# STAT330S18 - Mathematical Statistics

CLASSNOTES FOR SPRING 2018

by

Johnson Ng

BMath (Hons), Pure Mathematics major, Actuarial Science Minor University of Waterloo

# **Table of Contents**

	List	of Defi	nitions	7
	List	of Theo	orems	11
1	Lect	ure 1 N	Лау 1st 2018	17
	1.1		luction	, 17
	1.2		om Variable	-, 22
	1.3		ete Random Variable	 24
2	Lect	ure 2 N	May 03rd 2018	27
	2.1		nuous Random Variable	27
	2.2	Examı	ples of Discrete RVs	28
		2.2.1	Binomial Distribution	29
		2.2.2	Geometric Distribution	29
		2.2.3	Poisson Distribution	30
	2.3	Examp	ples of Continuous RVs	30
		2.3.1	Normal/Gaussian Distribution	30
		2.3.2	Uniform Distribution	31
		2.3.3	Exponential Distribution	31
		2.3.4	Gamma Distribution	32
	2.4	Functi	ions of Random Variables	32
		2.4.1	Discrete X and Discrete Y	33
		2.4.2	Continuous $X$ and Discrete $Y$	34
		2.4.3	Continuous $X$ and Continuous $Y$	35
		2.4.4	A Formula for the Continuous Case	36
3	Lect	ure 3 N	Aay 08th 2018	39
	3.1	Functi	ions of Random Variables (Continued)	39
		3.1.1	Special Cases	39
	3.2	Probal	bility Integral Transformation	40
	3.3	Locati	on-Scale Families	41
	3.4	Expec	tations	44

### 4 ≡ TABLE OF CONTENTS - ≡ TABLE OF CONTENTS

		3.4.1	Expectations	44
4	Lect	ure 4 N	1ay 10th 2018	47
	4.1	Expect	tations (Continued)	47
		4.1.1	Expectations (Continued)	47
		4.1.2	Moments and Variance	50
	4.2	Inequa	alities	52
		4.2.1	Markov/Chebyshev Style Inequalities	52
5	Lect	ure 5 M	1ay 15th 2018	55
	5.1	Inequa	alities (Continued)	55
		5.1.1	Markov/Chebyshev Style Inequalities (Continued)	55
	5.2	Mome	nt Generating Function	56
		5.2.1	MGF of a Linear Transformation	60
		5.2.2	Uniqueness of the MGF	61
6	Lect	ure 6 M	1ay 17th 2018	63
	6.1	Joint I	Distributions	63
		6.1.1	Introduction to Joint Distributions	63
		6.1.2	Joint and Marginal CDFs	63
		6.1.3	Joint Discrete RVs	67
		6.1.4	Independence of Discrete RVs	70
7	Lect	ure 7 M	Iay 24th 2018	73
	7.1	Joint I	Distributions (Continued)	73
		7.1.1	Independence of Discrete RVs (Continued)	73
		7.1.2	Joint Continuous RVs	73
		7.1.3	Marginal Distribution (Continuous)	77
		7.1.4	Independence of Continuous RVs	78
8	Lect	ure 8 M	Iay 29th 2018	81
	8.1	Joint I	Distributions (Continued 2)	81
		8.1.1	Independence of Continuous RVs (Continued)	81
		8.1.2	Conditional Distributions	83
		8.1.3	Joint Expectations	89
9	Lect	ure 9 N	lay 31st 2018	91
	9.1	Joint I	Distributions (Continued 3)	91
		9.1.1	Joint Expectations (Continued)	91
		9.1.2	Covariance	92
		9.1.3	Correlation	98
		014	Conditional Expectation	100

10	Lect	ure 10 J	Jun 05th 2018	103
	10.1	Joint D	Distribution (Continued 4)	103
		10.1.1	Conditional Expectation (Continued)	103
		10.1.2	Joint Moment Generating Functions	107
11	Lect		Jun 07th 2018	111
		11.0.1	Working with Multivariate Cases	111
12	Lect	ure 12 J	Jun 12th 2018	117
	12.1	Functi	ons of Random Variables	117
		12.1.1	Transformation of Two or More Random Variables	117
		12.1.2	One-to-One Bivariate Transformations	120
13	Lect	ure 13 J	Jun 14th 2018	125
	13.1	Functi	ons of Random Variables (Continued)	125
		13.1.1	One to One Bivariate Transformations (Continued)	125
		13.1.2	Moment Generating Function Method	126
14	Lect	ure 14 J	Jun 19th 2018	131
	14.1		ons of Random Variables (Continued 2)	
		14.1.1	Moment Generating Function Method (Continued)	131
15	Lect	ure 15 J	Jun 21st 2018	135
	15.1	Limiti	ng or Asymptotic Distributions	135
		15.1.1	Convergence in Distribution	135
		15.1.2	Convergence in Probability	139
16	Lect	ure 16 J	Jun 26th 2018	141
	16.1		ng or Asymptotic Distributions (Continued)	
		16.1.1	Convegence in Probability (Continued)	141
		16.1.2	Limit Theorems	146
17			Jun 28th 2018	157
	17.1	Estima	ation	157
		17.1.1	Maximum Likelihood Estimation	158
18	Lect	ure 18 J	Jul 3rd 2018	163
	18.1		ation (Continued)	
		18.1.1	Maximum Likelihood Estimation (Continued)	163
19	Lect	ure 19 J	July 5th 2018	171
	19.1	Estima	ation (Continued 2)	171
		19.1.1	Maximum Likelihood Estimation (Continued 2) .	171

20	Decture 20 Jul 10th 2018 179			179
	20.1	Estima	ation (Continued 3)	179
		20.1.1	$\label{lem:maximum Likelihood Estimation (Continued 3)} \ .$	179
		20.1.2	Asymptotic Properties of ML Estimators	182
21	Lect	ure 21 J	[ul 12th 2018	187
	21.1	Estima	tion (Continued 4)	187
		21.1.1	Asymptotic Properties of ML Estimators (Contin-	
			ued)	187
		21.1.2	Interval Estimators	188
22	Lect	ure 22 J	[ul 17th 2018	195
	22.1	Estima	te (Continued 5)	195
		22.1.1	Interval Estimators (Continued)	195
23	Lect	ure 23 J	[ul 19th 2018	201
	23.1	Hypot	hesis Testing	201
		23.1.1	Introduction	201
		23.1.2	Likelihood Ratio Tests for Simple Hypothesis	204
A	Reg	ularity	Conditions	211
В	Usei	ful Refe	erences	213
	В.1	Comm	only Used Distributions	213
Inc	lex			217
Bil	oliog	raphy		219

# List of Definitions

1	Sample Space	17
2	$\sigma$ -field	17
3	Measurable Space	18
4	Probability Measure	18
5	Conditional Probability	21
6	Independent Events	21
7	Random Variable	22
8	Cumulative Distribution Function	23
9	Discrete Random Variable	24
10	Continuous Random Variable	27
11	Binomial RV	29
12	Geometric RV	29
13	Poisson RV	30
14	Normal / Gaussian RV	30
15	Standard Normal Distribution	31
16	Uniform RV	31
17	Exponential RV	31
18	Gamma RV	32
19	Location Parameter and Family	42
20	Scale Parameter and Family	42
21	Location-Scale Family	12

22	Expectation of A Discrete RV
23	Expectation of A Continuous RV 45
24	Variance
25	Moments
26	Moment Generating Function
27	Indicator Function
28	Joint CDF
29	Marginal CDF
30	Joint Discrete RV 67
31	Marginal Distribution - Discrete Case 68
32	Independence of Discrete RVs
33	Joint Continuous RVs
34	Marginal PDF
35	Independence of Continuous RVs
36	Conditional Distributions 83
37	Beta Distribution
38	Joint Expectation
39	Covariance
40	Correlation Coefficient
41	Conditional Expectation
42	Joint Moment Generating Functions 107
43	Joint Moments and Marginal MGF 108
44	k-variate CDF112
45	k-variate Support Set
46	k-Variate Marginal Distribution
47	k-Variate Joint MGF
<sub>4</sub> 8	Mutlinomial Distribution
49	Bivariate Normal Distribution

50	Random Sample (IID)
51	One-to-One Bivariate Transformation
52	Convergence in Distribution
53	Degenerate Distribution
54	Convergence in Probability
55	Convergence in Probability to a Constant 143
56	Double Parameter Exponential Distribution 144
57	Statistic
58	Estimators and Estimates
59	Likelihood function
60	Score Function
51	Information Function
52	Fisher Information
53	Relative Likelihood
54	Likelihood Region & Likelihood Interval 172
55	Log Relative Likelihood
56	Score Vector
67	Information Matrix
58	Fisher Information Matrix
59	Likelihood Region - Multivariate
70	Interval Estimators
71	Pivotal Quantity
72	Confidence Interval
73	Likelihood Interval
<sup>7</sup> 4	Hypothesis
<sup>7</sup> 5	Null Hypothesis and Alternative Hypothesis 201
76	Test of Hypothesis 202
77	Test Statistic

#### 10 ≡ TABLE OF CONTENTS - ≡ TABLE OF CONTENTS

78	Significance Level and $p$ -Values 202
79	Type I and Type II Errors 203
80	Simple Hypothesis
81	Likelihood Ratio Statistic 204

# **S** List of Theorems

• Proposition 1	Properties of Probability Set Functions	19
• Proposition 2	Boole's Inequality	21
• Proposition 3	Bonferroni's Inequality	21
• Proposition 4	Continuity Property	22
• Proposition 5	Properties of the cdf	23
• Proposition 6	Properties of pmf	24
• Proposition 7	Properties of pdf	27
■ Theorem 8	One-to-One Transformation of a Random Vari	
able		36
■ Theorem 9	Probability Integral Transformation	40
■ Theorem 10	Converse of Probability Integral Transforma-	
tion		41
■ Theorem 11	Expectation from the cdf	47
■ Theorem 12	Expected Value of a Function of X	48
■ Theorem 13	Linearity of Expectation	49
■ Theorem 14	Variance of a Linear Function	51
Theorem 15	Markov's Inequality	52
■ Theorem 16	Markov's Inequality 2	53
■ Theorem 17	Chebyshev's Inequality	54
• Proposition 18	Properties of the MGF	58
■ Theorem 19	MGF of a Linear Transformation	60

Theorem 20	Uniqueness of the MGF	61
• Proposition 21	Properties of Joint CDF	64
♦ Proposition 22	Properties of Joint PMF	67
■ Theorem 23	Independence by PF	71
♦ Proposition 25	Properties of Joint PDF	74
■ Theorem 26	Factorization Theorem for Independence	82
• Proposition 27	Properties of Conditional Distributions	86
■ Theorem 28	Product Rule	86
♦ Proposition 29	Independence from Conditionality	87
Theorem 30	Linearity of Expectation in Bivariate Case .	91
■ Theorem 31 tation	Implication of Independence on Joint Expec-	91
■ Theorem 32  Joint Expectation	Generalized Implication of Independence on	92
Theorem 33	Variance of Linear Combinations	96
■ Theorem 34	Generalized Variance of Linear Combinations	s 97
♦ Proposition 35	Properties of the Correlation Coefficient	99
♦ Proposition 36	Independence on Conditional Expectation .	104
■ Theorem 37	Law of Total Expectation	104
■ Theorem 38	Law of Total Variance	106
♦ Proposition 39	Independence on Joint MGF	109
• Proposition 40	Law of Total Probability - Multivariate	111
• Proposition 41	Independence for Multivariate Cases	112
♦ Proposition 42	Properties of Multinomial Distribution	113
• Proposition 43	Properties of Bivariate Normal Distribution	115
■ Theorem 44	One-to-One Bivariate Transformations	121
• Proposition 45	Properties of the Jacobian	122
Theorem 46	Sums of MGF	126

Theorem 47	Gaussian Distribution 131
Theorem 48	Properties of the Gaussian Distribution 131
Theorem 49	F Distribution
Theorem 50	Taylor Series with Lagrange's Remainder . 135
Theorem 51	Generalized Limit Definition of $e$ 136
Corollary 52	Limit definition of $e$
Lemma 53	
Proposition 54 gence in Distrik	Convergence in Probability Implies Conver- oution
Proposition 55	Partial Converse of 6 Proposition 54 144
Proposition 56	Convergence in Distribution and MGF 146
■ Theorem 57	Central Limit Theorem
Proposition 58	Other Limit Theorems
Theorem 59	Generalized $\delta$ -Method 154
Corollary 60	$\delta$ -Method
Theorem 61	Invariance Property of the MLE 169
Theorem 62	Asymptotic Distribution of the ML Estimator 182
Proposition 63 otal Quantity .	MLE of a Location/Scale Parameter as a Piv-
Proposition 64	Asymptotic Confidence Intervals 196

## Foreword

The proofs in this set of notes will be more rigourous compared to the expectations of the course (at least, for the course this term). If you are not the author and is interested in reading the notes, you may skip the proofs should you have little interest in them. The rigour is required almost exclusively for the author himself, for his own practice, and because he transferred his STAT230 course from a class that is clean of proofs.

Also, many of the common mathematical notations will be heavily used both in the author's notes and proofs. The author cannot guarantee that his proofs are absolutely correct, but he tries, to the best of his abilities to minimize and assure of that. Should you be suspicious about the proofs, or should you notice errorneous ones, please do inform the author at https://github.com/japorized/TeX\_notes/issues.

You are also warned that the author is rather bummed with how the course is presented in the term that he is/was taking it, and so there may be sarcastic language (towards the lectures) mixed in his notes.

# 1 Lecture 1 May 1st 2018

#### 1.1 Introduction

#### Definition 1 (Sample Space)

A sample space, S of a random experiment is the set of all possible outcomes of the experiment.

#### Example 1.1.1

The following are some random experiments and their sample space.

- Flipping a coin  $S = \{H, T\}$  where H denotes head and T tail.
- Rolling a 6-faced dice twice

$$S = \{(x,y) : x,y \in \mathbb{N}, \ 1 \le x,y \le 6\}$$

• Measuring a patient's height

$$S = R^+ = \{x \in \mathbb{R} : x \ge 0\}$$

#### $\blacksquare$ Definition 2 ( $\sigma$ -field)

Let S be a sample space. The collection of sets  $\mathscr{B} \subseteq \mathbb{P}(S)^1$ , is called a  $\sigma$ -field (or  $\sigma$ -algebra) on S if:

1. 
$$\emptyset \in \mathcal{B}$$
 and  $S \in \mathcal{B}$ ;

2. 
$$\forall A \in \mathcal{B}$$
  $A^{C} \in \mathcal{B}$ ; <sup>2</sup> and

3. 
$$\forall n \in \mathbb{N} \quad \forall \{A_j\}_{j=1}^n \subseteq \mathscr{B} \quad \cup_{j=1}^n A_j \in \mathscr{B}.$$

<sup>1</sup> The **power set** of S,  $\mathbb{P}(S)$ , is defined as the set that contains all subsets of S.

 $^2$  We shall denote the compliment of a set by a superscript C in this set of notes. The supplemental notes provided in the class uses an overhead bar, e.g.  $\overline{A}$ , while lecture notes will use  $A^C$  and A' interchangably.

#### Definition 3 (Measurable Space)

Given that S is a non-empty set, and  $\mathcal{B}$  is a  $\sigma$ -field,  $(S,\mathcal{B})$  is a **measurable space**.<sup>3</sup>

<sup>3</sup> A measurable space is a basic object in measure theory.

#### Example 1.1.2

Consider  $S = \{1, 2, 3, 4\}$ . Check if  $\mathcal{B} = \{\emptyset, \{1, 2, 3, 4\}, \{1, 2\}, \{3, 4\}\}$  is a  $\sigma$ -field on S.

- 1. It is clear that  $\emptyset$ ,  $S \in \mathcal{B}$ .
- 2. Note that  $S^C = \emptyset$  and  $\{1,2\}^C = \{3,4\}$ .
- 3. Note that the largest possible result of any countable union of the elements of  $\mathcal{B}$  is  $\{1, 2, 3, 4\}$ , which is an element of  $\mathcal{B}$ .

BECAUSE  $(S, \mathcal{B})$  is a measurable space, we can define a measure on it.

#### Definition 4 (Probability Measure)

Suppose S is a sample space of a random experiment. Let  $\mathscr{B} = \{A_1, A_2, ...\} \subseteq \mathbb{P}(S)$  be the  $\sigma$ -field on S. The **probability set function** (or **probability measure**),  $P: \mathscr{B} \to [0,1]$ , is a function that satisfies the following:<sup>4</sup>

<sup>4</sup> These conditions are also known as Kolmogorov Axioms, or probability axioms.

- $\forall A \in \mathcal{B} \ P(A) > 0$ ;
- P(S) = 1;
- $\forall \{A_j\}_{j=1}^{\infty} \subseteq \mathscr{B} \ \forall i \neq j \in \mathbb{N} \ A_i \cap A_j = \emptyset \implies$

$$P\left(\bigcup_{j=1}^{\infty} A_j\right) = \sum_{j=1}^{\infty} P(A_j) \tag{1.1}$$

 $(S, \mathcal{B}, P)$  is called a probability space.

#### Example 1.1.3

Consider flipping a coin where  $S = \{H, T\}$ . Let P be defined as follows

$$P({H}) = \frac{1}{3} \quad P({T}) = \frac{2}{3} \quad P(\emptyset) = 0 \quad P(S) = 1$$

Conditions 1 and 2 of 🗐 Definition 4 are met. Notice that

$$P({H} \cup {T}) = P(S) = 1$$
 and  $P({H}) + P({T}) = \frac{1}{3} + \frac{2}{3} = 1$ .

Hence condition 3 is also fulfilled.

#### • Proposition 1 (Properties of Probability Set Functions)

Let P be a probability set function and A, B be any set in  $\mathcal{B}$ . Prove the following:5

1. 
$$P(A^C) = 1 - P(A)$$

2. 
$$P(\emptyset) = 0$$

3. 
$$P(A) \leq 1$$

4. 
$$P(A \cap B^C) = P(A) - P(A \cap B)$$

5. 
$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

6. 
$$A \subseteq B \implies P(A) \le P(B)$$

<sup>5</sup> Many among these properties illustrate that the probability is indeed a measure.

### Exercise 1.1.1

#### Proof

*Let S be the sample space for P.* 

#### 1. Note that

$$A \in \mathcal{B} \implies A \in \mathbb{P}(S) \iff A \subseteq S$$
  
 $A \in \mathcal{B} \iff A^C \in \mathcal{B} \implies A^C \subseteq S$ . Also, since  $A^C$  is the complement of  $A$ , it is clear that  $S = A \cup A^C$ .

$$\therefore P(S) = 1 \iff P(A \cup A^C) = 1 \iff P(A) + P(A^C) = 1$$

where 1 is by condition 3 in  $\blacksquare$  Definition 4 since  $A \cap A^C = \emptyset$  by definition of a complement of a set.

2. Note that  $S \cup \emptyset = S$  and  $S \cap \emptyset = \emptyset$ . Using a similar argument as above,

$$1 = P(S) = P(S \cup \emptyset) = P(S) + P(\emptyset) \implies P(\emptyset) = 0$$

- 3. By 1 from above,  $P(A) = 1 P(A^C)$ . Since  $0 \le P(A^C) \le 1$ , we have that P(A) is at most 1, as required.
- 4. Note that  $A = (A \cap B) \cup (A \cap B^C)$ . Clearly,  $(A \cap B) \cap (A \cap B^C) = \emptyset$ . Hence by condition 3 in  $\triangle$  Definition 4,

<sup>6</sup> This is an easy proof using the basic way of proving membership.

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

5. Consider  $P(A \cup B) + P(A \cap B)$ . By definition,

= P(A) + P(B)

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

where each of the sets in brackets are disjoint from each other<sup>7</sup>. By condition 3 of  $\blacksquare$  Definition 4, we would then have

$$P(A \cup B) + P(A \cap B)$$
  
=  $P(A \cap B^{C}) + P(A \cap B) + P(A^{C} \cap B) + P(A \cap B)$   
=  $2P(A \cap B) + P(A) - P(A \cap B) + P(B) - P(A \cap B)$  by 4

6. Note that  $B = B \cap S = B \cap (A^C \cup A) = (B \cap A^C) \cup A$ . Clearly,  $A \cap (B \cap A^C) \neq \emptyset$ . By condition 3 in  $\blacksquare$  Definition 4, we thus have that

$$P(B) = P(B \cap A^{C}) + P(A). \tag{\dagger}$$

Suppose  $A \subseteq B$ . Then  $B \cap A^C \neq \emptyset$ . I shall make the claim that  $B \cap A^C \in \mathcal{B}$ . Since  $A \subseteq B$  we have that

$$a \in (B \cap A^{C}) \iff a \in B \land a \in A^{C}$$
  
 $\iff a \in B \land a \notin A$   
 $\iff a \in (B \setminus A).$ 

But  $B \setminus A$  is a subset of B from the above steps<sup>8</sup>. Therefore,  $(B \cap A^C) \subseteq B \in \mathcal{B}$  as required.

With that done, by condition 1 in  $\blacksquare$  Definition 4,  $P(B \cap A^C) \ge 0$ . Hence from Equation (†), we have that

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

as required.

<sup>7</sup> Again, this is not hard to show

<sup>8</sup> This is rather obvious from the steps, since  $\forall a \in (B \cap A^C)$ ,  $a \in B$ .

#### Definition 5 (Conditional Probability)

Suppose S is a sample space of a random experiment, and  $A, B \subseteq S$ . The conditional probability of A given B is given by

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$
 provided  $P(B) > 0$ . (1.2)

#### Definition 6 (Independent Events)

Suppose S is a sample space of a random experiment, and A, B  $\subseteq$  S. A and B are said to be independent of each other if

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

#### • Proposition 2 (Boole's Inequality)

If  $\{A_j\}_{j=1}^{\infty}$  is a sequence of events, then

$$P\left(\bigcup_{j=1}^{\infty} A_j\right) \le \sum_{j=1}^{\infty} P(A_j)$$

#### Proof

Proof shall be provided later

#### • Proposition 3 (Bonferroni's Inequality)

If  $\{A_j\}_{j=1}^k$  is a set of events where  $k \in \mathbb{N}$ , then

$$P\left(\bigcap_{j=1}^{k} A_j\right) \ge 1 - \sum_{j=1}^{k} P(A_j^C)$$

Proof

Proof shall be provided later

#### • Proposition 4 (Continuity Property)

*If*  $A_1 \subset A_2 \subset ...$  *is a sequence where*  $A = \bigcup_{i=1}^n A_i$ , then

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

Proof

Proof shall be provided later

### 1.2 Random Variable

#### Definition 7 (Random Variable)

In a given probability space  $(S, \mathcal{B}, P)$ , the function  $X : S \to \mathbb{R}$  is called a random variable<sup>9</sup> if

$$P(X \le x) = P(\{\omega \in S : X(\omega) \le x\}) \tag{1.3}$$

is defined for all  $x \in \mathbb{R}^{10}$ .

<sup>9</sup> We shall use rv as shorthand for random variable in this set of notes.

 $^{10}$  *X* ≤ *x* is an abbreviation for { $\omega \in S$  :  $X(\omega) < x$ } ∈  $\mathcal{R}$ 

#### Example 1.2.1

In a coin flip experiment, we have that  $S = \{H, T\}$  where  $\mathbb{P}(S) = \{\emptyset, S, \{H\}, \{T\}\}$ . Define X: the number of heads in a flip, i.e.

$$X({H}) = 1$$
 and  $X({T}) = 0$ 

To prove why X is a random variable given this definition, notice that

$$x < 0 \implies P(X \le x) = P(\{\omega \in S : X(\omega) < 0\}) = P(\emptyset) = 0$$

$$x \ge 1 \implies P(X \le x) = P(\{\omega \in S : X(\omega) \le x\}) = P(\{H, T\})$$

$$= P(\{H\}) + P(\{T\}) = 1 \text{ by Independence}$$

$$0 \le x < 1 \implies P(X \le x) = P(\{\omega \in S : X(\omega) \le x\}) = P(T) \ge 0$$

which shows that P is defined for all  $x \in \mathbb{R}$ . Hence X is a random variable.

#### Definition 8 (Cumulative Distribution Function)

The cumulative distribution function (c.d.f) of a random variable X is defined as

$$\forall x \in \mathbb{R} \quad F(x) = P(X \le x)$$

#### 66 Note

Notice that F(x) is defined for all real numbers, and since it is a proba*bility, we have*  $0 \le F(x) \le 1$ .

#### • Proposition 5 (Properties of the cdf)

1. 
$$\forall x_1 < x_2 \in \mathbb{R}$$
  $F(x_1) \leq F(x_2)$ 

2. 
$$\lim_{x\to-\infty}=0 \wedge \lim_{x\to\infty}=1$$

3. 
$$\lim_{x\to a^+} F(x) = F(a)^{-11}$$

4.  $\forall a < b \in \mathbb{R} \ P(a < X \le b) = P(X \le b) - P(X \le a) =$ F(b) - F(a)

5. 
$$P(X = b) = F(b) - \lim_{a \to b^{-}} F(a)^{-12}$$

<sup>11</sup> *F* is a **right-continuous** function.

12 This is also called the magnitude of the jump.

#### Proof

Proof shall be provided later

#### 66 Note

The definition and properties of the cdf hold for the rv X regardless of whether S is discrete (finite or countable) or not.

#### 1.3 Discrete Random Variable

#### Definition 9 (Discrete Random Variable)

An  $rv\ X$  is a **discrete random variable** when its image is finite or countably infinite, i.e.  $X \in \{x_1, x_2, ...\}$ . The function

$$\forall x \in \mathbb{R} \quad f(x) := P(X = x) = F(x) - \lim_{\varepsilon \to 0^+} F(x - \varepsilon)$$

is its probability function, commonly known as the **probability mass** function (pmf). The set  $A := \{x : f(x) > 0\}$  is called the **support set** of X, and

$$\sum_{x \in A} f(x) = \sum_{i=1}^{\infty} f(x_i) = 1.$$
 (1.4)

### • Proposition 6 (Properties of pmf)

With the notation from 🗐 Definition 9, prove that

- $1. \ \forall x \in \mathbb{R} \ f(x) \ge 0$
- $2. \sum_{x \in A} f(x) = 1$

#### Proof

- 1. This result follows from the fact that f is a pdf, a probability, i.e.  $\forall x \in R$ , f(x) = 0 is  $x \notin S$  where S is the sample space, and  $0 \le f(x) \le 1$  if  $x \in S$ .
- 2. Since  $A = \{x : f(x) > 0\}$ , we know that

$$\sum_{x \in A} f(x) > 0.$$

*If we consider all the elements of A, we have that the events*  $(X = x_i)$ *,* for  $x_i \in A$ , constitutes the entire sample space. Therefore,

$$\sum_{x \in A} f(x) = \sum_{x \in A} P(X = x) = P(S) = 1.$$

#### Exercise 1.3.1

Consider an urn containing r red marbles and b black marbles. Find the pmf of the rv for the following:

- 1. X = number of red balls in n selections without replacement.
- 2. X = number of red balls in n selections with replacement.
- 3. X = number of black balls selected before obtaining the first red ball if sampling is done with replacement.
- 4. X = number of black balls selected before obtaining the kth red ball if sampling is done with replacement.

#### Solution

1. Let  $\overline{d} = \max\{n, r+b\}$ . The desired pmf is therefore the pmf from the hypergeometric distribution

$$\forall x \in \mathbb{Z}_{\leq r}^+ \quad f(x) = \frac{\binom{r}{x}\binom{b}{d-x}}{\binom{r+b}{d}}.$$

2.  $\forall x \in \mathbb{Z}^+$   $f(x) = \binom{n}{x} \left(\frac{r}{r+b}\right)^x \left(\frac{b}{r+b}\right)^{n-x}$ , which is the pmf of the binomial distribution.

3. 
$$\forall x \in \mathbb{Z}^+$$
  $f(x) = \left(\frac{b}{r+b}\right)^x \left(\frac{r}{r+b}\right)$ 

4. 
$$\forall x \in \mathbb{Z}^+$$
  $f(x) = \binom{x+k-1}{k-1} \left(\frac{b}{r+b}\right)^x \left(\frac{r}{r+b}\right)^k$ 

#### Example 1.3.1

Consider the function

$$f(x) = \begin{cases} \frac{C\mu^x}{x!} & x \in \mathbb{Z}^+, \, \mu > 0\\ 0 & otherwise \end{cases}$$

Find C such that f(x) is a pmf for the rv X.

#### Solution

We have that

$$1 = \sum_{x \in \mathbb{Z}^+} \frac{C\mu^x}{x!}$$
$$= C \sum_{x \in \mathbb{Z}^+} \frac{\mu^x}{x!}$$
$$= Ce^{\mu}$$

Thus  $C = e^{-\mu}$ .

#### Exercise 1.3.2

*Prove that the pdf of X*  $\sim$  Poi( $\mu$ ) *sums to 1 over all of its values.* 

#### Solution

$$\sum_{x \in \mathbb{N}} \frac{\mu^x e^{-\mu}}{x!} = e^{-\mu} \sum_{x \in \mathbb{N}} \frac{\mu^x}{x!}$$

$$= e^{-\mu} e^{\mu} \quad \because \sum_{x \in \mathbb{N}}^{\infty} \frac{k^x}{x!} = e^k$$

$$= 1$$

#### Exercise 1.3.3

If X is a random variable with pmf

$$f(x) = \frac{-(1-p)^x}{x \log p}$$
,  $x = 1, 2, ...$ ;  $0 ,$ 

show that

$$\sum_{x \in \mathbb{N}} f(x) = 1$$

#### Solution

$$\sum_{x \in \mathbb{N}} \frac{-(1-p)^x}{x \log p} = -\frac{1}{\log p} \sum_{x \in \mathbb{N}} \frac{(-1)^x (p-1)^x}{x}$$

$$= -\frac{1}{\log p} \underbrace{\left[ -(p-1) + \frac{(p-1)^2}{2} - \frac{(p-1)^3}{3} + \ldots \right]}_{\text{Taylor expansion of } -\log p}$$

$$= 1$$

This gives us that  $\forall x \in \mathbb{Z}^+$ ,  $f(x) = \frac{e^{-\mu}\mu^x}{x!}$ , and this is, of course, the pmf of the Poisson distribution.

#### 2.1 Continuous Random Variable

#### Definition 10 (Continuous Random Variable)

Suppose X is an rv with cdf F. If F is a continuous function for all  $x \in \mathbb{R}$  and F is differentiable except possibly at countably many points, then X is a **continuous** rv. The probability function, or more commonly known as the **probability density function** (pdf), of X is f(x) = F'(x) wherever F is differentiable on x and x otherwise.

The set  $A = \{x : f(x) > 0\}$  is called the **support set** of X and

$$\int_{x \in A} f(x) \, dx = 1$$

#### • Proposition 7 (Properties of pdf)

Let X be a random variable and f be its pdf.

1. 
$$\forall x \in \mathbb{R} \quad f(x) \ge 0$$

$$2. \int_{-\infty}^{\infty} f(x) \, dx = 1$$

3. 
$$f(x) = \lim_{h\to 0} \frac{F(x+h)-F(x)}{h} = \lim_{h\to 0} \frac{P(x\leq X\leq x+h)}{h}$$
 (if the limit exists)

4. 
$$\forall x \in \mathbb{R}$$
  $F(x) = \int_{-\infty}^{x} f(t) dt$ 

5. 
$$P(a < X \le b) = \int_a^b f(x) dx = F(b) - F(a)$$

6. 
$$P(X = b) = F(b) - \lim_{a \to b^{-}} F(a) = F(b) - F(b) = 0$$

#### Proof

- 1. The argument of this proof is similar to that provided in ♠ Proposition 6.
- 2. Same as above, except that the support set can now have complete intervals.
- 3. The first equation follows from the first principles of Calculus. The second equation follows by method of calculation using the cdf.
- 4.  $F(x) = P(X \le x) = \int_{-\infty}^{x} f(t) dt$ .
- 5. This follows immediately from the above property.
- 6. The first part of the equation is a way to interpret the above property. The limit equates to F(b) since F is continuous.

#### Example 2.1.1

Consider the function

$$f(x) = \begin{cases} \frac{\theta}{x^{\theta+1}} & x \ge 1\\ 0 & x < 1 \end{cases}$$

For what values of  $\theta$  is f a pdf?

#### Solution

If f is a pdf, then  $\theta \geq 0$ . In fact,  $\theta \neq 0$ ; otherwise f would be equivalently 0 for all  $x \in \mathbb{R}$ , which would imply that  $\int_{\mathbb{R}} f = 0$ , which is impossible. It remains to check if  $\theta > 0$  is a safe choice. Now

$$\int_{1}^{\infty} \frac{\theta}{x^{\theta+1}} dx = -\frac{1}{x^{\theta}} \Big|_{1}^{\infty} = 1$$

Note that the above integral is valid because  $\frac{1}{x^{\theta+1}} \leq \frac{1}{x}$ . Therefore the choice of  $\theta > 0$  is safe.

#### Binomial Distribution 2.2.1

#### Definition 11 (Binomial RV)

Consider X to be the number of successes in a sequence of n experiments where

- 1. experiments are independent;
- 2. the outcome of each experient is a binary (e.g. success or failure); and
- 3. has the **probability of success**, p for each singular experiment.

*X* is called a *Binomial rv*, and we write  $X \sim Bin(n, p)$  and its pmf is

$$P(X = x) = \begin{cases} \binom{n}{x} p^{x} (1 - p)^{n - x} & x = 0, 1, 2, ..., n \\ 0 & otherwise \end{cases}$$

#### Geometric Distribution 2.2.2

#### Definition 12 (Geometric RV)

Consider a sequence of independent success/failure (binary) experiments, each of which has a success probability of p. Let X be the number of failures before the first success is reached. We call X a Geometric rv, and we write  $X \sim \text{Geo}(p)$ , and its pmf is

$$P(X = x) = \begin{cases} (1 - p)^{x} p & x = 0, 1, 2, ..., n \\ 0 & otherwise \end{cases}$$

#### 66 Note

Some authors would define the Geometric rv as:

*Let X be the number of experiments until the first success.* 

But that really is just a play of words.

#### 2.2.3 Poisson Distribution

#### Definition 13 (Poisson RV)

Suppose X is defined to be the number of occurrences of an event in a given time period. If the process on which the events occur satisfies the following:

- 1. The number of occurrences in non-overlapping intervals are independent of each other;
- 2. The probability of the occurrence of an event in a short interval of length h is proportional to h;
- 3. For sufficiently short time periods of length h, the probability of 2 or more events occurring in the interval is negligible, i.e. almost zero;

then X is a Poisson rv, and we write  $X \sim Poi(\lambda)$ , with  $\lambda > 0$ , and the pmf is

$$P(X = x) = egin{cases} rac{e^{-\lambda} \lambda^x}{x!} & x = 0, 1, ... \ 0 & otherwise \end{cases}$$

### 2.3 Examples of Continuous RVs

#### 2.3.1 Normal/Gaussian Distribution

### Definition 14 (Normal / Gaussian RV)

The **Normal/Gaussian** Distribution is a continuous probability distribution that is symmetric about the mean, showing that data around the

mean is more frequent than data far from the mean. If X is a Normall-*Gaussian rv, we write*  $X \sim N(\mu, \sigma^2)$ , and its pdf is

$$f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$
 for  $x \in \mathbb{R}$ .

#### Definition 15 (Standard Normal Distribution)

The Standard Normal Distribution is the simplest case of a Normal Distribution. An rv Z is called the Standard Normal rv if  $\mu = 0$  and  $\sigma = 1$ . We write  $Z \sim N(0,1)$  and its pdf is

$$f(x) = \frac{1}{\sqrt{2\pi}}e^{-\frac{x^2}{2}}$$
 for  $x \in \mathbb{R}$ .

#### Uniform Distribution 2.3.2

#### Definition 16 (Uniform RV)

If X represents the result of drawing a real number from an interval (a,b), with a < b, such that all numbers in between are equally likely to be chosen, then X is called a **Uniform** rv, and we write  $X \sim \text{Unif}(a, b)$ , and its pdf is

$$f(x) = \begin{cases} \frac{1}{b-a} & x \in (a,b) \\ 0 & otherwise \end{cases}$$

#### 2.3.3 Exponential Distribution

#### Definition 17 (Exponential RV)

Let X show the time between two consecutive events in a Poisson process, i.e. the 3 conditions in Poisson Distribution are satisfied. Then X is called an Exponential rv, and we write  $X \sim \text{Exp}(\theta)$ , where  $\theta > 0$ , with its pdf

$$f(x) = \begin{cases} \frac{1}{\theta}e^{-\frac{x}{\theta}} & x > 0\\ 0 & otherwise \end{cases}$$

#### 2.3.4 *Gamma Distribution*

#### Definition 18 (Gamma RV)

Let X be the sum of n independent Exponential rvs with some fixed  $\theta$ . Then X is called a Gamma rv, in which we write  $X \sim \Gamma(n, \theta)$ , and its pdf is

$$f(x) = \begin{cases} \frac{x^{n-1}e^{-\frac{x}{\theta}}}{\Gamma(n)\theta^n} & x > 0 \land \theta, n > 0\\ 0 & otherwise \end{cases}$$

where  $\Gamma(n) = \int_0^\infty e^{-y} y^{n-1} dy = (n-1)!$ , where the last equality is true when n is an integer.

#### 66 Note

The Gamma Distribution is usually used for when we are looking for the probability of the occurrence of the n-th event in the desired waiting time.

#### 2.4 Functions of Random Variables

Consider the rv X with pdf/pmf f and cdf F. Given Y = h(X) where h is some real-valued function, we are interested in finding the pdf/pmf of Y.

The following are some possible scenarios:

1. *X* and *Y* are both discrete;

- 2. *X* is continuous and *Y* is discrete;
- 3. *X* and *Y* are both continuous

We may also define Y = h(X) for a continuous rv X such that Y is neither discrete nor continuous (e.g. discrete for some values of X while continuous for others).

#### 2.4.1 Discrete X and Discrete Y

If X and Y = h(X) are both discrete, we can derive P(Y = y) by mapping values in Y onto their corresponding value through h, i.e.

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

#### Exercise 2.4.1

*Let X have the following probability function:* 

$$f_X(x) = egin{cases} rac{e^{-1}}{x!} & x = 0,1,2,... \ 0 & otherwise \end{cases}$$

Find the pmf of  $Y = (X - 1)^2$ .

#### Solution

Note that since

Dom 
$$X = \{0, 1, 2, 3, 4, ...\},\$$

we have that

Dom 
$$Y = \{1, 0, 1, 4, 9, ...\}.$$

With that, note that

$$P(Y = 0) = P(X = 1) = \frac{e^{-1}}{1!}$$

$$P(Y = 1) = P(X = 0 \text{ or } 2) = P(X = 0) + P(X = 2)$$

$$= \frac{e^{-1}}{0!} + \frac{e^{-1}}{2!} = e^{-1} \left( 1 + \frac{1}{2} \right) = \frac{3}{2} e^{-1}$$

$$P(Y = 4) = P(X = 3) = \frac{e^{-1}}{3!}$$

$$P(Y = 9) = P(X = 4) = \frac{e^{-1}}{4!}.$$

Therefore, the pmf of  $Y = (X - 1)^2$  is

$$P(Y = y) = egin{cases} e^{-1} & y = 0 \ rac{3}{2}e^{-1} & y = 1 \ rac{e^{-1}}{(1+\sqrt{y})!} & y = 4,9,16,... \ 0 & otherwise \end{cases}$$

#### 2.4.2 *Continuous X and Discrete Y*

If X is continuous and Y is discrete, we can use the method that we have used in the previous subsection, and replace  $\Sigma$  by the integral sign  $\int$ , i.e. define  $A := \{x : h(x) = y\}$  such that we have

$$P(Y = y) = \int_{A} f(x) \, dx$$

#### Example 2.4.1 (Example 2.9)

Suppose X is a random variable with the following probability function

$$f_X(x) = egin{cases} 2e^{2x} & x > 0 \ 0 & otherwise \end{cases}.$$

Suppose Y = h(X) is defined as follows:

$$Y = \begin{cases} 1 & X < 1 \\ 2 & 1 \le X \le 2 \\ 3 & X > 2 \end{cases}$$

Find the probability function of Y.

#### Solution

*Note that*  $X \sim \text{Exp}(\frac{1}{2})$ . *So it is clear that* X *is a crv and since* Y = 1, 2, *or* 

3, we have that Y is discrete. Now

$$P(Y = 1) = P(X < 1) = \int_0^1 2e^{-2x} dx$$

$$= -e^{-2x} \Big|_0^1 = 1 - e^{-2}$$

$$P(Y = 2) = P(1 \le X \le 2) = \int_1^2 2e^{-2x} dx$$

$$= -e^{-2x} \Big|_1^2 = e^{-2} - e^{-4}$$

$$P(Y = 3) = P(X > 2) = \int_2^\infty 2e^{-2x} dx$$

$$= -e^{-2x} \Big|_2^\infty = e^{-4}$$

*Thus the pmf is* 

$$P(Y = y) = \begin{cases} 1 - e^{-2} & Y = 1 \\ e^{-2} - e^{-4} & Y = 2 \\ e^{-4} & Y = 3 \end{cases}$$

#### Continuous X and Continuous Y 2.4.3

If *X* and Y = h(X) are both continous, start with the definition of the cdf of Y, i.e.

$$F_Y(y) = P(Y \le y) = P(h(X) \le y)$$

solve the inequality for *X*, and then obtain the cdf of *Y*. We will then only need to differentiate the cdf wrt *y* to get the pdf that we desire.

#### Example 2.4.2 (Example 2.10)

*Let X have the following pdf:* 

$$f_X(x) = egin{cases} 2e^{-2x} & x \geq 0 \ 0 & otherwise \end{cases}$$

Find the pdf of  $Y = \sqrt{X}$ .

#### Solution

We have that the range of values where  $f_Y(y) \leq 0$  is  $y \geq 0$ . Now

$$F_Y(y) = P(Y \le y) = P(\sqrt{X} \le y) = P(X \le y^2)$$

$$= \int_0^{y^2} 2e^{-2x} dx$$

$$= -e^{-2x} \Big|_0^{y^2} = 1 - e^{-2y^2}$$

*Therefore, the pdf of Y is* 

$$f_Y(y) = egin{cases} rac{d}{dy} 1 - e^{-2y^2} = 4ye^{-2y^2} & y \leq 0 \ 0 & otherwise \end{cases}.$$

#### 2.4.4 A Formula for the Continuous Case

# **■** Theorem 8 (One-to-One Transformation of a Random Variable)

Suppose X is a continuous random variable with pdf  $f_X$  and support set  $A = \{x : f_X(x) > 0\}$  and Y = h(X) where h is a real-valued function. Let  $f_Y$  be the pdf of the rv Y and let  $B = \{y : f_Y(y) > 0\}$ . If h is a one-to-one function from A to B and if h' is continuous, then

$$f_Y(y) = f(h^{-1}(y)) \cdot \left| \frac{d}{dy} h^{-1}(y) \right|, \quad y \in B$$

#### Proof

Note that since h is one-to-one, it is monotonous. Suppose h is increasing. Then  $h^{-1}$  is also an increasing function. Note that the cdf of Y is

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

Then the cdf of Y is

$$f_Y(y) = \frac{d}{dy} F_X(h^{-1}(y)) = f_X(h^{-1}(y)) \cdot \frac{d}{dy} h^{-1}(y)$$

*If h is decreasing, then so is its inverse. Thus* 

$$F_Y(y) = P(Y \le y) = P(X \ge h^{-1}(y)) = 1 - F_X(h^{-1}(y))$$

Thus the cdf of Y is

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

Note that the pdf of Y is indeed positive since  $h^{-1}$  is decreasing.

Combining the two, we have that

$$f_Y(y) = f_X(h^{-1}(y)) \cdot \left| \frac{d}{dy} h^{-1}(y) \right|,$$

as required.

