-tossing context this can be thought of keeping track of the number of Heads and Tails out of an IID sequence of n tosses of a coin with probability θ of observing Heads. In the Quincunx context, this amounts to keeping track of the number of right and left turns made by the ball as it 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 Binomial(n, θ) \vec{RV} $(Y, n - Y)$ is the sum of n IID Bernoulli(θ) \vec{RV} 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 Bernoulli(θ) \vec{RV} . 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 Bernoulli($1/2$) \vec{RV} s) Let us understand the Quincunx construction of the Binomial($n, 1/2$) \vec{RV} X as the sum of n independent and identical Bernoulli($1/2$) \vec{RV} s by calling the interactive visual cognitive tool as follows:

```
>> guiMultinomial
```

Figure 3.20: Visual Cognitive Tool GUI: Quincunx & Septcunx.

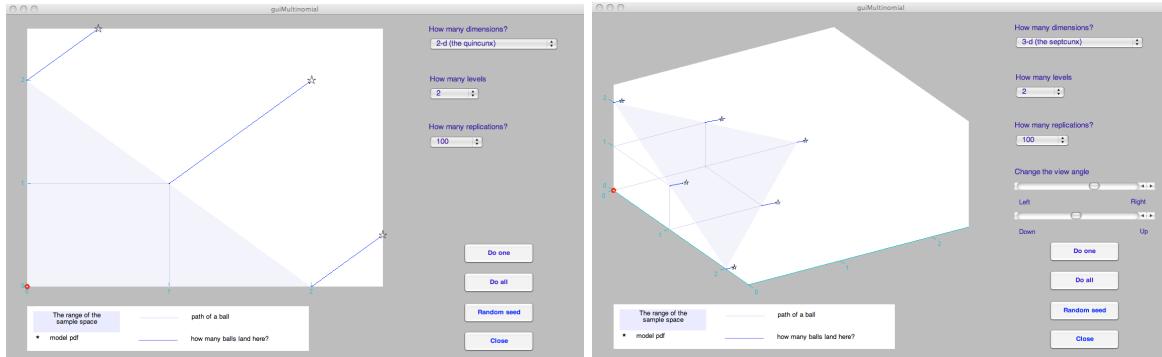
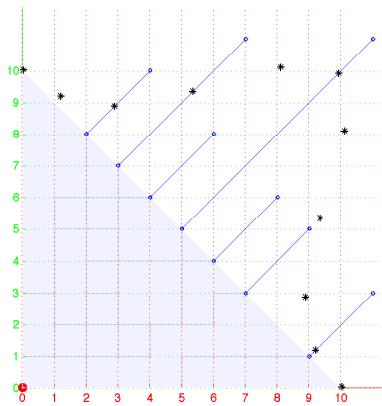
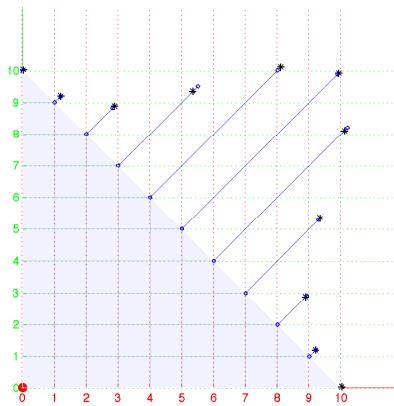


Figure 3.21: Quincunx on the Cartesian plane. Simulations of $\text{Binomial}(n = 10, \theta = 0.5)$ RV as the x-coordinate of the ordered pair resulting from the culmination of sample trajectories formed by the accumulating sum of $n = 10$ IID $\text{Bernoulli}(\theta = 0.5)$ random vectors over $\{(1, 0), (0, 1)\}$ 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.



(a) Ten samples



(b) Thousand samples

We are now ready to extend the $\text{Binomial}(n, \theta)$ RV or $R\vec{V}$ to its multivariate version called the $\text{Multinomial}(n, \theta_1, \theta_2, \dots, \theta_k)$ $R\vec{V}$. We develop this $R\vec{V}$ as the sum of n IID de Moivre($\theta_1, \theta_2, \dots, \theta_k$) $R\vec{V}$ that is defined next by extending 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 .

Model 16 (de Moivre($\theta_1, \theta_2, \dots, \theta_k$) $R\vec{V}$) The PMF of the de Moivre($\theta_1, \theta_2, \dots, \theta_k$) $R\vec{V}$ $X := (X_1, X_2, \dots, X_k)$ 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:

$$f(x; \theta_1, \theta_2, \dots, \theta_k) := P(X = x) = \sum_{i=1}^k \theta_i \mathbb{1}_{\{e_i\}}(x) = \begin{cases} \theta_1 & \text{if } x = e_1 := (1, 0, \dots, 0) \in \mathbb{R}^k \\ \theta_2 & \text{if } x = e_2 := (0, 1, \dots, 0) \in \mathbb{R}^k \\ \vdots \\ \theta_k & \text{if } x = e_k := (0, 0, \dots, 1) \in \mathbb{R}^k \\ 0 & \text{otherwise} \end{cases}$$

Of course, $\sum_{i=1}^k \theta_i = 1$.

When we add n IID de Moivre($\theta_1, \theta_2, \dots, \theta_k$) $R\vec{V}$ together, we get the $\text{Multinomial}(n, \theta_1, \theta_2, \dots, \theta_k)$ $R\vec{V}$ as defined below.

Model 17 ($\text{Multinomial}(n, \theta_1, \theta_2, \dots, \theta_k)$ $R\vec{V}$) We say that a $R\vec{V}$ $Y := (Y_1, Y_2, \dots, Y_k)$ obtained from the sum of n IID de Moivre($\theta_1, \theta_2, \dots, \theta_k$) $R\vec{V}$ 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:

$$f(y; n, \theta) := f(y; n, \theta_1, \theta_2, \dots, \theta_k) := P(Y = y; n, \theta_1, \theta_2, \dots, \theta_k) = \binom{n}{y_1, y_2, \dots, y_k} \prod_{i=1}^k \theta_i^{y_i},$$

where, the multinomial coefficient:

$$\binom{n}{y_1, y_2, \dots, y_k} := \frac{n!}{y_1! y_2! \cdots y_k!}.$$

Note that the marginal PMF of Y_j is $\text{Binomial}(n, \theta_j)$ for any $j = 1, 2, \dots, k$.

We can visualize the $\text{Multinomial}(n, \theta_1, \theta_2, \theta_3)$ process as a sum of n IID de Moivre($\theta_1, \theta_2, \theta_3$) $R\vec{V}$ s via a three dimensional extension of the Quincunx called the “Septcunx” and relate the number of paths that lead to a given trivariate sum (y_1, y_2, y_3) with $\sum_{i=1}^3 y_i = n$ as the multinomial coefficient $\frac{n!}{y_1! y_2! y_3!}$. In the Septcunx, 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, $\mathbb{Y} = \{(y_1, y_2, y_3) \in \mathbb{Z}_+^3 : \sum_{i=1}^3 y_i = n\}$. Once again, we can visualize that the sum of n IID de Moivre($\theta_1, \theta_2, \theta_3$) $R\vec{V}$ s constitute the $\text{Multinomial}(n, \theta_1, \theta_2, \theta_3)$ $R\vec{V}$.

Labwork 100 (Septcunx Sampler Demo – Sum of n IID de Moivre(1/3, 1/3, 1/3) \vec{RV} s)

Let us understand the Septcunx construction of the Multinomial($n, 1/3, 1/3, 1/3$) \vec{RV} X as the sum of n independent and identical de Moivre(1/3, 1/3, 13/) \vec{RV} s by calling the interactive visual cognitive tool as follows:

```
>> guiMultinomial
```

Multinomial distributions are at the very foundations of various machine learning algorithms, including, filtering junk email, learning from large knowledge-based resources like www, Wikipedia, word-net, etc.

Model 18 (Normal(μ, Σ) \vec{RV}) The univariate 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 matrix Σ . To begin, let

$$Z = \begin{pmatrix} Z_1 \\ Z_2 \\ \vdots \\ Z_m \end{pmatrix}$$

where, Z_1, Z_2, \dots, Z_m are jointly independent Normal(0, 1) RVs. Then the JPDF of Z is

$$f_Z(z) = f_{Z_1, Z_2, \dots, Z_m}(z_1, z_2, \dots, z_m) = \frac{1}{(2\pi)^{m/2}} \exp\left(-\frac{1}{2} \sum_{j=1}^m z_j^2\right) = \frac{1}{(2\pi)^{m/2}} \exp\left(-\frac{1}{2} z^T z\right)$$

We say that Z has a standard multivariate normal distribution and write $Z \sim \text{Normal}(0, I)$, where it is understood that 0 represents the vector of m zeros and I is the $m \times m$ identity matrix (with 1 along the diagonal entries and 0 on all off-diagonal entries).

More generally, a vector X has a multivariate normal distribution denoted by $X \sim \text{Normal}(\mu, \Sigma)$, if it has joint probability density function

$$f_X(x; \mu, \Sigma) = f_{X_1, X_2, \dots, X_m}(x_1, x_2, \dots, x_m; \mu, \Sigma) = \frac{1}{(2\pi)^{m/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu)\right)$$

where $|\Sigma|$ denotes the determinant of Σ , μ is a vector of length m and Σ is a $m \times m$ symmetric, positive definite matrix. Setting $\mu = 0$ and $\Sigma = I$ gives back the standard multivariate normal \vec{RV} .

When we have a non-zero mean vector

$$\mu = \begin{pmatrix} \mu_X \\ \mu_Y \end{pmatrix} = \begin{pmatrix} 6.49 \\ 5.07 \end{pmatrix}$$

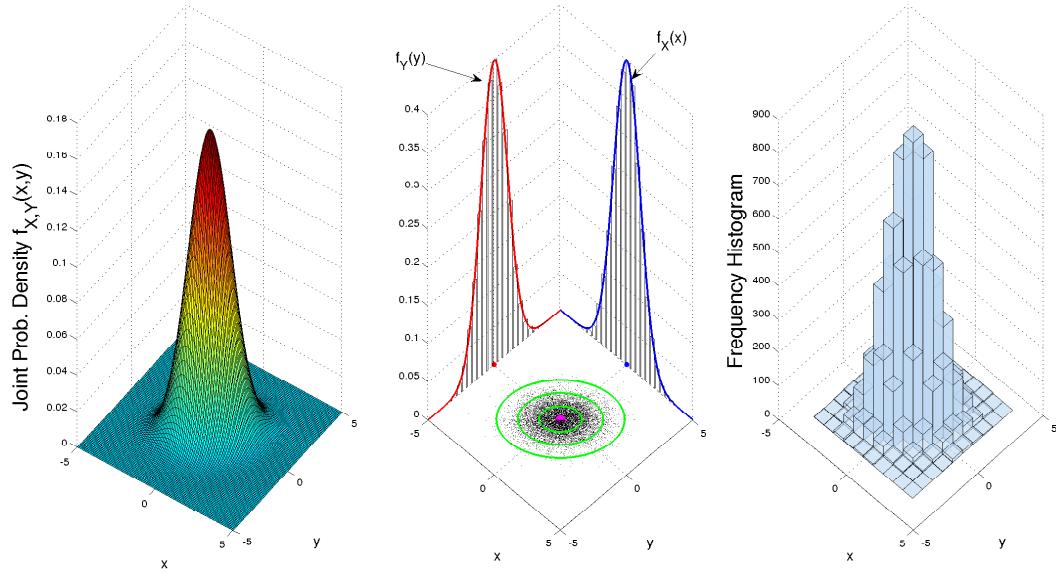
for the mean lengths and girths of cylindrical shafts from a manufacturing process with variance-covariance matrix

$$\Sigma = \begin{pmatrix} \text{Cov}(X, X) & \text{Cov}(X, Y) \\ \text{Cov}(Y, X) & \text{Cov}(Y, Y) \end{pmatrix} = \begin{pmatrix} V(X) & \text{Cov}(X, Y) \\ \text{Cov}(Y, X) & V(Y) \end{pmatrix} = \begin{pmatrix} 0.59 & 0.24 \\ 0.24 & 0.26 \end{pmatrix}$$

then the Normal(μ, Σ) \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:

Figure 3.22: JPDF, Marginal PDFs and Frequency Histogram of Bivariate Standard Normal RV.



```
>> 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 ??.

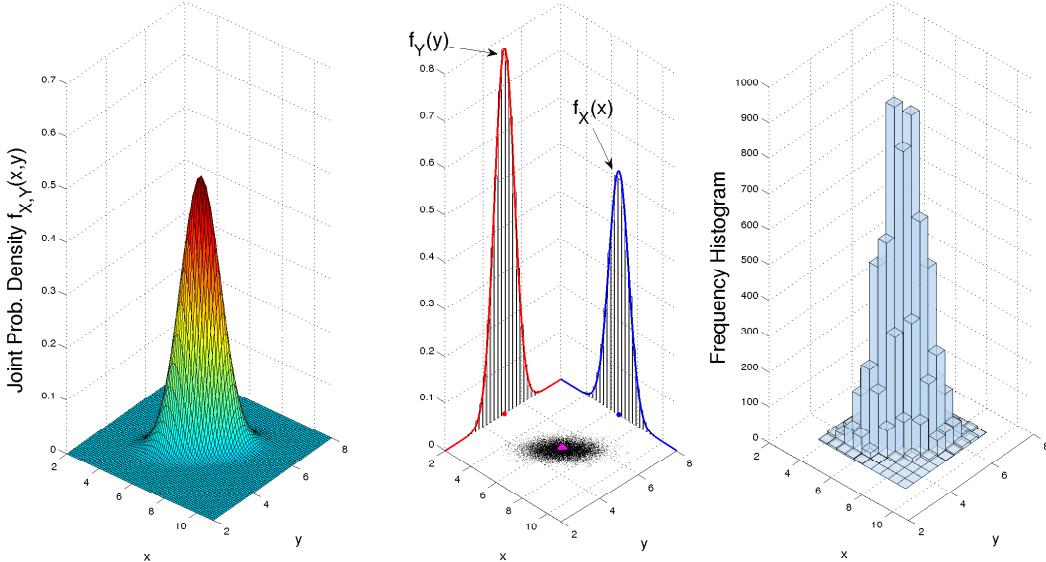
3.11 Exercises in Multivariate Random Variables

Ex. 3.38 — Find the probability that none of the three bulbs in a traffic signal, 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.

Figure 3.23: JPDF, Marginal PDFs and Frequency Histogram of a Bivariate Normal RV for lengths of girths of cylindrical shafts in a manufacturing process (in cm).



Ex. 3.39 — Let (X, Y) be a continuous RV with joint probability density function (JPDF)

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

Find the following:

1. the normalizing constant a which will ensure $P(\Omega) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f_{X,Y}(x,y) dx dy = 1$
2. $f_X(x) = \int_{-\infty}^{\infty} f_{X,Y}(x,y) dy$ called the marginal probability density function (MPDF) of X
3. $f_Y(y) = \int_{-\infty}^{\infty} f_{X,Y}(x,y) dx$ called the marginal probability density function (MPDF) of Y
4. Check if $f_X(x)f_Y(y) = f_{X,Y}(x,y)$ for every (x,y) and decide whether X and Y are independent random variables. Hint: X and Y are said to be independent if $f_X(x)f_Y(y) = f_{X,Y}(x,y)$ for every (x,y) .
5. $F_{X,Y}(x,y)$, the joint cumulative distribution function (JCDF) of (X, Y) for any $(x,y) \in (0, 1) \times (0, 1)$
6. the probability that $X > 0.5$ and $Y < 0.6$, i.e., $P(X > 0.5, Y < 0.6)$
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. $\text{Cov}(X, Y) = E(XY) - E(X)E(Y)$, 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 (RV) 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 RV (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 $\text{Cov}(X, Y) = E(XY) - E(X)E(Y)$]

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 $(x,y) \in \mathbb{R}^2$.]

Ex. 3.45 — Let (X, Y) be a continuous RV with joint probability density function (JPDF):

$$f_{X,Y}(x,y) = \begin{cases} e^{-x} & \text{if } x \in [0, \infty) \text{ and } y \in [2, 3] \\ 0 & \text{otherwise} . \end{cases}$$

Are X and Y independent?

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 JPDF of these variables is

$$f_{X_1, X_2, X_3, X_4}(x_1, x_2, x_3, x_4) = \begin{cases} 9 \times 10^{-12} e^{-0.001x_1 - 0.002x_2 - 0.0015x_3 - 0.003x_4} & \text{if } x_1 \geq 0, x_2 \geq 0, x_3 \geq 0, x_4 \geq 0 \\ 0 & \text{otherwise .} \end{cases}$$

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.]

Ex. 3.47 — Suppose the RVs Y_1, Y_2 and Y_3 represent the thickness in micrometers of a substrate, an active layer, and a coating layer of a chemical product. Assume Y_1, Y_2 and Y_3 are $\text{Normal}(10000, 250^2)$, $\text{Normal}(1000, 20^2)$ and $\text{Normal}(80, 4^2)$ RVs, respectively. Further suppose that they are independent. The required specifications for the thickness of the substrate, active layer and coating layer are $[9500, 10500]$, $[950, 1050]$ and $[75, 85]$, respectively. What proportion of chemical products meets all thickness specifications? [Hint: this is just $P(9500 < Y_1 < 10500, 950 < Y_2 < 1050, 75 < Y_3 < 85)$] Which one of the three thicknesses has the least probability of meeting specifications?

Ex. 3.48 — Soft drink cans are filled by an automated filling machine. Assume the fill volumes of the cans are independent $\text{Normal}(12.1, 0.01)$ RVs. What is the probability that the average volume of ten cans selected from this process is less than 12.01 fluid ounces?

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 6 of the bits received are *excellent*, 2 are *good*, 2 are *fair* and none are *poor* under the assumption that the classification of bits are independent events and that the probabilities of each bit being *excellent*, *good*, *fair* and *poor* are 0.6, 0.3, 0.08 and 0.02, respectively. [Hint: Think of $\text{Multinomial}(n = 10, \theta_1 = 0.6, \theta_2 = 0.3, \theta_3 = 0.08, \theta_4 = 0.02)$ as a model for bit classification in this digital channel.]

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 $\iota = \sqrt{-1}$. The function $\varphi_X(t) : \mathbb{R} \rightarrow \mathbb{C}$ defined by

$$\boxed{\varphi_X(t) := E(\exp(\iota tX)) = \begin{cases} \sum_x \exp(\iota tx) f_X(x) & \text{if } X \text{ is discrete RV} \\ \int_{-\infty}^{\infty} \exp(\iota tx) f_X(x) dx & \text{if } X \text{ is continuous RV} \end{cases}} \quad (3.70)$$

is called the **characteristic function** of X .

NOTE: $\varphi_X(t)$ exists for any $t \in \mathbb{R}$, because

$$\begin{aligned} \varphi_X(t) &= E(\exp(\iota tX)) \\ &= E(\cos(tX) + \iota \sin(tX)) \\ &= E(\cos(tX)) + \iota E(\sin(tX)) \end{aligned}$$

and the last two expected values are well-defined, because the sine and cosine functions are bounded by $[-1, 1]$.

For a continuous RV, $\int_{-\infty}^{\infty} \exp(-\iota tx) f_X(x) dx$ 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

$$E(X^k) = \frac{1}{\iota^k} \left[\frac{d^k \varphi_X(t)}{dt^k} \right]_{t=0} .$$

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 $t = 0$.

Proof: 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:

$$\frac{d^k \varphi_X(t)}{dt^k} = \frac{d^k}{dt^k} E(\exp(itX)) = E\left(\frac{d^k}{dt^k} \exp(itX)\right) = E\left((it)^k \exp(itX)\right) = it^k E\left(X^k \exp(itX)\right)$$

The RHS evaluated at $t = 0$ is

$$\left[\frac{d^k \varphi_X(t)}{dt^k}\right]_{t=0} = \left[it^k E\left(X^k \exp(itX)\right)\right]_{t=0} = it^k E\left(X^k\right)$$

This completes the sketch of the proof.

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.

Proposition 49 (Moments from CF.) Let X be a random variable and $\varphi_X(t)$ be its CF. If $\varphi_X(t)$ is k times differentiable at the point $t = 0$, then

1. if k is even, the n -th moment of X exists and is finite for any $0 \leq n \leq k$;
2. if k is odd, the n -th moment of X exists and is finite for any $0 \leq n \leq k - 1$.

In both cases,

$$E(X^k) = \frac{1}{i^k} \left[\frac{d^k \varphi_X(t)}{dt^k} \right]_{t=0}. \quad (3.71)$$

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 $t = 0$.

Proof: For proof see e.g., Ushakov, N. G. (1999) Selected topics in characteristic functions, VSP (p. 39).

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 $V(X) = E(X^2) - (E(X))^2$.

Solution:

Part 1

Recall the PMF for this discrete RV with parameter $\theta \in (0, 1)$ is

$$f_X(x; \theta) = \begin{cases} \theta & \text{if } x = 1 \\ 1 - \theta & \text{if } x = 0 \\ 0 & \text{otherwise.} \end{cases}$$

Let's first find the CF of X

$$\begin{aligned} \varphi_X(t) &= E(\exp(itX)) = \sum_x \exp(itx) f_X(x; \theta) \quad \text{By Defn. in Equation (3.70)} \\ &= \exp(it \times 0)(1 - \theta) + \exp(it \times 1)\theta = \exp(0)(1 - \theta) + \exp(it)\theta = 1 - \theta + \theta \exp(it) \end{aligned}$$

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:

Let us differentiate the CF to get moments using Equation (3.71) (CF has to be once and twice differentiable at $t = 0$ to get the first and second moments).

$$\begin{aligned}\frac{d}{dt}\varphi_X(t) &= \frac{d}{dt}\left(\frac{\lambda}{\lambda - it}\right) = \lambda\left(-1 \times (\lambda - it)^{-2} \times \frac{d}{dt}(\lambda - it)\right) \\ &= \lambda\left(\frac{-1}{(\lambda - it)^2} \times (-i)\right) = \frac{\lambda i}{(\lambda - it)^2}\end{aligned}$$

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:

$$E(X) = \frac{1}{i} \left[\frac{d}{dt}\varphi_X(t) \right]_{t=0} = \frac{1}{i} \left[\frac{\lambda i}{(\lambda - it)^2} \right]_{t=0} = \frac{1}{i} \left(\frac{\lambda i}{\lambda^2} \right) = \frac{1}{i} \left(\frac{i}{\lambda} \right) = \frac{1}{\lambda}$$

Let's pause and see if this makes sense.... Yes, because the expected value of Exponential(λ) RV is indeed $1/\lambda$ (recall from when we introduced this RV).

Similarly from Equation (3.71) we can get $E(X^2)$ as follows:

$$\begin{aligned}E(X^2) &= \frac{1}{i^2} \left[\frac{d^2}{dt^2}\varphi_X(t) \right]_{t=0} = \frac{1}{i^2} \left[\frac{d}{dt} \frac{d}{dt}\varphi_X(t) \right]_{t=0} = \frac{1}{i^2} \left[\frac{d}{dt} \frac{\lambda i}{(\lambda - it)^2} \right]_{t=0} \\ &= \frac{1}{i^2} \left[\lambda i \times \frac{d}{dt}(\lambda - it)^{-2} \right]_{t=0} = \frac{1}{i^2} \left[\lambda i \left(-2(\lambda - it)^{-3} \frac{d}{dt}(\lambda - it) \right) \right]_{t=0} \\ &= \frac{1}{i^2} \left[\lambda i (-2(\lambda - it)^{-3} \times (-i)) \right]_{t=0} = \frac{1}{i^2} \left[\frac{2\lambda i^2}{(\lambda - it)^3} \right]_{t=0} = \frac{1}{i^2} \left(\frac{2\lambda i^2}{\lambda^3} \right) = \frac{2}{\lambda^2}.\end{aligned}$$

Finally, from the first and second moments we can get the variance as follows:

$$V(X) = E(X^2) - (E(X))^2 = \frac{2}{\lambda^2} - \left(\frac{1}{\lambda}\right)^2 = \frac{2}{\lambda^2} - \frac{1}{\lambda^2} = \frac{2-1}{\lambda^2} = \frac{1}{\lambda^2}.$$

Let's check that this is what we had as variance for the Exponential(λ) RV when we first introduced it and directly computed using integrals for definition of expectation.

Characteristic functions can be used to characterize the distribution of a random variable.

Two RVs X and Y have the same DFs , i.e., $F_X(x) = F_Y(x)$ for all $x \in \mathbb{R}$, if and only if they have the same characteristic functions, i.e. $\varphi_X(t) = \varphi_Y(t)$ for all $t \in \mathbb{R}$ (for proof see Resnick, S. I. (1999) A Probability Path, Birkhauser).

Thus, if we can show that two RVs have the same CF then we know they are the same. This can be much more challenging or impossible to do directly with their DFs.

Let Z be Normal($0, 1$), the standard normal RV. We can find the CF for Z using couple of tricks as follows

$$\begin{aligned}\varphi_Z(t) &= E(e^{itZ}) \\ &= \int_{-\infty}^{\infty} e^{itz} f_Z(z) dz = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{itz} e^{-z^2/2} dz = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{itz - z^2/2} dz \\ &= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{-(t^2 + (z-it)^2)/2} dz = e^{-t^2/2} \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{-(z-it)^2/2} dz \\ &= e^{-t^2/2} \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{-y^2/2} dy \quad \text{substituting } y = z - it, dy = dz \\ &= e^{-t^2/2} \frac{1}{\sqrt{2\pi}} \sqrt{2\pi} \quad \text{using the normalizing constant in PDF of Normal}(0, 1) \text{ RV} \\ &= e^{-t^2/2}\end{aligned}$$

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

$$\boxed{\varphi_Y(t) = \exp(\imath at)\varphi_X(bt)} \quad (3.73)$$

Proof: This is easy to prove using the definition of CF as follows:

$$\begin{aligned} \varphi_Y(t) &= E(\exp(\imath tY)) = E(\exp(\imath t(a + bX))) = E(\exp(\imath ta + \imath tbX)) \\ &= E(\exp(\imath ta)\exp(\imath tbX)) = \exp(\imath ta)E(\exp(\imath tbX)) = \exp(\imath ta)\varphi_X(bt) \end{aligned}$$

Example 103 Let Y be a $\text{Normal}(\mu, \sigma^2)$ RV. Recall that Y is a linear transformation of Z , i.e., $Y = \mu + \sigma Z$ where Z is a $\text{Normal}(0, 1)$ RV. Using Equations (3.72) and (3.73) find the CF of Y .

Solution:

$$\begin{aligned} \varphi_Y(t) &= \exp(\imath \mu t)\varphi_Z(\sigma t), \quad \text{since } Y = \mu + \sigma Z \\ &= e^{\imath \mu t} e^{(-\sigma^2 t^2)/2}, \quad \text{since } \varphi_Z(t) = e^{-t^2/2} \\ &= e^{\imath \mu t - (\sigma^2 t^2)/2} \end{aligned}$$

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 $(1 - \theta + \theta e^t)$ then $Y = \sum_{i=1}^n X_i$ is the Binomial(n, θ) RV and by Eqn. (3.74) with $a_1 = a_2 = \cdots = 1$, we get

$$\varphi_Y(t) = \varphi_{X_1}(t) \times \varphi_{X_2}(t) \cdots \varphi_{X_n}(t) = \prod_{i=1}^n \varphi_{X_i}(t) = \prod_{i=1}^n (1 - \theta + \theta e^t) = (1 - \theta + \theta e^t)^n.$$

Example 105 Let Z_1 and Z_2 be independent $\text{Normal}(0, 1)$ RVs.

1. Use Eqn. (3.74) to find the CF of $Z_1 + Z_2$.
2. From the CF of $Z_1 + Z_2$ identify what RV it is.
3. Use Eqn. (3.74) to find the CF of $2Z_1$.
4. From the CF of $2Z_1$ identify what RV it is.
5. Try to understand the difference between the distributions of $Z_1 + Z_2$ and $2Z_1$ inspite of Z_1 and Z_2 having the same distribution.

Hint: from lectures we know that $\varphi_X(t) = e^{\mu t - (\sigma^2 t^2)/2}$ for a $\text{Normal}(\mu, \sigma^2)$ RV X .

Solution:

1. By Eqn. (3.74) we just multiply the characteristic functions of Z_1 and Z_2 , both of which are $e^{-t^2/2}$,

$$\varphi_{Z_1+Z_2}(t) = \varphi_{Z_1}(t) \times \varphi_{Z_2}(t) = e^{-t^2/2} \times e^{-t^2/2} = e^{-2t^2/2} = e^{-t^2} .$$

2. The CF of $Z_1 + Z_2$ is that of the $\text{Normal}(\mu, \sigma^2)$ RV with $\mu = 0$ and $\sigma^2 = 2$. Thus $Z_1 + Z_2$ is the $\text{Normal}(0, 2)$ RV with mean parameter $\mu = 0$ and variance parameter $\sigma^2 = 2$.
3. We can again use Eqn. (3.74) to find the CF of $2Z_1$ as follows

$$\varphi_{2Z_1} = \varphi_{Z_1}(2t) = e^{-2^2 t^2/2} .$$

4. The CF of $2Z_1$ is that of the $\text{Normal}(\mu, \sigma^2)$ RV with $\mu = 0$ and $\sigma^2 = 2^2 = 4$. Thus $2Z_1$ is the $\text{Normal}(0, 4)$ RV with mean parameter $\mu = 0$ and variance parameter $\sigma^2 = 4$.
5. $2Z_1$ has a bigger variance from multiplying the standard normal RV by 2 while $Z_1 + Z_2$ has a smaller variance from adding two independent standard normal RVs. Thus, the result of adding the same RV twice does not have the same distribution as that of multiplying it by 2. In other words $2 \times Z$ is not equal to $Z + Z$ in terms of its probability distribution!

3.12.2 Moment Generating Function

Moment generating functions are special cases of characteristic functions and we won't be explicitly using them here as it is more convenient to work in the complex plane.

3.13 Exercises in Characteristic Functions

Ex. 3.50 — Let X be a discrete random variable (RV) with probability mass function (PMF)

$$f_X(x) = \begin{cases} \frac{1}{3} & \text{if } x = 0 \\ \frac{1}{3} & \text{if } x = 1 \\ \frac{1}{3} & \text{if } x = 2 \\ 0 & \text{otherwise} . \end{cases}$$

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

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

1. Find the CF of X . (Hint: the sum of the infinite geometric series $\sum_{x=0}^{\infty} ar^x = \frac{a}{1-r}$.)
2. 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 $e^\alpha = \sum_{x=0}^{\infty} \frac{\alpha^x}{x!}$.

1. Find the CF of X .
2. 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.

Chapter 4

Statistics

4.1 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} :

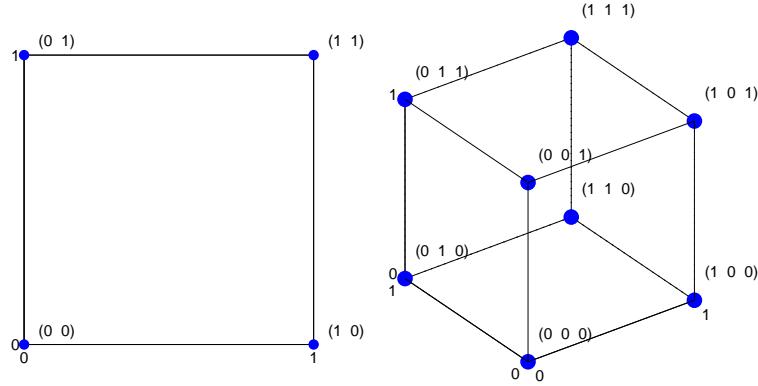
$$X(\omega) : \Omega \rightarrow \mathbb{X} .$$

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 .

Figure 4.1: Sample Space, Random Variable, Realisation, Data, and Data Space.

Example 106 (Tossing a coin n times) For some given parameter $\theta \in \Theta := [0, 1]$, consider n IID Bernoulli(θ) trials, i.e. $X_1, X_2, \dots, X_n \stackrel{\text{IID}}{\sim} \text{Bernoulli}(\theta)$. 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 hyper-cube, measures the outcomes of this experiment. A particular realisation of X , upon performance of this experi-

Figure 4.2: Data Spaces $\mathbb{X} = \{0, 1\}^2$ and $\mathbb{X} = \{0, 1\}^3$ for two and three Bernoulli trials, respectively.



ment, is the observation, 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 $(x_1, x_2, \dots, x_{n-1}, x_n) = (0, 0, \dots, 0, 1)$.

Definition 51 (Statistic) A **statistic** T is any function of the data:

$$T(x) : \mathbb{X} \rightarrow \mathbb{T} .$$

Thus, a statistic T is also an RV that takes values in the space \mathbb{T} . When $x \in \mathbb{X}$ is the realisation of an experiment, we let $T(x) = t$ denote the corresponding realisation of the statistic T . Sometimes we use $T_n(X)$ and \mathbb{T}_n to emphasise that X is an n -dimensional random vector, i.e. $\mathbb{X} \subset \mathbb{R}^n$

Classwork 107 (Is data a statistic?) Is the RV X , for which the realisation is the observed 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 $E(X_1)$ and $V(X_1)$, 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:

$$T_n((X_1, X_2, \dots, X_n)) = \bar{X}_n((X_1, X_2, \dots, X_n)) := \frac{1}{n} \sum_{i=1}^n X_i . \quad (4.1)$$

For brevity, we write

$$\bar{X}_n((X_1, X_2, \dots, X_n)) \quad \text{as} \quad \bar{X}_n ,$$

and its realisation

$$\bar{X}_n((x_1, x_2, \dots, x_n)) \quad \text{as} \quad \bar{x}_n .$$

Note that the expectation and variance of \bar{X}_n are:

$$\begin{aligned} E(\bar{X}_n) &= E\left(\frac{1}{n} \sum_{i=1}^n X_i\right) && [\text{by definition (4.1)}] \\ &= \frac{1}{n} \sum_{i=1}^n E(X_i) && [\text{by property (3.50)}] \end{aligned}$$

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 $E(X_1)$, then:

$$E(\bar{X}_n) = \frac{1}{n} \sum_{i=1}^n E(X_i) = \frac{1}{n} \sum_{i=1}^n E(X_1) = \frac{1}{n} n E(X_1) = E(X_1) . \quad (4.2)$$

Similarly, we can show that:

$$\begin{aligned} V(\bar{X}_n) &= V\left(\frac{1}{n} \sum_{i=1}^n X_i\right) && [\text{by definition (4.1)}] \\ &= \left(\frac{1}{n}\right)^2 V\left(\sum_{i=1}^n X_i\right) && [\text{by property (3.49)}] \end{aligned}$$

Furthermore, if the original sequence of RVs X_1, X_2, \dots is **independently** distributed then:

$$V(\bar{X}_n) = \left(\frac{1}{n}\right)^2 V\left(\sum_{i=1}^n X_i\right) = \frac{1}{n^2} \sum_{i=1}^n V(X_i) \quad [\text{by property (3.51)}]$$

Finally, if the original sequence of RVs X_1, X_2, \dots is **independently and identically** distributed with the same variance ($V(X_1)$ by convention) then:

$$V(\bar{X}_n) = \frac{1}{n^2} \sum_{i=1}^n V(X_i) = \frac{1}{n^2} \sum_{i=1}^n V(X_1) = \frac{1}{n^2} n V(X_1) = \frac{1}{n} V(X_1) . \quad (4.3)$$

Labwork 108 (Sample mean) After initializing the fundamental sampler, 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 `XsFromUni01Twstr101`.

```
>> rand('twister',101); % initialise the fundamental Uniform(0,1) sampler
>> XsFromUni01Twstr101=rand(1,5); % simulate n=5 IID samples from Uniform(0,1) RV
>> SampleMean=mean(XsFromUni01Twstr101);% find sample mean
>> 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
>> disp(SampleMean); % The Sample mean is :
    0.3945
```

We can thus use `mean` to obtain the sample mean \bar{x}_n of n sample points x_1, x_2, \dots, x_n .

We may also obtain the sample mean using the `sum` function and a division by sample size:

```
>> sum(XsFromUni01Twstr101) % take the sum of the elements of the XsFromUni01Twstr101 array
ans =      1.9723
>> sum(XsFromUni01Twstr101) / 5 % divide the sum by the sample size 5
ans =      0.3945
```

We can also obtain the sample mean via matrix product or multiplication as follows:

```
>> size(XsFromUni01Twstr101) % size(SomeArray) gives the size or dimensions of the arrar SomeArray
ans =      1      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
>> XsFromUni01Twstr101 * ones(5,1) % multiplying an 1 X 5 matrix with a 5 X 1 matrix of Ones
ans =    1.9723
>> XsFromUni01Twstr101 * ( ones(5,1) * 1/5) % multiplying an 1 X 5 matrix with a 5 X 1 matrix of 1/5 's
ans =    0.3945
```

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 simply the sample variance :

$$T_n((X_1, X_2, \dots, X_n)) = S_n^2((X_1, X_2, \dots, X_n)) := \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2 . \quad (4.4)$$

For brevity, we write $S_n^2((X_1, X_2, \dots, X_n))$ as S_n^2 and its realisation $S_n^2((x_1, x_2, \dots, x_n))$ as s_n^2 .

Sample standard deviation is simply the square root of sample variance:

$$S_n((X_1, X_2, \dots, X_n)) = \sqrt{S_n^2((X_1, X_2, \dots, X_n))} \quad (4.5)$$

For brevity, we write $S_n((X_1, X_2, \dots, X_n))$ as S_n and its realisation $S_n((x_1, x_2, \dots, x_n))$ as s_n .

Once again, if $X_1, X_2, \dots, X_n \stackrel{\text{IID}}{\sim} X_1$, the expectation of the sample variance is:

$$\mathbb{E}(S_n^2) = \mathbb{V}(X_1) .$$

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:

$$X_{([n])}((X_1, X_2, \dots, X_n)) := (X_{(1)}, X_{(2)}, \dots, X_{(n)}) , \text{ such that, } X_{(1)} \leq X_{(2)} \leq \dots \leq X_{(n)} . \quad (4.6)$$

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:

$$F^{[-1]}(q) := \inf \{x : F(x) > q\}, \quad \text{for some } q \in [0, 1] . \quad (4.7)$$

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:

1. The **mean** of RV $X \sim F$ is a function of the DF F :

$$T(F) = E(X) = \int x dF(x) .$$

2. The **variance** of RV $X \sim F$ is a function of the DF F :

$$T(F) = E(X - E(X))^2 = \int (x - E(X))^2 dF(x) .$$

3. The **value of DF at a given** $x \in \mathbb{R}$ of RV $X \sim F$ is also a function of DF F :

$$T(F) = F(x) .$$

Other functionals of F that depend on the quantile function $F^{[-1]}$ are:

1. The q^{th} **quantile** of RV $X \sim F$:

$$T(F) = F^{[-1]}(q) \quad \text{where } q \in [0, 1] .$$

2. The **first quartile** or the 0.25^{th} **quantile** of the RV $X \sim F$:

$$T(F) = F^{[-1]}(0.25) .$$

3. The **median** or the **second quartile** or the 0.50^{th} **quantile** of the RV $X \sim F$:

$$T(F) = F^{[-1]}(0.50) .$$

4. The **third quartile** or the 0.75^{th} **quantile** of the RV $X \sim F$:

$$T(F) = F^{[-1]}(0.75) .$$

Definition 56 (Empirical Distribution Function (EDF or ECDF)) Suppose we have n IID RVs, $X_1, X_2, \dots, X_n \stackrel{\text{IID}}{\sim} F$, where F is a DF from the set of all DFs over the real line. Then, the n -sample empirical distribution function (EDF or ECDF) is the discrete distribution function \hat{F}_n that puts a probability mass of $1/n$ at each sample or data point x_i :

$$\hat{F}_n(x) = \frac{\sum_{i=1}^n \mathbb{1}(X_i \leq x)}{n}, \quad \text{where} \quad \mathbb{1}(X_i \leq x) := \begin{cases} 1 & \text{if } x_i \leq x \\ 0 & \text{if } x_i > x \end{cases} \quad (4.8)$$

Labwork 111 (Plot of empirical CDF) Let us plot the ECDF for the five samples drawn from the $\text{Uniform}(0, 1)$ RV in Labwork 108 using the MATLAB function ECDF. Let us superimpose the samples and the true DF as depicted in Figure 4.3 with the following script:

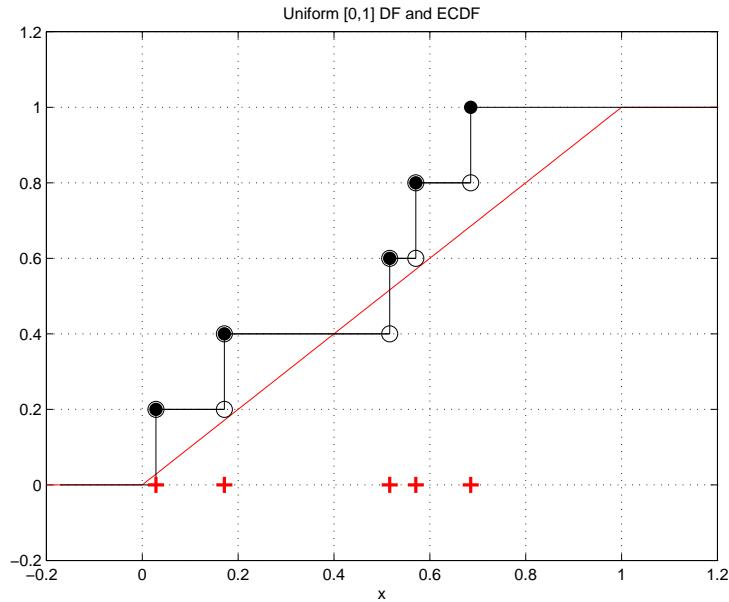
```
plotunifecdf.m
xs = -1:0.01:2; % vector xs from -1 to 2 with increment .05 for x values
% get the [0,1] uniform DF or cdf of xs in vector cdf
cdf=zeros(size(xs));% initialise cdf as zero
indices = find(xs>=1); cdf(indices) = 1; % set cdf as 1 when xs >= 1
indices = find(xs>=0 & xs<=1); cdf(indices)=xs(indices); % cdf=xs when 0 <= xs <= 1
plot(xs,cdf,'r') % plot the DF
hold on; title('Uniform [0,1] DF and ECDF'); xlabel('x'); axis([-0.2 1.2 -0.2 1.2])
x=[0.5164, 0.5707, 0.0285, 0.1715, 0.6853]; % five samples
plot(x,zeros(1,5),'r+','LineWidth',2,'MarkerSize',10)% plot the data as red + marks
hold on; grid on; % turn on grid
ECDF(x,.1,.6);% ECDF (type help ECDF) plot is extended to left and right by .2 and .4, respectively.
```

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 (4.8). The q^{th} **sample quantile** is defined as the statistic (statistical functional):

$$T(\hat{F}_n) = \hat{F}_n^{[-1]}(q) := \inf \{x : \hat{F}_n^{[-1]}(x) \geq q\} . \quad (4.9)$$

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.

Figure 4.3: Plot of the DF of Uniform(0, 1), five IID samples from it, and the ECDF \widehat{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.



Algorithm 1 q^{th} Sample Quantile of Order Statistics

1: *input:*

1. q in the q^{th} sample quantile, i.e. the argument q of $\widehat{F}_n^{[-1]}(q)$,
 2. order statistic $(x_{(1)}, x_{(2)}, \dots, x_{(n)})$, i.e. the sorted (x_1, x_2, \dots, x_n) , where $n > 0$.
- 2: *output:* $\widehat{F}_n^{[-1]}(q)$, the q^{th} sample quantile
- 3: $i \leftarrow \lfloor (n-1)q \rfloor$
 - 4: $\delta \leftarrow (n-1)q - i$
 - 5: **if** $i = n-1$ **then**
 - 6: $\widehat{F}_n^{[-1]}(q) \leftarrow x_{(i+1)}$
 - 7: **else**
 - 8: $\widehat{F}_n^{[-1]}(q) \leftarrow (1-\delta)x_{(i+1)} + \delta x_{(i+2)}$
 - 9: **end if**
 - 10: *return:* $\widehat{F}_n^{[-1]}(q)$
-

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`

```
>> disp(SortedXsFromUni01Twstr101)
    0.0285    0.1715    0.5164    0.5707    0.6853
>> rand('twister',420);
>> SortedXs=sort(rand(1,7));
>> disp(SortedXs)
    0.1089    0.2670    0.3156    0.3525    0.4530    0.6297    0.8682
>> for q=[0, 0.25, 0.5, 0.75, 1.0]
    disp([q, qthSampleQuantile(q,SortedXsFromUni01Twstr101) ...
            qthSampleQuantile(q,SortedXs)])
end
    0    0.0285    0.1089
    0.2500    0.1715    0.2913
    0.5000    0.5164    0.3525
    0.7500    0.5707    0.5414
    1.0000    0.6853    0.8682
```

4.1.1 Univariate Data

A **histogram** is a graphical representation of the frequency with which elements of a data array:

$$x = (x_1, x_2, \dots, x_n),$$

of real numbers fall within each of the m intervals or **bins** of some **interval partition**:

$$b := (b_1, b_2, \dots, b_m) := ([\underline{b}_1, \bar{b}_1], [\underline{b}_2, \bar{b}_2], \dots, [\underline{b}_m, \bar{b}_m])$$

of the **data range** of x given by the closed interval:

$$\mathcal{R}(x) := [\min\{x_1, x_2, \dots, x_n\}, \max\{x_1, x_2, \dots, x_n\}].$$

Elements of this partition b are called bins, their mid-points are called **bin centres**:

$$c := (c_1, c_2, \dots, c_m) := ((\underline{b}_1 + \bar{b}_1)/2, (\underline{b}_2 + \bar{b}_2)/2, \dots, (\underline{b}_m + \bar{b}_m)/2)$$

and their overlapping boundaries, i.e. $\bar{b}_i = \underline{b}_{i+1}$ for $1 \leq i < m$, are called **bin edges**:

$$d := (d_1, d_2, \dots, d_{m+1}) := (\underline{b}_1, \bar{b}_1, \dots, \underline{b}_{m-1}, \bar{b}_m, \bar{b}_m).$$

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:

$$((b_1, h_1), (b_2, h_2), \dots, (b_m, h_m)) , \text{ where } h_k := \#\{x_i : x_i \in b_k\} ,$$

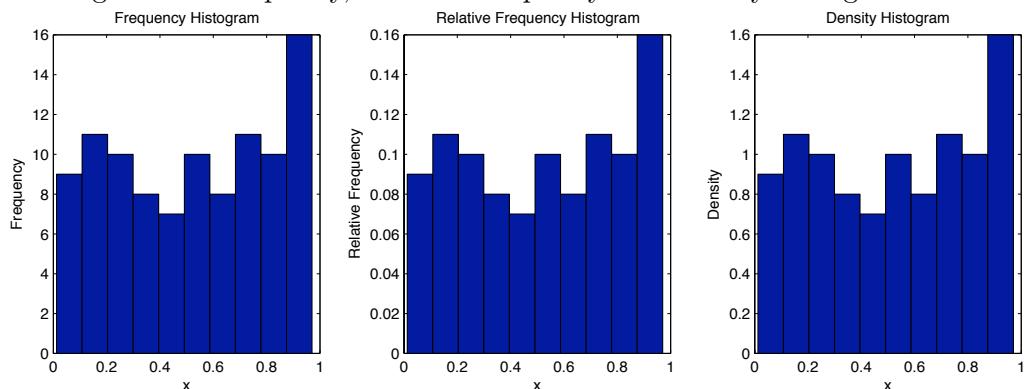
a **relative frequency histogram** is the histogram:

$$((b_1, h_1), (b_2, h_2), \dots, (b_m, h_m)) , \text{ where } h_k := n^{-1} \#\{x_i : x_i \in b_k\} ,$$

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 4.4: 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 4.4 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

```

Try making a density histogram with 1000 samples from `rand` with 15 bins. You can specify the number of bins by adding an extra argument to `hist`, for e.g. `[Fs, Cs] = hist(x,15)` will produce 15 bins of equal width over the data range $\mathcal{R}(x)$.

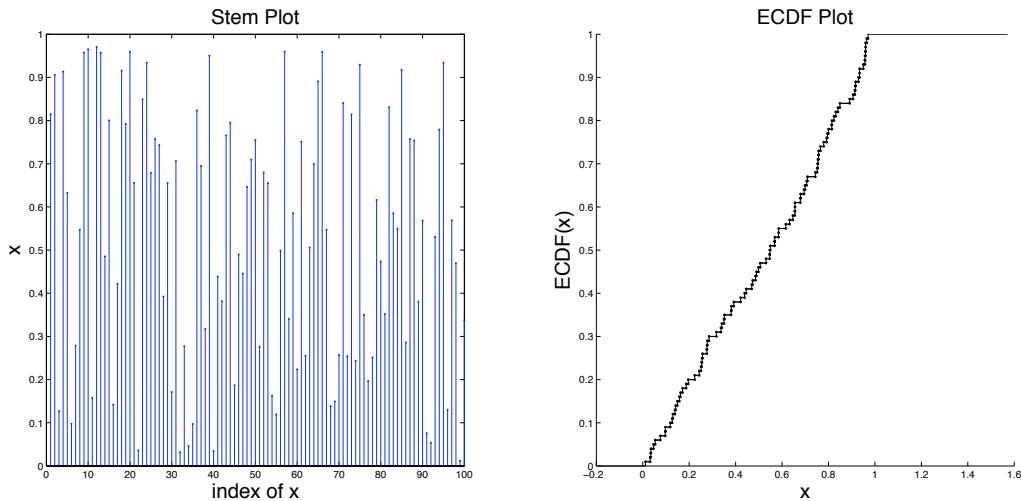
Labwork 114 (Stem plots and ECDF plots for univariate data) We can also visualise the 100 data points in the array `x` using stem plot and ECDF plot as shown in Figure 4.5 as follows:

```

>> rand('twister',5489);
>> x=rand(1,100); % produce 100 samples with rand
>> stem(x,'.') % make a stem plot of the 100 data points in x (the option '.' gives solid circles for x)
>>% ECDF (type help ECDF) plot is extended to left and right by .2 and .6, respectively
>>% (second parameter 6 makes the dots in the plot smaller).
>> ECDF(x,.2,.6);

```

Figure 4.5: Frequency, Relative Frequency and Density Histograms



We can also visually summarise univariate data using the **box plot** or **box-whisker plot** available in the Stats Toolbox of MATLAB. These family of plots display a set of sample quantiles, typically they include, the median, the first and third quartiles and the minimum and maximum values of our data array x .

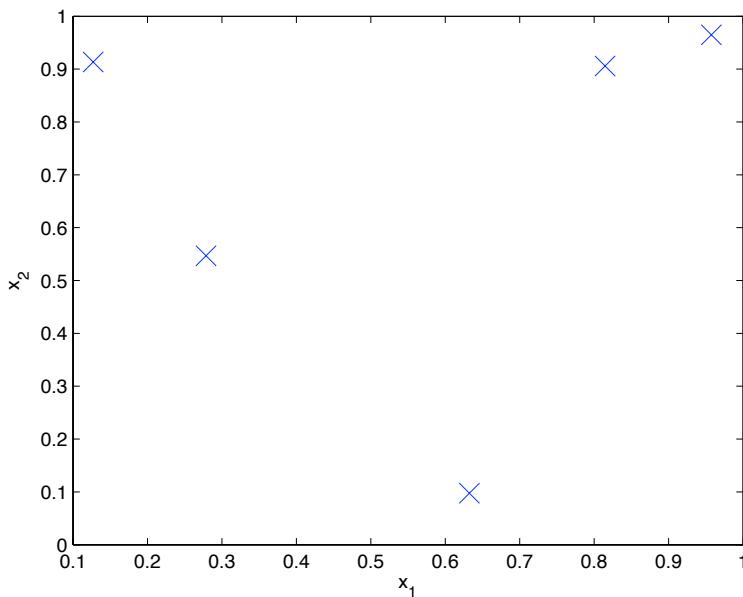
4.1.2 Bivariate Data

By bivariate data array x we mean a $2 \times n$ matrix of real numbers or equivalently n ordered pairs of points $(x_{1,i}, x_{2,i})$ as $i = 1, 2, \dots, n$. The most elementary visualisation 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 of 5 ordered pairs sampled uniformly at random over the unit square $[0, 1] \times [0, 1]$. We can make 2D scatter plot as shown in Figure 4.6 as follows:

```
>> rand('twister',5489);
>> x=rand(2,5)% create a sequence of 5 ordered pairs uniformly from unit square [0,1]X[0,1]
x =
    0.8147    0.1270    0.6324    0.2785    0.9575
    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 4.6: 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.

4.1.3 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 4.7 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,
```

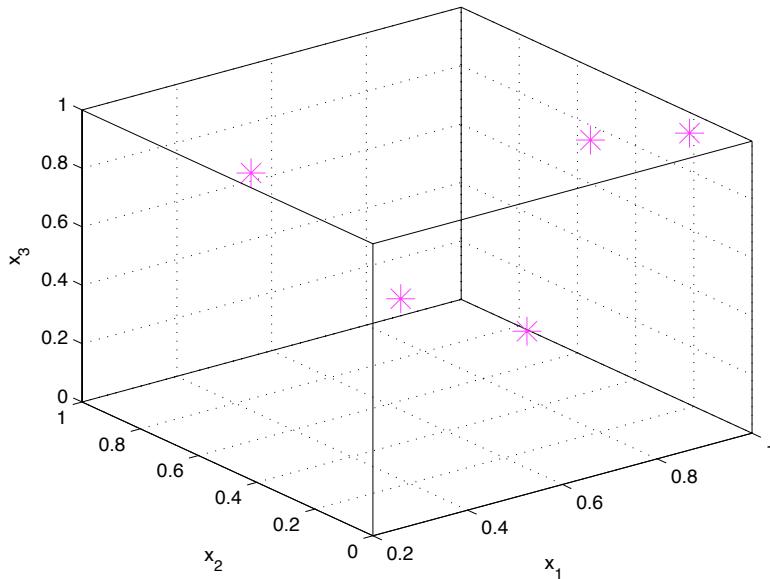
```

>>% 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?

Figure 4.7: 3D Scatter Plot



There are several other techniques for visualising trivariate data, including, iso-surface plots, moving surface or heat plots, and you will encounter some of them in the future.

4.1.4 Multivariate Data

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 $\{(i, j) \in \{1, 2, \dots, d\}^2 : i \neq j\}$ 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**.

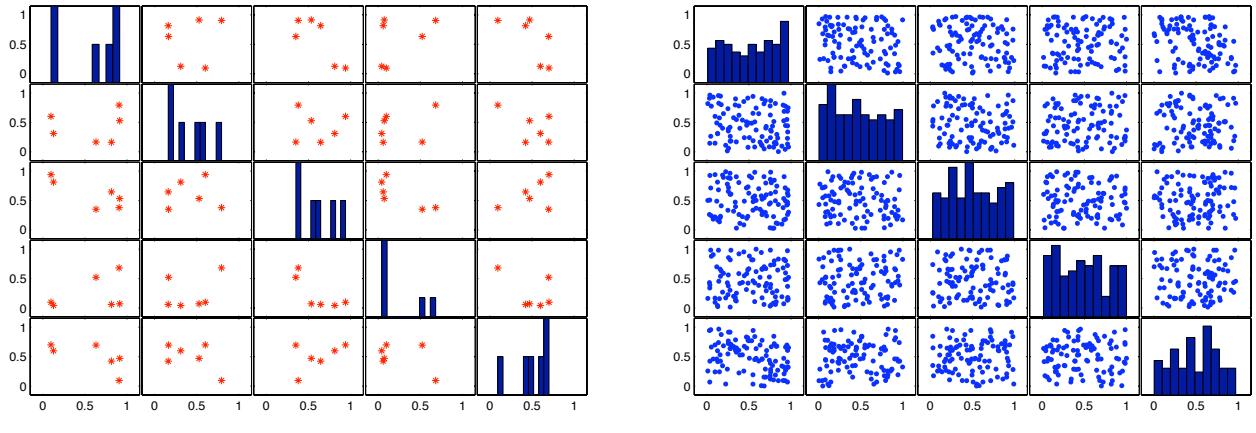
Labwork 117 Let us make matrix plots from a uniformly generated sequence of 100 points in 5D unit cube $[0, 1]^5$ as shown in Figure 4.8.

```

>> 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

```

Figure 4.8: Plot Matrix of uniformly generated data in $[0, 1]^5$



(a) First six samples

(b) All thousand samples

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				

4.1.5 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-Wesely, New York, 1977*].

4.1.6 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.

1. 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).
2. 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:

```
CUSP_ID,LAT,LONG,NZMGE,NZMGN,ORI_YEAR,ORI_MONTH,ORI_DAY,ORI_HOUR,ORI_MINUTE,ORI_SECOND,MAG,DEPTH
3481751,-43.55432,172.68898,2484890,5739375,2011,2,22,0,0,31.27814,3.79,5.8559,
3481760,-43.56579,172.70621,2486287,5738106,2011,2,22,0,0,43.70276,3.76,5.4045,
.
.
.
3469114,-43.58007,172.67126,2483470,5736509,2011,2,22,23,28,11.1014,3.117,3,
3469122,-43.55949,172.70396,2486103,5738805,2011,2,22,23,50,1.06171,3.136,12,
```

The thirteen columns correspond to fairly self-descriptive features of each measured earthquake 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.

3. The next set of commands show you how to load, manipulate and visually explore this data.

```
%% Load the data from the comma delimited text file 'NZ20110222earthquakes.csv' with
%% the following column IDs
%% CUSP_ID,LAT,LONG,NZMGE,NZMGN,ORI_YEAR,ORI_MONTH,ORI_DAY,ORI_HOUR,ORI_MINUTE,ORI_SECOND,MAG,DEPTH
%% Using MATLAB's dlmread command we can assign the data as a matrix to EQ;
%% note that the option 1,0 to dlmread skips first row of column descriptors
%
% the variable EQall is about to be assigned the data as a matrix
EQall = dlmread('NZ20110222earthquakes.csv', ',', 1, 0);
size(EQall) % report the dimensions or size of the matrix EQall
ans =
    145      14
```

4. In order to understand the syntax in detail get `help` from MATLAB !

```
>> help dlmread
DLMREAD Read ASCII delimited file.
.
.
.
```

5. When there 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 if there 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 earthquake in Christchurch on February 22 of 2011 as a scatter plot shown in Figure 4.9 (the axes labels were subsequently added from clicking <Edit> and <Figure Properties...> tabs of the output Figure Window).

