# STAT 6710Mathematical Statistics I Fall Semester 2000

# Dr. Jürgen Symanzik

Utah State University

Department of Mathematics and Statistics

3900 Old Main Hill

Logan, UT 84322-3900

Tel.: (435) 797–0696

FAX: (435) 797-1822

e-mail: symanzik@sunfs.math.usu.edu

# Contents

$\mathbf{A}$	ckno	wledgements	1				
1	Axi	Axioms of Probability					
	1.1	$\sigma$ –Fields	1				
	1.2	Manipulating Probability	5				
	1.3	Combinatorics and Counting	11				
	1.4	Conditional Probability and Independence	17				
2	Rar	ndom Variables	24				
	2.1	Measurable Functions	24				
	2.2	Probability Distribution of a Random Variable	27				
	2.3	Discrete and Continuous Random Variables	31				
	2.4	Transformations of Random Variables	36				
3	Mo	ments and Generating Functions	42				
	3.1	Expectation	42				
	3.2	Generating Functions	50				
	3.3	Complex–Valued Random Variables and Characteristic Functions	56				
	3.4	Probability Generating Functions	66				
	3.5	Moment Inequalities	68				
4	Rar	ndom Vectors	71				
	4.1	Joint, Marginal, and Conditional Distributions	71				
	4.2	Independent Random Variables	77				
	4.3	Functions of Random Vectors	82				
	4.4	Order Statistics	89				
	4.5	Multivariate Expectation	91				
	4.6	Multivariate Generating Functions	97				
	4.7	Conditional Expectation	102				
	4.8	Inequalities and Identities	106				

5	Particular Distributions			
	5.1	Multivariate Normal Distributions	. 112	
	5.2	Exponential Family of Distributions	. 119	
6	Lim	tit Theorems	121	
	6.1	Modes of Convergence	. 122	
	6.2	Weak Laws of Large Numbers	. 135	
In	$\mathbf{dex}$		139	

# Acknowledgements

I would like to thank my students, Hanadi B. Eltahir, Rich Madsen, and Bill Morphet, who helped during the Fall 1999 and Spring 2000 semesters in typesetting these lecture notes using LaTeX and for their suggestions how to improve some of the material presented in class. Thanks are also due to my Fall 2000 and Spring 2001 semester students, Weiping Deng, David B. Neal, Emily Simmons, Shea Watrin, Aonan Zhang, and Qian Zhao, for their comments that helped to improve these lecture notes.

In addition, I particularly would like to thank Mike Minnotte and Dan Coster, who previously taught this course at Utah State University, for providing me with their lecture notes and other materials related to this course. Their lecture notes, combined with additional material from Rohatgi (1976) and other sources listed below, form the basis of the script presented here.

The textbook required for this class is:

Rohatgi, V. K. (1976): An Introduction to Probability Theory and Mathematical Statistics, John Wiley and Sons, New York.

A Web page dedicated to this class is accessible at:

http://www.math.usu.edu/~symanzik/teaching/2000\_stat6710/stat6710.html

This course closely follows Rohatgi (1976) as described in the syllabus. Additional material originates from the lectures from Professors Hering, Trenkler, and Gather I have attended while studying at the Universität Dortmund, Germany, the collection of Masters and PhD Preliminary Exam questions from Iowa State University, Ames, Iowa, and the following textbooks:

- Bandelow, C. (1981): Einführung in die Wahrscheinlichkeitstheorie, Bibliographisches Institut, Mannheim, Germany.
- Casella, G., and Berger, R. L. (1990): *Statistical Inference*, Wadsworth & Brooks/Cole, Pacific Grove, CA.
- Fisz, M. (1989): Wahrscheinlichkeitsrechnung und mathematische Statistik, VEB Deutscher Verlag der Wissenschaften, Berlin, German Democratic Republic.
- Kelly, D. G. (1994): Introduction to Probability, Macmillan, New York, NY.
- Mood, A. M., and Graybill, F. A., and Boes, D. C. (1974): Introduction to the Theory of Statistics (Third Edition), McGraw-Hill, Singapore.
- Parzen, E. (1960): Modern Probability Theory and Its Applications, Wiley, New York, NY.

• Searle, S. R. (1971): Linear Models, Wiley, New York, NY.

Additional definitions, integrals, sums, etc. originate from the following formula collections:

- Bronstein, I. N. and Semendjajew, K. A. (1985): Taschenbuch der Mathematik (22. Auflage), Verlag Harri Deutsch, Thun, German Democratic Republic.
- Bronstein, I. N. and Semendjajew, K. A. (1986): Ergänzende Kapitel zu Taschenbuch der Mathematik (4. Auflage), Verlag Harri Deutsch, Thun, German Democratic Republic.
- Sieber, H. (1980): Mathematische Formeln Erweiterte Ausgabe E, Ernst Klett, Stuttgart, Germany.

Jürgen Symanzik, January 12, 2001

Lecture 02: We 08/30/00

# 1 Axioms of Probability

# 1.1 $\sigma$ -Fields

Let  $\Omega$  be the **sample space** of all possible outcomes of a chance experiment. Let  $\omega \in \Omega$  (or  $x \in \Omega$ ) be any outcome.

# Example:

Count # of heads in n coin tosses.  $\Omega = \{0, 1, 2, \dots, n\}$ .

Any subset A of  $\Omega$  is called an **event**.

For each event  $A \subseteq \Omega$ , we would like to assign a number (i.e., a probability). Unfortunately, we cannot always do this for every subset of  $\Omega$ .

Instead, we consider classes of subsets of  $\Omega$  called *fields* and  $\sigma$ -fields.

# Definition 1.1.1:

A class L of subsets of  $\Omega$  is called a **field** if  $\Omega \in L$  and L is *closed* under complements and finite unions, i.e., L satisfies

- (i)  $\Omega \in L$
- (ii)  $A \in L \Longrightarrow A^C \in L$
- (iii)  $A, B \in L \Longrightarrow A \cup B \in L$

Since  $\Omega^C = \emptyset$ , (i) and (ii) imply  $\emptyset \in L$ . Therefore, (i)':  $\emptyset \in L$  [can replace (i)].

#### Note: De Morgan's Laws

For any class A of sets, and sets  $A \in A$ , it holds:

$$\bigcup_{A\in\mathcal{A}}A=(\bigcap_{A\in\mathcal{A}}A^C)^C \text{ and } \bigcap_{A\in\mathcal{A}}A=(\bigcup_{A\in\mathcal{A}}A^C)^C.$$

#### Note:

So (ii), (iii) imply (iii)':  $A, B \in L \Longrightarrow A \cap B \in L$  [can replace (iii)].

#### Proof:

$$A,B\in L \xrightarrow{(ii)} A^C,B^C \in L \xrightarrow{(iii)} (A^C \cup B^C) \in L \xrightarrow{(ii)} (A^C \cup B^C)^C \in L \xrightarrow{DM} A \cap B \in L$$

#### Definition 1.1.2:

A class L of subsets of  $\Omega$  is called a  $\sigma$ -field (Borel field,  $\sigma$ -algebra) if it is a field and closed under countable unions, i.e.,

(iv) 
$$\{A_n\}_{n=1}^{\infty} \in L \Longrightarrow \bigcup_{n=1}^{\infty} A_n \in L.$$

# Note:

(iv) implies (iii) by taking  $A_n = \emptyset$  for  $n \geq 3$ .

# Example 1.1.3:

For some  $\Omega$ , let L contain all finite and all cofinite sets (A is *cofinite* if  $A^C$  is finite — for example, if  $\Omega = \mathbb{N}$ ,  $A = \{x \mid x \geq c\}$  is not finite but since  $A^C = \{x \mid x < c\}$  is finite, A is cofinite). Then L is a field. But L is a  $\sigma$ -field **iff** (if and only if)  $\Omega$  is finite.

For example, let  $\Omega = Z$ . Take  $A_n = \{n\}$ , each finite, so  $A_n \in L$ . But  $\bigcup_{n=1}^{\infty} A_n = Z^+ \notin L$ , since

the set is not finite (it is infinite) and also not cofinite  $((\bigcup_{n=1}^{\infty} A_n)^C = Z_0^-)$  is infinite, too).

Question: Does this construction work for  $\Omega = Z^+$ ??

# Note:

The largest  $\sigma$ -field in  $\Omega$  is the power set  $\mathcal{P}(\Omega)$  of all subsets of  $\Omega$ . The smallest  $\sigma$ -field is  $L = \{\emptyset, \Omega\}$ .

# Terminology:

A set  $A \in L$  is said to be "measurable L".

# Note:

We often begin with a class of sets, say a, which may not be a field or a  $\sigma$ -field.

# Definition 1.1.4:

The  $\sigma$ -field generated by a,  $\sigma(a)$ , is the smallest  $\sigma$ -field containing a, or the intersection of all  $\sigma$ -fields containing a.

# Note:

(i) Such  $\sigma$ -fields containing a always exist (e.g.,  $\mathcal{P}(\Omega)$ ), and (ii) the intersection of an arbitrary # of  $\sigma$ -fields is always a  $\sigma$ -field.

# Proof:

(ii) Suppose  $L = \bigcap_{\theta} L_{\theta}$ . We have to show that conditions (i) and (ii) of Def. 1.1.1 and (iv) of

Def. 1.1.2 are fulfilled:

- (i)  $\Omega \in L_{\theta} \ \forall \theta \Longrightarrow \Omega \in L$
- (ii) Let  $A \in L \Longrightarrow A \in L_{\theta} \ \forall \theta \Longrightarrow A^C \in L_{\theta} \ \forall \theta \Longrightarrow A^C \in L$

(iv) Let 
$$A_n \in L \ \forall n \Longrightarrow A_n \in L_\theta \ \forall \theta \ \forall n \Longrightarrow \bigcup_n A_n \in L_\theta \ \forall \theta \Longrightarrow \bigcup_n A_n \in L$$

Lecture 03: Fr 09/01/00

# Example 1.1.5:

$$\Omega = \{0,1,2,3\}, a = \{\{0\}\}, b = \{\{0\},\{0,1\}\}.$$

What is  $\sigma(a)$ ?

 $\sigma(a)$ : must include  $\Omega, \emptyset, \{0\}$ 

also:  $\{1, 2, 3\}$  by 1.1.1 (ii)

Since all unions are included, we have  $\sigma(a) = \{\Omega, \emptyset, \{0\}, \{1, 2, 3\}\}\$ 

What is  $\sigma(b)$ ?

 $\sigma(b)$ : must include  $\Omega, \emptyset, \{0\}, \{0, 1\}$ 

also:  $\{1, 2, 3\}, \{2, 3\}$  by 1.1.1 (ii)

 $\{0, 2, 3\}$  by 1.1.1 (iii)

{1} by 1.1.1 (ii)

Since all unions are included, we have  $\sigma(b) = \{\Omega, \emptyset, \{0\}, \{1\}, \{0, 1\}, \{2, 3\}, \{0, 2, 3\}, \{1, 2, 3\}\}$ 

#### Note:

If  $\Omega$  is finite or countable, we will usually use  $L = \mathcal{P}(\Omega)$ . If  $|\Omega| = n < \infty$ , then  $|L| = 2^n$ .

If  $\Omega$  is uncountable,  $\mathcal{P}(\Omega)$  may be too large to be useful and we may have to use some smaller  $\sigma$ -field.

#### Definition 1.1.6:

If  $\Omega = \mathbb{R}$ , an important special case is the **Borel**  $\sigma$ -field, i.e., the  $\sigma$ -field generated from all half-open intervals of the form (a, b], denoted  $\mathcal{B}$  or  $\mathcal{B}_1$ . The sets of  $\mathcal{B}$  are called **Borel sets**.

The Borel  $\sigma$ -field on  $\mathbb{R}^d$  ( $\mathcal{B}_d$ ) is the  $\sigma$ -field generated by d-dimensional rectangles of the form  $\{(x_1, x_2, \dots, x_d) \mid a_i < x_i \leq b_i; i = 1, 2, \dots, d\}.$ 

#### Note:

 $\mathcal{B}$  contains all points:  $\{x\} = \bigcap_{n=1}^{\infty} (x - \frac{1}{n}, x]$ 

closed intervals:  $[x, y] = (x, y] + \{x\} = (x, y] \cup \{x\}$ 

open intervals:  $(x, y) = (x, y] - \{y\} = (x, y] \cap \{y\}^C$ 

and semi–infinite intervals:  $(x, \infty) = \bigcup_{n=1}^{\infty} (x, x + n]$ 

# Note:

We now have a measurable space  $(\Omega, L)$ . We next define a probability measure  $P(\cdot)$  on  $(\Omega, L)$  to obtain a probability space  $(\Omega, L, P)$ .

# Definition 1.1.7: Kolmogorov Axioms of Probability

A probability measure (pm), P, on  $(\Omega, L)$  is a set function  $P: L \to \mathbb{R}$  satisfying

- (i)  $0 \le P(A) \ \forall A \in L$
- (ii)  $P(\Omega) = 1$
- (iii) If  $\{A_n\}_{n=1}^{\infty}$  are disjoint sets in L and  $\bigcup_{n=1}^{\infty} A_n \in L$ , then  $P(\bigcup_{n=1}^{\infty} A_n) = \sum_{n=1}^{\infty} P(A_n)$ .

#### Note:

 $\bigcup_{n=1}^{\infty} A_n \in L \text{ holds automatically if } L \text{ is a } \sigma\text{-field but it is needed as a precondition in the case that } L \text{ is just a field. Property (iii) is called } countable \ additivity.}$ 

# 1.2 Manipulating Probability

#### <u>Theorem 1.2.1:</u>

For P a pm on  $(\Omega, L)$ , it holds:

- (i)  $P(\emptyset) = 0$
- (ii)  $P(A^C) = 1 P(A) \quad \forall A \in L$
- (iii)  $P(A) < 1 \ \forall A \in L$
- (iv)  $P(A \cup B) = P(A) + P(B) P(A \cap B) \quad \forall A, B \in L$
- (v) If  $A \subseteq B$ , then P(A) < P(B).

# Proof:

(i) 
$$A_n \equiv \emptyset \quad \forall n \Rightarrow \bigcup_{n=1}^{\infty} A_n = \emptyset \in L.$$

 $A_i \cap A_j = \varnothing \cap \varnothing = \varnothing \ \, \forall i,j \Rightarrow A_n \text{ are disjoint } \forall n.$ 

$$P(\emptyset) = P(\bigcup_{n=1}^{\infty} A_n) \stackrel{Def 1.1.7(iii)}{=} \sum_{n=1}^{\infty} P(A_n) = \sum_{n=1}^{\infty} P(\emptyset)$$

This can only hold if  $P(\emptyset) = 0$ .

(ii) 
$$A_1 \equiv A, A_2 \equiv A^C, A_n \equiv \emptyset \ \forall n \geq 3.$$

$$\Omega = \bigcup_{n=1}^{\infty} A_n = A_1 \cup A_2 \cup \bigcup_{n=3}^{\infty} A_n = A_1 \cup A_2 \cup \emptyset.$$

 $A_1 \cap A_2 = A_1 \cap \emptyset = A_2 \cap \emptyset = \emptyset \Rightarrow A_1, A_2, \emptyset$  are disjoint.

$$1 = P(\Omega) = P(\bigcup_{n=1}^{\infty} A_n)$$

$$\stackrel{Def1.1.7(iii)}{=} \sum_{n=1}^{\infty} P(A_n)$$

$$= P(A_1) + P(A_2) + \sum_{n=3}^{\infty} P(A_n)$$

$$\stackrel{Th1.2.1(i)}{=} P(A_1) + P(A_2)$$

$$= P(A) + P(A^C)$$

$$\implies P(A^C) = 1 - P(A) \ \forall A \in L.$$

- (iii) By Th. 1.2.1 (ii)  $P(A) = 1 P(A^C)$  $\implies P(A) \le 1 \quad \forall A \in L \text{ since } P(A^C) \ge 0 \text{ by Def. 1.1.7 (i)}.$
- (iv)  $A \cup B = (A \cap B^C) \cup (A \cap B) \cup (B \cap A^C)$ . So,  $(A \cup B)$  can be written as a union of disjoint sets  $(A \cap B^C)$ ,  $(A \cap B)$ ,  $(B \cap A^C)$ .

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

$$\stackrel{Def.1.1.7(iii)}{=} P(A \cap B^C) + P(A \cap B) + P(B \cap A^C)$$

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

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

$$\stackrel{Def.1.1.7(iii)}{=} P(A) + P(B) - P(A \cap B)$$

(v)  $B = (B \cap A^C) \cup A$  where  $(B \cap A^C)$  and A are disjoint sets.

Lecture 04: We 09/06/00

$$\begin{split} P(B) &= P((B \cap A^C) \cup A) \overset{Def 1.1.7(iii)}{=} P(B \cap A^C) + P(A) \\ &\Longrightarrow P(A) = P(B) - P(B \cap A^C) \\ &\Longrightarrow P(A) \leq P(B) \text{ since } P(B \cap A^C) \geq 0 \text{ by Def. } 1.1.7 \text{ (i)}. \end{split}$$

# Theorem 1.2.2: Principle of Inclusion-Exclusion

Let  $A_1, A_2, \ldots, A_n \in L$ . Then

$$P(\bigcup_{k=1}^{n} A_k) = \sum_{k=1}^{n} P(A_k) - \sum_{k_1 < k_2}^{n} P(A_{k_1} \cap A_{k_2}) + \sum_{k_1 < k_2 < k_3}^{n} P(A_{k_1} \cap A_{k_2} \cap A_{k_3}) - \dots + (-1)^{n+1} P(\bigcap_{k=1}^{n} A_k)$$

#### Proof:

n=1 is trivial

n=2 is Theorem 1.2.1 (iv)

Use induction for higher n (Homework).

#### Note:

A proof by **induction** consists of two steps:

First, we have to establish the **induction base**. For example, if we state that something holds for all non-negative integers, then we have to show that it holds for n = 0. Similarly, if we state that something holds for all integers, then we have to show that it holds for n = 1. Formally, it is sufficient to verify a claim for the smallest valid integer only. However, to get some feeling how the proof from n to n + 1 might work, it is sometimes beneficial to verify a claim for 1, 2, or 3 as well.

In the second step, we have to establish the result in the **induction step**, showing that something holds for n + 1, using the fact that it holds for n (alternatively, we can show that it holds for n, using the fact that it holds for n - 1).

# Theorem 1.2.3: Bonferroni's Inequality

Let  $A_1, A_2, \ldots, A_n \in L$ . Then

$$\sum_{i=1}^{n} P(A_i) - \sum_{i < j} P(A_i \cap A_j) \le P(\bigcup_{i=1}^{n} A_i) \le \sum_{i=1}^{n} P(A_i)$$

Proof:

Right side: 
$$P(\bigcup_{i=1}^{n} A_i) \le \sum_{i=1}^{n} P(A_i)$$

# Induction Base:

For n = 1, the right side evaluates to  $P(A_1) \leq P(A_1)$ , which is true.

Formally, the next step is not required. However, it does not harm to verify the claim for n=2 as well. For n=2, the right side evaluates to  $P(A_1 \cup A_2) \leq P(A_1) + P(A_2)$ .

$$P(A_1 \cup A_2) \stackrel{Th.1.2.1(iv)}{=} P(A_1) + P(A_2) - P(A_1 \cap A_2) \le P(A_1) + P(A_2)$$
 since  $P(A_1 \cap A_2) \ge 0$  by Def. 1.1.7 (i).

This establishes the induction base for the right side.

#### Induction Step:

We assume the right side is true for n and show that it is true for n + 1:

$$P(\bigcup_{i=1}^{n+1} A_i) = P((\bigcup_{i=1}^{n} A_i) \cup A_{n+1})$$

$$\stackrel{Th.1.2.1(iv)}{=} P(\bigcup_{i=1}^{n} A_i) + P(A_{n+1}) - P((\bigcup_{i=1}^{n} A_i) \cap A_{n+1})$$

$$\stackrel{Def.1.1.7(i)}{\leq} P(\bigcup_{i=1}^{n} A_i) + P(A_{n+1})$$

$$\stackrel{I.B.}{\leq} \sum_{i=1}^{n} P(A_i) + P(A_{n+1})$$

$$= \sum_{i=1}^{n+1} P(A_i)$$

Left side: 
$$\sum_{i=1}^{n} P(A_i) - \sum_{i < j} P(A_i \cap A_j) \le P(\bigcup_{i=1}^{n} A_i)$$

# <u>Induction Base:</u>

For n = 1, the left side evaluates to  $P(A_1) \leq P(A_1)$ , which is true.

For n=2, the left side evaluates to  $P(A_1)+P(A_2)-P(A_1\cap A_2)\leq P(A_1\cup A_2)$ , which is true by Th. 1.2.1 (iv).

For n = 3, the left side evaluates to

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

This holds since

$$\begin{split} P(A_1 \cup A_2 \cup A_3) \\ &= P((A_1 \cup A_2) \cup A_3) \\ &\stackrel{Th.1.2.1(iv)}{=} P(A_1 \cup A_2) + P(A_3) - P((A_1 \cup A_2) \cap A_3) \\ &= P(A_1 \cup A_2) + P(A_3) - P((A_1 \cap A_3) \cup (A_2 \cap A_3)) \\ &\stackrel{Th.1.2.1(iv)}{=} P(A_1) + P(A_2) - P(A_1 \cap A_2) + P(A_3) - P(A_1 \cap A_3) - P(A_2 \cap A_3) \\ &+ P((A_1 \cap A_3) \cap (A_2 \cap A_3)) \\ &= P(A_1) + P(A_2) + P(A_3) - P(A_1 \cap A_2) - P(A_1 \cap A_3) - P(A_2 \cap A_3) + P(A_1 \cap A_2 \cap A_3) \\ &\stackrel{Def1.1.7(i)}{\geq} P(A_1) + P(A_2) + P(A_3) - P(A_1 \cap A_2) - P(A_1 \cap A_3) - P(A_2 \cap A_3) \end{split}$$

This establishes the induction base for the left side.

# Induction Step:

We assume the left side is true for n and show that it is true for n + 1:

$$P(\bigcup_{i=1}^{n+1} A_i) = P((\bigcup_{i=1}^{n} A_i) \cup A_{n+1})$$

$$= P((\bigcup_{i=1}^{n} A_i) + P(A_{n+1}) - P((\bigcup_{i=1}^{n} A_i) \cap A_{n+1})$$

$$\stackrel{left\ I.B.}{\geq} \sum_{i=1}^{n} P(A_i) - \sum_{i < j}^{n} P(A_i \cap A_j) + P(A_{n+1}) - P((\bigcup_{i=1}^{n} A_i) \cap A_{n+1})$$

$$= \sum_{i=1}^{n+1} P(A_i) - \sum_{i < j}^{n} P(A_i \cap A_j) - P(\bigcup_{i=1}^{n} (A_i \cap A_{n+1}))$$

$$\stackrel{Th.1.2.3\ right\ side}{\geq} \sum_{i=1}^{n+1} P(A_i) - \sum_{i < j}^{n} P(A_i \cap A_j) - \sum_{i=1}^{n} P(A_i \cap A_{n+1})$$

$$= \sum_{i=1}^{n+1} P(A_i) - \sum_{i < j}^{n+1} P(A_i \cap A_j)$$

# Theorem 1.2.4: Boole's Inequality

Let  $A, B \in L$ . Then

(i) 
$$P(A \cap B) \ge P(A) + P(B) - 1$$

(ii) 
$$P(A \cap B) \ge 1 - P(A^C) - P(B^C)$$

# Proof:

Homework

# <u>Definition 1.2.5:</u> Continuity of sets

For a sequence of sets  $\{A_n\}_{n=1}^{\infty}$ ,  $A_n \in L$  and  $A \in L$ , we say

(i) 
$$A_n \uparrow A$$
 if  $A_1 \subseteq A_2 \subseteq A_3 \subseteq \dots$  and  $A = \bigcup_{n=1}^{\infty} A_n$ .

(ii) 
$$A_n \downarrow A$$
 if  $A_1 \supseteq A_2 \supseteq A_3 \supseteq \dots$  and  $A = \bigcap_{n=1}^{\infty} A_n$ .

# Theorem 1.2.6:

If  $\{A_n\}_{n=1}^{\infty}$ ,  $A_n \in L$  and  $A \in L$ , then  $\lim_{n \to \infty} P(A_n) = P(A)$  if 1.2.5 (i) or 1.2.5 (ii) holds.

# Proof:

Part (i): Assume that 1.2.5 (i) holds.

Let 
$$B_1 = A_1$$
 and  $B_k = A_k - A_{k-1} = A_k \cap A_{k-1}^C \ \forall k \geq 2$ 

By construction,  $B_i \cap B_j = \emptyset$  for  $i \neq j$ 

It is 
$$A = \bigcup_{n=1}^{\infty} A_n = \bigcup_{n=1}^{\infty} B_n$$

and also 
$$A_n = \bigcup_{i=1}^n A_i = \bigcup_{i=1}^n B_i$$

$$\underline{\underline{P(A)}} = \underline{P(\bigcup_{k=1}^{\infty} B_k)}^{\text{Def. 1.1.7 (iii)}} \stackrel{\sum}{\underset{k=1}{\overset{\infty}{=}}} P(B_k) = \lim_{n \to \infty} [\sum_{k=1}^{n} P(B_k)]$$

$$\stackrel{\text{Def. 1.1.7 (iii)}}{=} \lim_{n \to \infty} [P(\bigcup_{k=1}^{n} B_k)] = \lim_{n \to \infty} [P(\bigcup_{k=1}^{n} A_k)] = \lim_{n \to \infty} P(A_n)$$

The last step is possible since  $A_n = \bigcup_{k=1}^n A_k$ 

Part (ii): Assume that 1.2.5 (ii) holds. Then, 
$$A_1^C \subseteq A_2^C \subseteq A_3^C \subseteq \dots$$
 and  $A^C = (\bigcap_{n=1}^{\infty} A_n)^{C} \stackrel{De\ Morgan}{=} \bigcup_{n=1}^{\infty} A_n^C$ 

$$P(A^C) \stackrel{\text{By Part (i)}}{=} \lim_{n \to \infty} P(A_n^C)$$

So 
$$1 - P(A^C) = 1 - \lim_{n \to \infty} P(A_n^C)$$

$$\Longrightarrow P(A) = \lim_{n \to \infty} (1 - P(A_n^C)) = \lim_{n \to \infty} P(A_n)$$

# Theorem 1.2.7:

Lecture 05: Fr 09/08/00

- (i) Countable unions of probability 0 sets have probability 0.
- (ii) Countable intersections of probability 1 sets have probability 1.

# Proof:

# Part (i):

Let 
$$\{A_n\}_{n=1}^{\infty} \in \mathcal{L}, P(A_n) = 0 \ \forall n$$

$$0 \stackrel{\text{Def. 1.7.7 (i)}}{\leq} P(\bigcup_{n=1}^{\infty} A_n) \stackrel{\text{Th. 1.2.3}}{\leq} \sum_{n=1}^{\infty} P(A_n) = \sum_{n=1}^{\infty} 0 = 0$$

Therefore 
$$P(\bigcup_{n=1}^{\infty} A_n) = 0$$
.

# Part (ii):

Let 
$$\{A_n\}_{n=1}^{\infty} \in \mathcal{L}, P(A_n) = 1 \ \forall n$$

$$\overset{\text{Th. 1.2.1 (ii)}}{\Longrightarrow} P(A_n^C) = 0 \quad \forall n \overset{\text{Th. 1.2.7 (ii)}}{\Longrightarrow} P(\bigcup_{n=1}^{\infty} A_n^C) = 0 \overset{\text{De Morgan}}{\Longrightarrow} P(\bigcap_{n=1}^{\infty} A_n) = 1$$

# 1.3 Combinatorics and Counting

For now, we restrict ourselves to sample spaces containing a finite number of points.

Let 
$$\Omega = \{\omega_1, \dots, \omega_n\}$$
 and  $L = \mathcal{P}(\Omega)$ . For any  $A \in L, P(A) = \sum_{\omega_j \in A} P(\omega_j)$ .

# Definition 1.3.1:

We say the elements of  $\Omega$  are **equally likely** (or occur with uniform probability) if  $P(\omega_j) = \frac{1}{n} \ \forall j = 1, \dots, n$ .

#### Note:

If this is true,  $P(A) = \frac{\text{number } \omega_j \text{ in } A}{\text{number } \omega_j \text{ in } \Omega}$ . Therefore, to calculate such probabilities, we just need to be able to count elements accurately.

# Theorem 1.3.2: Fundamental Theorem of Counting

If we wish to select one element  $(a_1)$  out of  $n_1$  choices, a second element  $(a_2)$  out of  $n_2$  choices, and so on for a total of k elements, there are

$$n_1 \times n_2 \times n_3 \times \ldots \times n_k$$

ways to do it.

<u>Proof:</u> (By Induction)

# Induction Base:

k = 1: trivial

k=2:  $n_1$  ways to choose  $a_1$ . For each,  $n_2$  ways to choose  $a_2$ .

Total # of ways = 
$$\underbrace{n_2 + n_2 + \ldots + n_2}_{n_1 \text{ times}} = n_1 \times n_2.$$

# Induction Step:

Suppose it is true for (k-1). We show that it is true for k=(k-1)+1.

There are  $n_1 \times n_2 \times n_3 \times \ldots \times n_{k-1}$  ways to select one element  $(a_1)$  out of  $n_1$  choices, a second element  $(a_2)$  out of  $n_2$  choices, and so on, up to the  $(k-1)^{th}$  element  $(a_{k-1})$  out of  $n_{k-1}$  choices. For each of these  $n_1 \times n_2 \times n_3 \times \ldots \times n_{k-1}$  possible ways, we can select the  $k^{th}$  element  $(a_k)$  out of  $n_k$  choices. Thus, the total # of ways =  $(n_1 \times n_2 \times n_3 \times \ldots \times n_{k-1}) \times n_k$ .

# <u>Definition 1.3.3:</u>

For positive integer n, we define **n factorial** as  $n! = n \times (n-1) \times (n-2) \times ... \times 2 \times 1 = n \times (n-1)!$  and 0! = 1.

#### Definition 1.3.4:

For nonnegative integers  $n \geq r$ , we define the **binomial coefficient** (read as n choose r) as

$$\binom{n}{r} = \frac{n!}{r!(n-r)!} = \frac{n \cdot (n-1) \cdot (n-2) \cdot \ldots \cdot (n-r+1)}{1 \cdot 2 \cdot 3 \cdot \ldots \cdot r}.$$

# Note:

A useful extension for the binomial coefficient for n < r is

$$\binom{n}{r} = \frac{n \cdot (n-1) \cdot \ldots \cdot 0 \cdot \ldots \cdot (n-r+1)}{1 \cdot 2 \cdot \ldots \cdot r} = 0.$$

# Note:

Most counting problems consist of drawing a fixed number of times from a set of elements (e.g.,  $\{1, 2, 3, 4, 5, 6\}$ ). To solve such problems, we need to know

- (i) the size of the set, n;
- (ii) the size of the sample, r;
- (iii) whether the result will be **ordered** (i.e., is  $\{1,2\}$  different from  $\{2,1\}$ ); and
- (iv) whether the draws are with replacement (i.e., can results like {1, 1} occur?).

# <u>Theorem 1.3.5:</u>

The number of ways to draw r elements from a set of n, if

- (i) ordered, without replacement, is  $\frac{n!}{(n-r)!}$ ;
- (ii) ordered, with replacement, is  $n^r$ ;
- (iii) unordered, without replacement, is  $\frac{n!}{r!(n-r)!} = \binom{n}{r}$ ;
- (iv) unordered, with replacement, is  $\frac{(n+r-1)!}{r!(n-1)!} = \binom{n+r-1}{r}$ .

#### Proof:

```
(i) n choices to select 1^{st} n-1 \text{ choices to select } 2^{nd} \vdots n-r+1 \text{ choices to select } r^{th} By Theorem 1.3.2, there are n \times (n-1) \times \ldots \times (n-r+1) = \frac{n \times (n-1) \times \ldots \times (n-r+1) \times (n-r)!}{(n-r)!} = \frac{n!}{(n-r)!} ways to do so.
```

# Corollary:

The number of permutations of n objects is n!.

```
(ii) n choices to select 1^{st}
n choices to select 2^{nd}
\vdots
n choices to select r^{th}

By Theorem 1.3.2, there are \underbrace{n \times n \times \ldots \times n}_{r \text{ times}} = n^r ways to do so.
```

- (iii) We know from (i) above that there are  $\frac{n!}{(n-r)!}$  ways to draw r elements out of n elements without replacement in the ordered case. However, for each unordered set of size r, there are r! related ordered sets that consist of the same elements. Thus, there are  $\frac{n!}{(n-r)!} \cdot \frac{1}{r!} = \binom{n}{r}$  ways to draw r elements out of n elements without replacement in the unordered case.
- (iv) There is no immediate direct way to show this part. We have to come up with some extra motivation. We assume that there are (n-1) walls that separate the n bins of possible outcomes and there are r markers. If we shake everything, there are (n-1+r)! permutations to arrange these (n-1) walls and r markers according to the Corollary. Since the r markers are indistinguishable and the (n-1) walls are also indistinguishable, we have to divide the number of permutations by r! to get rid of identical permutations where only the markers are changed and by (n-1)! to get rid of identical permutations where only the walls are changed. Thus, there are  $\frac{(n-1+r)!}{r!(n-1)!} = \binom{n+r-1}{r}$  ways to draw r elements out of n elements with replacement in the unordered case.

#### Theorem 1.3.6: The Binomial Theorem

If n is a non-negative integer, then

$$(1+x)^n = \sum_{r=0}^n \binom{n}{r} x^r$$

**Proof:** (By Induction)

#### Induction Base:

$$n = 0: \ 1 = (1+x)^0 = \sum_{r=0}^0 \binom{0}{r} x^r = \binom{0}{0} x^0 = 1$$
$$n = 1: \ (1+x)^1 = \sum_{r=0}^1 \binom{1}{r} x^r = \binom{1}{0} x^0 + \binom{1}{1} x^1 = 1 + x$$

# Induction Step:

Suppose it is true for k. We show that it is true for k+1.

$$(1+x)^{k+1} = (1+x)^k (1+x)$$

$$\stackrel{IB}{=} \left(\sum_{r=0}^k \binom{k}{r} x^r\right) (1+x)$$

$$= \sum_{r=0}^k \binom{k}{r} x^r + \sum_{r=0}^k \binom{k}{r} x^{r+1}$$

$$= \binom{k}{0} x^0 + \sum_{r=1}^k \left[\binom{k}{r} + \binom{k}{r-1}\right] x^r + \binom{k}{k} x^{k+1}$$

$$= \binom{k+1}{0} x^0 + \sum_{r=1}^k \left[\binom{k}{r} + \binom{k}{r-1}\right] x^r + \binom{k+1}{k+1} x^{k+1}$$

$$\stackrel{(*)}{=} \binom{k+1}{0} x^0 + \sum_{r=1}^k \binom{k+1}{r} x^r + \binom{k+1}{k+1} x^{k+1}$$

$$= \sum_{r=0}^{k+1} \binom{k+1}{r} x^r$$

(\*) Here we use Theorem 1.3.8 (i). Since the proof of Theorem 1.3.8 (i) only needs algebraic transformations without using the Binomial Theorem, part (i) of Theorem 1.3.8 can be applied here.

Corollary 1.3.7:

For a non-negative integer n, it holds:

Lecture 06: Mo 09/11/00

(i) 
$$\binom{n}{0} + \binom{n}{1} + \ldots + \binom{n}{n} = 2^n$$

(ii) 
$$\binom{n}{0} - \binom{n}{1} + \binom{n}{2} - \binom{n}{3} + \dots + (-1)^n \binom{n}{n} = 0$$

(iii) 
$$1 \cdot \binom{n}{1} + 2 \cdot \binom{n}{2} + 3 \cdot \binom{n}{3} + \ldots + n \cdot \binom{n}{n} = n2^{n-1}$$

(iv) 
$$1 \cdot \binom{n}{1} - 2 \cdot \binom{n}{2} + 3 \cdot \binom{n}{3} + \dots + (-1)^{n-1} n \cdot \binom{n}{n} = 0$$

Proof:

Use the Binomial Theorem:

(i) Let x = 1. Then

$$2^{n} = (1+1)^{n} \stackrel{Th.1.3.6}{=} \sum_{r=0}^{n} \binom{n}{r} 1^{r} = \sum_{r=0}^{n} \binom{n}{r}$$

(ii) Let x = -1. Then

$$0 = (1 + (-1))^n \stackrel{Th.1.3.6}{=} \sum_{r=0}^n \binom{n}{r} (-1)^r$$

(iii) 
$$\frac{d}{dx}(1+x)^n = \frac{d}{dx} \sum_{r=0}^n \binom{n}{r} x^r$$
$$\implies n(1+x)^{n-1} = \sum_{r=0}^n r \cdot \binom{n}{r} x^{r-1}$$

Substitute x = 1, then

$$n2^{n-1} = n(1+1)^{n-1} = \sum_{r=1}^{n} r \cdot \binom{n}{r}$$

(iv) Substitute x = -1 in (iii) above, then

$$0 = n(1 + (-1))^{n-1} = \sum_{r=1}^{n} r \cdot \binom{n}{r} (-1)^r = \sum_{r=1}^{n} r \cdot \binom{n}{r} (-1)^{r-1}$$

since for  $\sum a_i = 0$  also  $\sum (-a_i) = 0$ .

# Theorem 1.3.8:

For non-negative integers, n, m, r, it holds:

(i) 
$$\binom{n-1}{r} + \binom{n-1}{r-1} = \binom{n}{r}$$

(ii) 
$$\binom{n}{0}\binom{m}{r} + \binom{n}{1}\binom{m}{r-1} + \dots + \binom{n}{r}\binom{m}{0} = \binom{m+n}{r}$$

(iii) 
$$\binom{0}{r} + \binom{1}{r} + \binom{2}{r} + \ldots + \binom{n}{r} = \binom{n+1}{r+1}$$

# <u>Proof:</u>

Homework

# 1.4 Conditional Probability and Independence

So far, we have computed probability based only on the information that  $\Omega$  is used for a probability space  $(\Omega, L, P)$ . Suppose, instead, we know that event  $H \in L$  has happened. What statement should we then make about the chance of an event  $A \in L$ ?

# Definition 1.4.1:

Given  $(\Omega, L, P)$  and  $H \in L, P(H) > 0$ , and  $A \in L$ , we define

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

and call this the **conditional probability of** A **given** H.

#### Note:

$$P(A|H)$$
 is undefined if  $P(H) = 0$ .

# <u>Theorem 1.4.2:</u>

In the situation of Definition 1.4.1,  $(\Omega, L, P_H)$  is a probability space.

#### Proof:

If  $P_H$  is a probability measure, it must satisfy Def. 1.1.7.

(i) 
$$P(H)>0$$
 and by Def. 1.1.7 (i)  $P(A\cap H)\geq 0 \Longrightarrow P_H(A)=\frac{P(A\cap H)}{P(H)}\geq 0 \ \ \forall A\in L$ 

(ii) 
$$P_H(\Omega) = \frac{P(\Omega \cap H)}{P(H)} = \frac{P(H)}{P(H)} = 1$$

(iii) Let  $\{A_n\}_{n=1}^{\infty}$  be a sequence of disjoint sets. Then,

$$P_{H}(\bigcup_{n=1}^{\infty} A_{n}) \qquad \stackrel{Def.1.4.1}{=} \qquad \frac{P((\bigcup_{n=1}^{\infty} A_{n}) \cap H)}{P(H)}$$

$$= \qquad \frac{P(\bigcup_{n=1}^{\infty} (A_{n} \cap H))}{P(H)}$$

$$= \qquad \sum_{n=1}^{\infty} P(A_{n} \cap H)$$

$$= \qquad \sum_{n=1}^{\infty} P(H)$$

$$= \qquad \sum_{n=1}^{\infty} (\frac{P(A_{n} \cap H)}{P(H)})$$

$$\stackrel{Def.1.4.1}{=} \qquad \sum_{n=1}^{\infty} P_{H}(A_{n})$$

#### Note:

What we have done is to move to a new sample space  $\mathcal{H}$  and a new  $\sigma$ -field  $L_H = L \cap H$  of subsets  $A \cap H$  for  $A \in L$ . We thus have a new measurable space  $(\mathcal{H}, L_H)$  and a new probability space  $(\mathcal{H}, L_H, P_H)$ .

# Note:

From Definition 1.4.1, if  $A, B \in L, P(A) > 0$ , and P(B) > 0, then

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

which generalizes to the following Theorem.

