

# Math 3070/6070 Introduction to Probability

Mon/Wed/Fri 9:00am - 9:50am

Instructor: Dr. Xiang Ji, xji4@tulane.edu

## Lecture 1: Aug 21

### Today

- Introduction
- Introduce yourself
- Course logistics

### What is this course about?

This course will provide a calculus-based introduction to probability theory. Material covered will include fundamental axioms of probability, combinatorics, discrete and continuous random variables, multivariate distributions, expectation, and limit theorems, generally following Chapters 1-5 of the textbook. This course is a critical prerequisite for more advanced work in statistical theory and analysis.

### Prerequisite

- Calculus

### Why learn probability

- The subject of probability theory is the foundation upon which all of statistics is built.
- It provides you a tool to model
  - populations
  - experiments
  - almost anything else that could be considered a random phenomenon
  - example topics in [Data Analysis course](#)
- Through these models, statisticians are able to draw inferences about populations based on examination of only a part of the whole.
- A must have for any Data Scientists.

## What this course WILL NOT do for you

It will not help you:

- Beat the casino at blackjack (although it may convince you that it is better not to gamble, or that a casino is a great business).
- Answer your friends' silly questions such as "What are the chances it will rain tomorrow?" (although it might make you think of ways that you might model and compute it).

## Syllabus

Check course website frequently for updates and announcements.

<https://tulane-math-3070-2023.github.io/>

## HW submission

Students are required to submit hand-written homework in recitations to the TA. Homework assignments are expected every two weeks with 4-5 problems at a time.

## Last year comments

Your experience in this course:

- Exactly what I was looking for in a probability course - I got a really good grasp on the theory and this math has already come in useful in other areas and fields that I'm studying. Glad I took this class, and I appreciate the tests being more accessible and spaced out to provide less pressure - highly recommend.
- I really appreciate the lecture shift that Prof. Xiang Ji had after our midterm survey. He took suggestions seriously and dramatically improved how the content was presented to our class's needs. Having the opportunity to present and listen to classmates on related statistical topics was also fun and rewarding.
- Absolutely stupendous course. Fabulous structure, even more fabulous professor.
- Lectures were pretty disengaging. I would've rather had lecture notes written out to us rather than being read to us. Presentation extra credit opportunities were nice. Would've liked more communication and collaboration between TA and teacher.
- Once the lecture structure changed after midterms, I felt I learned much more during lectures. I think that the concepts and theories were explained clearly in class. I appreciate the generosity Professor Ji showed when grading exams, but I think that receiving more detailed feedback would have helped me learn the material better. The exams felt more like a test of our ability to make a formula sheet than a test of our understanding of course material. While this was nice in terms of my grade, I don't think this helped me with my understanding of the material. The recitations often

felt disconnected from the course material because we did practice with numerical applications and calculations instead of theory. Overall I feel prepared for the second half of this course next semester.

- I felt the course provided me with very little context for why we were learning about the things we were learning. For example, now I know about a lot of different distributions and their moment generating functions, but I have no idea when I might need to use a Gamma distribution or Poisson. Additionally, I did not find the textbook to be a particularly helpful resource. There also seemed to be little communication between the professor and TA, so the activities we did in recitation weren't always relevant to what we were doing in the class.
- This course and professor were great. I got introduced to an entirely new facet of mathematics and it excited me for my major. Professor X was great and fostered a great learning environment in the classroom.
- I appreciate the professor's willingness to adjust his teaching style to include more written-out derivations. I would have preferred not to have the 10-minute presentations in class on Fridays. I felt like especially at the beginning of the semester most of the content in the presentations was more complex than what we had learned and took class time away from the material. I think it would have been helpful for the TA and Professor to communicate about the level of instruction of the course. It was hard to understand where I stood from a knowledge perspective when the level of difficulty ranged so greatly from homework, quizzes, worksheets, and lectures. I understand that this class is a probability theory class, however, I think it would be helpful for math majors interested in applied math to have some way to learn more about the applications of probability since the Stats for scientists does not count toward the major.
- Wonderful experience! I learned a tremendous amount, and always left class wanting to know more. This course did not shy away from difficulty and I could not have appreciated it more. Far too often professors dumb down the material in order to cater to the students. Not this class, we learned intense probability theory and I could not be happier with it. I look forward to continuing my studies next semester with Statistical Inference. The rigor and complexity of the course demanded respect and, if given, the knowledge learned is powerful.
- Professor was good, but I think the subject matter was oftentimes too confusing
- The only reasons why I personally did not like the course is firstly because I do not like the subject matter and secondly, I do not like how the material was organized. I do not learn math the best when it is purely based off lecture notes. I did notice that the professor was trying to write more on the board during lectures which I did appreciate.
- I appreciate that Professor Ji mid-semester began to workout problems in class on the whiteboard as that kept me more engaged. Sometimes it was hard to follow the work on the board however as the steps didn't seem to be organized (they looked like they were written all over the place). I think it would be more easy to follow if the notes

on the whiteboard were more organized linearly. Also, we only did this like once, but I think it could be a good idea to also give students problems and have them come up to the whiteboard and solve them (maybe for bonus points, doesn't have to be) as that is another way to make class more interactive. I also liked having a practice test as it always nice to get more practice.

