

BIOS 660/BIOS 672 (3 Credits): Probability and Statistical Inference I

Jianwen Cai

<https://sakai.unc.edu/portal/site/bios660-bios672-3-credits>

Notes 1

Elementary Set Theory: Basic Notation	2
Sets	3
Common Notation	4
The Sample Space:	5
Events:	6
Elementary Set Theory: Set Operations	7
Union	8
Intersection	9
Complementation	10
Difference	11
Symmetric Difference	12
Disjoint Union	13
Properties of the Union/Intersection	14
Countable Unions/Intersections	15
Uncountable Unions/Intersections	16
DeMorgan's Rule:	17
Proof of DeMorgan's Rule (2 sets):	18
Proof of DeMorgan's Rule (continued):	19
Limit Sets	20
Limsup and Liminf of Sets	21
Limit of Sets	22
More Examples	23
Monotonicity	24

Sets

- Sets:
 - Sets are basic concepts in mathematics and probability.
 - Crudely defined as: “a collection of some elements”.
 - Usually denoted with a capital letter (e.g. A , B , S)
- Special sets:
 - \mathbb{N} = Natural numbers (1, 2, 3, 4, ...)
 - \mathbb{Z} = Integers (0, +1, -1, +2, -2, ...)
 - \mathbb{Q} = Rational numbers
 - \mathbb{R} = Real numbers

BIOS 660/BIOS 672 (3 Credits)

Notes 1 – 3 / 24

Common Notation

- \emptyset : empty set
- $w \in A$: w is an element of the set A
- $\{w\} \subset A$ or $\{w\} \subseteq A$: the set consisting of the singleton $w \in A$ is a subset of A
- (a, b) , where $a, b \in \mathbb{R}$, is the set of real numbers between (but not including) a and b
- $[a, b]$, where $a, b \in \mathbb{R}$, is the set of real numbers between and including a and b
- $[a, b)$, where $a, b \in \mathbb{R}$, is the set of real numbers between and including a but not b
- $(a, b]$, where $a, b \in \mathbb{R}$, is the set of real numbers between and including b but not a
- $\{w : \text{a statement}\}$: the set of elements w for which the statement holds. Example: the open interval (a, b) can be defined as $\{w : a < w < b\}$.
- $A = B$ if A and B contain exactly the same elements (this can be shown by showing (1) $A \subset B$ and (2) $B \subset A$)

BIOS 660/BIOS 672 (3 Credits)

Notes 1 – 4 / 24

The Sample Space:

- Need to define an “abstract space”, often denoted by Ω , as a non-empty set of all the elements concerned. These elements are called “points” and often denoted with lower case letter
- In probability we often use Ω to denote the **sample space**, which is the collection of all possible distinct realizations of a non-deterministic experiment. (Usually idealized)
- Choice of sample space is the first stem in formulating a probabilistic model for an experiment.

Examples of sample spaces:

1. Draw a card from a deck of 52 cards: $\Omega = \{1, 2, \dots, 52\}$, which is a **finite** sample space
2. Toss a coin until one gets two successive heads and record the number of tosses performed: $\Omega = \{2, 3, 4, \dots, \infty\}$, which is a **countably infinite** sample space.
3. Two components in an electrical system record their failure times: $\Omega = \{(x, y) : x \geq 0, y \geq 0\}$, which is an **uncountably infinite** sample space.

Events:

- An event is always a subset of Ω , “events”
- ... but not all subsets are necessarily events!
- If Ω is countable (i.e. finite or countably infinite), then any subset of Ω is an event
- If Ω is uncountable, we cannot handle all possible subsets (not all are events). Instead, we restrict events to be a “well-behaved” class of subsets. More on this later.
- Individual points in Ω are called “simple events”

Union

Definition: For two sets A and B , their **union** is denoted as $A \cup B = \{w : w \in A \text{ or } w \in B\}$

BIOS 660/BIOS 672 (3 Credits)

Notes 1 – 8 / 24

Intersection

Definition: For two sets A and B , their **intersection** is denoted as $A \cap B = \{w : w \in A \text{ and } w \in B\}$

BIOS 660/BIOS 672 (3 Credits)

Notes 1 – 9 / 24

Complementation

Definition: For a set A in Ω , the complement of A (w.r.t. Ω) is denoted by $A^c = \{w \in \Omega : w \notin A\}$

Note: One can show that under complementation, \subset and \supset are swapped.

BIOS 660/BIOS 672 (3 Credits)

Notes 1 – 10 / 24

Difference

Definition: For two sets A and B , their **difference** is denoted as

$$A - B = \{w : w \in A, w \notin B\} = A \cap B^c$$

BIOS 660/BIOS 672 (3 Credits)

Notes 1 – 11 / 24

Symmetric Difference

Definition: For two sets A and B , their **symmetric difference** is denoted as $A \Delta B = (A - B) \cup (B - A) = \{w : w \in \text{exactly one of } A \text{ and } B\}$

BIOS 660/BIOS 672 (3 Credits)

Notes 1 – 12 / 24

Disjoint Union

Definition: Two sets A and B are called **disjoint** if $A \cap B = \phi$.

Definition: For two disjoint sets A and B , their **disjoint union** is denoted as $A \cup B = A + B$

BIOS 660/BIOS 672 (3 Credits)

Notes 1 – 13 / 24

Properties of the Union/Intersection

For sets A_1, A_2, A_3 :

- Properties of Union:

1. Associative: $(A_1 \cup A_2) \cup A_3 = A_1 \cup (A_2 \cup A_3)$
2. Distributive: $(A_1 \cup A_2) \cap A_3 = (A_1 \cap A_3) \cup (A_2 \cap A_3)$
3. Commutative: $A_1 \cup A_2 = A_2 \cup A_1$

- Properties of Intersection:

1. Associative: $(A_1 \cap A_2) \cap A_3 = A_1 \cap (A_2 \cap A_3)$
2. Distributive: $(A_1 \cap A_2) \cup A_3 = (A_1 \cup A_3) \cap (A_2 \cup A_3)$
3. Commutative: $A_1 \cap A_2 = A_2 \cap A_1$

BIOS 660/BIOS 672 (3 Credits)

Notes 1 – 14 / 24

Countable Unions/Intersections

- Finite sequence of sets: Let $\{A_1, \dots, A_k\}$ be a finite sequence of k sets in Ω . We define

$$\sup_{1 \leq n \leq k} A_n = \bigcup_{n=1}^k A_n = \{w : w \in A_n \text{ for some } 1 \leq n \leq k\}$$

$$\inf_{1 \leq n \leq k} A_n = \bigcap_{n=1}^k A_n = \{w : w \in A_n \text{ for any } 1 \leq n \leq k\}$$

- Countable sequence of sets: Let $\{A_n\}$ be an infinite sequence of sets in Ω . We define

$$\sup_{n \geq 1} A_n = \bigcup_{n=1}^{\infty} A_n = \{w : w \in A_n \text{ for some } n\}$$

$$\inf_{n \geq 1} A_n = \bigcap_{n=1}^{\infty} A_n = \{w : w \in A_n \text{ for any } n \geq 1\}$$

BIOS 660/BIOS 672 (3 Credits)

Notes 1 – 15 / 24

Uncountable Unions/Intersections

We extend the intersection over a set of integers to any arbitrary set:

Definition: For $\{A_t, t \in T\}$, where T is a (possibly uncountable) index set,
 $\bigcup_{t \in T} A_t = \{w : w \in A_t \text{ for some } t \in T\}$. The definition for intersection is similar.

BIOS 660/BIOS 672 (3 Credits)

Notes 1 – 16 / 24

DeMorgan's Rule:

DeMorgan's Rule (for 2 sets):

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

and

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

DeMorgan's Rule (general):

$$\left(\bigcup_{t \in T} A_t \right)^c = \bigcap_{t \in T} A_t^c$$

and

$$\left(\bigcap_{t \in T} A_t \right)^c = \bigcup_{t \in T} A_t^c$$

where T is any index set (finite, countably infinite, uncountably infinite).

BIOS 660/BIOS 672 (3 Credits)

Notes 1 – 17 / 24

Proof of DeMorgan's Rule (2 sets):

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

Proof:

BIOS 660/BIOS 672 (3 Credits)

Notes 1 – 18 / 24

Proof of DeMorgan's Rule (continued):

The proof that $(A \cup B)^c = A^c \cap B^c$ is similar.

BIOS 660/BIOS 672 (3 Credits)

Notes 1 – 19 / 24

Limsup and Liminf of Sets

- **Definition:** Let $\{A_n\}$ be a sequence of sets in Ω . Define

$$A^* = \limsup_n A_n = \bigcap_{n=1}^{\infty} \bigcup_{k=n}^{\infty} A_k$$

and

$$A_* = \liminf_n A_n = \bigcup_{n=1}^{\infty} \bigcap_{k=n}^{\infty} A_k$$

these are the upper and lower limits of the sequence.

- **Theorem:** (a) $w \in \limsup_n A_n$ if and only if w is in infinitely many of the A_n .
(b) $w \in \liminf_n A_n$ if and only if there is an m such that $w \in A_n$ for all $n \geq m$ (i.e. w is in all but the first m) [Proof omitted]
- **Theorem:** $\liminf_n A_n \subset \limsup_n A_n$ (proof for homework)

BIOS 660/BIOS 672 (3 Credits)

Notes 1 – 21 / 24

Limit of Sets

If $\limsup A_n = \liminf A_n$ then we say that $\{A_n\}$ is convergent and write
 $\lim_{n \rightarrow \infty} A_n = \limsup A_n = \liminf A_n$

Example: let $\Omega = [0, 1]$, $A_n = [0, 1/n]$ if n is even and $A_n = [1 - 1/n, 1]$ if n is odd.

- Then by definition $A_* = \emptyset$ and $A^* = \{0\}$.
- Since $A_* \neq A^*$, then $\lim A_n$ does not exist.
- On the other hand, if we let $A_n = [0, 1/n]$ for all n , then $A_* = A^* = \{0\}$ and the limit exists.

BIOS 660/BIOS 672 (3 Credits)

Notes 1 – 22 / 24

More Examples

What is A_* and A^* for:

- Example: let $\Omega = \{0, 1\}$, then consider the sequence $\{0\}, \{1\}, \{0\}, \{1\}, \{0\}, \dots$
- Example: let $\Omega = \mathbb{Z}$, Consider the sequence of sets: $\{0\}, \{0, 1\}, \{0\}, \{0, 1\}, \{0\}, \{0, 1\}, \{0\}, \dots$
- Example: let $\Omega = \mathbb{Z}$, Consider the sequence of sets: $\{0\}, \{0, 1\}, \{0, 1, 2\}, \{0, 1, 2, 3\}, \{0, 1, 2, 3, 4\}, \dots$
- Example: let $\Omega = [0, 1]$, $A_n = [0, 1/n]$ if n is even and $A_n = [1 - 1/n, 1)$ if n is odd.

BIOS 660/BIOS 672 (3 Credits)

Notes 1 – 23 / 24

Monotonicity

- **Monotonicity:** A monotone sequence of sets is defined as:
 - $\{A_n\}$ is called monotone increasing iff $A_n \subset A_{n+1}$ for any n
 - $\{A_n\}$ is called monotone decreasing iff $A_n \supset A_{n+1}$ for any n
- **Theorem:** A monotone sequence of sets is convergent
Proof for increasing sequence:

BIOS 660/BIOS 672 (3 Credits)

Notes 1 – 24 / 24

BIOS 660/BIOS 672 (3 Credits): Probability and Statistical Inference I

Jianwen Cai

<https://sakai.unc.edu/portal/site/bios660-bios672-3-credits>

Notes 2

Introduction to Fields and σ-fields	2
Intro: More on Events	3
Intro: On Measures and Probability	4
Classes of Sets	5
Fields	6
Some Properties	7
σ -Fields	8
Example	9
σ -field Generated by a Class of Sets	10
Borel σ -fields (\mathcal{B})	11
Events... Again!	12

Intro: More on Events

- **Casella & Berger Definition:** An event is any collection of possible outcomes of an experiment, that is, any subset of Ω (including Ω itself). \leftarrow not strictly true!

Previously:

- An event is necessarily a collection of possible outcomes of a random experiment (a set of elements in Ω , i.e. a subset of Ω)
- For discrete (finite and countable) sample space, the set of possible events is the power set (2^Ω), i.e. any subset of Ω .
- For continuous sample spaces (uncountably infinite spaces), the set of possible events is NOT the power set, there are subsets that are not considered events.
- Strictly speaking, **events** are the subsets of the sample space for which a probability is defined.

BIOS 660/BIOS 672 (3 Credits)

Notes 2 – 3 / 12

Intro: On Measures and Probability

A **measure** is a set function that assigns a number $\mu(A)$ to each set A in a certain class of sets.

Examples:

- Length
- Area
- Volume
- Probability

Some structure must be imposed on the class of sets in which the set function μ is defined (i.e. cannot take probability of any set in 2^Ω).

BIOS 660/BIOS 672 (3 Credits)

Notes 2 – 4 / 12

Classes of Sets

- **Definition:** A **class** is a collection of sets (set of sets) that satisfy some conditions. Usually denoted with script characters (\mathcal{S} , \mathcal{X} , \mathcal{A} , etc.)
- **Definition:** A class of sets \mathcal{X} is closed under an operation (e.g. union, intersection, etc.) if when performed on any members of \mathcal{X} yields a set which also belongs to the class.
- Example: $\mathcal{X} = \{\emptyset, A, B, C, A \cup B, A \cup C, B \cup C, A \cup B \cup C\}$ is closed under Union operation.

BIOS 660/BIOS 672 (3 Credits)

Notes 2 – 5 / 12

Fields

Definition: A class \mathcal{X} of sets in Ω is called a **field** if

1. \mathcal{X} is non-empty
2. \mathcal{X} is closed under finite union
3. \mathcal{X} is closed under complementation.

Examples:

- $\mathcal{X} = \{\emptyset, \Omega\}$ (trivial class)
- $\mathcal{X} =$ all subsets of Ω
- Let $\Omega = (-\infty, \infty)$, $\mathcal{X}_3 =$ class of all finite intervals (a, b) where $a, b \in \mathbb{R}$.

BIOS 660/BIOS 672 (3 Credits)

Notes 2 – 6 / 12

Some Properties

Properties of a field \mathcal{X} :

1. \mathcal{X} is also closed under finite intersection
2. $\emptyset \in \mathcal{X}$ and $\Omega \in \mathcal{X}$

BIOS 660/BIOS 672 (3 Credits)

Notes 2 – 7 / 12

σ -Fields

- **Definition:** A class \mathcal{X} of sets is a σ -**field** if
 1. \mathcal{X} is non-empty
 2. \mathcal{X} is closed under countable unions
 3. \mathcal{X} is closed under complementation
- Examples: (1) $\mathcal{X} = \{\emptyset, \Omega\}$ (2) \mathcal{X} = all subsets of Ω are σ -fields
- Obviously, a σ -field is necessarily a field, but the converse does not hold: Consider class of all finite sets and sets whose complement is finite.
- **Theorem:** \mathcal{X}_1 and \mathcal{X}_2 are σ -fields, then $\mathcal{X}_1 \cap \mathcal{X}_2 = \{A : A \in \mathcal{X}_1 \text{ and } \mathcal{X}_2\}$ is also a σ -field.

BIOS 660/BIOS 672 (3 Credits)

Notes 2 – 8 / 12

Example

Let $\Omega = \mathbb{R}$ and $\mathcal{A} = \{A : A \text{ or } A^c \text{ is finite}\}$. Show that \mathcal{A} is a field, but that \mathcal{A} is not a σ -field. Note that \mathcal{A} is the finite-cofinite field.

- Let $A_1 = \{1\}, A_2 = \{2\}, A_3 = \{3\}, \dots$,
- Then $\bigcup_{n=1}^{\infty} A_n = \mathbb{N} \notin \Omega$.

BIOS 660/BIOS 672 (3 Credits)

Notes 2 – 9 / 12

σ -field Generated by a Class of Sets

- **Theorem:** Given a class of sets, \mathcal{S} , not necessarily a σ -field, there is a minimum σ -field (denoted $\sigma(\mathcal{S})$ containing) it, i.e. $\sigma(\mathcal{S})$ is a class of sets that:
 - is a σ -field
 - it contains \mathcal{S} : if $A \in \mathcal{S}$, then it is also in $\sigma(\mathcal{S})$.
 - if \mathcal{X} is another σ -field that contains \mathcal{S} , then \mathcal{X} contains $\sigma(\mathcal{S})$.
- $\sigma(\mathcal{S})$ is also called the **σ -field generated by \mathcal{S}** .
- $\sigma(\mathcal{S})$ also defined as the intersection of all σ -fields that contain \mathcal{S}
- Properties:
 - $\sigma(\mathcal{S})$ is itself a σ -field
 - $\mathcal{S} \subset \sigma(\mathcal{S})$
 - $\mathcal{S}_1 \subset \mathcal{S}_2$ implies $\sigma(\mathcal{S}_1) \subset \sigma(\mathcal{S}_2)$
 - if \mathcal{S} is itself a σ -field, then $\sigma(\mathcal{S}) = \mathcal{S}$

BIOS 660/BIOS 672 (3 Credits)

Notes 2 – 10 / 12

Borel σ -fields (\mathcal{B})

- Let $\Omega = \mathbb{R} = (-\infty, \infty)$. We have 4 types of finite intervals:
 1. $\mathcal{S}_1 = \{(a, b) : a < b, a, b \in \mathbb{R}\}$,
 2. $\mathcal{S}_2 = \{(a, b) : a < b, a, b \in \mathbb{R}\}$,
 3. $\mathcal{S}_3 = \{(a, b] : a < b, a, b \in \mathbb{R}\}$,
 4. $\mathcal{S}_4 = \{[a, b] : a < b, a, b \in \mathbb{R}\}$
- Let $\mathcal{S} = \bigcup_{i=1}^4 \mathcal{S}_i$ = a class of all finite intervals. Note that \mathcal{S} is neither a field nor a σ -field.
- We extend \mathcal{S} to a σ -field through the following definition:

Definition: The σ -field generated by \mathcal{S} is called the **Borel σ -field on \mathbb{R}** and is denoted by $\mathcal{B} = \sigma(\mathcal{S})$. Any set in \mathcal{B} is called a **Borel Set**.

- Important: Do not try to characterize \mathcal{B} .
- We will return to this σ -field when we start talking about Random Variables.

Events... Again!

- Why do we care about σ -fields?
They are key to understanding the definition of an event.
- **Definition:** A **measurable space** is a set Ω endowed with a σ -field \mathcal{F} of subsets of Ω denoted by the pair (Ω, \mathcal{F}) .
- Sets in \mathcal{F} are defined as **events**!
- In this class, we will often deal with the triplet (Ω, \mathcal{F}, P) where P is a probability measure on the space (Ω, \mathcal{F}) .
- For discrete (finite and countably infinite) Ω , $\mathcal{F} = 2^\Omega$, and for uncountable Ω , $\mathcal{F} = \mathcal{B}$.

BIOS 660/BIOS 672 (3 Credits): Probability and Statistical Inference I

Jianwen Cai

<https://sakai.unc.edu/portal/site/bios660-bios672-3-credits>

Notes 2.1

Examples for Sample Space and σ -fields	2
Examples (cont.)	3

Examples for Sample Space and σ -fields

Definition: A σ -field on a set Ω is a collection of subsets of Ω that includes the empty subset, is closed under complement, and is closed under countable unions and countable intersections.

Examples for Sample Space and σ -fields:

1. **Experiment:** Toss a coin with a "Head" and a "Tail".
Sample space Ω : {Head, Tail}
 σ -field on Ω : $\{\emptyset, \{\text{Head}\}, \{\text{Tail}\}, \{\text{Head}, \text{Tail}\}\}$.
2. **Experiment:** In a health survey among cancer patients, we ask each patient to report their perceived quality of life, categorized as {Very Good, Good, Poor}.
Sample space Ω : {Very Good, Good, Poor}
 σ -field on Ω : $\{\emptyset, \{\text{Very Good}\}, \{\text{Good}\}, \{\text{Poor}\}, \{\text{Very Good, Good}\}, \{\text{Very Good, Poor}\}, \{\text{Good, Poor}\}, \{\text{Very Good, Good, Poor}\}\}$.

Examples (cont.)

3. **Experiment:** Toss 2 coins and record the results from both tosses.
Sample space Ω : {HH, HT, TH, TT}
 σ -field on Ω : $\{\emptyset, \{\text{HH}\}, \{\text{HT}\}, \{\text{TH}\}, \{\text{HT}\}, \{\text{HH, HT}\}, \{\text{HH, TH}\}, \{\text{HH, TT}\}, \{\text{HT, TH}\}, \{\text{HT, TT}\}, \{\text{TH, TT}\}, \{\text{HH, HT, TH}\}, \{\text{HH, HT, TT}\}, \{\text{HH, TH, TT}\}, \{\text{HT, TH, TT}\}, \{\text{HH, HT, TH, TT}\}\}$.
4. **Experiment:** Toss 2 coins and record the number of heads.
Sample Space Ω : {0, 1, 2}
 σ -field on Ω : $\{\emptyset, \{0\}, \{1\}, \{2\}, \{0,1\}, \{0,2\}, \{1,2\}, \{0,1,2\}\}$.

BIOS 660/BIOS 672 (3 Credits): Probability and Statistical Inference I

Jianwen Cai

<https://sakai.unc.edu/portal/site/bios660-bios672-3-credits>

Notes 3

Axiomatic Probability	2
Intuitive probability	3
Axiomatic probability	4
Kolmogorov's Axioms of Probability	5
Axioms continued	6
Probability calculus	7
Monotone Sequence of Events	8
Boole's inequality (Union bound)	9
Counting: Preliminaries	10
Discrete Probability	11
Simplest Setting	12
Simple Examples	13
More Preliminaries	14
Counting: Ordered Samples	15
Ordered Samples	16
Sampling with Replacement	17
Sampling without Replacement	18
Examples	19
Examples (1)	20
Examples (2)	21
Examples (3)	22
Examples (4): Birthday Problem	23
Additional Reading	24
Additional Reading	25

Intuitive probability

- **Random experiment:** Can be repeated indefinitely but future outcomes cannot be exactly predicted.
- **Stabilization of relative frequency:** Suppose we perform “independent” repetitions of a random experiment and count how many times an “event” E occurs. Let $f_n(E)$ be the number of occurrences in n repetitions. Then the *relative frequency* of E

$$r_n(E) = \frac{f_n(E)}{n}$$

“converges” to some number as $n \rightarrow \infty$. We call this number $P(E)$, the probability of E .

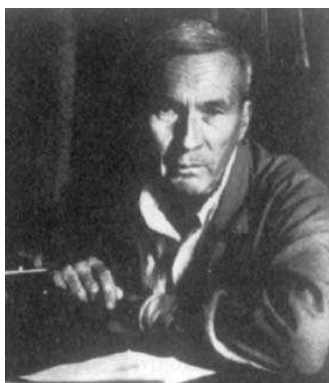
- Examples: rolling of a die, spin the wheel.
- But what is “event”, “independent”, “convergence”?

BIOS 660/BIOS 672 (3 Credits)

Notes 3 – 3 / 25

Axiomatic probability

Andrey Kolmogorov (1903–1987)



Grundbegriffe der Wahrscheinlichkeitsrechnung (1933)
 Basic Concepts of Probability Calculation
 Foundations of the Theory of Probability

BIOS 660/BIOS 672 (3 Credits)

Notes 3 – 4 / 25

Kolmogorov's Axioms of Probability

Suppose $E \subset \Omega$ and we want to assign a probability to E .

$P(E)$ is a real-valued function of the event E , called the probability of E , that satisfies the following:

Axioms

- i. **Regularity:** $P(\Omega) = 1$
- ii. **Non-negativity:** $P(E) \geq 0$
- iii. **Countable Additivity:** If E_1, E_2, \dots are mutually exclusive ($E_i E_j = \emptyset, i \neq j$), then
 $P(\bigcup_{i=1}^{\infty} E_i) = \sum_{i=1}^{\infty} P(E_i)$ (countable additivity)

The last axiom contains a hidden assumption, namely that if E_1, E_2, \dots is an infinite sequence of events, then $\bigcup_{i=1}^{\infty} E_i$ is also an event. This is guaranteed if the events E_1, E_2, \dots belong to a σ -field.

Axioms continued

Definition: A probability measure P defined on a σ -field of subsets of Ω is a real-valued set function satisfying Kolmogorov's Axioms.

Definition - formal: A probability space is denoted by (Ω, \mathcal{A}, P) where Ω is the sample space, \mathcal{A} refers to a σ -field of subsets of Ω , and P is a probability measure.

Probability calculus

If P is a probability function and A is any set in \mathcal{A} :

1. $P(\emptyset) = 0$
2. $P(A) \leq 1$
3. $P(A^c) = 1 - P(A)$

If P is a probability function and A, B are sets in \mathcal{A} :

1. $P(A - B) = P(A) - P(A \cap B)$ (Note: C&B uses \setminus instead of '-')
2. $P(A \cup B) = P(A) + P(B) - P(A \cap B)$
3. $A \subset B \Rightarrow P(A) \leq P(B)$

Law of Total Probability: If P is a probability function and $\{C_1, C_2, \dots\}$ is a partition of Ω :

$$P(A) = \sum_{i=1}^{\infty} P(A \cap C_i)$$

Monotone Sequence of Events

- (i) If $\{E_n\}$ is an increasing sequence of sets and $E_n \in \mathcal{A}$, then $P(\lim_n E_n) = \lim_n P(E_n)$;
(ii) If $\{E_n\}$ is a decreasing sequence of sets, then $P(\lim_n E_n) = \lim_n P(E_n)$.

Proof of (i):

Boole's inequality (Union bound)

Suppose E_1, E_2, \dots, E_n are events.

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

BIOS 660/BIOS 672 (3 Credits)

Notes 3 – 9 / 25

Counting: Preliminaries

10 / 25

Discrete Probability

- **Discrete Sample Space:** A sample space is called **discrete** if it contains only finitely many points or infinitely many points which can be arranged into a simple sequence e_1, e_2, \dots , i.e. a countable number of elements.
- $\mathcal{F} = 2^\Omega$.
- $P(\{e_1\}) + P(\{e_2\}) + \dots = P(\Omega) = 1$, and more generally, for a set A that contains the points a_1, a_2, \dots, a_k , then $P(A) = P(\{a_1\}) + P(\{a_2\}) + \dots + P(\{a_k\})$.

BIOS 660/BIOS 672 (3 Credits)

Notes 3 – 11 / 25

Simplest Setting

Simplest setting:

1. Finite sample space with n elements, $\Omega = \{e_1, e_2, \dots, e_n\}$
2. All sample points are equally probable.

In this case, $P(\{e_j\}) = \frac{1}{n}$, and for an event $\{a_1, \dots, a_k\} = A \subset \Omega$,
 $P(A) = kP(\{a_j\}) = \mu(A)/\mu(\Omega) = k/n$ where μ is the counting measure (counts the number of elements in a set).

$P(A)$ is simply the sum of the probabilities of each of the elements in A .

Counting and combinatorics will be concerned with identifying the number of individual elements within A .

Simple Examples

Examples: Suppose we have a standard deck of cards and we draw a single card. (probability of drawing any card is the same so they are equiprobable)

- Probability of drawing Jack of Hearts is $1/52$
- Probability of drawing a Club is $13/52 = 1/4$
- Probability of drawing an Ace is $4/52 = 1/13$
- Probability of drawing an Ace or Heart is $16/52 = 4/13$

More Preliminaries

- **Pairs:** With m elements a_1, \dots, a_m , and n elements b_1, \dots, b_n , it is possible to form nm pairs (a_j, b_k) containing one element from each group.
- **Multiplets:** Given r experiments that are to be performed such that the r^{th} experiment may result in any of n_r possible outcomes, then there is a total of $n_1 n_2 \dots n_r$ possible outcomes of the r experiments.
- Examples:
 1. College committee consists of 3 freshman, 4 sophomores, 5 juniors, and 2 seniors. A subcommittee consisting of 1 person from each class is to be chosen. How many subcommittees are possible?
 2. How many 7-place license plates are possible if the first 3 places are to be occupied by letters and the final 4 by numbers?

BIOS 660/BIOS 672 (3 Credits)

Notes 3 – 14 / 25

Counting: Ordered Samples

15 / 25

Ordered Samples

- Many probability problems use these same ideas of counting in settings where repeated selections are taken from the same set of objects. The number of possible selection patterns (and hence the probability of a particular pattern) depends on
 1. whether or not items are replaced after selection
 2. whether or not the order of selection matters
- Suppose we have a “population” of n elements a_1, \dots, a_n . Any ordered arrangement a_{j_1}, \dots, a_{j_r} of r symbols is called an *ordered sample of size r* .
 1. Sampling with replacement: each selection is made from the entire population (each of the r elements can be chosen in n ways. there are n^r possible samples)
 2. Sampling without replacement: an element, once chosen, is removed from the population. (sample cannot exceed size n)

BIOS 660/BIOS 672 (3 Credits)

Notes 3 – 16 / 25

Sampling with Replacement

<u>Draw</u>	<u>Set</u>
1	A_1, A_2, \dots, A_n
2	A_1, A_2, \dots, A_n
\vdots	
r	A_1, A_2, \dots, A_n

Theorem: For a population of n elements and a prescribed sample size r , there exist n^r different samples with replacement

Example: The Braille alphabet has 6 locations that can be raised or not raised. How many combinations of raised/not raised locations are there?

Why does this differ from the actual number of Braille letters?

Sampling without Replacement

<u>Draw</u>	<u>Set</u>	
1	A_1, A_2, \dots, A_n	
2	B_1, B_2, \dots, B_{n-1}	B'_i s are A'_i s not drawn
3	C_1, C_2, \dots, C_{n-2}	C'_i s are B'_i s not drawn
\vdots		
r	$X_1, X_2, \dots, X_{n-(r-1)}$	X'_i s are remaining items after first $r - 1$ draws

- **Theorem:** For a population of n elements and a prescribed sample size r , there exist $n(n-1)(n-2)\dots(n-r+1) = \frac{n!}{(n-r)!}$ different samples without replacement.
- An important special case occurs when $r = n$ (define $0!=1$):

$$\frac{n!}{(n-r)!} = \frac{n!}{(n-n)!} = n(n-1)(n-2)\dots 2 \cdot 1 = n!$$

This is the number of **permutations** of a set of size n

Examples

1. You want to take a photo of six of your friends. How many ways can you arrange them in a row?
2. A competition has 10 competitors of which 4 are Russian, 3 are from the USA, 1 is from Brazil, and 2 are from China. If the tournament result lists just the nationalities of the players in which they placed, how many outcomes are possible?
3. You want to invite those same 6 friends to dinner. How many ways can you arrange them at your round dining table?
4. 4 Martians, 3 Plutonians and 5 Jupiterians are to be seated in a row at an interplanetary conference. How many seating arrangements are possible?
5. With indistinguishable elements: How many different letter arrangements can be formed using the letters B-A-N-A-N-A?

Examples (1)

Suppose one has a population of size n . What is probability that a particular element will not be in sample of size r ?

Solution Without Replacement:

$$\frac{(n-1)!}{(n-1-r)!} = \text{No. of samples without particular element.}$$

$$\frac{n!}{(n-r)!} = \text{No. of Samples}$$

$$q = \text{Probability element is not in sample}$$

$$= \frac{(n-1)!/(n-(r+1))!}{n!/(n-r)!}$$

$$= \frac{(n-1)(n-2)\dots(n-r)}{n(n-1)(n-2)\dots(n-r+1)}$$

$$= \frac{n-r}{n} = 1 - \frac{r}{n}$$

$$p = \text{Probability element in the sample} = 1 - q = r/n$$

Examples (2)

Suppose one has a population of size n . What is probability that a particular element will not be in sample of size r ?

Solution With Replacement:

$$q = \frac{(n-1)^r}{n^r} = \left(1 - \frac{1}{n}\right)^r$$
$$p = 1 - q = 1 - \left(1 - \frac{1}{n}\right)^r$$

BIOS 660/BIOS 672 (3 Credits)

Notes 3 – 21 / 25

Examples (3)

Suppose one samples r items with replacement from a total of n . What is probability of no repetitions in the sample?

$$q = \frac{\text{No. of Samples with No Repetition}}{\text{No. of Samples}} = \frac{n!/(n-r)!}{n^r}$$
$$= \frac{n \cdot (n-1) \cdot (n-2) \dots (n-r+1)}{n \cdot n \cdot n \dots n}$$
$$= 1\left(1 - \frac{1}{n}\right)\left(1 - \frac{2}{n}\right) \dots \left(1 - \frac{(r-1)}{n}\right)$$
$$\sim \left(1 - \frac{1+2+\dots+r-1}{n}\right)$$
$$= 1 - \frac{r(r-1)}{2n}$$

because $1 + 2 + \dots + r - 1 = \frac{r(r-1)}{2}$

Note: $p = \text{Probability of repetitions} = 1 - q \sim \frac{(r-1)r}{2n}$

BIOS 660/BIOS 672 (3 Credits)

Notes 3 – 22 / 25

Examples (4): Birthday Problem

Application - the Birthday Problem. If there is a class of r people, what is the probability that some have birthdays on the same day? ($n = 365$)

$$p \sim \frac{r(r-1)}{2 \cdot 365}$$

Suppose $r = 25$

$$p \sim \frac{25 \cdot 24}{2 \cdot 365} = 0.82$$

Suppose $p = 1/2$. What is value of r ?

$$r(r-1)/730 = 1/2 \quad \text{Solution is } r \approx \frac{39}{2} \sim 20$$

BIOS 660/BIOS 672 (3 Credits)

Notes 3 – 23 / 25

Additional Reading

24 / 25

Additional Reading

See Chapter 1.2 in Casella and Berger.