# Theorem 1.4.3: Multiplication Rule

If  $A_1, \ldots, A_n \in L$  and  $P(\bigcap_{j=1}^{n-1} A_j) > 0$ , then

$$P(\bigcap_{j=1}^{n} A_j) = P(A_1) \cdot P(A_2|A_1) \cdot P(A_3|A_1 \cap A_2) \cdot \ldots \cdot P(A_n|\bigcap_{j=1}^{n-1} A_j).$$

# Proof:

Homework

# Definition 1.4.4:

A collection of subsets  $\{A_n\}_{n=1}^{\infty}$  of  $\Omega$  form a **partition** of  $\Omega$  if

(i) 
$$\bigcup_{n=1}^{\infty} A_n = \Omega$$
, and

(ii)  $A_i \cap A_j = \emptyset \ \forall i \neq j$ , i.e., elements are pairwise disjoint.

# Theorem 1.4.5: Law of Total Probability

If  $\{H_j\}_{j=1}^{\infty}$  is a partition of  $\Omega$ , and  $P(H_j) > 0 \ \forall j$ , then, for  $A \in L$ ,

$$P(A) = \sum_{j=1}^{\infty} P(A \cap H_j) = \sum_{j=1}^{\infty} P(H_j) P(A|H_j).$$

#### Proof:

By the Note preceding Theorem 1.4.3, the summands on both sides are equal  $\implies$  the right side of Th. 1.4.5 is true.

The left side proof:

 $H_j$  are disjoint  $\Rightarrow A \cap H_j$  are disjoint

$$A = A \cap \Omega \stackrel{Def_{1.4.4}}{=} A \cap (\bigcup_{j=1}^{\infty} H_j) = \bigcup_{j=1}^{\infty} (A \cap H_j)$$

$$\Longrightarrow P(A) = P(\bigcup_{j=1}^{\infty} (A \cap H_j)) \stackrel{Def_{1.1.7(iii)}}{=} \sum_{j=1}^{\infty} P(A \cap H_j)$$

# Theorem 1.4.6: Bayes' Rule

Let  $\{H_j\}_{j=1}^{\infty}$  be a partition of  $\Omega$ , and  $P(H_j) > 0 \ \forall j$ . Let  $A \in L$  and P(A) > 0. Then

$$P(H_j|A) = \frac{P(H_j)P(A|H_j)}{\sum_{n=1}^{\infty} P(H_n)P(A|H_n)} \quad \forall j.$$

Proofs

$$\overline{P(H_j \cap A)} \stackrel{Def_{1.4.1}}{=} P(A) \cdot P(H_j | A) = P(H_j) \cdot P(A | H_j)$$

$$\Longrightarrow P(H_j|A) = \frac{P(H_j) \cdot P(A|H_j)}{P(A)} \stackrel{Th.1.4.5}{=} \frac{P(H_j) \cdot P(A|H_j)}{\sum_{n=1}^{\infty} P(H_n) P(A|H_n)}.$$

# Definition 1.4.7:

For  $A, B \in L$ , A and B are **independent** iff  $P(A \cap B) = P(A)P(B)$ .

# Note:

- There are no restrictions on P(A) or P(B).
- If A and B are independent, then P(A|B) = P(A) (given that P(B) > 0) and P(B|A) = P(B) (given that P(A) > 0).
- If A and B are independent, then the following events are independent as well: A and  $B^C$ ;  $A^C$  and B;  $A^C$  and  $B^C$ .

# Definition 1.4.8:

Let  $\mathcal{A}$  be a collection of L-sets. The events of  $\mathcal{A}$  are **pairwise independent** iff for every distinct  $A_1, A_2 \in \mathcal{A}$  it holds  $P(A_1 \cap A_2) = P(A_1)P(A_2)$ .

#### Definition 1.4.9:

Let  $\mathcal{A}$  be a collection of L-sets. The events of  $\mathcal{A}$  are **mutually independent** (or completely independent) iff for every finite subcollection  $\{A_{i_1}, \ldots, A_{i_k}\}, A_{i_j} \in \mathcal{A}$ , it holds

$$P(\bigcap_{j=1}^{k} A_{i_j}) = \prod_{j=1}^{k} P(A_{i_j}).$$

#### Note:

To check for mutually independence of n events  $\{A_1, \ldots, A_n\} \in L$ , there are  $2^n - n - 1$  relations (i.e., all subcollections of size 2 or more) to check.

# Example 1.4.10:

Flip a fair coin twice.  $\Omega = \{HH, HT, TH, TT\}.$ 

 $A_1 = "H \text{ on 1st toss"}$ 

 $A_2 = "H \text{ on 2nd toss"}$ 

 $A_3 =$  "Exactly one H"

Obviously,  $P(A_1) = P(A_2) = P(A_3) = \frac{1}{2}$ .

Question: Are  $A_1, A_2$  and  $A_3$  pairwise independent and also mutually independent?

 $P(A_1 \cap A_2) = .25 = .5 \cdot .5 = P(A_1) \cdot P(A_2) \Rightarrow A_1, A_2 \text{ are independent.}$ 

 $P(A_1 \cap A_3) = .25 = .5 \cdot .5 = P(A_1) \cdot P(A_3) \Rightarrow A_1, A_3 \text{ are independent.}$ 

 $P(A_2 \cap A_3) = .25 = .5 \cdot .5 = P(A_2) \cdot P(A_3) \Rightarrow A_2, A_3 \text{ are independent.}$ 

Thus,  $A_1, A_2, A_3$  are pairwise independent.

 $P(A_1 \cap A_2 \cap A_3) = 0 \neq .5 \cdot .5 \cdot .5 = P(A_1) \cdot P(A_2) \cdot P(A_3) \Rightarrow A_1, A_2, A_3$  are not mutually independent.

# Example 1.4.11: (from Rohatgi, page 37, Example 5)

- r students. 365 possible birthdays for each student that are equally likely.
- One student at a time is asked for his/her birthday.
- If one of the other students hears this birthday and it matches his/her birthday, this other student has to raise his/her hand if at least one other student raises his/her hand, the procedure is over.

# • We are interested in

 $p_k = P(\text{procedure terminates at the } k \text{th student})$ = P(a hand is first risen when the k th student is asked for his/her birthday)

It is

 $p_1 = P(\text{at least 1 other from the } (r-1) \text{ students has a birthday on this particular day.})$  = 1 - P(all (r-1) students have a birthday on the remaining 364 out of 365 days)  $= 1 - \left(\frac{364}{365}\right)^{r-1}$ 

 $p_2 = P$ (no student has a birthday matching the first student and at least one of the other (r-2) students has a b-day matching the second student)

Let  $A \equiv \text{No}$  student has a b-day matching the 1<sup>st</sup> student Let  $B \equiv \text{At}$  least one of the other (r-2) has b-day matching 2<sup>nd</sup>

So 
$$p_2 = P(A \cap B)$$
  
=  $P(A) \cdot P(B|A)$   
=  $P(\text{no student has a matching b-day with the 1}^{\text{st}}$  student ) ×  
 $P(\text{at least one of the remaining students has a matching b-day with the second,}$   
given that no one matched the first.)  
=  $(1 - p_1)[1 - P(\text{all } (r - 2) \text{ students have a b-day on the remaining 363 out of 364 days})$   
=  $\left(\frac{364}{365}\right)^{r-1}\left[1 - \left(\frac{363}{364}\right)^{r-2}\right]$   
=  $\left(\frac{365 - 1}{365}\right)^{r-1}\left[1 - \left(\frac{363}{364}\right)^{r-2}\right]$  (\*)  
=  $\left(\frac{365}{365}\left(1 - \frac{1}{365}\right)^{r-1}\left[1 - \left(\frac{363}{364}\right)^{r-2}\right]$   
=  $\left(\frac{365}{365}P_{2-1}\right)\left(1 - \frac{2-1}{365}\right)^{r-2+1}\left[1 - \left(\frac{365-2}{365-2+1}\right)^{r-2}\right]$ 

Formally, we have to write this sequence of equalities in this order. However, it often might be easier to first work from both sides towards a particular result and combine partial results afterwards. Here, one might decide to stop at (\*) with the "forward" direction of the equalities

and first work "backwards" from the book, which makes things a lot simpler:

$$p_{2} = \left(\frac{\frac{365}{365}P_{2-1}}{(365)^{2-1}}\right) \left(1 - \frac{2-1}{365}\right)^{r-2+1} \left[1 - \left(\frac{365-2}{365-2+1}\right)^{r-2}\right]$$

$$= \frac{365}{365} \left(1 - \frac{1}{365}\right)^{r-1} \left[1 - \left(\frac{363}{364}\right)^{r-2}\right]$$

$$= \left(\frac{365-1}{365}\right)^{r-1} \left[1 - \left(\frac{363}{364}\right)^{r-2}\right]$$

We see that this is the same result as (\*).

Now let us consider  $p_3$ :

 $p_3 = P(\text{No one has same b-day as first and no one same as second, and at least one of the remaining <math>(r-3)$  has a matching b-day with the 3rd student)

Let  $A \equiv No$  one has the same b-day as the first student

Let  $B \equiv No$  one has the same b-day as the second student

Let  $C \equiv At$  least one of the other (r-3) has the same b-day as the third students

Now:

$$\begin{array}{ll} p_3 & = & P(A \cap B \cap C) \\ & = & P(A) \cdot P(B|A) \cdot P(C|A \cap B) \\ & = & \left(\frac{364}{365}\right)^{r-1} \left(\frac{363}{364}\right)^{r-2} \cdot \left[1 - P(\text{all }(r-3) \text{ students have a b-day on the remaining } 362 \text{ out of } 363 \text{ days}\right] \\ & = & \left(\frac{364}{365}\right)^{r-1} \left(\frac{363}{364}\right)^{r-2} \cdot \left[1 - \left(\frac{362}{363}\right)^{r-3}\right] \\ & = & \frac{(364)^{r-1}}{(365)^{r-1}} \cdot \frac{(363)^{r-2}}{(364)^{r-2}} \cdot \left[1 - \left(\frac{362}{363}\right)^{r-3}\right] \\ & = & \left(\frac{364^{r-1}}{364^{r-2}}\right) \left(\frac{363^{r-2}}{365^{r-1}}\right) \cdot \left[1 - \left(\frac{362}{363}\right)^{r-3}\right] \\ & = & \left(\frac{364}{365}\right) \left(\frac{363^{r-2}}{365^{r-2}}\right) \cdot \left[1 - \left(\frac{362}{363}\right)^{r-3}\right] \\ & = & \left(\frac{(365)(364)}{(365)^2}\right) \left(\frac{363}{365}\right)^{r-2} \left[1 - \left(\frac{362}{363}\right)^{r-3}\right] \end{array}$$

$$= \left(\frac{{}_{365}P_2}{(365)^2}\right) \left(1 - \frac{2}{365}\right)^{r-2} \left[1 - \left(\frac{362}{363}\right)^{r-3}\right]$$

$$= \left(\frac{{}_{365}P_{3-1}}{(365)^{3-1}}\right) \left(1 - \frac{3-1}{365}\right)^{r-3+1} \left[1 - \left(\frac{365-3}{365-3+1}\right)^{r-3}\right]$$

Once again, working "backwards" from the book should help to better understand these transformations.

For general  $p_k$  and restrictions on r and k see Homework.

Lecture 08: Fr 09/15/00

# 2 Random Variables

# 2.1 Measurable Functions

# Definition 2.1.1:

- A random variable (rv) is a set function from  $\Omega$  to  $\mathbb{R}$ .
- More formally: Let  $(\Omega, L, P)$  be any probability space. Suppose  $X : \Omega \to \mathbb{R}$  and that X is a measurable function, then we call X a random variable.
- More generally: If  $X : \Omega \to \mathbb{R}^k$ , we call X a **random vector**,  $\underline{X} = (X_1(\omega), X_2(\omega), \dots, X_k(\omega))$ .

What does it mean to say that a function is measurable?

#### Definition 2.1.2:

Suppose  $(\Omega, L)$  and (S, B) are two measurable spaces and  $X : \Omega \to S$  is a mapping from  $\Omega$  to S. We say that X is **measurable** L - B if  $X^{-1}(B) \in L$  for every set  $B \in B$ , where  $X^{-1}(B) = \{\omega \in \Omega : X(\omega) \in B\}$ .

# Example 2.1.3:

Record the opinion of 50 people: "yes" (y) or "no" (n).

 $\Omega = \{\text{All } 2^{50} \text{ possible sequences of y/n}\}$  — HUGE!

 $L = \mathcal{P}(\Omega)$ 

 $X: \Omega \to \mathcal{S} = \{\text{All } 2^{50} \text{ possible sequences of } 1 \text{ (=y) and } 0 \text{ (=n)} \}$ 

 $\mathcal{B} = \mathcal{P}(\mathcal{S})$ 

X is a random vector since each element in S has a corresponding element in  $\Omega$ , for  $B \in \mathcal{B}, X^{-1}(B) \in L = \mathcal{P}(\Omega)$ .

Consider  $X : \Omega \to \mathcal{S} = \{0, 1, 2, ..., 50\}$ , where  $X(\omega) = "\# \text{ of y's in } \omega" \text{ is a more manageable random variable.}$ 

A simple function, which takes only finite many values  $x_1, \ldots, x_k$  is measurable iff  $X^{-1}(x_i) \in L \ \forall x_i$ .

Here,  $X^{-1}(k) = \{\omega \in \Omega : \# \text{ 1's in sequence } \omega = k\}$  is a subset of  $\Omega$ , so it is in  $L = \mathcal{P}(\Omega)$ 

# Example 2.1.4:

Let  $\Omega$  = "infinite fair coin tossing space", i.e., infinite sequence of H's and T's.

Let  $L_n$  be a  $\sigma$ -field for the 1st n tosses.

Define 
$$L = \sigma(\bigcup_{n=0}^{\infty} L_n)$$
.

Let  $X_n: \Omega \to \mathbb{R}^{n=1}$  be  $X_n(\omega) =$  "proportion of H's in 1st n tosses".

For each  $n, X_n(\cdot)$  is simple (values  $\{0, \frac{1}{n}, \frac{2}{n}, \dots, n\}$ ) and  $X_n^{-1}(\frac{k}{n}) \in L_n \quad \forall k = 0, 1, \dots, n$ . Therefore,  $X_n^{-1}(\frac{k}{n}) \in L$ .

So every random variable  $X_n(\cdot)$  is measurable  $L-\mathcal{B}$ . Now we have a sequence of rv's  $\{X_n\}_{n=1}^{\infty}$ . We will show later that  $P(\{\omega: X_n(\omega) \to \frac{1}{2}\}) = 1$ , i.e., the *Strong Law of Large Numbers* (SLLN).

# Some Technical Points about Measurable Functions

# 2.1.5:

Suppose  $(\Omega, L)$  and (S, B) are measure spaces and that a collection of sets A generates B, i.e.,  $\sigma(A) = B$ . Let  $X : \Omega \to S$ . If  $X^{-1}(A) \in L$   $\forall A \in A$ , then X is measurable L - B.

This means we only have to check measurability on a basis collection  $\mathcal{A}$ . The usage is:  $\mathcal{B}$  on  $\mathbb{R}$  is generated by  $\{(-\infty, x] : x \in \mathbb{R}\}.$ 

#### 2.1.6:

If  $(\Omega, L), (\Omega', L')$ , and  $(\Omega'', L'')$  are measure spaces and  $X : \Omega \to \Omega'$  and  $Y : \Omega' \to \Omega''$  are measurable, then the composition  $(YX) : \Omega \to \Omega''$  is measurable L - L''.

# 2.1.7:

If  $f: \mathbb{R}^i \to \mathbb{R}^k$  is a continuous function, then f is measurable  $\mathcal{B}^i - \mathcal{B}^k$ .

#### 2.1.8:

If  $f_j: \Omega \to I\!\!R, j=1,\ldots k$  and  $g: I\!\!R^k \to I\!\!R$  are measurable, then  $g(f_1(\cdot),\ldots,f_k(\cdot))$  is measurable.

The usage is: g could be sum, average, difference, product, (finite) maximums and minimums of  $x_1, \ldots, x_k$ , etc.

# 2.1.9:

Limits: Extend the real line to  $[-\infty, \infty] = \mathbb{R} \cup \{-\infty, \infty\}$ . We say  $f: \Omega \to \mathbb{R}$  is measurable  $L - \mathcal{B}$  if

(i) 
$$f^{-1}(B) \in L \ \forall B \in \mathcal{B}$$
, and

(ii) 
$$f^{-1}(-\infty), f^{-1}(\infty) \in L$$
 also.

# <u>2.1.10:</u>

Suppose  $f_1, f_2, ...$  is a sequence of real-valued measurable functions  $(\Omega, L) \to (I\!\!R, \mathcal{B})$ . Then it holds:

(i)  $\sup_{n} f_n$ ,  $\inf_{n} f_n$ ,  $\lim_{n} \sup_{n} f_n$ ,  $\lim_{n} \inf_{n} f_n$ , are measurable.

(ii) If  $f = \lim_{n} f_n$  exists, then f is measurable.

(iii) The set  $\{\omega : f_n(\omega) \text{ converges}\} \in L$ .

(iv) If f is any measurable function, the set  $\{\omega: f_n(\omega) \to f(\omega)\} \in L$ .

# 2.2 Probability Distribution of a Random Variable

The definition of a random variable  $X:(\Omega,L)\to(\mathcal{S},\mathcal{B})$  makes no mention of P. We now introduce a probability measure on  $(\mathcal{S},\mathcal{B})$ .

# <u>Theorem 2.2.1:</u>

A random variable X on  $(\Omega, L, P)$  induces a probability measure on a space  $(\mathbb{R}, \mathcal{B}, Q)$  with the probability distribution Q of X defined by

$$Q(B) = P(X^{-1}(B)) = P(\{\omega : X(\omega) \in B\}) \quad \forall B \in \mathcal{B}.$$

# Note:

By the definition of a random variable,  $X^{-1}(B) \in L \ \forall B \in \mathcal{B}$ . Q is called **induced probability**.

#### Proof:

If X induces a probability measure Q on  $(I\!\!R,\mathcal{B})$ , then Q must satisfy the Kolmogorov Axioms of probability.

$$X:(\Omega,L)\to(S,\mathcal{B}).$$
 X is a rv  $\Rightarrow X^{-1}(B)=\{\omega:X(\omega)\in B\}=A\in L\ \forall B\in\mathcal{B}.$ 

(i) 
$$Q(B) = P(X^{-1}(B)) = P(\{\omega : X(\omega) \in B\}) = P(A) \stackrel{Def.1.1.7(i)}{\geq} 0 \ \forall B \in \mathcal{B}$$

(ii) 
$$Q(I\!\!R) = P(X^{-1}(I\!\!R)) \stackrel{X=rv}{=} P(\Omega) \stackrel{Def.1.1.7(ii)}{=} 1$$

(iii) Let  $\{B_n\}_{n=1}^{\infty} \in \mathcal{B}, B_i \cap B_j = \emptyset \ \forall i \neq j$ . Then,

$$Q(\bigcup_{n=1}^{\infty} B_n) = P(X^{-1}(\bigcup_{n=1}^{\infty} B_n)) \stackrel{(*)}{=} P(\bigcup_{n=1}^{\infty} (X^{-1}(B_n))) \stackrel{Def.1.1.7(iii)}{=} \sum_{n=1}^{\infty} P(X^{-1}(B_n)) = \sum_{n=1}^{\infty} Q(B_n)$$

(\*) holds since  $X^{-1}(\cdot)$  commutes with unions/intersections and preserves disjointedness.

# <u>Definition 2.2.2:</u>

A real–valued function F on  $(-\infty, \infty)$  that is non–decreasing, right–continuous, and satisfies

Lecture 09: Mo 09/18/00

$$F(-\infty) = 0, F(\infty) = 1$$

is called a **cumulative distribution function (cdf)** on  $\mathbb{R}$ .

# Note:

No mention of probability space or measure P in Definition 2.2.2 above.

#### Definition 2.2.3:

Let P be a probability measure on  $(IR, \mathcal{B})$ . The cdf associated with P is

$$F(x) = F_P(x) = P((-\infty, x]) = P(\{\omega : X(\omega) \le x\}) = P(X \le x)$$

for a random variable X defined on  $(\mathbb{R}, \mathcal{B}, P)$ .

#### Note:

 $F(\cdot)$  defined as in Definition 2.2.3 above indeed is a cdf.

# Proof (of Note):

(i) Let  $x_1 < x_2$ 

$$\implies (-\infty, x_1] \subset (-\infty, x_2]$$

$$\implies F(x_1) = P(\{\omega : X(\omega) \le x_1\}) \stackrel{Th.1.2.1(v)}{\le} P(\{\omega : X(\omega) \le x_2\}) = F(x_2)$$

Thus, since  $x_1 < x_2$  and  $F(x_1) \le F(x_2)$ , F(.) is non-decreasing.

(ii) Since F is non-decreasing, it is sufficient to show that F(.) is right-continuous if for any sequence of numbers  $x_n \to x+$  (which means that  $x_n$  is approaching x from the right) with  $x_1 > x_2 > \ldots > x_n > \ldots > x : F(x_n) \to F(x)$ .

Let  $A_n = \{\omega : X(\omega) \in (x, x_n]\} \in L$  and  $A_n \downarrow \emptyset$ . None of the intervals  $(x, x_n]$  contains x. As  $x_n \to x_n$ , the number of points  $\omega$  in  $A_n$  diminishes until the set is empty. Formally,

$$\lim_{n \to \infty} A_n = \lim_{n \to \infty} \bigcap_{i=1}^n A_i = \bigcap_{n=1}^\infty A_n = \emptyset.$$

By Theorem 1.2.6 it follows that

$$\lim_{n \to \infty} P(A_n) = P(\lim_{n \to \infty} A_n) = P(\emptyset) = 0.$$

It is

$$P(A_n) = P(\{\omega : X(\omega) \le x_n\}) - P(\{\omega : X(\omega) \le x\}) = F(x_n) - F(x).$$

$$\Longrightarrow (\lim_{n\to\infty} F(x_n)) - F(x) = \lim_{n\to\infty} (F(x_n) - F(x)) = \lim_{n\to\infty} P(A_n) = 0$$

$$\implies \lim_{n \to \infty} F(x_n) = F(x)$$

 $\Longrightarrow F(x)$  is right-continuous.

(iii) 
$$F(-n) \stackrel{Def.}{=} \stackrel{2.2.3}{=} P(\{\omega : X(\omega) \le -n\})$$
 $\Longrightarrow$ 

$$F(-\infty) = \lim_{n \to \infty} F(-n)$$

$$= \lim_{n \to \infty} P(\{\omega : X(\omega) \le -n\})$$

$$= P(\lim_{n \to \infty} \{\omega : X(\omega) \le -n\})$$

$$= P(\emptyset)$$

$$= 0$$

(iv) 
$$F(n) \stackrel{Def.=2.2.3}{=} P(\{\omega : X(\omega) \le n\})$$
 $\Longrightarrow$ 

$$F(\infty) = \lim_{n \to \infty} F(n)$$

$$= \lim_{n \to \infty} P(\{\omega : X(\omega) \le n\})$$

$$= P(\lim_{n \to \infty} \{\omega : X(\omega) \le n\})$$

$$= P(\Omega)$$

$$= 1$$

Note that (iii) and (iv) implicitly use Theorem 1.2.6. In (iii), we use  $A_n = (-\infty, -n)$  where  $A_n \supset A_{n+1}$  and  $A_n \downarrow \emptyset$ . In (iv), we use  $A_n = (-\infty, n)$  where  $A_n \subset A_{n+1}$  and  $A_n \uparrow \mathbb{R}$ .

# Definition 2.2.4:

If a random variable  $X: \Omega \to \mathbb{R}$  has induced a probability measure  $P_X$  on  $(\mathbb{R}, \mathcal{B})$  with cdf F(x), we say

- (i) rv X is **continuous** if F(x) is continuous in x.
- (ii) rv X is **discrete** if F(x) is a step function in x.

# Note:

There are rvs that are mixtures of continuous and discrete rvs. One such example is a truncated failure time distribution. We assume a continuous distribution (e.g., exponential) up to a given truncation point x and assign the "remaining" probability to the truncation point. Thus, a single point has a probability > 0 and F(x) jumps at the truncation point x.

# <u>Definition 2.2.5:</u>

Two random variables X and Y are identically distributed iff

$$P_X(X \in A) = P_Y(Y \in A) \ \forall A \in L.$$

# Note:

Def. 2.2.5 does not mean that  $X(\omega) = Y(\omega) \ \forall \omega \in \Omega$ . For example,

X = # H in 3 coin tosses

Y = # T in 3 coin tosses

X,Y are both Bin(3,0.5), i.e., identically distributed, but for  $\omega=(H,H,T), X(\omega)=2\neq 1=Y(\omega)$ , i.e.,  $X\neq Y$ .

# Theorem 2.2.6:

The following two statements are equivalent:

- (i) X, Y are identically distributed.
- (ii)  $F_X(x) = F_Y(x) \quad \forall x \in \mathbb{R}.$

# Proof:

 $(i) \Rightarrow (ii)$ :

$$F_X(x) = P_X((-\infty, x])$$

$$= P(\{\omega : X(\omega) \in (-\infty, x]\})$$

$$\stackrel{byDef.2.2.5}{=} P(\{\omega : Y(\omega) \in (-\infty, x]\})$$

$$= P_Y((-\infty, x])$$

$$= F_Y(X)$$

(ii)  $\Rightarrow$  (i):

Requires extra knowledge from measure theory.

Lecture 10: We 09/20/00

## 2.3 Discrete and Continuous Random Variables

We now extend Definition 2.2.4 to make our definitions a little bit more formal.

## Definition 2.3.1:

Let X be a real-valued random variable with cdf F on  $(\Omega, L, P)$ . X is **discrete** if there exists a countable set  $E \subset \mathbb{R}$  such that  $P(X \in E) = 1$ , i.e.,  $P(\{\omega : X(\omega) \in E\}) = 1$ . The points of E which have positive probability are the **jump points** of the step function F, i.e., the cdf of X.

Define 
$$p_i = P(\{\omega : X(\omega) = x_i, x_i \in E\}) = P_X(X = x_i) \ \forall i \ge 1.$$
 Then,  $p_i \ge 0, \sum_{i=1}^{\infty} p_i = 1.$ 

We call  $\{p_i : p_i \geq 0\}$  the **probability mass function (pmf)** (also: probability frequency function) of X.

# Note:

Given any set of numbers  $\{p_n\}_{n=1}^{\infty}, p_n \geq 0 \ \forall n \geq 1, \sum_{n=1}^{\infty} p_n = 1, \{p_n\}_{n=1}^{\infty}$  is the pmf of some rv X.

## Note:

The issue of continuous rv's and probability density functions (pdfs) is more complicated. A rv  $X: \Omega \to \mathbb{R}$  always has a cdf F. Whether there exists a function f such that f integrates to F and F' exists and equals f (almost everywhere) depends on something stronger than just continuity.

## Definition 2.3.2:

A real-valued function F is **continuous** in  $x_0 \in \mathbb{R}$  iff

$$\forall \epsilon > 0 \ \exists \delta > 0 \ \forall x : \ |x - x_0| < \delta \Rightarrow |F(x) - F(x_0)| < \epsilon.$$

F is continuous iff F is continuous in all  $x \in R$ .

## Definition 2.3.3:

A real-valued function F defined on [a, b] is **absolutely continuous** on [a, b] iff

 $\forall \epsilon > 0 \ \exists \delta > 0 \ \forall$  finite subcollection of disjoint subintervals  $[a_i, b_i], i = 1, \ldots, n$ :

$$\sum_{i=1}^{n} (b_i - a_i) < \delta \Rightarrow \sum_{i=1}^{n} |F(b_i) - F(a_i)| < \epsilon.$$

Note:

Absolute continuity implies continuity.

## <u>Theorem 2.3.4:</u>

- (i) If F is absolutely continuous, then F' exists almost everywhere.
- (ii) A function F is an indefinite integral iff it is absolutely continuous. Thus, every absolutely continuous function F is the indefinite integral of its derivative F'.

## Definition 2.3.5:

Let X be a random variable on  $(\Omega, L, P)$  with cdf F. We say X is a **continuous** rv iff F is absolutely continuous. In this case, there exists a non-negative integrable function f, the **probability density function (pdf)** of X, such that

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

From this it follows that, if  $a, b \in \mathbb{R}$ , a < b, then

$$P_X(a < X \le b) = F(b) - F(a) = \int_a^b f(t)dt$$

exists and is well defined.

## Theorem 2.3.6:

Let X be a continuous random variable with pdf f. Then it holds:

- (i) For every Borel set  $B \in \mathcal{B}, P(B) = \int_{B} f(t)dt$ .
- (ii) If F is absolutely continuous and f is continuous at x, then  $F'(x) = \frac{dF(x)}{dx} = f(x)$ .

## Proof:

Part (i): From Definition 2.3.5 above.

Part (ii): By Fundamental Theorem of Calculus.

## Note:

As already stated in the Note following Definition 2.2.4, not every rv will fall into one of these two (or if you prefer – three –, i.e., discrete, continuous/absolutely continuous) classes. However, most rv which arise in practice will. We look at one example that is unlikely to occur in practice in the next Homework assignment.

However, note that every  $\operatorname{cdf} F$  can be written as

$$F(x) = aF_d(x) + (1-a)F_c(x), \ 0 \le a \le 1,$$

where  $F_d$  is the cdf of a discrete rv and  $F_c$  is a continuous (but not necessarily absolute continuous) cdf.

Some authors, such as Marek Fisz Wahrscheinlichkeitsrechnung und mathematische Statistik, VEB Deutscher Verlag der Wissenschaften, Berlin, 1989, are even more specific. There it is stated that every cdf F can be written as

$$F(x) = a_1 F_d(x) + a_2 F_c(x) + a_3 F_s(x), \quad a_1, a_2, a_3 \ge 0, a_1 + a_2 + a_3 = 1.$$

Here,  $F_d(x)$  and  $F_c(x)$  are discrete and absolute continuous cdfs.  $F_s(x)$  is called a **singular** cdf. Singular means that  $F_s(x)$  is continuous and its derivative F'(x) equals 0 almost everywhere (i.e., everywhere but in those points that belong to a Borel-measurable set of probability 0).

Question: Does "continuous" but "not absolutely continuous" mean "singular"? — We will (hopefully) see later...

#### Example 2.3.7:

Consider

$$F(x) = \begin{cases} 0, & x < 0 \\ 1/2, & x = 0 \\ 1/2 + x/2, & 0 < x < 1 \\ 1, & x \ge 1 \end{cases}$$

We can write F(x) as  $aF_d(x) + (1-a)F_c(x), 0 \le a \le 1$ . How?

Since F(x) has only one jump at x = 0, it is reasonable to get started with a pmf  $p_0 = 1$  and corresponding cdf

$$F_d(x) = \begin{cases} 0, & x < 0 \\ 1, & x \ge 0 \end{cases}$$

Since F(x) = 0 for x < 0 and F(x) = 1 for  $x \ge 1$ , it must clearly hold that  $F_c(x) = 0$  for x < 0 and  $F_c(x) = 1$  for  $x \ge 1$ . In addition F(x) increases linearly in 0 < x < 1. A good guess would be a pdf  $f_c(x) = 1 \cdot I_{(0,1)}(x)$  and corresponding cdf

$$F_c(x) = \begin{cases} 0, & x \le 0 \\ x, & 0 < x < 1 \\ 1, & x \ge 1 \end{cases}$$

Knowing that F(0) = 1/2, we have at least to multiply  $F_d(x)$  by 1/2. And, indeed, F(x) can be written as

$$F(x) = \frac{1}{2}F_d(x) + \frac{1}{2}F_c(x).$$

# Definition 2.3.8:

The two-valued function  $I_A(x)$  is called **indicator function** and it is defined as follows:

$$I_A(x) = 1$$
 if  $x \in A$  and  $I_A(x) = 0$  if  $x \notin A$  for any set  $A$ .

# An Excursion into Logic

When proving theorems we only used direct methods so far. We used induction proofs to show that something holds for arbitrary n. To show that a statement A implies a statement B, i.e.,  $A \Rightarrow B$ , we used proofs of the type  $A \Rightarrow A_1 \Rightarrow A_2 \Rightarrow \ldots \Rightarrow A_{n-1} \Rightarrow A_n \Rightarrow B$  where one step directly follows from the previous step. However, there are different approaches to obtain the same result.

Implication: (A implies B)

 $A \Rightarrow B$  is equivalent to  $\neg B \Rightarrow \neg A$  is equivalent to  $\neg A \lor B$ :

A	B	$A \Rightarrow B$	$\neg A$	$\neg B$	$\neg B \Rightarrow \neg A$	$\neg A \lor B$
1	1	1	0	0	1	1
1	0	0	0	1	0	0
0	1	1	1	0	1	1
0	0	1	1	1	1	1

Lecture 11: Fr 09/22/00

Equivalence: (A is equivalent to B)

 $A \Leftrightarrow B$  is equivalent to  $(A \Rightarrow B) \land (B \Rightarrow A)$  is equivalent to  $(\neg A \lor B) \land (A \lor \neg B)$ :

A	B	$A \Leftrightarrow B$	$A \Rightarrow B$	$B \Rightarrow A$	$(A \Rightarrow B) \land (B \Rightarrow A)$	$\neg A \lor B$	$A \vee \neg B$	$(\neg A \lor B) \land (A \lor \neg B)$
1	1	1	1	1	1	1	1	1
1	0	0	0	1	0	0	1	0
0	1	0	1	0	0	1	0	0
0	0	1	1	1	1	1	1	1

## Negations of Quantifiers:

 $\neg \forall x \in X \ : \ B(x) \text{ is equivalent to } \exists x \in X \ : \ \neg B(x)$ 

 $\neg \exists x \in X : B(x) \text{ is equivalent to } \forall x \in X : \neg B(x)$ 

 $\exists x \in X \ \forall y \in Y : B(x, y) \text{ implies } \forall y \in Y \ \exists x \in X : B(x, y)$ 

## 2.4 Transformations of Random Variables

Let X be a real-valued random variable on  $(\Omega, L, P)$ , i.e.,  $X : (\Omega, L) \to (\mathbb{R}, \mathcal{B})$ . Let g be any Borel-measurable real-valued function on  $\mathbb{R}$ . Then, by statement 2.1.6, Y = g(X) is a random variable.

## <u>Theorem 2.4.1:</u>

Given a random rariable X with known induced distribution and a Borel-measurable function g, then the distribution of the random variable Y = g(X) is determined.

## Proof:

$$\begin{split} F_Y(y) &= P_Y(Y \leq y) \\ &= P(\{\omega: g(X(\omega)) \leq y\}) \\ &= P(\{\omega: X(\omega) \in B_y\}) \quad \text{where } B_y = g^{-1}((-\infty, y]) \in \mathcal{B} \text{ since g is Borel-measureable.} \\ &= P(X^{-1}(B_y)) \end{split}$$

#### Note:

From now on, we restrict ourselves to real-valued (vector-valued) functions that are Borel-measurable, i.e., measurable with respect to  $(\mathbb{R}, \mathcal{B})$  or  $(\mathbb{R}^k, \mathcal{B}^k)$ .

More generally, 
$$P_Y(Y \in C) = P_X(X \in g^{-1}(C)) \ \forall C \in \mathcal{B}.$$

#### Example 2.4.2:

Suppose X is a discrete random variable. Let A be a countable set such that  $P(X \in A) = 1$  and  $P(X = x) > 0 \ \forall x \in A$ .

Let Y = g(X). Obviously, the sample space of Y is also countable. Then,

$$P_Y(Y = y) = \sum_{x \in g^{-1}(\{y\})} P_X(X = x) = \sum_{\{x:g(x) = y\}} P_X(X = x) \quad \forall y \in g(A).$$

# Example 2.4.3:

 $X \sim U(-1,1)$  so the pdf of X is  $f_X(x) = 1/2I_{[-1,1]}(x)$ , which, according to Definition 2.3.8, reads as  $f_X(x) = 1/2$  for  $-1 \le x \le 1$  and 0 otherwise.

Let 
$$Y = X^+ = \begin{cases} x, & x \ge 0 \\ 0, & \text{otherwise} \end{cases}$$

Then,

$$F_Y(y) = P_Y(Y \le y) = \begin{cases} 0, & y < 0 \\ 1/2, & y = 0 \\ 1/2 + y/2, & 0 < y < 1 \\ 1, & y \ge 1 \end{cases}$$

This is the mixed discrete/continuous distribution from Example 2.3.7.

#### Note:

We need to put some conditions on g to ensure g(X) is continuous if X is continuous and avoid cases as in Example 2.4.3 above.

## Definition 2.4.4:

For a random variable X from  $(\Omega, L, P)$  to  $(I\!\!R, \mathcal{B})$ , the **support** of X (or P) is any set  $A \in L$  for which P(A) = 1. For a continuous random variable X with pdf f, we can think of the support of X as  $\mathcal{X} = X^{-1}(\{x : f_X(x) > 0\})$ .

## Definition 2.4.5:

Let f be a real-valued function defined on  $D \subseteq \mathbb{R}, D \in \mathcal{B}$ . We say:

f is (strictly) non-decreasing if  $x < y \Longrightarrow f(x)$  (<)  $\leq f(y) \ \forall x, y \in D$ 

f is (strictly) non-increasing if  $x < y \Longrightarrow f(x)$  (>)  $\geq f(y) \ \forall x, y \in D$ 

f is **monotonic** on D if f is either increasing or decreasing and write  $f \uparrow$  or  $f \downarrow$ .

## <u>Theorem 2.4.6:</u>

Lecture 12: Mo 09/25/00

Let X be a continuous rv with pdf  $f_X$  and support  $\mathcal{X}$ . Let y = g(x) be differentiable for all x and either (i) g'(x) > 0 or (ii) g'(x) < 0 for all x.

Then, Y = g(X) is also a continuous rv with pdf

$$f_Y(y) = f_X(g^{-1}(y)) \cdot |\frac{d}{dy}g^{-1}(y)| \cdot I_{g(X)}(y).$$

#### Proof:

Part (i):  $g'(x) > 0 \ \forall x \in \mathcal{X}$ 

So q is strictly increasing and continuous.

Therefore,  $x = g^{-1}(y)$  exists and it is also strictly increasing and also differentiable.

Then, from Rohatgi, page 9, Theorem 15:

$$\frac{d}{dy}g^{-1}(y) = \left(\frac{d}{dx}g(x)\mid_{x=g^{-1}(y)}\right)^{-1} > 0$$

We get  $F_Y(y) = P_Y(Y \le y) = P_Y(g(X) \le y) = P_X(X \le g^{-1}(y)) = F_X(g^{-1}(y))$  for  $y \in g(\mathcal{X})$  and, by differentiation,

$$f_Y(y) = F'_Y(y) = \frac{d}{dy} (F_X(g^{-1}(y))) \stackrel{\text{By Chain Rule}}{=} f_X(g^{-1}(y)) \cdot \frac{d}{dy} g^{-1}(y)$$

Part (ii):  $g'(x) < 0 \ \forall x \in \mathcal{X}$ 

So g is strictly decreasing and continuous.

Therefore,  $x = g^{-1}(y)$  exists and it is also strictly decreasing and also differentiable.

Then, from Rohatgi, page 9, Theorem 15:

$$\frac{d}{dy}g^{-1}(y) = \left(\frac{d}{dx}g(x)\mid_{x=g^{-1}(y)}\right)^{-1} < 0$$

We get  $F_Y(y) = P_Y(Y \le y) = P_Y(g(X) \le y) = P_X(X \ge g^{-1}(y)) = 1 - P_X(X \le g^{-1}(y)) = 1 - F_X(g^{-1}(y))$  for  $y \in g(\mathcal{X})$  and, by differentiation,

$$f_Y(y) = F_Y'(y) = \frac{d}{dy} (1 - F_X(g^{-1}(y))) \stackrel{\text{By Chain Rule}}{=} -f_X(g^{-1}(y)) \cdot \frac{d}{dy} g^{-1}(y) = f_X(g^{-1}(y)) \cdot \left( -\frac{d}{dy} g^{-1}(y) \right)$$

Since  $\frac{d}{dy}g^{-1}(y) < 0$ , the negative sign will cancel out, always giving us a positive value. Hence the need for the absolute value signs.

Combining parts (i) and (ii), we can therefore write

$$f_Y(y) = f_X(g^{-1}(y)) \cdot |\frac{d}{dy}g^{-1}(y)| \cdot I_{g(X)}(y).$$

Note:

In Theorem 2.4.6, we can also write

$$f_Y(y) = \frac{f(x)}{\left|\frac{dg(x)}{dx}\right|}\Big|_{x=g^{-1}(y)}, y \in g(\mathcal{X})$$

If g is monotonic over disjoint intervals, we can also get an expression for the pdf/cdf of Y = g(X) as stated in the following Theorem.