# 3 Lecture 3 May 08th 2018

# 3.1 Functions of Random Variables (Continued)

# 3.1.1 Special Cases

### Example 3.1.1

Recall Example 2.4.1. Suppose X is a rv with the following probability function

$$f_X(x) = egin{cases} 2e^{-2x} & x > 0 \ 0 & otherwise \end{cases}.$$

Define Y = h(X) as follows:

$$Y = \begin{cases} 1 & X < 1 \\ X & 1 \le X \le 2 \\ 3 & X > 2 \end{cases}$$

Find the cdf of Y.

#### Solution

Solution is given differently in the 2 sections. I am not happy with either solutions because some things don't add up. My opinion is that the definition of Y is badly given, along with a badly phrased question. As a result, there are more ways than one to interpret an already confusing information, and thus we have ourselves one hell of a mess.

# 3.2 Probability Integral Transformation

## **■** Theorem 9 (Probability Integral Transformation)

If X is a continuous rv with cdf F, then  $Y = F(X) \sim \text{Unif}(0,1)$ . Y = F(X) is called the **probability integral transformation**.

### 66 Note

The distribution of Y = F(X) can be proven.

#### Proof

Let X be a continuous rv and Y = F(X). Since F(X) is one-to-one and increasing (i.e. monotonous), there exists  $F^{-1}(Y)$  that is a real-valued and increasing function. Then

$$F_Y(y) = P(Y \le y) = P(F_X(X) \le y) = P(X \le F^{-1}(y))$$
  
=  $F(F^{-1}(y)) = y$ 

Note that  $F_Y(y) = y$  is the cdf of a Unif(0,1) rv, i.e. the standard uniform random variable. Thus  $Y \sim Unif(0,1)$ .

#### 66 Note

This theorem essentially states that any ro from a continuous distribution can be transformed into a standard uniform distribution.

### **Example 3.2.1 (Example 2.11)**

Suppose  $X \sim \text{Exp}(01)$ . We know that  $F_X(x) = 1 - e^{-10x}$  for all  $x \in \mathbb{R}a$ . By Probability Integral Transformation, we have that  $Y = F_X(X) = 1 - e^{-10X} \sim \text{Unif}(0,1)$ .

Note that the converse of Probability Integral Transformation is

true:

### ■ Theorem 10 (Converse of Probability Integral Transformation)

Suppose X is a continuous rv with cdf F such that  $F^{-1}$  exists. If  $U \sim$ Unif(0, 1), we have that  $Y = F^{-1}(U) \sim X$ .

### Proof

*Note that* 

$$F_Y(y) = P(Y \le y) = P(F^{-1}(U) \le y)$$
  
=  $P(U \le F_X(y)) = F_X(y)$ .

### Example 3.2.2 (Example 2.12)

Suppose  $X \sim \text{Unif}(0,1)$ . Find a transformation T such that  $T(X) \sim$  $\exp(\theta)$ .

### Solution

Let  $Y = T(X) \sim \text{Exp}(\theta)$ . Note that

$$F_Y(y) = 1 - e^{-\frac{y}{\theta}}, \quad y > 0$$

Observe that since

$$x = 1 - e^{-\frac{y}{\theta}} \implies y = -\theta \ln(1 - x)$$

we have that

$$F_Y^{-1}(X) = -\theta \ln(1-X).$$

By Converse of Probability Integral Transformation 10, we have that T =

# 3.3 Location-Scale Families

When we look into methods for constructing confidence intervals for an unknown parameter  $\theta$ . If the parameter  $\theta$  is either a *scale parameter* 

or *location parameter*, then a confidence interval is easier to construct.

### Definition 19 (Location Parameter and Family)

Suppose X is a continuous rv with pdf  $f(x; \mu)$ , where  $\mu$  is a parameter of the distribution of X. Let  $F_0(x) = F_X(x; \mu = 0)$ , where  $F_X$  is the cdf of X, and  $f_0(x) = f(x; \mu = 0)$ . The parameter  $\mu$  is called a **location** parameter of the distribution if

$$F_X(x;\mu) = F_0(x-\mu), \quad \mu \in \mathbb{R}$$

or equivalently,

$$f(x;\mu) = f_0(x-\mu), \quad \mu \in \mathbb{R}.$$

We say that F belongs to a **location family** of distributions.

### Definition 20 (Scale Parameter and Family)

Suppose X is a continuous rv with pdf  $f(x;\theta)$ , where  $\theta$  is a parameter of the distribution of X. Let  $F_1(x) = F_X(x;\theta=1)$ , where  $F_X$  is the cdf of X, and  $f_1(x) = (x;\theta=1)$ . The parameter  $\theta$  is called a **scale parameter** of the distribution if

$$F_X(x;\theta) = F_1(\frac{x}{\theta}). \quad \theta > 0$$

or equivalently,

$$f(x;\theta) = \frac{1}{\theta} f_0(\frac{x}{\theta}), \quad \theta > 0.$$

We say that F belongs to a **scale family** of distributions.

#### Definition 21 (Location-Scale Family)

Suppose X is an rv with  $cdf\ F(x;\mu,\sigma)$  where  $\mu\in\mathbb{R}$  and  $\sigma>0$  are the parameters of the distribution. Let  $Y=\frac{X-\mu}{\sigma}$ . If the distribution of Y does not depend on  $\mu$  and/or  $\sigma$ , then F is said to belong to a **location-scale family** of distributions, with **location parameter**  $\mu$  and **scale parameter**  $\sigma$ . In other words, F belongs to a location-scale family of distributions if

$$F(x;\mu,\theta) = F_0\left(\frac{x-\mu}{\theta}\right),$$

where  $F_0(x) = F(x; \mu = 0, \theta = 1)$ , or equivalently,

$$f(x;\mu,\theta) = \frac{1}{\theta} f_0\left(\frac{x-\mu}{\theta}\right),$$

where  $f_0(x) = f(x; \mu = 0, \theta = 1)$ .

### **Example 3.3.1 (Example 2.13)**

Consider  $X \sim G(\mu, \sigma)$ . Show that  $F_X$  belongs to a location-scale family of distributions.

We know that if  $\mu=0$  and  $\sigma=1$ , then  $Y=\frac{X-\mu}{\sigma}\sim G(0,1)$ , and we know that G(0,1) has no dependence on unknowns  $\mu$  and  $\sigma$ . Therefore,  $F_X$ belongs to the location-scale family of distributions, with location parameter  $\mu$  and scale parameter  $\sigma$ .

Another solution is to show that one of the equations in the definition is fulfilled. Observe that

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

So if we set  $\mu = 0$  and  $\sigma = 1$  to get  $f_0$ , we have that

$$f_0(x) = \frac{1}{\sqrt{2\pi}}e^{-\frac{x^2}{2}}.$$

Now, note that

$$f(x) = \frac{1}{\sigma} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{(x-\mu)}{\sigma}\right)^2}.$$

Let  $y = \frac{x-\mu}{\sigma}$ , and we have ourselves

$$f(x) = \frac{1}{\sigma} \frac{1}{\sqrt{2\pi}} e^{-\frac{y^2}{2}} = \frac{1}{\sigma} f_0(\frac{x-\mu}{\sigma})$$

#### Example 3.3.2 (Example 2.14)

Consider  $X \in G(\mu, 2)$  where  $\mu = E(X)$ . Show that  $\mu$  is a location parame-

We can use a similar approach as before and define  $Y = X - \mu$  which follows G(0,2). It is clear that we then have that  $F_X$ , the cdf of X, belongs to a location family of distributions.

### Example 3.3.3 (Example 2.15)

Consider  $X \sim \text{Exp}(\theta)$ . Show that  $F_X$  belongs to a scale family of distributions and find the scale parameter.

*Note that* 

$$f(x) = \begin{cases} \frac{1}{\theta}e^{-\frac{x}{\theta}} & x > 0\\ 0 & otherwise \end{cases}$$

Let  $Y = \frac{X}{\theta}$ . Then

$$F_Y(y) = P(Y \le y) = P(\frac{X}{\theta} \le y)$$

$$= P(X \le \theta y) = \int_0^{\theta y} \frac{1}{\theta} e^{-\frac{x}{\theta}} dx$$

$$= -e^{-\frac{x}{\theta}} \Big|_0^{\theta} y = 1 - e^{-y}$$

and we have

$$f_Y(y) = egin{cases} e^{-y} & y > 0 \ 0 & otherwise \end{cases}$$

*Note that if we set*  $\sigma = 1$  *to get*  $f_1$ *, we have* 

$$f_1(x) = egin{cases} e^{-x} & x > 0 \ 0 & otherwise \end{cases}.$$

Therefore,  $F_X$  belongs to a scale family of distributions.

#### 3.4 Expectations

### 3.4.1 Expectations

### Definition 22 (Expectation of A Discrete RV)

If X is a discrete rv with pmf f and support set A, then the **expectation** of X, or the **expected** value of X is defined by

$$E(X) = \sum_{x \in A} x f(x) \tag{3.1}$$

provided that the sum converges absolutely, i.e.

$$E(|X|) = \sum_{x \in A} |x| f(x) < \infty.$$

If E(|X|) does not converge, then we say that E(X) does not exist.

### Definition 23 (Expectation of A Continuous RV)

If X is a continuous rv with pdf f and support set A, then the expectation of X, or the expected value of X is defined by

$$E(X) = \int_{Y \in A} x f(x) \tag{3.2}$$

provided that the integral converges absolutely, i.e.

$$E(|X|) = \int_{x \in A} |x| f(x) < \infty.$$

If E(|X|) does not converge, then we say that E(X) does not exist.

### **Example 3.4.1 (Example 2.16)**

Suppose  $X \sim Poi(\lambda)$ . Calculate E(X).

### Solution

Note

$$f(x) = \begin{cases} \frac{e^{-\lambda}\lambda^x}{x!} & x = 0, 1, 2, \dots \\ 0 & otherwise \end{cases}.$$

Then

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

$$= 0 + \sum_{x=1}^{\infty} \frac{e^{-\lambda} \lambda^x}{(x-1)!}$$

$$= e^{-\lambda} \lambda \sum_{x=1}^{\infty} \frac{\lambda^{x-1}}{(x-1)!}$$

$$= e^{-\lambda} \lambda e^{\lambda} = \lambda$$

### **Example 3.4.2 (Example 2.18)**

Suppose X is an rv with

$$f(x) = \begin{cases} \frac{1}{x^2} & 1 < x < \infty \\ 0 & otherwise \end{cases}.$$

Calculate E(X).

### Solution

Observe that  $x \cdot \frac{1}{x^2} = \frac{1}{x}$  and the antiderivative of  $\frac{1}{x}$  is  $\ln x$ , which would need to be evaluated at  $\ln \infty$ . Thus, we should instead immediately check if the integral converges absolutely.

$$E(|X|) = \int_{1}^{\infty} |x| \frac{1}{x^{2}} dx$$

$$= \int_{1}^{\infty} |x| \frac{1}{|x| |x|} dx$$

$$= \int_{1}^{\infty} \frac{1}{|x|} dx$$

$$= \int_{1}^{\infty} \frac{1}{x} dx,$$

and we notice that the integral would not converge. Therefore, E(X) does not exist.

# 4 Lecture 4 May 10th 2018

# 4.1 Expectations (Continued)

### **4.1.1** *Expectations* (Continued)

### Theorem 11 (Expectation from the cdf)

Suppose X is a non-negative continuous rv with cdf F, and  $E(X) < \infty$ . Then

$$E(X) = \int_0^\infty [1 - F(x)] \, dx = \int_0^\infty P(X \ge x) \, dx \tag{4.1}$$

If X is a discrete rv with cdf F, and  $E(X) < \infty$ , then

$$E(X) = \sum_{x=0}^{\infty} [1 - F(x)] = \sum_{x=0}^{\infty} P(X \ge x)$$
 (4.2)

#### Proof

*Note that for a continuous rv X, we have* 

$$1 - F(x) = P(X \ge x) = \int_{x}^{\infty} f(t) dt$$

Therefore,

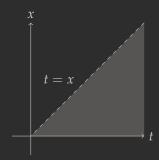
$$\int_0^\infty \left[1 - F(x)\right] dx = \int_0^\infty \int_x^\infty f(t) \, dt \, dx.$$

Since 1 - F(x) is a finite value, so is  $\int_0^\infty f(t) dt$ , and thus we can apply *Fubini's Theorem*<sup>1</sup>:

$$\int_0^\infty [1 - F(x)] \, dx = \int_0^\infty \int_x^\infty f(t) \, dt \, dx = \int_0^\infty \int_0^t f(t) \, dx \, dt$$

Note that the limits of the integral utilizes the following figure:

<sup>&</sup>lt;sup>1</sup> Condition for Fubini's Theorem to hold is that the integrand of the double integral must be absolutely convergent. See Wikipedia.



With that, note that

$$\int_0^t f(t) \, dx = x f(t) \Big|_0^t = t f(t)$$

Since t is just a dummy variable, we can indeed let t = x, and thus we have

$$\int_0^\infty \left[1 - F(x)\right] dx = \int_0^\infty x f(x) \, dx = E(X)$$

as required.

Work on the discrete case as an exer-

#### Exercise 4.1.1

For a non-negative discrete rv X with cdf F and  $E(X) < \infty$ , prove that

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

# **Example 4.1.1 (Example 2.20)**

Suppose  $X \sim \text{Exp}(\theta)$ . Use  $\blacksquare$  Theorem 11 to calculate E(X).

### Solution

*Note that X is a non-negative rv. The cdf of X* Im  $Exp(\theta)$  *is* 

$$F_X(x) = 1 - e^{-\frac{x}{\theta}}.$$

Then

$$E(X) = \int_0^\infty 1 - F_X(x) \, dx = \int_0^\infty e^{-\frac{x}{\theta}} \, dx$$
$$= -\theta e^{-\frac{x}{\theta}} \Big|_0^\infty = \theta$$

Theorem 12 (Expected Value of a Function of X)

Suppose h(x) is a real-valued function.

If X is a discrete rv with pmf f and support set A, then

$$E[h(x)] = \sum_{x \in A} h(x)f(x) \tag{4.3}$$

provided that the sum converges absolutely.

*If X is a continuous rv with pdf f, then* 

$$E[h(x)] = \int_{-\infty}^{\infty} h(x)f(x) dx \tag{4.4}$$

provided that the integral converges absolutely.

The proof is, unfortunately, not trivial. One would have to look into Lesbesgue integrals (or at the very least, Riemann-Stieltjes integrals) in order to prove this statement. This "theorem" is also called The Law of the Unconscious Statistician [Reference - Wikipedia]. An idea of the proof is given on Math SE.

### Example 4.1.2

*Suppose*  $X \sim \text{Unif}(0, \theta)$ . *Calculate*  $E(X^2)$ .

#### Solution

$$E(X^2) = \int_0^\theta \frac{x^2}{\theta} dx = \frac{1}{\theta} \frac{x^3}{3} \Big|_{x=0}^\theta = \frac{\theta^2}{3}$$

#### Exercise 4.1.2

Find the pdf of  $Y = X^2$  and find E(Y) by evaluating  $\int_{-\infty}^{\infty} y f_Y(y) dy$ 

### Theorem 13 (Linearity of Expectation)

Suppose X is an rv with pf f. Let  $a_i, b_i \in \mathbb{R}$ , for i = 1, ..., n, be constants, and  $g_i(x)$ , for i = 1, ..., n, are real-valued functions. Then

$$E\left[\sum_{i=1}^{n} (a_i g_i(X) + b_i)\right] = \sum_{i=1}^{n} (a_i E[g_i(X)] + b_i)$$
(4.5)

provided that  $E[g_i(X)] < \infty$  for i = 1, ..., n.

This theorem essentially states that the expectation is a linear operator.

### Proof

Suppose X is a discrete rv with support set A. Then

$$E\left[\sum_{i=1}^{n} (a_i g_i(X) + b_i)\right] = \sum_{x \in A} \left[\sum_{i=1}^{n} (a_i g_i(x) + b_i)\right] f(x) \quad \therefore 1 \text{ Theorem } 12$$

$$= \sum_{x \in A} \sum_{i=1}^{n} \left[a_i g_i(x) f(x) + b_i f(x)\right]$$

$$= \sum_{i=1}^{n} \sum_{x \in A} \left[a_i g_i(x) f(x) + b_i f(x)\right] \quad (*)$$

$$= \sum_{i=1}^{n} \left[a_i \sum_{x \in A} g_i(x) f(x) + b_i \sum_{x \in A} f(x)\right]$$

$$= \sum_{i=1}^{n} \left[a_i E[g_i(X)] + b_i\right]$$

where note that (\*) is valid because  $a_i, b_i$  are constants, and  $g_i(x), f(x)$  are finite real-valued functions.

### 66 Note

In general,  $E(g(X)) \neq g(E(X))$  unless if g is a linear function. For example, for  $a, b \in \mathbb{R}$ , we have

$$E(aX + b) = aE(X) + b$$

### 4.1.2 *Moments and Variance*

Since these concepts were introduced in STAT230 and were given little treatment in the lecture, we shall only cover over them briefly.

#### Definition 24 (Variance)

The expectation tof the squared deviation of an rv from its mean is called the variance, i.e. for an rv X with mean  $\mu = E(X)$ ,

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

### Definition 25 (Moments)

Let X be an rv with mean  $\mu$ .

The  $k^{th}$  moment about the origin is defined as:

$$E(X^k)$$

The k<sup>th</sup> moment about the mean is defined as:

$$E[(X-\mu)^k]$$

The  $k^{th}$  factorial moment is defined as:

$$E[X^{(k)}] = E[X(X-1)...(X-k+1)] = E\left[\frac{X!}{(X-k)!}\right]$$

### Theorem 14 (Variance of a Linear Function)

Suppose X is an rv with pf f and  $a, b \in \mathbb{R}$ . Then

$$Var(aX + b) = a^2 Var(X)$$

### Proof

Observe that

$$Var(aX + b) = E[(aX + b)^{2}] - E(aX + b)^{2}$$

$$= E[a^{2}X^{2} + 2abX + b^{2}] - (aE(X) + b)^{2}$$

$$= a^{2}E(X^{2}) + 2abE(X) + b^{2} - (a^{2}E(X)^{2} + 2abE(X) + b^{2})$$

$$= a^{2}E(X^{2}) - a^{2}E(X)^{2} = a^{2}Var(X)$$

### Example 4.1.3 (Example 2.22 (course notes - 2.6.10 (1)))

If 
$$X \sim \text{Poi}(\theta)$$
, then  $E[X^{(k)}] = \theta^k$  for  $k = 1, 2, ...$ 

### Solution

Note

$$f_X(x) = egin{cases} rac{e^{- heta} heta^x}{x!} & x = 0, 1, 2, ... \ 0 & otherwise \end{cases}$$

So

$$\begin{split} E[X^{(k)}] &= E(X(X-1)(X-2)\dots(X-k+1)) \\ &= \sum_{x=0}^{\infty} x(x-1)(x-2)\dots(x-k+1) \frac{e^{-\theta}\theta^x}{x!} \\ &= 0 + \sum_{x=k}^{\infty} x(x-1)(x-2)\dots(x-k+1) \frac{e^{-\theta}\theta^x}{x!} \quad (*) \\ &= \sum_{x=k}^{\infty} \frac{x!}{(x-k)!} \frac{e^{-\theta}\theta^x}{x!} \quad \because x(x-1)\dots(x-k+1) = \frac{x!}{(x-k)!} \\ &= e^{-\theta}\theta^k \sum_{x=k}^{\infty} \frac{\theta^{x-k}}{(x-k)!} \\ &= e^{-\theta}\theta^k \sum_{y=0}^{\infty} \frac{\theta^y}{y!} \qquad let \ y = x-k \\ &= e^{-\theta}\theta^k e^{\theta} = \theta^k \end{split}$$

where for (\*) we have that  $\sum_{x=0}^{k-1} x(x-1) \dots (x-k+1)A = 0$  for any  $A \in \mathbb{R}$ .

Note that it is not necessarily true that

$$x(x-1)\dots(x-k+1) = \frac{x!}{(x-k)!}$$

for  $0 \le x \le k - 1$ . And so we can only say that the equality is true for  $x \ge k$ , and hence we have the approach that we use in (\*).

# 4.2 Inequalities

### 4.2.1 Markov/Chebyshev Style Inequalities

#### Theorem 15 (Markov's Inequality)

If X is a non-negative rv and a > 0, then the probability that X is no less than a is no greater than the expectation of X divided by a, i.e.

$$P(X \ge a) \le \frac{E(X)}{a} \tag{4.6}$$

#### Proof

We shall prove for the discrete case. Suppose X is a non-negative discrete rv with pf f. Let  $A \subset S$ , where S is the sample space, such that  $A = \{ w \in S : X(w) \ge a \}.$ 

$$E(X) = \sum_{x \in S} xf(x)$$

$$= \sum_{x \in A} xf(x) + \sum_{x \notin A} xf(x)$$

$$\geq \sum_{x \in A} xf(x) \quad \therefore \sum_{x \notin A} xf(x) \geq 0$$

$$\geq \sum_{x \in A} af(x)$$

$$= a \sum_{x \in A} f(x) = a \cdot P(A)$$

$$= a \cdot P(\{w \in S : X(w) \geq a\}) = aP(X \geq a).$$

#### Exercise 4.2.1

Prove Markov's Inequality for a continuous

# ■ Theorem 16 (Markov's Inequality 2)

If X is a non-negative rv and a, k > 0, then the probability that X is no less than a is no greater than the expectation of X divided by a, i.e.

$$P(|X| \ge a) \le \frac{E(|X|^k)}{a^k} \tag{4.7}$$

#### Proof

We shall, again, prove for the discrete case. Suppose X is a non-negative discrete rv with pf f.  $A := \{w \in S : |X(w)| \ge a\} \subseteq S$ . Then

$$E(|X|^k) = \sum_{x \in S} |x|^k f(x)$$

$$= \sum_{x \in A} |x|^k f(x) + \sum_{x \notin A} |x|^k f(x)$$

$$\geq \sum_{x \in A} |x|^k f(x) \geq \sum_{x \in A} af(x)$$

$$= a^k P(A) = a^k P(|X| \geq a).$$

Question: Can we write

$$P(\{w \in S : |X(w)| > a\}) = P(|X| > a)$$
?

### Exercise 4.2.2

Prove for the continuous case.

### Theorem 17 (Chebyshev's Inequality)

Suppose X is an rv with finite mean  $\mu$  and finite variance  $\sigma^2$ . Then for any k > 0,

$$P(|X - \mu| \ge k\sigma) \le \frac{1}{k^2} \tag{4.8}$$

### Proof

By 💻 Theorem 16,

$$P(|X - \mu| \ge k\sigma) \le \frac{E(|X - \mu|^2)}{(k\sigma)^2} = \frac{1}{k^2}$$

since 
$$E(|X - \mu|^2) = Var(X) = \sigma^2$$
.

# **Example 4.2.1 (Example 2.23)**

A post office handles, on average, 10000 letters a day. What can be said about the probability that it will handle at least 15000 letters tomorrow?

### Solution

X:= number of letters handled in a day. Note that by its definition, X is a non-negative discrete v. Then, using  $\blacksquare$  Theorem 15, since E(X)=10000

$$P(X \ge 15000) \le \frac{10000}{15000} = \frac{2}{3}.$$

Thus, we know that there is less than two-third of chance that the post office will handle more than 15000 tomorrow.

# 5 Lecture 5 May 15th 2018

# 5.1 Inequalities (Continued)

### 5.1.1 Markov/Chebyshev Style Inequalities (Continued)

### Example 5.1.1 (Example 2.24)

A post office handles 10000 letters per day with a variance of 2000 letters. What can be said about the probability that this post office handles between 8000 and 12000 letters tomorrow? What about the probability that more than 15000 letters come in (use Parent 17)?

1. Probability that this post office handles between 8000 and 12000 letters tomorrow:

$$P(8000 < X < 12000)$$

$$= P(-2000 < X - 10000 < 2000)$$

$$= P(|X - 10000| < 2000) = 1 - P(|X - 10000| \ge 2000)$$

$$\ge 1 - \frac{1}{(\sqrt{2000})^2} \quad \because 1 \text{ Theorem } 17 \land k = \frac{2000}{\sigma} = \sqrt{2000}$$

$$= \frac{1999}{2000}$$

2. Probability that more than 15000 letters come in:

$$\begin{split} P(X > 15000) &= P(X - 10000 > 15000 - 10000) \\ &= P(X - 10000 > 5000) \\ &\leq P(X - 10000 > 5000) + P(X - 10000 < 5000) \\ &\leq P(|X - 10000| > 5000) \\ &\leq \frac{1}{\left(\frac{5000}{\sqrt{2000}}\right)^2} = \frac{2000}{5000^2} \end{split}$$

# 5.2 Moment Generating Function

Moment generating functions are important because they uniquely define the distribution of an rv.

### Definition 26 (Moment Generating Function)

If X is an rv, then  $M_X(t) = E(e^{tx})$  is called the moment generating function (mgf) of X provided this expectation exists for all  $t \in (-h, h)$  for some h > 0.

### 66 Note

When determining the mgf of an rv, the values of t for which the expectation exists must always be stated. The range of t where the expectation is defined is "essentially" the radius of convergence.

### Exercise 5.2.1 (Example 2.25 (2.9.2 (1) of the course notes))

Find the mgf of  $X \sim \Gamma(\alpha, \beta)$ . Make sure you specify the domain on which the mgf is defined.

### Solution

Note that the pdf of the Gamma distribution is:

$$f(x) = \begin{cases} \frac{1}{\beta^{\alpha}\Gamma(\alpha)}x^{\alpha-1}e^{-\frac{X}{\beta}} & x > 0\\ 0 & \text{otherwise} \end{cases}$$

*Therefore* 

$$\begin{split} M_X(t) &= E(e^{tx}) = \int_0^\infty e^{tx} \frac{1}{\beta^\alpha} x^{\alpha - 1} e^{-\frac{x}{\beta}} \, dx \\ &= \frac{1}{\beta^\alpha} \int_0^\infty \frac{1}{\Gamma(\alpha)} x^\alpha e^{-x\left(\frac{1}{\beta} - t\right)} \, dx \\ &= \frac{\left(\frac{\beta}{1 - t\beta}\right)^\alpha}{\beta^\alpha} \underbrace{\int_0^\infty \frac{1}{\left(\frac{\beta}{1 - t\beta}\right)^\alpha \Gamma(\alpha)} x^{\alpha - 1} e^{-\frac{x}{\frac{\beta}{1 - t\beta}}} \, dx}_{\text{sum over all values for pdf of } \Gamma(\alpha, \frac{\beta}{1 - t\beta} = 1)} \quad \text{for } \frac{1}{\beta} - t > 0 \end{split}$$

$$&= (1 - t\beta)^{-\alpha} \qquad \text{for } t < \frac{1}{\beta}$$

### Definition 27 (Indicator Function)

The function  $\mathbb{1}_A$  is called the **indicator function** of the set A, i.e.

$$\mathbb{1}_{A} = \begin{cases}
1 & \text{if A occurs} \\
0 & \text{if } A^{C} \text{ occurs}
\end{cases}$$
(5.1)

# Example 5.2.1

The pdf

$$f(x) = \begin{cases} \frac{1}{\theta} & 0 \le x \le \theta \\ 0 & otherwise \end{cases}$$

can be represented as

$$f(x) = \frac{1}{\theta} \mathbb{1}_{\{0 \le x \le \theta\}}$$

### **Example 5.2.2 (Example 2.26)**

Find the mgf of  $X \sim Poi(\lambda)$ . Make sure you specify the domain on which the mgf is defined.

### Solution

*Note that the pmf of X is* 

$$f_X(x) = \frac{e^{-\lambda} \lambda^x}{x!} \mathbb{1}_{\{0,1,2,\dots\}}$$

The mgf is thus

$$M_X(t) = E(e^{tX}) = \sum_{x=0}^{\infty} e^{tx} \frac{e^{-\lambda} \lambda^x}{x!}$$
$$= e^{-\lambda} \sum_{x=0}^{\infty} \frac{(e^t \lambda)^x}{x!} = e^{-\lambda} e^{e^t \lambda}$$
$$= e^{\lambda(e^t - 1)} \quad \forall t \in \mathbb{R}$$

# • Proposition 18 (Properties of the MGF)

Suppose X is an rv. Then

- 1.  $M_X(0) = 1$
- 2. Suppose the derivatives  $M_X^{(k)}(t)$ , for k = 1, 2, ..., exists for  $t \in (-h, h)$  for some h > 0, then the Maclaurin Series<sup>1</sup> of  $M_X(t)$  is

<sup>1</sup> The Maclaurin series is the Taylor expansion around 0.

$$M_X(t) = \sum_{k=0}^{\infty} \frac{M_X^{(k)}(t)\Big|_{t=0}}{k!} t^k$$

3. If the mgf exists, then the  $k^{th}$  moment of X is:

$$E(X^k) = \frac{d^k M_X(t)}{dt^k} \Big|_{t=0}$$

4. Putting 2 and 3 together, we have

$$M_X(t) = \sum_{k=0}^{\infty} \frac{E(X^k)}{k!} t^k$$

The final item shows why  $M_X(t)$  is called the moment generating function.

# Proof

1. 
$$M_X(t)\Big|_{t=0} = E(e^{tX})\Big|_{t=0} = E(e^0) = 1$$

2. This is simply a result of using the Maclaurin series.

3. Note that

$$E(e^{tX}) = E\left[1 + tX + \frac{1}{2}(tX)^2 + \frac{1}{3!}(tX)^3 + \ldots\right]$$
$$= 1 + tE(X) + \frac{t^2}{2}E(X^2) + \frac{t^3}{3!}E(X^3) + \ldots$$

So

$$\frac{d^k}{dt^k} E(e^{tX}) \Big|_{t=0} = \frac{k!}{k!} E(X^k) + \underbrace{\frac{k! \cdot t}{(k+1)!} E(X^{k+1}) + \dots}_{=0 \text{ when } t=0} \Big|_{t=0} = E(X^k)$$

### **Example 5.2.3 (Example 2.27)**

A discrete random variable X has the pmf

$$f(x) = \left(\frac{1}{2}\right)^{x+1} \mathbb{1}_{\{0,1,2,\dots\}}$$

Derive the mgf of X and use it calculate its mean and variance.

$$M_X(t) = \sum_{x=0}^{\infty} e^{tx} \left(\frac{1}{2}\right)^{x+1}$$

$$= \frac{1}{2} \cdot \sum_{x=0}^{\infty} \left(\frac{e^t}{2}\right)^x$$

$$= \frac{1}{2} \cdot \frac{1}{1 - \frac{e^t}{2}} \quad for \left|\frac{e^t}{2}\right| < 1 \text{ or } t < \ln 2$$

$$= \frac{1}{2 - e^t}$$

To get the first two moments,

$$E(X) = \frac{d}{dt} M_X(t) \Big|_{t=0}$$

$$= \frac{e^t}{(2 - e^t)^2} \Big|_{t=0}^{=} 1$$

$$E(X^2) = \frac{d^2}{dt^2} M_X(t) \Big|_{t=0}$$

$$= \frac{e^t}{(2 - e^t)^2} + \frac{2e^t}{(2 - e^t)^3} \Big|_{t=0}$$

$$= 1 + 2 = 3$$

Thus we have that the expected value and variance are

$$E(X) = 1$$
  
 $Var(X) = E(X^2) - E(X)^2 = 3 - 1 = 2$ 

respectively.

# **5.2.1** MGF of a Linear Transformation

### **■** Theorem 19 (MGF of a Linear Transformation)

Suppose the rv X has an mgf  $M_X(t)$  defined for  $t \in (-h,h)$  for some h > 0. Let Y = aX + b, where  $a, b \in \mathbb{R}$  and  $a \neq 0$ . Then the mgf of Y is

$$M_Y(t) = e^{bt} M_X(at), \quad |t| \le \frac{h}{|a|}.$$
 (5.2)

### Proof

Observe that

$$M_Y(t) = E(e^{tY}) = E(e^{t(aX+b)}) = E(e^{atX}e^{tb}) = e^{bt}M_X(at).$$

The range of t is

$$|at| < h \iff |t| < \frac{h}{|a|}$$

# Example 5.2.4 (Example 2.28)

Consider  $X \sim \text{Unif}(\theta_1, \theta_2)$ . Find the mgf of Y = 5X + 3.

### Solution

Note that

$$M_X(t) = \int_{\theta_1}^{\theta_2} \frac{e^{tx}}{\theta_2 - \theta_1} dx$$

$$= \begin{cases} \frac{e^{tx}}{t(\theta_2 - \theta_1)} \Big|_{\theta_1}^{\theta_2} & t \neq 0 \\ 1 & t = 0 \end{cases}$$

$$= \begin{cases} \frac{e^{t\theta_2} - e^{t\theta_1}}{t(\theta_2 - \theta_1)} & t \neq 0 \\ 1 & t = 0 \end{cases}$$

Thus by P Theorem 19,

$$M_Y(t) = e^{3t} M_X(5t) = egin{cases} e^{3t} rac{e^{5t heta_2} - e^{5t heta_1}}{5t( heta_2 - heta_1)} & t 
eq 0 \ 1 & t = 0 \end{cases}$$

#### 5.2.2 *Uniqueness of the MGF*

### Theorem 20 (Uniqueness of the MGF)

Suppose the rv X has mgf  $M_X(t)$  and the rv Y has mgf  $M_Y(t)$ . Suppose also that  $M_X(y) = M_Y(t)$  for all  $t \in (-h,h)$  for some h > 0. Then X and Y have the same distribution, that is,  $\forall s \in \mathbb{R}$ ,

$$P(X < s) = F_X(s) = F_Y(s) = P(Y < s)$$

### Proof

The proof of this theorem is not trivial. See this comment on Math SE for information. It appears that the 2nd bullet point points to a material that I might be able to understand. If I can find that material, and understand it, I may change this proof section to become my own notes.

### Example 5.2.5 (Example 2.29)

Suppose  $X \sim \text{Unif}(0,1)$ . Define  $Y = -2 \log X$ , and use the mgf method to show that  $Y \sim \chi_2^2$ .

( Hint: Find mgf of  $\chi_2$  and show that Y has the same mgf )

### Solution

Let  $Z = \chi_2^2$ . The pdf of Z is therefore

$$f_Z(z) = \frac{1}{2}e^{-\frac{z}{2}}\mathbb{1}_{\{z>0\}}.$$

Then

$$\begin{split} M_Z(t) &= E(e^{tZ}) = \int_0^\infty e^{tz} \frac{1}{2} e^{-\frac{z}{2}} \, dz \\ &= \frac{1}{2} \int_0^\infty e^{(t - \frac{1}{2})z} \, dz \\ &= \begin{cases} \frac{1}{2} \frac{1}{t - \frac{1}{2}} e^{(t - \frac{1}{2})z} \Big|_{z = 0}^\infty & t \neq \frac{1}{2} \\ \infty & t = \frac{1}{2} \end{cases} \\ &= \frac{1}{2t - 1} \qquad t \neq \frac{1}{2} \end{split}$$

# 6 *Lecture 6 May 17th 2018*

# 6.1 Joint Distributions

# 6.1.1 Introduction to Joint Distributions

### 66 Note (Motivation)

Most studies collect information for multiple variables per subject rather than just one variable. Because these variables may interfere/interact with each other and hence give us results that may not be fully reliant on a single variable, it is in our interest to study the interaction of these variables.

To start off with the basics, we will first look at the bivariate case of a joint distribution.

# 6.1.2 Joint and Marginal CDFs

### **Definition 28 (Joint CDF)**

Suppose X and Y are rvs defined on a sample space S. The **joint cdf** of X and Y is given by

$$\forall (x,y) \in \mathbb{R}^2$$
  $F(x,y) = P(X \le x, Y \le y).$ 

### 66 Note

- Depending on whether X and Y are both discrete or both continuous, we can derive the joint pmf or joint pdf of (X, Y), respectively.
- Definition 28 only concerns two variables (a bivariate case), but we can certainly extend the idea to a k-dimensional joint cdf for the rvs  $X_1, X_2, ..., X_k$  as  $\forall (x_1, x_2, ..., x_k) \in \mathbb{R}^k$ ,

$$F(x_1, x_2, ..., x_k) = P(X_1 \le x_1, X_2 \le x_2, ..., X_k \le x_k).$$

### • Proposition 21 (Properties of Joint CDF)

Suppose X, Y are rvs, either both continuous or discrete, and has a joint cdf F. Then

- 1. F is non-decreasing in x for fixed y.
- 2. F is non-decreasing in y for fixed x.
- 3.  $\lim_{x\to-\infty} F(x,y) = 0$  and  $\lim_{y\to-infty} F(x,y) = 0$ .
- 4.  $\lim_{(x,y)\to(-\infty,-\infty)} F(x,y) = 0$  and  $\lim_{(x,y)\to(\infty,\infty)} F(x,y) = 1$

### Proof

1. Suppose not, i.e. that we have instead that F is decreasing for x. Then for  $x_1 < x_2 \in \mathbb{R}$ , we would have

$$F(x_1, y) > F(x_2, y)$$

$$\implies P(X \le x_1, Y \le y) > P(X \le x_2, Y \le y)$$

In other words,

$$P(\{(w,v): (w,v) \in S, X(w) \le x_1, Y(v) \le y\})$$
  
>  $P(\{(w,v): (w,v) \in S, X(w) \le x_2, Y(v) \le y\})$ 

However, note that for fixed y, since  $x_1 < x_2$ , we must have that

$$\{(w,v) \in S : X(w) \le x_1, Y(v) \le y\}$$
  
 
$$\subseteq \{(w,v) \in S : X(w) \le x_2, Y(v) \le y\}.$$

By ♠ Proposition 1, we have that

$$P(\{(w,v): (w,v) \in S, \ X(w) \le x_1, \ Y(v) \le y\})$$
  
 
$$\le P(\{(w,v): (w,v) \in S, \ X(w) \le x_2, \ Y(v) \le y\}).$$

This is clearly a contradiction.

- 2. The proof for this statement is similar to the above.
- 3. Note that

$$\lim_{x \to -\infty} F(x, y) = \lim_{x \to -\infty} P(X \le x, Y \le y)$$

$$= P(X \le -\infty, Y \le y)$$

$$= P([X \le -\infty] \cap [Y \le y])$$

$$= P(\emptyset \cup [Y \le y]) = P(\emptyset) = 0$$

*The proof for the case where*  $y \to -\infty$  *is similar.* 

*4.* This is simply a consequence of 3.

#### 66 Note

We say that F is a joint cdf if it satisfies all the conditions in 6 Proposition 21.1 Many literature actually claims this, and it does look like it will be assumed so for this class.

#### Example 6.1.1 (Example 3.1)

Consider the following joint cdf of two rvs  $(X_1, X_2)$ :

$$F(x_1, x_2) = \begin{cases} 0 & x_1 < 0 \lor x_2 < 0 \\ 0.49 & 0 \le x_1 < 1 \land 0 \le x_2 < 1 \\ 0.7 & 0 \le x_1 < 1 \land x_2 > 1 \\ 0.7 & x_1 \ge 1 \land 0 \le x_2 < 1 \\ 1 & x_1 \ge 1 \land x_2 \ge 1 \end{cases}$$

Flipping an unfair coin with  $P({H}) = 0.3$  twice independently, we define

*for* i = 1, 2

$$X_i = egin{cases} 1 & ext{if the } i^{th} ext{ flip is heads} \ 0 & ext{otherwise} \end{cases}$$

The joint cdf of  $(X_1, X_2)$  is the given F above. Verify that under this experiment, F is indeed a cdf.

### Solution

Note that conditions 3 and 4 of • Proposition 21 are automatically satisfied by the definition of F.

incomplete example

### Definition 29 (Marginal CDF)

For the rvs X, Y with joint cdf F, the marginal cdf of X is

$$F_X(x) = P(X \le x) = \lim_{y \to \infty} F(x, y) = F(x, \infty) \quad \forall x \in \mathbb{R}$$

and the marginal cdf of Y is

$$F_Y(y) = P(Y \le y) = \lim_{x \to \infty} F(x, y) = F(\infty, y) \quad \forall y \in \mathbb{R}$$

Note that the marginal cdf is defined for both discrete and continuous cases.

# Example 6.1.2

Based on Example 6.1.1, derive  $F_{X_i}(x_i)$  for i = 1, 2.

### Solution

$$F_{X_1}(x_1) = \lim_{x_2 \to \infty} F(x_1, x_2)$$

$$= \begin{cases} 0 & x_1 < 0 \\ 0.7 & 0 \le x_1 < 1 \\ 1 & x_1 \ge 1 \end{cases}$$

The solution for  $F_{X_2}(x_2)$  is similar.

#### Joint Discrete RVs 6.1.3

### Definition 30 (Joint Discrete RV)

Suppose X amd Y are rvs defined on a sample space S. If S is discrete then X and Y are discrete rvs. The joint pmf of X and Y is given by

$$\forall (x,y) \in \mathbb{R}^2 \quad f(x,y) = P(X=x,Y=y).$$

The set  $A = \{(x,y) : f(x,y) > 0\}$  is called the support set of (X,Y).

### • Proposition 22 (Properties of Joint PMF)

Suppose X, Y are discrete rvs with joint pmf f and support set A. Then

1. 
$$\forall (x,y) \in \mathbb{R}^2$$
  $f(x,y) \ge 0$ 

$$2. \sum_{(x,y)\in A} \sum f(x,y) = 1$$

3. 
$$\forall R \subset \mathbb{R}^2$$
,

$$P[(X,Y) \in R] = \sum_{(x,y) \in R} f(x,y)$$

The proof is analogous to the univariate case as seen in 6 Proposition 6

### Example 6.1.3 (Example 3.2)

Consider the following joint pmf where the numbers inside the table show P(X = x, Y = y). Find c. Then, calculate  $P(X + Y \le 2)$ .

	<i>x</i> = -2	x = 0	x = 2
<i>y</i> = <i>o</i>	0.05	0.1	0.15
<i>y</i> = 1	0.07	0.11	С
<i>y</i> = 2	0.02	0.25	0.05

### Solution

Since the sum of all the probabilities must be 1, thus

$$c = 1 - 0.05 - 0.07 - 0.02 - \dots - 0.15 - 0.05 = 0.2.$$

*Notice that the only cases where* X + Y > 2 *is when* 

- X = 2, Y = 1; and
- X = 2, Y = 2.

Thus

$$P(X + Y \le 2) = 1 - P(X = 2, Y = 1) - P(X = 2, Y = 2)$$
  
= 1 - 0.2 - 0.05 = 0.75

### Example 6.1.4 (Example 3.3)

A small college has 90 male and 30 female professors. An ad hoc committee of 5 is selected at random to write the vision and mission of the college. Let X and Y be the number of men and women in this committee, respectively. Derive the joint distribution of (X,Y).

#### Solution

Observe that the support set of this distribution is

$$A = \{(x,y) : x + y = 5, x, y = 0, 1, 2, 3, 4, 5\}.$$

We have that the distribution is

$$P(X = x, Y = y) = \begin{cases} \frac{\binom{90}{x}\binom{30}{y}}{\binom{120}{5}} & x, y = 0, 1, 2, 3, 4, 5\\ \frac{120}{5} & x + y = 5 \end{cases}$$

$$0 & otherwise$$

### Definition 31 (Marginal Distribution - Discrete Case)

Suppose X and Y are discrete rvs with joint pf f. Then the marginal pf of X is

$$\forall x \in \mathbb{R}^2 \quad f_X(x) = P(X = x) = \sum_{y \in \mathbb{R}} f(x, y),$$

and the **marginal pf** of Y is

$$\forall y \in \mathbb{R}^2$$
  $f_Y(y) = P(Y = Y) = \sum_{x \in \mathbb{R}} f(x, y).$ 

### Example 6.1.5 (Example 3.4)

Consider the joint pmf from Example 6.1.3. Find the marginal distributions,

i.e. marginal pmfs of X and Y.