BIOS 660/BIOS 672 (3 Credits)

Notes 3 – 25 / 25

BIOS 660/BIOS 672 (3 Credits): Probability and Statistical Inference I

Jianwen Cai

<https://sakai.unc.edu/portal/site/bios660-bios672-3-credits>

Notes 4

Counting: Unordered Samples	2
Combinations (Subpopulations)	3
Combinations (2)	4
Conventions	5
Examples	6
Examples(2)	7
Binomial Theorem	8
Pascal's Triangle (1)	9
Pascal's Triangle (2)	10
Partitions.	11
Occupancy Problem	12
Occupancy: Results	13
Examples	14
Stirling's Approximation	15

Combinations (Subpopulations)

How many ways can we select r objects from n paying no attention to order? Continue to suppose that the n items are distinct.

$n = 4$ A, B, C, D

$r = 2$ Possible ordered samples:
 $(A, B), (A, C), (A, D), (B, C), (B, D), (C, D)$
 $(B, A), (C, A), (D, A), (C, B), (D, B), (D, C)$

$$n!/(n-r)! = 4!/2! = 4 \cdot 3 = 12$$

6 = No. of ways of choosing 2 objects from among 4 objects with no ordering... Why?

There are $4 \cdot 3 = 12$ ways of choosing ordered pairs. Then there are $2! = 2 \cdot 1$ orderings that we don't really care about, so $12/2 = 6$.

BIOS 660/BIOS 672 (3 Credits)

Notes 4 – 3 / 15

Combinations (2)

No. of ordered samples = No. of Combinations \times No. of ways of permuting each ordered sample.

$$\begin{aligned} \frac{n!}{(n-r)!} &= c \cdot r! \\ c &= \frac{n!}{(n-r)!r!} = \binom{n}{r} = \binom{n}{n-r} \\ \binom{n}{r} &= \text{No. of ways of choosing } r \text{ objects from among } n \\ &= \text{we say “} n \text{ choose } r \text{”} \end{aligned}$$

Theorem: A population of n elements possesses $\binom{n}{r}$ different combinations (subpopulations) of size $r \leq n$, i.e. $\binom{n}{r}$ represents the number of possible combinations of n objects taken r at a time.

BIOS 660/BIOS 672 (3 Credits)

Notes 4 – 4 / 15

Conventions

Conventions:

$$0! = 1; \binom{n}{0} = 1; \binom{n}{r} = 0 \quad \text{if } r > n, r < 0$$

BIOS 660/BIOS 672 (3 Credits)

Notes 4 – 5 / 15

Examples

Applications of $\binom{n}{r} = \frac{n!}{r!(n-r)!}$

- From a group of 5 women and 7 men, how many different committees consisting of 2 women and 3 men can be formed?
- How many different poker hands with 5 cards are possible?
- What is the *probability* of 4 aces in a poker hand with 5 cards?

BIOS 660/BIOS 672 (3 Credits)

Notes 4 – 6 / 15

Examples(2)

- What is the *probability* of 4 of a kind (e.x. 4 Kings, 4 Queens, etc)?

- What is the *probability* of (exactly) 3 of a kind with no other pairs?

Size of Event:

Binomial Theorem

Suppose n is a positive integer. The Binomial Theorem says

$$(a + b)^n = \sum_{r=0}^n \binom{n}{r} a^r b^{n-r}$$

Proof: Homework, by induction.

Note that if $a + b = 1$ one has

$$\sum_{r=0}^n \binom{n}{r} a^r (1-a)^{n-r} = 1,$$

and that if $a = b = 1$,

$$\sum_{r=0}^n \binom{n}{r} = 2^n.$$

Note: the left hand side is the number of all possible subsets from $\{X_1, \dots, X_n\}$. The above equality tells us that the number of all possible subsets from $\{X_1, \dots, X_n\}$ is 2^n , which motivates the notation 2^Ω for indicating power set .

Pascal's Triangle (1)

Pascal's triangle is an illustration of the following result:

$$\binom{n+1}{r} = \binom{n}{r-1} + \binom{n}{r}$$

Proof: Homework.

BIOS 660/BIOS 672 (3 Credits)

Notes 4 – 9 / 15

Pascal's Triangle (2)

Pascal's triangle

$$\begin{array}{ccccccc} & & & & \binom{0}{0} = 1 & & \\ & & & \binom{1}{0} = 1 & & \binom{1}{1} = 1 & \\ & & \binom{2}{0} = 1 & & \binom{2}{1} = 2 & & \binom{2}{2} = 1 \\ \binom{3}{0} = 1 & & \binom{3}{1} = 3 & & \binom{3}{2} = 3 & & \binom{3}{3} = 1 \\ & & & 1 & & & \\ & & 1 & & 1 & & \\ & & & 1 & 2 & 1 & \\ & 1 & & 3 & & 3 & 1 \\ & & 1 & 4 & 6 & 4 & 1 \end{array}$$

BIOS 660/BIOS 672 (3 Credits)

Notes 4 – 10 / 15

Partitions

- Beyond a single subpopulation...
- **Theorem:** Let r_1, r_2, \dots, r_k be nonnegative integers such that

$$r_1 + r_2 + \dots + r_k = n$$

The number of ways in which a population of n elements can be divided into k ordered groups (partitioned into k -subpopulations) of which the first contains r_1 elements, the second r_2 elements, etc., is

$$\frac{n!}{r_1! r_2! \dots r_k!}$$

- Example: Ten children are being divided into an A team and a B team of 5 each. How many different divisions are possible?

Occupancy Problem

Suppose we have n indistinguishable balls which we wish to place in r distinguishable urns. How many different outcomes are possible?

- We can describe the outcome by a vector (x_1, \dots, x_r) where $x_j \geq 0$ is the number of balls in the j th urn and $\sum_{j=1}^r x_j = n$. We are therefore looking for the number of possible vectors.
- We can further think of this as taking n indistinguishable items and dividing them into r groups.
- First consider the situation that each urn should have at least 1 ball. In this situation, we are dealing with $x_j > 0$ ($j = 1, \dots, r$). There are $n - 1$ spaces between the objects, and we select $r - 1$ of them to choose our dividing points. Therefore, there are $\binom{n-1}{r-1}$ different partitions.
- Now suppose that some of the urns are allowed to be empty. In this case we can imagine that we have n objects plus $r - 1$ dividers. The number of possible partitions is therefore the number of (non-redundant) orderings of these $n + r - 1$ items which is equal to $(n + r - 1)! / [n!(r - 1)!] = \binom{n+r-1}{r-1}$

Occupancy: Results

- **Theorem:** there are $\binom{n-1}{r-1}$ distinct positive integer valued vectors (x_1, \dots, x_r) satisfying

$$x_1 + \dots + x_r = n, \quad x_j > 0, j = 1, \dots, r$$

- **Theorem:** there are $\binom{n+r-1}{r-1}$ distinct nonnegative integer-valued vectors (x_1, \dots, x_r) satisfying

$$x_1 + \dots + x_r = n, \quad x_j \geq 0, j = 1, \dots, r$$

BIOS 660/BIOS 672 (3 Credits)

Notes 4 – 13 / 15

Examples

1. How many distinct nonnegative integer-valued solutions of $x_1 + x_2 = 3$ are possible?
2. The state health department has 20 cases of flu vaccine to distribute among 4 possible hospitals. If all of the vaccine must be distributed, how many different ways can the 20 cases be distributed if not every hospital has to receive some? What if the department wants to have the flexibility to hold some in reserve and not every hospital has to receive some?

BIOS 660/BIOS 672 (3 Credits)

Notes 4 – 14 / 15

Stirling's Approximation

- **Stirling's Formula:**

$$n! \sim \sqrt{2\pi n} n^{n+\frac{1}{2}} e^{-n}$$

BIOS 660/BIOS 672 (3 Credits): Probability and Statistical Inference I

Jianwen Cai

<https://sakai.unc.edu/portal/site/bios660-bios672-3-credits>

Notes 5

Conditional Probability and Independence	2
Conditional Probability	3
cont.	4
Independence.	5
Independence of many events	6
Independence of many events (cont.)	7
Independence of many events (cont.)	8
Sequential conditioning	9
The Birthday Problem.	10
The Birthday Problem.	11
Turning around probabilities	12
Decomposition Formula (Total Probability).	13
(So called) Bayes' Theorem	14
Bayes and Screening	15
Papanicolaou Example	16
Relative risks and relative odds	17
Bayes and Case Control Studies	18
Additional Reading	19
Additional Reading	20

Conditional Probability

If $P(B) > 0$ we define the *conditional probability* of the event A given B as

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

Intuitively, conditioning on B means reducing the original sample space S to B , which becomes the new, reduced sample space. All probabilities are computed with respect to B . Notice:

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

Disjoint events: If $A \cap B = \emptyset$, then $P(A|B) = 0$ and $P(B|A) = 0$.

BIOS 660/BIOS 672 (3 Credits)

Notes 5 – 3 / 20

cont.

Conditional probability satisfies the axioms of probability:

1. $P(\Omega|B) = 1$
2. $P(A|B) \geq 0$
3. If A_1, A_2, \dots are mutually exclusive events, then $P(\bigcup_{i=1}^{\infty} A_i|B) = \sum_{i=1}^{\infty} P(A_i|B)$

and all the other properties:

1. $P(\emptyset|B) = 0$
2. $P(A|B) \leq 1$
3. $P(A^c|B) = 1 - P(A|B)$

etc.

BIOS 660/BIOS 672 (3 Credits)

Notes 5 – 4 / 20

Independence

Two events A and B are said to be *independent* if

$$P(A \cap B) = P(A)P(B). \quad (1)$$

Why? Because then

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

This could have been taken as the definition, but (1) is easier to generalize.

If A and B are independent, then so are

- A^c and B
- A and B^c
- A^c and B^c

* Can two disjoint events be independent and vice versa? (HW)

Independence of many events

The events A_1, A_2, \dots, A_n are mutually independent if for every subcollection A_{i_1}, \dots, A_{i_k} of size $k = 2, \dots, n$

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

Notice that is a very strong condition. But it is necessary to ensure that

$$P(A_j | A_{i_1}, \dots, A_{i_k}) = P(A_j)$$

for every j and every subcollection A_{i_1}, \dots, A_{i_k} that does not include A_j . Pairwise independence is not enough.

Independence of many events (cont.)

Example: Toss a coin three times. Sample space $\Omega = \{HHH, HHT, HTH, HTT, THH, THT, TTH, TTT\}$. Define the events:

- $H_1 = \{ \text{The outcome of the first toss is heads} \}$
 $= \{HHH, HHT, HTH, HTT\}$
- $H_2 = \{ \text{The outcome of the second toss is heads} \}$
 $= \{HHH, HHT, THH, THT\}$
- $H_3 = \{ \text{The outcome of the third toss is heads} \}$
 $= \{HHH, HTH, THH, TTH\}$

Suppose every outcome is equally likely. Then these events are independent.

$$H_1 \cap H_2 = \{HHH, HHT\}; H_1 \cap H_3 = \{HHH, HTH\}; H_2 \cap H_3 = \{HHH, THH\};$$

$$P(H_1) = P(H_2) = P(H_3) = 4/8 = 1/2; P(H_1 \cap H_2) = 2/8 = 1/4 = P(H_1)P(H_2);$$

$$H_1 \cap H_2 \cap H_3 = \{HHH\}; P(H_1 \cap H_2 \cap H_3) = 1/8 = P(H_1)P(H_2)P(H_3)$$

Independence of many events (cont.)

Now define the events:

- $A_{12} = \{ \text{The outcome of the first toss equals the second} \}$
- $A_{13} = \{ \text{The outcome of the first toss equals the third} \}$
- $A_{23} = \{ \text{The outcome of the second toss equals the third} \}$

These events are pairwise independent but not mutually independent.

$$(A_{12} = \{HHH, HHT, TTH, TTT\}; A_{13} = \{HHH, HTH, THT, TTT\}; A_{23} = \{HHH, HTT, THH, TTT\};$$

$$P(A_{12}) = P(A_{13}) = P(A_{23}) = 4/8 = 1/2;$$

$$A_{12} \cap A_{13} = \{HHH, TTT\};$$

$$P(A_{12} \cap A_{13}) = 2/8 = 1/4 = P(A_{12})P(A_{13}).$$

On the other hand,

$$A_{12} \cap A_{13} \cap A_{23} = \{HHH, TTT\};$$

$$P(A_{12} \cap A_{13} \cap A_{23}) = 2/8 = 1/4 \neq P(A_{12})P(A_{13})P(A_{23}).$$

Sequential conditioning

By the definition of conditional probability:

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

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

This is useful for computing probabilities of sequential events.

E.g: What is the probability of dealing two aces in a row?

More generally:

$$P(A_1 \cap \dots \cap A_n) = P(A_1)P(A_2|A_1)P(A_3|A_1A_2) \dots P(A_n|A_1 \dots A_{n-1})$$

BIOS 660/BIOS 672 (3 Credits)

Notes 5 – 9 / 20

The Birthday Problem

In a group of n students in a class, what is the probability that at least two have the same birthday?

Solution:

Suppose we order the n students in an arbitrary order. Let D_j be the event that the first j have different birthdays. Based on page 22 in Notes 3, we have

$$P(D_j) = \frac{\text{No. of Samples with No Repetition}}{\text{No. of Samples}} = \frac{365!/(365-j)!}{365^j}$$

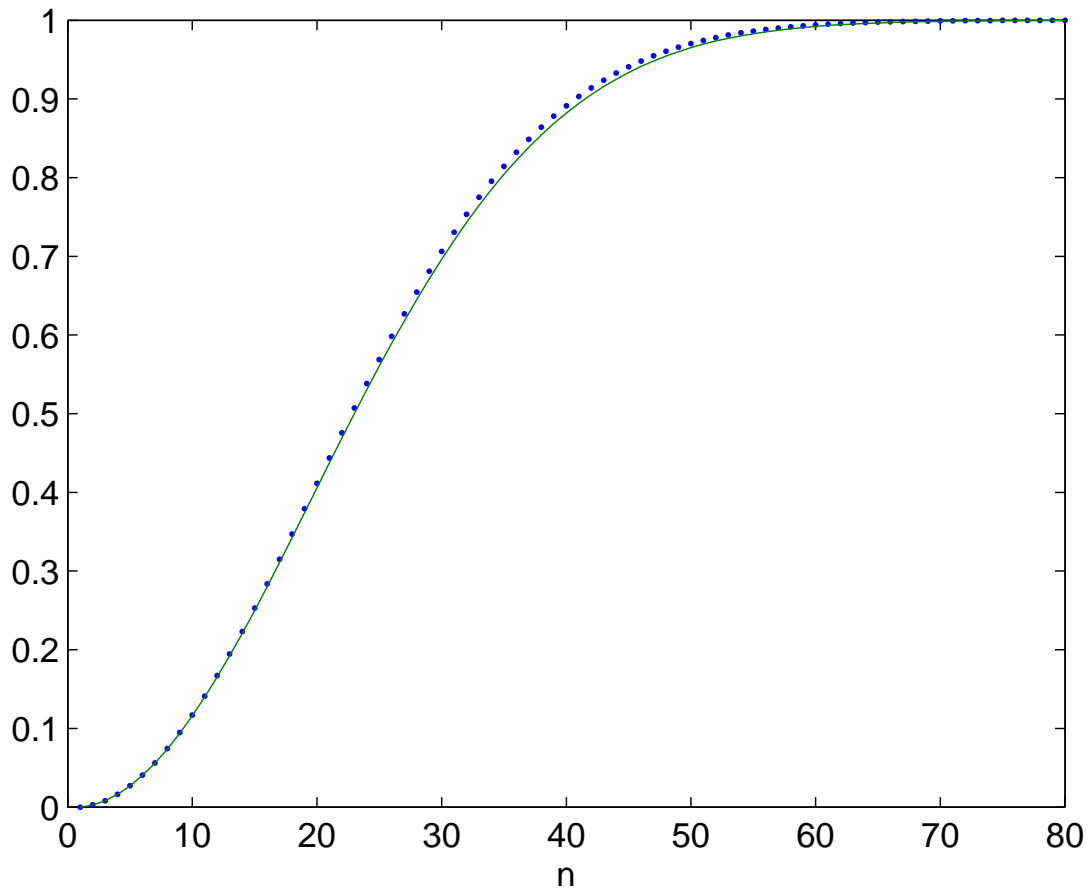
Let $D_j = \{A_1, A_2, A_3, \dots, A_j\}$, where A_1 is the birth day of the first person, A_2 is the birth day of the second person, etc., and all the A_i ($i = 1, 2, \dots, j$) are different. Based on sequential conditioning,

$$\begin{aligned} P(D_j) &= P(\{A_1, A_2, A_3, \dots, A_j\}) \\ &= P(A_1)P(A_2|A_1)P(A_3|A_1A_2) \dots P(A_j|A_1 \dots A_{j-1}) \\ &= \frac{365}{365} \frac{365-1}{365} \dots \frac{365-j+1}{365} = \frac{365!/(365-j)!}{365^j} \end{aligned}$$

BIOS 660/BIOS 672 (3 Credits)

Notes 5 – 10 / 20

The Birthday Problem



BIOS 660/BIOS 672 (3 Credits)

Notes 5 – 11 / 20

Turning around probabilities

Also by the definition of conditional probability:

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)} \quad (2)$$

This is useful for computing conditional probabilities when the reverse conditioning is easier to compute.

E.g.: Prob. that it will rain given that it is thundering vs. prob. that it thundered given that it is raining.

(2) is called sometimes **Bayes' rule**. It is often used in a context where we want to know the probability that a particular hypothesis is true. We have an *a priori* belief in whether or not the hypothesis is true, then update that probability by collecting data.

E.g.: Suppose a priori boys are equally likely to be born as girls. Say 90% of boys play with trucks. Baby X plays with trucks. What is the probability that Baby X is a boy?

BIOS 660/BIOS 672 (3 Credits)

Notes 5 – 12 / 20

Decomposition Formula (Total Probability)

Let $\{A_1, A_2, \dots\}$ be a partition of Ω . Let B be any subset in Ω . Then

$$P(B) = \sum_{i=1}^{\infty} P(A_i)P(B|A_i)$$

Proof:

(So called) Bayes' Theorem

Let $\{A_1, A_2, \dots\}$ be a partition of Ω . Let B be any subset in Ω . Then

$$P(A_j|B) = \frac{P(B|A_j)P(A_j)}{\sum_{i=1}^{\infty} P(B|A_i)P(A_i)}$$

Proof:

Bayes and Screening

An important application of Bayes' theorem is screening.

Notation:

Let D be the disease:

D means diseased

\overline{D} means "no disease"

and T be the diagnostic test:

T^+ means a positive test

T^- means a negative test.

Then we have that the **positive predictive value**

$$\begin{aligned} P(D|T^+) &= \frac{P(D) P(T^+|D)}{P(D) P(T^+|D) + P(\overline{D}) P(T^+|\overline{D})} \\ &\equiv \frac{\text{prevalence} \times \text{sensitivity}}{\text{prev.} \times \text{sens.} + (1 - \text{prev.}) \times (1 - \text{specificity})} \end{aligned}$$

BIOS 660/BIOS 672 (3 Credits)

Notes 5 – 15 / 20

Papanicolaou Example

Let D be cervical cancer.

$P(D)$ we'll take to be 1 in 21,000, which is the approximate annual incidence rate in the US (SEER 2002 estimate).

$$P(D) = .00004762$$

Let us take the sensitivity ($P(T^+|D)$) to be 0.71 and the specificity ($1 - P(T^+|\overline{D})$) to be 0.75.

Thence the positive predictive value is:

$$\begin{aligned} P(D|T^+) &= \frac{0.00004762 \times 0.71}{0.00004762 \times 0.71 + (1 - 0.00004762) \times (1 - 0.75)} \\ &= 0.000135 \end{aligned}$$

That means that for every 1,000,000 positive results, about 135 truly have cervical cancer. But for any particular patient, testing positive increases the probability of having the disease by a factor of 2.8!

BIOS 660/BIOS 672 (3 Credits)

Notes 5 – 16 / 20

Relative risks and relative odds

Given two conditions—smokers (S) and non-smokers (\bar{S}), say—then we say the *relative risk* of a disease (D)—lung cancer, say—due to smoking is:

$$RR = \frac{P(D|S)}{P(D|\bar{S})}$$

The *relative odds* of the disease (D) due to smoking is:

$$OR = \frac{\frac{P(D|S)}{1-P(D|S)}}{\frac{P(D|\bar{S})}{1-P(D|\bar{S})}}$$

Of course, if $P(D|S) \approx 0$ and $P(D|\bar{S}) \approx 0$, (*rare disease*) then

$$OR \approx \frac{\frac{P(D|S)}{1}}{\frac{P(D|\bar{S})}{1}} = RR$$

Bayes and Case Control Studies

Consider

$$OR(D|S) \equiv \frac{P(D|S)}{1 - P(D|S)} / \frac{P(D|\bar{S})}{1 - P(D|\bar{S})}$$

Now consider the numerator, which from Bayes theorem,

$$\frac{\frac{P(D)P(S|D)}{P(S)}}{\frac{P(\bar{D})P(S|\bar{D})}{P(S)}} = \frac{P(D)P(S|D)}{P(\bar{D})P(S|\bar{D})}$$

Do the same for the denominator, and get,

$$OR(D|S) = \frac{\frac{P(D)P(S|D)}{P(\bar{D})P(S|\bar{D})}}{\frac{P(D)P(\bar{S}|D)}{P(\bar{D})P(\bar{S}|\bar{D})}} = \frac{\frac{P(S|D)}{P(\bar{S}|D)}}{\frac{P(S|\bar{D})}{P(\bar{S}|\bar{D})}} = OR(S|D)$$

Additional Reading

See Chapter 1.2-1.3 in Casella and Berger.

BIOS 660/BIOS 672 (3 Credits): Probability and Statistical Inference I

Jianwen Cai

<https://sakai.unc.edu/portal/site/bios660-bios672-3-credits>

Notes 6

Random Variables	2
Random Variables	3
Random Variables (Formal Definition)	4
Conventions	5
Examples	6
Distribution Functions	7
Cumulative Distribution Functions	8
Some Properties of $F(y)$	9
Induced Probability Space	10
Identically distributed rvs	11
Types of Random Variables	12
Discrete random variables	13
Discrete random variables	14
Probability mass function	15
Properties of the pmf	16
Example	17
Continuous Random Variables	18
Continuous Random Variables	19
Absolute Continuity	20
cdf — density relation	21
Properties	22
Notes	23
Notes (cont.)	24
Notes (cont.)	25
Stochastic ordering	26

Random Variables

Suppose we start with a probability space (Ω, \mathcal{A}, P) . Instead of referring to outcomes and events observed from the sample space Ω , it is often convenient to assign a number to each possible outcome and record that instead.

Example: Flip a coin three times.

- $\Omega = \{HHH, HHT, HTH, HTT, THH, THT, TTH, TTT\}$
- $\mathcal{A} = 2^\Omega$
- Define a random variable Y to be the number of heads.
- $Y(HHH) = 3, Y(HHT) = 2, Y(HTT) = 1$, etc.

BIOS 660/BIOS 672 (3 Credits)

Notes 6 – 3 / 26

Random Variables (Formal Definition)

Definition: A random variable Y is a real-valued and measurable function defined on a probability space. That is, $Y : \Omega \rightarrow \mathbb{R}$.

Every point ω in Ω maps to a point in \mathbb{R} , namely $Y(\omega)$.

Conversely, we define the *inverse image* under Y of a subset B of \mathbb{R} as

$$Y^{-1}(B) = \{\omega : Y(\omega) \in B\}$$

The definition of a random variable requires that the inverse image of every Borel set $B \subset \mathbb{R}$ is an element of \mathcal{A} . This property allows us to assign probabilities to random variables. More precisely,

$$P\{Y \in B\} = P\{Y^{-1}(B)\}$$

BIOS 660/BIOS 672 (3 Credits)

Notes 6 – 4 / 26

Conventions

A Random Variable is a set function which takes values on the real line (for now). Often the argument is omitted and one writes Y instead of $Y(\omega)$.

Random variables are usually denoted by capital letters (e.g. Y).

Values which random variable can take on are denoted by lower case letters (e.g. y).

Example: Coin Tosses

$$\Omega : \{H, T\}$$
$$Y(H) = 1, \quad Y(T) = 0$$

If $P(\text{head}) = .5$, then $P(Y = 1) = .5$

BIOS 660/BIOS 672 (3 Credits)

Notes 6 – 5 / 26

Examples

Roll of a die: $\Omega = \{1, 2, 3, 4, 5, 6\}$

Define $Y(\omega) = \omega$

$$P\{Y(\omega) = \omega\} = P(\omega) = 1/6$$

An artificial example:

Suppose $\Omega = \{\omega : 0, \pm 1, \pm 2, \dots\}$

Define

$$Y(\omega) = a \quad \text{if} \quad \omega \leq 0$$

$$Y(\omega) = b \quad \text{if} \quad \omega > 0$$

$$P\{Y = a\} = \sum_{i=-\infty}^0 P(i) \quad ,$$

$$P\{Y = b\} = \sum_{i=1}^{\infty} P(i)$$

BIOS 660/BIOS 672 (3 Credits)

Notes 6 – 6 / 26

Cumulative Distribution Functions

Distribution Functions are used to describe the behavior of a rv.

Definition The *cumulative distribution function (cdf)* of a random variable Y is a real valued function $F(y)$ defined by

$$F(y) = P\{Y \leq y\} = P\{\omega : Y(\omega) \leq y\}$$

Example: cdf of a die: $F(y) = y/6$

Definition The *survival function* of Y is defined by

$$S(y) = 1 - F(y) = P(Y > y)$$

BIOS 660/BIOS 672 (3 Credits)

Notes 6 – 8 / 26

Some Properties of $F(y)$

1. $0 \leq F(y) \leq 1$
2. $\lim_{y \rightarrow -\infty} F(y) = 0$
3. $\lim_{y \rightarrow \infty} F(y) = 1$
4. F is nondecreasing: i.e. if $a < b$, then $F(a) \leq F(b)$
5. F is right continuous: that is, for any b and any decreasing sequence $b_n, n \geq 1$ that converges to b , $\lim_{n \rightarrow \infty} F(b_n) = F(b)$
6. $P\{a < Y \leq b\} = F(b) - F(a)$

These properties can all be proved using the properties of probability measures.

BIOS 660/BIOS 672 (3 Credits)

Notes 6 – 9 / 26

Induced Probability Space

All probability questions about a random variable can be answered via its *cdf*.

Every random variable defined on a probability space induces a probability space on \mathbb{R} :

$$(\Omega, \mathcal{A}, P) \longrightarrow Y(\omega) \longrightarrow (\mathbb{R}, \mathcal{B}, F(\cdot))$$

- Points in Ω are transformed to points on \mathbb{R} (Real line)
- Sets (events in \mathcal{A}) are mapped into intervals on real line, i.e., into members of the Borel sets, \mathcal{B} .
- P is replaced by $F(\cdot)$.

Because of this, the abstract notion of a sample space recedes, and attention is usually given primarily to random variables and their distributions.

We will sometimes refer to the ‘sample space’ of a random variable, which will be taken to be the values in \mathbb{R} that a random variable takes on.

Identically distributed rvs

The cdf does not contain information about the original sample space.

Example: Toss a coin n times. The number of heads and number of tails have the same distribution.

Definition: Two rvs X and Y are identically distributed if for every Borel set $A \subset \mathbb{R}$,
 $P(X \in A) = P(Y \in A)$.

Theorem C&B 1.5.10 The following two statements are equivalent:

- a. The rvs X and Y are identically distributed
- b. $F_X(x) = F_Y(x)$ for every x .

Note that two rvs can have the same distribution even if they are not equal to one another.

The distinction between two rvs being equal and having the same distribution will become important later in questions of convergence.

Types of Random Variables

A random variable Y can be

- *discrete* - Y takes on a finite or countably infinite number of values
- *continuous* - the range of Y consists of subsets of the real line.
- *mixed* - best to see this with an example

Example of a mixed random variable:

Consider a random variable with *cdf* given by:

$$F(x) = \begin{cases} 0 & x < 0 \\ x/2 & 0 \leq x < 1 \\ 2/3 & 1 \leq x < 2 \\ 11/12 & 2 \leq x < 3 \\ 1 & 3 \leq x \end{cases}$$

Discrete random variables

13 / 26

Discrete random variables

If a random variable Y takes only a finite or countable number of values, then the *cdf* can be expressed as:

$$F(y) = P\{Y \leq y\} = \sum_{z \leq y} P\{Y = z\}$$

If the sample space of Y is $\Omega = \{y_1, y_2, \dots\}$, then

$$F(y) = \sum_{y_i \leq y} P\{Y = y_i\}$$

Probability mass function

Definition The *prob. mass function (pmf)* or *frequency function* is a function $f(y)$ defined by

$$f(y) = P\{Y(\omega) = y\}$$

Thus, we can write $F(y) = P\{Y \leq y\} = \sum_{z \leq y} f(z)$.

If the sample space of Y is $\Omega = \{y_1, y_2, \dots\}$, then

$$f(y_i) = P(Y = y_i) = P(y_{i-1} < Y \leq y_i) = F(y_i) - F(y_{i-1})$$

Example: Suppose Y is a random variable that takes the values 0, 1 or 2 with probability .5, .3, and .2, respectively. Then

$$f(0) = 0.5, f(1) = .3, \text{ and } f(2) = .2.$$

Properties of the pmf

Definition: The *domain* of a random variable Y is the set of all values of y for which $f(y) > 0$.

Properties of the pmf:

1. $f(y) > 0$ for at most a countable number of values y . For all other values y , $f(y) = 0$.
2. Let $\{y_1, y_2, \dots\}$ denote the domain of Y . Then

$$\sum_{i=1}^{\infty} f(y_i) = 1$$

An obvious consequence is that $f(y) \leq 1$ over the domain.

Example: What is the pmf of a deterministic rv (a constant)?

$$f(x) = 1 \text{ for } x = k \text{ and } f(x) = 0 \text{ for } x \neq k.$$

Example

In many applications, a formula can be used to represent the *pmf* of a random variable. Suppose Y can take values $1, 2, \dots$ with *pmf*

$$f(y) = \begin{cases} \frac{1}{y(y+1)} & y = 1, 2, \dots \\ 0 & \text{otherwise} \end{cases}$$

How would we determine if this is an allowable *pmf*?

$$\begin{aligned} \sum_{y=1}^{\infty} \frac{1}{y(y+1)} &= \sum_{y=1}^{\infty} \left(\frac{1}{y} - \frac{1}{y+1} \right) \\ &= \left(1 - \frac{1}{2} \right) + \left(\frac{1}{2} - \frac{1}{3} \right) + \dots \\ &= 1 - \lim_{y \rightarrow \infty} \frac{1}{y} \\ &= 1 \end{aligned}$$

Continuous Random Variables

18 / 26

Continuous Random Variables

Recall that a random variable is a function that maps a probability space $\{\Omega, \mathcal{A}, P\}$ to the real line, thereby inducing a new probability space:

$$\{\Omega, \mathcal{A}, P\} \longrightarrow Y \longrightarrow (\mathbb{R}, \mathcal{F}, F(\cdot))$$

We have discussed the setting where Y is discrete (i.e. Y can take on a finite or countably infinite number of values).

A random variable Y is called *continuous* if its distribution function $F(y) = P(Y \leq y)$ is a continuous function.

Absolute Continuity

The distribution of a continuous random variable is characterized by the probability of falling in intervals, e.g. $P(Y \in (a, b])$.

We will focus on *absolutely continuous* random variables.

Definition: A function $F(y)$ is *absolutely continuous* if it can be written

$$F(y) = \int_{-\infty}^y f(x)dx,$$

where for now, you may interpret \int as the usual Riemann integral.

A random variable is said to be absolutely continuous if its distribution function is absolutely continuous.

Note: Absolute continuity is stronger than continuity but weaker than differentiability. An example of an absolutely continuous function is one that is differentiable everywhere except for a countable number of points.

cdf — density relation

If $F(y)$ is absolutely continuous, $f(y)$ is called the *probability density function (pdf)* of Y and

$$F'(y) = \frac{dF(y)}{dy} = f(y).$$

Building on this idea,

$$P(a < Y \leq b) = F(b) - F(a) = \int_a^b f(x)dx$$

More generally, for a set B ,

$$P(Y \in B) = \int_B f(x)dx$$

Note that of course B has to be an “allowable” subset of the real line \mathbb{R} , that is, a Borel set.

Properties

In general, a function $f(x)$ is a *pdf* iff

1. $f(x) \geq 0$
2. $\int_{-\infty}^{\infty} f(x)dx = 1$

Examples:

- Suppose $F(x) = 1 - e^{-\lambda x}$ for $x > 0$ and $F(x) = 0$ otherwise. Is $F(x)$ a cdf? What is the associated pdf?
 - $F(x)$ is continuous and nondecreasing, and $\lim_{x \rightarrow -\infty} F(x) = 0$, $\lim_{x \rightarrow +\infty} F(x) = 1$, so $F(x)$ is a cdf.
 - $f(x) = \frac{d}{dx} F(x) = \lambda e^{-\lambda x}$
- What about $f(x) = 1/x^r$ for $x > 1$ and $f(x) = 0$ otherwise?

$$\int_1^{\infty} \frac{dx}{x^r} = \left[\frac{x^{1-r}}{(1-r)} \right] = \frac{1}{r-1}$$

Then $f(x)$ is not a pdf but $(r-1)f(x)$ is a pdf.