```

>> EQall(any(isnan(EQall),2),:) = []; %Remove any rows containing NaNs from the matrix EQall
>> % report the size of EQall and see if it is different from before we removed and NaN containing rows
>> size(EQall)
ans = 145 14
>> % remove locations outside Chch and assign it to a new variable called EQ
>> EQ = EQall(-43.75<EQall(:,2) & EQall(:,2)<-43.45 ...
& 172.45<EQall(:,3) & EQall(:,3)<172.9 & EQall(:,12)>3, :);
>> % now report the size of the earthquakes in Christchurch in variable EQ
>> size(EQ)
ans = 124 14
>> % assign the four variables of interest
>> 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:

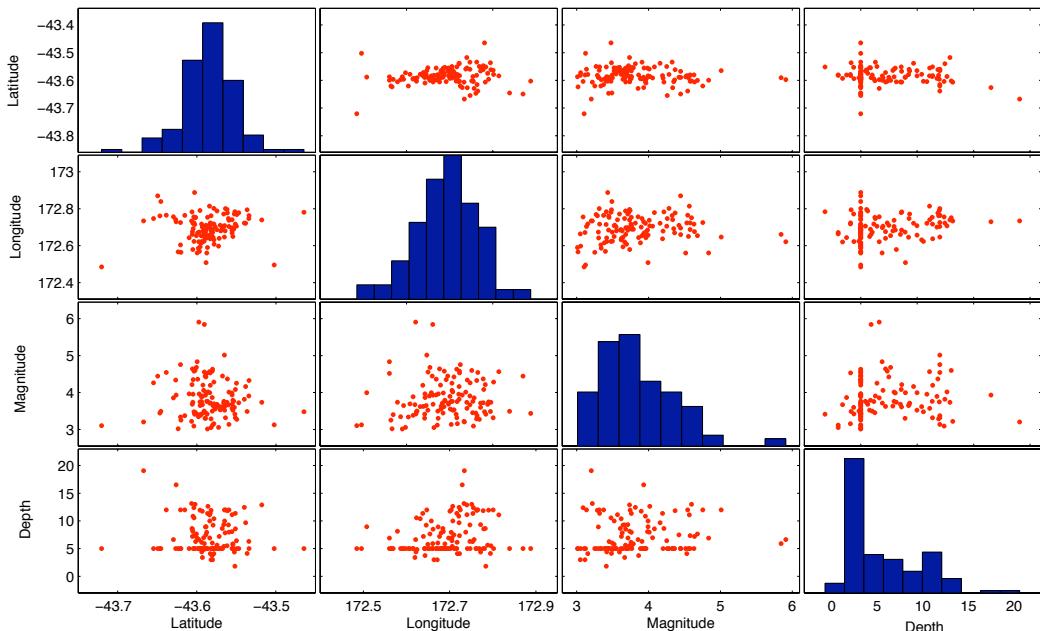
```

>> 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 4.9: 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
```

```
NZEQChCch20110222EDA.m
%% Load the data from the comma delimited text file 'NZ20110222earthquakes.csv'
% using the script M-file NZEQChCch20110222.m
NZEQChCch20110222
%% working with time stamps is tricky
%% time is encoded by columns 6 through 11
%% as origin of earthquake in year, month, day, hour, minute, sec:
%% ORI_YEAR,ORI_MONTH,ORI_DAY,ORI_HOUR,ORI_MINUTE,ORI_SECOND
%% datenum is Matlab's date encoding function see help datenum
TimeData=datenum(EQ(:,6:11)); % assign origin times of earth quakes in datenum coordinates
MaxD=max(TimeData); % get the latest time of observation in the data
MinD=min(TimeData); % % get the earliest time of observation in the data
datestr(MinD) % a nice way to conver to calendar time!
datestr(MaxD) % ditto

% recall that there four variables were assigned in NZEQChCch20110222.m
% LatData=EQ(:,2); LonData=EQ(:,3); MagData=EQ(:,12); DepData=EQ(:,13);

%clear any existing Figure windows
clf
plot(TimeData,MagData,'o-') % plot origin time against magnitude of each earth quake

figure % tell matlab you are about to make another figure
plotmatrix([LatData,LonData,MagData,DepData],'r.');

figure % tell matlab you are about to make another figure
scatter(LonData,LatData,'.') % plot the LONGitude Vs. LATtitude

figure % tell matlab you are about to make another figure
% relative frequency histogram of magnitudes from 0 to 12 on Richter Scale with 15 bins
hist(MagData,15)

%max(MagData)

figure % tell matlab you are about to make another figure
semilogx(DepData,MagData,'.') % see the depth in log scale

%%%%%
% more advanced topic - uncomment and read help if bored
%tri = delaunay(LatData,LonData);
%triplot(tri,LatData,LonData,DepData);
```

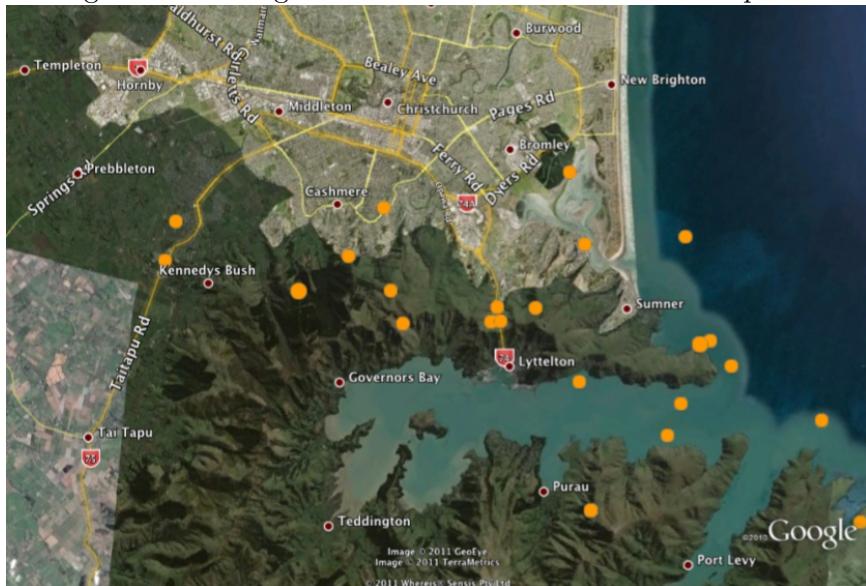
Geostatistical exploratory data analysis with Google Earth

A global search at <http://neic.usgs.gov/cgi-bin/epic/epic.cgi> with the following parameters:

Date Range: 2011 2 22 to 2011 2 22

Catalog: USGS/NEIC (PDE-Q)

Figure 4.10: Google Earth Visualisation of the earth quakes



produced 43 earth quakes world-wide, including those in Christchurch as shown in Figure 4.10. One can do a lot more than a mere visualisation with the USGS/NEIC database of earthquakes world-wide, the freely available **Google earth** software bundle <http://www.google.com/earth/index.html> and the freely available MATLAB package **googleearth** from http://www.mathworks.com/matlabcentral/fx_files/12954/4/content/googleearth/html/html_product_page.html.

4.1.7 Metereological Data

New Zealand's meteorological service NIWA provides weather data under its TERMS AND CONDITIONS FOR ACCESS TO DATA (See http://cliflo.niwa.co.nz/doc/terms_print.html). 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 4.11 is obtained by the script file `RainFallsInChch.m` by typing and following the comments:

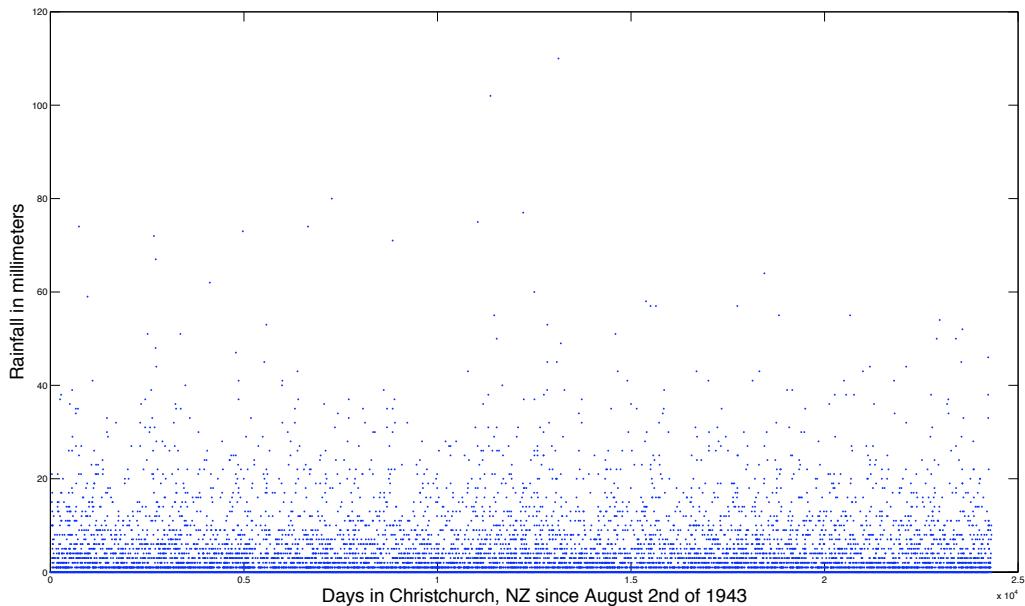
```
>> RainFallsInChch
RainFallsChch = [24312x1 int32]    [24312x1 double]
ans =      24312          2
FirstDayOfData =   19430802
LastDayOfData =   20100721
```

```
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)
```

Figure 4.11: Daily rainfalls in Christchurch since March 27 2010



Daily Temperatures in Christchurch

Labwork 121 Understand how Figure 4.12 is being generated by following the comments in the script file ChchTempsLoad.m by typing:

```
>> ChchTempsLoad
```

```
ChchTempsLoad.m
%% Load the data from the comma delimited text file 'NIWACliFloChchAeroClubStationTemps.txt',
%% with the following column IDs
%% Max_min: Daily Temperature in Christchurch New Zealand
%% Stationate(NZST),Tmax(C),Period(Hrs),Tmin(C),Period(Hrs),Tgmin(C),Period(Hrs),Tmean(C),RHmean(%),Period(Hrs)

% the matrix T is about to be assigned the data as a matrix; the option [27,1,20904,5] to
% specify the upper-left and lower-right corners of an imaginary rectangle
% over the text file 'NIWACliFloChchAeroClubStationTemps.txt'.
```

```

% here we start from line number 27 and end at the last line number 20904
% and we read only columns NZST,Tmax(C),Period(Hrs),Tmin(C),Period(Hrs)

T = dlmread('NIWACliFloChchAeroClubStationTemps.txt','','',[27,1,20904,5]);
% just keep column 1,2 and 4 named NZST,Tmax(C),Period(Hrs),Tmin(C),
% i.e. date in YYYYMMDD foramt, maximum temperature, minimum temperature
T = T(:,[1,2,4]); % just pull the time
% print size before removing missig data rows are removed
size(T) % report the dimensions or size of the matrix T

% This file has a lot of missing data points and they were replaced with
% NaN values - see the file for various manipulations that were done to the
% raw text file from NIWA (Copyright NIWA 2011 Subject to NIWA's Terms and
% Conditions. See: http://cliflo.niwa.co.nz/pls/niwp/doc/terms.html)
T(any(isnan(T),2),:) = [];% Remove any rows containing NaNs from a matrix

size(T) % if the matrix has a different size now then the data-less days now!

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

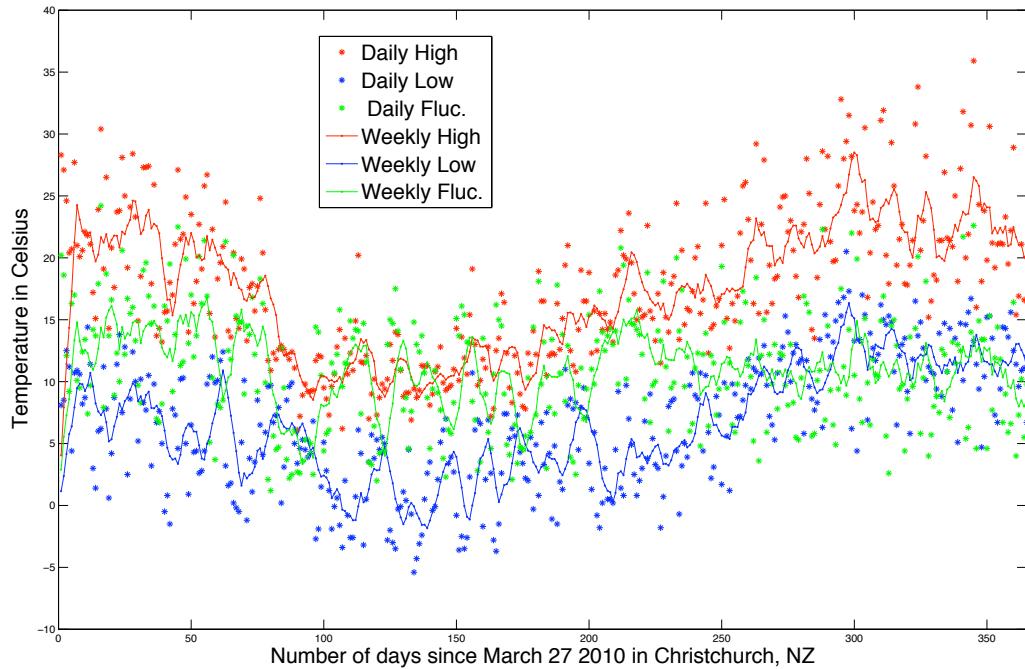
```

4.1.8 Textual Data

Processing and analysing textual data to make a decision is another important computational statistical experiment. An obvious example is machine translation and a less obvious one is exploratory data analysis of the textual content of

- a large document
- twitter messages within an online social network of interest
- etc.

Figure 4.12: Daily temperatures in Christchurch for one year since March 27 2010



An interesting document with a current affairs projection is the Joint Operating Environment 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 74 page document (JOE 2010 Report) reads:

ABOUT THIS STUDY The Joint Operating Environment is intended to inform joint concept development and experimentation throughout the Department of Defense. It provides a perspective on future trends, shocks, contexts, and implications for future 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 Mitscher Avenue, Suite 200, Norfolk, VA 23551-2488, (757) 836-6555.

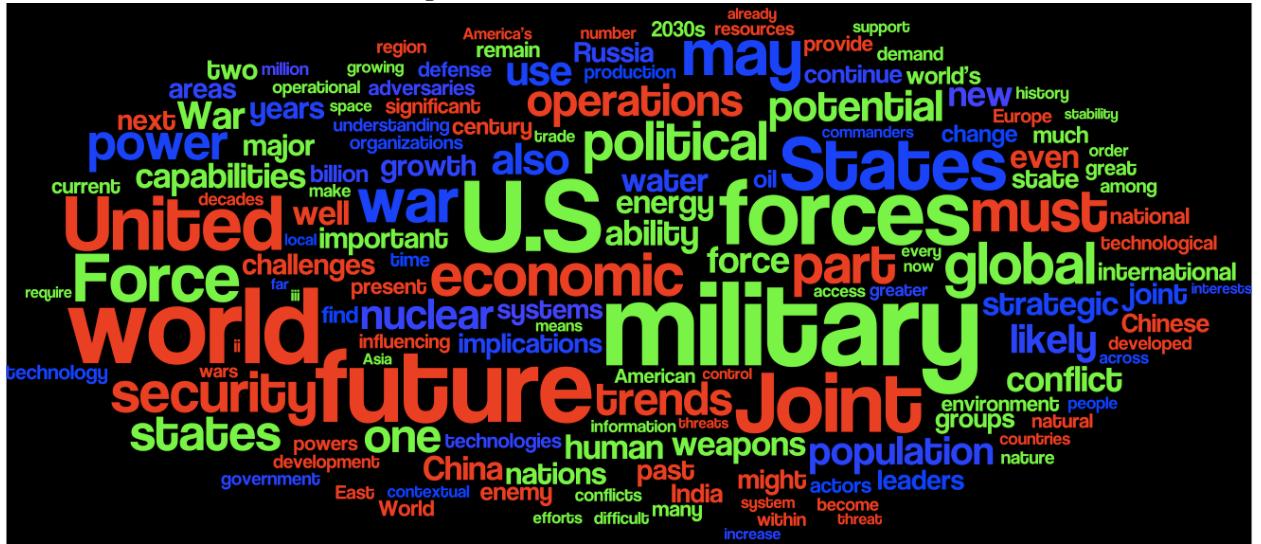
Distribution Statement A: Approved for Public Release

We can try to produce a statistic of this document by recording the frequency of words in its textual content. Then we can produce a “word histogram” or “word cloud” to explore the document visually at one of the coarsest possible resolutions of the textual content in the JOE 2010 Report. The “word cloud” shown in Figure 4.13 was produced by Phillip Wilson using *wordle* from <http://www.wordle.net/>. A description from the wordle URL says:

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 yours to use however you like. You can print them out, or save them to the Wordle gallery to share with your friends.

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

Figure 4.13: Wordle of JOE 2010



wordle from <http://www.wordle.net/>. Play with the aesthetic features to change colour, shapes, etc.

4.1.9 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 another 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

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 4.14. 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*.

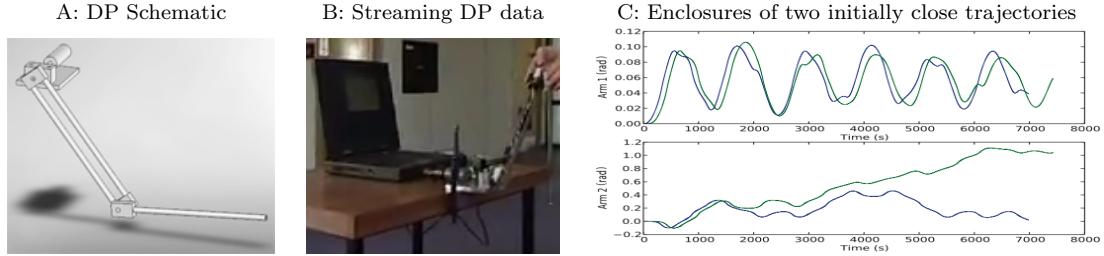


Figure 4.14: Double Pendulum

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 4.14(C) using time series data of trajectory 1 and trajectory 2 linked from the bottom of the above URL.

4.2 Exercises in Statistics

Ex. 4.1 — What is the sample mean and sample variance of the following dataset:

$$1, 3, 2, 1, 2, 3, 3$$

4.3 Fundamentals of Estimation

4.3.1 Introduction

Now that we have been introduced to two notions of convergence for RV sequences, we can begin to appreciate the basic limit theorems used in statistical inference. The problem of estimation is of fundamental importance in statistical inference and learning. We will formalise the general estimation problem here. There are two basic types of estimation. In point estimation we are interested in estimating a particular point of interest that is supposed to belong to a set of points. In (confidence) set estimation, we are interested in estimating a set with a particular form that has a specified probability of “trapping” the particular point of interest from a set of points. Here, a point should be interpreted as an element of a collection of elements from some space.

4.3.2 Point Estimation

Point estimation is any statistical methodology that provides one with a “single best guess” of some specific quantity of interest. Traditionally, we denote this **quantity of interest as θ^*** and its **point estimate as $\hat{\theta}$ or $\hat{\theta}_n$** . The subscript n in the point estimate $\hat{\theta}_n$ emphasises that our estimate is based on n observations or data points from a given statistical experiment to estimate θ^* . This quantity of interest, which is usually unknown, can be:

- an **integral** $\vartheta^* := \int_A h(x) dx \in \Theta$. If ϑ^* is finite, then $\Theta = \mathbb{R}$, or
- a **parameter** θ^* which is an element of the **parameter space** Θ , denoted $\theta^* \in \Theta$,
- a **distribution function (DF)** $F^* \in \mathbb{F} :=$ the set of all DFs

- a **density function (pdf)** $f \in \{ \text{"not too wiggly Sobolev functions"} \}$, or
- a **regression function** $g^* \in \mathbb{G}$, where \mathbb{G} is a class of regression functions in a regression experiment with model: $Y = g^*(X) + \epsilon$, such that $E(\epsilon) = 0$, from pairs of observations $\{(X_i, Y_i)\}_{i=1}^n$, or
- a **classifier** $g^* \in \mathbb{G}$, i.e. a regression experiment with discrete $Y = g^*(X) + \epsilon$, or
- a **prediction** in a regression experiment, i.e. when you want to estimate Y_i given X_i .

Recall that a statistic is an RV $T(X)$ that maps every data point x in the data space \mathbb{X} with $T(x) = t$ in its range \mathbb{T} , i.e. $T(x) : \mathbb{X} \rightarrow \mathbb{T}$ (Definition 51). Next, we look at a specific class of statistics whose range is the parameter space Θ .

Definition 58 (Point Estimator) A **point estimator** $\hat{\Theta}$ of some **fixed and possibly unknown** $\theta^* \in \Theta$ is a statistic that associates each data point $x \in \mathbb{X}$ with an estimate $\hat{\Theta}(x) = \hat{\theta} \in \Theta$,

$$\boxed{\hat{\Theta} := \hat{\Theta}(x) = \hat{\theta} : \mathbb{X} \rightarrow \Theta} .$$

If our data point $x := (x_1, x_2, \dots, x_n)$ is an n -vector or a point in the n -dimensional real space, i.e. $x := (x_1, x_2, \dots, x_n) \in \mathbb{X}_n \subset \mathbb{R}^n$, then we emphasise the dimension n in our point estimator $\hat{\Theta}_n$ of $\theta^* \in \Theta$.

$$\boxed{\hat{\Theta}_n := \hat{\Theta}_n(x := (x_1, x_2, \dots, x_n)) = \hat{\theta}_n : \mathbb{X}_n \rightarrow \Theta, \quad \mathbb{X}_n \subset \mathbb{R}^n} .$$

The typical situation for us involves point estimation of $\theta^* \in \Theta$ on the basis of one realisation $x \in \mathbb{X}_n \subset \mathbb{R}^n$ of an independent and identically distributed (IID) random vector $X = (X_1, X_2, \dots, X_n)$, such that $X_1, X_2, \dots, X_n \stackrel{\text{IID}}{\sim} X_1$ and the DF of X_1 is $F(x_1; \theta^*)$, i.e. the distribution of the IID RVs, X_1, X_2, \dots, X_n , is parameterised by $\theta^* \in \Theta$.

Example 124 (Coin Tossing Experiment) ($(X_1, \dots, X_n) \stackrel{\text{IID}}{\sim} \text{Bernoulli}(\theta^*)$) I tossed a coin that has an unknown probability θ^* of landing Heads independently and identically 10 times in a row. Four of my outcomes were Heads and the remaining six were Tails, with the actual sequence of Bernoulli outcomes (Heads $\rightarrow 1$ and Tails $\rightarrow 0$) being $(1, 0, 0, 0, 1, 1, 0, 0, 1, 0)$. I would like to estimate the probability $\theta^* \in \Theta = [0, 1]$ of observing Heads using the natural estimator $\hat{\Theta}_n((X_1, X_2, \dots, X_n))$ of θ^* :

$$\hat{\Theta}_n((X_1, X_2, \dots, X_n)) := \hat{\Theta}_n = \frac{1}{n} \sum_{i=1}^n X_i =: \bar{X}_n$$

For the coin tossing experiment I just performed ($n = 10$ times), the point estimate of the unknown θ^* is:

$$\begin{aligned} \hat{\theta}_{10} &= \hat{\Theta}_{10}((x_1, x_2, \dots, x_{10})) = \hat{\Theta}_{10}((1, 0, 0, 0, 1, 1, 0, 0, 1, 0)) \\ &= \frac{1+0+0+0+1+1+0+0+1+0}{10} = \frac{4}{10} = 0.40 . \end{aligned}$$

Labwork 125 (Bernoulli(38/75) Computer Experiment) Simulate one thousand IID samples from a $\text{Bernoulli}(\theta^* = 38/75)$ RV and store this data in an array called **Samples**. Use your student ID to initialise the fundamental sampler. Now, pretend that you don't know the true

θ^* and estimate θ^* using our estimator $\widehat{\Theta}_n = \overline{X}_n$ from the data array **Samples** for each sample size $n = 1, 2, \dots, 1000$. Plot the one thousand estimates $\widehat{\theta}_1, \widehat{\theta}_2, \dots, \widehat{\theta}_{1000}$ as a function of the corresponding sample size. Report your observations regarding the behaviour of our estimator as the sample size increases.

4.3.3 Some Properties of Point Estimators

Given that an estimator is merely a function from the data space to the parameter space, we need choose only the best estimators available. Recall that a point estimator $\widehat{\Theta}_n$, being a statistic or an RV of the data has a probability distribution over its range Θ . This distribution over Θ is called the **sampling distribution** of $\widehat{\Theta}_n$. Note that the sampling distribution not only depends on the statistic $\widehat{\Theta}_n := \widehat{\Theta}_n(X_1, X_2, \dots, X_n)$ but also on θ^* which in turn determines the distribution of the IID data vector (X_1, X_2, \dots, X_n) . The following definitions are useful for selecting better estimators from some lot of them.

Definition 59 (Bias of a Point Estimator) The bias_n of an estimator $\widehat{\Theta}_n$ of $\theta^* \in \Theta$ is:

$$\boxed{\text{bias}_n = \text{bias}_n(\widehat{\Theta}_n) := E_{\theta^*}(\widehat{\Theta}_n) - \theta^* = \int_{\mathbb{X}_n} \widehat{\Theta}_n(x) dF(x; \theta^*) - \theta^*} . \quad (4.10)$$

We say that the estimator $\widehat{\Theta}_n$ is **unbiased** if $\text{bias}_n(\widehat{\Theta}_n) = 0$ or if $E_{\theta^*}(\widehat{\Theta}_n) = \theta^*$ for every n . If $\lim_{n \rightarrow \infty} \text{bias}_n(\widehat{\Theta}_n) = 0$, we say that the estimator is **asymptotically unbiased**.

Since the expectation of the sampling distribution of the point estimator $\widehat{\Theta}_n$ depends on the unknown θ^* , we emphasise the θ^* -dependence by $E_{\theta^*}(\widehat{\Theta}_n)$.

Example 126 (Bias of our Estimator of θ^*) Consider the sample mean estimator $\widehat{\Theta}_n := \overline{X}_n$ of θ^* , from $X_1, X_2, \dots, X_n \stackrel{\text{IID}}{\sim} \text{Bernoulli}(\theta^*)$. That is, we take the sample mean of the n IID Bernoulli(θ^*) trials to be our point estimator of $\theta^* \in [0, 1]$. Then, **this estimator is unbiased** since:

$$E_{\theta^*}(\widehat{\Theta}_n) = E_{\theta^*} \left(n^{-1} \sum_{i=1}^n X_i \right) = n^{-1} E_{\theta^*} \left(\sum_{i=1}^n X_i \right) = n^{-1} \sum_{i=1}^n E_{\theta^*}(X_i) = n^{-1} n \theta^* = \theta^* .$$

Definition 60 (Standard Error of a Point Estimator) The standard deviation of the point estimator $\widehat{\Theta}_n$ of $\theta^* \in \Theta$ is called the **standard error**:

$$\boxed{\text{se}_n = \text{se}_n(\widehat{\Theta}_n) = \sqrt{V_{\theta^*}(\widehat{\Theta}_n)} = \sqrt{\int_{\mathbb{X}_n} (\widehat{\Theta}_n(x) - E_{\theta^*}(\widehat{\Theta}_n))^2 dF(x; \theta^*)}} . \quad (4.11)$$

Since the variance of the sampling distribution of the point estimator $\widehat{\Theta}_n$ depends on the fixed and possibly unknown θ^* , as emphasised by V_{θ^*} in (4.11), the se_n is also a possibly unknown quantity and may itself be estimated from the data. The estimated standard error, denoted by $\widehat{\text{se}}_n$, is calculated by replacing $V_{\theta^*}(\widehat{\Theta}_n)$ in (4.11) with its appropriate estimate.

Example 127 (Standard Error of our Estimator of θ^*) Consider the sample mean estimator $\widehat{\Theta}_n := \overline{X}_n$ of θ^* , from $X_1, X_2, \dots, X_n \stackrel{\text{IID}}{\sim} \text{Bernoulli}(\theta^*)$. Observe that the statistic:

$$T_n((X_1, X_2, \dots, X_n)) := n \widehat{\Theta}_n((X_1, X_2, \dots, X_n)) = \sum_{i=1}^n X_i$$

is the Binomial(n, θ^*) RV. The standard error se_n of this estimator is:

$$\text{se}_n = \sqrt{\text{V}_{\theta^*}(\hat{\Theta}_n)} = \sqrt{\text{V}_{\theta^*}\left(\sum_{i=1}^n \frac{X_i}{n}\right)} = \sqrt{\left(\sum_{i=1}^n \frac{1}{n^2} \text{V}_{\theta^*}(X_i)\right)} = \sqrt{\frac{n}{n^2} \text{V}_{\theta^*}(X_i)} = \sqrt{\theta^*(1-\theta^*)/n}.$$

Another reasonable property of an estimator is that it converge to the “true” parameter θ^* – here “true” means the supposedly fixed and possibly unknown θ^* , as we gather more and more IID data from a θ^* -specified DF $F(x; \theta^*)$. This property is stated precisely next.

Definition 61 (Asymptotic Consistency of a Point Estimator) A point estimator $\hat{\Theta}_n$ of $\theta^* \in \Theta$ is said to be **asymptotically consistent** if:

$$\boxed{\hat{\Theta}_n \xrightarrow{P} \theta^*} \quad \text{i.e., for any real } \epsilon > 0, \quad \lim_{n \rightarrow \infty} P(|\hat{\Theta}_n - \theta^*| > \epsilon) = 0.$$

Definition 62 (Mean Squared Error (MSE) of a Point Estimator) Often, the quality of a point estimator $\hat{\Theta}_n$ of $\theta^* \in \Theta$ is assessed by the **mean squared error** or MSE_n defined by:

$$\boxed{\text{MSE}_n = \text{MSE}_n(\hat{\Theta}_n) := \text{E}_{\theta^*}((\hat{\Theta}_n - \theta^*)^2) = \int_{\mathbb{X}} (\hat{\Theta}_n(x) - \theta^*)^2 dF(x; \theta^*)}. \quad (4.12)$$

The following proposition shows a simple relationship between the mean square error, bias and variance of an estimator $\hat{\Theta}_n$ of θ^* .

Proposition 63 (The $\sqrt{\text{MSE}_n}$: se_n : bias_n -Sided Right Triangle of an Estimator) Let $\hat{\Theta}_n$ be an estimator of $\theta^* \in \Theta$. Then:

$$\boxed{\text{MSE}_n(\hat{\Theta}_n) = (\text{se}_n(\hat{\Theta}_n))^2 + (\text{bias}_n(\hat{\Theta}_n))^2}. \quad (4.13)$$

Proof:

$$\begin{aligned}
& LHS \\
&= \text{MSE}_n(\hat{\Theta}_n) \\
&:= \mathbb{E}_{\theta^*}((\hat{\Theta}_n - \theta^*)^2), \quad \text{by definition of } \text{MSE}_n \text{ (4.12)} \\
&= \mathbb{E}_{\theta^*} \left(\left(\underbrace{\hat{\Theta}_n - \mathbb{E}_{\theta^*}(\hat{\Theta}_n)}_A + \underbrace{\mathbb{E}_{\theta^*}(\hat{\Theta}_n) - \theta^*}_B \right)^2 \right), \quad \text{by subtracting and adding the constant } \mathbb{E}_{\theta^*}(\hat{\Theta}_n) \\
&= \mathbb{E}_{\theta^*} \left(\underbrace{(\hat{\Theta}_n - \mathbb{E}_{\theta^*}(\hat{\Theta}_n))^2}_{A^2} + \underbrace{2(\hat{\Theta}_n - \mathbb{E}_{\theta^*}(\hat{\Theta}_n))(\mathbb{E}_{\theta^*}(\hat{\Theta}_n) - \theta^*)}_{2AB} + \underbrace{(\mathbb{E}_{\theta^*}(\hat{\Theta}_n) - \theta^*)^2}_{B^2} \right), \quad \because (A+B)^2 = A^2 + 2AB + B^2 \\
&= \mathbb{E}_{\theta^*}((\hat{\Theta}_n - \mathbb{E}_{\theta^*}(\hat{\Theta}_n))^2) + \mathbb{E}_{\theta^*}(2(\hat{\Theta}_n - \mathbb{E}_{\theta^*}(\hat{\Theta}_n))(\mathbb{E}_{\theta^*}(\hat{\Theta}_n) - \theta^*)) + \mathbb{E}_{\theta^*}((\mathbb{E}_{\theta^*}(\hat{\Theta}_n) - \theta^*)^2), \\
&= \mathbb{E}_{\theta^*}((\hat{\Theta}_n - \mathbb{E}_{\theta^*}(\hat{\Theta}_n))^2) + \underbrace{2(\mathbb{E}_{\theta^*}(\hat{\Theta}_n) - \theta^*)}_{C} \underbrace{\mathbb{E}_{\theta^*}((\hat{\Theta}_n - \mathbb{E}_{\theta^*}(\hat{\Theta}_n)))}_{D} + \mathbb{E}_{\theta^*}((\mathbb{E}_{\theta^*}(\hat{\Theta}_n) - \theta^*)^2), \quad \because C \text{ is constant} \\
&= \mathbb{E}_{\theta^*}((\hat{\Theta}_n - \mathbb{E}_{\theta^*}(\hat{\Theta}_n))^2) + 0 + \mathbb{E}_{\theta^*}((\mathbb{E}_{\theta^*}(\hat{\Theta}_n) - \theta^*)^2), \quad \because D := \mathbb{E}_{\theta^*}((\hat{\Theta}_n - \mathbb{E}_{\theta^*}(\hat{\Theta}_n))) = \mathbb{E}_{\theta^*}(\hat{\Theta}_n) - \mathbb{E}_{\theta^*}(\hat{\Theta}_n) = 0 \\
&= V_{\theta^*}(\hat{\Theta}_n) + \mathbb{E}_{\theta^*}((\mathbb{E}_{\theta^*}(\hat{\Theta}_n) - \theta^*)^2), \quad \because V_{\theta^*}(\hat{\Theta}_n) := \mathbb{E}_{\theta^*}((\hat{\Theta}_n - \mathbb{E}_{\theta^*}(\hat{\Theta}_n))^2), \text{ by definition of variance} \\
&= \left(\sqrt{V_{\theta^*}(\hat{\Theta}_n)} \right)^2 + \mathbb{E}_{\theta^*}((\text{bias}_n(\hat{\Theta}_n))^2), \quad \because \text{bias}_n(\hat{\Theta}_n) = \mathbb{E}_{\theta^*}(\hat{\Theta}_n) - \theta^*, \text{ by definition of bias}_n \text{ of an estimator } \hat{\Theta}_n \\
&= (\text{se}_n(\hat{\Theta}_n))^2 + \mathbb{E}_{\theta^*}((\text{bias}_n(\hat{\Theta}_n))^2), \quad \because \text{se}_n(\hat{\Theta}_n) := \sqrt{V_{\theta^*}(\hat{\Theta}_n)}, \text{ by definition (4.11)} \\
&= (\text{se}_n(\hat{\Theta}_n))^2 + (\text{bias}_n(\hat{\Theta}_n))^2, \quad \because \text{bias}_n(\hat{\Theta}_n) = \mathbb{E}_{\theta^*}(\hat{\Theta}_n) - \theta^* \text{ and } (\text{bias}_n(\hat{\Theta}_n))^2 \text{ are constants.} \\
&= RHS
\end{aligned}$$

Proposition 64 (Asymptotic consistency of a point estimator) Let $\hat{\Theta}_n$ be an estimator of $\theta^* \in \Theta$. Then, if $\text{bias}_n(\hat{\Theta}_n) \rightarrow 0$ and $\text{se}_n(\hat{\Theta}_n) \rightarrow 0$ as $n \rightarrow \infty$, the estimator $\hat{\Theta}_n$ is asymptotically consistent:

$$\hat{\Theta}_n \xrightarrow{P} \theta^*.$$

Proof: If $\text{bias}_n(\hat{\Theta}_n) \rightarrow 0$ and $\text{se}_n(\hat{\Theta}_n) \rightarrow 0$, then by (4.13), $\text{MSE}_n(\hat{\Theta}_n) \rightarrow 0$, i.e. that $\mathbb{E}_{\theta^*}((\hat{\Theta}_n - \theta^*)^2) \rightarrow 0$. This type of convergence of the RV $\hat{\Theta}_n$ to the Point Mass(θ^*) RV as $n \rightarrow \infty$ is called convergence in **quadratic mean** or **convergence in BL2** and denoted by $\hat{\Theta}_n \xrightarrow{qm} \theta^*$. Convergence in quadratic mean is a stronger notion of convergence than convergence in probability, in the sense that

$$\mathbb{E}_{\theta^*}((\hat{\Theta}_n - \theta^*)^2) \rightarrow 0 \quad \text{or} \quad \hat{\Theta}_n \xrightarrow{qm} \theta^* \implies \hat{\Theta}_n \xrightarrow{P} \theta^*.$$

Thus, if we prove the above implication we are done with the proof of our proposition. To show that convergence in quadratic mean implies convergence in probability for general sequence of RVs X_n converging to an RV X , we first assume that $X_n \xrightarrow{qm} X$. Now, fix any $\epsilon > 0$. Then by Markov's inequality (7.2),

$$P(|X_n - X| > \epsilon) = P(|X_n - X|^2 > \epsilon^2) \leq \frac{E(|X_n - X|^2)}{\epsilon^2} \rightarrow 0,$$

and we have shown that the definition of convergence in probability holds provided convergence in quadratic mean holds.

We want our estimator to be unbiased with small standard errors as the sample size n gets large. The **point estimator** $\hat{\Theta}_n$ will then produce a **point estimate** $\hat{\theta}_n$:

$$\hat{\Theta}_n((x_1, x_2, \dots, x_n)) = \hat{\theta}_n \in \Theta,$$

on the basis of the **observed data** (x_1, x_2, \dots, x_n) , that is close to the **true parameter** $\theta^* \in \Theta$.

Example 128 (Asymptotic consistency of our Estimator of θ^*) Consider the sample mean estimator $\widehat{\Theta}_n := \bar{X}_n$ of θ^* , from $X_1, X_2, \dots, X_n \stackrel{\text{IID}}{\sim} \text{Bernoulli}(\theta^*)$. Since $\text{bias}_n(\widehat{\Theta}_n) = 0$ for any n and $\text{se}_n = \sqrt{\theta^*(1-\theta^*)/n} \rightarrow 0$, as $n \rightarrow \infty$, by Proposition 64, $\widehat{\Theta}_n \xrightarrow{P} \theta^*$. That is $\widehat{\Theta}_n$ is an **asymptotically consistent estimator** of θ^* .

4.3.4 Confidence Set Estimation

As we saw in Section 4.3.2, the point estimate $\widehat{\theta}_n$ is a “single best guess” of the fixed and possibly unknown parameter $\theta^* \in \Theta$. However, if we wanted to make a statement about our confidence in an estimation procedure, then one possibility is to produce subsets from the parameter space Θ called **confidence sets** that “engulf” θ^* with a probability of at least $1 - \alpha$.

Formally, an $1 - \alpha$ **confidence interval** for the parameter $\theta^* \in \Theta \subset \mathbb{R}$, based on n observations or data points X_1, X_2, \dots, X_n , is an interval C_n that is a function of the data:

$$C_n := [\underline{C}_n, \bar{C}_n] = [\underline{C}_n(X_1, X_2, \dots, X_n), \bar{C}_n(X_1, X_2, \dots, X_n)] ,$$

such that:

$$P_{\theta^*}(\theta^* \in C_n := [\underline{C}_n, \bar{C}_n]) \geq 1 - \alpha .$$

Note that the confidence interval $C_n := [\underline{C}_n, \bar{C}_n]$ is a two-dimensional RV or a random vector in \mathbb{R}^2 that depends on the two statistics $\underline{C}_n(X_1, X_2, \dots, X_n)$ and $\bar{C}_n(X_1, X_2, \dots, X_n)$, as well as θ^* , which in turn determines the distribution of the data (X_1, X_2, \dots, X_n) . In words, C_n engulfs the true parameter $\theta^* \in \Theta$ with a probability of at least $1 - \alpha$. We call $1 - \alpha$ as the **coverage** of the confidence interval C_n .

Formally, a $1 - \alpha$ **confidence set** C_n for a vector-valued $\theta^* \in \Theta \subset \mathbb{R}^k$ is any subset of Θ such that $P_{\theta^*}(\theta^* \in C_n) \geq 1 - \alpha$. The typical forms taken by C_n are k -dimensional boxes or hyper-cuboids, hyper-ellipsoids and subsets defined by inequalities involving level sets of some estimator of θ^* .

Typically, we take $\alpha = 0.05$ because we are interested in the $1 - \alpha = 0.95$ or 95% confidence interval/set $C_n \subset \Theta$ of $\theta^* \in \Theta$ from an estimator $\widehat{\Theta}_n$ of θ^* .

Let us look at an example that makes use of the CLT next (Exercise in Prob. Theor.I).

Example 129 (Errors in computer code (Wasserman03, p. 78)) 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 modelling the number of errors in the i^{th} program is the *Poisson*($\lambda^* = 5$) for any $i = 1, 2, \dots, n$. Also suppose that they are independently distributed. In short, we suppose that:

$$X_1, X_2, \dots, X_n \stackrel{\text{IID}}{\sim} \text{Poisson}(\lambda^* = 5) .$$

Suppose we have $n = 125$ programs and want to make a probability statement about \bar{X}_n which is the average number of errors per program out of these 125 programs. Since $E(X_i) = \lambda^* = 5$ and $V(X_i) = \lambda^* = 5$, we may want to know how often our sample mean \bar{X}_{125} differs from the expectation of 5 errors per program. Using the CLT, we can approximate $P(\bar{X}_n < 5.5)$, for

instance, as follows:

$$\begin{aligned}
P(\bar{X}_n < 5.5) &= P\left(\frac{\sqrt{n}(\bar{X}_n - E(X_1))}{\sqrt{V(X_1)}} < \frac{\sqrt{n}(5.5 - E(X_1))}{\sqrt{V(X_1)}}\right) \\
&\approx P\left(Z < \frac{\sqrt{n}(5.5 - \lambda^*)}{\sqrt{\lambda^*}}\right) \quad [\text{by CLT, and } E(X_1) = V(X_1) = \lambda^*] \\
&= P\left(Z < \frac{\sqrt{125}(5.5 - 5)}{\sqrt{5}}\right) \quad [\text{Since, } \lambda^* = 5 \text{ and } n = 125 \text{ in this Example}] \\
&= P(Z \leq 2.5) = \Phi(2.5) = \int_{-\infty}^{2.5} \left(\frac{1}{\sqrt{2\pi}} \exp\left(-\frac{x^2}{2}\right) \right) dx \approx 0.993790334674224 .
\end{aligned}$$

To obtain the final number in this approximation, we need the following:

Labwork 130 The numerical approximation of $\Phi(2.5)$ was obtained via the call shown below to our erf-based `NormalCdf` function from [??](#). We could have also found it from a pre-computed Table for $\Phi(x)$.

```
>> format long
>> disp(NormalCdf(2.5,0,1))
0.993790334674224
```

The CLT says that if $X_1, X_2, \dots \stackrel{\text{IID}}{\sim} X_1$ then $Z_n := \sqrt{n}(\bar{X}_n - E(X_1))/\sqrt{V(X_1)}$ is approximately distributed as $\text{Normal}(0, 1)$. In Example 129, we knew $\sqrt{V(X_1)}$ since we assumed knowledge of $\lambda^* = 5$. However, in general, we may not know $\sqrt{V(X_1)}$. The next proposition says that we may estimate $\sqrt{V(X_1)}$ using the sample standard deviation S_n of X_1, X_2, \dots, X_n , according to (4.5), and still make probability statements about the sample mean \bar{X}_n using a Normal distribution, **provided n is not too small**, for e.g. $n \geq 30$.

Proposition 65 (CLT based on Sample Variance) Let $X_1, X_2, \dots \stackrel{\text{IID}}{\sim} X_1$ and suppose $E(X_1)$ and $V(X_1)$ exists, then:

$$\frac{\sqrt{n}(\bar{X}_n - E(X_1))}{S_n} \rightsquigarrow \text{Normal}(0, 1) . \quad (4.14)$$

The following property of an estimator makes it easy to obtain confidence intervals.

Definition 66 (Asymptotic Normality of Estimators) An estimator $\hat{\Theta}_n$ of a fixed and possibly unknown parameter $\theta^* \in \Theta$ is **asymptotically normal** if:

$$\frac{\hat{\Theta}_n - \theta^*}{\text{se}_n} \rightsquigarrow \text{Normal}(0, 1) . \quad (4.15)$$

That is, $\hat{\Theta}_n \rightsquigarrow \text{Normal}(\theta^*, \text{se}_n^2)$. By a further estimation of $\text{se}_n := \sqrt{V_{\theta^*}(\hat{\Theta}_n)}$ by $\hat{\text{se}}_n$, we can see that $\hat{\Theta}_n \rightsquigarrow \text{Normal}(\theta^*, \hat{\text{se}}_n^2)$ on the basis of (4.14).

Proposition 67 (Normal-based Asymptotic Confidence Interval) Suppose an estimator $\hat{\Theta}_n$ of parameter $\theta^* \in \Theta \subset \mathbb{R}$ is asymptotically normal:

$$\hat{\Theta}_n \rightsquigarrow \text{Normal}(\theta^*, \hat{\text{se}}_n^2) .$$

Let the RV $Z \sim \text{Normal}(0, 1)$ have DF Φ and inverse DF Φ^{-1} . Let:

$$z_{\alpha/2} = \Phi^{-1}(1 - (\alpha/2)), \quad \text{that is, } P(Z > z_{\alpha/2}) = \alpha/2 \quad \text{and} \quad P(-z_{\alpha/2} < Z < z_{\alpha/2}) = 1 - \alpha.$$

Then:

$$P_{\theta^*}(\theta^* \in C_n) = P\left(\theta^* \in [\hat{\Theta}_n - z_{\alpha/2}\hat{s}\hat{e}_n, \hat{\Theta}_n + z_{\alpha/2}\hat{s}\hat{e}_n]\right) \rightarrow 1 - \alpha.$$

Therefore:

$$C_n := [\underline{C}_n, \bar{C}_n] = [\hat{\Theta}_n - z_{\alpha/2}\hat{s}\hat{e}_n, \hat{\Theta}_n + z_{\alpha/2}\hat{s}\hat{e}_n]$$

is the $1 - \alpha$ Normal-based asymptotic confidence interval that relies on the asymptotic normality of the estimator $\hat{\Theta}_n$ of $\theta^* \in \Theta \subset \mathbb{R}$.

Proof: Define the centralised and scaled estimator as $Z_n := (\hat{\Theta}_n - \theta^*)/\hat{s}\hat{e}_n$. By assumption, $Z_n \rightsquigarrow Z \sim \text{Normal}(0, 1)$. Therefore,

$$\begin{aligned} P_{\theta^*}(\theta^* \in C_n) &= P_{\theta^*}\left(\theta^* \in [\hat{\Theta}_n - z_{\alpha/2}\hat{s}\hat{e}_n, \hat{\Theta}_n + z_{\alpha/2}\hat{s}\hat{e}_n]\right) \\ &= P_{\theta^*}\left(\hat{\Theta}_n - z_{\alpha/2}\hat{s}\hat{e}_n \leq \theta^* \leq \hat{\Theta}_n + z_{\alpha/2}\hat{s}\hat{e}_n\right) \\ &= P_{\theta^*}\left(-z_{\alpha/2}\hat{s}\hat{e}_n \leq \hat{\Theta}_n - \theta^* \leq z_{\alpha/2}\hat{s}\hat{e}_n\right) \\ &= P_{\theta^*}\left(-z_{\alpha/2} \leq \frac{\hat{\Theta}_n - \theta^*}{\hat{s}\hat{e}_n} \leq z_{\alpha/2}\right) \\ &\rightarrow P_{\theta^*}(-z_{\alpha/2} \leq Z \leq z_{\alpha/2}) \\ &= 1 - \alpha \end{aligned}$$

Figure 4.15: Density and Confidence Interval of the Asymptotically Normal Point Estimator

For 95% confidence intervals, $\alpha = 0.05$ and $z_{\alpha/2} = z_{0.025} = 1.96 \approx 2$. This leads to the **approximate 95% confidence interval** of $\hat{\theta}_n \pm 2\hat{s}\hat{e}_n$, where $\hat{\theta}_n = \hat{\Theta}_n(x_1, x_2, \dots, x_n)$ and x_1, x_2, \dots, x_n are the data or observations of the RVs X_1, X_2, \dots, X_n .

Example 131 (Confidence interval for θ^* from n Bernoulli(θ^*) trials) Let $X_1, X_2, \dots, X_n \stackrel{\text{IID}}{\sim} \text{Bernoulli}(\theta^*)$ for some fixed but unknown parameter $\theta^* \in \Theta = [0, 1]$. Consider the following point estimator of θ^* :

$$\hat{\Theta}_n((X_1, X_2, \dots, X_n)) = n^{-1} \sum_{i=1}^n X_i.$$

That is, we take the sample mean of the n IID Bernoulli(θ^*) trials to be our point estimator of $\theta^* \in [0, 1]$. Then, we already saw that **this estimator is unbiased**

We already saw that the standard error $s\hat{e}_n$ of this estimator is:

$$s\hat{e}_n = \sqrt{\theta^*(1 - \theta^*)/n}.$$

Since θ^* is unknown, we obtain the estimated standard error \widehat{s}_n from the point estimate $\widehat{\theta}_n$ of θ^* on the basis of n observed data points $x = (x_1, x_2, \dots, x_n)$ of the experiment:

$$\widehat{s}_n = \sqrt{\widehat{\theta}_n(1 - \widehat{\theta}_n)/n}, \quad \text{where,} \quad \widehat{\theta}_n = \widehat{\Theta}_n((x_1, x_2, \dots, x_n)) = n^{-1} \sum_{i=1}^n x_i.$$

By the central limit theorem, $\widehat{\Theta}_n \rightsquigarrow \text{Normal}(\theta^*, \widehat{s}_n)$, i.e. $\widehat{\Theta}_n$ is asymptotically normal. Therefore, an asymptotically (for large sample size n) approximate $1 - \alpha$ normal-based confidence interval is:

$$\widehat{\theta}_n \pm z_{\alpha/2} \widehat{s}_n = \widehat{\theta}_n \pm z_{\alpha/2} \sqrt{\frac{\widehat{\theta}_n(1 - \widehat{\theta}_n)}{n}} := \left[\widehat{\theta}_n - z_{\alpha/2} \sqrt{\frac{\widehat{\theta}_n(1 - \widehat{\theta}_n)}{n}}, \widehat{\theta}_n + z_{\alpha/2} \sqrt{\frac{\widehat{\theta}_n(1 - \widehat{\theta}_n)}{n}} \right]$$

We also saw that $\widehat{\Theta}_n$ is an **asymptotically consistent estimator** of θ^* due to Proposition 64.

The confidence Interval for the coin tossing experiment in Example 124 with the observed sequence of Bernoulli outcomes (Heads $\rightarrow 1$ and Tails $\rightarrow 0$) being $(1, 0, 0, 0, 1, 1, 0, 0, 1, 0)$. We estimated the probability θ^* of observing Heads with the **unbiased, asymptotically consistent estimator** $\widehat{\Theta}_n((X_1, X_2, \dots, X_n)) = n^{-1} \sum_{i=1}^n X_i$ of θ^* . The point estimate of θ^* was:

$$\begin{aligned} \widehat{\theta}_{10} &= \widehat{\Theta}_{10}((x_1, x_2, \dots, x_{10})) = \widehat{\Theta}_{10}((1, 0, 0, 0, 1, 1, 0, 0, 1, 0)) \\ &= \frac{1+0+0+0+1+1+0+0+1+0}{10} = \frac{4}{10} = 0.40. \end{aligned}$$

The normal-based confidence interval for θ^* may not be a valid approximation here with just $n = 10$ samples. Nevertheless, we will compute a 95% normal-based confidence interval:

$$C_{10} = 0.40 \pm 1.96 \sqrt{\frac{0.40(1 - 0.40)}{10}} = 0.40 \pm 0.3036 = [0.0964, 0.7036]$$

with a width of 0.6072. When I increased the sample size n of the experiment from 10 to 100 by tossing the same coin another 90 times, I discovered that a total of 57 trials landed as Heads. Thus my point estimate and confidence interval for θ^* are:

$$\widehat{\theta}_{100} = \frac{57}{100} = 0.57 \quad \text{and} \quad C_{100} = 0.57 \pm 1.96 \sqrt{\frac{0.57(1 - 0.57)}{100}} = 0.57 \pm 0.0495 = [0.5205, 0.6195]$$

with a much smaller width of 0.0990. Thus our confidence interval shrank considerably from a width of 0.6072 after an additional 90 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.

4.4 Parameter Estimation and Likelihood

Now that we have been introduced to point and set estimation of the population mean and the population proportion using the notion of convergence in distribution for sequences of RVs as well as concentration inequalities, we can begin to appreciate the art of estimation in a more general setting. Parameter estimation is the basic problem in statistical inference and machine learning. We will formalize the general estimation problem here.

As we have already seen, when estimating the population mean or population proportion, there are two basic types of estimators. In point estimation, as seen in Definition 58, we are interested in estimating a particular point of interest that is supposed to belong to a set of points. In (confidence) set estimation, we are interested in estimating a set with a particular form that has a specified probability of “trapping” the particular point of interest from a set of points. Here, a point should be interpreted as an element of a collection of elements from some space.

4.4.1 Point and Set Estimation – A General Likelihood Approach

Point estimation is any statistical methodology that provides one with a “**single best guess**” of some specific quantity of interest. Traditionally, we denote this **quantity of interest as θ^*** and its **point estimate as $\hat{\theta}$ or $\hat{\theta}_n$** . The subscript n in the point estimate $\hat{\theta}_n$ emphasizes that our estimate is based on n observations or data points from a given statistical experiment to estimate θ^* . This quantity of interest, which is usually unknown, can be:

- a **parameter** θ^* which is an element of the **parameter space** Θ , i.e. $\theta^* \in \Theta$ such that θ^* specifies the “law” of the observations (realizations or samples) of the $\vec{RV}(X_1, \dots, X_n)$ modeled by JPDF or JPMF $f_{X_1, \dots, X_n}(x_1, \dots, x_n; \theta^*)$, or
- a **regression function** $\theta^* \in \Theta$, where Θ is a class of regression functions in a regression experiment with model: $Y = \theta^*(X) + \epsilon$, such that $e(\epsilon) = 0$ and θ^* specifies the “law” of pairs of observations $\{(X_i, Y_i)\}_{i=1}^n$, for e.g., fitting parameters in noisy ODE or PDEs from observed data — one can always do a **prediction** in a regression experiment, i.e. when you want to estimate Y_i given X_i , or
- a **classifier** $\theta^* \in \Theta$, i.e. a regression experiment with discrete $Y = \theta^*(X) + \epsilon$, for e.g. training an scrub-nurse robot to assist a human surgeon, or
- an **integral** $\theta^* := \int_A h(x) dx \in \Theta$. If θ^* is finite, then $\Theta = \mathbb{R}$, for e.g. θ^* could be the volume of a high-dimensional irregular polyhedron, a traffic congestion measure on a network of roadways, the expected profit from a new brew of beer, or the probability of an extreme event such as the collapse of a dam in the Southern Alps in the next 150 years.

Set estimation is any statistical methodology that provides one with a “**best smallest set**”, such as an interval, rectangle, ellipse, etc. that contains θ^* with a high probability $1 - \alpha$.

Recall that a statistic is a RV or $\vec{RV} T(X)$ that maps every data point x in the data space \mathbb{X} with $T(x) = t$ in its range \mathbb{T} , i.e. $T(x) : \mathbb{X} \rightarrow \mathbb{T}$ (Definition 51). Next, we look at a specific class of estimators based on the likelihood of the data.

4.4.2 Likelihood

We take a look at **likelihood** — one of the most fundamental concepts in Statistics.

Definition 68 (Likelihood Function) Suppose (X_1, X_2, \dots, X_n) is a \vec{RV} with JPDF or JPMF $f_{X_1, X_2, \dots, X_n}(x_1, x_2, \dots, x_n; \theta)$ specified by parameter $\theta \in \Theta$. Let the observed data be (x_1, x_2, \dots, x_n) . Then the **likelihood** function given by $L_n(\theta)$ is merely the joint probability of the data, with the exception that we see it as a function of the parameter:

$$L_n(\theta) := L_n(x_1, x_2, \dots, x_n; \theta) = f_{X_1, X_2, \dots, X_n}(x_1, x_2, \dots, x_n; \theta) . \quad (4.16)$$

The **log-likelihood** function is defined by:

$$\ell_n(\theta) := \log(L_n(\theta)) \quad (4.17)$$

Example 132 (Likelihood of the IID Bernoulli(θ^*) experiment) Consider our IID Bernoulli experiment:

$$X_1, X_2, \dots, X_n \stackrel{IID}{\sim} \text{Bernoulli}(\theta^*), \text{ with PDF } f_{X_i}(x_i; \theta) = \theta^{x_i}(1-\theta)^{1-x_i} \mathbb{1}_{\{0,1\}}(x_i), \text{ for } i \in \{1, 2, \dots, n\} .$$

Let us understand the likelihood function for one observation first. There are two possibilities for the first observation.

If we only have one observation and it happens to be $x_1 = 1$, then our likelihood function is:

$$L_1(\theta) = L_1(x_1; \theta) = f_{X_1}(x_1; \theta) = \theta^1(1-\theta)^{1-1} \mathbb{1}_{\{0,1\}}(1) = \theta(1-\theta)^0 1 = \theta$$

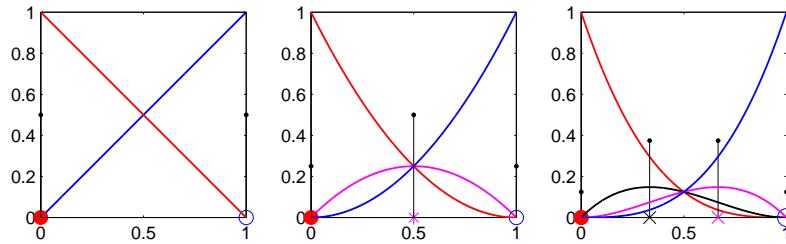
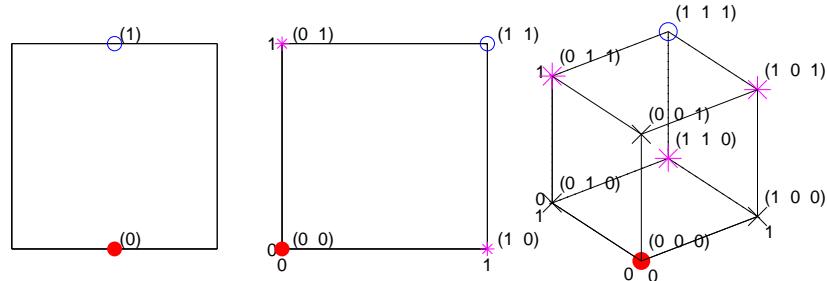
If we only have one observation and it happens to be $x_1 = 0$, then our likelihood function is:

$$L_1(\theta) = L_1(x_1; \theta) = f_{X_1}(x_1; \theta) = \theta^0(1-\theta)^{1-0} \mathbb{1}_{\{0,1\}}(0) = 1(1-\theta)^1 1 = 1 - \theta$$

If we have n observations (x_1, x_2, \dots, x_n) , i.e. a vertex point of the unit hyper-cube $\{0, 1\}^n$ (see top panel of Figure 4.16 when $n \in \{1, 2, 3\}$), then our likelihood function (see bottom panel of Figure 4.16) is obtained by multiplying the densities due to our IID assumption:

$$\begin{aligned} L_n(\theta) &:= L_n(x_1, x_2, \dots, x_n; \theta) = f_{X_1, X_2, \dots, X_n}(x_1, x_2, \dots, x_n; \theta) \\ &= f_{X_1}(x_1; \theta) f_{X_2}(x_2; \theta) \cdots f_{X_n}(x_n; \theta) := \prod_{i=1}^n f_{X_i}(x_i; \theta) \\ &= \theta^{\sum_{i=1}^n x_i} (1-\theta)^{n-\sum_{i=1}^n x_i} \end{aligned} \quad (4.18)$$

Figure 4.16: Data Spaces $\mathbb{X}_1 = \{0, 1\}$, $\mathbb{X}_2 = \{0, 1\}^2$ and $\mathbb{X}_3 = \{0, 1\}^3$ for one, two and three IID Bernoulli trials, respectively and the corresponding likelihood functions.



Definition 69 (Maximum Likelihood Estimator (MLE)) Let the model for the data be

$$(X_1, \dots, X_n) \sim f_{X_1, X_2, \dots, X_n}(x_1, \dots, x_n; \theta^*) .$$

Then the maximum likelihood estimator (MLE) $\hat{\Theta}_n$ of the fixed and possibly unknown parameter $\theta^* \in \Theta$ is the value of θ that maximizes the likelihood function:

$$\boxed{\hat{\Theta}_n := \hat{\Theta}_n(X_1, X_2, \dots, X_n) := \operatorname{argmax}_{\theta \in \Theta} L_n(\theta) ,}$$

Equivalently, MLE is the value of θ that maximizes the log-likelihood function (since $\log = \log_e = \ln$ is a monotone increasing function):

$$\boxed{\hat{\Theta}_n := \operatorname{argmax}_{\theta \in \Theta} \ell_n(\theta) ,}$$

Useful Properties of the Maximum Likelihood Estimator

1. The ML Estimator is *asymptotically consistent* (gives the “true” θ^* as sample size $n \rightarrow \infty$):

$$\boxed{\hat{\Theta}_n \rightsquigarrow \text{Point Mass}(\theta^*)}$$

2. The ML Estimator is asymptotically normal (has a normal distribution concentrating on θ^* as $n \rightarrow \infty$):

$$\boxed{\hat{\Theta}_n \rightsquigarrow \text{Normal}(\theta^*, (\hat{s}\hat{e}_n)^2)}$$

or equivalently:

$$\boxed{(\hat{\Theta}_n - \theta^*)/\hat{s}\hat{e}_n \rightsquigarrow \text{Normal}(0, 1)}$$

where $\hat{s}\hat{e}_n$ is the **estimated standard error**, i.e. the standard deviation of $\hat{\Theta}_n$, and it is given by the square-root of the inverse negative curvature of $\ell_n(\theta)$ at $\hat{\theta}_n$:

$$\boxed{\hat{s}\hat{e}_n = \sqrt{\left(\left[-\frac{d^2 \ell_n(\theta)}{d\theta^2} \right]_{\theta=\hat{\theta}_n} \right)^{-1}}}$$

3. Because of the previous two properties, the $1 - \alpha$ confidence interval is:

$$\boxed{\hat{\Theta}_n \pm z_{\alpha/2} \hat{s}\hat{e}_n}$$

MLE is a general methodology for parameter estimation in an essentially arbitrary parameter space Θ that is defining or indexing the laws in a parametric family of models, although we are only seeing it in action when $\Theta \subset \mathbb{R}$ for simplest parametric family of models involving IID product experiments here. When $\Theta \subset \mathbb{R}^d$ with $2 \leq d < \infty$ then MLE $\hat{\Theta}_n \rightsquigarrow \text{Point Mass}(\theta^*)$, where $\theta^* = (\theta_1^*, \theta_2^*, \dots, \theta_d^*)^T$ is a column vector in $\Theta \subset \mathbb{R}^d$ and $\hat{\Theta}_n \rightsquigarrow \text{Normal}(\theta^*, \widehat{\Sigma(s\hat{e})}_n)$, a multivariate Normal distribution with mean vector θ^* and variance-covariance matrix of standard errors given by the *Hessian* (a $d \times d$ matrix of mixed partial derivatives) of $\ell_n(\theta_1, \theta_2, \dots, \theta_d)$. The ideas in the case of dimension $d = 1$ naturally generalize to an arbitrary, but finite, dimension d .