	<i>x</i> = -2	x = 0	x = 2
y = o	0.05	0.1	0.15
<i>y</i> = 1	0.07	0.11	0.2
<i>y</i> = 2	0.02	0.25	0.05

### Solution

Using the definition, we have that

$$f_X(x) = \sum_{y \in \mathbb{R}} f(x, y) = \begin{cases} 0.14 & x = -2 \\ 0.46 & x = 0 \\ 0.40 & x = 2 \end{cases}$$

and

$$f_Y(y) = \sum_{x \in \mathbb{R}} f(x, y) = \begin{cases} 0.3 & y = 0 \\ 0.38 & y = 1 \\ 0.32 & y = 2 \end{cases}$$

### Example 6.1.6 (Example 3.5)

Suppose that a penny and a nickel are each tossed 10 times so that every pair of sequences of tosses (n tosses in each sequence) is equally likely to occur. Let X be the number of heads obtained with the penny, and Y be the number of heads obtained with the nickel. It can be shown that (show it!) the joint pmf of X and Y is as follows.

$$P(X = x, Y = y) = \begin{cases} \binom{10}{x} \binom{10}{y} \left(\frac{1}{2}\right)^{20} & x, y = 0, ..., 10 \\ 0 & otherwise \end{cases}$$

### Solution

*Note that the support set of X and Y are the same, i.e.* 

$$A_X = A_Y = \{0, 1, ..., 10\}.$$

We may assume that the penny and the nickel are fair coins, i.e. if we let  $p_x$  and  $p_y$  be the probability of getting a head for a penny and nickel, respectively, then  $p_x = p_y = \frac{1}{2}$ . Since there are 10 ways to get x heads with the penny, and similarly so for the nickel, we have that

$$P(X = x, Y = y) = \begin{cases} \binom{10}{x} \binom{10}{y} \left(\frac{1}{2}\right)^{10} & x, y = 0, 1, ..., 10 \\ 0 & otherwise \end{cases}$$
$$= \begin{cases} \binom{10}{x} \binom{10}{y} \left(\frac{1}{2}\right)^{20} & x, y = 0, 1, ..., 10 \\ 0 & otherwise \end{cases}$$

as required.

### 66 Note

6.1.4

It is interesting to observe that the two rvs in the last example have seemingly no relationship with one another in terms of the experiment conducted, since they do not affect each other. This leads us to introducing the next concept.

### Independence of Discrete RVs

### Definition 32 (Independence of Discrete RVs)

Two rvs X and Y with joint cdf F are said to be **independent** if and only if

$$\forall x, y \in \mathbb{R}$$
  $F(x,y) = F_X(x)F_Y(y)$ 

### Theorem 23 (Independence by PF)

Suppose X and Y are rvs with joint cdf F, joint pf f, marginal cdf  $F_X$  and  $F_Y$  respectively, and marginal pf  $f_X$  and  $f_Y$  respectively. Also, suppose that  $A_X = \{x : f_X(x) > 0\}$  is the support set of X and  $A_Y = \{y : f_X(x) > 0\}$  $f_Y(y) > 0$  is the support set of Y. Then X and Y are independent ros if and only if either

$$\forall (x,y) \in A_X \times A_Y \quad f(x,y) = f_X(x)f_Y(y)$$

holds, or

$$\forall x, y \in \mathbb{R}$$
  $F(x, y) = F_X(x)F_Y(y)$ 

### Proof

The ( $\Longrightarrow$ ) direction is simply a result of Clairaut's Theorem<sup>2</sup>. While the  $(\Leftarrow)$  direction is a direct result of applying double integrals. 

### Example 6.1.7 (Example 3.6)

Suppose X and Y are discrete rvs with joint pf

$$f(x,y) = \frac{\theta^{x+y}e^{-2\theta}}{x!\nu!} \mathbb{1}_{\{x,y=0,1,\dots\}}.$$

Are X and Y independent of each other?

### Solution

Note that we may write f as

$$f(x,y) = \left(\frac{\theta^x e^{-\theta}}{x!} \cdot \frac{\theta^y e^{-\theta}}{y!}\right) \mathbb{1}_{\{x,y=0,1,\dots\}}$$

and so this suggests that we can indeed break down f into two parts, each only affected by x and y respectively, "indeppdent" of each other. Indeed,

I am not certain as to why this is presented as a theorem that repeats the definition. As so, the prove for the 2nd equation will not be shown.

- <sup>2</sup> Work needs to be done to show that our statement actually satisfies the condition for Clairaut's Theorem to apply. Clairaut's Theorem states that:
  - Theorem 24 (Clairaut's Theorem) If  $(x_0, y_0)$  is a point in the domain of a function f with
  - f is defined on all points in an open disk centered at  $(x_0, y_0)$ ;
  - the first partial derivatives,  $f_{xy}$  and  $f_{yx}$  are all continuous for all points in the open disk.

since

$$f_X(x) = \sum_{y=0}^{\infty} \frac{\theta^{x+y}e^{-\theta}}{x!y!} \mathbb{1}_{\{x,y=0,1,\dots\}}$$

$$= \sum_{y=0}^{\infty} \left(\frac{\theta^x e^{-\theta}}{x!} \cdot \frac{\theta^y e^{-\theta}}{y!}\right) \mathbb{1}_{\{x=0,1,\dots\}}$$

$$= \frac{\theta^x e^{-\theta}}{x!} \sum_{y=0}^{\infty} \frac{\theta^y e^{-\theta}}{y!}$$

$$= \frac{\theta^x e^{-\theta}}{x!}$$

$$= \frac{\theta^x e^{-\theta}}{x!}$$

$$= \frac{\theta^x e^{-\theta}}{x!}$$

Similarly, we can obtain

$$f_Y(y) = \frac{\theta^y e^{-\theta}}{y!}$$

Multiplying  $f_X(x)$  and  $f_Y(y)$  together, we indeed get back to the original joint pmf.

# 7 *Lecture 7 May 24th 2018*

# 7.1 *Joint Distributions (Continued)*

# 7.1.1 Independence of Discrete RVs (Continued)

# Example 7.1.1 (Example 3.7)

Consider the joint pmf below from Example 6.1.3. Are X and Y independent? Prove or disprove.

	<i>x</i> = -2	x = 0	x = 2	P(Y=y)
y = 0	0.05	0.1	0.15	0.3
<i>y</i> = 1	0.07	0.11	0.2	0.38
<i>y</i> = 2	0.02	0.25	0.05	0.32
P(X=x)	0.14	0.46	0.4	

# Solution

*Note that* 

$$P(X = -2, Y = 0) = 0.5 \ but$$
 
$$P(X = -2)P(Y = 0) = 0.14 \cdot 0.3 = 0.042 \neq 0.5.$$

Thus X and Y are not independent.

# 7.1.2 *Joint Continuous RVs*

Two random variables X and Y are said to be **jointly continuous** if there exists a function f(x,y) such that the joint cdf of X and Y can be written as

$$\forall (x,y) \in \mathbb{R}^2 \quad F(x,y) = \int_{-\infty}^x \int_{-\infty}^y f(t_1,t_2) d_{t_2} d_{t_1}.$$

The function f is called the **joint density function** of X and Y. It follows from the above defintiion that when the second partial derivative exists, we have

$$f(x,y) = \frac{\partial^2}{\partial x \partial y} F(x,y)$$

The set  $\{(x,y): f(x,y) > 0\}$  is called the support set of (X,Y).

# **66** Note (Convention)

Define f(x,y) = 0 when  $\frac{\partial^2}{\partial x \partial y} F(x,y)$  does not exist.

## Example 7.1.2 (Example 3.8)

Suppose X and Y have joint pdf  $f(x,y) = \mathbb{1}_{\{0 < x,y < 1\}} = \mathbb{1}_{\{0 < x < 1,0 < y < 1\}}$ . Calculate the joint cdf of X and Y.

## Solution

$$F(x,y) = \begin{cases} 0 & x \le 0, \forall y \le 0 \\ \int_0^x \int_0^y 1 \, ds \, dt = xy & 0 < x < 1 \, \land 0 < y < 1 \\ \int_0^1 \int_0^y 1 \, ds \, dt = y & x \ge 1 \, \land 0 < y < 1 \\ \int_0^x \int_0^1 1 \, ds \, dt = x & 0 < x < 1 \, \land y \ge 1 \\ \int_0^1 \int_0^1 1 \, ds \, dt = 1 & x \ge 1 \, \land y \ge 1 \end{cases}$$

## • Proposition 25 (Properties of Joint PDF)

1. 
$$\forall (x,y) \in \mathbb{R}^2 \quad f(x,y) \ge 0$$

$$2. \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(x, y) \, dx \, dy = 1$$

3. 
$$\forall B \subset \mathbb{R}^2$$
, 
$$P[(X,Y) \in B] = \int\limits_{(x,y) \in B} \int f(x,y) \, dx \, dy$$



# Example 7.1.3 (Example 3.9)

Suppose that  $f(x,y) = Kxy \cdot \mathbb{1}_{\{0 < x, y < 1\}}$  for some constant K > 0. Find Kso that f is a valid joint pdf. If X and Y have the joint density f, calculate P(X > Y).

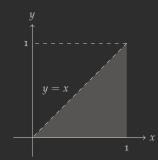
## Solution

Note that

$$1 = \int_0^1 \int_0^1 Kxy \, dx \, dy = \frac{K}{4}.$$

Thus K = 4. To solve the next part, observe that for X > Y, we have the diagram to the right to show the support set of the joint distribution. The shaded region is the support set. We then have

$$P(X > Y) = \int_0^1 \int_0^x 4xy \, dy \, dx = \int_0^1 2xy^2 \Big|_0^x \, dx$$
$$= \int_0^1 2x^3 \, dx = \frac{1}{2}x^3 \Big|_0^1 = \frac{1}{2}$$



# Example 7.1.4 (Example 3.10)

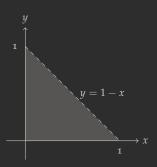
Suppose that

$$f(x,y) = \begin{cases} Cxy & 0 < x, y < 1, x + y < 1 \\ 0 & otherwise \end{cases}$$

Find C so that f(x,y) is a valid joint probability density function, and calculate  $P(Y^2 < X)$ .

# Solution

Note that the diagram on the right shows the support set of (X,Y). To find

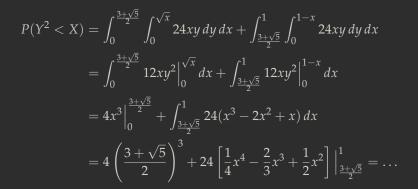


С,

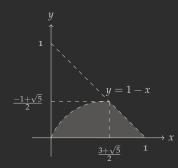
$$\begin{split} 1 &= \int_0^1 \int_0^{1-x} Cxy \, dy \, dx = \int_0^1 \frac{C}{2} xy^2 \Big|_0^{1-x} \, dx \\ &= C \int_0^1 \frac{1}{2} x(x^2 - 2x + 1) \, dx = C \int_0^1 \frac{1}{2} (x^3 - 2x^2 + x) \, dx \\ &= C \left( \frac{1}{8} x^4 - \frac{1}{3} x^3 + \frac{1}{4} x^2 \right) \Big|_0^1 = C \left( \frac{3}{24} - \frac{8}{24} + \frac{6}{24} \right) = \frac{C}{24}. \end{split}$$

And so C = 24.

To calculate  $P(Y^2 < X)$ , note the diagram to the right. Then



We shall not proceed to get the final solution since it is a messy process and the result is not important.



Solve for y = 1 - x and  $y^2 = x$  to get the intersection.

#### Marginal Distribution (Continuous) 7.1.3

# Definition 34 (Marginal PDF)

Suppose X and Y are continuous rvs with joint pdf f. Then the marginal pdf of X is given by

$$\forall x \in \mathbb{R} \quad f_X(x) = \int_{-\infty}^{\infty} f \, dy,$$

and the marginal pdf of Y is

$$\forall y \in \mathbb{R} \quad f_Y(y) = \int_{-\infty}^{\infty} f \, dx.$$

## **Example 7.1.5 (Example 3.11)**

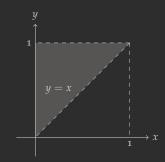
Suppose X and Y have joint pdf  $f(x,y) = K(x+y)\mathbb{1}_{0 \le x \le y \le 1}$  for some constant K. Find K. Then, calculate the marginal density of X.

# Solution

A diagram showing the region of the support set is on the right.

To get K,

$$1 = \int_0^1 \int_x^1 K(x+y) \, dy \, dx = \int_0^1 \left( Kxy + \frac{1}{2} Ky^2 \right) \Big|_x^1 \, dx$$
$$= \int_0^1 Kx + \frac{K}{2} - Kx^2 - \frac{1}{2} Kx^2 \, dx$$
$$= \frac{K}{2} \left( x^2 + x - x^3 \right) \Big|_0^1 = \frac{K}{2}$$



Thus K = 2.

To get the marginal density of X, note that our joint pdf is now the following:

$$f(x,y) = 2(x+y)\mathbb{1}_{\{0 \le x < y \le 1\}}$$

Thus

$$\int_{x}^{1} 2(x+y) \, dy = 2xy + y^{2} \Big|_{x}^{1} = 2x + 1 - 3x^{2}$$

And hence

$$f_X(x) = \begin{cases} -3x^2 + 2x + 1 & 0 \le x \le 1\\ 0 & \text{otherwise} \end{cases}$$

# 7.1.4 Independence of Continuous RVs

# Definition 35 (Independence of Continuous RVs)

Two random variables X and Y with joint cdf F and joint pdf f are independent iff

$$\forall x, y \in \mathbb{R} \quad F(x, y) = F_X(x)F_Y(y)$$

 $or^1$ 

$$\forall x, y \in \mathbb{R} \quad f(x, y) = f_X(x) f_Y(y).$$

<sup>1</sup> It's really an "AND"

## 66 Note

A necessary, but insufficient, condition for X and Y to be independent is that

$$\mathrm{supp}(X,Y) = \mathrm{supp}(X) \times \mathrm{supp}(Y)$$

# **Example 7.1.6 (Example 3.12)**

Are random variables X and Y introduced in Example 7.1.5 independent? Explain.

## Solution

Recall that the pdf was given as

$$f(x,y) = 2(x+y) \mathbb{1}_{\{1 \le x < y \le 1\}}.$$

We derived the marginal pdf of X in the earlier example:

$$f_X(x) = (-3x^2 + 2x + 1) \mathbb{1}_{\{0 \le x \le 1\}}.$$

To get the marginal pdf of Y, note

$$\int_0^y 2(x+y) \, dx = x^2 + 2xy \Big|_0^y = 3y^2.$$

Thus

$$f_Y(y) = egin{cases} 3y^2 & 0 \leq y \leq 1 \ 0 & otherwise. \end{cases}$$

*Note that* 

$$f_X(x)f_Y(y) = -9x^2y^2 + 6xy^2 + 3y^2$$
  $0 \le x < y \le 1$ 

which is not equal to f. Thus, X and Y are not independent.

# 8 Lecture 8 May 29th 2018

# 8.1 Joint Distributions (Continued 2)

# 8.1.1 Independence of Continuous RVs (Continued)

# Example 8.1.1 (Example 3.12 (3.4.8 course note))

Suppose X and Y are continuous with joint pdf

$$f(x,y) = \frac{3}{2}y(1-x^2)\mathbb{1}_{\{-1 \le x \le 1\}}\mathbb{1}_{\{0 \le y \le 1\}}$$

Are X and Y independent?

# Solution

The marginal pdf of X is

$$f_X(x) = \int_0^1 \frac{3}{2} y(1 - x^2) \, dy = \frac{3}{4} y^2 (1 - x^2) \Big|_0^1$$

$$= \begin{cases} \frac{3}{4} (1 - x^2) & -1 \le x \le 1 \\ 0 & \text{otherwise} \end{cases}$$

The marginal pdf of Y is

$$f_Y(x) = \int_{-1}^{1} \frac{3}{2} y (1 - x^2) dx = \frac{3}{2} y \left( x - \frac{1}{3} x^3 \right) \Big|_{-1}^{1}$$

$$= \begin{cases} 2y & 0 \le y \le 1 \\ 0 & otherwise. \end{cases}$$

Clearly, we have

$$f_X(x)f_Y(y) = \frac{3}{2}y(1-x^2) = f(x,y) - 1 \le x \le 1, 0 \le y \le 1.$$

Thus X and Y are independent.

# Theorem 26 (Factorization Theorem for Independence)

Suppose X and Y are rvs with joint pf f, and marginal pf  $f_X$  and  $f_Y$ , respectively. Suppose also that

$$A = \{(x,y) : f(x,y) > 0\}$$
 is the support set of  $(X,Y)$   
 $A_X = \{x : f_X(x) > 0\}$  is the support set of  $X$ , and  
 $A_Y = \{y : f_Y(y) > 0\}$  is the support set of  $Y$ 

Then X and Y are independent rvs iff  $A = A_X \times A_Y$  and there exist non-negative functions g and h such that

$$f(x,y) = g(x)h(y)$$

for all  $(x,y) \in A_X \times A_Y$ .

## Proof

The  $\implies$  direction is straightforward: Since X and Y are independent, we have that  $f = f_X f_Y$ , and so clearly,  $A = A_X \times A_Y$  and so  $\forall (x, y) \in A = A_X \times A_Y$ , we have that  $f_X$  and  $f_Y$  are non-negative.

*For the* ← *direction, note that* 

$$f_Y(y) = \int_{x \in A_X} g(x)h(y) dx = h(y) \int_{x \in A_X} g(x) dx$$
  
$$f_X(x) = \int_{y \in A_Y} g(x)h(y) dy = g(x) \int_{y \in A_Y} h(x) dy.$$

Thus,

$$f_X(x)f_Y(y) = g(x)h(y) \int_{x \in A_X} g(x) \, dx \int_{y \in A_Y} h(y) \, dy$$
  
=  $g(x)h(x) \int_{x \in A_X} \int_{y \in A_Y} g(x)h(y) \, dy \, dx = g(x)h(y)$ 

where line 2 is by linearity of integration. Thus  $f(x,y) = f_X(x)f_Y(y)$ . Thus X and Y are independent.  $\Box$ 

### 66 Note

1. If  $\blacksquare$  Theorem 26 holds, then  $f_X$  will be proportional to g and  $f_Y$  will be proportional to h. Clearly so, since

$$g(x) \cdot h(y) = f_X(x) f_Y(y)$$
  
$$g(x) \propto f_X(x) \wedge h(y) \propto f_Y(y)$$

2. The definitions and theorems can be easily extended to the random vector  $(X_1, X_2, ..., X_n)$ . Indeed, if we apply mathematical induction on the proof above, we will be able to get our desired result.1

<sup>1</sup> I wonder if this statement is equivalent to the Fisher-Neyman Factorization Theorem.

#### Conditional Distributions 8.1.2

## Definition 36 (Conditional Distributions)

Suppose X and Y are rvs with joint pf f, and marginal pfs  $f_X$  and  $f_Y$ , respectively. Suppose also that  $A = \{(x,y) : f(x,y) > 0\}$ . The **conditional pf** of X given Y = y is given by

$$f_X(x|y) = \frac{f(x,y)}{f_Y(y)}$$

for  $(x,y) \in A$  provided that  $f_Y(y) \neq 0$ . The **conditional pf** of Y given X = x is given by

$$f_Y(y|x) = \frac{f(x,y)}{f_X(x)}$$

for  $(x,y) \in A$  provided that  $f_X(x) \neq 0$ .

#### Remark

If X and Y are discrete rvs then

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

and

$$\sum_{x} f_X(x|y) = \sum_{x} \frac{f(x,y)}{f_Y(y)} = \frac{1}{f_Y(y)} \sum_{x} f(x,y) = \frac{f_Y(y)}{f_Y(y)} = 1,$$

and similarly so for  $f_Y(y|x)$ . Similarly, if X and Y are both continuous rvs, then

$$\int_{-\infty}^{\infty} f_X(x|y) \, dx = \int_{-\infty}^{\infty} \frac{f(x,y)}{f_Y(y)} \, dx = \frac{1}{f_Y(y)} \int_{-\infty}^{\infty} f(x,y) \, dx = \frac{f_Y(y)}{f_Y(y)} = 1$$

Now consider if X is a continuous rv such that  $f_X(x) \neq P(X = x)$  and P(X = x) = 0 for all x. Then to justify the definition of the conditional pdf of Y given X = x, when X and Y are both continuous rvs, we consider  $P(Y \leq y|X = x)$  using the limit approach:

$$\begin{split} P(Y \leq y | X = x) &= \lim_{h \to 0} P(Y \leq y | x \leq X \leq x + h) \\ &= \lim_{h \to 0} \frac{\int_{x}^{x+h} \int_{-\infty}^{y} f(u,v) \, dv \, du}{\int_{x}^{x+h} f_{X}(u) \, du} \\ &= \lim_{h \to 0} \frac{\frac{d}{dh} \int_{x}^{x+h} \int_{-\infty}^{y} f(u,v) \, dv \, du}{\frac{d}{dh} \int_{x}^{x+h} f_{X}(u) \, du} \quad \text{by L'Hôpital's Rule} \\ &= \lim_{h \to 0} \frac{\int_{-\infty}^{y} \frac{d}{dh} \int_{x}^{x+h} f(u,v) \, du \, dv}{\frac{d}{dh} \int_{x}^{x+h} f_{X}(u) \, du} \quad (1) \\ &= \lim_{h \to 0} \frac{\int_{-\infty}^{y} f(x+h,v) \, dv}{f_{X}(x+h)} \quad (2) \\ &= \frac{\int_{-\infty}^{y} f(x,v) \, dv}{f_{X}(x)} \end{split}$$

where (1) is by assuming that the integrands are all convergent so that we may interchange the integral signs and the differential operator, and (2) by the Fundamental Theorem of Calculus. If we differentiate the last line with respect to y, by the Fundamental Theorem of Calculus, we have

$$\frac{d}{dy}P(Y \le y|X = x) = \frac{f(x,y)}{f_X(x)}$$

which justifies the using of our definition

$$f_Y(y|x) = \frac{f(x,y)}{f_X(x)}.$$

# Example 8.1.2

A fair coin is flipped 10 times independently.

- 1. What is the distribution of Y, the number of heads in 10 flips?
- 2. Suppose the first 4 flips have all landed on tails. What is the distribution of Y given this information?

#### Solution

- 1. Clearly, we know that  $Y \sim Bin(10, \frac{1}{2})$ .
- 2. Since each flip is independent of each other and the first four flips have already been determined, the range of values for Y changes from  $\{0,...,10\}$ to  $\{0,...,6\}$ . Since the experiment is still essentially the same, we have that

 $Y \mid first \mid 4 flips are tails \sim Bin \left(6, \frac{1}{2}\right)$ .

# Example 8.1.3 (Example 3.13)

Consider the experiment carried out in Example 8.1.2. Let

X := number of heads in the first 4 flips

Y := number of heads in 10 flips

Derive the conditional distribution of Y given that the first 4 flips landed on heads, i.e. derive the distribution for Y|X=4.

#### Solution

Let W be the number of heads in the last 6 flips. Then W has the same distribution as in part 2 of our earlier example. Also,  $X \sim \text{Bin}\left(4,\frac{1}{2}\right)$ Clearly, Y = X + W. We proceed to derive the joint pf of X and Y:

$$P(X = x, Y = y) = P(X = x, X + W = y) = P(X = x, W = y - x)$$

$$= P(X = x)P(W = y - x) \quad \text{by Independence}$$

$$= {4 \choose x} \left(\frac{1}{2}\right)^x \left(\frac{1}{2}\right)^{4-x} \cdot {6 \choose y - x} \left(\frac{1}{2}\right)^{y-x} \left(\frac{1}{2}\right)^{6-y+x}$$

$$= {4 \choose x} {6 \choose y - x} \left(\frac{1}{2}\right)^{10}$$

Then

$$\begin{split} P(Y|X=4) &= \frac{P(X=4, Y=y)}{P(X=4)} \\ &= \frac{\binom{4}{4}\binom{6}{y-4}\left(\frac{1}{2}\right)^{10}}{\binom{4}{4}\left(\frac{1}{2}\right)^{4}\left(\frac{1}{2}\right)^{0}} \\ &= \binom{6}{y-4}\left(\frac{1}{2}\right)^{6} \quad y \in \{4, 5, ..., 10\}. \end{split}$$

We may also re-label the conditional distribution to have

$$\begin{pmatrix} 6 \\ y^* \end{pmatrix} \left(\frac{1}{2}\right)^6 \quad y^* \in \{0, ..., 6\}$$

# **Example 8.1.4 (Example 3.14)**

From Example 7.1.5, we had that

$$f(x,y) = 2(x+y) \mathbb{1}_{\P 0 \le x \le y \le 1}$$

and the marginal density of X is

$$f_X(x) = (2x - 3x^2 + 1)\mathbb{1}_{\{0 \le x < 1\}}.$$

Derive the conditional distribution of  $Y|X = \frac{1}{2}$ .

### Solution

Observe that

$$f(y|X = \frac{1}{2}) = \frac{f\left(\frac{1}{2}, y\right)}{f_X\left(\frac{1}{2}\right)} = \frac{2\left(\frac{1}{2} + y\right)}{2\left(\frac{1}{2}\right) - 3\left(\frac{1}{2}\right)^2 + 1} = \frac{8}{13}(1 + 2y)$$

for  $\frac{1}{2} < y \le 1$ .

# • Proposition 27 (Properties of Conditional Distributions)

Let X and Y be rvs. If both X and Y are discrete, then

- $\sum_{x} f(x|y) = 1$ ;
- $F(x|y) = \sum_{\{w:w \le x\}} f(w|y)$ ; and
- $f(x|y) = F(x|y) F(x^-|y)$ .

If X and Y are both continuous, then

- $\int_{\mathcal{X}} f(x|y) dx = 1$ ;
- $F(x|y) = \int_{-\infty}^{x} f(t|y) dt$ ; and
- $f(x|y) = \frac{\partial}{\partial x} F(x|y)$

Exercise 8.1.1

Prove • Proposition 27.

#### ■ Theorem 28 (Product Rule)

Suppose X and Y are rvs with joint pf f, marginal pfs  $f_X(x)$  and  $f_Y(y)$  respectively, and conditional pfs  $f_X(x|y)$  and  $f_Y(y|x)$  respectively. Then

$$f(x,y) = f_X(x|y)f_Y(y) = f_Y(y|x)f_X(x).$$

# Proof

Notice once and for all that by rearranging the definition of conditional distribution

$$f_X(x|y)f_Y(y) = f(x,y) = f_Y(y|x)f_X(x)$$

# • Proposition 29 (Independence from Conditionality)

Suppose X and Y are rvs with marginal pfs  $f_X(x)$  and  $f_Y(y)$  respectively, and conditional pfs  $f_X(x|y)$  and  $f_Y(y|x)$  respectively. Let  $A_X = \{x : x \in A_X = \{x : x \in A_X = x \in$  $f_X(x) > 0$  and  $A_Y = \{y : f_Y(y) > 0\}$ . X and Y are independent rvs iff either of the following holds:

$$\forall x \in A_X \quad f_X(x|y) = f_X(x)$$

or

$$\forall y \in A_Y \quad f_Y(y|x) = f_Y(y).$$

# Proof

Suppose that X and Y are independent rvs. Then

$$f(x,y) = f_X(x)f_Y(y).$$

Then

$$f_X(x|y) = \frac{f(x,y)}{f_Y(y)} = \frac{f_X(x)f_Y(y)}{f_Y(y)} = f_X(x)$$
$$f_Y(y|x) = \frac{f(x,y)}{f_X(x)} = \frac{f_X(x)f_Y(y)}{f_Y(y)} = f_Y(y).$$

for  $x \in A_X$  and  $y \in A_Y$ .

Suppose WLOG that  $f_X(x|y) = f_X(x)$ . Thus by  $\blacksquare$  Theorem 28

$$f(x,y) = f_X(x|y)f_Y(y) = f_X(x)f_Y(y)$$
 for  $x \in A_X$ ,  $y \in A_Y$ .

# **Example 8.1.5 (Example 3.15)**

In a game of chance, a random number is generated from  $P \sim \text{Beta}(\alpha, \beta)$ . Given P = p, a coin with  $P(\{H\}) = p$  is flipped independently n times, where the player is rewarded the same amount of dollars as the number of heads in n. Calculate the probability that a random player earns at least \$1 in this game.

#### Solution

Let X be the number of heads that appear in n flips, which equates to the total amount of \$1 earned. Then

$$X | P = p \sim Bin(n, p)$$

However, note that

$$P(X \ge 1) = P(earn \ at \ least \ \$1) = 1 - P(earn \ nothing) = 1 - P(X = 0)$$

To get P(X = x), we need to do the following: note that the support set of P is from 0 to 1, then

$$Pr(X = x) = \int_0^1 Pr(X = x, P = p) dp$$

$$= \int_0^1 Pr(X = x | P = p) Pr(P = p) dp \qquad by 1 \text{ Theorem 28}$$

$$= \int_0^1 \binom{n}{x} p^x (1 - p)^{n-x} \frac{p^{\alpha - 1} (1 - x)^{\beta - 1}}{B(\alpha, \beta)} dp$$

$$= \binom{n}{x} \frac{1}{B(\alpha, \beta)} \int_0^1 p^{\alpha + x - 1} (1 - p)^{\beta + n - x - 1} dp$$

$$= \binom{n}{x} \frac{B(\alpha + x, \beta + n - x)}{B(\alpha, \beta)} \int_0^1 \underbrace{\frac{p^{\alpha + x - 1} (1 - p)^{\beta + n - x - 1}}{B(\alpha + x, \beta + n - x)}}_{pdf \text{ of Beta}(\alpha + x, \beta + n - x)} dp$$

$$= \binom{n}{x} \frac{B(\alpha + x, \beta + n - x)}{B(\alpha, \beta)}$$

Therefore,

$$\begin{split} P(X \geq 1) &= 1 - P(X = 0) = 1 - \binom{n}{0} \frac{\Gamma(\alpha)\Gamma(\beta + n)}{\Gamma(\alpha + \beta + n)} \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \\ &= 1 - \frac{\Gamma(\alpha + \beta)\Gamma(\beta + n)}{\Gamma(\alpha + \beta + n)\Gamma(\beta)} \end{split}$$

# Definition 37 (Beta Distribution)

*If*  $X \sim \text{Beta}(\alpha, \beta)$ , then

$$f(x) = \frac{x^{\alpha - 1}(1 - x)^{\beta - 1}}{B(\alpha, \beta)}$$

where

$$B(\alpha, \beta) = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha + \beta)}$$

# Definition 38 (Joint Expectation)

Suppose h(x,y) is a real-valued function. If X and Y are discrete rvs with joint pf f and support A, then

$$E[h(x,y)] = \sum_{(x,y)\in A} \sum h(x,y) f(x,y).$$

provided that the joint sum converges absolutely.

If X and Y are continuous rvs with joint pf f, then

$$E[h(x,y)] = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} h(x,y) f(x,y) \, dx \, dy$$

provided that the joint integral converges absolutely.

This is also known as the Law of the Unconscious Statistician for two rvs.

# **Example 8.1.6 (Example 3.16)**

Consider X and Y with the following joint probability distribution. Calculate E(XY).

	x = -2	x = o	x = 2
<i>y</i> = 0	0.05	0.1	0.15
<i>y</i> = 1	0.07	0.1	0.2
<i>y</i> = 2	0.02	0.25	0.05

#### Solution

$$E(XY) = \sum_{x} \sum_{y} xy f(x, y)$$

$$= -2(1)(0.07) - 2(2)(0.02) + 2(1)(0.2) + 2(2)(0.05)$$

$$= 0.38$$

# 9 Lecture 9 May 31st 2018

# 9.1 *Joint Distributions* (Continued 3)

# *9.1.1 Joint Expectations (Continued)*

# **■** Theorem 30 (Linearity of Expectation in Bivariate Case)

Suppose X and Y are two rvs with joint pf f,  $a_i$ ,  $b_i$ , for i = 1, ..., n, are constants, and  $g_i(x, y)$ , for i = 1, ..., n, are real-valued functions. Then

$$E\left[\sum_{i=1}^{n} (a_{i}g_{i}(X,Y) + b_{i})\right] = \sum_{i=1}^{n} (a_{i}E[g_{i}(X,Y)]) + \sum_{i=1}^{n} b_{i}$$

provided that  $E[g_i(X,Y)]$  is finite for i = 1,...,n.

## Proof

*This is simply an extension of* **P** *Theorem 13.* 

# **■** Theorem 31 (Implication of Independence on Joint Expectation)

If X and Y are independent rvs with joint pf f, and g(x) and h(y) are real valued functions, then

$$E[g(X)h(Y)] = E[g(X)]E[h(Y)]$$

#### Proof

We shall prove for the discrete case and leave the continuous case for future exercises.

Observe that

$$E[g(X)h(Y)] = \sum_{x} \sum_{y} g(x)h(y)f(x,y) \quad \therefore \Phi \text{ Definition 38}$$

$$= \sum_{x} \sum_{y} g(x)h(y)f_{X}(x)f_{Y}(y)$$

$$= \sum_{x} g(x)f_{X}(x) \sum_{y} h(y)f_{Y}(y)$$

$$= E[g(X)]E[h(Y)]$$

where  $f_X(x)$  and  $f_Y(y)$  are the marginal pfs of X and Y respectively.  $\square$ 

We may repeatedly apply the above proof for n rvs through induction and get the following result.

# **■** Theorem 32 (Generalized Implication of Independence on Joint Expectation)

If  $X_1, X_2, ..., X_n$ , for some  $n \in \mathbb{N}$ , are independent rvs and  $h_1, h_2, ..., h_n$  are real valued functions, then

$$E\left[\prod_{i=1}^n h_i(X_i)\right] = \prod_{i=1}^n E[h_i(X_i)].$$

### 9.1.2 Covariance

INDEPENDENCE of two rvs X and Y implies that knowledge of the value of X does not provide any information whatsoever about the distribution of Y. Essentially, we can say that there is no "relationship" between X and Y. In statistics, **linear relationships** are often the subject of interest. The strength of a linear relationship is related to **covariance** and measured by the **correlation coefficient**, usually denoted by  $\rho$ .

#### Exercise 9.1.1

*Prove* **P** *Theorem 31 for the continuous case.* 

It can be shown that when *X* and *Y* have no linear relationship iff their covariance is 0.

On a related thought, does covariance relate to independence? If so, how?

# Definition 39 (Covariance)

The covariance of rvs X and Y is given by

$$Cov(X,Y) = E[(X - \mu_X)(Y - \mu_Y)] = E(XY) - \mu_X \mu_Y.$$
 (9.1)

where  $\mu_X = E[X]$  and  $\mu_Y = E[Y]$ .

If Cov(X, Y) = 0, then X and Y are called uncorrelated rvs.

#### 66 Note

Note that the 2nd and 3rd term are equivalent in Equation (9.1) since

$$\begin{split} E[(X - \mu_X)(Y - \mu_Y)] &= E[XY - \mu_X Y - \mu_Y X + \mu_X \mu_Y] \\ &= E[XY] - \mu_X E[Y] - \mu_Y E[X] + \mu_X \mu_Y :: 1 \ \, \textit{Theorem 30} \\ &= E[XY] - \mu_X \mu_Y - \mu_X \mu_Y + \mu_X \mu_Y = E[XY] - \mu_X \mu_Y \end{split}$$

From here, it is easy to see from  $\blacksquare$  Theorem 31, that since E[XY] = $E[X]E[Y] = \mu_X \mu_Y$ , we have that the independence of X from Y will *imply that* Cov(X, Y) = 0.

However, the converse of the above is **not true**.

## Example 9.1.1

Source: Stats SE

Let X be an rv that it is -1 or 1 with probability 0.5. Then let Y be an rv such that Y = 0 if X = -1, and Y is randomly -1 or 1 with probability 0.5 if X = 1.

Clearly X and Y are highly dependent (since knowing Y allows me to

perfectly know X). They both have zero mean:

$$E[X] = -1\left(\frac{1}{2}\right) + 1\left(\frac{1}{2}\right) = 0$$
$$E[Y] = -1\left(\frac{1}{2}\right) + 1\left(\frac{1}{2}\right) = 0$$

and

$$E[XY] = (-1) \cdot 0P(X = -1, Y = 0) + 1(-1) \cdot P(X = 1, Y = -1)$$
$$+ 1(1)P(X = 1, Y = 1)$$
$$= -\frac{1}{4} + \frac{1}{4} = 0$$

Thus Cov(X, Y) = 0

Or more generally, take any distribution P(X) and any P(Y|X) such that P(Y=a|X)=P(Y=-a|X) for all X (i.e., a joint distribution that is symmetric around the x axis), and you will always have zero covariance. But you will have non-independence whenever  $P(Y|X) \neq P(Y)$ , i.e., the conditionals are not all equal to the marginal, and vice versa for symmetry around the y axis.

## 66 Note

- If Cov(X, Y) = 0, then X and Y are called uncorrelated rvs.
- By definition, Cov(X, X) = Var(X), since

$$Cov(X, X) = E[(X - \mu_X)^2] = Var(X)$$

## Example 9.1.2 (Example 3.17)

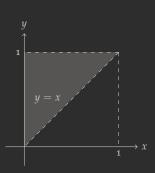
Consider the rvs X and Y with the joint pdf

$$f(x,y) = \begin{cases} 2 & 0 \le x \le y \le 1 \\ 0 & otherwise \end{cases}$$

Calculate Cov(X, Y).

# Solution

Observe the diagram of the support set of X and Y to our right.



Then we can calculate

$$E[XY] = \int_0^1 \int_x^1 2xy \, dy \, dx = \int_0^1 xy^2 \Big|_x^1 \, dx$$

$$= \int_0^1 x - x^3 \, dx = \frac{1}{2} - \frac{1}{4} = \frac{1}{4}$$

$$f_X(x) = \int_x^1 2 \, dy = \begin{cases} 2 - 2x & 0 \le x \le 1 \\ 0 & \text{otherwise} \end{cases}$$

$$E[X] = \int_0^1 x(2 - 2x) \, dx = 1 - \frac{2}{3} = \frac{1}{3}$$

$$f_Y(y) = \int_0^y 2 \, dx = \begin{cases} 2y & 0 \le y \le 1 \\ 0 & \text{otherwise} \end{cases}$$

$$E[Y] = \int_0^1 2y^2 \, dy = \frac{2}{3}$$

Thus

$$Cov(X,Y) = \frac{1}{4} - \frac{1}{3} \left(\frac{2}{3}\right) = \frac{1}{36}$$

We observe that the covariance is positive. This implies a positive linear relationship. However, we cannot tell from this value the strength of the relationship between X and Y.

## **Example 9.1.3 (Example 3.18)**

Consider rvs X and Y with the joint pf

$$f(x,y) = \frac{xy}{7} \mathbb{1}_{\{y=1,2\}} \mathbb{1}_{\{x=0,\dots,y\}}.$$

Calculate Cov(X, Y).

## Solution

The following table captures all the probabilities that can be found from the given pf

	x = 0	<i>x</i> = 1	<i>x</i> = 2
<i>y</i> = 1	0	<u>1</u> 7	0
<i>y</i> = 2	o	<u>2</u> 7	$\frac{4}{7}$

Observe that we thus have

$$f_X(x) = \begin{cases} \frac{3}{7} & x = 1\\ \frac{4}{7} & x = 2 \end{cases}$$

$$E[X] = \frac{3}{7} + 2\left(\frac{4}{7}\right) = \frac{11}{7}$$

$$f_Y(y) = \begin{cases} \frac{1}{7} & y = 1\\ \frac{6}{7} & y = 2 \end{cases}$$

$$E[Y] = \frac{1}{7} + 2\left(\frac{6}{7}\right) = \frac{13}{7}$$

Also

$$E[XY] = \frac{1}{7} + 2\left(\frac{2}{7}\right) + 4\left(\frac{4}{7}\right) = 3$$

Therefore,

$$Cov(X,Y) = E[XY] - E[X]E[Y] = 3 - \frac{11}{7} \frac{13}{7} = \frac{4}{49}$$

## Theorem 33 (Variance of Linear Combinations)

Suppose X and Y are rvs and a, b, c are real constants. Then

$$Var(aX + bY + c) = a^{2} Var(X) + b^{2} Var(Y) + 2ab Cov(X, Y)$$

## Proof

Let 
$$E[aX + bY + c] = \mu$$
. Observe that

$$Var[aX + bY + c]$$

$$= E \left[ [(aX + bY + c) - \mu]^2 \right]$$

$$= E \left[ (aX + bY + c)^2 - 2\mu(aX + bY + c) + \mu^2 \right]$$

$$= E \left[ a^2X^2 + abXY + acX + abXY + b^2Y^2 + bcY + c^2 - 2\mu(aX + bY + c) + \mu^2 \right]$$

$$= a^2E[X^2] + 2abE[XY] + acE[X] + b^2E[Y^2] + bcE[Y] + c^2 - \mu^2$$

*Note that* 

$$\mu^{2} = E [aX + bY + c]^{2}$$

$$= (aE[X] + bE[Y] + c)^{2}$$

$$= a^{2}E[X]^{2} + b^{2}E[Y]^{2} + 2abE[X]E[Y] + acE[X] + bcE[Y] + c^{2}$$

Therefore, we have that

$$\begin{aligned} & \operatorname{Var}[aX + bY + c] \\ &= a^2 E[X^2] + a^2 E[X]^2 + b^2 E[Y^2] - b^2 E[Y]^2 + 2ab E[XY] - 2ab E[X] E[Y] \\ &= a^2 \left( E[X^2] - E[X]^2 \right) + b^2 \left( E[Y^2] - E[Y]^2 \right) + 2ab \left( E[XY] - E[X] E[Y] \right) \\ &= a^2 \operatorname{Var}(X) + b^2 \operatorname{Var} Y + 2ab \operatorname{Cov}(X, Y) \end{aligned}$$

as required.

By applying P Theorem 33 repeatedly, we have the following generalized theorem.

# ■ Theorem 34 (Generalized Variance of Linear Combinations)

Suppose  $X_1, X_2, ..., X_n$  are rvs with  $Var(X_i) = \sigma_i^2$ , and  $a_1, a_2, ..., a_n$  are real constants. Then

$$\operatorname{Var}\left(\sum_{i=1}^{n} a_{i} X_{i}\right) = \sum_{i=1}^{n} a_{i}^{2} \sigma_{i}^{2} + 2 \sum_{i=1}^{n-1} \sum_{j=i+1}^{n} a_{i} a_{j} \operatorname{Cov}(X_{i}, X_{j})$$

Note that the prove the above, we have to also use the fact that

$$Cov(X_i, X_i) = Cov(X_i, X_i)$$

#### 66 Note

Note that in  $\blacksquare$  Theorem 34, if the rvs are independent rvs, then  $Cov(X_i, X_i) =$ 0 for  $i \neq j$ , thus wiping off the 2nd term in the equation, leaving us with

$$\operatorname{Var}\left(\sum_{i=1}^{n} a_i X_i\right) = \sum_{i=1}^{n} a_i^2 \sigma_i^2$$

## Example 9.1.4 (Example 3.19)

To build a ship engine piece, suppose two pole-shaped components A and B are attached at one end to each other to make one long pole-shaped component C. Suppose the length of part A is an rv with a mean of 3 inches and a variance of 0.25 inch<sup>2</sup>. Similarly, the length of component B is an rv with a mean of 25 inches and a variance of 0.5 inch<sup>2</sup>.

Find the mean and the variance of the length of part C if

- 1. the lengths of A and B are independent;
- 2. the covariance between lengths of A and B is -0.3 inch<sup>2</sup>.

#### Solution

*Note that we are given that* C = A + B

1. We have that

$$E(C) = E(A + B) = E(A) + E(B) = 3 + 25 = 28.$$

For variance, since A and B are independent, Cov(A, B) = 0, thus

$$Var(C) = Var(A + B) = Var(A) + Var(B) + 2Cov(A, B)$$
  
= 0.25 + 0.5 + 0 = 0.75

2. Since C is a linear equation, the covariance does not affect the expectation and thus we still have

$$E(C) = 28.$$

*Now, given that* Cov(A, B) = -0.3*, we have* 

$$Var(C) = Var(A) + Var(B) + 2Cov(A, B) = 0.75 - 0.6 = 0.15.$$

#### 9.1.3 Correlation

The **covariance** is a real number which depends on the units of measurement of *X* and *Y*. The information part of a covariance is its **sign**, unless if it is used as the context.