Notes

- $f(x)$ is not the probability that $Y = x$. In fact, if Y is an absolutely continuous random variable with density function $f(x)$, then $P(Y = x) = 0$. Why?

$$\begin{aligned} P(Y = x) &= \lim_{h \rightarrow 0} \int_{x-h}^{x+h} f(u) du \\ &= \lim_{h \rightarrow 0} [F(x+h) - F(x-h)] \\ &= F(x+) - F(x-) \\ &= 0 \end{aligned}$$

Notes (cont.)

- More generally, if B is a subset of \mathbb{R} with

$$\int_B dx = 0,$$

then if Y is an absolutely continuous random variable defined on \mathbb{R} , then $P(Y \in B) = 0$ also.

- Because $P(Y = a) = 0$, all the following are equivalent:

$$P(a \leq Y \leq b), \quad P(a \leq Y < b) \quad \text{and} \quad P(a < Y < b)$$

- Also, note that $f(x)$ can exceed one!

BIOS 660/BIOS 672 (3 Credits)

Notes 6 – 24 / 26

Notes (cont.)

- $f(x)$ can be interpreted as the *relative* probability that Y takes the value x . Why? By the mean value theorem, we can say

$$P(x < Y \leq x + \Delta) \approx f(x)\Delta$$

Thus

$$P(Y \in \text{interval of width } \Delta \text{ centered at } a) \approx f(a)\Delta$$

and

$$P(Y \in \text{interval of width } \Delta \text{ centered at } b) \approx f(b)\Delta$$

Hence, if $f(b) > f(a)$, we can say that it is more likely for Y to take the values near b rather than near a .

BIOS 660/BIOS 672 (3 Credits)

Notes 6 – 25 / 26

Stochastic ordering

Suppose Y is a rv and define $X = Y + 2$. Then $X > Y$ always.

Now suppose

$$X \sim F_X(t) = (1 - e^{-t}) \mathbf{1}(t > 0)$$

$$Y \sim F_Y(t) = (1 - e^{-2t}) \mathbf{1}(t > 0)$$

Then X is not always greater than Y , but it is likely to be.

Definition: X is *stochastically greater* than Y if

$$F_X(t) \leq F_Y(t) \text{ for all } t$$

$$F_X(t) < F_Y(t) \text{ for some } t$$

or equivalently

$$P(X > t) \geq P(Y > t) \text{ for all } t$$

$$P(X > t) > P(Y > t) \text{ for some } t$$

BIOS 660/BIOS 672 (3 Credits): Probability and Statistical Inference I

Jianwen Cai

<https://sakai.unc.edu/portal/site/bios660-bios672-3-credits>

Notes 7

Transformations of Random Variables	2
Functions of Random Variables	3
Simple example	4
Probability mapping	5
Discrete RVs	6
Continuous RVs	7
Continuous RVs	8
Examples	9
Linear Transformation	10
Normal Distribution	11
Square root of an exponential RV	12
Probability Integral Transform	13
Probability Integral Transform (cont.)	14
Inverse Probability Integral Transform	15
Inverse Probability Integral Transform (cont.)	16
Example: Cauchy Distribution	17
Non-monotone transformations	18
One-to-many	19
Quadratic transformation	20
Example	21
General Result—Theorem C-B 2.1.8	22
Example: A wrapped distribution	23

Functions of Random Variables

(C-B Chap 2.1 & Gut I.2)

If X is a rv with sample space $\mathcal{X} \subset \mathbb{R}$ and cdf $F_X(x)$ then any function of X , say $Y = g(X)$ is also a random variable. The new random variable Y has a new sample space $\mathcal{Y} = g(\mathcal{X}) \subset \mathbb{R}$. The objective is to find the cdf $F_Y(y)$ of Y .

Example: Suppose X is an exponential random variable with parameter 1, i.e. $F_X(x) = 1 - e^{-x}$, $f_X(x) = e^{-x}$. What is the distribution of $Y = X/\lambda$?

$$\begin{aligned} F_Y(y) &= P(Y \leq y) = P(X/\lambda \leq y) = P(X \leq \lambda y) \\ &= F_X(\lambda y) = 1 - e^{-\lambda y} \end{aligned}$$

Y is distributed $\text{exp}(\lambda)$ with density $f_Y(y) = \lambda e^{-\lambda y}$.

BIOS 660/BIOS 672 (3 Credits)

Notes 7 – 3 / 23

Simple example

Example: Change units – Fahrenheit to Celsius

$$C = \frac{5}{9}(F - 32)$$

Ranges:

$$20^\circ C < C < 30^\circ C \Leftrightarrow 68^\circ F < F < 86^\circ F$$

$$(\mathcal{Y}, \mathcal{B}, F_Y) \leftarrow g(\cdot) \leftarrow (\mathcal{X}, \mathcal{B}, F_X)$$

BIOS 660/BIOS 672 (3 Credits)

Notes 7 – 4 / 23

Probability mapping

For any Borel set A :

$$\begin{aligned}P(Y \in A) &= P(g(X) \in A) \\&= P(\{x \in \mathcal{X} : g(x) \in A\}) \\&= P(X \in g^{-1}(A)).\end{aligned}$$

where we have defined

$$g^{-1}(A) = \{x \in \mathcal{X} : g(x) \in A\}$$

Notice that $g^{-1}(A)$ is well defined even if $g(\cdot)$ is not bijective (one-to-one).

Example: Let $g(x) = x^2$.

Then

$$g^{-1}([-1, 1]) = [-1, 1]$$

But

$$g(g^{-1}([-1, 1])) = [0, 1]$$

Discrete RVs

Suppose that X is a discrete random variable with probability mass function $p(x) = P(X = x)$.

Then, the *pmf* of a 1-1 transformation $Y = g(X)$ is given by

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

In practice, one never sees many general results about transformations of discrete random variables because the results are so simple!

Continuous RVs

Consider the transformation $Y = g(X)$ where $g(x)$ is strictly increasing (consequently a one-to-one transformation), and suppose g is differentiable. This means that we can also define the *inverse function*, $g^{-1}(y)$.

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

The *pdf* of Y is thus,

$$f_Y(y) = \frac{d}{dy} F_Y(y) = F'_X[g^{-1}(y)] \frac{dg^{-1}(y)}{dy} = f_X(g^{-1}(y)) \frac{dx}{dy}$$

since

$$x = g^{-1}(y), \text{ so that } \frac{dx}{dy} = \frac{dg^{-1}(y)}{dy}$$

Continuous RVs

Suppose $Y = g(x)$ is still one-to-one, but decreasing instead of increasing.

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

and

$$\begin{aligned} f_Y(y) &= -F'_X(g^{-1}(y)) \frac{dg^{-1}(y)}{dy} = -f_X(g^{-1}(y)) \frac{dx}{dy} \\ &= f_X(g^{-1}(y)) \left| \frac{dx}{dy} \right| \end{aligned}$$

The last step follows because $\frac{dx}{dy}$ is negative.

Therefore, regardless of whether $Y = g(x)$ is increasing or decreasing, so long as it is monotonic, we have

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

Linear Transformation

Given X with pdf $f_X(x)$, let

$$Y = a + bX, \quad \frac{dy}{dx} = b$$

Then

$$f_Y(y) = f_X[g^{-1}(y)] \left| \frac{dx}{dy} \right| = f_X\left(\frac{y-a}{b}\right) \frac{1}{|b|}$$

This transformation is often used when X has mean 0 and standard deviation 1. The linear transformation above creates a rv Y with a distribution that has the same shape as that of X but has mean a and standard deviation b .

Conversely, if Y has mean a and standard deviation b , then $X = (Y - a)/b$ has mean 0 and standard deviation 1. This is called sometimes the “*Studentized*” transform.

BIOS 660/BIOS 672 (3 Credits)

Notes 7 – 10 / 23

Normal Distribution

Let $X \sim N(0, 1)$:

$$f_X(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}, \quad -\infty < x < \infty$$

The transformation

$$Y = \mu + \sigma X, \quad X = \frac{Y - \mu}{\sigma}$$

yields

$$f_Y(y) = f_X\left(\frac{y - \mu}{\sigma}\right) \frac{1}{\sigma} = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(y-\mu)^2}{2\sigma^2}}$$

More generally, a distribution is a member of the class of *location-scale distributions* if the distribution of a linear transformation of a random variable with that distribution has the same distribution, but with different parameters.

BIOS 660/BIOS 672 (3 Credits)

Notes 7 – 11 / 23

Square root of an exponential RV

We have already seen that a constant times an exponential random variable leads to another exponential random variable. Suppose $X \sim \exp(\lambda)$, so that

$$f_X(x) = \lambda e^{-\lambda x} \mathbf{1}(x \geq 0)$$

and consider the distribution of $Y = \sqrt{X}$.

The transformation

$$y = g(x) = \sqrt{x}, \quad x \geq 0$$

is one-to-one and has an inverse $x = y^2$ with $dx/dy = 2y$. Thus

$$f_Y(y) = f_X(y^2)2y = 2\lambda y e^{-\lambda y^2}, \quad y \geq 0$$

This distribution is a particular form of the Rayleigh distribution and is a special case of the χ , Rice and Weibull distributions.

Probability Integral Transform

Let $X \sim F_X(x)$. Define the transformation

$$Y = F_X(X) \in [0, 1], \quad X = F_X^{-1}(Y)$$

Here

$$\begin{aligned} \frac{dy}{dx} &= F'_X(x) = f_X(x) \\ f_Y(y) &= f_X[F_X^{-1}(y)] \frac{1}{f_X[F_X^{-1}(y)]} = 1 \end{aligned}$$

i.e. Y is uniform over $[0, 1]$.

Another way to see it is through the distribution function:

$$\begin{aligned} F_Y(y) &= P(Y \leq y) = P(F_X(X) \leq y) \\ &= P(X \leq F_X^{-1}(y)) = F_X(F_X^{-1}(y)) = y \end{aligned}$$

Probability Integral Transform (cont.)

The probability integral transform is useful in statistics for checking goodness of fit of a distribution to a set of data.

Example:

$$X \sim \exp(\lambda) \Rightarrow Y = 1 - \exp(-\lambda X) \sim U[0, 1]$$

If one has data X_1, \dots, X_n , one could compute the transformed data $Y_i = 1 - \exp(-\lambda X_i)$, $i = 1, \dots, n$, and check whether the Y_i 's appear uniformly distributed over the interval $[0, 1]$.

Inverse Probability Integral Transform

We can also start from the uniform distribution and do the inverse procedure.

Suppose $X \sim U[0, 1]$, so that $f_X(x) = 1$ and $F_X(x) = x$ for $x \in [0, 1]$. Let

$$Y = F^{-1}(X), \quad X = F(Y)$$

where $F(\cdot)$ is a non-decreasing absolutely continuous function $F : \mathbb{R} \rightarrow [0, 1]$, $F(y) = \int_{-\infty}^y f(x) dx$. Then

$$\frac{dx}{dy} = F'(y) = f(y) \Rightarrow f_Y(y) = f(y)$$

i.e. Y has the *pdf* corresponding to F .

Another way to see it is through the distribution function:

$$\begin{aligned} F_Y(y) &= P(Y \leq y) = P(F^{-1}(X) \leq y) \\ &= P(X \leq F(y)) = F(y) \end{aligned}$$

Inverse Probability Integral Transform (cont.)

The Inverse Probability Integral Transform is used extensively in simulation of random variables.

Example: Most random number generators generate random numbers uniformly in the interval $[0, 1]$. Suppose we want to generate random numbers from an exponential distribution with parameter λ .

Let $X \sim U[0, 1]$. We want Y with distribution function $F(y) = 1 - \exp(-\lambda y)$. Then we need the transformation

$$Y = F^{-1}(X) = \frac{-1}{\lambda} \log(1 - X)$$

The recipe is

1. Generate random numbers X_i uniformly over $[0, 1]$.
2. Compute $Y_i = -\log(1 - X_i)/\lambda$.

Example: Cauchy Distribution

Let θ be distributed uniformly between $(-\pi/2, \pi/2)$:

$$f(\theta) = \frac{1}{\pi}, \quad -\pi/2 < \theta < \pi/2$$

Consider $Y = \tan \theta$.

$$\begin{aligned} \frac{dy}{d\theta} = \sec^2 \theta &= 1 + \tan^2 \theta = 1 + y^2 \\ f_Y(y) = \frac{1}{\pi} \left| \frac{d\theta}{dy} \right| &= \frac{1}{\pi} \frac{1}{(1 + y^2)} \quad -\infty < y < \infty \end{aligned}$$

The distribution with this density is known as the Cauchy distribution.

One-to-many

What if the transformation is not 1-1? The trick is to start with the cdf of the transformed random variable.

Example: Let $Y = |X|$, and assume X is continuous.

$$\begin{aligned} F_Y(y) &= P\{Y \leq y\} = P\{-y \leq X \leq y\} = F_X(y) - F_X(-y) \\ f_Y(y) &= F'_X(y) - F'_X(-y)(-1) = f_X(y) + f_X(-y) \end{aligned}$$

Suppose

$$X \sim N(0, 1), \quad f_X(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}$$

Then

$$f_Y(y) = \frac{2}{\sqrt{2\pi}} e^{-y^2/2}, \quad 0 < y < \infty$$

BIOS 660/BIOS 672 (3 Credits)

Notes 7 – 19 / 23

Quadratic transformation

Let

$$Y = X^2, \quad \frac{dy}{dx} = 2x, \quad \left| \frac{dy}{dx} \right| = 2\sqrt{y}$$

Then

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

$$\begin{aligned} f_Y(y) &= F'_X(\sqrt{y}) \left(\frac{1}{2} y^{-\frac{1}{2}} \right) - F'_X(-\sqrt{y}) \left(-\frac{1}{2} y^{-\frac{1}{2}} \right) \\ &= \frac{1}{2\sqrt{y}} [f_X(\sqrt{y}) + f_X(-\sqrt{y})], \quad y > 0 \end{aligned}$$

BIOS 660/BIOS 672 (3 Credits)

Notes 7 – 20 / 23

Example

Suppose $X \sim N(0, 1)$, $Y = X^2$:

$$\begin{aligned}f_X(x) &= \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}x^2} \quad -\infty < x < \infty \\f_Y(y) &= \frac{1}{2\sqrt{y}} \left[\frac{2}{\sqrt{2\pi}} e^{-\frac{1}{2}y} \right] \\&= \frac{(y/2)^{-\frac{1}{2}} e^{-\frac{y}{2}}}{2\sqrt{\pi}}, \quad y > 0\end{aligned}$$

This is the density of a χ^2 distribution with 1 degree of freedom.

Result: If $X \sim N(0, 1)$, then $X^2 \sim \chi^2(1)$.

General Result–Theorem C-B 2.1.8

Suppose $Y = g(X)$ is not 1-1, but there are disjoint sets A_1, \dots, A_k that span the domain (sample space) of X such that $g(\cdot) = g_j(\cdot)$ is continuous and 1-1 on each A_j . This means that the inverse, $x = g_j^{-1}(y)$ exists on each A_j . Then

$$f_Y(y) = \sum_{j=1}^k f(g_j^{-1}(y)) \left| \frac{dg_j^{-1}(y)}{dy} \right|$$

Example: A wrapped distribution

Suppose $X \in \mathbb{R}$ with density $f_X(x)$ represents a random angle of rotation (in radians) from the x -axis on the unit circumference. The observed angle is

$$\Theta = X \bmod 2\pi, \quad \Theta \in [0, 2\pi)$$

because it is impossible to tell if the rotation involved more than one full turn. In this case

$$f_{\Theta}(\theta) = \sum_{j=-\infty}^{\infty} f(\theta + 2\pi j), \quad 0 \leq \theta < 2\pi$$

If $X \sim N(0, \sigma^2)$, then

$$f_{\Theta}(\theta) = \sum_{j=-\infty}^{\infty} \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(\theta + 2\pi j)^2}{2\sigma^2}\right)$$

BIOS 660/BIOS 672 (3 Credits): Probability and Statistical Inference I

Jianwen Cai

<https://sakai.unc.edu/portal/site/bios660-bios672-3-credits>

Notes 8

Expected Values	2
Expectation	3
Expectation of some continuous variables	4
Expectation of some discrete variables	5
cont.	6
Mean = area under survival curve	7
Survival curves	8
Properties of expectation	9
Method of indicators	10
Minimal property of the mean	11
Moments	12
Central moments	13
Variance	14
Skewness	15
Central and non-central moments	16
Method of indicators (again)	17

Expectation

- The *expected value* or *mean* of a rv X , denoted $E(X)$ or EX is

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

Provided the integral or summation exists.

- This is generalized for a function of a random variable $g(X)$ as

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

Notice we could also find the pdf or pmf of Y and use the first definition. Both give the same answer (HW).

BIOS 660/BIOS 672 (3 Credits)

Notes 8 – 3 / 17

Expectation of some continuous variables

- Let $X \sim U[a, b]$. Then

$$EX = \int_a^b x \frac{1}{b-a} dx = \frac{a+b}{2}$$

- Let $X \sim \text{Exp}(1)$, $f_X(x) = e^{-x} 1(x > 0)$. Then

$$EX = \int_0^{\infty} x e^{-x} dx = -x e^{-x} \Big|_0^{\infty} + \int_0^{\infty} e^{-x} dx = 1$$

- Let $X \sim N(0, 1)$. Then

$$EX = \int_{-\infty}^{\infty} x \frac{1}{\sqrt{2\pi}} e^{-x^2/2} dx = 0$$

(since the above integral is an odd function)

BIOS 660/BIOS 672 (3 Credits)

Notes 8 – 4 / 17

Expectation of some discrete variables

- Let $X = 1(A)$, where A is a Borel set. Then

$$EX = 0 \cdot P(A^c) + 1 \cdot P(A) = P(A)$$

- Let $X \sim \text{Binomial}(n, p)$ for n positive integer and $0 < p < 1$ (n is the number of independent identical binary trials and p is the probability of success). Then

$$f_X(x) = \binom{n}{x} p^x (1-p)^{n-x} \text{ for } x = 0, \dots, n$$

and

$$\begin{aligned} EX &= \sum_{x=0}^n x \frac{n!}{(n-x)!x!} p^x (1-p)^{n-x} \\ &= np \sum_{x=1}^n \frac{(n-1)!}{(n-x)!(x-1)!} p^{x-1} (1-p)^{n-x} \\ (\text{let } y = x-1) &= np \sum_{y=0}^{n-1} \frac{(n-1)!}{(n-1-y)!y!} p^y (1-p)^{n-1-y} \\ &= np \end{aligned}$$

An easier way later.

cont.

- Let $X \sim \text{Geom}(p)$ (e.x. toss till first success), $f_X(x) = pq^{x-1}$, $x = 1, 2, \dots$, where $q = 1 - p$. Then

$$\begin{aligned} EX &= \sum_{x=1}^{\infty} x \cdot pq^{x-1} = \sum_{x=1}^{\infty} p \frac{d}{dq} (q^x) \\ &= p \frac{d}{dq} \sum_{x=1}^{\infty} q^x = p \frac{d}{dq} \left(\frac{1}{1-q} - 1 \right) \\ &= p \frac{1}{(1-q)^2} = \frac{1}{p} \end{aligned}$$

Note: The differentiation operator can be moved outside the summation sign because the geometric series converges uniformly.

Mean = area under survival curve

For a non-negative random variable X (i.e. $f(x) = 0$ for $x < 0$),

$$E(X) = \begin{cases} \int_0^\infty (1 - F(x)) dx, & X \text{ continuous} \\ \sum_{x=0}^\infty (1 - F(x)), & X \text{ discrete} \end{cases}$$

Proof: Homework! (Exercise 2.14)

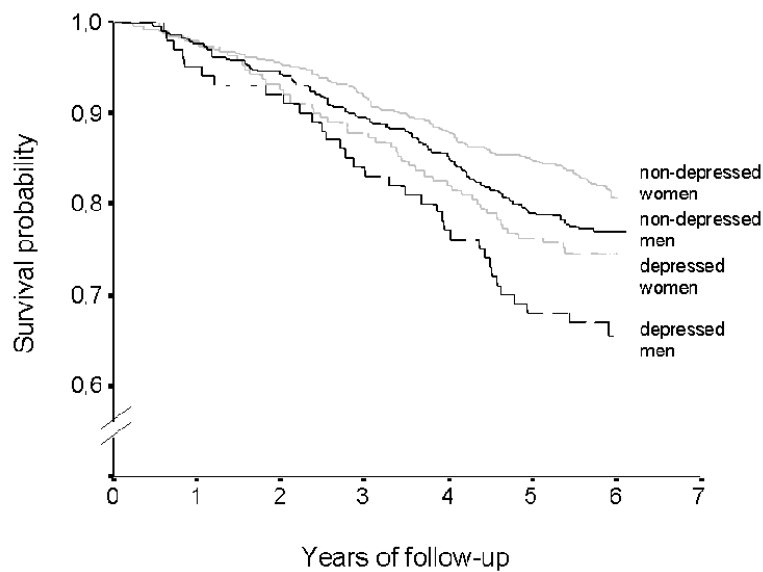
Example: Let $X \sim \text{Exp}(\lambda)$ with

$$F_X(x) = (1 - e^{-\lambda x}) 1(x > 0)$$

Then

$$EX = \int_0^\infty e^{-\lambda x} dx = \lambda$$

Survival curves



Properties of expectation

$$\mathbb{E}[ag(X) + c] = a\mathbb{E}[g(X)] + c, \quad a, c \text{ constants}$$

$$\mathbb{E}[g_1(X) + g_2(X)] = \mathbb{E}[g_1(X)] + \mathbb{E}[g_2(X)]$$

$$\text{If } g_1(x) \geq 0 \forall x, \quad \text{then } \mathbb{E}g_1(X) \geq 0$$

$$\text{If } a \leq g(X) \leq b \forall x, \quad \text{then } a \leq \mathbb{E}g(X) \leq b$$

(Proofs are immediate page 57)

Note: The linearity property above applies to two random variables X and Y that have the same distribution:

$$\mathbb{E}[g_1(X) + g_2(Y)] = \mathbb{E}[g_1(X)] + \mathbb{E}[g_2(Y)]$$

We will see later that this is true even if X and Y do not have the same distribution.

Method of indicators

An example of how the above properties are useful.

Let $X \sim \text{Binomial}(n, p)$ for n positive integer and $0 < p < 1$ (n is the number of independent identical binary trials and p is the probability of success). We can write

$$X = \sum_{i=1}^n I_i$$

where I_i is the indicator that the i^{th} trial is a success. We have

$$\mathbb{E}I_i = p$$

Therefore

$$\mathbb{E}X = \sum_{i=1}^n \mathbb{E}I_i = \sum_{i=1}^n p = np$$

Minimal property of the mean

Assuming all integrals exist,

$$\min_b \mathbf{E}(X - b)^2 = \mathbf{E}(X - \mathbf{E}X)^2 = \text{Var}X$$

Two proofs:

- By differentiation with respect to b (homework). Requires more assumptions.
- By sum-of-squares decomposition.

Moments

For a random variable X , the expectation of the polynomials $g(X) = X^r$, $r = 0, 1, 2, \dots$ are called the *moments* of X :

$$m_r = \mathbf{E}(X^r), \quad r = 0, 1, 2, \dots$$

These are sometimes called *non-central* moments or *moments about the origin*.

Notes:

- $m_0 = 1$.
- m_1 is the *mean*, usually denoted by $m_1 = \mu$.

Central moments

The r^{th} central moment of X is

$$\mu_r = E[X - EX]^r = E[X - \mu]^r, \quad r = 0, 1, 2, \dots$$

These are sometimes called *moments about the mean*.

Notes:

- $\mu_0 = 1$.
- $\mu_1 = 0$.
- μ_2 is the *variance*.
- μ_3 is related to the *skewness*.
- μ_4 is related to the *kurtosis*.

Variance

The second central moment μ_2 is called the *variance of X* and is usually denoted by σ^2 :

$$\sigma^2 = \mu_2 = \text{Var}(X) = E[X - \mu]^2 = E(X^2) - \mu^2$$

Notes:

- Useful property: For $a, b \in \mathbb{R}$,
$$\text{Var}(aX + b) = a^2 \text{Var}(X)$$
- Notice that $E(X^2) \neq [E(X)]^2$. In fact, $E(X^2) \geq [E(X)]^2$ because

$$\text{Var}(X) = E(X^2) - [E(X)]^2 \geq 0$$

Skewness

The *skewness* of a rv X is defined as

$$\alpha_3 = \frac{\mu_3}{(\mu_2)^{3/2}}, \quad \mu_3 = \mathbf{E}[(X - \mu)^3]$$

Skewness measures the symmetry of a distribution.

$$\mu_3 = 0 \Rightarrow \text{symmetric}$$

$$\mu_3 \geq 0 \Rightarrow \text{right skew}$$

$$\mu_3 \leq 0 \Rightarrow \text{left skew}$$

Central and non-central moments

Central moments can be written as a function of non-central moments and vice-versa:

$$\mu_2 = m_2 - m_1^2$$

$$\mu_3 = m_3 - 3m_1m_2 + 2m_1^3$$

etc.

Method of indicators (again)

Recall that $X \sim \text{Binomial}(n, p)$ can be written as $X = \sum_{i=1}^n I_i$ where I_i is the indicator that the i^{th} trial out of n is a success. We have

$$\mathbb{E}I_i = p, \quad \mathbb{E}I_i^2 = p, \quad \text{Var}I_i = p(1 - p)$$

Therefore

$$\begin{aligned} \mathbb{E}X^2 &= \mathbb{E}\left(\sum_{i=1}^n I_i\right)^2 = \mathbb{E}\sum_{i=1}^n \sum_{j=1}^n I_i I_j \\ &= \mathbb{E}\left(\sum_{i=1}^n I_i + \sum_{i \neq j} I_i I_j\right) = \sum_{i=1}^n p + \sum_{i \neq j} p^2 = np + (n^2 - n)p^2 \end{aligned}$$

and

$$\text{Var}X = \mathbb{E}X^2 - \mathbb{E}^2X = np + (n^2 - n)p^2 - (np)^2 = np(1 - p)$$

BIOS 660/BIOS 672 (3 Credits): Probability and Statistical Inference I

Jianwen Cai

<https://sakai.unc.edu/portal/site/bios660-bios672-3-credits>

Notes 9

Moment Generating Functions	2
Moment Generating Function	3
cont.	4
Example: Continuous	5
Example: Discrete	6
Linear transformations	7
Example: Continuous	8
Existence of moments	9
Can moments not exist?	10
Cont.	11
Can mgf's not exist?	12
Characterizing distributions	13
Convergence of mgfs	14
Application	15
Relationship to other transforms.	16
Characteristic functions	17
Characteristic Function	18
Examples	19
Properties	20
Characterizing distributions	21
Relationship to other transforms.	22

Moment Generating Function

(C-B 2.3, Gut III.3)

The *moment generating function (mgf)* of the rv X is defined to be

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

provided that the expectation exists in a neighbourhood $(-h, h)$ of $t = 0$.

Theorem: Suppose the mgf $M_X(t)$ of X exists for $t \in (-h, h)$ for some $h > 0$. Then for any positive integer n ,

$$\mathbf{E}(X^n) = M_X^{(n)}(0) = \left. \frac{d^n}{dt^n} M_X(t) \right|_{t=0}$$

Notice that $M_X(0) = 1$ always.

BIOS 660/BIOS 672 (3 Credits)

Notes 9 – 3 / 22

cont.

Proof: Assuming that we can interchange expectation and differentiation,

$$\begin{aligned} \frac{d}{dt} M_X(t) &= \frac{d}{dt} \mathbf{E}(e^{tX}) = \mathbf{E}\left(\frac{d}{dt} e^{tX}\right) = \mathbf{E}(X e^{tX}) \\ \Rightarrow \left. \frac{d}{dt} M_X(t) \right|_{t=0} &= \mathbf{E}(X) \\ \left. \frac{d^2}{dt^2} M_X(t) \right|_{t=0} &= \mathbf{E}(X^2 e^{tX}) \Big|_{t=0} = \mathbf{E}(X^2) \end{aligned}$$

Another way to see this is

$$M_X(t) = \mathbf{E}(e^{tX}) = \mathbf{E}\left(\sum_{n=0}^{\infty} \frac{t^n}{n!} X^n\right) = \sum_{n=0}^{\infty} \frac{t^n}{n!} \mathbf{E}X^n$$

so the moments can be obtained from a Taylor expansion of $M_X(t)$ around $t = 0$.

BIOS 660/BIOS 672 (3 Credits)

Notes 9 – 4 / 22

Example: Continuous

Mgf of an exponential rv: Let $f_X(x) = \lambda e^{-\lambda x} 1(x > 0)$. Then

$$\begin{aligned} M_X(t) &= \mathbb{E}(e^{tX}) = \int_0^\infty e^{tx} \cdot \lambda e^{-\lambda x} dx = \lambda \int_0^\infty e^{-(\lambda-t)x} dx \\ &= \lambda \frac{-1}{\lambda-t} e^{-(\lambda-t)x} \Big|_0^\infty = \begin{cases} \frac{\lambda}{\lambda-t} & \text{if } t < \lambda \\ \infty & \text{otherwise} \end{cases} \end{aligned}$$

This is fine as we only need the mgf to be defined near zero. To obtain the moments, assume $|t| < \lambda$:

$$M_X(t) = \frac{1}{1-t/\lambda} = \sum_{n=0}^{\infty} \frac{t^n}{\lambda^n} = \sum_{n=0}^{\infty} \frac{t^n}{n!} \mathbb{E}(X^n) \Rightarrow \mathbb{E}(X^n) = \frac{n!}{\lambda^n}$$

In particular,

$$\begin{aligned} \mathbb{E}X &= 1/\lambda \\ \text{Var}X &= \mathbb{E}X^2 - \mathbb{E}^2X = 2/\lambda^2 - 1/\lambda^2 = 1/\lambda^2 \end{aligned}$$

Example: Discrete

Mgf of a geometric rv: Let $f_X(x) = pq^{x-1}$, $x = 1, 2, \dots$ ($q = 1 - p$).

$$\begin{aligned} M_X(t) &= \mathbb{E}(e^{tX}) = \sum_{x=1}^{\infty} e^{tx} \cdot pq^{x-1} = \frac{p}{q} \sum_{x=1}^{\infty} (qe^t)^x \\ &= \frac{p}{q} \left(\frac{1}{1-qe^t} - 1 \right) = \frac{pe^t}{1-qe^t} \end{aligned}$$

The sum converges if $e^t q < 1$, that is, $t < \log(1/q)$.

The moments can be obtained by differentiation:

$$\mathbb{E}X = \frac{d}{dt} \left[\frac{p}{e^{-t} - q} \right]_{t=0} = \frac{pe^{-t}}{(e^{-t} - q)^2} \Big|_{t=0} = \frac{1}{p}$$

Linear transformations

For any constants a, b , the mgf of the rv $g(X) = aX + b$ is

$$M_{aX+b}(t) = e^{bt} M_X(at)$$

Proof:

$$\begin{aligned} M_{aX+b}(t) &= \mathbf{E}(e^{t(aX+b)}) = \mathbf{E}(e^{taX} e^{bt}) \\ &= e^{bt} \mathbf{E}(e^{(at)X}) = e^{bt} M_X(at) \end{aligned}$$

Assign as HW

Example: Continuous

Mgf of a Gaussian: Let $X \sim N(0, 1)$. Then

$$M_X(t) = \int_{-\infty}^{\infty} e^{tx} \cdot \frac{e^{-x^2/2}}{\sqrt{2\pi}} dx = e^{t^2/2} \int_{-\infty}^{\infty} \frac{e^{-(x-t)^2/2}}{\sqrt{2\pi}} dx = e^{t^2/2}$$

Also,

$$e^{t^2/2} = \sum_{n=0}^{\infty} \frac{1}{n!} \left(\frac{t^2}{2} \right)^n = \sum_{n=0}^{\infty} \frac{t^{2n}}{2^n n!} = \sum_{m: \text{ even}} \frac{t^m}{m!} \frac{m!}{2^{m/2} (m/2)!}$$

Matching coefficients of $t^m/m!$ we get that all the odd moments are zero and that the even moments are $E(X^m) = m!/(2^{m/2} (m/2)!)$.

Now let $Y = \mu + \sigma X$ so that $Y \sim N(\mu, \sigma^2)$,

$$M_Y(t) = e^{\mu t} M_X(\sigma t) = \exp \left(\mu t + \frac{\sigma^2 t^2}{2} \right)$$

Existence of moments

Too hard. May eliminate. **Theorem:** Suppose the mgf $M_X(t)$ of X exists for $t \in (-h, h)$ for some $h > 0$. Then all moments exist: $|EX^r| < \infty$ for all $r > 0$.

Proof: Fix $t \in (-h, h)$. There exists $C > 0$ such that

$$|x^r| \leq Ce^{|tx|}, \quad \forall x \in \mathbb{R}$$

(What is C ?)

so

$$\begin{aligned} |EX^r| &\leq E|X^r| \\ &\leq CE^{|tX|} \\ &\leq CE^{tX + e^{-tX}} \\ &\leq C[M_X(t) + M_X(-t)] < \infty \end{aligned}$$

Can moments not exist?

Example: Cauchy distribution

$$f_X(x) = \frac{1}{\pi} \frac{1}{1+x^2}, \quad x \in \mathbb{R}$$

The mean is

$$\begin{aligned} EX &= \int_{-\infty}^{\infty} x \frac{1}{\pi} \frac{1}{1+x^2} dx \\ &= \int_{-\infty}^0 \frac{1}{\pi} \frac{x}{1+x^2} dx + \int_0^{\infty} \frac{1}{\pi} \frac{x}{1+x^2} dx \\ &= \int_0^{\infty} \frac{1}{\pi} \frac{x}{1+x^2} dx - \int_0^{\infty} \frac{1}{\pi} \frac{x}{1+x^2} dx = ?? \end{aligned}$$

because

$$\int_0^{\infty} \frac{1}{\pi} \frac{x}{1+x^2} dx = \int_0^{\infty} \frac{1}{2\pi} \frac{dy}{1+y} = \frac{1}{2\pi} \log(1+y) \Big|_0^{\infty} = \infty$$

Cont.