## Theorem 2.4.7:

Let Y = g(X) where X is a rv with pdf  $f_X(x)$  on support  $\mathcal{X}$ . Suppose there exists a partition  $A_0, A_1, \ldots, A_k$  of  $\mathcal{X}$  such that  $P(X \in A_0) = 0$  and  $f_X(x)$  is continuous on each  $A_i$ . Suppose there exist functions  $g_1(x), \ldots, g_k(x)$  defined on  $A_1$  through  $A_k$ , respectively, satisfying

- (i)  $g(x) = g_i(x) \quad \forall x \in A_i$ ,
- (ii)  $g_i(x)$  is monotonic on  $A_i$ ,
- (iii) the set  $\mathcal{Y} = g_i(A_i) = \{y : y = g_i(x) \text{ for some } x \in A_i\}$  is the same for each  $i = 1, \ldots, k$ , and
- (iv)  $g_i^{-1}(y)$  has a continuous derivative on  $\mathcal{Y}$  for each  $i=1,\ldots,k$ .

Then,

$$f_Y(y) = \sum_{i=1}^k f_X(g_i^{-1}(y)) \cdot |\frac{d}{dy}g_i^{-1}(y)| \cdot I_{\mathcal{Y}}(y)$$

Note:

Rohatgi, page 73, Theorem 4, removes condition (iii) by defining n = n(y) and  $x_1(y), \ldots, x_n(y)$ .

Example 2.4.8:

Let X be a rv with pdf  $f_X(x) = \frac{2x}{\pi^2} \cdot I_{(0,\pi)}(x)$ .

Let  $Y = \sin(X)$ . What is  $f_Y(y)$ ?

Since sin is not monotonic on  $(0, \pi)$ , Theorem 2.4.6 cannot be used to determine the pdf of Y.

Two possible approaches:

## Method 1: cdfs

For 0 < y < 1 we have

$$F_Y(y) = P_Y(Y \le y)$$

$$= P_X(\sin X \le y)$$

$$= P_X([0 \le X \le \sin^{-1}(y)] \text{ or } [\pi - \sin^{-1}(y) \le X \le \pi])$$

$$= F_X(\sin^{-1}(y)) + (1 - F_X(\pi - \sin^{-1}(y)))$$

since  $[0 \le X \le \sin^{-1}(y)]$  and  $[\pi - \sin^{-1}(y) \le X \le \pi]$  are disjoint sets. Then,

$$f_Y(y) = F_Y'(y)$$

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

$$= \frac{1}{\sqrt{1 - y^2}} \left( f_X(\sin^{-1}(y)) + f_X(\pi - \sin^{-1}(y)) \right)$$

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

$$= \frac{1}{\pi^2 \sqrt{1 - y^2}} 2\pi$$

$$= \frac{2}{\pi \sqrt{1 - y^2}} \cdot I_{(0,1)}(y)$$

# Method 2: Use of Theorem 2.4.7

Let  $A_1 = (0, \frac{\pi}{2}), A_2 = (\frac{\pi}{2}, \pi), \text{ and } A_0 = {\frac{\pi}{2}}.$ 

Let  $g_1^{-1}(y) = \sin^{-1}(y)$  and  $g_2^{-1}(y) = \pi - \sin^{-1}(y)$ .

It is 
$$\frac{d}{dy}g_1^{-1}(y) = \frac{1}{\sqrt{1-y^2}} = -\frac{d}{dy}g_2^{-1}(y)$$
 and  $\mathcal{Y} = (0,1)$ .

Thus, by use of Theorem 2.4.7, we get

$$f_{Y}(y) = \sum_{i=1}^{2} f_{X}(g_{i}^{-1}(y)) \cdot \left| \frac{d}{dy} g_{i}^{-1}(y) \right| \cdot I_{Y}(y)$$

$$= \frac{2 \sin^{-1}(y)}{\pi^{2}} \frac{1}{\sqrt{1 - y^{2}}} \cdot I_{(0,1)}(y) + \frac{2(\pi - \sin^{-1}(y))}{\pi^{2}} \frac{1}{\sqrt{1 - y^{2}}} \cdot I_{(0,1)}(y)$$

$$= \frac{2\pi}{\pi^{2}} \frac{1}{\sqrt{1 - y^{2}}} \cdot I_{(0,1)}(y)$$

$$= \frac{2}{\pi \sqrt{1 - y^{2}}} \cdot I_{(0,1)}(y)$$

Obviously, both results are identical.

Lecture 13: We 09/27/00

# <u>Theorem 2.4.9:</u>

Let X be a rv with a continuous cdf  $F_X(x)$  and let  $Y = F_X(X)$ . Then,  $Y \sim U(0,1)$ .

## Proof:

We have to consider two possible cases:

- (a)  $F_X$  is strictly increasing, i.e.,  $F_X(x_1) < F_X(x_2)$  for  $x_1 < x_2$ , and
- (b)  $F_X$  is non-decreasing, i.e., there exists  $x_1 < x_2$  and  $F_X(x_1) = F_X(x_2)$ . Assume that  $x_1$  is the infimum and  $x_2$  the supremum of those values for which  $F_X(x_1) = F_X(x_2)$  holds.

In (a),  $F_X^{-1}(y)$  is uniquely defined. In (b), we define  $F_X^{-1}(y) = \inf\{x : F_X(x) \ge y\}$ 

Without loss of generality:

$$F_X^{-1}(1) = +\infty \text{ if } F_X(x) < 1 \ \forall x \in \mathbb{R} \text{ and } F_X^{-1}(0) = -\infty \text{ if } F_X(x) > 0 \ \forall x \in \mathbb{R}.$$

For  $Y = F_X(X)$  and 0 < y < 1, we have

$$P(Y \le y) = P(F_X(X) \le y)$$

$$F_X^{-1} \uparrow P(F_X^{-1}(F_X(X)) \le F_X^{-1}(y))$$

$$\stackrel{(*)}{=} P(X \le F_X^{-1}(y))$$

$$= F_X(F_X^{-1}(y))$$

$$= y$$

At the endpoints, we have  $P(Y \le y) = 1$  if  $y \ge 1$  and  $P(Y \le y) = 0$  if  $y \le 0$ .

But why is (\*) true? — In (a), if  $F_X$  is strictly increasing and continuous, it is certainly  $x = F_X^{-1}(F_X(x))$ .

In (b), if  $F_X(x_1) = F_X(x_2)$  for  $x_1 < x < x_2$ , it may be that  $F_X^{-1}(F_X(x)) \neq x$ . But by definition,  $F_X^{-1}(F_X(x)) = x_1 \quad \forall x \in [x_1, x_2]$ . (\*) holds since on  $[x_1, x_2]$ , it is  $P(X \leq x) = P(X \leq x_1) \quad \forall x \in [x_1, x_2]$ . The flat cdf denotes  $F_X(x_2) - F_X(x_1) = P(x_1 < X \leq x_2) = 0$  by definition.

## Note:

This proof also holds if there exist multiple intervals with  $x_i < x_j$  and  $F_X(x_i) = F_X(x_j)$ , i.e., if the support of X is split in more than just 2 disjoint intervals.

# 3 Moments and Generating Functions

# 3.1 Expectation

## Definition 3.1.1:

Let X be a real-valued rv with cdf  $F_X$  and pdf  $f_X$  if X is continuous (or pmf  $f_X$  and support  $\mathcal{X}$  if X is discrete). The **expected value** (**mean**) of a measurable function  $g(\cdot)$  of X is

$$E(g(X)) = \begin{cases} \int_{-\infty}^{\infty} g(x) f_X(x) dx, & \text{if } X \text{ is continuous} \\ \sum_{x \in \mathcal{X}} g(x) f_X(x), & \text{if } X \text{ is discrete} \end{cases}$$

if  $E(|g(X)|) < \infty$ ; otherwise E(g(X)) is undefined, i.e., it does not exist.

## Example:

$$X \sim \text{Cauchy}, f_X(x) = \frac{1}{\pi(1+x^2)}, -\infty < x < \infty$$
:

$$E(|X|) = \frac{2}{\pi} \int_0^\infty \frac{x}{1+x^2} dx = \frac{1}{\pi} [\log(1+x^2)]_0^\infty = \infty$$

So, E(X) does not exist for the Cauchy distribution.

## Theorem 3.1.2:

If E(X) exists and a and b are finite constants, then E(aX + b) exists and equals aE(X) + b.

#### Proof:

Continuous case only:

# Existence:

$$E(|aX + b|) = \int_{-\infty}^{\infty} |ax + b| f_X(x) dx$$

$$\leq \int_{-\infty}^{\infty} (|a| \cdot |x| + |b|) f_X(x) dx$$

$$= |a| \int_{-\infty}^{\infty} |x| f_X(x) dx + |b| \int_{-\infty}^{\infty} f_X(x) dx$$

$$= |a| E(|X|) + |b|$$

$$< \infty$$

## Numerical Result:

$$E(aX + b) = \int_{-\infty}^{\infty} (ax + b) f_X(x) dx$$
$$= a \int_{-\infty}^{\infty} x f_X(x) dx + b \int_{-\infty}^{\infty} f_X(x) dx$$
$$= aE(X) + b$$

# Theorem 3.1.3:

Lecture 14: Fr 09/29/00

If X is bounded (i.e., there exists a M,  $0 < M < \infty$ , such that P(|X| < M) = 1), then E(X) exists.

## Definition 3.1.4:

The  $k^{th}$  moment of X, if it exists, is  $m_k = E(X^k)$ .

The  $k^{th}$  central moment of X, if it exists, is  $\mu_k = E((X - E(X))^k)$ .

# <u>Definition 3.1.5:</u>

The **variance** of X, if it exists, is the second central moment of X, i.e.,

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

## Theorem 3.1.6:

$$Var(X) = E(X^2) - (E(X))^2.$$

# Proof:

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

$$= E(X^{2} - 2XE(X) + (E(X))^{2})$$

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

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

## <u>Theorem 3.1.7:</u>

If Var(X) exists and a and b are finite constants, then Var(aX + b) exists and equals  $a^2Var(X)$ .

# Proof:

# Existence & Numerical Result:

$$Var(aX+b) = E\left(((aX+b)-E(aX+b))^2\right) \text{ exists if } E\left(|\ ((aX+b)-E(aX+b))^2\ |\right) \text{ exists.}$$

It holds that

$$E\left(|\ ((aX+b)-E(aX+b))^2\ |\right)$$

$$= E\left(((aX+b)-E(aX+b))^2\right)$$

$$= Var(aX+b)$$

$$\stackrel{Th.3.1.6}{=} E((aX+b)^2) - (E(aX+b))^2$$

$$\stackrel{Th.3.1.2}{=} E(a^2X^2 + 2abX + b^2) - (aE(X)+b)^2$$

$$\stackrel{Th.3.1.2}{=} a^2E(X^2) + 2abE(X) + b^2 - a^2(E(X))^2 - 2abE(X) - b^2$$

$$= a^2(E(X^2) - (E(X))^2)$$

$$\stackrel{Th.3.1.6}{=} a^2Var(X)$$

$$< \infty \text{ since } Var(X) \text{ exists}$$

## Theorem 3.1.8:

If the  $t^{th}$  moment of a rv X exists, then all moments of order 0 < s < t exist.

#### Proof:

Continuous case only:

$$E(|X|^{s}) = \int_{|x| \le 1} |x|^{s} f_{X}(x) dx + \int_{|x| > 1} |x|^{s} f_{X}(x) dx$$

$$\leq \int_{|x| \le 1} 1 \cdot f_{X}(x) dx + \int_{|x| > 1} |x|^{t} f_{X}(x) dx$$

$$\leq P(|X| \le 1) + E(|X|^{t})$$

$$< \infty$$

## <u>Theorem 3.1.9:</u>

If the  $t^{th}$  moment of a rv X exists, then

$$\lim_{n \to \infty} n^t P(\mid X \mid > n) = 0.$$

#### Proof:

Continuous case only:

$$\infty > \int_{\mathbb{R}} |x|^t f_X(x) dx = \lim_{n \to \infty} \int_{|x| \le n} |x|^t f_X(x) dx$$

$$\implies \lim_{n \to \infty} \int_{|x| > n} |x|^t f_X(x) dx = 0$$
But, 
$$\lim_{n \to \infty} \int_{|x| > n} |x|^t f_X(x) dx \ge \lim_{n \to \infty} n^t \int_{|x| > n} f_X(x) dx = \lim_{n \to \infty} n^t P(|X| > n) = 0$$

Lecture 16: We 10/04/00

# Note:

The inverse is not necessarily true, i.e., if  $\lim_{n\to\infty} n^t P(\mid X\mid > n) = 0$ , then the  $t^{th}$  moment of a rv X does not necessarily exist. We can only approach t up to some  $\delta > 0$  as the following Theorem 3.1.10 indicates.

## Theorem 3.1.10:

Let X be a rv with a distribution such that  $\lim_{n\to\infty} n^t P(\mid X\mid >n)=0$  for some t>0. Then,

$$E(\mid X\mid^s) < \infty \ \forall \ 0 < s < t.$$

## Note:

To prove this Theorem, we need Lemma 3.1.11 and Corollary 3.1.12.

## <u>Lemma 3.1.11:</u>

Let X be a non-negative rv with cdf F. Then,

$$E(X) = \int_0^\infty (1 - F_X(x)) dx$$

(if either side exists).

#### <u>Proof:</u>

Continuous case only:

To prove that the left side implies that the right side is finite and both sides are identical, we assume that E(X) exists. It is

$$E(X) = \int_0^\infty x f_X(x) dx = \lim_{n \to \infty} \int_0^n x f_X(x) dx$$

Replace the expression for the right side integral using integration by parts.

Let u = x and  $dv = f_X(x)dx$ , then

$$\int_0^n x f_X(x) dx = (xF(x)) \mid_0^n - \int_0^n F_X(x) dx$$

$$= nF_X(n) - 0F_X(0) - \int_0^n F_X(x) dx$$

$$= nF_X(n) - n + n - \int_0^n F_X(x) dx$$

$$= nF_X(n) - n + \int_0^n [1 - F_X(x)] dx$$

$$= n[F_X(n) - 1] + \int_0^n [1 - F_X(x)] dx$$

$$= -n[1 - F_X(n)] + \int_0^n [1 - F_X(x)] dx$$

$$= -nP(X > n) + \int_0^n [1 - F_X(x)] dx$$

$$\stackrel{X \ge 0}{=} -n^1 P(|X| > n) + \int_0^n [1 - F_X(x)] dx$$

$$\implies E(X^1) = \lim_{n \to \infty} [-n^1 P(|X| > n) + \int_0^n [1 - F_X(x)] dx$$

$$\stackrel{Th. 3.1.9}{=} 0 + \lim_{n \to \infty} \int_0^n [1 - F_X(x)] dx$$

$$= \int_0^\infty [1 - F_X(x)] dx$$

Thus, the existence of E(X) implies that  $\int_0^\infty [1 - F_X(x)] dx$  is finite and that both sides are identical.

We still have to show the converse implication:

If  $\int_0^\infty [1 - F_X(x)] dx$  is finite, then E(X) exists, i.e.,  $E(|X|) = E(X) < \infty$ , and both sides are identical. It is

$$\int_0^n x f_X(x) dx \stackrel{X \ge 0}{=} \int_0^n |x| f_X(x) dx = -n[1 - F_X(n)] + \int_0^n [1 - F_X(x)] dx$$

as seen above.

Since  $-n[1 - F_X(n)] \le 0$ , we get

$$\int_{0}^{n} |x| f_{X}(x) dx \le \int_{0}^{n} [1 - F_{X}(x)] dx \le \int_{0}^{\infty} [1 - F_{X}(x)] dx < \infty \quad \forall n$$

Thus,

$$\lim_{n \to \infty} \int_0^n |x| f_X(x) = \int_0^\infty |x| f_X(x) dx \le \int_0^\infty [1 - F_X(x)] dx < \infty$$

 $\Longrightarrow E(X)$  exists and is identical to  $\int_0^\infty [1 - F_X(x)] dx$  as seen above.

Corollary 3.1.12:

$$E(|X|^s) = s \int_0^\infty y^{s-1} P(|X| > y) dy$$

Proof:

$$E(\mid X\mid^{s}) \stackrel{Lemma}{=} {}^{3.1.11} \int_{0}^{\infty} [1 - F_{|X|^{s}}(z)] dz = \int_{0}^{\infty} P(\mid X\mid^{s} > z) dz$$

Let  $z=y^s$ . Then  $\frac{dz}{dy}=sy^{s-1}$  and  $dz=sy^{s-1}dy$ . Therefore,

$$\int_0^\infty P(\mid X\mid^s > z)dz = \int_0^\infty P(\mid X\mid^s > y^s)sy^{s-1}dy$$

$$= s \int_0^\infty y^{s-1}P(\mid X\mid^s > y^s)dy$$

$$\stackrel{monotonic \uparrow}{=} s \int_0^\infty y^{s-1}P(\mid X\mid > y)dy$$

# Proof (of Theorem 3.1.10):

For any given  $\epsilon > 0$ , choose N such that the tail probability  $P(|X| > n) < \frac{\epsilon}{n^t} \quad \forall n \geq N$ .  $E(|X|^s) \stackrel{Cor.}{=} \frac{3 \cdot 1 \cdot 12}{s} s \int_0^\infty y^{s-1} P(|X| > y) dy$ 

$$= s \int_0^N y^{s-1} P(|X| > y) dy + s \int_N^\infty y^{s-1} P(|X| > y) dy$$

$$\leq \int_0^N s y^{s-1} \cdot 1 \ dy + s \int_N^\infty y^{s-1} \frac{\epsilon}{y^t} dy$$

$$= y^s \mid_0^N + s\epsilon \int_N^\infty y^{s-1} \frac{1}{y^t} dy$$

$$= N^s + s\epsilon \int_N^\infty y^{s-1-t} dy$$

It is

$$\int_{N}^{\infty} y^{c} dy = \begin{cases} \frac{1}{c+1} y^{c+1} |_{N}^{\infty}, & c \neq -1 \\ \ln y |_{N}^{\infty}, & c = -1 \end{cases}$$
$$= \begin{cases} \infty, & c \geq -1 \\ -\frac{1}{c+1} N^{c+1} < \infty, & c < -1 \end{cases}$$

Thus, for  $E(\mid X\mid^s) < \infty$ , it must hold that s-1-t<-1, or equivalently, s< t. So  $E(\mid X\mid^s) < \infty$ , i.e., it exists, for every s with 0 < s < t for a rv X with a distribution such that  $\lim_{n\to\infty} n^t P(\mid X\mid > n) = 0$  for some t>0.

Lecture 17: Fr 10/06/00

<u>Theorem 3.1.13:</u>

Let X be a rv such that

$$\lim_{k \to \infty} \frac{P(\mid X \mid > \alpha k)}{P(\mid X \mid > k)} = 0 \quad \forall \alpha > 1.$$

Then, all moments of X exist.

Proof:

• For  $\epsilon > 0$ , we select some  $k_0$  such that

$$\frac{P(\mid X\mid >\alpha k)}{P(\mid X\mid >k)}<\epsilon \ \forall k\geq k_0.$$

- Select  $k_1$  such that  $P(\mid X \mid > k) < \epsilon \ \forall k \geq k_1$ .
- Select  $N = \max(k_0, k_1)$ .
- If we have some fixed positive integer r:

$$\frac{P(|X| > \alpha^{r}k)}{P(|X| > k)} = \frac{P(|X| > \alpha k)}{P(|X| > k)} \cdot \frac{P(|X| > \alpha^{2}k)}{P(|X| > \alpha k)} \cdot \frac{P(|X| > \alpha^{3}k)}{P(|X| > \alpha^{2}k)} \cdot \dots \cdot \frac{P(|X| > \alpha^{r}k)}{P(|X| > \alpha^{r-1}k)}$$

$$= \frac{P(|X| > \alpha k)}{P(|X| > k)} \cdot \frac{P(|X| > \alpha \cdot (\alpha k))}{P(|X| > 1 \cdot (\alpha k))} \cdot \frac{P(|X| > \alpha \cdot (\alpha^{2}k))}{P(|X| > 1 \cdot (\alpha^{2}k))} \cdot \dots \cdot \frac{P(|X| > \alpha \cdot (\alpha^{r-1}k))}{P(|X| > 1 \cdot (\alpha^{r-1}k))}$$

- Note: Each of these r terms on the right side is  $<\epsilon$  by our original statement of selecting some  $k_0$  such that  $\frac{P(|X|>\alpha k)}{P(|X|>k)}<\epsilon$   $\forall k\geq k_0$  and since  $\alpha>1$  and therefore  $\alpha^n k\geq k_0$ .
- Now we get for our entire expression that  $\frac{P(|X|>\alpha^r k)}{P(|X|>k)} \le \epsilon^r$  for  $k \ge N$  (since in this case also  $k \ge k_0$ ) and  $\alpha > 1$ .
- Overall, we have  $P(|X| > \alpha^r k) \le \epsilon^r P(|X| > k) \le \epsilon^{r+1}$  for  $k \ge N$  (since in this case also  $k \ge k_1$ ).
- For a fixed positive integer n:

$$\mathrm{E}(\mid X\mid^{n}) \stackrel{Cor.3.1.12}{=} n \cdot \int\limits_{0}^{\infty} x^{n-1} \mathrm{P}(\mid X\mid > x) dx \ = \ n \int\limits_{0}^{N} x^{n-1} \mathrm{P}(\mid X\mid > x) dx + n \int\limits_{N}^{\infty} x^{n-1} \mathrm{P}(\mid X\mid > x) dx$$

• We know that:

$$n\int_{0}^{N} x^{n-1} P(|X| > x) dx \le \int_{0}^{N} n x^{n-1} dx = x^{n} \mid_{0}^{N} = N^{n} < \infty$$

but is

$$n\int_{N}^{\infty} x^{n-1} P(\mid X\mid > x) dx < \infty ?$$

• To check the second part, we use:

$$\int_{N}^{\infty} x^{n-1} P(|X| > x) dx = \sum_{r=1}^{\infty} \int_{\alpha^{r-1} N}^{\alpha^{r} N} x^{n-1} P(|X| > x) dx$$

• We know that:

$$\int_{\alpha^{r-1}N}^{\alpha^{r}N} x^{n-1} P(\mid X \mid > x) dx \leq \epsilon^{r} \int_{\alpha^{r-1}N}^{\alpha^{r}N} x^{n-1} dx$$

This step is possible since  $\epsilon^r \geq P(\mid X \mid \geq \alpha^{r-1}N) \geq P(\mid X \mid > x) \geq P(\mid X \mid \geq \alpha^r N)$  $\forall x \in (\alpha^{r-1}N, \alpha^r N) \text{ and } N = \max(k_0, k_1).$ 

• Since  $(\alpha^{r-1}N)^{n-1} \le x^{n-1} \le (\alpha^r N)^{n-1} \ \forall x \in (\alpha^{r-1}N, \alpha^r N)$ , we get:

$$\epsilon^r \int\limits_{\alpha^{r-1}N}^{\alpha^r N} x^{n-1} dx \;\; \leq \;\; \epsilon^r (\alpha^r N)^{n-1} \int\limits_{\alpha^{r-1}N}^{\alpha^r N} 1 dx \;\; \leq \;\; \epsilon^r (\alpha^r N)^{n-1} (\alpha^r N) \;\; = \;\; \epsilon^r (\alpha^r N)^n$$

• Now we go back to our original inequality:

$$\int_{N}^{\infty} x^{n-1} P(|X| > x) dx \leq \sum_{r=1}^{\infty} \epsilon^{r} \int_{\alpha^{r-1}N}^{\alpha^{r}N} x^{n-1} dx \leq \sum_{r=1}^{\infty} \epsilon^{r} (\alpha^{r}N)^{n} = N^{n} \sum_{r=1}^{\infty} (\epsilon \cdot \alpha^{n})^{r}$$

$$= \frac{N^{n} \epsilon \alpha^{n}}{1 - \epsilon \alpha^{n}} \text{ if } \epsilon \alpha^{n} < 1 \text{ or, equivalently, if } \epsilon < \frac{1}{\alpha^{n}}$$

 • Since  $\frac{N^n \epsilon \alpha^n}{1 - \epsilon \alpha^n}$  is finite, all moments  $E(\mid X \mid^n)$  exist.

Lecture 18: Mo 10/09/00

# 3.2 Generating Functions

## Definition 3.2.1:

Let X be a rv with cdf  $F_X$ . The **moment generating function** (**mgf**) of X is defined as

$$M_X(t) = E(e^{tX})$$

provided that this expectation exists in an (open) interval around 0, i.e., for -h < t < h for some h > 0.

# Theorem 3.2.2:

If a rv X has a mgf  $M_X(t)$  that exists for -h < t < h for some h > 0, then

$$E(X^n) = M_X^{(n)}(0) = \frac{d^n}{dt^n} M_X(t) \mid_{t=0}$$
.

## Proof:

We assume that we can differentiate under the integral sign. If, and when, this really is true will be discussed later in this section.

$$\frac{d}{dt}M_X(t) = \frac{d}{dt} \int_{-\infty}^{\infty} e^{tx} f_X(x) dx$$

$$= \int_{-\infty}^{\infty} (\frac{\partial}{\partial t} e^{tx} f_X(x)) dx$$

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

$$= E(Xe^{tX})$$

Evaluating this at t = 0, we get:  $\frac{d}{dt}M_X(t)|_{t=0} = E(X)$ 

By induction, we get for  $n \geq 2$ :

$$\frac{d^n}{dt^n} M_X(t) = \frac{d}{dt} \left( \frac{d^{n-1}}{dt^{n-1}} M_X(t) \right) 
= \frac{d}{dt} \left( \int_{-\infty}^{\infty} x^{n-1} e^{tx} f_X(x) dx \right) 
= \int_{-\infty}^{\infty} \left( \frac{\partial}{\partial t} x^{n-1} e^{tx} f_X(x) \right) dx 
= \int_{-\infty}^{\infty} x^n e^{tx} f_X(x) dx 
= E(X^n e^{tX})$$

Evaluating this at t=0, we get:  $\frac{d^n}{dt^n}M_X(t)\mid_{t=0}=E(X^n)$ 

# Example 3.2.3:

 $\overline{X} \sim U(a, b)$ , where a < b;  $f_X(x) = \frac{1}{b-a} \cdot I_{[a,b]}(x)$ Then,

$$M_X(t) = \int_a^b \frac{e^{tx}}{b-a} dx = \frac{e^{tb} - e^{ta}}{t(b-a)}$$

$$M_X(0) = \frac{0}{0}$$

$$\text{L'Hospital} \quad \frac{be^{tb} - ae^{ta}}{b-a} \Big|_{t=0}$$

$$= 1$$

So  $M_X(0) = 1$  and since  $\frac{e^{tb} - e^{ta}}{t(b-a)}$  is continuous, it also exists in an open interval around 0 (in fact, it exists for every  $t \in \mathbb{R}$ ).

$$\begin{array}{lll} M_X'(t) & = & \frac{(be^{tb}-ae^{ta})t(b-a)-(e^{tb}-e^{ta})(b-a)}{t^2(b-a)^2} \\ & = & \frac{t(be^{tb}-ae^{ta})-(e^{tb}-e^{ta})}{t^2(b-a)} \\ \\ \Longrightarrow E(X) & = & M_X'(0) = \frac{0}{0} \\ \\ \text{L'Hospital} & \frac{be^{tb}-ae^{ta}+tb^2e^{tb}-ta^2e^{ta}-be^{tb}+ae^{ta}}{2t(b-a)} \Big|_{t=0} \\ & = & \frac{tb^2e^{tb}-ta^2e^{ta}}{2(b-a)} \Big|_{t=0} \\ & = & \frac{b^2e^{tb}-a^2e^{ta}}{2(b-a)} \Big|_{t=0} \\ & = & \frac{b^2-a^2}{2(b-a)} \\ & = & \frac{b+a}{2} \end{array}$$

## Note:

In the previous example, we made use of **L'Hospital's rule**. This rule gives conditions under which we can resolve indefinite expressions of the type " $\frac{\pm 0}{\pm 0}$ " and " $\frac{\pm \infty}{\pm \infty}$ ".

- (i) Let f and g be functions that are differentiable in an open interval around  $x_0$ , say in  $(x_0 \delta, x_0 + \delta)$ , but not necessarily differentiable in  $x_0$ . Let  $f(x_0) = g(x_0) = 0$  and  $g'(x) \neq 0 \ \forall x \in (x_0 \delta, x_0 + \delta) \{x_0\}$ . Then,  $\lim_{x \to x_0} \frac{f'(x)}{g'(x)} = A$  implies that also  $\lim_{x \to x_0} \frac{f(x)}{g(x)} = A$ . The same holds for the cases  $\lim_{x \to x_0} f(x) = \lim_{x \to x_0} g(x) = \infty$  and  $x \to x_0^+$  or  $x \to x_0^-$ .
- (ii) Let f and g be functions that are differentiable for x>a (a>0). Let  $\lim_{x\to\infty}f(x)=\lim_{x\to\infty}g(x)=0$  and  $\lim_{x\to\infty}g'(x)\neq0$ . Then,  $\lim_{x\to\infty}\frac{f'(x)}{g'(x)}=A$  implies that also  $\lim_{x\to\infty}\frac{f(x)}{g(x)}=A$ .
- (iii) We can iterate this process as long as the required conditions are met and derivatives exist, e.g., if the first derivatives still result in an indefinite expression, we can look at the second derivatives, then at the third derivatives, and so on.
- (iv) It is recommended to keep expressions as simple as possible. If we have identical factors in the numerator and denominator, we can exclude them from both and continue with the simpler functions.
- (v) Indefinite expressions of the form " $0 \cdot \infty$ " can be handled by rearranging them to " $\frac{0}{1/\infty}$ " and  $\lim_{x \to -\infty} \frac{f(x)}{g(x)}$  can be handled by use of the rules for  $\lim_{x \to \infty} \frac{f(-x)}{g(-x)}$ .

## Note:

Lecture 19: We 10/11/00

The following Theorems provide us with rules that tell us when we can differentiate under the integral sign. Theorem 3.2.4 relates to finite integral bounds  $a(\theta)$  and  $b(\theta)$  and Theorems 3.2.5 and 3.2.6 to infinite bounds.

## Theorem 3.2.4: Leibnitz's Rule

If  $f(x,\theta), a(\theta)$ , and  $b(\theta)$  are differentiable with respect to  $\theta$  (for all x) and  $-\infty < a(\theta) < b(\theta) < \infty$ , then

$$\frac{d}{d\theta} \int_{a(\theta)}^{b(\theta)} f(x,\theta) dx = f(b(\theta),\theta) \frac{d}{d\theta} b(\theta) - f(a(\theta),\theta) \frac{d}{d\theta} a(\theta) + \int_{a(\theta)}^{b(\theta)} \frac{\partial}{\partial \theta} f(x,\theta) dx.$$

The first 2 terms are vanishing if  $a(\theta)$  and  $b(\theta)$  are constant in  $\theta$ .

#### Proof:

Uses the Fundamental Theorem of Calculus and the chain rule.

# Theorem 3.2.5: Lebesque's Dominated Convergence Theorem

Let g be an integrable function such that  $\int_{-\infty}^{\infty} g(x)dx < \infty$ . If  $|f_n| \le g$  almost everywhere (i.e., except for a set of Borel-measure 0) and if  $f_n \to f$  almost everywhere, then  $f_n$  and f are integrable and

 $\int_{-\infty}^{\infty} f_n(x) dx \to \int_{-\infty}^{\infty} f(x) dx.$ 

# Note:

If f is differentiable with respect to  $\theta$ , then

$$\frac{\partial}{\partial \theta} f(x, \theta) = \lim_{\delta \to 0} \frac{f(x, \theta + \delta) - f(x, \theta)}{\delta}$$

and

$$\int_{-\infty}^{\infty} \frac{\partial}{\partial \theta} f(x, \theta) dx = \int_{-\infty}^{\infty} \lim_{\delta \to 0} \frac{f(x, \theta + \delta) - f(x, \theta)}{\delta} dx$$

while

$$\frac{\partial}{\partial \theta} \int_{-\infty}^{\infty} f(x,\theta) dx = \lim_{\delta \to 0} \int_{-\infty}^{\infty} \frac{f(x,\theta+\delta) - f(x,\theta)}{\delta} dx$$

#### Theorem 3.2.6:

Let  $f_n(x, \theta_0) = \frac{f(x, \theta_0 + \delta_n) - f(x, \theta_0)}{\delta_n}$  for some  $\theta_0$ . Suppose there exists an integrable function g(x) such that  $\int_{-\infty}^{\infty} g(x) dx < \infty$  and  $|f_n(x, \theta)| \le g(x) \ \forall x$ , then

$$\left[\frac{d}{d\theta} \int_{-\infty}^{\infty} f(x,\theta) dx\right]\Big|_{\theta=\theta_0} = \int_{-\infty}^{\infty} \left[\frac{\partial}{\partial \theta} f(x,\theta) \mid_{\theta=\theta_0}\right] dx.$$

Usually, if f is differentiable for all  $\theta$ , we write

$$\frac{d}{d\theta} \int_{-\infty}^{\infty} f(x,\theta) dx = \int_{-\infty}^{\infty} \frac{\partial}{\partial \theta} f(x,\theta) dx.$$

## Corollary 3.2.7:

Let  $f(x,\theta)$  be differentiable for all  $\theta$ . Suppose there exists an integrable function  $g(x,\theta)$  such that  $\int_{-\infty}^{\infty} g(x,\theta) dx < \infty$  and  $\left| \frac{\partial}{\partial \theta} f(x,\theta) \right|_{\theta=\theta_0} \le g(x,\theta) \quad \forall x \; \forall \theta_0 \text{ in some } \epsilon\text{-neighborhood of } \theta$ , then

$$\frac{d}{d\theta} \int_{-\infty}^{\infty} f(x,\theta) dx = \int_{-\infty}^{\infty} \frac{\partial}{\partial \theta} f(x,\theta) dx.$$

# More on Moment Generating Functions

Consider

$$\left| \frac{\partial}{\partial t} e^{tx} f_X(x) \right|_{t=t'} = |x| e^{t'x} f_X(x) \text{ for } |t'-t| \le \delta_0.$$

Choose  $t, \delta_0$  small enough such that  $t + \delta_0 \in (-h, h)$  and  $t - \delta_0 \in (-h, h)$ , or, equivalently,  $|t + \delta_0| < h$  and  $|t - \delta_0| < h$ . Then,

$$\left| \frac{\partial}{\partial t} e^{tx} f_X(x) \right|_{t=t'} \le g(x,t)$$

where

$$g(x,t) = \begin{cases} |x| e^{(t+\delta_0)x} f_X(x), & x \ge 0 \\ |x| e^{(t-\delta_0)x} f_X(x), & x < 0 \end{cases}$$

To verify  $\int g(x,t)dx < \infty$ , we need to know  $f_X(x)$ .

Suppose  $\operatorname{mgf} M_X(t)$  exists for  $|t| \leq h^*$  for some  $h^* > 1$ , where  $h^* - 1 \geq h$ . Then  $|t + \delta_0 + 1| < h^*$  and  $|t - \delta_0 - 1| < h^*$ . Since  $|x| \leq e^{|x|} \ \forall x$ , we get

$$g(x,t) \le \begin{cases} e^{(t+\delta_0+1)x} f_X(x), & x \ge 0\\ e^{(t-\delta_0-1)x} f_X(x), & x < 0 \end{cases}$$

Then, 
$$\int_0^\infty g(x,t)dx \le M_X(t+\delta_0+1) < \infty$$
 and  $\int_{-\infty}^0 g(x,t)dx \le M_X(t-\delta_0-1) < \infty$  and, therefore,  $\int_{-\infty}^\infty g(x)dx < \infty$ .

Together with Corollary 3.2.7, this establishes that we can differentiate under the integral in the Proof of Theorem 3.2.2.

If  $h^* \leq 1$ , we may need to check more carefully to see if the condition holds.

## Note:

If  $M_X(t)$  exists for  $t \in (-h,h)$ , then we have an infinite collection of moments.

Does a collection of integer moments  $\{m_k : k = 1, 2, 3, ...\}$  completely characterize the distribution, i.e., cdf, of X? — Unfortunately not, as Example 3.2.8 shows.

#### Example 3.2.8:

Let  $X_1$  and  $X_2$  be rv's with pdfs

$$f_{X_1}(x) = \frac{1}{\sqrt{2\pi}} \frac{1}{x} \exp(-\frac{1}{2} (\log x)^2) \cdot I_{(0,\infty)}(x)$$

and

$$f_{X_2}(x) = f_{X_1}(x) \cdot (1 + \sin(2\pi \log x)) \cdot I_{(0,\infty)}(x)$$

It is  $E(X_1^r) = E(X_2^r) = e^{r^2/2}$  for r = 0, 1, 2, ... as you have to show in the Homeworks.

Two different pdfs/cdfs have the same moment sequence! What went wrong? In this example,  $M_{X_1}(t)$  does not exist as shown in the Homeworks!

# Theorem 3.2.9: Let V and V be 2 with edfin E and E for which all moments exist Fr 10/13/00

Let X and Y be 2 rv's with cdf's  $F_X$  and  $F_Y$  for which all moments exist.

- (i) If  $F_X$  and  $F_Y$  have bounded support, then  $F_X(u) = F_Y(u) \ \forall u \text{ iff } E(X^r) = E(Y^r) \text{ for } r = 0, 1, 2, \dots$
- (ii) If both mgf's exist, i.e.,  $M_X(t) = M_Y(t)$  for t in some neighborhood of 0, then  $F_X(u) = F_Y(u) \ \forall u$ .

# Note:

The existence of moments is not equivalent to the existence of a mgf as seen in Example 3.2.8 above and some of the Homework assignments.

## <u>Theorem 3.2.10:</u>

Suppose rv's  $\{X_i\}_{i=1}^{\infty}$  have mgf's  $M_{X_i}(t)$  and that  $\lim_{i\to\infty} M_{X_i}(t) = M_X(t) \ \forall t\in (-h,h)$  for some h>0 and that  $M_X(t)$  itself is a mgf. Then, there exists a cdf  $F_X$  whose moments are determined by  $M_X(t)$  and for all continuity points x of  $F_X(x)$  it holds that  $\lim_{i\to\infty} F_{X_i}(x) = F_X(x)$ , i.e., the convergence of mgf's implies the convergence of cdf's.

# Proof:

Uniqueness of Laplace transformations, etc.

#### <u>Theorem 3.2.11:</u>

For constants a and b, the mgf of Y = aX + b is

$$M_Y(t) = e^{bt} M_X(at),$$

given that  $M_X(t)$  exists.

# Proof:

$$M_Y(t) = E(e^{(aX+b)t})$$

$$= E(e^{aXt}e^{bt})$$

$$= e^{bt}E(e^{Xat})$$

$$= e^{bt}M_X(at)$$

# 3.3 Complex-Valued Random Variables and Characteristic Functions

Recall the following facts regarding complex numbers:

$$i^0 = +1; i = \sqrt{-1}; i^2 = -1; i^3 = -i; i^4 = +1;$$
 etc.

in the planar Gauss'ian number plane it holds that i = (0, 1)

$$z = a + ib = r(\cos\phi + i\sin\phi)$$

$$r = |z| = \sqrt{a^2 + b^2}$$

$$\tan \phi = \frac{b}{a}$$

Euler's Relation:  $z = r(\cos \phi + i \sin \phi) = re^{i\phi}$ 

# Mathematical Operations on Complex Numbers:

$$z_1 \pm z_2 = (a_1 \pm a_2) + i(b_1 \pm b_2)$$

$$z_1 \cdot z_2 = r_1 r_2 e^{i(\phi_1 + \phi_2)} = r_1 r_2 (\cos(\phi_1 + \phi_2) + i \sin(\phi_1 + \phi_2))$$

$$\frac{z_1}{z_2} = \frac{r_1}{r_2} e^{i(\phi_1 - \phi_2)} = \frac{r_1}{r_2} (\cos(\phi_1 - \phi_2) + i\sin(\phi_1 - \phi_2))$$

Moivre's Theorem:  $z^n = (r(\cos\phi + i\sin\phi))^n = r^n(\cos(n\phi) + i\sin(n\phi))$ 

$$\sqrt[n]{z} = \sqrt[n]{a+ib} = \sqrt[n]{r} \left( \cos(\frac{\phi+k\cdot 360^0}{n}) + i\sin(\frac{\phi+k\cdot 360^0}{n}) \right)$$
 for  $k = 0, 1, \dots, (n-1)$  and the main value for  $k = 0$ 

 $\ln z = \ln(a+ib) = \ln(|z|) + i\phi \pm i2n\pi$  where  $\phi = \arctan \frac{b}{a}$  and the main value for n=0

# Conjugate Complex Numbers:

For z = a + ib, we define the **conjugate complex number**  $\overline{z} = a - ib$ . It holds:

$$\overline{\overline{z}} = z$$

$$z = \overline{z} \text{ iff } z \in \mathbb{R}$$

$$\overline{z_1 \pm z_2} = \overline{z_1} \pm \overline{z_2}$$

$$\overline{z_1 \cdot z_2} = \overline{z_1} \cdot \overline{z_2}$$

$$\overline{\left(\frac{z_1}{z_2}\right)} = \overline{\left(\frac{\overline{z_1}}{\overline{z_2}}\right)}$$

$$z\cdot \overline{z}=a^2+b^2$$

$$Re(z) = a = \frac{1}{2}(z + \overline{z})$$

$$Im(z) = b = \frac{1}{2i}(z - \overline{z})$$

$$\mid z \mid = \sqrt{a^2 + b^2} = \sqrt{z \cdot \overline{z}}$$

## Definition 3.3.1:

Let  $(\Omega, L, P)$  be a probability space and X and Y real-valued rv's, i.e.,  $X, Y : (\Omega, L) \to (I\!\!R, \mathcal{B})$ 

(i)  $Z = X + iY : (\Omega, L) \to (\mathcal{C}, \mathcal{B}_{\mathcal{C}})$  is called a **complex-valued random variable**  $(\mathcal{C}-rv)$ .