Remark 70 In order to use MLE for parameter estimation we need to ensure that the following two conditions hold:

1. The *support* of the data, i.e. the set of possible values of (X_1, X_2, \dots, X_n) must not depend on θ for every $\theta \in \Theta$ — of course the probabilities do depend on θ in an *identifiable* manner, i.e. for every θ and ϑ in Θ , if $\theta \neq \vartheta$ then $f_{X_1, X_2, \dots, X_n}(x_1, x_2, \dots, x_n; \theta) \neq f_{X_1, X_2, \dots, X_n}(x_1, x_2, \dots, x_n; \vartheta)$ at least for some $(x_1, x_2, \dots, x_n) \in \mathbb{X}$.
2. If the parameter space Θ is bounded then θ^* must not belong to the boundaries of Θ .

Maximum Likelihood Estimation Method in Six Easy Steps

Background: We have observed data:

$$(x_1, x_2, \dots, x_n)$$

which is modeled as a sample or realization from the random vector:

$$(X_1, X_2, \dots, X_n) \sim f_{X_1, X_2, \dots, X_n}(x_1, x_2, \dots, x_n; \theta^*), \quad \theta^* \in \Theta .$$

Objective: We want to obtain an estimator $\hat{\Theta}_n$ that will give:

1. the point estimate $\hat{\theta}_n$ of the “true” parameter θ^* and
2. the $(1 - \alpha)$ confidence interval for θ^* .

Steps of MLE:

- Step 1: Find the expression for the log likelihood function:

$$\ell_n(\theta) = \log(L_n(\theta)) = \log(f_{X_1, X_2, \dots, X_n}(x_1, x_2, \dots, x_n; \theta)) .$$

Note that if the model assumes that (X_1, X_2, \dots, X_n) is jointly independent, i.e. we have an independent and identically distributed (IID) experiment, then $\ell_n(\theta)$ simplifies further as follows:

$$\ell_n(\theta) = \log(L_n(\theta)) = \log(f_{X_1, X_2, \dots, X_n}(x_1, x_2, \dots, x_n; \theta)) = \log\left(\prod_{i=1}^n f_{X_i}(x_i; \theta)\right) .$$

- Step 2: Obtain the derivative of $\ell_n(\theta)$ with respect to θ :

$$\frac{d}{d\theta} (\ell_n(\theta)) .$$

- Step 3: Set the derivative equal to zero, solve for θ and let $\hat{\theta}_n$ equal to this solution.
- Step 4: Check if this solution is indeed a maximum of $\ell_n(\theta)$ by checking if:

$$\frac{d^2}{d\theta^2} (\ell_n(\theta)) < 0 .$$

- Step 5: If $\frac{d^2}{d\theta^2} (\ell_n(\theta)) < 0$ then you have found the maximum likelihood estimate $\hat{\theta}_n$.

- Step 6: If you also want the $(1 - \alpha)$ confidence interval then get it from

$$\hat{\theta}_n \pm z_{\alpha/2} \widehat{\text{se}}_n \quad , \text{ where } \widehat{\text{se}}_n = \sqrt{\left(\left[-\frac{d^2 \ell_n(\theta)}{d\theta^2} \right]_{\theta=\hat{\theta}_n} \right)^{-1}} .$$

Let us apply this method in some examples.

Example 133 (Maximum likelihood estimation for IID Exponential(λ^*) trials) Find (or derive) the maximum likelihood estimate $\hat{\lambda}_n$ and the $(1 - \alpha)$ confidence interval of the fixed and possibly unknown parameter λ^* for the IID experiment:

$$X_1, \dots, X_n \stackrel{IID}{\sim} \text{Exponential}(\lambda^*), \quad \lambda^* \in \mathbf{A} = (0, \infty) .$$

Note that \mathbf{A} is the parameter space.

We first obtain the log-likelihood function $\ell_n(\theta)$ given data (x_1, x_2, \dots, x_n) .

$$\begin{aligned} \ell_n(\lambda) &:= \log(L(x_1, x_2, \dots, x_n; \lambda)) = \log \left(\prod_{i=1}^n f_{X_i}(x_i; \lambda) \right) = \log \left(\prod_{i=1}^n \lambda e^{-\lambda x_i} \right) \\ &= \log \left(\lambda e^{-\lambda x_1} \cdot \lambda e^{-\lambda x_2} \cdots \lambda e^{-\lambda x_n} \right) = \log \left(\lambda^n e^{-\lambda \sum_{i=1}^n x_i} \right) = \log(\lambda^n) + \log \left(e^{-\lambda \sum_{i=1}^n x_i} \right) \\ &= \boxed{\log(\lambda^n) - \lambda \sum_{i=1}^n x_i} \end{aligned}$$

Now, let us take the derivative with respect to λ ,

$$\begin{aligned} \frac{d}{d\lambda} (\ell_n(\lambda)) &:= \frac{d}{d\lambda} \left(\log(\lambda^n) - \lambda \sum_{i=1}^n x_i \right) = \frac{d}{d\lambda} (\log(\lambda^n)) - \frac{d}{d\lambda} \left(\lambda \sum_{i=1}^n x_i \right) = \frac{1}{\lambda^n} \frac{d}{d\lambda} (\lambda^n) - \sum_{i=1}^n x_i \\ &= \frac{1}{\lambda^n} n \lambda^{n-1} - \sum_{i=1}^n x_i = \boxed{\frac{n}{\lambda} - \sum_{i=1}^n x_i} . \end{aligned}$$

Next, we set the derivative to 0, solve for λ , and let the solution equal to the ML estimate $\hat{\lambda}_n$.

$$0 = \frac{d}{d\lambda} (\ell_n(\lambda)) \iff 0 = \frac{n}{\lambda} - \sum_{i=1}^n x_i \iff \sum_{i=1}^n x_i = \frac{n}{\lambda} \iff \lambda = \frac{n}{\sum_{i=1}^n x_i} \quad \text{and let } \boxed{\hat{\lambda}_n = \frac{1}{\bar{x}_n}} .$$

Next, we find the second derivative and check if it is negative.

$$\frac{d^2}{d\lambda^2} (\ell_n(\lambda)) = \frac{d}{d\lambda} \left(\frac{d}{d\lambda} (\ell_n(\lambda)) \right) = \frac{d}{d\lambda} \left(\frac{n}{\lambda} - \sum_{i=1}^n x_i \right) = \boxed{-n\lambda^{-2}} .$$

Since $\lambda > 0$ and $n \in \mathbb{N}$, $\boxed{-n\lambda^{-2} = -n/\lambda^2 < 0}$, so we have found the maximum likelihood estimate:

$$\boxed{\hat{\lambda}_n = \frac{1}{\bar{x}_n}} .$$

Now, let us find the estimated standard error:

$$\begin{aligned}\widehat{\text{se}}_n &= \sqrt{\left(\left[-\frac{d^2\ell_n(\lambda)}{d\lambda^2}\right]_{\lambda=\widehat{\lambda}_n}\right)^{-1}} = \sqrt{\left(\left[-\left(-\frac{n}{\lambda^2}\right)\right]_{\lambda=\widehat{\lambda}_n}\right)^{-1}} = \sqrt{\left(\frac{n}{\widehat{\lambda}_n^2}\right)^{-1}} = \sqrt{\frac{\widehat{\lambda}_n^2}{n}} = \frac{\widehat{\lambda}_n}{\sqrt{n}} \\ &= \frac{1}{\bar{x}_n\sqrt{n}}.\end{aligned}$$

And finally, the $(1 - \alpha)$ confidence interval is

$$\widehat{\lambda}_n \pm z_{\alpha/2} \widehat{\text{se}}_n = \left[\frac{1}{\bar{x}_n} \pm z_{\alpha/2} \frac{1}{\bar{x}_n\sqrt{n}} \right].$$

Since we have worked “hard” to get the maximum likelihood estimate for a general IID model $X_1, X_2, \dots, X_n \stackrel{IID}{\sim} \text{Exponential}(\lambda^*)$. Let us kill two birds with the same stone by applying it to two datasets:

1. Orbiter waiting times and
2. Time between measurable earthquakes in New Zealand over a few months.

Therefore, the ML estimate $\widehat{\lambda}_n$ of the unknown rate parameter $\lambda^* \in \Lambda$ on the basis of n IID observations $x_1, x_2, \dots, x_n \stackrel{IID}{\sim} \text{Exponential}(\lambda^*)$ is $1/\bar{x}_n$ and the ML estimator $\widehat{\Lambda}_n = 1/\bar{X}_n$.

Example 134 (Orbiter Waiting Times) Let us apply this ML estimator of the rate parameter for the supposedly exponentially distributed waiting times at the on-campus Orbiter bus-stop.

Joshua Fenemore and Yiran Wang collected data on waiting times between buses at an Orbiter bus-stop close to campus. They collected a sample of size $n = 132$ with sample mean $\bar{x}_{132} = 9.0758$.

```
% Joshu Fenemores Data from 2007 on Waiting Times at Orbiter Bust Stop by Balgay Street
%The raw data -- the waiting times to nearest minute between Orbiter buses
>> orbiterTimes=[8 3 7 18 18 3 7 9 9 25 0 0 25 6 10 0 10 8 16 9 1 5 16 6 4 1 3 21 0 28 3 8 ...
6 6 11 8 10 15 0 8 7 11 10 9 12 13 8 10 11 8 7 11 5 9 11 14 13 5 8 9 12 10 13 6 11 13 ...
0 0 11 1 9 5 14 16 2 10 21 1 14 2 10 24 6 1 14 14 0 14 4 11 15 0 10 2 13 2 22 ...
10 5 6 13 1 13 10 11 4 7 9 12 8 16 15 14 5 10 12 9 8 0 5 13 13 6 8 4 13 15 7 11 6 23 1];
>> mean(orbiterTimes)
ans =
9.0758
```

From our work in Example 133 we can now easily obtain the maximum likelihood estimate of λ^* and the 95% confidence interval for it, under the assumption that the waiting times X_1, \dots, X_{132} are IID $\text{Exponential}(\lambda^*)$ RVs as follows:

$$\widehat{\lambda}_{132} = 1/\bar{x}_{132} = 1/9.0758 = 0.1102 \quad (0.1102 \pm 1.96 \times 0.1102/\sqrt{132}) = (0.0914, 0.1290),$$

and thus the estimated mean waiting time is

$$1/\widehat{\lambda}_{132} = 9.0763 \text{ minutes.}$$

The estimated mean waiting time for a bus to arrive is well within the 10 minutes promised by the Orbiter bus company. This data and its maximum likelihood analysis is presented visually in Figure 4.17.

Figure 4.17: Plot of $\log(L(\lambda))$ as a function of the parameter λ and the MLE $\hat{\lambda}_{132}$ of 0.1102 for Fenemore-Wang Orbiter Waiting Times Experiment from STAT 218 S2 2007. The density or PDF and the DF at the MLE of 0.1102 are compared with a histogram and the empirical DF.

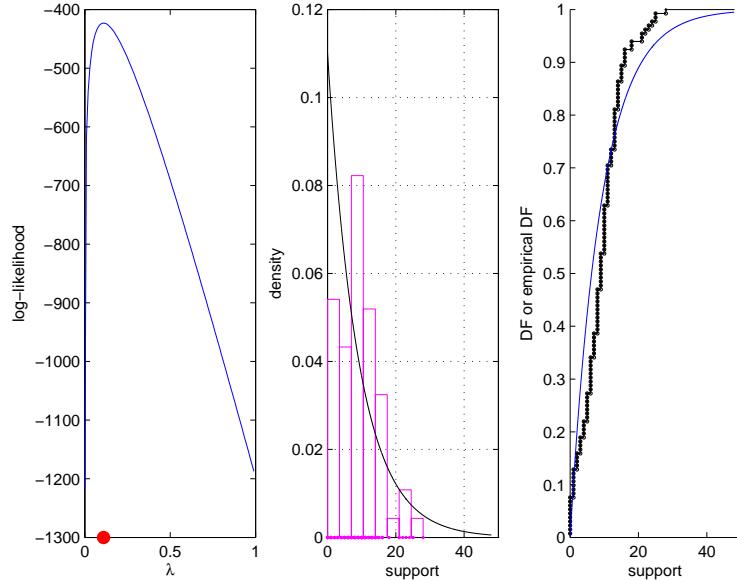
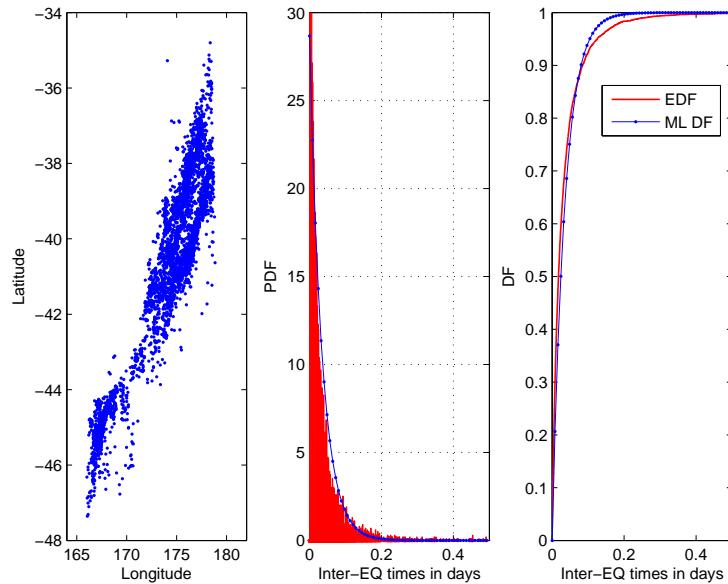


Figure 4.18: Comparing the Exponential($\hat{\lambda}_{6128} = 28.6694$) PDF and DF with a histogram and empirical DF of the times (in units of days) between earth quakes in NZ. The epicenters of 6128 earth quakes are shown in left panel.



The following script was used to generate the Figure 4.17: Notice how the exponential PDF $f(x; \hat{\lambda}_{132} = 0.1102)$ and the DF $F(x; \hat{\lambda}_{132} = 0.1102)$ based on the MLE fits with the histogram and the empirical DF, respectively.

Example 135 (Waiting Times between Earth Quakes in NZ) Once again from our work in Example 133 we can now easily obtain the maximum likelihood estimate of λ^* and the 95% confidence interval for it, under the assumption that the waiting times (in days) between the 6128 measurable earth-quakes in NZ from 18-Jan-2008 02:23:44 to 18-Aug-2008 19:29:29 are IID Exponential(λ^*) RVs as follows:

$$\hat{\lambda}_{6128} = 1/\bar{x}_{6128} = 1/0.0349 = 28.6694 \quad (28.6694 \pm 1.96 \times 28.6694/\sqrt{6128}) = (27.95, 29.39) ,$$

and thus the estimated mean time in days and minutes between earth quakes (somewhere in NZ over the first 8 months in 2008), as processed in Labwork 136, is

$$1/\hat{\lambda}_{6128} = \bar{x}_{6128} = 0.0349 \text{ days} = 0.0349 * 24 * 60 = 50.2560 \text{ minutes} .$$

This data and its maximum likelihood analysis is presented visually in Figure 4.18. The PDF and DF corresponding to the $\hat{\lambda}_{6128}$ (blue curves in Figure 4.18) are the best fitting PDF and DF from the parametric family of PDFs in $\{\lambda e^{-\lambda x} : \lambda \in (0, \infty)\}$ and DFs in $\{1 - e^{-\lambda x} : \lambda \in (0, \infty)\}$ to the density histogram and the empirical distribution function given by the data, respectively. Clearly, there is room for improving beyond the model of IID Exponential(λ) RVs, but the fit with just one real-valued parameter is not too bad either. Finally, with the best fitting PDF $28.6694e^{-28.6694x}$ we can get probabilities of events and answer questions like: “what is the probability that there will be three earth quakes somewhere in NZ within the next hour?”, etc.

Labwork 136 (Inter Earth Quake Time Processing) To process the data to get the times between earth quakes, we can compute as in the following script:

```
NZSEarthQuakesExponentialMLE.m
%% The columns in earthquakes.csv file have the following headings
%%CUSP_ID,LAT,LONG,NZMGE,NZMGN,ORI_YEAR,ORI_MONTH,ORI_DAY,ORI_HOUR,ORI_MINUTE,ORI_SECOND,MAG,DEPTH
EQ=dlmread('earthquakes.csv',','); % load the data in the matrix EQ
size(EQ) % report the size of the matrix EQ

% Read help datenum -- converts time stamps into numbers in units of days
MaxD=max(datenum(EQ(:,6:11))); % maximum datenum
MinD=min(datenum(EQ(:,6:11))); % minimum datenum
disp('Earth Quakes in NZ between')
disp(strcat(datestr(MinD), ' and ', datestr(MaxD))) % print MaxD and MinD as a date string

% get an array of sorted time stamps of EQ events starting at 0
Times=sort(datenum(EQ(:,6:11))-MinD);
TimeDiff=diff(Times); % inter-EQ times = times between successive EQ events
clf % clear any current figures
%figure
%plot(TimeDiff) % plot the inter-EQ times
subplot(1,3,1)
plot(EQ(:,3),EQ(:,2),'.')
axis([164 182 -48 -34])
xlabel('Longitude'); ylabel('Latitude');

subplot(1,3,2) % construct a histogram estimate of inter-EQ times
histogram(TimeDiff',1,[min(TimeDiff),max(TimeDiff)],'r',2);
SampleMean=mean(TimeDiff) % find the sample mean
```

```
% the MLE of LambdaStar if inter-EQ times are IID Exponential(LambdaStar)
MLELambdaHat=1/SampleMean
hold on;
TIMES=linspace(0,max(TimeDiff),100);
plot(TIMES,MLELambdaHat*exp(-MLELambdaHat*TIMES), 'b.-')
axis([0 0.5 0 30])
xlabel('Inter-EQ times in days'); ylabel('PDF');

subplot(1,3,3)
[x y]=ECDF(TimeDiff,0,0,0); % get the coordinates for empirical DF
stairs(x,y,'r','linewidth',1) % draw the empirical DF
hold on; plot(TIMES,ExponentialCdf(TIMES,MLELambdaHat), 'b.-');% plot the DF at MLE
axis([0 0.5 0 1])
xlabel('Inter-EQ times in days'); ylabel('DF'); legend('EDF', 'ML DF')
```

We first load the data in the text file `earthquakes.csv` into a matrix `EQ`. Using the `datenum` function in MATLAB we transform the time stamps into a number starting at zero. These transformed time stamps are in units of days. Then we find the times between consecutive events and estimate a histogram. We finally compute the ML estimate of λ^* and super-impose the PDF of the Exponential($\hat{\lambda}_{6128} = 28.6694$) upon the histogram.

```
>> NZSEarthQuakesExponentialMLE
ans =          6128          13

Earth Quakes in NZ between
18-Jan-2008 02:23:44 and 18-Aug-2008 19:29:29

SampleMean =    0.0349
MLELambdaHat =  28.6694
```

Thus, the average time between earth quakes is $0.0349 * 24 * 60 = 50.2560$ minutes.

Example 137 (ML Estimation for the IID Bernoulli(θ^*) experiment) Let us do maximum likelihood estimation for the coin-tossing experiment of Example 124 with likelihood derived in Example 132 to obtain the maximum likelihood estimate $\hat{\theta}_n$ of the unknown parameter $\theta^* \in \Theta = [0, 1]$ and the $(1 - \alpha)$ confidence interval for it.

From Equation (4.18) the log likelihood function is

$$\ell_n(\theta) = \log(L_n(\theta)) = \log\left(\theta^{\sum_{i=1}^n x_i} (1-\theta)^{n-\sum_{i=1}^n x_i}\right) = \left[\left(\sum_{i=1}^n x_i\right) \log(\theta) + \left(n - \sum_{i=1}^n x_i\right) \log(1-\theta)\right],$$

Next, we take the derivative with respect to the parameter θ :

$$\frac{d}{d\theta} (\ell_n(\theta)) = \frac{d}{d\theta} \left(\left(\sum_{i=1}^n x_i\right) \log(\theta) \right) + \frac{d}{d\theta} \left(\left(n - \sum_{i=1}^n x_i\right) \log(1-\theta) \right) = \left[\frac{\sum_{i=1}^n x_i}{\theta} - \frac{n - \sum_{i=1}^n x_i}{1-\theta} \right].$$

Now, set $\frac{d}{d\theta} \log(L_n(\theta)) = 0$, solve for θ and set the solution equal to $\hat{\theta}_n$:

$$\begin{aligned} \frac{d}{d\theta} (\ell_n(\theta)) = 0 &\iff \frac{\sum_{i=1}^n x_i}{\theta} = \frac{n - \sum_{i=1}^n x_i}{1-\theta} \iff \frac{1-\theta}{\theta} = \frac{n - \sum_{i=1}^n x_i}{\sum_{i=1}^n x_i} \\ &\iff \frac{1}{\theta} - 1 = \frac{n}{\sum_{i=1}^n x_i} - 1 \quad \text{let } \hat{\theta}_n = \frac{\sum_{i=1}^n x_i}{n} \end{aligned}$$

Next, we find the second derivative and check if it is negative.

$$\frac{d^2}{d\theta^2}(\ell_n(\theta)) = \frac{d}{d\theta} \left(\frac{\sum_{i=1}^n x_i}{\theta} - \frac{n - \sum_{i=1}^n x_i}{1-\theta} \right) = \boxed{-\frac{\sum_{i=1}^n x_i}{\theta^2} - \frac{n - \sum_{i=1}^n x_i}{(1-\theta)^2}}$$

Since each term in the numerator and the denominator of the two fractions in the above box are non-negative, $\frac{d^2}{d\theta^2}(\ell_n(\theta)) < 0$ and therefore we have found the maximum likelihood estimate

$$\hat{\theta}_n = \frac{1}{n} \sum_{i=1}^n x_i = \bar{x}_n$$

We already knew this to be a point estimate for $E(X_i) = \theta^*$ from LLN and CLT. But now we know that MLE also agrees. Now, let us find the estimated standard error:

$$\begin{aligned} \widehat{\text{se}}_n &= \sqrt{\left(\left[-\frac{d^2 \ell_n(\theta)}{d\theta^2} \right]_{\theta=\hat{\theta}_n} \right)^{-1}} = \sqrt{\left(\left[-\left(-\frac{\sum_{i=1}^n x_i}{\theta^2} - \frac{n - \sum_{i=1}^n x_i}{(1-\theta)^2} \right) \right]_{\theta=\hat{\theta}_n} \right)^{-1}} \\ &= \sqrt{\left(\frac{\sum_{i=1}^n x_i}{\hat{\theta}_n^2} + \frac{n - \sum_{i=1}^n x_i}{(1-\hat{\theta}_n)^2} \right)^{-1}} = \sqrt{\left(\frac{n\bar{x}_n}{\bar{x}_n^2} + \frac{n - n\bar{x}_n}{(1-\bar{x}_n)^2} \right)^{-1}} = \sqrt{\left(\frac{n}{\bar{x}_n} + \frac{n}{(1-\bar{x}_n)} \right)^{-1}} \\ &= \sqrt{\left(\frac{n(1-\bar{x}_n) + n\bar{x}_n}{\bar{x}_n(1-\bar{x}_n)} \right)^{-1}} = \sqrt{\frac{\bar{x}_n(1-\bar{x}_n)}{n((1-\bar{x}_n) + \bar{x}_n)}} = \boxed{\sqrt{\frac{\bar{x}_n(1-\bar{x}_n)}{n}}} . \end{aligned}$$

And finally, the $(1 - \alpha)$ confidence interval is

$$\hat{\theta}_n \pm z_{\alpha/2} \widehat{\text{se}}_n = \boxed{\bar{x}_n \pm z_{\alpha/2} \sqrt{\frac{\bar{x}_n(1-\bar{x}_n)}{n}}} .$$

For the coin tossing experiment that was performed ($n = 10$ times) in Example 124, the maximum likelihood estimate of θ^* and the 95% confidence interval for it, under the model that the tosses are IID Bernoulli(θ^*) RVs, are as follows:

$$\hat{\theta}_{10} = \bar{x}_{10} = \frac{4}{10} = 0.40 \quad \text{and} \quad \left(0.4 \pm 1.96 \times \sqrt{\frac{0.4 \times 0.6}{10}} \right) = (0.0964, 0.7036) .$$

See Figures 4.19 and 4.20 to completely understand parameter estimation for IID Bernoulli experiments.

Figure 4.19: Plots of the log likelihood $\ell_n(\theta) = \log(L(1, 0, 0, 0, 1, 1, 0, 0, 1, 0; \theta))$ as a function of the parameter θ over the parameter space $\Theta = [0, 1]$ and the MLE $\hat{\theta}_{10}$ of 0.4 for the coin-tossing experiment shown in standard scale (left panel) and log scale for x -axis (right panel).

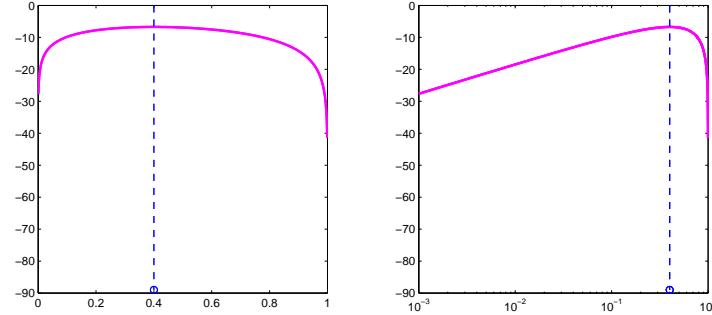
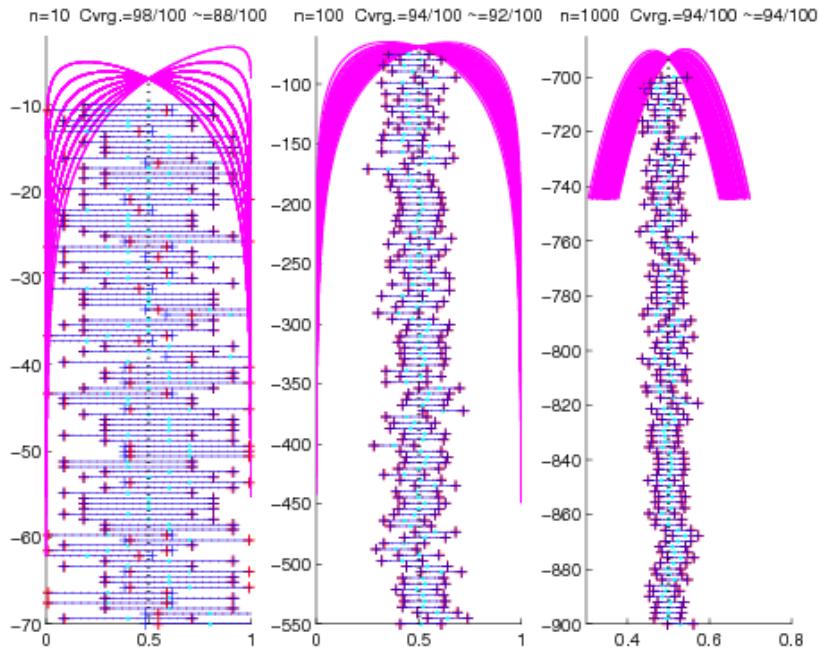


Figure 4.20: 100 realizations of 95% confidence intervals based on samples of size $n = 10, 100$ and 1000 simulated from IID Bernoulli($\theta^* = 0.5$) RVs. The MLE $\hat{\theta}_n$ (cyan dot) and the log-likelihood function (magenta curve) for each of the 100 replications of the experiment for each sample size n are depicted. The approximate normal-based 95% confidence intervals with blue boundaries are based on the exact $\text{se}_n = \sqrt{\theta^*(1 - \theta^*)/n} = \sqrt{1/4}$, while those with red boundaries are based on the estimated $\widehat{\text{se}}_n = \sqrt{\hat{\theta}_n(1 - \hat{\theta}_n)/n} = \sqrt{\bar{x}_n(1 - \bar{x}_n)/n}$. The fraction of times the true parameter $\theta^* = 0.5$ was contained by the exact and approximate confidence interval (known as *empirical coverage*) over the 100 replications of the simulation experiment for each of the three sample sizes are given by the numbers after $\text{Cvrg.} =$ and $\sim =$, above each sub-plot, respectively.



Exercise 4.2 (Likelihoods of tiny Bernoulli trials) Find and plot the likelihood function of the following observations (x_1, x_2, \dots, x_n) from the following IID sequence of $\text{Bernoulli}(\theta)$ RVs:

1. $(x_1) = (1)$
2. $(x_1) = (0)$
3. $(x_1, x_2) = (0, 0)$
4. $(x_1, x_2) = (1, 1)$
5. $(x_1, x_2) = (1, 0)$
6. $(x_1, x_2) = (0, 1)$
7. $(x_1, x_2, x_3) = (1, 1, 0)$
8. $(x_1, x_2, x_3) = (0, 0, 1)$

[Hint: your x-axis is θ with values in $[0, 1]$, the parameter space, and y-axis is $L_n(\theta; x_1, x_2, \dots, x_n) = \prod_{i=1}^n f_{X_i}(x_i; \theta)$, where $f_{X_i}(x_i; \theta)$ is the PMF of $\text{Bernoulli}(\theta)$ RV X_i]

Exercise 4.3 (MLE Exercises) Assume that an independent and identically distributed sample, X_1, X_2, \dots, X_n is drawn from the distribution of X with PDF $f(x; \theta^*)$ for a fixed and unknown parameter θ^* and derive the maximum likelihood estimate of θ^* (you only need to do Steps 1–5 from **Steps of MLE** in Lecture Notes on pages 126–127). Consider the following PDFs:

1. The parameter θ is a real number in $(0, \infty)$ and the PDF is given by

$$f(x; \theta) = \begin{cases} \theta x^{\theta-1} & \text{if } 0 < x < 1 \\ 0 & \text{otherwise} . \end{cases}$$

2. The parameter θ is a real number in $(0, \infty)$ and the PDF is given by

$$f(x; \theta) = \begin{cases} \frac{1}{\theta} x^{(1-\theta)/\theta} & \text{if } 0 < x < 1 \\ 0 & \text{otherwise} . \end{cases}$$

3. The parameter θ is a real number in $(0, \infty)$ and the PDF is given by

$$f(x; \theta) = \begin{cases} \frac{1}{2\theta^3} x^2 e^{-x/\theta} & \text{if } 0 < x < \infty \\ 0 & \text{otherwise} . \end{cases}$$

4. The parameter θ is a real number in $(0, \infty)$ and the PDF is given by

$$f(x; \theta) = \begin{cases} \frac{x}{\theta^2} e^{-\frac{1}{2}(x/\theta)^2} & \text{if } 0 < x < \infty \\ 0 & \text{otherwise} . \end{cases}$$

4.4.3 Moment Estimator (MME)

See notes from class.

4.5 Practical Excursion in One-dimensional Optimisation

Numerically maximising a log-likelihood function of one parameter is a useful technique. This can be used for models with no analytically known MLE. A fairly large field of maths, called optimisation, exists for this sole purpose. Conventionally, in optimisation, one is interested in minimisation. Therefore, the basic algorithms are cast in the “find the minimiser and the minimum” of a target function $f : \mathbb{R} \rightarrow \mathbb{R}$. Since we are interested in maximising our target, which is the likelihood or log-likelihood function, say $\log(L(x_1, \dots, x_n; \theta)) : \Theta \rightarrow \mathbb{R}$, we will simply apply the standard optimisation algorithms directly to $-\log(L(x_1, \dots, x_n; \theta)) : \Theta \rightarrow \mathbb{R}$.

The algorithm implemented in `fminbnd` is based on the golden section search and an inverse parabolic interpolation, and attempts to find the minimum of a function of one variable within a given fixed interval. Briefly, the golden section search proceeds by successively **bracketing** the minimum of the target function within an acceptably small interval inside the given starting interval [see Section 8.2 of Forsythe, G. E., M. A. Malcolm, and C. B. Moler, 1977, *Computer Methods for Mathematical Computations*, Prentice-Hall]. MATLAB’s `fminbnd` also relies on Brent’s inverse parabolic interpolation [see Chapter 5 of Brent, Richard. P., 1973, *Algorithms for Minimization without Derivatives*, Prentice-Hall, Englewood Cliffs, New Jersey]. Briefly, additional smoothness conditions are assumed for the target function to aid in a faster bracketing strategy through polynomial interpolations of past function evaluations. MATLAB’s `fminbnd` has several limitations, including:

- The likelihood function must be continuous.
- Only local MLE solutions, i.e. those inside the starting interval, are given.
- One needs to know or carefully guess the starting interval that contains the MLE.
- MATLAB’s `fminbnd` exhibits slow convergence when the solution is on a boundary of the starting interval.

Labwork 138 (Coin-tossing experiment) The following script was used to study the coin-tossing experiment in MATLAB. The plot of the log-likelihood function and the numerical optimisation of MLE are carried out using MATLAB’s built-in function `fminbnd` (See Figure 4.19).

BernoulliMLE.m

```
% To simulate n coin tosses, set theta=probability of heads and n
% Then draw n IID samples from Bernoulli(theta) RV
% theta=0.5; n=20; x=floor(rand(1,n) + theta);
% enter data from a real coin tossing experiment
x=[1 0 0 0 1 1 0 0 1 0]; n=length(x);
t = sum(x); % statistic t is the sum of the x_i values
% display the outcomes and their sum
display(x)
display(t)

% Analytically MLE is t/n
MLE=t/n
% l is the log-likelihood of data x as a function of parameter theta
l=@(theta)log(theta ^ t * (1-theta)^(n-t));
ThetaS=[0:0.001:1]; % sample some values for theta

% plot the log-likelihood function and MLE in two scales
subplot(1,2,1);
plot(ThetaS,arrayfun(l,ThetaS),'m','LineWidth',2);
hold on; stem([MLE],[-89],'b--'); % plot MLE as a stem plot
```

```

subplot(1,2,2);
semilogx(ThetaS,arrayfun(l,ThetaS),'m','LineWidth',2);
hold on; stem([MLE],[-89],'b--'); % plot MLE as a stem plot

% Now we will find the MLE by finding the minimiser or argmin of -l
% negative log-likelihood function of parameter theta
negl=@(theta)-(log(theta.^t * (1-theta).^(n-t)));
% read help fminbnd
% you need to supply the function to be minimised and its search interval
% NumericalMLE = fminbnd(negl,0,1)
% to see the iteration in the numerical minimisation
NumericalMLE = fminbnd(negl,0,1,optimset('Display','iter'))

```

```

>> BernoulliMLE
x = 1 0 0 0 1 1 0 0 1 0
t = 4
MLE = 0.4000
Func-count x f(x) Procedure
1 0.381966 6.73697 initial
2 0.618034 7.69939 golden
3 0.236068 7.3902 golden
4 0.408979 6.73179 parabolic
5 0.399339 6.73013 parabolic
6 0.400045 6.73012 parabolic
7 0.400001 6.73012 parabolic
8 0.399968 6.73012 parabolic
Optimisation terminated:
the current x satisfies the termination criteria using OPTIONS.TolX of 1.000000e-04
NumericalMLE = 0.4000

```

Labwork 139 (Numerical MLE of λ from n IID Exponential(λ) RVs) Joshua Fenemore and Yiran Wang collected data on waiting times between buses at an Orbiter bus-stop close to campus and modelled the waiting times as IID Exponential(λ^*) RVs (<http://www.math.canterbury.ac.nz/~r.sainudiin/courses/STAT218/projects/Stat218StudentProjects2007.pdf>). We can use their data `sampleTimes` to find the MLE of λ^* under the assumption that the waiting times X_1, \dots, X_{132} are IID Exponential(λ^*). We find the ML estimate $\hat{\lambda}_{132} = 0.1102$ and thus the estimated mean waiting time is $1/\hat{\lambda}_{132} = 9.0763$ minutes. The estimated mean waiting time for a bus to arrive is well within the 10 minutes promised by the Orbiter bus company. The following script was used to generate the Figure 4.17:

ExponentialMLEOrbiter.m

```

% Joshu Fenemore's Data from 2007 on Waiting Times at Orbiter Bust Stop
%The raw data -- the waiting times i minutes for each direction
antiTimes=[8 3 7 18 18 3 7 9 9 25 0 0 25 6 10 0 10 8 16 9 1 5 16 6 4 1 3 21 0 28 3 8 ...
6 6 11 8 10 15 0 8 7 11 10 9 12 13 8 10 11 8 7 11 5 9 11 14 13 5 8 9 12 10 13 6 11 13];
clockTimes=[0 0 11 1 9 5 14 16 2 10 21 1 14 2 10 24 6 1 14 14 0 14 4 11 15 0 10 2 13 2 22 ...
10 5 6 13 1 13 10 11 4 7 9 12 8 16 15 14 5 10 12 9 8 0 5 13 13 6 8 4 13 15 7 11 6 23 1];
sampleTimes=[antiTimes clockTimes];% pool all times into 1 array
% L = Log Likelihood of data x as a function of parameter lambda
L=@(lambda)sum(log(lambda*exp(-lambda * sampleTimes)));
LAMBDA=0.0001:0.01:1; % sample some values for lambda
clf;
subplot(1,3,1);
plot(LAMBDA,arrayfun(L,LAMBDA)); % plot the Log Likelihood function
% Now we will find the Maximum Likelihood Estimator by finding the minimizer of -L
MLE = fminbnd(@(lambda)-sum(log(lambda*exp(-lambda * sampleTimes))),0.0001,1)
MeanEstimate=1/MLE
hold on; % plot the MLE
plot([MLE],[-1300],'r.','MarkerSize',25); ylabel('log-likelihood'); xlabel('\lambda');

```

```

subplot(1,3,2); % plot a histogram estimate
histogram(sampleTimes,1,[min(sampleTimes),max(sampleTimes)],'m',2);
hold on; TIMES=[0.00001:0.01:max(sampleTimes)+20]; % points on support
plot(TIMES,MLE*exp(-MLE * TIMES ),'k-') % plot PDF at MLE to compare with histogram
% compare the empirical DF to the best fitted DF
subplot(1,3,3)
ECDF(sampleTimes,5,0.0,20); hold on
plot(TIMES,ExponentialCdf(TIMES,MLE),'b-')
ylabel('DF or empirical DF'); xlabel('support');

```

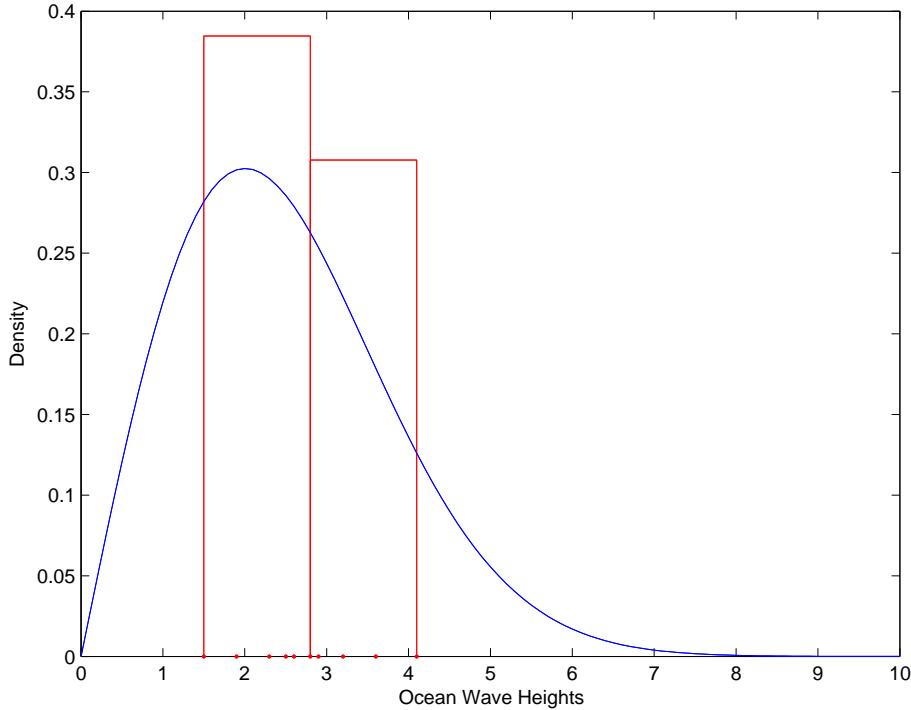
The script output the following in addition to the plot:

```

>> ExponentialMLEOrbiter
MLE =      0.1102
MeanEstimate =    9.0763

```

Figure 4.21: The ML fitted Rayleigh($\hat{\alpha}_{10} = 2$) PDF and a histogram of the ocean wave heights.



Example 140 (6.7, p. 275 of Ang & Tang) The distribution of ocean wave heights, H , may be modeled with the Rayleigh(α) RV with parameter α and probability density function,

$$f(h; \alpha) = \frac{h}{\alpha^2} \exp\left(-\frac{1}{2}(h/\alpha)^2\right), \quad h \in \mathbb{H} := [0, \infty) .$$

The parameter space for α is $\mathbb{A} = (0, \infty)$. Suppose that the following measurements h_1, h_2, \dots, h_{10} of wave heights in meters were observed to be

$$1.50, 2.80, 2.50, 3.20, 1.90, 4.10, 3.60, 2.60, 2.90, 2.30 ,$$

respectively. Under the assumption that the 10 samples are IID realisations from a Rayleigh(α^*) RV with a fixed and unknown parameter α^* , find the ML estimate $\hat{\alpha}_{10}$ of α^* .

We first obtain the log-likelihood function of α for the data $h_1, h_2, \dots, h_n \stackrel{IID}{\sim} \text{Rayleigh}(\alpha)$.

$$\begin{aligned}\ell(\alpha) &:= \log(L(h_1, h_2, \dots, h_n; \alpha)) = \log \left(\prod_{i=1}^n f(h_i; \alpha) \right) = \sum_{i=1}^n \log(f(h_i; \alpha)) \\ &= \sum_{i=1}^n \log \left(\frac{h_i}{\alpha^2} e^{-\frac{1}{2}(h_i/\alpha)^2} \right) = \sum_{i=1}^n \left(\log(h_i) - 2 \log(\alpha) - \frac{1}{2}(h_i/\alpha)^2 \right) \\ &= \sum_{i=1}^n (\log(h_i)) - 2n \log(\alpha) - \sum_{i=1}^n \left(\frac{1}{2} h_i^2 \alpha^{-2} \right)\end{aligned}$$

Now, let us take the derivative with respect to α ,

$$\begin{aligned}\frac{\partial}{\partial \alpha} (\ell(\alpha)) &:= \frac{\partial}{\partial \alpha} \left(\sum_{i=1}^n (\log(h_i)) - 2n \log(\alpha) - \sum_{i=1}^n \left(\frac{1}{2} h_i^2 \alpha^{-2} \right) \right) \\ &= \frac{\partial}{\partial \alpha} \left(\sum_{i=1}^n (\log(h_i)) \right) - \frac{\partial}{\partial \alpha} (2n \log(\alpha)) - \frac{\partial}{\partial \alpha} \left(\sum_{i=1}^n \left(\frac{1}{2} h_i^2 \alpha^{-2} \right) \right) \\ &= 0 - 2n \frac{1}{\alpha} - \sum_{i=1}^n \left(\frac{1}{2} h_i^2 (-2\alpha^{-3}) \right) = -2n\alpha^{-1} + \alpha^{-3} \sum_{i=1}^n (h_i^2)\end{aligned}$$

Next, we set the derivative to 0, solve for α , and set the solution equal to the ML estimate $\hat{\alpha}_n$.

$$\begin{aligned}0 = \frac{\partial}{\partial \alpha} (\ell(\alpha)) &\iff 0 = -2n\alpha^{-1} + \alpha^{-3} \sum_{i=1}^n h_i^2 \iff 2n\alpha^{-1} = \alpha^{-3} \sum_{i=1}^n h_i^2 \\ &\iff 2n\alpha^{-1}\alpha^3 = \sum_{i=1}^n h_i^2 \iff \alpha^2 = \frac{1}{2n} \sum_{i=1}^n h_i^2 \iff \hat{\alpha}_n = \sqrt{\frac{1}{2n} \sum_{i=1}^n h_i^2}\end{aligned}$$

Therefore, the ML estimate of the unknown $\alpha^* \in \mathbb{A}$ on the basis of our 10 observations h_1, h_2, \dots, h_{10} of wave heights is

$$\begin{aligned}\hat{\alpha}_{10} &= \sqrt{\frac{1}{2 * 10} \sum_{i=1}^{10} h_i^2} \\ &= \sqrt{\frac{1}{20} (1.50^2 + 2.80^2 + 2.50^2 + 3.20^2 + 1.90^2 + 4.10^2 + 3.60^2 + 2.60^2 + 2.90^2 + 2.30^2)} \approx 2\end{aligned}$$

We use the following script file to compute the MLE $\hat{\alpha}_{10}$ and plot the PDF at $\hat{\alpha}_{10}$ in Figure 4.21.

```
RayleighOceanHeightsMLE.m
OceanHeights=[1.50, 2.80, 2.50, 3.20, 1.90, 4.10, 3.60, 2.60, 2.90, 2.30];% data
histogram(OceanHeights,1,[min(OceanHeights),max(OceanHeights)],'r',2); % make a histogram
Heights=0:0.1:10; % get some heights for plotting
AlphaHat=sqrt(sum(OceanHeights.^2)/(2*length(OceanHeights))) % find the MLE
hold on; % superimpose the PDF at the MLE
plot(Heights,(Heights/AlphaHat.^2).* exp(-((Heights/AlphaHat).^2)/2))
xlabel('Ocean Wave Heights'); ylabel('Density');
```

```
>> RayleighOceanHeightsMLE
AlphaHat = 2.0052
```

4.6 More Properties of the Maximum Likelihood Estimator

Next, we list some nice properties of the ML Estimator $\hat{\Theta}_n$ for the fixed and possibly unknown $\theta^* \in \Theta$.

1. The ML Estimator is asymptotically consistent, i.e. $\hat{\Theta}_n \xrightarrow{P} \theta^*$.
2. The ML Estimator is asymptotically normal, i.e. $(\hat{\Theta}_n - \theta^*)/\hat{s}\hat{e}_n \rightsquigarrow \text{Normal}(0, 1)$.
3. The estimated standard error of the ML Estimator, $\hat{s}\hat{e}_n$, can usually be computed analytically using the **Fisher Information**.
4. Because of the previous two properties, the $1 - \alpha$ confidence interval can also be computed analytically as $\hat{\Theta}_n \pm z_{\alpha/2}\hat{s}\hat{e}_n$.
5. The ML Estimator is **equivariant**, i.e. $\hat{\psi}_n = g(\hat{\theta}_n)$ is the ML Estimate of $\psi^* = g(\theta^*)$, for some smooth function $g(\theta) = \psi : \Theta \rightarrow \Psi$.
6. We can also obtain the estimated standard error of the estimator $\hat{\Psi}_n$ of $\psi^* \in \Psi$ via the **Delta Method**.
7. The ML Estimator is **asymptotically optimal** or **efficient**. This means that the MLE has the smallest variance among the well-behaved class of estimators as the sample size gets larger.
8. ML Estimator is close to the Bayes estimator (obtained in the Bayesian inferential paradigm).

4.7 Fisher Information

Let $X_1, X_2, \dots, X_n \stackrel{IID}{\sim} f(X_1; \theta)$. Here, $f(X_1; \theta)$ is the probability density function (pdf) or the probability mass function (pmf) of the RV X_1 . Since all RVs are identically distributed, we simply focus on X_1 without loss of generality.

Definition 71 (Fisher Information) The **score function** of an RV X for which the density is parameterised by θ is defined as:

$$\mathcal{S}(X; \theta) := \frac{\partial}{\partial \theta} \log f(X; \theta), \quad \text{and} \quad E_\theta(\mathcal{S}(X; \theta)) = 0.$$

The expectation of the score function is 0 if X is distributed according to $f(x; \theta)$ and under regularity conditions that allow for the interchange of differentiation and integration or summation, as shown below:

$$\begin{aligned} E_\theta(\mathcal{S}(X; \theta)) &= E_\theta \left(\frac{\partial}{\partial \theta} \log f(X; \theta) \right) = \int \frac{\frac{\partial}{\partial \theta} f(x; \theta)}{f(x; \theta)} dF(x; \theta) \\ &= \begin{cases} \sum_{x \in \mathbb{X}} \left(\frac{\frac{\partial}{\partial \theta} f(x; \theta)}{f(x; \theta)} \right) f(x; \theta) = \frac{\partial}{\partial \theta} \left(\sum_{x \in \mathbb{X}} f(x; \theta) \right) = \frac{\partial}{\partial \theta}(1) = 0 & \text{for discrete } X \\ \int_{x \in \mathbb{X}} \left(\frac{\frac{\partial}{\partial \theta} f(x; \theta)}{f(x; \theta)} \right) f(x; \theta) dx = \frac{\partial}{\partial \theta} \left(\int_{x \in \mathbb{X}} f(x; \theta) dx \right) = \frac{\partial}{\partial \theta}(1) = 0 & \text{for continuous } X \end{cases} \end{aligned}$$

The **Fisher Information** is simply the variance of the score function.

$$I_n := V_\theta \left(\sum_{i=1}^n S(X_i; \theta) \right) = \sum_{i=1}^n V_\theta(S(X_i; \theta)) = n I_1(\theta), \quad (4.19)$$

where I_1 is the Fisher Information of just one of the RVs X_i , say for e.g. X :

$$\begin{aligned} I_1(\theta) &:= V_\theta(S(X; \theta)) = E_\theta(S^2(X, \theta)) - (E_\theta S(X, \theta))^2 = E_\theta(S^2(X, \theta)) - 0^2 \\ &= E_\theta(S^2(X, \theta)) = \begin{cases} \sum_{x \in \mathbb{X}} \left(\frac{\partial}{\partial \theta} \log(f(x; \theta)) \right)^2 f(x; \theta) & \text{for discrete } X \\ \int_{x \in \mathbb{X}} \left(\frac{\partial}{\partial \theta} \log(f(x; \theta)) \right)^2 f(x; \theta) dx & \text{for continuous } X \end{cases} \end{aligned} \quad (4.20)$$

If $\log(f(x; \theta))$ is twice differentiable with respect to θ and under further regularity conditions, we can also express Fisher Information as the negative of the expected curvature of the log likelihood function:

$$I_1(\theta) = -E_\theta \left(\frac{\partial^2 \log f(X; \theta)}{\partial \theta^2} \right) = \begin{cases} -\sum_{x \in \mathbb{X}} \left(\frac{\partial^2 \log f(x; \theta)}{\partial \theta^2} \right) f(x; \theta) & \text{for discrete } X \\ -\int_{x \in \mathbb{X}} \left(\frac{\partial^2 \log f(x; \theta)}{\partial \theta^2} \right) f(x; \theta) dx & \text{for continuous } X \end{cases} \quad (4.21)$$

Note that (4.21) is due to taking expectation, $E_\theta(\cdot)$, of the LHS and RHS of the following equalities:

$$\begin{aligned} \left(\frac{\partial^2}{\partial \theta^2} \log(f(x; \theta)) \right) &= \frac{\partial}{\partial \theta} \left(\frac{\partial}{\partial \theta} \log(f(x; \theta)) \right) = \frac{\partial}{\partial \theta} \left(\frac{\frac{\partial}{\partial \theta} f(x; \theta)}{f(x; \theta)} \right) \\ &= \frac{\frac{\partial^2}{\partial \theta^2} f(x; \theta)}{f(x; \theta)} - \left(\frac{\frac{\partial}{\partial \theta} f(x; \theta)}{f(x; \theta)} \right)^2 = \frac{\frac{\partial^2}{\partial \theta^2} f(x; \theta)}{f(x; \theta)} - \left(\frac{\partial}{\partial \theta} \log(f(x; \theta)) \right)^2 \end{aligned}$$

and due to the E_θ of the second-last term above being 0 and that of the last term being I_1 as per (4.20).

Next, we give a **general method** for obtaining:

1. The standard error $se_n(\hat{\Theta}_n)$ of **any** maximum likelihood estimator $\hat{\Theta}_n$ of the possibly unknown and fixed parameter of interest $\theta^* \in \Theta$, and
2. The $1 - \alpha$ confidence interval for θ^* .

Proposition 72 (Asymptotic Normality of the ML Estimator & Confidence Intervals)

Let $\hat{\Theta}_n$ be the maximum likelihood estimator of $\theta^* \in \Theta$ with standard error $se_n := \sqrt{V_{\theta^*}(\hat{\Theta}_n)}$. Under appropriate regularity conditions, the following propositions are true:

1. The standard error se_n can be approximated by the side of a square whose area is the inverse Fisher Information at θ^* , and the distribution of $\hat{\Theta}_n$ approaches that of the $\text{Normal}(\theta^*, se_n^2)$ distribution as the samples size n gets larger. In other terms:

$$se_n \approx \sqrt{1/I_n(\theta^*)} \quad \text{and} \quad \frac{\hat{\Theta}_n - \theta^*}{se_n} \rightsquigarrow \text{Normal}(0, 1)$$

2. The approximation holds even if we substitute the ML Estimate $\hat{\theta}_n$ for θ^* and use the estimated standard error $\hat{s}\hat{e}_n$ instead of $s\hat{e}_n$. Let $\hat{s}\hat{e}_n = \sqrt{1/I_n(\hat{\theta}_n)}$. Then:

$$\frac{\hat{\Theta}_n - \theta^*}{\hat{s}\hat{e}_n} \rightsquigarrow \text{Normal}(0, 1)$$

3. Using the fact that $\hat{\Theta}_n \rightsquigarrow \text{Normal}(\theta^*, \hat{s}\hat{e}_n^2)$, we can construct the estimate of an approximate Normal-based $1 - \alpha$ confidence interval as:

$$C_n = [C_n, \bar{C}_n] = [\hat{\theta}_n - z_{\alpha/2} \hat{s}\hat{e}_n, \hat{\theta}_n + z_{\alpha/2} \hat{s}\hat{e}_n] = \hat{\theta}_n \pm z_{\alpha/2} \hat{s}\hat{e}_n$$

Now, let us do an example.

Example 141 (MLE and Confidence Interval for the IID Poisson(λ) experiment) Suppose the fixed parameter $\lambda^* \in \Lambda = (0, \infty)$ is unknown. Let $X_1, X_2, \dots, X_n \stackrel{IID}{\sim} \text{Poisson}(\lambda^*)$. We want to find the ML Estimate $\hat{\lambda}_n$ of λ^* and produce a $1 - \alpha$ confidence interval for λ^* .

The MLE can be obtained as follows:

The likelihood function is:

$$L(\lambda) := L(x_1, x_2, \dots, x_n; \lambda) = \prod_{i=1}^n f(x_i; \lambda) = \prod_{i=1}^n e^{-\lambda} \frac{\lambda^{x_i}}{x_i!}$$

Hence, the log-likelihood function is:

$$\begin{aligned} \ell(\theta) := \log(L(\lambda)) &= \log \left(\prod_{i=1}^n e^{-\lambda} \frac{\lambda^{x_i}}{x_i!} \right) = \sum_{i=1}^n \log \left(e^{-\lambda} \frac{\lambda^{x_i}}{x_i!} \right) = \sum_{i=1}^n (\log(e^{-\lambda}) + \log(\lambda^{x_i}) - \log(x_i!)) \\ &= \sum_{i=1}^n (-\lambda + x_i \log(\lambda) - \log(x_i!)) = \sum_{i=1}^n -\lambda + \sum_{i=1}^n x_i \log(\lambda) - \sum_{i=1}^n \log(x_i!) \\ &= n(-\lambda) + \log(\lambda) \left(\sum_{i=1}^n x_i \right) - \sum_{i=1}^n \log(x_i!) \end{aligned}$$

Next, take the derivative of $\ell(\lambda)$:

$$\frac{\partial}{\partial \lambda} \ell(\lambda) = \frac{\partial}{\partial \lambda} \left(n(-\lambda) + \log(\lambda) \left(\sum_{i=1}^n x_i \right) - \sum_{i=1}^n \log(x_i!) \right) = n(-1) + \frac{1}{\lambda} \left(\sum_{i=1}^n x_i \right) + 0$$

and set it equal to 0 to solve for λ , as follows:

$$0 = n(-1) + \frac{1}{\lambda} \left(\sum_{i=1}^n x_i \right) + 0 \iff n = \frac{1}{\lambda} \left(\sum_{i=1}^n x_i \right) \iff \lambda = \frac{1}{n} \left(\sum_{i=1}^n x_i \right) = \bar{x}_n$$

Finally, the ML Estimator of λ^* is $\hat{\lambda}_n = \bar{X}_n$ and the ML estimate is $\hat{\lambda}_n = \bar{x}_n$.

Now, we want an $1 - \alpha$ confidence interval for λ^* using the $\hat{s}\hat{e}_n \approx \sqrt{1/I_n(\hat{\lambda}_n)}$ that is based on the Fisher Information $I_n(\lambda) = nI_1(\lambda)$ given in (4.19). We need I_1 given in (4.21). Since $X_1, X_2, \dots, X_n \sim \text{Poisson}(\lambda)$, we have discrete RVs:

$$I_1 = - \sum_{x \in \mathbb{X}} \left(\frac{\partial^2 \log(f(x; \lambda))}{\partial \lambda^2} \right) f(x; \lambda) = - \sum_{x=0}^{\infty} \left(\frac{\partial^2 \log(f(x; \lambda))}{\partial \lambda^2} \right) f(x; \lambda)$$

First find

$$\begin{aligned}\frac{\partial^2 \log(f(x; \lambda))}{\partial \lambda^2} &= \frac{\partial}{\partial \lambda} \left(\frac{\partial}{\partial \lambda} \log(f(x; \lambda)) \right) = \frac{\partial}{\partial \lambda} \left(\frac{\partial}{\partial \lambda} \log \left(e^{-\lambda} \frac{\lambda^x}{x!} \right) \right) \\ &= \frac{\partial}{\partial \lambda} \left(\frac{\partial}{\partial \lambda} (-\lambda + x \log(\lambda) - \log(x!)) \right) = \frac{\partial}{\partial \lambda} \left(-1 + \frac{x}{\lambda} - 0 \right) = -\frac{x}{\lambda^2}\end{aligned}$$

Now, substitute the above expression into the right-hand side of I_1 to obtain:

$$I_1 = - \sum_{x=0}^{\infty} \left(-\frac{x}{\lambda^2} \right) f(x; \lambda) = \frac{1}{\lambda^2} \sum_{x=0}^{\infty} (x) f(x; \lambda) = \frac{1}{\lambda^2} \sum_{x=0}^{\infty} (x) e^{-\lambda} \frac{\lambda^x}{x!} = \frac{1}{\lambda^2} \mathbb{E}_{\lambda}(X) = \frac{1}{\lambda^2} \lambda = \frac{1}{\lambda}$$

In the third-to-last step above, we recognise the sum as the expectation of the Poisson(λ) RV X , namely $\mathbb{E}_{\lambda}(X) = \lambda$. Therefore, the estimated standard error is:

$$\hat{s}\hat{e}_n \approx \sqrt{1/I_n(\hat{\lambda}_n)} = \sqrt{1/(nI_1(\hat{\lambda}_n))} = \sqrt{1/(n(1/\hat{\lambda}_n))} = \sqrt{\hat{\lambda}_n/n}$$

and the approximate $1 - \alpha$ confidence interval is

$$\hat{\lambda}_n \pm z_{\alpha/2} \hat{s}\hat{e}_n = \hat{\lambda}_n \pm z_{\alpha/2} \sqrt{\hat{\lambda}_n/n}$$

Thus, using the MLE and the estimated standard error via the Fisher Information, we can carry out point estimation and confidence interval construction in **most** parametric families of RVs encountered in typical engineering applications.