One intuitive explanation for this is that the Cauchy distribution has infinite variance:

$$\begin{aligned} EX^2 &= \int_{-\infty}^{\infty} \frac{1}{\pi} \frac{x^2}{1+x^2} dx = \int_{-\infty}^{\infty} \frac{1}{\pi} \left(1 - \frac{1}{1+x^2} \right) dx \\ &= \frac{1}{\pi} \int_{-\infty}^{\infty} dx - 1 = \infty \end{aligned}$$

BIOS 660/BIOS 672 (3 Credits)

Notes 9 – 11 / 22

Can mgf's not exist?**Example: Cauchy distribution**

If the 1st and 2nd moments do not exist, certainly the mgf does not either!
(How would you prove this statement is true?)

Example: Log-normal distribution

If $X \sim N(0, 1)$ then $Y = e^X$ is called log-normal.

$$f_Y(y) = \frac{1}{y\sqrt{2\pi}} e^{-(\log y)^2/2}, \quad y > 0$$

For $n = 0, 1, 2, \dots$ the moments exist, but the mgf does not, i.e. the integral $E(e^{tY})$ does not converge. (Homework).

BIOS 660/BIOS 672 (3 Credits)

Notes 9 – 12 / 22

Characterizing distributions

For rvs with unbounded support, the moments do not specify the distribution: there exist distributions with different pdfs and yet have all the same moments (see Example C-B 2.3.10).

However, moments uniquely identify distributions when the rvs have bounded support.

Also, mgfs uniquely identify distributions when the mgfs exist.

Theorem: Let $F_X(x)$ and $F_Y(y)$ be cdfs all of whose moments exist.

1. If X and Y have bounded support, then $F_X(u) = F_Y(u)$ for all u iff $EX^n = EY^n$ for all $n = 0, 1, 2, \dots$
2. If the mgfs exist and $M_X(t) = M_Y(t)$ for all t in a neighborhood of 0, then $F_X(u) = F_Y(u)$ for all u .

BIOS 660/BIOS 672 (3 Credits)

Notes 9 – 13 / 22

Convergence of mgfs

Convergence of mgfs implies convergence of cdfs.

Theorem 2.3.12 C-B: Let X_1, X_2, \dots be a sequence of rvs with corresponding mgfs $M_{X_1}(t), M_{X_2}(t), \dots$ such that

$$\lim_{n \rightarrow \infty} M_{X_n}(t) = M_X(t), \quad \forall t \in (-h, h), \quad h > 0.$$

Then $\exists!$ a unique cdf $F_X(t)$ whose moments are given by $M_X(t)$ and for all x where $F_X(x)$ is continuous,

$$\lim_{n \rightarrow \infty} F_{X_n}(x) = F_X(x).$$

BIOS 660/BIOS 672 (3 Credits)

Notes 9 – 14 / 22

Application

Example: Normal approximation to Poisson.

Let $X \sim \text{Poisson}(\lambda)$, then $M_X(t) = \exp[\lambda(e^t - 1)]$

(Exercise C-B 2.33), with $\mathbf{E}X = \lambda$, $\text{Var}X = \lambda$.

Let $Y = (X - \lambda)/\sqrt{\lambda}$. Then

$$M_Y(t) = \mathbf{E}(e^{tY}) = \mathbf{E}(e^{t(X - \lambda)/\sqrt{\lambda}}) = e^{-\sqrt{\lambda}t} M_X(t/\sqrt{\lambda}).$$

Hence,

$$\begin{aligned} \log(M_Y(t)) &= -t\sqrt{\lambda} + \lambda(e^{t/\sqrt{\lambda}} - 1) \\ (\text{when } \lambda \text{ is large}) &= -t\sqrt{\lambda} + \lambda\left(\frac{t}{\sqrt{\lambda}} + \frac{t^2}{2\lambda} + \frac{t^3}{3!\lambda^{3/2}} + \cdots\right) \\ &= \frac{t^2}{2} + \frac{t^3}{3!\lambda^{1/2}} + \cdots \end{aligned}$$

Therefore,

$$\lim_{\lambda \rightarrow \infty} M_Y(t) = e^{t^2/2},$$

which is the mgf of a $N(0, 1)$ variable.

Relationship to other transforms

For continuous rvs:

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

is similar to the two-sided *Laplace transform* of the function $f_X(x)$.

The transform

$$S_X(t) = \log M_X(t) = \log \mathbf{E}e^{tX}$$

is called *cumulant generating function*. The derivatives at $t = 0$ are called *cumulants*. In particular, (Homework)

$$S_X(0) = 0, \quad S_X^{(1)}(0) = \mathbf{E}X, \quad S_X^{(2)}(0) = \text{Var}X$$

E.g. $X \sim N(\mu, \sigma^2)$,

$$S_X(t) = \log(e^{\mu t + \sigma^2 t^2/2}) = \mu t + \sigma^2 t^2/2$$

Characteristic Function

(Gut III.4)

The *characteristic function (cf)* of the rv X is defined as

$$\phi_X(t) = \mathbb{E}e^{itX} = \mathbb{E}[\cos(tX) + i \sin(tX)]$$

where $i^2 = -1$.

- The cf is complex-valued, $\phi_X(t) \in \mathbb{C}$.
- The cf always exists because

$$|\mathbb{E}e^{itX}| \leq \mathbb{E}|e^{itX}| = \mathbb{E}1 = 1$$

- For calculations, the cf can often be obtained from the mgf replacing the argument t by it (as long as the mgf exists!).

BIOS 660/BIOS 672 (3 Credits)

Notes 9 – 18 / 22

Examples

$X \sim \text{Exp}(\lambda)$:

$$\phi_X(t) = \frac{\lambda}{\lambda - it}$$

$X \sim N(\mu, \sigma^2)$:

$$\phi_X(t) = \exp\left(i\mu t - \frac{\sigma^2 t^2}{2}\right)$$

$X \sim \text{Geom}(p)$:

$$\phi_X(t) = \frac{pe^{it}}{1 - qe^{it}}$$

And the range is $t \in \mathbb{R}$ in all cases.

Also, for X with the Cauchy distribution

$$f_X(x) = \frac{1}{\pi(1+x^2)} \quad \Rightarrow \quad \phi_X(t) = e^{-|t|}$$

Exists!

BIOS 660/BIOS 672 (3 Credits)

Notes 9 – 19 / 22

Properties

1. $|\phi_X(t)| \leq \phi_X(0) = 1$
2. Complex conjugate: $\overline{\phi_X(t)} = \phi_X(-t)$ (Homework)
3. Linear transformations:

$$\phi_{aX+b}(t) = e^{ibt} \phi_X(at)$$

4. The distribution is symmetric about 0, $f_X(x) = f_X(-x)$, iff $\phi_X(t)$ is real.
5. Moment generation:

$$\phi_X^{(n)}(0) = \left. \frac{d^n}{dt^n} M_X(t) \right|_{t=0} = i^n \mathbf{E}(X^n)$$

BIOS 660/BIOS 672 (3 Credits)

Notes 9 – 20 / 22

Characterizing distributions

Characteristic functions uniquely identify distributions, always.

Theorem: If X and Y are rvs with cfs $\phi_X(t) = \phi_Y(t)$, then $F_X(u) = F_Y(u)$ for all u .

BIOS 660/BIOS 672 (3 Credits)

Notes 9 – 21 / 22

Relationship to other transforms

For continuous rvs:

$$\phi_X(t) = \int_{-\infty}^{\infty} e^{itx} f_X(x) dx$$

is similar to the *Fourier transform* of the function $f_X(x)$.

For discrete rvs:

$$\phi_X(t) = \sum_{x=-\infty}^{\infty} e^{itx} f_X(x)$$

is similar to the *discrete Fourier transform* of the sequence $f_X(x)$ (when x is sampled at equal intervals).

Thus Fourier transform tables can be helpful for finding cfs.

BIOS 660/BIOS 672 (3 Credits): Probability and Statistical Inference I

Jianwen Cai

<https://sakai.unc.edu/portal/site/bios660-bios672-3-credits>

Notes 10

Common Discrete Distributions	2
Why parametric models?	3
Discrete uniform	4
Binomial Distribution	5
Bernoulli distribution	6
Binomial Distribution	7
Example: Coin Tossing	8
Binomial distribution	9
Binomial cont.	10
Binomial example	11
Poisson Distribution	12
Poisson Distribution	13
Binomial vs Poisson	14
Binomial to Poisson via cfs	15
Poisson example	16
Uses of the Poisson	17
Hypergeometric Distribution	18
Example: Capture-Recapture Method	19
Hypergeometric Distribution	20
Hypergeometric Moments	21
Random sampling.	22
cont.	23
Hypergeometric vs. Binomial	24
cont.	25
cont.	26
Summary	27

Why parametric models?

- *Parametric models* or *distribution families* have a specific form but can change according to a fixed number of parameters.
- The objective is to model a population. Parametric models are often appropriate in common situations with similar mechanisms.
- Parametric models have many known and useful properties and are easy to work with. When fitting a population, only a few parameters need to be estimated: *parametric inference*.
- Sometimes one does not want to make parametric assumptions and would rather work with non-parametric models. But non-parametric models can be infinite dimensional. E.g. $f_X(x)$, $x = 0, 1, 2, \dots$ or $F_X(x)$, $x \in \mathbb{R}$.
- In this course we emphasize parametric models.

BIOS 660/BIOS 672 (3 Credits)

Notes 10 – 3 / 27

Discrete uniform

X has the *discrete uniform*(1, N) distribution if X is equally likely to be one of $\{1, 2, \dots, N\}$.

sample space: $\{1, 2, \dots, N\}$

pmf:

$$f_X(x) = \frac{1}{N}, \quad x = 1, 2, \dots, N$$

cdf:

$$F_X(x) = P(X \leq x) = \frac{x}{N}, \quad x = 1, 2, \dots, N$$

moments:

$$EX = \frac{N+1}{2}$$

This definition can be extended to the range N_0, \dots, N_1 (consecutive integers starting at any integer N_0 and ending with N_1) with $f_X(x) = 1/(N_1 - N_0 + 1)$.

BIOS 660/BIOS 672 (3 Credits)

Notes 10 – 4 / 27

Bernoulli distribution

Consider an experiment where outcomes are binary (say, Success or Failure) and the probability of success is p .

Define the following random variable

$$Y = \begin{cases} 1 & \text{outcome is success} \\ 0 & \text{outcome is failure} \end{cases}$$

Then, Y has a Bernoulli Distribution.

sample space: $\{0, 1\}$

pmf: $P(Y = 1) = p$ and $P(Y = 0) = 1 - p$. We can write this as:

$$f(y) = P(Y = y) = \begin{cases} p^y(1-p)^{(1-y)} & y = 0, 1 \\ 0 & \text{otherwise} \end{cases}$$

What are the cdf, mean and variance?

Binomial Distribution

Now consider a *series of n Bernoulli trials* where

1. trials are independent
2. Prob of success or failure is the same for each trial, i.e.

$$P(S_i) = p \text{ and } P(F_i) = q = 1 - p \text{ for the } i^{th} \text{ trial.}$$

More concisely, consider n *iid* (independent, identically distributed) Bernoulli rvs Y_i .

A *binomial(n, p)* random variable X is defined as the number of successes in n iid Bernoulli trials, each with probability p of success:

$$X = \sum_{i=1}^n Y_i$$

Example: Coin Tossing

Toss 3 coins ($n = 3$) $P(H) = p$, $P(T) = q$

$$\begin{aligned}P(HHH) &= p^3 & P(TTT) &= q^3 \\P(THH) &= p^2q & P(HTT) &= pq^2 \\P(HTH) &= p^2q & P(THT) &= pq^2 \\P(HHT) &= p^2q & P(TTH) &= pq^2\end{aligned}$$

The binomial distribution is concerned with the distribution of the **number** of successes (heads):

$$\begin{aligned}P(0H) &= P(3T) = q^3 \\P(1H) &= P(2T) = 3pq^2 \\P(2H) &= P(1T) = 3p^2q \\P(3H) &= P(0T) = p^3\end{aligned}$$

Binomial distribution

For any particular sequence of s successes and $n - s$ failures

$$P(SFS \dots F) = p^s q^{n-s}$$

However there are $\binom{n}{s}$ ways to get s successes from n trials.

Formally, the *binomial distribution* has:

sample space: $\{0, 1, \dots, n\}$

pmf:

$$f_Y(s) = \begin{cases} \binom{n}{s} p^s q^{n-s} & s = 0, 1, \dots, n \\ 0 & \text{otherwise} \end{cases}$$

Is it easy to check that the pmf sums to 1?

Binomial cont.

cdf:

$$F_Y(y) = \sum_{s=0}^y \binom{n}{s} p^s q^{n-s}$$

i.e.

$$F(y) = 0 \quad \text{for } y < 0$$

$$F(0) = q^n$$

$$F(1) = q^n + npq^{n-1}$$

$$F(2) = q^n + npq^{n-1} + \frac{n(n-1)}{2} p^2 q^{n-2}$$

\vdots

$$F(y) = 1 \quad \forall y \geq n$$

BIOS 660/BIOS 672 (3 Credits)

Notes 10 – 10 / 27

Binomial example

A physician treats $n = 10$ people having a particular disease where $P(\text{success}) = 1/4$. What are probabilities of outcomes?

$$f(s) = \binom{10}{s} p^s q^{10-s}, \quad p = 1/4$$

s	$f(s)$	$F(s)$	s	$f(s)$	$F(s)$
0	.056		6	.016	
1	.188		7	.003	
2	.282		8	.0004	
3	.250		9	.00003	
4	.146		10	.00000009	
5	.059				

BIOS 660/BIOS 672 (3 Credits)

Notes 10 – 11 / 27

Poisson Distribution

The Poisson distribution was derived by the French mathematician Poisson in 1837 as a limiting version of the binomial distribution.

Suppose Y has a binomial(n, p) distribution, but consider what happens when n becomes large, but p is small enough so that np stays constant, and equal to a fixed value λ .

$$\lim_{n \rightarrow \infty} \binom{n}{y} p^y q^{n-y} = \frac{e^{-\lambda} \lambda^y}{y!}$$

Proof:

$$\begin{aligned} f_Y(y) &= \frac{n!}{y!(n-y)!} p^y (1-p)^{n-y} \\ &= \frac{n!}{y!(n-y)!} \left(\frac{\lambda}{n}\right)^y \left(1 - \frac{\lambda}{n}\right)^{n-y} \\ &= \frac{\lambda^y}{y!} \left(1 - \frac{\lambda}{n}\right)^n \left(1 - \frac{\lambda}{n}\right)^{-y} \frac{(n-y+1) \dots n}{n^y} \end{aligned}$$

BIOS 660/BIOS 672 (3 Credits)

Notes 10 – 13 / 27

Binomial vs Poisson

Comparison of binomial and Poisson *pmf's*

$n = 5, p = 1/5 \ (\lambda = 1)$

y	binomial	poisson
0	.328	.368
1	.410	.368
2	.205	.184
3	.051	.061
4	.006	.015
5	.000	.003
6+	0	.001

BIOS 660/BIOS 672 (3 Credits)

Notes 10 – 14 / 27

Binomial to Poisson via cfs

Another way to see the convergence of a binomial to a Poisson is via convergence of the cf.

The cf of $X \sim \text{Binomial}(n, p)$ is

$$\phi_X(t) = (pe^{it} + 1 - p)^n$$

Let $p = \lambda/n$:

$$\phi_X(t) = \left(\frac{\lambda}{n} e^{it} + 1 - \frac{\lambda}{n} \right)^n = \left(1 + \frac{\lambda(e^{it} - 1)}{n} \right)^n$$

As $n \rightarrow \infty$,

$$\phi_X(t) \rightarrow \exp[\lambda(e^{it} - 1)]$$

which is the cf of a Poisson. Since cfs characterize distributions, the distribution for the rv X converges to a Poisson.

Poisson example

The death rate among females (Age 15-44) from pulmonary embolism is 4 per million. In a city of one million females (Age 15-44), what is the probability distribution of number of cases; i.e. $\lambda = 4$.

Note that $\lambda = np$ when $p = \frac{4}{1,000,000}$ is probability for a woman $\lambda = np = 10^6 \cdot 4/10^6 = 4$.

s	$e^{-4}4^s/s!$	s	$e^{-4}4^s/s!$	s	$e^{-4}4^s/s!$
0	.018	4	.20	8	.03
1	.07	5	.16	9	.01
2	.15	6	.10		
3	.20	7	.06		

Uses of the Poisson

The Poisson distribution is often used to describe:

- incidence of rare events in time or space
- failure of equipment
- incidence of rare diseases
- mortality
- pixel intensity in CT and PET imaging
- readings of molecular binding experiments (e.g. gene transcription)

Also arises in queueing theory (cashiers and internet), survival analysis, etc.

Hypergeometric Distribution

18 / 27

Example: Capture-Recapture Method

Suppose we wish to estimate how many fish there are in a lake.

Consider the following technique:

Capture a certain number of fish, 100 say, mark them, and return them to the lake.

Wait a sufficient amount of time for them to intermix again with the other fish.

Capture a new batch, 120 say, and see how many of these were previously marked.

Suppose there are 10 marked. Then the argument goes: the proportion marked that were recaptured should be in the same proportion as those amongst the non-captured. And so we can estimate the total number of fish in the lake to be the solution N to:

$$\frac{100 - 10}{N - 120} = \frac{10}{120} \quad \Rightarrow \quad N =$$

This is a situation where we might apply a hypergeometric distribution.

Hypergeometric Distribution

Suppose a population of N entities is made up of two types, and there are M of the first type; and so $N - M$ of the second type.

Suppose we take a sample of size K , and we wish to know X , the number in the sample of the first type.

The probability mass function of X is given by:

$$f_X(x) = \frac{\binom{M}{x} \binom{N-M}{K-x}}{\binom{N}{K}}$$

for $x = \max(0, M - N + K), \dots, \min(M, K)$.

The sample space is defined so that all binomial coefficients are valid.

We must have:

$$0 \leq x \leq K, \quad 0 \leq x \leq M, \quad 0 \leq K - x \leq N - M$$

Often $K < M$ and $K < N - M$ so the range becomes $0 \leq x \leq K$.

BIOS 660/BIOS 672 (3 Credits)

Notes 10 – 20 / 27

Hypergeometric Moments

Mean:

$$EX = \sum_{x=0}^K x \frac{\binom{M}{x} \binom{N-M}{K-x}}{\binom{N}{K}} = \sum_{x=1}^K \text{ditto}$$

We need the identity

$$x \binom{M}{x} = x \frac{M!}{x!(M-x)!} = M \frac{(M-1)!}{(x-1)!(M-x)!} = M \binom{M-1}{x-1}$$

or

$$\binom{M}{x} = \frac{M}{x} \binom{M-1}{x-1}$$

assuming everything is legit, i.e. all the numbers are positive integers, etc.. So

$$EX = \sum_{x=1}^K \frac{M}{N} \frac{\binom{M-1}{x-1} \binom{N-M}{K-1-(x-1)}}{\binom{N-1}{K-1}} = \frac{MK}{N}$$

BIOS 660/BIOS 672 (3 Credits)

Notes 10 – 21 / 27

Random sampling

Example: Vaccines are manufactured in batches of size N . Suppose n vials are sampled.

Decision Rule: If no vials are defective, batch is accepted.

What is the probability that batch with M defective vials is accepted?

Example: Experience suggests that a treatment for liver cancer should be considered effective if 20% of treated patients respond. A hospital plans to run a trial of a new treatment in 12 patients, and will consider the drug ineffective (no better than standard) if less than 2 patients respond. What is the probability that a drug with a true efficacy rate of 30% is classified as ineffective?

BIOS 660/BIOS 672 (3 Credits)

Notes 10 – 22 / 27

cont.

If we think of N being very large compared with n , then it makes sense to approximate this probability by thinking of sampling **with** replacement, in which case, we have

$$P(\text{Accept batch}) = (1 - M/N)^n$$

Table of $P(A)$ (approximate)

		M/N		
		.001	.01	.1
n	5	.995	.951	.59
	10	.990	.904	.349
	50	.95	.605	.005

BIOS 660/BIOS 672 (3 Credits)

Notes 10 – 23 / 27

Hypergeometric vs. Binomial

We can show that the limiting form of the hypergeometric pmf is the binomial pmf

$$\begin{aligned}
 P(s) &= \frac{\binom{M}{s} \binom{N-M}{n-s}}{\binom{N}{n}} \\
 &= \frac{\frac{M!}{s!(M-s)!} \frac{(N-M)!}{(n-s)!(N-M-n+s)!}}{\frac{N!}{n!(N-n)!}} \\
 &= \frac{\frac{n!}{s!(n-s)!} \frac{M!}{(M-s)!} \frac{(N-M)!}{(N-M-n+s)!}}{\frac{N!}{(N-n)!}}
 \end{aligned}$$

BIOS 660/BIOS 672 (3 Credits)

Notes 10 – 24 / 27

cont.

Note

$$\begin{aligned}
 \frac{M!}{(M-s)!} &= \frac{M(M-1)(M-2)\dots(M-s)!}{(M-s)!} \\
 &= M^s \left[1 \left(1 - \frac{1}{M}\right) \dots \left(1 - \frac{s-1}{M}\right) \right] \\
 \frac{N!}{(N-n)!} &= N^n \left[1 \left(1 - \frac{1}{N}\right) \dots \left(1 - \frac{n-1}{N}\right) \right] \\
 \frac{(N-M)!}{[(N-M)-(n-s)]!} &= (N-M)^{n-s} \\
 &\quad \left[1 \cdot \left(1 - \frac{1}{N-M}\right) \dots \left(1 - \frac{n-s-1}{N-M}\right) \right]
 \end{aligned}$$

BIOS 660/BIOS 672 (3 Credits)

Notes 10 – 25 / 27

cont.

Letting $N \rightarrow \infty, M \rightarrow \infty, \frac{M}{N} \rightarrow p,$

$$\begin{aligned} P(s) &= \frac{\binom{M}{s} \binom{N-M}{n-s}}{\binom{N}{n}} \\ &\sim \binom{n}{s} \frac{M^s (N-M)^{n-s}}{N^n} \\ &= \binom{n}{s} \left(\frac{M}{N}\right)^s \left(1 - \frac{M}{N}\right)^{n-s} \\ &\rightarrow \binom{n}{s} p^s (1-p)^{n-s} \quad \text{Binomial Distribution} \end{aligned}$$

BIOS 660/BIOS 672 (3 Credits)

Notes 10 – 26 / 27

Summary

Hypergeometric \rightarrow Binomial \rightarrow Poisson

$N \rightarrow \infty,$

$n \rightarrow \infty$

$\lambda = np$

$M \rightarrow \infty$

$p \rightarrow 0$

$\frac{M}{N} \rightarrow p$

$np \rightarrow \lambda$

BIOS 660/BIOS 672 (3 Credits)

Notes 10 – 27 / 27

BIOS 660/BIOS 672 (3 Credits): Probability and Statistical Inference I

Jianwen Cai

<https://sakai.unc.edu/portal/site/bios660-bios672-3-credits>

Notes 11

Geometric and Negative Binomial Distributions	2
Geometric Distribution	3
Memoryless property	4
Note	5
Example	6
Negative Binomial Distribution	7
Notes	8
Negative binomial sampling	9
Other parametrizations	10
Negative binomial vs. Poisson	11
cont.	12
Continuous Random Variables	13
Uniform Distribution	14
Uniform Distribution (cont.)	15
Notes	16
Exponential Distribution	17
Exponential Distribution	18
Interpretation	19
Notes	20
Shifted exponential	21
Double Exponential	22

Geometric Distribution

Consider a series of iid Bernoulli Trials with p =probability of success in each trial. Define a random variable X representing the number of trials until first success. *Note: X includes the trial at which the success occurs.* Then, X has a geometric distribution.

sample space: $\{1, 2, \dots\}$

pmf:

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

cdf:

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

Moments:

$$\begin{aligned} E(X) &= 1/p \\ \text{Var}(X) &= (1-p)/p^2 \end{aligned}$$

Memoryless property

Suppose $k > i$, then

$$P(X > k | X > i) = P(X > k - i)$$

Proof:

$$\begin{aligned} P(X > k | X > i) &= \frac{P(X > k)}{P(X > i)} = \frac{(1-p)^k}{(1-p)^i} \\ &= (1-p)^{k-i} = P(X > k - i) \end{aligned}$$

Example: Suppose X is # years you live

$$\begin{aligned} P(\text{survive two more years}) &= P(X > \text{current age} + 2 | X > \text{current age}) \\ &= P(X > 2) \end{aligned}$$

This model is clearly too simple for human populations (since we do age).

Note

Some texts use the geometric distribution to describe the distribution of the number of trials **before** the first success. In this case, the pdf changes to:

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

The moments also change. Be careful!

BIOS 660/BIOS 672 (3 Credits)

Notes 11 – 5 / 22

Example

In studies of human fertility, researchers are often interested in exposures that increase the time it takes for a couple to successfully conceive. If a woman has a $p = 25\%$ chance of conceiving in a particular menstrual cycle, what is the probability that she becomes pregnant on the 3rd cycle?

What is the expected number of cycles to pregnancy?

What is her probability of not becoming pregnant within 12 cycles? (this is the definition of clinical infertility)

BIOS 660/BIOS 672 (3 Credits)

Notes 11 – 6 / 22

Negative Binomial Distribution

Still in the context of iid Bernoulli trials, define a random variable corresponding to the number of trials required to have s successes. We say $X \sim \text{Negbin}(s, p)$.

sample space: $\{s, (s+1), \dots\}$

pmf: for $x = s, s+1, s+2, \dots$,

$$\begin{aligned} f(x) &= \binom{x-1}{s-1} p^{s-1} q^{x-s} \cdot p \\ &= \binom{x-1}{s-1} p^s q^{x-s} \end{aligned}$$

cdf: no closed form.

expectation: $E(X) = s/p$

Variance: $\text{Var}(X) = s(1-p)/p^2$

Notes

- Why the name? See C-B p.95.
- $X \sim \text{Negbin}(1, p)$ is the same as $X \sim \text{geometric}(p)$
- $\text{Negbin}(n, p)$ is the same as the sum of n $\text{geometric}(p)$ random variables
- Note that the outcome *SFFSFSSFS* may have been generated from a $\text{Binomial}(10, p)$ or from a $\text{NegBinom}(6, p)$. Need to know the experimental design to compute probabilities.

Negative binomial sampling

Example: Suppose a proportion p of the population possess a certain characteristic (e.g., have a certain disease). How many people should we expect to sample in order to collect r people with that characteristic? (Suppose the population is big enough so that you can assume sampling with replacement)

How many should we sample in order to be 95% sure of getting r people with the characteristic?

BIOS 660/BIOS 672 (3 Credits)

Notes 11 – 9 / 22

Other parametrizations

As with the geometric, some books define a negative binomial random variable as the number of **failures before the s^{th} success**. This is equal to the previous definition minus s .

sample space: $\{0, 1, 2, \dots\}$

pmf:

$$f(x) = \binom{s+x-1}{x} p^s q^x, \quad x = 0, 1, 2, \dots$$

cdf: no closed form.

expectation: $E(X) = s(1-p)/p$

Variance: $\text{Var}(X) = s(1-p)/p^2$

BIOS 660/BIOS 672 (3 Credits)

Notes 11 – 10 / 22

Negative binomial vs. Poisson

The negative binomial distribution is often good for modeling count data as an alternative to the Poisson.

In the previous parametrization, define

$$\lambda = \frac{s(1-p)}{p} \Leftrightarrow p = \frac{s}{s+\lambda}$$

Then we have

$$EX = \lambda$$

$$\text{Var}X = \frac{\lambda}{p} = \lambda \left(1 + \frac{\lambda}{s}\right) = \lambda + \frac{\lambda^2}{s}$$

For the Poisson we had that the variance equals the mean.

For the negative binomial, the variance is equal to the mean plus a quadratic term. Thus the negative binomial can capture overdispersion in count data.

cont.

In the previous parametrization, the pmf becomes

$$\begin{aligned} f(x) &= \binom{s+x-1}{x} p^s q^x = \frac{(s+x-1)!}{x!(s-1)!} \left(\frac{s}{s+\lambda}\right)^s \left(\frac{\lambda}{s+\lambda}\right)^x \\ &= \frac{\lambda^x}{x!} \frac{s(s+1)\dots(s+x-1)}{(s+\lambda)^x} \left(1 + \frac{\lambda}{s}\right)^{-s} \end{aligned}$$

Letting $s \rightarrow \infty$ we get that

$$f(x) \rightarrow \frac{\lambda^x}{x!} e^{-\lambda}$$

So for large s , the negative binomial can be approximated by a Poisson with parameter $\lambda = s(1-p)/p$.

We can see that also from convergence of the moments (Homework).

Uniform Distribution

A random variable X having a *pdf*

$$f(x) = \begin{cases} 1 & \text{for } 0 < x \leq 1 \\ 0 & \text{elsewhere} \end{cases}$$

is said to have a *uniform distribution* over the interval $(0, 1)$.

The *cdf* is:

$$F(y) = \int_{-\infty}^y f(x)dx = \int_0^y dx = \begin{cases} 0 & \text{for } y \leq 0 \\ y & \text{for } 0 \leq y \leq 1 \\ 1 & \text{for } y > 1 \end{cases}$$

BIOS 660/BIOS 672 (3 Credits)

Notes 11 – 14 / 22

Uniform Distribution (cont.)

Uniform: $Y \sim U[a, b]$:

sample space $[a, b]$

pdf:

$$f(y) = \begin{cases} \frac{1}{b-a} & \text{for } a < y \leq b \\ 0 & \text{elsewhere} \end{cases}$$

cdf:

$$F(y) = \int_a^y \frac{1}{b-a} dx = \begin{cases} 0 & \text{for } y < a \\ \frac{y-a}{b-a} & \text{for } a < y \leq b \\ 1 & \text{for } y > b \end{cases}$$

moments

$$\begin{aligned} E(Y) &= (a+b)/2 \\ \text{Var}(Y) &= \frac{(b-a)^2}{12} \end{aligned}$$

BIOS 660/BIOS 672 (3 Credits)

Notes 11 – 15 / 22

Notes

- The uniform extends to the continuous case the idea of equally likely outcomes.
- If $Y \sim U[0, 1]$, then $a + (b - a)Y \sim U[a, b]$.
- Useful for settings where individuals arrive at a destination, enter a study, etc. randomly over time.

Example: Suppose that during rush hour buses leave every 10 minutes, starting at 7:00am. Suppose you arrive at the bus stop according to a uniform distribution between 7am and 7:30am. What is the probability that you have to wait more than 5 minutes? What about 2 minutes?

Exponential Distribution

17 / 22

Exponential Distribution

Denoted $X \sim \text{Exp}(\lambda)$:

sample space: $y \geq 0$

pdf:

$$f(y) = \begin{cases} \lambda e^{-\lambda y} & \text{for } y \geq 0 \\ 0 & \text{elsewhere} \end{cases}$$

cdf:

$$F(y) = \int_0^y \lambda e^{-\lambda x} dx = \begin{cases} 1 - e^{-\lambda y} & \text{for } y \geq 0 \\ 0 & \text{for } y < 0 \end{cases}$$

moments:

$$\begin{aligned} E(Y) &= 1/\lambda \\ \text{Var}(Y) &= 1/\lambda^2 \\ M_X(t) &= \lambda/(\lambda - t) \end{aligned}$$

Interpretation

The exponential can be derived as the waiting time between Poisson events. Suppose that the number of events in a unit interval of time follows a $\text{Poisson}(\lambda)$ distribution. Then, let Y be the time until the first event.

$$P(Y > t) = P(0 \text{ events in } [0, t])$$

But, the number of events in $[0, t]$ follows a Poisson distribution with parameter λt . Hence,

$$P(Y > t) = e^{-\lambda t}$$

The *cdf* corresponding to Y is

$$F(t) = 1 - P(Y > t) = 1 - e^{-\lambda t}$$

and hence the density is

$$f(t) = \lambda e^{-\lambda t}$$

Notes

- Many books write the density as

$$f(y) = \begin{cases} \frac{1}{\theta} e^{-y/\theta} & \text{for } y \geq 0 \\ 0 & \text{elsewhere} \end{cases}$$

so that $E(Y) = \theta$ and $\text{Var}(Y) = \theta^2$

- The exponential has a *memorylessness property*, just like the geometric.

$$P(Y > s + t | Y > t) = P(Y > s)$$

Same interpretation as the geometric for continuous time: the probability of an event in a time interval depends only on the length of the interval, not the absolute time of the interval.

The underlying Poisson process is stationary: the rate λ is constant. (In the geometric, the prob. p of getting an event in every discrete time unit is constant.)

Shifted exponential

Let $X \sim \text{Exp}(\lambda)$ and $Y = X + v, v \in \mathbb{R}$.

Y has the *shifted exponential* distribution with pdf:

$$f(y) = \begin{cases} \lambda e^{-(y-v)\lambda} & \text{for } y \geq v \\ 0 & \text{elsewhere} \end{cases}$$

Double Exponential

The *double exponential* distribution is formed by reflecting the exponential distribution around zero. It has pdf:

$$f(x) = 0.5\lambda e^{-\lambda|x|}, \quad x \in \mathbb{R}$$

Suppose X has the above distribution with $\lambda = 1$.