(ii) If E(X) and E(Y) exist, then E(Z) is defined as  $E(Z) = E(X) + iE(Y) \in \mathcal{C}$ .

Note:

E(Z) exists iff E(|X|) and E(|Y|) exist. It also holds that if E(Z) exists, then  $|E(Z)| \le E(|Z|)$  (see Homework).

## Definition 3.3.2:

Let X be a real-valued rv on  $(\Omega, L, P)$ . Then,  $\Phi_X(t) : \mathbb{R} \to \mathbb{C}$  with  $\Phi_X(t) = E(e^{itX})$  is called the **characteristic function** of X.

Note:

(i)  $\Phi_X(t) = \int_{-\infty}^{\infty} e^{itx} f_X(x) dx = \int_{-\infty}^{\infty} \cos(tx) f_X(x) dx + i \int_{-\infty}^{\infty} \sin(tx) f_X(x) dx$  if X is continuous.

(ii)  $\Phi_X(t) = \sum_{x \in \mathcal{X}} e^{itx} P(X = x) = \sum_{x \in \mathcal{X}} \cos(tx) P(X = x) + i \sum_{x \in \mathcal{X}} \sin(tx) P(X = x)(x)$ 

if X is discrete and  $\mathcal{X}$  is the support of X.

(iii)  $\Phi_X(t)$  exists for all real-valued rv's X since  $|e^{itx}| = 1$ .

# Theorem 3.3.3:

Let  $\Phi_X$  be the characteristic function of a real-valued rv X. Then it holds:

- (i)  $\Phi_X(0) = 1$ .
- (ii)  $|\Phi_X(t)| \leq 1 \quad \forall t \in \mathbb{R}$ .
- (iii)  $\Phi_X$  is uniformly continuous, i.e.,  $\forall \epsilon > 0 \ \exists \delta > 0 \ \forall t_1, t_2 \in \mathbb{R} : |t_1 t_2| < \delta \Rightarrow |\Phi(t_1) \Phi(t_2)| < \epsilon$ .
- (iv)  $\Phi_X$  is a positive definite function, i.e.,  $\forall n \in \mathbb{I} N \ \forall \alpha_1, \dots, \alpha_n \in \mathbb{C} \ \forall t_1, \dots, t_n \in \mathbb{R} :$   $\sum_{l=1}^n \sum_{j=1}^n \alpha_l \overline{\alpha_j} \Phi_X(t_l t_j) \geq 0.$

- (v)  $\Phi_X(t) = \overline{\Phi_X(-t)}$ .
- (vi) If X is symmetric around 0, i.e., if X has a pdf that is symmetric around 0, then  $\Phi_X(t) \in \mathbb{R} \ \forall t \in \mathbb{R}$ .
- (vii)  $\Phi_{aX+b}(t) = e^{itb}\Phi_X(at)$ .

# Proof:

See Homework for parts (i), (ii), (iv), (v), (vi), and (vii).

# Part (iii):

Known conditions:

- (i) Let  $\epsilon > 0$ .
- (ii)  $\exists a > 0 : P(-a < X < +a) > 1 \frac{\epsilon}{4} \text{ and } P(|X| > a) < \frac{\epsilon}{4}$
- (iii)  $\exists \ \delta > 0 : \ |e^{i(t'-t)x} 1| < \frac{\epsilon}{2} \ \forall x \text{ s.t.} \ |x| < a \text{ and } \forall (t'-t) \text{ s.t. } 0 < (t'-t) < \delta.$

This third condition holds since  $|e^{i0}-1|=0$  and the exponential function is continuous. Therefore, if we select (t'-t) and x small enough,  $|e^{i(t'-t)x}-1|$  will be  $<\frac{\epsilon}{2}$  for a given  $\epsilon$ .

Let  $t, t' \in \mathbb{R}$ , t < t', and  $t' - t < \delta$ . Then,

$$\begin{split} |\Phi_{X}(t') - \Phi_{X}(t)| &= |\int_{-\infty}^{+\infty} e^{\imath t'x} f_{X}(x) dx - \int_{-\infty}^{+\infty} e^{\imath tx} f_{X}(x) dx | \\ &= |\int_{-\infty}^{+\infty} (e^{\imath t'x} - e^{\imath tx}) f_{X}(x) dx | \\ &= |\int_{-\infty}^{-a} (e^{\imath t'x} - e^{\imath tx}) f_{X}(x) dx + \int_{-a}^{+a} (e^{\imath t'x} - e^{\imath tx}) f_{X}(x) dx + \int_{+a}^{+\infty} (e^{\imath t'x} - e^{\imath tx}) f_{X}(x) dx | \\ &\leq |\int_{-\infty}^{-a} (e^{\imath t'x} - e^{\imath tx}) f_{X}(x) dx | + |\int_{-a}^{+a} (e^{\imath t'x} - e^{\imath tx}) f_{X}(x) dx | \\ &+ |\int_{+a}^{+\infty} (e^{\imath t'x} - e^{\imath tx}) f_{X}(x) dx | \end{split}$$

Lecture 21: Mo 10/16/00

We now take a closer look at the first and third of these absolute integrals. It is:

$$| \int_{-\infty}^{-a} (e^{it'x} - e^{itx}) f_X(x) dx | = | \int_{-\infty}^{-a} e^{it'x} f_X(x) dx - \int_{-\infty}^{-a} e^{itx} f_X(x) dx |$$

$$\leq | \int_{-\infty}^{-a} e^{it'x} f_X(x) dx | + | \int_{-\infty}^{-a} e^{itx} f_X(x) dx |$$

$$\leq \int_{-\infty}^{-a} | e^{it'x} | f_X(x) dx + \int_{-\infty}^{-a} | e^{itx} | f_X(x) dx$$

$$\stackrel{(A)}{=} \int_{-\infty}^{-a} 1 f_X(x) dx + \int_{-\infty}^{-a} 1 f_X(x) dx$$

$$= \int_{-\infty}^{-a} 2 f_X(x) dx.$$

(A) holds due to Note (iii) that follows Definition 3.3.2.

Similarly,

$$\mid \int_{+a}^{+\infty} (e^{\imath t'x} - e^{\imath tx}) f_X(x) dx \mid \leq \int_{+a}^{+\infty} 2 f_X(x) dx$$

Returning to the main part of the proof, we get

$$\begin{split} |\Phi_{X}(t') - \Phi_{X}(t)| &\leq \int_{-\infty}^{-a} 2f_{X}(x)dx + |\int_{-a}^{+a} (e^{\imath t'x} - e^{\imath tx})f_{X}(x)dx| + \int_{+a}^{+\infty} 2f_{X}(x)dx \\ &= 2\left(\int_{-\infty}^{-a} f_{X}(x)dx + \int_{+a}^{+\infty} f_{X}(x)dx\right) + |\int_{-a}^{+a} (e^{\imath t'x} - e^{\imath tx})f_{X}(x)dx| \\ &= 2P(|X| > a) + |\int_{-a}^{+a} (e^{\imath t'x} - e^{\imath tx})f_{X}(x)dx| \\ &\stackrel{Condition\ (ii)}{\leq} 2\frac{\epsilon}{4} + |\int_{-a}^{+a} (e^{\imath t'x} - e^{\imath tx})f_{X}(x)dx| \\ &= \frac{\epsilon}{2} + |\int_{-a}^{+a} e^{\imath tx} (e^{\imath (t'-t)x} - 1)f_{X}(x)dx| \\ &\leq \frac{\epsilon}{2} + \int_{-a}^{+a} |e^{\imath tx}| (e^{\imath (t'-t)x} - 1) |f_{X}(x)dx| \\ &\leq \frac{\epsilon}{2} + \int_{-a}^{+a} |e^{\imath tx}| |e^{\imath (t'-t)x} - 1| |f_{X}(x)dx| \end{split}$$

$$\stackrel{(B)}{\leq} \frac{\epsilon}{2} + \int_{-a}^{+a} 1 \frac{\epsilon}{2} f_X(x) dx$$

$$\stackrel{(B)}{\leq} \frac{\epsilon}{2} + \int_{-\infty}^{+\infty} \frac{\epsilon}{2} f_X(x) dx$$

$$= \frac{\epsilon}{2} + \frac{\epsilon}{2}$$

$$= \epsilon$$

(B) holds due to Note (iii) that follows Definition 3.3.2 and due to condition (iii).

## Theorem 3.3.4: Bochner's Theorem

Let  $\Phi : \mathbb{R} \to \mathbb{C}$  be any function with properties (i), (ii), (iii), and (iv) from Theorem 3.3.3. Then there exists a real-valued rv X with  $\Phi_X = \Phi$ .

## Theorem 3.3.5:

Let X be a real-valued rv and  $E(X^k)$  exists for an integer k. Then,  $\Phi_X$  is k times differentiable and  $\Phi_X^{(k)}(t) = i^k E(X^k e^{itX})$ . In particular for t = 0, it is  $\Phi_X^{(k)}(0) = i^k m_k$ .

#### Theorem 3.3.6:

Let X be a real-valued rv with characteristic function  $\Phi_X$  and let  $\Phi_X$  be k times differentiable, where k is an even integer. Then the  $k^{th}$  moment of X,  $m_k$ , exists and it is  $\Phi_X^{(k)}(0) = i^k m_k$ .

# Theorem 3.3.7: Levy's Theorem

Let X be a real-valued rv with cdf  $F_X$  and characteristic function  $\Phi_X$ . Let  $a, b \in \mathbb{R}$ , a < b. If P(X = a) = P(X = b) = 0, i.e.,  $F_X$  is continuous in a and b, then

$$F(b) - F(a) = \frac{1}{2\pi} \int_{-\infty}^{\infty} \frac{e^{-ita} - e^{-itb}}{it} \Phi_X(t) dt.$$

## <u>Theorem 3.3.8:</u>

Let X and Y be a real-valued rv with characteristic functions  $\Phi_X$  and  $\Phi_Y$ . If  $\Phi_X = \Phi_Y$ , then X and Y are identically distributed.

## Theorem 3.3.9:

Let X be a real-valued rv with characteristic function  $\Phi_X$  such that  $\int_{-\infty}^{\infty} |\Phi_X(t)| dt < \infty$ . Then X has pdf

$$f_X(x) = \frac{1}{2\pi} \int_{-\infty}^{\infty} e^{-itx} \Phi_X(t) dt.$$

# <u>Theorem 3.3.10:</u>

Let X be a real-valued rv with mgf  $M_X(t)$ , i.e., the mgf exists. Then  $\Phi_X(t) = M_X(it)$ .

## <u>Theorem 3.3.11:</u>

Suppose real-valued rv's  $\{X_i\}_{i=1}^{\infty}$  have cdf's  $\{F_{X_i}\}_{i=1}^{\infty}$  and characteristic functions  $\{\Phi_{X_i}(t)\}_{i=1}^{\infty}$ . If  $\lim_{i\to\infty}\Phi_{X_i}(t)=\Phi_X(t)\ \, \forall t\in(-h,h)$  for some h>0 and  $\Phi_X(t)$  is itself a characteristic function (of a rv X with cdf  $F_X$ ), then  $\lim_{i\to\infty}F_{X_i}(x)=F_X(x)$  for all continuity points x of  $F_X(x)$ , i.e., the convergence of characteristic functions implies the convergence of cdf's.

## Theorem 3.3.12:

Characteristic functions for some well-known distributions:

Distribution	$\Phi_X(t)$
$(i)$ $X \sim Dirac(c)$	$e^{itc}$
(ii) $X \sim Bin(1, p)$	$1 + p(e^{it} - 1)$
(iii) $X \sim \text{Poisson}(c)$	$\exp(c(e^{it}-1))$
(iv) $X \sim U(a, b)$	$rac{e^{itb}-e^{ita}}{(b-a)it}$
(v) $X \sim N(0, 1)$	$\exp(-t^2/2)$
(vi) $X \sim N(\mu, \sigma^2)$	$e^{it\mu} \exp(-\sigma^2 t^2/2)$
(vii) $X \sim \Gamma(p, q)$	$(1-\frac{it}{q})^{-p}$
(viii) $X \sim Exp(c)$	$(1-\frac{it}{c})^{-1}$
(ix) $X \sim \chi_n^2$	$(1-2it)^{-n/2}$

# Proof:

(i) 
$$\Phi_X(t) = E(e^{itX}) = e^{itc}P(X = c) = e^{itc}$$

(ii) 
$$\Phi_X(t) = \sum_{k=0}^{1} e^{itk} P(X=k) = e^{it0} (1-p) + e^{it1} p = 1 + p(e^{it}-1)$$

(iii) 
$$\Phi_X(t) = \sum_{n \in \mathbb{N}_0} e^{itn} \cdot \frac{c^n}{n!} e^{-c} = e^{-c} \sum_{n=0}^{\infty} \frac{1}{n!} (c \cdot e^{it})^n = e^{-c} \cdot e^{c \cdot e^{it}} = e^{c(e^{it} - 1)}$$

$$\text{since } \sum_{n=0}^{\infty} \frac{x^n}{n!} = e^x$$

(iv) 
$$\Phi_X(t) = \frac{1}{b-a} \int_a^b e^{itx} dx = \frac{1}{b-a} \left[ \frac{e^{itx}}{it} \right]_a^b = \frac{e^{itb} - e^{ita}}{(b-a)it}$$

Lecture 22: We 10/18/00

(v)  $X \sim N(0, 1)$  is symmetric around 0

 $\Longrightarrow \Phi_X(t)$  is real since there is no imaginary part according to Theorem 3.3.3 (vi)

$$\Longrightarrow \Phi_X(t) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} \cos(tx) e^{\frac{-x^2}{2}} dx$$

Since E(X) exists,  $\Phi_X(t)$  is differentiable according to Theorem 3.3.5 and the following holds:

$$\begin{split} \Phi_X'(t) &= Re(\Phi_X'(t)) \\ &= Re\left(\int_{-\infty}^{\infty} ix \underbrace{e^{itx}}_{\cos(tx)+i\sin(tx)} \frac{1}{\sqrt{2\pi}} e^{\frac{-x^2}{2}} dx\right) \\ &= Re\left(\int_{-\infty}^{\infty} ix \cos(tx) \frac{1}{\sqrt{2\pi}} e^{\frac{-x^2}{2}} dx + \int_{-\infty}^{\infty} -x \sin(tx) \frac{1}{\sqrt{2\pi}} e^{\frac{-x^2}{2}} dx\right) \\ &= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} \underbrace{(-\sin(tx))}_{u} \underbrace{xe^{\frac{-x^2}{2}}}_{v'} dx \quad | \quad u' = -t \cos(tx) \text{ and } v = -e^{\frac{-x^2}{2}} \\ &= \underbrace{\frac{1}{\sqrt{2\pi}}}_{=0} \sin(tx) e^{\frac{-x^2}{2}} \Big|_{-\infty}^{\infty} - \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} (-t \cos(tx)) (-e^{\frac{-x^2}{2}}) dx \\ &= 0 \text{ since sin is odd} \\ &= -t \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} \cos(tx) e^{\frac{-x^2}{2}} dx \\ &= -t \Phi_X(t) \end{split}$$

Thus,  $\Phi_X'(t) = -t\Phi_X(t)$ . It follows that  $\frac{\Phi_X'(t)}{\Phi_X(t)} = -t$  and by integrating both sides, we get  $\ln |\Phi_X(t)| = -\frac{1}{2}t^2 + c$  with  $c \in \mathbb{R}$ .

For t=0, we know that  $\Phi_X(0)=1$  by Theorem 3.3.3 (i) and  $\ln |\Phi_X(0)|=0$ . It follows that 0=0+c. Therefore, c=0 and  $|\Phi_X(t)|=e^{-\frac{1}{2}t^2}$ .

If we take t=0, then  $\Phi_X(0)=1$  by Theorem 3.3.3 (i). Since  $\Phi_X$  is uniformly continuous,  $\Phi_X$  must take the value 0 before it can eventually take a negative value. However, since  $e^{-\frac{1}{2}t^2}>0 \ \forall t\in \mathbb{R}, \ \Phi_X$  cannot take 0 as a possible value and therefore cannot pass into the negative numbers. So, it must hold that  $\Phi_X(t)=e^{-\frac{1}{2}t^2} \ \forall t\in \mathbb{R}$ .

(vi) For  $\sigma > 0, \mu \in \mathbb{R}$ , we know that if  $X \sim N(0,1)$ , then  $\sigma X + \mu \sim N(\mu, \sigma^2)$ . By Theorem 3.3.3 (vii) we have

$$\Phi_{\sigma X+\mu}(t) = e^{it\mu}\Phi_X(\sigma t) = e^{it\mu}e^{-\frac{1}{2}\sigma^2t^2}.$$

(vii)

$$\begin{split} \Phi_{X}(t) &= \int_{0}^{\infty} e^{itx} \gamma(p,q,x) dx \\ &= \int_{0}^{\infty} \frac{q^{p}}{\Gamma(p)} x^{p-1} e^{-(q-it)x} dx \\ &= \frac{q^{p}}{\Gamma(p)} (q-it)^{-p} \int_{0}^{\infty} ((q-it)x)^{p-1} e^{-(q-it)x} (q-it) dx \mid u = (q-it)x, \ du = (q-it) dx \\ &= \frac{q^{p}}{\Gamma(p)} (q-it)^{-p} \underbrace{\int_{0}^{\infty} (u)^{p-1} e^{-u} du}_{=\Gamma(p)} \\ &= q^{p} (q-it)^{-p} \\ &= \frac{(q-it)^{-p}}{q^{-p}} \\ &= \left(\frac{q-it}{q}\right)^{-p} \\ &= \left(1 - \frac{it}{q}\right)^{-p} \end{split}$$

(viii) Since an Exp(c) distribution is a  $\Gamma(1,c)$  distribution, we get for  $X \sim Exp(c) = \Gamma(1,c)$ :  $\operatorname{Fr} \frac{10}{20}/00$ 

$$\Phi_X(t) = (1 - \frac{it}{c})^{-1}$$

(ix) Since a  $\chi^2_n$  distribution (for  $n \in I\!\!N$ ) is a  $\Gamma(\frac{n}{2},\frac{1}{2})$  distribution, we get for  $X \sim \chi^2_n =$  $\Gamma(\frac{n}{2},\frac{1}{2})$ :

$$\Phi_X(t) = \left(1 - \frac{it}{1/2}\right)^{-n/2} = (1 - 2it)^{-n/2}$$

# Example 3.3.13:

Since we know that  $m_1 = E(X)$  and  $m_2 = E(X^2)$  exist for  $X \sim Bin(1, p)$ , we can determine these moments according to Theorem 3.3.5 using the characteristic function.

It is

$$\Phi_X(t) = 1 + p(e^{it} - 1)$$

$$\Phi_X'(t) = pie^{it}$$

$$\Phi_X'(0) = pi$$

$$\Rightarrow m_1 = \frac{\Phi_X'(0)}{i} = \frac{pi}{i} = p = E(X)$$

$$\Phi_X''(t) = pi^2e^{it}$$

$$\Phi_X''(0) = pi^2$$

$$\Rightarrow m_2 = \frac{\Phi_X''(0)}{i^2} = \frac{pi^2}{i^2} = p = E(X^2)$$

$$\Rightarrow Var(X) = E(X^2) - (E(X))^2 = p - p^2 = p(1 - p)$$

Note:

The restriction  $\int_{-\infty}^{\infty} |\Phi_X(t)| dt < \infty$  in Theorem 3.3.9 works in such a way that we don't end up with a (non-existing) pdf if X is a discrete rv. For example,

•  $X \sim \text{Dirac}(c)$ :

$$\int_{-\infty}^{\infty} |\Phi_X(t)| dt = \int_{-\infty}^{\infty} |e^{itc}| dt$$
$$= \int_{-\infty}^{\infty} 1 dt$$
$$= x \Big|_{-\infty}^{\infty}$$

which is undefined.

• Also for  $X \sim Bin(1, p)$ :

$$\int_{-\infty}^{\infty} | \Phi_{X}(t) | dt = \int_{-\infty}^{\infty} | 1 + p(e^{it} - 1) | dt 
= \int_{-\infty}^{\infty} | pe^{it} - (p - 1) | dt 
\ge \int_{-\infty}^{\infty} | | pe^{it} | - | (p - 1) | | dt 
\ge \int_{-\infty}^{\infty} | pe^{it} | dt - \int_{-\infty}^{\infty} | (p - 1) | dt$$

$$= p \int_{-\infty}^{\infty} 1 dt - (1 - p) \int_{-\infty}^{\infty} 1 dt$$
$$= (2p - 1) \int_{-\infty}^{\infty} 1 dt$$
$$= (2p - 1)x \Big|_{-\infty}^{\infty}$$

which is undefined for  $p \neq 1/2$ .

If p = 1/2, we have

$$\int_{-\infty}^{\infty} |pe^{it} - (p-1)| dt = 1/2 \int_{-\infty}^{\infty} |e^{it} + 1| dt$$

$$= 1/2 \int_{-\infty}^{\infty} |\cos t + i\sin t + 1| dt$$

$$= 1/2 \int_{-\infty}^{\infty} \sqrt{(\cos t + 1)^2 + (\sin t)^2} dt$$

$$= 1/2 \int_{-\infty}^{\infty} \sqrt{\cos^2 t + 2\cos t + 1 + \sin^2 t} dt$$

$$= 1/2 \int_{-\infty}^{\infty} \sqrt{2 + 2\cos t} dt$$

which also does not exist.

• Otherwise,  $X \sim N(0, 1)$ :

$$\int_{-\infty}^{\infty} |\Phi_X(t)| dt = \int_{-\infty}^{\infty} \exp(-t^2/2) dt$$

$$= \sqrt{2\pi} \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} \exp(-t^2/2) dt$$

$$= \sqrt{2\pi}$$

$$< \infty$$

# 3.4 Probability Generating Functions

## Definition 3.4.1:

Let X be a discrete rv which only takes non-negative integer values, i.e.,  $p_k = P(X = k)$ , and  $\sum_{k=0}^{\infty} p_k = 1$ . Then, the **probability generating function** (**pgf**) of X is defined as

$$G(s) = \sum_{k=0}^{\infty} p_k s^k.$$

## Theorem 3.4.2:

G(s) converges for |s| < 1.

## Proofs

$$\overline{\mid G(s)\mid} \leq \sum_{k=0}^{\infty} \mid p_k s^k \mid \leq \sum_{k=0}^{\infty} \mid p_k \mid = 1$$

# <u>Theorem 3.4.3:</u>

Let X be a discrete rv which only takes non-negative integer values and has pgf G(s). Then it holds:

$$P(X = k) = \frac{1}{k!} \frac{d^k}{ds^k} G(s) |_{s=0}$$

## <u>Theorem 3.4.4:</u>

Let X be a discrete rv which only takes non-negative integer values and has pgf G(s). If E(X) exists, then it holds:

$$E(X) = \frac{d}{ds}G(s)\mid_{s=1}$$

# Definition 3.4.5:

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

$$E[X(X-1)(X-2)\cdot\ldots\cdot(X-k+1)]$$

if this expectation exists.

# <u>Theorem 3.4.6:</u>

Let X be a discrete rv which only takes non-negative integer values and has pgf G(s). If  $E[X(X-1)(X-2)\cdot\ldots\cdot(X-k+1)]$  exists, then it holds:

$$E[X(X-1)(X-2)\cdot\ldots\cdot(X-k+1)] = \frac{d^k}{ds^k}G(s)|_{s=1}$$

# Note:

Similar to the Cauchy distribution for the continuous case, there exist discrete distributions where the mean (or higher moments) do not exist. See Homework.

Lecture 24: Mo 10/23/00

# 3.5 Moment Inequalities

## <u>Theorem 3.5.1:</u>

Let h(X) be a non-negative Borel-measurable function of a rv X. If E(h(X)) exists, then it holds:

$$P(h(X) \ge \epsilon) \le \frac{E(h(X))}{\epsilon} \quad \forall \epsilon > 0$$

# Proof:

Continuous case only:

$$\begin{split} E(h(X)) &= \int_{-\infty}^{\infty} h(x) f_X(x) dx \\ &= \int_A h(x) f_X(x) dx + \int_{A^C} h(x) f_X(x) dx \quad | \text{ where } A = \{x : h(x) \ge \epsilon\} \\ &\ge \int_A h(x) f_X(x) dx \\ &\ge \int_A \epsilon f_X(x) dx \\ &= \epsilon P(h(X) \ge \epsilon) \quad \forall \epsilon > 0 \end{split}$$

Therefore,  $P(h(X) \ge \epsilon) \le \frac{E(h(X))}{\epsilon} \quad \forall \epsilon > 0.$ 

# Corollary 3.5.2: Markov's Inequality

Let  $h(X) = |X|^r$  and  $\epsilon = k^r$  where r > 0 and k > 0. If  $E(|X|^r)$  exists, then it holds:

$$P(\mid X \mid \ge k) \le \frac{E(\mid X \mid^r)}{k^r}$$

# Proof:

Since  $P(\mid X \mid \geq k) = P(\mid X \mid^r \geq k^r)$  for k > 0, it follows using Theorem 3.5.1:

$$P(\mid X \mid \geq k) = P(\mid X \mid^{r} \geq k^{r}) \stackrel{Th.3.5.1}{\leq} \frac{E(\mid X \mid^{r})}{k^{r}}$$

# Corollary 3.5.3: Chebychev's Inequality

Let  $h(X) = (X - \mu)^2$  and  $\epsilon = k^2 \sigma^2$  where  $E(X) = \mu$ ,  $Var(X) = \sigma^2 < \infty$ , and k > 0. Then it holds:

$$P(\mid X - \mu \mid > k\sigma) \le \frac{1}{k^2}$$

# Proof:

Since  $P(|X - \mu| > k\sigma) = P(|X - \mu|^2 > k^2\sigma^2)$  for k > 0, it follows using Theorem 3.5.1:

$$P(\mid X - \mu \mid > k\sigma) = P(\mid X - \mu \mid^{2} > k^{2}\sigma^{2}) \stackrel{Th.3.5.1}{\leq} \frac{E(\mid X - \mu \mid^{2})}{k^{2}\sigma^{2}} = \frac{Var(X)}{k^{2}\sigma^{2}} = \frac{\sigma^{2}}{k^{2}\sigma^{2}} = \frac{1}{k^{2}}$$

# Theorem 3.5.4: Lyapunov's Inequality

Let  $0 < \beta_n = E(|X|^n) < \infty$ . For arbitrary k such that  $2 \le k \le n$ , it holds that

$$(\beta_{k-1})^{\frac{1}{k-1}} \le (\beta_k)^{\frac{1}{k}},$$

i.e., 
$$(E(|X|^{k-1}))^{\frac{1}{k-1}} \le (E(|X|^k))^{\frac{1}{k}}$$
.

# Proof:

Continuous case only:

Let  $Q(u,v) = E\left(\left(u \mid X \mid \frac{k-1}{2} + v \mid X \mid \frac{k+1}{2}\right)^2\right)$ . Obviously,  $Q(u,v) \geq 0 \ \forall u,v \in \mathbb{R}$ . Also,

$$Q(u,v) = \int_{-\infty}^{\infty} (u \mid x \mid^{\frac{k-1}{2}} + v \mid x \mid^{\frac{k+1}{2}})^{2} f_{X}(x) dx$$

$$= u^{2} \int_{-\infty}^{\infty} |x|^{k-1} f_{X}(x) dx + 2uv \int_{-\infty}^{\infty} |x|^{k} f_{X}(x) dx + v^{2} \int_{-\infty}^{\infty} |x|^{k+1} f_{X}(x) dx$$

$$= u^{2} \beta_{k-1} + 2uv \beta_{k} + v^{2} \beta_{k+1}$$

$$\geq 0 \quad \forall u, v \in \mathbb{R}$$

Using the fact that  $Ax^2 + 2Bxy + Cy^2 \ge 0 \quad \forall x, y \in \mathbb{R}$  iff A > 0 and  $AC - B^2 > 0$  (see Rohatgi, page 6, Section P2.4), we get with  $A = \beta_{k-1}, B = \beta_k$ , and  $C = \beta_{k+1}$ :

$$\beta_{k-1}\beta_{k+1} - \beta_k^2 \ge 0$$

$$\Longrightarrow \beta_k^2 \le \beta_{k-1}\beta_{k+1}$$

$$\Longrightarrow \beta_k^{2k} \le \beta_{k-1}^k \beta_{k+1}^k$$

This means that  $\beta_1^2 \leq \beta_0 \beta_2$ ,  $\beta_2^4 \leq \beta_1^2 \beta_3^2$ ,  $\beta_3^6 \leq \beta_2^3 \beta_4^3$ , and so on. Multiplying these, we get:

$$\prod_{j=1}^{k-1} \beta_j^{2j} \leq \prod_{j=1}^{k-1} \beta_{j-1}^j \beta_{j+1}^j 
= (\beta_0 \beta_2) (\beta_1^2 \beta_3^2) (\beta_2^3 \beta_4^3) (\beta_3^4 \beta_5^4) \dots (\beta_{k-3}^{k-2} \beta_{k-1}^{k-2}) (\beta_{k-2}^{k-1} \beta_k^{k-1}) 
= \beta_0 \beta_{k-1}^{k-2} \beta_k^{k-1} \prod_{j=1}^{k-2} \beta_j^{2j}$$

Dividing both sides by  $\prod_{j=1}^{k-2} \beta_j^{2j}$ , we get:

$$\beta_{k-1}^{2k-2} \leq \beta_0 \beta_k^{k-1} \beta_{k-1}^{k-2}$$

$$\stackrel{\beta_0=1}{\Longrightarrow} \beta_{k-1}^k \leq \beta_k^{k-1}$$

$$\Longrightarrow \beta_{k-1}^1 \leq \beta_k^{\frac{k-1}{k}}$$

$$\Longrightarrow \beta_{k-1}^{\frac{1}{k-1}} \leq \beta_k^{\frac{1}{k}}$$

70

Lecture 25: We 10/25/00

# 4 Random Vectors

# 4.1 Joint, Marginal, and Conditional Distributions

#### Definition 4.1.1:

The vector  $\underline{X} = (X_1, \dots, X_n)'$  on  $(\Omega, L, P) \to \mathbb{R}^n$  defined by  $\underline{X}(\omega) = (X_1(\omega), \dots, X_n(\omega))', \omega \in \Omega$ , is an n-dimensional random vector (n-rv) if  $\underline{X}^{-1}(I) = \{\omega : X_1(\omega) \le a_1, \dots, X_n(\omega) \le a_n\} \in L$  for all n-dimensional intervals  $I = \{(x_1, \dots, x_n) : -\infty < x_i \le a_i, a_i \in \mathbb{R} \ \forall i = 1, \dots, n\}$ .

# Note:

It follows that if  $X_1, \ldots, X_n$  are any n rv's on  $(\Omega, L, P)$ , then  $\underline{X} = (X_1, \ldots, X_n)'$  is an n-rv on  $(\Omega, L, P)$  since for any I, it holds:

$$\underline{X}^{-1}(I) = \{\omega : (X_1(\omega), \dots, X_n(\omega)) \in I\}$$

$$= \{\omega : X_1(\omega) \le a_1, \dots, X_n(\omega) \le a_n\}$$

$$= \bigcap_{k=1}^{n} \underbrace{\{\omega : X_k(\omega) \le a_k\}}_{\in L}$$

# Definition 4.1.2:

For an n-rv  $\underline{X}$ , a function F defined by

$$F(x) = P(X \le x) = P(X_1 \le x_1, \dots, X_n \le x_n) \ \forall x \in \mathbb{R}^n$$

is the joint cumulative distribution function (joint cdf) of X.

#### Note:

(i) F is non-decreasing and right-continuous in each of its arguments  $x_i$ .

(ii) 
$$\lim_{\underline{x} \to \underline{\infty} \atop \underline{I\!\!R}} F(\underline{x}) = \lim_{x_1 \to \infty, \dots x_n \to \infty} F(\underline{x}) = 1 \text{ and } \lim_{x_k \to -\infty} F(\underline{x}) = 0 \quad \forall x_1, \dots, x_{k-1}, x_{k+1}, \dots, x_n \in \underline{I\!\!R}$$

However, conditions (i) and (ii) together are not sufficient for F to be a joint cdf. Instead we need the conditions from the next Theorem.

#### Theorem 4.1.3:

A function  $F(\underline{x}) = F(x_1, \dots, x_n)$  is the joint cdf of some n-rv  $\underline{X}$  iff

(i) F is non-decreasing and right-continuous with respect to each  $x_i$ ,

(ii) 
$$F(-\infty, x_2, ..., x_n) = F(x_1, -\infty, x_3, ..., x_n) = ... = F(x_1, ..., x_{n-1}, -\infty) = 0$$
 and  $F(\infty, ..., \infty) = 1$ , and

(iii)  $\forall \underline{x} \in \mathbb{R}^n \ \forall \epsilon_i > 0, i = 1, \dots, n$ , the following inequality holds:

$$F(\underline{x} + \underline{\epsilon}) - \sum_{i=1}^{n} F(x_1 + \epsilon_1, \dots, x_{i-1} + \epsilon_{i-1}, x_i, x_{i+1} + \epsilon_{i+1}, \dots, x_n + \epsilon_n)$$

$$+ \sum_{1 \le i < j \le n} F(x_1 + \epsilon_1, \dots, x_{i-1} + \epsilon_{i-1}, x_i, x_{i+1} + \epsilon_{i+1}, \dots, x_n + \epsilon_n)$$

$$x_{j-1} + \epsilon_{j-1}, x_j, x_{j+1} + \epsilon_{j+1}, \dots, x_n + \epsilon_n)$$

$$\mp \dots$$

$$+ (-1)^n F(\underline{x})$$

$$\geq 0$$

Note:

We won't prove this Theorem but just see why we need condition (iii) for n=2:

$$P(x_1 < X \le x_2, y_1 < Y \le y_2) = P(X \le x_2, Y \le y_2) - P(X \le x_1, Y \le y_2) - P(X \le x_2, Y \le y_1) + P(X \le x_1, Y \le y_1) \ge 0$$

Note:

We will restrict ourselves to n=2 for most of the next Definitions and Theorems but those can be easily generalized to n>2. The term **bivariate**  $\mathbf{rv}$  is often used to refer to a 2-rv and **multivariate**  $\mathbf{rv}$  is used to refer to an n-rv,  $n \geq 2$ .

#### Definition 4.1.4:

A 2-rv (X,Y) is **discrete** if there exists a countable collection  $\mathcal{X}$  of pairs  $(x_i,y_i)$  that has probability 1. Let  $p_{ij} = P(X = x_i, Y = y_j) > 0 \quad \forall (x_i,y_j) \in \mathcal{X}$ . Then,  $\sum_{i,j} p_{ij} = 1$  and  $\{p_{ij}\}$  is the **joint probability mass function (joint pmf)** of (X,Y).

#### Definition 4.1.5:

Let (X, Y) be a discrete 2-rv with joint pmf  $\{p_{ij}\}$ . Define

$$p_{i\cdot} = \sum_{j=1}^{\infty} p_{ij} = \sum_{j=1}^{\infty} P(X = x_i, Y = y_j) = P(X = x_i)$$

and

$$p_{\cdot j} = \sum_{i=1}^{\infty} p_{ij} = \sum_{i=1}^{\infty} P(X = x_i, Y = y_j) = P(Y = y_j).$$

Then  $\{p_{i\cdot}\}$  is called the **marginal probability mass function** (marginal pmf) of X and  $\{p_{\cdot j}\}$  is called the marginal probability mass function of Y.

# Definition 4.1.6:

A 2-rv (X,Y) is **continuous** if there exists a non-negative function f such that

$$F(x,y) = \int_{-\infty}^{x} \int_{-\infty}^{y} f(u,v) \ dv \ du \ \forall (x,y) \in \mathbb{R}^2$$

where F is the joint cdf of (X, Y). We call f the **joint probability density function (joint pdf)** of (X, Y).

#### Note:

If F is continuous at (x, y), then

$$\frac{d^2F(x,y)}{dx\ dy} = f(x,y).$$

# Definition 4.1.7:

Let (X,Y) be a continuous 2-rv with joint pdf f. Then  $f_X(x) = \int_{-\infty}^{\infty} f(x,y) dy$  is called the marginal probability density function (marginal pdf) of X and  $f_Y(y) = \int_{-\infty}^{\infty} f(x,y) dx$  is called the marginal probability density function of Y.

# Note:

(i)

$$\int_{-\infty}^{\infty} f_X(x)dx = \int_{-\infty}^{\infty} \left( \int_{-\infty}^{\infty} f(x,y)dy \right) dx = F(\infty,\infty) = 1 = 0$$

$$\int_{-\infty}^{\infty} \left( \int_{-\infty}^{\infty} f(x,y)dx \right) dy = \int_{-\infty}^{\infty} f_Y(y)dy$$

and  $f_X(x) \ge 0 \ \forall x \in \mathbb{R}$  and  $f_Y(y) \ge 0 \ \forall y \in \mathbb{R}$ .

(ii) Given a 2-rv (X,Y) with joint cdf F(x,y), how do we generate a marginal cdf  $F_X(x) = P(X \le x)$ ? — The answer is  $P(X \le x) = P(X \le x, -\infty < Y < \infty) = F(x,\infty)$ .

#### Definition 4.1.8:

If  $F_{\underline{X}}(x_1,\ldots,x_n)=F_{\underline{X}}(\underline{x})$  is the joint cdf of an n-rv  $\underline{X}=(X_1,\ldots,X_n)$ , then the **marginal** cumulative distribution function (marginal cdf) of  $(X_{i_1},\ldots,X_{i_k})$ ,  $1 \leq k \leq n-1$ ,  $1 \leq i_1 < i_2 < \ldots < i_k \leq n$ , is given by

$$\lim_{x_i \to \infty, i \neq i_1, \dots, i_k} F_{\underline{X}}(\underline{x}) = F_{\underline{X}}(\infty, \dots, \infty, x_{i_1}, \infty, \dots, \infty, x_{i_2}, \infty, \dots, \infty, x_{i_k}, \infty, \dots, \infty).$$

#### Note:

In Definition 1.4.1, we defined conditional probability distributions in some probability space  $(\Omega, L, P)$ . This definition extends to conditional distributions of 2-rv's (X, Y).

# Definition 4.1.9:

Let (X, Y) be a discrete 2-rv. If  $P(Y = y_j) = p_{\cdot j} > 0$ , then the **conditional probability** mass function (conditional pmf) of X given  $Y = y_j$  (for fixed j) is defined as

$$p_{i|j} = P(X = x_i \mid Y = y_j) = \frac{P(X = x_i, Y = y_j)}{P(Y = y_j)} = \frac{p_{ij}}{p_{ij}}.$$

#### Note:

For a continuous 2-rv (X,Y) with pdf f,  $P(X \le x \mid Y = y)$  is not defined. Let  $\epsilon > 0$  and suppose that  $P(y - \epsilon < Y \le y + \epsilon) > 0$ . For every x and every interval  $(y - \epsilon, y + \epsilon]$ , consider the conditional probability of  $X \le x$  given  $Y \in (y - \epsilon, y + \epsilon]$ . We have

$$P(X \le x \mid y - \epsilon < Y \le y + \epsilon) = \frac{P(X \le x, y - \epsilon < Y \le y + \epsilon)}{P(y - \epsilon < Y \le y + \epsilon)}$$

which is well-defined if  $P(y - \epsilon < Y \le y + \epsilon) > 0$  holds.

So, when does

$$\lim_{\epsilon \to 0^+} P(X \le x \mid Y \in (y - \epsilon, y + \epsilon])$$

exist? See the next definition.