Example 142 (Fisher Information of the Bernoulli Experiment) Suppose $X_1, X_2, \dots, X_n \stackrel{IID}{\sim} \text{Bernoulli}(\theta^*)$. Also, suppose that $\theta^* \in \Theta = [0, 1]$ is unknown. We have already shown in Example 137 that the ML estimator of θ^* is $\hat{\theta}_n = \bar{X}_n$. Using the identity:

$$\hat{s}\hat{e}_n = \frac{1}{\sqrt{I_n(\hat{\theta}_n)}}$$

(1) we can compute $\hat{s}\hat{e}_n(\hat{\theta}_n)$, the estimated standard error of the unknown parameter θ^* as follows:

$$\hat{s}\hat{e}_n(\hat{\theta}_n) = \frac{1}{\sqrt{I_n(\hat{\theta}_n)}} = \frac{1}{\sqrt{nI_1(\hat{\theta}_n)}}.$$

So, we need to first compute $I_1(\theta)$, the Fisher Information of one sample. Due to (4.21) and the fact that the $\text{Bernoulli}(\theta^*)$ distributed RV X is discrete with probability mass function $f(x; \theta) = \theta^x (1 - \theta)^{1-x}$, for $x \in \mathbb{X} := \{0, 1\}$, we have,

$$I_1(\theta) = -\mathbb{E}_{\theta} \left(\frac{\partial^2 \log f(X; \theta)}{\partial \theta^2} \right) = - \sum_{x \in \mathbb{X}=\{0,1\}} \left(\frac{\partial^2 \log (\theta^x (1 - \theta)^{1-x})}{\partial \theta^2} \right) \theta^x (1 - \theta)^{1-x}$$

Next, let us compute,

$$\begin{aligned}\frac{\partial^2 \log (\theta^x (1 - \theta)^{1-x})}{\partial \theta^2} &:= \frac{\partial}{\partial \theta} \left(\frac{\partial}{\partial \theta} (\log (\theta^x (1 - \theta)^{1-x})) \right) = \frac{\partial}{\partial \theta} \left(\frac{\partial}{\partial \theta} (x \log(\theta) + (1 - x) \log(1 - \theta)) \right) \\ &= \frac{\partial}{\partial \theta} (x\theta^{-1} + (1 - x)(1 - \theta)^{-1}(-1)) = \frac{\partial}{\partial \theta} (x\theta^{-1} - (1 - x)(1 - \theta)^{-1}) \\ &= x(-1)\theta^{-1-1} - (1 - x)(-1)(1 - \theta)^{-1-1}(-1) = -x\theta^{-2} - (1 - x)(1 - \theta)^{-2}\end{aligned}$$

Now, we compute the expectation I_1 , i.e. the sum over the two possible values of $x \in \{0, 1\}$,

$$\begin{aligned} I_1(\theta) &= - \sum_{x \in \mathbb{X}=\{0,1\}} \left(\frac{\partial^2 \log (\theta^x (1-\theta)^{1-x})}{\partial \theta^2} \right) \theta^x (1-\theta)^{1-x} \\ &= - ((-0 \theta^{-2} - (1-0)(1-\theta)^{-2}) \theta^0 (1-\theta)^{1-0} + (-1 \theta^{-2} - (1-1)(1-\theta)^{-2}) \theta^1 (1-\theta)^{1-1}) \\ &= - ((0 - 1(1-\theta)^{-2}) 1 (1-\theta)^1 + (-\theta^{-2} - 0) \theta^1 1) = (1-\theta)^{-2} (1-\theta)^1 + \theta^{-2} \theta^1 \\ &= (1-\theta)^{-1} + \theta^{-1} = \frac{1}{1-\theta} + \frac{1}{\theta} = \frac{\theta}{\theta(1-\theta)} + \frac{1-\theta}{\theta(1-\theta)} = \frac{\theta + (1-\theta)}{\theta(1-\theta)} = \frac{1}{\theta(1-\theta)} \end{aligned}$$

Therefore, the desired estimated standard error of our estimator, can be obtained by substituting the ML estimate $\hat{\theta}_n = \bar{x}_n := n^{-1} \sum_{i=1}^n x_i$ of the unknown θ^* as follows:

$$\widehat{\text{se}}_n(\hat{\theta}_n) = \frac{1}{\sqrt{I_n(\hat{\theta}_n)}} = \frac{1}{\sqrt{n I_1(\hat{\theta}_n)}} = \sqrt{\frac{1}{n \frac{1}{\hat{\theta}_n(1-\hat{\theta}_n)}}} = \sqrt{\frac{\hat{\theta}_n(1-\hat{\theta}_n)}{n}} = \sqrt{\frac{\bar{x}_n(1-\bar{x}_n)}{n}}.$$

(2) Using $\widehat{\text{se}}_n(\hat{\theta}_n)$ we can construct an approximate 95% confidence interval C_n for θ^* , due to the asymptotic normality of the ML estimator of θ^* , as follows:

$$C_n = \hat{\theta}_n \pm 1.96 \sqrt{\frac{\hat{\theta}_n(1-\hat{\theta}_n)}{n}} = \bar{x}_n \pm 1.96 \sqrt{\frac{\bar{x}_n(1-\bar{x}_n)}{n}}$$

Recall that C_n is the realisation of a random set based on your observed samples or data x_1, x_2, \dots, x_n . Furthermore, C_n 's construction procedure ensures the engulfing of the unknown θ^* with probability approaching 0.95 as the sample size n gets large.

Example 143 ([Fisher Information of the Exponential Experiment]) Let us get our hands dirty with a continuous RV next. Let $X_1, X_2, \dots, X_n \stackrel{IID}{\sim} \text{Exponential}(\lambda^*)$. We saw that the ML estimator of $\lambda^* \in \Lambda = (0, \infty)$ is $\hat{\lambda}_n = 1/\bar{X}_n$ and its ML estimate is $\hat{\lambda}_n = 1/\bar{x}_n$, where x_1, x_2, \dots, x_n are our observed data.

(1) Let us obtain the Fisher Information I_n for this experiment to find the standard error:

$$\widehat{\text{se}}_n(\hat{\lambda}_n) = \frac{1}{\sqrt{I_n(\hat{\lambda}_n)}} = \frac{1}{\sqrt{n I_1(\hat{\lambda}_n)}}$$

and construct an approximate 95% confidence interval for λ^* using the asymptotic normality of its ML estimator $\hat{\lambda}_n$.

So, we need to first compute $I_1(\theta)$, the Fisher Information of one sample. Due to (4.21) and the fact that the $\text{Exponential}(\lambda^*)$ distributed RV X is continuous with probability density function $f(x; \lambda) = \lambda e^{-\lambda x}$, for $x \in \mathbb{X} := [0, \infty)$, we have,

$$I_1(\theta) = -E_\theta \left(\frac{\partial^2 \log f(X; \theta)}{\partial \theta^2} \right) = - \int_{x \in \mathbb{X}=[0, \infty)} \left(\frac{\partial^2 \log (\lambda e^{-\lambda x})}{\partial \lambda^2} \right) \lambda e^{-\lambda x} dx$$

Let us compute the above integrand next.

$$\begin{aligned} \frac{\partial^2 \log (\lambda e^{-\lambda x})}{\partial \lambda^2} &:= \frac{\partial}{\partial \lambda} \left(\frac{\partial}{\partial \lambda} \left(\log (\lambda e^{-\lambda x}) \right) \right) = \frac{\partial}{\partial \lambda} \left(\frac{\partial}{\partial \lambda} \left(\log(\lambda) + \log(e^{-\lambda x}) \right) \right) \\ &= \frac{\partial}{\partial \lambda} \left(\frac{\partial}{\partial \lambda} (\log(\lambda) - \lambda x) \right) = \frac{\partial}{\partial \lambda} (\lambda^{-1} - x) = -\lambda^{-2} - 0 = -\frac{1}{\lambda^2} \end{aligned}$$

Now, let us evaluate the integral by recalling that the expectation of the constant 1 is 1 for any RV X governed by some parameter, say θ . For instance when X is a continuous RV, $E_\theta(1) = \int_{x \in \mathbb{X}} 1 f(x; \theta) = \int_{x \in \mathbb{X}} f(x; \theta) = 1$. Therefore, the Fisher Information of one sample is

$$\begin{aligned} I_1(\theta) &= - \int_{x \in \mathbb{X}=[0, \infty)} \left(\frac{\partial^2 \log(\lambda e^{-\lambda x})}{\partial \lambda^2} \right) \lambda e^{-\lambda x} dx = - \int_0^\infty \left(-\frac{1}{\lambda^2} \right) \lambda e^{-\lambda x} dx \\ &= - \left(-\frac{1}{\lambda^2} \right) \int_0^\infty \lambda e^{-\lambda x} dx = \frac{1}{\lambda^2} 1 = \frac{1}{\lambda^2} \end{aligned}$$

Now, we can compute the desired estimated standard error, by substituting in the ML estimate $\hat{\lambda}_n = 1/(\bar{x}_n) := 1/(\sum_{i=1}^n x_i)$ of λ^* , as follows:

$$\widehat{\text{se}}_n(\hat{\lambda}_n) = \frac{1}{\sqrt{I_n(\hat{\lambda}_n)}} = \frac{1}{\sqrt{n I_1(\hat{\lambda}_n)}} = \frac{1}{\sqrt{n \frac{1}{\hat{\lambda}_n^2}}} = \frac{\hat{\lambda}_n}{\sqrt{n}} = \frac{1}{\sqrt{n} \bar{x}_n}$$

Using $\widehat{\text{se}}_n(\hat{\lambda}_n)$ we can construct an approximate 95% confidence interval C_n for λ^* , due to the asymptotic normality of the ML estimator of λ^* , as follows:

$$C_n = \hat{\lambda}_n \pm 1.96 \frac{\hat{\lambda}_n}{\sqrt{n}} = \frac{1}{\bar{x}_n} \pm 1.96 \frac{1}{\sqrt{n} \bar{x}_n} .$$

Let us compute the ML estimate and the 95% confidence interval for the rate parameter for the waiting times at the Orbiter bus-stop (see labwork 139). The sample mean $\bar{x}_{132} = 9.0758$ and the ML estimate is:

$$\hat{\lambda}_{132} = 1/\bar{x}_{132} = 1/9.0758 = 0.1102 ,$$

and the 95% confidence interval is:

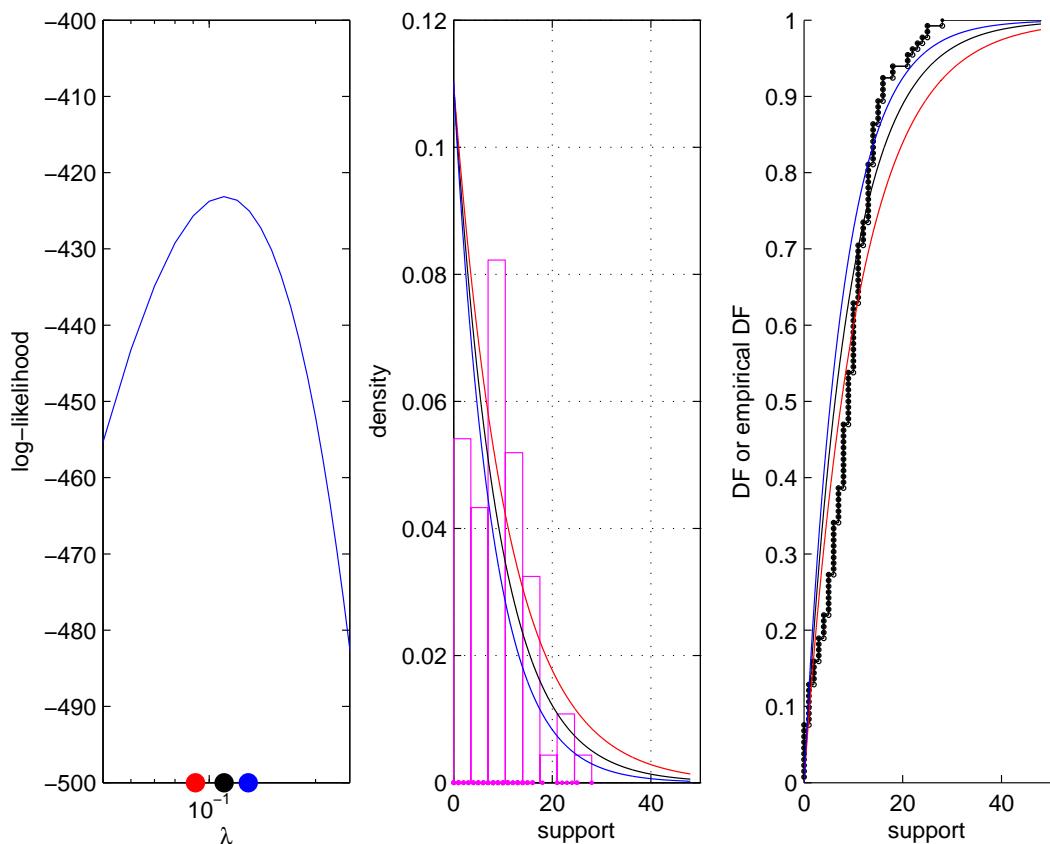
$$C_n = \hat{\lambda}_{132} \pm 1.96 \frac{\hat{\lambda}_{132}}{\sqrt{132}} = \frac{1}{\bar{x}_{132}} \pm 1.96 \frac{1}{\sqrt{132} \bar{x}_{132}} = 0.1102 \pm 1.96 \cdot 0.0096 = [0.0914, 0.1290] .$$

Notice how poorly the exponential PDF $f(x; \hat{\lambda}_{132} = 0.1102)$ and the DF $F(x; \hat{\lambda}_{132} = 0.1102)$ based on the MLE fits with the histogram and the empirical DF, respectively, in Figure 4.22, despite taking the the confidence interval into account. This is a further indication of the inadequacy of our parametric model.

Labwork 144 (Maximum likelihood estimation for Orbiter bus-stop) The above analysis was undertaken with the following M-file:

```
ExponentialMLECIOrbiter.m
OrbiterData; % load the Orbiter Data sampleTimes
% L = Log Likelihood of data x as a function of parameter lambda
L=@(lambda)sum(log(lambda*exp(-lambda * sampleTimes)));
LAMBDA=[0.01:0.01:1]; % sample some values for lambda
clf;
subplot(1,3,1);
semilogx(LAMBDA,arrayfun(L,LAMBDA)); % plot the Log Likelihood function
axis([0.05 0.25 -500 -400])
SampleMean = mean(sampleTimes) % sample mean
MLE = 1/SampleMean % ML estimate is 1/mean
n=length(sampleTimes) % sample size
StdErr=1/(sqrt(n)*SampleMean) % Standard Error
```

Figure 4.22: Plot of $\log(L(\lambda))$ as a function of the parameter λ , the MLE $\hat{\lambda}_{132} = 0.1102$ and 95% confidence interval $C_n = [0.0914, 0.1290]$ for Fenemore-Wang Orbiter Waiting Times Experiment from STAT 218 S2 2007. The PDF and the DF at (1) the MLE 0.1102 (black), (2) lower 95% confidence bound 0.0914 (red) and (3) upper 95% confidence bound 0.1290 (blue) are compared with a histogram and the empirical DF.



```

MLE95CI=[MLE-(1.96*StdErr), MLE+(1.96*StdErr)] % 95 % CI
hold on; % plot the MLE
plot([MLE], [-500], 'k.', 'MarkerSize', 25);
plot([MLE95CI(1)], [-500], 'r.', 'MarkerSize', 25);
plot([MLE95CI(2)], [-500], 'b.', 'MarkerSize', 25);
ylabel('log-likelihood'); xlabel('\lambda');
subplot(1,3,2); % plot a histogram estimate
histogram(sampleTimes, 1, [min(sampleTimes), max(sampleTimes)], 'm', 2);
hold on; TIMES=[0.00001:0.01:max(sampleTimes)+20]; % points on support
% plot PDF at MLE and 95% CI to compare with histogram
plot(TIMES,MLE*exp(-MLE*TIMES), 'k-');
plot(TIMES,MLE*exp(-MLE95CI(1)*TIMES), 'r-'); plot(TIMES,MLE*exp(-MLE95CI(2)*TIMES), 'b-')
% compare the empirical DF to the best fitted DF at MLE and 95% CI
subplot(1,3,3)
ECDF(sampleTimes,5,0.0,20); hold on; plot(TIMES,ExponentialCdf(TIMES,MLE), 'k-');
plot(TIMES,ExponentialCdf(TIMES,MLE95CI(1)), 'r-'); plot(TIMES,ExponentialCdf(TIMES,MLE95CI(2)), 'b-')
ylabel('DF or empirical DF'); xlabel('support');

```

A call to the script generates Figure 4.22 and the following output of the sample mean, MLE, sample size, standard error and the 95% confidence interval.

```

>> ExponentialMLECIOrbiter
SampleMean =      9.0758
MLE =      0.1102
n =      132
StdErr =     0.0096
MLE95CI =    0.0914    0.1290

```

Labwork 145 (Maximum likelihood estimation for your bus-stop) Recall labwork 161 where you modeled the arrival of buses using $\text{Exponential}(\lambda^* = 0.1)$ distributed inter-arrival time with a mean of $1/\lambda^* = 10$ minutes. Using the data of these seven inter-arrival times at your ID-seeded bus stop and pretending that you do not know the true λ^* , report (1) the ML estimate of λ^* , (2) 95% confidence interval for it and (3) whether the true value $\lambda^* = 1/10$ is engulfed by your confidence interval.

4.8 Delta Method

A more general estimation problem of interest concerns some function of the parameter $\theta \in \Theta$, say $g(\theta) = \psi : \Theta \rightarrow \Psi$. So, $g(\theta) = \psi$ is a function from the parameter space Θ to Ψ . Thus, we are not only interested in estimating the fixed and possibly unknown $\theta^* \in \Theta$ using the ML estimator $\hat{\Theta}_n$ and its ML estimate $\hat{\theta}_n$, but also in estimating $\psi^* = g(\theta^*) \in \Psi$ via an estimator $\hat{\Psi}_n$ and its estimate $\hat{\psi}_n$. We exploit the equivariance property of the ML estimator $\hat{\Theta}_n$ of θ^* and use the Delta method to find the following analytically:

1. The ML estimator of $\psi^* = g(\theta^*) \in \Psi$ is

$$\hat{\Psi}_n = g(\hat{\Theta}_n)$$

and its point estimate is

$$\hat{\psi}_n = g(\hat{\theta}_n)$$

2. Suppose $g(\theta) = \psi : \Theta \rightarrow \Psi$ is **any** smooth function of θ , i.e. g is differentiable, and $g'(\theta) := \frac{\partial}{\partial \theta} g(\theta) \neq 0$. Then, the distribution of the ML estimator $\hat{\Psi}_n$ is asymptotically

$\text{Normal}(\psi^*, \widehat{\text{se}}_n(\widehat{\Psi}_n)^2)$, i.e.:

$$\frac{\widehat{\Psi}_n - \psi^*}{\widehat{\text{se}}_n(\widehat{\Psi}_n)} \rightsquigarrow \text{Normal}(0, 1)$$

where the standard error $\widehat{\text{se}}_n(\widehat{\Psi}_n)$ of the ML estimator $\widehat{\Psi}_n$ of the unknown quantity $\psi^* \in \Psi$ can be obtained from the standard error $\widehat{\text{se}}_n(\widehat{\Theta}_n)$ of the ML estimator $\widehat{\Theta}_n$ of the parameter $\theta^* \in \Theta$, as follows:

$$\widehat{\text{se}}_n(\widehat{\Psi}_n) = |g'(\widehat{\theta}_n)| \widehat{\text{se}}_n(\widehat{\Theta}_n)$$

3. Using $\text{Normal}(\psi^*, \widehat{\text{se}}_n(\widehat{\Psi}_n)^2)$, we can construct the estimate of an approximate Normal-based $1 - \alpha$ confidence interval for $\psi^* \in \Psi$:

$$C_n = [\underline{C}_n, \bar{C}_n] = \widehat{\psi}_n \pm z_{\alpha/2} \widehat{\text{se}}_n(\widehat{\psi}_n)$$

Let us do an example next.

Example 146 Let $X_1, X_2, \dots, X_n \stackrel{IID}{\sim} \text{Bernoulli}(\theta^*)$. Let $\psi = g(\theta) = \log(\theta/(1-\theta))$. Suppose we are interested in producing a point estimate and confidence interval for $\psi^* = g(\theta^*)$. We can use the Delta method as follows:

First, the estimated standard error of the ML estimator of θ^* , as shown in Example 142, is

$$\widehat{\text{se}}_n(\widehat{\Theta}_n) = \sqrt{\frac{\widehat{\theta}_n(1 - \widehat{\theta}_n)}{n}} .$$

The ML estimator of ψ^* is:

$$\widehat{\Psi}_n = \log(\widehat{\Theta}_n / (1 - \widehat{\Theta}_n))$$

and the ML estimate of ψ^* is:

$$\widehat{\psi}_n = \log(\widehat{\theta}_n / (1 - \widehat{\theta}_n)) .$$

Since, $g'(\theta) = 1/(\theta(1-\theta))$, by the Delta method, the estimated standard error of the ML estimator of ψ^* is:

$$\widehat{\text{se}}_n(\widehat{\Psi}_n) = |g'(\widehat{\theta}_n)| (\widehat{\text{se}}_n(\widehat{\Theta}_n)) = \frac{1}{\widehat{\theta}_n(1 - \widehat{\theta}_n)} \sqrt{\frac{\widehat{\theta}_n(1 - \widehat{\theta}_n)}{n}} = \frac{1}{\sqrt{n\widehat{\theta}_n(1 - \widehat{\theta}_n)}} = \frac{1}{\sqrt{n\bar{x}_n(1 - \bar{x}_n)}} .$$

An approximate 95% confidence interval for $\psi^* = \log(\theta^*/(1-\theta^*))$ is:

$$\widehat{\psi}_n \pm \frac{1.96}{\sqrt{n\widehat{\theta}_n(1 - \widehat{\theta}_n)}} = \log(\widehat{\theta}_n / (1 - \widehat{\theta}_n)) \pm \frac{1.96}{\sqrt{n\widehat{\theta}_n(1 - \widehat{\theta}_n)}} = \log(\bar{x}_n / (1 - \bar{x}_n)) \pm \frac{1.96}{\sqrt{n\bar{x}_n(1 - \bar{x}_n)}} .$$

Example 147 (Delta Method for a Normal Experiment) Let us try the Delta method on a continuous RV. Let $X_1, X_2, \dots, X_n \stackrel{IID}{\sim} \text{Normal}(\mu^*, \sigma^{*2})$. Suppose that μ^* is known and σ^* is unknown. Let us derive the ML estimate $\widehat{\psi}_n$ of $\psi^* = \log(\sigma^*)$ and a 95% confidence interval for it in 6 steps.

(1) First let us find the log-likelihood function $\ell(\sigma)$

$$\begin{aligned}
\ell(\sigma) &:= \log(L(\sigma)) := \log(L(x_1, x_2, \dots, x_n; \sigma)) = \log \left(\prod_{i=1}^n f(x_i; \sigma) \right) = \sum_{i=1}^n \log(f(x_i; \sigma)) \\
&= \sum_{i=1}^n \log \left(\frac{1}{\sigma \sqrt{2\pi}} \exp \left(-\frac{1}{2\sigma^2}(x_i - \mu)^2 \right) \right) \quad \because f(x_i; \sigma) \text{ in (3.39) is pdf of } \text{Normal}(\mu, \sigma^2) \text{ RV with known } \mu \\
&= \sum_{i=1}^n \left(\log \left(\frac{1}{\sigma \sqrt{2\pi}} \right) + \log \left(\exp \left(-\frac{1}{2\sigma^2}(x_i - \mu)^2 \right) \right) \right) \\
&= \sum_{i=1}^n \log \left(\frac{1}{\sigma \sqrt{2\pi}} \right) + \sum_{i=1}^n \left(-\frac{1}{2\sigma^2}(x_i - \mu)^2 \right) = n \log \left(\frac{1}{\sigma \sqrt{2\pi}} \right) + \left(-\frac{1}{2\sigma^2} \right) \sum_{i=1}^n (x_i - \mu)^2 \\
&= n \left(\log \left(\frac{1}{\sqrt{2\pi}} \right) + \log \left(\frac{1}{\sigma} \right) \right) - \left(\frac{1}{2\sigma^2} \right) \sum_{i=1}^n (x_i - \mu)^2 \\
&= n \log \left(\sqrt{2\pi}^{-1} \right) + n \log(\sigma^{-1}) - \left(\frac{1}{2\sigma^2} \right) \sum_{i=1}^n (x_i - \mu)^2 \\
&= -n \log \left(\sqrt{2\pi} \right) - n \log(\sigma) - \left(\frac{1}{2\sigma^2} \right) \sum_{i=1}^n (x_i - \mu)^2
\end{aligned}$$

(2) Let us find its derivative with respect to the unknown parameter σ next.

$$\begin{aligned}
\frac{\partial}{\partial \sigma} \ell(\sigma) &:= \frac{\partial}{\partial \sigma} \left(-n \log \left(\sqrt{2\pi} \right) - n \log(\sigma) - \left(\frac{1}{2\sigma^2} \right) \sum_{i=1}^n (x_i - \mu)^2 \right) \\
&= \frac{\partial}{\partial \sigma} \left(-n \log \left(\sqrt{2\pi} \right) \right) - \frac{\partial}{\partial \sigma} (n \log(\sigma)) - \frac{\partial}{\partial \sigma} \left(\left(\frac{1}{2\sigma^2} \right) \sum_{i=1}^n (x_i - \mu)^2 \right) \\
&= 0 - n \frac{\partial}{\partial \sigma} (\log(\sigma)) - \left(\frac{1}{2} \sum_{i=1}^n (x_i - \mu)^2 \right) \frac{\partial}{\partial \sigma} (\sigma^{-2}) \\
&= -n\sigma^{-1} - \left(\frac{1}{2} \sum_{i=1}^n (x_i - \mu)^2 \right) (-2\sigma^{-3}) = -n\sigma^{-1} + \sigma^{-3} \sum_{i=1}^n (x_i - \mu)^2
\end{aligned}$$

(3) Now, let us set the derivative equal to 0 and solve for σ .

$$\begin{aligned}
0 = -n\sigma^{-1} + \sigma^{-3} \sum_{i=1}^n (x_i - \mu)^2 &\iff n\sigma^{-1} = \sigma^{-3} \sum_{i=1}^n (x_i - \mu)^2 \iff n\sigma^{-1}\sigma^{+3} = \sum_{i=1}^n (x_i - \mu)^2 \\
&\iff n\sigma^{-1+3} = \sum_{i=1}^n (x_i - \mu)^2 \iff n\sigma^2 = \sum_{i=1}^n (x_i - \mu)^2 \\
&\iff \sigma^2 = \left(\sum_{i=1}^n (x_i - \mu)^2 \right) / n \iff \sigma = \sqrt{\sum_{i=1}^n (x_i - \mu)^2 / n}
\end{aligned}$$

Finally, we set the solution, i.e. the maximiser of the concave-down log-likelihood function of σ with a known and fixed μ^* as our ML estimate $\hat{\sigma}_n = \sqrt{\sum_{i=1}^n (x_i - \mu^*)^2 / n}$. Analogously, the ML estimator of σ^* is $\hat{\Sigma}_n = \sqrt{\sum_{i=1}^n (X_i - \mu^*)^2 / n}$. Don't confuse Σ , the upper-case sigma, with $\sum_{i=1}^n \bigcirc_i$, the summation over some \bigcirc_i 's. This is usually clear from the context.

(4) Next, let us get the estimated standard error $\hat{s}e_n$ for the estimator of σ^* via Fisher Information. The Log-likelihood function of σ , based on one sample from the $\text{Normal}(\mu, \sigma^2)$ RV with known μ is,

$$\log f(x; \sigma) = \log \left(\frac{1}{\sigma \sqrt{2\pi}} \exp \left(-\frac{1}{2\sigma^2} (x - \mu)^2 \right) \right) = -\log(\sqrt{2\pi}) - \log(\sigma) - \left(\frac{1}{2\sigma^2} \right) (x - \mu)^2$$

Therefore, in much the same way as in part (2) earlier,

$$\begin{aligned} \frac{\partial^2 \log f(x; \sigma)}{\partial \sigma^2} &:= \frac{\partial}{\partial \sigma} \left(\frac{\partial}{\partial \sigma} \left(-\log(\sqrt{2\pi}) - \log(\sigma) - \left(\frac{1}{2\sigma^2} \right) (x - \mu)^2 \right) \right) \\ &= \frac{\partial}{\partial \sigma} (-\sigma^{-1} + \sigma^{-3}(x - \mu)^2) = \sigma^{-2} - 3\sigma^{-4}(x - \mu)^2 \end{aligned}$$

Now, we compute the Fisher Information of one sample as an expectation of the continuous RV X over $\mathbb{X} = (-\infty, \infty)$ with density $f(x; \sigma)$,

$$\begin{aligned} I_1(\sigma) &= - \int_{x \in \mathbb{X} = (-\infty, \infty)} \left(\frac{\partial^2 \log f(x; \sigma)}{\partial \lambda^2} \right) f(x; \sigma) dx = - \int_{-\infty}^{\infty} (\sigma^{-2} - 3\sigma^{-4}(x - \mu)^2) f(x; \sigma) dx \\ &= \int_{-\infty}^{\infty} -\sigma^{-2} f(x; \sigma) dx + \int_{-\infty}^{\infty} 3\sigma^{-4}(x - \mu)^2 f(x; \sigma) dx \\ &= -\sigma^{-2} \int_{-\infty}^{\infty} f(x; \sigma) dx + 3\sigma^{-4} \int_{-\infty}^{\infty} (x - \mu)^2 f(x; \sigma) dx \\ &= -\sigma^{-2} + 3\sigma^{-4}\sigma^2 \quad \because \sigma^2 = V(X) = E(X - E(X))^2 = \int_{-\infty}^{\infty} (x - \mu)^2 f(x; \sigma) dx \\ &= -\sigma^{-2} + 3\sigma^{-4+2} = -\sigma^{-2} + 3\sigma^{-2} = 2\sigma^{-2} \end{aligned}$$

Therefore, the estimated standard error of the estimator of the unknown σ^* is

$$\hat{s}e_n(\hat{\Sigma}_n) = \frac{1}{\sqrt{I_n(\hat{\sigma}_n)}} = \frac{1}{\sqrt{nI_1(\hat{\sigma}_n)}} = \frac{1}{\sqrt{n2\sigma^{-2}}} = \frac{\sigma}{\sqrt{2n}} .$$

(5) Given that $\psi = g(\sigma) = \log(\sigma)$, we derive the estimated standard error of $\psi^* = \log(\sigma^*)$ via the Delta method as follows:

$$\hat{s}e_n(\hat{\Psi}_n) = |g'(\sigma)| \hat{s}e_n(\hat{\Sigma}_n) = \left| \frac{\partial}{\partial \sigma} \log(\sigma) \right| \frac{\sigma}{\sqrt{2n}} = \frac{1}{\sigma} \frac{\sigma}{\sqrt{2n}} = \frac{1}{\sqrt{2n}} .$$

(6) Finally, the 95% confidence interval for ψ^* is $\hat{\psi}_n \pm 1.96 \hat{s}e_n(\hat{\Psi}_n) = \log(\hat{\sigma}_n) \pm 1.96 \frac{1}{\sqrt{2n}}$.

Chapter 5

Maximum Likelihood Estimation for Multiparameter Models

5.1 Introduction

When two or more parameters index a statistical experiment we want to estimate the vector-valued parameter $\theta^* := (\theta_1^*, \dots, \theta_k^*)$. Here we will find the maximum likelihood estimates of vector-valued parameters.

The maximum likelihood estimator (MLE) of a possibly unknown but fixed parameter $\theta^* := (\theta_1^*, \dots, \theta_k^*)$ in a multi-parametric experiment, i.e. $\theta^* \in \Theta \subset \mathbb{R}^k$ with $1 < k < \infty$ is defined analogously to Definition 69 with the exception that we allow the parameter to be a vector. We take an excursion in multi-dimensional optimisation before finding the MLE of a parametric experiment involving two parameters.

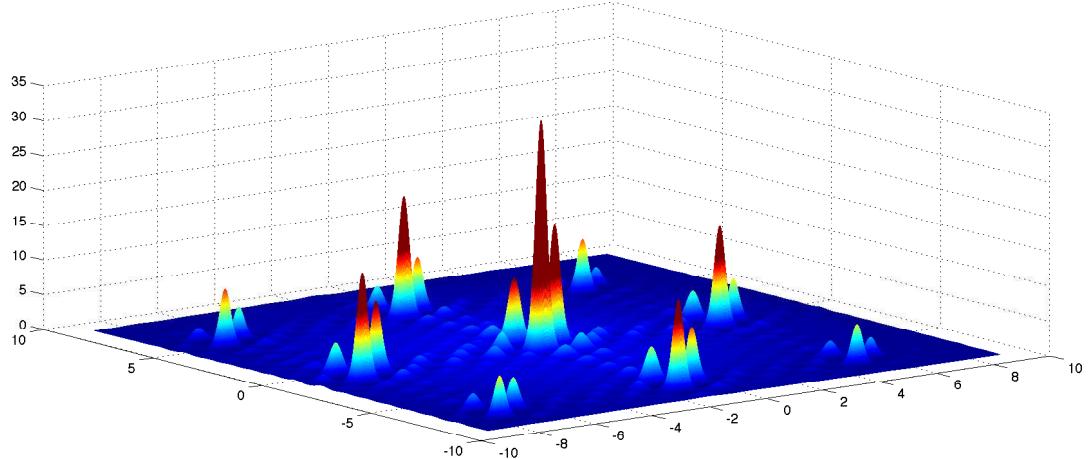
5.2 Practical Excursion in Multi-dimensional Optimisation

The basic idea involves multi-dimensional iterations that attempt to converge on a local maximum close to the starting vector $\theta^{(0)} \in \Theta$ (our initial guess). We can employ MATLAB's built-in function `fminsearch` to find the MLE of vector-valued parameters such as in the Lognormal model with two parameters, i.e. $\theta = (\lambda, \zeta) \in \Theta \subset \mathbb{R}^2$. The function `fminsearch` is similar to `fminbnd` except that it handles a given function of many variables, and the user specifies a starting vector $\theta^{(0)}$ rather than a starting interval. Thus, `fminsearch` tries to return a vector $\theta^{(*)}$ that is a local minimiser of, $-\log(L(x_1, x_2, \dots, x_n; \theta))$, the negative log-likelihood function of the vector-valued parameter θ , near this starting vector $\theta^{(0)}$. We illustrate the use of `fminsearch` on a more challenging target called the Levy density:

$$f(x, y) = \exp \left(-\frac{1}{50} \left(\left(\sum_{i=1}^5 i \cos((i-1)x + i) \right) \left(\sum_{j=1}^5 j \cos((j+1)y + j) \right) + (x + 1.42513)^2 + (y + 0.80032)^2 \right) \right) \quad (5.1)$$

`fminsearch` uses the simplex search method [Nelder, J.A., and Mead, R. 1965, Computer Journal, vol. 7, p. 308-313]. For an animation of the method and more details, please visit http://en.wikipedia.org/wiki/Nelder-Mead_method. An advantage of the method is that it does not use numerical (finite differencing) or analytical (closed-form expressions) gradients but relies on a direct search method. Briefly, the simplex algorithm tries to “tumble and shrink” a simplex towards the local valley of the function to be minimised. If k is the dimension of

Figure 5.1: Plot of Levy density as a function of the parameter $(x, y) \in [-10, 10]^2$ scripted in Labwork ??.



the parameter space or domain of the function to be optimised, a k -dimensional simplex is specified by its $k + 1$ distinct vertices each of dimension k . Thus, a simplex is a triangle in a two-dimensional space and a pyramid in a three-dimensional space. At each iteration of the algorithm:

1. A new point inside or nearby the current simplex is proposed.
2. The function's value at the newly proposed point is compared with its values at the vertices of the simplex.
3. One of the vertices is typically replaced by the proposed point, giving rise to a new simplex.
4. The first three steps are repeated until the diameter of the simplex is less than the specified tolerance.

A major limitation of `fminsearch`, as demonstrated with the Levy target (encoded in Labwork ??) is that it can only give local solutions. The **global maximiser** of the Levy function $f(x, y)$ is $(-1.3069, -1.4249)$ and the **global maximum** is $f(-1.3069, -1.4249) = 33.8775$ For instance, if we start the search close to, say $(x^{(0)}, y^{(0)}) = (-1.3, -1.4)$, as shown below, then the simplex algorithm converges as desired to the solution $(-1.3068, -1.4249)$.

```
>> [params, fvalue, exitflag, output] = fminsearch('NegLevyDensity', [-1.3 -1.4], options)
params =    -1.3068    -1.4249
fvalue =   -33.8775
exitflag =      1
output =
    iterations: 24
    funcCount: 46
    algorithm: 'Nelder-Mead simplex direct search'
    message: [1x194 char]
```

However, if we start the search further away, say $(x^{(0)}, y^{(0)}) = (1.3, 1.4)$, as shown below, then the algorithm converges to the **local maximiser** $(1.1627, 1.3093)$ with a **local maximum** value of $f(1.1627, 1.3093) = 0.9632$, which is clearly smaller than the global maximum of 33.8775.

```
>> [params, fvalue, exitflag, output] = fminsearch('NegLevyDensity',[1.3 1.4],options)
params = 1.1627    1.3093
fvalue = -0.9632
exitflag = 1
output =
    iterations: 29
    funcCount: 57
    algorithm: 'Nelder-Mead simplex direct search'
    message: [1x194 char]
```

Therefore, we have to be extremely careful when using point-valued, iterative, local optimisation algorithms, implemented in floating-point arithmetic to find the global maximum. Other examples of such algorithms include:

- **Conjugate Gradient Method:**
http://en.wikipedia.org/wiki/Conjugate_gradient_method
- **Broyden-Fletcher-Goldfarb-Shanno (BFGS) method:**
http://en.wikipedia.org/wiki/BFGS_method
- **Simulated Annealing:**
http://en.wikipedia.org/wiki/Simulated_annealing

In general, we have no guarantee that the output of such local optimisation routines will indeed be the global optimum. In practice, you can start the search at several distinct starting points and choose the best local maximum from the lot.

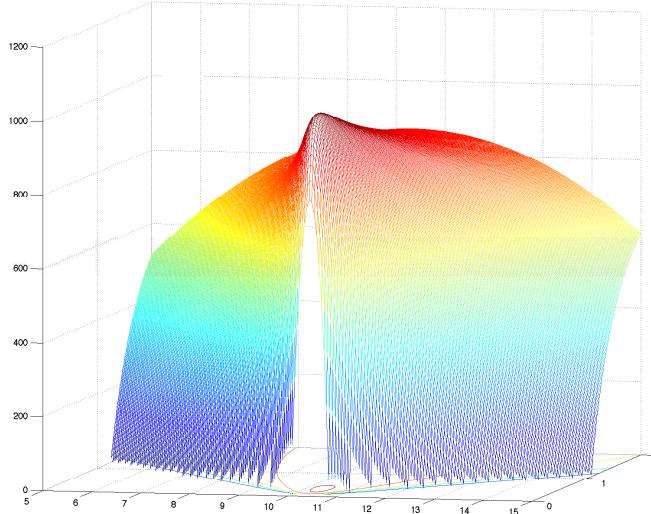
When the target function is “well-behaved,” i.e. uni-modal or single-peaked and not too spiky, the optimisation routine can be expected to perform well. Log-likelihood functions are often well-behaved. Let us generate 100 samples from an RV $C \sim \text{Lognormal}(\lambda^* = 10.36, \zeta^* = 0.26)$ by exponentiating the samples from the $\text{Normal}(10.36, 0.26^2)$ RV, and then compute the corresponding MMEs and MLEs for parameters (λ, ζ) using the formulae in Table 5.1.

```
>> rand('twister',001); % set the fundamental sampler
>> % draw 100 samples from the Lognormal(10.36,0.26) RV
>> Cs = exp(arrayfun(@(u)(Sample1NormalByNewRap(u,10.36,0.26^2)),rand(1,100)));
>> MLElambdahat = mean(log(Cs)) % maximum likelihood estimate of lambda
MLElambdahat = 10.3397
>> MLEzetahat = sqrt(mean((log(Cs)-MLElambdahat) .^ 2)) % max. lkl. estimate of zeta
MLEzetahat = 0.2744
>> MMEzetaahat = sqrt(log(var(Cs)/(mean(Cs)^2) + 1)) % moment estimate of zeta
MMEzetaahat = 0.2624
>> MMElambdahat = log(mean(Cs))-(0.5*MMEzetaahat^2) % moment estimate of lambda
MMElambdahat = 10.3417
```

Let us try to apply the simplex algorithm to find the MLE numerically. We first encode the negative log-likelihood function of the parameters $(\lambda, \zeta) \in (0, \infty)^2$ for the given data x , as follows:

```
function l = NegLogNormalLogLkl(x,params)
% Returns the -log likelihood of [lambda zeta]=exp(params)
% for observed data vector x=(x_1,...,x_n) ~ IID LogNormal(lambda, zeta).
% We define lambda and zeta as exp(params) to allow for unconstrained
```

Figure 5.2: Plot of the “well-behaved” (uni-modal and non-spiky) $\log(L((x_1, x_2, \dots, x_{100}); \lambda, \zeta))$, based on 100 samples $(x_1, x_2, \dots, x_{100})$ drawn from the Lognormal($\lambda^* = 10.36, \zeta^* = 0.26$) as per Labwork ??.



```
% minimisation by fminsearch and respect the positive domain constraints
% for Lambda and zeta. So in the end we re-transform, i.e. [lambda zeta]=exp(params)
% lambda=params(1); zeta=params(1);
lambda=exp(params(1)); zeta=exp(params(2));
% minus Log-likelihood function
l = -sum(log((1 ./ (sqrt(2*pi)*zeta) .* x) .* exp((-1/(2*zeta^2))*(log(x)-lambda).^2)));
```

Here is how we can call `fminsearch` and find the MLE after the re-transformation.

```
>> [params, fvalue, exitflag, output] = ...
fminsearch(@(params)(NegLogNormalLogLkl(Cs,params)),[log(5), log(1)])
params =    2.3360   -1.2931
fvalue =  -1.0214e+03
exitflag =      1
output =
    iterations: 74
    funcCount: 131
    algorithm: 'Nelder-Mead simplex direct search'
    message: [1x194 char]
>> % But we want exp(params) since we defined lambda and zeta as exp(params)
exp(params)
ans =    10.3397    0.2744
```

Note that the MLEs $(\hat{\lambda}_{100}, \hat{\zeta}_{100}) = (10.3397, 0.2744)$ from 74 iterations or “tumbles” of the ‘Nelder-Mead simplex (triangle)’ and the MLEs agree well with the direct evaluations `MLElambdahat` and `MLEzetahat` based on the formulae in Table 5.1.

Summarizing Table of Point Estimators

Using the sample mean \bar{X}_n and sample standard deviation S_n defined in (4.1) and (4.5), respectively, we summarise the two point estimators of the parameters of some common distributions

below. For some cases, the MLE is the same as the MME (method of moments) and can be solved analytically.

Table 5.1: Summary of the Method of Moment Estimator (MME) and the Maximum Likelihood Estimator (MLE) for some IID Experiments.

Statistical Experiment	MLE	MME
$X_1, X_2, \dots, X_n \stackrel{IID}{\sim} \text{Bernoulli}(\theta)$	$\hat{\theta} = \bar{X}_n$	same as MLE
$X_1, X_2, \dots, X_n \stackrel{IID}{\sim} \text{Exponential}(\lambda)$	$\hat{\lambda} = 1/\bar{X}_n$	same as MLE
$X_1, X_2, \dots, X_n \stackrel{IID}{\sim} \text{Normal}(\mu, \sigma^2)$	$\hat{\mu} = \bar{X}_n, \hat{\sigma} = \sqrt{\frac{n-1}{n} S_n^2}$	$\hat{\mu} = \bar{X}_n, \hat{\sigma} = S_n$
$X_1, X_2, \dots, X_n \stackrel{IID}{\sim} \text{Lognormal}(\lambda, \zeta)$	$\hat{\lambda} = \frac{1}{n} \sum_{i=1}^n \log(X_i)$ $\hat{\zeta} = \sqrt{\frac{1}{n} \sum_{i=1}^n (\log(X_i) - \hat{\lambda})^2}$	$\hat{\lambda} = \log(\bar{X}_n) - \frac{1}{2} \hat{\zeta}^2$ $\hat{\zeta} = \sqrt{\log(S_n^2/\bar{X}_n^2 + 1)}$

5.3 Confidence Sets for Multiparameter Models

We will extend the Fisher Information and Delta method to models with more than one parameter:

$$X_1, X_2, \dots, X_n \stackrel{IID}{\sim} f(x; \theta^*), \quad \theta^* := (\theta_1^*, \theta_2^*, \dots, \theta_k^*) \in \Theta \subset \mathbb{R}^k.$$

Let, the ML estimator of the fixed and possibly unknown vector-valued parameter θ^* be:

$$\hat{\Theta}_n := (\hat{\Theta}_{1,n}, \hat{\Theta}_{2,n}, \dots, \hat{\Theta}_{k,n}), \quad \hat{\Theta}_n := \hat{\Theta}_n(X_1, X_2, \dots, X_n) : \mathbb{X}_n \rightarrow \Theta$$

and the ML estimate based on n observations x_1, x_2, \dots, x_n be:

$$\hat{\theta}_n := (\hat{\theta}_{1,n}, \hat{\theta}_{2,n}, \dots, \hat{\theta}_{k,n}), \quad \hat{\theta}_n := \hat{\theta}_n(x_1, x_2, \dots, x_n) \in \Theta.$$

Let the log-likelihood function and its Hessian matrix $H = (H_{i,j})_{i,j=1,2,\dots,k}$ of partial derivatives be:

$$\ell_n(\theta) := \ell_n(\theta_1, \theta_2, \dots, \theta_k) := \sum_{i=1}^n \log(f(x_i; (\theta_1, \theta_2, \dots, \theta_k))), \quad H_{i,j} := \frac{\partial}{\partial \theta_i} \frac{\partial}{\partial \theta_j} \ell_n(\theta_1, \theta_2, \dots, \theta_k),$$

respectively, provided the log-likelihood function is sufficiently smooth.

Definition 73 (Fisher Information Matrix) The Fisher Information matrix is:

$$I_n(\theta) := I_n(\theta_1, \theta_2, \dots, \theta_k) = - \begin{bmatrix} E_\theta(H_{1,1}) & E_\theta(H_{1,2}) & \cdots & E_\theta(H_{1,k}) \\ E_\theta(H_{2,1}) & E_\theta(H_{2,2}) & \cdots & E_\theta(H_{2,k}) \\ \vdots & \vdots & \ddots & \vdots \\ E_\theta(H_{k,1}) & E_\theta(H_{k,2}) & \cdots & E_\theta(H_{k,k}) \end{bmatrix} \quad (5.2)$$

and its matrix inverse is denoted by $I_n^{-1}(\theta)$.

Proposition 74 (Asymptotic Normality of MLE in Multiparameter Models) Let

$$X_1, X_2, \dots, X_n \stackrel{IID}{\sim} f(x_1; \theta_1^*, \theta_2^*, \dots, \theta_k^*), \quad \theta^* = (\theta_1^*, \theta_2^*, \dots, \theta_k^*) \in \Theta \subset \mathbb{R}^k,$$

for some fixed and possibly unknown $\theta^* \in \Theta \subset \mathbb{R}^k$. Then, under appropriate regularity conditions:

$$\widehat{\Theta}_n := (\widehat{\theta}_{1,n}, \widehat{\theta}_{2,n}, \dots, \widehat{\theta}_{k,n}) \rightsquigarrow \text{Normal}(\theta^*, I_n^{-1})$$

In other words, the vector-valued estimator $\widehat{\Theta}_n$ converges in distribution to the multivariate Normal distribution centred at the unknown parameter θ^* with the variance-covariance matrix given by inverse Fisher Information matrix I_n^{-1} . Furthermore, let $I_n^{-1}(j,j)$ denote the j^{th} diagonal entry of I_n^{-1} . In this case:

$$\frac{\widehat{\theta}_{j,n} - \theta_j^*}{\sqrt{I_n^{-1}(j,j)}} \rightsquigarrow \text{Normal}(0, 1)$$

and the approximate covariance of $\widehat{\Theta}_{i,n}$ and $\widehat{\Theta}_{j,n}$ is:

$$\text{Cov}(\widehat{\Theta}_{i,n}, \widehat{\Theta}_{j,n}) \approx I_n^{-1}(i,j).$$

Now, let us look at a way of obtaining ML estimates and confidence sets for functions of θ . Suppose the real-valued function $g(\theta) = \psi : \Theta \rightarrow \Psi$ maps points in the k -dimensional parameter space $\Theta \subset \mathbb{R}^k$ to points in $\Psi \subset \mathbb{R}$. Let the gradient of g be

$$\nabla g(\theta) := \nabla g(\theta_1, \theta_2, \dots, \theta_k) = \begin{pmatrix} \frac{\partial}{\partial \theta_1} g(\theta_1, \theta_2, \dots, \theta_k) \\ \frac{\partial}{\partial \theta_2} g(\theta_1, \theta_2, \dots, \theta_k) \\ \vdots \\ \frac{\partial}{\partial \theta_k} g(\theta_1, \theta_2, \dots, \theta_k) \end{pmatrix}.$$

Proposition 75 (Multiparameter Delta Method) Suppose:

1. $X_1, X_2, \dots, X_n \stackrel{IID}{\sim} f(x_1; \theta_1^*, \theta_2^*, \dots, \theta_k^*), \quad \theta^* = (\theta_1^*, \theta_2^*, \dots, \theta_k^*) \in \Theta \subset \mathbb{R}^k,$
2. Let $\widehat{\Theta}_n$ be a ML estimator of $\theta^* \in \Theta$ and let $\widehat{\theta}_n$ be its ML estimate, and
3. Let $g(\theta) = \psi : \Theta \rightarrow \Psi \subset \mathbb{R}$ be a smooth function such that $\nabla g(\widehat{\theta}_n) \neq 0$.

Then:

1. $\widehat{\Psi}_n = g(\widehat{\Theta}_n)$ is the ML estimator and $\widehat{\psi}_n = g(\widehat{\theta}_n)$ is the ML estimate of $\psi^* = g(\theta^*) \in \Psi$,
2. The standard error of the ML estimator of ψ^* is:

$$\widehat{\text{se}}_n(\widehat{\Psi}_n) = \sqrt{\left(\nabla g(\widehat{\theta}_n) \right)^T I_n^{-1}(\widehat{\theta}_n) \left(\nabla g(\widehat{\theta}_n) \right)},$$

3. The ML estimator of ψ^* is asymptotically normal, i.e.:

$$\frac{\widehat{\Psi}_n - \psi^*}{\widehat{\text{se}}_n(\widehat{\Psi}_n)} \rightsquigarrow \text{Normal}(0, 1),$$

4. And a $1 - \alpha$ confidence interval for ψ^* is:

$$\hat{\psi}_n \pm z_{\alpha/2} \widehat{\text{se}}_n(\hat{\Psi}_n)$$

Let us put the theory to practice in the problem of estimating the coefficient of variation from samples of size n from an RV.

Example 148 (Estimating the Coefficient of Variation of a $\text{Normal}(\mu^*, \sigma^{*2})$ RV) Let

$$\psi^* = g(\mu^*, \sigma^*) = \sigma^*/\mu^*, \quad X_1, X_2, \dots, X_n \stackrel{IID}{\sim} \text{Normal}(\mu^*, \sigma^{*2}).$$

We do not know the fixed parameters (μ^*, σ^*) and are interested in estimating the coefficient of variation ψ^* based on n IID samples x_1, x_2, \dots, x_n . We have already seen that the ML estimates of μ^* and σ^* are:

$$\hat{\mu}_n = \bar{x}_n := \frac{1}{n} \sum_{i=1}^n x_i, \quad \hat{\sigma}_n = s_n := \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu}_n)^2}.$$

Thus, the ML estimate of $\psi^* = \sigma^*/\mu^*$ is:

$$\hat{\psi}_n = \frac{\hat{\sigma}_n}{\hat{\mu}_n} = \frac{s_n}{\bar{x}_n}$$

We can now derive the standard error of the ML estimator $\hat{\Psi}_n$ by first computing $I_n(\mu, \sigma)$, $I_n^{-1}(\mu, \sigma)$, and $\nabla g(\mu, \sigma)$. A careful computation shows that:

$$I_n(\mu, \sigma) = \begin{bmatrix} \frac{n}{\sigma^2} & 0 \\ 0 & \frac{2n}{\sigma^2} \end{bmatrix}, \quad I_n^{-1}(\mu, \sigma) = \frac{1}{n} \begin{bmatrix} \sigma^2 & 0 \\ 0 & \frac{\sigma^2}{2} \end{bmatrix}, \quad \nabla g(\mu, \sigma) = \begin{pmatrix} -\frac{\sigma}{\mu^2} \\ \frac{1}{\mu} \end{pmatrix}.$$

Therefore, the standard error of interest is:

$$\widehat{\text{se}}_n(\hat{\Psi}_n) = \sqrt{\left(\nabla g(\hat{\theta}_n) \right)^T I_n^{-1}(\hat{\theta}_n) \left(\nabla g(\hat{\theta}_n) \right)} = \frac{1}{\sqrt{n}} \sqrt{\frac{1}{\hat{\mu}_n^4} + \frac{\hat{\sigma}_n^2}{2\hat{\mu}_n^2}}$$

and the 95% confidence interval for the unknown coefficient of variation ψ^* is:

$$\hat{\psi}_n \pm z_{\alpha/2} \widehat{\text{se}}_n(\hat{\Psi}_n) = \frac{s_n}{\bar{x}_n} \pm z_{\alpha/2} \left(\frac{1}{\sqrt{n}} \sqrt{\frac{1}{\hat{\mu}_n^4} + \frac{\hat{\sigma}_n^2}{2\hat{\mu}_n^2}} \right)$$

Let us get our hands dirty in the machine with Labwork 149 next.

Labwork 149 (Computing the coefficient of variation of a $\text{Normal}(\mu^*, \sigma^{*2})$ RV) Let us apply these results to $n = 100$ simulated samples from $\text{Normal}(100, 10^2)$ as follows.

```
n=100; % sample size
Mustar=100; % true mean
Sigmapstar=10; % true standard deviation
rand('twister',67345); Us=rand(1,100); % draw some Uniform(0,1) samples
x=arrayfun(@(u)(Sample1NormalByNewRap(u,Mustar,Sigmapstar^2)),Us); % get normal samples
Muhat=mean(x) % sample mean is MLE of Mustar
Sigmahat=std(x) % sample standard deviation is MLE for Sigmapstar
Pihat=Sigmahat/Muhat % MLE of coefficient of variation std/mean
Sehat = sqrt((1/Muhat^4)+(Sigmahat^2/(2*Muhat^2)))/sqrt(n) % standar error estimate
ConfInt95=[Pihat-1.96*Sehat, Pihat+1.96*Sehat] % 1.96 since 1-alpha=0.95
```

```
>> CoeffOfVarNormal  
Muhat = 100.3117  
Sigmahat = 10.9800  
Psihat = 0.1095  
Sehat = 0.0077  
ConfInt95 = 0.0943 0.1246
```

Chapter 6

Simulation

“Anyone who considers arithmetical methods of producing random digits is, of course, in a state of sin.” — John von Neumann (1951)

6.1 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 fourty 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).

6.2 Pseudo-Random Number Generators

Our probability model and the elementary continuous $\text{Uniform}(0, 1)$ 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) $\text{Uniform}(0, 1)$ RVs are generated. Second, appropriate transformations are made to these imitations of IID $\text{Uniform}(0, 1)$ 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 $\text{Uniform}(0, 1)$ random variates.

In the following chapters, we focus on transforming IID Uniform(0, 1) variates to other non-uniform variates. In this chapter, we focus on the art of imitating IID Uniform(0, 1) variates using simple deterministic rules.

6.2.1 Linear Congruential Generators

The following procedure introduced by D. H. Lehmer in 1949 [*Proc. 2nd Symp. on Large-Scale Digital Calculating Machinery, Harvard Univ. Press, Cambridge, Mass., 1951, 141–146*] gives the simplest popular PRNG that can be useful in many statistical situations if used wisely.

Algorithm 2 Linear Congruential Generator (LCG)

- 1: *input:* five suitable integers:
 1. m , the modulus; $0 < m$
 2. a , the multiplier; $0 \leq a < m$
 3. c , the increment; $0 \leq c < m$
 4. x_0 , the seed; $0 \leq x_0 < m$
 5. n , the number of desired pseudo-random numbers
 - 2: *output:* $(x_0, x_1, \dots, x_{n-1})$, the linear congruential sequence of length n
 - 3: **for** $i = 1$ to $n - 1$ **do**
 - 4: $x_i \leftarrow (ax_{i-1} + c) \bmod m$
 - 5: **end for**
 - 6: *return:* (x_1, x_2, \dots, x_n)
-

In order to implement LCGs we need to be able to do high precision exact integer arithmetic in MATLAB. We employ the Module `vpi` to implement variable precision integer arithmetic. You need to download this module for the next Labwork.

Labwork 150 (Generic Linear Congruential Sequence) Let us implement Algorithm 2 in MATLAB as follows.

LinConGen.m

```

function x = LinConGen(m,a,c,x0,n)
% Returns the linear congruential sequence
% Needs variable precision integer arithmetic in MATLAB!!!
% Usage: x = LinConGen(m,a,c,x0,n)
% Tested:
% Knuth3.3.4Table1.Line1: LinConGen(100000001,23,0,01234,10)
% Knuth3.3.4Table1.Line5: LinConGen(256,137,0,01234,10)
% Knuth3.3.4Table1.Line20: LinConGen(2147483647,48271,0,0123456,10)
% Knuth3.3.4Table1.Line21: LinConGen(2147483399,40692,0,0123456,10)

x=zeros(1,n); % initialize an array of zeros
X=vpi(x0); % X is a variable precision integer seed
x(1) = double(X); % convert to double
A=vpi(a); M=vpi(m); C=vpi(c); % A,M,C as variable precision integers
for i = 2:n % loop to generate the Linear congruential sequence
    % the linear congruential operation in variable precision integer
    % arithmetic
    % comment out the next ';' to get integer output
    X=mod(A * X + C, M);

```

```

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) .

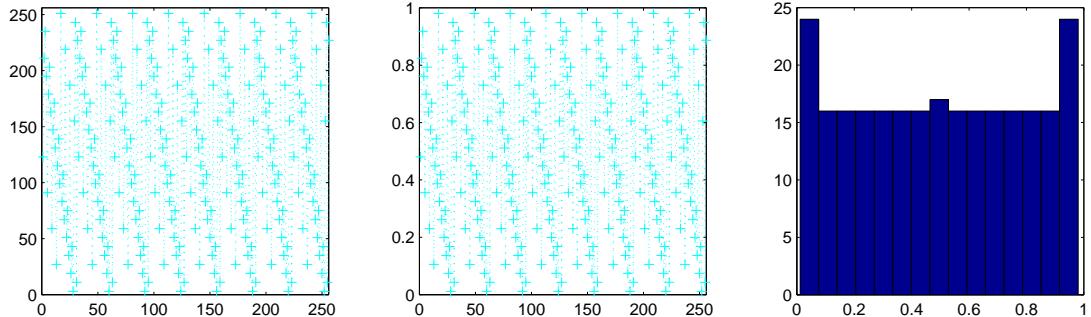
Labwork 151 (LCG with period length of 32) Consider the linear congruential sequence with $(m, a, c, x_0, n) = (256, 137, 0, 123, 257)$ with period length of only $32 < m = 256$. We can visualise the sequence as plots in Figure 6.1 after calling the following M-file.

```

LinConGenKnuth334T1L5Plots.m
-----
LCGSeq=LinConGen(256,137,0,123,257)
subplot(1,3,1)
plot(LCGSeq,'+')
axis([0 256 0 256]); axis square
LCGSeqIn01=LCGSeq ./ 256
subplot(1,3,2)
plot(LCGSeqIn01,'+')
axis([0 256 0 1]); axis square
subplot(1,3,3)
hist(LCGSeqIn01,15)
axis square

```

Figure 6.1: 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*:

$$x_i = f(x_{i-1}), \quad f : \{0, 1, 2, \dots, m-1\} \rightarrow \{0, 1, 2, \dots, m-1\} \text{ and } f(x_{i-1}) = (ax_{i-1} + c) \pmod{m}.$$

Since f maps the finite set $\{1, 2, \dots, m-1\}$ into itself, such systems are bound to have a repeating cycle of numbers called the **period**. In Labwork 150, the generator $\text{LinConGen}(13, 12, 11, 10, 12)$

has period $(10, 1)$ of length 2, the generator `LinConGen(13, 10, 9, 8, 12)` has period $(8, 11, 2, 3, 0, 9)$ of length 6 and the generator `LinConGen(256, 137, 0, 123, 257)` has a period of length 32. All these generators have a non-maximal period length less than their modulus m . A good generator should have a maximal period of m . Let us try to implement a generator with a maximal period of $m = 256$.

The period of a general LCG is at most m , and for some choices of a the period can be much less than m as shown in the examples considered earlier. The LCG will have a full period if and only if:

1. c and m are relatively prime,
2. $a - 1$ is divisible by all prime factors of m ,
3. $a - 1$ is a multiple of 4 if m is a multiple of 4

Labwork 152 (LCG with maximal 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 `LinConGen` and repeat Labwork 151 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 151? Next produce a Figure to visualise the sequence as done in Figure 6.1.

A useful sequence should clearly have a relatively long period, say at least 2^{30} . Therefore, the **modulus m has to be rather large** because the **period** cannot have more than m elements. Moreover, the quality of pseudo-random numbers of a LCG is extremely sensitive to the choice of m , a and c even if the maximal period is attained. The next example illustrates this point.

Labwork 153 (The infamous RANDU) RANDU is an infamous LCG, which has been used since the 1960s. It is widely considered to be one of the most ill-conceived random number generators designed. Notably, it fails the **spectral test** badly for dimensions greater than 2. The following commands help visualise the sequence of first 5001/3 triplets (x_i, x_{i+1}, x_{i+2}) seeded from $x_0 = 1$ (Figure 6.2). Read `help reshape` and `help plot3`.