Now let $Y = \sigma X + \mu, \mu \in \mathbb{R}$ (shifting) and $\sigma > 0$ (scaling).

Then Y has the *double exponential distribution* with pdf:

$$f_Y(y) = \frac{1}{2\sigma} \exp\left(-\frac{|y - \mu|}{\sigma}\right)$$

and moments

$$\mathbf{E}Y = \mu, \quad \text{Var}Y = 2\sigma^2$$

The double exponential distribution provides an alternative to the normal for centered data with fatter tails but all finite moments.

BIOS 660/BIOS 672 (3 Credits): Probability and Statistical Inference I

Jianwen Cai

<https://sakai.unc.edu/portal/site/bios660-bios672-3-credits>

Notes 12

Normal Distribution	2
Normal Distribution	3
Normal Moments	4
Standartization	5
Density integrates to 1	6
cont.	7
Notes	8
χ^2 distribution	9
Student's t and F distributions	10
More Distributions	11
Gamma distribution	12
Notes	13
Weibull distribution	14
Cauchy distribution	15
Beta distribution	16
cont.	17
Larger Families of Distributions	18
Location and Scale families	19
Group families	20
Group families: Examples	21

Normal Distribution

Introduced by De Moivre (1667 - 1754) in 1733 as an approximation to the binomial. Later studied by Laplace and others as part of the Central Limit Theorem. Gauss derived the normal as a suitable distribution for outcomes that could be thought of as sums of many small deviations.

sample space: $R = (-\infty, \infty)$

pdf: For $Y \sim N(\mu, \sigma^2)$,

$$f(y) = \frac{1}{\sqrt{2\pi}\sigma} e^{-(y-\mu)^2/2\sigma^2} \quad -\infty < y < \infty$$

If $\mu = 0$ and $\sigma^2 = 1$ is referred to as the *standard normal*.

cdf: There is no closed form. The notation $\Phi(x)$ is often used for $F(x) = P(Y \leq x)$ for the standard normal case. Many books have tables of its values for $x > 0$. Values for $x < 0$ can be obtained by the formula $\Phi(-x) = 1 - \Phi(x)$.

BIOS 660/BIOS 672 (3 Credits)

Notes 12 – 3 / 21

Normal Moments

Mean:

$$EY = \mu$$

Variance:

$$\text{Var}(Y) = E(Y - \mu)^2 = \sigma^2$$

Higher central moments:

$$E(Y - \mu)^m = \begin{cases} \frac{m!}{2^{m/2}(m/2)!} \sigma^m & m \text{ even} \\ 0 & m \text{ odd} \end{cases}$$

In particular:

$$\mu_3 = E(Y - \mu)^3 = 0 \quad (\text{Skewness})$$

$$\mu_4 = E(Y - \mu)^4 = 3\sigma^4$$

BIOS 660/BIOS 672 (3 Credits)

Notes 12 – 4 / 21

Standartization

Standartization:

$$Y \sim N(\mu, \sigma^2) \Leftrightarrow Z = \frac{Y - \mu}{\sigma} \sim N(0, 1)$$

Shifting and scaling:

$$Z \sim N(0, 1) \Leftrightarrow Y = \sigma Z + \mu \sim N(\mu, \sigma^2)$$

Easy to prove using the mgf.

BIOS 660/BIOS 672 (3 Credits)

Notes 12 – 5 / 21

Density integrates to 1

Theorem:

$$\int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-y^2/2} dy = 1$$

Proof:

This is not as easy as one might think. Call the integral I . Then,

$$\begin{aligned} I^2 &= \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-y^2/2} dy \cdot \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-x^2/2} dx \\ &= \frac{1}{2\pi} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} e^{-\frac{(x^2+y^2)}{2}} dx dy \end{aligned}$$

Now make a change of variables to polar coordinates, i.e. put

$$y = r \sin \theta \quad x = r \cos \theta, \quad 0 < \theta \leq 2\pi, \quad 0 < r < \infty$$

BIOS 660/BIOS 672 (3 Credits)

Notes 12 – 6 / 21

cont.

Now, $dx dy \rightarrow r dr d\theta$, because

$$\text{Jacobian} = J = \begin{vmatrix} \frac{\partial x}{\partial r} & \frac{\partial y}{\partial r} \\ \frac{\partial x}{\partial \theta} & \frac{\partial y}{\partial \theta} \end{vmatrix} = \begin{vmatrix} \cos \theta & \sin \theta \\ -r \sin \theta & r \cos \theta \end{vmatrix} = r$$

$$\begin{aligned} I^2 &= \int_{r=0}^{\infty} \int_{\theta=0}^{2\pi} \frac{e^{-r^2/2}}{2\pi} r dr d\theta \\ &= \int_0^{\infty} e^{-r^2/2} r dr = -e^{-r^2/2} \Big|_0^{\infty} = 0 - (-1) = 1 \end{aligned}$$

BIOS 660/BIOS 672 (3 Credits)

Notes 12 – 7 / 21

Notes

- Normal distribution useful in many practical settings
- Plays an important role in *sampling distributions in large samples*, since the Central Limit theorem says that sums of independent identically distributed random variables are approximately normal
- There are many important distributions that can be derived from functions of normal random variables (e.g. χ^2 , t , F). We will see much more on this later, for now, we will briefly present the *pdf's* and sample spaces of these distributions.

BIOS 660/BIOS 672 (3 Credits)

Notes 12 – 8 / 21

χ^2 distribution

If $Z \sim N(0, 1)$, then $X = Z^2$ has the χ^2 *distribution* with 1 degree of freedom.

More generally, we have the χ^2 *distribution* with ν degrees of freedom with pdf:

$$f(x) = \frac{(x/2)^{\frac{\nu}{2}-1} e^{-x/2}}{2\Gamma(\nu/2)}, \quad x > 0$$

where $\Gamma(a)$ is the complete gamma function,

$$\Gamma(a) = \int_0^\infty x^{a-1} e^{-x} dx$$

Note that if a is an integer, $\Gamma(a) = (a-1)!$.

The $\chi^2(\nu)$ distribution is a special case of the gamma distribution, so it is easier to derive its properties from the gamma.

Student's t and F distributions

Y has a t_k distribution (t with ν degrees of freedom) if its *pdf* can be written as:

$$f(y) = \frac{\Gamma[(\nu+1)/2]}{\sqrt{\nu\pi}\Gamma(\nu/2)} \frac{1}{(1+y^2/\nu)^{(\nu+1)/2}}, \quad -\infty < y < \infty$$

Y has an $F(\nu_1, \nu_2)$ distribution if its *pdf* can be written as:

$$f(y) = \frac{(\nu_1/\nu_2)\Gamma[(\nu_1+\nu_2)/2](\nu_1 y/\nu_2)^{\nu_1/2-1}}{\Gamma(\nu_1/2)\Gamma(\nu_2/2)(1+\nu_1 y/\nu_2)^{(\nu_1+\nu_2)/2}}, \quad 0 \leq y < \infty$$

There are many important properties and relationships between these three distributions (e.g. χ_k^2 is the distribution of the sum of the squares of k independent standard normals). We'll come back to these in a few weeks when we do *sampling distributions and transformations of the normal distribution*.

Gamma distribution

Notation: $Y \sim \text{gamma}(a, \lambda)$.

pdf:

$$f(y) = \frac{\lambda e^{-\lambda y} (\lambda y)^{a-1}}{\Gamma(a)}, \quad y \geq 0$$

where $\Gamma(a)$ is the *complete gamma function*,

$$\Gamma(a) = \int_0^{\infty} x^{a-1} e^{-x} dx$$

Note that if a is an integer, $\Gamma(a) = (a-1)!$.

cdf: In general, there is no closed form, unless a is an integer.

moments

$$\begin{aligned} \mathbf{E}(Y) &= a/\lambda \\ \mathbf{Var}(Y) &= a/\lambda^2 \end{aligned}$$

BIOS 660/BIOS 672 (3 Credits)

Notes 12 – 12 / 21

Notes

- The special case $a = 1$ corresponds to an exponential(λ)
- Can be thought of as a flexible generalization of the exponential (a can be interpreted as a *shape parameter*)
- The special case $\text{gamma}(n/2, 1/2)$, for integer n , corresponds to the χ^2 distribution with n degrees of freedom.
- We will see later in the class that the gamma distribution can be derived as the sum of a independent exponential(λ) distributions
- When a is an integer, the $\text{gamma}(a, \lambda)$ distribution can be derived as the distribution of time until the occurrence of the a^{th} event in a Poisson process.

BIOS 660/BIOS 672 (3 Credits)

Notes 12 – 13 / 21

Weibull distribution

This is another useful generalization of the exponential. It is useful to begin with the *cdf* instead of the *pdf*:

sample space: $[v, \infty]$

cdf:

$$F(y) = 1 - \exp \left[- \left(\frac{y - v}{\alpha} \right)^\beta \right], \quad y \geq v$$

It follows that the *pdf* is:

$$f(y) = \frac{\beta}{\alpha} \left(\frac{y - v}{\alpha} \right)^{(\beta-1)} \exp \left[- \left(\frac{y - v}{\alpha} \right)^\beta \right], \quad y \geq v$$

The usual case is $v = 0$.

If $\beta = 1$ we get an exponential with parameter $\lambda = 1/\alpha$.

Cauchy distribution

This is a famous distribution to mathematical statisticians, since it often serves as a useful counterexample.

pdf

$$f(y) = \frac{1}{\pi} \frac{1}{[1 + (y - \mu)^2/\sigma^2]} \quad \text{for } -\infty < y < \infty$$

The Cauchy with $\mu = 0$, $\sigma = 1$, corresponds to the t -distribution with 1 degree of freedom.

While the moments of the Cauchy are not defined, its quantiles are (HW).

The Cauchy is not just a pathological case. We'll see later that the ratio of two standard normals is Cauchy. So ratios of observations can be problematic (e.g. BMI).

Beta distribution

Notation: $Y \sim \text{beta}(a, b)$.

sample space: $[0, 1]$

pdf:

$$f(y) = \frac{y^{a-1}(1-y)^{b-1}}{B(a, b)}, \quad 0 \leq y \leq 1$$

where $B(a, b)$ is the (complete) Beta function,

$$B(a, b) = \int_0^1 x^{a-1}(1-x)^{b-1} dx = \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)},$$

where $\Gamma(a)$ is the complete gamma function. The normalizing constant is required so that $\int_0^1 f(x) dx = 1$.

Note that if a and b are integers, then $B(a, b)$ can be calculated in closed form.

cont.

cdf: In general, there is no closed form, except if a and b are integers.

moments

$$\begin{aligned} E(Y) &= \frac{a}{a+b} \\ \text{Var}(Y) &= \frac{ab}{(a+b)^2(a+b+1)} \end{aligned}$$

The beta distribution is very flexible, and can take a wide variety of shapes by varying its parameters.

Special case: $\text{beta}(1, 1) = U(0, 1)$.

* Read C-B Section 3.3 (normal, beta, Cauchy, lognormal and double exponential)

Location and Scale families

Let $f(x)$ be any pdf. Then the family of pdfs

$$f_{\mu,\sigma}(x) = \frac{1}{\sigma} f\left(\frac{x - \mu}{\sigma}\right)$$

for $\mu \in \mathbb{R}$, $\sigma > 0$, is called a location-scale family.

If $\mu = 0$ we get a scale family; if $\sigma = 1$ we get a location family.

Examples: Normal, Laplace, Cauchy, exponential.

Properties: Let $Z \sim f(z)$ and $X = \sigma Z + \mu$. Then

1. X has pdf $f_{\mu,\sigma}(x)$.
- 2.

$$E(X) = \sigma E(Z) + \mu, \quad \text{Var}(X) = \sigma^2 \text{Var}(Z)$$

Group families

Let \mathcal{G} be a class of 1-to-1 functions $g : \mathbb{R} \rightarrow \mathbb{R}$. The class of transformations \mathcal{G} is called a *transformation group* if

1. \mathcal{G} is closed under composition: $g_1, g_2 \in \mathcal{G}$ implies $g_2 \circ g_1 \in \mathcal{G}$.
2. \mathcal{G} is closed under inversion: $g \in \mathcal{G}$ implies $g^{-1} \in \mathcal{G}$.

Given a rv Z with cdf $F(z)$, the class

$$\{X = g(Z), g \in \mathcal{G}\}$$

is a group family.

Group families: Examples

- *Parametric*: Location-scale families, $g(z) = \sigma z + \mu$, $\sigma > 0$, $\mu \in \mathbb{R}$.
- *Non-parametric*: Let \mathcal{G} is the class of all continuous strictly increasing functions $g(z)$ such that

$$\lim_{z \rightarrow -\infty} g(z) = -\infty, \quad \lim_{z \rightarrow \infty} g(z) = \infty$$

Let Z be a rv supported on $(-\infty, \infty)$. Then the class $\{X = g(Z), g \in \mathcal{G}\}$ is the class of all rvs supported on $(-\infty, \infty)$ whose cdfs are continuous and strictly increasing.

- *Non-parametric*: Same as before with the additional restriction that Z has a symmetric distribution about 0 and g is odd: $g(-z) = -g(z)$. The generated rvs are now the class of all rvs with symmetric distributions about 0.

BIOS 660/BIOS 672 (3 Credits): Probability and Statistical Inference I

Jianwen Cai

<https://sakai.unc.edu/portal/site/bios660-bios672-3-credits>

Notes 13

Multinomial Distribution	2
Multiple Possible Outcomes of a Trial	3
continued	4
Multinomial distribution	5
More on the Poisson	6
The Poisson distribution again	7
continued	8
Poisson: Examples and generalizations	9
continued	10
Example: Failure of Equipment	11
Continued	12
Example: Leukemia	13
Example: Safety Testing of Vaccine	14
Example: Leukemia in Woburn, MA	15
continued	16
More Negative Binomial	17
Negative Binomial	18
Example	19
Example (continued)	20
Comments	21
Exponential Families	22
Exponential Families	23
Example: Binomial	24
Example: Gaussian	25
Theorem C-B 3.4.2	26
Example: Binomial	27
Indicator function	28
More examples	29
Natural parameters	30

Example: Gaussian	31
Curved exponential families	32
Probability Inequalities	33
Chebychev Inequality	34
Application	35
Normal tail bound	36
Multiple Random Variables	37
Multiple random measurements	38
Random Vectors	39
Example: Bivariate	40
Discrete Bivariate RVs	41
Example	42
Bivariate cdfs	43
Marginal distributions	44
Joint probabilities	45
Continuous Bivariate Random Variables	46
Continuous Bivariate RVs	47
Properties of the bivariate pdf	48
Example 1	49
Example 2	50
Conditional Distributions and Independence	51
Conditional Distributions - Discrete	52
Example: Discrete	53
Conditional Distributions - Continuous	54
Example 1	55
Example 2	56

Multiple Possible Outcomes of a Trial

Suppose: Result of drug trial is Failure, Partial Success and Success. Let

$$P(F) = p_1, \quad P(PS) = p_2, \quad P(S) = p_3, \quad (p_1 + p_2 + p_3) = 1$$

Suppose in a sample of size n

s_1 = Number of Failures

s_2 = Number of Partial Successes

s_3 = Number of Successes

where $s_1 + s_2 + s_3 = n$,

$$P(s_1, s_2, s_3) = \frac{n!}{s_1!s_2!s_3!} p_1^{s_1} p_2^{s_2} p_3^{s_3}$$

These are *multinomial probabilities*.

BIOS 660/BIOS 672 (3 Credits)

Notes 13 – 3 / 56

continued

The probability of any ordered arrangement resulting in s_1 “F”, s_2 “PS” and s_3 “S” is

$$p_1^{s_1} p_2^{s_2} p_3^{s_3}$$

However there are $\frac{n!}{s_1!s_2!s_3!}$ such ordered arrangements. Therefore

$$P(s_1, s_2, s_3) = \frac{n!}{s_1!s_2!s_3!} p_1^{s_1} p_2^{s_2} p_3^{s_3}$$

BIOS 660/BIOS 672 (3 Credits)

Notes 13 – 4 / 56

Multinomial distribution

The generalization to k classes gives us the *Multinomial Distribution*

$$p(s_1, s_2, \dots, s_k) = \frac{n!}{s_1! s_2! \dots s_k!} p_1^{s_1} p_2^{s_2} \dots p_k^{s_k}$$

where $\sum_{i=1}^k s_i = n$ and $\sum_{i=1}^k p_i = 1$.

BIOS 660/BIOS 672 (3 Credits)

Notes 13 – 5 / 56

More on the Poisson

6 / 56

The Poisson distribution again

sample space: $\{0, 1, 2, \dots\}$

pmf:

$$P(s) = \begin{cases} e^{-\lambda} \lambda^s / s! & s = 0, 1, 2, \dots \\ 0 & \text{otherwise} \end{cases}$$

cdf:

$$F(y) = \sum_{s=0}^y e^{-\lambda} \lambda^s / s!$$

expectation:

$$E(Y) = \lambda$$

Variance:

$$\text{Var}(Y) = \lambda$$

BIOS 660/BIOS 672 (3 Credits)

Notes 13 – 7 / 56

continued

Note: we can write

$$\frac{P(X = i + 1)}{P(X = i)} = \frac{\lambda}{i + 1}$$

BIOS 660/BIOS 672 (3 Credits)

Notes 13 – 8 / 56

Poisson: Examples and generalizations

Last time, we had an example concerning pulmonary embolism among young women. The rate is 4 per million, and we had looked at the pdf for the number of cases in a city with 1,000,000 women. But, suppose we are interested in a city that only has 100,000 women. How does the probability distribution change?

$$p = \frac{4}{1,000,000}, \quad n = 100,000, \quad np = \frac{4}{10} = .4 = \lambda$$

$$P(s) = e^{-.4}(.4)^s/s!$$

s	$p(s)$
0	.67
1	.27
2	.05
3	.007
4	.0007

BIOS 660/BIOS 672 (3 Credits)

Notes 13 – 9 / 56

continued

It is often useful to write

$$P(s) = e^{-\lambda n} (\lambda n)^s / s! \quad \begin{array}{l} \lambda = \text{rate/unit population} \\ n = \text{size of population} \end{array}$$

where λn is the expected number of events.

Note: "population" can be in units of 10, 100, 1000, etc. λ is the mean number of events per unit population.

Example: Failure of Equipment

A computer has a failure rate of 1 failure per 1,000 hours of use. How many failures would be expected in 500 hours of use?

What is the probability distribution of the number of failures?

s	$P(s) = e^{-.5} (.5)^s / s!$
0	.61
1	.30
2	.08
3	.01

More generally, the distribution may be written

$$P(s) = e^{-\lambda t} (\lambda t)^s / s!$$

where λt = number of failures (events) in t units of time.

Now suppose there are n computers, each being observed for t time units and each having a failure rate of λ , what is $P(s)$?

Continued

$$P(s) = e^{-\lambda nt} (\lambda nt)^s / s!$$

λ = rate per unit time per unit individual

n = number of units (individuals)

t = time frame

λnt = Expected number of events for n units
and time t

Example: Leukemia

Suppose the incidence rate for childhood leukemia is 3.9 per 100,000 per year for children less than 4 years old. In a population of 5,000 children observed for 10 years, what would the probability distribution be for number of cases?

$$\lambda nt = \frac{3.9}{100,000} \times 5,000 \times 10 = 1.95$$

$$P(s) = e^{-1.95} (1.95)^s / s!$$

s	$P(s)$
0	.14
1	.28
2	.27
3	.18
4	.08
5	.03
6	.01
7	.003

Example: Safety Testing of Vaccine

Suppose a vaccine contains m live virus per cm^3 . Suppose a sample of $v \text{ cm}^3$ of vaccine is tested. The expected number of virus in $v \text{ cm}^3$ is thus equal to mv . What is probability that vaccine tested will be free of a virus?

$$P(s) = e^{-mv} (mv)^s / s!$$

$$P(0) = e^{-mv}$$

e.g. Suppose $m = .005$ and $v = 600\text{cc}$. Then, $mv = 3$ and $P(0) = e^{-3} = .05$

BIOS 660/BIOS 672 (3 Credits)

Notes 13 – 14 / 56

Example: Leukemia in Woburn, MA

During a 19 year period 15 leukemias were observed. Is this an unusual event?

<u>Age</u>	<u>Population (n)</u>	<u>Rate per $10^5/\text{yr}(\lambda)$</u>	<u>Expected Number</u>
0 – 4	2120	6.27	.133
5 – 9	2191	3.09	.068
10 – 14	2969	2.04	.061
15 – 19	3592	2.19	.079
	10,872		.341

During a 19 year period, we expect $(19)(.341) = 6.5$

$$\text{Therefore, } P(s) = e^{-6.5} (6.5)^s / s! \text{ and } \sum_{s=15}^{\infty} e^{-6.5} (6.5)^s / s! = .007$$

Probability of observing 15 or more leukemias = .007.

BIOS 660/BIOS 672 (3 Credits)

Notes 13 – 15 / 56

continued

Note: For this example, we are glossing over the fact here that we are really combining several different populations that have different Poisson rates. We will see that if X_1 and X_2 are two independent Poisson random variables with parameters λ_1 and λ_2 , then their sum is Poisson with parameter $\lambda_1 + \lambda_2$.

More Negative Binomial

17 / 56

Negative Binomial

In the context of iid Bernoulli trials, define a random variable corresponding to the number of trials required to have s successes. We say $Y \sim \text{Negbin}(s, p)$:

sample space: $\{s, (s + 1), \dots\}$

pmf: for $y = s, s + 1, s + 2, \dots$,

$$f(y) = \binom{y-1}{s-1} p^s q^{y-s}$$

Why? If y is number of trials, then first $(y - 1)$ trials resulted in $(s - 1)$ successes and the last trial is success.

cdf: no closed form.

expectation: $E(Y) = s/p$

Variance: $\text{Var}(Y) = s(1 - p)/p^2$

Recall that the negative binomial is the sum of s independent geometrics with parameter p such that these follow immediately.

Example

The Red Sox and the Atlanta Braves are playing in the world series. The winning team is the first one to win 4 games. Suppose each game is independent of the others, and that the Red Sox win a game with probability p . What is the probability that the Red Sox win?

$$\begin{aligned} P(\text{Red Sox wins}) = & \binom{3}{3} p^4 + \\ & \binom{4}{3} p^4 (1-p) + \\ & \binom{5}{3} p^4 (1-p)^2 + \\ & \binom{6}{3} p^4 (1-p)^3 \end{aligned}$$

BIOS 660/BIOS 672 (3 Credits)

Notes 13 – 19 / 56

Example (continued)

What is the probability that the series goes to 7 games?

If the Red Sox wins, then we need to have three Braves wins and three Red Sox wins (in any order) followed by a Red Sox win:

$$\binom{6}{3} p^4 (1-p)^3.$$

Similarly, the probability that the Braves wins in 7 games is

$$\binom{6}{3} p^3 (1-p)^4.$$

Hence the total probability is

$$\binom{6}{3} [p^4 (1-p)^3 + p^3 (1-p)^4].$$

BIOS 660/BIOS 672 (3 Credits)

Notes 13 – 20 / 56

Comments

- One of the most important things to get out of this class is to understand the different distributions and when/where you would choose to use them.
- You need to be very familiar with all of these distributions as well as simple transformations of these distributions e.g. what happens if you scale an exponential? what about a gamma?
- You will also need to be very familiar with what happens if there are multiple random variables (next class) and need to understand what happens if they are combined (transformed), e.g. what happens if you add Poissons, normals, exponentials? what happens if you take the ratio of a binomial to a Poisson?

BIOS 660/BIOS 672 (3 Credits)

Notes 13 – 21 / 56

5

Exponential Families

22 / 56

Exponential Families

A family of pdfs or pmfs with vector parameter θ is called an *exponential family* if it can be expressed as

$$f(x|\theta) = h(x) c(\theta) \exp \left(\sum_{j=1}^k w_j(\theta) t_j(x) \right), \quad x \in S \subset \mathbb{R}$$

where

- S is not defined in terms of θ (i.e., **the support of the distribution does not depend on the unknown parameter**)
- $h(x), c(\theta) \geq 0$ and the functions are just functions of the parameters specified; i.e. h is free of θ , $c(\theta)$ is free of x , etc...

BIOS 660/BIOS 672 (3 Credits)

Notes 13 – 23 / 56

Example: Binomial

Let $X \sim \text{Binom}(n, p)$, $0 < p < 1$.

$$\begin{aligned} f(x|p) &= \binom{n}{x} p^x (1-p)^{n-x} = \binom{n}{x} (1-p)^n \left[\frac{p}{1-p} \right]^x \\ &= \binom{n}{x} (1-p)^n \exp \left[\log \left(\frac{p}{1-p} \right) x \right] \end{aligned}$$

Thus,

$$\begin{aligned} h(x) &= \binom{n}{x}, \quad x = 0, \dots, n & w_1(p) &= \log \left(\frac{p}{1-p} \right) \\ c(p) &= (1-p)^n, \quad 0 < p < 1 & t_1(x) &= x \end{aligned}$$

Note that this works when p is considered the parameter, while n is fixed. If n is not fixed, then the support depends on the unknown parameter. Also, p cannot be 0 or 1.

Example: Gaussian

Let $X \sim N(\mu, \sigma^2)$.

$$\begin{aligned} f_X(x) &= \frac{1}{\sqrt{2\pi}\sigma} \exp \left(-\frac{(x-\mu)^2}{2\sigma^2} \right) \\ &= \frac{1}{\sqrt{2\pi}\sigma} \exp \left(-\frac{\mu^2}{2\sigma^2} \right) \exp \left(-\frac{x^2}{2\sigma^2} + \frac{\mu x}{\sigma^2} \right) \end{aligned}$$

Thus

$$\begin{aligned} h(x) &= \frac{1}{\sqrt{2\pi}} & c(\mu, \sigma) &= \frac{1}{\sigma} \exp \left(-\frac{\mu^2}{2\sigma^2} \right) \\ w_1(\mu, \sigma) &= -\frac{1}{2\sigma^2} & w_2(\mu, \sigma) &= \frac{\mu}{\sigma^2} \\ t_1(x) &= x^2 & t_2(x) &= x \end{aligned}$$

The parameter space is $(\mu, \sigma^2) \in \mathbb{R} \times (0, \infty)$.

Theorem C-B 3.4.2

If X is a rv from the exponential family, then

$$\begin{aligned} \mathbb{E} \left(\sum_{i=1}^k \frac{\partial w_i(\boldsymbol{\theta})}{\partial \theta_j} t_i(X) \right) &= -\frac{\partial}{\partial \theta_j} \log c(\boldsymbol{\theta}) \\ \text{Var} \left(\sum_{i=1}^k \frac{\partial w_i(\boldsymbol{\theta})}{\partial \theta_j} t_i(X) \right) &= -\frac{\partial^2}{\partial \theta_j^2} \log c(\boldsymbol{\theta}) \\ &\quad - \mathbb{E} \left(\sum_{i=1}^k \frac{\partial^2 w_i(\boldsymbol{\theta})}{\partial \theta_j^2} t_i(X) \right) \end{aligned}$$

Proof: Homework

BIOS 660/BIOS 672 (3 Credits)

Notes 13 – 26 / 56

Example: Binomial

Recall:

$$w_1(p) = \log \frac{p}{1-p}, \quad c(p) = (1-p)^n$$

Relevant derivatives:

$$\begin{aligned} \frac{\partial}{\partial p} w_1(p) &= \frac{\partial}{\partial p} \log \frac{p}{1-p} = \frac{1}{p(1-p)} \\ \frac{\partial}{\partial p} \log c(p) &= \frac{\partial}{\partial p} n \log(1-p) = \frac{-n}{1-p} \end{aligned}$$

So

$$\mathbb{E} \left[\frac{1}{p(1-p)} X \right] = \frac{n}{1-p} \Rightarrow \mathbb{E}(X) = np$$

BIOS 660/BIOS 672 (3 Credits)

Notes 13 – 27 / 56

Indicator function

Definition 3.4.5. The *indicator function* of a set A , most often denoted by $I_A(x)$, is the function

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

Also denoted as $I(x \in A)$, $1_A(x)$, or $1(x \in A)$.

Note on exponential family:

- The set of x values for which $f(x|\theta) > 0$ cannot depend on θ in an exponential family.
- The entire definition of the pdf or pmf must be incorporated into the form for the exponential family.
- Incorporate the range of x into the expression for $f(x|\theta)$ through the use of an indicator function.

Example. Normal pdf:

$$f_X(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right) I_{(-\infty, \infty)}(x).$$

More examples

Other exponential families are: Poisson, negative binomial, gamma, beta.

Some densities that are not exponential families: t , F .

Uniform: $X \sim U(0, \theta)$

$$f_X(x) = \theta^{-1} I(0 < x < \theta)$$

Truncated exponential:

$$f_X(x) = \theta^{-1} \exp(1 - x/\theta) I(\theta, \infty)$$

What about $X \sim U(0, 1)$?

Natural parameters

An exponential family can be reparametrized as

$$f(x|\boldsymbol{\eta}) = h(x) c^*(\boldsymbol{\eta}) \exp \left(\sum_{j=1}^k \eta_j t_j(x) \right), \quad x \in S \subset \mathbb{R}$$

where $\boldsymbol{\eta}$ is called the natural parameter vector.

This parametrization is often more useful. We have the following property:

$$\begin{aligned} \mathbb{E} [t_j(X)] &= -\frac{\partial}{\partial \eta_j} \log c^*(\boldsymbol{\eta}) \\ \text{Var} [t_j(X)] &= -\frac{\partial^2}{\partial \eta_j^2} \log c^*(\boldsymbol{\eta}) \end{aligned}$$

Proof: Homework

Example: Gaussian

Let $X \sim N(\mu, \sigma^2)$. Recall:

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

and

$$w_1(\mu, \sigma) = -\frac{1}{2\sigma^2}, \quad w_2(\mu, \sigma) = \frac{\mu}{\sigma^2}$$

Define

$$\eta_1 = -\frac{1}{2\sigma^2} < 0, \quad \eta_2 = \frac{\mu}{\sigma^2} \in \mathbb{R}$$

then

$$\sigma^2 = -\frac{1}{2\eta_1} > 0, \quad \mu = -\frac{\eta_2}{2\eta_1} \in \mathbb{R}$$

The parameter space is now $(\eta_1, \eta_2) \in (-\infty, 0) \times \mathbb{R}$.

Curved exponential families

Let d be the dimension of the parameter space of the exponential family with k terms. The exp. family is called

$$\begin{array}{ll} \text{full} & \text{if } d = k \\ \text{curved} & \text{if } d < k \end{array}$$

Example: $X \sim N(\mu, \sigma^2)$. Suppose $\sigma^2 = \mu^2$, i.e. the coefficient of variation is constant equal to 1. The parameter space $(\mu, \sigma^2) = (\mu, \mu^2)$ is now a parabola. For the natural parameters:

$$\eta_1 = -\frac{1}{2\mu^2}, \quad \eta_2 = \frac{1}{\mu} \quad \Rightarrow \quad \eta_1 = -\frac{\eta_2^2}{2}$$

Example: Let X_1, \dots, X_n be an iid sample from $Po(\lambda)$. Let $\bar{X} = \sum_{i=1}^n X_i/n$. Then for large n (by the Central Limit Theorem (you will learn this in bios 661/673)),

$$\frac{\bar{X} - \lambda}{\sqrt{\lambda}} \Rightarrow N(0, 1) \quad \text{so} \quad X \dot{\sim} N(\lambda, \lambda)$$

This is a curved exp. family with parameter space $(\mu, \sigma^2) = (\mu, \mu)$.

BIOS 660/BIOS 672 (3 Credits)

Notes 13 – 32 / 56

Probability Inequalities

33 / 56

Chebychev Inequality

Let X be a random variable and let $g(x)$ be a non-negative function. Then for any $r > 0$,

$$P[g(X) \geq r] \leq \frac{Eg(X)}{r}$$

Proof:

$$\begin{aligned} Eg(X) &= \int_{-\infty}^{\infty} g(x) f_X(x) dx \\ &\geq \int_{\{x: g(x) \geq r\}} g(x) f_X(x) dx \\ &\geq r \int_{\{x: g(x) \geq r\}} f_X(x) dx \\ &= r P\{g(X) \geq r\} \end{aligned}$$

BIOS 660/BIOS 672 (3 Credits)

Notes 13 – 34 / 56

Application

Let

$$g(x) = \frac{(x - \mu)^2}{\sigma^2}$$

where $\mu = E(X)$, $\sigma^2 = \text{Var}(X)$.

Let $r = t^2$, then

$$P\left[\frac{(x - \mu)^2}{\sigma^2} \geq t^2\right] \leq \frac{1}{t^2} E\left[\frac{(x - \mu)^2}{\sigma^2}\right]$$
$$P[|X - \mu| \geq t\sigma] \leq \frac{1}{t^2}$$

The probability that X is more than $t\sigma$ away from μ cannot be more than $1/t^2$, no matter what the distribution of X . E.g. $t = 2$.

Normal tail bound

Let $Z \sim N(0, 1)$:

$$P[Z \geq t] = \frac{1}{\sqrt{2\pi}} \int_t^\infty e^{-x^2/2} dx$$
$$\leq \frac{1}{\sqrt{2\pi}} \int_t^\infty \frac{x}{t} e^{-x^2/2} dx = \frac{1}{\sqrt{2\pi}} \frac{e^{-t^2/2}}{t}$$

and so

$$P[|Z| \geq t] = 2P[Z \geq t] \leq \sqrt{\frac{2}{\pi}} \frac{e^{-t^2/2}}{t}$$

For $t = 2$, the bound is 0.054.