#### Definition 4.1.10:

The conditional cumulative distribution function (conditional cdf) of a rv X given that Y = y is defined to be

$$F_{X|Y}(x \mid y) = \lim_{\epsilon \to 0^+} P(X \le x \mid Y \in (y - \epsilon, y + \epsilon])$$

provided that this limit exists. If it does exist, the **conditional probability density func**tion (conditional pdf) of X given that Y = y is any non-negative function  $f_{X|Y}(x \mid y)$  satisfying

$$F_{X|Y}(x \mid y) = \int_{-\infty}^{x} f_{X|Y}(t \mid y) dt \quad \forall x \in \mathbb{R}.$$

Note:

For fixed 
$$y$$
,  $f_{X|Y}(x \mid y) \ge 0$  and  $\int_{-\infty}^{\infty} f_{X|Y}(x \mid y) dx = 1$ . So it is really a pdf.

#### <u>Theorem 4.1.11:</u>

Let (X, Y) be a continuous 2-rv with joint pdf  $f_{X,Y}$ . It holds that at every point (x, y) where f is continuous and the marginal pdf  $f_Y(y) > 0$ , we have

$$F_{X|Y}(x \mid y) = \lim_{\epsilon \to 0+} \frac{P(X \le x, Y \in (y - \epsilon, y + \epsilon])}{P(Y \in (y - \epsilon, y + \epsilon])}$$

$$= \lim_{\epsilon \to 0+} \left( \frac{\frac{1}{2\epsilon} \int_{-\infty}^{x} \int_{y - \epsilon}^{y + \epsilon} f_{X,Y}(u, v) dv \ du}{\frac{1}{2\epsilon} \int_{y - \epsilon}^{y + \epsilon} f_{Y}(v) dv} \right)$$

$$= \frac{\int_{-\infty}^{x} f_{X,Y}(u, y) du}{f_{Y}(y)}$$

$$= \int_{-\infty}^{x} \frac{f_{X,Y}(u, y)}{f_{Y}(y)} du.$$

Thus,  $f_{X|Y}(x \mid y)$  exists and equals  $\frac{f_{X,Y}(x,y)}{f_Y(y)}$ , provided that  $f_Y(y) > 0$ . Furthermore, since

$$\int_{-\infty}^{x} f_{X,Y}(u,y)du = f_Y(y)F_{X\mid Y}(x\mid y),$$

we get the following marginal cdf of X:

$$F_X(x) = \int_{-\infty}^{\infty} \left( \int_{-\infty}^x f_{X,Y}(u, y) du \right) dy = \int_{-\infty}^{\infty} f_Y(y) F_{X|Y}(x \mid y) dy$$

# Example 4.1.12:

Consider

Lecture 26: Fr 10/27/00

$$f_{X,Y}(x,y) = \begin{cases} 2, & 0 < x < y < 1 \\ 0, & \text{otherwise} \end{cases}$$

We calculate the marginal pdf's  $f_X(x)$  and  $f_Y(y)$  first:

$$f_X(x) = \int_{-\infty}^{\infty} f_{X,Y}(x,y)dy = \int_{x}^{1} 2dy = 2(1-x) \text{ for } 0 < x < 1$$

and

$$f_Y(y) = \int_{-\infty}^{\infty} f_{X,Y}(x,y) dx = \int_{0}^{y} 2dx = 2y \text{ for } 0 < y < 1$$

The conditional pdf's  $f_{Y|X}(y\mid x)$  and  $f_{X|Y}(x\mid y)$  are calculated as follows:

$$f_{Y|X}(y \mid x) = \frac{f_{X,Y}(x,y)}{f_X(x)} = \frac{2}{2(1-x)} = \frac{1}{1-x}$$
 for  $x < y < 1$  (where  $0 < x < 1$ )

and

$$f_{X|Y}(x \mid y) = \frac{f_{X,Y}(x,y)}{f_Y(y)} = \frac{2}{2y} = \frac{1}{y} \text{ for } 0 < x < y \text{ (where } 0 < y < 1)$$

Thus, it holds that  $Y \mid X = x \sim U(x,1)$  and  $X \mid Y = y \sim U(0,y)$ , i.e., both conditional pdf's are related to uniform distributions.

# 4.2 Independent Random Variables

Example 4.2.1: (from Rohatgi, page 119, Example 1)

Let  $f_1, f_2, f_3$  be 3 pdf's with cdf's  $F_1, F_2, F_3$  and let  $|\alpha| \leq 1$ . Define

$$f_{\alpha}(x_1, x_2, x_3) = f_1(x_1)f_2(x_2)f_3(x_3) \cdot (1 + \alpha(2F_1(x_1) - 1)(2F_2(x_2) - 1)(2F_3(x_3) - 1)).$$

We can show

- (i)  $f_{\alpha}$  is a pdf for all  $\alpha \in [-1, 1]$ .
- (ii)  $\{f_{\alpha}: -1 \leq \alpha \leq 1\}$  all have marginal pdf's  $f_1, f_2, f_3$ .

See book for proof and further discussion — but when do the marginal distributions uniquely determine the joint distribution?

# <u>Definition 4.2.2:</u>

Let  $F_{X,Y}(x,y)$  be the joint cdf and  $F_X(x)$  and  $F_Y(y)$  be the marginal cdf's of a 2-rv (X,Y). X and Y are **independent** iff

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

#### Lemma 4.2.3:

If X and Y are independent,  $a, b, c, d \in \mathbb{R}$ , and a < b and c < d, then

$$P(a < X < b, c < Y < d) = P(a < X < b)P(c < Y < d).$$

Proof:

$$\begin{split} P(a < X \le b, c < Y \le d) &= F_{X,Y}(b, d) - F_{X,Y}(a, d) - F_{X,Y}(b, c) + F_{X,Y}(a, c) \\ &= F_X(b)F_Y(d) - F_X(a)F_Y(d) - F_X(b)F_Y(c) + F_X(a)F_Y(c) \\ &= (F_X(b) - F_X(a))(F_Y(d) - F_Y(c)) \\ &= P(a < X < b)P(c < Y < d) \end{split}$$

# Definition 4.2.4:

A collection of rv's  $X_1, \ldots, X_n$  with joint cdf  $F_{\underline{X}}(\underline{x})$  and marginal cdf's  $F_{X_i}(x_i)$  are **mutually** (or completely) independent iff

$$F_{\underline{X}}(\underline{x}) = \prod_{i=1}^{n} F_{X_i}(x_i) \ \forall \underline{x} \in I\!\!R^n.$$

#### Note:

We often simply say that the rv's  $X_1, \ldots, X_n$  are **independent** when we really mean that they are mutually independent.

# Theorem 4.2.5: Factorization Theorem

(i) A necessary and sufficient condition for discrete rv's  $X_1, \ldots, X_n$  to be independent is that

$$P(\underline{X} = \underline{x}) = P(X_1 = x_1, \dots, X_n = x_n) = \prod_{i=1}^n P(X_i = x_i) \ \forall \underline{x} \in \mathcal{X}$$

where  $\mathcal{X} \subset \mathbb{R}^n$  is the countable support of  $\underline{X}$ .

(ii) For an absolutely continuous n-rv  $\underline{X} = (X_1, \dots, X_n), X_1, \dots, X_n$  are independent iff

$$f_{\underline{X}}(\underline{x}) = f_{X_1,...,X_n}(x_1,...,x_n) = \prod_{i=1}^n f_{X_i}(x_i),$$

where  $f_{\underline{X}}$  is the joint pdf and  $f_{X_1}, \ldots, f_{X_n}$  are the marginal pdfs of  $\underline{X}$ .

# Proof:

(i) Discrete case:

Let  $\underline{X}$  be a random vector whose components are independent random variables of the discrete type with  $P(\underline{X} = \underline{b}) > 0$ . Lemma 4.2.3 extends to:

$$\lim_{\underline{a}\uparrow\underline{b}} P(\underline{a} < \underline{X} \leq \underline{b}) = \lim_{\substack{a_i\uparrow b_i \ \forall i \in \{1,\dots,n\}}} P(a_1 < X_1 \leq b_1, \ a_2 < X_2 \leq b_2, \ \dots, \ a_n < X_n \leq b_n)$$

$$= P(X_1 = b_1, \ X_2 = b_2, \ \dots, \ X_n = b_n)$$

$$= P(\underline{X} = \underline{b})$$

$$\stackrel{indep.}{=} \prod_{i=1}^n P(X_i = b_i)$$

Before considering the converse factorization of the joint cdf of an n-rv, recall that independence allows each value in the support of a component to combine with each of all possible combinations of the other component values: a 3-dimensional vector where the first component has support  $\{x_{1a}, x_{1b}, x_{1c}\}$ , the second component has support  $\{x_{2a}, x_{2b}, x_{2c}\}$ , and the third component has support  $\{x_{3a}, x_{3b}, x_{3c}\}$ , has  $3 \times 3 \times 3 = 27$  points in its support. Due to independence, all these vectors can be arranged into three sets: the first set having  $x_{1a}$ , the second set having  $x_{1b}$ , and the third set having  $x_{1c}$ , with each set having 9 combinations of  $x_2$  and  $x_3$ .

$$\Rightarrow F(x_{1}, x_{2}, x_{3})$$

$$= P(X_{1} = x_{1a})F(x_{2}, x_{3}) + P(X_{1} = x_{1b})F(x_{2}, x_{3}) + P(X_{1} = x_{1c})F(x_{2}, x_{3})$$

$$= F(x_{1})F(x_{2}, x_{3})$$

$$= F(x_{1})[P(X_{2} = x_{2a})F(x_{3}) + P(X_{2} = x_{2b})F(x_{3}) + P(X_{2} = x_{2c})F(x_{3})]$$

$$= F(x_{1})[P(X_{2} = x_{2a}) + P(X_{2} = x_{2b}) + P(X_{2} = x_{2c})]F(x_{3})$$

$$= F(x_{1})F(x_{2})F(x_{3})$$

More generally, for n dimensions, let  $\underline{x}_i = (x_{i1}, x_{i2}, \dots, x_{in}), B = \{\underline{x}_i : x_{i1} \leq x_1, x_{i2} \leq x_2, \dots, x_{in} \leq x_n\}, B \in \mathcal{X}$ . Then it holds:

$$F_{\underline{X}}(\underline{x}) = \sum_{\underline{x}_i \in B} P(\underline{X} = \underline{x}_i)$$

$$= \sum_{\underline{x}_i \in B} P(X_1 = x_{i1}, X_2 = x_{i2}, \dots, X_n = x_{in})$$

$$\stackrel{indep.}{=} \sum_{\underline{x}_i \in B} P(X_1 = x_{i1}) P(X_2 = x_{i2}, \dots, X_n = x_{in})$$

$$= \sum_{x_{i1} \le x_1, \dots, x_{in} \le x_n} P(X_1 = x_{i1}) P(X_2 = x_{i2}, \dots, X_n = x_{in})$$

$$= \sum_{x_{i1} \le x_1} P(X_1 = x_{i1}) \sum_{x_{i2} \le x_2, \dots, x_{in} \le x_n} P(X_2 = x_{i2}, \dots, X_n = x_{in})$$

$$= F_{X_1}(x_1) \sum_{x_{i2} \le x_2, \dots, x_{in} \le x_n} P(X_2 = x_{i2}, \dots, X_n = x_{in})$$

$$\stackrel{indep.}{=} F_{X_1}(x_1) \sum_{x_{i2} \le x_2} P(X_2 = x_{i2}) \sum_{x_{i3} \le x_3, \dots, x_{in} \le x_n} P(X_3 = x_{i3}, \dots, X_n = x_{in})$$

$$= F_{X_1}(x_1) F_{X_2}(x_2) \sum_{x_{i3} \le x_3, \dots, x_{in} \le x_n} P(X_3 = x_{i3}, \dots, X_n = x_{in})$$

$$= \dots$$

$$= F_{X_1}(x_1) F_{X_2}(x_2) \dots F_{X_n}(x_n)$$

$$= \prod_{i=1}^n F_{X_i}(x_i)$$

(ii) Continuous case: Homework

# <u>Theorem 4.2.6:</u>

Lecture 27: Mo 10/30/00

 $X_1, \ldots, X_n$  are independent iff  $P(X_i \in A_i, i = 1, \ldots, n) = \prod_{i=1}^n P(X_i \in A_i) \ \forall \text{ Borel sets } A_i \in \mathcal{B}$  (i.e., rv's are independent iff all events involving these rv's are independent).

# Proof:

Lemma 4.2.3 and definition of Borel sets.

# Theorem 4.2.7:

Let  $X_1, \ldots, X_n$  be independent rv's and  $g_1, \ldots, g_n$  be Borel-measurable functions. Then  $g_1(X_1), g_2(X_2), \ldots, g_n(X_n)$  are independent.

# Proof:

$$F_{g_{1}(X_{1}),g_{2}(X_{2}),...,g_{n}(X_{n})}(h_{1},h_{2},...,h_{n}) = P(g_{1}(X_{1}) \leq h_{1},g_{2}(X_{2}) \leq h_{2},...,g_{n}(X_{n}) \leq h_{n})$$

$$\stackrel{(*)}{=} P(X_{1} \in g_{1}^{-1}((-\infty,h_{1}]),...,X_{n} \in g_{n}^{-1}((-\infty,h_{n}]))$$

$$\stackrel{Th.4.2.6}{=} \prod_{i=1}^{n} P(X_{i} \in g_{i}^{-1}((-\infty,h_{i}]))$$

$$= \prod_{i=1}^{n} P(g_{i}(X_{i}) \leq h_{i})$$

$$= \prod_{i=1}^{n} F_{g_{i}(X_{i})}(h_{i})$$

(\*) holds since 
$$g_1^{-1}((-\infty, h_1]) \in \mathcal{B}, \dots, g_n^{-1}((-\infty, h_n]) \in \mathcal{B}$$

# Theorem 4.2.8:

If  $X_1, \ldots, X_n$  are independent, then also every subcollection  $X_{i_1}, \ldots, X_{i_k}, k = 2, \ldots, n-1, 1 \le i_1 < i_2 \ldots < i_k \le n$ , is independent.

# Definition 4.2.9:

A set (or a sequence) of rv's  $\{X_n\}_{n=1}^{\infty}$  is independent iff every finite subcollection is independent.

#### Note:

Recall that X and Y are identically distributed iff  $F_X(x) = F_Y(x) \quad \forall x \in \mathbb{R}$  according to Definition 2.2.5 and Theorem 2.2.6.

# Definition 4.2.10:

We say that  $\{X_n\}_{n=1}^{\infty}$  is a set (or a sequence) of **independent identically distributed** (iid) rv's if  $\{X_n\}_{n=1}^{\infty}$  is independent and all  $X_n$  are identically distributed.

# Note:

Recall that X and Y being identically distributed does not say that X = Y with probability 1. If this happens, we say that X and Y are **equivalent** rv's.

#### Note:

We can also extend the defintion of independence to 2 random vectors  $\underline{X}^{n\times 1}$  and  $\underline{Y}^{n\times 1}$ :  $\underline{X}$  and  $\underline{Y}$  are independent iff  $F_{\underline{X},\underline{Y}}(\underline{x},\underline{y}) = F_{\underline{X}}(\underline{x})F_{\underline{Y}}(\underline{y}) \ \ \forall \underline{x},\underline{y} \in I\!\!R^n$ .

This does not mean that the components  $X_i$  of  $\underline{X}$  or the components  $Y_i$  of  $\underline{Y}$  are independent. However, it does mean that each pair of components  $(X_i, Y_i)$  are independent, any subcollections  $(X_{i_1}, \ldots, X_{i_k})$  and  $(Y_{j_1}, \ldots, Y_{j_l})$  are independent, and any Borel-measurable functions  $f(\underline{X})$  and  $g(\underline{Y})$  are independent.

Corollary 4.2.11: (to Factorization Theorem 4.2.5)

If X and Y are independent rv's, then

$$F_{X|Y}(x \mid y) = F_X(x) \ \forall x,$$

and

$$F_{Y|X}(y \mid x) = F_Y(y) \ \forall y.$$

# 4.3 Functions of Random Vectors

# <u>Theorem 4.3.1:</u>

If X and Y are rv's on  $(\Omega, L, P) \to \mathbb{R}$ , then

- (i)  $X \pm Y$  is a rv.
- (ii) XY is a rv.
- (iii) If  $\{\omega : Y(\omega) = 0\} = \emptyset$ , then  $\frac{X}{Y}$  is a rv.

# Theorem 4.3.2:

Let  $X_1, \ldots, X_n$  be rv's on  $(\Omega, L, P) \to I\!\!R$ . Define

$$MAX_n = \max\{X_1, \dots, X_n\} = X_{(n)}$$

by

$$MAX_n(\omega) = \max\{X_1(\omega), \dots, X_n(\omega)\} \ \forall \omega \in \Omega$$

and

$$MIN_n = \min\{X_1, \dots, X_n\} = X_{(1)} = -\max\{-X_1, \dots, -X_n\}$$

by

$$MIN_n(\omega) = \min\{X_1(\omega), \dots, X_n(\omega)\} \ \forall \omega \in \Omega.$$

Then,

- (i)  $MIN_n$  and  $MAX_n$  are rv's.
- (ii) If  $X_1, \ldots, X_n$  are independent, then

$$F_{MAX_n}(z) = P(MAX_n \le z) = P(X_i \le z \ \forall i = 1, ..., n) = \prod_{i=1}^n F_{X_i}(z)$$

and

$$F_{MIN_n}(z) = P(MIN_n \le z) = 1 - P(X_i > z \ \forall i = 1, \dots, n) = 1 - \prod_{i=1}^n (1 - F_{X_i}(z)).$$

(iii) If  $\{X_i\}_{i=1}^n$  are iid rv's with common cdf  $F_X$ , then

$$F_{MAX_n}(z) = F_X^n(z)$$

and

$$F_{MIN_n}(z) = 1 - (1 - F_X(z))^n.$$

If  $F_X$  is absolutely continuous with pdf  $f_X$ , then the pdfs of  $MAX_n$  and  $MIN_n$  are

$$f_{MAX_n}(z) = n \cdot F_X^{n-1}(z) \cdot f_X(z)$$

and

$$f_{MIN_n}(z) = n \cdot (1 - F_X(z))^{n-1} \cdot f_X(z)$$

for all continuity points of  $F_X$ .

# Note:

Using Theorem 4.3.2, it is easy to derive the joint cdf and pdf of  $MAX_n$  and  $MIN_n$  for iid rv's  $\{X_1, \ldots, X_n\}$ . For example, if the  $X_i$ 's are iid with cdf  $F_X$  and pdf  $f_X$ , then the joint pdf of  $MAX_n$  and  $MIN_n$  is

$$f_{MAX_n,MIN_n}(x,y) = \begin{cases} 0, & x \le y \\ n(n-1) \cdot (F_X(x) - F_X(y))^{n-2} \cdot f_X(x) f_X(y), & x > y \end{cases}$$

However, note that  $MAX_n$  and  $MIN_n$  are not independent. See Rohatgi, page 129, Corollary, for more details.

#### Note:

The previous transformations are special cases of the following Theorem 4.3.3.

# <u>Theorem 4.3.3:</u>

If  $g: \mathbb{R}^n \to \mathbb{R}^m$  is a Borel-measurable function (i.e.,  $\forall B \in \mathcal{B}^m: g^{-1}(B) \in \mathcal{B}^n$ ) and if  $\underline{X} = (X_1, \dots, X_n)$  is an n-rv, then  $g(\underline{X})$  is an m-rv.

# Proof:

If 
$$B \in \mathcal{B}^m$$
, then  $\{\omega : g(\underline{X}(\omega)) \in B\} = \{\omega : \underline{X}(\omega) \in g^{-1}(B)\} \in \mathcal{B}^n$ .

Question: How do we handle more general transformations of  $\underline{X}$ ?

# <u>Discrete Case:</u>

Let  $\underline{X} = (X_1, \dots, X_n)$  be a discrete n-rv and  $\mathcal{X} \subset \mathbb{R}^n$  be the countable support of  $\underline{X}$ , i.e.,  $P(\underline{X} \in \mathcal{X}) = 1$  and  $P(\underline{X} = \underline{x}) > 0 \ \forall \underline{x} \in \mathcal{X}$ .

Define  $u_i = g_i(x_1, \ldots, x_n), i = 1, \ldots, n$  to be 1-to-1-mappings of  $\mathcal{X}$  onto B. Let  $\underline{u} = (u_1, \ldots, u_n)'$ . Then

$$P(\underline{U} = \underline{u}) = P(g_1(\underline{X}) = u_1, \dots, g_n(\underline{X}) = u_n) = P(X_1 = h_1(\underline{u}), \dots, X_n = h_n(\underline{u})) \quad \forall \underline{u} \in B$$

where  $x_i = h_i(\underline{u}), i = 1, ..., n$ , is the inverse transformation (and  $P(\underline{U} = \underline{u}) = 0 \ \forall \underline{u} \notin B$ ).

The joint marginal pmf of any subcollection of  $u_i$ 's is now obtained by summing over the other remaining  $u_i$ 's.

 $\frac{\text{Example 4.3.4:}}{\text{Let }X,Y \text{ be iid }} \sim Bin(n,p), 0$ 

Then X = UV and Y = V - 1. So the joint pmf of U, V is

$$P(U = u, V = v) = \binom{n}{uv} p^{uv} (1-p)^{n-uv} \binom{n}{v-1} p^{v-1} (1-p)^{n+1-v}$$
$$= \binom{n}{uv} \binom{n}{v-1} p^{uv+v-1} (1-p)^{2n+1-uv-v}$$

for  $v \in \{1, 2, \dots, n+1\}$  and  $uv \in \{0, 1, \dots, n\}$ .

# Continuous Case:

Let  $\underline{X} = (X_1, \dots, X_n)$  be a continuous n–rv with joint cdf  $F_{\underline{X}}$  and joint pdf  $f_{\underline{X}}$ .

Lecture 28: We 11/01/00

Let

$$\underline{U} = \begin{pmatrix} U_1 \\ \vdots \\ U_n \end{pmatrix} = g(\underline{X}) = \begin{pmatrix} g_1(\underline{X}) \\ \vdots \\ g_n(\underline{X}) \end{pmatrix},$$

i.e.,  $U_i = g_i(\underline{X})$ , be a mapping from  $\mathbb{R}^n$  into  $\mathbb{R}^n$ .

If  $B \in \mathcal{B}^n$ , then

$$P(\underline{U} \in B) = P(\underline{X} \in g^{-1}(B)) = \int \dots \int_{g^{-1}(B)} f_{\underline{X}}(\underline{x}) d(\underline{x}) = \int \dots \int_{g^{-1}(B)} f_{\underline{X}}(\underline{x}) \prod_{i=1}^{n} dx_i$$

where  $g^{-1}(B) = \{ \underline{x} = (x_1, \dots, x_n) \in \mathbb{R}^n : g(\underline{x}) \in B \}.$ 

Suppose we define B as the half-infinite n-dimensional interval

$$B_{\underline{u}} = \{(u'_1, \dots, u'_n) : -\infty < u'_i < u_i \ \forall i = 1, \dots, n\}$$

for any  $\underline{u} \in \mathbb{R}^n$ . Then the joint cdf of  $\underline{U}$  is

$$G(\underline{u}) = P(\underline{U} \in B_{\underline{u}}) = P(g_1(\underline{X}) \le u_1, \dots, g_n(\underline{X}) \le u_n) = \int_{q^{-1}(B_u)}^{\int \dots \int} f_{\underline{X}}(\underline{x}) d(\underline{x}).$$

If G happens to be absolutely continuous, the joint pdf of  $\underline{U}$  will be given by  $f_{\underline{U}}(\underline{u}) = \frac{\partial^n G(\underline{u})}{\partial u_1 \partial u_2 ... \partial u_n}$  at every continuity point of  $f_{\underline{U}}$ .

Under certain conditions, we can write  $f_{\underline{U}}$  in terms of the original pdf  $f_{\underline{X}}$  of  $\underline{X}$  as stated in the next Theorem:

#### Theorem 4.3.5: Multivariate Transformation

Let  $\underline{X} = (X_1, \dots, X_n)$  be a continuous n-rv with joint pdf  $f_{\underline{X}}$ .

(i) Let

$$\underline{U} = \begin{pmatrix} U_1 \\ \vdots \\ U_n \end{pmatrix} = g(\underline{X}) = \begin{pmatrix} g_1(\underline{X}) \\ \vdots \\ g_n(\underline{X}) \end{pmatrix},$$

(i.e.,  $U_i = g_i(\underline{X})$ ) be a 1-to-1-mapping from  $\mathbb{R}^n$  into  $\mathbb{R}^n$ , i.e., there exist inverses  $h_i$ ,  $i = 1, \ldots, n$ , such that  $x_i = h_i(\underline{u}) = h_i(u_1, \ldots, u_n), i = 1, \ldots, n$ , over the range of the transformation g.

- (ii) Assume both g and h are continuous.
- (iii) Assume partial derivatives  $\frac{\partial x_i}{\partial u_j} = \frac{\partial h_i(\underline{u})}{\partial u_j}$ ,  $i, j = 1, \ldots, n$ , exist and are continuous.
- (iv) Assume that the Jacobian of the inverse transformation

$$J = \det\left(\frac{\partial(x_1, \dots, x_n)}{\partial(u_1, \dots, u_n)}\right) = \det\left(\begin{array}{ccc} \frac{\partial x_1}{\partial u_1} & \dots & \frac{\partial x_1}{\partial u_n} \\ \vdots & & \vdots \\ \frac{\partial x_n}{\partial u_1} & \dots & \frac{\partial x_n}{\partial u_n} \end{array}\right) = \left(\begin{array}{ccc} \frac{\partial x_1}{\partial u_1} & \dots & \frac{\partial x_1}{\partial u_n} \\ \vdots & & \vdots \\ \frac{\partial x_n}{\partial u_1} & \dots & \frac{\partial x_n}{\partial u_n} \end{array}\right)$$

is different from 0 for all  $\underline{u}$  in the range of g.

Then the n-rv  $\underline{U} = g(\underline{X})$  has a joint absolutely continuous cdf with corresponding joint pdf

$$f_U(\underline{u}) = |J| f_X(h_1(\underline{u}), \dots, h_n(\underline{u})).$$

Proof:

Let  $\underline{u} \in I\!\!R^n$  and

$$B_{\underline{u}} = \{(u'_1, \dots, u'_n) : -\infty < u'_i < u_i \ \forall i = 1, \dots, n\}.$$

Then,

$$G_{\underline{U}}(\underline{u}) = \int \dots \int \int f_{\underline{u}} f_{\underline{x}}(\underline{x}) d(\underline{x})$$

$$= \int \dots \int \int f_{\underline{x}}(h_1(\underline{u}), \dots, h_n(\underline{u})) \mid J \mid d(\underline{u})$$

The result follows from differentiation of  $G_U$ .

For additional steps of the proof see Rohatgi (page 135 and Theorem 17 on page 10) or a book on multivariate calculus.

# Theorem 4.3.6:

Let  $\underline{X} = (X_1, \dots, X_n)$  be a continuous n–rv with joint pdf  $f_{\underline{X}}$ .

(i) Let

$$\underline{U} = \begin{pmatrix} U_1 \\ \vdots \\ U_n \end{pmatrix} = g(\underline{X}) = \begin{pmatrix} g_1(\underline{X}) \\ \vdots \\ g_n(\underline{X}) \end{pmatrix},$$

(i.e.,  $U_i = g_i(\underline{X})$ ) be a mapping from  $\mathbb{R}^n$  into  $\mathbb{R}^n$ .

- (ii) Let  $\mathcal{X} = \{\underline{x} : f_{\underline{X}}(\underline{x}) > 0\}$  be the support of  $\underline{X}$ .
- (iii) Suppose that for each  $\underline{u} \in B = \{\underline{u} \in \mathbb{R}^n : \underline{u} = g(\underline{x}) \text{ for some } \underline{x} \in \mathcal{X} \}$  there is a finite number  $k = k(\underline{u})$  of inverses.
- (iv) Suppose we can partition  $\mathcal{X}$  into  $\mathcal{X}_0, \mathcal{X}_1, \dots, \mathcal{X}_k$  s.t.
  - (a)  $P(X \in \mathcal{X}_0) = 0$ .
  - (b)  $\underline{U} = g(\underline{X})$  is a 1-to-1-mapping from  $\mathcal{X}_l$  onto B for all  $l = 1, \ldots, k$ , with inverse transformation  $h_l(\underline{u}) = \begin{pmatrix} h_{l1}(\underline{u}) \\ \vdots \\ h_{ln}(\underline{u}) \end{pmatrix}$ ,  $\underline{u} \in B$ , i.e., for each  $\underline{u} \in B$ ,  $h_l(\underline{u})$  is the unique  $\underline{x} \in \mathcal{X}_l$  such that  $\underline{u} = g(\underline{x})$ .
- (v) Assume partial derivatives  $\frac{\partial x_i}{\partial u_j} = \frac{\partial h_{li}(\underline{u})}{\partial u_j}$ ,  $l = 1, \ldots, k, i, j = 1, \ldots, n$ , exist and are continuous.
- (vi) Assume the Jacobian of each of the inverse transformations

$$J_{l} = \det \begin{pmatrix} \frac{\partial x_{1}}{\partial u_{1}} & \cdots & \frac{\partial x_{1}}{\partial u_{n}} \\ \vdots & & \vdots \\ \frac{\partial x_{n}}{\partial u_{1}} & \cdots & \frac{\partial x_{n}}{\partial u_{n}} \end{pmatrix} = \det \begin{pmatrix} \frac{\partial h_{l1}}{\partial u_{1}} & \cdots & \frac{\partial h_{l1}}{\partial u_{n}} \\ \vdots & & \vdots \\ \frac{\partial h_{ln}}{\partial u_{1}} & \cdots & \frac{\partial h_{ln}}{\partial u_{n}} \end{pmatrix} = \begin{vmatrix} \frac{\partial h_{l1}}{\partial u_{1}} & \cdots & \frac{\partial h_{l1}}{\partial u_{n}} \\ \vdots & & \vdots \\ \frac{\partial h_{ln}}{\partial u_{1}} & \cdots & \frac{\partial h_{ln}}{\partial u_{n}} \end{vmatrix}, \quad l = 1, \dots, k,$$

is different from 0 for all  $\underline{u}$  in the range of g.

Then the joint pdf of  $\underline{U}$  is given by

$$f_{\underline{U}}(\underline{u}) = \sum_{l=1}^{k} |J_l| f_{\underline{X}}(h_{l1}(\underline{u}), \dots, h_{ln}(\underline{u})).$$

# Example 4.3.7:

Let X, Y be iid  $\sim N(0, 1)$ . Define

$$U = g_1(X, Y) = \begin{cases} \frac{X}{Y}, & Y \neq 0\\ 0, & Y = 0 \end{cases}$$

and

$$V = g_2(X, Y) = |Y|.$$

 $\mathcal{X} = \mathbb{R}^2$ , but U, V are not 1-to-1 mappings since (U, V)(x, y) = (U, V)(-x, -y), i.e., conditions do not apply for the use of Theorem 4.3.5. Let

$$\mathcal{X}_0 = \{(x,y) : y = 0\}$$
  
 $\mathcal{X}_1 = \{(x,y) : y > 0\}$   
 $\mathcal{X}_2 = \{(x,y) : y < 0\}$ 

Then  $P((X,Y) \in \mathcal{X}_0) = 0$ .

Let 
$$B = \{(u, v) : v > 0\} = g(\mathcal{X}_1) = g(\mathcal{X}_2).$$

Inverses:

$$B \to \mathcal{X}_1 : x = h_{11}(u, v) = uv$$

$$y = h_{12}(u, v) = v$$

$$B \to \mathcal{X}_2 : x = h_{21}(u, v) = -uv$$

$$y = h_{22}(u, v) = -v$$

$$J_1 = \begin{vmatrix} v & u \\ 0 & 1 \end{vmatrix} \Rightarrow |J_1| = |v|$$

$$J_2 = \begin{vmatrix} -v & -u \\ 0 & -1 \end{vmatrix} \Rightarrow |J_2| = |v|$$

$$f_{X,Y}(x, y) = \frac{1}{2\pi} e^{-x^2/2} e^{-y^2/2}$$

$$f_{U,V}(u, v) = |v| \frac{1}{2\pi} e^{-(uv)^2/2} e^{-v^2/2} + |v| \frac{1}{2\pi} e^{-(-uv)^2/2} e^{-(-v)^2/2}$$

$$= \frac{v}{\pi} e^{\frac{-(u^2+1)v^2}{2}}, \quad -\infty < u < \infty, \quad 0 < v < \infty$$

Marginal:

$$f_U(u) = \int_0^\infty \frac{v}{\pi} e^{\frac{-(u^2+1)v^2}{2}} dv \quad | \quad z = \frac{(u^2+1)v^2}{2}, \quad \frac{dz}{dv} = (u^2+1)v$$
$$= \int_0^\infty \frac{1}{\pi(u^2+1)} e^{-z} dz$$

$$= \frac{1}{\pi(u^2+1)}(-e^{-z})\Big|_0^{\infty}$$
$$= \frac{1}{\pi(1+u^2)}, \quad -\infty < u < \infty$$

Thus, the ratio of two iid N(0,1) rv's is a rv that has a Cauchy distribution.

Lecture 29: Fr 11/03/00

# 4.4 Order Statistics

# Definition 4.4.1:

Let  $(X_1, ..., X_n)$  be an n-rv. The  $k^{th}$  **order statistic**  $X_{(k)}$  is the  $k^{th}$  smallest of the  $X_i's$ , i.e.,  $X_{(1)} = \min\{X_1, ..., X_n\}, \ X_{(2)} = \min\{\{X_1, ..., X_n\} - X_{(1)}\}, ..., \ X_{(n)} = \max\{X_1, ..., X_n\}.$  It is  $X_{(1)} \leq X_{(2)} \leq ... \leq X_{(n)}$  and  $\{X_{(1)}, X_{(2)}, ..., X_{(n)}\}$  is the set of order statistics for  $(X_1, ..., X_n)$ .

# Note:

As shown in Theorem 4.3.2,  $X_{(1)}$  and  $X_{(n)}$  are rv's. This result will be extended in the following Theorem:

# Theorem 4.4.2:

Let  $(X_1, \ldots, X_n)$  be an n-rv. Then the  $k^{th}$  order statistic  $X_{(k)}, k = 1, \ldots, n$ , is also an rv.

# <u>Theorem 4.4.3:</u>

Let  $X_1, \ldots, X_n$  be continuous iid rv's with pdf  $f_X$ . The joint pdf of  $X_{(1)}, \ldots, X_{(n)}$  is

$$f_{X_{(1)},...,X_{(n)}}(x_1,...,x_n) = \begin{cases} n! & \prod_{i=1}^n f_X(x_i), & x_1 \le x_2 \le ... \le x_n \\ 0, & \text{otherwise} \end{cases}$$

# Proof:

For the case n = 3, look at the following scenario how  $X_1, X_2$ , and  $X_3$  can be possibly ordered to yield  $X_{(1)} < X_{(2)} < X_{(3)}$ . Columns represent  $X_{(1)}, X_{(2)}$ , and  $X_{(3)}$ . Rows represent  $X_1, X_2$ , and  $X_3$ :

$$X_1 < X_2 < X_3$$
 :  $\begin{vmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{vmatrix}$ 
 $X_1 < X_3 < X_2$  :  $\begin{vmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{vmatrix}$ 
 $X_2 < X_1 < X_3$  :  $\begin{vmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{vmatrix}$ 

$$X_2 < X_3 < X_1$$
 :  $egin{bmatrix} 0 & 0 & 1 \ 1 & 0 & 0 \ 0 & 1 & 0 \ \end{bmatrix}$ 
 $X_3 < X_1 < X_2$  :  $egin{bmatrix} 0 & 1 & 0 \ 0 & 0 & 1 \ 1 & 0 & 0 \ \end{bmatrix}$ 
 $X_3 < X_2 < X_1$  :  $egin{bmatrix} 0 & 0 & 1 \ 0 & 1 & 0 \ 1 & 0 & 0 \ \end{bmatrix}$ 

For n = 3, there are 3! = 6 possible arrangements.

In general, there are n! arrangements of  $X_1, \ldots, X_n$  for each  $(X_{(1)}, \ldots, X_{(n)})$ . This mapping is not 1-to-1. For each mapping, we have a  $n \times n$  matrix J that results from an  $n \times n$  identity matrix through the rearrangement of rows. Therefore, |J| = 1. By Theorem 4.3.6, we get

$$f_{X_{(1)},...,X_{(n)}}(x_{(1)},...,x_{(n)}) = n! f_{X_1,...,X_n}(x_{(k_1)},x_{(k_2)},...,x_{(k_n)})$$

$$= n! \prod_{i=1}^n f_{X_i}(x_{(k_i)})$$

$$= n! \prod_{i=1}^n f_X(x_i)$$

Theorem 4.4.4: Let  $X_1, \ldots, X_n$  be continuous iid rv's with pdf  $f_X$  and cdf  $F_X$ . Then the following holds:

(i) The marginal pdf of  $X_{(k)}$ , k = 1, ..., n, is

$$f_{X_{(k)}}(x) = \frac{n!}{(k-1)!(n-k)!} (F_X(x))^{k-1} (1 - F_X(x))^{n-k} f_X(x).$$

(ii) The joint pdf of  $X_{(j)}$  and  $X_{(k)}, 1 \leq j < k \leq n$ , is

$$f_{X_{(j)},X_{(k)}}(x_j,x_k) = \frac{n!}{(j-1)!(k-j-1)!(n-k)!} \times (F_X(x_j))^{j-1} (F_X(x_k) - F_X(x_j))^{k-j-1} (1 - F_X(x_k))^{n-k} f_X(x_j) f_X(x_k)$$

if  $x_j < x_k$  and 0 otherwise.

# 4.5 Multivariate Expectation

In this section, we assume that  $\underline{X} = (X_1, \dots, X_n)$  is an n-rv and  $g : \mathbb{R}^n \to \mathbb{R}^n$  is a Borel-measurable function.

# Definition 4.5.1:

If n = 1, i.e., g is univariate, we define the following:

- (i) Let  $\underline{X}$  be discrete with joint pmf  $p_{i_1,...,i_n} = P(X_1 = x_{i_1}, \ldots, X_n = x_{i_n})$ . If  $\sum_{\substack{i_1,...,i_n \ \text{and this value exists.}}} p_{i_1,...,i_n} \cdot |g(x_{i_1},\ldots,x_{i_n})| < \infty$ , we define  $E(g(\underline{X})) = \sum_{\substack{i_1,...,i_n \ \text{and this value exists.}}} p_{i_1,...,i_n} \cdot g(x_{i_1},\ldots,x_{i_n})$
- (ii) Let  $\underline{X}$  be continuous with joint pdf  $f_{\underline{X}}(\underline{x})$ . If  $\int_{\mathbb{R}^n} |g(\underline{x})| f_{\underline{X}}(\underline{x}) d\underline{x} < \infty$ , we define  $E(g(\underline{X})) = \int_{\mathbb{R}^n} g(\underline{x}) f_{\underline{X}}(\underline{x}) d\underline{x}$  and this value exists.

# Note:

The above can be extended to vector-valued functions g (n > 1) in the obvious way. For example, if g is the identity mapping from  $\mathbb{R}^n \to \mathbb{R}^n$ , then

$$E(\underline{X}) = \begin{pmatrix} E(X_1) \\ \vdots \\ E(X_n) \end{pmatrix} = \begin{pmatrix} \mu_1 \\ \vdots \\ \mu_n \end{pmatrix}$$

provided that  $E(\mid X_i \mid) < \infty \ \forall i = 1, ..., n$ .

Similarly, provided that all expectations exist, we get for the variance-covariance matrix:

$$Var(\underline{X}) = \Sigma_{\underline{X}} = E((\underline{X} - E(\underline{X})) (\underline{X} - E(\underline{X}))')$$

with  $(i, j)^{th}$  component

$$E((X_i - E(X_i)) (X_j - E(X_j))) = Cov(X_i, X_j)$$

and with  $(i, i)^{th}$  component

$$E((X_i - E(X_i)) (X_i - E(X_i))) = Var(X_i) = \sigma_i^2.$$

Joint higher-order moments can be defined similarly when needed.

#### Note:

We are often interested in (weighted) sums of rv's or products of rv's and their expectations.

This will be addressed in the next two Theorems.