```
>> x=reshape( (LinConGen(2147483648,65539,0,1,5001)./ 2147483648) ,3,[]);
>> plot3(x(1,:),x(2,:),x(3,:),'.'
```

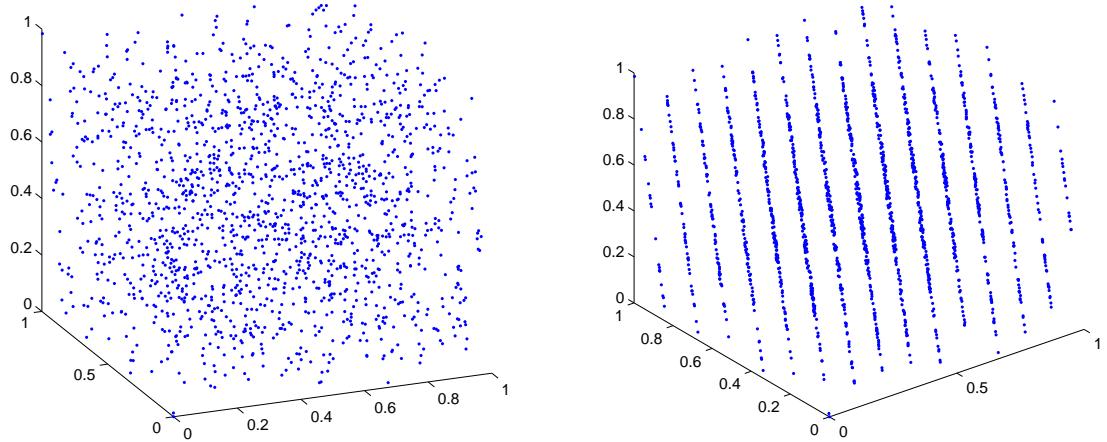
Labwork 154 (Fishman20 and Lecuyer21 LCGs) The following two LCGs are recommended in Knuth's Art of Computer Programming, vol. 2, for generating pseudo-random numbers for simple simulation tasks.

```
>> LinConGen(2147483647,48271,0,08787458,10) ./ 2147483647
ans =    0.0041    0.5239    0.0755    0.7624    0.6496    0.0769    0.9030    0.4259    0.9948    0.8868

>> LinConGen(2147483399,40692,0,01234567,10) ./ 2147483399
ans =    0.0006    0.3934    0.4117    0.7893    0.3913    0.6942    0.6790    0.3337    0.2192    0.1883
```

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.

Figure 6.2: 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 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.

6.2.2 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 155 (Calling PRNG in MATLAB) In MATLAB **rand** is basic PRNG command.

```
>> rand(1,10) % generate a 1 X 10 array of PRNs
ans =
    0.8147    0.9058    0.1270    0.9134    0.6324    0.0975    0.2785    0.5469    0.9575    0.9649
>> rand(1,10) % generate another 1 X 10 array of PRNs
ans =
```

```

0.1576    0.9706    0.9572    0.4854    0.8003    0.1419    0.4218    0.9157    0.7922    0.9595
>> rand('twister',5489) % reset the PRNG to default state Mersenne Twister with seed=5489
>> rand(1,10) % reproduce the first array
ans =
0.8147    0.9058    0.1270    0.9134    0.6324    0.0975    0.2785    0.5469    0.9575    0.9649
>> rand(1,10) % reproduce the second array
ans =
0.1576    0.9706    0.9572    0.4854    0.8003    0.1419    0.4218    0.9157    0.7922    0.9595

```

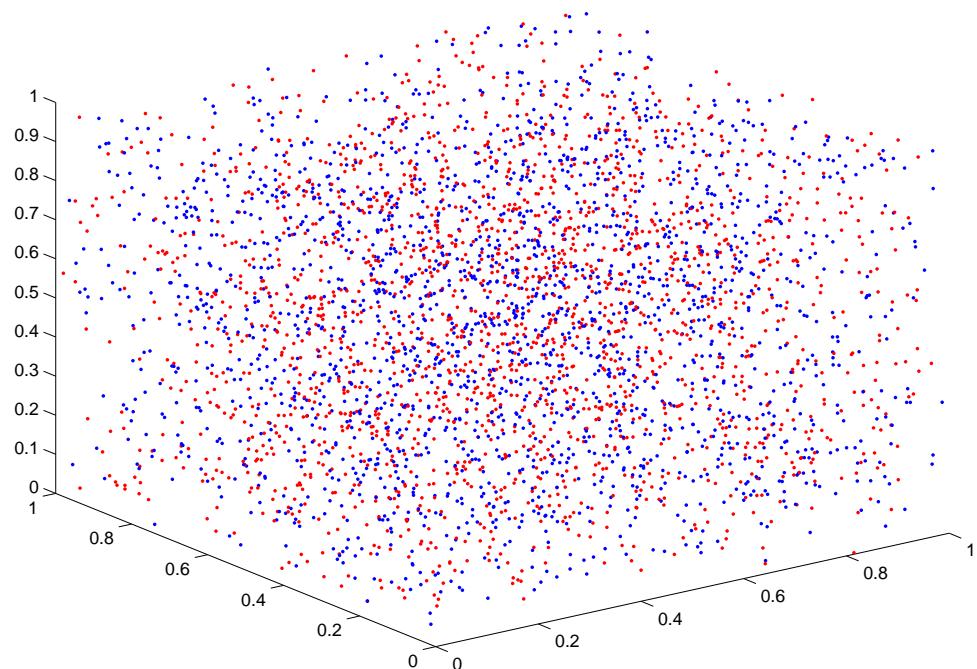
In general, you can use any seed value to initiate your PRNG. You may use the `clock` command to set the seed:

```

>> SeedFromClock=sum(100*clock); % save the seed from clock
>> rand('twister',SeedFromClock) % initialize the PRNG
>> rand(1,10)
ans =
0.3696    0.3974    0.6428    0.6651    0.6961    0.7311    0.8982    0.6656    0.6991    0.8606
>> rand(2,10)
ans =
0.3432    0.9511    0.3477    0.1007    0.8880    0.0853    0.6067    0.6976    0.4756    0.1523
0.5827    0.5685    0.0125    0.1555    0.5551    0.8994    0.2502    0.5955    0.5960    0.5700
>> rand('twister',SeedFromClock) % initialize the PRNG to same SeedFromClock
>> rand(1,10)
ans =
0.3696    0.3974    0.6428    0.6651    0.6961    0.7311    0.8982    0.6656    0.6991    0.8606

```

Figure 6.3: Triplet point clouds from the “Mersenne Twister” with two different seeds (see Labwork 156). .



Labwork 156 (3D plots of triplets generated by the “Mersenne Twister”) Try to find any correlation between triplets generated by the “Mersenne Twister” by rotating the 3D plot generated by the following code:

```
>> rand('twister',1234)
>> x=rand(3,2000); % store PRNs in a 3X2000 matrix named x
>> plot3(x(1,:),x(2,:),x(3,:),'.'
```

Compare this with the 3D plot of triplets from RANDU of Labwork 153. Which of these two PRNGs do you think is “more random” looking? and why?

Change the seed value to the recommended default by the authors and look at the point cloud (in red) relative to the previous point cloud (in blue). Rotate the plots to visualise from multiple angles. Are they still random looking?

```
>> rand('twister',1234)% same seed as before
>> x=rand(3,2000); % store PRNs in a 3X2000 matrix named x
>> rand('twister',5489)% the recommended default seed
>> y=rand(3,2000);% store PRNs seeded by 5489 in a 3X2000 matrix named y
>> plot3(x(1,:),x(2,:),x(3,:),'b.') % plot triplets as blue dots
>> hold on;
>> plot3(y(1,:),y(2,:),y(3,:),'r.') % plot triplets as red dots
```

6.3 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.

6.3.1 Inversion Sampler for Continuous Random Variables

Proposition 76 (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 .

Proof: The “one-line proof” of the proposition is due to the following equalities:

$$\mathbb{P}(F^{[-1]}(U) \leq x) = \mathbb{P}(\inf\{y : F(y) = U\} \leq x) = \mathbb{P}(U \leq F(x)) = F(x), \quad \text{for all } x \in \mathbb{R}.$$

Algorithm 3 Inversion Sampler or Inverse (C)DF Sampler

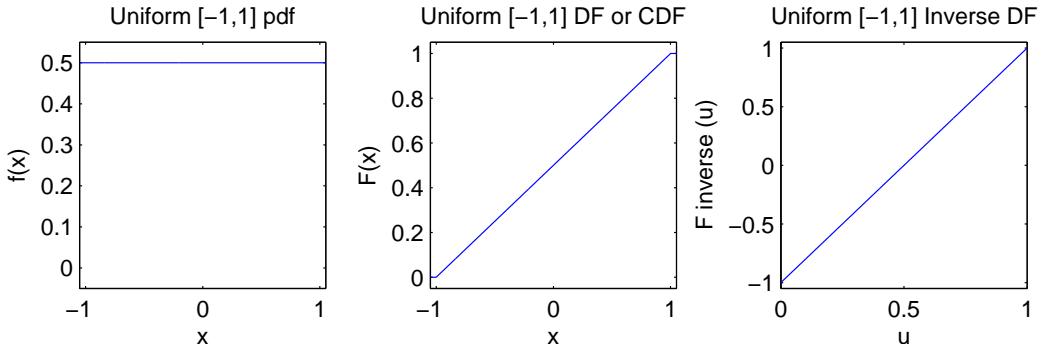
- 1: *input*: (1) $F^{[-1]}(x)$, inverse of the DF of the target RV X , (2) the fundamental sampler
 - 2: *initialise*: set the seed, if any, for the fundamental sampler
 - 3: *output*: a sample from X distributed according to F
 - 4: *draw* $u \sim \text{Uniform}(0, 1)$
 - 5: *return*: $x = F^{[-1]}(u)$
-

This yields the inversion sampler or the inverse (C)DF sampler, where we (i) *generate* $u \sim \text{Uniform}(0, 1)$ and (ii) *return* $x = F^{[-1]}(u)$, as formalised by the following algorithm.

This algorithm emphasises 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 universal 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 $\text{Uniform}(\theta_1, \theta_2)$ RV of Model 9 with the following PDF, DF and inverse DF. Let us simulate from it using the inversion sampler.

Figure 6.4: A plot of the PDF, DF or CDF and inverse DF of the $\text{Uniform}(-1, 1)$ RV X .



Simulation 157 ($\text{Uniform}(\theta_1, \theta_2)$) To simulate from $\text{Uniform}(\theta_1, \theta_2)$ RV X using the Inversion Sampler, we first need to find $F^{[-1]}(u)$ by solving for x in terms of $u = F(x; \theta_1, \theta_2)$:

$$u = \frac{x - \theta_1}{\theta_2 - \theta_1} \iff x = (\theta_2 - \theta_1)u + \theta_1 \iff F^{[-1]}(u; \theta_1, \theta_2) = \theta_1 + (\theta_2 - \theta_1)u$$

Here is a simple implementation of the Inversion Sampler for the $\text{Uniform}(\theta_1, \theta_2)$ RV in MATLAB:

```
>> rand('twister', 786); % initialise the fundamental sampler for Uniform(0,1)
>> theta1=-1; theta2=1; % declare values for parameters theta1 and theta2
>> u=rand; % rand is the Fundamental Sampler and u is a sample from it
>> x=theta1+(theta2 - theta1)*u; % sample from Uniform(-1,1) RV
>> disp(x); % display the sample from Uniform[-1,,1] RV
0.5134
```

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

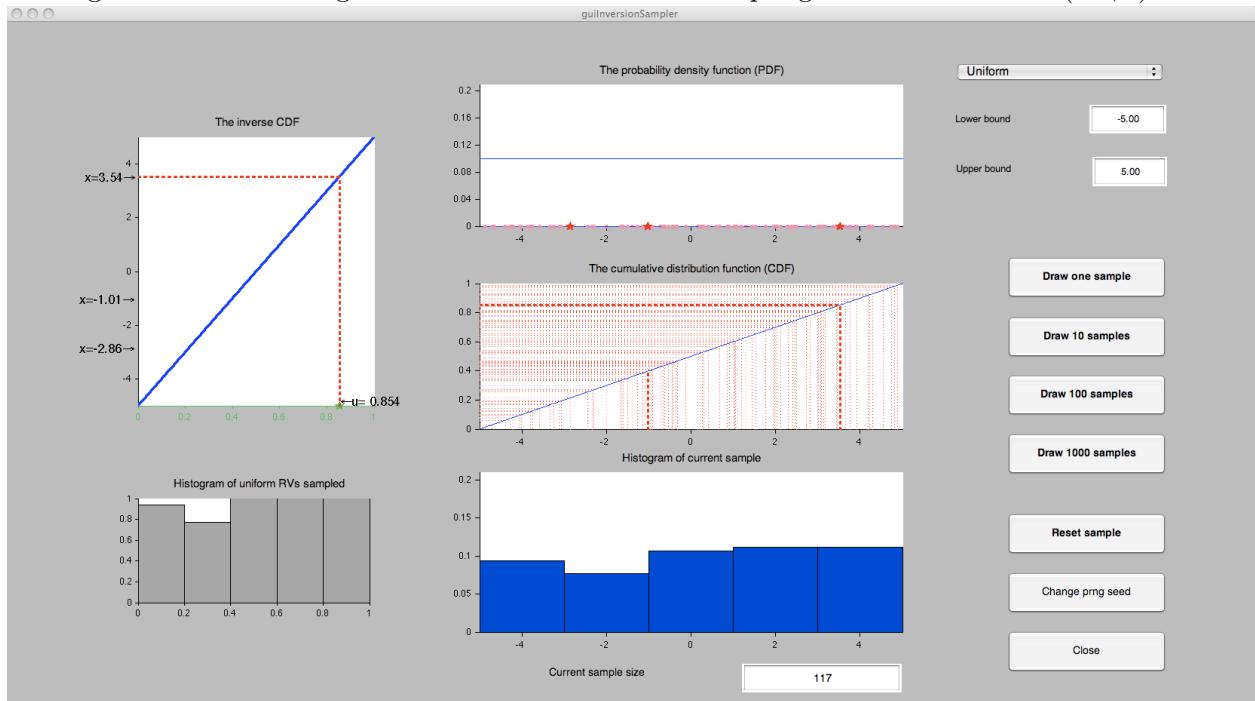
```

Labwork 158 (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 6.5. The default target distribution is 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 Uniform(θ_1, θ_2) RV.

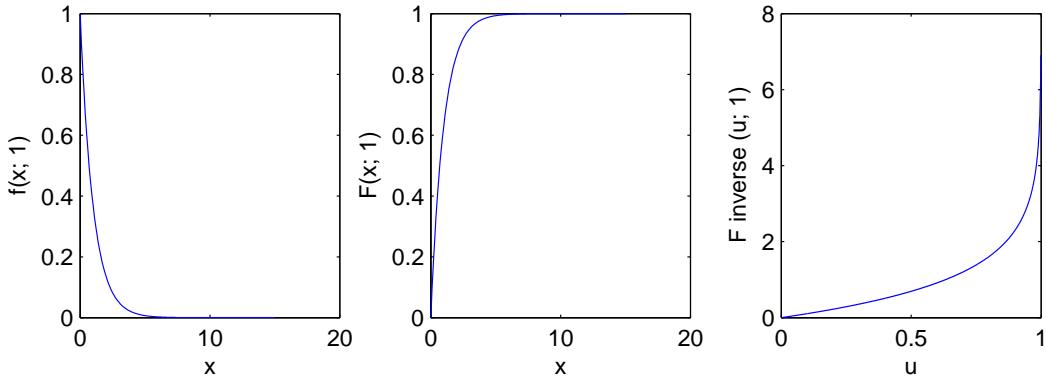
Figure 6.5: Visual Cognitive Tool GUI: Inversion Sampling from $X \sim \text{Uniform}(-5, 5)$.



Recall the Exponential(λ) RV of Model 8. Let us simulate from it using the inversion sampler.

Let us consider the problem of simulating from an Exponential(λ) 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.

Figure 6.6: The PDF f , DF F , and inverse DF $F^{[-1]}$ of the Exponential($\lambda = 1.0$) RV.



Simulation 159 (Exponential(λ)) For a given $\lambda > 0$, an Exponential(λ) RV has the following PDF f , DF F and inverse DF $F^{[-1]}$:

$$f(x; \lambda) = \lambda e^{-\lambda x} \quad F(x; \lambda) = 1 - e^{-\lambda x} \quad F^{[-1]}(u; \lambda) = \frac{-1}{\lambda} \log_e(1 - u) \quad (6.1)$$

We write the natural logarithm \log_e as log for notational simplicity. An implementation of the Inversion Sampler for Exponential(λ) as a function in the M-file:

```
function x = ExpInvCDF(u,lambda);
% Return the Inverse CDF of Exponential(lambda) RV X
% Call Syntax: x = ExpInvCDF(u,lambda);
%             ExpInvCDF(u,lambda);
% Input      : lambda = rate parameter,
%                 u = array of numbers in [0,1]
% Output     : x
x=-(1/lambda) * log(1-u);
```

We can simply call the function to draw a sample from, say the Exponential($\lambda = 1.0$) RV by:

```
lambda=1.0; % some value for lambda
u=rand; % rand is the Fundamental Sampler
ExpInvCDF(u,lambda) % sample from Exponential(1) RV via function in ExpInvCDF.m
```

Because of the following (recall Example 65):

$$U \sim \text{Uniform}(0, 1) \implies -U \sim \text{Uniform}(-1, 0) \implies 1 - U \sim \text{Uniform}(0, 1),$$

we could save a subtraction operation in the above algorithm by replacing $-(1/\lambda) * \log(1-u)$ by $-(1/\lambda) * \log(u)$. Recall that the transformation of $U \sim \text{Uniform}(0, 1)$ by $X = -(1/\lambda) \log(U)$ is exactly how we defined X as the Exponential(λ) RV in Model 8. This is implemented as the following function.

```
function x = ExpInvSam(u,lambda);
% Return the Inverse CDF based Sample from Exponential(lambda) RV X
% Call Syntax: x = ExpInvSam(u,lambda);
%             or ExpInvSam(u,lambda);
% Input      : lambda = rate parameter,
%                 u = array of numbers in [0,1] from Uniform[0,1] RV
% Output     : x
x=-(1/lambda) * log(u);
```

```

>> 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

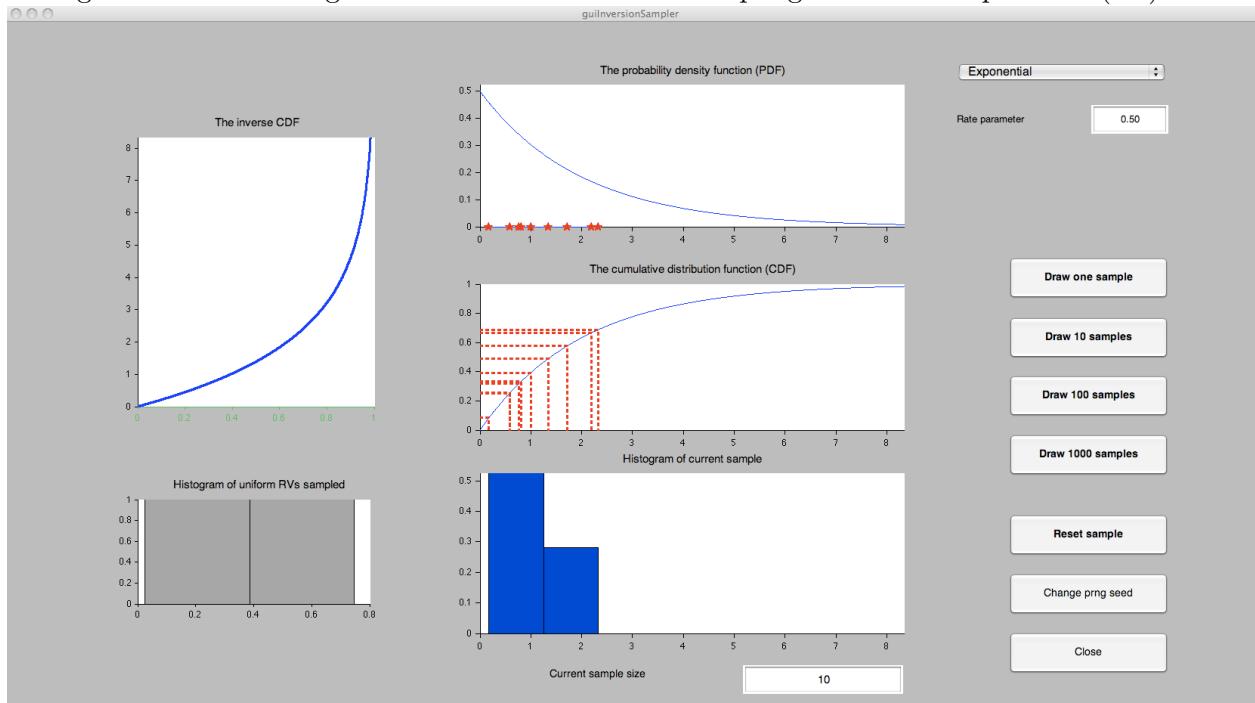
```

Labwork 160 (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 6.7. 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 6.7: 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 $\text{Exponential}(\lambda = 1.0)$ 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
0.4204  0.1219  2.2182  3.6692  0.9654  0.0093  1.7126
2.1427  0.1281  0.8500  1.4065  0.1160  0.1324  0.2635
0.6620  1.1729  0.6301  0.6375  0.3793  0.6525  0.8330
>> %sum up the outcomes of the sequence of 5 draws in each replicate
>> s=sum(a); disp(s);
5.2216   2.1856   6.6378   7.6490   2.0168   1.3073   7.9990

```

Labwork 161 (Next seven buses at your bus-stop) Consider the problem of modelling the arrival of buses at a bus stop. Suppose that the time between arrivals is an $\text{Exponential}(\lambda = 0.1)$ RV X with a mean inter-arrival time of $1/\lambda = 10$ minutes. Suppose you go to your bus stop and zero a stop-watch. Simulate the times of arrival for the next seven buses as indicated by your stop-watch. Seed the fundamental sampler by your Student ID (eg. if your ID is 11424620 then type `rand('twister', 11424620)`; just before the simulation). Hand in the code with the arrival times of the next seven buses at your ID-seeded bus stop.

The support of the $\text{Exponential}(\lambda)$ RV is $\mathbb{R}_+ := [0, \infty)$. Let us consider a RV built by mirroring the $\text{Exponential}(\lambda)$ RV about the origin with the entire real line as its support.

Model 19 (Laplace(λ) or Double Exponential(λ) RV) If a RV X is equally likely to be either positive or negative with an exponential density, then the Laplace(λ) or Double Exponential(λ) RV, with the rate parameter $\lambda > 0, \lambda \in \mathbb{R}$, may be used to model it. The density function for the Laplace(λ) RV given by $f(x; \lambda)$ is

$$f(x; \lambda) = \frac{\lambda}{2} e^{-\lambda|x|} = \begin{cases} \frac{\lambda}{2} e^{\lambda x} & \text{if } x < 0 \\ \frac{\lambda}{2} e^{-\lambda x} & \text{if } x \geq 0 \end{cases}. \quad (6.2)$$

Let us define the sign of a real number x by

$$\text{sign}(x) = \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{if } x = 0 \\ -1 & \text{if } x < 0 \end{cases}.$$

Then, the DF of the Laplace(λ) RV X is

$$F(x; \lambda) = \int_{-\infty}^x f(y; \lambda) dy = \frac{1}{2} \left(1 + \text{sign}(x) \left(1 - e^{-\lambda|x|} \right) \right), \quad (6.3)$$

and its inverse DF is

$$F^{[-1]}(u; \lambda) = -\frac{1}{\lambda} \text{sign} \left(u - \frac{1}{2} \right) \log \left(1 - 2 \left| u - \frac{1}{2} \right| \right), \quad u \in [0, 1] \quad (6.4)$$

Mean and Variance of Laplace(λ) RV X : Show that the mean of a Laplace(λ) RV X is

$$\mathbb{E}(X) = \int_0^\infty x f(x; \lambda) dx = \int_0^\infty x \frac{\lambda}{2} e^{-\lambda|x|} dx = 0,$$

and the variance is

$$\text{V}(X) = \left(\frac{1}{\lambda} \right)^2 + \left(\frac{1}{\lambda} \right)^2 = 2 \left(\frac{1}{\lambda} \right)^2.$$

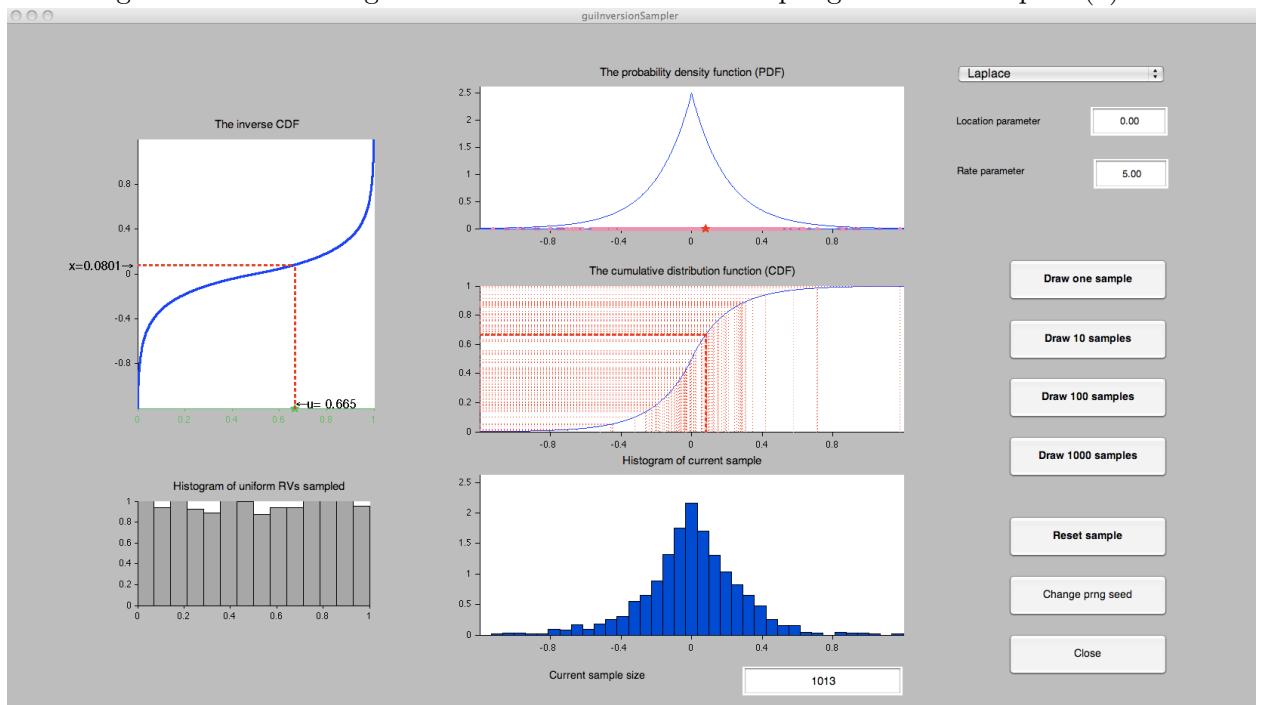
Note that the mean is 0 due to the symmetry of the density about 0 and the variance is twice that of the $\text{Exponential}(\lambda)$ RV.

Labwork 162 (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 6.8. 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 (6.2) with an additional location parameter μ .

Figure 6.8: Visual Cognitive Tool GUI: Inversion Sampling from $X \sim \text{Laplace}(5)$.



Simulation 163 (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));
```

We can simply call the function to draw a sample from, say the Laplace($\lambda = 1.0$) RV by

```
>> lambda=1.0; % some value for lambda
>> rand('twister',6567); % initialize the fundamental sampler
>> u=rand(1,5); % draw 5 IID samples from Uniform(0,1) RV
>> disp(u); % display the samples in u
0.6487 0.9003 0.3481 0.6524 0.8152

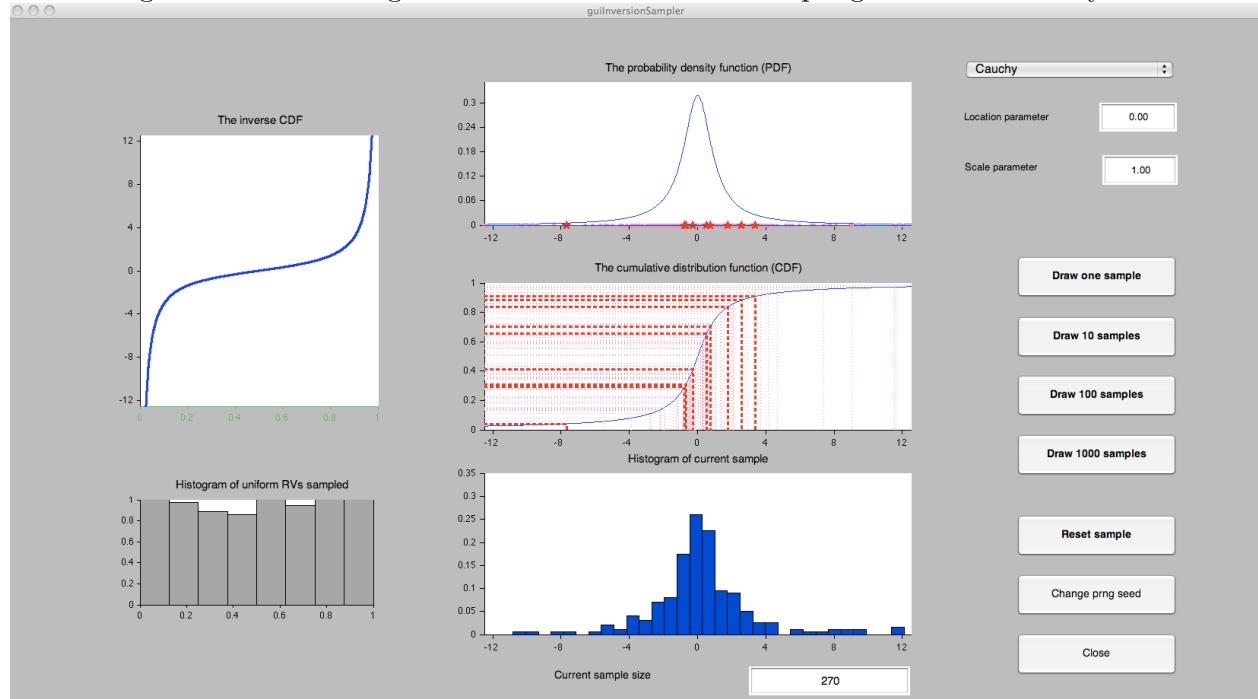
>> x=LaplaceInvCDF(u,lambda); % draw 5 samples from Laplace(1) RV using inverse CDF
>> disp(x); % display the samples
0.3530 1.6127 -0.3621 0.3637 0.9953
```

Labwork 164 (Inversion Sampler Demo – Cauchy) Let us comprehend 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 6.9. 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 comprehend the simulation process. You can also press “Draw 10 samples” several times and see the density histogram of the generated samples. 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 .

Figure 6.9: Visual Cognitive Tool GUI: Inversion Sampling from $X \sim \text{Cauchy}$.



Simulation 165 (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
```

6.3.2 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 77 (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:

$$\Pr(F^{[-1]}(U) \leq x) = \Pr(U \leq F(x)) = F(x) := \Pr(X \leq x)$$

Simulation 166 (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.

Simulation 167 (Point Mass(θ)) Let us simulate a sample from the Point Mass(θ) RV X . Since this RV produces the same realisation θ we can implement it via the following M-file:

```

function x = Sim1PointMass(u,theta)
% Returns one sample from the Point Mass(theta) RV X
% Call Syntax: x = Sim1PointMass(u,theta);
% Input      : u = one uniform random number eg. rand()
%                 theta = a real number (scalar)
% Output     : x = sample from X
x=theta;

```

Here is call to the function.

```

>> Sim1PointMass(rand(),2)
ans =
    2
>> % we can use arrayfun to apply Sim1PointMass to any array of Uniform(0,1) samples
>> arrayfun(@(u)(Sim1PointMass(u,17)),rand(2,10))
ans =
    17    17    17    17    17    17    17    17    17    17
    17    17    17    17    17    17    17    17    17    17

```

Note that it is not necessary to have input IID samples from Uniform(0, 1) RV via `rand` in order to draw samples from the Point Mass(θ) RV. For instance, an input matrix of zeros can do the job:

```

>> arrayfun(@(u)(Sim1PointMass(u,17)),zeros(2,8))
ans =
    17    17    17    17    17    17    17    17
    17    17    17    17    17    17    17    17

```

Next we simulate from de Moivre($\theta_1, \theta_2, \dots, \theta_k$) RV X of Model 14 via its inverse DF

$$F^{[-1]} : [0, 1] \rightarrow [k] := \{1, 2, \dots, k\},$$

given by:

$$F^{[-1]}(u; \theta_1, \theta_2, \dots, \theta_k) = \begin{cases} 1 & \text{if } 0 \leq u < \theta_1 \\ 2 & \text{if } \theta_1 \leq u < \theta_1 + \theta_2 \\ 3 & \text{if } \theta_1 + \theta_2 \leq u < \theta_1 + \theta_2 + \theta_3 \\ \vdots & \vdots \\ k & \text{if } \theta_1 + \theta_2 + \dots + \theta_{k-1} \leq u < 1 \end{cases} \quad (6.5)$$

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 6.10.

First we simulate from an equi-probable special case of the de Moivre($\theta_1, \theta_2, \dots, \theta_k$) RV, with $\theta_1 = \theta_2 = \dots = \theta_k = 1/k$.

Simulation 168 (de Moivre($1/k, 1/k, \dots, 1/k$)) The equi-probable de Moivre($1/k, 1/k, \dots, 1/k$) RV X with a discrete uniform distribution over $[k] = \{1, 2, \dots, k\}$ can be efficiently sampled using the ceiling function. Recall that $\lceil y \rceil$ is the smallest integer larger than or equal to y , eg. $\lceil 13.1 \rceil = 14$. Algorithm 4 produces samples from the de Moivre($1/k, 1/k, \dots, 1/k$) RV X .

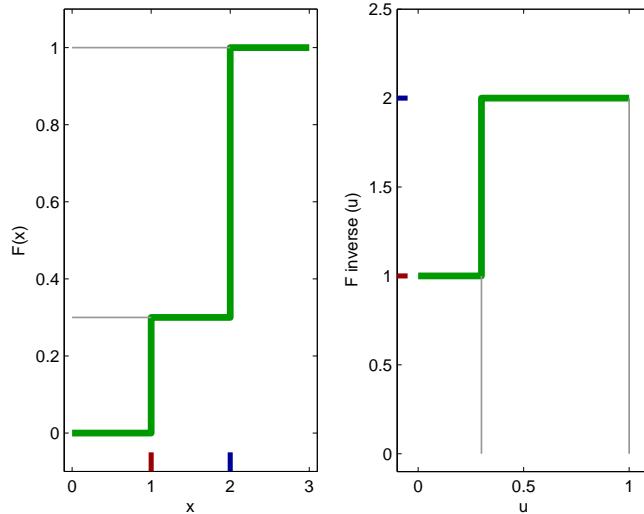
The M-file implementing Algorithm 4 is:

```

function x = SimdeMoivreEqui(u,k);
% return samples from de Moivre(1/k,1/k,...,1/k) RV X
% Call Syntax: x = SimdeMoivreEqui(u,k);

```

Figure 6.10: 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:* $x \leftarrow \lceil ku \rceil$
-

```
% 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 $1/k$.

Algorithm 5 Inversion Sampler for de Moivre($\theta_1, \theta_2, \dots, \theta_k$) RV X

1: *input:*1. parameter vector $(\theta_1, \theta_2, \dots, \theta_k)$ of de Moivre($\theta_1, \theta_2, \dots, \theta_k$) RV X .2. $u \sim \text{Uniform}(0, 1)$ 2: *output:* a sample from X 3: *initialise:* $F \leftarrow \theta_1, i \leftarrow 1$ 4: **while** $u > F$ **do**5: $i \leftarrow i + 1$ 6: $F \leftarrow F + \theta_i$ 7: **end while**8: *return:* $x \leftarrow i$

Simulation 169 (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 when $(\theta_1, \theta_2, \dots, \theta_k)$ are specifiable as an input vector via the following algorithm.

The M-file implementing Algorithm 5 is:

SimdeMoivreOnce.m

```
function x = SimdeMoivreOnce(u,thetas)
% Returns a sample from the de Moivre(thetas=(theta_1,...,theta_k)) RV X
% Call Syntax: x = SimdeMoivreOnce(u,thetas);
%           deMoivreEqui(u,thetas);
% Input      : u = a uniform random number eg. rand
%             thetas = an array of probabilities thetas=[theta_1 ... theta_k]
% Output     : x = sample from X
x=1; % initial index is 1
cum_theta=thetas(x);
while u > cum_theta;
    x=x+1;
    cum_theta = cum_theta + thetas(x);
end
```

Let us use the function `deMoivreEqui` to draw five samples from a fair seven-faced dice.

```
>> k=7; % number of faces of the fair dice
>> n=5; % number of trials
>> rand('twister',78657); % initialise the fundamental sampler
>> Us=rand(1,n); % draw n samples from Uniform(0,1)
>> disp(Us);
    0.8330    0.6819    0.6468    0.6674    0.2577
>> % 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)
>> f=[1/7 1/7 1/7 1/7 1/7 1/7 1/7];
>> disp(f);
    0.1429    0.1429    0.1429    0.1429    0.1429    0.1429
>> % use funarray to apply function-handled SimdeMoivreOnce to
>> % each element of array Us and save it in array outcomes2
>> outcomes2=arrayfun(@(u)(SimdeMoivreOnce(u,f)),Us);
>> disp(outcomes2);
    6      5      5      5      2
>> disp(SimdeMoivreEqui(u,k)); % same result using the previous algorithm
    6      5      5      5      2
```

Clearly, Algorithm 5 may be used to sample from any de Moivre($\theta_1, \dots, \theta_k$) RV X . We demonstrate this by producing five samples from a randomly generated PMF `f2`.

```

>> 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 170 (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

```

Labwork 171 (PMF versus relative frequency histogram of simulated Geometric(θ) RV)

It is a good idea to make a relative frequency histogram of a simulation algorithm and compare that to the PDF of the discrete RV we are simulating from. We use the following script to create Figure 3.7:

```

theta=0.5;
SampleSize=1000;
% simulate 1000 samples from Geometric(theta) RV
Samples=floor(log(rand(1,SampleSize))/ log (1-theta));
Xs = 0:10; % get some values for x
RelFreqs=hist(Samples,Xs)/SampleSize; % relative frequencies of Samples
stem(Xs,theta*((1-theta) .^ Xs),'*')% PDF of Geometric(theta) over Xs
hold on;
plot(Xs,RelFreqs,'o')% relative frequency histogram
RelFreqs100=hist(Samples(1:100),Xs)/100; % Relative Frequencies of first 100 samples
plot(Xs,RelFreqs100,'x')
legend('PDF of Geometric(0.5)', 'Relative freq. hist. (1000 samples)', ...
'Relative freq. hist. (100 samples)')

```

Let us simulate from the Binomial(n, θ) RV of Model 5.

Labwork 172 (Binomial coefficient) The MATLAB function `BinomialCoefficient` can be used to compute:

$$\binom{n}{x} = \frac{n!}{x!(n-x)!} = \frac{n(n-1)(n-2)\dots(n-x+1)}{x(x-1)(x-2)\cdots(2)(1)} = \frac{\prod_{i=(n-x+1)}^n i}{\prod_{i=2}^x i},$$

with the following M-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
NminusX = n-x;
NumeratorPostCancel = prod(n:-1:(max([NminusX,x])+1));
DenominatorPostCancel = prod(2:min([NminusX, 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 = BinomialPdf(x,n,theta)
% f is the prob mass function for the Binomial(x;n,theta) RV
% and x can be array of samples.
% Values of x are integers in [0,n] and theta is a number in [0,1]
fx = zeros(size(x));
fx = arrayfun(@(xi)(BinomialCoefficient(n,xi)),x);
fx = fx .* (theta .^ x) .* (1-theta) .^ (n-x);

```

For example, we can compute the desired PDF for an array of samples \mathbf{x} from $\text{Binomial}(8, 0.5)$ RV X , as follows:

```
>> 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 173 ($\text{Binomial}(n, \theta)$ as $\sum_{i=1}^n \text{Bernoulli}(\theta)$) Since the $\text{Binomial}(n, \theta)$ RV X is the sum of n IID $\text{Bernoulli}(\theta)$ RVs we can also simulate from X by first simulating n IID $\text{Bernoulli}(\theta)$ 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      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 173, the number of IID $\text{Bernoulli}(\theta)$ RVs needed to simulate one sample from the $\text{Binomial}(n, \theta)$ RV is exactly n . Thus, as n increases, the amount of time needed to simulate from $\text{Binomial}(n, \theta)$ is $O(n)$, i.e. linear in n . We can simulate more efficiently by exploiting a simple relationship between the $\text{Geometric}(\theta)$ RV and the $\text{Binomial}(n, \theta)$ RV.

The $\text{Binomial}(n, \theta)$ RV X is related to the IID $\text{Geometric}(\theta)$ RV Y_1, Y_2, \dots : X is the number of successful $\text{Bernoulli}(\theta)$ outcomes (outcome is 1) that occur in a total of n $\text{Bernoulli}(\theta)$ trials, with the number of trials between consecutive successes distributed according to IID $\text{Geometric}(\theta)$ RV.

Simulation 174 ($\text{Binomial}(\theta)$ from IID $\text{Geometric}(\theta)$ RVs) By this principle, we can simulate from the $\text{Binomial}(\theta)$ X by Step 1: generating IID $\text{Geometric}(\theta)$ 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

```

Here is a call to simulate 12 samples from $\text{Binomial}(n = 10, \theta = 0.5)$ RV:

```

>> theta=0.5; % declare theta
>> n=10; % say n=10
>> SampleSize=12;% say you want to simulate 12 samples
>> rand('twister',10001) % seed the fundamental sampler
>> Samples=arrayfun(@(T)Sim1BinomByGeoms(n,T),theta*ones(1,SampleSize))
Samples =    7      5      8      8      4      1      4      8      2      4      6      5

```

Figure 3.8 depicts a comparison of the PDF of $\text{Binomial}(n = 10, \theta = 0.5)$ RV and a relative frequency histogram based on 100,000 simulations from it.

Let us simulate from the $\text{Poisson}(\lambda)$ RV of Model 6 as shown in Figure 3.18.

Simulation 175 (Poisson(λ) from IID Exponential(λ) RVs) By this principle, we can simulate from the $\text{Poisson}(\lambda)$ X by Step 1: generating IID Exponential(λ) RVs Y_1, Y_2, \dots , Step 2: stopping as soon as $\sum_{i=1}^k Y_i \geq 1$ and Step 3: setting $x \leftarrow k - 1$.

We implement the above algorithm via the following M-file:

function x = Sim1Poisson(lambda)
% Simulate one sample from Poisson(lambda) via Exponentials
YSum=0; k=0; % initialise
while (YSum < 1),
 YSum = YSum + -(1/lambda) * log(rand);
 k=k+1;
end
x=k-1; % return x

Sim1Poisson.m

Here is a call to simulate 10 samples from $\text{Poisson}(\lambda = 10.0)$ and $\text{Poisson}(\lambda = 0.1)$ RVs:

```

>> arrayfun(@(lambda)Sim1Poisson(lambda),10.0*ones(1,10)) % lambda=10.0
ans =    14      7     10     13     11      3      6      5      8      5
>> arrayfun(@(lambda)Sim1Poisson(lambda),0.1*ones(1,10)) % lambda=0.1
ans =      2      0      0      0      0      0      0      0      0      0

```

Figure 3.18 depicts a comparison of the PDF of $\text{Poisson}(\lambda = 10)$ RV and a relative frequency histogram based on 1000 simulations from it.

Simulating from a $\text{Poisson}(\lambda)$ RV is also a special case of simulating from the following more general RV.

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 is defined as follows:

$$f(X = x; \theta_0, \theta_1, \dots) = \begin{cases} 0, & \text{if } x \notin \{0, 1, 2, \dots\} \\ \theta_0, & \text{if } x = 0 \\ \theta_1, & \text{if } x = 1 \\ \vdots & \end{cases}$$

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:*

1. θ_0 and $\{C(i) = \theta_i / \theta_{i-1}\}$ for any $i \in \{1, 2, 3, \dots\}$.
2. $u \sim \text{Uniform}(0, 1)$

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 176 ($\text{Binomial}(n, \theta)$) To simulate from a $\text{Binomial}(n, \theta)$ 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 $\text{Geometric}(\theta)$ and $\text{Poisson}(\lambda)$ RVs.

Labwork 177 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 `@`) that can be passed as a parameter to the function `MyGenDiscInvSampler`].
2. Show that your code works for drawing samples from a $\text{Binomial}(n, p)$ 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 $\text{Binomial}(n = 20, p = 0.5)$ 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 $\text{Binomial}(n, p)$ RV is given above].
3. Show that your code works for drawing samples from a $\text{Geometric}(p)$ RV by doing the following:

- (a) Seed the fundamental sampler by your Student ID.
- (b) Set the variable `Mytheta=rand`.
- (c) Draw 100 samples from the `Geometric(Mytheta)` 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 shown].

6.3.3 von Neumann Rejection Sampler (RS)

Rejection sampling [John von Neumann, 1947, in *Stanislaw Ulam 1909-1984*, a special issue of Los Alamos Science, Los Alamos National Lab., 1987, p. 135-136] is a Monte Carlo method to draw independent samples from a target RV X with probability density $f(x)$, where $x \in \mathbb{X} \subset \mathbb{R}^k$. Typically, the target density f is only known up to a constant and therefore the (normalised) density f itself may be unknown and it is difficult to generate samples directly from X .

Suppose we have another density or mass function g for which the following are true:

- (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 $a > 1$ exists, such that:

$$f(x) \leq ag(x). \quad (6.6)$$

for any $x \in \mathbb{X}$, the support of X . Then x can be generated from Algorithm 7.

Algorithm 7 Rejection Sampler (RS) of von Neumann

1: *input*:

- (1) a target density $f(x)$,
- (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** $u \leq \frac{f(y)}{ag(y)}$

6: **return**: $x \leftarrow y$

Proposition 78 (Fundamental Theorem of Simulation) The von Neumann rejection sampler of Algorithm 7 produces a sample x from the random variable X with density $f(x)$.

Proof: We shall prove the result for the continuous case. For any real number t :

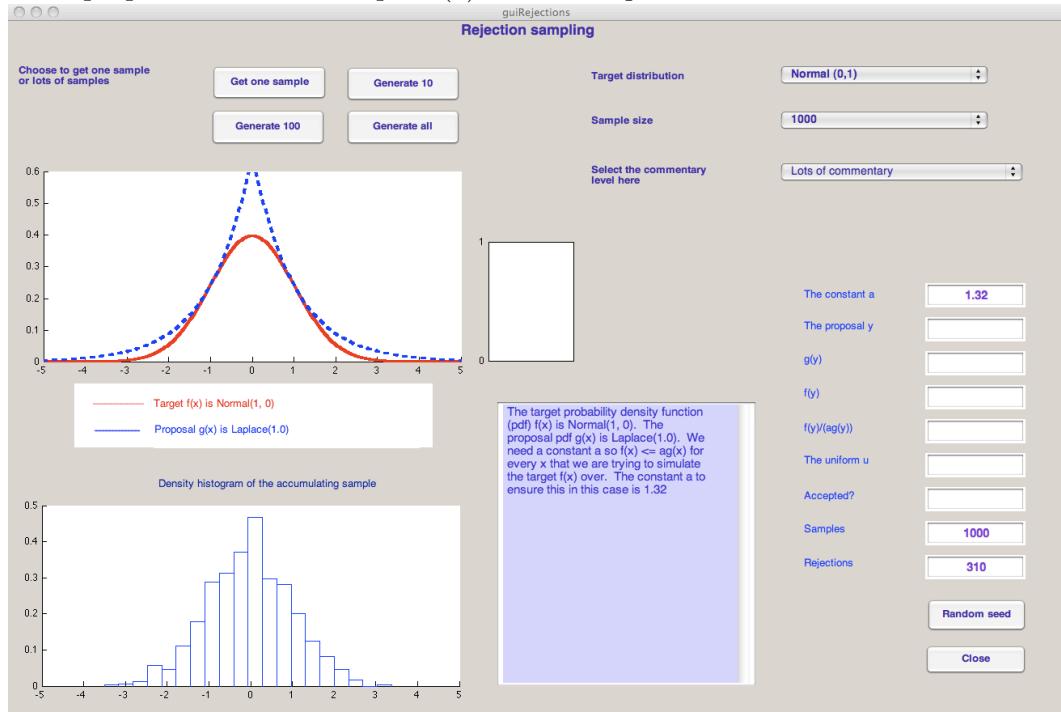
$$\begin{aligned} F(t) &= P(X \leq t) = P\left(Y \leq t \mid U \leq \frac{f(Y)}{ag(Y)}\right) = \frac{P\left(Y \leq t, U \leq \frac{f(Y)}{ag(Y)}\right)}{P\left(U \leq \frac{f(Y)}{ag(Y)}\right)} \\ &= \frac{\int_{-\infty}^t \left(\int_0^{f(y)/ag(y)} 1 du \right) g(y) dy}{\int_{-\infty}^{\infty} \left(\int_0^{f(y)/ag(y)} 1 du \right) g(y) dy} = \frac{\int_{-\infty}^t \left(\frac{f(y)}{ag(y)} \right) g(y) dy}{\int_{-\infty}^{\infty} \left(\frac{f(y)}{ag(y)} \right) g(y) dy} \\ &= \int_{-\infty}^t f(y) dy \end{aligned}$$

Labwork 178 (Rejection Sampler Demo) Let us understand the rejection sampler by calling the interactive visual cognitive tool:

```
>> guiRejections
```

The M-file `guiRejections.m` will bring a graphical user interface (GUI) as shown in Figure 6.11. 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 6.11: 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 179 (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 163). 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:

$$\frac{f(y)}{ag(y)} = \frac{f(y)}{g(y)a} = \sqrt{\frac{2}{\pi}} \exp\left(|y| - \frac{y^2}{2}\right) \frac{1}{\sqrt{\frac{2}{\pi}} \exp\left(\frac{1}{2}\right)} = \exp\left(|y| - \frac{y^2}{2} - \frac{1}{2}\right)$$

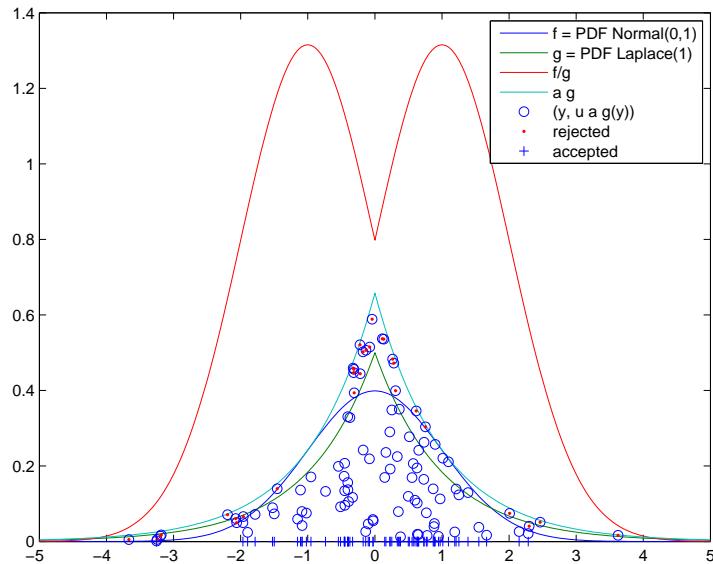
Let us implement a rejection sampler as a function in the M-file `RejectionNormalLaplace.m` by reusing the function in `LaplaceInvCDF.m`.

```
RejectionNormalLaplace.m
function x = RejectionNormalLaplace()
Accept = 0; % a binary variable to indicate whether a proposed point is accepted
while ~Accept % ~ is the logical NOT operation
    y = LaplaceInvCDF(rand(),1); % sample Laplace(1) RV
    Bound = exp( abs(y) - (y*y+1)/2 );
    u = rand();
    if u <= Bound
        x = y;
        Accept = 1;
    end % if
end % while
```

We may obtain a large number of samples and plot them as a histogram using the following commands:

```
>> % use funarray to convert 1000 zeros into samples from the Normal(0,1)
>> y=arrayfun(@(x)(RejectionNormalLaplace()),zeros(1,1000));
>> hist(y,20) % histogram with 20 bins
```

Figure 6.12: Rejection Sampling from $X \sim \text{Normal}(0, 1)$ with PDF f based on 100 proposals from $Y \sim \text{Laplace}(1)$ with PDF g .



Classwork 180 (A note on the proposal's tail in rejection sampling) The condition $f(x) \leq ag(x)$ is equivalent to $f(x)/g(x) \leq a$, which says that $f(x)/g(x)$ must be bounded; therefore, g must have higher tails than f . The rejection method cannot be used to generate from a Cauchy distribution using a normal distribution, because the latter has lower tails than the former.

The next result tells us how many iterations of the algorithm are needed, on average, to get a sample value from a RV with PDF f .

Proposition 79 (Acceptance Probability of RS) The expected number of iterations of the rejection algorithm to get a sample x is the constant a .

Proof: For the continuous case:

$$P(\text{'accept } y') = P\left(u \leq \frac{f(y)}{ag(y)}\right) = \int_{-\infty}^{\infty} \left(\int_0^{f(y)/ag(y)} du \right) g(y) dy = \int_{-\infty}^{\infty} \frac{f(y)}{ag(y)} g(y) dy = \frac{1}{a}.$$

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:

$$\tilde{f}(x) \leq \tilde{a}\tilde{g}(x), \quad (6.7)$$

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 181 (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.