To use the covariance as the context, and to quantitatively measure the strength of a linear relationship, which we have discussed and desired before, we use the **correlation coefficient**.

# Definition 40 (Correlation Coefficient)

The correlation coefficient of rvs X and Y is given by

$$\rho(X,Y) = \frac{\mathrm{Cov}(X,Y)}{\sigma_X \sigma_Y}$$

where  $\sigma_X = \sqrt{\operatorname{Var}(X)}$  and  $\sigma_Y = \sqrt{\operatorname{Var}(Y)}$ .

Note that this is the definition of the Pearson Correlation Coefficient. There are other correlation coefficients but we will be using only Pearson, at it seems.

## • Proposition 35 (Properties of the Correlation Coefficient)

Let X and Y be rvs, and  $\rho(X,Y)$  the correlation coefficient of X and Y. Then

- 1.  $|\rho(X,Y)| \leq 1$ ;
- 2. (perfect positive linear relationship)  $\rho(X,Y) = 1 \iff Y = aX + b \text{ for some } a > 0;$
- 3. (perfect inverse linear relationship)  $\rho(X,Y) = -1 \iff Y = aX + b \text{ for some } a < 0.$

# Proof

1. This is somewhat beyond the scope of what we can cover now but we shall use this result presented on Wikipedia:

$$|\text{Cov}(X,Y)| \leq \sqrt{\text{Var}(X)\,\text{Var}(Y)}.$$

*Then given the formula of*  $\rho(X,Y)$ *, the proof is complete:* 

$$|\rho(X,Y)| = \left| \frac{\operatorname{Cov}(X,Y)}{\sqrt{\operatorname{Var}(X)\operatorname{Var}(Y)}} \right| \le \frac{\sqrt{\operatorname{Var}(X)\operatorname{Var}(Y)}}{\sqrt{\operatorname{Var}(X)\operatorname{Var}(Y)}} = 1$$

Proving 2 and 3 is currently outside of my abilities. Refer to this Math SE Q&A for a hint on how to prove this statement.

Consider rvs X and Y with the joint pdf

$$f(x,y) = \begin{cases} 2 & 0 \le x \le y \le 1 \\ 0 & otherwise \end{cases}$$

Calculate  $\rho(X, Y)$ .

#### Solution

The diagram to the right is an illustration of the region of support for X and Y. We now calculate the following values:

$$E(XY) = \int_0^1 \int_x^1 2xy \, dy \, dx = \int_0^1 x - x^3 \, dx = \frac{1}{4}$$

$$f_X(x) = \int_x^1 2 \, dy = 2 - 2x \quad 0 \le x \le 1$$

$$f_Y(y) = \int_0^y 2 \, dx = 2y \quad 0 \le y \le 1$$

$$E(X) = \int_0^1 2x - 2x^2 \, dx = \left(x^2 - \frac{2}{3}x^3\right) \Big|_0^1 = \frac{1}{3}$$

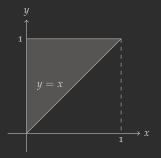
$$E(X^2) = \int_0^1 2x^2 - 2x^3 \, dx = \left(\frac{2}{3}x^3 - \frac{1}{2}x^4\right) \Big|_0^1 = \frac{1}{6}$$

$$E(Y) = \int_0^1 2y^2 \, dy = \frac{2}{3}$$

$$E(Y^2) = \int_0^1 2y^3 \, dy = \frac{1}{2}$$

$$Var(X) = E(X^2) - E(X)^2 = \frac{1}{6} - \left(\frac{1}{3}\right)^2 = \frac{1}{18}$$

$$Var(Y) = E(Y^2) - E(Y)^2 = \frac{1}{2} - \left(\frac{2}{3}\right)^2 = \frac{1}{18}$$



Therefore we have that

$$Cov(X,Y) = E(XY) - E(X)E(Y) = \frac{1}{4} - \frac{2}{9} = \frac{1}{36}$$
$$\sqrt{Var(X) Var(Y)} = \frac{1}{18}$$

and so

$$\rho(X,Y) = \frac{\text{Cov}(X,Y)}{\sqrt{\text{Var}(X)\,\text{Var}(Y)}} = \frac{\frac{1}{36}}{\frac{1}{18}} = \frac{1}{2}$$

# Definition 41 (Conditional Expectation)

Let G be a real-valued function. The **conditional expectation** of g(Y)given X = x, denoted as  $g(Y) \mid (X = x)$  is given by

$$E[g(Y) | x] = \sum_{y} g(y) f_Y(y | x)$$

if Y|(X = x) is a discrete rv and

$$E[g(Y) | x] = \int_{\mathcal{Y}} g(y) f_Y(y | x)$$

if Y|(X=x) is a continuous rv. This definition only holds provided that the sum and the integral converges absolutely. The conditional expectation of h(X) given Y = y, where h is a real-valued function, is defined in a similar manner.

We also call E[Y | X = x] the **conditional mean**, which may be denoted as E(Y | x), and Var(Y | X = x) the conditional variance, which may be denoted as Var(Y | x).

#### 66 Note

Note that there is also the notation E(Y | X), which is an rv, and hence different from E(Y | x).

#### Example 9.1.6 (Example 3.21)

Consider  $f(x,y) = 8xy\mathbb{1}_{\{0 < x < y < 1\}}$ . Calculate the conditional mean and the conditional variance of  $X \mid (Y = \frac{1}{2})$ .

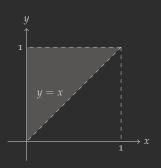
### Solution

The diagram to the right illustrates the region of support for X and Y. To derive the conditional distribution, we first need  $f_Y(y)$ .

$$f_Y(y) = \int_0^y 8xy \, dx = 4x^2y \Big|_0^y = 4y^3$$

Thus

$$f_X\left(X \mid Y = \frac{1}{2}\right) = \frac{f\left(x, \frac{1}{2}\right)}{f_Y\left(\frac{1}{2}\right)} = \frac{4x}{\frac{1}{2}} = 8x \quad 0 < x < \frac{1}{2}$$



Therefore the conditional mean is

$$E\left[X \mid \frac{1}{2}\right] = \int_0^{\frac{1}{2}} 8x^2 dx = \frac{8}{3} \cdot \frac{1}{2^3} = \frac{1}{3},$$

the second conditional moment is

$$E\left[X^2 \mid \frac{1}{2}\right] = \int_0^{\frac{1}{2}} 8x^3 \, dx = 2 \cdot \frac{1}{2^4} = \frac{1}{8},$$

and so the conditional variance is

$$Var\left(X^{2} \mid \frac{1}{2}\right) = E\left[X^{2} \mid \frac{1}{2}\right] + E\left[X \mid \frac{1}{2}\right]^{2} = \frac{1}{8} - \frac{1}{9} = \frac{1}{72}$$

# 10 Lecture 10 Jun 05th 2018

# 10.1 Joint Distribution (Continued 4)

# 10.1.1 Conditional Expectation (Continued)

# Example 10.1.1 (Example 3.22)

Given the joint distribution below, calculate Var(X | Y = 1) and compare it to Var(X).

	<i>x</i> = -2	x = 0	x = 2	P(Y=y)
y = o	0.05	0.1	0.15	0.3
<i>y</i> = 1	0.07	0.11	0.2	0.38
<i>y</i> = 2	0.02	0.25	0.05	0.32
P(X=x)	0.14	0.46	0.4	

# Solution

*Note that* 

$$Var(X) = E(X^{2}) - E(X)^{2}$$

$$= 4 \cdot 0.14 + 4 \cdot 0.4 - (-2 \cdot 0.14 + 2 \cdot 0.4)^{2} = 1.8896$$

To get Var(X | Y = 1), we first need

$$f(X | Y = 1) = \frac{P(X = x, Y = 1)}{P(Y = 1)} = \begin{cases} \frac{0.07}{0.38} = 0.1842 & x = -2\\ \frac{0.11}{0.38} = 0.2895 & x = 0\\ \frac{0.2}{0.38} = 0.5263 & x = 2 \end{cases}$$

Thus

$$Var(X | Y = 1) = E[X^{2} | Y = 1] - E[X | Y = 1]^{2}$$
$$= \frac{14}{19} + \frac{40}{19} - (\frac{13}{19})^{2} = \frac{857}{361} = 2.3740$$

# • Proposition 36 (Independence on Conditional Expectation)

If X and Y are independent rvs then E[g(Y) | x] = E[g(Y)] and E[h(X) | y] = E[h(X)].

## Proof

We shall prove one of the above for the other will follow a similar argument. Also, we shall prove the continuous case and leave the discrete case as an exercise.

Observe that

$$E[g(Y) | X = x] = \int_{y} g(y) \frac{f(x,y)}{f_X(x)} dy$$

$$= \int_{y} g(y) \frac{f_X(x) f_Y(y)}{f_X(x)} dy \quad \because \text{ independence}$$

$$= \int_{y} g(y) f_Y(y) dy = E[g(Y)]$$

Exercise 10.1.1

*Prove the discrete case for* **♦** *Proposition 36.* 

# Theorem 37 (Law of Total Expectation)

Suppose X and Y are rvs, then

$$E(E[g(Y) \mid X]) = E[g(Y)]$$

*If* g *is the identity function, we have* E(E[Y | X]) = E(Y).

#### Proof

We shall prove for the discrete case and leave the continuous case as an exercise. Observe that

$$E[g(Y) | X] = \sum_{y} [g(y) \cdot P(Y = y | X)]$$

$$E[E[g(Y) | X]] = \sum_{x} \left[ \sum_{y} [g(y) \cdot P(Y = y | X)] \right] P(X = x)$$

$$= \sum_{x} \sum_{y} g(y) \cdot P(X = x, Y = y)$$

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

$$= \sum_{y} [g(y) \cdot P(Y = y)] = E[g(Y)]$$

#### Exercise 10.1.2

Theorem 37.

# Example 10.1.2 (Example 3.23 - A Classical Example)

A man is lost in a mine, and 3 paths are in front of him. If he takes path 1, after 3 hours, he will be back at his current place. If he takes path 2, the time to get out of the mine (in hours) follows an Exp(1) distribution. If he takes the 3rd path, he will be back to his current place after 2 hours. Suppose that the man cannot recognize which path he took every time he comes back to the original spot (after going through either path 1 or 3), and so he randomly chooses a path every time he comes back to this original spot. What is the expected time that he will take to get out of the mine?

This is a very classical example to illustrate the power of the Law of Total Expectation.

#### Solution

Let X an rv that represents the path number, i.e. X = 1, 2 or 3, and let Y represent the total time that the man takes to exit the mine. We are given

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

We are also given that

$$E[Y | X = 1] = 3 + E[Y]$$
  
 $E[Y | X = 2] = 1 \quad \because Y | (X = 2) \sim Exp(1)$   
 $E[Y | X = 3] = 2 + E[Y].$ 

Therefore, to get the expected time, by P Theorem 37,

$$E[Y] = E[E[Y \mid X]]$$

$$= \frac{1}{3} \cdot E[Y \mid X = 1] + \frac{1}{3} \cdot E[Y \mid X = 2] + \frac{1}{3} \cdot E[Y \mid X = 3]$$

$$= \frac{1}{3} \cdot (3 + E[Y]) + \frac{1}{3} \cdot 1 + \frac{1}{3} (2 + E[Y])$$

$$= 2 + \frac{2}{3} E[Y]$$

and hence

$$E[Y] = 6$$

# **■** Theorem 38 (Law of Total Variance)

Suppose X and Y are rvs. Then

$$Var(Y) = E[Var(Y \mid X)] + Var[E(Y \mid X)]$$

# Proof

Note that

$$Var(Y|X) = E(Y^{2}|X) - E(Y|X)^{2}$$

$$E[Var(Y|X)] = E[E(Y^{2}|X) - E(Y|X)^{2}]$$

$$= E[Y^{2}] - E[E(Y|X)^{2}]$$

$$= E[Y^{2}] - \left[Var(E(Y|X)) + E(E(Y|X))^{2}\right]$$

$$= Var(Y) - Var(E(Y|X))$$

By rearranging the above, we get

$$Var(Y) = E[Var(Y|X)] + Var[E(Y|X)]$$

# Example 10.1.3 (Example 3.24 (Course Note 3.7.11))

Suppose  $P \sim \text{Unif}(0,0.1)$  and  $Y \mid P = p \sim \text{Bin}(10,p)$ . Find E(Y) and Var(Y).

## Solution

$$E[Y] = E[E[Y|P]] = E[10P] = 10E[P]$$

$$= 10 \cdot \int_{0}^{0.1} \frac{p}{0.1} dp = \frac{10}{0.2} p^{2} \Big|_{0}^{0.1} = 0.05$$

$$Note: E[P^{2}] = \int_{0}^{0.1} \frac{p^{2}}{0.1} dp = \frac{p^{3}}{0.3} \Big|_{0}^{0.1} = \frac{0.001}{0.3} = \frac{1}{300}$$

$$Var(Y) = E[Var(V|P)] + Var[E[Y|P]]$$

$$= E[10P(1-P)] + Var[10P]$$

$$= 10E[P] - 10E[P^{2}] + 100 Var(P)$$

$$= 0.05 - \frac{1}{30} + 100 \left[ E[P^{2}] - E[P]^{2} \right]$$

$$= \frac{1}{60} + 100 \left[ \frac{1}{300} - 0.05^{2} \right] = \frac{1}{60} + 100 \left[ \frac{1}{300} - \frac{1}{400} \right]$$

$$= \frac{1}{60} + \frac{1}{12} = \frac{1}{10}$$

# Joint Moment Generating Functions

## Definition 42 (Joint Moment Generating Functions)

The joint moment generating function of two rvs X and Y is defined as

$$M(t_1, t_2) = E\left(e^{t_1 X + t_2 Y}\right)$$

if the expectation exists for all  $t_1 \in (-h_1, h_1)$  and  $t_2 \in (-h_2, h_2)$  for some  $h_1, h_2 > 0$ .

More generally, if  $X_1, X_2, ..., X_n$  are rvs, then

$$M(t_1, t_2, ..., t_n) = E\left[\exp\left(\sum_{i=1}^n t_i X_i\right)\right]$$

is called the **joint mgf** of  $X_1, X_2, ..., X_n$  is the expectation exists for all  $t_i \in (-h_i, h_i)$  for some  $h_i > 0$ , where i = 1, ..., n.

# Definition 43 (Joint Moments and Marginal MGF)

Given the joint  $mgf M(t_1, t_2)$ , we can calculate the joint moments. In particular,

$$E\left(X^{j}Y^{k}\right) = \frac{\partial^{j+k}}{\partial t_{1}^{j}\partial t_{2}^{k}}M(t_{1},t_{2})\Big|_{(t_{1},t_{2})=(0,0)}$$

If  $M(t_1, t_2)$  exists for all  $t_1 \in (-h_1, h_1)$  and  $t_2 \in (-h_2, h_2)$  for some  $h_1, h_2 > 0$ , then the mdf of X is given by

$$M_X(t) = E\left(e^{tX}\right) = M(t,0) \quad t \in (-h_1, h_1)$$

and the mgf of Y is given by

$$M_Y(t) = E\left(e^{tY}\right) = M(0,t) \quad t \in (-h_2,h_2).$$

# Example 10.1.4 (Example 3.25)

Given the joint distribution below, calculate the joint  $mgf M(t_1, t_2)$ , the first joint moment, E[XY], from the joint mgf, and the marginal mgf of X and that of Y.

$$x = -1$$
  $x = 1$ 
 $y = 1$   $0.5$   $0.3$ 
 $y = 2$   $0.1$   $0.1$ 

#### Solution

Since all probabilities are provided,

$$\begin{split} M(t_1,t_2) &= E\left(e^{t_1X+t_2Y}\right) = \sum_{x} \sum_{y} e^{t_1x+t_2y} P(X=x,Y=y) \\ &= 0.5e^{-t_1+t_2} + 0.3e^{t_1+t_2} + 0.1e^{-t_1+2t_2} + 0.1e^{t_1+2t_2} \\ E(XY) &= \frac{\partial^2}{\partial t_1 \partial t_2} M(t_1,t_2) \Big|_{t_1=0,t_2=0} \\ &= \frac{\partial}{\partial t_1} \left[ 0.5e^{-t_1+t_2} + 0.3e^{t_1+t_2} + 0.2e^{-t_1+2t_2} + 0.2e^{t_1+2t_2} \right] \Big|_{t_1=0,t_2=0} \\ &= -0.5e^{-t_1+t_2} + 0.3e^{t_1+t_2} - 0.2e^{-t_1+2t_2} + 0.2e^{t_1+2t_2} \Big|_{t_1=0,t_2=0} \\ &= -0.2 \\ M_X(t_1) &= M(t_1,0) = 0.5e^{-t_1} + 0.3e^{t_1} + 0.1e^{-t_1} + 0.1^{t_1} \\ &= 0.6e^{-t_1} + 0.4e^{t_1} \\ M_Y(t_2) &= M(0,t_2) = 0.5e^{t_2} + 0.3e^{t_2} + 0.1e^{2t_2} + 0.1e^{2t_2} \\ &= 0.8e^{t_2} + 0.2e^{2t_2} \end{split}$$

# • Proposition 39 (Independence on Joint MGF)

Suppose X and Y are rvs with joint mgf  $M(t_1, t_2)$  which exists  $\forall t_1 \in$  $(-h_1, h_1), t_2 \in (-h_2, h_2), \text{ for some } h_1, h_2 > 0. \text{ Then } X \text{ and } Y \text{ are } Y \text{ are } Y \text{ and } Y \text{ are } Y \text{$ independent rvs iff

$$\forall t_1 \in (-h_1,h_1), t_2 \in (-h_2,h_2) \quad M(t_1,t_2) = M_X(t_1)M_Y(t_2)$$
 where  $M_X(t_1) = M(t_1,0)$  and  $M_Y(t_2) = M(0,t_2)$ .

#### Proof

to be proven later

#### Example 10.1.5 (Example 3.26 (Course Note 3.8.5))

Suppose X and Y are continous rvs with joint pdf

$$f(x,y) = e^{-y} \quad 0 < x < y < \infty$$

Find the joint mdf of X and Y. Are X and Y independent rvs? What is the

marginal mgf of X and Y?

#### Solution

$$\begin{split} M(t_1,t_2) &= E[e^{t_1X+t_2Y}] = \int_0^\infty \int_0^y e^{t_1x+t_2y}e^{-y}\,dx\,dy \\ &= \int_0^\infty \frac{1}{t_1}e^{t_1x+t_2y-y}\Big|_0^y\,dy \\ &= \int_0^\infty \frac{1}{t_1}\left[e^{y(t_2-1)} - e^{y(t_1+t_2-1)}\right]\,dy \\ &= \frac{1}{t_1}\left[\frac{1}{t_2-1}e^{y(t_2-1)} - \frac{1}{t_1+t_2-1}e^{y(t_1+t_2-1)}\right]\Big|_0^\infty \\ &= \frac{1}{t_1}\left[\frac{t_1}{(t_2-1)(t_1+t_2-1)}\right] \\ &= \frac{1}{(t_2-1)(t_1+t_2-1)} \quad t_2 < 1 \wedge t_1 + t_2 < 1 \\ M_X(t_1) &= M(t_1,0) = \frac{1}{t_1-1} \quad t_1 < 1 \\ M_Y(t_2) &= M(0,t_2) = \frac{1}{(t_2-1)^2} \quad t_2 < 1 \end{split}$$

Observe that

$$M_X(t_1)M_Y(t_2) = \frac{1}{(t_1 - 1)(t_2 - 1)^2} \neq M(t_1, t_2)$$

and so by **♦** Proposition 39, X and Y are not independent.

#### Example 10.1.6 (Example 3.27)

Investigate the independence of X and Y in Example 10.1.4 using the mgf method.

## Solution

We had that

$$M_X(t_1) = 0.6^{-t_1} + 0.4e^{t_1}$$
  $t_1 \in \mathbb{R}$   
 $M_Y(t_2) = 0.8e^{t_2} + 0.2e^{2t_2}$   $t_2 \in \mathbb{R}$ .

Since

$$M_X\left(\frac{1}{2}\right)M_Y\left(\frac{1}{2}\right) \neq M\left(\frac{1}{2},\frac{1}{2}\right)$$

we have that X and Y are not independent.

# 11 Lecture 11 Jun 07th 2018

# 11.0.1 Working with Multivariate Cases

Almost everything that has been introduced above can be extended to cases where we have more than just 2 rvs. For example:

# Definition 44 (k-variate CDF)

The *k*-variate CDF, k > 2, rvs  $X_1, ..., X_k$  is defined as

$$F(x_1,...,x_k) = P(X_1 \le x_1, X_2 \le x_2,..., X_k \le x_k).$$

*In the continuous case, we may write* 

$$f(x_1,...,x_k) = \frac{\partial^k}{\partial x_1 \dots \partial x_k} F(x_1,...,x_k).$$

# Definition 45 (k-variate Support Set)

The support set of the distribution for  $X_1, X_2, ..., X_k$  is

$$\{(x_1,...,x_k): f(x_1,...,x_k) > 0\}$$

We also have the following:

# • Proposition 40 (Law of Total Probability - Multivariate)

If  $X_1, ..., X_k$  are continuous rvs, then

$$\int_{x_1} \dots \int_{x_k} f(x_1, ..., x_k) dx_1 \dots dx_k = 1.$$

Should they by discrete, then

$$\sum_{x_1} \dots \sum_{x_k} f(x_1, ..., x_k) = 1$$

# Definition 46 (k-Variate Marginal Distribution)

To get the marginal distribution of a subset of m variables from  $X_1, ..., X_k$  ( $1 \le m \le k$ ), we will sum or integrate over the other ones if they are discrete or continuous, respectively. For example,

$$f(x_1, x_2, x_3) = \int_{x_4} \dots \int_{x_k} f(x_1, ..., x_k) dx_4 \dots dx_k$$

# Definition 47 (k-Variate Joint MGF)

The joint mgf of  $X_1, ..., X_k$  is defined as

$$M(t_1, t_2, ..., t_k) = E\left(e^{t_1X_1 + ... + t_kX_k}\right)$$

#### • Proposition 41 (Independence for Multivariate Cases)

If  $X_1, ..., X_k$  are independent, then

$$f(x_1,...,x_k) = \prod_{i=1}^k f_{X_i}(x_i) \qquad F(x_1,...,x_k) = \prod_{i=1}^k F_{X_i}(x_i)$$
$$M(t_1,...,t_k) = \prod_{i=1}^k M_{X_i}(t_i)$$

THERE ARE many different examples of multivariate distributions. We shall discuss two:

• Multinomial Distribution

#### • Multivariate Normal Distribution

The multinomial distribution is an extension of the binomial distribution to cases where there are more categories than two results. For a multinomial distribution, we have that

- the experiment involves *n* trials, each with *k* categories
- the outcome of trials are independent of each other
- the probability of each category,  $p_i$ , remains the same across ntrials
- $X = (X_1, ..., X_k) \sim \text{Mult}(n, p_1, ..., p_k)$  counts the number of elements in each category among the n trials.

# Definition 48 (Mutlinomial Distribution)

Suppose  $X_1, ..., X_k$  are discrete rvs with joint pf

$$f(x_1,...,x_k) = \frac{n!}{x_1!x_2!...x_k!x_{k+1}!} p_1^{x_1} p_2^{x_2} ... p_k^{x_k} p_{k+1}^{x_{k+1}}$$

where 
$$x_i = 0, ..., n_i, x_{k+1} = n - \sum_{i=1}^k x_i, 0 < p_i < 1, p_{k+1} = 1 - \sum_{i=1}^k p_i$$
, for  $i = 1, ..., k + 1$ .

*Under these conditions,*  $(X_1, ..., X_k)$  *is said to have a multinomial* distribution, and we write  $(X_1,...,X_k) \sim \text{Mult}(n, p_1,...,p_k)$ .

#### 66 Note

*Observe that* Bin(n, p) = Mult(n, p, p).

#### • Proposition 42 (Properties of Multinomial Distribution)

Suppose  $(X_1,...,X_k) \sim \text{Mult}(n,p_1,...,p_k)$ , then

1.  $\forall (t_1,...,t_k) \in \mathbb{R}^k$ , the random vector  $(X_1,...,X_k)$  has joint mgf

$$M(t_1,...,t_k) = E\left(e^{t_1X_1+...t_kX_k}\right) = (p_1e^{t_1}+...+p_ke^{t_k}+p_{k+1})^n$$

2. Any subset of  $X_1, ..., X_{k+1}$  also has a multinomial distribution. In particular,  $X_i \sim \text{Bin}(n, p_i)$ , i = 1, ..., k + 1.

3. If 
$$T = X_i + X_j$$
, for  $i \neq j$ , then  $T \sim \text{Bin}(n, p_i + p_j)$ 

- 4.  $Cov(X_i, X_j) = -np_ip_j$ , for  $i \neq j$
- 5. The conditional distribution of any subset of  $(X_1, ..., X_{k+1})$  given the rest of the coordinates is a multinomial distribution. In particular, the conditional pf of  $X_i$  given  $X_j = x_j$ ,  $i \neq j$ , is

$$X_i \mid X_j = x_j \sim \operatorname{Bin}\left(n - x_j, \frac{p_i}{1 - p_j}\right)$$

6. The conditional distribution of  $X_i$  given  $T = X_i + X_j = t$ , for  $i \neq j$ , is

$$X_i \mid X_i + X_j = t \sim \operatorname{Bin}\left(t, \frac{p_i}{p_i + p_j}\right)$$

WE SHALL look at the bivariate normal distribution so that it is easier to be explained. The same idea can be extended to a multivariate Normal distribution.

# Definition 49 (Bivariate Normal Distribution)

Let  $X_1$  and  $X_2$  be rvs with joint pdf

$$f(x_1, x_2) = \frac{1}{2\pi |\Sigma|^{\frac{1}{2}}} \exp\left\{-\frac{1}{2}(x - \mu)^T \Sigma^{-1}(x - \mu)\right\}$$

$$= \frac{1}{2\pi\sigma_1\sigma_2\sqrt{1 - \rho^2}} \exp\left\{-\frac{1}{2(1 - \rho^2)} \left[\frac{(x_1 - \mu_1)^2}{\sigma_1^2} + \frac{(x_2 - \mu_2)^2}{\sigma_2^2} - \frac{2\rho(x_1 - \mu_1)(x_2 - \mu_2)}{\sigma_1\sigma_2}\right]\right\}$$

where  $(x_1, x_2) \in \mathbb{R}^2$  and

$$x = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$$
,  $\mu = \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}$ ,  $\Sigma = \begin{bmatrix} \sigma_1^2 & \rho \sigma_1 \sigma_2 \\ \rho \sigma_1 \sigma_2 & \sigma_2^2 \end{bmatrix}$ 

where  $\Sigma$  is a nonsingular matrix. Then  $X = (X_1, X_2)^T$  is said to have a bivariate normal distribution, and we write  $X \sim \text{BVN}(\mu, \Sigma)$ .

### • Proposition 43 (Properties of Bivariate Normal Distribution)

Suppose  $X \sim BVN(\mu, \Sigma)$ , where

$$\mu = egin{pmatrix} \mu_1 \ \mu_2 \end{pmatrix}$$
 ,  $\Sigma = egin{bmatrix} \sigma_1^2 & 
ho\sigma_1\sigma_2 \ 
ho\sigma_1\sigma_2 & \sigma_2^2 \end{bmatrix}$ 

1. X has a joint mgf of

$$M(t_1, t_2) = E[\exp\left(t^T X\right)] = E\left(e^{t_1 X_1 + t_2 X_2}\right) = \exp\left(\mu^T t + \frac{1}{2} t^T \Sigma t\right)$$
  
 $\forall (t_1, t_2) \in \mathbb{R}^2.$ 

- 2.  $X_1 \sim N(\mu_1, \sigma_1^2)$  and  $X_2 \sim N(\mu_2, \sigma_2^2)$
- 3.  $Cov(X_1, X_2) = \rho \sigma_1 \sigma_2$  and  $Corr(X_1, X_2) = \rho$  where  $|rho| \le 1$
- 4.  $X_1$  and  $X_2$  are independent rvs iff  $\rho = 0$
- 5.  $c = (c_1, c_2)^T$  is a nonzero vector of constants  $\Longrightarrow$

$$c^T X = \sum_{i=1}^2 c_i X_i \sim N\left(c^T \mu, c^T \Sigma c\right).$$

6. If A is a  $2 \times 2$  nonsingular matrix and b is a  $2 \times 1$  vector, then Y = $AX + b \sim \text{BVN}(A\mu + b, A\Sigma A^T).$ 

$$X_2 \mid X_1 = x_1 \sim N\left(\mu_2 + \rho \frac{\sigma_2}{\sigma_1}(x_1 - \mu_1), \sigma_2^2(1 - \rho^2)\right)$$
$$X_1 \mid X_2 = x_2 \sim N\left(\mu_1 + \rho \frac{\sigma_1}{\sigma_2}(x_2 - \mu_2), \sigma_1^2(1 - \rho^2)\right)$$

8. 
$$(X - \mu)^T \Sigma^{-1} (X - \mu) \sim \chi^2(2)$$

#### Definition 50 (Random Sample (IID))

Rvs  $X_1, ..., X_2$  are said to form a **simple random sample** or are said to be independent and identically distributed (IID) if  $X_1, ..., X_n$  are independent, and  $f_{X_i} = f_{X_i}$ ,  $\forall i \neq j$ .

# Example 11.0.1 (Example 3.28)

Let  $X_1, ..., X_{10}$  be a random sample of standard normal distribution. Let  $Y_1$  denote the number of these variables that are between -1 and 1, let  $Y_2$  denote the number that have absolute value between 1 and 2, and let  $Y_3$  denote the number that have absolute value larger than 2. Calculate:

- 1.  $P(Y_1 \le 2)$
- 2.  $E[Y_2 | Y_1 = 5]$

# 12 Lecture 12 Jun 12th 2018

# 12.1 Functions of Random Variables

# 12.1.1 Transformation of Two or More Random Variables

In earlier lectures we discussed about basic transformations from one random variable to another, for example, from a continuous rv X to Y = g(X). In particular, we methods were presented:

- The direct method, i.e.  $P(Y \le y) = P(g(X) \le y)$ , and taking the derivative of  $P(Y \le y)$  with respect to y.
- Using the MGF of *Y*, and then translate it as the mgf of *X*.

#### 66 Note (Recall)

In Section 2.4.4, we used the following idea to obtain the result that we desire: for rvs X and Y = g(X) where g is some continuous and injective function

$$P(Y \le y) = P(g(X) \le y) = P(X \le g^{-1}(y)) = F_X(g^{-1}(y))$$
$$f_Y(y) = \frac{d}{dy} F_X(g^{-1}(y)) = (g^{-1}(y))' f_X(g^{-1}(y))$$

In this chapter, we will now study the case where we have more than one rv involved. In particular, let X and Y be two continuous rvs with joint pdf f(x, y). Our questions are:

- 1. What is the distribution of  $U = h_1(X, Y)$ ?
- 2. What is the joint distribution of  $U = h_1(X, Y)$  and  $V = h_2(X, Y)$ ?

To answer the first question, we can actually still employ the direct method:

#### Example 12.1.1 (Example 4.1 (Course Notes 4.1.1))

Suppose X and Y are continuous rvs with joint pdf

$$f(x,y) = 3y \mathbb{1}_{0 < x < y < 1}$$

Find the pdf of T = XY.

#### Solution

First, note that<sup>1</sup>

$$P(T \le t) = P(XY \le t) = P\left(Y \le \frac{t}{X}\right)$$

The diagram to the right shows us the support of the joint probability. We observe that if  $t \le 0$ , then  $P(T \le t) = 0$ , and if  $t \ge 1$ , then  $P(T \le t) = 1$ . Now if 0 < t < 1, the region that we are looking for is the shaded region with the label A, and so we consider

$$P(T \le t) = 1 - P(B) = 1 - \int_{B}^{\infty} \int_{B}^{\infty} f(x, y) \, dx \, dy$$

$$\stackrel{(1)}{=} 1 - \int_{\sqrt{t}}^{1} \int_{\frac{t}{y}}^{\sqrt{t}} f(x, y) \, dx \, dy - \int_{\sqrt{t}}^{1} \int_{\sqrt{t}}^{y} f(x, y) \, dx \, dy$$

$$\stackrel{(2)}{=} 1 - \int_{\sqrt{t}}^{1} \int_{\frac{t}{y}}^{y} f(x, y) \, dx \, dy$$

$$= 1 - \int_{\sqrt{t}}^{1} \int_{\frac{t}{y}}^{y} 3y \, dx \, dy$$

$$= 1 - \int_{\sqrt{t}}^{1} 3y \left( y - \frac{t}{y} \right) dy \quad \because FTC$$

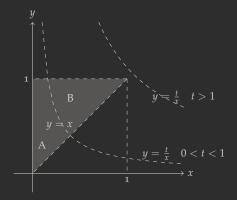
$$= 1 - \frac{\sqrt{t}}{1} \left( 3y^{2} - 3t \right) dy$$

$$= 1 - \left[ y^{3} - 3ty \right] \Big|_{\sqrt{t}}^{1} = 1 - (1 - 3t - \sqrt{t}^{3} + 3t\sqrt{t})$$

$$= 3t - 2t\sqrt{t} \text{ for } 0 < t < 1$$

where for step (1), we broke B into two parts, in particular at  $x = \sqrt{t}$  where y = x and  $y = \frac{t}{x}$  coincide, and step (2) is true by linearity of integration.

¹ I wonder if the last step is actually valid. The support of X definitely includes 0, so the division would not make sense with  $\frac{t}{X}$ . We can, however, still make sense of the event in the 2nd term, which we would have  $P(0 \le t)$ . Should we be concerned about X = 0, or can we neglect that single point given that X is a continuous rv?



With that, i.e. with the CDF of T, we can then obtain

$$f_T(t) = egin{cases} 3 - 3\sqrt{t} & 0 < t < 1 \ 0 & otherwise \end{cases}$$

## Example 12.1.2 (Example 4.2 (Course Note 4.1.2))

*Using the info in Example 12.1.1, find the pdf of T* =  $\frac{X}{Y}$ .

#### Solution

The diagram to the right shows the support of X, Y and the function  $t = \frac{x}{u}$ . We observe that if t = 0, then x = 0, and we would have the line on the axis, and so  $P(T \le t) = 0$ . If t < 0, we would have y = mx where  $m=\frac{1}{t}<0$ , which, regardless of what t<0 is, will not interact with the support of X and Y. So for t < 0,  $P(T \le t) = 0$ . Now if t > 0, we have

$$P(T \le t) = P\left(\frac{X}{Y} \le t\right) = P\left(Y \ge \frac{1}{t}X\right).$$

Consider the case where  $t \geq 1$ , we have that the event would still cover the entire support set of X and Y, and so  $P(T \le t) = 1$  for  $t \ge 1$ . With that, the only remaining case is when 0 < t < 1. In this case,

$$P(T \le t) = \int_0^1 \int_0^{ty} 3y \, dx \, dy = \int_0^1 3ty^2 \, dy = t$$

and so the pdf of T is

$$f_T(t) = egin{cases} 1 & 0 < t < 1 \ 0 & otherwise \end{cases}$$

## Example 12.1.3 (Example 4.3 (Course Note 4.1.3) - Order Statistics)

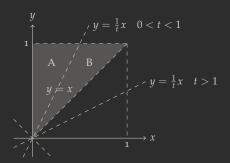
Suppose  $X_1, ..., X_n$  are IID samples, each from a continuous distribution, and with pdf f and cdf F. Find the pdf of

1. 
$$T = \min(X_1, ..., X_n) = X_{(1)}$$

2. 
$$Y = \max(X_1, ..., X_n) = X_{(n)}$$

## Solution

1. For T, we have that its cdf is<sup>2</sup>



<sup>&</sup>lt;sup>2</sup> The use the Law of Total Probability here so that we can use the following argument: "the smallest rv is larger than t, and so rest of the rvs must be the same."

120 Lecture 12 Jun 12th 2018 - Functions of Random Variables

$$\begin{split} P(T \leq t) &= 1 - P(T > t) = 1 - P(\min(X_1, ..., X_n) > t) \\ &= 1 - P(X_1 > t, X_2 > t, ..., X_n > t) \\ &= 1 - \prod_{i=1}^n P(X_i > t) \qquad \because independence \\ &= 1 - \prod_{i=1}^n P(X_1 > t) \qquad \because identical \ distribution \\ &= 1 - P(X_1 > t)^n = 1 - [1 - F_{X_1}(t)]^n \end{split}$$

and so its pdf is

$$f_T(t) = -\frac{d}{dt} [1 - F_{X_1}(t)]^n = -n(-F_{X_1}(t))' [1 - F_{X_1}(t)]^{n-1}$$
$$= nf_{X_1}(t) [1 - F_{X_1}(t)]^{n-1}.$$

Since T relies entirely on  $X_1$  (due to IID), and since we did not have to condition on the values of t, we have that

$$supp(T) = supp(X_1) = supp(X_i)$$
 for  $i = 1, ..., n$ .

2. For Y, we have that its cdf is<sup>3</sup>

$$P(Y \le t) = P(\max(X_1, ..., X_n) \le y) = P(X_1 \le y, ..., X_n \le y)$$

$$= \prod_{i=1}^{n} P(X_i \le y) \quad \because independence$$

$$= \prod_{i=1}^{n} P(X_1 \le y) \quad \because identical \ distribution$$

$$= P(X_1 \le t)^n = F_{X_1}(t)^n$$

and therefore its pdf is

$$f_Y(y) = \frac{d}{dy} F_{X_1}(y)^n = n f_{X_1}(y) F_{X_1}(y)^{n-1}.$$

#### Exercise 12.1.1

From Example 12.1.3, find the joint distribution of  $X_{(1)}$  and  $X_{(n)}$ .

<sup>3</sup> This time, we do not have to employ the Law of Total Probability, because we simply have that "the *largest* rv is smaller than *t*, and so must the rest of the rvs."

12.1.2 One-to-One Bivariate Transformations

Let X and Y be rvs, and  $R_{XY} = \text{supp}[(X,Y)] \in \mathbb{R}^2$ . We define

$$U = h_1(X, Y)$$
  $V = h_2(X, Y)$   
 $S: R_{XY} \to \mathbb{R}^2 \ by \ (x, y) \mapsto (h_1(x, y), h_2(x, y))$ 

The mapping S is called a one-to-one mapping if and only if  $\forall (u,v) \in$  $R_{UV}$ ,  $\exists !(x,y) \in R_{XY}$ ,  $\exists w_1, w_2$  that are functions such that

<sup>4</sup> There is nothing magnificent about this definition, since this is simply the definition of a one-to-one function.

$$x = w_1(u, v)$$
  $y = w_2(u, v)$ 

i.e.  $\exists S^{-1}: R_{UV} \to R_{XY}$  such that  $(u, v) \mapsto (x, y)$ . The Jacobian of the transformation  $S^{-1}$  is

$$\frac{\partial(x,y)}{\partial(u,v)} = \begin{vmatrix} \frac{\partial x}{\partial u} & \frac{\partial x}{\partial v} \\ \frac{\partial y}{\partial u} & \frac{\partial y}{\partial v} \end{vmatrix} = \left| \frac{\partial(u,v)}{\partial(x,y)} \right|^{-1}$$

where  $\frac{\partial(u,v)}{\partial(x,y)}$  is the Jacobian of the transformation S.

# Theorem 44 (One-to-One Bivariate Transformations)

Let X and Y be continuous rvs with joint pdf  $f_{XY}$  and let  $R_{XY} =$  $\{(x,y): f(x,y) > 0\}$  be the support set of (X,Y), and  $R_{UV}$  be the support set of (U,V). Suppose the transformation  $S:R_{XY}\to R_{UV}$ defined by

$$U = h_1(X, Y)$$
  $V = h_2(X, Y)$ 

is a one-to-one transformation, with inverse transformation

$$X = w_1(U, V)$$
  $Y = w_2(U, V)$ .

Then g(u, v), the joint pdf of U and V, is given by

$$\forall (u,v) \in R_{UV} \quad g(u,v) = f(w_1(u,v),w_2(u,v)) \left| \frac{\partial(x,y)}{\partial(u,v)} \right|.$$

This is a generalization of Theorem 8.

## Proof

Let the inverse transformation be labelled as  $S^{-1}: R_{UV} \supset B \rightarrow A \subset$ 

 $R_{XY}$ . Then

$$\int_{B} \int g(u,v) \, dv \, du = P[(U,V) \in B] = P[(X,Y) \in A]$$

$$= \int_{A} \int f(x,y) \, dx \, dy$$

$$= \int_{B} \int f(w_{1}(u,v), w_{2}(u,v)) \left| \frac{\partial(x,y)}{\partial(u,v)} \right| du \, dv$$

where the last step is by the Change of Variables Theorem. And so by comparing integrands, we have

$$\forall (u,v) \in R_{UV} \quad g(u,v) = f(w_1(u,v), w_2(u,v)) \left| \frac{\partial(x,y)}{\partial(u,v)} \right|$$

as required.

## • Proposition 45 (Properties of the Jacobian)

Given the setup in 🗗 Definition 51, we have that

- 1. if S is a linear transformation, i.e.  $\exists a_1, b_1, c_1, a_2, b_2, c_2 \in \mathbb{R}$  such that  $u(x,y) = a_1x + b_1y + c_1$  and  $v(x,y) = a_2x + b_2y + c_2$ , then the Jacobian is a constant;
- 2. if S is a one-to-one transformation, then  $\left|\frac{\partial(x,y)}{\partial(u,v)}\right| \neq 0$

## Proof

1. We have

$$\frac{\partial u}{\partial x} = a_1$$
  $\frac{\partial u}{\partial y} = b_1$   
 $\frac{\partial v}{\partial x} = a_2$   $\frac{\partial v}{\partial y} = b_2$ 

and so

$$|J|=rac{\partial(u,v)}{\partial(x,y)}=egin{array}{c|c} a_1&b_1\ a_2&b_2 \end{array} =a_1b_2-a_2b_1$$

which is a constant, as required.

2. I have no idea how to prove this.

# Example 12.1.4 (Example 4.4 (Course Notes 4.2.4))

Suppose  $X \sim \text{Gam}(a,1)$  and  $Y \sim \text{Gam}(b,1)$  independently. Find the joint pdf of U = X + Y and  $V = \frac{X}{X+Y}$ . Show that  $U \sim \text{Gam}(a+b,1)$  and  $V \sim \text{Beta}(a, b)$ , independently. Find E(V).