# Theorem 4.5.2:

Let  $X_i, i = 1, ..., n$ , be rv's such that  $E(|X_i|) < \infty$ . Let  $a_1, ..., a_n \in \mathbb{R}$  and define  $S = \sum_{i=1}^n a_i X_i$ . Then it holds that  $E(|S|) < \infty$  and

$$E(S) = \sum_{i=1}^{n} a_i E(X_i).$$

# Proof:

Continuous case only:

$$E(|S|) = \int_{\mathbb{R}^{n}} |\sum_{i=1}^{n} a_{i}x_{i}| f_{\underline{X}}(\underline{x}) d\underline{x}$$

$$\leq \int_{\mathbb{R}^{n}} \sum_{i=1}^{n} |a_{i}| \cdot |x_{i}| f_{\underline{X}}(\underline{x}) d\underline{x}$$

$$= \sum_{i=1}^{n} |a_{i}| \int_{\mathbb{R}} |x_{i}| \left( \int_{\mathbb{R}^{n-1}} f_{\underline{X}}(\underline{x}) dx_{1} \dots dx_{i-1} dx_{i+1} \dots dx_{n} \right) dx_{i}$$

$$= \sum_{i=1}^{n} |a_{i}| \int_{\mathbb{R}} |x_{i}| f_{X_{i}}(x_{i}) dx_{i}$$

$$= \sum_{i=1}^{n} |a_{i}| E(|X_{i}|)$$

$$\leq \infty$$

It follows that  $E(S) = \sum_{i=1}^{n} a_i E(X_i)$  by the same argument without using the absolute values  $|\cdot|$ .

#### Note:

If  $X_i$ , i = 1, ..., n, are iid with  $E(X_i) = \mu$ , then

$$E(\overline{X}) = E(\frac{1}{n} \sum_{i=1}^{n} X_i) = \sum_{i=1}^{n} \frac{1}{n} E(X_i) = \mu.$$

# Lecture 30:

# <u>Theorem 4.5.3:</u>

Let  $X_i, i=1,\ldots,n$ , be independent rv's such that  $E(\mid X_i\mid)<\infty$ . Let  $g_i, i=1,\ldots,n$ , be Borel-measurable functions. Then

$$E(\prod_{i=1}^{n} g_i(X_i)) = \prod_{i=1}^{n} E(g_i(X_i))$$

if all expectations exist.

# Proof:

By Theorem 4.2.5,  $f_{\underline{X}}(\underline{x}) = \prod_{i=1}^{n} f_{X_i}(x_i)$ , and by Theorem 4.2.7,  $g_i(X_i)$ ,  $i = 1, \ldots, n$ , are also independent. Therefore,

$$E(\prod_{i=1}^{n} g_{i}(X_{i})) = \int_{\mathbb{R}^{n}} \prod_{i=1}^{n} g_{i}(x_{i}) f_{\underline{X}}(\underline{x}) d\underline{x}$$

$$Th \stackrel{4.2.5}{=} \int_{\mathbb{R}^{n}} \prod_{i=1}^{n} (g_{i}(x_{i}) f_{X_{i}}(x_{i}) dx_{i})$$

$$= \int_{\mathbb{R}} \dots \int_{\mathbb{R}} \prod_{i=1}^{n} g_{i}(x_{i}) \prod_{i=1}^{n} f_{X_{i}}(x_{i}) \prod_{i=1}^{n} dx_{i}$$

$$Th \stackrel{4.2.7}{=} \int_{\mathbb{R}} g_{1}(x_{1}) f_{X_{1}}(x_{1}) dx_{1} \int_{\mathbb{R}} g_{2}(x_{2}) f_{X_{2}}(x_{2}) dx_{2} \dots \int_{\mathbb{R}} g_{n}(x_{n}) f_{X_{n}}(x_{n}) dx_{n}$$

$$= \prod_{i=1}^{n} \int_{\mathbb{R}} g_{i}(x_{i}) f_{X_{i}}(x_{i}) dx_{i}$$

$$= \prod_{i=1}^{n} E(g_{i}(X_{i}))$$

# Corollary 4.5.4:

If X, Y are independent, then Cov(X, Y) = 0.

# Theorem 4.5.5:

Two rv's X, Y are independent iff for all pairs of Borel-measurable functions  $g_1$  and  $g_2$  it holds that  $E(g_1(X) \cdot g_2(Y)) = E(g_1(X)) \cdot E(g_2(Y))$  if all expectations exist.

#### <u>Proof:</u>

" $\Longrightarrow$ ": It follows from Theorem 4.5.3 and the independence of X and Y that

$$E(g_1(X)g_2(Y)) = E(g_1(X)) \cdot E(g_2(Y)).$$

" $\Leftarrow$ ": From Theorem 4.2.6, we know that X and Y are independent iff

$$P(X \in A_1, Y \in A_2) = P(X \in A_1) P(Y \in A_2) \ \forall \text{ Borel sets } A_1 \text{ and } A_2.$$

How do we relate Theorem 4.2.6 to  $g_1$  and  $g_2$ ? Let us define two Borel-measurable functions  $g_1$  and  $g_2$  as:

$$g_1(x) = I_{A_1}(x) = \begin{cases} 1, & x \in A_1 \\ 0, & \text{otherwise} \end{cases}$$

$$g_2(y) = I_{A_2}(y) = \begin{cases} 1, & y \in A_2 \\ 0, & \text{otherwise} \end{cases}$$

Then,

$$E(g_1(X)) = 0 \cdot P(X \in A_1^c) + 1 \cdot P(X \in A_1) = P(X \in A_1),$$
  

$$E(g_2(Y)) = 0 \cdot P(Y \in A_2^c) + 1 \cdot P(Y \in A_2) = P(Y \in A_2)$$

and

$$E(g_1(X) \cdot g_2(Y)) = P(X \in A_1, Y \in A_2).$$

$$\implies P(X \in A_1, Y \in A_2) = E(g_1(X) \cdot g_2(Y))$$

$$\stackrel{given}{=} E(g_1(X)) \cdot E(g_2(Y))$$

$$= P(X \in A_1)P(Y \in A_2)$$

 $\implies X$ , Y independent by Theorem 4.2.6.

#### Definition 4.5.6:

The  $(i_1^{th}, i_2^{th}, \dots, i_n^{th})$  multi-way moment of  $\underline{X} = (X_1, \dots, X_n)$  is defined as

$$m_{i_1 i_2 \dots i_n} = E(X_1^{i_1} X_2^{i_2} \dots X_n^{i_n})$$

if it exists.

The  $(i_1^{th}, i_2^{th}, \dots, i_n^{th})$  multi-way central moment of  $\underline{X} = (X_1, \dots, X_n)$  is defined as

$$\mu_{i_1 i_2 \dots i_n} = E(\prod_{j=1}^n (X_j - E(X_j))^{i_j})$$

if it exists.

#### Note:

If we set  $i_r = i_s = 1$  and  $i_j = 0 \ \forall j \neq r, s$  in Definition 4.5.6 for the multi-way central moment, we get

$$\mu_{0...0} \quad 1 \quad 0...0 \quad 1 \quad 0...0 \quad = \mu_{rs} = Cov(X_r, X_s).$$

# Theorem 4.5.7: Cauchy-Schwarz Inequality

Let X, Y be 2 rv's with finite variance. Then it holds:

- (i) Cov(X,Y) exists.
- (ii)  $(E(XY))^2 \le E(X^2)E(Y^2)$ .
- (iii)  $(E(XY))^2 = E(X^2)E(Y^2)$  iff there exists an  $(\alpha, \beta) \in \mathbb{R}^2 \{(0, 0)\}$  such that  $P(\alpha X + \beta Y = 0) = 1$ .

# Proof:

Assumptions: Var(X),  $Var(Y) < \infty$ . Then also  $E(X^2)$ , E(X),  $E(Y^2)$ ,  $E(Y) < \infty$ .

Result used in proof:

$$0 \le (a-b)^2 = a^2 - 2ab + b^2 \Longrightarrow ab \le \frac{a^2 + b^2}{2}$$
$$0 \le (a+b)^2 = a^2 + 2ab + b^2 \Longrightarrow -ab \le \frac{a^2 + b^2}{2}$$
$$\Longrightarrow |ab| \le \frac{a^2 + b^2}{2} \ \forall \ a, b \in \mathbb{R} \quad (*)$$

(i)

$$E(|XY|) = \int_{\mathbb{R}^{2}} |xy| f_{X,Y}(x,y) dx dy$$

$$\stackrel{(*)}{\leq} \int_{\mathbb{R}^{2}} \frac{x^{2} + y^{2}}{2} f_{X,Y}(x,y) dx dy$$

$$= \int_{\mathbb{R}^{2}} \frac{x^{2}}{2} f_{X,Y}(x,y) dx dy + \int_{\mathbb{R}^{2}} \frac{y^{2}}{2} f_{X,Y}(x,y) dy dx$$

$$= \int_{\mathbb{R}} \frac{x^{2}}{2} f_{X}(x) dx + \int_{\mathbb{R}} \frac{y^{2}}{2} f_{Y}(y) dy$$

$$= \frac{E(X^{2}) + E(Y^{2})}{2}$$

$$< \infty$$

$$\Longrightarrow E(XY)$$
 exists

$$\implies Cov(X,Y) = E(XY) - E(X)E(Y)$$
 exists

(ii) 
$$0 \le E((\alpha X + \beta Y)^2) = \alpha^2 E(X^2) + 2\alpha \beta E(XY) + \beta^2 E(Y^2) \ \forall \ \alpha, \beta \in \mathbb{R}$$
 (A)

If  $E(X^2) = 0$ , then X has a degenerate 1-point Dirac distribution and the inequality trivially is true. Therefore, we can assume that  $E(X^2) > 0$ . As (A) is true for all  $\alpha, \beta \in \mathbb{R}$ , we can choose  $\alpha = \frac{-E(XY)}{E(X^2)}$ ,  $\beta = 1$ .

$$\implies \frac{(E(XY))^2}{E(X^2)} - 2\frac{(E(XY))^2}{E(X^2)} + E(Y^2) \ge 0$$

$$\implies -(E(XY))^2 + E(Y^2)E(X^2) \ge 0$$

$$\implies (E(XY))^2 < E(X^2) E(Y^2)$$

(iii) When are the left and right sides of the inequality in (ii) equal?

Assume that  $E(X^2) > 0$ .  $(E(XY))^2 = E(X^2)E(Y^2)$  holds iff  $E((\alpha X + \beta Y)^2) = 0$  based on (ii). It is therefore sufficient to show that  $E((\alpha X + \beta Y)^2) = 0$  iff  $P(\alpha X + \beta Y) = 0 = 1$ :

" $\Longrightarrow$ ":

Let  $Z = \alpha X + \beta Y$ . Since  $E((\alpha X + \beta Y)^2) = E(Z^2) = Var(Z) + (E(Z))^2 = 0$  and  $Var(Z) \ge 0$  and  $(E(Z))^2 \ge 0$ , it follows that E(Z) = 0 and Var(Z) = 0.

This means that Z has a degenerate 1-point Dirac(0) distribution with  $P(Z=0)=P(\alpha X+\beta Y=0)=1.$ 

"⇐=":

If  $P(\alpha X + \beta Y = 0) = P(Y = -\frac{\alpha}{\beta}X) = 1$  for some  $(\alpha, \beta) \in \mathbb{R}^2 - \{(0, 0)\}$ , i.e., Y is linearly dependent on X with probability 1, this implies:

$$(E(XY))^{2} = (E(X \cdot \frac{-\alpha X}{\beta}))^{2} = (\frac{\alpha}{\beta})^{2} (E(X^{2}))^{2} = E(X^{2})(\frac{\alpha}{\beta})^{2} E(X^{2}) = E(X^{2})E(Y^{2})$$

Lecture 31: We 11/08/00

# 4.6 Multivariate Generating Functions

# Definition 4.6.1:

Let  $\underline{X} = (X_1, \dots, X_n)$  be an n-rv. We define the **multivariate moment generating function** (**mmgf**) of  $\underline{X}$  as

$$M_{\underline{X}}(\underline{t}) = E(e^{\underline{t}'\underline{X}}) = E\left(\exp\left(\sum_{i=1}^{n} t_i X_i\right)\right)$$

if this expectation exists for  $|\underline{t}| = \sqrt{\sum_{i=1}^{n} t_i^2} < h$  for some h > 0.

# Definition 4.6.2:

Let  $\underline{X} = (X_1, \dots, X_n)$  be an n-rv. We define the *n*-dimensional characteristic function  $\Phi_X : \mathbb{R}^n \to \mathcal{C}$  of  $\underline{X}$  as

$$\Phi_{\underline{X}}(\underline{t}) = E(e^{i\underline{t}'\underline{X}}) = E\left(\exp\left(i\sum_{j=1}^n t_j X_j\right)\right).$$

Note:

- (i)  $\Phi_X(\underline{t})$  exists for any real-valued n-rv.
- (ii) If  $M_{\underline{X}}(\underline{t})$  exists, then  $\Phi_{\underline{X}}(\underline{t}) = M_{\underline{X}}(i\underline{t})$ .

<u>Theorem 4.6.3:</u>

- (i) If  $M_{\underline{X}}(\underline{t})$  exists, it is unique and uniquely determines the joint distribution of  $\underline{X}$ .  $\Phi_{\underline{X}}(\underline{t})$  is also unique and uniquely determines the joint distribution of  $\underline{X}$ .
- (ii)  $M_{\underline{X}}(\underline{t})$  (if it exists) and  $\Phi_{\underline{X}}(\underline{t})$  uniquely determine all marginal distributions of  $\underline{X}$ , i.e.,  $M_{X_i}(t_i) = M_{\underline{X}}(\underline{0}, t_i, \underline{0})$  and and  $\Phi_{X_i}(t_i) = \Phi_{\underline{X}}(\underline{0}, t_i, \underline{0})$ .
- (iii) Joint moments of all orders (if they exist) can be obtained as

$$m_{i_1...i_n} = \frac{\partial^{i_1+i_2+...+i_n}}{\partial t_1^{i_1} \partial t_2^{i_2} \dots \partial t_n^{i_n}} M_{\underline{X}}(\underline{t}) \bigg|_{\underline{t}=\underline{0}} = E(X_1^{i_1} X_2^{i_2} \dots X_n^{i_n})$$

if the mmgf exists and

$$m_{i_1...i_n} = \frac{1}{i^{i_1+i_2+...+i_n}} \frac{\partial^{i_1+i_2+...+i_n}}{\partial t_1^{i_1} \partial t_2^{i_2} \dots \partial t_n^{i_n}} \Phi_{\underline{X}}(\underline{0}) = E(X_1^{i_1} X_2^{i_2} \dots X_n^{i_n}).$$

(iv)  $X_1, \ldots, X_n$  are independent rv's iff

$$M_{\underline{X}}(t_1,\ldots,t_n) = M_{\underline{X}}(t_1,\underline{0}) \cdot M_{\underline{X}}(0,t_2,\underline{0}) \cdot \ldots \cdot M_{\underline{X}}(\underline{0},t_n) \quad \forall t_1,\ldots,t_n \in \mathbb{R},$$

given that  $M_{\underline{X}}(\underline{t})$  exists.

Similarly,  $X_1, \ldots, X_n$  are independent rv's iff

$$\Phi_X(t_1,\ldots,t_n) = \Phi_X(t_1,\underline{0}) \cdot \Phi_X(0,t_2,\underline{0}) \cdot \ldots \cdot \Phi_X(\underline{0},t_n) \quad \forall t_1,\ldots,t_n \in \mathbb{R}.$$

# Proof:

Rohatgi, page 162: Theorem 7, Corollary, Theorem 8, and Theorem 9 (for mmgf and the case n=2).

# <u>Theorem 4.6.4:</u>

Let  $X_1, \ldots, X_n$  be independent rv's.

(i) If mgf's  $M_{X_1}(t), \ldots, M_{X_n}(t)$  exist, then the mgf of  $Y = \sum_{i=1}^n a_i X_i$  is

$$M_Y(t) = \prod_{i=1}^n M_{X_i}(a_i t)$$

on the common interval where all individual mgf's exist.

(ii) The characteristic function of  $Y = \sum_{j=1}^{n} a_j X_j$  is

$$\Phi_Y(t) = \prod_{j=1}^n \Phi_{X_j}(a_j t)$$

(iii) If mgf's  $M_{X_1}(t), \dots, M_{X_n}(t)$  exist, then the mmgf of  $\underline{X}$  is

$$M_{\underline{X}}(\underline{t}) = \prod_{i=1}^n M_{X_i}(t_i)$$

on the common interval where all individual mgf's exist.

(iv) The n-dimensional characteristic function of  $\underline{X}$  is

$$\Phi_{\underline{X}}(\underline{t}) = \prod_{j=1}^{n} \Phi_{X_j}(t_j).$$

#### Proof:

Homework (parts (ii) and (iv) only)

# <u>Theorem 4.6.5:</u>

Let  $X_1, \ldots, X_n$  be independent discrete rv's on the non-negative integers with pgf's  $G_{X_1}(s), \ldots, G_{X_n}(s)$ .

The pgf of 
$$Y = \sum_{i=1}^{n} X_i$$
 is

$$G_Y(s) = \prod_{i=1}^n G_{X_i}(s).$$

# Proof:

# Version 1:

$$G_{X_i}(s) = E(s^{X_i})$$

$$G_Y(s) = E(s^Y)$$

$$= E(s^{\sum_{i=1}^n X_i})$$

$$\stackrel{indep.}{=} \prod_{i=1}^n E(s^{X_i})$$

$$= \prod_{i=1}^n G_{X_i}(s)$$

Version 2: (case n=2 only)

$$\begin{split} G_Y(s) &= P(Y=0) + P(Y=1)s + P(Y=2)s^2 + \dots \\ &= P(X_1=0,X_2=0) + \\ &\quad (P(X_1=1,X_2=0) + P(X_1=0,X_2=1)) \, s + \\ &\quad (P(X_1=2,X_2=0) + P(X_1=1,X_2=1) + P(X_1=0,X_2=2)) \, s^2 + \dots \\ &\stackrel{indep.}{=} P(X_1=0)P(X_2=0) + \\ &\quad (P(X_1=1)P(X_2=0) + P(X_1=0)P(X_2=1)) \, s + \\ &\quad (P(X_1=2)P(X_2=0) + P(X_1=1)P(X_2=1) + P(X_1=0)P(X_2=2)) \, s^2 + \dots \\ &= \left(P(X_1=0) + P(X_1=1)s + P(X_1=2)s^2 + \dots\right) \cdot \\ &\quad \left(P(X_2=0) + P(X_2=1)s + P(X_2=2)s^2 + \dots\right) \\ &= G_{X_1}(s) \cdot G_{X_2}(s) \end{split}$$

A generalized proof for  $n \geq 3$  needs to be done by induction on n.

#### Theorem 4.6.6:

Let  $X_1, \ldots, X_N$  be iid discrete rv's on the non-negative integers with common pgf  $G_X(s)$ . Let N be a discrete rv on the non-negative integers with pgf  $G_N(s)$ . Let N be independent of the  $X_i$ 's. Define  $S_N = \sum_{i=1}^N X_i$ . The pgf of  $S_N$  is

$$G_{S_N}(s) = G_N(G_X(s)).$$

Proof:

$$P(S_N = k) = \sum_{n=0}^{\infty} P(S_N = k | N = n) \cdot P(N = n)$$

$$\Longrightarrow G_{S_N}(s) = \sum_{k=0}^{\infty} \sum_{n=0}^{\infty} P(S_N = k | N = n) \cdot P(N = n) \cdot s^k$$

$$= \sum_{n=0}^{\infty} P(N = n) \sum_{k=0}^{\infty} P(S_N = k | N = n) \cdot s^k$$

$$= \sum_{n=0}^{\infty} P(N = n) \sum_{k=0}^{\infty} P(S_n = k) \cdot s^k$$

$$= \sum_{n=0}^{\infty} P(N = n) \sum_{k=0}^{\infty} P(\sum_{i=1}^{n} X_i = k) \cdot s^k$$

$$Th.4.6.5 \sum_{n=0}^{\infty} P(N = n) \prod_{i=1}^{n} G_{X_i}(s)$$

$$\stackrel{iid}{=} \sum_{n=0}^{\infty} P(N = n) \cdot (G_X(s))^n$$

$$= G_N(G_X(s))$$

#### Example 4.6.7:

Starting with a single cell at time 0, after one time unit there is probability p that the cell will have split (2 cells), probability q that it will survive without splitting (1 cell), and probability r that it will have died (0 cells). It holds that  $p, q, r \ge 0$  and p + q + r = 1. Any surviving cells have the same probabilities of splitting or dying. What is the pgf for the # of cells at time 2?

$$G_X(s) = G_N(s) = ps^2 + qs + r$$
  
 $G_{S_N}(s) = p(ps^2 + ps + r)^2 + q(ps^2 + ps + r) + r$ 

# <u>Theorem 4.6.8:</u>

Let  $X_1, \ldots, X_N$  be iid rv's with common mgf  $M_X(t)$ . Let N be a discrete rv on the non-negative integers with mgf  $M_N(t)$ . Let N be independent of the  $X_i$ 's. Define  $S_N = \sum_{i=1}^N X_i$ . The mgf of  $S_N$  is

$$M_{S_N}(t) = M_N(\ln M_X(t)).$$

# Proof:

Consider the case that the  $X_i$ 's are non-negative integers:

We know that

$$G_X(s) = E(s^X) = E(e^{\ln s^X}) = E(e^{(\ln s) \cdot X}) = M_X(\ln s)$$

$$\Longrightarrow M_X(s) = G_X(e^s)$$

$$\Longrightarrow M_{S_N}(t) = G_{S_N}(e^t) \stackrel{Th.4.6.6}{=} G_N(G_X(e^t)) = G_N(M_X(t)) = M_N(\ln M_X(t))$$

In the general case, i.e., if the  $X_i's$  are not non-negative integers, we need results from Section 4.7 (conditional expectation) to proof this Theorem.

Lecture 32: Fr 11/10/00

# 4.7 Conditional Expectation

In Section 4.1, we established that the conditional pmf of X given  $Y = y_j$  (for  $P_Y(y_j) > 0$ ) is a pmf. For continuous rv's X and Y, when  $f_Y(y) > 0$ ,  $f_{X|Y}(x \mid y) = \frac{f_{X,Y}(x,y)}{f_Y(y)}$ , and  $f_{X,Y}(x \mid y)$  are continuous, then  $f_{X|Y}(x \mid y)$  is a pdf and it is the conditional pdf of X given Y = y.

#### Definition 4.7.1:

Let X, Y be rv's on  $(\Omega, L, P)$ . Let h be a Borel-measurable function. Assume that E(h(X)) exists. Then the **conditional expectation** of h(X) given Y, i.e.,  $E(h(X) \mid Y)$ , is a rv that takes the value  $E(h(X) \mid y)$ . It is defined as

$$E(h(X) \mid y) = \begin{cases} \sum_{x \in \mathcal{X}} h(x) P(X = x \mid Y = y), & \text{if } (X, Y) \text{ is discrete and } P(Y = y) > 0 \\ \\ \int_{-\infty}^{\infty} h(x) f_{X|Y}(x \mid y) dx, & \text{if } (X, Y) \text{ is continuous and } f_{Y}(y) > 0 \end{cases}$$

Note:

- (i) The rv E(h(X) | Y) = g(Y) is a function of Y as a rv.
- (ii) The usual properties of expectations apply to the conditional expectation:
  - (a)  $E(c \mid Y) = c \ \forall c \in \mathbb{R}$ .
  - (b)  $E(aX + b | Y) = aE(X | Y) + b \ \forall a, b \in \mathbb{R}$ .
  - (c) If  $g_1, g_2$  are Borel-measurable functions and if  $E(g_1(X)), E(g_2(X))$  exist, then  $E(a_1g_1(X) + a_2g_2(X) | Y) = a_1E(g_1(X) | Y) + a_2E(g_2(X) | Y) \quad \forall a_1, a_2 \in \mathbb{R}.$
  - (d) If  $X \geq 0$  then  $E(X \mid Y) \geq 0$ .
  - (e) If  $X_1 \ge X_2$  then  $E(X_1 \mid Y) \ge E(X_2 \mid Y)$ .
- (iii) Moments are defined in the usual way. If  $E(|X|^r) < \infty$ , then  $E(X^r \mid Y)$  exists and is the  $r^{th}$  conditional moment of X given Y.

#### Example 4.7.2:

Recall Example 4.1.12:

$$f_{X,Y}(x,y) = \begin{cases} 2, & 0 < x < y < 1 \\ 0, & \text{otherwise} \end{cases}$$

The conditional pdf's  $f_{Y\mid X}(y\mid x)$  and  $f_{X\mid Y}(x\mid y)$  have been calculated as:

$$f_{Y|X}(y \mid x) = \frac{1}{1-x}$$
 for  $x < y < 1$  (where  $0 < x < 1$ )

and

$$f_{X|Y}(x \mid y) = \frac{1}{y}$$
 for  $0 < x < y$  (where  $0 < y < 1$ ).

So,

$$E(X \mid y) = \int_0^y \frac{x}{y} dx = \frac{y}{2}$$

and

$$E(Y \mid x) = \int_{x}^{1} \frac{1}{1-x} \ y \ dy = \frac{1}{1-x} \frac{y^{2}}{2} \Big|_{x}^{1} = \frac{1}{2} \frac{1-x^{2}}{1-x} = \frac{1+x}{2}.$$

Therefore, we get the rv's  $E(X \mid Y) = \frac{Y}{2}$  and  $E(Y \mid X) = \frac{1+X}{2}$ .

#### <u>Theorem 4.7.3:</u>

If E(h(X)) exists, then

$$E_Y(E_{X|Y}(h(X) \mid Y)) = E(h(X)).$$

#### Proof:

Continuous case only:

$$E_Y(E(h(X) \mid Y)) = \int_{-\infty}^{\infty} E(h(x) \mid y) f_Y(y) dy$$

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

$$= \int_{-\infty}^{\infty} h(x) \int_{-\infty}^{\infty} f_{X,Y}(x, y) dy dx$$

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

$$= E(h(X))$$

#### Theorem 4.7.4:

If  $E(X^2)$  exists, then

$$Var_Y(E(X \mid Y)) + E_Y(Var(X \mid Y)) = Var(X).$$

#### Proof:

$$Var_{Y}(E(X \mid Y)) + E_{Y}(Var(X \mid Y)) = E_{Y}((E(X \mid Y))^{2}) - (E_{Y}(E(X \mid Y)))^{2}$$

$$+ E_{Y}(E(X^{2} \mid Y) - (E(X \mid Y))^{2})$$

$$= E_{Y}((E(X \mid Y))^{2}) - (E(X))^{2} + E(X^{2}) - E_{Y}((E(X \mid Y))^{2})$$

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

$$= Var(X)$$

#### Note:

If  $E(X^2)$  exists, then  $Var(X) \ge Var_Y(E(X \mid Y))$ .  $Var(X) = Var_Y(E(X \mid Y))$  iff X = g(Y). The inequality directly follows from Theorem 4.7.4.

For equality, it is necessary that  $E_Y(Var(X \mid Y)) = E_Y((X - E(X \mid Y))^2 \mid Y) = 0$  which holds if  $X = E(X \mid Y) = g(Y)$ .

If X, Y are independent,  $F_{X|Y}(x \mid y) = F_X(x) \ \forall x$ .

Thus, if E(h(X)) exists, then  $E(h(X) \mid Y) = E(h(X))$ .

#### Proof: (of Theorem 4.6.8)

$$M_{S_N}(t) \stackrel{Def.}{=} E(e^{tS_N})$$

$$= E\left(\underbrace{\exp(t\sum_{i=1}^N X_i)}_{\text{``}h(X)''}\right)$$
 $Th = \underbrace{4.7.3}_{E_N} E_N\left(\underbrace{E\sum_{i=1}^N X_i \mid N}_{i=1} \left(\exp(t\sum_{i=1}^N X_i) \mid N\right)\right)$ 

First consider

$$E_{\sum_{i=1}^{N} X_{i} \mid N} \left( \exp\left(t \sum_{i=1}^{N} X_{i}\right) \mid n \right) = \sum_{i=1}^{n} X_{i} \left( \exp\left(t \sum_{i=1}^{n} X_{i}\right) \right)$$

$$Th. \stackrel{4.6.4(i)}{=} \prod_{i=1}^{n} M_{X_{i}}(t)$$

$$X_{i} \stackrel{iid}{=} \prod_{i=1}^{n} M_{X}(t)$$

$$= (M_{X}(t))^{n}$$

$$\Longrightarrow M_{S_{N}}(t) = E_{N} \left( \prod_{i=1}^{N} M_{X}(t) \right)$$

$$= E_{N}((M_{X}(t))^{N})$$

$$= E_{N}(\exp(N \ln M_{X}(t)))$$

$$\stackrel{(*)}{=} M_{N}(\ln M_{X}(t))$$

(\*) holds since  $M_N(k) = E_N(\exp(N \cdot k))$ .

Lecture 33: Mo 11/13/00

# 4.8 Inequalities and Identities

#### Lemma 4.8.1:

Let a, b be positive numbers and p, q > 1 such that  $\frac{1}{p} + \frac{1}{q} = 1$  (i.e., pq = p + q and  $q = \frac{p}{p-1}$ ). Then it holds that

$$\frac{1}{p}a^p + \frac{1}{q}b^q \ge ab$$

with equality iff  $a^p = b^q$ .

#### Proof:

Fix b. Let

$$g(a) = \frac{1}{p}a^p + \frac{1}{q}b^q - ab$$

$$\implies g'(a) = a^{p-1} - b \stackrel{!}{=} 0$$

$$\implies b = a^{p-1}$$

$$\implies b^q = a^{(p-1)q} = a^p$$

$$g''(a) = (p-1)a^{p-2} > 0$$

Since g''(a) > 0, this is really a minimum. The minimum value is obtained for  $b = a^{p-1}$  and it is

$$\frac{1}{p}a^p + \frac{1}{q}(a^{p-1})^q - aa^{p-1} = \frac{1}{p}a^p + \frac{1}{q}a^p - a^p = a^p(\frac{1}{p} + \frac{1}{q} - 1) = 0.$$

Since g''(a) > 0, the minimum is unique and  $g(a) \ge 0$ . Therefore,

$$g(a) + ab = \frac{1}{p}a^p + \frac{1}{q}b^q \ge ab.$$

#### Theorem 4.8.2: Hölder's Inequality

Let X, Y be 2 rv's. Let p, q > 1 such that  $\frac{1}{p} + \frac{1}{q} = 1$  (i.e., pq = p + q and  $q = \frac{p}{p-1}$ ). Then it holds that

$$E(\mid XY\mid) \le (E(\mid X\mid^p))^{\frac{1}{p}} (E(\mid Y\mid^q))^{\frac{1}{q}}.$$

#### Proof:

In Lemma 4.8.1, let

$$a = \frac{\mid X \mid}{\left(E(\mid X \mid^p)\right)^{\frac{1}{p}}} \quad \text{and} \quad b = \frac{\mid Y \mid}{\left(E(\mid Y \mid^q)\right)^{\frac{1}{q}}}.$$

$$\xrightarrow{Lem\underline{ma}} \stackrel{4.8.1}{\longrightarrow} \frac{1}{p} \frac{\mid X \mid^p}{E(\mid X \mid^p)} + \frac{1}{q} \frac{\mid Y \mid^q}{E(\mid Y \mid^q)} \ge \frac{\mid XY \mid}{(E(\mid X \mid^p))^{\frac{1}{p}}(E(\mid Y \mid^q))^{\frac{1}{q}}}$$

Taking expectations on both sides of this inequality, we get

$$1 = \frac{1}{p} + \frac{1}{q} \ge \frac{E(|XY|)}{(E(|X|^p))^{\frac{1}{p}}(E(|Y|^q))^{\frac{1}{q}}}$$

The result follows immediately when multiplying both sides with  $(E(\mid X\mid^p))^{\frac{1}{p}}(E(\mid Y\mid^q))^{\frac{1}{q}}$ .

#### Note:

Note that Theorem 4.5.7 (ii) (Cauchy–Schwarz Inequality) is a special case of Theorem 4.8.2 with p = q = 2.

#### Theorem 4.8.3: Minkowski's Inequality

Let X, Y be 2 rv's. Then it holds for  $1 \le p < \infty$  that

$$(E(|X+Y|^p))^{\frac{1}{p}} \le (E(|X|^p))^{\frac{1}{p}} + (E(|Y|^p))^{\frac{1}{p}}.$$

#### Proof:

$$\begin{split} E(\mid X+Y\mid^{p}) &= E(\mid X+Y\mid\cdot\mid X+Y\mid^{p-1}) \\ &\leq E((\mid X\mid+\mid Y\mid)\cdot\mid X+Y\mid^{p-1}) \\ &= E(\mid X\mid\cdot\mid X+Y\mid^{p-1}) + E(\mid Y\mid\cdot\mid X+Y\mid^{p-1}) \\ &\stackrel{Th.4.8.2}{\leq} \left( E(\mid X\mid^{p}) \right)^{\frac{1}{p}} \cdot \left( E(\mid X\mid Y\mid^{p-1})^{q} \right) \right)^{\frac{1}{q}} \\ &+ \left( E(\mid Y\mid^{p}) \right)^{\frac{1}{p}} \cdot \left( E(\mid X\mid Y\mid^{p-1})^{q} \right) \right)^{\frac{1}{q}} \quad (A) \\ &\stackrel{(*)}{=} \left( \left( E(\mid X\mid^{p}) \right)^{\frac{1}{p}} + \left( E(\mid Y\mid^{p}) \right)^{\frac{1}{p}} \right) \cdot \left( E(\mid X\mid Y\mid^{p}) \right)^{\frac{1}{q}} \end{split}$$

Divide the left and right side of this inequality by  $(E(|X+Y|^p))^{\frac{1}{q}}$ .

The result on the left side is  $(E(\mid X+Y\mid^p))^{1-\frac{1}{q}}=(E(\mid X+Y\mid^p))^{\frac{1}{p}},$  and the result on the right side is  $(E(\mid X\mid^p))^{\frac{1}{p}}+(E(\mid Y\mid^p))^{\frac{1}{p}}.$  Therefore, Theorem 4.8.3 holds.

In (A), we define  $q = \frac{p}{p-1}$ . Therefore,  $\frac{1}{p} + \frac{1}{q} = 1$  and (p-1)q = p.

(\*) holds since  $\frac{1}{p} + \frac{1}{q} = 1$ , a condition to meet Hölder's Inequality.

#### Definition 4.8.4:

A function g(x) is **convex** if

$$g(\lambda x + (1 - \lambda)y) \le \lambda g(x) + (1 - \lambda)g(y) \quad \forall x, y \in \mathbb{R} \quad \forall 0 < \lambda < 1.$$

#### Note:

- (i) Geometrically, a convex function falls above all of its tangent lines. Also, a connecting line between any pairs of points (x, g(x)) and (y, g(y)) in the 2-dimensional plane always falls above the curve.
- (ii) A function g(x) is **concave** iff -g(x) is convex.

#### Theorem 4.8.5: Jensen's Inequality

Let X be a rv. If g(x) is a convex function, then

$$E(g(X)) \ge g(E(X))$$

given that both expectations exist.

#### <u>Proof:</u>

Construct a tangent line l(x) to g(x) at the (constant) point  $x_0 = E(X)$ :

$$l(x) = ax + b$$
 for some  $a, b \in \mathbb{R}$ 

The function g(x) is a convex function and falls above the tangent line l(x)

$$\implies g(x) \ge ax + b \ \forall x \in IR$$

$$\Longrightarrow E(g(X)) \geq E(aX+b) = aE(X) + b = l(E(X)) \stackrel{\text{\tiny tangent point}}{=} g(E(X))$$

Therefore, Theorem 4.8.5 holds.

#### Note:

Typical convex functions g are:

(i) 
$$q_1(x) = |x| \implies E(|X|) > |E(X)|$$
.

(ii) 
$$q_2(x) = x^2 \implies E(X^2) > (E(X))^2 \implies Var(X) > 0.$$

(iii) 
$$g_3(x) = \frac{1}{x^p}$$
 for  $x > 0, p > 0 \implies E(\frac{1}{X^p}) \ge \frac{1}{(E(X))^p}$ ; for  $p = 1$ :  $E(\frac{1}{X}) \ge \frac{1}{E(X)}$ 

(iv) Other convex functions are  $x^p$  for  $x > 0, p \ge 1$ ;  $\theta^x$  for  $\theta > 1$ ;  $-\ln(x)$  for x > 0; etc.

- (v) Recall that if g is convex and differentiable, then  $g''(x) \geq 0 \ \forall x$ .
- (vi) If the function g is concave, the direction of the inequality in Jensen's Inequality is reversed, i.e.,  $E(g(X)) \leq g(E(X))$ .

Lecture 35/1: Fr 11/17/00

### Example 4.8.6:

Given the real numbers  $a_1, a_2, \ldots, a_n > 0$ , we define

$$\begin{aligned} \textbf{arithmetic mean} & : & a_A = \frac{1}{n}(a_1 + a_2 + \ldots + a_n) = \frac{1}{n}\sum_{i=1}^n a_i \\ \textbf{geometric mean} & : & a_G = (a_1 \cdot a_2 \cdot \ldots \cdot a_n)^{\frac{1}{n}} = \left(\prod_{i=1}^n a_i\right)^{\frac{1}{n}} \\ \textbf{harmonic mean} & : & a_H = \frac{1}{\frac{1}{n}\left(\frac{1}{a_1} + \frac{1}{a_2} + \ldots + \frac{1}{a_n}\right)} = \frac{1}{\frac{1}{n}\sum_{i=1}^n \frac{1}{a_i}} \end{aligned}$$

Let X be a rv that takes values  $a_1, a_2, \ldots, a_n > 0$  with probability  $\frac{1}{n}$  each.

(i)  $a_A \geq a_G$ :

$$\ln(a_A) = \ln\left(\frac{1}{n}\sum_{i=1}^n a_i\right) \\
= \ln(E(X)) \\
\ln concave \\
\geq E(\ln(X)) \\
= \sum_{i=1}^n \frac{1}{n}\ln(a_i) \\
= \frac{1}{n}\sum_{i=1}^n \ln(a_i) \\
= \frac{1}{n}\ln(\prod_{i=1}^n a_i) \\
= \ln(a_G)$$

Taking the anti-log of both sides gives  $a_A \geq a_G$ .

(ii)  $a_A \geq a_H$ :

$$\frac{1}{a_A} \qquad = \qquad \frac{1}{\frac{1}{n} \sum_{i=1}^n a_i}$$

$$= \frac{1}{E(X)}$$

$$^{1/X} \stackrel{convex}{\leq} E(\frac{1}{X})$$

$$= \sum_{i=1}^{n} \frac{1}{n} \frac{1}{a_i}$$

$$= \frac{1}{n} \left(\frac{1}{a_1} + \frac{1}{a_2} + \dots \frac{1}{a_n}\right)$$

$$= \frac{1}{a_H}$$

Inverting both sides gives  $a_A \geq a_H$ .

(iii)  $a_G \geq a_H$ :

$$-\ln(a_{H}) = \ln(\frac{1}{a_{H}})$$

$$= \ln\left(\frac{1}{n}(\frac{1}{a_{1}} + \frac{1}{a_{2}} + \dots + \frac{1}{a_{n}})\right)$$

$$= \ln(E(\frac{1}{X}))$$

$$= \ln(E(\frac{1}{X}))$$

$$= \sum_{i=1}^{n} \frac{1}{n} \ln(\frac{1}{a_{i}})$$

$$= \sum_{i=1}^{n} \ln(\frac{1}{a_{i}})$$

$$= \frac{1}{n} \sum_{i=1}^{n} -\ln a_{i}$$

$$= -\frac{1}{n} \ln(\prod_{i=1}^{n} a_{i})$$

$$= -\ln a_{G}$$

Multiplying both sides with -1 gives  $\ln a_H \leq \ln a_G$ . Then taking the anti-log of both sides gives  $a_H \leq a_G$ .

In summary,  $a_H \leq a_G \leq a_A$ . Note that it would have been sufficient to prove steps (i) and (iii) only to establish this result. However, step (ii) has been included to provide another example how to apply Theorem 4.8.5.

# Theorem 4.8.7: Covariance Inequality

Let X be a rv with finite mean  $\mu$ .

(i) If g(x) is non-decreasing, then

$$E(g(X)(X - \mu)) \ge 0$$

if this expectation exists.

(ii) If g(x) is non-decreasing and h(x) is non-increasing, then

$$E(g(X)h(X)) \le E(g(X))E(h(X))$$

if all expectations exist.

(iii) If g(x) and h(x) are both non-decreasing or if g(x) and h(x) are both non-increasing, then

$$E(g(X)h(X)) \ge E(g(X))E(h(X))$$

if all expectations exist.