```

>> randn('state',67678); % initialise the seed at 67678 and method as Ziggurat -- TYPE help randn
>> randn % produce 1 sample from Normal(0,1) RV
ans =
    1.5587
>> randn(2,8) % produce an 2 X 8 array of samples from Normal(0,1) RV
ans =
    1.2558    0.7834    0.6612    0.3247    0.1407    1.0562    0.8034    1.2970
   -0.5317    0.0417   -0.3454    0.6182   -1.4162    0.4796   -1.5015    0.3718

```

If we want to produce samples from $X \sim \text{Normal}(\mu, \sigma^2)$ with some user-specified μ and σ , then we can use the following relationship between X and $Z \sim \text{Normal}(0, 1)$:

$$X \leftarrow \mu + \sigma Z, \quad Z \sim \text{Normal}(0, 1).$$

Suppose we want samples from $X \sim \text{Normal}(\mu = \pi, \sigma^2 = 2)$, then we can do the following:

```

>> randn('state',679); % initialise the seed at 679 and method as Ziggurat -- TYPE help randn
>> mu=pi % set the desired mean parameter mu
mu =
    3.1416
>> sigma=sqrt(2) % set the desired standard deviation parameter sigma
sigma =
    1.4142
>> mu + sigma * randn(2,8) % produces a 2 X 8 array of samples from Normal(3.1416,1.4.42)
ans =
    1.3955    1.7107    3.9572    3.2618    6.1652    2.6971    2.4940    4.5928
    0.8442    4.7617    3.5397    5.0282    1.6139    5.0977    2.0477    2.3286

```

Labwork 182 (Sampling from truncated normal distributions) [Christian P. Robert, Simulation of truncated normal variables, Statistics and Computing (1995) 5, 121-125] Let $N_+(\mu, \tau, \sigma^2)$ denote the left-truncated normal distribution with truncation point τ and density given by

$$f(x|\mu, \tau, \sigma^2) = \frac{\exp(-(x-\mu)^2/2\sigma^2)}{\sqrt{2\pi}\sigma[1 - \Phi((\tau-\mu)/\sigma)]} \mathbf{1}_{x \geq \tau}.$$

When $\tau < \mu$, the rejection sampler can readily be used to simulate from $N_+(\mu, \tau, \sigma^2)$ by simulating from $\text{Normal}(\mu, \sigma^2)$ until a number larger than τ is obtained. When $\tau > \mu$, however, this can be inefficient and increasingly 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 as the proposal distribution:

$$g(y|\lambda, \tau) = \lambda \exp(-\lambda(y-\tau)) \mathbf{1}_{y \geq \tau}.$$

1. Show that for simulating from $N_+(\mu = 0, \tau, \sigma^2 = 1)$ when $\tau \geq 0$, the best choice of λ that maximizes the expected acceptance probability for the rejection sampler is given by

$$\lambda = \frac{\tau + \sqrt{\tau^2 + 4}}{2}$$

2. Find the maximum expected acceptance probabilities for the following truncation points, $\tau = 0, 0.5, 1, 1.5, 2, 2.5$ and 3 . What can you conclude about efficiency as τ gets further out into the right tail?
3. Describe how samples from $N_+(\mu, \tau, \sigma^2)$ can be obtained by simulating from $N_+(\mu = 0, \tau, \sigma^2 = 1)$ and using location-scale transformation.

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 .

6.4 Exercises in Simulation

Ex. 6.1 — 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.

Chapter 7

Limit Laws of Statistics

7.1 Convergence of Random Variables

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 (Ω, \mathcal{F}, P) .

$$\{X_i\}_{i=1}^n := X_1, X_2, X_3, \dots, X_{n-1}, X_n \quad \text{as } n \rightarrow \infty.$$

From a statistical or decision-making viewpoint, as you will see in Inference 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 .

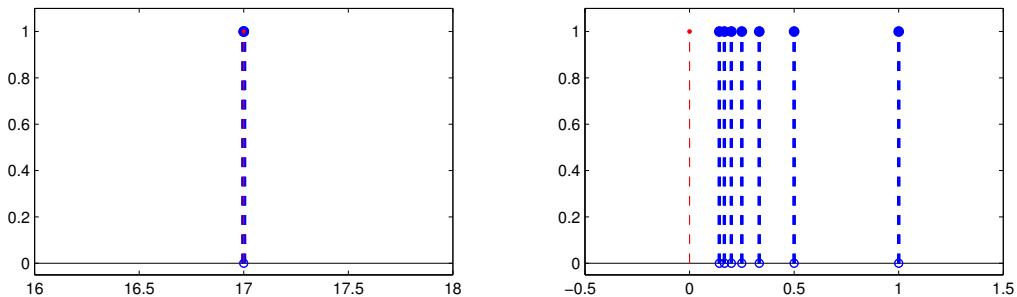
We need different notions of convergence to characterize such a behavior: two simplest behaviors are that the sequence eventually takes a constant value θ , i.e. X_n approaches $X \sim \text{Point Mass}(\theta)$ RV, or that values in the sequence continue to change but can be described by an unchanging probability distribution, i.e., X_n approaches $X \sim F(x)$. See https://en.wikipedia.org/wiki/Convergence_of_random_variables.

Let us first refresh ourselves with notions of convergence, limits and continuity in the real line (Sec. 1.8.1) before proceeding further.

Can the sequences of $\{\text{Point Mass}(\theta_i = 17)\}_{i=1}^{\infty}$ and $\{\text{Point Mass}(\theta_i = 1/i)\}_{i=1}^{\infty}$ RVs be the same as the two sequences of real numbers $\{x_i\}_{i=1}^{\infty} = 17, 17, 17, \dots$ and $\{x_i\}_{i=1}^{\infty} = \frac{1}{1}, \frac{1}{2}, \frac{1}{3}, \dots$ we saw in Examples 12 and 13?

Yes why not – just move to space of distributions over the reals! See Figure 7.1.

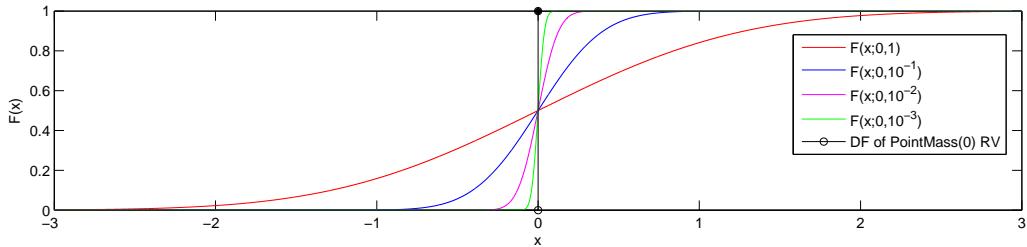
Figure 7.1: Sequence of $\{\text{Point Mass}(17)\}_{i=1}^{\infty}$ RVs (left panel) and $\{\text{Point Mass}(1/i)\}_{i=1}^{\infty}$ RVs (only the first seven are shown on right panel) and their limiting RVs in red.



Classwork 183 (Convergence of $X_i \sim \text{Normal}(0, 1/i)$) Suppose you are given an independent sequence of RVs $\{X_i\}_{i=1}^n$, where $X_i \sim \text{Normal}(0, 1/i)$. How would you talk about the convergence of $X_n \sim \text{Normal}(0, 1/n)$ as n approaches ∞ ? Take a look at Figure 7.2 for insight. The probability mass of X_n increasingly concentrates about 0 as n approaches ∞ and the variance $1/n$ approaches 0, as depicted in Figure 7.2. Based on this observation, can we expect $\lim_{n \rightarrow \infty} X_n = X$, where the limiting RV $X \sim \text{Point Mass}(0)$?

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 large. In other words, a continuous RV, such as X_n , has 0 probability of realizing any single real number in its support, such as 0.

Figure 7.2: 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 80 (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:

$$\begin{aligned} \lim_{n \rightarrow \infty} P(\{\omega : X_n(\omega) \leq t\}) &= P(\{\omega : X(\omega) \leq t\}), \\ \text{i.e. } P(\{\omega : X_n(\omega) \leq t\}) &\rightarrow P(\{\omega : X(\omega) \leq t\}), \quad \text{as } n \rightarrow \infty. \end{aligned}$$

Let us revisit the problem of convergence in Classwork 183 armed with our new notions of convergence.

Example 184 (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 . We can formalize our observation in Classwork 183 that X_n is concentrating about 0 as $n \rightarrow \infty$ by the statement:

$$X_n \text{ is converging in distribution to } X, \text{ ie, } X_n \rightsquigarrow X.$$

Proof: To check that the above statement is true we need to verify that the definition of convergence 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 , $\lim_{n \rightarrow \infty} F_n(t) = F(t)$. First note that

$$X_n \sim \text{Normal}(0, 1/n) \implies Z := \sqrt{n}X_n \sim \text{Normal}(0, 1),$$

and thus

$$F_n(t) = P(X_n < t) = P(\sqrt{n}X_n < \sqrt{nt}) = P(Z < \sqrt{nt}).$$

The only discontinuous point of F is 0 where F jump from 0 to 1.

When $t < 0$, $F(t)$, being the constant 0 function over the interval $(-\infty, 0)$, is continuous at t . Since $\sqrt{nt} \rightarrow -\infty$, as $n \rightarrow \infty$,

$$\lim_{n \rightarrow \infty} F_n(t) = \lim_{n \rightarrow \infty} P(Z < \sqrt{nt}) = 0 = F(t).$$

And, when $t > 0$, $F(t)$, being the constant 1 function over the interval $(0, \infty)$, is again continuous at t . Since $\sqrt{nt} \rightarrow \infty$, as $n \rightarrow \infty$,

$$\lim_{n \rightarrow \infty} F_n(t) = \lim_{n \rightarrow \infty} P(Z < \sqrt{nt}) = 1 = F(t).$$

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: $\lim_{n \rightarrow \infty} F_n(t) = F(t)$.

However, note that

$$F_n(0) = \frac{1}{2} \neq F(0) = 1,$$

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 .

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 $\mathbb{1}_{(0,1)}(x)(1 - \cos(2\pi nx))$. 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 7.3.

Proposition 81 (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.

Proof: We will state this without a Proof here as Proof of the Theorem requires measure theory 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 $n \rightarrow \infty$.

¹See https://en.wikipedia.org/wiki/Scheff%C3%A9%27s_lemma.

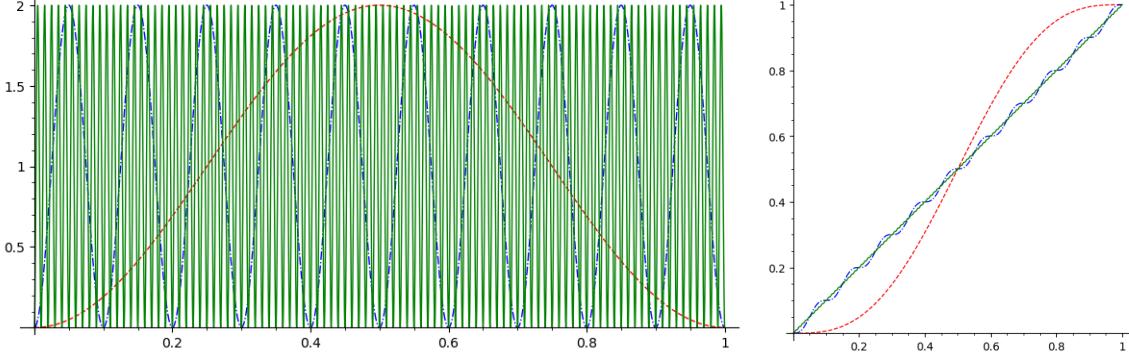


Figure 7.3: PDF $f_{X_n}(x) := \mathbb{1}_{(0,1)}(x)(1 - \cos(2\pi n x))$ of the RV X_n [the left sub-figure] and its DF $F_n(x) := \int_{-\infty}^x \mathbb{1}_{(0,1)}(v)(1 - \cos(2\pi n v)) dv$ [the right sub-figure], for $n = 1$ [red '---'], $n = 10$ [blue '-.-'], and $n = 100$ [green '-'], respectively. One can see clear convergence of the DFs F_n to $\mathbb{1}_{(0,1)}(x)x$, the DF of the Uniform(0, 1) RV, while the corresponding PDFs $f_n(x)$ keep oscillating wildly with n across $[0, 2]$ about $\mathbb{1}_{(0,1)}(x)$, 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 $F(x) = P(X \leq x)$, 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 $\lim_{n \rightarrow \infty} P(X_n = x) = P(X = x)$ 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 185 ($\text{Binomial}(n, \lambda/n) \rightsquigarrow \text{Poisson}(\lambda)$) 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 $\lim_{n \rightarrow \infty} P(X_n = x) = P(Y = x) = e^{-\lambda} \lambda^x / x!$ for any

$x \in \{0, 1, 2, 3, \dots, n\}$.

$$\begin{aligned}
P(X = x) &= \binom{n}{x} \left(\frac{\lambda}{n}\right)^x \left(1 - \frac{\lambda}{n}\right)^{n-x} \\
&= \frac{n(n-1)(n-2)\cdots(n-x+1)}{x(x-1)(x-2)\cdots(2)(1)} \left(\frac{\lambda^x}{n^x}\right) \left(1 - \frac{\lambda}{n}\right)^n \left(1 - \frac{\lambda}{n}\right)^{-x} \\
&= \overbrace{\left(\frac{n}{n}\right) \left(\frac{n-1}{n}\right) \left(\frac{n-2}{n}\right) \cdots \left(\frac{n-x+1}{n}\right)} \overbrace{\left(\frac{\lambda^x}{x!}\right)} \underbrace{\left(1 - \frac{\lambda}{n}\right)^n}_{\text{underbrace}} \underbrace{\left(1 - \frac{\lambda}{n}\right)^{-x}}_{\text{underbrace}}
\end{aligned} \tag{7.1}$$

As $n \rightarrow \infty$, the expression below the first overbrace $\rightarrow 1$, while that below 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:

$$\lim_{n \rightarrow \infty} P(X = x) = \frac{e^{-\lambda} \lambda^x}{x!}.$$

The second notion of convergence of RVs is convergence in probability.

Definition 82 (Convergence in Probability) 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 probability, and write:

$$X_n \xrightarrow{P} X$$

if for every real number $\epsilon > 0$,

$$\lim_{n \rightarrow \infty} P(|X_n - X| > \epsilon) = 0 \quad [\text{in the sense of Definition 4}].$$

Once again, the above limit, by (3.1) in our Definition 18 of a RV, can be equivalently expressed as follows:

$$\lim_{n \rightarrow \infty} P(\{\omega : |X_n(\omega) - X(\omega)| > \epsilon\}) = 0, \quad \text{ie,} \quad P(\{\omega : |X_n(\omega) - X(\omega)| > \epsilon\}) \rightarrow 0, \quad \text{as } n \rightarrow \infty.$$

For the same sequence of RVs in Classwork 183 and Example 184 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 need some elementary inequalities in Probability to help us answer this question. We visit these inequalities next.

Proposition 83 (Markov's Inequality) Let (Ω, \mathcal{F}, P) be a probability triple and let $X = X(\omega)$ be a non-negative RV. Then,

$$P(X \geq \epsilon) \leq \frac{E(X)}{\epsilon}, \quad \text{for any } \epsilon > 0. \tag{7.2}$$

Proof:

$$\begin{aligned}
X &= X \mathbf{1}_{\{y:y \geq \epsilon\}}(x) + X \mathbf{1}_{\{y:y < \epsilon\}}(x) \\
&\geq X \mathbf{1}_{\{y:y \geq \epsilon\}}(x) \\
&\geq \epsilon \mathbf{1}_{\{y:y \geq \epsilon\}}(x)
\end{aligned} \tag{7.3}$$

Finally, taking expectations on both sides of the above inequality and then using the fact that the expectation of an indicator function of an event is simply the probability of that event (3.46), we get the desired result:

$$E(X) \geq \epsilon E(\mathbf{1}_{\{y:y \geq \epsilon\}}(x)) = \epsilon P(X \geq \epsilon).$$

Let us look at some immediate consequences of Markov's inequality.

Proposition 84 (Chebychev's Inequality) For any RV X and any $\epsilon > 0$,

$$P(|X| > \epsilon) \leq \frac{E(|X|)}{\epsilon} \quad (7.4)$$

$$P(|X| > \epsilon) = P(X^2 \geq \epsilon^2) \leq \frac{E(X^2)}{\epsilon^2} \quad (7.5)$$

$$P(|X - E(X)| \geq \epsilon) = P((X - E(X))^2 \geq \epsilon^2) \leq \frac{E(X - E(X))^2}{\epsilon^2} = \frac{V(X)}{\epsilon^2} \quad (7.6)$$

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 183 and Example 184.

Example 186 (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 $X_n \xrightarrow{P} X$?

To find out if $X_n \xrightarrow{P} X$, we need to show that for any $\epsilon > 0$, $\lim_{n \rightarrow \infty} P(|X_n - X| > \epsilon) = 0$.

Let ϵ be any real number greater than 0, then

$$\begin{aligned} P(|X_n| > \epsilon) &= P(|X_n|^2 > \epsilon^2) \\ &\leq \frac{E(X_n^2)}{\epsilon^2} \quad [\text{by Markov's Inequality (7.2)}] \\ &= \frac{\frac{1}{n}}{\epsilon^2} \rightarrow 0, \quad \text{as } n \rightarrow \infty \quad [\text{in the sense of Definition 4}]. \end{aligned}$$

Hence, we have shown that for any $\epsilon > 0$, $\lim_{n \rightarrow \infty} P(|X_n - X| > \epsilon) = 0$ and therefore by Definition 82, $X_n \xrightarrow{P} X$ or $X_n \xrightarrow{P} 0$.

Convention: When X has a Point Mass(θ) distribution and $X_n \xrightarrow{P} X$, we simply write $X_n \xrightarrow{P} \theta$.

Definition 85 (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 (Ω, \mathcal{F}, P) , as denoted by

$$X_n \xrightarrow{a.s.} X$$

means that

$$P\left(\left\{\lim_{n \rightarrow \infty} X_n = X\right\}\right) = 1 \iff P\left(\left\{\omega \in \Omega : \lim_{n \rightarrow \infty} X_n(\omega) = X(\omega)\right\}\right) = 1.$$

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.

7.1.1 Properties of Convergence of RVs**

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 implication relations exist between the three notions.

- Convergence almost surely implies convergence in probability²

$$X_n \xrightarrow{a.s.} X \implies X_n \xrightarrow{P} X .$$

- By the Borel-Cantelli Lemma³, convergence in probability does not imply almost sure convergence in the discrete case⁴
- Convergence in probability implies convergence in distribution⁵

$$X_n \xrightarrow{P} X \implies X_n \rightsquigarrow X .$$

- Convergence in distribution to a constant θ implies convergence in probability to θ :⁶

$$X_n \rightsquigarrow \text{Point Mass}(\theta) \implies X_n \xrightarrow{P} \text{Point Mass}(\theta) .$$

- In general, convergence in distribution does not imply convergence in probability.

7.2 Law of Large Numbers

Proposition 86 (Law of Large Numbers (LLN): $\bar{X}_n \xrightarrow{P} E(X_1)$) If we are given a sequence of independent and identically distributed RVs, $X_1, X_2, \dots \stackrel{\text{IID}}{\sim} X_1$ and if $E(X_1)$ exists, as per (3.45), i.e., $E(\text{abs}(X_1)) < \infty$, and the variance is finite, i.e., $V(X_1) < \infty$, then the sample mean \bar{X}_n converges in probability to the expectation of any one of the IID RVs, say $E(X_1)$ by convention. More formally, we write:

$$\text{If } X_1, X_2, \dots \stackrel{\text{IID}}{\sim} X_1 \text{ and if } E(X_1) \text{ exists, then } \bar{X}_n \xrightarrow{P} E(X_1) .$$

Proof: Because $V(X_1) < \infty$, we have:

$$\begin{aligned} P(|\bar{X}_n - E(\bar{X}_n)| \geq \epsilon) &= \frac{V(\bar{X}_n)}{\epsilon^2} && [\text{by applying Chebychev's inequality (7.6) to the RV } \bar{X}_n] \\ &= \frac{\frac{1}{n} V(X_1)}{\epsilon^2} && [\text{by the IID assumption of } X_1, X_2, \dots \text{ we can apply (4.3)}] \end{aligned}$$

²https://en.wikipedia.org/wiki/Proofs_of_convergence_of_random_variables#Convergence_almost_surely_implies_convergence_in_probability

³https://en.wikipedia.org/wiki/Borel%20%93Cantelli_lemma

⁴https://en.wikipedia.org/wiki/Proofs_of_convergence_of_random_variables#Convergence_in_probability_does_not_imply_almost_sure_convergence_in_the_discrete_case

⁵https://en.wikipedia.org/wiki/Proofs_of_convergence_of_random_variables#Convergence_in_probability_implies_convergence_in_distribution

⁶https://en.wikipedia.org/wiki/Proofs_of_convergence_of_random_variables#Convergence_in_distribution_to_a_constant_implies_convergence_in_probability

Therefore, for any given $\epsilon > 0$,

$$\begin{aligned} \mathbb{P}(|\bar{X}_n - \mathbb{E}(X_1)| \geq \epsilon) &= \mathbb{P}(|\bar{X}_n - \mathbb{E}(\bar{X}_n)| \geq \epsilon) \quad [\text{by the IID assumption of } X_1, X_2, \dots, \mathbb{E}(\bar{X}_n) = \mathbb{E}(X_1), \text{ as per (4.2)}] \\ &= \frac{\frac{1}{n} V(X_1)}{\epsilon^2} \rightarrow 0, \quad \text{as } n \rightarrow \infty, \end{aligned}$$

or equivalently, $\lim_{n \rightarrow \infty} \mathbb{P}(|\bar{X}_n - \mathbb{E}(X_1)| \geq \epsilon) = 0$. And the last statement is the definition of the claim made by the law of large numbers (LLN), namely that $\bar{X}_n \xrightarrow{P} \mathbb{E}(X_1)$.

Proposition 87 (Weak Law of Large Numbers (WLLN): $\bar{X}_n \rightsquigarrow \text{Point Mass}(\mathbb{E}(X_1))$) 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 $\mathbb{E}(X_1)$ exists, i.e. $\mathbb{E}(\text{abs}(X)) < \infty$, then the sample mean \bar{X}_n converges in distribution to the expectation of any one of the IID RVs, say $\text{Point Mass}(\mathbb{E}(X_1))$ by convention. More formally, we write:

If $X_1, X_2, \dots \stackrel{\text{IID}}{\sim} X_1$ and if $\mathbb{E}(X_1)$ exists, then $\bar{X}_n \rightsquigarrow \text{Point Mass}(\mathbb{E}(X_1))$ 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 $\text{Point Mass}(\mathbb{E}(X_1))$ is

$$\mathbb{E}(e^{it\mathbb{E}(X_1)}) = e^{it\mathbb{E}(X_1)},$$

since $\mathbb{E}(X_1)$ is just a constant, i.e., a Point Mass RV that puts all of its probability mass at $\mathbb{E}(X_1)$.

Second, the CF of \bar{X}_n is

$$\begin{aligned} \mathbb{E}(e^{it\bar{X}_n}) &= \mathbb{E}\left(e^{it\frac{1}{n}\sum_{k=1}^n X_k}\right) = \mathbb{E}\left(\prod_{k=1}^n e^{itX_k/n}\right) = \prod_{k=1}^n \mathbb{E}\left(e^{itX_k/n}\right) = \prod_{k=1}^n \varphi_{X_k}(t/n) \\ &= \prod_{k=1}^n \varphi_{X_1}(t/n) = (\varphi_{X_1}(t/n))^n. \end{aligned}$$

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:

$$\varphi_X(t) = 1 + it\mathbb{E}(X) + o(t).$$

Hence,

$$\varphi_{X_1}(t/n) = 1 + it\frac{1}{n}\mathbb{E}(X_1) + o\left(\frac{t}{n}\right)$$

and

$$E\left(e^{it\bar{X}_n}\right) = \left(1 + it\frac{1}{n}\mathbb{E}(X_1) + o\left(\frac{t}{n}\right)\right)^n \rightarrow e^{it\mathbb{E}(X_1)} \text{ as } n \rightarrow \infty.$$

⁷https://en.wikipedia.org/wiki/L%C3%A9vy%27s_continuity_theorem

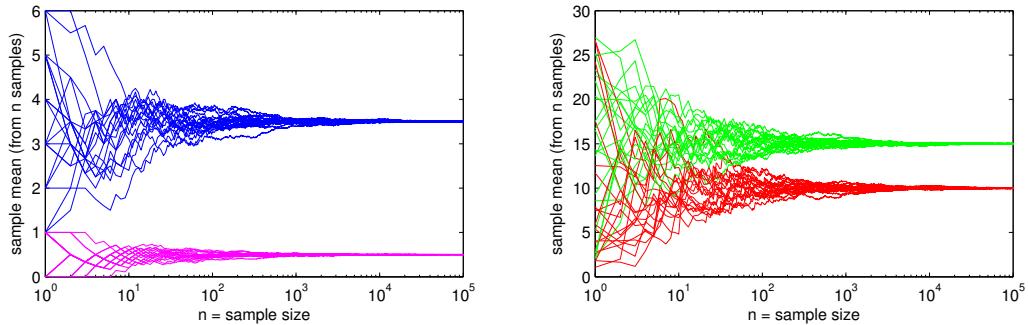
For the last limit we have used $\left(1 + \frac{x}{n}\right)^n \rightarrow e^x$ as $n \rightarrow \infty$.

Finally, we have shown that $E\left(e^{it\bar{X}_n}\right)$, the CF of the n -sample mean RV \bar{X}_n , converges to $E(e^{itE(X_1)}) = e^{itE(X_1)}$, the CF of the Point Mass($E(X_1)$) RV, as the sample size n tends to infinity.

Heuristic Interpretation of LLN

The distribution of the sample mean RV \bar{X}_n obtained from an independent and identically distributed sequence of RVs X_1, X_2, \dots [i.e. all the RVs X_i 's are independent of one another and have the same distribution function, and thereby the same expectation, variance and higher moments], concentrates around the expectation of any one of the RVs in the sequence, say that of the first one $E(X_1)$ [without loss of generality], as n approaches infinity. See Figure 7.4 for examples of 20 replicates of the sample mean of IID sequences from four RVs. All the sample mean trajectories converge to the corresponding population mean.

Figure 7.4: 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.



Example 187 (Bernoulli WLLN and Galton's Quincunx) We can appreciate the WLLN for $\bar{X}_n = n^{-1}S_n = \sum_{i=1}^n X_i$, where $X_1, X_2, \dots, X_n \stackrel{\text{IID}}{\sim} \text{Bernoulli}(p)$ using the paths of balls dropped into Galton's Quincunx of Sec. 3.2.3.

Cauchy whose expectations does not exist has no Law of Large Numbers

Recall that the mean of the Cauchy RV X does not exist since $\int |x| dF(x) = \infty$ (3.55). We will investigate this in Labwork 188.

Labwork 188 (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:

PlotStandardCauchyRunningMean.m

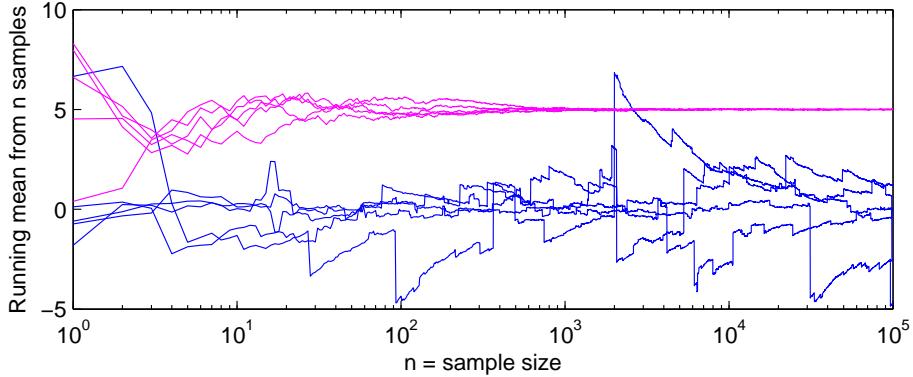
```
% script to plot the oscillating running mean of Std Cauchy samples
% relative to those for the Uniform(0,10) samples
rand('twister',25567); % initialize the fundamental sampler
for i=1:5
    N = 10^5; % maximum sample size
    u=rand(1,N); % draw N IID samples from Uniform(0,1)
```

```

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 7.5: 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 7.5. 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.

7.2.1 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 $X_1, X_2, \dots, X_n \stackrel{IID}{\sim} X_1$.

Example 189 Let $X_1, X_2, \dots, X_n \stackrel{IID}{\sim} X_1$, where X_1 is an $\text{Exponential}(\lambda^*)$ RV, i.e., let

$$X_1, X_2, \dots, X_n \stackrel{IID}{\sim} \text{Exponential}(\lambda^*) .$$

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

$$\bar{X}_n \rightsquigarrow \text{Point Mass}(E(X_1)) ,$$

and therefore, we can use the sample mean \bar{X}_n as a point estimator of $E(X_1) = 1/\lambda^*$.

Now, suppose you model seven waiting times in nearest minutes between Orbiter buses at Balgay street as follows:

$$X_1, X_2, \dots, X_7 \stackrel{IID}{\sim} \text{Exponential}(\lambda^*) ,$$

and have the following realization as your observed data:

$$(x_1, x_2, \dots, x_7) = (2, 12, 8, 9, 14, 15, 11) .$$

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 $E(X_1) = 1/\lambda^*$. By the rearrangement $\lambda^* = 1/E(X_1)$, we can also obtain a point estimate of the “true” parameter λ^* from $1/\bar{x}_7 = 7/71 \approx 0.0986$.

Remark 88 (Point estimates are realizations of the Point Estimator) We say the statistic \bar{X}_n , which is a random variable that depends on the data $\vec{RV}(X_1, X_2, \dots, X_n)$, is a **point estimator** of $E(X_1)$. But once we have a realization of the data \vec{RV} , i.e., our observed data vector (x_1, x_2, \dots, x_n) and its corresponding realization as observed sample mean \bar{x}_n , we say \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 the random variable \bar{X}_n called the point estimator of $E(X_1)$. Therefore, when we observe a new data vector $(x'_1, x'_2, \dots, x'_n)$ that is different from our first data vector (x_1, x_2, \dots, x_n) , our point estimator of $E(X_1)$ is still \bar{X}_n but the point estimate $n^{-1} \sum_{i=1}^n x'_i$ may be different from the first point estimate $n^{-1} \sum_{i=1}^n x_i$. The sample means from n samples for 20 replications (repeats of the experiment) are typically distinct especially for small n as shown in Figure 7.4.

Example 190 Let $X_1, X_2, \dots, X_n \stackrel{IID}{\sim} X_1$, where X_1 is an $\text{Bernoulli}(\theta^*)$ RV, i.e., let

$$X_1, X_2, \dots, X_n \stackrel{IID}{\sim} \text{Bernoulli}(\theta^*) .$$

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 LLN, we know that

$$\bar{X}_n \rightsquigarrow \text{Point Mass}(E(X_1)) ,$$

and therefore, we can use the sample mean \bar{X}_n as a point estimator of $E(X_1) = \theta^*$.

Now, suppose you model seven coin tosses (encoding **Heads** as 1 with probability θ^* and **Tails** as 0 with probability $1 - \theta^*$) as follows:

$$X_1, X_2, \dots, X_7 \stackrel{IID}{\sim} \text{Bernoulli}(\theta^*) ,$$

and have the following realization as your observed data:

$$(x_1, x_2, \dots, x_7) = (0, 1, 1, 0, 0, 1, 0) .$$

Then you can use the observed sample mean $\bar{x}_7 = (0 + 1 + 1 + 0 + 0 + 1 + 0)/7 = 3/7 \approx 0.4286$ as a **point estimate** of the population mean $E(X_1) = \theta^*$. Thus, our “single best guess” for $E(X_1)$ which is the same as the probability of **Heads** is $\bar{x}_7 = 3/7$.

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 orbiter buses at a different bus-stop or on a different day we may get a different point estimate for $E(X_1)$. See the intersection of the twenty magenta sample mean 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 7.4 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 $\text{error} = |\bar{X}_n - E(X_1)|$ in our estimate of $E(X_1)$ is within a required tolerance $= \epsilon$ we need to know the full distribution of $\bar{X}_n - E(X_1)$ itself. The Central Limit Theorem (CLT) helps us here.

7.3 Central Limit Theorem

What if we scale the sum of X_i ‘s by \sqrt{n} instead of n ?

Exercise 7.1 (What if we scale by \sqrt{n}) After reading Sec. 7.1 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 89 (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 $E(X_1)$ 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:

$$\begin{aligned} \text{If } X_1, X_2, \dots \stackrel{\text{IID}}{\sim} X_1 \text{ and if } E(X_1) < \infty, V(X_1) < \infty \\ \text{then } \bar{X}_n \rightsquigarrow \text{Normal}\left(E(X_1), \frac{V(X_1)}{n}\right) \text{ as } n \rightarrow \infty , \end{aligned} \quad (7.7)$$

or equivalently after standardization:

$$\begin{aligned} \text{If } X_1, X_2, \dots \stackrel{\text{IID}}{\sim} X_1 \text{ and if } E(X_1) < \infty, V(X_1) < \infty \\ \text{then } \frac{\bar{X}_n - E(X_1)}{\sqrt{V(X_1)/n}} \rightsquigarrow Z \sim \text{Normal}(0, 1) \text{ as } n \rightarrow \infty . \end{aligned} \quad (7.8)$$

Proof: Our proof is based on the convergence of characteristic functions (CFs). We will prove the standardized form of the CLT in Equation (7.8) 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 $\text{Normal}(0, 1)$ RV. First, note from Equation (3.72) that the CF of $Z \sim \text{Normal}(0, 1)$ is:

$$\varphi_Z(t) = E(e^{itZ}) = e^{-t^2/2} .$$

Second,

$$U_n := \frac{\bar{X}_n - E(X_1)}{\sqrt{V(X_1)/n}} = \frac{\sum_{k=1}^n X_k - nE(X_1)}{\sqrt{nV(X_1)}} = \frac{1}{\sqrt{n}} \sum_{k=1}^n \left(\frac{X_k - E(X_1)}{\sqrt{V(X_1)}} \right) .$$

Therefore, the CF of U_n is

$$\begin{aligned} \varphi_{U_n}(t) &= E(\exp(itU_n)) = E\left(\exp\left(i\frac{t}{\sqrt{n}} \sum_{k=1}^n \frac{X_k - E(X_1)}{\sqrt{V(X_1)}}\right)\right) = \prod_{k=1}^n E\left(\exp\left(i\frac{t}{\sqrt{n}} \frac{X_k - E(X_1)}{\sqrt{V(X_1)}}\right)\right) \\ &= \left(E\left(\exp\left(i\frac{t}{\sqrt{n}} \frac{X_1 - E(X_1)}{\sqrt{V(X_1)}}\right)\right)\right)^n . \end{aligned}$$

Now, if we let

$$Y = \frac{X_1 - E(X_1)}{\sqrt{V(X_1)}}$$

then

$$E(Y) = 0, \quad E(Y^2) = 1, \text{ and } V(Y) = 1 .$$

So, the CF of U_n is

$$\varphi_{U_n}(t) = \left(\varphi_Y\left(\frac{t}{\sqrt{n}}\right)\right)^n ,$$

and since we can Taylor expand $\varphi_Y(t)$ as follows:

$$\varphi_Y(t) = 1 + itE(Y) + i^2 \frac{t^2}{2} E(Y^2) + o(t^2) ,$$

which implies

$$\varphi_Y\left(\frac{t}{\sqrt{n}}\right) = 1 + \frac{it}{\sqrt{n}} E(Y) + \frac{i^2 t^2}{2n} E(Y^2) + o\left(\frac{t^2}{n}\right) ,$$

we finally get

$$\varphi_{U_n}(t) = \left(\varphi_Y\left(\frac{t}{\sqrt{n}}\right)\right)^n = \left(1 + \frac{it}{\sqrt{n}} \times 0 + \frac{i^2 t^2}{2n} \times 1 + o\left(\frac{t^2}{n}\right)\right)^n = \left(1 - \frac{t^2}{2n} + o\left(\frac{t^2}{n}\right)\right)^n \rightarrow e^{-t^2/2} = \varphi_Z(t) .$$

For the last limit we have used $(1 + \frac{x}{n})^n \rightarrow e^x$ as $n \rightarrow \infty$. Thus, we have proved Equation (7.8) which is equivalent to Equation (7.7) by a standardization argument that if $W \sim \text{Normal}(\mu, \sigma^2)$ then $Z = \frac{W-\mu}{\sigma} \sim \text{Normal}(0, 1)$ through the linear transformation $W = \sigma Z + \mu$ of Example 67.

7.3.1 Application: Tolerating Errors in our estimate of $E(X_1)$

Recall that we wanted to ensure the **error** = $|\bar{X}_n - E(X_1)|$ in our estimate of $E(X_1)$ is within a required **tolerance** = ϵ and make the following probability statement:

$$P(\text{error} < \text{tolerance}) = P(|\bar{X}_n - E(X_1)| < \epsilon) = P(-\epsilon < \bar{X}_n - E(X_1) < \epsilon) = 1 - \alpha .$$

To be able to do this we needed to know the full distribution of $\bar{X}_n - E(X_1)$ itself.

Due to the Central Limit Theorem (CLT) we now know that (assuming n is large)

$$\begin{aligned} P(-\epsilon < \bar{X}_n - E(X_1) < \epsilon) &\approx P\left(-\frac{\epsilon}{\sqrt{V(X_1)/n}} < \frac{\bar{X}_n - E(X_1)}{\sqrt{V(X_1)/n}} < \frac{\epsilon}{\sqrt{V(X_1)/n}}\right) \\ &= P\left(-\frac{\epsilon}{\sqrt{V(X_1)/n}} < Z < \frac{\epsilon}{\sqrt{V(X_1)/n}}\right), \end{aligned}$$

where $Z \sim \text{Normal}(0, 1)$.

Example 191 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,

$$P(\text{error} < 0.1) \approx P\left(-\frac{0.1}{\sqrt{4/80}} < Z < \frac{0.1}{\sqrt{4/80}}\right) = P(-0.447 < Z < 0.447) = 0.345.$$

Suppose you want the error to be less than tolerance $= \epsilon$ 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

$$n = \left(\frac{\sqrt{V(X_1)} z_{\alpha/2}}{\epsilon} \right)^2 \tag{7.9}$$

for the needed sample size that will ensure that our error is less than our tolerance $= \epsilon$ with probability $1 - \alpha$. Of course, if n given by Equation (7.9) 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 192 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 (7.9) we see that the needed sample size is

$$n = \left(\frac{\sqrt{4} \times 1.96}{0.1} \right)^2 \approx 1537$$

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.

7.3.2 Application: Set Estimation of $E(X_1)$

A useful byproduct of the CLT is the $(1 - \alpha)$ **confidence interval**, a random interval (or bivariate RV) that contains $E(X_1)$, the quantity of interest, with probability $1 - \alpha$:

$$(\bar{X}_n \pm z_{\alpha/2} \sqrt{V(X_1)/n}) := (\bar{X}_n - z_{\alpha/2} \sqrt{V(X_1)/n}, \bar{X}_n + z_{\alpha/2} \sqrt{V(X_1)/n}) . \quad (7.10)$$

We can easily see how Equation (7.10) is derived from CLT as follows:

$$\begin{aligned} P(-z_{\alpha/2} < Z < z_{\alpha/2}) &= 1 - \alpha \\ P\left(-z_{\alpha/2} < \frac{\bar{X}_n - E(X_1)}{\sqrt{V(X_1)/n}} < z_{\alpha/2}\right) &= 1 - \alpha \\ P\left(-\bar{X}_n - z_{\alpha/2} \sqrt{V(X_1)/n} < -E(X_1) < -\bar{X}_n + z_{\alpha/2} \sqrt{V(X_1)/n}\right) &= 1 - \alpha \\ P\left(\bar{X}_n + z_{\alpha/2} \sqrt{V(X_1)/n} > E(X_1) > \bar{X}_n - z_{\alpha/2} \sqrt{V(X_1)/n}\right) &= 1 - \alpha \\ P\left(\bar{X}_n - z_{\alpha/2} \sqrt{V(X_1)/n} < E(X_1) < \bar{X}_n + z_{\alpha/2} \sqrt{V(X_1)/n}\right) &= 1 - \alpha \\ P(E(X_1) \in (\bar{X}_n - z_{\alpha/2} \sqrt{V(X_1)/n}, \bar{X}_n + z_{\alpha/2} \sqrt{V(X_1)/n})) &= 1 - \alpha . \end{aligned}$$

Remark 90 (Heuristic interpretation of the $(1 - \alpha)$ confidence interval) If we repeatedly produced samples of size n to contain $E(X_1)$ within a $(\bar{X}_n \pm z_{\alpha/2} \sqrt{V(X_1)/n})$, say 100 times, then on average, $(1 - \alpha) \times 100$ repetitions will actually contain $E(X_1)$ within the random interval and $\alpha \times 100$ repetitions will fail to contain $E(X_1)$.

So far, we have assumed we know the population variance $V(X_1)$ in an IID experiment with n samples and tried to estimate the population mean $E(X_1)$. But in general, we will not know $V(X_1)$. We can still get a point estimate of $E(X_1)$ from the sample mean due to LLN but we won't be able to get a confidence interval for $E(X_1)$. Fortunately, a more elaborate form of the CLT tells us that even when we substitute the sample variance $S_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2$ for the population variance $V(X_1)$ the following $1 - \alpha$ confidence interval for $E(X_1)$ works!

$$(\bar{X}_n \pm z_{\alpha/2} S_n / \sqrt{n}) := (\bar{X}_n - z_{\alpha/2} S_n / \sqrt{n}, \bar{X}_n + z_{\alpha/2} S_n / \sqrt{n}) , \quad (7.11)$$

where, $S_n = \sqrt{S_n^2}$ is the sample standard deviation.

Let's return to our two examples again.

Example 193 We model the waiting times between Orbiter buses with unknown $E(X_1) = 1/\lambda^*$ as

$$X_1, X_2, \dots, X_n \stackrel{IID}{\sim} \text{Exponential}(\lambda^*)$$

and observed the following data, sample mean, sample variance and sample standard deviation:

$$(x_1, x_2, \dots, x_7) = (2, 12, 8, 9, 14, 15, 11), \bar{x}_7 = 10.143, s_7^2 = 19.143, s_7 = 4.375 ,$$

respectively. Our point estimate and $1 - \alpha = 95\%$ confidence interval for $E(X_1)$ are:

$$\bar{x}_7 = 10.143 \quad \text{and} \quad (\bar{x}_7 \pm z_{\alpha/2} s_7 / \sqrt{7}) = (10.143 \pm 1.96 \times 4.375 / \sqrt{7}) = (6.9016, 13.3841) ,$$

respectively. So with 95% probability the true population mean $E(X_1) = 1/\lambda^*$ is contained in $(6.9016, 13.3841)$ and since the mean waiting time of 10 minutes promised by the Orbiter bus company is also within $(6.9016, 13.3841)$ we can be fairly certain that the company sticks to its promise.

Example 194 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 91 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.

7.4 Exercises in Limit Laws of Statistics

Ex. 7.2 — 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. 7.3 — 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

$$X_1, X_2, \dots, X_n \stackrel{\text{IID}}{\sim} \text{Poisson}(\lambda = 5).$$

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. 7.4 — What is the distribution of $\sum_{i=1}^n X_i/n$ as $n \rightarrow \infty$ when $X_i \stackrel{\text{IID}}{\sim} \text{Uniform}(-10, 10)$?

Ex. 7.5 — What is the distribution of $\sum_{i=1}^n X_i/\sqrt{V(X_i)n}$ as $n \rightarrow \infty$ when $X_i \stackrel{\text{IID}}{\sim} \text{Uniform}(-10, 10)$?

Chapter 8

Finite Markov Chains

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.

8.1 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 92 (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 $P(H_{t-1} \cap \{X_t = x\}) > 0$, the following **Markov property** is satisfied:

$$P(X_{t+1} = y | H_{t-1} \cap \{X_t = x\}) = P(X_{t+1} = y | X_t = x) =: P(x, y) . \quad (8.1)$$

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 $|\mathbb{X}| \times |\mathbb{X}|$ matrix P is enough to obtain the state transitions since the x -th row of P is the probability distribution $P(x, \cdot) := (P(x, y))_{y \in \mathbb{X}}$. For this reason P is called a **stochastic matrix**, i.e.,

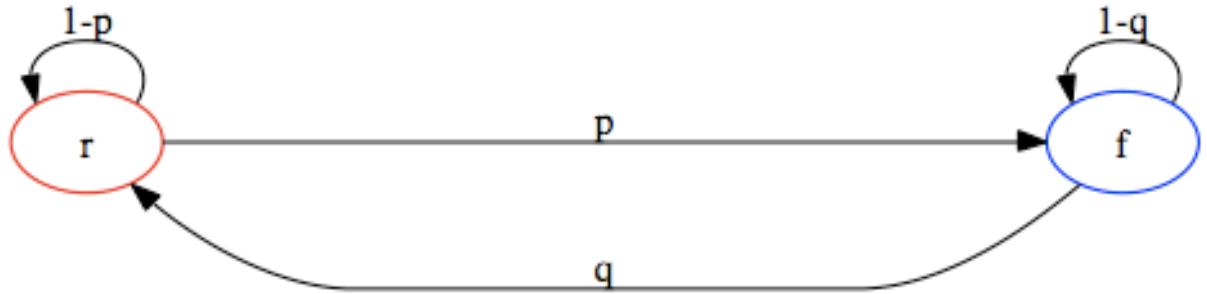
$$P(x, y) \geq 0 \quad \text{for all } (x, y) \in \mathbb{X}^2 \quad \text{and} \quad \sum_{y \in \mathbb{X}} P(x, y) = 1 \quad \text{for all } x \in \mathbb{X} . \quad (8.2)$$

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 195 (Flippant Freddy) Freddy the flippant frog lives in an enchanted pond with only two lily pads, *rollopia* and *flipopia*. A wizard gave a die and a silver coin to help flippant Freddy decide where to jump next. Freddy left the die on *rollopia* and the coin on *flipopia*. When Freddy got restless in *rollopia* he would roll the die and if the die landed odd he would leave the die behind and jump to *flipopia*, otherwise he would stay put. When Freddy got restless in *flipopia* he would flip the coin and if it landed Heads he would leave the coin behind and jump to *rollopia*, otherwise he would stay put.

Let the state space $\mathbb{X} = \{r, f\}$, and let (X_0, X_1, \dots) be the sequence of lily pads occupied by Freddy after his restless moments. Say the die on *rollopia* r has probability p of turning up odd and the coin on *flipopia* f has probability q of turning up heads. We can visualise the rules of Freddy's jumps by the following **transition diagram**:

Figure 8.1: Transition Diagram of Flippant Freddy's Jumps.



Then Freddy's sequence of jumps (X_0, X_1, \dots) is a Markov chain on \mathbb{X} with transition matrix:

$$P = \begin{pmatrix} r & f \\ r & f \end{pmatrix} = \begin{pmatrix} P(r, r) & P(r, f) \\ P(f, r) & P(f, f) \end{pmatrix} = \begin{pmatrix} r & f \\ f & r \end{pmatrix} \begin{pmatrix} 1-p & p \\ q & 1-q \end{pmatrix}. \quad (8.3)$$

Suppose we first see Freddy in *rollopia*, i.e., $X_0 = r$. When he gets restless for the first time we know from the first row of P that he will leave to *flipopia* with probability p and stay with probability $1 - p$, i.e.,

$$P(X_1 = f | X_0 = r) = p, \quad P(X_1 = r | X_0 = r) = 1 - p. \quad (8.4)$$

What happens when he is restless for the second time? By considering the two possibilities for

X_1 , Definition of conditional probability and the Markov property, we see that,

$$\begin{aligned}
P(X_2 = f | X_0 = r) &= P(X_2 = f, X_1 = f | X_0 = r) + P(X_2 = f, X_1 = r | X_0 = r) \\
&= \frac{P(X_2 = f, X_1 = f, X_0 = r)}{P(X_0 = r)} + \frac{P(X_2 = f, X_1 = r, X_0 = r)}{P(X_0 = r)} \\
&= P(X_2 = f | X_1 = f, X_0 = r) \frac{P(X_1 = f, X_0 = r)}{P(X_0 = r)} \\
&\quad + P(X_2 = f | X_1 = r, X_0 = r) \frac{P(X_1 = r, X_0 = r)}{P(X_0 = r)} \\
&= P(X_2 = f | X_1 = f, X_0 = r) P(X_1 = f | X_0 = r) \\
&\quad + P(X_2 = f | X_1 = r, X_0 = r) P(X_1 = r | X_0 = r) \\
&= P(X_2 = f | X_1 = f) P(X_1 = f | X_0 = r) \\
&\quad + P(X_2 = f | X_1 = r) P(X_1 = r | X_0 = r) \\
&= P(f, f) P(r, f) + P(r, f) P(r, r) \\
&= (1 - q)p + p(1 - p)
\end{aligned} \tag{8.5}$$

Similarly,

$$P(X_2 = r | X_0 = r) = P(f, r) P(r, f) + P(r, r) P(r, r) = qp + (1 - p)(1 - p) \tag{8.6}$$

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:

$$\mu_t := (P(X_t = r | X_0 = r), P(X_t = f | X_0 = r)) ,$$

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 (8.4) from $\mu_1 = \mu_0 P$ and the 2-step state probabilities in (8.5) and (8.6) 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:

```

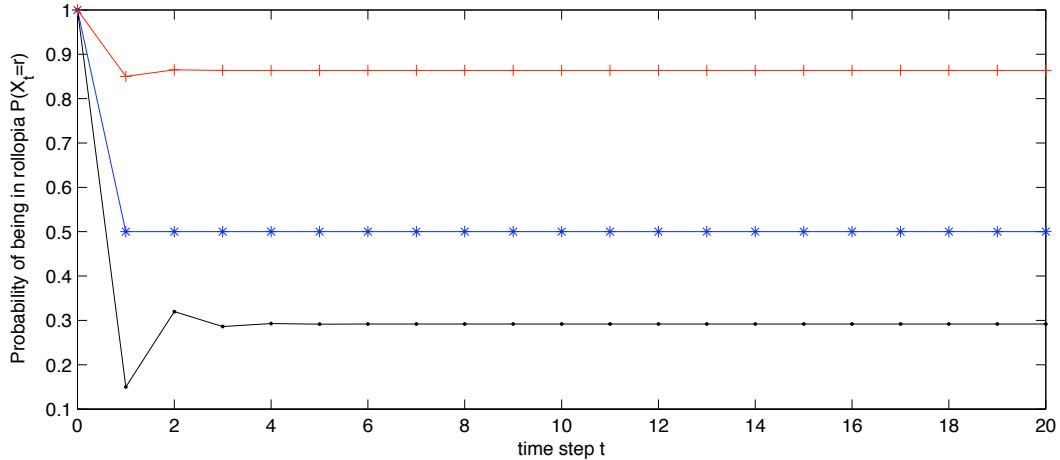
>> p=0.85; q=0.35; P = [1-p p; q 1-q] % assume an unfair coin and an unfair die
P =
    0.1500    0.8500
    0.3500    0.6500
>> mu0 = [1, 0] % initial state vector since Freddy started in rollovia
mu0 =
    1    0
>> mu0*P^0    % initial state distribution at t=0 is just mu0
ans =
    1    0
>> mu0*P^1    % 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

```

Now, let us compute and look at the probability of being in rollovia after having started there for three values of p and q according to the following script:

```
FlippantFreddyRolloviaProbs.m
p=0.5; q=0.5; P = [1-p p; q 1-q]; % assume a fair coin and a fair die
mu0 = [1, 0]; % initial state vector since Freddy started in rollovia
for t = 1: 1: 21, mu(t,:)= mu0*P^(t-1); end
t=0:1:20; % vector of time steps t
plot(t,mu(:,1)', 'b*-')
hold on;
p=0.85; q=0.35; P = [1-p p; q 1-q]; % assume an unfair coin and an unfair die
for t = 1: 1: 21, mu(t,:)= mu0*P^(t-1); end
t=0:1:20; % vector of time steps t
plot(t,mu(:,1)', 'k.-')
p=0.15; q=0.95; P = [1-p p; q 1-q]; % assume another unfair coin and another unfair die
for t = 1: 1: 21, mu(t,:)= mu0*P^(t-1); end
t=0:1:20; % vector of time steps t
plot(t,mu(:,1)', 'r+-')
xlabel('time step t'); ylabel('Probability of being in rollovia $P(X_{t=r})$')
xlabel('time step t'); ylabel('Probability of being in rollovia $P(X_{t=r})$')
```

Figure 8.2: The probability of being back in rollovia 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).



It is evident from Figure 8.2 that as $t \rightarrow \infty$, μ_t approaches a distribution, say π , that depends on p and q in P . Such a limit distribution is called the **stationary distribution** and must satisfy the fixed point condition:

$$\pi P = \pi ,$$

that gives the solution:

$$\pi(r) = \frac{q}{p+q}, \quad \pi(f) = \frac{p}{p+q} .$$

In Figure 8.2 we see that $P(X_t = r)$ approaches $\pi(r) = \frac{q}{p+q}$ for the three cases of p and q :

$$\begin{aligned} \text{(i)} \quad p = 0.50, q = 0.50, \quad P(X_t = r) \rightarrow \pi(r) = \frac{q}{p+q} = \frac{0.50}{0.50+0.50} = 0.5000, \\ \text{(ii)} \quad p = 0.85, q = 0.35, \quad P(X_t = r) \rightarrow \pi(r) = \frac{q}{p+q} = \frac{0.35}{0.85+0.35} = 0.2917, \\ \text{(iii)} \quad p = 0.15, q = 0.95, \quad P(X_t = r) \rightarrow \pi(r) = \frac{q}{p+q} = \frac{0.95}{0.15+0.95} = 0.8636. \end{aligned}$$

Now let us generalise the lessons learned from Example 195.

Proposition 93 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 $\mu_0(s_i) = P(X_0 = s_i)$, 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 $\mu_t(s_i) = P(X_t = s_i)$, satisfies:

$$\mu_t = \mu_0 P^t. \quad (8.7)$$

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

$$\begin{aligned} \mu_1(s_j) &= P(X_1 = s_j) = \sum_{i=1}^k P(X_1 = s_j, X_0 = s_i) \\ &= \sum_{i=1}^k P(X_1 = s_j | X_0 = s_i) P(X_0 = s_i) \\ &= \sum_{i=1}^k P(s_i, s_j) \mu_0(s_i) \\ &= (\mu_0 P)(s_j), \quad \text{the } j\text{-th entry of the row vector } (\mu_0 P). \end{aligned}$$

Hence, $\mu_1 = \mu_0 P$. Now, we will fix m and suppose that (8.7) holds for $t = m$ and prove that (8.7) also holds for $t = m + 1$. For each $j \in \{1, 2, \dots, k\}$, we get

$$\begin{aligned} \mu_{m+1}(s_j) &= P(X_{m+1} = s_j) = \sum_{i=1}^k P(X_{m+1} = s_j, X_m = s_i) \\ &= \sum_{i=1}^k P(X_{m+1} = s_j | X_m = s_i) P(X_m = s_i) \\ &= \sum_{i=1}^k P(s_i, s_j) \mu_m(s_i) \\ &= (\mu_m P)(s_j), \quad \text{the } j\text{-th entry of the row vector } (\mu_m P). \end{aligned}$$

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.

Thus, multiplying a row vector μ_0 by P^t on the right takes you from current distribution over the state space to the distribution in t steps of the chain.

Since we will be interested in Markov chains on $\mathbb{X} = \{s_1, s_2, \dots, s_k\}$ with the same transition matrix P but different initial distributions, we introduce P_μ and E_μ for probabilities and expectations given that the initial distribution is μ , respectively. When the initial distribution is concentrated at a single initial state x given by:

$$\mathbf{1}_{\{x\}}(y) := \begin{cases} 1 & \text{if } y = x \\ 0 & \text{if } y \neq x \end{cases}$$

we represent it by e_x , the $1 \times k$ ortho-normal basis row vector with a 1 in the x -th entry and a 0 elsewhere. We simply write P_x for $P_{\mathbf{1}_{\{x\}}}$ or P_{e_x} and E_x for $E_{\mathbf{1}_{\{x\}}}$ or E_{e_x} . Thus, Proposition 93 along with our new notations means that:

$$P_x(X_t = y) = (e_x P^t)(y) = P^t(x, y) .$$

In words, the probability of going to y from x in t steps is given by the (x, y) -th entry of P^t , the **t -step transition matrix**. We refer to the x -th row and the x -th column of P by $P(x, \cdot)$ and $P(\cdot, x)$, respectively.

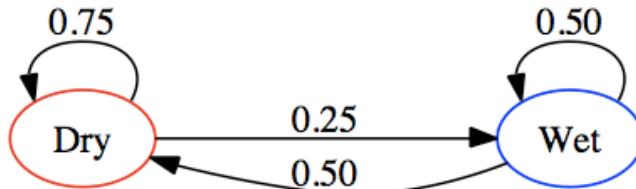
Let the function $f(x) : \mathbb{X} \rightarrow \mathbb{R}$ be represented by the column vector $f := (f(s_1), f(s_2), \dots, f(s_k)) \in \mathbb{R}^{k \times 1}$. Then the x -th entry of $P^t f$ is:

$$P^t f(x) = \sum_y P^t(x, y) f(y) = \sum_y f(y) P_x(X_t = y) = E_x(f(X_t)) .$$

This is the expected value of f under the distribution of states in t steps given that we start at state x . Thus multiplying a column vector f by P^t from the left takes you from a function on the state space to its expected value in t steps of the chain.

Example 196 (Dry-Wet Christchurch Weather) Consider a toy weather model for dry or wet days in Christchurch using a Markov chain with state space $\{d, w\}$. Let the transition diagram in Figure 8.3 give the transition matrix P for our dry-wet Markov chain. Using (8.7)

Figure 8.3: Transition Diagram of Dry and Wet Days in Christchurch.



we can find that the probability of being dry on the day after tomorrow is 0.625 given that it is wet today as follows:

```

>> P=[0.75 0.25; 0.5 0.5] % Transition Probability Matrix
P =
    0.7500    0.2500
    0.5000    0.5000
>> mu0=[0 1] % it is wet today gives the initial distribution
mu0 =
     0     1
>> mu0 * P^2 % the distribution in 2 days from today
ans =    0.6250    0.3750

```

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 197 (Freddy discovers a gold coin) Flippant Freddy of Example 195 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, rollophia 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 rollophia he will jump to flipopia but if he rolls an even number then he will stay in rollophia 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 rollophia.

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 198 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

$$P = \begin{matrix} & a & b & c \\ a & \left(\begin{matrix} 0 & 1 & 0 \end{matrix} \right) \\ b & \left(\begin{matrix} 0 & 0 & 1 \end{matrix} \right) \\ c & \left(\begin{matrix} 1 & 0 & 0 \end{matrix} \right) \end{matrix} .$$

For each t , define $Y_t = \mathbf{1}_{\{b,c\}}(X_t)$. Show that $(Y_t)_{t \in \mathbb{Z}_+}$ is not a Markov chain.

Exercise 199 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 92, 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.

Definition 94 (Inhomogeneous finite Markov chain) Let P_1, P_2, \dots be a sequence of $k \times k$ stochastic matrices satisfying the conditions in Equation 8.2. 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 , if for all pairs of states $(x, y) \in \mathbb{X} \times \mathbb{X}$, all integers $t \geq 1$, and all probable historical events $H_{t-1} := \bigcap_{n=0}^{t-1} \{X_n = x_n\}$ with $P(H_{t-1} \cap \{X_t = x\}) > 0$, the following **Markov property** is satisfied:

$$P(X_{t+1} = y | H_{t-1} \cap \{X_t = x\}) = P(X_{t+1} = y | X_t = x) =: P_{t+1}(x, y). \quad (8.8)$$

Proposition 95 For a finite inhomogeneous 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 $\mu_0(s_i) = P(X_0 = s_i)$, and transition matrices

$$(P_1, P_2, \dots), \quad P_t := (P_t(s_i, s_j))_{(s_i, s_j) \in \mathbb{X} \times \mathbb{X}}, \quad t \in \{1, 2, \dots\}$$

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 $\mu_t(s_i) = P(X_t = s_i)$, satisfies:

$$\mu_t = \mu_0 P_1 P_2 \cdots P_t. \quad (8.9)$$

Proof: Left as Exercise 200.

Exercise 200 Prove Proposition 95 using induction as done for Proposition 93.

Example 201 (a more sophisticated dry-wet chain) Let us make a more sophisticated version of the dry-wet chain of Example 196 with state space $\{d, w\}$. In order to take some seasonality into account in our weather model for dry and wet days in Christchurch, let us have two transition matrices for hot and cold days:

$$P_{\text{hot}} = \begin{matrix} d & w \\ \begin{pmatrix} 0.95 & 0.05 \\ 0.75 & 0.25 \end{pmatrix} \end{matrix}, \quad P_{\text{cold}} = \begin{matrix} d & w \\ \begin{pmatrix} 0.65 & 0.35 \\ 0.45 & 0.55 \end{pmatrix} \end{matrix}.$$

We say that a day is hot if its maximum temperature is more than 20° Celsius, otherwise it is cold. We use the transition matrix for today to obtain the state probabilities for tomorrow. If today is dry and hot and tomorrow is supposed to be cold then what is the probability that the day after tomorrow will be wet? We can use (8.9) to obtain the answer as 0.36:

```
>> Phot = [0.95 0.05; 0.75 0.25] % Transition Probability Matrix for hot day
Phot =
    0.9500    0.0500
    0.7500    0.2500
>> Pcold = [0.65 0.35; 0.45 0.55] % Transition Probability Matrix for cold day
Pcold =
    0.6500    0.3500
    0.4500    0.5500
>> mu0 = [1 0] % today is dry
mu0 =
    1    0
>> 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 202 For the Markov chain in Example 201 compute the probability that the day after tomorrow is wet if today is dry and hot but tomorrow is supposed to be cold.

8.2 Random Mapping Representation and Simulation

In order to simulate (x_0, x_1, \dots, x_n) , a sequential realisation or sequence of states visited by a Markov chain, say the sequence of lily pads that Flippant Freddy visits on his jumps, we need a random mapping representation of a Markov chain and its computer implementation.

Definition 96 (Random mapping representation (RMR)) A **random mapping representation (RMR)** of a transition matrix $P := (P(x, y))_{(x,y) \in \mathbb{X}^2}$ is a function

$$\rho(x, w) : \mathbb{X} \times \mathbb{W} \rightarrow \mathbb{X} , \quad (8.10)$$

along with the auxiliary \mathbb{W} -valued random variable W , satisfying

$$P(\{\rho(x, W) = y\}) = P(x, y), \quad \text{for each } (x, y) \in \mathbb{X}^2 . \quad (8.11)$$

Proposition 97 (Markov chain from RMR) If $W_1, W_2, \dots \stackrel{IID}{\sim} W$, the auxiliary RV in a RMR of a transition matrix $P := (P(x, y))_{(x,y) \in \mathbb{X}^2}$, and $X_0 \sim \mu_0$, then $(X_t)_{t \in \mathbb{Z}_+}$ defined by

$$X_t = \rho(X_{t-1}, W_t) , \quad \text{for all } t \geq 1$$

is a Markov chain with transition matrix P and initial distribution μ_0 on state space \mathbb{X} .

Proof: Left as Exercise 203.

Exercise 203 Do the proof of Proposition 97 by using the necessary Definitions.

Example 204 (An RMR for Flippant Freddy) Reconsider the Markov chain of Flippant Freddy with fair dice and fair coin on state space $\mathbb{X} = \{r, f\}$ with transition matrix

$$P = \begin{matrix} & \begin{matrix} r & f \end{matrix} \\ \begin{matrix} r \\ f \end{matrix} & \begin{pmatrix} 1/2 & 1/2 \\ 1/2 & 1/2 \end{pmatrix} \end{matrix}.$$

Let the auxiliary RV W have sample space $\mathbb{W} = \{0, 1\}$. Then an RMR $\rho : \mathbb{X} \times \mathbb{W} \rightarrow \mathbb{X}$ for this P is given by

$$\rho(x, w) : \{r, f\} \times \{0, 1\} \rightarrow \{r, f\}, \quad \rho(r, 0) = r, \quad \rho(r, 1) = f, \quad \rho(f, 0) = f, \quad \rho(f, 1) = r,$$

with $P(W = 0) = P(W = 1) = 1/2$. Now let us check that our ρ and W satisfy Equation 8.11:

$$\begin{aligned} \begin{matrix} r & f \\ r & f \end{matrix} \left(\begin{matrix} P(\{\rho(r, W) = r\}) & P(\{\rho(r, W) = f\}) \\ P(\{\rho(f, W) = r\}) & P(\{\rho(f, W) = f\}) \end{matrix} \right) &= \begin{matrix} r & f \\ r & f \end{matrix} \left(\begin{matrix} P(W = 0) & P(W = 1) \\ P(W = 1) & P(W = 0) \end{matrix} \right) \\ &= \begin{matrix} r & f \\ r & f \end{matrix} \left(\begin{matrix} 1/2 & 1/2 \\ 1/2 & 1/2 \end{matrix} \right) = P. \end{aligned}$$

Thus, by Proposition 97 we can obtain Freddy's Markov chain $(X_t)_{t \in \mathbb{Z}_+}$ by initialising $X_0 \sim \mu_0 = (1, 0)$, i.e., setting $X_0 = r$ since Freddy starts at r , and defining

$$X_t = \rho(X_{t-1}, W_t), \quad \text{for all } t \geq 1, \text{ where, } W_1, W_2, \dots \stackrel{IID}{\sim} \text{Bernoulli}(1/2) \text{ RV}.$$

In other words, we can simulate a sequence of states or lily pads visited by Freddy by merely doing independent Bernoulli(1/2) trials and use the mapping ρ . A MATLAB implementation of this RMR ρ as a MATLAB function is:

RMR10fFairFreddy.m

```

function y = RMR10fFairFreddy(x,w)
% Random Mapping Representation Number 1 of P=[1/2 1/2; 1/2 12/]
% input: character x as 'r' or 'f' and w as 0 or 1
% output: character y as 'r' or 'f'
if (x =='r')
    if (w==0)
        y = 'r';
    elseif (w==1)
        y = 'f';
    else
        y = Nan;
        print "when x = 'r' w is neither 0 nor 1!";
    end
elseif (x =='f')
    if (w==0)
        y = 'f';
    elseif (w==1)
        y = 'r';
    else
        y = Nan;
        print "when x='f' w is neither 0 nor 1!";
    end
else
end

```

```

y = Nan;
print "x is neither 'r' nor 'f'";
end

```

We can simulate one realisation of the first two states (x_0, x_1) visited by (X_0, X_1) as follows:

```

>> % set PRNG to be twister with seed 19731511
>> RandStream.setDefaultStream(RandStream('mt19937ar','seed',19731511));
>> x0 = 'r' % set x_0 = 'r'
x0 = r
>> w1 = floor( rand + 0.5 ) % a Bernoulli(0.5) trial
w1 =
0
>> x1 = RMR10fFairFreddy(x0,w1) % x_1 = rho(x_0,w1) is the state at time t=1
x1 = r

```

We can simulate one realisation of the first 10 states (x_0, x_1, \dots, x_9) visited by (X_0, X_1, \dots, X_9) using a for loop as follows:

```

>> % set PRNG to be twister with seed 19731511
>> RandStream.setDefaultStream(RandStream('mt19937ar','seed',19731511));
>> x0 = 'r' % set x_0 = 'r'
x0 = r
>> xt = x0; % current state x_t is x_0
>> Visited = x0; % initialise the variable Visited to hold the visited states
>> for t = 1:9 % start a for loop for t = 1,2,...,9
xt = RMR10fFairFreddy(xt, floor(rand+0.5) ); % update the current state at t
Visited = strcat(Visited,',',xt); % store the visited state in string Visited
end
>> Visited % disclose the string of visited state separated by commas
Visited = r,r,f,f,r,r,f,f,r,f,r

```

If we change the seed to some other number and repeat the code above, we will get another realisation of visits (x_0, x_1, \dots, x_9) of (X_0, X_1, \dots, X_9) . However, there are many distinct RMRs of the same transition matrix P . For example, we can define a new RMR ρ' from our first RMR ρ for P by $\rho'(x, w) = \rho(x, 1-w)$. The reader should check that ρ' also satisfies Equation 8.11 with $W \sim \text{Bernoulli}(1/2)$. But note that even for the same seed and the same PRNG the sequence of states (x_0, x_1, \dots, x_9) visited by (X_0, X_1, \dots, X_9) under the new RMR ρ' is different from that of the original RMR ρ :

```

>> % set PRNG to be twister with seed 19731511
>> RandStream.setDefaultStream(RandStream('mt19937ar','seed',19731511));
>> x0 = 'r' % set x_0 = 'r'
x0 = r
>> xt = x0; % current state x_t is x_0
>> Visited = x0; % initialise the variable Visited to hold the visited states
>> for t = 1:9 % start a for loop for t = 1,2,...,9
xt = RMR10fFairFreddy(xt, 1-floor(rand+0.5) ); % update the current state at t with new RMR rho'
Visited = strcat(Visited,',',xt); % store the visited state in string Visited
end
>> Visited % disclose the string of visited state separated by commas under new RMR rho'
Visited = r,f,f,r,r,f,f,r,f,f

```

Proposition 98 (Existence and non-uniqueness of RMR) Every transition matrix P on a finite state space \mathbb{X} has a random mapping representation (RMR) that is not necessarily unique.

Proof: Let $\mathbb{X} = \{s_1, s_2, \dots, s_k\}$ be sequentially accessible by $\psi(i) = s_i : \{1, 2, \dots, k\} \rightarrow \mathbb{X}$. We will prove the proposition constructively via the inversion sampler for \mathbb{X} -valued family of ψ -transformed de Moivre RVs. Let the auxiliary RV W be Uniform(0, 1) with $\mathbb{W} = [0, 1]$ and let $\rho(x, w) : \mathbb{X} \times \mathbb{W} \rightarrow \mathbb{X}$ be given by $F^{[-1]}(u; \theta_1, \theta_2, \dots, \theta_k)$ of Equation 6.5, the inverse DF of the de Moivre($\theta_1, \theta_2, \dots, \theta_k$) RV, as follows:

$$\rho(x, w) = \psi \left(F^{[-1]}(w; P(x, s_1), P(x, s_2), \dots, P(x, s_k)) \right), \quad \text{for each } x \in \mathbb{X}.$$

Then, by construction, this ρ is indeed an RMR of P since

$$P(\{\rho(x, W) = y\}) = P(x, y) \quad \text{for each } (x, y) \in \mathbb{X}^2.$$

Non-uniqueness is established by constructing another RMR for P as $\rho'(x, w) = \rho(x, 1 - w)$.

Labwork 205 (Markov chain from $\{\text{de Moivre}(P(x, .))\}_{x \in \mathbb{X}}$ RVs) Let us implement a function that will take a transition matrix P as input and produce a sequence of n states $(x_0, x_1, \dots, x_{n-1})$ visited by the corresponding Markov chain (X_0, X_1, \dots, X_n) using the function in the following M-file.