Multiple random measurements

Multiple endpoints in health studies:

- Cancer: survival, quality of life, toxicity
- Reproductive health: time to pregnancy, birth defects
- AIDS: time from infection to AIDS, time from AIDS to death
- Carcinogenicity studies: time to tumor, time to death
- Health care: cost, hospital duration of stay

Multiple time points/spatial locations:

- Environmental monitoring: daily temperature, humidity, CO₂/ozone concentration, geographically located monitoring stations.
- Finance: daily stock prices, portfolios.

Massively multivariate:

- Genomics and other omics: thousands of gene expression values, SNPs, protein concentrations, metabolite concentrations.
- Imaging: thousands of voxels (volume pixels) measuring brain activity (fMRI), tumor metabolism (PET); satellite remote sensing.

BIOS 660/BIOS 672 (3 Credits)

Notes 13 – 38 / 56

Random Vectors

Suppose we start with a probability space (Ω, \mathcal{A}, P) .

Definition: An n -dimensional random vector $\mathbf{X} = (X_1, \dots, X_n)$ is a function from a sample space Ω into \mathbb{R}^n .

- Each coordinate X_i is a random variable.
- The random vector is associated with a probability space $(\mathbb{R}^n, \mathcal{B}(\mathbb{R}^n), F)$.
- For every Borel set B ,

$$P\{\mathbf{X} \in B\} = P\{\mathbf{X}^{-1}(B)\}$$

where

$$\mathbf{X}^{-1}(B) = \{\omega : \mathbf{X}(\omega) \in B\}$$

BIOS 660/BIOS 672 (3 Credits)

Notes 13 – 39 / 56

Example: Bivariate

A fair coin is flipped 3 times. Define the random vector (X, Y) where X represents the number of heads on the last toss and Y the total number of heads. Then, the probabilities of various outcomes are given in the following table:

Outcome	(x, y)	$P(\text{outcome})$
(H,H,H)	(1,3)	1/8
(H,H,T)	(0,2)	1/8
(H,T,H)	(1,2)	1/8
(H,T,T)	(0,1)	1/8
(T,H,H)	(1,2)	1/8
(T,H,T)	(0,1)	1/8
(T,T,H)	(1,1)	1/8
(T,T,T)	(0,0)	1/8

Discrete Bivariate RVs

Two random variables X and Y are said to be jointly *discrete* if there is an associated *joint probability mass function*,

$$f_{X,Y}(x, y) = P\{X = x, Y = y\}$$

which sums to 1 over a finite or possibly countable combinations of x and y for which $f_{X,Y}(x, y) > 0$, i.e.,

$$\sum_{x,y} f_{X,Y}(x, y) = 1$$

From this, one can also obtain the marginal pmfs of X and Y as follows:

$$f_X(x) = P(X = x) = \sum_y f_{X,Y}(x, y)$$

$$f_Y(y) = P(Y = y) = \sum_x f_{X,Y}(x, y)$$

Example

Back to the fair coin example again. From the definition we can construct the joint pdf of X and Y :

		Y			
		0	1	2	3
X	0	1/8	1/4	1/8	0
	1	0	1/8	1/4	1/8

The marginal distributions of X and Y are also easy to find.

Note: Marginals do not determine joint pmf.

Bivariate cdfs

Regardless of whether they are discrete or continuous or some combination of the two, we can always define the *joint cumulative distribution function*.

For $n = 2$, the *bivariate cumulative distribution function* is

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

Properties:

- $F_{X,Y}(x, y) \geq 0$
- $F_{X,Y}(\infty, \infty) = 1$
- $F_{X,Y}(-\infty, y) = F(x, -\infty) = 0$
- $F_{X,Y}(-\infty, -\infty) = 0$
- F is non-decreasing and right-continuous in each variable separately

Marginal distributions

From $F_{X,Y}$, we can derive the univariate distribution functions for X and Y . These are generally called *marginal distributions*.

$$F_X(x) = P\{X \leq x\} = P\{X \leq x, Y < \infty\} = F_{X,Y}(x, \infty)$$

$$F_Y(y) = P\{Y \leq y\} = P\{X < \infty, Y \leq y\} = F_{X,Y}(\infty, y)$$

Note: Although we can obtain $F_X(x)$ and $F_Y(y)$ from the joint *cdf*, we cannot do the reverse.

Joint probabilities

All joint probability statements about X and Y can be answered in terms of their joint *cdf*. For example,

$$P(X > x, Y > y) = 1 - F_X(x) - F_Y(y) + F_{X,Y}(x, y)$$

More generally,

$$P(x_1 < X \leq x_2, y_1 < Y \leq y_2) = \\ F_{X,Y}(x_2, y_2) + F_{X,Y}(x_1, y_1) - F_{X,Y}(x_1, y_2) - F_{X,Y}(x_2, y_1)$$

Continuous Bivariate Random Variables

46 / 56

Continuous Bivariate RVs

The random variables X and Y are said to be *jointly (absolutely) continuous* if there exists a function $f_{X,Y}(x, y)$, such that for any Borel set B of 2-tuples in \mathbb{R}^2 ,

$$P\{(X, Y) \in B\} = \int \int_{(x,y) \in B} f_{X,Y}(x, y) dx dy$$

The function $f_{X,Y}(x, y)$ is called the *joint probability density function* for X and Y .

It follows in this case that

$$F_{X,Y}(x, y) = \int_{-\infty}^x \int_{-\infty}^y f_{X,Y}(s, t) ds dt,$$

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

BIOS 660/BIOS 672 (3 Credits)

Notes 13 – 47 / 56

Properties of the bivariate pdf

- $f_{X,Y}(x, y) \geq 0$
- $\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f_{X,Y}(x, y) dx dy = 1$
- $f_{X,Y}(x, y)$ is **not a probability**, but can be thought of as a relative probability of (X, Y) falling into a small rectangle located at (x, y) :

$$P\{x < X \leq x + dx, y < Y \leq y + dy\} \approx f(x, y) dx dy$$

- The *marginal probability density functions* for X and Y can be obtained as

$$f_X(x) = \int_{-\infty}^{\infty} f_{X,Y}(x, y) dy$$

$$f_Y(y) = \int_{-\infty}^{\infty} f_{X,Y}(x, y) dx$$

BIOS 660/BIOS 672 (3 Credits)

Notes 13 – 48 / 56

Example 1

$$F_{XY}(x, y) = xy \quad 0 < x \leq 1, \quad 0 < y \leq 1$$

$$f_{XY}(x, y) = \frac{\partial^2 F_{XY}(x, y)}{\partial x \partial y} = 1$$

$$f_X(x) = \int_0^1 dy = 1$$

$$f_Y(y) = \int_0^1 dx = 1$$

BIOS 660/BIOS 672 (3 Credits)

Notes 13 – 49 / 56

Example 2

$$F_{XY}(x, y) = x - x \log\left(\frac{x}{y}\right) \quad 0 < x \leq y \leq 1$$

$$f_{XY}(x, y) = \frac{\partial^2 F_{XY}}{\partial x \partial y} = \frac{\partial}{\partial x} \left[-x \left(\frac{y}{x} \right) \left(-\frac{x}{y^2} \right) \right] = \frac{\partial}{\partial x} \frac{x}{y} = \frac{1}{y}$$

$$f_X(x) = \int_x^1 \frac{dy}{y} = -\log(x)$$

$$f_Y(y) = \int_0^y \frac{dx}{y} = 1$$

Note: Once we have $f_X(y)$ and $f_Y(y)$, we can obtain $F_X(x)$ and $F_Y(y)$ directly.

Double check:

$$\begin{aligned} F_X(x) &= F_{X,Y}(x, 1) = x - x \log(x); \\ &\quad \frac{d}{dx} [x - x \log(x)] = -\log(x). \\ F_Y(y) &= F_{X,Y}(y, y) = y; \\ &\quad \frac{d}{dy} y = 1. \end{aligned}$$

BIOS 660/BIOS 672 (3 Credits)

Notes 13 – 50 / 56

Conditional Distributions - Discrete

Recall if A and B are two events, the probability of A conditional on B is:

$$P(A|B) = \frac{P(A \cap B)}{P(B)} \quad \text{or} \quad \frac{P(AB)}{P(B)}$$

Defining the events $A = \{Y = y\}$ and $B = \{X = x\}$, it follows that

$$\begin{aligned} P\{Y = y|X = x\} &= \frac{P(X = x, Y = y)}{P(X = x)} \\ &= \frac{f_{X,Y}(x, y)}{f_X(x)} \\ &= f_{Y|X}(y|x) \end{aligned}$$

This is called the **conditional probability mass function** of Y given X .

BIOS 660/BIOS 672 (3 Credits)

Notes 13 – 52 / 56

Example: Discrete

A fair coin is flipped 3 times. Define the random vector (X, Y) where X represents the number of heads on the last toss and Y the total number of heads. From the joint pmf of X and Y we can derive all the conditional pmfs: $f_{Y|X}(y|x) = f_{X,Y}(x, y)/f_X(x)$ and $f_{X|Y}(x|y) = f_{X,Y}(x, y)/f_Y(y)$.

Examples: $f_{Y|X}(0|0) = f_{X,Y}(0, 0)/f_X(0) = \frac{1/8}{1/2} = 1/4$;

$f_{X|Y}(1|2) = f_{X,Y}(1, 2)/f_Y(2) = \frac{1/4}{3/8} = 2/3$.

		Y				Sum
		0	1	2	3	
X	0	1/8	1/4	1/8	0	1/2
	1	0	1/8	1/4	1/8	1/2
	Sum	1/8	3/8	3/8	1/8	1

BIOS 660/BIOS 672 (3 Credits)

Notes 13 – 53 / 56

Conditional Distributions - Continuous

If $F(x, y)$ is absolutely continuous, we define the conditional density of Y given X as:

$$f_{X|Y}(x|y) = \frac{f_{X,Y}(x, y)}{f_Y(y)} \quad \text{if } f_Y(y) > 0$$

Because Y is continuous, we cannot directly evaluate this probability, since the denominator will be zero. Instead, think of small dx, dy :

$$\begin{aligned} & Pr(x \leq X < x + dx \mid y \leq Y < y + dy) \\ &= \frac{Pr(x \leq X < x + dx, y \leq Y < y + dy)}{Pr(y \leq Y < y + dy)} \\ &\approx \frac{f(x, y)dx dy}{f_Y(y)dy} \\ &= f_{X|Y}(x|y)dx \end{aligned}$$

Show that it satisfies the conditions for a density.

Example 1

$$F_{XY}(x, y) = xy \quad 0 < x \leq 1, \quad 0 < y \leq 1$$

$$f_{XY}(x, y) = 1 \quad 0 < x < 1, \quad 0 < y < 1$$

$$f_X(x) = 1 \quad 0 < x < 1$$

$$f_Y(y) = 1 \quad 0 < y < 1$$

$$f_{X|Y}(x|y) = \frac{f_{XY}(x, y)}{f_Y(y)} = 1 \quad 0 < x < 1 \quad (0 < y < 1)$$

$$f_{Y|X}(y|x) = \frac{f_{XY}(x, y)}{f_X(x)} = 1 \quad 0 < y < 1 \quad (0 < x < 1)$$

In this particular case, we get that the conditional densities are the same as the marginals. This means X and Y are independent.

Example 2

$$F_{XY}(x, y) = x - x \log \frac{x}{y} \quad 0 < x \leq y \leq 1$$

$$f_{XY}(x, y) = 1/y \quad 0 < x \leq y \leq 1$$

$$f_X(x) = -\log x \quad 0 < x \leq 1$$

$$f_Y(y) = 1 \quad 0 < y \leq 1$$

$$f_{X|Y}(x|y) = \frac{f_{XY}(x, y)}{f_Y(y)} = 1/y \quad 0 < x \leq y \quad (0 < y \leq 1)$$

$$f_{Y|X}(y|x) = \frac{f_{XY}(x, y)}{f_X(x)} = -\frac{1}{y \log x} \quad x \leq y \leq 1 \quad (0 < x \leq 1)$$

- Y is marginally uniform, but not conditionally
- X is conditionally uniform, but not marginally

BIOS 660/BIOS 672 (3 Credits): Probability and Statistical Inference I

Jianwen Cai

<https://sakai.unc.edu/portal/site/bios660-bios672-3-credits>

Notes 14

Independent Random Variables	2
Independence	3
Checking independence	4
Example	5
Example: Buffon's Needle	6
Example (cont.)	7
Expectations of Independent RVs	8
Proof	9
Bivariate Transformations	10
Functions of random vectors	11
Discrete RVs	12
Sum of two independent Poissons	13
cont.	14
Bivariate Transformations of Continuous RVs	15
Continuous RVs	16
Rotation of a bivariate normal vector	17
cont.	18
Functions of independent random variables	19
Extensions of previous example	20
Ratio of two independent normals	21
cont.	22
Sum of Two Independent RVs	23
Sum of Two Independent RVs	24
Sum of two independent Poissons	25
Characteristic Function	26
Sum of two independent Poissons	27
Sum of Two Independent Normals	28
Sum of two independent gammas	29

Independence

The random variables X and Y are said to be *independent* if for any two Borel sets A and B ,

$$P(X \in A, Y \in B) = P(X \in A) P(Y \in B)$$

All events defined in terms of X are independent of all events defined in terms of Y .

Using the Kolmogorov axioms of probability, it can be shown that X and Y are independent if and only if $\forall(x, y)$ (*except possibly for sets of prob. 0*)

$$F_{X,Y}(x, y) = F_X(x)F_Y(y)$$

or in terms of *pmfs* (discrete) and *pdf's* (continuous)

$$f_{X,Y}(x, y) = f_X(x)f_Y(y)$$

* Check previous examples.

BIOS 660/BIOS 672 (3 Credits)

Notes 14 – 3 / 29

Checking independence

1. A necessary condition for independence of X and Y is that their joint pdf/pmf has positive probability on a rectangular domain.
2. If the domain is rectangular, one can try to write the joint pdf/pmf as a product of functions of x and y only.

Lemma C-B 4.2.7: Let (X, Y) be a bivariate random vector with joint pdf or pmf $f(x, y)$. Then X and Y are independent if and only if there exist functions $g(x)$ and $h(y)$ such that for every $x \in \mathbb{R}$ and $y \in \mathbb{R}$,

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

Proof:

BIOS 660/BIOS 672 (3 Credits)

Notes 14 – 4 / 29

Example

Two points are selected randomly on a line of length a so as to be on opposite sides of the mid-point of the line. Find the probability that the distance between them is less than $a/3$.

Solution:

Let X be the coordinate of a point selected randomly in $[0, a/2]$ and Y the coordinate of a point selected randomly in $[a/2, a]$. Assume X and Y are independent and uniform over its interval. The joint density is:

$$f_{X,Y}(x, y) = 4/a^2, \quad 0 \leq x \leq a/2, \quad a/2 \leq y \leq a$$

Hence, the solution is

$$P(Y - X < a/3) = \int_{a/6}^{a/2} \int_{a/2}^{a/3+x} \frac{4}{a^2} dy dx = 2/9$$

Example: Buffon's Needle

A table is ruled with lines distance 1 unit apart. A needle of length $L \leq 1$ is thrown randomly on the table. What is the probability that the needle intersects a line?

Solution: Define two random variables:

- X : distance from low end of the needle to the nearest line above
- θ : angle from the vertical to the needle.

Example (cont.)

By “random”, we assume X and θ are independent, and

$$X \sim U(0, 1) \quad \text{and} \quad \theta \sim U[-\pi/2, \pi/2].$$

This means that

$$f_{X,\Theta}(x, \theta) = 1/\pi, \quad 0 \leq x \leq 1, \quad -\pi/2 \leq \theta \leq \pi/2$$

For the needle to intersect a line, we need $X < L \cos(\theta)$. So,

$$\begin{aligned} P(\text{needle intersects a line}) &= P(X < L \cos(\theta)) \\ &= \int_{-\pi/2}^{\pi/2} \int_0^{L \cos \theta} \frac{1}{\pi} dx d\theta \\ &= \frac{2L}{\pi} \end{aligned}$$

Expectations of Independent RVs

Theorem C-B 4.2.10 Let X and Y be independent rvs.

- For any $A \subset \mathcal{R}$ and $B \subset \mathcal{R}$,

$$P(X \in A, Y \in B) = P(X \in A)P(Y \in B)$$

i.e. the events $\{X \in A\}$ and $\{Y \in B\}$ are independent.
(C-B write it as a theorem, we took it as definition.)

- Let $g(x)$ be a function only of x and $h(y)$ be a function only of y . Then

$$E(g(X)h(Y)) = (Eg(X))(Eh(Y))$$

Example: X, Y indep.

$$E(X^2Y^3) = (EX^2)(EY^3)$$

$$E(Y^2Y^3) \neq (EY^2)(EY^3)$$

Proof

$$\begin{aligned} E(g(X)h(Y)) &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} g(x)h(y)f_{X,Y}(x,y)dx dy \\ &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} g(x)h(y)f_X(x)f_Y(y)dx dy \\ &= \left(\int_{-\infty}^{\infty} g(x)f_X(x)dx \right) \left(\int_{-\infty}^{\infty} h(y)f_Y(y)dy \right) \\ &= (E g(X))(E h(Y)) \end{aligned}$$

BIOS 660/BIOS 672 (3 Credits)

Notes 14 – 9 / 29

Bivariate Transformations

10 / 29

Functions of random vectors

Let (X, Y) be a bivariate rv with known distribution. Define (U, V) by

$$U = g_1(X, Y), \quad V = g_2(X, Y)$$

Probability mapping: For any Borel set $B \subset \mathbb{R}^2$,

$$P[(U, V) \in B] = P[(X, Y) \in A]$$

where A is the inverse mapping of B , i.e.

$$A = \{(x, y) \in \mathbb{R}^2 : (g_1(x, y), g_2(x, y)) \in B\}$$

The inverse is well defined even if the mapping is not bijective.

Example: Let $g_1(x, y) = x$, $g_2(x, y) = x^2 + y^2$.

BIOS 660/BIOS 672 (3 Credits)

Notes 14 – 11 / 29

Discrete RVs

Suppose that (X, Y) is a discrete rv, i.e. the pmf is positive on a countable set \mathcal{A} . Then (U, V) is also discrete and takes values on a countable set \mathcal{B} . Define

$$A_{uv} = \{(x, y) \in \mathcal{A} : g_1(x, y) = u, g_2(x, y) = v\}$$

Then

$$f_{U,V}(u, v) = P(U = u, V = v) = \sum_{(x,y) \in A_{uv}} f_{X,Y}(x, y)$$

BIOS 660/BIOS 672 (3 Credits)

Notes 14 – 12 / 29

Sum of two independent Poissons

Let $X \sim Po(\lambda_1)$, $Y \sim Po(\lambda_2)$, independent, and define

$$U = X + Y, \quad V = Y$$

- (X, Y) takes values in $\mathcal{A} = \{0, 1, 2, \dots\}^2$.
- (U, V) takes values on $\mathcal{B} = \{(u, v) : v = 0, 1, 2, \dots, u = v, v + 1, v + 2, \dots\}$.
- For a particular (u, v) , $A_{uv} = \{(x, y) \in \mathcal{A} : x + y = u, y = v\} = (u - v, u)$.

The joint pmf of U and V is

$$f_{U,V}(u, v) = f_{X,Y}(u - v, v) = \frac{e^{-\lambda_1} \lambda_1^{u-v}}{(u - v)!} \frac{e^{-\lambda_2} \lambda_2^v}{v!}$$

BIOS 660/BIOS 672 (3 Credits)

Notes 14 – 13 / 29

cont.

The distribution of $U = X + Y$ is the marginal

$$\begin{aligned} f_U(u) &= \sum_{v=0}^u \frac{e^{-\lambda_1} \lambda_1^{u-v}}{(u-v)!} \frac{e^{-\lambda_2} \lambda_2^v}{v!} \\ &= \frac{e^{-(\lambda_1+\lambda_2)}}{u!} \sum_{v=0}^u \binom{u}{v} \lambda_1^{u-v} \lambda_2^v \\ &= \frac{e^{-(\lambda_1+\lambda_2)}}{u!} (\lambda_1 + \lambda_2)^u \end{aligned}$$

We obtain that U is Poisson with parameter $\lambda = \lambda_1 + \lambda_2$.

Bivariate Transformations of Continuous RVs

15 / 29

Continuous RVs

Suppose (X, Y) is continuous and the joint transformation

$$u = g_1(x, y), \quad v = g_2(x, y)$$

is one-to-one and differentiable. Define the inverse mapping

$$x = h_1(u, v), \quad y = h_2(u, v)$$

Then

$$f_{UV}(u, v) = f_{XY}(h_1(u, v), h_2(u, v)) |J(u, v)|$$

where $J(u, v)$ is the Jacobian of the transformation $(x, y) \rightarrow (u, v)$ given by

$$J(u, v) = \frac{\partial(x, y)}{\partial(u, v)} = \begin{vmatrix} \frac{\partial x}{\partial u} & \frac{\partial x}{\partial v} \\ \frac{\partial y}{\partial u} & \frac{\partial y}{\partial v} \end{vmatrix}$$

Rotation of a bivariate normal vector

Let $X \sim N(0, 1)$, $Y \sim N(0, 1)$, independent. Define the rotation

$$U = X \cos \theta - Y \sin \theta$$

$$V = X \sin \theta + Y \cos \theta$$

for fixed θ . Then $U \sim N(0, 1)$, $V \sim N(0, 1)$, independent.

Proof: The range of (X, Y) is \mathbb{R}^2 . The range of (U, V) is \mathbb{R}^2 . Need the inverse transformation

$$X = U \cos \theta + V \sin \theta$$

$$Y = -U \sin \theta + V \cos \theta$$

with Jacobian

$$J(u, v) = \begin{vmatrix} \frac{\partial x}{\partial u} & \frac{\partial x}{\partial v} \\ \frac{\partial y}{\partial u} & \frac{\partial y}{\partial v} \end{vmatrix} = \begin{vmatrix} \cos \theta & \sin \theta \\ -\sin \theta & \cos \theta \end{vmatrix} = 1$$

cont.

The joint pdf of (X, Y) is

$$f_{XY}(x, y) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2} \cdot \frac{1}{\sqrt{2\pi}} e^{-y^2/2} = \frac{1}{2\pi} e^{-(x^2+y^2)/2}$$

The joint pdf of (U, V) is

$$\begin{aligned} f_{UV}(u, v) &= \frac{1}{2\pi} e^{-[(u \cos \theta + v \sin \theta)^2 + (-u \sin \theta + v \cos \theta)^2]/2} \cdot |1| \\ &= \frac{1}{2\pi} e^{-(u^2+v^2)/2} = \frac{1}{\sqrt{2\pi}} e^{-u^2/2} \cdot \frac{1}{\sqrt{2\pi}} e^{-v^2/2} \end{aligned}$$

so $U \sim N(0, 1)$, $V \sim N(0, 1)$, independent.

Functions of independent random variables

Theorem C-B 4.3.5: Let X and Y be independent rvs. Let $g : \mathbb{R} \rightarrow \mathbb{R}$ and $h : \mathbb{R} \rightarrow \mathbb{R}$ be functions. Then the random variables $U = g(X)$ and $V = h(Y)$ are independent.

Proof:

Extensions of previous example

- Suppose $X \sim N(0, \sigma^2)$, $Y \sim N(0, \sigma^2)$, independent

$$U = a(X \cos \theta - Y \sin \theta)$$

$$V = a(X \sin \theta + Y \cos \theta)$$

Then $U \sim N(0, a^2 \sigma^2)$, $V \sim N(0, a^2 \sigma^2)$, independent.

- Above, take $\theta = \pi/4$, $a = \sqrt{2}$:

$$U = \sqrt{2}(X/\sqrt{2} - Y/\sqrt{2}) = X - Y$$

$$V = \sqrt{2}(X/\sqrt{2} + Y/\sqrt{2}) = X + Y$$

We get $U \sim N(0, 2\sigma^2)$, $V \sim N(0, 2\sigma^2)$, independent.

Ratio of two independent normals

Let $X \sim N(0, 1)$, $Y \sim N(0, 1)$, independent. The ratio X/Y has the Cauchy distribution.

Proof: Define the variables

$$U = X/Y, \quad V = Y$$

with inverse

$$X = UV, \quad Y = V$$

The Jacobian is

$$J(u, v) = \begin{vmatrix} \frac{\partial x}{\partial u} & \frac{\partial x}{\partial v} \\ \frac{\partial y}{\partial u} & \frac{\partial y}{\partial v} \end{vmatrix} = \begin{vmatrix} v & u \\ 0 & 1 \end{vmatrix} = v$$

The range of (X, Y) is \mathbb{R}^2 . The range of (U, V) is \mathbb{R}^2 .

BIOS 660/BIOS 672 (3 Credits)

Notes 14 – 21 / 29

cont.

The joint pdf of (X, Y) is

$$f_{XY}(x, y) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2} \cdot \frac{1}{\sqrt{2\pi}} e^{-y^2/2} = \frac{1}{2\pi} e^{-(x^2+y^2)/2}$$

The joint pdf of (U, V) is

$$f_{UV}(u, v) = \frac{1}{2\pi} e^{-[(uv)^2+v^2]/2} \cdot |v| = \frac{|v|}{2\pi} e^{-(u^2+1)v^2/2}$$

The marginal of U is

$$\begin{aligned} f_U(u) &= \int_{-\infty}^{\infty} f_{UV}(u, v) dv = 2 \int_0^{\infty} \frac{v}{2\pi} e^{-(u^2+1)v^2/2} dv \\ &= \frac{1}{\pi} \int_0^{\infty} e^{-(u^2+1)z} dz = \frac{1}{\pi(u^2+1)} \end{aligned}$$

BIOS 660/BIOS 672 (3 Credits)

Notes 14 – 22 / 29

Sum of Two Independent RVs

Suppose X and Y are independent. What is the distribution of $Z = X + Y$?

In general:

$$F_Z(z) = P(X + Y \leq z) = P(\{(x, y) \text{ such that } x + y \leq z\})$$

Various approaches:

- Bivariate transformation method (continuous and discrete)
- Discrete convolution

$$f_Z(z) = \sum_{x+y=z} f_X(x)f_Y(y) = \sum_x f_X(x)f_Y(z-x)$$

- Continuous convolution (C-B Section 5.2)
- Mgf/cf method (continuous and discrete)

BIOS 660/BIOS 672 (3 Credits)

Notes 14 – 24 / 29

Sum of two independent Poissons

Define X, Y to be two independent random variables having Poisson distributions with parameters λ_i , $i = 1, 2$. Then:

$$f_{X,Y}(x, y) = \frac{e^{-\lambda_1} \lambda_1^x}{x!} \frac{e^{-\lambda_2} \lambda_2^y}{y!} \quad x, y = 0, 1, 2, \dots$$

The distribution of $S = X + Y$ is

$$\begin{aligned} f_S(s) &= \sum_{x=0}^s \frac{e^{-\lambda_1} \lambda_1^x}{x!} \frac{e^{-\lambda_2} \lambda_2^{s-x}}{(s-x)!} \\ &= \frac{e^{-(\lambda_1+\lambda_2)}}{s!} \sum_{x=0}^s \binom{s}{x} \lambda_1^x \lambda_2^{s-x} \\ &= \frac{e^{-(\lambda_1+\lambda_2)}}{s!} (\lambda_1 + \lambda_2)^s \end{aligned}$$

Thus S is Poisson with parameter $\lambda = \lambda_1 + \lambda_2$.

BIOS 660/BIOS 672 (3 Credits)

Notes 14 – 25 / 29

Characteristic Function

Theorem 4.2.12 Let X and Y be independent rvs with characteristic functions $\phi_X(\cdot)$ and $\phi_Y(\cdot)$, respectively. Then the characteristic function of $Z = X + Y$ is

$$\phi_Z(\theta) = \phi_X(\theta) \phi_Y(\theta)$$

Proof:

$$\begin{aligned}\phi_Z(\theta) &= \mathbb{E} \exp(iZ\theta) &&= \mathbb{E} \exp[i(X + Y)\theta] \\ &= \mathbb{E} \exp(iX\theta) \exp(iY\theta) &&= \mathbb{E} \exp(iX\theta) \mathbb{E} \exp(iY\theta) \\ &= \phi_X(\theta) \phi_Y(\theta)\end{aligned}$$

Corollary If X and Y independent and $Z = X - Y$,

$$\phi_Z(\theta) = \phi_X(\theta) \phi_Y(-\theta)$$

Sum of two independent Poissons

Suppose $X \sim \text{Poisson}(\lambda_X)$ and $Y \sim \text{Poisson}(\lambda_Y)$ and put $Z = X + Y$. Then, $Z \sim \text{Poisson}(\lambda_X + \lambda_Y)$.

Proof:

$$\begin{aligned}\phi_Z(\theta) &= \exp[\lambda_X(e^\theta - 1)] \exp[\lambda_Y(e^\theta - 1)] \\ &= \exp[(\lambda_X + \lambda_Y)(e^\theta - 1)]\end{aligned}$$

Sum of Two Independent Normals

Suppose $X \sim N(\mu_x, \sigma_x^2)$ and $Y \sim N(\mu_y, \sigma_y^2)$ and X and Y are independent and $Z = X + Y$ then

$$Z \sim N(\mu_x + \mu_y, \sigma_x^2 + \sigma_y^2)$$

Proof

$$\begin{aligned}\phi_Z(\theta) &= \exp\left(i\mu_x\theta - \frac{1}{2}\sigma_x^2\theta^2\right) \exp\left(i\mu_y\theta - \frac{1}{2}\sigma_y^2\theta^2\right) \\ &= \exp\left[i(\mu_x + \mu_y)\theta - \frac{1}{2}(\sigma_x^2 + \sigma_y^2)\theta^2\right]\end{aligned}$$

BIOS 660/BIOS 672 (3 Credits)

Notes 14 – 28 / 29

Sum of two independent gammas

Suppose $X \sim \Gamma(\alpha_x, \beta)$ and independently $Y \sim \Gamma(\alpha_y, \beta)$ and put $Z = X + Y$. Then,
 $Z \sim \Gamma((\alpha_x + \alpha_y), \beta)$

Proof:

$$\begin{aligned}\phi_Z(\theta) &= \left(\frac{1}{1 - \beta\theta}\right)^{\alpha_x} \left(\frac{1}{1 - \beta\theta}\right)^{\alpha_y} \\ &= \left(\frac{1}{1 - \beta\theta}\right)^{\alpha_x + \alpha_y}\end{aligned}$$

Remember that

- If $\alpha = 1$ we have an exponential with parameter β .
- If $\alpha = n/2$ and $\beta = 2$, we have a $\chi^2(n)$ (with n d.f.). The above result states that $\chi^2(n_1) + \chi^2(n_2) = \chi^2(n_1 + n_2)$.

BIOS 660/BIOS 672 (3 Credits)

Notes 14 – 29 / 29

BIOS 660/BIOS 672 (3 Credits): Probability and Statistical Inference I

Jianwen Cai

<https://sakai.unc.edu/portal/site/bios660-bios672-3-credits>

Notes 15

Conditional Expectation and Hierarchical Models	2
Conditional Expectation	3
Example	4
Iterative expectation formula.	5
Example	6
Hierarchical Model	7
cont.	8
Three-layer hierarchical model	9
Conditional Variance.	10
Conditional variance hierarchical formula.	11
cont.	12
Example	13
Covariance and Correlation	14
Definitions.	15
Properties of Covariance	16
Linear Combinations.	17
Correlation Coefficient	18
Bivariate Normal	19
Standard Bivariate Normal	20
Bivariate Normal.	21
Properties.	22
Multivariate Distributions	23
Multivariate Distributions	24
Marginals and Conditionals	25
Multinomial Distribution	26
Multinomial Theorem 4.6.4.	27
Multinomial Distribution: Properties	28
Multivariate Independence	29

Independent Random Vectors	30
cont.	31
Specific Characteristic Functions	32
Multivariate moments and multivariate normal(Gut, Chapter V)	33
Multivariate moments	34
Multivariate covariance.	35
Bivariate normal	36
Linear functions	37
Positive definiteness	38
Multivariate linear transformations	39
Multivariate linear transformations	40
Multivariate normal	41
Construction	42
Construction (cont.)	43
Properties	44

Conditional Expectation

Suppose we have discrete rvs X and Y with conditional pmf $f_{Y|X}(y|x)$. The *conditional expectation* of $g(Y)$ given $X = x$ is

$$E[g(Y)|X = x] = \sum_y g(y) f_{Y|X}(y|x)$$

Notice that this is a function $h(x)$ of x . We may define the random variable

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

In particular, if $X = x$ then $E[g(Y)|X] = E[g(Y)|x]$.

For continuous rvs:

$$h(x) = E[g(Y)|x] = \int_{-\infty}^{\infty} g(y) f_{Y|X}(y|x) dy$$

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

BIOS 660/BIOS 672 (3 Credits)

Notes 15 – 3 / 44

Example

Let $X \sim \text{Exp}(\lambda)$, $Z \sim \text{Exp}(\lambda)$, independent, be the times to mitosis of a cell and its daughter. Define the total time $Y = X + Z$.

Given $X = x > 0$, $Y = x + Z$ is a shifted exponential rv

$$f_{Y|X}(y|x) = \lambda e^{-\lambda(y-x)}, \quad y > x,$$

The conditional expectation of Y given $X = x$ is

$$E[Y|X = x] = \int_x^{\infty} y \lambda e^{-\lambda(y-x)} dy = x + 1/\lambda$$

yielding the rv

$$E[Y|X] = X + 1/\lambda$$

More directly

$$E(Y|X) = E(X + Z|X) = E(X|X) + E(Z|X) = X + 1/\lambda$$

BIOS 660/BIOS 672 (3 Credits)

Notes 15 – 4 / 44

Iterative expectation formula

If X and Y are any two random variables then

$$EX = E(E(X|Y))$$

Proof:

$$\begin{aligned} E(E(X|Y)) &= \int E(X|y)f(y)dy \\ &= \int \left[\int xf(x|y)dx \right] f(y)dy \\ &= \int \int xf(x,y)dxdy \\ &= EX \end{aligned}$$

Notice that the expectations are with respect to different variables and densities. Careful!