#### Solution

Given U = X + Y and  $V = \frac{X}{X+Y}$ , rearranging variables, we have

$$X = UV$$
 and  $Y = U(1 - V)$ 

In order for U to have a Gamma distribution, we need U to be non-negative, which we do since both X and Y have Gamma distributions. Note that the transformation is indeed one to one, since  $\forall (u_1, v_1), (u_2, v_2) \in R_{UV}$  with

$$u_1 = x_1 + y_1$$
  $v_1 = \frac{x_1}{x_1 + y_1}$   
 $u_2 = x_2 + y_2$   $v_2 = \frac{x_2}{x_2 + y_2}$ 

we have that, if we let  $\phi$  denote the transformation,

$$\phi(u_1, v_1) = \phi(u_2, v_2)$$

$$\implies \left(x_1 + y_1, \frac{x_1}{x_1 + y_1}\right) = \left(x_2 + y_2, \frac{x_2}{x_2 + y_2}\right)$$

which then

$$x_1 + y_1 = x_2 + y_2$$

$$\frac{x_1}{x_2 + y_2} = \frac{x_2}{x_2 + y_2}$$

$$\stackrel{Equation (12.1)}{\Longrightarrow} x_1 = x_2$$

$$\stackrel{Equation (12.1)}{\Longrightarrow} y_1 = y_2.$$

We shall now get the Jacobian so that we may use P Theorem 44, so that we may consequently get the distributions for U and V.

$$J = \begin{vmatrix} \frac{\partial x}{\partial u} & \frac{\partial x}{\partial v} \\ \frac{\partial y}{\partial u} & \frac{\partial y}{\partial v} \end{vmatrix} = \begin{vmatrix} v & u \\ 1 - v & -u \end{vmatrix} = -vu - u(1 - v) = -u$$
$$|J| = u$$

By P Theorem 44, and since X and Y are independent, we have

$$f_{U,V}(u,v) = f_{X,Y}(x,y) \cdot |J| = f_X(x)f_Y(y) \cdot |J|$$

$$= \frac{x^{a-1}e^{-x}}{\Gamma(a)} \frac{y^{b-1}e^{-y}}{\Gamma(b)} \cdot |J|$$

$$= \frac{e^{-a}e^{-b}}{\Gamma(a)\Gamma(b)} (uv)^{a-1} u^{b-1} (1-v)^{b-1} u$$

$$= \underbrace{\frac{u^{a+b-1}e^{-(a+b)}}{\Gamma(a+b)}}_{pdf \ of \ Gam(a+b,1)} \cdot \underbrace{\frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} v^{a-1} (1-v)^{b-1}}_{pdf \ of \ Beta(a,b)}$$

We have already shown that U is a non-negative v, and so  $U \sim \text{Gam}(a+b,1)$  as required. Note that for V, we have  $X+Y>X>0 \implies 1> \frac{X}{X+Y}>0 \quad \because X+Y\neq 0$  and so 0< V<1. Therefore  $V\sim \text{Beta}(a,b)$ .

With that, we can look for E[V].

$$E[V] = \int_0^1 v \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} v^{a-1} (1-v)^{b-1} dv$$

$$= \frac{\Gamma(a+b)\Gamma(a+1)}{\Gamma(a)\Gamma(a+b+1)} \int_0^1 \underbrace{\frac{\Gamma(a+b+1)}{\Gamma(a+1)\Gamma(b)} v^a (1-v)^{b-1}}_{pdf \text{ of Beta}(a+1,b)} dv$$

$$= \frac{a\Gamma(a)\Gamma(a+b)}{(a+b)\Gamma(a)\Gamma(a+b)} = \frac{a}{a+b}$$

# 13 Lecture 13 Jun 14th 2018

# 13.1 Functions of Random Variables (Continued)

# 13.1.1 One to One Bivariate Transformations (Continued)

#### Example 13.1.1 (Example 4.5 (Course Notes 4.2.8))

Suppose X and Y are continuous rvs with joint pdf

$$f(x,y) = e^{-x-y} \mathbb{1}_{\{0 < x,y < \infty\}}.$$

Let U = X + Y and V = X - Y. Find the joint pdf of U and V. (Note: Be sure to specify the support of (U, V).) Then, find the marginal pdf of U and V respectively.

# Solution

*Note that because*  $0 < x, y < \infty$ *,* 

$$u = x + y \implies 0 < u < \infty$$
  
 $v = x - y \implies -\infty < v < \infty$ 

The Jacobian is

$$\left| \frac{\partial(u,v)}{\partial(x,y)} \right| = \left| \begin{vmatrix} 1 & 1 \\ 1 & -1 \end{vmatrix} \right| = |-2| = 2$$

and so

$$\left| \frac{\partial(x,y)}{\partial(u,v)} \right| = \left| \frac{\partial(u,v)}{\partial(x,y)} \right|^{-1} = \frac{1}{2}.$$

*Note that* 

$$X = \frac{U - V}{2}$$
  $Y = \frac{U - V}{2}$ 

and so

$$x = \frac{u+v}{2} > 0 \implies u > -v$$
$$y = \frac{u-v}{2} > 0 \implies u > v.$$

The diagram to the right shows the support of U, V. With that,

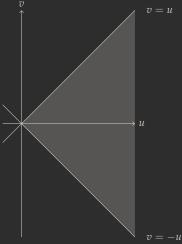
$$g(u,v) = f(\frac{u+v}{2}, \frac{u-v}{2}) \cdot \frac{1}{2}$$
  
=  $\frac{1}{2}e^{-(\frac{u+v}{2})-(\frac{u-v}{2})} = \frac{1}{2}e^{-u}$ .

With that, the marginal pdf of V is1

$$v < 0 \implies f_V(v) = \int_{-v}^{\infty} \frac{1}{2} e^{-u} du = \frac{1}{2} e^v$$
  
 $v \ge 0 \implies f_V(v) = \int_{v}^{\infty} \frac{1}{2} e^{-v} du = \frac{1}{2} e^{-v}.$ 

For the marignal pdf of U, we have

$$f_{U}(u) = \int_{-u}^{u} \frac{1}{2} e^{-u} dv = u e^{-u}$$
 for  $0 < u \le \infty$ .



 $^{\scriptscriptstyle \mathrm{I}}$  V is also called the Double Exponential Distribution.

# 13.1.2 Moment Generating Function Method

This method is particularly useful in finding distributions of sums of independent rvs.

#### Theorem 46 (Sums of MGF)

Suppose  $X_1, ..., X_n$  are independent rvs and  $X_i$  has  $mgf M_i(t)$  which exists for  $t \in (-h,h)$  for some H > 0. The mgf of  $Y = \sum_{i=1}^n X_i$  is given by

$$M_Y(t) = \prod_{i=1}^n M_i(t)$$

for  $t \in (-h, h)$ .

# Proof

Observe that

$$M_Y(t) = E\left[e^{tY}\right] = E\left[e^{t\sum_{i=1}^n X_i}\right] = E\left[\prod_{i=1}^n e^{tX_i}\right]$$

$$= \prod_{i=1}^n E\left[e^{tX_i}\right] = \prod_{i=1}^n M_i(t)$$

and  $t \in (-h, h)$  is preserved.

#### 66 Note

1. If  $X_i$ 's are IID rvs each with mgf M(t) then Y has mgf

$$M_Y(t) = [M(t)]^n$$
 for  $t \in (-h,h)$ .

2. Used in conjunction with the Uniqueness Theorem for mgfs, this theorem can be used to find the distribution of Y.

#### Exercise 13.1.1

*Show the following results:* 

- 1. If  $X \sim \text{Gam}(\alpha, \beta)$ , where  $\alpha \in \mathbb{N}$ , then  $\frac{2X}{\beta} \sim \chi^2(2\alpha)$ .
- 2. If  $X_i \sim \text{Gam}(\alpha_i, \beta)$ , i = 1, ..., n independently, then

$$\sum_{i=1}^{n} X_i \sim \operatorname{Gam}\left(\sum_{i=1}^{n} \alpha_i, \beta\right).$$

3. If  $X_i \sim \text{Gam}(1, \beta) = \text{Exp}(\beta)$ , i = 1, ..., n independently, then

$$\sum_{i=1}^{n} X_i \sim \operatorname{Gam}(n, \beta).$$

4. If  $X_i \sim \operatorname{Gam}\left(\frac{k_i}{2}, 2\right) = \chi^2(k_i)$ , i = 1, ..., n independently, then

$$\sum_{i=1}^{n} X_i \sim \chi^2 \left( \sum_{i=1}^{n} k_i \right).$$

5. If  $X_i \sim N(\mu, \sigma^2)$ , i = 1, ..., n independently, then

$$\sum_{i=1}^{n} \left( \frac{X_i - \mu}{\sigma} \right)^2 \sim \chi^2(n).$$

6. If  $X_i \sim \text{Poi}(\mu_i)$ , i = 1, ..., n independently, then

$$\sum_{i=1}^{n} X_i \sim \operatorname{Poi}\left(\sum_{i=1}^{n} \mu_i\right).$$

7. If  $X_i \sim \text{Bin}(n_i, p)$ , i = 1, ..., n independently, then

$$\sum_{i=1}^{n} X_i \sim \operatorname{Bin}\left(\sum_{i=1}^{n} n_i, p\right).$$

8. If  $X_i \sim NB(k_i, p)$ , i = 1, ..., n independently, then

$$\sum_{i=1}^{n} X_i \sim NB\left(\sum_{i=1}^{n} k_i, p\right).$$

### Solution

1. Using the mgf method,

$$\begin{split} M_{Y}(t) &= E\left[e^{tY}\right] = E\left[e^{t\cdot\frac{2x}{\beta}}\right] = \int_{0}^{\infty} e^{\frac{2tx}{\beta}} \frac{1}{\beta^{\alpha}\Gamma(a)} x^{\alpha-1} e^{-\frac{x}{\beta}} dx \\ &= \frac{1}{\beta^{\alpha}} \int_{0}^{\infty} \frac{1}{\Gamma(a)} x^{\alpha-1} e^{-\frac{x(1-2t)}{\beta}} dx \\ &\stackrel{(*)}{=} \frac{1}{\beta^{\alpha}} \left(\frac{\beta}{1-2t}\right)^{\alpha} \int_{0}^{\infty} \underbrace{\frac{1}{\left(\frac{\beta}{1-2t}\right)\Gamma(\alpha)} x^{\alpha-1} e^{-\frac{x}{\left(\frac{\beta}{1-2t}\right)}}}_{pdf \ of \ Gam\left(\alpha, \frac{\beta}{1-2t}\right)} dx \\ &= \left(\frac{1}{1-2t}\right)^{\alpha} = (1-2t)^{\alpha} \end{split}$$

where in (\*) we note that 1-2t>0 and so  $t>\frac{1}{2}$ . Observe that the mgf of Y is the mgf of  $\chi^2(2\alpha)$  and so by the  $\blacksquare$  Theorem 20,  $Y=\frac{2X}{\beta}\sim \chi^2(2\alpha)$ .

2. Using the mgf method, let  $Y = \sum_{i=1}^{n} X_i$ 

$$M_Y(t) = E\left[e^{tY}\right] = E\left[\prod_{i=1}^n e^{tX_i}\right] = \prod_{i=1}^n E[e^{tX_i}]$$
$$= \prod_{i=1}^n M_{X_i}(t) = \prod_{i=1}^n (1 - \beta t)^{-\alpha_1} = (1 - \beta t)^{-\sum_{i=1}^n \alpha_1}$$

and we observe that the last term is the mgf of Gam  $(\sum_{i=1}^{n} \alpha_i, \beta)$  and so the result follows.

- 3. From (2) is  $\alpha_i=1$  for i=1,...,n, then  $\sum_{i=1}^n \alpha_1=\sum_{i=1}^n 1=n$ . The result follows.
- 4. By (2), we have  $\sum_{i=1}^n X_i \sim \text{Gam}\left(\sum_{i=1}^n \frac{k_1}{2}, 2\right)$ . Then by (1), we have

$$\sum_{i=1}^{n} X_i \sim \chi^2 \left( \sum_{i=1}^{n} k_i \right)$$

as required.

5. We know that  $Y_i = \frac{X_i - \mu}{\sigma} \sim Z(0, 1)$ . Let  $W_i = Y_i^2$ . Then

$$\begin{split} M_{W_i}(t) &= E\left[e^{tY_i}\right] = \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-\frac{y^2}{2}} e^{ty^2} \, dy \\ &= \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{y^2}{2(1-2t)^{-1}}\right\} dy \\ &= (1-2t)^{-\frac{1}{2}} \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}(1-2t)^{-\frac{1}{2}}} \exp\left\{\frac{y^2}{2(1-2t)^{-1}}\right\} dy \\ &= (1-2t)^{-\frac{1}{2}}, \end{split}$$

which is the mgf of  $\chi^2(1)$ . Let  $W = \sum_{i=1}^n W_i$ . Then

$$M_W(t) = \prod_{i=1}^n (1-2t)^{-\frac{1}{2}} = (1-2t)^{-\frac{n}{2}}$$

which is the mgf of  $\chi^2(n)$ .

6. Let  $Y = \sum_{i=1}^{n} X_i$ . Then

$$M_Y(t) = \prod_{i=1}^n e^{\mu_i(e^t - 1)} = \exp\left[(e^t - 1)\sum_{i=1}^n \mu_i\right]$$

is the mgf of Poi  $(\sum_{i=1}^{n} \mu_i)$ .

7. Let  $Y = \sum_{i=1}^{n} X_i$ . Then

$$M_Y(t) = \prod_{i=1}^n (pe^t + q)^{n_i} = (pe^t + q)^{\sum_{i=1}^n n_i}$$

is the mgf of Bin  $(\sum_{i=1}^{n} n_i, p)$ .

8. Again, a similar approach: let  $Y = \sum_{i=1}^{n} X_i$ . Then

$$M_Y(t) = \prod_{i=1}^n \left(\frac{1-p}{1-pe^t}\right)^{k_i} = \left(\frac{1-p}{1-pe^t}\right)^{\sum_{i=1}^n k_i}$$

is the mgf of NB  $(\sum_{i=1}^{n} k_i, p)$ .

# 14 Lecture 14 Jun 19th 2018

# 14.1 Functions of Random Variables (Continued 2)

# **14.1.1** Moment Generating Function Method (Continued)

## **■** Theorem 47 (Gaussian Distribution)

*If*  $X_i \sim N(\mu_i, \sigma_i^2)$ , with i = 1, ..., n independently, then

$$\sum_{i=1}^n a_i X_i \sim N\left(\sum_{i=1}^n a_i \mu_i, \sum_{i=1}^n a_i^2 \sigma_i^2\right).$$

# Proof

Let  $Y_i = a_i X_i$ . Then by  $\blacksquare$  Theorem 19, we have

$$M_{Y_i}(t) = M_{X_i}(a_i t) = \exp\left[\mu_i a_i t + \frac{\sigma_i^2 a_i^2 t^2}{2}\right],$$

which implies that  $Y_i \sim N(a_i\mu_i, a_i^2\sigma_i^2)$  by  $\blacksquare$  Theorem 20. Then let  $Y = \sum_{i=1}^n Y_i$ . Then

$$M_Y(t) = \prod_{i=1}^n \exp\left[\mu_i a_i t + \frac{\sigma_i^2 a_i^2 t_2}{2}\right] = \exp\left[t \sum_{i=1}^n a_i \mu_i + \frac{t^2}{2} \sum_{i=1}^n a_i^2 \sigma_i^2\right],$$

which implies that  $Y \sim N\left(\sum_{i=1}^{n} a_i \mu_i, \sum_{i=1}^{n} a_i^2 \sigma_i^2\right)$ .

#### **■** Theorem 48 (Properties of the Gaussian Distribution)

Asumme  $X_1,...,X_n \sim N(\mu,\sigma_2)$ , independently, where  $\overline{X} = \frac{1}{n}\sum_{i=1}^n X_i$ , and  $S^2 = \frac{\sum_{i=1}^n (X_i - \overline{X})^2}{n-1}$ . Then

- 1.  $\overline{X} \sim N\left(\mu, \frac{\sigma^2}{n}\right)$ .
- 2. (Cochran's Theorem)  $\frac{(n-1)S^2}{\sigma^2} = \frac{\sum_{i=1}^n (X_i \overline{X})^2}{\sigma^2} \sim \chi^2(n-1)$ .
- 3. (*t-test*)  $\frac{\overline{X}-\mu}{S/\sqrt{n}} \sim t(n-1)$ .

#### Proof

- 1. Let  $Y_i = \frac{X_i}{n}$ .  $\blacksquare$  Theorem 47 implies that  $Y_i \sim N\left(\frac{\mu}{n}, \frac{\sigma^2}{n^2}\right)$ . Let  $Y = \sum_{i=1}^n Y_i$ . Then  $\blacksquare$  Theorem 47 implies that  $Y \sim N\left(\sum_{i=1}^n \frac{\mu}{n}, \sum_{i=1}^n \frac{\sigma^2}{n^2}\right) = N\left(\mu, \frac{\sigma^2}{n}\right)$ .
- 2. The proof of Cochran's Theorem is beyond the scope of this course, for it uses knowledge from linear algebra and involving Fourier Transforms. [Reference Wikipedia]
- 3. The proof of this statement is non-trivial: [Reference Wikipedia] [Reference Math SE]

#### Theorem 49 (F Distribution)

Suppose  $X_1, ..., X_n$  is a random sample from the  $N(\mu_1, \sigma_1^2)$  distribution and independently  $Y_1, ..., Y_m$  is a random sample from the  $N(\mu_2, \sigma_2^2)$  distribution. Let

$$S_1^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \overline{X})^2$$
 and  $S_2^2 = \frac{1}{m-1} \sum_{i=1}^m (Y_i - \overline{Y})^2$ .

Then

$$\frac{S_1^2/\sigma_1^2}{S_2^2/\sigma_2^2} \sim F(n-1, m-1).$$

## Proof

This is, once again, an incomplete proof. Rigor shall be appended where it is due or they shall be stated with reference, or if it is at a level that I cannot understand and that I cannot find a reference, I shall state that as well.

Note that the definition of the F distribution is as follows:

$$F_{\gamma_1,\gamma_2} = \frac{\chi^2(\gamma_1)/\gamma_1}{\chi^2(\gamma_2)/\gamma_2}$$

where  $\chi^2(\gamma_1)$  and  $\chi^2(\gamma_2)$  are independent.

Now we are given that  $X_1, ..., X_n$  and  $Y_1, ..., Y_m$  are independent of one another. It can likely be shown, using induction, that the following two are independent of one another:

$$\sum_{i=1}^{n} \left( \frac{X_i - \overline{X}}{\sigma_1} \right)^2 \sim \chi^2(n-1)$$

$$\sum_{j=1}^{m} \left( \frac{Y_j - \overline{Y}}{\sigma_2} \right)^2 \sim \chi^2(m-1).$$

Then

$$\frac{\sum\limits_{i=1}^{n}\left(\frac{X_1-\overline{X}}{\sigma_1}\right)^2}{\sum\limits_{j=1}^{m}\left(\frac{Y_j-\overline{Y}}{\sigma_2}\right)^2} = \frac{(n-1)S_1^2/\sigma_1^2}{(m-1)S_2^2/\sigma_2^2}.$$

With that, we have that

$$\frac{S_1^2/\sigma_1^2}{S_2^2/\sigma_2^2} = \frac{\frac{(n-1)S_1^2}{\sigma_1^2}/(n-1)}{\frac{(m-1)S_2^2}{\sigma_2^2}/(m-1)}.$$

The result follows.

# 15 Lecture 15 Jun 21st 2018

# 15.1 Limiting or Asymptotic Distributions

# 15.1.1 Convergence in Distribution

## **Definition** 52 (Convergence in Distribution)

Let  $X_1, X = 2, ..., X_n, ...$  be a sequence of rvs such that  $X_i$  has cdf  $F_i$ , for  $i \in \mathbb{N}$ . Let X be an rv with cdf F. We say that  $X_i$  converges in distribution to X and we write

$$X_n \stackrel{D}{\to} X$$

if

$$\lim_{i\to\infty} F_i(x) = F(x)$$

at all points x at which F is continuous. We call F the **limiting** or asymptotic distribution of  $X_i$ .

The following theorem, that we shall not prove, will be useful for the remainder of the notes.

# **■** Theorem 50 (Taylor Series with Lagrange's Remainder)

Suppose that  $f : [a,b] \to \mathbb{R}$  is infinitely differentiable, and  $x \in [a,b]$ . Then  $\forall x \in [a,b]$  and  $\forall k \in \mathbb{N}$ ,

$$f(x) = \sum_{i=0}^{k} \frac{f^{(i)}(c)(x-c)^{i}}{i!} + \frac{f^{(k+1)}(\zeta_X)(x-c)^{k+1}}{(k+1)!}$$

for some  $\zeta_X \in [c, x]$ .

The proof for the above involves the use of the Mean Value Theorem for Integrals.

#### Theorem 51 (Generalized Limit Definition of e)

*If*  $b,c \in \mathbb{R}$  *are constants and*  $\lim_{n \to \infty} \psi(n) = 0$ , then

$$\lim_{n\to\infty}\left[1+\frac{b}{n}+\frac{\psi(n)}{n}\right]^{cn}=e^{bc}.$$

# Proof

The above equation can be rewritten as

$$\lim_{n\to\infty}e^{cn\log\left(1+\frac{b}{n}+\frac{\psi(n)}{n}\right)}=e^{bc}$$

and so it suffices to prove that

$$\lim_{n \to \infty} cn \log \left( 1 + \frac{b}{n} + \frac{\psi(n)}{n} \right) = bc$$

Note that the Taylor Expansion with Lagrange's Remainders ( $\blacksquare$  Theorem 50) of log(1+x), where we pick c=0 and k=1 for convenience,

$$\log(1+x) = \frac{\log(1)x^0}{0!} + \frac{\left(\frac{1}{1-0}\right)x^1}{1!} + \frac{-\frac{1}{(1+\zeta_x)^2}x^2}{2!} = x - \frac{x^2}{2(1+\zeta_x)}$$

where  $\zeta \in [0, x]$ . Then

$$cn\log\left[1+\frac{b}{n}+\frac{\psi(n)}{n}\right]=cb+c\psi(n)-\frac{c(b+\psi(n))^2}{2n(1+\zeta)^2}$$

where  $\zeta \in \left[0, \frac{b+\psi(n)}{n}\right]$ . Note that by L'Hôpital's Rule, we have that

$$\lim_{n\to\infty}\frac{b+\psi(n)}{n}=\lim_{n\to\infty}\frac{\psi'(n)}{1}=0$$

and so the possible highest value for  $\zeta$  goes to 0 as  $n \to \infty$ . Then

$$\lim_{n\to\infty}\frac{c(b+\psi(n))^2}{2n(1+\zeta)^2}=0.$$

As a result, we have

$$\lim_{n \to \infty} cn \log \left( 1 + \frac{b}{n} + \frac{\psi(n)}{n} \right) = bc$$

as required.

#### **►** Corollary 52 (Limit definition of *e*)

*If*  $b, c \in \mathbb{R}$  *are constants, then* 

$$\lim_{n \to \infty} \left( 1 + \frac{b}{n} \right)^{cn} = e^{bc}$$

#### Example 15.1.1 (Example 5.1 (Course Notes 5.1.4))

Let  $X_i \sim \text{Exp}(1)$ , where  $i \in \mathbb{N}$ , independently so. Consider the sequence of rvs  $Y_1, Y_2, ..., Y_n, ...,$  where  $Y_n = \max(X_1, ..., X_n) - \log n$ . Find the limiting distribution of  $Y_n$ .

#### Solution

Firstly, observe that to find the support set of  $Y_n$ , note

$$0 < x_1, ..., x_n < \infty$$
$$0 < \max(x_1, ..., x_n) < \infty$$
$$-\log n < \max(x_1, ..., x_n) - \log n < \infty$$

$$supp(Y) = (-\log n, \infty)$$

Now the pf of  $Y_n$  is

$$F_n(y) = P(Y_n \le y) = P(\max(X_1, ..., X_n) - \log n \le y)$$

$$= P(\max(X_1, ..., X_n) \le y + \log n)$$

$$= \prod_{i=1}^n P(X_1 \le y + \log n) \quad \because X_1, ..., X_n \text{ are } IID$$

$$= \prod_{i=1}^n \left[ 1 - e^{-y - \log n} \right] = \prod_{i=1}^n \left[ 1 - \frac{e^- y}{n} \right]^n.$$

Thus the limiting distribution of  $Y_i$  is

$$\lim_{n\to\infty} F_n(y) = \lim_{n\to\infty} \left[1 - \frac{e^-y}{n}\right]^n = e^{-e^{-y}} \text{ for } y \in (-\log n, \infty).$$

#### Example 15.1.2 (Example 5.2 (Course Notes 5.1.5))

Let  $X_i \sim \text{Unif}(0,\theta)$ ,  $i \in \mathbb{N}$ , independently so. Consider the sequence of random variables  $Y_1, Y_2, ..., Y_n, ...,$  where  $Y_n = \max(X_1, ..., X_n)$ . Find the limiting distribution of  $Y_n$ .

#### Solution

*Clearly, support of*  $Y_n = (0, \theta)$ *. Note that* 

$$f_{X_i}(x) = \frac{1}{\theta} \mathbb{1}_{0 < x < \theta}.$$

Now since  $X_i$  are IID,

$$F_n(y) = P(Y_n \le y) = \prod_{i=1}^n P(X_1 \le y) = \prod_{i=1}^n \frac{y}{\theta} = \left(\frac{y}{\theta}\right)^n \quad y \in (0, \theta)$$

Then the limiting distribution is

$$\lim_{n\to\infty} F_n(y) = \lim_{n\to\infty} \left(\frac{y}{\theta}\right)^n = \begin{cases} 0 & y < \theta \\ 1 & y \ge \theta \end{cases}$$

*Define* Y *to have a distribution such that*  $P(Y = \theta) = 1$ *. Then* 

$$Y_n \stackrel{D}{\rightarrow} Y$$
.

# **Definition** 53 (Degenerate Distribution)

A function F is the cdf of a degenerate distribution at value y = c if

$$F(y) = \begin{cases} 0 & y < c \\ 1 & y \ge c \end{cases}.$$

In other words, F is the CDF of a discrete distribution where

$$P(Y = y) = egin{cases} 1 & y = c \ 0 & otherwise \end{cases}$$

We have that the earlier example gives us that the limiting distribution is a degenerate distribution.

#### 15.1.2 Convergence in Probability

## Definition 54 (Convergence in Probability)

A sequence of rvs  $X_1, X_2, ..., X_n, ...$  converges in probability to an rv X*if, for every*  $\varepsilon > 0$ *,* 

$$\lim_{n\to\infty} P(|X_n - X| \ge \varepsilon) = 0$$

or equivalently

$$\lim_{n\to\infty} P(|X_n - X| < \varepsilon) = 1$$

We write

$$X_n \stackrel{P}{\to} X$$
.

## Example 15.1.3 (Example 5.3)

Consider  $X \sim \text{Bernoulli}(0.3)$ . Define the sequence  $X_n = \left(1 + \frac{1}{n}\right) X$ ,  $n \in \mathbb{N}$ . Show that  $X_n \stackrel{P}{\rightarrow} X$ .

#### Solution

*Note that* 

$$P(X = x) = \begin{cases} 0.3 & x = 1 \\ 0.7 & x = 0 \end{cases}$$

Since  $X_n = \left(1 + \frac{1}{n}\right) X$ , we have

$$X_1 = 2X$$
,  $X_2 = \frac{3}{2}X$ ,  $X_3 = \frac{4}{3}X$ , ....

140 Lecture 15 Jun 21st 2018 - Limiting or Asymptotic Distributions

Note that 
$$|X_n - X| = \left| \left( 1 - \frac{1}{n} \right) X - X \right| = \left| \frac{1}{n} X \right| = \frac{1}{n} X$$
, so
$$P(|X_n - X| < \varepsilon) = P\left( \frac{1}{n} X < \varepsilon \right) = P(X < n\varepsilon)$$

$$= \begin{cases} P(X = 0) = 0.7 & n < \frac{1}{\varepsilon} \\ P(X = 1) + P(X = 0) = 1 & n \ge \frac{1}{\varepsilon} \end{cases}$$

Therefore

$$\lim_{n\to\infty} P(|X_n - X| < \varepsilon) = 1$$

and so

$$X_n \stackrel{P}{\to} X$$
.

We shall look into the following proposition in the next lecture.

# • Proposition

Suppose  $\forall n \in \mathbb{N}$ ,  $\{X_n\}$  is a sequence of rvs. Then

$$X_n \stackrel{P}{\to} X \implies X_n \stackrel{D}{\to} X$$

where X is an rv.

# 16 Lecture 16 Jun 26th 2018

# 16.1 Limiting or Asymptotic Distributions (Continued)

# 16.1.1 Convegence in Probability (Continued)

Before proving the proposition introduced in last class, we require the following lemma.

# ♣ Lemma 53

Let X, Y be rvs,  $a \in \mathbb{R}$ , and  $\varepsilon > 0$ . Then

$$P(Y \le a) \le P(X \le a + \varepsilon) + P(|Y - X| > \varepsilon)$$

#### Proof

*Note that* 

$$P(Y \le a) = P(Y \le a, X \le a + \varepsilon) + P(Y \le a, X > a + \varepsilon)$$
 (16.1)

$$\leq P(X \leq a + \varepsilon) + P(Y - X \leq a - X, a - X < -\varepsilon)$$
 (16.2)

$$\leq P(X \leq a + \varepsilon) + P(Y - X < -\varepsilon) \tag{16.3}$$

$$\leq P(X \leq a + \varepsilon) + P(Y - X < -\varepsilon) + P(Y - X > \varepsilon)$$
 (16.4)

$$= P(X \le a + \varepsilon) + P(|Y - X| > \varepsilon)$$

where Equation (16.1) is by Law of Total Probability, Equation (16.2) and Equation (16.3) are by the fact that for non-empty sets A and B,  $P(A \cap B) \leq P(A)$ , and Equation (16.4) is because  $P(Y - X > \varepsilon) \geq 0$ .  $\square$ 

# • Proposition 54 (Convergence in Probability Implies Convergence in Distribution)

Suppose  $\forall n \in \mathbb{N}$ ,  $\{X_n\}$  is a sequence of rvs. Then

$$X_n \stackrel{P}{\to} X \implies X_n \stackrel{D}{\to} X$$

where X is an rv.

This only states for the case of scalar random variables.

#### Proof

By Lemma 53, we have that 1

$$P(X_n \le a) \le P(X \le a + \varepsilon) + P(|X_n - X| > \varepsilon)$$
  
$$P(X \le a - \varepsilon) \le P(X_n \le a) + P(|X_n - X| > \varepsilon)$$

Then

$$P(X \le a - \varepsilon) - P(|X_n - X| > \varepsilon)$$

$$\le P(X_n \le a)$$

$$\le P(X \le a + \varepsilon) + P(|X_n - X| > \varepsilon)$$

As  $n \to \infty$ , by assumption that  $X_n \stackrel{P}{\to} X$ , we have that  $P(|X_n - X| > \epsilon) \to 0$ , and so

$$P(X \le a - \varepsilon) \le \lim_{n \to \infty} P(X_n \le a) \le P(X \le a + \varepsilon)$$

Now as  $\varepsilon \to 0^+$ , note that we must have the cdf of X be continuous on a by assumption, and so

$$P(X \le a) \le \lim_{n \to \infty} P(X_n \le a) \le P(X \le a)$$

and so by the Squeeze Theorem,

$$\lim_{n\to\infty} P(X_n \le a) = P(X \le a)$$

which then

$$X_n \stackrel{D}{\to} X$$

as required.

<sup>1</sup> We choose to use Lemma 53 to have  $P(X_n \le a)$  in two inequalities. Note that our goal is to show that

$$\lim_{n\to\infty} P(X_n \le a) = P(X \le a)$$

We already know that  $\varepsilon \to 0$ , so we should use that to our advantage in the case of X. Also, in hindsight, it is clear why we choose this method as the **Squeeze Theorem** is suitable to help us reach our conclusion.

# Definition 55 (Convergence in Probability to a Constant)

A sequence of rvs  $\{X_i\}_{i\in\mathbb{N}}$  converges in probability to a constant  $b\in\mathbb{R}$ if  $\forall \varepsilon > 0$ ,

$$\lim_{n\to\infty} P(|X_n - b| \ge \varepsilon) = 0$$

or equivalently

$$\lim_{n\to\infty} P(|X_n-b|<\varepsilon)=1.$$

We write

$$X_n \stackrel{P}{\rightarrow} b$$
.

# Example 16.1.1 (Example 5.4 - Weak Law of Large Numbers)

Let  $X_1, ..., X_n, ...$  be a sequence of IID rvs, each having mean  $\mu$  and variance

$$\overline{X}_n = \frac{1}{n} \sum_{i=1}^n X_i.$$

Then  $\overline{X}_n \stackrel{P}{\to} \mu$ .

# Solution

Observe that

$$E\left[\overline{X}_n\right] = E\left[\frac{1}{n}\sum_{i=1}^n X_i\right] = \frac{1}{n}\sum_{i=1}^n E[X_i] = \frac{n\mu}{n} = \mu,$$

Also, we have that

$$\operatorname{Var}\left(\overline{X}_{n}\right) = \operatorname{Var}\left(\frac{1}{n}\sum_{i=1}^{n}X_{i}\right) = \frac{1}{n^{2}}\sum_{i=1}^{n}\operatorname{Var}(X_{i}) = \frac{1}{n^{2}}\cdot n\sigma^{2} = \frac{\sigma^{2}}{2}.$$

Then by  $\blacksquare$  Theorem 16, we have that  $\forall \varepsilon > 0$ ,

$$P(|\overline{X}_n - \mu| \ge \varepsilon) \le \frac{\operatorname{Var}(\overline{X}_n)}{\varepsilon^2} = \frac{\sigma^2}{n\varepsilon^2}.$$

Therefore, we have

$$0 \le P(\left|\overline{X}_n - \mu\right| \ge \varepsilon) \le \min\left(1, \frac{\sigma^2}{n\varepsilon^2}\right)$$

*As*  $n \to \infty$ *, by Squeeze Theorem, we have that* 

$$\lim_{n\to\infty} P(\left|\overline{X}_n - \mu\right| \ge \varepsilon) = 0$$

i.e.  $\overline{X}_n \stackrel{P}{\rightarrow} \mu$ .

The Weak Law of Large Numbers is also known as Bernoulli's Theorem. Note that the converse of **6** Proposition 54 is not generally true. **(Example required)** However, with an additional condition, the converse becomes true.

# **♦** Proposition 55 (Partial Converse of **♦** Proposition 54)

Suppose  $\forall n \in \mathbb{N}$ ,  $\{X_n\}$  is a sequence of rvs, each with cdf  $F_n$ . If

$$\lim_{n \to \infty} F_n(x) = \lim_{n \to \infty} P(X_n \le x) = \begin{cases} 0 & x < c \\ & \end{cases},$$

$$1 & x > c$$

then

$$X_n \stackrel{P}{\to} c$$

#### Proof

*Note that for*  $\varepsilon > 0$ *, we have* 

$$P(|X_n - c| > \varepsilon) = P(X_n - c < -\varepsilon) + P(X_n - c > \varepsilon)$$

$$= P(X_n < c - \varepsilon) + P(X_n > c + \varepsilon)$$

$$= 1 - P(X_n < c + \varepsilon) + P(X_n < c - \varepsilon)$$

By assumption, we have that

$$\lim_{n \to \infty} P(X_n \le c + \varepsilon) = 1$$
$$\lim_{n \to \infty} P(X_n \le c - \varepsilon) = 0$$

So

$$\lim_{n\to\infty} P(|X_n-c|>\varepsilon)=1-1-0=0$$

and so by definition, we have  $X_n \stackrel{P}{\rightarrow} c$ .

# Definition 56 (Double Parameter Exponential Distribution)

We say that an rv  $X \sim \text{Exp}(\lambda, \theta)$  when X has pdf

$$f(x) = e^{-\frac{x-\theta}{\lambda}}$$
 for  $x \in (\theta, \infty)$ 

where  $\lambda$  is the scale parameter, and  $\theta$  the location parameter.

#### Example 16.1.2 (Example 5.5 (Course Notes 5.2.5))

Let  $X_i \sim \text{Exp}(1,\theta)$ , i = 1, 2, ..., independently. Consider the sequence of rvs  $Y_1, Y_2, ...$  where  $Y_n = \min(X_1, X_2, ..., X_n), n = 1, 2, ...$  Show that  $Y_n \stackrel{P}{\rightarrow} \theta$ .

#### Solution

Since we want to show that  $Y_n \stackrel{P}{\to} \theta$ , and  $\theta$  is a constant, we can use

• Proposition 55. Thus we need to show that

$$\lim_{n \to \infty} F_{Y_n}(y) = \lim_{n \to \infty} P(Y_n \le y) = egin{cases} 0 & y < \theta \\ 1 & y > \theta \end{cases}$$

Since  $Y_n$  is defined as the minimum of n of the first  $X_i$  rvs, we need to use the Law of Total Probability in order to be able to make sense of the order statistics, i.e.

$$P(Y_n \le y) = 1 - P(Y_n > y)$$

Now

$$P(Y_n > y) = \prod_{i=1}^{n} P(X_i > y) = \prod_{i=1}^{n} \left[ \int_{y}^{\infty} e^{-(x-\theta)} dx \right]$$
$$= \prod_{i=1}^{n} \left[ e^{-(y-\theta)} \right] = e^{-n(y-\theta)}$$

Thus

$$\lim_{n \to \infty} P(Y_n \le y) = \begin{cases} 1 - \lim_{n \to \infty} P(Y_n > y) & y > \theta \\ 0 & y \le \theta \end{cases}$$
$$= \begin{cases} 1 - \lim_{n \to \infty} e^{-n(y-\theta)} = 1 & y > \theta \\ 0 & y \le \theta \end{cases}$$

*since if*  $y < \theta$ 

$$P(Y_n \le y) = P(\min(X_1, X_2, ..., X_n) \le y) = 0.$$

*Note that if*  $y = \theta$ *, then* 

$$e^{-n(y-\theta)} = 1.$$

*The proof is complete with* • *Proposition 55.* 

# 16.1.2 Limit Theorems

# • Proposition 56 (Convergence in Distribution and MGF)

Let  $X_1, X_2, ..., X_n, ...$  be a sequence of rvs such that  $X_n$  has  $mgf M_n(t)$ . Let X be an rv with mgf M(t).

$$X_n \stackrel{D}{\to} X \iff \left[ \exists h > 0 \ \forall t \in (-h, h) \quad \lim_{n \to \infty} M_n(t) = M(t) \right]$$

#### Proof

*Note that* 

$$\lim_{n\to\infty} M_n(t) = \lim_{n\to\infty} \int_{\text{supp}(X_n)} e^{tx} \frac{d}{dx} F_{X_n}(x) dx$$

and

$$M(t) = \int_{\text{supp}(X)} e^{tx} \frac{d}{dx} F_X(x) dx$$

The result follows assuming that the integral converges<sup>2</sup>.

#### Example 16.1.3

Consider the sequence  $Y_1, Y_2, ..., where Y_i \sim Bin(n, \frac{\mu}{n})$ , for i = 1, 2, ... Find the limiting distribution of  $Y_n$ .

# Solution (Example 5.6)

*Note that*  $Y_n \sim \text{Bin}\left(n, \frac{\mu}{n}\right) \implies$ 

$$M_{Y_n}(t) = \left(\frac{\mu}{n}e^t + 1 - \frac{\mu}{n}\right)^n = \left(1 + \frac{\mu\left(e^t - 1\right)}{n}\right)^n.$$

So

$$\lim_{n\to\infty} M_{Y_n}(t) = \lim_{n\to\infty} \left(1 + \frac{\mu\left(e^t - 1\right)}{n}\right)^n = e^{\mu\left(e^t - 1\right)}$$

<sup>&</sup>lt;sup>2</sup> This allows us to "swap" the limit, the integral sign, and the differential operator.

which is the mgf of a Poisson Distribution. Thus by • Proposition 56, we have that

$$Y_n \stackrel{D}{\rightarrow} Y \sim Poi(\mu)$$
.

## Example 16.1.4 (Example 5.7)

Let  $Y_1, Y_2, ... \sim G(\mu, \sigma)$ . Then

$$\frac{\overline{Y}_n - \mu}{\sigma / \sqrt{n}} \sim G(0, 1)$$

where  $\overline{Y}_n = \frac{1}{n} \sum_{i=1}^n Y_i$ .

# Solution

Note that by P Theorem 19, we have that

$$M_{\overline{Y}_n}(t) = \prod_{i=1}^n M_{Y_i}\left(\frac{t}{n}\right) = \prod_{i=1}^n e^{\frac{\mu t}{n} + \frac{(\sigma t)^2}{2n^2}} = e^{\mu t + \frac{\sigma^2 t^2}{n}}.$$

which is the mgf of  $G\left(\mu, \frac{\sigma}{\sqrt{n}}\right)$ . Let  $X_n = \frac{\overline{Y}_n - \mu}{\sigma/\sqrt{n}}$ .

Then again, by 里 Theorem 19, we have

$$M_{X_n}(t) = e^{-\frac{\mu t}{\sigma/\sqrt{n}}} M_{\overline{Y}_n}(t) \left(\frac{t}{\sigma/\sqrt{n}}\right)$$
$$= e^{-\frac{\mu t}{\sigma/\sqrt{n}}} e^{\frac{\mu t}{\sigma/\sqrt{n}} + \frac{\sigma^2/n}{2} \left(\frac{t}{\sigma/\sqrt{n}}\right)^2} = e^{\frac{t^2}{2}}$$

which is the mgf of G(0,1).

## **Example 16.1.5 (Example 5.8)**

Let  $Y_n \sim \text{Gam}(n, 2)$ . Find the limiting distribution of

$$W_N = \frac{1}{2\sqrt{n}}(Y_n - 2n).$$

#### Solution

By 🖳 Theorem 19, we have

$$\begin{split} M_{W_n}(t) &= e^{-t\sqrt{n}} M_{Y_n}(t) \left(\frac{t}{2\sqrt{n}}\right) = e^{-t\sqrt{n}} \left(1 - 2\left(\frac{t}{2\sqrt{n}}\right)\right)^{-n} \\ &= \left[e^{\frac{t}{\sqrt{n}}} \left(1 - \frac{t}{\sqrt{n}}\right)\right]^{-n} = \left[\left(1 - \frac{t}{\sqrt{n}}\right) \sum_{k=0}^{\infty} \frac{(t/\sqrt{n})^k}{k!}\right]^{-n} \\ &= \left[\left(1 - \frac{t}{\sqrt{n}}\right) \left(1 + \frac{t}{\sqrt{n}} + \frac{t^2}{2n} + \frac{t^3}{3\sqrt{n^3}} + \dots\right)\right]^{-n} \\ &= \left(1 - \frac{t}{\sqrt{n}} + \frac{t}{\sqrt{n}} - \frac{t^2}{n} + \frac{t^2}{2n} - \frac{t^3}{2\sqrt{n^3}} + \frac{t^3}{3\sqrt{n^3}} - \frac{t^4}{3n^2} + \dots\right)^{-n} \\ &= \left(1 - \frac{t^2}{2n} - \underbrace{\frac{t^3}{6\sqrt{n^3}} - \frac{t^4}{3n^2} + \dots}_{\to 0 \text{ gs } n \to \infty}\right)^{-n} \end{split}$$

So we have, by P Theorem 51,

$$\lim_{n\to\infty} M_{W_n}(t) = \lim_{n\to\infty} \left(1 - \frac{t^2/2}{n} - \frac{t^3}{6\sqrt{n^3}} - \frac{t^4}{3n^2} + \ldots\right)^{-n} = e^{-\left(-\frac{t^2}{2}\right)} = e^{\frac{t^2}{2}},$$

which is the mgf of G(0,1). So we know that

$$W_n \stackrel{D}{\to} W \sim G(0,1)$$

for some limiting distribution W that has a G(0,1) distribution.

## Theorem 57 (Central Limit Theorem)

Suppose  $X_1, X_2, ...$  is a sequence of IID rvs with  $E(X_i) = \mu$  and  $Var(X_i) = \sigma^2 < \infty$ , for i = 1, 2, ... Then

$$\frac{\overline{X}_n - \mu}{\sigma / \sqrt{n}} = \frac{\sqrt{n}(\overline{X}_n - \mu)}{\sigma} \xrightarrow{D} Z \sim N(0, 1)$$

where  $\overline{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$ .