- Class sessions were a little slow paced and repetitive for me. However, towards the end as the Professor took some feedback from students and started writing out equations on the board it was easier to follow along and I was more actively learning. I think the weekly presentations were nice but took away a chunk of class time that could've prevented us from getting behind. The weekly reviews were really helpful and I think those are great for Fridays. Our Professor, Dr. Ji, was super accommodating and very adaptive towards creating the best learning environment for students. I appreciated the mid-semester surveys and his changes based off that immensely, and his efforts to supplement our grades with bonus presentations (I just think they could be shorter or maybe uploaded on a discussion post, rather than take up 1 of every 3 classes). I liked that we had a class notes document, rather than a textbook, and this also made the class very accessible when I couldn't be there in person.
- Teaching improved significantly through the semester, he was open to feedback so that helped some. Lectures were still almost entirely him reading off of a pdf though, which did not teach me much

#### strongest aspects of this course

- Very good in-depth class, useful theory knowledge and strong lectures.
- The strongest aspect of this course is how it provides a very good background which is needed for future statistics courses.
- Loved the professor, absolutely no complaints. Promote Xiang Ji
- Lecture notes helpful. Loved the TA and our lab sessions.
- I really appreciated that Professor Ji asked for student feedback in the middle of the semester and adjusted the lectures based on the results. Once we started doing more derivations on the board, I was able to understand them better.
- I love you as a person, but it is seriously hard to digest anything you say in class. Reading directly from the textbook is not teaching, it is just reading. The tests do not test us on our knowledge of the material AT ALL. They simply did you write down the right thing on your cheat sheet.
- The professor and TA were very flexible when it was clear that the class didn't understand something, and were always open for feedback.
- I really enjoyed having the lecture notes typed and uploaded ahead of class. This allowed me to add to these notes and not have to get a notetaker. I also enjoyed

having a notes sheet for the exams especially since there are so many formulas and distributions.

- The rigor and complexity of the course is the strongest aspect. Both Professor Zhao and Professor Ji made the material approachable and were always happy to explain and explain again the difficult concepts and proofs. I am excited to take Statistical Inference for the following semester. Additionally, the course has changed my way of thinking when assessing probabilities in all aspects of life, and there have been many instances over the semester when the knowledge imparted to me has been of service.
- The professor made himself available which I noticed and appreciated. I liked him as a person a lot and I noticed his deep knowledge on the topic.
- Switching to writing on the whiteboard vs pure lecture from typed notes. Enjoyed our bonus presentations on various math topics as it opened my eyes to how versatile statistics can be.
- I listed those above. But the professors attitude and flexibility, and use of technology to sum it up.
- professor was open to student feedback which was helpful
- I liked the presentations a lot. They were a good change of pace and way to understand how this is all applied

#### RateMyProfessor

- (4.0 Quality / 2.0 Difficulty) Dr Ji has a dry wit and is receptive to student feedback. He is a generous grader and offered an opportunity for generous extra credit. For the tests, he allowed a cheat sheet, and the final was take-home. I also think the tests were easy compared to how complicated they could have been. Beware the class is super theory-based similar to analysis.
- (4.0 Quality / 3.0 Difficulty) At the start of the semester I struggled with Dr. Ji's lecturing style, but after the midterm he asked for our feedback and made adjustments to his class so it was easier to follow his lessons. He didn't always explain things in great detail the first time, but if you ask questions he is always willing to clarify. Exams are also graded generously.
- (3.0 Quality / 3.0 Difficulty) Lectures are based off a pdf document which helps when you need to study, but can be terribly difficult to pay attention to in class. Tests account for about 65% of your grade but he goes pretty easy on the grading.
- (2.0 Quality / 2.0 Difficulty) Don't take this class if you're actually trying to learn the class content. I've been in here a whole semester and genuinely cannot tell you one thing I have retained. Does not communicate with the TA so recitation is not helpful either. Although, homework is graded for completion and quizzes are easy so at least its not a hard grade.

- (4.0 Quality / 5.0 Difficulty) Prof. Xi is hardcore. 3070 is definitely a theory-heavy class for people who really want to get into the underlying technical parts of probability, but if you go into it with that mindset it's really well structured and informative. Book is useful, but you have to be serious and commit time/effort into this class to do well.

## Lecture 2:Aug 23

### Last time

- Introduction
- Introduce yourself
- Course logistics

### Today

- Set theory (1.1)
- Axiomatic Foundations (1.2)

### Set Theory

One of the main objectives of a statistician is to draw conclusions about a population of objects by conducting an experiment. The first step in this endeavor is to identify the possible outcomes or, in statistical terminology, the *sample space*.