```
MCSimBydeMoivre.m
function VisitedStateIdxs = MCSimBydeMoivre(idx0, P, n)
% input: idx0 = index of initial state x_0, psi(idx0) = x_0
%         P = transition probability matrix (has to be stochastic matrix)
%         n = number of time steps to simulate, n >= 0
% output: VisitedStateIdxs = idx0, idx1, ..., idxn
VisitedStateIdxs = zeros(1,n);
VisitedStateIdxs(1, 0+1) = idx0; % initial state index is the input idx0
for i=1:n-1
    CurrentState = VisitedStateIdxs(1, i); % current state
    Thetas = P(CurrentState,:);
    VisitedStateIdxs(1, i+1) = SimdeMoivreOnce(rand, Thetas); % next state
end
end
```

Simulation 206 (Another simulation of Freddy's jumps) Let us simulate a sequence of 10 jumps of Flippant Freddy with fair dice and coin by using the function `MCSimBydeMoivre` defined in Labwork 205 as follows:

```
>> % set PRNG to be twister with seed 19731511
>> RandStream.setDefaultStream(RandStream('mt19937ar','seed',19731511));
>> MCSimBydeMoivre(1,[0.5 0.5; 0.5 0.5], 10)
ans =
     1     1     2     1     2     1     2     1     1     2
```

Here we need to further transform the output by $\psi : \{1, 2\} \rightarrow \{r, f\}$ with $\psi(1) = r$ and $\psi(2) = f$.

Labwork 207 (Markov chain from $\{\text{de Moivre}(P(x, .))\}_{x \in \mathbb{X}}$ RVs by Recursion) Let us implement a recursive function that will take a transition matrix P as input and produce a sequence of n states $(x_0, x_1, \dots, x_{n-1})$ visited by the corresponding Markov chain (X_0, X_1, \dots, X_n) using the function in the following M-file.

```
MCSimBydeMoivreRecurse.m
function VisitedStateIdxs = MCSimBydeMoivreRecurse(VisitedStateIdxs, P, n)
% input: VisitedStateIdxs = array of indexes of states visited so far
%         P = transition probability matrix (has to be stochastic matrix)
%         n = number of time steps to simulate, n >= 0
```

```
% output: VisitedStateIdxs = idx0, idx1, ..., idxn
i = length(VisitedStateIdxs);
if i < n
    CurrentState = VisitedStateIdxs(1, i); % current state
    Thetas = P(CurrentState,:);
    % recursion
    VisitedStateIdxs= MCSimBydeMoivreRecurse([VisitedStateIdxs SimdeMoivreOnce(rand,Thetas)],P,n); % next state
end
end
```

Now, let us compare this recursive function to the function `MCSimBydeMoivre` defined in Lab-work 205 as follows:

```
format compact
P=[1/3 2/3;1/4 4/5]
initial = 2
visited = [initial];
n = 12;

s = RandStream('mt19937ar','Seed', 5489);
RandStream.setDefaultStream(s) % reset the PRNG to default state Mersenne Twister with seed=5489

VisitByMethod1 = MCSimBydeMoivre(initial, P, n)

s = RandStream('mt19937ar','Seed', 5489);
RandStream.setDefaultStream(s) % reset the PRNG to default state Mersenne Twister with seed=5489

VisitByMethod2 = MCSimBydeMoivreRecurse(visited, P, n)
```

```
>> CompareMCSimBydeMoivreMethods
P =
    0.3333    0.6667
    0.2500    0.8000
initial =
    2
VisitByMethod1 =
    2    2    2    1    2    2    1    1    2    2    2    2    1
VisitByMethod2 =
    2    2    2    1    2    2    1    1    2    2    2    2    1
```

Therefore, both methods produce the same output. The recursive version of the function is more versatile and useful in the sequel.

Simulation 208 Using the function `MCSimBydeMoivre` of Labwork 205 simulate twenty states visited by the Markov chain in Exercise 197.

Simulation 209 (Drunkard's walk around the block) Consider the Markov chain $(X_t)_{t \in \mathbb{Z}_+}$ on $\mathbb{X} = \{0, 1, 2, 3\}$ with initial distribution $\mathbf{1}_{\{3\}}(x)$ and transition matrix

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

Draw the transition diagram for this Markov chain. Do you see why this chain can be called the “drunkard’s walk around the block”? Using the function `MCSimBydeMoivre` of Labwork 205 simulate a sequence of ten states visited by the drunkard (don’t forget to subtract 1 from the output of `MCSimBydeMoivre` since $\psi(i) = i - 1$ here).

There are many distinct and interesting RMRs of any given transition matrix P beyond that constructed in the proof above. Good RMRs will typically simplify the simulation of a Markov chain. Let us consider examples of Markov chains that can be simulated by simpler methods.

Example 210 (Jukes & Cantor Model of DNA mutation) The “blueprint” of organisms on earth are typically given by a long sequence of deoxyribonucleic acid or DNA. A DNA sequence of length n can be thought of as a string made up of n alphabets from the set of four nucleotides $\{a, c, g, t\}$. For example a DNA sequence of length 3 is *agg* and another is *act*. When an organism goes through time to “stay alive” it has to copy its DNA. This copying process is not perfect and mistakes or mutations are made. We can look at a particular position of a DNA sequence and keep track of its mutations using a simple Markov chain due to Jukes and Cantor [Jukes TH and Cantor CR (1969) Evolution of protein molecules. In Munro HN, editor, Mammalian Protein Metabolism, pp. 21-132, Academic Press, New York.] with the following transition probability matrix:

$$P = \begin{pmatrix} & a & c & g & t \\ a & 0 & \frac{1}{3} & \frac{1}{3} & \frac{1}{3} \\ c & \frac{1}{3} & 0 & \frac{1}{3} & \frac{1}{3} \\ g & \frac{1}{3} & \frac{1}{3} & 0 & \frac{1}{3} \\ t & \frac{1}{3} & \frac{1}{3} & \frac{1}{3} & 0 \end{pmatrix}.$$

Suppose you initially observe the particular position of a DNA sequence at state c and want to simulate a sequence of states visited due to mutation under this Markov chain model. We can achieve this by improvising the inversion sampler for the equi-probable de Moivre($1/3, 1/3, 1/3$) RV (Algorithm 4) in the following RMR:

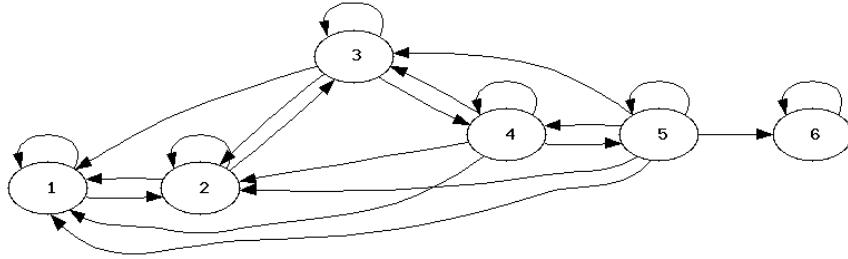
$$\rho(x, U) : \{a, c, g, t\} \times [0, 1] \rightarrow \{a, c, g, t\}, \quad \rho(x, U) = \psi_x(\lceil 3U \rceil), \quad U \sim \text{Uniform}(0, 1),$$

with any fixed bijection $\psi_x(i) : \{1, 2, 3\} \rightarrow \{a, c, g, t\} \setminus \{x\}$ for each $x \in \{a, c, g, t\}$. Then we can produce a sequence of visited states as follows:

$$X_0 \leftarrow c, \quad X_i \leftarrow \rho(X_{i-1}, U_i), \quad i = 1, 2, \dots.$$

Example 211 (Six Lounges) Suppose there are six lounges with doors that allow you to go only in one direction. These lounges are labelled by 1, 2, 3, 4, 5 and 6 and form our state space \mathbb{X} with one-way-doors as shown in Figure 8.4. Every hour an alarm rings and it can be heard in all six lounges. In each lounge $i \in \{1, 2, 3, 4, 5\}$ there is a fair i -sided polyhedral cylinder whose i faces are marked with lounge numbers $1, 2, \dots, i$ but in lounge 6 there is a hexagonal cylinder with all six faces marked by 6. Suppose you start from lounge number 1. When the hourly alarm rings you toss the polyhedral cylinder in the current lounge over the floor. When the cylinder comes to rest, you note the number on the face that touches the floor and go to the lounge labelled by this number. This scheme of lounge hopping can be formalised as a Markov

Figure 8.4: Transition diagram over six lounges (without edge probabilities).



chain starting at lounge number 1 and evolving according to the transition matrix P :

$$P = \begin{pmatrix} 1 & 2 & 3 & 4 & 5 & 6 \\ 1 & \frac{1}{2} & \frac{1}{2} & 0 & 0 & 0 \\ 2 & \frac{1}{3} & \frac{1}{3} & \frac{1}{3} & 0 & 0 \\ 3 & \frac{1}{4} & \frac{1}{4} & \frac{1}{4} & \frac{1}{4} & 0 \\ 4 & \frac{1}{5} & \frac{1}{5} & \frac{1}{5} & \frac{1}{5} & \frac{1}{5} \\ 5 & \frac{1}{6} & \frac{1}{6} & \frac{1}{6} & \frac{1}{6} & \frac{1}{6} \\ 6 & 0 & 0 & 0 & 0 & 0 \end{pmatrix}.$$

The inversion samplers for the family of equi-probable $\{\text{de Moivre}(1/i, 1/i, \dots, 1/i)\}_{i \in \{1, 2, \dots, 5\}}$ RVs (Algorithm 4) and the Point Mass(6) RV (Simulation 167) can be combined in the random mapping representation:

$$\rho(i, U) : \mathbb{X} \times [0, 1] \rightarrow \mathbb{X}, \quad \rho(i, U) = \lceil iU \rceil \mathbf{1}_{\{1, 2, 3, 4, 5\}}(i) + 6\mathbf{1}_{\{6\}}(i), \quad U \sim \text{Uniform}(0, 1) ,$$

in order to simulate a sequence of states from this markov chain as follows:

$$X_0 \leftarrow 1, \quad X_i \leftarrow \rho(X_{i-1}, U_i), \quad i = 1, 2, \dots . \quad (8.12)$$

Simulation 212 (Trapped in lounge 6) Implement the Algorithm described in Equation 8.12 in a MATLAB program to simulate the first ten states visited by the Markov chain in Example 211. Recall the “Hotel California” character of lounge 6 – *you can check out anytime you like, but you can never leave!* Repeat this simulation 1000 times and find the fraction of times your are not trapped in lounge 6 by the tenth time step.

Exercise 213 (Drunkard’s walk around a polygonal block with k corners) Can you think of another way to simulate the “drunkard’s walk around a polygonal block with k corners” labelled by $0, 1, \dots, k - 1$ that is more efficient than using the `MCSimBydeMoivre` function which relies on the `SimdeMoivreOnce` function that implements Algorithm 5 with an average-case efficiency that is linear in k ?

Hint: think of the drunkard tossing a fair coin to make his decision of where to go next from each corner and arithmetic mod k .

8.3 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 99 (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 100 (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 8.3. If $(p, q) \in (0, 1)^2$, then the corresponding Markov chain is irreducible because we can go from rolloplia 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 rolloplia 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. For instance, from the transition diagram in Figure 8.4 of the lounge-hopping Markov chain of Example 211, it is clear that if you start at state 6 you cannot find any arrow going to any other state. Therefore, the chain is reducible since $6 \not\rightarrow i$ for any $i \in \{1, 2, 3, 4, 5\}$.

Exercise 214 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.

Definition 101 (Return times and period) Let $\mathbb{T}(x) := \{t \in \mathbb{N} : P^t(x, x) > 0\}$ be the set of **possible return times** to the starting state x . The **period** of state x is defined to be $\gcd(\mathbb{T}(x))$, the greatest common divisor of $\mathbb{T}(x)$. When the period of a state x is 1, i.e., $\gcd(\mathbb{T}(x)) = 1$, then x is said to be an **aperiodic state**.

Proposition 102 If the Markov chain $(X_t)_{t \in \mathbb{Z}_+}$ with transition matrix P on state space \mathbb{X} is irreducible then $\gcd(\mathbb{T}(x)) = \gcd(\mathbb{T}(y))$ for any $(x, y) \in \mathbb{X}^2$.

Proof: Fix any pair of states $(x, y) \in \mathbb{X}^2$. Since, P is irreducible, $x \leftrightarrow y$ and therefore there exists natural numbers $\eta(x, y)$ and $\eta(y, x)$ such that $P^{\eta(x,y)}(x, y) > 0$ and $P^{\eta(y,x)}(y, x) > 0$. Let $\eta' = \eta(x, y) + \eta(y, x)$ and observe that $\eta' \in \mathbb{T}(x) \cap \mathbb{T}(y)$, $\mathbb{T}(x) \subset \mathbb{T}(y) - \eta' := \{t - \eta' : t \in \mathbb{T}(y)\}$ and

$\gcd(\mathbb{T}(y))$ divides all elements in $\mathbb{T}(x)$. Thus, $\gcd(\mathbb{T}(y)) \leq \gcd(\mathbb{T}(x))$. By a similar argument we can also conclude that $\gcd(\mathbb{T}(x)) \leq \gcd(\mathbb{T}(y))$. Therefore $\gcd(\mathbb{T}(x)) = \gcd(\mathbb{T}(y))$.

Definition 103 (Aperiodic) A Markov chain $(X_t)_{t \in \mathbb{Z}_+}$ with transition matrix P on state space \mathbb{X} is said to be aperiodic if all of its states are aperiodic, i.e., $\gcd(\mathbb{T}(x)) = 1$ for every $x \in \mathbb{X}$. If a chain is not aperiodic, we call it **periodic**.

We have already seen example of irreducible Markov chains that were either periodic or aperiodic. For instance, Freddy's Markov chain with $(p, q) \in (0, 1)^2$ is aperiodic since the period of either of its two states is given by $\gcd(\{1, 2, 3, \dots\}) = 1$. However, the Markov chain model for a drunkard's walk around a block over the state space $\{0, 1, 2, 3\}$ (Simulation 209) is periodic because you can only return to the starting state in an even number of time steps and

$$\gcd(\mathbb{T}(0)) = \gcd(\mathbb{T}(1)) = \gcd(\mathbb{T}(2)) = \gcd(\mathbb{T}(3)) = \gcd(\{2, 4, 6, \dots\}) = 2 \neq 1 .$$

Exercise 215 Show that the Markov chain corresponding to a drunkard's walk around a polygonal block with k corners is irreducible for any integer $k > 1$. Show that it is aperiodic only when k is odd and has period 2 when k is even.

Proposition 104 Let $A = \{a_1, a_2, \dots\} \subset \mathbb{N}$ that satisfies the following two conditions:

1. A is a **nonlattice**, meaning that $\gcd(A) = 1$ and
2. A is closed under addition, meaning that if $(a, a') \in A^2$ then $a + a' \in A$.

Then there exists a positive integer $\eta < \infty$ such that $n \in A$ for all $n \geq \eta$.

Proof: See Proofs of Lemma 1.1, Lemma 1.2 and Theorem 1.1 in Appendix of *Pierre Brémaud, Markov Chains, Gibbs Fields, Monte Carlo Simulation, and Queues, Springer, 1999*.

Proposition 105 If the Markov chain $(X_t)_{t \in \mathbb{Z}_+}$ with transition matrix P on state space \mathbb{X} is irreducible and aperiodic then there is an integer τ such that $P^t(x, x) > 0$ for all $t \geq \tau$ and all $x \in \mathbb{X}$.

Proof: TBD

Proposition 106 If the Markov chain $(X_t)_{t \in \mathbb{Z}_+}$ with transition matrix P on state space \mathbb{X} is irreducible and aperiodic then there is an integer τ such that $P^t(x, y) > 0$ for all $t \geq \tau$ and all $(x, y) \in \mathbb{X}^2$.

Proof: TBD

Exercise 216 (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 217 (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:

$$P(W_t = w) = \begin{cases} \frac{1}{8} & \text{if } w \in \{(1, 1), (1, 0), (1, -1), (0, -1), (-1, -1), (-1, 0), (-1, 1), (0, 1)\}, \\ 0 & \text{otherwise.} \end{cases}$$

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.

8.4 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 sattes that the distribution of X_t over \mathbb{X} settles down as $t \rightarrow \infty$, provided the Markov chain is irreducible and aperiodic.

Definition 107 (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 108 (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**

$$T(x, y) = \min\{t \geq 1 : X_t = y\} .$$

and let $T(x, y) = \min\{\} = \infty$ if the Markov chain never visits y after having started from x . Let the **mean hitting time**

$$\tau(x, y) := E(T(x, y)),$$

be the expected time taken to reach y after having started at x . Note that $\tau(x, x)$ is the **mean return time** to state x .

Proposition 109 (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 $(x, y) \in \mathbb{X}^2$,

$$\mathbb{P}(T(x, y) < \infty) = 1 ,$$

and the mean hitting time is finite, i.e.,

$$\tau(x, y) < \infty .$$

Proposition 110 (Existence of Stationary distribution) For any irreducible and aperiodic Markov chain there exists at least one stationary distribution.

Proof: TBD

Definition 111 (Total variation distance) If $\nu_1 := (\nu_1(x))_{x \in \mathbb{X}}$ and $\nu_2 := (\nu_2(x))_{x \in \mathbb{X}}$ are elements of $\mathcal{P}(\mathbb{X})$, the set of all probability distributions on $\mathbb{X} := \{s_1, s_2, \dots, s_k\}$, then we define the **total variation distance** between ν_1 and ν_2 as

$$d_{TV}(\nu_1, \nu_2) := \frac{1}{2} \sum_{x \in \mathbb{X}} \text{abs}(\nu_1(x) - \nu_2(x)), \quad d_{TV} : \mathcal{P}(\mathbb{X}) \times \mathcal{P}(\mathbb{X}) \rightarrow [0, 1] . \quad (8.13)$$

If ν_1, ν_2, \dots and ν are probability distributions on \mathbb{X} , then we say that ν_t **converges in total variation** to ν as $n \rightarrow \infty$ and write $\nu_t \xrightarrow{TV} \nu$, if

$$\lim_{t \rightarrow \infty} d_{TV}(\nu_t, \nu) = 0 .$$

Observe that if $d_{TV}(\nu_1, \nu_2) = 0$ then $\nu_1 = \nu_2$. The constant $1/2$ in Equation 8.13 ensures that the range of d_{TV} is in $[0, 1]$. If $d_{TV}(\nu_1, \nu_2) = 1$ then ν_1 and ν_2 have disjoint supports, i.e., we can partition \mathbb{X} into \mathbb{X}_1 and \mathbb{X}_2 , i.e., $\mathbb{X} = \mathbb{X}_1 \cup \mathbb{X}_2$ and $\mathbb{X}_1 \cap \mathbb{X}_2 = \emptyset$, such that $\sum_{x \in \mathbb{X}_1} \nu_1(x) = 1$ and $\sum_{x \in \mathbb{X}_2} \nu_2(x) = 1$. The total variation distance gets its name from the following natural interpretation:

$$d_{TV}(\nu_1, \nu_2) = \max_{A \subset \mathbb{X}} \text{abs}(\nu_1(A) - \nu_2(A)) .$$

This interpretation means that the total variation distance between ν_1 and ν_2 is the maximal difference in probabilities that the two distributions assign to any one event $A \in \sigma(\mathbb{X}) = 2^\mathbb{X}$.

In words, Proposition 112 says that if you run the chain for a sufficiently long enough time t , then, regardless of the initial distribution μ_0 , the distribution at time t will be close to the stationary distribution π . This is referred to as the Markov chain **approaching equilibrium** or **stationarity** as $t \rightarrow \infty$.

Proposition 112 (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

$$\mu_t \xrightarrow{TV} \pi . \quad (8.14)$$

Proof: TBD

Proposition 113 (Uniqueness of stationary distribution) Any irreducible aperiodic Markov chain has a unique stationary distribution.

Proof: TBD

Exercise 218 Consider the Markov chain on $\{1, 2, 3, 4, 5, 6\}$ with the following transition matrix:

$$P = \begin{matrix} & \begin{matrix} 1 & 2 & 3 & 4 & 5 & 6 \end{matrix} \\ \begin{matrix} 1 \\ 2 \\ 3 \\ 4 \\ 5 \\ 6 \end{matrix} & \left(\begin{matrix} \frac{1}{2} & \frac{1}{2} & 0 & 0 & 0 & 0 \\ \frac{1}{2} & \frac{1}{2} & 0 & 0 & 0 & 0 \\ 0 & 0 & \frac{1}{4} & \frac{3}{4} & 0 & 0 \\ 0 & 0 & \frac{3}{4} & \frac{1}{4} & 0 & 0 \\ 0 & 0 & 0 & 0 & \frac{3}{4} & \frac{1}{4} \\ 0 & 0 & 0 & 0 & \frac{1}{4} & \frac{3}{4} \end{matrix} \right) \end{matrix} .$$

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 219 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 220 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.

8.5 Reversibility

We introduce another specific property called reversibility. This property will assist in conjuring Markov chains with a desired stationary distribution.

Definition 114 (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 $(x, y) \in \mathbb{X}^2$:

$$\pi(x)P(x, y) = \pi(y)P(y, x) . \tag{8.15}$$

A Markov chain that has a reversible distribution is said to be a reversible Markov chain.

In words, $\pi(x)P(x, y) = \pi(y)P(y, x)$ says that if you start the chain at the reversible distribution π , i.e., $\mu_0 = \pi$, then the probability of going from x to y is the same as that of going from y to x .

Proposition 115 (A reversible π is a stationary π) Let $(X_t)_{t \in \mathbb{Z}_+}$ be a Markov chain on $\mathbb{X} = \{s_1, s_2, \dots, s_k\}$ with transition matrix P . If π is a reversible distribution for $(X_t)_{t \in \mathbb{Z}_+}$ then π is a stationary distribution for $(X_t)_{t \in \mathbb{Z}_+}$.

Proof: Suppose π is a reversible distribution for $(X_t)_{t \in \mathbb{Z}_+}$ then π is a probability distribution on \mathbb{X} and $\pi(x)P(x, y) = \pi(y)P(y, x)$ for each $(x, y) \in \mathbb{X}^2$. We need to show that for any $y \in \mathbb{X}$ we have

$$\pi(y) = \sum_{x \in \mathbb{X}} \pi(y)P(y, x) .$$

Fix a $y \in \mathbb{X}$,

$$\begin{aligned} LHS &= \pi(y) = \pi(y)1 = \pi(y) \sum_{x \in \mathbb{X}} P(y, x), \text{ since } P \text{ is a stochastic matrix} \\ &= \sum_{x \in \mathbb{X}} \pi(y)P(y, x) = \sum_{x \in \mathbb{X}} \pi(x)P(x, y), \text{ by reversibility} \\ &= RHS . \end{aligned}$$

Definition 116 (Graph) A **Graph** $\mathbb{G} := (\mathbb{V}, \mathbb{E})$ consists of a **vertex set** $\mathbb{V} := \{v_1, v_2, \dots, v_k\}$ together with an **edge set** $\mathbb{E} := \{e_1, e_2, \dots, e_l\}$. Each edge connects two of the vertices in \mathbb{V} . An edge e_h connecting vertices v_i and v_j is denoted by $\langle v_i, v_j \rangle$. Two vertices are **neighbours** if they share an edge. The **neighbourhood** of a vertex v_i denoted by $\text{nbhd}(v_i) := \{v_j : \langle v_i, v_j \rangle \in \mathbb{E}\}$ is the set of neighbouring vertices of v_i . The number of neighbours of a vertex v_i in an undirected graph is called its **degree** and is denoted by $\deg(v_i)$. Note that $\deg(v_i) = \#\text{nbhd}(v_i)$. In a graph 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 preserve the order of the pair of vertices they connect or they can be **undirected**. An edge can be **weighted** by being associated with a real number called its weight. We can represent a directed graph by its **adjacency matrix** given by:

$$A := (A(v_i, v_j))_{(v_i, v_j) \in \mathbb{V} \times \mathbb{V}}, \quad A(v_i, v_j) = \begin{cases} 1 & \text{if } \langle v_i, v_j \rangle \in \mathbb{E} \\ 0 & \text{otherwise} . \end{cases}$$

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 $\text{ideg}(v_i)$ and $\text{odeg}(v_i)$ respectively. Note that a transition diagram of a Markov chain is a weighted directed graph and is represented by the transition probability matrix.

Model 21 (Random Walk on an Undirected Graph) A random walk on an undirected graph $\mathbb{G} = (\mathbb{V}, \mathbb{E})$ is a Markov chain with state space $\mathbb{V} := \{v_1, v_2, \dots, v_k\}$ 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 neighbours of v_i at time $t + 1$. If $\deg(v_i)$ is the degree of v_i then the transition probabilities of this Markov chain is

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

Proposition 117 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 8.15 for each $(v_i, v_j) \in \mathbb{V}^2$. Fix a pair of states $(v_i, v_j) \in \mathbb{V}^2$, then

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

By Proposition 115 π is also the stationary distribution.

Exercise 221 Prove Proposition 117 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 222 (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 117, 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 223 (Triangulated Quadrangle) The random walk on the undirected graph

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

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

$$A = \begin{pmatrix} 1 & 2 & 3 & 4 \\ 0 & 1 & 1 & 0 \\ 1 & 0 & 1 & 1 \\ 1 & 1 & 0 & 1 \\ 0 & 1 & 1 & 0 \end{pmatrix}, \quad P = \begin{pmatrix} 1 & 2 & 3 & 4 \\ 0 & \frac{1}{2} & \frac{1}{2} & 0 \\ \frac{1}{3} & 0 & \frac{1}{3} & \frac{1}{3} \\ \frac{1}{3} & \frac{1}{3} & 0 & \frac{1}{3} \\ 0 & \frac{1}{2} & \frac{1}{2} & 0 \end{pmatrix}, \quad \text{Graph: } \begin{array}{c} \text{1} \text{---} \text{2} \text{---} \text{3} \text{---} \text{4} \\ | \qquad \qquad \qquad | \\ \text{1} \text{---} \text{3} \end{array}.$$

By Proposition 117, the stationary distribution of the random walk on \mathbb{G} is

$$\pi = \left(\frac{\deg(v_1)}{d}, \frac{\deg(v_2)}{d}, \frac{\deg(v_3)}{d}, \frac{\deg(v_4)}{d} \right) = \left(\frac{2}{10}, \frac{3}{10}, \frac{3}{10}, \frac{2}{10} \right).$$

Exercise 224 Show that the Drunkard's walk around the block from Simulation 209 is a random walk on the undirected graph $\mathbb{G} = (\mathbb{V}, \mathbb{E})$ with $\mathbb{V} = \{0, 1, 2, 3\}$ and $\mathbb{E} = \{\langle 0, 1 \rangle, \langle 1, 2 \rangle, \langle 2, 3 \rangle, \langle 0, 3 \rangle\}$. What is its reversible distribution?

Example 225 (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 $\mathbf{1}_{\{3\}}(x)$ and transition matrix

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

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 $\pi = (1/4, 1/4, 1/4, 1/4)$ (verify $\pi P = \pi$).

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 (aperiodicity is not necessary for uniqueness of π) π is the unique stationary distribution. Due to Proposition 115, π 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,

$$\pi(0)P(0, 1) = \frac{1}{4} \times \frac{1}{3} = \frac{1}{12} < \frac{1}{6} = \frac{1}{4} \times \frac{2}{3} = \pi(1)P(1, 0).$$

Exercise 226 Find the stationary distribution of the Markov chain in Exercise 217.

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

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

Example 227 (Directed Triangulated Quadrangle) The random walk on the directed graph

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

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

$$A = \begin{pmatrix} 1 & 2 & 3 & 4 \\ 1 & 0 & 1 & 0 & 0 \\ 2 & 0 & 0 & 1 & 1 \\ 3 & 1 & 0 & 0 & 0 \\ 4 & 0 & 0 & 1 & 0 \end{pmatrix}, \quad P = \begin{pmatrix} 1 & 2 & 3 & 4 \\ 1 & 0 & 1 & 0 & 0 \\ 2 & 0 & 0 & \frac{1}{2} & \frac{1}{2} \\ 3 & 1 & 0 & 0 & 0 \\ 4 & 0 & 0 & 1 & 0 \end{pmatrix}, \quad \text{Diagram: } \begin{array}{c} \text{1} \rightarrow \text{2} \rightarrow \text{4} \rightarrow \text{3} \rightarrow \text{1} \\ \text{1} \leftarrow \text{2} \leftarrow \text{4} \leftarrow \text{3} \end{array}.$$

Exercise 228 Show that there is no reversible distribution for the Markov chain in Example 227.

Example 229 (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 at random again. Google's PageRank formula uses a random surfer model who gets bored after several clicks and switches to a random page. The PageRank value of a page reflects the chance that the random surfer will land on that page by clicking on a link. The stationary distribution of the random surfer on the world wide web is a very successful model for ranking pages.

Model 23 (Lazy Random Walk) You can convert a random walk on an undirected graph $\mathbb{G} = (\mathbb{V}, \mathbb{E})$ into a **lazy random walk** on \mathbb{G} by the following steps:

- Add loops to each vertex in $\mathbb{V} = \{v_1, v_2, \dots, v_k\}$ to obtain a new set of edges $\mathbb{E}' = \mathbb{E} \cup \{\langle v_1, v_1 \rangle, \langle v_2, v_2 \rangle, \dots, \langle v_k, v_k \rangle\}$.
- Construct the lazy graph $\mathbb{G}' = (\mathbb{V}, \mathbb{E}')$.
- Do a random walk on the undirected graph \mathbb{G}' .

The lazy random walk allows us to introduce aperiodicity quite easily.

Exercise 230 (Lazy Random Walk on the Triangulated Quadrangle) Consider the random walk of Example 223 on the undirected graph

$$\mathbb{G} = (\{1, 2, 3, 4\}, \{\langle 1, 2 \rangle, \langle 3, 1 \rangle, \langle 2, 3 \rangle, \langle 2, 4 \rangle, \langle 4, 3 \rangle\}) .$$

Construct the lazy random walk on \mathbb{G} , obtain its transition probability matrix and state transition diagram. Show that the stationary distribution of this lazy random walk on \mathbb{G} is

$$\pi = \left(\frac{3}{14}, \frac{4}{14}, \frac{4}{14}, \frac{3}{14} \right) .$$

Model 24 (Random Walks on Groups) Under 

Model 25 (Birth-Death chains) Under 

8.6 Metropolis-Hastings Markov chain

Definition 118 (Metropolis-Hastings Markov chain) If we are given an irreducible Markov chain $(Y_t)_{t \in \mathbb{Z}_+}$ called the **base chain** or the **proposal chain** on a finite state space $\mathbb{X} = \{s_1, s_2, \dots, s_k\}$ with transition probability matrix $Q = (Q(x, y))_{(x,y) \in \mathbb{X}^2}$ and some probability distribution π on \mathbb{X} of interest that may only be known up to a normalizing constant as $\tilde{\pi}$, i.e., $\pi(x) = (\sum_{z \in \mathbb{X}} \tilde{\pi}(z))^{-1} \tilde{\pi}(x)$ for each $x \in \mathbb{X}$, then we can construct a new Markov chain $(X_t)_{t \in \mathbb{Z}_+}$ called the **Metropolis-Hastings** chain on \mathbb{X} with the following transition probabilities:

$$P(x, y) = \begin{cases} Q(x, y)a(x, y) & \text{if } x \neq y \\ 1 - \sum_{z \in \{z \in \mathbb{X}: z \neq x\}} Q(x, z)a(x, z) & \text{if } x = y \end{cases} , \quad (8.16)$$

where the acceptance probability is

$$a(x, y) := \min \left\{ 1, \frac{\pi(y)}{\pi(x)} \frac{Q(y, x)}{Q(x, y)} \right\} . \quad (8.17)$$

Note that we only need to know π up to ratios. Thus, $\pi(y)/\pi(x)$ in $a(x, y)$ can be replaced by $\tilde{\pi}(y)/\tilde{\pi}(x)$ since

$$\frac{\pi(y)}{\pi(x)} = \frac{(\sum_{z \in \mathbb{X}} \tilde{\pi}(z))^{-1} \tilde{\pi}(y)}{(\sum_{z \in \mathbb{X}} \tilde{\pi}(z))^{-1} \tilde{\pi}(x)} = \frac{\tilde{\pi}(y)}{\tilde{\pi}(x)} .$$

Algorithm 9 describes how to simulate samples from a Metropolis-Hastings Markov chain.

Proposition 119 (Stationarity of the Metropolis-Hastings chain) The Metropolis-Hastings chain constructed according to Definition 118 has π as its stationary distribution.

Proof: It suffices to show that π is the reversible distribution for $(X_t)_{t \in \mathbb{Z}_+}$, i.e., for each $(x, y) \in \mathbb{X}^2$, $\pi(x)P(x, y) = \pi(y)P(y, x)$. Fix a pair $(x, y) \in \mathbb{X}^2$ and suppose $x \neq y$. Then,

$$\begin{aligned} \pi(x)P(x, y) &= \pi(x)Q(x, y)a(x, y) \\ &= \pi(x)Q(x, y) \min \left\{ 1, \frac{\pi(y) Q(y, x)}{\pi(x) Q(x, y)} \right\} \\ &= \min \left\{ \pi(x)Q(x, y), \pi(x)Q(x, y) \frac{\pi(y) Q(y, x)}{\pi(x) Q(x, y)} \right\} \\ &= \min \{ \pi(x)Q(x, y), \pi(y)Q(y, x) \} \\ &= \min \{ \pi(y)Q(y, x), \pi(x)Q(x, y) \} \\ &= \min \left\{ \pi(y)Q(y, x), \pi(y)Q(y, x) \frac{\pi(x) Q(x, y)}{\pi(y) Q(y, x)} \right\} \\ &= \pi(y)Q(y, x) \min \left\{ 1, \frac{\pi(x) Q(x, y)}{\pi(y) Q(y, x)} \right\} \\ &= \pi(y)P(y, x) . \end{aligned}$$

When $x = y$, reversibility is trivially satisfied since $\pi(x)P(x, y) = \pi(y)P(y, x) = \pi(x)P(x, x)$.

Definition 120 If the base chain $(Y_t)_{t \in \mathbb{Z}_+}$ in the Metropolis-Hastings Markov chain of Definition 118 has a symmetric transition matrix Q with $Q(x, y) = Q(y, x)$ for each $(x, y) \in \mathbb{X}^2$ then the acceptance probability in Equation 8.17 simplifies to

$$a(x, y) = \min \left\{ 1, \frac{\pi(y)}{\pi(x)} \right\} ,$$

and the corresponding Metropolis-Hastings chain $(X_t)_{t \in \mathbb{Z}_+}$ is called the **Metropolis chain**.

Suppose you know neither the vertex set \mathbb{V} nor the edge set \mathbb{E} entirely for an undirected graph $\mathbb{G} = (\mathbb{V}, \mathbb{E})$ but you are capable of walking locally on \mathbb{G} . In other words, if you are currently at vertex x you are able to make a move to one of the neighbouring vertices of x . However, you do not know every single vertex in \mathbb{V} or the entire set of edges \mathbb{E} as an adjacency matrix for instance. Several real-world problems fall in this class. Some examples include the random surfer on www to rank web pages (Example 229), social network analyses in facebook or twitter, exact tests for contingency tables, etc.

Model 26 (Metropolis-Hastings Random Walk on Graph) Let $\mathbb{G} = (\mathbb{V}, \mathbb{E})$ be an undirected graph and let $(Y_t)_{t \in \mathbb{Z}_+}$ with transition matrix Q be an irreducible random walk on \mathbb{G}

Algorithm 9 Metropolis-Hastings Markov chain

1: *input:*

- (1) shape of a target density $\tilde{\pi}(x) = (\sum_{x \in \mathbb{X}} \tilde{\pi}(x)) \pi(x)$,
- (2) sampler for the base chain that can produce samples $y \sim Q(x, \cdot)$.

2: *output:* a sequence of samples x_0, x_1, \dots, x_n from the Metropolis-Hastings Markov chain $(X_t)_{t \in \mathbb{Z}_+}$ with stationary distribution π

3: Choose initial state $x_0 \in \mathbb{X}$ according to μ_0

4: **repeat**

5: At iteration t ,

6: Generate $y \sim Q(x_{t-1}, \cdot)$ and $u \sim \text{Uniform}(0, 1)$,

7: Compute *acceptance probability*

$$a(x_{t-1}, y) = \min \left\{ 1, \frac{\tilde{\pi}(y)}{\tilde{\pi}(x_{t-1})} \frac{Q(y, x_{t-1})}{Q(x_{t-1}, y)} \right\},$$

8: **If** $u \leq a(x_{t-1}, y)$ **then** $x_t \leftarrow y$, **else** $x_t \leftarrow x_{t-1}$

9: **until** desired number of samples n are obtained from $(X_t)_{t \in \mathbb{Z}_+}$

and let π be a probability distribution on $\mathbb{V} = \{v_1, v_2, \dots, v_k\}$ that is known upto a normalizing constant as $\tilde{\pi}$. The **Metropolis-Hastings random walk** on \mathbb{G} is the Metropolis-Hastings Markov chain $(X_t)_{t \in \mathbb{Z}_+}$ on \mathbb{V} with base chain $(Y_t)_{t \in \mathbb{Z}_+}$ and the following transition probabilities:

$$P(x, y) = \begin{cases} \frac{1}{\deg(v_i)} \min \left\{ 1, \frac{\tilde{\pi}(v_j)}{\tilde{\pi}(v_i)} \frac{\deg(v_i)}{\deg(v_j)} \right\} & \text{if } \langle v_i, v_j \rangle \in \mathbb{E} \\ 1 - \sum_{v_l \in \text{nbhd}(v_i)} \left(\frac{1}{\deg(v_i)} \min \left\{ 1, \frac{\tilde{\pi}(v_l)}{\tilde{\pi}(v_i)} \frac{\deg(v_i)}{\deg(v_l)} \right\} \right) & \text{if } v_i = v_j \\ 0 & \text{otherwise} \end{cases}.$$

By Proposition 119, $(X_t)_{t \in \mathbb{Z}_+}$ has π as its stationary distribution. This Markov chain can be simulated as follows:

- Suppose $x_t = v_i$ at time t
- Propose v_j uniformly at random from $\text{nbhd}(v_i)$
- Sample u from $\text{Uniform}(0, 1)$
- If $u < \min\{1, \pi(v_j) \deg(v_i) / \pi(v_i) \deg(v_j)\}$ then $x_{t+1} = v_j$ else $x_{t+1} = x_t$

Model 27 (Metropolis chain on a regular graph) Consider the random walk $(Y_t)_{t \in \mathbb{Z}_+}$ on a regular graph $\mathbb{G} = (\mathbb{V}, \mathbb{E})$ with $\deg(v_i) = \delta$ for every vertex $v_i \in \mathbb{V} = \{v_1, v_2, \dots, v_k\}$ and the symmetric transition matrix

$$Q(v_i, v_j) = \begin{cases} \frac{1}{\delta} & \text{if } \langle v_i, v_j \rangle \in \mathbb{E} \\ 0 & \text{otherwise} \end{cases}.$$

You can sample from a given distribution π on \mathbb{V} by constructing the Metropolis chain with stationary distribution π from the base chain given by $(Y_t)_{t \in \mathbb{Z}_+}$.

Model 28 (sampling from a uniform distribution over an irregular graph) A graph $\mathbb{G} = (\mathbb{V}, \mathbb{E})$ that is not regular is said to be irregular. Clearly, the stationary distribution of a random walk on \mathbb{G} is not uniform. Suppose you want to sample uniformly from \mathbb{V} according to $\pi(v_i) = (\#\mathbb{V})^{-1}$ for each $v_i \in \mathbb{V}$. We can accomplish this by constructing a Metropolis-Hastings Markov chain with the random walk on \mathbb{G} as the base chain and the following transition probabilities:

$$P(v_i, v_j) = \begin{cases} \frac{1}{\deg(v_i)} \min \left\{ 1, \frac{\deg(v_i)}{\deg(v_j)} \right\} & \text{if } \langle v_i, v_j \rangle \in \mathbb{E} \\ 1 - \sum_{v_l \in \text{nbhd}(v_i)} \left(\frac{1}{\deg(v_i)} \min \left\{ 1, \frac{\deg(v_i)}{\deg(v_l)} \right\} \right) & \text{if } v_i = v_j \\ 0 & \text{otherwise} \end{cases}.$$

Thus the Metropolis-Hastings walk on \mathbb{G} is biased against visiting higher degree vertices and thereby samples uniformly from \mathbb{V} at stationarity.

Example 231 (Stochastic Optimization) Let $f : \mathbb{V} \rightarrow \mathbb{R}$ and $\mathbb{G} = (\mathbb{V}, \mathbb{E})$ be an undirected graph. Let the global maximum be

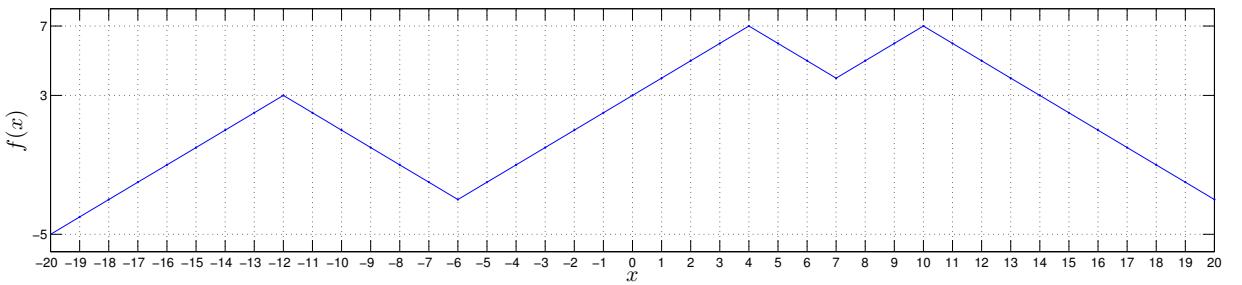
$$f^* := \max_{y \in \mathbb{V}} f(y) ,$$

and the set of maximizers of f be

$$\mathbb{V}^* := \operatorname{argmax}_{x \in \mathbb{V}} f(x) = \{x \in \mathbb{V} : f(x) = f^*\} .$$

In many problems such as maximum likelihood estimation, minimizing a cost function by maximizing its negative, etc, one is interested in $\mathbb{V}^* \subset \mathbb{V}$. This global maximization problem is difficult when $\#\mathbb{V}$ is huge. A deterministic hill-climbing or gradient ascent algorithm that iteratively moves from the current state v_i to a neighbouring state v_j if $f(v_j) > f(v_i)$ can easily get trapped in a local peak of f and thereby miss the global peak attained by elements in \mathbb{V}^* .

Figure 8.5: Stochastic Optimization with Metropolis chain.



For example consider the global maximization problem shown in Figure 8.5 with

$$f^* = 7 \text{ and } \mathbb{V}^* = \{4, 10\} \subset \mathbb{V} = \{-20, -19, \dots, 19, 20\} .$$

The deterministic hill-climbing algorithm will clearly miss \mathbb{V}^* and terminate at the local maximum of 3 at -12 if initialised at any element in $\{-20, -19, \dots, -8, -7\}$. Also, this algorithm will not find both elements in \mathbb{V}^* even when initialised more appropriately.

We will construct a Markov chain to solve this global maximization problem. For a fixed parameter $\lambda \in \mathbb{R}_{>0}$, let

$$\pi_\lambda(x) = \frac{\lambda^{f(x)}}{\sum_{z \in \mathbb{V}} \lambda^{f(z)}} .$$

Since $\pi_\lambda(x)$ is increasing in $f(x)$, $\pi_\lambda(x)$ favours vertices with large $f(x)$. First form a graph $\mathbb{G} = (\mathbb{V}, \mathbb{E})$ by adding edges between the vertices in \mathbb{V} so that you can get from any vertex to any other vertex in \mathbb{V} by following a sequence of edges in \mathbb{E} . Now using the random walk on \mathbb{G} as the base chain let us construct a Metropolis-Hastings chain $(X_t)_{t \in \mathbb{Z}_+}$ on \mathbb{G} with π_λ on \mathbb{V} as its stationary distribution.

For simplicity, let us suppose that \mathbb{G} is a regular graph with a symmetric transition matrix Q for the base chain and thereby making $(X_t)_{t \in \mathbb{Z}_+}$ a Metropolis chain. For instance, in the Example from Figure 8.5 with $\mathbb{V} = \{-20, -19, \dots, 19, 20\}$, we can obtain a Metropolis chain on \mathbb{V} with stationary distribution π_λ by taking \mathbb{E} in $\mathbb{G} = (\mathbb{V}, \mathbb{E})$ to be

$$\mathbb{E} = \{\langle -20, -19 \rangle, \langle -19, -18 \rangle, \langle -18, -17 \rangle, \dots, \langle 17, 18 \rangle, \langle 18, 19 \rangle, \langle 19, 20 \rangle\} .$$

Then, if $f(y) < f(x)$, the Metropolis chain accepts a transition from x to y with probability

$$\frac{\pi_\lambda(y)}{\pi_\lambda(x)} = \frac{\lambda^{f(y)}}{\lambda^{f(x)}} = \lambda^{f(y)-f(x)} = \lambda^{-(f(x)-f(y))} .$$

As $\lambda \rightarrow \infty$, the Metropolis chain approaches the deterministic hill-climbing algorithm and yields a uniform distribution over \mathbb{V}^* as follows:

$$\lim_{\lambda \rightarrow \infty} \pi_\lambda(x) = \lim_{\lambda \rightarrow \infty} \frac{\lambda^{f(x)} / \lambda^{f^*}}{\#\mathbb{V}^* + \sum_{z \in \mathbb{V} \setminus \mathbb{V}^*} \lambda^{f(x)} / \lambda^{f^*}} = \frac{\mathbb{1}_{\mathbb{V}^*}(x)}{\#\mathbb{V}^*} .$$

8.7 Glauber Dynamics

Let \mathbb{S} be a finite set of states. Let \mathbb{V} be a set of vertices. Typically, \mathbb{S} contains characters or colours that can be taken by each site or vertex in \mathbb{V} . Let $x \in \mathbb{S}^\mathbb{V}$ be a configuration, i.e., a function from \mathbb{V} to \mathbb{S} . A configuration can be thought of as a labelling of vertices in \mathbb{V} with elements in \mathbb{S} .

Definition 121 (Glauber dynamics for π) Let \mathbb{V} and \mathbb{S} be finite sets and let $\mathbb{X} \subset \mathbb{S}^\mathbb{V}$ which forms the support of the probability distribution π on $\mathbb{S}^\mathbb{V}$, i.e.,

$$\mathbb{X} = \{x \in \mathbb{S}^\mathbb{V} : \pi(x) > 0\} .$$

The **Glauber dynamics** or **Gibbs sampler** for π is a reversible Markov chain on \mathbb{X} with stationary distribution π under the following transition mechanism. Let the current state at time t be x . To obtain the state at time $t+1$ first choose a vertex v uniformly at random from \mathbb{V} and then choose a new state according to π conditioned on the set of states equal to x at all vertices other than v . We give the details of this transition mechanism next.

For $x \in \mathbb{X}$ and $v \in \mathbb{V}$, define the set of states identical to x everywhere except possibly at v as

$$\mathbb{X}(x, v) := \{y \in \mathbb{X} : y(w) = x(w) \text{ for all } w \neq v\} .$$

Now let

$$\pi^{x,v}(y) := \pi(y|\mathbb{X}(x,v)) = \begin{cases} \left(\sum_{z \in \mathbb{X}(x,v)} \pi(z)\right)^{-1} \pi(y) & \text{if } y \in \mathbb{X}(x,v) \\ 0 & \text{if } y \notin \mathbb{X}(x,v) \end{cases}$$

be the distribution π conditioned on $\mathbb{X}(x,v)$. Therefore the rule for updating the current state x is:

- pick a vertex v uniformly at random from \mathbb{V} ,
- choose a new configuration by sampling from $\pi^{x,v}$.

Proposition 122 (Stationarity of Glauber dynamics) The Glauber dynamics for π on $\mathbb{X} \subset \mathbb{S}^{\mathbb{V}}$ has π as its reversible and stationary distribution.

Proof: Exercise.

Model 29 (Hard-core model) Let $\mathbb{G} = (\mathbb{V}, \mathbb{E})$ be an undirected graph. An assignment of elements of $\mathbb{S} = \{0, 1\}$ to vertices in \mathbb{V} is called a configuration. Thus, the configuration x is a function $x : \mathbb{V} \rightarrow \mathbb{S}$ and $x \in \mathbb{S}^{\mathbb{V}}$. The vertices v of a configuration x with $x(v) = 1$ are said to be occupied and those with $x(v) = 0$ are said to be vacant. Thus a configuration models a placement of particles on the vertices of \mathbb{V} . A hard-core configuration is a configuration in which no two neighbouring vertices are occupied. More formally, a configuration x is called hard-core if $\sum_{\langle v_i, v_j \rangle \in \mathbb{E}} x(v_i)x(v_j) = 0$. Let the set of hard-core configurations be \mathbb{X} and let π be the uniform distribution on \mathbb{X} , given by

$$\pi(x) = \begin{cases} \frac{1}{\#\mathbb{X}} & \text{if } x \in \mathbb{X} \\ 0 & \text{otherwise} \end{cases}.$$

The Glauber dynamics $(X_t)_{t \in \mathbb{Z}_+}$ for the uniform distribution π on hard-core configurations can be simulated as follows:

- initialize with vacant vertices, i.e., $X_0(w) = 0$ for each $w \in \mathbb{V}$,
- let the current hard-core configuration be $x_t : \mathbb{V} \rightarrow \{0, 1\}$ at time t ,
- choose a vertex v uniformly at random from \mathbb{V} ,
- if any neighbour of v is occupied then v is left vacant, i.e., $x_{t+1}(v) = 0$
- if every neighbour of v is vacant then v is occupied with probability $1/2$, i.e., $x_{t+1}(v) = 1$,
- leave the values at all other vertices unchanged, i.e., $x_{t+1}(w) = x_t(w)$ for each $w \neq v$,
- the possibly modified configuration x_{t+1} is the updated hard-core configuration at time $t + 1$.

Proposition 123 The Glauber dynamics of Model 29 does indeed have π as its stationary distribution.

Proof: First we need to verify that $(X_t)_{t \in \mathbb{Z}_+}$, the Markov chain given by the Glauber dynamics for π in Model 29, is irreducible and aperiodic. Clearly $(X_t)_{t \in \mathbb{Z}_+}$ is aperiodic since we can get from any hard-core configuration $x \in \mathbb{X}$ to itself in one time step by choosing a vertex with at least one occupied neighbour and leaving the chosen vertex unchanged or by choosing a vertex with no occupied neighbours and leaving the chosen vertex unchanged with probability 1/2. Next we need to establish irreducibility, i.e., we need to show that we can get from any hardcore configuration x to any other hardcore configuration x' in finitely many steps. Let the vacant configuration be \tilde{x} , i.e., $\tilde{x}(v) = 0$ for every vertex $v \in \mathbb{V}$. In finitely many steps, we can go from any x to \tilde{x} and from \tilde{x} to x' . If x has $s(x) := \sum_{v \in \mathbb{V}} x(v)$ occupied sites or vertices then we can go to the vacant configuration \tilde{x} with $s(\tilde{x}) = 0$ in $s(x)$ time steps by picking one of the currently occupied sites and making it vacant as follows:

$$P(X_{t+s(x)} = \tilde{x} | X_t = x) = \prod_{i=0}^{s(x)-1} \frac{(s(x) - i)}{\#\mathbb{V}} \frac{1}{2} > 0 .$$

Similarly, we can go from \tilde{x} to any other configuration x' with $s(x')$ many occupied sites in $s(x')$ time steps with the following positive probability:

$$P(X_{t+s(x')} = x' | X_t = \tilde{x}) = \prod_{i=0}^{s(x')-1} \frac{(s(x') - i)}{\#\mathbb{V}} \frac{1}{2} > 0 .$$

Note that this is not the shortest possible number of steps to go from x to x' but just a finite number of steps. Thus we have established that $x \leftrightarrow x'$ for every $(x, x') \in \mathbb{X}$ and thereby established irreducibility of the chain $(X_t)_{t \in \mathbb{Z}_+}$.

If we now show that π is reversible for $(X_t)_{t \in \mathbb{Z}_+}$ then by Proposition 115 π is also stationary for $(X_t)_{t \in \mathbb{Z}_+}$ and finally π is the unique stationary distribution due to irreducibility and aperiodicity. Let $P(x, y)$ be the probability of going from x to y in one time step of $(X_t)_{t \in \mathbb{Z}_+}$. We need to show that for any pair of hardcore configurations $(x, y) \in \mathbb{X}^2$ the following equality holds:

$$\pi(x)P(x, y) = \pi(y)P(y, x), \quad \pi(x) = \frac{1}{\#\mathbb{X}} .$$

Let the number of vertices at which x and y differ be $d(x, y) := \sum_{v \in \mathbb{V}} \text{abs}(x(v) - y(v))$. Let us consider three cases of $(x, y) \in \mathbb{X}^2$.

Case i: When $d(x, y) = 0$ the two configurations are identical, i.e., $x = y$, and therefore we have the trivial equality:

$$\pi(x)P(x, y) = \pi(x)P(x, x) = \pi(y)P(y, x) .$$

Case ii: When $d(x, y) > 1$ the two configurations differ at more than one vertex and therefore $P(x, y) = 0$ and we have the trivial equality:

$$\pi(x)P(x, y) = \pi(x)0 = 0 = \pi(y)P(y, x) .$$

Case iii: When $d(x, y) = 1$ the two configurations differ at exactly one vertex v and therefore all neighbouring vertices of v must be vacant, i.e., take the value 0, in both x and y with $P(x, y) = P(y, x) = \frac{1}{\#\mathbb{V}} \frac{1}{2}$. Thus,

$$\pi(x)P(x, y) = \frac{1}{\#\mathbb{X}} \left(\frac{1}{\#\mathbb{V}} \frac{1}{2} \right) = \pi(y)P(y, x) .$$

We have established that $pi(x) = 1/\#\mathbb{X}$ for each $x \in \mathbb{X}$ is the reversible distribution and thereby also the unique stationarity distribution for $(X_t)_{t \in \mathbb{Z}_+}$, the Markov chain given by the Glauber dynamics for π in Model 29.

Exercise 232 (1-D hardcore model) Let \mathbb{X}_n be the set of hardcore configurations on a path graph with n vertices. Recall that a path graph $\mathbb{G}_n = (\mathbb{V}_n, \mathbb{E}_n)$ has n vertices and $n - 1$ edges, as follows:

$$\mathbb{V}_n = \{v_1, v_2, \dots, v_n\}, \quad \mathbb{E}_n = \{\langle v_1, v_2 \rangle, \langle v_2, v_3 \rangle, \dots, \langle v_{n-1}, v_n \rangle\} .$$

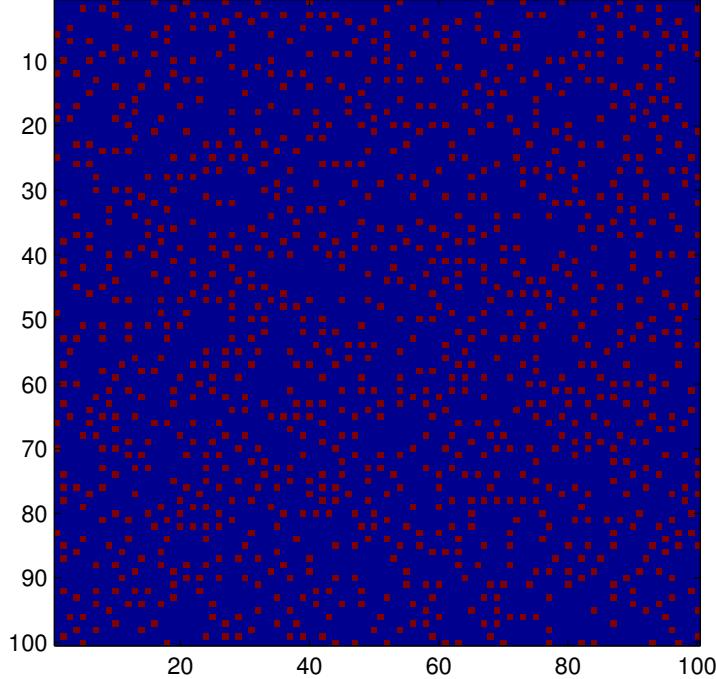
Draw all five hardcore configurations in \mathbb{X}_3 . Show that for any positive integer n ,

$$\#\mathbb{X}_n = \text{fibo}(n + 1) ,$$

the $(n + 1)$ -th Fibonacci number, that is defined recursively as follows:

$$\text{fibo}(0) := \text{fibo}(1) := 1, \quad \text{fibo}(n) = \text{fibo}(n - 1) + \text{fibo}(n - 2), \quad n \geq 1 .$$

Figure 8.6: The sample at time step 10^6 from the Glauber dynamics for the hardcore model on 100×100 regular torus grid. A red site is occupied while a blue site is vacant.



Simulation 233 (Glauber dynamics for the hardcore model on a 2D regular torus)
 Let us implement a program in MATLAB that will simulate Glauber dynamics to sample uniformly from the hardcore configurations on the undirected regular torus graph. We can report the sample mean of the fraction of occupied sites on this graph from the simulated sequence and make a movie of the simulations (last frame is shown in Figure 8.6).