**?** We actually also need the condition that all moments of  $\frac{X_i - \mu}{\sigma}$  exists and is bounded.

#### Proof

Since the first and second moment exist, the mgf must be well-defined, and is at least of class  $C^{2}$  3. We will use  $\bullet$  Proposition 56 and

- <sup>3</sup> You should remember this from your calculus courses. Otherwise, [Reference
- Wiki

💻 Theorem 50 in this proof. Firstly note that

$$E\left[\sum_{i=1}^{n} X_i\right] = n\mu \text{ and } \operatorname{Var}\left(\sum_{i=1}^{n} X_i\right) = n\sigma^2$$

Let  $\Sigma X_n = \sum_{i=1}^n X_i$  and  $Z = \frac{\sqrt{n}(\overline{X}_n - \mu)}{\sigma}$ . Note that

$$Z = \frac{\sqrt{n}(\overline{X}_n - \mu)}{\sigma} = \frac{n\overline{X}_n - \mu n}{\sigma\sqrt{n}} = \frac{\Sigma X_n - n\mu}{\sigma\sqrt{n}} = \sum_{i=1}^n \frac{X_i - \mu}{\sigma\sqrt{n}}.$$

Let  $Y_i = \frac{X_i - \mu}{\sigma}$  so that

$$Z = \sum_{i=1}^{n} \frac{1}{\sqrt{n}} Y_i.$$

Now the mgf of Z, by ■ Theorem 19 and 6 Proposition 39, is

$$M_Z(t) = \prod_{i=1}^n M_{Y_i} \left( \frac{t}{\sqrt{n}} \right) = M_{Y_i} \left( \frac{t}{\sqrt{n}} \right)^n.$$
 (16.5)

*Note that* 

$$E(Y_i) = E\left(\frac{X_i - \mu}{\sigma}\right) = \frac{1}{\sigma}E(X) - \frac{\mu}{\sigma} = 0$$
$$Var(Y_i) = Var\left(\frac{X_i - \mu}{\sigma}\right) = \frac{1}{\sigma^2}Var(X_i) = 1.$$

So by P Theorem 50 up to the second derivative, we have

$$M_{Y_i}\left(\frac{t}{\sqrt{n}}\right) = 1 + \frac{t^2}{2n} + \frac{M_{Y_i}^{(3)}(\zeta)\left(\frac{t}{\sqrt{n}}\right)^3}{3!}$$

where  $\zeta \in \left[0, \frac{t}{\sqrt{n}}\right]$ . Note that the last term in the above equation goes to 0 faster than the second term as  $n \to \infty$ , since we would have  $n^{\frac{3}{2}}$  in the denominator, which is larger than n in the second term. With that, continuing with Equation (16.5) and now taking the limit as  $n \to \infty$ , we have

$$\lim_{n \to \infty} M_Z(t) = \lim_{n \to \infty} \left[ 1 + \frac{t^2}{2n} + \phi(n) \right]^n$$

where

$$\phi(n) = \frac{M_{Y_i}^{(3)}(\zeta) \left(\frac{t}{\sqrt{n}}\right)^3}{3!}.$$

Then by P Theorem 51, we have

$$\lim_{n\to\infty} M_Z(t) = e^{\frac{t^2}{2}}$$

which is the mgf of N(0,1). This completes the proof.

## Example 16.1.6 (Example 5.9)

Revisit Example 16.1.5 using the 🖳 Central Limit Theorem. Given

$$W_n = \frac{Y - 2n}{2\sqrt{n}}$$

where  $Y \sim \text{Gam}(n,2)$ , find sequences  $a_n$  and  $b_n$  such that

$$a_n(Y-b_n) \stackrel{D}{\rightarrow} Z \sim G(0,1).$$

## Solution

Let  $X_1, X_2, ...$  be a sequence of IID rvs with distribution Exp(2). Note that we then have  $E(X_i) = 2$  and  $\text{Var}(X_i) = 4$ . Let  $Y_n = \sum_{i=1}^n X_i$ . Note that

$$M_{Y_n}(t) = [M_{X_i}(t)]^n = (1-2t)^{-n}$$

which is the mgf of Gam(n, 2). By CLT, we have that

$$\frac{\frac{1}{n}Y_n - 2}{2/\sqrt{n}} \stackrel{D}{\to} G(0,1).$$

Note that

$$\frac{\frac{1}{n}Y_n-2}{2/\sqrt{n}}=\frac{\sqrt{n}}{2n}(Y_n-2n)$$

and so

$$a_n = \frac{1}{2}\sqrt{n}$$
 and  $b_n = 2n$ 

#### Example 16.1.7 (Example 5.10 (Course Notes 5.3.6))

Suppose  $Y_n \sim \chi^2(n)$ , n = 1, 2, ... Consider the sequence of rvs  $Z_1, Z_2, ...$  where  $Z_n = (Y_n - n)/\sqrt{2n}$ . Show that

$$Z_n = \frac{Y_n - n}{\sqrt{2n}} \stackrel{D}{\rightarrow} Z \sim G(0,1).$$

#### Solution

Let  $X_1, X_2, ... \sim \chi^2(1)$  be IID rvs, and so  $E[X_i] = 1$  and  $Var(X_i) = 2$  for  $1 \le i \le n$ , and let  $Y_n = \sum_{i=1}^n X_i$ . Note that

$$M_{Y_n}(t) = [M_{X_i}(t)]^n = (1 - 2t)^{-\frac{n}{2}}$$

This involves the knowledge that the sum of n exponential distributions with mean  $\mu$  results in a gamma distribution with parameter n and  $\mu$ . Since this has never been proven in this set of notes, it shall be proven here in this solution.

Once again our class pulls the card of "not explaining or showing stuff that they are supposed to". We will be using the fact that a sum of n Chi-Squared Distirbutions, each with degree of freedom k, will result in a gamma distribution with parameter  $\frac{nk}{2}$  and 2. This will be proven in this exercise.

is the mgf of  $\chi^2(n)$ . By CLT, we have that

$$\frac{\frac{1}{n}\sum_{i=1}^{n}X_{i}-1}{\sqrt{2}/\sqrt{2}} \stackrel{D}{\rightarrow} Z \sim G(0,1).$$

Note that

$$\frac{\frac{1}{n}\sum_{i=1}^{n}X_{i}-1}{\sqrt{2}/\sqrt{n}} = \frac{Y_{n}-n}{\sqrt{2}n/\sqrt{n}} = \frac{Y_{n}-n}{\sqrt{2n}}.$$

This completes our example.

# Example 16.1.8 (Example 5.11 (Course Notes 5.3.7))

Suppose  $Y_n \sim Bin(n, p)$ , n = 1, 2, ... Consider the sequence of rvs  $Z_1, Z_2, ...,$  where  $Z_n = \frac{Y_n - np}{\sqrt{np(1-p)}}$ . Show that

$$Z_n = rac{Y_n - np}{\sqrt{np(1-p)}} \stackrel{D}{
ightarrow} Z \sim \mathrm{N}(0,1).$$

#### Solution

Let  $X_1, X_2, ...$  be a sequence of IID Bernoulli trials with success probability  $0 \le p \le 1$ . Note that for  $i \in \mathbb{N}$ , we have that  $E[X_i] = p$  and  $Var(X_i) = p(1-p)$ . Now let  $Y_n = \sum_{i=1}^n X_i$ . We note that this satisfies our assumption, i.e.  $Y_n \sim Bin(n, p)$  since

$$M_{Y_n}(t) = [M_{X_1}]^n = (pe^t + 1 - p)^n$$

is the mgf of Bin(n, p). By CLT, we have that

$$\frac{\frac{1}{n}\sum_{i=1}^{n}X_{i}-p}{\sqrt{p(1-p)}/\sqrt{n}} \xrightarrow{D} Z \sim N(0,1).$$

*Note that* 

$$\frac{\frac{1}{n}\sum_{i=1}^{n}X_{i}-p}{\sqrt{p(1-p)/\sqrt{n}}} = \frac{Y_{n}-np}{\sqrt{np(1-p)}}$$

and so this completes our proof.

# • Proposition 58 (Other Limit Theorems)

- 1.  $X_n \stackrel{P}{\to} a \land g$  continuous at  $x = a \implies g(X_n) \stackrel{P}{\to} g(a)$ .
- 2.  $X_n \stackrel{P}{\rightarrow} a \wedge Y_n \stackrel{P}{\rightarrow} b \wedge g$  continous at  $(x,y) = (a,b) \implies$  $g(X_n, Y_n) \stackrel{P}{\to} g(a, b).$
- 3. (Slutsky's Theorem)  $X_n \stackrel{D}{\rightarrow} X \wedge Y_n \stackrel{P}{\rightarrow} b \wedge g$  continous at (x,y) =(x,b) for all  $x \in \text{supp}(X) \implies g(X_n,Y_n) \stackrel{D}{\rightarrow} g(X,b)$ .

4. (Continuous Mapping Theorem)  $X_n \stackrel{D}{\to} X \wedge g$  continous at all  $x \in \text{supp}(X) \implies g(X_n) \stackrel{D}{\to} g(X)$ .

#### Proof

1. Since g is continuous at a, we have that

$$\forall \varepsilon > 0 \ \exists \delta > 0 \ \forall x \in \text{supp}(X) \ (|x - a| < \delta \implies |g(x) - g(a)| < \varepsilon)$$

It follows that

$$P[|g(X_n) - g(a)| < \varepsilon] \ge P[|X_n - a| < \delta]$$

since  $P(B) \ge P(A)$  whenever  $A \subset B^4$ . Since  $X_n \xrightarrow{P} a$ , we have that for any  $\delta > 0$ , we have

$$1 \ge \lim_{n \to \infty} P[|g(X_n) - g(a)| < \varepsilon] \ge \lim_{n \to \infty} P[|X_n - a| < \delta] = 1.$$

Thus

$$\lim_{n\to\infty} P[|g(X_n) - g(a)| < \varepsilon] = 1$$

and so

$$g(X_n) \stackrel{P}{\to} g(a)$$

as required.

- 2. This is simply a more general case than (1).
- 3. See this Wikipedia article for a proof. Requires knowledge of measure theory (in particular, convergence of measures).
- 4. See this Wikipedia article for a proof. Requires knowledge of measure theory (in particular, convergence of measures).

#### Example 16.1.9

If  $X_n \stackrel{P}{\to} 10$  and  $Y_n \stackrel{P}{\to} 2$ , then  $X_n Y_n \stackrel{P}{\to} 20$  since g(x,y) = xy is continuous at (10,2).

If  $Z_n \stackrel{D}{\to} Z \sim G(0,1)$  and  $a_n \stackrel{P}{\to} a$  where a is a constant, then  $a_n Z_n \stackrel{D}{\to} aZ \sim (0,a)$  since g(a,z) = az is continous at (a,z) for all  $z \in \text{supp}(Z)$ ..

Example 16.1.10 (Example 5.12 (Course Notes 5.3.10))

<sup>4</sup> While I do not have a rigorous proof of this for our case here, it is a sensible result seeing that  $\delta$  depends on  $\varepsilon$ . My hunch was right: thinking about the two probability measures using the definition of a random variable, we see that the  $w \in S$  that works for  $|X_n - a| < \delta$  will definitely work for  $|g(X_n) - g(a)| < \varepsilon$ , but the converse is not necessarily true. Therefore, the events covered in the left term is a larger set that contains the event on the right.

If  $X_n \stackrel{P}{\to} a > 0$ ,  $Y_n \stackrel{P}{\to} b \neq 0$  and  $Z_n \stackrel{D}{\to} Z \sim G(0,1)$ , the find the limiting distributions of each of the following:

- 1.  $\sqrt{X_n}$
- 2.  $X_n + Y_n$
- 3.  $X_nZ_n$

#### Solution

- 1. Since a > 0, we have the function  $g(x) = \sqrt{x}$  is continuous on a, and so by  $\bullet$  Proposition 58, we have that  $\sqrt{X_n} \stackrel{P}{\to} \sqrt{a}$ .
- 2. Since the function g(x,y) = x + y is continous on the real number 2-tuple (a,b), we have that  $X_n + Y_n \stackrel{P}{\rightarrow} a + b$ .
- 3. Since the function g(x,z) = xz is continuous on the real number 2-tuple (x,z) for all  $x \in \text{supp}(X_n)$  and  $z \in \text{supp}(Z_n)$ , by Slutsky's Theorem, we have  $X_n Z_n \stackrel{D}{\rightarrow} aZ \sim G(0, a)$ .
- 4. Note that the function  $g(z) = \frac{1}{z}$  is continuous for all  $z \in \operatorname{supp}(X)$ except at z = 0. But since  $Z \sim G(0,1)$  is a continuous distribution, at a single point z = 0, P(Z = 0) = 0, and in a distribution it is of negligible value<sup>5</sup>. Therefore, we have that

$$\frac{1}{Z_n} \stackrel{D}{\to} \frac{1}{Z}$$

## Example 16.1.11 (Example 5.13 (Course Notes 5.3.11))

Suppose  $X_1, X_2, ... \sim \text{Poi}(\mu)$  are IID rvs. Define  $Z_n = \frac{\sqrt{n}(\overline{X}_n - \mu)}{\sqrt{\overline{X}_n}}$ . Find the *limiting distribution of*  $Z_n$ .

## Solution

Firstly, note that the parameter for a Poisson distribution is positive. Note that by CLT, we have

$$\frac{\overline{X}_n - \mu}{\sqrt{\mu}/\sqrt{n}} \sim Z \sim G(0, 1)$$

*Note that since*  $\sqrt{m}$  *is a constant, we have that*  $\sqrt{m} \stackrel{P}{\to} \sqrt{m}$  *is trivially true.* Thus by our limit theorems, we have

$$\sqrt{n}(\overline{X}_n - \mu) \sim \sqrt{m}Z \sim G(0, \sqrt{m}).$$

<sup>5</sup> This is a painful thing to write down

Now by the Weak Law of Large Numbers, we have

$$\overline{X}_n \stackrel{P}{\to} \mu$$
.

Since a function  $g(x) = \sqrt{x}$  is continuous for  $x = \mu > 0$ , we have that

$$\sqrt{\overline{X}}_n \stackrel{P}{\to} \sqrt{m}$$
.

Then by the Continuous Mapping Theorem, we have that

$$Z_n = \frac{\sqrt{n}(\overline{X}_n - \mu)}{\sqrt{\overline{X}_n}} \sim \frac{\sqrt{\mu}Z}{\sqrt{\mu}} \sim G(0, 1).$$

# **P** Theorem 59 (Generalized $\delta$ -Method)

Let  $X_1, X_2, ...$  be a sequence of rvs such that

$$n^b(X_n-a)\stackrel{D}{\to} X$$

for some b > 0. Suppose the function g(x) is differentiable at a and  $g'(a) \neq 0$ . Then

$$n^b[g(X_n)-g(a)] \stackrel{D}{\to} g'(a)X.$$

## Proof

Using the Mean Value Theorem, we have that

$$g(X_n) = g(a) + g'(c)(X_n - a)$$
 (16.6)

for some c in between  $X_n$  and a. Since  $X_n \stackrel{P}{\to} a^6$ , and, WLOG,  $X_n < c < a$ ,  $c \stackrel{P}{\to} a$ . Then by the **Continuous Mapping Theorem**, we have  $g(c) \stackrel{P}{\to} g(a)$ .

Now by rearranging Equation (16.6) and multiplying both sides by  $n^b$ , we get

$$n^{b}[g(X_{n})-g(a)]=g'(c)n^{b}[X_{n}-a]\stackrel{D}{\to}g'(a)X$$

by Slutsky's Theorem.

<sup>6</sup> I have no idea why...

# $\blacktriangleright$ Corollary 60 ( $\delta$ -Method)

Let  $X_1, X_2, ...$  be a sequence of IID rvs with mean  $\mu$  and variance  $\sigma^2$ .

Suppose the function g(x) is differentiable at  $\mu$  and  $g'(\mu) \neq 0$ . Then

$$\sqrt{n}[g(\overline{X}_n) - g(\mu)] \stackrel{D}{\to} Z \sim N\left(0, [g'(\mu)]^2 \sigma^2\right)$$

#### Proof

Note that by the same working as in Example 16.1.11 for  $\sqrt{n}(\overline{X}_n - \mu)$ , we have that<sup>7</sup>

$$\sqrt{n}(\overline{X}_n - \mu) \stackrel{D}{\to} Y \sim N(0, \sigma^2).$$

Then by P Theorem 59, we have that

$$\sqrt{n}(g(\overline{X}_n) - g(\mu)) \stackrel{D}{\rightarrow} g'(\mu) Y \sim N\left(0, [g'(\mu)]^2 \sigma^2\right)$$

as required.

<sup>7</sup> This is actually stated as an assumption on the  $\delta$ -Method Wikipedia article -[Reference].

# Example 16.1.12 (Example 5.14)

Let  $X_1, X_2, ... \sim \text{Geo}(p)$  is a sequence of IID rvs, each with support  $supp(X_i) = \{0, 1, 2, ...\}$ . Derive the limiting distribution of

$$W_n = \sqrt{n} \left( \frac{1}{\overline{X}_n} - \frac{p}{1 - p} \right)$$

#### Solution

Note that by CLT, we have that

$$\frac{\overline{X}_n - \frac{1-p}{p}}{\sqrt{\frac{1-p}{p^2}}/\sqrt{n}} \stackrel{D}{\to} Z \sim N(0,1)$$

Then using **b** Proposition 58, we have

$$\sqrt{n}\left(\overline{\mathrm{X}}_n - rac{1-p}{p}
ight) \stackrel{D}{
ightarrow} \sqrt{rac{1-p}{p^2}} Z \sim N\left(0,rac{1-p}{p^2}
ight)$$

**We took**  $g(x) = \frac{1}{x}$  and the given solution in class somehow circumvents the fact that x = 0 is a case.

## Example 16.1.13 (Example 5.15)

Suppose that  $X_1, X_2, ... \sim \text{Exp}(\theta)$  is a sequence of IID rvs. Find constants  $a_n$  and  $b_n$  such that

$$W_n = b_n(\overline{X}_n^2 - a_n)$$

has a non-denegerate limiting distribution.

# Solution

By CLT, we have

$$\frac{\overline{X}_n - \theta}{\theta / \sqrt{n}} \stackrel{D}{\to} Z \sim N(0, 1)$$

which then by **b** Proposition 58, we have

$$\sqrt{n}(\overline{X}_n - \theta) \stackrel{D}{\to} \theta Z \sim N(0, \theta^2).$$

Let  $g(x) = x^2$ . Then g is continous on any  $x \in \text{supp}(X_n)$  and on  $\theta$ . Then using  $\blacktriangleright$  Corollary 60, we have

$$\sqrt{n}(\overline{X}_n^2 - \theta^2) \sim \theta^2 Z \sim N(0, 4\theta^2).$$

So we just need to pick  $b_n = \sqrt{n}$  and  $a_n = \theta^2$ .

# 17.1 Estimation

Suppose  $X_1,...,X_n \sim f(x;\theta)$  is an IID sequence of rvs, where  $f(x;\theta)$  is the pf of the  $X_i$ 's. The joint distribution of  $X_1,...,X_n$  is

$$\prod_{i=1}^{n} f(x_i; \theta)$$

where the unknown parameter  $\theta$  can either be a scalar in  $\Omega$ , where  $\Omega$  is the parameter space or the set of possible values of  $\theta$ , or a vector,

i.e. 
$$\theta = \begin{pmatrix} & & & & \\ \theta_1 & \theta_2 & \dots & \theta_n \end{pmatrix}^T$$
.

We are interested in making inferences about the unknown parameter  $\theta$ , i.e. we want to find **estimators** (point and interval) of  $\theta$  and we want to test our hypothesis about  $\theta$ .

Before proceeding the the rest of this chapter (actual chapter, ahem...), we require the following definitions.

# Definition 57 (Statistic)

A statistic,  $T = T(X) = T(X_1, ..., X_n)$ , is a function of the data which does not depend on any unknown parameter(s).

#### Example 17.1.1

Suppose  $X_1, ..., X_n$  is a random sample from a distribution with  $E[X_i] = \mu$  and  $Var(X_i) = \sigma^2$ , where  $\mu$  and  $\sigma^2$  are unknown. The sample mean  $\overline{X} = \frac{1}{n} \sum_{i=1}^n X_i$  and the sample variance  $S^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \overline{X})^2$  are statistics,

while  $\frac{\overline{X}-\mu}{\sigma/\sqrt{n}}$  is not a statistic.

## Definition 58 (Estimators and Estimates)

A statistic  $T = T(X) = T(X_1, ..., X_n)$  that is used to estimate  $\tau(\theta)$ , a function of  $\theta$ , is called an **estimator** of  $\tau(\theta)$ , and an observed value of the statistic  $t = t(x) = t(x_1, ..., x_n)$  is called an **estimate** of  $\tau(\theta)$ .

## **Example 17.1.2**

Suppose  $X_1,...,X_n$  are IID rvs with  $E[X_i] = \mu$ . The rv  $\overline{X}$  is an estimator of  $\mu$ . For a given set of observations,  $x_1,...,x_n$ , the number  $\overline{x} = \frac{1}{n} \sum_{i=1}^n x_i$  is an estimate of  $\mu$ .

There are certain "properties" of an estimator that we look for, in particular, for an estimator,  $\tilde{\theta}$ , of an unknown parameter  $\theta$ , we want that the estimator

- is not biased, aka an **unbiased estimator**, i.e.  $E[\tilde{\theta}]$ ;
- has small variance;
- is consistent, i.e.  $\tilde{\theta} \stackrel{P}{\to} \theta$ .

#### 17.1.1 Maximum Likelihood Estimation

Suppose X is a discrete rv with pf  $P(X=x;\theta)=f(x;\theta), \theta\in\Omega$  where the scalar parameter  $\theta$  is unknown. Suppose x is an observed value of the rv X. Then the probability of observing this value is

$$P(X = x; \theta) = f(x; \theta).$$

With the observed value of x substituted into  $f(x;\theta)$ , we have a function of the parameter  $\theta$ , referred to as the **likelihood function**, and denoted  $L(\theta)$ . In the absense of any other information, it seems logical (temptingly so) that we should estimate the parameter  $\theta$  using a value that is "the most compatible" with the data. E.g., we might choose the value of  $\theta$  of which it maximizes the probability of the observed data, or equivalently, the value of  $\theta$  which maximizes the likelihood function  $L(\theta)$ .

But first, let us formally state the definition of a likelihood function.

# Definition 59 (Likelihood function)

Suppose X is an rv with pf  $f(x;\theta)$ , where  $\theta \in \Omega$  is a scalar. If x is the observed data, then the **likelihood function** for  $\theta$  based on x is

$$L(\theta) = P(X = x; \theta) = f(x; \theta)$$
 for  $\theta \in \Omega$ .

Similarly so, suppose  $X_1, ..., X_n$  is a random sample from a distribution, each with pf  $f(x;\theta)$ , and let  $x_1,...,x_n$  be the observed data. Then the **likelihood function** for  $\theta$  based on  $x_1, ..., x_n$  is

$$L(\theta) = P(X_1 = x_1, ..., X_n = x_n; \theta) = \prod_{i=1}^n f(x_i; \theta)$$
 for  $\theta \in \Omega$ .

#### 66 Note

As discussed earlier, a natural estimate of  $\theta$  is the value which maximizes the probability of the observed sample, and we denote this notion as

$$\hat{\theta} = \argmax_{\theta \in \Omega} L(\theta)$$

and call  $\hat{\theta}$  the maximum likelihood estimate (ML estimate). In practice, it is often convenient to work, instead, with the natural logarithm of the likelihood function, in which we call the Log-likelihood:

$$\ell(\theta) = \ln L(\theta).$$

Note that the maximum likelihood estimate for both the likelihood function and the log-likelihood are the same since log is a strictly increasing function.

## Example 17.1.3 (Example 6.1)

Consider flipping a coin repeatedly, where for  $i \in \mathbb{N}$ 

$$X_i = egin{cases} 1 & ext{if the } i^{th} ext{ flip lands on heads} \ 0 & ext{otherwise} \end{cases}$$

Based on 4 independent flips, we are given that  $X_1, X_2, X_3, X_4 \sim \text{Bernoulli}(p)$  are IID rvs. The sample has been observed as  $x_1 = 1$ ,  $x_2 = 1$   $x_3 = 0$ ,  $x_4 = 1$ . Write the probability of this sample as a function of p. Then, compute the likelihood function and find the ML estimate.

#### Solution

Since the  $X_i$ 's are IID rvs, we have

$$f(x_1, x_2, x_3, x_4; p) = \prod_{i=1}^4 P(X_i = x_i; p)$$

$$= \prod_{i=1}^4 p^{x_i} (1-p)^{1-x_i} = p^{\sum_{i=1}^4 x_i} (1-p)^{\sum_{i=1}^4 (1-x_i)}.$$

Note that the likelihood function for p under the observed values is therefore

$$L(p) = p^3(1-p).$$

To find the ML estimate, we do

$$0 = \frac{dL(p)}{dp} = 3p^2 - 4p^3 = p^2(3 - 4p)$$
 (17.1)

and we observe that p=0 or  $\frac{3}{4}$ . Clearly,  $p=\frac{3}{4}$  maximizes the likelihood function.

#### 66 Note

We have been talking almost solely about the discrete case, so what about the continuous case? We have that P(X=x)=0, and so we consider a small neighbourhood of radius  $\delta>0$ . Then for a small  $\delta>0$  around any point  $x\in \operatorname{supp}(X)$ , we have that

$$P(x - \delta < X < x + \delta; \theta) = \int_{x - \delta}^{x + \delta} f(t; \theta) dt \approx 2\delta \cdot f(x; \theta).$$
 (17.2)

And so for an observed value x, since  $\delta$  is fixed in Equation (17.2), we have that the value of  $\theta$  that maximizes  $f(t;\theta)$  also maximizes  $2\delta \cdot f(t;\theta)$ .

#### \* Warning

We shall clarify the following notations: we use

- $\tilde{\theta}$  as the estimator, which is an rv; and
- $\hat{\theta}$  as an estimate, which is a fixed value.

#### Example 17.1.4 (Example 6.2 (Course Notes 6.2.4))

Recall the coin flip example in Example 17.1.3. Suppose  $X_1, ..., X_n \sim$ Bernoulli(p) is a sequence of IID rvs. Calculate the ML estimate of p.

#### Solution

Since  $X_1, ..., X_n$  are IID, let  $\vec{x} = (x_1 \ x_2 \ ... \ x_n)$ , then the joint distribution

$$f(\vec{x}; p) = \prod_{i=1}^{n} f(x_i; p)$$

$$= \prod_{i=1}^{n} p^{x_i} (1-p)^{1-x_i} = p^{\sum_{i=1}^{n} x_i} (1-p)^{\sum_{i=1}^{n} (1-x_i)}$$

To get an ML estimator, we shall use the log-likelihood:

$$\ell(p) = \left(\sum_{i=1}^{n} x_i\right) \ln p + \left(n - \sum_{i=1}^{n} x_i\right) \ln(1-p)$$

and so to get the ML estimator

$$0 = \frac{\partial}{\partial p} \ell(p) \Big|_{p=\hat{p}} = \left( \sum_{i=1}^{n} x_i \right) \left( \frac{1}{\hat{p}} \right) - \left( n - \sum_{i=1}^{n} x_i \right) \left( \frac{1}{1 - \hat{p}} \right)$$

$$= \frac{\sum x_i}{\hat{p}} - \frac{n - \sum x_i}{1 - \hat{p}} = \frac{(1 - \hat{p}) \sum x_i - n\hat{p} + \hat{p} \sum x_i}{\hat{p}(1 - \hat{p})}$$

$$= \frac{\sum x_i - n\hat{p}}{\hat{p}(1 - \hat{p})}$$

where we represent  $\sum_{i=1}^{n} x_i$  by  $\sum x_i$  for sanity. Thus we have that

$$n\hat{p} = \sum_{i=1}^{n} x_i \implies \hat{p} = \frac{1}{n} \sum_{i=1}^{n} x_i = \bar{x}$$

## Example 17.1.5 (Example 6.3 (Course Notes 6.2.5))

Suppose we have the data  $\vec{x} = (x_1 \ x_2 \ \dots \ x_n)$  for the sequence of IID rvs  $X_1, X_2, ..., X_n \sim Poi(\theta)$ . Find the likelihood function, log-likelihood, the ML

Note that when looking for an ML estimate, we should also check for the case when  $\left. \frac{\partial^2}{\partial p^2} \ell(p) \right|_{p=\hat{p}} < 0$  to ensure maximality (instead of minimality).

Note that the ML estimator in this case would be represented as

$$\tilde{p} = \overline{X}$$
.

estimate  $\hat{\theta}$ , and the ML estimator  $\tilde{\theta}$ .

# Solution

Since each of the  $X_i$ 's are IID, we have that their joint pmf, and in particular, their likelihood function, is

$$L(\theta) = \prod_{i=1}^{n} \frac{e^{-\theta} \theta^{x}}{x!} = e^{-n\theta} \theta^{\sum_{i=1}^{n} x_i} \left( \prod_{i=1}^{n} x! \right)^{-1}.$$

Then the log-likelihood is

$$\ell(\theta) = -n\theta + \left(\sum_{i=1}^{n} x_i\right) \log \theta - \sum_{i=1}^{n} \log(x_i!).$$

To get the ML estimate, note that

$$0 = \frac{\partial}{\partial \theta} \ell(\theta) \Big|_{\theta = \hat{\theta}} = -n + \frac{1}{\hat{\theta}} \sum_{i=1}^{n} x_i$$

and so

$$\hat{\theta} = \frac{1}{n} \sum_{i=1}^{n} x_i = \bar{x}.$$

Consequently, we have that the ML estimator is

$$\tilde{\theta} = \overline{X}$$
.

# 18 Lecture 18 Jul 3rd 2018

# **18.1** Estimation (Continued)

# 18.1.1 Maximum Likelihood Estimation (Continued)

# Definition 60 (Score Function)

The score function is defined as

$$S(\theta) = S(\theta; x) = \frac{d}{d\theta} \ell(\theta) = \frac{d}{d\theta} \ln L(\theta) \quad \theta \in \Omega.$$

## 66 Note

Notice that to find the ML estimate, we usually set the score function to zero and solve for  $S(\theta) = 0$ .

## **Definition 61 (Information Function)**

The *information function* is defined as

$$I(\theta) = I(\theta; x) = -\frac{d^2}{d\theta^2} \ell(\theta) = -\frac{d^2}{d\theta^2} \ln L(\theta) \quad \theta \in \Omega$$

If  $\hat{\theta}$  is the ML estimate of  $\theta$ , then  $l(\hat{\theta})$  is called the **observed information** 

# Example 18.1.1

Going back to our example in Example 17.1.4, we had that the likelihood function was

$$L(p) = p^{\sum_{i=1}^{n} x_i} (1-p)^{n-\sum_{i=1}^{n} x_i}$$

the log-likelihood was

$$\ell(p) = \left(\sum_{i=1}^{n} x_i\right) \ln p + \left(n - \sum_{i=1}^{n} x_i\right) \ln(1-p).$$

So the score function is

$$S(p) = \frac{d}{dp}\ell(p) = \frac{1}{p}\sum_{i=1}^{n}x_i - \frac{1}{1-p}\left(n - \sum_{i=1}^{n}x_i\right)$$

Consequently, the information function is

$$I(p) = -\frac{d^2}{dp^2}\ell(p) = -\frac{d}{dp}S(p)$$
$$= \frac{1}{p^2} \sum_{i=1}^n x_i + \frac{1}{(1-p)^2} \left(n - \sum_{i=1}^n x_i\right)$$

## 66 Note

Recall from the calculus knowledge of the 2nd derivative being the "curvature" of the original function, and so  $I(\theta)$  is a function about the curvature of the log-likelihood. In particular, we have that  $I(\hat{\theta})$  tells us about the convacity of the log-likelihood function at the ML estimate.

Note that the information function is a function that has two variables: the unknown parameter  $\theta$ , and the data  $\vec{X} = (X_1 \dots X_n)$ .

In a later lecture, we shall see how the observed information  $I(\hat{\theta})$  can be used to construct approximate confidence intervals for the unknown parameter  $\theta$ .

#### Definition 62 (Fisher Information)

If  $\theta$  is a scalar, then the **expected information**, or **Fisher information** (function) is given by

$$J(\theta) = E[I(\theta; \vec{X})] = E\left[-\frac{\partial^2}{\partial \theta^2}\ell(\theta; \vec{X})\right] \quad \theta \in \Omega$$

#### 66 Note

Just to take away the layers of definitions and compare the Fisher information with our pf for the rv(s), if  $X_1, ..., X_n$  is a random sample (i.e. IID rvs), each with pf  $f(x;\theta)$ , then

$$J(\theta) = E\left[-\frac{\partial^2}{\partial \theta^2}\ell(\theta; \vec{X})\right] = nE\left[-\frac{\partial^2}{\partial \theta^2}\ln f(\vec{X}; \theta)\right]$$

where  $\vec{X} = (X_1 \dots X_n)$ .

# Example 18.1.2 (Example 6.4 (Course Notes 6.2.10))

Suppose  $X_1, ..., X_n \sim Bernoulli(p)$  is a sequence of IID rvs. We have showed in Example 17.1.4 that the ML estimator of p is  $\tilde{p} = \overline{X}$ . Calculate the Fisher information and compare it with the variance of the ML estimator of p.

## Solution

Recall that from Example 18.1.1, we had

$$\ell(p) = \left(\sum_{i=1}^{n} x_i\right) \ln p + \left(n - \sum_{i=1}^{n} x_i\right) \ln(1-p)$$

$$S(p) = \frac{1}{p} \sum_{i=1}^{n} x_i - \frac{1}{1-p} \left(n - \sum_{i=1}^{n} x_i\right)$$

$$I(p) = \frac{1}{p^2} \sum_{i=1}^{n} x_i + \frac{1}{(1-p)^2} \left(n - \sum_{i=1}^{n} x_i\right),$$

So the Fisher information is

$$E[I(p)] = E\left[\frac{1}{p^2} \sum_{i=1}^n X_i + \frac{1}{(1-p)^2} \left(n - \sum_{i=1}^n X_i\right)\right]$$

$$= \frac{1}{p^2} \sum_{i=1}^n E[X_i] + \frac{1}{(1-p)^2} \left(n - \sum_{i=1}^n E[X_i]\right)$$

$$= \frac{1}{p^2} (np) + \frac{1}{(1-p)^2} (n-np)$$

$$= \frac{n}{p} + \frac{n}{1-p} = \frac{n}{p(1-p)}.$$

On the other hand, note that the variance of the ML estimator is

$$\operatorname{Var}(\tilde{p}) = \operatorname{Var}(\overline{X}) = \frac{1}{n} \operatorname{Var}(X_i) = \frac{p(1-p)}{n}.$$

## Example 18.1.3 (Example 6.5 (Course Notes 6.2.10))

Suppose  $X_1, ..., X_n \sim \text{Poi}(\theta)$  is a sequence of IID rvs. We showed in Example 17.1.5 the ML estimator of  $\theta$  is  $\tilde{\theta} = \overline{X}$ . Calculate the Fisher information and compare it with the variance of the ML estimator of  $\theta$ .

## Solution

We had the that log-likelihood is

$$\ell(\theta) = -n\theta + \left(\sum_{i=1}^{n} x_i\right) \log \theta - \sum_{i=1}^{n} \log(x_i!).$$

So the score function is

$$S(\theta) = -n + \frac{1}{\theta} \sum_{i=1}^{n} x_i,$$

and so the information function is

$$I(\theta; x) = \frac{1}{\theta^2} \sum_{i=1}^n x_i.$$

Then the Fisher information is

$$E[I(\theta; \vec{X})] = \frac{n}{\theta^2} E[X_i] = \frac{n}{\theta}.$$

*The variance of the ML estimator of*  $\theta$  *is* 

$$\operatorname{Var}(\tilde{\theta}) = \frac{1}{n} \operatorname{Var}(X_i) = \frac{\theta}{n}.$$

## 66 Note

We observe from the above two examples that

$$J(\theta) = \frac{1}{\operatorname{Var}\tilde{\theta}}$$

where  $\tilde{\theta} = \overline{X}$ . It is tempting to verify if this is the case in general, but there does not seem to be reliable sources that points to this case (at least, online searches have yet to yield me results). The lecture claims that this is only true for certain families of distributions.

## Example 18.1.4 (Example 6.6 (Course Notes 6.2.12))

Suppose  $X_1, ..., X_n$  is a random sample, each from a distribution with pdf

$$f(x;\theta) = \theta x^{\theta-1}$$
  $0 \le x \le 1$ ,  $\theta > 0$ .

Find the score function, the ML estimator of  $\theta$ , the information function and the observed information.

## Solution

Since  $X_1, ..., X_n$  are IID, we have

$$L(\theta) = \prod_{i=1}^{n} \theta x_i^{\theta-1} = \theta^n \left( \prod_{i=1}^{n} x_i \right)^{\theta-1}.$$

So the log-likelihood is

$$\ell(\theta) = n \ln \theta + (\theta - 1) \ln \left( \prod_{i=1}^{n} x_i \right) = n \ln \theta + (\theta - 1) \sum_{i=1}^{n} \ln x_i.$$

The score function is therefore

$$S(\theta) = \frac{n}{\theta} + \sum_{i=1}^{n} \ln x_i,$$

and so the ML estimator of  $\theta$  is

$$\tilde{\theta} = -\frac{n}{\sum_{i=1}^{n} \ln x_i}.$$

From the score function, the information function is

$$I(\theta) = \frac{n}{\theta^2},$$

and so the observed information is

$$I(\hat{\theta}) = \frac{n}{\left(-\frac{n}{\sum_{i=1}^{n} \ln x_i}\right)^2} = \frac{1}{n} \left(\sum_{i=1}^{n} \ln x_i\right)^2$$

# 66 Note

In the above example, note that the Fisher information is

$$J(\theta) = E[I(\theta; \vec{X})] = \frac{n}{\theta^2}.$$

## Example 18.1.5 (Example 6.7 (Course Notes 6.2.13))

Suppose  $X_1, ..., X_n \sim \text{Unif}(0, \theta)$  is a random sample. Find the ML estimator of  $\theta$ .

## Solution

Note that  $f_{X_i}(x_i;\theta) = \frac{1}{\theta}$  for  $0 \le x_i \le \theta$ . So the likelihood function is

$$L(\theta) = \prod_{i=1}^{n} \frac{1}{\theta} \mathbb{1}_{\{0 \le x_i \le \theta\}}$$

$$= \theta^{-n} \mathbb{1}_{\{0 \le x_1 \le \theta, 0 \le x_2 \le \theta, \dots, 0 \le x_n \le \theta\}}$$

$$= \theta^{-n} \mathbb{1}_{\{0 \le x_{(n)} \le \theta\}}$$

where we use the order statistics notation to simplify the equation. Notice that if we take the derivative of the likelihood function from this point and try to get the ML estimate/estimator, we would run into the following equation:

$$-\frac{n}{\theta^{n+1}}=0$$

in which we would have trouble getting the MLE.

# **Example 18.1.6 (Example 6.8)**

Suppose  $X_1, ..., X_n$  is a random sample, each from the Unif $(\theta, \theta + 1)$  distribution. Find the ML estimator of  $\theta$ .

## Solution

This time, we have that the pdf for each  $X_i$  is

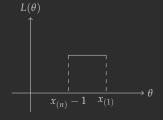
$$f_{X_i}(x_i;\theta) = \mathbb{1}_{\{\theta < x_i < \theta + 1\}}.$$

Then the likelihood function is

$$L(\theta) = \prod_{i=1}^{n} f(x_i; \theta) = \mathbb{1}_{\{\theta \le x_1 \le \theta + 1, \dots, \theta \le x_n \le \theta + 1\}}$$
$$= \mathbb{1}_{\{\theta \le x_{(1)}\}} \mathbb{1}_{\{x_{(n)} \le \theta + 1\}}$$
$$= \mathbb{1}_{\{\theta \le x_{(1)}\}} \mathbb{1}_{\{x_{(n)} - 1 \le \theta\}} = \mathbb{1}_{\{x_{(n)} - 1 \le \theta \le x_{(1)}\}}$$

where we, once again, use the order statistics notation. Consequently, we have the graph of  $\theta$  versus  $L(\theta)$  on the right.

Thus we observe that the MLE is not unique, since it can be any number between  $x_{(n)} - 1$  and  $x_{(1)}$ .



# **■** Theorem 61 (Invariance Property of the MLE)

Suppose  $\tau = h(\theta)$  is an injective function of  $\theta$ . Suppose also that  $\hat{\theta}$  is the *ML* estimator of  $\theta$ . Then  $\hat{\tau} = h(\theta)$  is the *ML* estimator of  $\tau$ .

# \* Warning

We are short on certain tools to actually prove this theorem.

# **Example 18.1.7 (Example 6.9)**

Suppose  $X_1,...,X_n \sim f(x;\theta) = \theta x^{\theta-1} \mathbb{1}_{\{0 < x < 1\}}$ , where  $\theta > 0$ . Find the MLE of the median of the distribution.

# Solution

Recall from Example 18.1.4 that we had

$$\hat{\theta} = -\frac{n}{\sum\limits_{i=1}^{n} \ln x_i}.$$

Let m be the median. The goal is to use 🖳 Theorem 61. Something feels off...

# 19 Lecture 19 July 5th 2018

# 19.1 Estimation (Continued 2)

## 19.1.1 Maximum Likelihood Estimation (Continued 2)

# Definition 63 (Relative Likelihood)

Suppose  $X_1, ..., X_n \sim f(x; \theta)$  is an IID sequence of rvs where the likelihood function is  $L(\theta)$  and the MLE of  $\theta$  is  $\hat{\theta}$ . The **relative likelihood** function is defined by

$$R(\theta) = R(\theta; x) = \frac{L(\theta)}{L(\hat{\theta})}, \quad \theta \in \Omega$$

## 66 Note

Note that since  $L(\hat{\theta})$  is maximum since  $\hat{\theta}$  is the maximum likelihood estimate, we have that  $L(\theta) \leq L(\hat{\theta})$ , and so<sup>1</sup>

$$0 \le R(\theta) = \frac{L(\theta)}{L(\hat{\theta})} \le 1.$$

 $^1$  Note that the likelihood function  $L(\theta)=f(x;\theta)$  is nonnegative, since the cdf of a distribution is nondecreasing, and so the pf must be nonnegative.

## 66 Note

Given the data  $\theta = (x_1 \dots x_n)$ , we say that

- if  $R(\theta)$  is small (usual threshold is 0.1), then we say that  $\theta$  is **implausible**;
- if  $R(\theta)$  is large (usual threshold is 0.5), then we say that  $\theta$  is **plausible**

# Definition 64 (Likelihood Region & Likelihood Interval)

The set of  $\theta$  values for which  $R(\theta) \geq p$  is called a 100p% **likelihood** region for  $\theta$ . If the region is an interval of real values, then it is called a 100p% **likelihood interval** (LI) for  $\theta$ .

#### 66 Note

*Using the definition, we can extend on the previous note:* 

- Values inside the 10% LI are referred to as **plausible** and values outside of this interval as **implausible**;
- Values inside a 50% LI are very plausible; and
- Values outside a 1% LI are very implausible in light of the data.

As how we had the log-likelihood from the likelihood function, we can have the log relative likelihood.

# Definition 65 (Log Relative Likelihood)

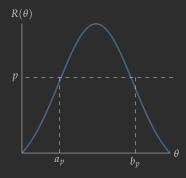
The *log relative likelihood* is the natural logarithm of the relative *likelihood function*:

$$r(\theta) = r(\theta; x) = \ln[R(\theta)] = \log[L(\theta)] - \log[L(\hat{\theta})] = \ell(\theta) - \ell(\hat{\theta})$$

*for*  $\theta \in \Omega$ .

## 66 Note (Unimodality)

If for  $\theta \in \Omega$ ,  $R(\theta)$  is unimodal as shown in the graph to the right, then  $R(\theta) \geq p$  will give us the likelihood interval, in which this case we have that the interval is  $[a_p, b_p]$ .



## **Example 19.1.1 (Example 6.10)**

Let  $X_1,...,X_{100} \sim Poi(\theta)$  be an IID sequence of rvs. Based on the observed values, you are given that  $\sum_{i=1}^{100} x_i = 980$ . Find a 10% and 50% likelihood *intervals for*  $\theta$ *.* 

#### Solution

From Example 17.1.5, we had

$$L(\theta) = e^{-n\theta} \theta^{n\bar{x}} \left( \prod_{i=1}^n x_i! \right)^{-1}$$
 and  $\hat{\theta} = \bar{x}$ .