Example

Back to the cells example:

$$E(Y|X) = X + 1/\lambda$$

so

$$E[E(Y|X)] = E[X + 1/\lambda] = EX + 1/\lambda = 1/\lambda + 1/\lambda = 2/\lambda$$

This should be equal to EY , which is

$$EY = E(X + Z) = EX + EZ = 1/\lambda + 1/\lambda = 2/\lambda$$

Hierarchical Model

Example: Suppose an insect lays eggs according to a $\text{Poisson}(\lambda)$ and each egg survives with probability p . Assume that the survival of eggs is independent of each other, then what is the average number of eggs surviving?

Let's say $Y \sim \text{Poisson}(\lambda)$ and $X|Y \sim \text{Binomial}(Y, p)$ where X is the total number of eggs surviving.

$$\begin{aligned} P(X = x) &= \sum_{y=0}^{\infty} P(X = x, Y = y) \\ &= \sum_{y=0}^{\infty} P(X = x|Y = y)P(Y = y) \\ &= \sum_{y=x}^{\infty} \left[\binom{y}{x} p^x (1-p)^{y-x} \right] \left[\frac{e^{-\lambda} \lambda^y}{y!} \right] \\ &= \frac{(\lambda p)^x e^{-\lambda}}{x!} \sum_{y=x}^{\infty} \frac{((1-p)\lambda)^{y-x}}{(y-x)!} \end{aligned}$$

BIOS 660/BIOS 672 (3 Credits)

Notes 15 – 7 / 44

cont

$$\begin{aligned} &= \frac{(\lambda p)^x e^{-\lambda}}{x!} \exp((1-p)\lambda) \\ &= \frac{(\lambda p)^x}{x!} \exp(-\lambda p) \end{aligned}$$

So $X \sim \text{Poisson}(\lambda p)$ and $EX = \lambda p$.

Using the iterative expectation formula

$$EX = E(E(X|Y)) = E(Yp) = \lambda p$$

BIOS 660/BIOS 672 (3 Credits)

Notes 15 – 8 / 44

Three-layer hierarchical model

Example: In the previous example, suppose that there are several insect mothers, each with different average number of eggs. Model λ of the Poisson as being an exponential rv Λ with mean parameter β . Now what is the expected number of eggs surviving?

$$\begin{aligned} EX &= E(E(X|Y)) \\ &= E(pY) \\ &= E(E(pY|\Lambda)) \\ &= E(p\Lambda) \\ &= p\beta \end{aligned}$$

BIOS 660/BIOS 672 (3 Credits)

Notes 15 – 9 / 44

Conditional Variance

Suppose we have rvs X and Y . Recall that the marginal variance of $g(Y)$ is

$$\text{Var}[g(Y)] = E[g(Y) - E(g(Y))]^2$$

The *conditional variance* of $g(Y)$ given X is

$$\text{Var}[g(Y)|X] = E\{[g(Y) - E(g(Y)|X)]^2|X\}$$

where both expectations are taken with respect to the conditional pmf or pdf $f_{Y|X}(y)$.

Like the conditional expectation, the conditional variance of Y given X is a random variable whose value depends on the rv X .

BIOS 660/BIOS 672 (3 Credits)

Notes 15 – 10 / 44

Conditional variance hierarchical formula

For any two random variables X and Y ,

$$\text{Var}X = E(\text{Var}(X|Y)) + \text{Var}(E(X|Y))$$

Proof

$$\begin{aligned}\text{Var}X &= E\{[X - EX]^2\} \\ &= E\{[X - E(X|Y) + E(X|Y) - EX]^2\} \\ &= E\{[X - E(X|Y)]^2\} + E\{[E(X|Y) - EX]^2\} \\ &\quad + 2E\{[X - E(X|Y)][E(X|Y) - EX]\}\end{aligned}$$

Study the three terms separately. 1st term:

$$\begin{aligned}E\{[X - E(X|Y)]^2\} &= E\left(E\{[X - E(X|Y)]^2|Y\}\right) \\ &= E(\text{Var}(X|Y))\end{aligned}$$

cont

2nd term:

$$\begin{aligned}E\{[E(X|Y) - EX]^2\} &= E\{[E(X|Y) - E(E(X|Y))]^2\} \\ &= \text{Var}(E(X|Y))\end{aligned}$$

3rd term:

$$\begin{aligned}E\{[X - E(X|Y)][E(X|Y) - EX]\} &= E\left(E\{[X - E(X|Y)][E(X|Y) - EX]|Y\}\right) \\ &= E\left([E(X|Y) - EX] E\{[X - E(X|Y)]|Y\}\right) \\ &= E\left([E(X|Y) - EX] \{E[X|Y] - E(X|Y)\}\right) \\ &= 0\end{aligned}$$

Example

In the the Poisson-Binomial hierarchical model, we had

$$Y \sim \text{Poisson}(\lambda), \quad X|Y \sim \text{Binomial}(Y, p)$$

and showed that $EX = \lambda p$. Using the conditional variance formula,

$$\begin{aligned}\text{Var}X &= E(\text{Var}(X|Y)) + \text{Var}(E(X|Y)) \\ &= E(Yp(1-p)) + \text{Var}(Yp) \\ &= \lambda p(1-p) + \lambda p^2 \\ &= \lambda p\end{aligned}$$

This is consistent with the result that $X \sim \text{Poisson}(\lambda p)$.

Covariance and Correlation

14 / 44

Definitions

Let X and Y be two random variables with respective means μ_X, μ_Y and variances $\sigma_X^2 > 0$ and $\sigma_Y^2 > 0$, all assumed to exist.

- The *covariance* of X and Y is

$$\text{Cov}(X, Y) = E[(X - \mu_X)(Y - \mu_Y)] = \sigma_{XY}$$

- The *correlation* between X and Y is

$$\begin{aligned}\rho_{XY} &= \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}X \text{Var}Y}} = \frac{\sigma_{XY}}{\sigma_X \sigma_Y} \\ &= E\left[\left(\frac{X - \mu_X}{\sigma_X}\right)\left(\frac{Y - \mu_Y}{\sigma_Y}\right)\right]\end{aligned}$$

- X and Y are called *uncorrelated* iff

$$\text{Cov}(X, Y) = 0 \quad \text{or equivalently} \quad \rho_{XY} = 0$$

Properties of Covariance

1. $\text{Cov}(X, X) = \text{Var}(X)$
2. $\text{Cov}(X, Y) = \text{Cov}(Y, X)$
3. For any constant c , $\text{Cov}(X, c) = 0$
4. $\text{Cov}(X, Y) = E(XY) - E(X)E(Y)$
5. If X and Y are independent and $\text{Cov}(X, Y)$ exists, then $\text{Cov}(X, Y) = 0$.

Note: If X and Y are uncorrelated, this does not imply that they are independent. Example:

$$X \sim U[-1, 1], \quad Y = \begin{cases} X, & \text{prob. } 1/2 \\ -X, & \text{prob. } 1/2 \end{cases}$$

BIOS 660/BIOS 672 (3 Credits)

Notes 15 – 16 / 44

Linear Combinations

If X , Y , and Z are rvs each with a variance, and a and b are constants, then

$$\text{Cov}(aX + bY, Z) = a\text{Cov}(X, Z) + b\text{Cov}(Y, Z)$$

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

$$\text{Corr}(aX + b, cY + d) = \text{sign}(ac)\text{Corr}(X, Y)$$

Proof:

BIOS 660/BIOS 672 (3 Credits)

Notes 15 – 17 / 44

Correlation Coefficient

Theorem C-B 4.5.7 For any rvs X and Y ,

1. $-1 \leq \rho_{XY} \leq 1$
2. $|\rho_{XY}| = 1$ if and only if $\exists a \neq 0$ and b such that

$$P(Y = aX + b) = 1.$$

If $\rho_{XY} = 1 \Rightarrow a > 0$, and if $\rho_{XY} = -1 \Rightarrow a < 0$.

Proof: Let $\tilde{X} = (X - \mu_X)/\sigma_X$, $\tilde{Y} = (Y - \mu_Y)/\sigma_Y$, so that $\rho_{XY} = E(\tilde{X}\tilde{Y})$.

1.
$$0 \leq E(\tilde{X} - \tilde{Y})^2 = 1 + 1 - 2E(\tilde{X}\tilde{Y}) \Rightarrow E(\tilde{X}\tilde{Y}) \leq 1$$
$$0 \leq E(\tilde{X} + \tilde{Y})^2 = 1 + 1 + 2E(\tilde{X}\tilde{Y}) \Rightarrow -1 \leq E(\tilde{X}\tilde{Y})$$
2.
$$\rho_{XY} = 1 \text{ iff } P(\tilde{Y} = \tilde{X}) = 1 \Rightarrow a > 0$$
$$\rho_{XY} = -1 \text{ iff } P(\tilde{Y} = -\tilde{X}) = 1 \Rightarrow a < 0$$

Bivariate Normal

19 / 44

Standard Bivariate Normal

Given a number $-1 \leq \rho \leq 1$, define the *standard bivariate normal density* of $(X, Y) \in \mathbb{R}^2$ by

$$f_{XY}(x, y) = \frac{1}{2\pi\sqrt{1-\rho^2}} \exp \left[-\frac{x^2 - 2\rho xy + y^2}{2(1-\rho^2)} \right]$$

Properties:

1. The marginal distribution of X is $N(0, 1)$.
2. The marginal distribution of Y is $N(0, 1)$.
3. The correlation of X and Y is ρ .
4. The conditional distributions are normal:

$$Y|X \sim N(\rho X, 1 - \rho^2), \quad X|Y \sim N(\rho Y, 1 - \rho^2)$$

The means are the *regression lines* of Y on X and X on Y respectively.

Bivariate Normal

Let \tilde{X} and \tilde{Y} have a standard bivariate normal distribution with correlation ρ . Let

$$\begin{aligned}X &= \mu_X + \sigma_X \tilde{X}, & \mu_X &\in \mathbb{R}, \sigma_X > 0 \\Y &= \mu_Y + \sigma_Y \tilde{Y}, & \mu_Y &\in \mathbb{R}, \sigma_Y > 0\end{aligned}$$

Then (X, Y) has the *bivariate normal density*

$$\begin{aligned}f_{XY}(x, y) &= \left(2\pi\sigma_X\sigma_Y\sqrt{1-\rho^2}\right)^{-1} \\&\times \exp \left\{ -\frac{1}{2(1-\rho^2)} \left[\left(\frac{x-\mu_X}{\sigma_X}\right)^2 \right. \right. \\&\quad \left. \left. - 2\rho \left(\frac{x-\mu_X}{\sigma_X}\right) \left(\frac{y-\mu_Y}{\sigma_Y}\right) + \left(\frac{y-\mu_Y}{\sigma_Y}\right)^2 \right] \right\}\end{aligned}$$

for $-\infty < x < \infty$ and $-\infty < y < \infty$.

BIOS 660/BIOS 672 (3 Credits)

Notes 15 – 21 / 44

Properties

1. The marginal distribution of X is $N(\mu_X, \sigma_X^2)$.
2. The marginal distribution of Y is $N(\mu_Y, \sigma_Y^2)$.
3. The correlation between X and Y is ρ .
4. The conditional distributions are normal:

$$Y|X \sim N \left[\mu_Y + \rho(\sigma_Y/\sigma_X)(x - \mu_X), \sigma_Y^2(1 - \rho^2) \right]$$

The mean is the *regression line* of Y on X .

5. For any constants a and b , the distribution of $aX + bY$ is

$$N(a\mu_X + b\mu_Y, a^2\sigma_X^2 + b^2\sigma_Y^2 + 2ab\rho\sigma_X\sigma_Y)$$

Proof: HW

(Suggestion: standardize the variables first)

BIOS 660/BIOS 672 (3 Credits)

Notes 15 – 22 / 44

Multivariate Distributions

The n -dimensional random vector $\mathbf{X} = (X_1, X_2, \dots, X_n)^\top$ defined on the triplet (S, \mathcal{B}, P) takes values on the hyperspace \mathcal{R}^n .

If the variables X_1, \dots, X_n are *discrete* then we have a discrete random vector, if the X 's are *continuous*, then we have a continuous random vector.

If \mathbf{X} is discrete then

$$P(\mathbf{X} \in A) = \sum_{\mathbf{x} \in A} f(\mathbf{x})$$

where $f(\mathbf{x})$ denotes the *joint pmf*.

If \mathbf{X} is continuous then

$$P(\mathbf{X} \in A) = \int \dots \int_A f(x_1, \dots, x_n) dx_1 \dots dx_n$$

where $f(\mathbf{x})$ denotes the *joint pdf*.

BIOS 660/BIOS 672 (3 Credits)

Notes 15 – 24 / 44

Marginals and Conditionals

Definition The *marginal pdf* or *pmf* of any subset of the coordinates of (X_1, \dots, X_n) can be computed by integrating or summing the joint pdf or pmf over all possible values of the other coordinates.

Definition The *conditional pdf* or *pmf* of a subset of the coordinates of (X_1, \dots, X_n) given the values of the remaining coordinates is obtained by dividing the full joint pdf or pmf by the joint pdf or pmf of the conditioning variates:

$$f(x_{k+1}, \dots, x_n | x_1, \dots, x_k) = \frac{f(x_1, \dots, x_n)}{f(x_1, \dots, x_k)}$$

BIOS 660/BIOS 672 (3 Credits)

Notes 15 – 25 / 44

Multinomial Distribution

Let n and m be positive integers and let p_1, \dots, p_n be probabilities summing to one. Then the random vector (X_1, \dots, X_n) has a *multinomial distribution with m trials and cell probabilities p_1, \dots, p_n* if its joint pmf is

$$\begin{aligned} f(x_1, \dots, x_n) &= \binom{m}{x_1 \dots x_n} p_1^{x_1} \dots p_n^{x_n} \\ &= \frac{m!}{x_1! \dots x_n!} p_1^{x_1} \dots p_n^{x_n} = m! \prod_{j=1}^n \frac{p_j^{x_j}}{x_j!} \end{aligned}$$

for $x_i = 0, \dots, m$, $i = 1, \dots, n$, and $x_1 + \dots + x_n = m$.

The proof that this is a pmf is called the Multinomial Theorem.

BIOS 660/BIOS 672 (3 Credits)

Notes 15 – 26 / 44

Multinomial Theorem 4.6.4

Let m and n be positive integers. Let \mathcal{A} be the set of all vectors $\mathbf{x} = (x_1, \dots, x_n)$ which are such that the sum of their nonnegative integer components is m , i.e. $\sum_{j=1}^n x_j = m$ and $x_j \geq 0$. Then, for any real numbers p_1, \dots, p_n

$$(p_1 + \dots + p_n)^m = \sum_{\mathbf{x} \in \mathcal{A}} \frac{m!}{x_1! \dots x_n!} p_1^{x_1} \dots p_n^{x_n}$$

E.g. for $n = 2$ we have the binomial theorem.

E.g. for $n = 3$:

$$(p_1 + p_2 + p_3)^m = \sum_{x_1=0}^m \sum_{x_2=0}^{m-x_1} \frac{m!}{x_1! x_2! (m-x_1-x_2)!} p_1^{x_1} p_2^{x_2} p_3^{m-x_1-x_2}$$

BIOS 660/BIOS 672 (3 Credits)

Notes 15 – 27 / 44

Multinomial Distribution: Properties

- Marginals are multinomials. E.g.

$$\begin{aligned} X_1 &\sim \text{Binomial}(m, p_1) \\ (X_1, X_2, m - X_1 - X_2) &\sim \text{Multinomial}(m, p_1, p_2, 1 - p_1 - p_2) \end{aligned}$$

- Conditionals are multinomials. E.g.

$$(X_1, \dots, X_{n-1}) | [X_n = x_n] \sim \text{Multinomial}\left(m - x_n, \frac{p_1}{1 - p_n}, \dots, \frac{p_{n-1}}{1 - p_n}\right)$$

- Variance and covariance:

$$\text{Cov}(X_j, X_k) = E[(X_j - mp_j)(X_k - mp_k)] = \begin{cases} mp_j(1 - p_j), & j = k \\ -mp_j p_k, & j \neq k \end{cases}$$

Multivariate Independence

29 / 44

Independent Random Vectors

Let $\mathbf{X}_1, \dots, \mathbf{X}_n$ be random vectors with joint pdf or pmf $f(\mathbf{x}_1, \dots, \mathbf{x}_n)$. Let $f_{\mathbf{X}_j}(\mathbf{x}_j)$ denote the marginal pdf or pmf of \mathbf{X}_j . Then $\mathbf{X}_1, \dots, \mathbf{X}_n$ are called *mutually independent random vectors* if, for every $(\mathbf{x}_1, \dots, \mathbf{x}_n)$,

$$f(\mathbf{x}_1, \dots, \mathbf{x}_n) = \prod_{j=1}^n f_{\mathbf{X}_j}(\mathbf{x}_j)$$

cont.

Theorem C-B 4.6.6 (Generalization of 4.2.10) Let X_1, \dots, X_n be independent rvs. Let g_1, \dots, g_n be real-valued functions such that $g(x_j)$ is only a function of x_j . Then

$$E[g_1(X_1) \dots g_n(X_n)] = E g_1(X_1) \dots E g_n(X_n)$$

Theorem C-B 4.6.7 (Generalization of 4.2.12) Let X_1, \dots, X_n be mutually independent rvs with characteristic functions $\phi_{X_1}(\theta), \dots, \phi_{X_n}(\theta)$. Let $Z = X_1 + \dots + X_n$. Then the characteristic function of Z is

$$\phi_Z(\theta) = \prod_{j=1}^n \phi_{X_j}(\theta)$$

Note simplification if $\phi_{X_j}(\theta) = \phi_X(\theta)$.

Specific Characteristic Functions

	<u>mgf</u>	<u>cf</u>
Bernoulli(p)	$pe^t + q$	$pe^{it} + q$
Binomial(n, p)	$(pe^t + q)^n$	$(pe^{it} + q)^n$
Poisson(λ)	$e^{\lambda(e^t - 1)}$	$e^{\lambda(e^{it} - 1)}$
Geometric(p)	$pe^t / (1 - qe^t)$	$pe^{it} / (1 - qe^{it})$
Negbin(n, p)	$\left[\frac{pe^t}{1 - qe^t} \right]^n$	$\left[\frac{pe^{it}}{1 - qe^{it}} \right]^n$
Uniform(a, b)	$\frac{e^{tb} - e^{ta}}{t(b-a)}$	$\frac{e^{itb} - e^{ita}}{it(b-a)}$
Normal(μ, σ^2)	$e^{t\mu + \frac{1}{2}\sigma^2 t^2}$	$e^{it\mu - \frac{1}{2}\sigma^2 t^2}$
Exponential(λ)	$\frac{\lambda}{\lambda - t} = (1 - \frac{t}{\lambda})^{-1}$	$(1 - \frac{it}{\lambda})^{-1}$
Gamma(a, λ)	$(1 - \frac{t}{\lambda})^{-a}$	$(1 - \frac{it}{\lambda})^{-a}$

Multivariate moments and multivariate normal (Gut, Chapter V)

33 / 44

Multivariate moments

Let $\mathbf{X} = (X_1, \dots, X_n)^\top$ be a random vector.
The *mean vector* $\boldsymbol{\mu}$ is

$$\boldsymbol{\mu} = E(\mathbf{X}) = \begin{pmatrix} EX_1 \\ \vdots \\ EX_n \end{pmatrix} = \begin{pmatrix} \mu_1 \\ \vdots \\ \mu_n \end{pmatrix}$$

We can also define the second moment matrix

$$E(\mathbf{X}\mathbf{X}^\top) = \begin{pmatrix} EX_1^2 & EX_1X_2 & \cdots & EX_1X_n \\ EX_2X_1 & EX_2^2 & \cdots & EX_2X_n \\ \vdots & \vdots & \ddots & \vdots \\ EX_nX_1 & EX_nX_2 & \cdots & EX_n^2 \end{pmatrix}$$

BIOS 660/BIOS 672 (3 Credits)

Notes 15 – 34 / 44

Multivariate covariance

The *variance-covariance matrix* Σ is defined as

$$\begin{aligned} \Sigma &= \text{Cov} [\mathbf{X} - \boldsymbol{\mu}, (\mathbf{X} - \boldsymbol{\mu})^\top] \\ &= E [(\mathbf{X} - \boldsymbol{\mu})(\mathbf{X} - \boldsymbol{\mu})^\top] \\ &= \begin{pmatrix} \text{Var}(X_1) & \text{Cov}(X_1, X_2) & \cdots & \text{Cov}(X_1, X_n) \\ \text{Cov}(X_2, X_1) & \text{Var}(X_2) & \cdots & \text{Cov}(X_2, X_n) \\ \vdots & \vdots & \ddots & \vdots \\ \text{Cov}(X_n, X_1) & \text{Cov}(X_n, X_2) & \cdots & \text{Var}(X_n) \end{pmatrix} \end{aligned}$$

Notice that Σ is a symmetric matrix.

BIOS 660/BIOS 672 (3 Credits)

Notes 15 – 35 / 44

Bivariate normal

Let X and Y are bivariate normal, then

$$\boldsymbol{\mu} = \begin{pmatrix} \mu_X \\ \mu_Y \end{pmatrix}, \quad \boldsymbol{\Sigma} = \begin{pmatrix} \sigma_X^2 & \rho\sigma_X\sigma_Y \\ \rho\sigma_X\sigma_Y & \sigma_Y^2 \end{pmatrix}$$

The joint density of the vector $\mathbf{X} = (X, Y)^\top$ can be written as

$$f_{\mathbf{X}}(\mathbf{X}) = \frac{1}{2\pi\sqrt{\det(\boldsymbol{\Sigma})}} \exp \left[-\frac{1}{2}(\mathbf{X} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{X} - \boldsymbol{\mu}) \right]$$

Proof:

BIOS 660/BIOS 672 (3 Credits)

Notes 15 – 36 / 44

Linear functions

Let $\mathbf{X} = (X_1, \dots, X_n)^\top$ be a random vector with mean $\boldsymbol{\mu}$ and covariance $\boldsymbol{\Sigma}$. For a vector $\mathbf{c} = (c_1, c_2, \dots, c_n)^\top \in \mathbb{R}^n$, define

$$Y = \mathbf{c}^\top \mathbf{X} = \sum_{i=1}^n c_i X_i$$

Then

$$\begin{aligned} \mathbb{E}(Y) &= \mathbf{c}^\top \boldsymbol{\mu} \\ \text{Var}(Y) &= \mathbf{c}^\top \boldsymbol{\Sigma} \mathbf{c} \end{aligned}$$

Proof:

BIOS 660/BIOS 672 (3 Credits)

Notes 15 – 37 / 44

Positive definiteness

Definition:

- An $n \times n$ symmetric matrix $\mathbf{\Lambda}$ is called *positive semi-definite* or *nonnegative definite* iff for every vector $\mathbf{c} \in \mathbb{R}^n$, $\mathbf{c}^\top \mathbf{\Lambda} \mathbf{c} \geq 0$.
- An $n \times n$ symmetric matrix $\mathbf{\Lambda}$ is called *positive definite* iff for every vector $\mathbf{c} \in \mathbb{R}^n$, $\mathbf{c}^\top \mathbf{\Lambda} \mathbf{c} > 0$.

Properties:

- A positive definite matrix:
 - is invertible
 - its determinant is positive
 - all its eigenvalues are positive
- The variance-covariance matrix is positive semi-definite.

BIOS 660/BIOS 672 (3 Credits)

Notes 15 – 38 / 44

Multivariate linear transformations

Let $\mathbf{X} = (X_1, X_2, \dots, X_n)^\top$ have joint pdf $f_{\mathbf{X}}(\mathbf{X})$ and let

$$\begin{aligned} Y_1 &= a_{11}X_1 + a_{12}X_2 + \dots a_{1n}X_n \\ Y_2 &= a_{21}X_1 + a_{22}X_2 + \dots a_{2n}X_n \\ &\vdots \\ Y_n &= a_{n1}X_1 + a_{n2}X_2 + \dots a_{nn}X_n \end{aligned}$$

Using vector and matrix notation, we can write this transformation as

$$\mathbf{Y} = \mathbf{A}\mathbf{X}$$

where

$$\mathbf{A} = \begin{pmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \dots & a_{nn} \end{pmatrix}$$

BIOS 660/BIOS 672 (3 Credits)

Notes 15 – 39 / 44

Multivariate linear transformations

The Jacobian of the transformation

$$\mathbf{Y} = \mathbf{A}\mathbf{X}$$

is

$$\frac{\partial \mathbf{Y}}{\partial \mathbf{X}^\top} = \begin{pmatrix} \partial y_1 / \partial x_1 & \dots & \partial y_1 / \partial x_n \\ \vdots & \ddots & \vdots \\ \partial y_n / \partial x_1 & \dots & \partial y_n / \partial x_n \end{pmatrix} = \begin{pmatrix} a_{11} & \dots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{n1} & \dots & a_{nn} \end{pmatrix} = \mathbf{A}$$

with determinant $\det(\mathbf{A})$.

The inverse transformation is

$$\mathbf{X} = \mathbf{A}^{-1}\mathbf{Y}$$

with Jacobian determinant $\det(\mathbf{A}^{-1}) = 1/\det(\mathbf{A})$.

BIOS 660/BIOS 672 (3 Credits)

Notes 15 – 40 / 44

Multivariate normal

Let $\boldsymbol{\mu}$ be a vector in \mathbb{R}^n and $\boldsymbol{\Sigma}$ be an $n \times n$ symmetric positive definite matrix. The vector $\mathbf{X} = (X_1, \dots, X_n)$ has a *multivariate normal distribution* with parameters $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$ if it has joint density

$$f_{\mathbf{X}}(\mathbf{X}) = \frac{1}{(2\pi)^{n/2} \sqrt{\det(\boldsymbol{\Sigma})}} \exp \left[-\frac{1}{2} (\mathbf{X} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{X} - \boldsymbol{\mu}) \right]$$

In particular:

$$\begin{aligned} \boldsymbol{\mu} &= \mathbf{E}\mathbf{X} \\ \boldsymbol{\Sigma} &= \mathbf{E}[(\mathbf{X} - \boldsymbol{\mu})(\mathbf{X} - \boldsymbol{\mu})^\top] \end{aligned}$$

BIOS 660/BIOS 672 (3 Credits)

Notes 15 – 41 / 44

Construction

Let $\mathbf{Z} = (Z_1, \dots, Z_n)$ be i.i.d. $N(0, 1)$ variables. Then

$$f_{\mathbf{Z}}(z_1, \dots, z_n) = \prod_{i=1}^n \frac{e^{-z_i^2/2}}{\sqrt{2\pi}} = \frac{e^{-\sum_{i=1}^n z_i^2/2}}{(2\pi)^{n/2}} = \frac{e^{-\mathbf{Z}^\top \mathbf{Z}/2}}{(2\pi)^{n/2}}$$

Now define

$$\mathbf{X} = \mathbf{A}\mathbf{Z} + \boldsymbol{\mu}$$

where $\boldsymbol{\mu} \in \mathbb{R}^n$ is a vector and $\mathbf{A} \in \mathbb{R}^{n \times n}$ is an invertible matrix.

The inverse transformation is

$$\mathbf{Z} = \mathbf{A}^{-1}(\mathbf{X} - \boldsymbol{\mu})$$

with Jacobian

$$J = \frac{\partial \mathbf{Z}}{\partial \mathbf{X}^\top} = \mathbf{A}^{-1}$$

and determinant $|J| = 1/|\mathbf{A}|$.

BIOS 660/BIOS 672 (3 Credits)

Notes 15 – 42 / 44

Construction (cont.)

The joint pdf of \mathbf{X} is

$$f_{\mathbf{X}}(\mathbf{X}) = \frac{e^{-(\mathbf{X}-\boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{X}-\boldsymbol{\mu})/2}}{(2\pi)^{n/2} |\boldsymbol{\Sigma}|^{1/2}}$$

where $\boldsymbol{\Sigma} = \mathbf{A}\mathbf{A}^\top$. We can confirm that

$$\mathbb{E}(\mathbf{X}) = \mathbb{E}(\mathbf{A}\mathbf{Z} + \boldsymbol{\mu}) = \boldsymbol{\mu}$$

and

$$\begin{aligned} \mathbb{E}[(\mathbf{X} - \boldsymbol{\mu})(\mathbf{X} - \boldsymbol{\mu})^\top] &= \mathbb{E}[(\mathbf{A}\mathbf{Z})(\mathbf{A}\mathbf{Z})^\top] = \mathbb{E}(\mathbf{A}\mathbf{Z}\mathbf{Z}^\top \mathbf{A}^\top) \\ &= \mathbf{A}\mathbb{E}(\mathbf{Z}\mathbf{Z}^\top) \mathbf{A}^\top = \mathbf{A}\mathbf{I}\mathbf{A}^\top = \boldsymbol{\Sigma} \end{aligned}$$

BIOS 660/BIOS 672 (3 Credits)

Notes 15 – 43 / 44

Properties

Theorem:

1. Let X and Y be jointly normal random variables. Then X and Y are independent if and only if they are uncorrelated.
2. Let X_1, \dots, X_n be jointly normal random variables. Then they are mutually independent if and only if they are pairwise uncorrelated.

Theorem: If \mathbf{X} has a multivariate normal distribution, then

1. All marginal distributions are normal.
2. All conditional distributions are normal.
3. For any constants $\mathbf{A} = (a_1, \dots, a_n)^\top$ the random variable $Y = \mathbf{A}^\top \mathbf{X}$ has a normal distribution.

BIOS 660/BIOS 672 (3 Credits): Probability and Statistical Inference I

Jianwen Cai

<https://sakai.unc.edu/portal/site/bios660-bios672-3-credits>

Notes 16

Inequalities	2
Recall: Chebychev Inequality	3
Corollaries	4
Special cases	5
Functional inequalities	6
Convex functions	7
Jensen's Inequality	8
Jensen's Inequality (proof)	9
Example	10
Young's Inequality	11
Hölder's inequality	12
Corollaries	13
Application of Cauchy-Schwartz:	14
Minkowski's inequality	15
Order Statistics (C-B, Section 5.4; Gut, Chapter IV)	16
Distribution of the Maximum	17
Distribution of the Minimum	18
Example	19
Order Statistics	20
r^{th} order statistic	21
cont.	22
Example	23
Distribution of the median	24
Joint distribution of $Y_{(r)}, Y_{(s)}, r < s$	25
cont.	26
Joint distribution of all order statistics	27
Distribution of the range	28
Example:	29
Convergence	30
Agenda	31

Modes of Convergence.	32
Notes	33
Convergence in Distribution	34
Convergence in Distrib. cont.	35
Other modes of convergence	36
Example 1.	37
Example 2.	38
Example 3.	39
Example 4.	40
Relations	41
Relationships among convergence modes.	42
In probability and distribution	43
cont.	44
cont.	45
r th moment and in Probability	46
Convergence properties	47
Convergence in probability	48
Slutsky's Theorem	49
Example	50
Example	51
Convergence in distribution	52
Example	53
The Delta Method	54
Approximate mean and variance	55
Example	56
The Delta method.	57
Example	58
Second-order Delta method	59

Recall: Chebychev Inequality

Let X be a random variable and let $g(x)$ be a non-negative function. Then for any $r > 0$,

$$P[g(X) \geq r] \leq \frac{Eg(X)}{r}$$

Proof:

$$\begin{aligned} Eg(X) &= \int_{-\infty}^{\infty} g(x) f_X(x) dx \\ &\geq \int_{\{x: g(x) \geq r\}} g(x) f_X(x) dx \\ &\geq r \int_{\{x: g(x) \geq r\}} f_X(x) dx \\ &= r P\{g(X) \geq r\} \end{aligned}$$

BIOS 660/BIOS 672 (3 Credits)

Notes 16 – 3 / 59

Corollaries

1. Suppose X is a non-negative and g is a positive, non-decreasing function, with $E[g(X)] < \infty$. Then

$$P\{X \geq a\} \leq \frac{E(g(X))}{g(a)}$$

2. Suppose g is a non-negative symmetric function, increasing on \mathbb{R}^+ , with $E[g(X)] < \infty$. Then

$$P\{|X| \geq a\} \leq \frac{E[g(X)]}{g(a)}$$

Proof: $P\{X \geq a\} = P\{g(X) \geq g(a)\}$, so the results follow from the inequality on the previous slides.

BIOS 660/BIOS 672 (3 Credits)

Notes 16 – 4 / 59

Special cases

Provided that the appropriate expectations exist, for $a > 0$:

$$X \geq 0 : \quad P\{X \geq a\} \leq \frac{E(e^{tX})}{e^{ta}} \quad (1)$$

$$X \in \mathbb{R} : \quad P\{|X| \geq a\} \leq \frac{E(|X|)}{a} \quad (2)$$

$$X \in \mathbb{R}, p > 0 : \quad P\{|X| \geq a\} \leq \frac{E(|X|^p)}{a^p} \quad (3)$$

$$\sigma^2 = \text{Var}(X) : \quad P\{|X - EX| \geq a\sigma\} \leq \frac{1}{a^2} \quad (4)$$

Note:

- (1) is called Chernoff bound, useful when the mgf is easier to compute than the cdf.
- (3) is sometimes called Markov's inequality.
- (4) is sometimes called Chebychev's inequality.