```
>> Hardcore2D
Avg1s = 0.1128
```

The simulation was implemented in the following M-file:

```
HardCore2D.m
% simulation of Glauber dynamics for the hardcore model on
% 2D regular torus grid
clf; %clear; clc; % clear current settings
Seed=347632321; rand('twister',Seed); % set seed for PRNG
MaxSteps=1000000; % number of time steps to simulate
DisplayStepSize=10000; % display interval
Steps=0; % initialize time-step to 0
StepsM=1; % index for movie frame
Rows=100; % number of rows
Cols=100; % number of columns
CC = zeros(Rows,Cols,'int8'); %initialize all sites to be vacant
Delta=[-1,0,+1]; % neighbourhood of indices along one coordinate
Avg1s=0.0;%initialise the Average Fraction of occupied sites
while(Steps <= MaxSteps)
    % find a random site with 0 for possible swap
    I=ceil(Rows*rand); J=ceil(Cols*rand);
    % Get the Nbhd of CC(I,J)
    RowNbhd = mod((I-1)+Delta,Rows)+1;
    ColNbhd = mod((J-1)+Delta,Cols)+1;
    Nbhd=CC(RowNbhd, ColNbhd);
    To1Is=find(Nbhd); % find the 1s in Nbhd of CC(I,J)
    Num1s=length(To1Is); % total number of 1s in Nbhd
    if(Num1s > 0)
        CC(I,J)=0; % set site to be vacant
    elseif(rand < 0.5)
        CC(I,J)=1; % set site to be occupied
    else
        CC(I,J)=0; % set site to be vacant
    end
    Steps=Steps+1; % increment time step
    Frac1s=sum(sum(CC))/(Rows*Cols); % fraction of occupied sites
    Avg1s = Avg1s + (Frac1s - Avg1s)/Steps; % online sample mean
    if(mod(Steps,DisplayStepSize)==0)
        A(StepsM)=getframe; % get the frame into A
        imagesc(CC)
        axis square
        StepsM=StepsM+1;
    end
end
Avg1s % print the sample mean of fraction of occupied sites
movie(A,5) % make a movie
```

Model 30 (Ising model) Let $\mathbb{G} = (\mathbb{V}, \mathbb{E})$ be an undirected graph. The Ising model is a probability distribution on $\mathbb{X} = \{-1, +1\}^{\mathbb{V}}$, i.e., a way of randomly assigning elements from the set $\{-1, +1\}$ to vertices of \mathbb{G} . The physical interpretation of the model is that each vertex is the position of an atom in a ferromagnetic material and $+1$'s or -1 's denote the two possible spin orientations of the atoms. There is a parameter β in the model called inverse temperature and $\beta \in [0, \infty)$. Associated with each spin configuration $x \in \mathbb{X}$ is its energy

$$H(x) = - \sum_{\langle u,v \rangle \in \mathbb{E}} x(u)x(v)$$

where $x(u)$ and $x(v)$ give the spin orientations of the atoms at vertices u and v , respectively. So, each edge $\langle u, v \rangle$ adds 1 to the energy $H(x)$ if its neighbouring vertices have opposite spins

and subtracts 1 from $H(x)$ otherwise. Thus, lower energy is equivalent to a higher agreement in spins between neighbouring vertices.

The Ising model on \mathbb{G} at inverse temperature β means a random spin configuration X with

$$P(X = x) = \pi_{\mathbb{G}, \beta}(x) = \frac{1}{Z_{\mathbb{G}, \beta}} \exp(-\beta H(x)) = \frac{1}{Z_{\mathbb{G}, \beta}} \exp\left(\beta \sum_{\langle u, v \rangle \in \mathbb{E}} x(u)x(v)\right),$$

where $Z_{\mathbb{G}, \beta} = \sum_{x \in \mathbb{X}} \exp(-\beta H(x))$ is the normalising constant.

Labwork 234 (Glauber dynamics for the Ising model on a 2D regular torus) Implement a program in MATLAB to simulate from the Ising model on the undirected regular torus graph.

Let us explore the physical interpretation of the Ising model further. If the inverse temperature $\beta = 0$ then we are at infinite temperature and therefore every configuration in \mathbb{X} is equally likely, i.e., $\pi_{\mathbb{G}, 0} = 1/\#\mathbb{X}$. At the other extreme, if $\beta \rightarrow \infty$ then we are approaching zero temperature and the probability over \mathbb{X} under $\pi_{\mathbb{G}, \infty}$ is equally split between “all +1” configuration and “all -1” configuration. However, if $\beta > 0$, then we are at some temperature $1/\beta$ that is neither absolutely hot or absolutely cold and therefore the model will favour configurations with lower energy as opposed to higher energy. Such favourable low energy configurations tend to have neighbouring clumps of identical spins. We say that there is a phase transition in β since the Ising model’s qualitative behaviour depends on whether β is above or below a critical threshold β_c .

Model 31 (Proper q -colourings) A proper q -colouring of an undirected graph $\mathbb{G} = (\mathbb{V}, \mathbb{E})$ is an assignment of q colours labelled $\{1, 2, \dots, q\}$ to vertices in \mathbb{V} , subject to the constraint that neighbouring vertices do not receive the same colour. Let \mathbb{X} denote the set of all proper q -colourings of \mathbb{G} . If \mathbb{V} is large then \mathbb{X} can be a large and complicated subset of $\{1, 2, \dots, q\}^{\mathbb{V}}$. Note that proper q colourings are a natural generalisation of the hardcore model.

8.7.1 Random Walks on \mathbb{Z} and the reflection principle



8.8 Coupling from the past

MCMC algorithms make it easy to implement a Markov chain that has a given distribution as its stationary distribution. When used on their own, however, MCMC algorithms can only provide sample values that approximate a desired distribution. To obtain sample values that have a desired distribution *exactly* or *perfectly*, MCMC algorithms must be used in conjunction with ideas that make clever use of coupling.

MCMC convergence diagnostics based on *multiple* independent or *coupled* Markov chains running *forward* in time have been suggested, but are not completely reliable. The chains are coupled if the same sequence of random numbers is used to propagate all of them. By adopting a different perspective - running multiple coupled chains from the past or *backward coupling* - Propp & Wilson (1996) developed the *coupling from the past (CFTP)* algorithm, which allowed exact sample values to be obtained from the stationary distribution of an ergodic Markov chain with *finite* state space.

Let us first appreciate the trouble with MCMC algorithms such as Metropolis-Hastings chain, Metropolis chain and Glauber dynamics. Firstly, no matter how large we make time t

to be we cannot avoid the discrepancy between the t -step distribution μ_t and the stationary distribution π . Consider the following transition probability matrix:

$$P = \begin{matrix} & s_1 & s_2 \\ s_1 & \frac{3}{4} & \frac{1}{4} \\ s_2 & \frac{1}{4} & \frac{3}{4} \end{matrix}$$

We can prove by induction that

$$\mu_t = \left(\frac{1}{2} (1 + 2^{-t}), \frac{1}{2} (1 - 2^{-t}) \right)$$

for every $t \in \mathbb{Z}_+$. The stationary distribution is $\pi = (1/2, 1/2)$. So, as t approaches infinity $\mu_t \xrightarrow{\text{TV}} \pi$, however for any t the total variation distance between $d_{\text{TV}}(\mu_t, \pi) = 2^{-t}$ is strictly positive. Even in this simple example μ_t may never equal π for any finite t , however large. Thus, we have to settle for an approximation to π with some acceptable error ϵ . Secondly, to make the approximation error measured by $d_{\text{TV}}(\mu_t, \pi)$ smaller than ϵ we have to find the ϵ -burnin time τ_ϵ by which $d_{\text{TV}}(\mu_{\tau_\epsilon}, \pi) < \epsilon$. Determining τ_ϵ is nontrivial except in special cases and constitutes an active field of research.

The following material is under 

Demonstration 235 (Applet – Perfect sampling.) The CFTP algorithm starts multiple Markov chains, one for each possible state, at some time $t_0 < 0$ in the past, and uses coupled transitions to propagate them to time 0. If all the chains *coalesce*, (i.e. end up having the same state, at or before time 0), then they will have “forgotten” their starting values and will evolve as a single chain from that point onwards. The common state at time zero ($X^{(0)}$) is an exact sample value from the stationary distribution. Intuitively, if coalescence occurs at some finite time, $t^* < 0$, then if the chains had been started in the infinite past, coupling with the same sequence of random numbers will ensure that they coalesce at t^* , and the common chain at time 0 must be stationary because it had been running for an infinitely long time. Thus, the existence of a finite coalescence time can give a stationary sample value in finite time. The use of coupling is essential to induce coalescence in a finite length of time.

Consider a Markov chain with finite state space, $S = 1, 2, \dots, K$. The CFTP algorithm starts K Markov chains, one from each state in S , at some time $t_0 < 0$ in the past. A sequence of t_0 random vectors, $R^{t+1}, R^{t+2}, \dots, R^0$, is generated and used to propagate all K Markov chains to time 0. Let $X^{t,k(t_0)}$ represent the state of the Markov chain at time t , starting from state $k \in S$ at time $t_0 < t$, and let ϕ be the update function of the Markov chain, such that:

$$X^{(t+1,k(t_0))} = \phi(X^{(t,k(t_0))}, R^{(t+1)}) \quad (8.18)$$

8.8.1 Algorithm – Coupling from the past.

Set $t_0 = 0$.

Repeat

- Set $t_0 = t_0 - 1$, (take 1 time-step back)
- Generate $R^{(t_0+1)}$,
- For $k = 1, 2, \dots, K$, (for each state)
 - Set $X^{(t_0,k(t_0))} = k$, (start chain in that state)

For $t = t_0, t_0 + 1, \dots, -1$, (propagate chain to time 0)

$$\text{Set } X^{(t+1,k(t_0))} = \phi(X^{(t,k(t_0))}, R^{(t+1)}).$$

Until $X^{(0,1(t_0))} = X^{(0,2(t_0))} = \Lambda = X^{(0,K(t_0))}$. (check for coalescence at time 0)

Return $X^{(0)}$.

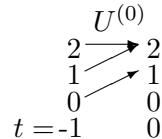
Example 236 Suppose that the Markov chain has the state space, $S = 0, 1, 2$, and a transition matrix:

$$Q = \begin{pmatrix} 0.6 & 0.3 & 0.1 \\ 0.4 & 0.4 & 0.2 \\ 0.3 & 0.4 & 0.3 \end{pmatrix}$$

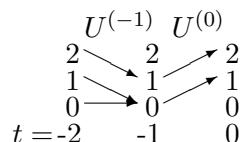
where the (i, j) -element is the conditional probability, $P(X^{(t+1)} = j | X^{(t)} = i)$. The matrix of conditional cumulative probabilities is

$$C = \begin{pmatrix} 0.6 & 0.9 & 1 \\ 0.4 & 0.8 & 1 \\ 0.3 & 0.7 & 1 \end{pmatrix}$$

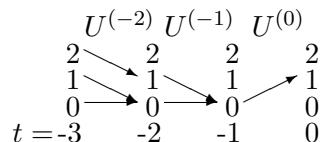
where the (i, j) -element is the probability, $P(X^{(t+1)} = j | X^{(t)} = i)$. Beginning at $t_0 = -1$, three chains are started at 0, 1 and 2. A uniform $(0, 1)$ random number, $U^{(0)}$, is generated (in this example, $R^{(0)} = U^{(0)}$) and used to propagate all three chains to time 0. Suppose that $U^{(0)} \in (0.8, 0.9)$. Then the three chains are updated as shown:



The chains have not coalesced at $t = 0$, so we need to move one time-step back to $t_0 = -2$, start three chains at 0, 1 and 2, generate a second uniform $(0, 1)$ random number, $U^{(-1)}$ and use it along with the previous $U^{(0)}$ to propagate the chains to time 0. Suppose that $U^{(-1)} \in (0.3, 0.4)$. The three chains then evolve as shown:



The chains have still not coalesced at $t = 0$, so we must move another time-step back to $t_0 = -3$ and start again, generating a third uniform $(0, 1)$ random number, $U^{(-2)}$. Suppose that $U^{(-2)} \in (0.3, 0.4)$. This is used with $U^{(-1)}$ and $U^{(0)}$ from before, giving the following transitions:



All three chains have now coalesced at $t = 0$ and so $X^{(0)} = 1$ is accepted as a sample value from the stationary distribution. The whole process is repeated to get another independent sample value. It is important to note that even though the chains have coalesced at $t = 1$, with the common value $X^{(-1)} = 0$; this value at the time of coalescence is not accepted as being from the stationary distribution. This is because the time of coalescence is a random time that depends only on the sequence of random numbers, $U^{(0)}, U^{(-1)}, \dots$; while the time at which a coalesced state has the required stationary distribution must be a fixed time. In the CFTP algorithm, this *fixed* time has been arbitrarily specified to be $t = 0$.

Example 237 To see that the state at the time of coalescence does not have the stationary distribution, suppose that the state space is $S = 1, 2$ and the transition matrix is:

$$Q = \begin{pmatrix} 0.5 & 0.5 \\ 1 & 0 \end{pmatrix}.$$

Since $Q(2, 1) = 1$, the two coupled chains must be in state 1 at the time of coalescence. However, the stationary distribution of this Markov chain is $f(1) = 2/3$ and $f(2) = 1/3$, and so the state at the time of coalescence cannot be from the stationary distribution.

Instead of taking a single step back when the two bounding chains fail to coalesce, any decreasing sequence of time-steps may be used. The “double-until-overshoot” choice of $t_0 = -2^0, -2^1, -2^2, \dots$ is optimal in the sense that it minimises the worst-case number of steps and almost minimises the expected number of steps for coalescence.

Exercise 238 Implement the CFTP algorithm for the Markov chain in Example 2.5.3 and use it to generate 1000 sample points from the stationary distribution of the chain. The stationary distribution can be shown to be:

x	0	1	2
$f(x)$	0.4789	0.3521	0.1690

Compare the relative frequencies of the generated sample with the true stationary probabilities.

Exercise 239 2.6.19 Consider a Markov chain with a state space $S = 0, 1, 2, 3$ and the transition matrix:

$$Q = \begin{pmatrix} 0.6 & 0.4 & 0 & 0 \\ 0.4 & 0.2 & 0.4 & 0 \\ 0.2 & 0.4 & 0 & 0.4 \\ 0 & 0.2 & 0.4 & 0.4 \end{pmatrix}.$$

Let $f = (f_0, f_1, f_2, f_3)$ be a row vector containing the stationary probabilities of the chain.

(a) By solving $fQ = f$ and $f_0 + f_1 + f_2 + f_3 = 1$ simultaneously, show that the stationary distribution of the chain is $f = (14/35, 11/35, 6/35, 4/35)$.

(b) Implement the “double-until-overshoot” version of the CFTP algorithm to generate from the stationary distribution, and use it to obtain 1000 sample points. Compare the relative frequencies of the generated sample with the true stationary probabilities.

8.9 Non-parametric DF Estimation

So far, we have been interested in some estimation problems involved in parametric experiments. In parametric experiments, the parameter space Θ can have many dimensions, but these are finite. For example, in the n IID Bernoulli(θ^*) and the n IID Exponential(λ^*) experiments:

$$\begin{aligned} X_1, \dots, X_n &\stackrel{\text{IID}}{\sim} \text{Bernoulli}(\theta^*), & \theta^* \in \Theta = [0, 1] \subset \mathbb{R}^1, \\ X_1, \dots, X_n &\stackrel{\text{IID}}{\sim} \text{Exponential}(\lambda^*), & \lambda^* \in \Lambda = (0, \infty) \subset \mathbb{R}^1, \end{aligned}$$

the parameter spaces Θ and Λ are of dimension 1. Similarly, in the n IID Normal(μ, σ^2) and the n IID Lognormal(λ, ζ), experiments:

$$\begin{aligned} X_1, \dots, X_n &\stackrel{\text{IID}}{\sim} \text{Normal}(\mu, \sigma^2), & (\mu, \sigma^2) \in \Theta = (-\infty, +\infty) \times (0, +\infty) \subset \mathbb{R}^2 \\ X_1, \dots, X_n &\stackrel{\text{IID}}{\sim} \text{Lognormal}(\lambda, \zeta), & (\lambda, \zeta) \in \Theta = (0, +\infty) \times (0, +\infty) \subset \mathbb{R}^2 \end{aligned}$$

the parameter space is of dimension 2.

An experiment with an infinite dimensional parameter space Θ is said to be **non-parametric**. Next we consider a non-parametric experiment in which n IID samples are drawn according to some fixed and possibly unknown DF F^* from the space of **All Distribution Functions**:

$X_1, X_2, \dots, X_n \stackrel{\text{IID}}{\sim} F^*, \quad F^* \in \Theta = \{\text{All DFs}\} := \{F(x; F) : F \text{ is a DF}\}$

where the DF $F(x; F)$ is indexed or parameterised by itself. Thus, the parameter space $\Theta = \{\text{All DFs}\}$ is the **infinite dimensional** space of **All DFs**. In this section, we look at estimation problems in non-parametric experiments with an infinite dimensional parameter space. That is, we want to estimate the DF F^* from which our IID data are drawn.

The next proposition is often referred to as the **fundamental theorem of statistics** and is at the heart of non-parametric inference, empirical processes, and computationally intensive bootstrap techniques. Recall Definition 56 of the n -sample empirical distribution function (EDF or ECDF) \widehat{F}_n that assigns a probability mass of $1/n$ at each data point x_i :

$$\widehat{F}_n(x) = \frac{\sum_{i=1}^n \mathbf{1}(X_i \leq x)}{n}, \quad \text{where} \quad \mathbf{1}(X_i \leq x) := \begin{cases} 1 & \text{if } x_i \leq x \\ 0 & \text{if } x_i > x \end{cases}$$

Proposition 124 (Gilvenko-Cantelli Theorem) Let $X_1, X_2, \dots, X_n \stackrel{\text{IID}}{\sim} F^*$. Then:

$$\sup_x |\widehat{F}_n(x) - F^*(x)| \xrightarrow{P} 0.$$

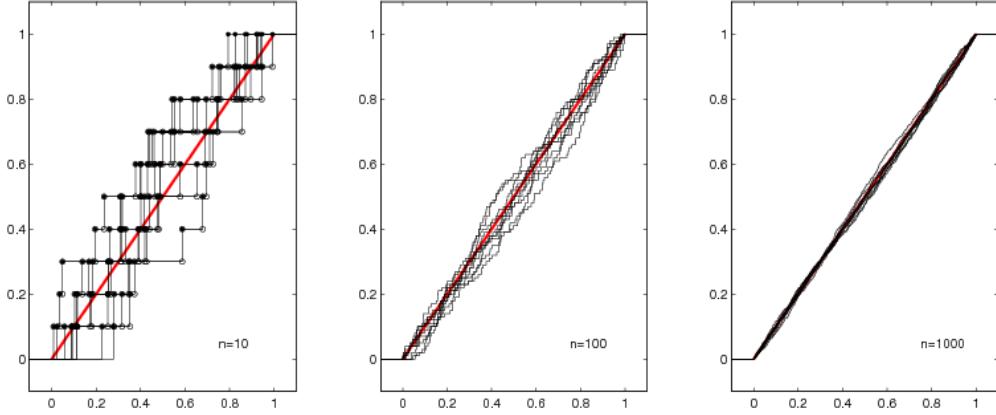
Heuristic Interpretation of the Gilvenko-Cantelli Theorem: As the sample size n increases, the empirical distribution function \widehat{F}_n converges to the true DF F^* in probability, as shown in Figure 8.7.

Proposition 125 (The Dvoretzky-Kiefer-Wolfowitz (DKW) Inequality) Let $X_1, X_2, \dots, X_n \stackrel{\text{IID}}{\sim} F^*$. Then, for any $\epsilon > 0$:

$$P \left(\sup_x |\widehat{F}_n(x) - F^*(x)| > \epsilon \right) \leq 2 \exp(-2n\epsilon^2) \tag{8.19}$$

Recall that $\sup(A)$ or supremum of a set $A \subset \mathbb{R}$ is the least upper bound of every element in A .

Figure 8.7: Plots of ten distinct ECDFs \hat{F}_n based on 10 sets of n IID samples from Uniform(0, 1) RV X , as n increases from 10 to 100 to 1000. The DF $F(x) = x$ over $[0, 1]$ is shown in red. The script of Labwork ?? was used to generate this plot.



8.9.1 Estimating DF

Let $X_1, X_2, \dots, X_n \stackrel{\text{IID}}{\sim} F^*$, where F^* is some particular DF in the space of all possible DFs, i.e. the experiment is non-parametric. Then, based on the data sequence X_1, X_2, \dots, X_n we want to estimate F^* .

For any fixed value of x , the expectation and variance of the empirical DF (4.8) are:

$$E(\hat{F}_n(x)) = F^*(x) \implies \text{bias}_n(\hat{F}_n(x)) = 0 \quad (8.20)$$

$$V(\hat{F}_n(x)) = \frac{F^*(x)(1 - F^*(x))}{n} \implies \lim_{n \rightarrow \infty} \text{se}_n(\hat{F}_n(x)) = 0 \quad (8.21)$$

Therefore, by Proposition 64, the empirical DF evaluated at x , i.e. $\hat{F}_n(x)$ is an asymptotically consistent estimator of the DF evaluated at x , i.e. $F^*(x)$. More formally, (8.20) and (8.21), by Proposition 64, imply that for any fixed value of x :

$$\hat{F}_n(x) \xrightarrow{P} F^*(x).$$

We are interested in a point estimate of the entire DF F^* , i.e. $F^*(x)$ over all x . A point estimator $T_n = T_n(X_1, X_2, \dots, X_n)$ of a fixed and possibly unknown $F \in \{\text{All DFs}\}$ is the empirical DF \hat{F}_n . This estimator has an asymptotically desirable property:

$$\sup_x |\hat{F}_n(x) - F^*(x)| \xrightarrow{P} 0$$

because of the Gilvenko-Cantelli theorem in Proposition 124. Thus, we can simply use \hat{F}_n , based on the realized data (x_1, x_2, \dots, x_n) , as a point estimate of F^* .

On the basis of the DKW inequality (8.19), we can obtain a $1 - \alpha$ confidence set or **confidence band** $C_n(x) := [\underline{C}_n(x), \bar{C}_n(x)]$ about our point estimate of F^* :

$$\begin{aligned} \underline{C}_n(x) &= \max\{\hat{F}_n(x) - \epsilon_n, 0\}, \\ \bar{C}_n(x) &= \min\{\hat{F}_n(x) + \epsilon_n, 1\}, \\ \epsilon_n &= \sqrt{\frac{1}{2n} \log \left(\frac{2}{\alpha} \right)}. \end{aligned} \quad (8.22)$$

It follows from (8.19) that for any fixed and possibly unknown F^* :

$$P(\underline{C}_n(x) \leq F^*(x) \leq \bar{C}_n(x)) \geq 1 - \alpha .$$

Let us look at a simple example next.

Labwork 240 (Estimating the DF of Uniform(0, 1) RV) Consider the problem of estimating the DF of Uniform(0, 1) RV U on the basis of $n=10$ samples. We use the function ECDF of Labwork ?? and MATLAB's built-in function stairs to render the plots. Figure 8.8 was generated by PlotUniformECDFsConfBands.m given below.

```
% script PlotUniformECDFsConfBands.m to plot the ECDF from 10 and 100 samples
% from Uniform(0,1) RV
rand('twister',76534); % initialize the Uniform(0,1) Sampler
N = 3; % 10^N is the maximum number of samples from Uniform(0,1) RV
u = rand(1,10^N); % generate 1000 samples from Uniform(0,1) RV U

% plot the ECDF from the first 10 samples using the function ECDF
for i=1:N
    SampleSize=10^i;
    subplot(1,N,i)
    % Get the x and y coordinates of SampleSize-based ECDF in x1 and y1 and
    % plot the ECDF using the function ECDF
    if (i==1) [x1 y1] = ECDF(u(1:SampleSize),2,0.2,0.2);
    else
        [x1 y1] = ECDF(u(1:SampleSize),0,0.1,0.1);
        stairs(x1,y1,'k');
    end
    % Note PlotFlag is 1 and the plot range of x-axis is
    % incremented by 0.1 or 0.2 on either side due to last 2 parameters to ECDF
    % being 0.1 or 0.2
    Alpha=0.05; % set alpha to 5% for instance
    Epsn = sqrt((1/(2*SampleSize))*log(2/Alpha)); % epsilon_n for the confidence band
    hold on;
    stairs(x1,max(y1-Epsn,zeros(1,length(y1))), 'g'); % lower band plot
    stairs(x1,min(y1+Epsn,ones(1,length(y1))), 'g'); % upper band plot
    axis([-0.1 1.1 -0.1 1.1]);
    axis square;
    x=[0:0.001:1];
    plot(x,x,'r'); % plot the DF of Uniform(0,1) RV in red
    LabelString=['n=' num2str(SampleSize)];
    text(0.75,0.05,LabelString)
    hold off;
end
```

Next we look at a more interesting example involving real-world data.

Labwork 241 (Non-parametric Estimation of the DF of Times Between Earth Quakes) Suppose that the 6,128 observed times between Earth quakes in NZ between 18-Jan-2008 02:23:44 and 18-Aug-2008 19:29:29 are:

$$X_1, \dots, X_{6128} \stackrel{IID}{\sim} F^*, \quad F^* \in \{\text{all DFs}\} .$$

Then the non-parametric point estimate of the unknown F^* is \hat{F}_{6128} , the ECDF of the inter earth quake times. We plot the non-parametric point estimate as well as the 95% confidence bands for F^* in Figure 8.9.

Figure 8.8: The empirical DFs $\hat{F}_n^{(1)}$ from sample size $n = 10, 100, 1000$ (black), is the point estimate of the fixed and known DF $F(x) = x, x \in [0, 1]$ of Uniform(0, 1) RV (red). The 95% confidence band for each \hat{F}_n are depicted by green lines.

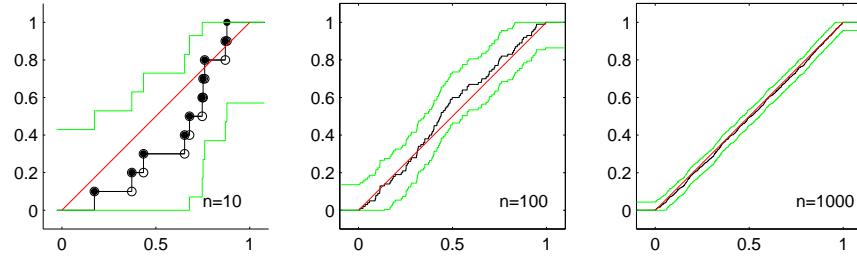
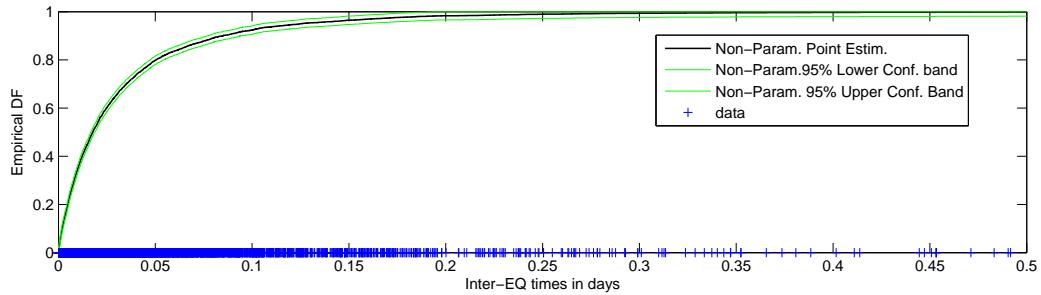


Figure 8.9: The empirical DF \hat{F}_{6128} for the inter earth quake times and the 95% confidence bands for the non-parametric experiment.



```

NZSIEQTimesECDFsConfBands.m
%% The columns in earthquakes.csv file have the following headings
%% CUSP_ID, LAT, LONG, NZMGE, NZMGN, ORI_YEAR, ORI_MONTH, ORI_DAY, ORI_HOUR, ORI_MINUTE, ORI_SECOND, MAG, DEPTH
EQ=dlmread('earthquakes.csv',','); % load the data in the matrix EQ
size(EQ) % report thr size of the matrix EQ
% Read help datenum -- converts time stamps into numbers in units of days
MaxD=max(datenum(EQ(:,6:11)));% maximum datenum
MinD=min(datenum(EQ(:,6:11)));% minimum datenum
% get an array of sorted time stamps of EQ events starting at 0
Times=sort(datenum(EQ(:,6:11))-MinD);
TimeDiff=diff(Times); % inter-EQ times = times between successtive EQ events
n=length(TimeDiff); %sample size
clf % clear any current figures
%% Non-parametric Estimation X_1,X_2,...,X_132 ~ IID F
[x y] = ECDF(TimeDiff,0,0,0); % get the coordinates for empirical DF
stairs(x,y,'k','linewidth',1) % draw the empirical DF
hold on;
% get the 5% non-parametric confidence bands
Alpha=0.05; % set alpha to 5% for instance
Epsn = sqrt((1/(2*n))*log(2/Alpha)); % epsilon_n for the confidence band
stairs(x,max(y-Epsn,zeros(1,length(y))), 'g') % non-parametric 95% lower confidence band
stairs(x,min(y+Epsn,ones(1,length(y))), 'g') % non-parametric 95% upper confidence band
plot(TimeDiff,zeros(1,n),'+')
axis([0 0.5 0 1])
xlabel('Inter-EQ times in days'); ylabel('Empirical DF');
legend('Non-Param. Point Estim.', 'Non-Param. 95% Lower Conf. band', 'Non-Param. 95% Upper Conf. Band', 'data')

```

Recall the poor fit of the Exponential PDF at the MLE for the Orbiter waiting time data. We can attribute the poor fit to coarse resolution of the waiting time measurements in minutes

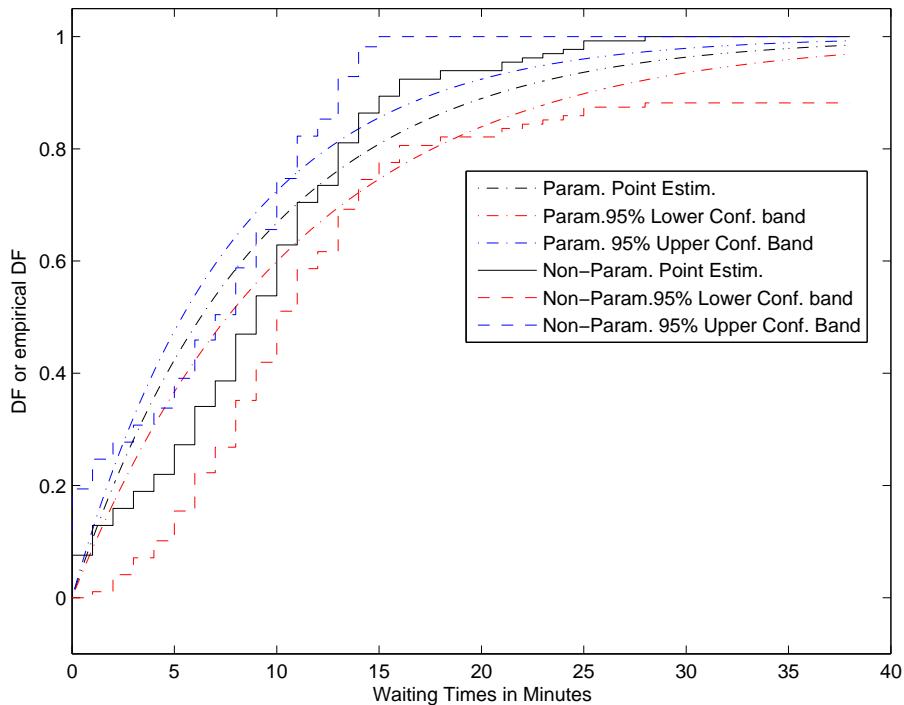
and the rigid decaying form of the exponential PDFs. Let us revisit the Orbiter waiting time problem with our non-parametric estimator.

Labwork 242 (Non-parametric Estimation of Orbiter Waiting Times DF) Suppose that the waiting times at the Orbiter bus stop are:

$$X_1, \dots, X_{132} \stackrel{IID}{\sim} F^*, \quad F^* \in \{\text{all DFs}\}$$

Then the non-parametric point estimate of F^* is \hat{F}_{132} , the ECDF of the 132 Orbiter waiting times. We compute and plot the non-parametric point estimate as well as the 95% confidence

Figure 8.10: The empirical DF \hat{F}_{132} for the Orbiter waiting times and the 95% confidence bands for the non-parametric experiment.



bands for the unknown DF F^* beside the parametric estimate and 95% confidence bands from Labwork 144. Clearly, the non-parametric estimate is preferable to the parametric one for this example. Notice how the non-parametric confidence bands do not contain the parametric estimate of the DF.

```
OrbiterData; % load the Orbiter Data sampleTimes
clf; % clear any current figures
%% Parametric Estimation X_1,X_2,...,X_132 ~ IID Exponential(lambda)
SampleMean = mean(sampleTimes) % sample mean
MLE = 1/SampleMean % ML estimate is 1/mean
n=length(sampleTimes) % sample size
StdErr=1/(sqrt(n)*SampleMean) % Standard Error
MLE95CI=[MLE-(1.96*StdErr), MLE+(1.96*StdErr)] % 95 % CI
TIMES=[0.00001:0.01:max(sampleTimes)+10]; % points on support
plot(TIMES,ExponentialCdf(TIMES,MLE),'k-.'); hold on; % Parametric Point Estimate
```

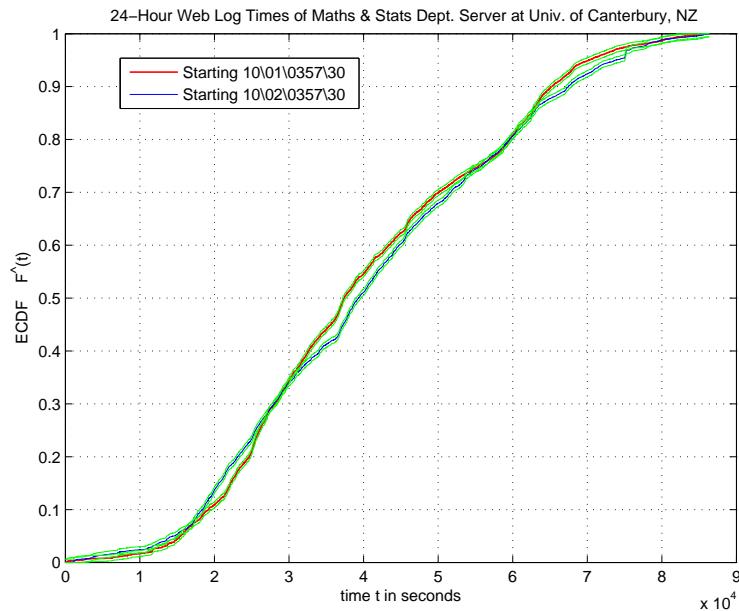
```

plot(TIMES,ExponentialCdf(TIMES,MLE95CI(1)), 'r-.');% Normal-based Parametric 95% lower C.I.
plot(TIMES,ExponentialCdf(TIMES,MLE95CI(2)), 'b-.');% Normal-based Parametric 95% upper C.I.
ylabel('DF or empirical DF'); xlabel('Waiting Times in Minutes');
%% Non-parametric Estimation X_1,X_2,...,X_132 ~ IID F
[x1 y1] = ECDF(sampleTimes,0,0.0,10); stairs(x1,y1,'k');% plot the ECDF
% get the 5% non-parametric confidence bands
Alpha=0.05; % set alpha to 5% for instance
Epsn = sqrt((1/(2*n))*log(2/Alpha)); % epsilon_n for the confidence band
stairs(x1,max(y1-Epsn,zeros(1,length(y1))), 'r--'); % non-parametric 95% lower confidence band
stairs(x1,min(y1+Epsn,ones(1,length(y1))), 'b--'); % non-parametric 95% upper confidence band
axis([0 40 -0.1 1.05]);
legend('Param. Point Estim.', 'Param. 95% Lower Conf. band', 'Param. 95% Upper Conf. Band',...
'Non-Param. Point Estim.', 'Non-Param. 95% Lower Conf. band', 'Non-Param. 95% Upper Conf. Band')

```

Example 243 First take a look at Data ?? to understand how the web login times to our Maths & Stats Department's web server (or requests to our WWW server) were generated. Figure 8.11 shows the login times in units of seconds over a 24 hour period starting at 0357 hours and 30 seconds (just before 4:00AM) on October 1st, 2007 (red line) and on October 2nd, 2007 (magenta). If we assume that some fixed and unknown DF $F^{(1)}$ specifies the distribution of login

Figure 8.11: The empirical DFs $\hat{F}_{n_1}^{(1)}$ with $n_1 = 56485$, for the web log times starting October 1, and $\hat{F}_{n_2}^{(2)}$ with $n_2 = 53966$, for the web log times starting October 2. Their 95% confidence bands are indicated by the green.



times for October 1st data and another DF $F^{(2)}$ for October 2nd data, then the non-parametric point estimates of $F^{(1)}$ and $F^{(2)}$ are simply the empirical DFs $\hat{F}_{n_1}^{(1)}$ with $n_1 = 56485$ and $\hat{F}_{n_2}^{(2)}$ with $n_2 = 53966$, respectively, as depicted in Figure 8.11. See the script of `WebLogDataProc.m` in Data ?? to appreciate how the ECDF plots in Figure 8.11 were made.

8.10 Plug-in Estimators of Statistical Functionals

Recall from Chapter 4 that a **statistical functional** is simply any function of the DF F . For example, the median $T(F) = F^{[-1]}(1/2)$ is a statistical functional. 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. The idea behind the plug-in estimator for a statistical functional is simple: just plug-in the point estimate \hat{F}_n instead of the unknown DF F^* to estimate the statistical functional of interest.

Definition 126 (Plug-in estimator) Suppose, $X_1, \dots, X_n \stackrel{IID}{\sim} F^*$. The plug-in estimator of a statistical functional of interest, namely, $T(F^*)$, is defined by:

$$\hat{T}_n := \hat{T}_n(X_1, \dots, X_n) = T(\hat{F}_n) .$$

Definition 127 (Linear functional) If $T(F) = \int r(x)dF(x)$ for some function $r(x) : \mathbb{X} \rightarrow \mathbb{R}$, then T is called a **linear functional**. Thus, T is linear in its arguments:

$$T(aF + a'F') = aT(F) + a'T(F') .$$

Proposition 128 (Plug-in Estimator of a linear functional) The plug-in estimator for a linear functional $T = \int r(x)dF(x)$ is:

$$T(\hat{F}_n) = \int r(x)d\hat{F}_n(x) = \frac{1}{n} \sum_{i=1}^n r(X_i) .$$

Some specific examples of statistical linear functionals we have already seen include:

1. The **mean** of RV $X \sim F$ is a function of the DF F :

$$T(F) = \mathbb{E}(X) = \int x dF(x) .$$

2. The **variance** of RV $X \sim F$ is a function of the DF F :

$$T(F) = \mathbb{E}(X - \mathbb{E}(X))^2 = \int (x - \mathbb{E}(X))^2 dF(x) .$$

3. The **value of DF at a given $x \in \mathbb{R}$** of RV $X \sim F$ is also a function of DF F :

$$T(F) = F(x) .$$

4. The q^{th} **quantile** of RV $X \sim F$:

$$T(F) = F^{[-1]}(q) \text{ where } q \in [0, 1] .$$

5. The **first quartile** or the 0.25^{th} **quantile** of the RV $X \sim F$:

$$T(F) = F^{[-1]}(0.25) .$$

6. The **median** or the **second quartile** or the 0.50^{th} **quantile** of the RV $X \sim F$:

$$T(F) = F^{[-1]}(0.50) .$$

7. The **third quartile** or the 0.75^{th} **quantile** of the RV $X \sim F$:

$$T(F) = F^{[-1]}(0.75) .$$

Labwork 244 (Plug-in Estimate for Median of Web Login Data) Compute the plug-in estimates for the median for each of the data arrays:

WebLogSeconds20071001035730 and WebLogSeconds20071002035730

that can be loaded into memory by following the commands in the first 13 lines of the script file `WebLogDataProc.m` of Data ??.

Labwork 245 (Plug-in Estimates of Times Between Earth Quakes) Compute the plug-in estimates for the median and mean time in minutes between earth quakes in NZ using the data in `earthquakes.csv`.

```
%>> NZSIEQTimesPlugInEstimates.m
%% The columns in earthquakes.csv file have the following headings
%% CUSP_ID,LAT,LONG,NZMGE,NZMGN,ORI_YEAR,ORI_MONTH,ORI_DAY,ORI_HOUR,ORI_MINUTE,ORI_SECOND,MAG,DEPTH
EQ=dlmread('earthquakes.csv',','); % load the data in the matrix EQ
% Read help datenum -- converts time stamps into numbers in units of days
MaxD=max(datenum(EQ(:,6:11)));% maximum datenum
MinD=min(datenum(EQ(:,6:11)));% minimum datenum
% get an array of sorted time stamps of EQ events starting at 0
Times=sort(datenum(EQ(:,6:11))-MinD);
TimeDiff=diff(Times); % inter-EQ times = times between successive EQ events
n=length(TimeDiff); %sample size
PlugInMedianEstimate=median(TimeDiff) % plug-in estimate of median
PlugInMedianEstimateMinutes=PlugInMedianEstimate*24*60 % median estimate in minutes
PlugInMeanEstimate=mean(TimeDiff) % plug-in estimate of mean
PlugInMeanEstimateMinutes=PlugInMeanEstimate*24*60 % mean estimate in minutes
```

```
>> NZSIEQTimesPlugInEstimates
PlugInMedianEstimate =    0.0177
PlugInMedianEstimateMinutes =   25.5092
PlugInMeanEstimate =    0.0349
PlugInMeanEstimateMinutes =   50.2278
```

Note that any statistical functional can be estimated using the plug-in estimator. However, to produce a $1 - \alpha$ confidence set for the plug-in point estimate, we need bootstrap methods. The subject of next chapter.

8.11 Bootstrap

The **bootstrap** is a statistical method for estimating standard errors and confidence sets of statistics, such as estimators.

8.11.1 Non-parametric Bootstrap for Confidence Sets

Let $T_n := T_n((X_1, X_2, \dots, X_n))$ be a statistic, i.e. any function of the data $X_1, X_2, \dots, X_n \stackrel{\text{IID}}{\sim} F^*$. Suppose we want to know its variance $V_{F^*}(T_n)$, which clearly depends on the fixed and possibly unknown DF F^* .

If our statistic T_n is one with an analytically unknown variance, then we can use the bootstrap to estimate it. The bootstrap idea has the following two basic steps:

Step 1: Estimate $V_{F^*}(T_n)$ with $V_{\hat{F}_n}(T_n)$.

Step 2: Approximate $V_{\hat{F}_n}(T_n)$ using simulated data from the “Bootstrap World.”

For example, if $T_n = \bar{X}_n$, in Step 1, $V_{\hat{F}_n}(T_n) = s_n^2/n$, where $s_n^2 = n^{-1} \sum_{i=1}^n (x_i - \bar{x}_n)^2$ is the sample variance and \bar{x}_n is the sample mean. In this case, Step 1 is enough. However, when the statistic T_n is more complicated (e.g. $T_n = \tilde{X}_n = F^{[-1]}(0.5)$), the sample median, then we may not be able to find a simple expression for $V_{\hat{F}_n}(T_n)$ and may need Step 2 of the bootstrap.

Real World Data come from $F^* \implies X_1, X_2, \dots, X_n \implies T_n((X_1, X_2, \dots, X_n)) = t_n$
 Bootstrap World Data come from $\hat{F}_n \implies X_1^\bullet, X_2^\bullet, \dots, X_n^\bullet \implies T_n((X_1^\bullet, X_2^\bullet, \dots, X_n^\bullet)) = t_n^\bullet$

Observe that drawing an observation from the ECDF \hat{F}_n is equivalent to drawing one point at random from the original data (think of the indices $[n] := \{1, 2, \dots, n\}$ of the original data X_1, X_2, \dots, X_n being drawn according to the equi-probable de Moivre($1/n, 1/n, \dots, 1/n$) RV on $[n]$). Thus, to simulate $X_1^\bullet, X_2^\bullet, \dots, X_n^\bullet$ from \hat{F}_n , it is enough to draw n observations with replacement from X_1, X_2, \dots, X_n .

In summary, the algorithm for Bootstrap Variance Estimation is:

Step 1: Draw $X_1^\bullet, X_2^\bullet, \dots, X_n^\bullet \sim \hat{F}_n$

Step 2: Compute $t_n^\bullet = T_n((X_1^\bullet, X_2^\bullet, \dots, X_n^\bullet))$

Step 3: Repeat Step 1 and Step 2 B times, for some large B , say $B > 1000$, to get $t_{n,1}^\bullet, t_{n,2}^\bullet, \dots, t_{n,B}^\bullet$

Step 4: Several ways of estimating the bootstrap confidence intervals are possible:

(a) The $1 - \alpha$ Normal-based bootstrap confidence interval is:

$$C_n = [T_n - z_{\alpha/2} \hat{s}e_{boot}, T_n + z_{\alpha/2} \hat{s}e_{boot}] ,$$

where the bootstrap-based standard error estimate is:

$$\hat{s}e_{boot} = \sqrt{v_{boot}} = \sqrt{\frac{1}{B} \sum_{b=1}^B \left(t_{n,b}^\bullet - \frac{1}{B} \sum_{r=1}^B t_{n,r}^\bullet \right)^2}$$

(b) The $1 - \alpha$ percentile-based bootstrap confidence interval is:

$$C_n = [\widehat{G}^{\bullet -1}_n(\alpha/2), \widehat{G}^{\bullet -1}_n(1 - \alpha/2)],$$

where \widehat{G}^{\bullet}_n is the empirical DF of the bootstrapped $t_{n,1}^{\bullet}, t_{n,2}^{\bullet}, \dots, t_{n,B}^{\bullet}$ and $\widehat{G}^{\bullet -1}_n(q)$ is the q^{th} sample quantile (4.9) of $t_{n,1}^{\bullet}, t_{n,2}^{\bullet}, \dots, t_{n,B}^{\bullet}$.

Labwork 246 (Confidence Interval for Median Estimate of Inter Earth Quake Times)

Let us find the 95% Normal-based bootstrap confidence interval as well as the 95% percentile-based bootstrap confidence interval for our plug-in estimate of the median of inter earth quake times from Labwork 245 using the following script:

```
%> NZSIEQTimesMedianBootstrap.m
%% The columns in earthquakes.csv file have the following headings
%% CUSP_ID,LAT,LONG,NZMGE,NZMGN,ORI_YEAR,ORI_MONTH,ORI_DAY,ORI_HOUR,ORI_MINUTE,ORI_SECOND,MAG,DEPTH
EQ=dlmread('earthquakes.csv',','); % load the data in the matrix EQ
% Read help datenum -- converts time stamps into numbers in units of days
MaxD=max(datenum(EQ(:,6:11)));% maximum datenum
MinD=min(datenum(EQ(:,6:11)));% minimum datenum
% get an array of sorted time stamps of EQ events starting at 0
Times=sort(datenum(EQ(:,6:11))-MinD);
TimeDiff=diff(Times); % inter-EQ times = times between successive EQ events
n=length(TimeDiff) %sample size
Medianhat=median(TimeDiff)*24*60 % plug-in estimate of median in minutes
B= 1000 % Number of Bootstrap replications
% REPEAT B times: PROCEDURE of sampling n indices uniformly from 1,...,n with replacement
BootstrappedDataSet = TimeDiff([ceil(n*rand(n,B))]);
size(BootstrappedDataSet) % dimension of the BootstrappedDataSet
BootstrappedMedians=median(BootstrappedDataSet)*24*60; % get the statistic in Bootstrap world
% 95% Normal based Confidence Interval
SehatBoot = std(BootstrappedMedians); % std of BootstrappedMedians
% 95% C.I. for median from Normal approximation
ConfInt95BootNormal = [Medianhat-1.96*SehatBoot, Medianhat+1.96*SehatBoot]
% 95% Percentile based Confidence Interval
ConfInt95BootPercentile = ...
    [qthSampleQuantile(0.025,sort(BootstrappedMedians)),...
     qthSampleQuantile(0.975,sort(BootstrappedMedians))]
```

We get the following output when we call the script file.

```
>> NZSIEQTimesMedianBootstrap
n =       6127
Medianhat =   25.5092
B =       1000
ans =      6127      1000
ConfInt95BootNormal =   24.4383   26.5800
ConfInt95BootPercentile =   24.4057   26.4742
```

Labwork 247 (Confidence Interval for Median Estimate of Web Login Data) Find the 95% Normal-based bootstrap confidence interval as well as the 95% percentile-based bootstrap confidence interval for our plug-in estimate of the median for each of the data arrays:

WebLogSeconds20071001035730 and WebLogSeconds20071002035730 .

Once again, the arrays can be loaded into memory by following the commands in the first 13 lines of the script file `WebLogDataProc.m` of Section ???. Produce four intervals (two for each data-set). Do the confidence intervals for the medians for the two days intersect?

```

>> WebLogDataProc % load in the data
>> Medianhat = median(WebLogSeconds20071001035730) % plug-in estimate of median
Medianhat =
      37416
>> % store the length of data array
>> K=length(WebLogSeconds20071001035730)
K =
      56485
>> B= 1000 % Number of Bootstrap replications
B =
      1000
>> BootstrappedDataSet = WebLogSeconds20071001035730([ceil(K*rand(K,B))]);
>> size(BootstrappedDataSet) % dimension of the BootstrappedDataSet
ans =
      56485      1000
>> BootstrappedMedians=median(BootstrappedDataSet); % get the statistic in Bootstrap world
>> % 95% Normal based Confidence Interval
>> SehatBoot = std(BootstrappedMedians); % std of BootstrappedMedians
>> % 95% C.I. for median from Normal approximation
>> ConfInt95BootNormal = [Medianhat-1.96*SehatBoot, Medianhat+1.96*SehatBoot]
ConfInt95BootNormal =
      37242      37590
>> % 95% Percentile based Confidence Interval
ConfInt95BootPercentile =
      ... [qthSampleQuantile(0.025,sort(BootstrappedMedians)),...
      qthSampleQuantile(0.975,sort(BootstrappedMedians))]
ConfInt95BootPercentile =
      37239      37554

```

Labwork 248 (Confidence interval for correlation) Here is a classical data set used by Bradley Efron (the inventor of bootstrap) to illustrate the method. The data are LSAT (Law School Admission Test in the U.S.A.) scores and GPA of fifteen individuals.

Thus, we have bivariate data of the form (Y_i, Z_i) , where $Y_i = \text{LSAT}_i$ and $Z_i = \text{GPA}_i$. For example, the first individual had an LSAT score of $y_1 = 576$ and a GPA of $z_1 = 3.39$ while the fifteenth individual had an LSAT score of $y_{15} = 594$ and a GPA of $z_{15} = 3.96$. We suppose that the bivariate data $(Y_i, Z_i) \stackrel{IID}{\sim} F^*$, such that $F^* \in \{\text{all bivariate DFs}\}$. This is a bivariate non-parametric experiment. The bivariate data are plotted in Figure .

The law school is interested in the correlation between the GPA and LSAT scores:

$$\theta^* = \frac{\int \int (y - \mathbb{E}(Y))(z - \mathbb{E}(Z)) dF(y, z)}{\sqrt{\int (y - \mathbb{E}(Y))^2 dF(y) \int (z - \mathbb{E}(Z))^2 dF(z)}}$$

The plug-in estimate of the population correlation θ^* is the sample correlation:

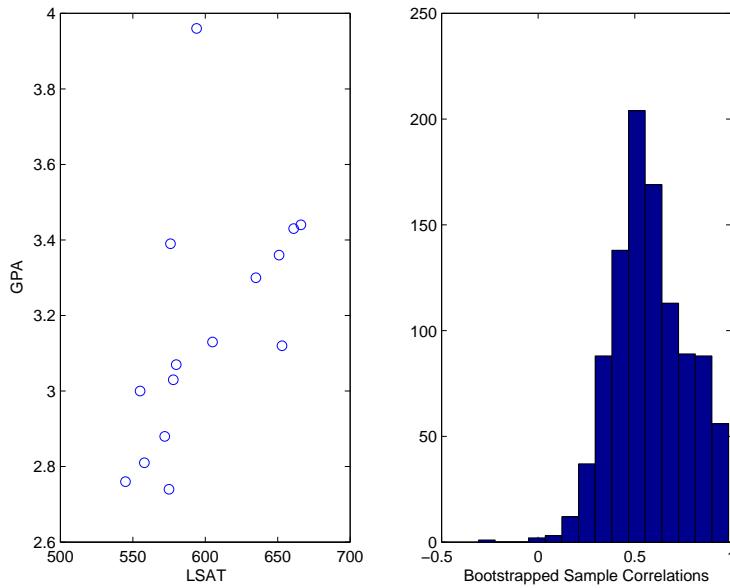
$$\widehat{\Theta}_n = \frac{\sum_{i=1}^n (Y_i - \bar{Y}_n)(Z_i - \bar{Z}_n)}{\sqrt{\sum_{i=1}^n (Y_i - \bar{Y}_n)^2 \sum_{i=1}^n (Z_i - \bar{Z}_n)^2}}$$

```

%% Data from Bradley Efron's LSAT,GPA correlation estimation
LSAT=[576 635 558 578 666 580 555 661 651 605 653 575 545 572 594]; % LSAT data
GPA=[3.39 3.30 2.81 3.03 3.44 3.07 3.00 3.43 3.36 3.13 3.12 2.74 2.76 2.88 3.96]; % GPA data
subplot(1,2,1); plot(LSAT,GPA, 'o'); xlabel('LSAT'); ylabel('GPA') % make a plot of the data
CC=corrcoef(LSAT,GPA); % use built-in function to compute sample correlation coefficient matrix
SampleCorrelation=CC(1,2) % plug-in estimate of the correlation coefficient
%% Bootstrap
B = 1000; % Number of Bootstrap replications
BootstrappedCCs=zeros(1,B); % initialise a vector of zeros
N = length(LSAT); % sample size
rand('twister',767671); % initialise the fundamental sampler
for b=1:B
    Indices=ceil(N*rand(N,1)); % uniformly sample random indices from 1 to 15 with replacement
    BootstrappedLSAT = LSAT([Indices]); % bootstrapped LSAT data
    BootstrappedGPA = GPA([Indices]); % bootstrapped GPA data

```

Figure 8.12: Data from Bradley Efrons LSAT and GPA scores for fifteen individuals (left). The confidence interval of the sample correlation, the plug-in estimate of the population correlation, is obtained from the sample correlation of one thousand bootstrapped data sets (right).



```

CCB=corrcoef(BootstrappedLSAT,BootstrappedGPA);
BootstrappedCCs(b)=CCB(1,2); % sample correlation of bootstrapped data
end
%plot the histogram of Bootstrapped Sample Correlations with 15 bins
subplot(1,2,2);hist(BootstrappedCCs,15);xlabel('Bootstrapped Sample Correlations')

% 95% Normal based Confidence Interval
SehatBoot = std(BootstrappedCCs); % std of BootstrappedMedians
% 95% C.I. for median from Normal approximation
ConfInt95BootNormal = [SampleCorrelation-1.96*SehatBoot, SampleCorrelation+1.96*SehatBoot]
% 95% Percentile based Confidence Interval
ConfInt95BootPercentile = ...
[qthSampleQuantile(0.025,sort(BootstrappedCCs)),...
qthSampleQuantile(0.975,sort(BootstrappedCCs))]
```

We get the following output when we call the script file.

```

>> LSATGPACorrBootstrap
SampleCorrelation =      0.5459
ConfInt95BootNormal =    0.1770    0.9148
ConfInt95BootPercentile =  0.2346    0.9296
```

8.11.2 Parametric Bootstrap for Confidence Sets

The **bootstrap** may also be employed for estimating standard errors and confidence sets of statistics, such as estimators, even in a parametric setting. This is much easier than the variance calculation based on Fisher Information and/or the Delta method.

The only difference in the **parametric bootstrap** as opposed to the **non-parametric bootstrap** we saw earlier is that our statistic of interest $T_n := T_n((X_1, X_2, \dots, X_n))$ is a function of the data:

$$X_1, X_2, \dots, X_n \stackrel{\text{IID}}{\sim} F(x; \theta^*) .$$

That is, our data come from a parametric distribution $F(x; \theta^*)$ and we want to know the variance of our statistic T_n , i.e. $V_{\theta^*}(T_n)$.

The parametric bootstrap concept has the following two basic steps:

Step 1: Estimate $V_{\theta^*}(T_n)$ with $V_{\hat{\theta}_n}(T_n)$, where $\hat{\theta}_n$ is an estimate of θ^* based on maximum likelihood or the method of moments.

Step 2: Approximate $V_{\hat{\theta}_n}(T_n)$ using simulated data from the “Bootstrap World.”

For example, if $T_n = \bar{X}_n$, the sample mean, then in **Step 1**, $V_{\hat{\theta}_n}(T_n) = n^{-1} \sum_{i=1}^n (x_i - \bar{x}_n)$ is the sample variance. Thus, in this case, **Step 1** is enough. However, when the statistic T_n is more complicated, say $T_n = \tilde{X}_n = F^{[-1]}(0.5)$, the sample median, then we may not be able to write down a simple expression for $V_{\hat{\theta}_n}(T_n)$ and may need **Step 2** of the bootstrap.

$$\begin{aligned} \text{Real World Data come from } F(\theta^*) &\implies X_1, X_2, \dots, X_n \implies T_n((X_1, X_2, \dots, X_n)) = t_n \\ \text{Bootstrap World Data come from } F(\hat{\theta}_n) &\implies X_1^\bullet, X_2^\bullet, \dots, X_n^\bullet \implies T_n((X_1^\bullet, X_2^\bullet, \dots, X_n^\bullet)) = t_n^\bullet \end{aligned}$$

To simulate $X_1^\bullet, X_2^\bullet, \dots, X_n^\bullet$ from $F(\hat{\theta}_n)$, we must have a simulation algorithm that allows us to draw IID samples from $F(\theta)$, for instance the inversion sampler. In summary, the algorithm for Bootstrap Variance Estimation is:

Step 1: Draw $X_1^\bullet, X_2^\bullet, \dots, X_n^\bullet \sim F(\hat{\theta}_n)$

Step 2: Compute $t_n^\bullet = T_n((X_1^\bullet, X_2^\bullet, \dots, X_n^\bullet))$

Step 3: Repeat **Step 1** and **Step 2** B times, for some large B , say $B \geq 1000$, to get $t_{n,1}^\bullet, t_{n,2}^\bullet, \dots, t_{n,B}^\bullet$

Step 4: We can estimate the bootstrap confidence intervals in several ways:

(a) The $1 - \alpha$ normal-based bootstrap confidence interval is:

$$C_n = [T_n - z_{\alpha/2} \hat{s}e_{boot}, T_n + z_{\alpha/2} \hat{s}e_{boot}] ,$$

where the bootstrap-based standard error estimate is:

$$\hat{s}e_{boot} = \sqrt{v_{boot}} = \sqrt{\frac{1}{B} \sum_{b=1}^B \left(t_{n,b}^\bullet - \frac{1}{B} \sum_{r=1}^B t_{n,r}^\bullet \right)^2}$$

(b) The $1 - \alpha$ percentile-based bootstrap confidence interval:

$$C_n = [\widehat{G}_n^\bullet^{-1}(\alpha/2), \widehat{G}_n^\bullet^{-1}(1 - \alpha/2)],$$

where \widehat{G}_n^\bullet is the empirical DF of the bootstrapped $t_{n,1}^\bullet, t_{n,2}^\bullet, \dots, t_{n,B}^\bullet$ and $\widehat{G}_n^\bullet(q)$ is the q^{th} sample quantile (4.9) of $t_{n,1}^\bullet, t_{n,2}^\bullet, \dots, t_{n,B}^\bullet$.

Let us apply the bootstrap method to the previous problem of estimating the standard error of the coefficient of variation from $n = 100$ samples from $\text{Normal}(100, 10^2)$ RV. The confidence intervals from bootstrap-based methods are similar to those from the Delta method.

```

n=100; Mustar=100; Sigmastar=10; % sample size, true mean and standard deviation
rand('twister',67345);
x=arrayfun(@(u)(Sample1NormalByNewRap(u,Mustar,Sigmatstar^2)),rand(n,1)); % normal samples
Muhat=mean(x) Sigmahat=std(x) Psihat=Sigmahat/Muhat % MLE of Mustar, Sigmastar and Psistar
Sehat = sqrt((1/Muhat^4)+(Sigmahat^2/(2*Muhat^2))/sqrt(n)) % standard error estimate
% 95% Confidence interval by Delta Method
ConfInt95DeltaMethod=[Psihat-1.96*Sehat, Psihat+1.96*Sehat] % 1.96 since 1-alpha=0.95
B = 1000; % B is number of bootstrap replications
% Step 1: draw n IID samples in Bootstrap World from Normal(Muhat,Sigmahat^2)
xBoot = arrayfun(@(u)(Sample1NormalByNewRap(u,Muhat,Sigmahat^2)),rand(n,B));
% Step 2: % Compute Bootstrapped Statistic Psihat
PsihatBoot = std(xBoot) ./ mean(xBoot);
% 95% Normal based Confidence Interval
SehatBoot = std(PsihatBoot); % std of PsihatBoot
ConfInt95BootNormal = [Psihat-1.96*SehatBoot, Psihat+1.96*SehatBoot] % 1-alpha=0.95
% 95% Percentile based Confidence Interval
ConfInt95BootPercentile = ...
[qthSampleQuantile(0.025,sort(PsihatBoot)),qthSampleQuantile(0.975,sort(PsihatBoot))]

```

```

>> CoeffOfVarNormal
Muhat = 100.3117
Sigmahat = 10.9800
Psihat = 0.1095
Sehat = 0.0077
ConfInt95DeltaMethod = 0.0943 0.1246
ConfInt95BootNormal = 0.0943 0.1246
ConfInt95BootPercentile = 0.0946 0.1249

```