The relative likelihood is therefore

$$R(\theta) = \frac{L(\theta)}{L(\hat{\theta})} = \frac{e^{-n\theta}\theta^{n\bar{x}}}{e^{-n\hat{\theta}\hat{\theta}^{n\bar{x}}}} \left(\prod_{i=1}^n x_i!\right)^{-1} \left(\prod_{i=1}^n x_i!\right) = e^{-n(\theta-\bar{x})} \left(\frac{\theta}{\bar{x}}\right)^{n\bar{x}}.$$

We were given that

$$\bar{x} = \frac{1}{100} \sum_{i=1}^{100} x_i = \frac{980}{100} = 9.8,$$

and so get

$$R(\theta) = e^{-100(\theta - 9.8)} \left(\frac{\theta}{9.8}\right)^{980}.$$

Using Wolfram Alpha, we get that the 50% likelihood interval for  $\theta$  is [9.436, 10.173]. Unfortunately so, I cannot get the solution for the 10% likelihood interval since the allowed free computation time is exceeded, and so I shall leave the likelihood interval as

$$\left\{ \theta \in \Omega : e^{-100(\theta - 9.8)} \left( \frac{\theta}{9.8} \right)^{980} > 0.1 \right\}.$$

Example 19.1.2 (Example 6.11 (Course Notes 6.2.24))

Suppose  $X_1, ..., X_n$  is a random sample from  $Exp(1, \theta)$ . Plot the relative likelihood function for  $\theta$  if n = 20 and  $x_{(1)} = \min\{x_1, ..., x_n\} = 1$ . Find the 10% and 50% likelihood intervals for  $\theta$ .

#### Solution

Recall the definition of the Double Parameter Exponential Distribution. We have that the pdf of each  $X_i$  is

$$f_{X_i}(x_i) = e^{-\frac{x_i - \theta}{1}} = e^{-(x_i - \theta)}$$

Since the  $X_i$ 's form a random sample, we have that

$$L(\theta) = \prod_{i=1}^{n} e^{-(x_i - \theta)} \mathbb{1}_{\{x_i \ge \theta\}} = e^{-\sum_{i=1}^{n} x_i + n\theta} \mathbb{1}_{\{x_{(1)} \ge \theta\}}.$$

Firstly, note that  $L(\theta)$  is an increasing function with respect to  $\theta$ . Now, note that if the indicator function is 1, we would have that  $\hat{\theta}_0 = \bar{x}$ , where we denote the "supposed" ML estimate as  $\hat{\theta}_0$ . However, since the condition for the indicator function to be 1 is that  $\theta \leq x_{(1)}$ , we must then have that  $\hat{\theta} = x_{(1)}$ , since  $x_{(1)} \leq \bar{x}$  as  $x_{(1)} = \min\{x_1, ..., x_n\}$ .

With that, we can calculate the relative likelihood function:

$$R(\theta) = \frac{L(\theta)}{L(\hat{\theta})} = e^{n(\theta - x_{(1)})} \mathbb{1}_{\{x_{(1)} \ge \theta\}}.$$

Given n = 20 and  $x_{(1)} = 1$ , we have that

$$R(\theta) = e^{20(\theta-1)} \mathbb{1}_{\{\theta < 1\}}.$$

To get the 50% likelihood interval, we have

$$e^{20(\theta-1)} \ge 0.5$$
  
 $\theta \ge 1 + \frac{1}{20} \ln 0.5$ .

Note that  $\ln 0.5 < 0$ , so we have that the 50% likelihood interval is  $[1 + 0.05 \ln 0.5, 1]$ . For the 10% likelihood interval, we have

$$\begin{split} e^{20(\theta-1)} &\geq 0.1 \\ \theta &\geq 1 + \frac{1}{20} \ln 0.1 \; . \end{split}$$

Again, note that  $\ln 0.1 < 0$ , and so we have that the 10% likelihood interval is  $[1 + 0.05 \ln 0.1, 1]$ .

We can extend our methods analogously for the multiparameter case. Assume  $\theta = (\theta_1 \dots \theta_k)^T$ . Then the likelihood function,  $L(\theta)$ , is a function of k parameters. Then the log-likelihood function is  $\ell(\theta) = \ln L(\theta)$ . From there, we can then use the same approach to get the ML estimate of  $\theta$ ,

$$\hat{\theta} = (\hat{\theta}_1 \dots \hat{\theta}_k)^T = \underset{\theta_1, \dots, \theta_k}{\operatorname{arg max}} \ell(\theta).$$

We shall also note that we can extend **P** Theorem 61 to the multi-

variate case.

Furthermore, we will restate some of the definitions that were stated in univariate context for the multivariate case.

## Definition 66 (Score Vector)

Suppose  $X_1, ..., X_n \sim f(x; \theta)$  is a sequence of IID rvs, where  $\theta =$  $(\theta_1 \ldots \theta_k)^T$ . The **score vector** is defined as

$$S(\theta) = S(\theta; x) = \left[\frac{\partial \ell}{\partial \theta_1} \dots \frac{\partial \ell}{\partial \theta_k}\right]^T$$

## Definition 67 (Information Matrix)

Suppose  $X_1, ..., X_n \sim f(x; \theta)$  be a sequence of IID rvs, where  $\theta =$  $(\theta_1 \ldots \theta_k)^T$ . The information matrix,  $I(\theta) = I(\theta; x)$ , is a  $k \times k$ symmetric matrix whose (i, j) entry is given by

$$-\frac{\partial^2}{\partial \theta_i \partial \theta_j} \ell(\theta).$$

The constant matrix  $I(\hat{\theta})$  is called the observed information matrix.

#### Definition 68 (Fisher Information Matrix)

Suppose  $X_1, ..., X_n \sim f(x; \theta)$  is a sequence of IID rvs, where  $\theta =$  $(\theta_1 \ldots \theta_k)^T$ . The **Fisher information matrix**,  $J(\theta)$  is a  $k \times k$  symmetric matrix whose (i, j) entry is given by

$$E\left[-\frac{\partial^2}{\partial \theta_i \partial \theta_j}\ell(\theta;X)\right]$$
,

where  $X = (X_1 \dots X_k)^T$ .

## Definition 69 (Likelihood Region - Multivariate)

Suppose  $X_1, ..., X_n \sim f(x; \theta)$  is a sequence of IID rvs, where  $\theta =$  $(\theta_1 \ldots \theta_k)^T$ . The set of  $\theta = (\theta_1 \ldots \theta_k)^T$  for which  $R(\theta) \geq p$ , for some  $0 \le p \le 1$ , is called a 100p% **likelihood region** for  $\theta$ , which is a region in a k-dimensional space  $\mathbb{R}^k$ .

#### 66 Note

Recall that in the univariate case, we can certainly check if our choice of the  $\hat{\theta}$  is indeed a maximizing value by taking the second derivative of  $\ell(\theta)$  or  $L(\theta)$ , and checking that  $I(\hat{\theta}) > 0$ . The analogous approach for the multivariate case is to ensure that the **observed information matrix**,  $I(\hat{\theta})$  is positive definite<sup>2</sup>.

We have two methods to verify that a matrix M is positive definite:

- $\det M > 0$ ;
- eigenvalues of M are all positive.

<sup>2</sup> A matrix M is called a **positive definite** if the scalar  $z^T M z > 0$  for every non-zero column vector z.

#### Example 19.1.3 (Example 6.12 Part 1)

Suppose  $X_1,...,X_n \sim N(\mu,\sigma^2)$  is a sequence of IID rvs.

- 1. Derive the MLE of the parameters  $\mu$  and  $\sigma^2$ .
- 2. Derive the MLE of the coefficient of variation  $CV = \frac{\sigma}{\mu}$ .

#### Solution

1. Since the  $X_i$ 's are IID, we have

$$L(\mu, \sigma^2) = (2\pi)^{-\frac{n}{2}} (\sigma^2)^{-\frac{n}{2}} \exp\left\{-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2\right\}.$$

Then the log-likelihood is

$$\ell(\mu, \sigma^2) = -\frac{n}{2} \ln(2\pi) - \frac{n}{2} \log \sigma^2 - \frac{1}{2\sigma^2} \sum_{i=1}^{n} (x_i - \mu)^2.$$

*Now, note that the first partial derivatives of*  $\ell(\mu, \sigma^2)$  *are* 

$$\frac{\partial \ell}{\partial \mu} = \frac{1}{\sigma^2} \sum_{i=1}^n (x_i - \mu) = \frac{1}{\sigma^2} \left( \sum_{i=1}^n x_i - n\mu \right)$$
$$\frac{\partial \ell}{\partial \sigma^2} = -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_{i=1}^n (x_i - \mu)^2.$$

So the score vector is

$$S\left(\mu,\sigma^{2}\right) = \begin{bmatrix} \frac{\partial\ell}{\partial\mu} \\ \frac{\partial\ell}{\partial\sigma^{2}} \end{bmatrix} = \begin{bmatrix} \frac{1}{\sigma^{2}} \left(\sum_{i=1}^{n} x_{i} - n\mu\right) \\ -\frac{n}{2\sigma^{2}} + \frac{1}{2\sigma^{4}} \sum_{i=1}^{n} (x_{i} - \mu)^{2} \end{bmatrix}.$$

So the possible candidates for  $\hat{\mu}$  and  $\hat{\sigma}^2$  is

$$\hat{\mu} = \bar{x}$$
  $\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2$ .

It remains to verify that  $\hat{n}u$  and  $\hat{\sigma}^2$  are indeed the "maximizers", note that

$$\frac{\partial^2 \ell}{\partial \mu^2} = -\frac{n}{\sigma^2}$$

$$\frac{\partial^2 \ell}{\partial (\sigma^2)^2} = \frac{n}{2\sigma^4} - \frac{1}{\sigma^8} \sum_{i=1}^n (x_i - \mu)^2$$

$$\frac{\partial^2 \ell}{\partial \mu \partial \sigma^2} = \frac{n\mu}{\sigma^4} - \frac{1}{\sigma^4} \sum_{i=1}^n x_i,$$

and so the information matrix is

$$I(\mu, \sigma^{2}) = \begin{bmatrix} \frac{n}{\sigma^{2}} & \frac{n\mu}{\sigma^{4}} - \frac{1}{\sigma^{4}} \sum_{i=1}^{n} x_{i} \\ \frac{n\mu}{\sigma^{4}} - \frac{1}{\sigma^{4}} \sum_{i=1}^{n} x_{i} & \frac{1}{\sigma^{8}} \sum_{i=1}^{n} (x_{i} - \mu)^{2} - \frac{n}{2\sigma^{4}} \end{bmatrix}.$$

Then the observed matrix is

$$I(\hat{\mu}, \hat{\sigma}^{2}) = \begin{bmatrix} \frac{n}{\hat{\sigma}^{2}} & \frac{n\bar{x}}{(\hat{\sigma}^{2})^{2}} - \frac{1}{(\hat{\sigma}^{2})^{2}} \sum_{i=1}^{n} x_{i} \\ \frac{n\bar{x}}{(\hat{\sigma}^{2})^{2}} - \frac{1}{(\hat{\sigma}^{2})^{2}} \sum_{i=1}^{n} x_{i} & \frac{1}{(\hat{\sigma}^{2})^{4}} \sum_{i=1}^{n} (x_{i} - \bar{x})^{2} - \frac{n}{2(\hat{\sigma}^{2})^{2}} \end{bmatrix}$$

$$= \begin{bmatrix} \frac{n}{\hat{\sigma}^{2}} & 0 \\ 0 & \frac{1}{(\hat{\sigma}^{2})^{2}} (\hat{\sigma}^{2}) - \frac{n}{2(\hat{\sigma}^{2})^{2}} \end{bmatrix}$$

$$= \begin{bmatrix} \frac{n}{\hat{\sigma}^{2}} & 0 \\ 0 & \frac{n}{(\hat{\sigma}^{2})^{2}} \end{bmatrix}.$$

We have that the determinant of the observed information matrix is

$$\det I\left(\hat{\mu}, \hat{\sigma}^2\right) = \frac{n}{\hat{\sigma}^2} \cdot \frac{n}{\left(\hat{\sigma}^2\right)^2} > 0$$

since  $\hat{\sigma}^2 > 0$ .

2. By  $\blacksquare$  Theorem 61, we have that the MLE of the coefficient of variation, CV,  $is^3$ 

$$\hat{CV} = \frac{\hat{\sigma}}{\hat{\mu}} = \frac{1}{\bar{x}} \sqrt{\frac{\sum_{i=1}^{n} (x_i - \bar{x})^2}{n}} = \frac{\hat{\sigma}}{\sqrt{n}\bar{x}}$$

³ Note that the function  $f(x,y) = \frac{x}{y}$  is not injective if x and y are not coprime. Since both  $\hat{\sigma}$  and  $\hat{\mu}$  has a common multiple of  $\sqrt{\frac{1}{n}}$ , they are not coprime, so **I** do not think we can actually invoke ■ Theorem 61. I shall, however, leave this answer here from the lecture.

# 20 Lecture 20 Jul 10th 2018

# **20.1** *Estimation* (Continued 3)

## **20.1.1** *Maximum Likelihood Estimation (Continued 3)*

# Example 20.1.1 (Example 6.12 Part 2)

Continuing with Example 19.1.3, answer the following questions:

• Derive the Fisher information for the parameter vector  $(\mu, \sigma^2)$ .

#### Solution

Recall that the information matrix was

$$I(\mu, \sigma^{2}) = \begin{bmatrix} \frac{n}{\sigma^{2}} & \frac{n\mu}{\sigma^{4}} - \frac{1}{\sigma^{4}} \sum_{i=1}^{n} x_{i} \\ \frac{n\mu}{\sigma^{4}} - \frac{1}{\sigma^{4}} \sum_{i=1}^{n} x_{i} & \frac{1}{\sigma^{8}} \sum_{i=1}^{n} (x_{i} - \mu)^{2} - \frac{n}{2\sigma^{4}} \end{bmatrix}.$$

Thus the Fisher information matrix is

$$J(\mu, \sigma^2) = \begin{bmatrix} E\left(\frac{n}{\sigma^2}\right) & E\left[\frac{n\mu}{\sigma^4} - \frac{1}{\sigma^4}\sum_{i=1}^n x_i\right] \\ E\left[\frac{n\mu}{\sigma^4} - \frac{1}{\sigma^4}\sum_{i=1}^n x_i\right] & E\left[\frac{1}{\sigma^8}\sum_{i=1}^n (x_i - \mu)^2 - \frac{n}{2\sigma^4}\right] \end{bmatrix}$$
$$= \begin{bmatrix} \frac{n}{\sigma^2} & \frac{n\mu}{\sigma^4} - \frac{1}{\sigma^4}\sum_{i=1}^n E[X_i] = 0 \\ 0 & \frac{n}{\sigma^6} - \frac{n}{2\sigma^4} \end{bmatrix}.$$

#### 66 Note

Notice that we have, yet again, a similar scenario as when we first intro-

duced the Fisher information in the univariate case: notice that

$$\tilde{\mu} = \overline{X} \implies \operatorname{Var}(\tilde{\mu}) = \operatorname{Var}(\overline{X}) = \frac{\sigma^2}{n} = \frac{1}{J(\mu, \sigma^2)_{11}}$$

and

$$\tilde{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2 \implies \text{Var}\left(\tilde{\sigma}^2\right) = \frac{1}{n^2} \text{Var}\left((n-1)S^2\right)$$
$$= \frac{\sigma^4}{n^2} \text{Var}\left(\frac{n-1}{\sigma^2}S^2\right) = \frac{\sigma^4 2(n-1)}{n^2}$$

by Cochran's Theorem. Note that  $Var(\hat{\sigma}^2) \approx J(\mu, \sigma^2)_{22}$ .

#### Example 20.1.2 (Example 6.13 (Course Notes 6.3.8))

Suppose  $Y_i \sim N(\alpha + \beta x_i, \sigma^2)$ , i = 1, ..., n, is a sequence of IID rvs. Show that the ML estimators of  $\alpha$ ,  $\beta$  and  $\sigma^2$  are given by

$$\tilde{\alpha} = \overline{Y} - \tilde{\beta}\bar{x}$$

$$\tilde{\beta} = \frac{\sum_{i=1}^{n} (x_i - \bar{x})(Y_i - \overline{Y})}{\sum_{i=1}^{n} (x_i - \bar{x})^2}$$

$$\tilde{\sigma}^2 = \frac{1}{n} \sum_{i=1}^{n} (Y_i - \tilde{\alpha} - \tilde{\beta}x_i)^2$$

#### Solution

Since the  $Y_i$ 's constitute a random sample, the likelihood function is

$$L(\alpha, \beta, \sigma) = \prod_{i=1}^{n} \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{(y_i - \alpha - \beta x_i)^2}{2\sigma^2}\right\}$$
$$= (2\pi)^{-\frac{n}{2}} \left(\sigma^2\right)^{-\frac{n}{2}} \exp\left\{-\frac{1}{2\sigma^2} \sum_{i=1}^{n} (y_i - \alpha - \beta x_i)^2\right\}$$

and consequently the log-likelihood is

$$\ell(\alpha, \beta, \sigma) = -\frac{n}{2} \ln(2\pi) - \frac{n}{2} \ln \sigma^2 - \frac{1}{2\sigma^2} \sum_{i=1}^{n} (y_i - \alpha - \beta x_i)^2.$$

Note that  $\tilde{\alpha}$  and  $\tilde{\beta}$  are also the **least** squares estimators of  $\alpha$  and  $\beta$ .

The first partial derivatives are

$$\begin{split} \frac{\partial \ell}{\partial \alpha} &= \frac{1}{\sigma^2} \sum_{i=1}^n (y_i - \alpha - \beta x_i) = \frac{n}{\sigma^2} \left[ \bar{y} - \alpha - \beta \bar{x} \right] \\ \frac{\partial \ell}{\partial \beta} &= \frac{1}{\sigma^2} \sum_{i=1}^n \left[ x_i (y_i - \alpha - \beta x_i) \right] = \frac{1}{\sigma^2} \left[ \sum_{i=1}^n x_i y_i - \alpha \sum_{i=1}^n x_i - \beta \sum_{i=1}^n x_i^2 \right] \\ \frac{\partial \ell}{\partial \sigma} &= -\frac{n}{\sigma} + \frac{1}{\sigma^3} \sum_{i=1}^n (y_i - \alpha - \beta x_i)^2 \,. \end{split}$$

The score vector is therefore

$$S(\alpha, \beta, \sigma) = \begin{bmatrix} \frac{n}{\sigma^2} \left[ \bar{y} - \alpha - \beta \bar{x} \right] \\ \frac{1}{\sigma^2} \left[ \sum_{i=1}^n x_i y_i - \alpha \sum_{i=1}^n x_i - \beta \sum_{i=1}^n x_i^2 \right] \\ \frac{1}{\sigma^3} \sum_{i=1}^n (y_i - \alpha - \beta x_i)^2 - \frac{n}{\sigma} \end{bmatrix}.$$

To get the candidates for respective MLE,

$$\hat{\alpha} = \bar{y} - \hat{\beta}\bar{x}$$

$$\hat{\beta} = \frac{\sum_{i=1}^{n} x_{i} y_{i} - \hat{\alpha} \sum_{i=1}^{n} x_{i}}{\sum_{i=1}^{n} x_{i}^{2}} = \frac{\sum_{i=1}^{n} x_{i} y_{i} - (\bar{y} - \hat{\beta}\bar{x}) \sum_{i=1}^{n} x_{i}}{\sum_{i=1}^{n} x_{i}^{2}}$$

$$= \frac{\sum_{i=1}^{n} x_{i} y_{i} - n\bar{x}\bar{y} + n\hat{\beta}\bar{x}^{2}}{\sum_{i=1}^{n} x_{i}^{2} - n\bar{x}\bar{y}}$$

$$= \frac{\sum_{i=1}^{n} x_{i} y_{i} - n\bar{x}\bar{y}}{\sum_{i=1}^{n} x_{i}^{2}} = \frac{\sum_{i=1}^{n} x_{i} y_{i} - n\bar{x}\bar{y}}{\sum_{i=1}^{n} x_{i}^{2} - n\bar{x}^{2}}$$

$$\hat{\sigma}^{2} = \frac{1}{n} \sum_{i=1}^{n} (y_{i} - \alpha - \beta x_{i})^{2}$$

For  $\hat{\beta}$ , note that

$$\sum_{i=1}^{n} (x_i - \bar{x})(y_i - \bar{y}) = \sum_{i=1}^{n} [x_i y_i - \bar{x} y_i - \bar{y} x_i + \bar{x} \bar{y}]$$

$$= \sum_{i=1}^{n} x_i y_i - \bar{x} \sum_{i=1}^{n} y_i - \bar{y} \sum_{i=1}^{n} x_i + n \bar{x} \bar{y}$$

$$= \sum_{i=1}^{n} x_i y_i - \bar{x} (n \bar{y}) - \bar{y} (n \bar{x}) + n \bar{x} \bar{y}$$

$$= \sum_{i=1}^{n} x_i y_i - n \bar{x} \bar{y}$$

and similarly so we have

$$\sum_{i=1}^{n} (x_i - \bar{x})^2 = \sum_{i=1}^{n} x_i^2 - n\bar{x}^2.$$

**Therefore** 

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

Thus the potential ML esimators are

$$\tilde{\alpha} = \overline{Y} - \tilde{\beta}\overline{x}$$

$$\tilde{\beta} = \frac{\sum_{i=1}^{n} (x_i - x)(Y_i - \overline{Y})}{\sum_{i=1}^{n} (x_i - \overline{x})^2}$$

$$\tilde{\sigma}^2 = \frac{1}{n} \sum_{i=1}^{n} (Y_i - \tilde{\alpha} - \tilde{\beta}x_i)^2.$$

To verify that these are indeed maximizing estimates, we need the following second order partial derivatives:

$$\begin{split} \frac{\partial^2 \ell}{\partial \alpha^2} &= -\frac{n}{\sigma^2} \\ \frac{\partial^2 \ell}{\partial \beta \partial \alpha} &= -\frac{n\bar{x}}{\sigma^2} \\ \frac{\partial^2 \ell}{\partial \sigma \partial \alpha} &= -\frac{2n}{\sigma^3} \left[ \bar{y} - \alpha - \beta \bar{x} \right] \\ \frac{\partial^2 \ell}{\partial \alpha \partial \beta} &= -\frac{n\bar{x}}{\sigma^2} \\ \frac{\partial^2 \ell}{\partial \beta^2} &= -\frac{1}{\sigma^2} \sum_{i=1}^n x_i^2 \\ \frac{\partial^2 \ell}{\partial \sigma \partial \beta} &= -\frac{2}{\sigma^3} \left( \sum x_i y_i - \alpha \sum x_i - \beta \sum x_i^2 \right) \\ \frac{\partial^2 \ell}{\partial \alpha \partial \sigma} &= -\frac{2n}{\sigma^3} \left[ \bar{y} - \alpha - \beta \bar{x} \right] \\ \frac{\partial^2 \ell}{\partial \beta \partial \sigma} &= -\frac{2}{\sigma^3} \left( \sum x_i y_i - \alpha \sum x_i - \beta \sum x_i^2 \right) \\ \frac{\partial^2 \ell}{\partial \beta \partial \sigma} &= -\frac{2}{\sigma^3} \left( \sum x_i y_i - \alpha \sum x_i - \beta \sum x_i^2 \right) \\ \frac{\partial^2 \ell}{\partial \sigma^2} &= \frac{n}{\sigma^2} - \frac{3}{\sigma^4} \sum_{i=1}^n \left( y_i - \alpha - \beta x_i \right)^2 . \end{split}$$

It remains to show that the observed information matrix is positive definite, which will not be shown since the working is too tedious.

# 20.1.2 Asymptotic Properties of ML Estimators

# ■ Theorem 62 (Asymptotic Distribution of the ML Estimator)

Suppose  $X = (X_1 ... X_n)$  be a random sample from  $f(x;\theta)$ . Let  $\tilde{\theta}_n = \tilde{\theta}_n(X_1,...,X_n)$  be the ML estimator of  $\theta$  based on X. Then under certain regularity conditions<sup>1</sup>:

1. Consistency:

$$\tilde{\theta}_n \stackrel{P}{\to} \theta.$$
 (20.1)

2. Asymptotic Normality:

$$\sqrt{J(\theta)}(\tilde{\theta}_n - \theta) \stackrel{D}{\to} Z \sim N(0, 1).$$
 (20.2)

3. Asymptotic Distribution of Relative Likelihood:

$$-2\ln R(\theta;X) = 2[\ell(\tilde{\theta}_n;X) - \ell(\theta;X)] \xrightarrow{D} W \sim \chi^2(1)$$

where  $\theta$  is the unknown true parameter. If **consistency** holds, then we call  $\hat{\theta}_n$  a consistent estimator of  $\theta$ .

For a proof of the above theorem see Casella & Berger 2002<sup>2</sup>.

<sup>2</sup> R. L. Casella, G. Berger. Statistical *Inference*. Thomson Learning, 2002

#### 66 Note

The above theorem implies that for sufficiently large n,  $\hat{\theta}_n$  has, approximately, an N  $(\theta, [J(\theta)]^{-1})$  distribution (using methods from

• Proposition 58).  $[I(\theta)]^{-1}$  is called the asymptotic variance of  $\hat{\theta}_n$ . Consequently, for sufficiently large n, we have that

$$\operatorname{Var}(\hat{\theta}_n) \approx [J(\theta)]^{-1}.$$

Note that  $J(\theta)$  is unknown since  $\theta$  is unknown.<sup>3</sup> Since  $\tilde{\theta}_n \stackrel{P}{\to} \theta$ , by

• Proposition 58, we have  $J(\tilde{\theta}_n) \stackrel{P}{\to} J(\theta)$  4. Then for large n, we have

$$\operatorname{Var}(\tilde{\theta}_n) \approx [J(\theta)]^{-1} \approx [J(\hat{\theta}_n)]^{-1}$$

and so  $\tilde{\theta}_n \sim N(\theta, [J(\theta)]^{-1}) \approx N(\theta, [J(\tilde{\theta})]^{-1})$ , and therefore we get

$$[J(\tilde{\theta}_n)]^{\frac{1}{2}}(\tilde{\theta}_n-\theta)\stackrel{D}{\to} Z\sim N(0,1).$$

Now, note that the definition of the information function is

$$I(\theta;X) = -\sum_{i=1}^{n} \frac{d^2}{d\theta^2} \ell(\theta;X_i).$$

According to Geyer 5, by the Law of Large Numbers,

$$\frac{1}{n}I(\tilde{\theta}) \stackrel{P}{\to} \frac{1}{n}J(\theta).$$

- <sup>3</sup> Comments from here on are my own after discussion with the lecturer. The explanation here will try to stay within the syllabus.
- <sup>4</sup> This is if I is continuous on both  $\theta$  and

<sup>5</sup> C. J. Geyer. Stat 5102 notes: Fisher information and confidence intervals using maximum likelihood. https: //www.stat.umn.edu/geyer/s06/ 5102/notes/fish.pdf, March 2007

Therefore, by Limit Theorems, we have

$$\frac{I(\tilde{ heta})}{J( heta)} \stackrel{P}{ o} 1$$

Then by the Continuous Mapping Theorem, we have

$$[I(\tilde{\theta})]^{\frac{1}{2}}(\tilde{\theta}-\theta) \stackrel{D}{\rightarrow} Z \sim N(0,1).$$

This result implies, by Limit Theorems, that for sufficiently large n, we have

$$Var(\tilde{\theta}_n) \approx [I(\hat{\theta})]^{-1}$$
.

In the next section, we shall study how these results can be used to construct approximate **confidence intervals** for  $\theta$ .

# Example 20.1.3 (Example 6.14 (Course Notes 6.4.2))

Suppose  $X_1, ..., X_n$  is a random sample from the Wei $(\theta, 2)$  distribution. Verify that the consistency and Normal distribution assumption hold for this distribution. Also, show that

$$[I(\tilde{\theta}_n; X)]^{\frac{1}{2}}(\tilde{\theta}_n - \theta) \stackrel{D}{\rightarrow} Z \sim N(0, 1).$$

Note that if  $X \sim Wei(\theta, 2)$ , then  $E(X^k) = \theta^k \Gamma\left(\frac{k}{2} + 1\right)$  and its pf is

$$f_{X_i}(x_i) = \frac{2}{\theta^2} x_i e^{-\left(\frac{x_i}{\theta}\right)^2}$$

#### Solution

Since the  $X_i$ 's form a random sample, we have that the likelihood function is

$$L(\theta) = \prod_{i=1}^{n} \frac{2}{\theta^{2}} x_{i} e^{-\left(\frac{x_{i}}{\theta}\right)^{2}} = \left(\frac{2}{\theta^{2}}\right)^{n} e^{-\frac{1}{\theta^{2}} \sum_{i=1}^{n} x_{i}^{2}} \prod_{i=1}^{n} x_{i}.$$

Then the log-likelihood is

$$\ell(\theta) = n \ln 2 - 2n \ln \theta - \frac{1}{\theta^2} \sum_{i=1}^n x_i^2 + \sum_{i=1}^n \ln x_i.$$

To get a candidate for the MLE, we have

$$\frac{d\ell}{d\theta} = -\frac{2n}{\theta} + \frac{2}{\theta^3} \sum_{i=1}^n x_i^2,$$

and so

$$\hat{\theta}^2 = \frac{1}{n} \sum_{i=1}^n x_i \text{ and } \tilde{\theta}^2 = \frac{1}{n} \sum_{i=1}^n X_i^2.$$

To verify that this candidate is indeed a maximum, note that we have

$$\frac{d^2\ell}{d\theta^2} = \frac{2n}{\theta^2} - \frac{6}{\theta^4} \sum_{i=1}^n x_i^2.$$
 (20.3)

Evaluating the above at  $\hat{\theta}^2 = \frac{1}{n} \sum_{i=1}^n x_i^2$ , we have

$$\frac{2n}{\frac{1}{n}\sum_{i=1}^{n}x_{i}^{2}} - \frac{6}{\left(\frac{1}{n}\right)^{2}\left(\sum_{i=1}^{n}x_{i}^{2}\right)^{2}}\sum_{i=1}^{n}x_{i}^{2} = -\frac{4n^{2}}{\sum_{i=1}^{n}x_{i}^{2}} < 0.$$

Thus  $\hat{\theta}$  is indeed the maximum likelihood estimate. Now to show consistency, note that we have

$$E(X_i^2) = \theta^2 \Gamma(1+1) = \theta^2$$

and so

$$E\left(\tilde{\theta}^2\right) = E\left(\frac{1}{n}\sum_{i=1}^n X_i^2\right) = \frac{1}{n}\sum_{i=1}^n E\left(X_i^2\right) = \theta^2.$$

From here, to simplify notations, we shall write<sup>6</sup>

<sup>6</sup> Note that  $\overline{X^2} \neq \overline{X}^2$ .

$$\overline{X^2} = \frac{1}{n} \sum_{i=1}^n X_i^2 \text{ and } \sum x_i^2 = \sum_{i=1}^n x_i^2.$$

By the Weak Law of Large Numbers, we have that  $\overline{X^2} \stackrel{P}{\to} \theta^2$ . Since  $\overline{X^2} \ge$ 0, the function  $g(x) = \sqrt{x}$  is continuous on its support, and so using • Proposition 58, we have

$$\tilde{\theta} = \sqrt{\overline{X^2}} \stackrel{P}{\to} \theta.$$

To show normality, note that the Fisher information is

$$J(\theta) = E\left[-\frac{d^2}{d\theta^2}\ell(\theta)\right] = \frac{6}{\theta^4} \sum_{i=1}^n E(X_i^2) - \frac{2n}{\theta^2} = \frac{4n}{\theta^2}$$
 (20.4)

where we shall note the use of Equation (20.3). We want to invoke CLT, and so we also need  $Var(X^2)$ .

$$\operatorname{Var}(\overline{X^2}) = \frac{1}{n^2} \sum_{i=1}^n \operatorname{Var}(X_i^2) = \frac{1}{n^2} \sum_{i=1}^n \left[ E\left(X_i^4\right) - \left[ E\left(X_i^2\right) \right]^2 \right] = \frac{\theta^4}{n}.$$

Then by the CLT, we have

$$\frac{\sqrt{n}(\overline{X^2} - \theta^2)}{\theta^2} \stackrel{D}{\to} Z \sim N(0, 1).$$

By ♠ Proposition 58,

$$\sqrt{n}(\overline{X^2} - \theta^2) \stackrel{D}{\rightarrow} \theta^2 Z \sim N\left(0, \theta^4\right).$$

Using the same  $g(x) = \sqrt{x}$  as above, by  $\blacktriangleright$  Corollary 60, we have

$$\sqrt{n}( ilde{ heta}- heta)\overset{D}{
ightarrow}rac{ heta}{2}Z\sim N\left(0,rac{ heta^2}{4}
ight)$$

Since Equation (20.4), using 6 Proposition 58 again, we have

$$[J(\theta)]^{\frac{1}{2}}(\tilde{\theta} - \theta) = \frac{2\sqrt{n}}{\theta}(\tilde{\theta} - \theta) \stackrel{D}{\rightarrow} \frac{2}{\theta} \cdot \frac{\theta}{2}Z = Z \sim N(0, 1)$$

Now from Equation (20.3), we have that the information function is

$$I(\theta; X) = \frac{6}{\theta^4} \sum_{i=1}^{n} X_i^2 - \frac{2n}{\theta^2}.$$

Then

$$I(\tilde{\theta};X) = \frac{6n^2}{\left(\sum_{i=1}^n X_i^2\right)^2} \sum_{i=1}^n X_i^2 - \frac{2n^2}{\sum_{i=1}^n X_i^2} = \frac{4n^2}{\sum_{i=1}^n X_i^2} = \frac{4n}{\tilde{\theta}^2}.$$

Since  $\tilde{\theta} \stackrel{P}{\to} \theta$ , again by Limit Theorems, we have  $\frac{\theta^2}{\tilde{\theta}^2} \stackrel{P}{\to} 1$ , and so

$$rac{\sqrt{I( ilde{ heta};X)}}{\sqrt{J( heta)}} = rac{ heta^2}{ ilde{ heta}^2} \stackrel{P}{
ightarrow} 1.$$

Then

$$[I(\tilde{\theta};X)]^{\frac{1}{2}}(\tilde{\theta}-\theta) = \frac{\sqrt{I(\tilde{\theta};X)}}{\sqrt{J(\theta)}}[J(\theta)]^{\frac{1}{2}}(\tilde{\theta}-\theta) \xrightarrow{D} 1 \cdot Z = Z \sim N(0,1).$$

# 21 Lecture 21 Jul 12th 2018

# **21.1** *Estimation (Continued 4)*

# 21.1.1 Asymptotic Properties of ML Estimators (Continued)

# Example 21.1.1 (Example 6.15 (Course Notes 6.4.3))

Suppose  $X_1, ..., X_n$  is a random sample from the  $Unif(0, \theta)$  distribution. We showed in Example 18.1.5 that  $\tilde{\theta}_n = X_{(n)}$  is the ML estimator of  $\theta$ . Since the support of  $X_i$  depends on  $\theta$ ,  $\blacksquare$  Theorem 62 does not hold. Show, however, that  $X_{(n)}$  is still a consistent estimator of  $\theta$ .

# Solution

Since  $\theta$  is a fixed (but unknown) value, we have the option of use

• Proposition 55. We have

$$F_{X_{(n)}}(t) = P(X_{(n)} \le t) = P(X_1 \le t, X_2 \le t, ..., X_n \le t)$$

$$= \prod_{i=1}^{n} P(X_i \le t) \quad \because \text{ independence}$$

$$= \prod_{i=1}^{n} \frac{t}{\theta} = \left(\frac{t}{\theta}\right)^n$$

So the full expression for the CDF of  $X_{(n)}$  is

$$F_{\mathrm{X}_{(n)}}(t) = egin{cases} 0 & t \leq 0 \ \left(rac{t}{ heta}
ight)^n & 0 < t < heta \ 1 & t \geq heta \end{cases}.$$

Observe that

$$\lim_{n o\infty}F_{X_{(n)}}(t)=egin{cases} 0 & t< heta\ 1 & t\geq heta \end{cases}$$

and so as  $n \to \infty$ , we have that  $X_{(n)}$  has a denegerate distribution, and so the proof for consistency follows from  $\bullet$  Proposition 55.

# 21.1.2 Interval Estimators

# Definition 70 (Interval Estimators)

Suppose X is an rv whose distribution depends on  $\theta$ . Suppose that A(x) and B(x) are functions such that  $A(x) \leq B(x)$  for all  $x \in \text{supp}(X)$  and  $\theta \in \Omega$ . Let x be the observed data. Then (A(x), B(x)) is an **interval** estimate for  $\theta$ . The interval (A(X), B(X)) is an **interval** estimator for  $\theta$ 

#### 66 Note

Two interval estimators that we already know<sup>1</sup>:

- 1. Confidence intervals
- 2. Likelihood intervals

WE NOW CONSIDER a general approach for constructing confidence intervals based on pivotal quantities.

# Definition 71 (Pivotal Quantity)

Suppose X is an rv whose distribution depends on  $\theta$ . The rv  $Q(X;\theta)$  is called a **pivotal quantity** if the distribution of Q does not depend on  $\theta$ .  $Q(X;\theta)$  is called an **asymptotic pivotal quantity** if the limiting distribution of Q as  $n \to \infty$  does not depend on  $\theta$ .

<sup>1</sup> Well, one from STAT231, that is.

# Example 21.1.2 (Example 6.16)

Suppose  $X_1, ..., X_n$  is a random sample from the  $N(\mu, \sigma^2)$  distribution. The following rvs are pivotal quantities.:

$$\frac{\overline{X} - \mu}{\sigma / \sqrt{n}}, \frac{\overline{X} - \mu}{S / \sqrt{n}}, \frac{\sum\limits_{i=1}^{n} (X_i - \mu)^2}{\sigma^2} \text{ and } \frac{(n-1)S^2}{\sigma^2} = \frac{\sum\limits_{i=1}^{n} (X_i - \overline{X})^2}{\sigma^2}$$

# Example 21.1.3 (Example 6.17 (Course Notes 6.5.4))

Suppose  $X_1,...,X_n$  is a random sample from the  $Poi(\theta)$  distribution. Show that

$$\sqrt{n} \frac{(\overline{X}_n - \theta)}{\sqrt{\overline{X}_n}}$$

is an asymptotic pivotal quantity.

# Solution

By P Theorem 57, we have that

$$\sqrt{n} \frac{\overline{X}_n - \theta}{\sqrt{\theta}} \stackrel{D}{\to} Z \sim N(0, 1).$$

*By* ♠ *Proposition* 58, we have

$$\sqrt{n}(\overline{X}_n - \theta) \stackrel{D}{\to} \sqrt{\theta}Z \sim N(0, \theta).$$

Now by the Weak Law of Large Numbers, we have that

$$\overline{X}_n \stackrel{P}{\to} \theta$$
.

Since  $\overline{X}_n$  takes only non-negative values, and  $\theta > 0$ , the function g(x) = $\sqrt{x}$  is well-defined on these "points", and so by  $\bullet$  Proposition 58, we have

$$\sqrt{\overline{X}_n} \stackrel{P}{\to} \sqrt{\theta}.$$

Then by Slutsky's Theorem, we have

$$\sqrt{n} \frac{\overline{X}_n - \theta}{\sqrt{\overline{X}_n}} \stackrel{D}{\to} Z \sim N(0, 1).$$

We observe that the limiting distribution of  $\sqrt{n}(\overline{X}_n - \theta) / \sqrt{\overline{X}_n}$  depends not on the unknown parameter  $\theta$ . Therefore it is a asymptotic pivotal quantity.

#### Exercise 21.1.1

Show that the given rvs are pivotal quanti-

This may be a little late but at this point, it should be clear how Limit Theorems work when it comes to playing with the Normal distribution as the asymptotic distribution. If this is not clear to you, read back from

• Proposition 58.

Suppose A(X) and B(X) are statistics, and are called the **lower** and **upper confidence limits**, respectively. If  $P[A(X) < \theta < B(X)] = p$ , where 0 , then <math>(a(x), b(x)) is called a 100p% **confidence interval** (CI) for  $\theta$ .

Pivotal quantities can be used for constructing CIs in the following way: since the distribution of  $Q(X;\theta)$  is known, we can write down a probability statement of the form

$$P(q_1 \leq Q(X;\theta) \leq q_2) = p.$$

If Q is a monotone function of  $\theta$ , then this statement can be rewritten as

$$P[A(X) \le \theta \le B(X)] = p$$

and the "realized" interval [a(x), b(x)] is a 100p% CI.

#### Example 21.1.4 (Example 6.18 (Course Notes 6.5.6))

Suppose  $X_1, ..., X_n$  is a random sample from the  $N(\mu, \sigma^2)$  distribution. Use the pivotal quantities in Example 21.1.2 to find:

- 1. a 100p% CI for  $\mu$  if  $\sigma^2$  is unknown;
- 2. a 100p% CI for  $\sigma^2$  is  $\mu$  is known.

# Solution

1. Since  $\sigma^2$  is unknown, we cannot use

$$\frac{\overline{X}_n - \mu}{\sigma / \sqrt{n}} \sim N(0, 1)$$

as our pivotal quantity. We shall in fact, use

$$\frac{\overline{X}_n - \mu}{S/\sqrt{n}} \sim t(n-1)$$

since it contains the unknown parameter whose confidence interval is that of what we want to construct, where we know that the rv follows a t-Distribution with degrees of freedom n-1 from  $\blacksquare$  Theorem 47. Then since the t-Distribution is symmetric<sup>2</sup>, we have

$$P\left[-a \le \frac{\overline{X}_n - \mu}{S/\sqrt{n}} \le a\right] - p$$

$$P\left[\overline{X}_n - (S/\sqrt{n})a \le \mu \le \overline{X}_n + (S/\sqrt{n})a\right] = p.$$

<sup>2</sup> See Wikipedia for a graph of the *t*-Distribution.

And so the 100p% CI for  $\mu$  when  $\sigma^2$  is unknown is

$$[\overline{X}_n - (S/\sqrt{n})a, \overline{X}_n + (S/\sqrt{n})a]$$

2. Since  $\mu$  is known, we shall use the rv

$$\frac{\sum\limits_{i=1}^{n}(X_i-\mu)^2}{\sigma^2}\sim \chi^2(n-1),$$

where the following of the chi-squared Distribution is from Cochran's Theorem. To construct the CI for  $\sigma^2$ ,

$$P\left(a \le \frac{\sum\limits_{i=1}^{n} (X_i - \mu)^2}{\sigma^2} \le b\right) = p$$

$$P\left(\frac{1}{b} \sum\limits_{i=1}^{n} (X_i - \mu)^2 \le \sigma^2 \le \frac{1}{a} \sum\limits_{i=1}^{n} (X_i - \mu)^2\right) = p.$$

Thus the CI for  $\sigma^2$  is

$$\left[\frac{1}{b}\sum_{i=1}^{n}(X_{i}-\mu)^{2},\,\frac{1}{a}\sum_{i=1}^{n}(X_{i}-\mu)^{2}\right].$$

# • Proposition 63 (MLE of a Location/Scale Parameter as a Pivotal Quantity)

Let  $X = (X_1, ..., X_n)$  be a random sample from  $f(x; \theta)$  and let  $\tilde{\theta} = \tilde{\theta}(X)$ be the ML estimator of the scalar parameter  $\theta$  based on X.

- 1. If  $\theta$  is a location parameter, then  $Q = \tilde{\theta} \theta$  is a pivotal quantity.
- 2. If  $\theta$  is a scale parameter, then  $Q = \tilde{\theta}/\theta$  is a pivotal quantity.

There is a proof out there and I wish to hunt for it but I do not have the time.