Functional inequalities

First a couple of useful items:

- *L^p spaces:*
The space called L^p consists of all random variables whose p^{th} absolute power is integrable, i.e., $E(|X|^p) < \infty$.
- *Triangle inequality:*
For two real or complex numbers a and b ,

$$|a + b| \leq |a| + |b|$$

Proof: HW

Convex functions

Definition: Let I be an interval on \mathbb{R} . A function $g : I \rightarrow \mathbb{R}$ is *convex* on I if for any $\lambda \in [0, 1]$, and any points x and y in I

$$g[\lambda x + (1 - \lambda)y] \leq \lambda g(x) + (1 - \lambda)g(y)$$

Properties:

- A differentiable function g is convex iff it lies above all its tangents.
- A twice differentiable function g is convex iff its second derivative is non-negative.

Definition: Let I be an interval on \mathbb{R} . A function $g : I \rightarrow \mathbb{R}$ is *concave* on I if $-g$ is convex on I .

Examples:

- $g(x) = x^2$ is a convex function for all x .
- $g(x) = \log(x)$ is concave for $x > 0$.

Jensen's Inequality

Let $X \in L^1$ and $g(x)$ be a convex function where $E[g(X)]$ exists. Then,

$$E[g(X)] \geq g[EX]$$

with equality if and only if for every line $a + bx$ tangent to $g(x)$ at $x = EX$, $P[g(X) = a + bX] = 1$.

Examples:

$$\begin{aligned} g(x) = x^2 &\rightarrow E(X^2) \geq E^2(X) \\ g(x) = 1/x, x > 0 &\rightarrow E(1/X) \geq 1/E(X), X > 0 \end{aligned}$$

Note: The direction of the inequality is reversed if g is concave.

Jensen's Inequality (proof)

Let $l(x) = a + bx$ be the tangent line to $g(x)$ at $g(\mathbf{E}X)$. Then

$$\begin{aligned}\mathbf{E}g(X) &\geq \mathbf{E}(a + bX) \\ &= a + b\mathbf{E}X \\ &= l(\mathbf{E}X) \\ &= g(\mathbf{E}X)\end{aligned}$$

BIOS 660/BIOS 672 (3 Credits)

Notes 16 – 9 / 59

Example

Let $a_1, \dots, a_n > 0$. Then

$$\left(\frac{1}{n} \sum_{i=1}^n \frac{1}{a_i}\right)^{-1} \leq \left(\prod_{i=1}^n a_i\right)^{1/n} \leq \frac{1}{n} \sum_{i=1}^n a_i$$

Proof: Let X be a rv such that $P(X = a_i) = 1/n$. Since $\log(x)$ is concave,

$$\begin{aligned}\log\left(\prod_{i=1}^n a_i\right)^{1/n} &= \frac{1}{n} \sum_{i=1}^n \log(a_i) \\ &= \mathbf{E}(\log(X)) \\ &\leq \log(\mathbf{E}(X)) \\ &= \log\left(\frac{1}{n} \sum_{i=1}^n a_i\right)\end{aligned}$$

The proof of the second inequality is similar.

BIOS 660/BIOS 672 (3 Credits)

Notes 16 – 10 / 59

Young's Inequality

Let $a, b > 0$ and $p, q > 1$ with $1/p + 1/q = 1$. Then

$$\frac{a^p}{p} + \frac{b^q}{q} \geq ab$$

With equality only if $a^p = b^q$.

Proof: Consider

$$g(a) = \frac{1}{p}a^p + \frac{1}{q}b^q - ab$$

To minimize $g(a)$, differentiate and set equal to 0:

$$\frac{d}{da}g(a) = 0 \rightarrow a^{p-1} - b = 0 \rightarrow a = b^{1/(p-1)}.$$

Since $g(b^{1/(p-1)}) = 0$, the result follows.

Hölder's inequality

Suppose $X \in L^p, Y \in L^q$ where $p, q > 1$ and $1/p + 1/q = 1$. Then

$$\mathbb{E}|XY| \leq [\mathbb{E}|X|^p]^{1/p} [\mathbb{E}|Y|^q]^{1/q}$$

with equality if $X^p = cY^q$ for some $c \in \mathbb{R}$.

Proof: Let

$$a = \frac{|X|}{(\mathbb{E}|X|^p)^{1/p}} \quad \text{and} \quad b = \frac{|Y|}{(\mathbb{E}|Y|^q)^{1/q}}$$

By Young's Inequality,

$$\frac{|X|^p}{p\mathbb{E}|X|^p} + \frac{|Y|^q}{q\mathbb{E}|Y|^q} \geq \frac{|XY|}{(\mathbb{E}|X|^p)^{1/p}(\mathbb{E}|Y|^q)^{1/q}}$$

The result follows by taking the expected value of both sides and noting that the expected value of the left-hand side is 1.

Corollaries

- **Cauchy-Schwartz inequality:** Special case where $p = q = 2$.

$$E|XY| \leq [E|X|^2]^{1/2} [E|Y|^2]^{1/2} = \sqrt{E[X^2]E[Y^2]}$$

with equality if $X = cY$ for some $c \in \mathbb{R}$.

- **Lyapunov's inequality:** For $1 \leq r \leq s$ and $X \in L^s$,

$$[E|X|^r]^{1/r} \leq [E|X|^s]^{1/s}$$

Proof:

Apply Hölder's inequality to $|X|^r$ with $Y = 1$ and $p = s/r$.

Application of Cauchy-Schwartz:

Let ρ represent the correlation between two rvs X and Y , ie,

$$\rho = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X) \text{Var}(Y)}}.$$

Then, $|\rho| \leq 1$, with equality iff $Y - \mu_Y = c(X - \mu_X)$ for some $c \in \mathbb{R}$.

Proof: By the Cauchy-Schwartz Inequality,

$$E|(X - \mu_X)(Y - \mu_Y)| \leq \{E(X - \mu_X)^2\}^{\frac{1}{2}} \{E(Y - \mu_Y)^2\}^{\frac{1}{2}}.$$

Squaring both sides, we obtain

$$(\text{Cov}(X, Y))^2 \leq \sigma_X^2 \sigma_Y^2.$$

Thus, $|\rho| \leq 1$.

Minkowski's inequality

Suppose $X, Y \in L^p$, $p \geq 1$. Then $(X + Y) \in L^p$ and

$$[E|X + Y|^p]^{1/p} \leq [E|X|^p]^{1/p} + [E|Y|^p]^{1/p}$$

Proof:

For $p = 1$, the proof follows almost immediately from the triangle inequality (HW).

The case $p > 1$ is more complicated. Consider

$$\begin{aligned} E|X + Y|^p &= E(|X + Y| |X + Y|^{p-1}) \\ &\leq E(|X| |X + Y|^{p-1}) + E(|Y| |X + Y|^{p-1}) \\ &\leq [E|X|^p]^{1/p} [E|X + Y|^{(p-1)q}]^{1/q} \\ &\quad + [E|Y|^p]^{1/p} [E|X + Y|^{(p-1)q}]^{1/q} \end{aligned}$$

where the last row follows by Hölder's inequality with $1/p + 1/q = 1$.

Order Statistics (C-B, Section 5.4; Gut, Chapter IV)

16 / 59

Distribution of the Maximum

The *cdf* of $Z = \max(Y_1, \dots, Y_n)$ is

$$\begin{aligned} F_Z(z) &= Pr\{Z \leq z\} \\ &= Pr\{Y_1 \leq z, Y_2 \leq z, \dots, Y_n \leq z\} \\ &= \prod_{j=1}^n Pr\{Y_j \leq z\} \quad (\text{indep}) \\ &= F_Y(z)^n \quad (\text{ident. distrib.}) \end{aligned}$$

and thus the density (or pmf) is:

$$f_Z(z) = nF_Y(z)^{n-1}f_Y(z)$$

Distribution of the Minimum

Similarly, consider $W = \min(Y_1, Y_2, \dots, Y_n)$.

$$\begin{aligned} 1 - F_W(w) &= Pr\{W > w\} \\ &= Pr\{Y_1 > w, Y_2 > w, \dots, Y_n > w\} \\ &= \prod_{j=1}^n Pr\{Y_j > w\} \quad (\text{indep}) \\ &= (1 - F_Y(w))^n \quad (\text{ident. distrib}) \end{aligned}$$

Thus

$$F_W(w) = 1 - (1 - F_Y(w))^n$$

and the corresponding density (or pmf) is:

$$f_W(w) = n(1 - F_Y(w))^{n-1} f_Y(w)$$

Example

Suppose $Y_i \sim \exp(\lambda)$:

$$f_Y(y) = \lambda e^{-\lambda y} \quad \text{for } y > 0, \quad 1 - F(y) = e^{-\lambda y}$$

Maximum:

$$f_Z(z) = n(1 - e^{-\lambda z})^{n-1} \lambda e^{-\lambda z} = n\lambda e^{-\lambda z} (1 - e^{-\lambda z})^{n-1}$$

Minimum:

$$f_W(w) = n(e^{-\lambda w})^{n-1} \lambda e^{-\lambda w} = (n\lambda) e^{-n\lambda w}$$

\Rightarrow exponential with parameter $n\lambda$

The next obvious statistic is the range defined as the difference of the maximum and the minimum, but to get its distribution we need the joint distribution of the maximum and the minimum.

Order Statistics

Let Y_1, Y_2, \dots, Y_n be *iid* with *pdf* $f_Y(x)$.
Order the observations; i.e.

$$Y_{(1)} \leq Y_{(2)} \leq \dots \leq Y_{(n)}$$

The $Y_{(i)}$ are called *order statistics*. For example, the *minimum* is $Y_{(1)}$ and the *maximum* is $Y_{(n)}$.

We are interested in finding the distribution of an arbitrary $Y_{(i)}$, as well as the joint distributions of sets of $Y_{(i)}$ s and $Y_{(j)}$ s.

e.g. Range = $Y_{(n)} - Y_{(1)}$
or interquartile range, or joint of median and interquartile range, etc....

r^{th} order statistic

We need to find the density of $Y_{(r)}$ at a value y :

$$\frac{\overbrace{\quad}^{r-1}}{y} \quad \overbrace{\quad}^1 \quad \overbrace{\quad}^{n-r}}{y + dy}$$

Consider 3 intervals $(-\infty, y)$, $[y, y + dy)$, $[y + dy, \infty)$. The number of observations in each of the intervals follows the tri-nomial distribution

$$f(s_1, s_2, s_3) = \frac{n!}{s_1! s_2! s_3!} p_1^{s_1} p_2^{s_2} p_3^{s_3}$$

The event that $y \leq Y_{(r)} < y + dy$ is the event that we have

$(r - 1)$ observations are less than y

$(n - r)$ observations are greater than y

1 observation is in interval; $y, y + dy$

In the trinomial distribution, this corresponds to

$$s_1 = r - 1, \quad s_2 = 1, \quad s_3 = n - r$$

$$p_1 = F_Y(y), \quad p_2 = f_Y(y)dy, \quad p_3 = 1 - F_Y(y + dy)$$

cont.

Taking the limit as $dy \rightarrow 0$, we get:

$$\begin{aligned} f_{Y_{(r)}}(y) &= \frac{n!}{(r-1)!(n-r)!} F_Y(y)^{r-1} [1 - F_Y(y)]^{n-r} f_Y(y) \\ &= \frac{F_Y(y)^{r-1} [1 - F_Y(y)]^{n-r} f_Y(y)}{B(r, n-r+1)} \end{aligned}$$

Gut has a more formal derivation based on deriving the joint density of the order statistics, then integrating out all but the r^{th} order statistic. (see also Casella and Berger, p.228).

BIOS 660/BIOS 672 (3 Credits)

Notes 16 – 22 / 59

Example

$F_Y(y) = y$, that is, $Y \sim \text{Uniform}(0, 1)$

$$f_{Y_{(r)}}(y) = \frac{y^{r-1}(1-y)^{n-r}}{B(r, n-r+1)}$$

hence, $Y_{(r)}$ follows a *Beta Distribution* with parameters r and $n-r+1$.

Note: $E[Y_{(r)}] = \frac{r}{n+1}$

BIOS 660/BIOS 672 (3 Credits)

Notes 16 – 23 / 59

Distribution of the median

To simplify, suppose the sample size is odd, $n = 2m + 1$, so that the median corresponds to the $(m + 1)^{\text{th}}$ order statistic.

Setting $r = m + 1$ and $n = 2m + 1$ into the formula derived earlier

$$f_{\text{med}}(y) = f_{Y_{(m+1)}}(y) = \frac{F_Y(y)^m (1 - F_Y(y))^m f_Y(y)}{B(m + 1, m + 1)}$$

If the density $f_Y(y)$ is symmetric around zero, so that $EY = 0$, then

$$F_Y(-y) = 1 - F_Y(y)$$

and so the density of the median is also symmetric around zero, so that

$$E[\text{med}(Y_1, \dots, Y_n)] = 0$$

Joint distribution of $Y_{(r)}, Y_{(s)}, r < s$

	<u>Interval</u>	<u>Prob.</u>	<u># obs = s_i</u>
1.	$(-\infty, u]$	$p_1 = F_Y(u)$	$r - 1$
2.	$(u, u + du]$	$p_2 = f_Y(u)du$	1
3.	$(u + du, v]$	$p_3 = F_Y(v) - F_Y(u + du)$	$s - r - 1$
4.	$(v, v + dv]$	$p_4 = f_Y(v)dv$	1
5.	$(v + dv, \infty)$	$p_5 = 1 - F_Y(v + dv)$	$n - s$

This is a multinomial with 5 cells:

$$f(s_1, \dots, s_5) = \frac{n!}{\prod s_i!} \prod_{i=1}^5 p_i^{s_i}$$

cont.

Taking limits as du and dv approach 0, we get

$$f_{Y_{(r)}, Y_{(s)}}(u, v) = \frac{n!}{(r-1)!(s-r-1)!(n-s)!} F_Y(u)^{r-1} \\ \times [F_Y(v) - F_Y(u)]^{s-r-1} (1 - F_Y(v))^{n-s} f_Y(u) f_Y(v)$$

Example: Suppose $F_Y(x) = x$ (Uniform)

$$f_{Y_{(r)}, Y_{(s)}}(u, v) = \frac{n!}{(r-1)!(s-r-1)!(n-s)!} \\ \times u^{r-1} (v-u)^{s-r-1} (1-v)^{n-s}$$

for $u < v$.

BIOS 660/BIOS 672 (3 Credits)

Notes 16 – 26 / 59

Joint distribution of all order statistics

Multinomial with $2n + 1$ cells, where we have one observation in each interval $[u_i, u_i + du_i)$, $i = 1, \dots, n$, and zero on the others.

$$f_{Y_{(1)}, Y_{(2)}, \dots, Y_{(n)}}(u_1, \dots, u_n) = n! \prod_{i=1}^n f_Y(u_i)$$

for $u_1 < \dots < u_n$.

Example: Suppose $F_Y(x) = x$ (Uniform)

$$f_{X_{(1)}, X_{(2)}, \dots, X_{(n)}}(u_1, \dots, u_n) = n! \quad u_1 < \dots < u_n$$

BIOS 660/BIOS 672 (3 Credits)

Notes 16 – 27 / 59

Distribution of the range

Setting $r = 1$ and $s = n$ in the joint dist. of the r^{th} and s^{th} order statistics gives the joint dist. of the *min* and *max*:

$$f_{Y_{(1)}, Y_{(n)}}(u, v) = \frac{n!}{(n-2)!} [F_Y(v) - F_Y(u)]^{n-2} f_Y(u) f_Y(v)$$

Now, do a transformation to $R = Y_{(n)} - Y_{(1)}$ and $W = Y_{(1)}$. Note that the Jacobian is 1. What is the range?

Hence,

$$f_{W,R}(w, r) = n(n-1) [F_Y(w+r) - F_Y(w)]^{n-2} f_Y(w) f_Y(w+r)$$

The density of R can be obtained by integrating out W :

$$f_R(r) = \int_{-\infty}^{\infty} f_{W,R}(w, r) dw$$

Example:

Suppose $Y \sim U[0, 1]$, i.e. $F_Y(x) = x$

$$\begin{aligned} f_R(r) &= \int_0^{1-r} n(n-1) r^{n-2} dw \\ &= n(n-1) r^{n-2} (1-r) \end{aligned}$$

Note that R has a Beta distribution.

$$\begin{aligned} E(R) &= n(n-1) \int_0^1 r \cdot r^{n-2} (1-r) dr \\ &= n(n-1) \left[\frac{1}{n} - \frac{1}{n+1} \right] \\ &= \frac{n-1}{n+1} \end{aligned}$$

What happens when $n = 2$ and $n \rightarrow \infty$?

Agenda

In the last part of the course, we discuss

- Convergence of random variables. Several different kinds
 - Convergence in probability
 - Almost sure convergence
 - Convergence in distribution
 - Convergence in L^p
 - Complete convergence
- Weak law of large numbers
- Strong law of large numbers
- Central limit theorems

The moment inequalities will be useful in proving these results.

Material is in *C-B*, Section 5.5, and *Gut*, Chapter VI.

BIOS 660/BIOS 672 (3 Credits)

Notes 16 – 31 / 59

Modes of Convergence

There are five modes of convergence. If $X_n \rightarrow X$ by any of these modes, then X is unique (see Section 2, Chapter VI of *Gut*).

1. *Convergence in Probability* $X_n \xrightarrow{P} X$

For any $\epsilon > 0$, $\lim_{n \rightarrow \infty} P\{|X_n - X| < \epsilon\} = 1$

Or equivalently,

for any $\epsilon > 0$, $\lim_{n \rightarrow \infty} P\{|X_n - X| > \epsilon\} = 0$

2. *Convergence “almost surely” (a.s.)*, denoted $X_n \xrightarrow{a.s.} X$.

Also called *Convergence with Prob. 1*

For any $\epsilon > 0$, $P\{\lim_{n \rightarrow \infty} |X_n - X| < \epsilon\} = 1$

Or

$$P\left\{\lim_{n \rightarrow \infty} X_n = X\right\} = 1$$

Or

$$P\left\{\omega : \lim_{n \rightarrow \infty} X_n(\omega) = X(\omega)\right\} = 1$$

BIOS 660/BIOS 672 (3 Credits)

Notes 16 – 32 / 59

Notes

The distinction between *Convergence almost surely* and *Convergence in probability* is subtle.

We'll see how the Markov inequality and Chebychev's inequality can often be used to establish convergence in probability. Establishing convergence a.s. is often more difficult.

Almost sure convergence is stronger than convergence in probability. Or equivalently,

$$X_n \xrightarrow{a.s.} X \Rightarrow X_n \xrightarrow{p} X.$$

Convergence in Distribution

3. *Convergence in Distribution* $X_n \xrightarrow{d} X$

Also called *convergence in law* or *weak convergence*.

$\lim_{n \rightarrow \infty} F_{X_n}(x) = F_X(x)$ for all continuity points of $F_X(x)$

Notes:

- Recall that cdfs can have at most a countable number of discontinuities.
- Theorem (no proof):

$$X_n \xrightarrow{d} X \Leftrightarrow \forall \text{ bounded continuous functions } g, \\ \mathbb{E}g(X_n) \rightarrow \mathbb{E}g(X) \text{ as } n \rightarrow \infty.$$

Convergence in Distrib. cont.

- Convergence in distribution is the weakest convergence and does not imply the other modes.
E.g. $X_n \sim N(0, 1)$ and $Y = -X$.
- An exception is the following important special case:

Suppose $X_n \xrightarrow{d} X$ where X has the degenerate distribution (i.e. $P\{X = a\} = 1$). Then,
 $X_n \xrightarrow{d} X \Rightarrow X_n \xrightarrow{P} X$.

Proof:

$$\begin{aligned} P\{|X_n - a| < \epsilon\} &= F_{X_n}(a + \epsilon) - F_{X_n}(a - \epsilon) \\ \lim_{n \rightarrow \infty} P\{|X_n - a| < \epsilon\} &= F(a + \epsilon) - F(a - \epsilon) \\ &= 1 - 0 = 1 \end{aligned}$$

Other modes of convergence

4. *Convergence in r^{th} mean* ($r \geq 1$) $X_n \xrightarrow{r} X$
If $E|X_n|^r < \infty$ for all n and

$$\lim_{n \rightarrow \infty} (E|X_n - X|^r) = 0.$$

Also called *convergence in L^r* and sometimes referred to as *convergence in the L^r norm*.
(Some books use L^p)

5. *Complete convergence* (see Chapter VI, section 4 in *Gut*), defined as,

$$\sum_{n=1}^{\infty} P(|X_n - X| > \epsilon) < \infty \quad \forall \epsilon > 0,$$

is slightly stronger than a.s. convergence, but much easier to verify.

Hence, it sometimes provides a relatively easy way to establish a.s. convergence. Some books use this as the definition of a.s. convergence.

Example 1

1. Let X_n be *gamma*(n, n) Then $X_n \xrightarrow{p} 1$.

Proof:

Since $E(X_n) = 1$ and $\text{Var}(X_n) = 1/n$, we can apply Chebychev's inequality to obtain

$$P(|X_n - 1| > \epsilon) \leq \frac{1}{n\epsilon^2} \rightarrow 0 \text{ as } n \rightarrow \infty.$$

Therefore $X_n \xrightarrow{p} 1$.

Example 2

2. Suppose $X_n \sim \text{binom}(n, \lambda/n)$. Then $X_n \xrightarrow{d} X$, where X has a *Poisson*(λ) distribution.

Proof:

$$F_{X_n}(x) = \sum_{y=0}^x \binom{n}{y} \left(\frac{\lambda}{n}\right)^y \left(1 - \frac{\lambda}{n}\right)^{n-y} \rightarrow \sum_{y=0}^x e^{-\lambda} \frac{\lambda^y}{y!},$$

as $n \rightarrow \infty$. (We saw this before). The RHS is the distribution function of the Poisson with parameter λ .

Example 3

3. Let X_2, X_3, \dots be a sequence of binary random variables defined by

$$P(X_n = 1) = 1 - \frac{1}{n} \quad \text{and} \quad P(X_n = n) = \frac{1}{n}$$

If we choose an ϵ smaller than 1, then

$$P(|X_n - 1| > \epsilon) = P(X_n = n) = 1/n \rightarrow 0,$$

hence $X_n \xrightarrow{P} 1$.

It turns out that X_n does not converge to 1 almost surely (see page 156 in *Gut*).

BIOS 660/BIOS 672 (3 Credits)

Notes 16 – 39 / 59

Example 4

4. Let X_2, X_3, \dots be a binary random variables defined by

$$P(X_n = 1) = 1 - \frac{1}{n^2} \quad \text{and} \quad P(X_n = n) = \frac{1}{n^2}$$

Now, for ϵ small enough,

$$\sum_{n=1}^{\infty} P(|X_n - 1| > \epsilon) = \sum_{n=1}^{\infty} \frac{1}{n^2}$$

which converges (the series $\sum_{n=1}^{\infty} \frac{1}{n^k}$ converges for $k > 1$). I.e., $X_n \rightarrow X$ in complete convergence, with $X = 1$. Hence $X_n \xrightarrow{a.s.} 1$.

BIOS 660/BIOS 672 (3 Credits)

Notes 16 – 40 / 59

Relationships among convergence modes

$$\begin{array}{ccc}
 X_n \xrightarrow{Compl} X \Rightarrow X_n \xrightarrow{a.s.} X & \searrow & \\
 & & X_n \xrightarrow{P} X \Rightarrow X_n \xrightarrow{d} X \\
 X_n \xrightarrow{r} X & \nearrow &
 \end{array}$$

A silly mnemonic is *All Probabilists Drink*.

Also: If $r \geq s \geq 1$

$$X_n \xrightarrow{r} X \implies X_n \xrightarrow{s} X.$$

(Try to prove this one).

BIOS 660/BIOS 672 (3 Credits)

Notes 16 – 42 / 59

In probability and distribution

Theorem: If $X_n \xrightarrow{P} X \Rightarrow X_n \xrightarrow{d} X$

Proof:

For any $\epsilon > 0$:

$$\begin{aligned}
 F_{X_n}(x) = P\{X_n \leq x\} &= P\{X_n \leq x \cap |X - X_n| \leq \epsilon\} \\
 &\quad + P\{X_n \leq x \cap |X_n - X| > \epsilon\} \\
 &\leq P\{X \leq x + \epsilon\} + P\{|X_n - X| > \epsilon\}
 \end{aligned}$$

because

$$\begin{aligned}
 \{|X - X_n| \leq \epsilon\} &= \{-\epsilon \leq X - X_n \leq \epsilon\} \\
 &\subset \{X - X_n \leq \epsilon\} = \{X \leq X_n + \epsilon\}
 \end{aligned}$$

and $\{X_n \leq x\} \cap \{X \leq X_n + \epsilon\} \subset \{X \leq x + \epsilon\}$.

BIOS 660/BIOS 672 (3 Credits)

Notes 16 – 43 / 59

cont.

Hence,

$$F_{X_n}(x) \leq F_X(x + \epsilon) + P\{|X - X_n| > \epsilon\} \quad (5)$$

and as $n \rightarrow \infty$, this implies

$$\limsup_{n \rightarrow \infty} F_{X_n}(x) \leq F_X(x + \epsilon)$$

Now, interchange the roles of X_n and X in (5) and repeat to get

$$F_X(x) \leq F_{X_n}(x + \epsilon) + P\{|X_n - X| > \epsilon\}$$

but apply inequality to $x = x - \epsilon$ instead of x , yielding

$$F_X(x - \epsilon) \leq F_{X_n}(x) + P\{|X_n - X| > \epsilon\}$$

or

$$F_X(x - \epsilon) - P\{|X_n - X| > \epsilon\} \leq F_{X_n}(x)$$

cont.

As $n \rightarrow \infty$, this implies

$$F_X(x - \epsilon) \leq \liminf_{n \rightarrow \infty} F_{X_n}(x)$$

Putting these together, we have

$$F_X(x - \epsilon) \leq \liminf_{n \rightarrow \infty} F_{X_n}(x) \leq \limsup_{n \rightarrow \infty} F_{X_n}(x) \leq F_X(x + \epsilon)$$

for all $\epsilon > 0$. Therefore, for all x where $F_X(x)$ is continuous,

$$\lim_{n \rightarrow \infty} F_{X_n}(x) = F_X(x).$$

Note: If $F_X(x)$ is not continuous at x , then all we can claim is

$$F_X(x-) \leq \liminf_{n \rightarrow \infty} F_{X_n}(x) \leq \limsup_{n \rightarrow \infty} F_{X_n}(x) \leq F_X(x)$$

r th moment and in Probability

Theorem: $X_n \xrightarrow{r} X \Rightarrow X_n \xrightarrow{P} X$

Proof: By Markov's inequality,

$$\begin{aligned} P(|X_n - X| > \epsilon) &= P(|X_n - X|^r > \epsilon^r) \\ &\leq \frac{E(|X_n - X|^r)}{\epsilon^r} \rightarrow 0 \end{aligned}$$

Example: Let Y_1, \dots, Y_n be iid with common mean μ and variance σ^2 . Let $\bar{Y}_n = \sum_{i=1}^n Y_i/n$. Then

$$E(\bar{Y}_n - \mu)^2 = \text{var}(\bar{Y}) = \frac{\sigma^2}{n} \rightarrow 0$$

Therefore $\bar{Y}_n \xrightarrow{r=2} \mu$ and so $\bar{Y}_n \xrightarrow{P} \mu$.

This result is the *weak law of large numbers*.

Convergence properties

47 / 59

Convergence in probability

Theorem: If $X_n \xrightarrow{P} X$ and $Y_n \xrightarrow{P} Y$, then

1. $X_n + Y_n \xrightarrow{P} X + Y$
2. $X_n Y_n \xrightarrow{P} XY$
3. If $g(x)$ is a continuous function, then $g(X_n) \xrightarrow{P} g(X)$

Proof:

Slutsky's Theorem

Also known as *Cramer's Theorem* - **VERY VERY USEFUL**

Theorem: If $X_n \xrightarrow{d} X$ and $Y_n \xrightarrow{P} a$, where a is a constant, then

1. $X_n + Y_n \xrightarrow{d} X + a$
2. $Y_n X_n \xrightarrow{d} aX$

Proof: Homework

BIOS 660/BIOS 672 (3 Credits)

Notes 16 – 49 / 59

Example

Let X_1, \dots, X_n be iid with mean μ , variance σ^2 , and finite moments up to fourth order. The Central Limit Theorem (CLT) says that

$$\sqrt{n} \frac{\bar{X}_n - \mu}{\sigma} \xrightarrow{d} N(0, 1)$$

But the empirical variance is a consistent estimator of σ^2 , i.e. $S_n^2 \xrightarrow{P} \sigma^2$. Therefore

$$\sqrt{n} \frac{\bar{X}_n - \mu}{S_n} \xrightarrow{d} N(0, 1)$$

by Slutsky's Theorem.

This is useful for constructing confidence intervals for μ .

BIOS 660/BIOS 672 (3 Credits)

Notes 16 – 50 / 59

Example

Show that a t -distribution with n degrees of freedom converges in distribution to the standard normal as $n \rightarrow \infty$.

Proof: Let $Y_n \sim \text{ChiSquare}(n)$. Then $Y_n/n \rightarrow 1$ by the Weak Law of Large Number (WLLN).

By Slutsky's Theorem, if $X \sim \text{Normal}(0, 1)$, then

$$\frac{X}{\sqrt{\frac{Y_n}{n}}} \xrightarrow{d} \text{Normal}(0, 1)$$

Convergence in distribution

Theorem: Suppose $X_n \xrightarrow{d} X$ and $Y_n \xrightarrow{d} Y$. Suppose further that X_n and Y_n are independent for all n , and that X and Y are independent. Then

$$X_n + Y_n \xrightarrow{d} X + Y$$

Proof: Omitted (use characteristic functions)

Example

Let $X_n \sim \text{Bin}(n_x, p_x)$ with $n_x p_x \rightarrow \lambda_x$ as $n_x \rightarrow \infty$.

Let $Y_n \sim \text{Bin}(n_y, p_y)$ with $n_y p_y \rightarrow \lambda_y$ as $n_y \rightarrow \infty$, indep. of X_n .

Then

$$X_n \xrightarrow{d} \text{Po}(\lambda_x), \quad Y_n \xrightarrow{d} \text{Po}(\lambda_y)$$

and

$$X_n + Y_n \xrightarrow{d} \text{Po}(\lambda_x + \lambda_y)$$

The Delta Method

54 / 59

Approximate mean and variance

Suppose we know the distribution of X and want to get the distribution of $Y = g(X)$. The general method is to use the Jacobian transformation. But if the distribution of X is “well concentrated” around its mean $\mu = EX$, we can approximate the mean and variance of Y as follows.

The Taylor expansion of $g(X)$ around μ is

$$g(X) = g(\mu) + g'(\mu)(X - \mu) + g''(\mu)(X - \mu)^2 + \dots$$

Therefore

$$\begin{aligned} E[g(X)] &= g(\mu) + E[g''(\mu)(X - \mu)^2] + \dots \\ &\approx g(\mu) \end{aligned}$$

Similarly

$$\begin{aligned} \text{Var}[g(X)] &\approx E[g(X) - g(\mu)]^2 = E[g'(\mu)(X - \mu)]^2 \\ &= E[g'(\mu)]^2 \text{Var} X \end{aligned}$$

Example

Let $X \sim N(\mu, \sigma^2)$ and $Y = \exp(X)$. The exact mean and variance of Y are

$$EY = e^{\mu + \sigma^2/2}, \quad \text{Var}Y = (e^{\sigma^2} - 1)e^{2\mu + \sigma^2}$$

The first order Taylor expansion gives

$$E(Y) = e^{\mu}, \quad \text{Var}(Y) = e^{2\mu}\sigma^2$$

The Delta method

Let Y_n be a sequence of rvs such that $\sqrt{n}(Y_n - \theta) \xrightarrow{d} N(0, \sigma^2)$. For a given function g and a specific value of θ , suppose $g'(\theta)$ exists and is nonzero. Then

$$\sqrt{n}[g(Y_n) - g(\theta)] \xrightarrow{d} N(0, \sigma^2[g'(\theta)]^2)$$

Proof: The Taylor expansion of $g(Y_n)$ around $Y_n = \theta$ is

$$g(Y_n) = g(\theta) + g'(\theta)(Y_n - \theta) + R_2$$

where $R_2 \rightarrow 0$ as $Y_n \rightarrow \theta$. Apply Slutsky's Theorem to

$$\sqrt{n}[g(Y_n) - g(\theta)] = g'(\theta)\sqrt{n}(Y_n - \theta)$$

Example

Let X_i iid with mean μ and variance σ^2 . The CLT gives that

$$\sqrt{n}(\bar{X}_n - \mu) \xrightarrow{d} N(0, \sigma^2)$$

Now let $g(x) = e^x$, where $g'(x) = e^x > 0$ for all x . The Delta method gives that

$$\sqrt{n}(e^{\bar{X}_n} - e^\mu) \xrightarrow{d} N(0, \sigma^2 e^{2\mu})$$

Let $Y_i = \exp(X_i)$, then $e^{\bar{X}_n}$ is the geometric average of the Y_i . So we have an approximation for the distribution of the geometric average.

Second-order Delta method

Let Y_n be a sequence of rvs such that $\sqrt{n}(Y_n - \theta) \xrightarrow{d} N(0, \sigma^2)$. For a given function g and a specific value of θ , suppose $g'(\theta) = 0$ and $g''(\theta)$ exists and is nonzero. Then

$$n[g(Y_n) - g(\theta)] \xrightarrow{d} \frac{\sigma^2 g''(\theta)}{2} \chi_1^2$$

Proof: See C&B.