#### Proof:

Homework

Lecture 34: We 11/15/00

# 5 Particular Distributions

#### 5.1 Multivariate Normal Distributions

#### Definition 5.1.1:

A rv X has a (univariate) Normal distribution, i.e.,  $X \sim N(\mu, \sigma^2)$  with  $\mu \in \mathbb{R}$  and  $\sigma > 0$ , iff it has the pdf

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

X has a standard Normal distribution iff  $\mu = 0$  and  $\sigma^2 = 1$ , i.e.,  $X \sim N(0, 1)$ .

#### Note:

If  $X \sim N(\mu, \sigma^2)$ , then  $E(X) = \mu$  and  $Var(X) = \sigma^2$ .

If  $X_1 \sim N(\mu_1, \sigma_1^2)$ ,  $X_2 \sim N(\mu_2, \sigma_2^2)$ ,  $X_1$  and  $X_2$  independent, and  $c_1, c_2 \in \mathbb{R}$ , then

$$Y = c_1 X_1 + c_2 X_2 \sim N(c_1 \mu_1 + c_2 \mu_2, c_1^2 \sigma_1^2 + c_2^2 \sigma_2^2).$$

#### Definition 5.1.2:

A 2-rv (X,Y) has a **bivariate Normal distribution** iff there exist constants  $a_{11}, a_{12}, a_{21}, a_{22}, \mu_1, \mu_2 \in \mathbb{R}$  and iid N(0,1) rv's  $Z_1$  and  $Z_2$  such that

$$X = \mu_1 + a_{11}Z_1 + a_{12}Z_2, \quad Y = \mu_2 + a_{21}Z_1 + a_{22}Z_2.$$

If we define

$$A = \begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{pmatrix}, \quad \underline{\mu} = \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}, \quad \underline{X} = \begin{pmatrix} X \\ Y \end{pmatrix}, \quad \underline{Z} = \begin{pmatrix} Z_1 \\ Z_2 \end{pmatrix},$$

then we can write

$$\underline{X} = A\underline{Z} + \underline{\mu}.$$

#### Note:

- (i) Recall that for  $X \sim N(\mu, \sigma^2)$ , X can be defined as  $X = \sigma Z + \mu$ , where  $Z \sim N(0, 1)$ .
- (ii)  $E(X) = \mu_1 + a_{11}E(Z_1) + a_{12}E(Z_2) = \mu_1$  and  $E(Y) = \mu_2 + a_{21}E(Z_1) + a_{22}E(Z_2) = \mu_2$ . The marginal distributions are  $X \sim N(\mu_1, a_{11}^2 + a_{12}^2)$  and  $Y \sim N(\mu_2, a_{21}^2 + a_{22}^2)$ . Thus, X and Y have (univariate) Normal marginal densities or degenerate marginal densities (which correspond to Dirac distributions) if  $a_{i1} = a_{i2} = 0$ .

(iii) There exists another (equivalent) formulation of the previous defintion using the joint pdf (see Rohatgi, page 227).

#### Example:

Let  $X \sim N(\mu_1, \sigma_1^2)$ ,  $Y \sim N(\mu_2, \sigma_2^2)$ , X and Y independent.

What is the distribution of  $\begin{pmatrix} X \\ Y \end{pmatrix}$ ?

Since  $X \sim N(\mu_1, \sigma_1^2)$ , it follows that  $X = \mu_1 + \sigma_1 Z_1$ , where  $Z_1 \sim N(0, 1)$ .

Since  $Y \sim N(\mu_2, \sigma_2^2)$ , it follows that  $Y = \mu_2 + \sigma_2 Z_2$ , where  $Z_2 \sim N(0, 1)$ .

Therefore,

$$\left(\begin{array}{c} X \\ Y \end{array}\right) = \left(\begin{array}{cc} \sigma_1 & 0 \\ 0 & \sigma_2 \end{array}\right) \left(\begin{array}{c} Z_1 \\ Z_2 \end{array}\right) + \left(\begin{array}{c} \mu_1 \\ \mu_2 \end{array}\right).$$

Since X, Y are independent, it follows by Theorem 4.2.7 that  $Z_1 = \frac{X - \mu_1}{\sigma_1}$  and  $Z_2 = \frac{X - \mu_2}{\sigma_2}$  are independent. Thus,  $Z_1, Z_2$  are iid N(0, 1). It follows from Definition 5.1.2 that  $\begin{pmatrix} X \\ Y \end{pmatrix}$  has a bivariate Normal distribution.

## Theorem 5.1.3:

Define  $\underline{g}: \mathbb{R}^2 \to \mathbb{R}^2$  as  $\underline{g}(\underline{x}) = C\underline{x} + \underline{d}$ . If  $\underline{X}$  is a bivariate Normal rv, then  $\underline{g}(\underline{X})$  also is a bivariate Normal rv.

#### Proof:

$$\begin{array}{lll} g(\underline{X}) & = & C\underline{X} + \underline{d} \\ & = & C(A\underline{Z} + \underline{\mu}) + \underline{d} \\ & = & \underbrace{(CA)}_{\text{another matrix}} \underline{Z} + \underbrace{(C\underline{\mu} + \underline{d})}_{\text{another vector}} \\ & = & \tilde{A}\underline{Z} + \tilde{\mu} & \text{which represents another bivariate Normal distribution} \end{array}$$

Note:

$$\rho \sigma_1 \sigma_2 = Cov(X, Y) = Cov(a_{11}Z_1 + a_{12}Z_2, a_{21}Z_1 + a_{22}Z_2) 
= a_{11}a_{21}Cov(Z_1, Z_1) + (a_{11}a_{22} + a_{12}a_{21})Cov(Z_1, Z_2) + a_{12}a_{22}Cov(Z_2, Z_2) 
= a_{11}a_{21} + a_{12}a_{22}$$

since  $Z_1, Z_2$  are iid N(0, 1) rv's.

#### Definition 5.1.4:

The variance–covariance matrix of (X, Y) is

$$\Sigma = AA' = \begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{pmatrix} \begin{pmatrix} a_{11} & a_{21} \\ a_{12} & a_{22} \end{pmatrix} = \begin{pmatrix} a_{11}^2 + a_{12}^2 & a_{11}a_{21} + a_{12}a_{22} \\ a_{11}a_{21} + a_{12}a_{22} & a_{21}^2 + a_{22}^2 \end{pmatrix} = \begin{pmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 \end{pmatrix}.$$

Note:

$$det(\Sigma) = |\Sigma| = \sigma_1^2 \sigma_2^2 - \rho^2 \sigma_1^2 \sigma_2^2 = \sigma_1^2 \sigma_2^2 (1 - \rho^2),$$
$$\sqrt{|\Sigma|} = \sigma_1 \sigma_2 \sqrt{1 - \rho^2}$$

and

$$\Sigma^{-1} = \frac{1}{|\Sigma|} \cdot \begin{pmatrix} \sigma_2^2 & -\rho\sigma_1\sigma_2 \\ -\rho\sigma_1\sigma_2 & \sigma_1^2 \end{pmatrix} = \begin{pmatrix} \frac{1}{\sigma_1^2(1-\rho^2)} & \frac{-\rho}{\sigma_1\sigma_2(1-\rho^2)} \\ \frac{-\rho}{\sigma_1\sigma_2(1-\rho^2)} & \frac{1}{\sigma_2^2(1-\rho^2)} \end{pmatrix}$$

#### Theorem 5.1.5:

Assume that  $\sigma_1 > 0, \sigma_2 > 0$  and  $|\rho| < 1$ . Then the joint pdf of  $\underline{X} = (X, Y) = A\underline{Z} + \underline{\mu}$  (as defined in Definition 5.1.2) is

$$\begin{split} f_{\underline{X}}(\underline{x}) &= \frac{1}{2\pi\sqrt{|\Sigma|}}\exp\left(-\frac{1}{2}(\underline{x}-\underline{\mu})'\Sigma^{-1}(\underline{x}-\underline{\mu})\right) \\ &= \frac{1}{2\pi\sigma_1\sigma_2\sqrt{1-\rho^2}}\exp\left(-\frac{1}{2(1-\rho^2)}\left(\left(\frac{x-\mu_1}{\sigma_1}\right)^2 - 2\rho\left(\frac{x-\mu_1}{\sigma_1}\right)\left(\frac{y-\mu_2}{\sigma_2}\right) + \left(\frac{y-\mu_2}{\sigma_2}\right)^2\right)\right) \end{split}$$

#### Proof:

Since  $\Sigma$  is positive definite and symmetric, A is invertible.

The mapping  $\underline{Z} \to \underline{X}$  is 1-to-1:

$$\underline{X} = A\underline{Z} + \underline{\mu}$$

$$\Longrightarrow \underline{Z} = A^{-1}(\underline{X} - \underline{\mu})$$

$$J = |A^{-1}| = \frac{1}{|A|}$$

$$|A| = \sqrt{|A|^2} = \sqrt{|A| \cdot |A^T|} = \sqrt{|AA^T|} = \sqrt{|\Sigma|} = \sqrt{\sigma_1^2 \sigma_2^2 - \rho^2 \sigma_1^2 \sigma_2^2} = \sigma_1 \sigma_2 \sqrt{1 - \rho^2}$$

We can use this result to get to the second line of the theorem:

$$f_{\underline{Z}}(\underline{z}) = \frac{1}{\sqrt{2\pi}} e^{-\frac{z_1^2}{2}} \frac{1}{\sqrt{2\pi}} e^{-\frac{z_2^2}{2}}$$
$$= \frac{1}{2\pi} e^{-\frac{1}{2}(\underline{z}^T \underline{z})}$$

As already stated, the mapping from  $\underline{Z}$  to  $\underline{X}$  is 1-to-1, so we can apply Theorem 4.3.5:

$$f_{\underline{X}}(\underline{x}) = \frac{1}{2\pi\sqrt{|\Sigma|}} \exp(-\frac{1}{2}(\underline{x} - \underline{\mu})^T \underbrace{(A^{-1})^T A^{-1}}_{\Sigma^{-1}}(\underline{x} - \underline{\mu}))$$

$$\stackrel{(*)}{=} \frac{1}{2\pi\sqrt{|\Sigma|}} \exp(-\frac{1}{2}(\underline{x} - \underline{\mu})^T \Sigma^{-1}(\underline{x} - \underline{\mu}))$$

This proves the 1<sup>st</sup> line of the Theorem. Step (\*) holds since

$$(A^{-1})^T A^{-1} = (A^T)^{-1} A^{-1} = (AA^T)^{-1} = \Sigma^{-1}.$$

The second line of the Theorem is based on the transformations stated in the Note following Definition 5.1.4:

$$f_{\underline{X}}(\underline{x}) = \frac{1}{2\pi\sigma_1\sigma_2\sqrt{1-\rho^2}} \exp\left(-\frac{1}{2(1-\rho^2)} \left(\left(\frac{x-\mu_1}{\sigma_1}\right)^2 - 2\rho\left(\frac{x-\mu_1}{\sigma_1}\right)\left(\frac{y-\mu_2}{\sigma_2}\right) + \left(\frac{y-\mu_2}{\sigma_2}\right)^2\right)\right)$$

Note:

In the situation of Theorem 5.1.5, we say that  $(X,Y) \sim N(\mu_1, \mu_2, \sigma_1^2, \sigma_2^2, \rho)$ .

#### Theorem 5.1.6:

If (X, Y) has a non-degenerate  $N(\mu_1, \mu_2, \sigma_1^2, \sigma_2^2, \rho)$  distribution, then the conditional distribution of X given Y = y is

$$N(\mu_1 + \rho \frac{\sigma_1}{\sigma_2}(y - \mu_2), \sigma_1^2(1 - \rho^2)).$$

Proof:

Homework

#### Example 5.1.7:

Let rv's  $(X_1, Y_1)$  be N(0, 0, 1, 1, 0) with pdf  $f_1(x, y)$  and  $(X_2, Y_2)$  be  $N(0, 0, 1, 1, \rho)$  with pdf  $f_2(x, y)$ . Let (X, Y) be the rv that corresponds to the pdf

$$f_{X,Y}(x,y) = \frac{1}{2}f_1(x,y) + \frac{1}{2}f_2(x,y).$$

(X, Y) is a bivariate Normal rv iff  $\rho = 0$ . However, the marginal distributions of X and Y are always N(0, 1) distributions. See also Rohatgi, page 229, Remark 2.

Lecture 35/2: Fr 11/17/00

#### Theorem 5.1.8:

The mgf  $M_{\underline{X}}(\underline{t})$  of a bivariate Normal rv  $\underline{X} = (X, Y)$  is

$$M_{\underline{X}}(\underline{t}) = M_{X,Y}(t_1, t_2) = \exp(\underline{\mu}' \underline{t} + \frac{1}{2} \underline{t}' \underline{\Sigma} \underline{t}) = \exp\left(\mu_1 t_1 + \mu_2 t_2 + \frac{1}{2} (\sigma_1^2 t_1^2 + \sigma_2^2 t_2^2 + 2\rho \sigma_1 \sigma_2 t_1 t_2)\right).$$

#### Proof:

The mgf of a univariate Normal rv  $X \sim N(\mu, \sigma^2)$  will be used to develop the mgf of a bivariate Normal rv X = (X, Y):

$$\begin{split} M_X(t) &= E(\exp(tX)) \\ &= \int_{-\infty}^{\infty} \exp(tx) \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x-\mu)^2\right) \, dx \\ &= \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}[-2\sigma^2tx \, + \, (x-\mu)^2]\right) \, dx \\ &= \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}[-2\sigma^2tx \, + \, (x^2-2\mu x \, + \, \mu^2)]\right) \, dx \\ &= \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}[x^2-2(\mu \, + \, \sigma^2t)x \, + \, (\mu + t\sigma^2)^2 \, - \, (\mu + t\sigma^2)^2 \, + \, \mu^2]\right) \, dx \\ &= \exp\left(-\frac{1}{2\sigma^2}[-(\mu + t\sigma^2)^2 + \mu^2]\right) \underbrace{\int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}[x - (\mu + t\sigma^2)]^2\right) \, dx}_{\text{pdf of } N(\mu + t\sigma^2, \sigma^2), \text{ that integrates to 1}} \\ &= \exp\left(-\frac{1}{2\sigma^2}[-\mu^2 - 2\mu t\sigma^2 - t^2\sigma^4 + \mu^2]\right) \\ &= \exp\left(\frac{-2\mu t\sigma^2 - t^2\sigma^4}{-2\sigma^2}\right) \\ &= \exp\left(\mu t + \frac{1}{2}\sigma^2t^2\right) \end{split}$$

#### Bivariate Normal mgf:

$$M_{X,Y}(t_1, t_2) = E(exp(t_1X + t_2Y))$$

$$= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \exp(t_1x + t_2y) \ f_{X,Y}(x, y) \ dx \ dy$$

$$= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \exp(t_1x) \ \exp(t_2y) \ f_{X}(x) \ f_{Y|X}(y \mid x) \ dy \ dx$$

$$= \int_{-\infty}^{\infty} \left( \int_{-\infty}^{\infty} \exp(t_2y) \ f_{Y|X}(y \mid x) \ dy \right) \ \exp(t_1x) \ f_{X}(x) \ dx$$

$$\begin{array}{ll} \stackrel{(A)}{=} & \int_{-\infty}^{\infty} \left( \int_{-\infty}^{\infty} \frac{\exp(t_2 y)}{\sigma_2 \sqrt{1 - \rho^2} \sqrt{2\pi}} \, \exp\left( \frac{-(y - \beta_X)^2}{2\sigma_2^2 (1 - \rho^2)} \right) \, dy \right) \, \exp(t_1 x) \, f_X(x) \, dx \\ & | \beta_X = \mu_2 + \rho \frac{\sigma_2}{\sigma_1} (x - \mu_1) \right) \\ \stackrel{(B)}{=} & \int_{-\infty}^{\infty} \, \exp\left( \beta_X t_2 + \frac{1}{2} \sigma_2^2 (1 - \rho^2) t_2^2 \right) \, \exp(t_1 x) \, f_X(x) \, dx \\ & = \int_{-\infty}^{\infty} \, \exp\left( \left[ \mu_2 + \rho \frac{\sigma_2}{\sigma_1} (x - \mu_1) \right] t_2 + \frac{1}{2} \sigma_2^2 (1 - \rho^2) t_2^2 + t_1 x \right) \, f_X(x) \, dx \\ & = \int_{-\infty}^{\infty} \, \exp\left( \mu_2 t_2 + \rho \frac{\sigma_2}{\sigma_1} t_2 x - \rho \frac{\sigma_2}{\sigma_1} \mu_1 t_2 + \frac{1}{2} \sigma_2^2 (1 - \rho^2) t_2^2 + t_1 x \right) \, f_X(x) \, dx \\ & = \exp\left( \frac{1}{2} \sigma_2^2 (1 - \rho^2) t_2^2 + t_2 \mu_2 - \rho \frac{\sigma_2}{\sigma_1} \mu_1 t_2 \right) \int_{-\infty}^{\infty} \exp\left( \left( t_1 + \rho \frac{\sigma_2}{\sigma_1} t_2 \right) x \right) \, f_X(x) \, dx \\ \stackrel{(C)}{=} \, \exp\left( \frac{1}{2} \sigma_2^2 (1 - \rho^2) t_2^2 + t_2 \mu_2 - \rho \frac{\sigma_2}{\sigma_1} \mu_1 t_2 \right) \cdot \exp\left( \mu_1 (t_1 + \rho \frac{\sigma_2}{\sigma_1} t_2) + \frac{1}{2} \sigma_1^2 (t_1 + \rho \frac{\sigma_2}{\sigma_1} t_2)^2 \right) \\ & = \exp\left( \frac{1}{2} \sigma_2^2 t_2^2 - \frac{1}{2} \rho^2 \sigma_2^2 t_2^2 + \mu_2 t_2 - \mu_1 \rho \frac{\sigma_2}{\sigma_1} t_2 + \mu_1 t_1 + \mu_1 \rho \frac{\sigma_2}{\sigma_1} t_2 + \frac{1}{2} \sigma_1^2 t_1^2 + \rho \sigma_1 \sigma_2 t_1 t_2 + \frac{1}{2} \rho^2 \sigma_2^2 t_2^2 \right) \\ & = \exp\left( \mu_1 t_1 + \mu_2 t_2 + \frac{\sigma_1^2 t_1^2 + \sigma_2^2 t_2^2 + 2 \rho \sigma_1 \sigma_2 t_1 t_2}{2} \right) \end{array}$$

(A) follows from Theorem 5.1.6 since  $Y \mid X \sim N(\beta_X, \sigma_2^2(1-\rho^2))$ . (B) follows when we apply our calculations of the mgf of a  $N(\mu, \sigma^2)$  distribution to a  $N(\beta_X, \sigma_2^2(1-\rho^2))$  distribution. (C) holds since the integral represents  $M_X(t_1 + \rho \frac{\sigma_2}{\sigma_1}t_2)$ .

#### Corollary 5.1.9:

Let (X, Y) be a bivariate Normal rv. X and Y are independent iff  $\rho = 0$ .

Lecture 36: Mo 11/20/00

#### Definition 5.1.10:

Let  $\underline{Z}$  be a k-rv of k iid N(0,1) rv's. Let  $A \in \mathbb{R}^{k \times k}$  be a  $k \times k$  matrix, and let  $\underline{\mu} \in \mathbb{R}^k$  be a k-dimensional vector. Then  $\underline{X} = A\underline{Z} + \underline{\mu}$  has a **multivariate Normal distribution** with mean vector  $\underline{\mu}$  and variance-covariance matrix  $\Sigma = AA'$ .

#### Note:

(i) If  $\Sigma$  is non-singular,  $\underline{X}$  has the joint pdf

$$f_{\underline{X}}(\underline{x}) = \frac{1}{(2\pi)^{k/2} (|\Sigma|)^{1/2}} \exp\left(-\frac{1}{2} (\underline{x} - \underline{\mu})' \Sigma^{-1} (\underline{x} - \underline{\mu})\right).$$

If  $\Sigma$  is singular, the joint pdf does exist but it cannot be written explicitly.

- (ii) If  $\Sigma$  is singular, then  $\underline{X} \underline{\mu}$  takes values in a linear subspace of  $\mathbb{R}^k$  with probability 1.
- (iii) If  $\Sigma$  is non-singular, then  $\underline{X}$  has mgf

$$M_{\underline{X}}(\underline{t}) = \exp(\underline{\mu}'\underline{t} + \frac{1}{2}\underline{t}'\Sigma\underline{t}).$$

(iv)  $\underline{X}$  has characteristic function

$$\Phi_{\underline{X}}(\underline{t}) = \exp(i\underline{\mu}'\underline{t} - \frac{1}{2}\underline{t}'\underline{\Sigma}\underline{t})$$

(no matter whether  $\Sigma$  is singular or non-singular).

(v) See Searle, S. R. (1971): "Linear Models", Chapter 2.7, for more details on singular Normal distributions.

#### <u>Theorem 5.1.11:</u>

The components  $X_1, \ldots, X_k$  of a normally distributed k-rv  $\underline{X}$  are independent iff  $Cov(X_i, X_j) = 0 \ \forall i, j = 1, \ldots, k, \ i \neq j$ .

#### <u>Theorem 5.1.12:</u>

Let  $\underline{X} = (X_1, \dots, X_k)'$ .  $\underline{X}$  has a k-dimensional Normal distribution iff every linear function of  $\underline{X}$ , i.e.,  $\underline{X}'\underline{t} = t_1X_1 + t_2X_2 + \dots + t_kX_k$ , has a univariate Normal distribution.

#### Proof:

The Note following Definition 5.1.1 states that any linear function of two Normal rv's has a univariate Normal distribution. By induction on k, we can show that every linear function of  $\underline{X}$ , i.e.,  $\underline{X't}$ , has a univariate Normal distribution.

Conversely, if  $\underline{X't}$  has a univariate Normal distribution, we know from Theorem 5.1.8 that

$$\begin{split} M_{\underline{X'}\underline{t}}(s) &= \exp\left(E(\underline{X'}\underline{t}) \cdot s + \frac{1}{2}Var(\underline{X'}\underline{t}) \cdot s^2\right) \\ &= \exp\left(\underline{\mu'}\underline{t}s + \frac{1}{2}\underline{t'}\Sigma\underline{t}s^2\right) \\ \Longrightarrow M_{\underline{X'}\underline{t}}(1) &= \exp\left(\underline{\mu'}\underline{t} + \frac{1}{2}\underline{t'}\Sigma\underline{t}\right) \\ &= M_{\underline{X}}(\underline{t}) \end{split}$$

By uniqueness of the mgf and Note (iii) that follows Definition 5.1.10,  $\underline{X}$  has a multivariate Normal distribution.

#### 5.2 Exponential Family of Distributions

#### Definition 5.2.1:

Let  $\vartheta$  be an interval on the real line. Let  $\{f(\cdot;\theta):\theta\in\vartheta\}$  be a family of pdf's (or pmf's). We assume that the set  $\{\underline{x}:f(\underline{x};\theta)>0\}$  is independent of  $\theta$ , where  $\underline{x}=(x_1,\ldots,x_n)$ .

We say that the family  $\{f(\cdot;\theta):\theta\in\vartheta\}$  is a **one-parameter exponential family** if there exist real-valued functions  $Q(\theta)$  and  $D(\theta)$  on  $\vartheta$  and Borel-measurable functions  $T(\underline{X})$  and  $S(\underline{X})$  on  $\mathbb{R}^n$  such that

$$f(\underline{x}; \theta) = \exp(Q(\theta)T(\underline{x}) + D(\theta) + S(\underline{x})).$$

#### Note:

We can also write  $f(\underline{x}; \theta)$  as

$$f(\underline{x}; \eta) = h(\underline{x})c(\eta) \exp(\eta T(\underline{x}))$$

where  $h(\underline{x}) = \exp(S(\underline{x}))$ ,  $\eta = Q(\theta)$ , and  $c(\eta) = \exp(D(Q^{-1}(\eta)))$ , and call this the **exponential family in canonical form** for a natural parameter  $\eta$ .

#### Definition 5.2.2:

Let  $\underline{\vartheta} \subseteq \mathbb{R}^k$  be a k-dimensional interval. Let  $\{f(\cdot;\underline{\theta}):\underline{\theta}\in\underline{\vartheta}\}$  be a family of pdf's (or pmf's). We assume that the set  $\{\underline{x}:f(\underline{x};\underline{\theta})>0\}$  is independent of  $\underline{\theta}$ , where  $\underline{x}=(x_1,\ldots,x_n)$ .

We say that the family  $\{f(\cdot;\underline{\theta}):\underline{\theta}\in\underline{\vartheta}\}$  is a k-parameter exponential family if there exist real-valued functions  $Q_1(\underline{\theta}),\ldots Q_k(\underline{\theta})$  and  $D(\underline{\theta})$  on  $\underline{\vartheta}$  and Borel-measurable functions  $T_1(\underline{X}),\ldots,T_k(\underline{X})$  and  $S(\underline{X})$  on  $\mathbb{R}^n$  such that

$$f(\underline{x}; \underline{\theta}) = \exp\left(\sum_{i=1}^{k} Q_i(\underline{\theta}) T_i(\underline{x}) + D(\underline{\theta}) + S(\underline{x})\right).$$

#### Note:

Similar to the Note following Definition 5.2.1, we can express the k-parameter **exponential** family in canonical form for a natural  $k \times 1$  parameter vector  $\eta = (\eta_1, \dots, \eta_k)'$ .

Example 5.2.3: Let  $X \sim N(\mu, \sigma^2)$  with both parameters  $\mu$  and  $\sigma^2$  unknown. We have:

$$\begin{split} f(x;\underline{\theta}) &= \frac{1}{\sqrt{2\pi\sigma^2}} \, \exp\left(-\frac{1}{2\sigma^2}(x-\mu)^2\right) = \exp\left(-\frac{1}{2\sigma^2}x^2 + \frac{\mu}{\sigma^2}x - \frac{\mu^2}{2\sigma^2} - \frac{1}{2}\ln(2\pi\sigma^2)\right) \\ \underline{\theta} &= (\mu,\sigma^2) \\ \underline{\theta} &= \{(\mu,\sigma^2) : \mu \in I\!\!R, \sigma^2 > 0\} \end{split}$$

Therefore,

$$Q_{1}(\underline{\theta}) = -\frac{1}{2\sigma^{2}}$$

$$T_{1}(x) = x^{2}$$

$$Q_{2}(\underline{\theta}) = \frac{\mu}{\sigma^{2}}$$

$$T_{2}(x) = x$$

$$D(\underline{\theta}) = -\frac{\mu^{2}}{2\sigma^{2}} - \frac{1}{2}\ln(2\pi\sigma^{2})$$

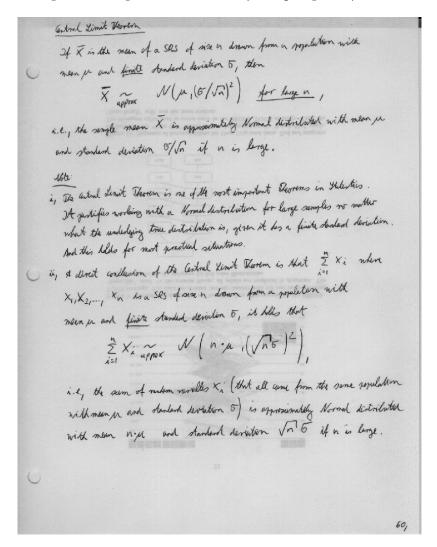
$$S(x) = 0$$

Thus, this is a 2-parameter exponential family.

# 6 Limit Theorems

#### Motivation:

I found this slide from my Stat 250, Section 003, "Introductory Statistics" class (an undergraduate class I taught at George Mason University in Spring 1999):



What does this mean at a more theoretical level???

Lecture 37: Mo 11/27/00

#### Modes of Convergence 6.1

#### Definition 6.1.1:

Let  $X_1, \ldots, X_n$  be iid rv's with common cdf  $F_X(x)$ . Let  $\underline{T} = \underline{T}(\underline{X})$  be any **statistic**, i.e., a Borel-measurable function of X that does not involve the population parameter(s)  $\vartheta$ , defined on the support  $\mathcal{X}$  of  $\underline{X}$ . The induced probability distribution of  $\underline{T}(\underline{X})$  is called the **sampling** distribution of  $\underline{T}(\underline{X})$ .

#### Note:

(i) Commonly used statistics are:

Sample Mean: 
$$\overline{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$$

Sample Variance: 
$$S_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \overline{X}_n)^2$$

Sample Median, Order Statistics, Min, Max, etc.

- (ii) Recall that if  $X_1, \ldots, X_n$  are iid and if E(X) and Var(X) exist, then  $E(\overline{X}_n) = \mu$  $E(X), E(S_n^2) = \sigma^2 = Var(X), \text{ and } Var(\overline{X}_n) = \frac{\sigma^2}{n}.$
- (iii) Recall that if  $X_1, \ldots, X_n$  are iid and if X has mgf  $M_X(t)$  or characteristic function  $\Phi_X(t)$ then  $M_{\overline{X}_n}(t) = (M_X(\frac{t}{n}))^n$  or  $\Phi_{\overline{X}_n}(t) = (\Phi_X(\frac{t}{n}))^n$ .

Note: Let  $\{X_n\}_{n=1}^{\infty}$  be a sequence of rv's on some probability space  $(\Omega, L, P)$ . Is there any meaning behind the expression  $\lim_{n\to\infty} X_n = X$ ? Not immediately under the usual definitions of limits. We first need to define modes of convergence for rv's and probabilities.

#### Definition 6.1.2:

Let  $\{X_n\}_{n=1}^{\infty}$  be a sequence of rv's with cdf's  $\{F_n\}_{n=1}^{\infty}$  and let X be a rv with cdf F. If  $F_n(x) \to F(x)$  at all continuity points of F, we say that  $X_n$  converges in distribution to X  $(X_n \stackrel{d}{\longrightarrow} X)$  or  $X_n$  converges in law to X  $(X_n \stackrel{L}{\longrightarrow} X)$ , or  $F_n$  converges weakly to F $(F_n \stackrel{w}{\longrightarrow} F).$ 

Example 6.1.3: Let  $X_n \sim N(0, \frac{1}{n})$ . Then

$$F_n(x) = \int_{-\infty}^x \frac{\exp\left(-\frac{1}{2}nt^2\right)}{\sqrt{\frac{2\pi}{n}}} dt$$

$$= \int_{-\infty}^{\sqrt{n}x} \frac{\exp(-\frac{1}{2}s^2)}{\sqrt{2\pi}} ds$$
$$= \phi(\sqrt{n}x)$$

$$\implies F_n(x) \to \begin{cases} \phi(\infty) = 1, & \text{if } x > 0 \\ \phi(0) = \frac{1}{2}, & \text{if } x = 0 \\ \phi(-\infty) = 0, & \text{if } x < 0 \end{cases}$$

If  $F_X(x) = \begin{cases} 1, & x \ge 0 \\ 0, & x < 0 \end{cases}$  the only point of discontinuity is at x = 0. Everywhere else,  $\phi(\sqrt{n}x) = F_n(x) \to F_X(x)$ .

So,  $X_n \stackrel{d}{\longrightarrow} X$ , where P(X=0)=1, or  $X_n \stackrel{d}{\longrightarrow} 0$  since the limiting rv here is degenerate, i.e., it has a Dirac(0) distribution.

#### Example 6.1.4:

In this example, the sequence  $\{F_n\}_{n=1}^{\infty}$  converges pointwise to something that is not a cdf: Let  $X_n \sim \text{Dirac}(n)$ , i.e.,  $P(X_n = n) = 1$ . Then,

$$F_n(x) = \begin{cases} 0, & x < n \\ 1, & x \ge n \end{cases}$$

It is  $F_n(x) \to 0 \ \forall x$  which is not a cdf. Thus, there is no rv X such that  $X_n \stackrel{d}{\longrightarrow} X$ .

#### Example 6.1.5:

Let  $\{X_n\}_{n=1}^{\infty}$  be a sequence of rv's such that  $P(X_n=0)=1-\frac{1}{n}$  and  $P(X_n=n)=\frac{1}{n}$  and let  $X \sim \text{Dirac}(0)$ , i.e., P(X=0)=1.

It is

$$F_n(x) = \begin{cases} 0, & x < 0 \\ 1 - \frac{1}{n}, & 0 \le x < n \\ 1, & x \ge n \end{cases}$$

$$F_X(x) = \begin{cases} 0, & x < 0 \\ 1, & x \ge 0 \end{cases}$$

It holds that  $F_n \xrightarrow{w} F_X$  but

$$E(X_n^k) = n^{k-1} \not\to E(X^k) = 0.$$

Thus, convergence in distribution does not imply convergence of moments/means.

#### Note:

Convergence in distribution does not say that the  $X_i$ 's are close to each other or to X. It only means that their cdf's are (eventually) close to some cdf F. The  $X_i$ 's do not even have to be defined on the same probability space.

#### Example 6.1.6:

Let X and  $\{X_n\}_{n=1}^{\infty}$  be iid N(0,1). Obviously,  $X_n \stackrel{d}{\longrightarrow} X$  but  $\lim_{n \to \infty} X_n \neq X$ .

#### <u>Theorem 6.1.7:</u>

Let X and  $\{X_n\}_{n=1}^{\infty}$  be discrete rv's with support  $\mathcal{X}$  and  $\{\mathcal{X}_n\}_{n=1}^{\infty}$ , respectively. Define the countable set  $A = \mathcal{X} \cup \bigcup_{n=1}^{\infty} \mathcal{X}_n = \{a_k : k = 1, 2, 3, \ldots\}$ . Let  $p_k = P(X = a_k)$  and  $p_{nk} = P(X_n = a_k)$ . Then it holds that  $p_{nk} \to p_k \ \forall k \text{ iff } X_n \stackrel{d}{\longrightarrow} X$ .

#### Theorem 6.1.8:

Let X and  $\{X_n\}_{n=1}^{\infty}$  be continuous rv's with pdf's f and  $\{f_n\}_{n=1}^{\infty}$ , respectively. If  $f_n(x) \to f(x)$  for almost all x as  $n \to \infty$  then  $X_n \stackrel{d}{\longrightarrow} X$ .

#### Theorem 6.1.9:

Let X and  $\{X_n\}_{n=1}^{\infty}$  be rv's such that  $X_n \stackrel{d}{\longrightarrow} X$ . Let  $c \in \mathbb{R}$  be a constant. Then it holds:

- (i)  $X_n + c \xrightarrow{d} X + c$ .
- (ii)  $cX_n \stackrel{d}{\longrightarrow} cX$ .
- (iii) If  $a_n \to a$  and  $b_n \to b$ , then  $a_n X_n + b_n \stackrel{d}{\longrightarrow} aX + b$ .

#### Proof:

#### Part (iii):

Suppose that a > 0,  $a_n > 0$ . (If a < 0,  $a_n < 0$ , the result follows via (ii) and c = -1.) Let  $Y_n = a_n X_n + b_n$  and Y = aX + b. It is

$$F_Y(y) = P(Y < y) = P(aX + b < y) = P(X < \frac{y - b}{a}) = F_X(\frac{y - b}{a}).$$

Likewise,

$$F_{Y_n}(y) = F_{X_n}(\frac{y - b_n}{a_n}).$$

If y is a continuity point of  $F_Y$ ,  $\frac{y-b}{a}$  is a continuity point of  $F_X$ . Since  $a_n \to a, b_n \to b$  and  $F_{X_n}(x) \to F_X(x)$ , it follows that  $F_{Y_n}(y) \to F_Y(y)$  for every continuity point y of  $F_Y$ . Thus,  $a_n X_n + b_n \xrightarrow{d} aX + b$ .

#### Definition 6.1.10:

Lecture 38: We 11/29/00

Let  $\{X_n\}_{n=1}^{\infty}$  be a sequence of rv's defined on a probability space  $(\Omega, L, P)$ . We say that  $X_n$  converges in probability to a rv X  $(X_n \xrightarrow{p} X, P-\lim_{n\to\infty} X_n = X)$  if

$$\lim_{n \to \infty} P(\mid X_n - X \mid > \epsilon) = 0 \quad \forall \epsilon > 0.$$

#### Note:

The following are equivalent:

$$\lim_{n \to \infty} P(|X_n - X| > \epsilon) = 0$$

$$\iff \lim_{n \to \infty} P(|X_n - X| \le \epsilon) = 1$$

$$\iff \lim_{n \to \infty} P(\{\omega : |X_n(\omega) - X(\omega)| > \epsilon)) = 0$$

If X is degenerate, i.e., P(X=c)=1, we say that  $X_n$  is **consistent** for c. For example, let  $X_n$  such that  $P(X_n=0)=1-\frac{1}{n}$  and  $P(X_n=1)=\frac{1}{n}$ . Then

$$P(\mid X_n \mid > \epsilon) = \begin{cases} \frac{1}{n}, & 0 < \epsilon < 1\\ 0, & \epsilon \ge 1 \end{cases}$$

Therefore,  $\lim_{n\to\infty} P(|X_n| > \epsilon) = 0 \quad \forall \epsilon > 0$ . So  $X_n \xrightarrow{p} 0$ , i.e.,  $X_n$  is consistent for 0.

#### Theorem 6.1.11:

(i) 
$$X_n \xrightarrow{p} X \iff X_n - X \xrightarrow{p} 0$$
.

(ii) 
$$X_n \xrightarrow{p} X, X_n \xrightarrow{p} Y \Longrightarrow P(X = Y) = 1.$$

(iii) 
$$X_n \xrightarrow{p} X, X_m \xrightarrow{p} X \Longrightarrow X_n - X_m \xrightarrow{p} 0 \text{ as } n, m \to \infty.$$

(iv) 
$$X_n \xrightarrow{p} X, Y_n \xrightarrow{p} Y \Longrightarrow X_n \pm Y_n \xrightarrow{p} X \pm Y$$
.

(v) 
$$X_n \xrightarrow{p} X, k \in \mathbb{R}$$
 a constant  $\Longrightarrow kX_n \xrightarrow{p} kX$ .

(vi) 
$$X_n \xrightarrow{p} k, k \in \mathbb{R}$$
 a constant  $\Longrightarrow X_n^r \xrightarrow{p} k^r \ \forall r \in \mathbb{N}$ .

(vii) 
$$X_n \xrightarrow{p} a, Y_n \xrightarrow{p} b, \ a, b \in \mathbb{R} \Longrightarrow X_n Y_n \xrightarrow{p} ab.$$

(viii) 
$$X_n \xrightarrow{p} 1 \Longrightarrow X_n^{-1} \xrightarrow{p} 1$$
.

$$\text{(ix)} \ \ X_n \overset{p}{\longrightarrow} a, Y_n \overset{p}{\longrightarrow} b, \ \ a \in I\!\!R, b \in I\!\!R - \{0\} \Longrightarrow \frac{X_n}{Y_n} \overset{p}{\longrightarrow} \frac{a}{b}.$$

(x)  $X_n \xrightarrow{p} X, Y$  an arbitrary rv  $\Longrightarrow X_n Y \xrightarrow{p} XY$ .

(xi) 
$$X_n \xrightarrow{p} X, Y_n \xrightarrow{p} Y \Longrightarrow X_n Y_n \xrightarrow{p} XY$$
.

#### Proof:

See Rohatgi, page 244–245 for partial proofs.

#### <u>Theorem 6.1.12:</u>

Let  $X_n \stackrel{p}{\longrightarrow} X$  and let g be a continuous function on  $\mathbb{R}$ . Then  $g(X_n) \stackrel{p}{\longrightarrow} g(X)$ .

# <u>Pro</u>of:

Preconditions:

1.) 
$$X \text{ rv} \Longrightarrow \forall \epsilon > 0 \ \exists k = k(\epsilon) : P(|X| > k) < \frac{\epsilon}{2}$$

2.) g is continuous on IR

 $\implies$  g is also uniformly continuous on [-k, k] (see Definition of uniformly continuous in Theorem 3.3.3 (iii))

$$\implies \exists \delta = \delta(\epsilon, k) : |X| \le k, |X_n - X| < \delta \implies |g(X_n) - g(X)| < \epsilon$$

Let

$$A = \{|X| \le k\} = \{\omega : |X(\omega)| \le k\}$$

$$B = \{|X_n - X| < \delta\} = \{\omega : |X_n(\omega) - X(\omega)| < \delta\}$$

$$C = \{|g(X_n) - g(X)| < \epsilon\} = \{\omega : |g(X_n(\omega)) - g(X(\omega))| < \epsilon\}$$

If  $\omega \in A \cap B$ 

$$\overset{2.)}{\Longrightarrow} \omega \in C$$

$$\Longrightarrow A \cap B \subseteq C$$

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

$$\Longrightarrow P(C^C) < P(A^C \cup B^C) < P(A^C) + P(B^C)$$

Now:

$$P(|g(X_n) - g(X)| \ge \epsilon) \le \underbrace{P(|X| > k)}_{\le \frac{\epsilon}{2} \text{ by } 1.)} + \underbrace{P(|X_n - X| \ge \delta)}_{\le \frac{\epsilon}{2} \text{ for } n \ge n_0(\epsilon, \delta, k) \text{ since } X_n \stackrel{p}{\longrightarrow} X}_{\le \epsilon}$$

$$\le \epsilon \quad \text{for } n \ge n_0(\epsilon, \delta, k)$$

#### Corollary 6.1.13:

Let  $X_n \xrightarrow{p} c, c \in \mathbb{R}$  and let g be a continuous function on  $\mathbb{R}$ . Then  $g(X_n) \xrightarrow{p} g(c)$ .