# Example 21.1.5 (Example 6.19)

Suppse  $X_1,...,X_n \sim \text{Exp}(\theta)$  is a random sample. Find a 100p% equal tails CI for  $\theta$ . For the data n=15 and  $\sum_{i=1}^{15} x_i=36$ , find a 95% equal tail CI for

# Solution

Referring to Example 3.3.3, we know that  $Exp(\theta)$  belongs to a scale parameter family of distributions. We shall proceed with finding the MLE of  $\theta$  and then use  $\bullet$  Proposition 63 to get a CI for  $\theta$ . Since the  $X_i$ 's are IID,

$$L(\theta) = \prod_{i=1}^{n} \frac{1}{\theta} e^{-\frac{x_i}{\theta}} = \theta^{-n} \exp\left\{-\frac{1}{\theta} \sum_{i=1}^{n} x_i\right\}$$
$$\ell(\theta) = -n \ln \theta - \frac{1}{\theta} \sum_{i=1}^{n} x_i$$
$$\frac{d\ell}{d\theta} = -\frac{n}{\theta} + \frac{1}{\theta} \sum_{i=1}^{n} x_i$$
$$\implies \hat{\theta} = \frac{1}{n} \sum_{i=1}^{n} x_i = \bar{x}.$$

Thus the MLE of  $\theta$  is  $\tilde{\theta} = \overline{X}$ . To get the 100p% CI for  $\theta$ , we have by  $\Phi$  Proposition 63,  $Y = \overline{X}/\theta$  is a pivotal quantity. Then

$$P\left[a \le \frac{\overline{X}}{\theta} \le b\right] = p$$

$$P\left[\frac{1}{b}\overline{X} \le \theta \le \frac{1}{a}\overline{X}\right] = p.$$

*And so the* 100p% *CI for*  $\theta$  *is* 

$$\left[\frac{1}{b}\overline{X},\,\frac{1}{a}\overline{X}\right].$$

To solve for a and b, we would need to know what the distribution of  $Y = \overline{X}/\theta$  is:

$$M_Y(t) = E\left(\exp\left\{\frac{t}{\theta n}\sum_{i=1}^n X_i\right\}\right) = E\left(e^{\frac{tX_1}{n\theta}}\right)^n :: IID$$

$$= M_{X_1}\left(\frac{t}{n\theta}\right)^n = \left(1 - \frac{t}{n}\right)^{-n},$$

which we observe that then  $Y \sim Gam\left(n, \frac{1}{n}\right)$ .

For n=15,  $\sum_{i=1}^{15} x_i=36$  (and so  $\bar{x}=36/15=2.4$ ), and to find the 95% equal tail CI for  $\theta$ , we have that

$$P\left[0.559 \le \frac{2.4}{\theta} \le 1.566\right] = 0.95$$
  
 $P\left[2.4/1.566 \le \theta \le 2.4/0.559\right] = 0.95.$ 

where the values of a and b are verified on R using the following:

Listing 21.1: R interactive code to get a and b

Thus

$$\left[\frac{2.4}{1.566}, \frac{2.4}{0.559}\right]$$
[1.533, 4.293]

# Example 21.1.6 (Example 6.20 (Course Notes 6.5.11))

Suppose  $X_1,...,X_n \sim \text{Unif}(0,\theta)$  is a random sample. Find an appropriate pivotal quantity. Determine a such that

$$[\hat{\theta}, a\hat{\theta}]$$

is a 100p% CI for  $\theta$ .

# Solution

Since the  $X_i$ 's are IID, we have that

$$L(\theta) = \prod_{i=1}^{n} \frac{1}{\theta} \mathbb{1}_{\{0 \le x_i \le \theta\}} = \theta^{-n} \mathbb{1}_{\{0 \le x_{(1)} \le x_{(n)} \le \theta\}}$$
$$\ell(\theta) = -n \ln \theta + \ln \mathbb{1}_{\{0 \le x_{(1)} \le x_{(n)} \le \theta\}}.$$

*Using a similar argument as in Example 19.1.2, notice that*  $-n \ln \theta$  *is a* decreasing function, and so to get the maximum likelihood estimate, we want  $\theta$  to be as close to 0 as possible, but the closest  $\theta$  can get to is  $x_{(n)}$ . Note that if the indicator function turns out to be 1, then  $\theta$  must also be 1, but this does not maximize the likelihood function.

So we have  $\hat{\theta} = x_{(n)}$  and  $\tilde{\theta} = X_{(n)}$ . Note that  $\text{Unif}(0,\theta)$  is from a scale parameter family of distributions since

$$f_{X_i}(x_i) = \frac{1}{\theta} \mathbb{1}_{\{0 \le x_i \le \theta\}} = \frac{1}{\theta} f_{X_i}\left(\frac{x_i}{\theta}\right).$$

Then by  $\bullet$  Proposition 63,  $Y = \tilde{\theta}/\theta = X_{(n)}/\theta$  is a pivotal quantity. Note, then, that we want to solve for

$$P\left[\hat{\theta} \leq \theta \leq a\hat{\theta}\right] = p$$
 $P\left[\frac{1}{a} \leq \frac{\hat{\theta}}{\theta} \leq 1\right] = p.$ 

It suffices, from here, to find the CDF of Y:

$$F_Y(y) = P(Y \le y) = P\left(\frac{X_{(n)}}{\theta} \le y\right) = P(X_{(n)} \le y\theta)$$

$$= P(X_1 \le y\theta)^n \quad \therefore IID$$

$$= \left(\int_0^{y\theta} \frac{1}{\theta} dx\right)^n = y^n$$

Thus

$$F(1) - F\left(\frac{1}{a}\right) = p$$
$$1 - \left(\frac{1}{a}\right)^n = p$$
$$a = \frac{1}{\sqrt[n]{p-1}}$$

# 22 Lecture 22 Jul 17th 2018

# 22.1 Estimate (Continued 5)

# 22.1.1 Interval Estimators (Continued)

# Example 22.1.1 (Example 6.21)

Suppose  $X_1,...,X_n \sim \operatorname{Exp}(1,\theta)$  is a random sample. Find an appropriate pivotal quantity. Determine a such that

$$[\hat{\theta} - a, \hat{\theta}]$$

is a 100p% CI for  $\theta$ .

# Solution

Recall from Example 19.1.2 that the ML estimate is  $\hat{\theta} = x_{(1)}$  and and the estimator is  $\tilde{\theta} = X_{(1)}$ . Note that  $\text{Exp}(1,\theta)$  is from a location parameter family of distributions, since

$$f(x;1,0) = e^{-x} = f(x-\theta;1,\theta).$$

We can then use  $\bullet$  Proposition 63, and use  $Y = \tilde{\theta} - \theta$  as our pivotal quantity. Note that

$$P(\hat{\theta} - a \le \theta \le \hat{\theta}) = p$$

$$P(-\hat{\theta} \le -\theta \le a - \hat{\theta}) = p$$

$$P(0 \le \hat{\theta} - \theta \le a) = p$$

$$P(0 \le Y \le a) = p$$

It remains, again, to find the CDF.

$$F_Y(y) = P(X_{(1)} - \theta \le y) = 1 - P(X_{(1)} > y + \theta)$$

$$= 1 - P(X_1 > y + \theta)^n \quad \because IID$$

$$= 1 - \left[ \int_{y+\theta}^{\infty} e^{-(x+\theta)} dx \right]^n$$

$$= 1 - \left[ -e^{-(x+\theta)} \Big|_{y+\theta}^{\infty} \right]^n$$

$$= 1 - e^{-ny}$$

Thus

$$F(a) - F(0) = p$$

$$1 - e^{-na} - 1 + 1 = p$$

$$a = -\frac{1}{n} \ln(1 - p)$$

In cases where we cannot construct an exact pivotal quantity, we can use the limiting distribution of the ML estimator  $\tilde{\theta} = \tilde{\theta}(X_1, ..., X_n)$  and construct approximate CIs.

# • Proposition 64 (Asymptotic Confidence Intervals)

Since1

$$[J(\tilde{\theta})]^{\frac{1}{2}}(\tilde{\theta}-\theta) \stackrel{D}{\rightarrow} Z \sim N(0,1),$$

then  $[J(\tilde{\theta})]^{\frac{1}{2}}(\tilde{\theta}-\theta)$  is an asymptotic pivotal quantity. An asymptotic 100p% CI<sup>2</sup> based on this asymptotic pivotal quantity is given by

$$[\hat{\theta} - a[J(\hat{\theta})]^{-\frac{1}{2}}, \ \hat{\theta} + a[J(\hat{\theta})]^{-\frac{1}{2}}]$$

where P(-a < Z < a) = p and  $Z \sim N(0,1)$ .

# ¹ This is based on ■ Theorem 62, which we could not prove, and its consequences are more than just fishy. I am also unsure whether this is presented as a proof or definition in the lecture, since there are no indications of that whatsoever both during the lectures and in the course notes.

<sup>2</sup> Also called **approximate** 100*p*% **CI**.

# 66 Note

Similarly, since

$$[I(\tilde{\theta})]^{\frac{1}{2}}(\tilde{\theta}-\theta) \stackrel{D}{\to} Z \sim N(0,1),$$

then  $[I(\tilde{\theta})]^{\frac{1}{2}}(\tilde{\theta}-\theta)$  is an asymptotic pivotal quantity. An asymptotic 100p% CI based on this asymptotic pitoval quantity is given by

$$[\hat{\theta} - a[I(\hat{\theta})]^{-\frac{1}{2}}, \ \hat{\theta} + a[I(\hat{\theta})]^{-\frac{1}{2}}]$$

where  $I(\hat{\theta})$  is the observed information.

# Example 22.1.2 (Example 6.22)

Suppose  $X \sim Bin(n, \theta)$ . Show how you would construct an approximate 100p% CI for  $\theta$ .

# Solution

Since this is just a single rv, we have

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

$$\ell(\theta) = \ln \binom{n}{x} + x \ln \theta + (n - x) \ln(1 - \theta)$$

$$\frac{d\ell}{d\theta} = \frac{x}{\theta} - \frac{n - x}{1 - \theta}$$

$$\frac{d^2 \ell}{d\theta} = -\frac{x}{\theta^2} - \frac{n - x}{(1 - \theta)^2} < 0$$

where we note that the ML estimate is  $\hat{\theta} = \frac{x}{n}$  and the estimator  $\tilde{\theta} = \frac{X}{n}$ , and we know that this will be the maximum since the second order derivative is always negative<sup>3</sup>. The information function is

$$I(\theta; X) = \frac{X}{\theta} + \frac{n - X}{(1 - \theta)^2}$$

and so the Fisher information is

$$J(\theta) = \frac{E(X)}{\theta^2} + \frac{n - E(X)}{(1 - \theta)^2} = \frac{n}{\theta} + \frac{n}{1 - \theta} = \frac{n}{\theta(1 - \theta)}.$$

Thus

$$J(\hat{\theta}) = \frac{n^3}{x(n-x)}.$$

*Then by*  $\bullet$  *Proposition 64, we have that the asymptotic* 100p% CI for  $\theta$  *is* 

$$\left[\frac{x}{n}-zn\sqrt{\frac{x(n-x)}{n}}, \frac{x}{n}+zn\sqrt{\frac{x(n-x)}{n}}\right],$$

where  $P(-z \le Z \le z) = p$  where  $Z \sim N(0,1)$ .

# Example 22.1.3 (Example 6.23 (Course Notes 6.5.15))

Suppose  $X_1,...,X_n$  is a random sample from the  $Poi(\theta)$  distribution. Show how you would construct an asymptotic 100p% CI for  $\theta$ .

<sup>3</sup> In fact, by rearranging  $\frac{d\ell}{d\theta}$ , we can show that it is a linear function.

# Solution

Since the  $X_i$ 's are IID, we have

$$L(\theta) = \prod_{i=1}^{n} \frac{e^{-\theta} \theta^{x_i}}{x_i!} = e^{-n\theta} \theta^{\sum_{i=1}^{n} x_i} \prod_{i=1}^{n} \frac{1}{x_i!}$$

$$\ell(\theta) = -n\theta + \left(\sum_{i=1}^{n} x_i\right) \ln \theta - \sum_{i=1}^{n} \ln x_i!$$

$$\frac{d\ell}{d\theta} = -n + \frac{1}{\theta} \sum_{i=1}^{n} x_i$$

$$\frac{d^2 \ell}{d\theta^2} = -\frac{1}{\theta^2} \sum_{i=1}^{n} x_i.$$

The candidate for MLE from equating the score function to 0 is  $\hat{\theta} = \frac{1}{n} \sum_{i=1}^{n} x_i = \bar{x}$ . We verify that it is indeed a maximum point since

$$\left. \frac{d^2 \ell}{d\theta^2} \right|_{\theta = \hat{\theta}} = -\frac{1}{\bar{x}} \sum_{i=1}^n = -\frac{1}{n} < 0.$$

The information function is

$$I(\theta; X) = \frac{1}{\theta^2} \sum_{i=1}^n X_i$$

and so the Fisher information is

$$J(\theta) = E[I(\theta; X)] = \frac{1}{\theta^2} n\theta = \frac{n}{\theta}.$$

Thus

$$J(\hat{\theta}) = \frac{n}{\bar{x}}$$

Assuming that our random sample satisfies the regularity conditions, we have that the asymptotic 100p% CI for  $\theta$ , by  $\bullet$  Proposition 64, is

$$\left[\bar{x}-z\sqrt{\frac{\bar{x}}{n}},\ \bar{x}+z\sqrt{\frac{\bar{x}}{n}}\right]$$

# Definition 73 (Likelihood Interval)

Given a random sample  $X_1, ..., X_n$  from  $f(x; \theta)$ , the 100p% likelihood interval for  $\theta$  is defined by

$$\{\theta: R(\theta) > p\}.$$

# 66 Note

Recall the 3rd property from 💻 Theorem 62 that

 $-2\log R(\theta;X) \stackrel{D}{\rightarrow} W \sim \chi^2(1).$ 

Since the  $R(\theta)$  is a unimodal graph, we have

$$P[R(\theta; X) \ge p] = P[-2 \ln R(\theta; X) \le -2 \ln p]$$

$$\approx P(W \le -2 \ln p)$$

$$= P(Z^2 \le -2 \ln p)$$

$$= P\left(-\sqrt{-2 \ln p} \le Z \le \sqrt{-2 \ln p}\right)$$

$$= 2P\left(Z \le \sqrt{-2 \ln p}\right) - 1.$$

Consequently, if p = 0.15, then

$$P[R(\theta; X) \ge 0.15] \approx 0.95$$

i.e. a 15% likelihood interval is, approximately, a 95% CI.

The amount of things built on this shady theorem is insane...

# Example 22.1.4 (Example 6.24)

In Example 19.1.1, based on  $X_1,...,X_{100} \sim Poi(\theta)$  is a random sample, where  $\sum_{i=1}^{100} x_i = 980$ , we derived the 10% and 50% likelihood intervals for θ. Compare the approximate 95% confidence interval based on the asymptotic distribution of the MLE to the 15% likelihood interval for  $\theta$ .

# Solution

Recall that

$$R(\theta) = e^{-n(\theta - \bar{x})} \left(\frac{\theta}{\bar{x}}\right)^{n\bar{x}} = e^{-100(\theta - 9.8)} \left(\frac{\theta}{980}\right)^{980}$$

and that the MLE is  $\hat{\theta} = \bar{x}$ . Using R 4, one can find that the 15% likelihood interval is

<sup>4</sup> This is actually a non-trivial process...

$$9.2 < \theta < 10.4$$
.

Again, assuming that the regularity conditions are satisfied, by ♠ Proposition 64,

and the previous example, the 95% CI for  $\theta$  is

$$\left[9.8 - 1.96 \frac{\sqrt{980}}{100}, 9.8 + 1.96 \frac{\sqrt{980}}{100}\right]$$
$$\left[9.186, 10.414\right]$$

# <sup>23</sup> Lecture 23 Jul 19th 2018

# 23.1 Hypothesis Testing

# 23.1.1 Introduction

# Definition 74 (Hypothesis)

In inferential statistics, given a model  $f(x;\theta)$ , where  $\theta = (\theta_1 \dots \theta_n)^T \in \Omega$ , and  $\Omega$  is a parameter space, a **hypothesis** is a general statement about the parameters for the model, denoted by

$$H: \theta \in \Omega_0 \subseteq \Omega$$
.

# Definition 75 (Null Hypothesis and Alternative Hypothesis)

A *null hypothesis* is a hypothesis of which its default position is that there is no relationship between the two phenomena measured by the model  $f(x;\theta)$ . We usually denote the null hypothesis by

$$H_0: \theta \in \Omega_0$$
.

The **alternative hypothesis** is the other case of the null hypothesis; it is a hypothesis that states that there is a relationship between the two phenomena measured by the model  $f(x;\theta)$ . We usually denote the alternative hypothesis as

$$H_a: \theta \notin \Omega_0$$
.

# Definition 76 (Test of Hypothesis)

A **test of hypothesis** is a proedure for evaluating the strength of the <u>evidence provided</u> by the data against the null hypothesis.

To measure the evidence against  $H_0$  based on the observed data, we use a test statistic or a discrepancy measure.

# Definition 77 (Test Statistic)

A **test statistic** is a statistic that is used for measuring the evidence against  $H_0$  using the observed data.

#### 66 Note

- 1. A small observed value of the test statistic shows close agreement between the observed data and the null hypothesis; while
- 2. A large observed value of the test statistic shows poor agreement.

# $\blacksquare$ Definition 78 (Significance Level and p-Values)

Assuming that  $H_0$  is true, we compute the probability of observing a value of the test statistic being greater than or equal to the actual observed value. This probability is called the **significance level** or p-value of the data in relation with the null hypothesis.

#### 66 Note

The p-value is the probability of observing a poor agreement between the null hypothesis and the data.

- A small p-value indicates that the probability that our observed test statistic occurs under the null hypothesis is low, i.e. an evidence against the null hypothesis;
- A large p-value indicates that the probability that our observed test statistic occurs under the null hypothesis is high, i.e. we have little evidence against the null hypothesis.

This also means that the smaller the p-value, the stronger the evidence against the null hypothesis.

There are two types of errors that we may run into while performing hypothesis testing:

# Definition 79 (Type I and Type II Errors)

A Type I Error is the case where we reject a true null hypothesis. A Type II Error is the case where we fail to reject a false null hypothesis.

$$H_0: \theta \in \Omega_0 \qquad H_a: \theta \notin \Omega_0$$

$$Reject \ H_0 \qquad \textbf{Type I Error} \qquad \checkmark \ (True \ Negatives)$$

$$Fail \ to \ Reject \ H_0 \qquad \checkmark \ (True \ Positives) \qquad \textbf{Type II Error}$$

#### 66 Note

The common practice is to reduce the occurrence of a Type I Error. It is analogous to the Presumption of Innocence, where the "null hypothesis" is "correct until proven wrong". In other words, the goal of this practice is to reduce the occurrence of false positives.

#### Example 23.1.1 (Example 7.11 (Course Notes 7.1.1))

Suppose  $X_1, ..., X_n \sim N(\mu, \sigma^2)$  is a random sample, where  $\sigma^2$  is known. Suppose n=25,  $\sigma=1$ ,  $\bar{x}=0.5$ , and  $\mu_0=0$ . Perform a test of hypothesis on the null hypothesis  $H_0: \mu = \mu_0$  and state the conclusion.

# Solution

By the CLT, we have

$$Z = \frac{\overline{X} - \mu}{\sigma / \sqrt{n}} \sim N(0, 1)$$

which contains all the parameters of our interest and is a statistic<sup>1</sup>. We shall use this as a test statistic. Under  $H_0$ , we have  $\mu = \mu_0$ , i.e.

<sup>1</sup> Note that it is also a pivotal quantity.

$$\frac{\overline{X} - \mu_0}{\sigma / \sqrt{n}} \sim N(0, 1).$$

Then, the observed test statistic is

$$z_0 = \frac{\bar{x} - \mu_0}{\sigma / \sqrt{n}} = \frac{0.5 - 0}{1 / \sqrt{25}} = 2.5,$$

and so the p-value is

$$p = 2P(Z \ge 2.5) \approx 0.01241923$$

by the Z-table. Therefore, we observe that the p-value is small, and so we have evidence against the null hypothesis  $H_0: \mu_0 = 0$ .

# 23.1.2 Likelihood Ratio Tests for Simple Hypothesis

# Definition 80 (Simple Hypothesis)

Suppose  $X_1, ..., X_n$  is a random sample from  $f(x; \theta)$ , where  $\theta = (\theta_1 ... \theta_k)^T$ . Suppose we wish to test  $H : \theta = \theta_0$ , where  $\theta_0$  is a completely known  $k \times 1$  vector. This hypothesis is called a **simple hypothesis**, which specifies the values of all unknown parameters in the mode.

#### Definition 81 (Likelihood Ratio Statistic)

The likelihood ratio (LR) statistic is defined as

$$\Lambda( heta_0) = -2\ln R( heta_0;X) = -2\ln \left[rac{L( heta_0;X)}{L( ilde{ heta};X)}
ight] = 2\left[\ell( ilde{ heta};X) - \ell( heta_0;X)
ight] \; ,$$

where  $X = (X_1, ..., X_n)$  are the data and  $\tilde{\theta} = \tilde{\theta}(X_1, ..., X_n)$  is the ML estimator of  $\theta$ . The observed likelihood ratio statistic is thus

$$\lambda(\theta_0) = -2\ln R(\theta_0; x)$$

where  $x = (x_1, ..., x_n)$ .

Assuming  $H_0$  is true, under certain regularity conditions<sup>2</sup>, we have

<sup>2</sup> This probably involves **Theorem** 62

$$\Lambda(\theta_0) = -2 \ln R(\theta_0, X) \stackrel{D}{\to} W \sim \chi^2(k).$$

where k is the number of parameters under the general hypothesis minus the number of parameters under  $H_0$ . Thus an approximate

p-value is given by

$$p \approx P(W \ge -2 \ln R(\theta_0; x)])$$
,

where  $x = (x_1, ..., x_n)$  is the observed data.

# Example 23.1.2 (Example 7.2 (Course Notes 7.2.1))

Continuing with Example 23.1.1, find the LR statistic and compare it to the test statistic used in the example.

# Solution

Let  $H_0: \mu = \mu_0$  and  $H_a: \mu \neq \mu_0$ . We need the MLE of  $\mu$ . Since the  $X_i$ 's forms a random sample, we have

$$L(\mu) = \prod_{i=1}^{n} \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x_i - \mu)^2}{2\sigma^2}} = (2\pi)^{-\frac{n}{2}} \sigma^{-n} \exp\left\{-\frac{1}{2\sigma^2} \sum_{i=1}^{n} (x_i - \mu)^2\right\}$$
$$\ell(\mu) = -\frac{n}{2} \ln(2\pi) - n \ln \sigma - \frac{1}{2\sigma^2} \sum_{i=1}^{n} (x_i - \mu)^2$$
$$\frac{d\ell}{d\mu} = \frac{1}{\sigma^2} \sum_{i=1}^{n} (x_i - \mu) = \frac{1}{\sigma^2} (\sum_{i=1}^{n} x_i - n\mu)$$
$$\frac{d^2\ell}{d\mu^2} = -\frac{n}{\sigma^2} < 0$$

and so the ML estimate for  $\mu$  is

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^{n} x_i = \bar{x}.$$

Thus the likelihood ratio statistic is

$$\begin{split} \Lambda(\mu_0) &= -2 \ln R(\mu_0, X) = 2 [\ell(\tilde{\mu}; X) - \ell(\mu_0; X)] \\ &= \frac{1}{\sigma^2} \left( \sum_{i=1}^n (X_i - \mu_0)^2 - \sum_{i=1}^n (X_i - \tilde{\mu})^2 \right) \\ &= \frac{n}{\sigma^2} \left[ 2\bar{x}(\tilde{\mu} - \mu_0) + \mu_0^2 - \tilde{\mu}^2 \right] \end{split}$$

By our argument before this example, we have that

$$k = number of unknown parameters$$
 $- number of unknown parameters under H_0$ 
 $= 1 - 0 = 1.$ 

Thus

$$\Lambda(\theta_0) \stackrel{D}{\to} W \sim \chi^2(1).$$

The observed likelihood ratio statistic is

$$\begin{split} \lambda(\mu_0) &= \frac{n}{\sigma^2} [2\bar{x}(\bar{x} - \mu_0) + \mu_0^2 - \bar{x}^2] = \frac{n(\bar{x}^2 - 2\mu_0\bar{x} + \mu_0^2)}{\sigma^2} \\ &= \frac{n(\bar{x} - \mu_0)^2}{\sigma^2}. \end{split}$$

The p-value is therefore

$$p pprox P(W \ge \lambda(\mu_0)) = P\left(W \ge \frac{n(\bar{x} - \mu_0)^2}{\sigma^2}\right).$$

Using the values given in Example 23.1.1, we have

$$p \approx P(W \ge 6.25) \approx 0.01241923$$
,

where our value is obtained by the R code:

We can now conclude that the two tests are similar.

# Example 23.1.3 (Example 7.3 (Course Notes 7.2.4))

Suppose  $X_1, ..., X_n$  is a random sample from the  $Poi(\theta)$  distribution. Find the LR statistic for testing  $H : \theta = \theta_0$ . Another test statistic which could be used is

$$T = \frac{\left| \overline{X} - \theta_0 \right|}{\sqrt{\theta_0 / n}}.$$

What is the approximate distribution of

$$\frac{X - \theta_0}{\sqrt{\theta_0/n}}$$

for large n if  $H_0$  is true, and how could you use this to find an approximate p-value?

# Solution

Let  $H_0: \theta \in \Omega_0$  and  $H_a: \theta \notin \Omega_0$ . By Example 22.1.3, we have that  $\hat{\theta} = \bar{x}$ , and the log-likelihood is

$$\ell(\theta) = -n\theta + \left(\sum_{i=1}^{n} x_i\right) \ln \theta - \sum_{i=1}^{n} \ln x_i!.$$

Thus the likelihood ratio statistic is

$$\Lambda( heta_0) = 2[n\overline{X}\ln ilde{ heta} - n\overline{X}\ln heta_0] = 2n\overline{X}\ln\left(rac{ ilde{ heta}}{ heta_0}
ight)$$

and so the observed likelihood ratio statistic is

$$\lambda(\theta_0) = 2n\bar{x} \ln\left(\frac{\bar{x}}{\theta_0}\right).$$

Thus the p-value under the assumption of  $H_0$  is

$$p \approx P(W \ge \lambda(\theta_0)).$$

On the other hand, by the CLT, we have that

$$\frac{\overline{X} - \theta_0}{\sqrt{\theta_0/n}} \stackrel{D}{\to} Z \sim N(0,1),$$

which then the p-value under this method is

$$p pprox 2P\left(Z \ge \left| rac{ar{x} - heta_0}{\sqrt{ heta_0/n}} \right| 
ight).$$

However, it is not necessary that

$$P(W \ge \lambda(\theta_0)) = 2P\left(Z \ge \left| \frac{\bar{x} - \theta_0}{\sqrt{\theta_0/n}} \right| \right)$$

# Example 23.1.4 (Example 7.4 (Course Notes 7.2.6))

The following table gives the observed frequencies of the six faces in 100 rolls of a die:

Are these observations consistent with the hypothesis that the die is fair?

#### Solution

Let  $\theta_i$  be the probability that we get face i, and let

$$H_0: \theta_1=\theta_2=\ldots=\theta_6=\frac{1}{6}$$

 $H_a$ : at least one of the  $\theta_i \neq \theta_j$  for  $i \neq j$ ,  $0 \leq i, j \leq 6$ 

Note that the above experiment is of a multinomial model. Let  $\theta = (\theta_1, ..., \theta_6)$ .

Thus we have

$$L(\theta) = \frac{100!}{\prod_{i=1}^{6} x_i!} \theta_1^{x_1} \theta_2^{x_2} \dots \theta_6^{x_6}.$$

With the constraint that  $\sum_{i=1}^{6} \theta_i = 1$ , we can use Lagrange's Multiplier<sup>3</sup> to

<sup>3</sup> A step of that can be found here:

find that

$$\hat{\theta}_i = \frac{x_i}{100}.$$

Then the likelihood ratio statistic is

$$\begin{split} \Lambda(\theta_0) &= -2\ln\left[\frac{L(\theta_0)}{L(\hat{\theta})}\right] = -2\ln\left[\frac{\prod_{i=1}^6(\frac{1}{6})^{X_i}}{\prod_{i=1}^6(\frac{X_i}{100})^{X_i}}\right] \\ &= 2\sum_{i=1}^6 X_i \ln\frac{6X_i}{100} \overset{D}{\to} W \sim \chi^2(5) \end{split}$$

where  $\hat{\theta} = (\hat{\theta}_1, ..., \hat{\theta}_6)$ , and the degree of freedom is because  $\theta_6 = 1 - \sum_{i=1}^5 \theta_i$ . The observed likelihood ratio statistic is

$$\lambda(\theta_0) = 2\sum_{i=1}^{6} x_i \ln \frac{6x_i}{100}$$

$$= -1.3063 - 3.1608 - 4.8819 + 7.2929 + 12.2158 - 6.4600$$

$$= 3.6997.$$

And so the p-value, by

pchisq(3.6997, 5, lower.tail=FALSE)

is

$$p \approx P(W \ge 3.6997) \approx 0.5934$$
.

This shows that we do not have sufficient evidence to reject  $H_0$ , i.e. we have a fair die.

# Example 23.1.5 (Example 7.5)

Suppose  $X_1, ..., X_n \sim N(\mu_0, \sigma^2)$  is a random sample. Test the hypothesis  $H_0: \sigma^2 = \sigma_0^2$ , where the mean  $\mu_0$  is known.

# Solution

Let  $H_a: \sigma^2 \neq \sigma_0^2$ . Since the  $X_i$ 's are IID, we have

$$L(\theta) = \prod_{i=1}^{n} \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{(x_i - \mu_0)^2}{2\sigma^2}\right\}$$

$$= (2\pi)^{-\frac{n}{2}} (\sigma^2)^{-\frac{n}{2}} \exp\left\{-\frac{1}{2\sigma^2} \sum_{i=1}^{n} (x_i - \mu_0)^2\right\}$$

$$\ell(\sigma^2) = -\frac{n}{2} \ln 2\pi - \frac{n}{2} \sigma^2 - \frac{1}{2\sigma^2} \sum_{i=1}^{n} (x_i - \mu_0)^2$$

$$\frac{d\ell}{d\sigma^2} = -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_{i=1}^{n} (x_i - \mu_0)^2$$

$$\implies \hat{\sigma} = \frac{1}{n} \sum_{i=1}^{n} (x_i - \mu_0)^2$$

$$\frac{d^2\ell}{d(\sigma^2)^2} = \frac{n}{2\sigma^4} - \frac{1}{\sigma^6} \sum_{i=1}^{n} (x_i - \mu_0)^2$$

$$\frac{d^2\ell}{d(\sigma^2)^2} \Big|_{\sigma^2 = \hat{\sigma}^2} = \frac{n^3}{\left(\sum_{i=1}^{n} (x_i - \mu_0)^2\right)^2} (\frac{1}{2} - n^3) < 0$$

Then, under the null hypothesis, note that

$$\begin{split} R(\sigma^2) &= \frac{L(\sigma_0^2)}{L(\tilde{\sigma}^2)} = \frac{\left(\sigma_0^2\right)^{-\frac{n}{2}}}{\left(\tilde{\sigma}^2\right)^{-\frac{n}{2}}} \exp\left\{ \left(\frac{1}{2\tilde{\sigma}^2} - \frac{1}{2\sigma_0^2}\right) \sum_{i=1}^n (X_i - \mu_0)^2 \right\} \\ &= \left(\frac{\tilde{\sigma}^2}{\sigma_0^2}\right)^{\frac{n}{2}} \exp\left\{ \frac{n}{2} \left(1 - \frac{\tilde{\sigma}^2}{\sigma_0^2}\right) \right\} \end{split}$$

and thus

$$\begin{split} \Lambda(\sigma_0^2) &= -2 \ln R(\sigma_0^2) = -n \ln \left( \frac{\tilde{\sigma}^2}{\sigma_0^2} \right) - n \left( 1 - \frac{\tilde{\sigma}^2}{\sigma_0^2} \right) \\ &= n \left[ \frac{\tilde{\sigma}^2}{\sigma_0^2} - 1 - \ln \left( \frac{\tilde{\sigma}^2}{\sigma_0^2} \right) \right] \overset{D}{\to} W \sim \chi^2(1). \end{split}$$

The observed likelihood ratio statistic is therefore

$$\lambda(\sigma_0^2) = n \left[ \frac{\hat{\sigma}^2}{\sigma_0^2} - 1 - \ln \left( \frac{\hat{\sigma}^2}{\sigma_0^2} \right) \right].$$

Then the p-value can be obtained by evaluating

$$p pprox P\left(W \ge n\left[\frac{\hat{\sigma}^2}{\sigma_0^2} - 1 - \ln\left(\frac{\hat{\sigma}^2}{\sigma_0^2}\right)\right]\right)$$

# A Regularity Conditions

Throughout this course, the phrase "regularity conditions" have been regularly coined but never formally introduced. To clarify that, I have sought out for the instructor and confirmed that these are the regularity conditions that are used. This is taken out from *Casella & Berger 2002*<sup>1</sup>.

"(R)egularity conditions" are typically very technical, rather boring, and usually satisfied in most reasonable problems. But they are a necessary evil...

<sup>1</sup> R. L. Casella, G. Berger. *Statistical Inference*. Thomson Learning, 2002

#### Furthermore,

These conditions mainly relate to differentiability of the density<sup>2</sup> and the ability to interchange differentiation and integration.

The following are the regularity conditions:

- 1.  $X_1, ..., X_n \sim f(x; \theta)$  are IID;
- 2. The unknown parameter is **identifiable**, i.e. if  $\theta \neq \theta'$ , then  $f(x;\theta) \neq f(x;\theta')$ ;
- 3. The pf's  $f(x;\theta)$  have common support, and  $f(x;\theta)$  is differentiable in  $\theta$ ;
- 4. The parameter space  $\Omega$  contains an open set  $\omega \subset \Omega$  of which the true parameter value  $\theta_0$  is an interior point<sup>3</sup>.
- 5.  $\forall x \in \text{supp}(X)$ , the pf  $f(x;\theta)$  is three times differentiable with respect to  $\theta$ , and the third derivative is continuous on  $\theta$ , and  $\int f(x;\theta)$  can be differentiated three times under the integral sign.
- 6.  $\forall \theta_0 \in \Omega, \exists c \in \mathbb{R}_{>0}, \exists a \text{ function } M(x) \text{ }^4 \text{ such that } \forall x \in \text{supp}(X),$

$$|\theta - \theta_0| < c \implies \left| \frac{\partial^3}{\partial \theta^3} \ln f(x; \theta) \right| \le M(x)$$

<sup>2</sup> probability function

 $^{_{3}}\omega$  is an open set, and so we can really only talk about interior points. The laymen definition of an interior point, is a point that is in a set but is not the limit point of the set.

<sup>4</sup> Note that both c and M(x) can depend on  $\theta_0$ , by the way that the statement is laid out.

with 
$$E_{\theta_0}[M(X)] \leq \infty$$
 5

 $^{5}$   $E_{\theta_0}$  is the expected value operator that acts on the variable  $\theta_0$ .

# B Useful References

# B.1 Commonly Used Distributions

Distribution	pf	Mean	Var	mgf			
Binomial Distribution : $X \sim Bin(n, p)$							
$x \in \mathbb{N} \cup \{0\}$							
$n \in \mathbb{N}$	$\binom{n}{x}p^x(1-p)^{n-x}$	пр	np(1-p)	$(pe^t + 1 - p)^n$			
0 < p < 1							
Geometric Distribution : $X \sim \text{Geo}(p)$							
$x \in \mathbb{N} \cup \{0\}$	$p(1-p)^x$	$\frac{1-p}{p}$	$\frac{1-p}{p^2}$	$\frac{p}{1-(1-p)e^t}$			
$0$				$t < -\log(1-p)$			
Poisson Distribution : $X \sim \text{Poi}(\mu)$							
$x \in \mathbb{N} \cup \{0\}$	$\frac{e^{-\mu}\mu^x}{x!}$	μ	$\mu$	$e^{\mu(e^t-1)}$			
$\mu > 0$							
Uniform Distribution : $X \sim \text{Unif}(a, b)$							
$a \le x \le b$	$\frac{1}{b-a}$	$\frac{a+b}{2}$	$\frac{(b-a)^2}{12}$	$\frac{e^{bt}-e^{at}}{(b-a)^t}$			
$a < b \in \mathbb{R}$				$t \neq 0$			
Normal Distribution : $X \sim N(\mu, \sigma)$							
$x \in \mathbb{R}$							
	$\frac{1}{\sqrt{2\pi\sigma^2}}e^{-\frac{(x-\mu)^2}{2\sigma^2}}$	μ	$\sigma^2$	$e^{\mu t + rac{\sigma^2 t^2}{2}}$			

$$\mu \in \mathbb{R}$$

Gamma Distribution :  $X \sim \Gamma(\alpha, \beta)$ 

$$x \in \mathbb{R}_{>0}$$

$$\alpha > 0$$

$$\frac{1}{\beta^{\alpha}\Gamma(\alpha)}x^{\alpha-1}e^{-\frac{x}{\beta}}$$

$$\alpha\beta^2$$

 $\theta^2$ 

$$(1-\beta t)^{-\alpha}$$

$$t < \frac{1}{\beta}$$

Exponential Distribution :  $X \sim \text{Exp}(\theta)$ 

$$x \in \mathbb{R}_{>0}$$

$$\frac{1}{\theta}e^{-\frac{x}{\theta}}$$

$$\theta > 0$$

$$t < \frac{1}{\theta}$$

Negative Binomial Distribution :  $X \sim NB(r, p)$ 

$$x \in \mathbb{N} \cup \{0\}$$

$$\binom{x+r-1}{x} \cdot (1-p)^r p^x$$
  $\frac{pr}{1-p}$ 

$$\frac{pr}{(1-p)^2}$$

$$\left(\frac{1-p}{1-pe^t}\right)^r$$

$$t < -\log p$$

 $p \in (0,1)$ 

r > 0

Beta Distribution :  $X \sim \text{Beta}(\alpha, \beta)$ 

 $x \in [0,1]$  or

$$x \in (0,1)$$

$$\frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+b)}x^{\alpha-1}(1-x)^{\beta-1}$$

$$\frac{\alpha}{(\alpha+\beta)^2(\alpha+\beta+1)}$$
  $\frac{\alpha\beta}{(\alpha+\beta)^2(\alpha+\beta+1)}$ 

$$\frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(a+b)}x^{\alpha-1}(1-x)^{\beta-1} \qquad \frac{\alpha}{\alpha+\beta} \qquad \frac{\alpha\beta}{(\alpha+\beta)^2(\alpha+\beta+1)} \qquad 1+\sum_{k=1}^{\infty} \left(\prod_{r=0}^{k-1} \frac{\alpha+r}{\alpha+\beta+r}\right) \frac{t^k}{k!}$$

Chi-Squared Distribution :  $X \sim \chi^2(k)$ 

 $k \in \mathbb{N}$ 

$$\frac{1}{2^{\frac{k}{2}}\Gamma(\frac{k}{2})}x^{\frac{k}{2}-1}e^{-\frac{x}{2}}$$

$$(1-2t)^{-\frac{t}{2}}$$

if k = 1

$$\frac{1}{2^{\frac{k}{2}}\Gamma(\frac{k}{2})}x^{\frac{\kappa}{2}-1}e^{-\frac{k}{2}}$$

$$t < \frac{1}{2}$$

 $x \in (0, \infty)$ 

otherwise

$$x \in [0, \infty]$$

# Index

 $\delta$ -Method, 154  $\sigma$ -algebra, 17  $\sigma$ -field, 17 p-Value, 202

Alternative Hypothesis, 201
Asymptotic Confidence Intervals, 196
Asymptotic Distribution, 135
Asymptotic Distribution of Relative
Likelihood, 183
Asymptotic Normality, 183
Asymptotic Pivotal Quantity, 188
asymptotic variance, 183

Bernoulli's Theorem, 143
Beta Distribution, 88
Binomial Distribution, 29
Bivariate Normal Distribution, 114
Bonferroni's Inequality, 11, 21
Boole's Inequality, 11, 21

Central Limit Theorem, 148
Chebyshev's Inequality, 54
Clairaut's Theorem, 71
Cochran's Theorem, 132
Conditional Distributions, 83
Conditional Expectation, 101
conditional mean, 101
Conditional PF, 83
Conditional Probability, 21
conditional variance, 101
Confidence Interval, 189
Consistency, 183
consistent estimator, 183
Continuity Property, 11, 22

Continuous Mapping Theorem, 152 Continuous Random Variable, 27 Convergence in Distribution, 135 Convergence in Probability, 139 Correlation Coefficient, 99 Covariance, 93 Cumulative Distribution Function, 23

Degenerate Distribution, 138
Discrete Random Variable, 24
Double Parameter Exponential Distribution, 144

Estimator, 158
Estimator, 158
expected value, 44, 45
Exponential Distribution, 31

F Distribution, 132 factorial moment, 51 Factorization Theorem, 82 Fisher Information, 164 Fisher Information Matrix, 175

Gamma Distribution, 32 Gaussian Distribution, 30, 131 Geometric Distribution, 29

Hypothesis, 201

implausible, 171
Independence, 70, 78
Independent Events, 21
Indicator Function, 57
Information Function, 163

Information Matrix, 175 Interval Estimate, 188 Interval Estimators, 188 Invariance Property of MLE, 169

Jacobian, 121, 122
Joint CDF, 63
Joint Continuous Random Variables, 74
joint density function, 74
Joint Discrete Random Variables, 67
Joint Expectation, 89
Joint Moment Generating Functions, 107
Joint Moments, 108
Joint PMF, 67

k-variate CDF, 111 k-Variate Joint MGF, 112 k-variate Support Set, 111 Kolmogorov Axioms, 18

Law of the Unconscious Statistician, 49, 89

Law of Total Expectation, 104

Law of Total Variance, 106

Likelihood function, 159

Likelihood Interval, 172, 198

Likelihood Ratio Statistic, 204

Likelihood Region, 172, 175

Limiting Distribution, 135

Linearity - Expectation, 49, 91

Location Family, 42

Location-Scale Family, 42

Log Relative Likelihood, 172 Log-likelihood, 159 Lower Confidence Limit, 189

Marginal CDF, 66
Marginal Distribution, 68
Marginal MGF, 108
Marginal Probability Density Function, 77
Markov's Inequality, 52
Markov's Inequality 2, 53
maximum likelihood estimate, 159
Measurable Space, 18
Moment Generating Function, 56
Moments, 51
Mutlinomial Distribution, 113

Normal Distribution, 30 Null Hypothesis, 201

observed information, 163 observed information matrix, 175 One-to-One Bivariate Transformations, 121 One-to-One Transformation, 120 Pearson Correlation Coefficient, 99
Pivotal Quantity, 188
plausible, 171
Poisson Distribution, 30
power set, 17
probability axioms, 18
probability density function, 27
Probability Integral Transformation,

probability mass function, 24
Probability Measure, 18
probability set function, 18
probability space, 18
Product Rule, 86
Properties of pdf, 27
Properties of pmf, 24
Properties of the cdf, 11, 23

Random Sample, 115 Random Variable, 22 Relative Likelihood, 171 right-continuous, 23

Sample Space, 17 Scale Family, 42 Scale Parameter, 42 Score Function, 163
Score Vector, 175
Significance Level, 202
Simple Hypothesis, 204
Slutsky's Theorem, 151
Standard Normal Distribution, 31
Statistic, 157
support set, 24, 67, 74

t-test, 132
Taylor Series with Lagrange's Remainder, 135
Test of Hypothesis, 202
Test Statistic, 202
Type I Error, 203
Type II Error, 203

unbiased estimator, 158 uncorrelated, 93 Uniform Distribution, 31 Upper Confidence Limit, 189

Variance, 50

Weak Law of Large Numbers, 143

# Bibliography

- R. L. Casella, G. Berger. Statistical Inference. Thomson Learning, 2002.
- C. J. Geyer. Stat 5102 notes: Fisher information and confidence intervals using maximum likelihood. https://www.stat.umn.edu/geyer/s06/5102/notes/fish.pdf, March 2007.