#### <u>Theorem 6.1.14:</u>

$$\overline{X_n \xrightarrow{p} X \Longrightarrow X_n \xrightarrow{d} X}.$$

#### Proof:

$$X_n \xrightarrow{p} X \Leftrightarrow P(|X_n - X| > \epsilon) \to 0 \text{ as } n \to \infty \quad \forall \epsilon > 0$$

It holds:

$$P(X \le x - \epsilon) = P(X \le x - \epsilon, |X_n - X| \le \epsilon) + P(X \le x - \epsilon, |X_n - X| > \epsilon)$$

$$\stackrel{(A)}{\le} P(X_n \le x) + P(|X_n - X| > \epsilon)$$

(A) holds since  $X \leq x - \epsilon$  and  $X_n$  within  $\epsilon$  of X, thus  $X_n \leq x$ .

Similarly, it holds:

$$P(X_n \le x) = P(X_n \le x, |X_n - X| \le \epsilon) + P(X_n \le x, |X_n - X| > \epsilon)$$
  
$$\le P(X \le x + \epsilon) + P(|X_n - X| > \epsilon)$$

Combining the 2 inequalities from above gives:

$$P(X \le x - \epsilon) - \underbrace{P(|X_n - X| > \epsilon)}_{\to 0 \text{ as } n \to \infty} \le \underbrace{P(X_n \le x)}_{=F_n(x)} \le P(X \le x + \epsilon) + \underbrace{P(|X_n - X| > \epsilon)}_{\to 0 \text{ as } n \to \infty}$$

Therefore,

$$P(X \le x - \epsilon) \le F_n(x) \le P(X \le x + \epsilon) \text{ as } n \to \infty.$$

Since the cdf's  $F_n(\cdot)$  are not necessarily left continuous, we get the following result for  $\epsilon \downarrow 0$ :

$$P(X < x) < F_n(x) < P(X < x) = F_X(x)$$

Let x be a continuity point of F. Then it holds:

$$F(x) = P(X < x) \le F_n(x) \le F(x)$$

$$\Longrightarrow F_n(x) \to F(x)$$

$$\Longrightarrow X_n \stackrel{d}{\longrightarrow} X$$

#### <u>Theorem 6.1.15:</u>

Let  $c \in \mathbb{R}$  be a constant. Then it holds:

$$X_n \xrightarrow{d} c \iff X_n \xrightarrow{p} c.$$

### Example 6.1.16:

In this example, we will see that

$$X_n \xrightarrow{d} X \not\Longrightarrow X_n \xrightarrow{p} X$$

for some rv X. Let  $X_n$  be identically distributed rv's and let  $(X_n, X)$  have the following joint distribution:

$$\begin{array}{c|cccc} X_n & 0 & 1 \\ \hline X & & & \\ \hline 0 & 0 & \frac{1}{2} & \frac{1}{2} \\ 1 & \frac{1}{2} & 0 & \frac{1}{2} \\ \hline & \frac{1}{2} & \frac{1}{2} & 1 \\ \hline \end{array}$$

Obviously,  $X_n \stackrel{d}{\longrightarrow} X$  since all have exactly the same cdf, but for any  $\epsilon \in (0,1)$ , it is

$$P(|X_n - X| > \epsilon) = P(|X_n - X| = 1) = 1 \ \forall n,$$

so 
$$\lim_{n\to\infty} P(|X_n - X| > \epsilon) \neq 0$$
. Therefore,  $X_n \not\stackrel{p}{\longrightarrow} X$ .

#### <u>Theorem 6.1.17:</u>

Let  $\{X_n\}_{n=1}^{\infty}$  and  $\{Y_n\}_{n=1}^{\infty}$  be sequences of rv's and X be a rv defined on a probability space  $(\Omega, L, P)$ . Then it holds:

$$Y_n \xrightarrow{d} X$$
,  $|X_n - Y_n| \xrightarrow{p} 0 \Longrightarrow X_n \xrightarrow{d} X$ .

#### Proof:

Similar to the proof of Theorem 6.1.14. See also Rohatgi, page 253, Theorem 14.

### Theorem 6.1.18: Slutsky's Theorem

Let  $\{X_n\}_{n=1}^{\infty}$  and  $\{Y_n\}_{n=1}^{\infty}$  be sequences of rv's and X be a rv defined on a probability space  $(\Omega, L, P)$ . Let  $c \in \mathbb{R}$  be a constant. Then it holds:

(i) 
$$X_n \xrightarrow{d} X, Y_n \xrightarrow{p} c \Longrightarrow X_n + Y_n \xrightarrow{d} X + c$$
.

(ii) 
$$X_n \xrightarrow{d} X, Y_n \xrightarrow{p} c \Longrightarrow X_n Y_n \xrightarrow{d} cX$$
.  
If  $c = 0$ , then also  $X_n Y_n \xrightarrow{p} 0$ .

(iii) 
$$X_n \xrightarrow{d} X, Y_n \xrightarrow{p} c \Longrightarrow \frac{X_n}{Y_n} \xrightarrow{d} \frac{X}{c}$$
 if  $c \neq 0$ .

Proof:

(i) 
$$Y_n \xrightarrow{p} c \xrightarrow{Th.6.1.11(i)} Y_n - c \xrightarrow{p} 0$$
  
 $\Longrightarrow Y_n - c = Y_n + (X_n - X_n) - c = (X_n + Y_n) - (X_n + c) \xrightarrow{p} 0$  (A)  
 $X_n \xrightarrow{d} X \xrightarrow{Th.6.1.9(i)} X_n + c \xrightarrow{d} X + c$  (B)

Combining (A) and (B), it follows from Theorem 6.1.17:

$$X_n + Y_n \stackrel{d}{\longrightarrow} X + c$$

(ii) Case c = 0:

 $\forall \epsilon > 0 \ \forall k > 0$ , it is

$$P(|X_n Y_n| > \epsilon) = P(|X_n Y_n| > \epsilon, Y_n \le \frac{\epsilon}{k}) + P(|X_n Y_n| > \epsilon, Y_n > \frac{\epsilon}{k})$$

$$\le P(|X_n \frac{\epsilon}{k}| > \epsilon) + P(Y_n > \frac{\epsilon}{k})$$

$$\le P(|X_n| > k) + P(|Y_n| > \frac{\epsilon}{k})$$

Since  $X_n \stackrel{d}{\longrightarrow} X$  and  $Y_n \stackrel{p}{\longrightarrow} 0$ , it follows

$$\overline{\lim_{n\to\infty}} P(|X_n Y_n| > \epsilon) \le P(|X_n| > k) \to 0 \text{ as } k \to \infty.$$

Therefore,  $X_n Y_n \xrightarrow{p} 0$ .

Case  $c \neq 0$ :

Since  $X_n \xrightarrow{d} X$  and  $Y_n \xrightarrow{p} c$ , it follows from (ii), case c = 0, that  $X_n Y_n - c X_n = X_n (Y_n - c) \xrightarrow{p} 0$ .

$$\Longrightarrow X_n Y_n \xrightarrow{p} cX_n$$

$$\xrightarrow{Th.6.1.14} X_n Y_n \xrightarrow{d} cX_n$$

Since  $cX_n \xrightarrow{d} cX$  by Theorem 6.1.9 (ii), it follows from Theorem 6.1.17:

$$X_n Y_n \stackrel{d}{\longrightarrow} cX$$

(iii) Let  $Z_n \stackrel{p}{\longrightarrow} 1$  and let  $Y_n = cZ_n$ .

$$\begin{array}{c} \overset{c\neq 0}{\Longrightarrow} \frac{1}{Y_n} = \frac{1}{Z_n} \cdot \frac{1}{c} \\ \\ \overset{Th.6.1.11(v,viii)}{\Longrightarrow} \frac{1}{Y_n} \stackrel{p}{\longrightarrow} \frac{1}{c} \end{array}$$

With part (ii) above, it follows:

$$X_n \xrightarrow{d} X$$
 and  $\frac{1}{Y_n} \xrightarrow{p} \frac{1}{c}$ 

$$\Longrightarrow \frac{X_n}{Y_n} \xrightarrow{d} \frac{X}{c}$$

#### Definition 6.1.19:

Let  $\{X_n\}_{n=1}^{\infty}$  be a sequence of rv's such that  $E(|X_n|^r) < \infty$  for some r > 0. We say that  $X_n$  converges in the  $r^{th}$  mean to a rv X  $(X_n \xrightarrow{r} X)$  if  $E(|X|^r) < \infty$  and

$$\lim_{n\to\infty} E(\mid X_n - X\mid^r) = 0.$$

#### Example 6.1.20:

Let  $\{X_n\}_{n=1}^{\infty}$  be a sequence of rv's defined by  $P(X_n=0)=1-\frac{1}{n}$  and  $P(X_n=1)=\frac{1}{n}$ . It is  $E(\mid X_n\mid^r)=\frac{1}{n} \ \forall r>0$ . Therefore,  $X_n\stackrel{r}{\longrightarrow} 0 \ \forall r>0$ .

#### Note:

The special cases r=1 and r=2 are called **convergence in absolute mean** for r=1  $(X_n \xrightarrow{1} X)$  and **convergence in mean square** for r=2  $(X_n \xrightarrow{ms} X \text{ or } X_n \xrightarrow{2} X)$ .

#### <u>Theorem 6.1.21:</u>

Assume that  $X_n \stackrel{r}{\longrightarrow} X$  for some r > 0. Then  $X_n \stackrel{p}{\longrightarrow} X$ .

#### Proof:

Using Markov's Inequality (Corollary 3.5.2), it holds for any  $\epsilon > 0$ :

$$\frac{E(\mid X_n - X\mid^r)}{\epsilon^r} \ge P(\mid X_n - X\mid \ge \epsilon)$$

$$X_n \xrightarrow{r} X \Longrightarrow \lim_{n \to \infty} E(|X_n - X|^r) = 0$$

$$\Longrightarrow \lim_{n \to \infty} P(|X_n - X| \ge \epsilon) \le \lim_{n \to \infty} \frac{E(|X_n - X|^r)}{\epsilon^r} = 0$$

$$\Longrightarrow X_n \xrightarrow{p} X$$

#### Example 6.1.22:

Let  $\{X_n\}_{n=1}^{\infty}$  be a sequence of rv's defined by  $P(X_n=0)=1-\frac{1}{n^r}$  and  $P(X_n=n)=\frac{1}{n^r}$  for some r>0.

For any  $\epsilon > 0$ ,  $P(|X_n| > \epsilon) \to 0$  as  $n \to \infty$ ; so  $X_n \stackrel{p}{\longrightarrow} 0$ .

For 0 < s < r,  $E(|X_n|^s) = \frac{1}{n^{r-s}} \to 0$  as  $n \to \infty$ ; so  $X_n \xrightarrow{s} 0$ . But  $E(|X_n|^r) = 1 \not\to 0$  as  $n \to \infty$ ; so  $X_n \not\stackrel{r}{\longrightarrow} 0$ .

#### Theorem 6.1.23:

If  $X_n \stackrel{r}{\longrightarrow} X$ , then it holds:

(i) 
$$\lim_{n\to\infty} E(\mid X_n\mid^r) = E(\mid X\mid^r)$$
; and

(ii) 
$$X_n \stackrel{s}{\longrightarrow} X$$
 for  $0 < s < r$ .

#### Proof:

(i) For  $0 < r \le 1$ , it holds:

$$E(||X_{n}||^{r}) = E(||X_{n} - X + X||^{r}) \stackrel{(*)}{\leq} E(||X_{n} - X||^{r} + ||X||^{r})$$

$$\Longrightarrow E(||X_{n}||^{r}) - E(||X||^{r}) \leq E(||X_{n} - X||^{r})$$

$$\Longrightarrow \lim_{n \to \infty} E(||X_{n}||^{r}) - \lim_{n \to \infty} E(||X||^{r}) \leq \lim_{n \to \infty} E(||X_{n} - X||^{r}) = 0$$

$$\Longrightarrow \lim_{n \to \infty} E(||X_{n}||^{r}) \leq E(||X||^{r}) \quad (A)$$

(\*) holds due to Bronstein/Semendjajew (1986), page 36 (see Handout)

Similarly,

$$E(\mid X\mid^{r}) = E(\mid X - X_{n} + X_{n}\mid^{r}) \leq E(\mid X_{n} - X\mid^{r} + \mid X_{n}\mid^{r})$$

$$\Longrightarrow E(\mid X\mid^{r}) - E(\mid X_{n}\mid^{r}) \leq E(\mid X_{n} - X\mid^{r})$$

$$\Longrightarrow \lim_{n \to \infty} E(\mid X\mid^{r}) - \lim_{n \to \infty} E(\mid X_{n}\mid^{r}) \leq \lim_{n \to \infty} E(\mid X_{n} - X\mid^{r}) = 0$$

$$\Longrightarrow E(\mid X\mid^{r}) \leq \lim_{n \to \infty} E(\mid X_{n}\mid^{r}) \quad (B)$$

Combining (A) and (B) gives

$$\lim_{n \to \infty} E(\mid X_n \mid^r) = E(\mid X \mid^r)$$

For r > 1, it follows from Minkowski's Inequality (Theorem 4.8.3):

$$\begin{split} &[E(\mid X - X_n + X_n \mid^r)]^{\frac{1}{r}} \leq [E(\mid X - X_n \mid^r)]^{\frac{1}{r}} + [E(\mid X_n \mid^r)]^{\frac{1}{r}} \\ &\Longrightarrow [E(\mid X \mid^r)]^{\frac{1}{r}} - [E(\mid X_n \mid^r)]^{\frac{1}{r}} \leq [E(\mid X - X_n \mid^r)]^{\frac{1}{r}} \\ &\Longrightarrow [E(\mid X \mid^r)]^{\frac{1}{r}} - \lim_{n \to \infty} [E(\mid X_n \mid^r)]^{\frac{1}{r}} \leq \lim_{n \to \infty} [E(\mid X_n - X \mid^r)]^{\frac{1}{r}} = 0 \text{ since } X_n \xrightarrow{r} X \\ &\Longrightarrow [E(\mid X \mid^r)]^{\frac{1}{r}} \leq \lim_{n \to \infty} [E(\mid X_n \mid^r)]^{\frac{1}{r}} \quad (C) \end{split}$$

Similarly,

$$\begin{split} &[E(\mid X_{n}-X+X\mid^{r})]^{\frac{1}{r}} \leq [E(\mid X_{n}-X\mid^{r})]^{\frac{1}{r}} + [E(\mid X\mid^{r})]^{\frac{1}{r}} \\ &\Longrightarrow \lim_{n \to \infty} [E(\mid X_{n}\mid^{r})]^{\frac{1}{r}} - \lim_{n \to \infty} [E(\mid X\mid^{r})]^{\frac{1}{r}} \leq \lim_{n \to \infty} [E(\mid X_{n}-X\mid^{r})]^{\frac{1}{r}} = 0 \text{ since } X_{n} \xrightarrow{r} X \end{split}$$

$$\Longrightarrow \lim_{n \to \infty} [E(\mid X_n \mid^r)]^{\frac{1}{r}} \le [E(\mid X \mid^r)]^{\frac{1}{r}} \quad (D)$$

Combining (C) and (D) gives

$$\lim_{n \to \infty} [E(\mid X_n \mid^r)]^{\frac{1}{r}} = [E(\mid X \mid^r)]^{\frac{1}{r}}$$

$$\implies \lim_{n \to \infty} E(\mid X_n \mid^r) = E(\mid X \mid^r)$$

(ii) For  $1 \le s < r$ , it follows from Lyapunov's Inequality (Theorem 3.5.4):

Lecture 42/1: Fr 12/08/00

$$\begin{aligned} &[E(\mid X_n - X\mid^s)]^{\frac{1}{s}} \leq [E(\mid X_n - X\mid^r)]^{\frac{1}{r}} \\ &\Longrightarrow E(\mid X_n - X\mid^s) \leq [E(\mid X_n - X\mid^r)]^{\frac{s}{r}} \\ &\Longrightarrow \lim_{n \to \infty} E(\mid X_n - X\mid^s) \leq \lim_{n \to \infty} [E(\mid X_n - X\mid^r)]^{\frac{s}{r}} = 0 \quad \text{since } X_n \stackrel{r}{\longrightarrow} X \\ &\Longrightarrow X_n \stackrel{s}{\longrightarrow} X \end{aligned}$$

An additional proof is required for 0 < s < r < 1.

#### Definition 6.1.24:

Let  $\{X_n\}_{n=1}^{\infty}$  be a sequence of rv's on  $(\Omega, L, P)$ . We say that  $X_n$  converges almost surely to a rv X  $(X_n \xrightarrow{a.s.} X)$  or  $X_n$  converges with probability 1 to X  $(X_n \xrightarrow{w.p.1} X)$  or  $X_n$  converges strongly to X iff

$$P(\{\omega: X_n(\omega) \to X(\omega) \text{ as } n \to \infty\}) = 1.$$

#### Note:

An interesting characterization of convergence with probability 1 and convergence in probability can be found in Parzen (1960) "Modern Probability Theory and Its Applications" on page 416 (see Handout).

#### Example 6.1.25:

Let  $\Omega = [0, 1]$  and P a uniform distribution on  $\Omega$ . Let  $X_n(\omega) = \omega + \omega^n$  and  $X(\omega) = \omega$ .

For  $\omega \in [0, 1)$ ,  $\omega^n \to 0$  as  $n \to \infty$ . So  $X_n(\omega) \to X(\omega) \ \forall \omega \in [0, 1)$ .

However, for  $\omega = 1$ ,  $X_n(1) = 2 \neq 1 = X(1) \ \forall n$ , i.e., convergence fails at  $\omega = 1$ .

Anyway, since  $P(\{\omega: X_n(\omega) \to X(\omega) \text{ as } n \to \infty\}) = P(\{\omega \in [0,1)\}) = 1$ , it is  $X_n \xrightarrow{a.s.} X$ .

#### Theorem 6.1.26:

$$X_n \xrightarrow{a.s.} X \Longrightarrow X_n \xrightarrow{p} X.$$

#### Proof:

Choose  $\epsilon > 0$  and  $\delta > 0$ . Find  $n_0 = n_0(\epsilon, \delta)$  such that

$$P\left(\bigcap_{n=n_0}^{\infty} \{\mid X_n - X \mid \leq \epsilon\}\right) \geq 1 - \delta.$$

Since  $\bigcap_{n=0}^{\infty} \{|X_n - X| \le \epsilon\} \subseteq \{|X_n - X| \le \epsilon\} \quad \forall n \ge n_0, \text{ it is}$ 

$$P\left(\left\{\mid X_n - X \mid \leq \epsilon\right\}\right) \geq P\left(\bigcap_{n=n_0}^{\infty} \left\{\mid X_n - X \mid \leq \epsilon\right\}\right) \geq 1 - \delta \ \forall n \geq n_0.$$

Therefore,  $P(\{|X_n - X| \le \epsilon\}) \to 1 \text{ as } n \to \infty.$  Thus,  $X_n \stackrel{p}{\longrightarrow} X$ .

$$\frac{\text{Example 6.1.27:}}{X_n \xrightarrow{p} X \not\Longrightarrow X_n \xrightarrow{a.s.} X:}$$

Let  $\Omega = (0, 1]$  and P a uniform distribution on  $\Omega$ .

Define  $A_n$  by

$$A_1 = (0, \frac{1}{2}], A_2 = (\frac{1}{2}, 1]$$

$$A_3 = (0, \frac{1}{4}], A_4 = (\frac{1}{4}, \frac{1}{2}], A_5 = (\frac{1}{2}, \frac{3}{4}], A_6 = (\frac{3}{4}, 1]$$

$$A_7 = (0, \frac{1}{8}], A_8 = (\frac{1}{8}, \frac{1}{4}], \dots$$

Let  $X_n(\omega) = I_{A_n}(\omega)$ .

It is  $P(|X_n - 0| \ge \epsilon) \to 0 \ \forall \epsilon > 0$  since  $X_n$  is 0 except on  $A_n$  and  $P(A_n) \downarrow 0$ . Thus  $X_n \stackrel{p}{\longrightarrow} 0$ .

But  $P(\{\omega: X_n(\omega) \to 0\}) = 0$  (and not 1) because any  $\omega$  keeps being in some  $A_n$  beyond any  $n_0$ , i.e.,  $X_n(\omega)$  looks like  $0 \dots 010 \dots 010 \dots 010 \dots$ , so  $X_n \stackrel{q.s.}{\longrightarrow} 0$ .

$$\frac{\text{Example 6.1.28:}}{X_n \xrightarrow{r} X \not\Longrightarrow X_n \xrightarrow{a.s.} X:}$$

Let  $X_n$  be independent rv's such that  $P(X_n = 0) = 1 - \frac{1}{n}$  and  $P(X_n = 1) = \frac{1}{n}$ .

It is 
$$E(\mid X_n - 0 \mid^r) = E(\mid X_n \mid^r) = E(\mid X_n \mid) = \frac{1}{n} \to 0$$
 as  $n \to \infty$ , so  $X_n \stackrel{r}{\longrightarrow} 0 \quad \forall r > 0$ .

But

$$P(X_n = 0 \ \forall m \le n \le n_0) = \prod_{n=m}^{n_0} (1 - \frac{1}{n}) = (\frac{m-1}{m}) (\frac{m}{m+1}) (\frac{m+1}{m+2}) \dots (\frac{n_0-2}{n_0-1}) (\frac{n_0-1}{n_0}) = \frac{m-1}{n_0}$$

As 
$$n_0 \to \infty$$
, it is  $P(X_n = 0 \ \forall m \le n \le n_0) \to 0 \ \forall m$ , so  $X_n \xrightarrow{q_s \cdot s} 0$ .

$$\frac{\text{Example 6.1.29:}}{X_n \xrightarrow{a.s.} X \not\Longrightarrow X_n \xrightarrow{r} X:}$$

Let  $\Omega = [0, 1]$  and P a uniform distribution on  $\Omega$ .

Let 
$$A_n = [0, \frac{1}{\ln n}].$$

Let  $X_n(\omega) = nI_{A_n}(\omega)$  and  $X(\omega) = 0$ .

It holds that  $\forall \omega > 0 \quad \exists n_0 : \frac{1}{\ln n_0} < \omega \Longrightarrow X_n(\omega) = 0 \quad \forall n > n_0 \text{ and } P(\omega = 0) = 0.$  Thus,  $X_n \xrightarrow{a.s.} 0.$ 

But 
$$E(|X_n - 0|^r) = \frac{n^r}{\ln n} \to \infty \quad \forall r > 0$$
, so  $X_n \not\stackrel{r}{\longrightarrow} X$ .

Lecture 39: Fr 12/01/00

# 6.2 Weak Laws of Large Numbers

### Theorem 6.2.1: WLLN: Version I

Let  $\{X_i\}_{i=1}^{\infty}$  be a sequence of iid rv's with mean  $E(X_i) = \mu$  and variance  $Var(X_i) = \sigma^2 < \infty$ . Let  $\overline{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$ . Then it holds

$$\lim_{n \to \infty} P(|\overline{X}_n - \mu| \ge \epsilon) = 0,$$

i.e.,  $\overline{X}_n \stackrel{p}{\longrightarrow} \mu$ .

#### Proof:

By Markov's Inequality, it holds for all  $\epsilon > 0$ :

$$P(|\overline{X}_n - \mu| \ge \epsilon) \le \frac{E((\overline{X}_n - \mu)^2)}{\epsilon^2} = \frac{Var(\overline{X}_n)}{\epsilon^2} = \frac{\sigma^2}{n\epsilon^2} \longrightarrow 0 \text{ as } n \to \infty$$

#### Note:

For iid rv's with finite variance,  $\overline{X}_n$  is consistent for  $\mu$ .

A more general way to derive a "WLLN" follows in the next Definition.

#### Definition 6.2.2:

Let  $\{X_i\}_{i=1}^{\infty}$  be a sequence of rv's. Let  $T_n = \sum_{i=1}^n X_i$ . We say that  $\{X_i\}$  obeys the WLLN with respect to a sequence of **norming constants**  $\{B_i\}_{i=1}^{\infty}$ ,  $B_i > 0$ ,  $B_i \uparrow \infty$ , if there exists a sequence of **centering constants**  $\{A_i\}_{i=1}^{\infty}$  such that

$$B_n^{-1}(T_n - A_n) \stackrel{p}{\longrightarrow} 0.$$

#### Theorem 6.2.3:

Let  $\{X_i\}_{i=1}^{\infty}$  be a sequence of pairwise uncorrelated rv's with  $E(X_i) = \mu_i$  and  $Var(X_i) = \sigma_i^2$ ,  $i \in \mathbb{N}$ . If  $\sum_{i=1}^{n} \sigma_i^2 \to \infty$  as  $n \to \infty$ , we can choose  $A_n = \sum_{i=1}^{n} \mu_i$  and  $B_n = \sum_{i=1}^{n} \sigma_i^2$  and get

$$\frac{\sum_{i=1}^{n} (X_i - \mu_i)}{\sum_{i=1}^{n} \sigma_i^2} \stackrel{p}{\longrightarrow} 0.$$

#### Proof:

By Markov's Inequality, it holds for all  $\epsilon > 0$ :

$$P(|\sum_{i=1}^{n} X_{i} - \sum_{i=1}^{n} \mu_{i}| > \epsilon \sum_{i=1}^{n} \sigma_{i}^{2}) \le \frac{E((\sum_{i=1}^{n} (X_{i} - \mu))^{2})}{\epsilon^{2} (\sum_{i=1}^{n} \sigma_{i}^{2})^{2}} = \frac{1}{\epsilon^{2} \sum_{i=1}^{n} \sigma_{i}^{2}} \longrightarrow 0 \text{ as } n \to \infty$$

#### Note:

To obtain Theorem 6.2.1, we choose  $A_n = n\mu$  and  $B_n = n\sigma^2$ .

#### Theorem 6.2.4:

Let  $\{X_i\}_{i=1}^{\infty}$  be a sequence of rv's. Let  $\overline{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$ . A necessary and sufficient condition for  $\{X_i\}$  to obey the WLLN with respect to  $B_n = n$  is that

$$E\left(\frac{\overline{X}_n^2}{1+\overline{X}_n^2}\right) \to 0$$

as  $n \to \infty$ .

#### Proof:

Rohatgi, page 258, Theorem 2.

#### Example 6.2.5:

Let  $(X_1, \ldots, X_n)$  be jointly Normal with  $E(X_i) = 0$ ,  $E(X_i^2) = 1$  for all i, and  $Cov(X_i, X_j) = \rho$  if |i-j| = 1 and  $Cov(X_i, X_j) = 0$  if |i-j| > 1. Then,  $T_n \sim N(0, n + 2(n-1)\rho) = N(0, \sigma^2)$ .

$$\begin{split} E\left(\frac{\overline{X}_n^2}{1+\overline{X}_n^2}\right) &= E\left(\frac{T_n^2}{n^2+T_n^2}\right) \\ &= \frac{2}{\sqrt{2\pi}\sigma} \int_0^\infty \frac{x^2}{n^2+x^2} e^{-\frac{x^2}{2\sigma^2}} dx \quad | \quad y = \frac{x}{\sigma}, \quad dy = \frac{dx}{\sigma} \\ &= \frac{2}{\sqrt{2\pi}} \int_0^\infty \frac{\sigma^2 y^2}{n^2+\sigma^2 y^2} e^{-\frac{y^2}{2}} dy \\ &= \frac{2}{\sqrt{2\pi}} \int_0^\infty \frac{(n+2(n-1)\rho)y^2}{n^2+(n+2(n-1)\rho)y^2} e^{-\frac{y^2}{2}} dy \\ &\leq \frac{n+2(n-1)\rho}{n^2} \underbrace{\int_0^\infty \frac{2}{\sqrt{2\pi}} y^2 e^{-\frac{y^2}{2}} dy}_{=1, \text{ since Var of } N(0,1) \text{ distribution}} \\ &\to 0 \quad \text{as } n \to \infty \\ \Longrightarrow \overline{X}_n \stackrel{p}{\longrightarrow} 0 \end{split}$$

#### Note:

We would like to have a WLLN that just depends on means but does not depend on the existence of finite variances. To approach this, we consider the following:

Let  $\{X_i\}_{i=1}^{\infty}$  be a sequence of rv's. Let  $T_n = \sum_{i=1}^n X_i$ . We truncate each  $X_i$  at c > 0 and get

$$X_i^c = \begin{cases} X_i, & |X_i| \le c \\ 0, & \text{otherwise} \end{cases}$$

Let 
$$T_n^c = \sum_{i=1}^n X_i^c$$
 and  $m_n = \sum_{i=1}^n E(X_i^c)$ .

#### Lemma 6.2.6:

For  $T_n$ ,  $T_n^c$  and  $m_n$  as defined in the Note above, it holds:

$$P(|T_n - m_n| > \epsilon) \le P(|T_n^c - m_n| > \epsilon) + \sum_{i=1}^n P(|X_i| > c) \ \forall \epsilon > 0$$

#### Proof:

It holds for all  $\epsilon > 0$ :

$$P(\mid T_{n} - m_{n} \mid > \epsilon) = P(\mid T_{n} - m_{n} \mid > \epsilon \text{ and } \mid X_{i} \mid \leq c \quad \forall i \in \{1, \dots, n\}) +$$

$$P(\mid T_{n} - m_{n} \mid > \epsilon \text{ and } \mid X_{i} \mid > c \text{ for at least one } i \in \{1, \dots, n\}\}$$

$$\stackrel{(*)}{\leq} P(\mid T_{n}^{c} - m_{n} \mid > \epsilon) + P(\mid X_{i} \mid > c \text{ for at least one } i \in \{1, \dots, n\})$$

$$\leq P(\mid T_{n}^{c} - m_{n} \mid > \epsilon) + \sum_{i=1}^{n} P(\mid X_{i} \mid > c)$$

(\*) holds since 
$$T_n^c = T_n$$
 when  $|X_i| \le c \ \forall i \in \{1, ..., n\}$ .

#### Note:

If the  $X_i$ 's are identically distributed, then

$$P(\mid T_n - m_n \mid > \epsilon) \le P(\mid T_n^c - m_n \mid > \epsilon) + nP(\mid X_1 \mid > c) \quad \forall \epsilon > 0.$$

If the  $X_i$ 's are iid, then

$$P(|T_n - m_n| > \epsilon) \le \frac{nE((X_1^c)^2)}{\epsilon^2} + nP(|X_1| > c) \quad \forall \epsilon > 0 \quad (*).$$

Note that  $P(|X_i|>c)=P(|X_1|>c) \ \forall i\in \mathbb{N}$  if the  $X_i$ 's are identically distributed and that  $E((X_i^c)^2)=E((X_1^c)^2) \ \forall i\in \mathbb{N}$  if the  $X_i$ 's are iid.

Lecture 42/2: Fr 12/08/00

#### Theorem 6.2.7: Khintchine's WLLN

Let  $\{X_i\}_{i=1}^{\infty}$  be a sequence of iid rv's with finite mean  $E(X_i) = \mu$ . Then it holds:

$$\overline{X}_n = \frac{1}{n} T_n \stackrel{p}{\longrightarrow} \mu$$

#### Proof:

If we take c = n and replace  $\epsilon$  by  $n\epsilon$  in (\*) in the Note above, we get

$$P\left(\left|\frac{T_n - m_n}{n}\right| > \epsilon\right) = P(|T_n - m_n| > n\epsilon) \le \frac{E((X_1^n)^2)}{n\epsilon^2} + nP(|X_1| > n).$$

Since  $E(\mid X_1\mid)<\infty$ , it is  $nP(\mid X_1\mid>n)\to 0$  as  $n\to\infty$  by Theorem 3.1.9. From Corollary 3.1.12 we know that  $E(\mid X\mid^{\alpha})=\alpha\int_0^\infty x^{\alpha-1}P(\mid X\mid>x)dx$ . Therefore,

$$\begin{split} E((X_1^n)^2) &= 2\int_0^n x P(\mid X_1^n\mid > x) dx \\ &= 2\int_0^A x P(\mid X_1^n\mid > x) dx + 2\int_A^n x P(\mid X_1^n\mid > x) dx \\ \overset{(+)}{\leq} K + \delta\int_A^n dx \\ &< K + n\delta \end{split}$$

In (+), A is chosen sufficiently large such that  $xP(\mid X_1^n\mid>x)<\frac{\delta}{2} \ \forall x\geq A$  for an arbitrary constant  $\delta>0$  and K>0 a constant.

Therefore,

$$\frac{E((X_1^n)^2)}{n\epsilon^2} \le \frac{K}{n\epsilon^2} + \frac{\delta}{\epsilon^2}$$

Since  $\delta$  is arbitrary, we can make the right hand side of this last inequality arbitrarily small for sufficiently large n.

Since 
$$E(X_i) = \mu \ \forall i$$
, it is  $\frac{m_n}{n} \to \mu$  as  $n \to \infty$ .

#### Note:

Theorem 6.2.7 meets the previously stated goal of not having a finite variance requirement.

# Merry Xmas and a Happy New Year!

## Index

 $\sigma$ -Algebra, 2 Consistent, 125  $\sigma$ -Field, 2 Continuity of Sets, 9  $\sigma$ -Field, Borel, 4 Continuous, 29, 31, 32  $\sigma$ -Field, Generated, 2 Continuous RV, 73 k-Parameter Exponential Family, 119 Continuous, Absolutely, 32  $k^{th}$  Central Moment, 43 Convergence, Almost Sure, 132  $k^{th}$  Factorial Moment, 66 Convergence, In  $r^{th}$  Mean, 130  $k^{th}$  Moment, 43 Convergence, In Absolute Mean, 130 Convergence, In Distribution, 122 n-Dimensional Characteristic Function, 97 n-RV, 71 Convergence, In Law, 122 Convergence, In Mean Square, 130 Absolutely Continuous, 32 Convergence, In Probability, 125 Additivity, Countable, 4 Convergence, Strong, 132 Arithmetic Mean, 109 Convergence, Weak, 122 Convergence, With Probability 1, 132 Bayes' Rule, 19 Convex, 108 Binomial Coefficient, 12 Countable Additivity, 4 Binomial Theorem, 14 Counting, Fundamental Theorem of, 11 Bivariate Normal Distribution, 112 Covariance Inequality, 111 Bivariate RV, 72 Cumulative Distribution Function, 27 Bochner's Theorem, 60 Cumulative Distribution Function, Joint, 71 Bonferroni's Inequality, 7 Boole's Inequality, 9 Discrete, 29, 31 Borel  $\sigma$ -Field, 4 Discrete RV, 72 Borel Field, 2 Distribution, Sampling, 122 Borel Sets, 4 Equally Likely, 11 Canonical Form, 119 Equivalence, 35 Cauchy-Schwarz Inequality, 95 Equivalent RV, 81 CDF, 27 Euler's Relation, 56 CDF, Conditional, 75 Event, 1 CDF, Joint, 71 Expectation, Conditional, 102 CDF, Marginal, 74 Expectation, Multivariate, 91 Centering Constants, 135 Expected Value, 42 Characteristic Function, 57 Exponential Family in Canonical Form, 119 Characteristic Function, n-Dimensional, 97 Exponential Family, k-Parameter, 119 Chebychev's Inequality, 68 Exponential Family, One-Parameter, 119 Cofinite, 2 Completely Independent, 20, 77 Factorial, 12 Complex Numbers, 56 Factorization Theorem, 78 Complex Numbers, Conjugate, 56 Field, 1 Complex-Valued Random Variable, 57 Fundamental Theorem of Counting, 11 Concave, 108 Geometric Mean, 109 Conditional CDF, 75 Conditional Cumulative Distribution Function, 75 Hölder's Inequality, 106 Conditional Expectation, 102 Harmonic Mean, 109 Conditional PDF, 75 Conditional PMF, 74 Identically Distributed, 30 Conditional Probability, 17 IID, 81 Conditional Probability Density Function, 75 Implication, 35 Conditional Probability Mass Function, 74 Inclusion-Exclusion, Principle, 6 Conjugate Complex Numbers, 56

Independent, 19, 77, 78

Independent Identically Distributed, 81 Measurable Space, 4 Independent, Completely, 20, 77 MGF, 50 Independent, Mutually, 20, 77 MIN, 82 Independent, Pairwise, 19 Minkowski's Inequality, 107 MMGF, 97 Indicator Function, 34 Induced Probability, 27 Moivre's Theorem, 56 Induction, 6 Moment Generating Function, 50 Induction Base, 6 Moment Generating Function, Multivariate, 97 Moment,  $k^{\,th}\,,~43$ Induction Step, 7 Moment, k<sup>th</sup> Central, 43 Inequality, Bonferroni's, 7 Moment,  $k^{th}$  Factorial, 66 Inequality, Boole's, 9 Inequality, Cauchy-Schwarz, 95 Moment, Multi-Way, 94 Inequality, Chebychev's, 68 Moment, Multi-Way Central, 94 Inequality, Covariance, 111 Monotonic, 37 Inequality, Hölder's, 106 Multi-Way Central Moment, 94 Inequality, Jensen's, 108 Multi-Way Moment, 94 Inequality, Lyapunov's, 69 Multiplication Rule, 18 Inequality, Markov's, 68 Multivariate Expectation, 91 Inequality, Minkowski's, 107 Multivariate Moment Generating Function, 97 Multivariate Normal Distribution, 117 Jensen's Inequality, 108 Multivariate RV, 72 Joint CDF, 71 Multivariate Transformation, 85 Joint Cumulative Distribution Function, 71 Mutually Independent, 20, 77 Joint PDF, 73 Joint PMF, 72 Non-Decreasing, 37 Joint Probability Density Function, 73 Non-Decreasing, Strictly, 37 Joint Probability Mass Function, 72 Non-Increasing, 37 Jump Points, 31 Non-Increasing, Strictly, 37 Normal Distribution, 112 Khintchine's Weak Law of Large Numbers, 138 Normal Distribution, Bivariate, 112 Kolmogorov Axioms of Probability, 4 Normal Distribution, Multivariate, 117 Normal Distribution, Standard, 112 L'Hospital's Rule, 52 Normal Distribution, Univariate, 112 Law of Total Probability, 18 Norming Constants, 135 Lebesque's Dominated Convergence Theorem, 53 Leibnitz's Rule, 52 One-Parameter Exponential Family, 119 Levy's Theorem, 60 Order Statistic, 89 Logic, 35 Pairwise Independent, 19 Lyapunov's Inequality, 69 Partition, 18 Marginal CDF, 74 PDF, 32 Marginal Cumulative Distribution Function, 74 PDF, Conditional, 75 Marginal PDF, 73 PDF, Joint, 73 Marginal PMF, 73 PDF, Marginal, 73 Marginal Probability Density Function, 73 PGF, 66 Marginal Probability Mass Function, 73 PM, 4 Markov's Inequality, 68 PMF, 31 MAX, 82 PMF, Conditional, 74 Mean, 42 PMF, Joint, 72 Mean, Arithmetic, 109 PMF, Marginal, 73 Mean, Geometric, 109 Power Set, 2 Mean, Harmonic, 109 Principle of Inclusion-Exclusion, 6 Measurable, 2, 24

Measurable  $L-\mathcal{B}$ , 24

Measurable Functions, 25

Probability Density Function, 32

Probability Mass Function, 31

Probability Generating Function, 66

Probability Measure, 4 Probability Space, 4 Probability, Conditional, 17 Probability, Induced, 27 Probability, Total, 18

#### Quantifier, 35

Random Variable, 24
Random Variable, Complex-Valued, 57
Random Vector, 24
Random Vector, n-Dimensional, 71
RV, 24
RV, Bivariate, 72
RV, Continuous, 73
RV, Discrete, 72
RV, Equivalent, 81
RV, Multivariate, 72

Sample Space, 1 Sampling Distribution, 122 Sets, Continuity of, 9 Singular, 33 Slutsky's Theorem, 128 Standard Normal Distribution, 112 Statistic, 122 Support, 37

Total Probability, 18 Transformation, 36 Transformation, Multivariate, 85

Univariate Normal Distribution, 112

Variance, 43 Variance-Covariance Matrix, 114

Weak Law Of Large Numbers, 135 Weak Law Of Large Numbers, Khintchine's, 138