**Definition** The set,  $S$ , of all possible outcomes of a particular experiment is called the *sample space* for the experiment.

**Example** The sample space of

- tossing a coin just once, contains two outcomes, heads and tails

$$S = \{H, T\}$$

- observing reported SAT scores of randomly selected students at a certain university

$$S = \{200, 210, 220, \dots, 780, 790, 800\}$$

- an experiment where the observation is reaction time to a certain stimulus

$$S = (0, \infty)$$

**Definition** An *event* is any collection of possible outcomes of an experiment, that is, any subset of  $S$  (including  $S$  itself).

Let  $A$  be an event,

- $A$  is a subset of  $S$ ,
- event  $A$  occurs if the outcome of the experiment is in the set  $A$ ,
- we generally speak of the probability of an event, rather than a set.

Set operations:

- Containment:

$$A \subset B \iff x \in A \implies x \in B$$

- Equality:

$$A = B \iff A \subset B \text{ and } B \subset A$$

- Union: the union of  $A$  and  $B$ , written as  $A \cup B$ , is the set of elements that belong to either  $A$  or  $B$  or both

$$A \cup B = \{x : x \in A \text{ or } x \in B\}.$$

- Intersection: the intersection of  $A$  and  $B$ , written  $A \cap B$ , is the set of elements that belong to both  $A$  and  $B$ :

$$A \cap B = \{x : x \in A \text{ and } x \in B\}.$$

- Complementation: the complement of  $A$ , written  $A^c$ , is the set of all elements that are not in  $A$ :

$$A^c = \{x : x \notin A\}.$$



## Lecture 3: Aug 25

### Last time

- Set theory (1.1)
- Axiomatic Foundations (1.2)

### Today

- Axiomatic Foundations (1.2)
- Calculus of Probabilities (1.2)
- Conditional Probability (1.3)

**Theorem** For any three events,  $A$ ,  $B$ , and  $C$ , defined on a sample space  $S$ ,

1. Commutativity

$$A \cup B = B \cup A,$$
$$A \cap B = B \cap A;$$

2. Associativity

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

3. Distributive Laws

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

4. DeMorgan's Laws

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

We show the proof of  $A \cap (B \cup C) = (A \cap B) \cup (A \cap C)$  in the distributive laws. Caution: Venn diagrams are helpful in visualization, but they do not constitute a formal proof. To prove that two sets are equal, we need to show that each set contains the other.

*proof:*

- $A \cap (B \cup C) \subset (A \cap B) \cup (A \cap C)$ :  
Let  $x \in (A \cap (B \cup C))$ . By definition of intersection,  $x \in (B \cup C)$  that is, either  $x \in B$  or  $x \in C$ . Since  $x$  also must be in  $A$ , we have that either  $x \in (A \cap B)$  or  $x \in (A \cap C)$ ; therefore,  $x \in ((A \cap B) \cup (A \cap C))$ .
- $(A \cap B) \cup (A \cap C) \subset A \cap (B \cup C)$ :  
Let  $x \in ((A \cap B) \cup (A \cap C))$ . This implies that  $x \in (A \cap B)$  or  $x \in (A \cap C)$ . If  $x \in (A \cap B)$ , then  $x$  is in both  $A$  and  $B$ . Since  $x \in B$ , then  $x \in (B \cup C)$  and thus

$x \in (A \cap (B \cup C))$ . It follows the same argument when  $x \in (A \cap C)$ , we still have  $x \in (A \cap (B \cup C))$ .

**Definition** Two events  $A$  and  $B$  are *disjoint* (or *mutually exclusive*) if  $A \cap B = \emptyset$ . The events  $A_1, A_2, \dots$  are *pairwise disjoint* (or *mutually exclusive*) if  $A_i \cap A_j = \emptyset$  for all  $i \neq j$ .

**Definition** If  $A_1, A_2, \dots$  are pairwise disjoint and  $\cup_{i=1}^{\infty} A_i = A_1 \cup A_2 \cup \dots = S$ , then the collection of  $A_1, A_2, \dots$  forms a *partition* of  $S$ .

**Example** The sets  $A_i = [i, i + 1), i = 0, 1, 2, \dots$  form a partition of  $[0, \infty)$ .

## Basics of Probability Theory

When an experiment is performed, the realization of the experiment is an outcome in the sample space. If the experiment is performed a number of times, then

- different outcomes may occur each time
- some outcomes may repeat
- the “frequency of occurrence” of an outcome can be thought of as a probability

However, we **do not** define probabilities in terms of frequencies but instead take the mathematically simpler axiomatic approach. The axiomatic approach is not concerned with the interpretations of probabilities, but is concerned only that the probabilities are defined by a function satisfying the axioms. Interpretations of the probabilities are quite another matter:

- The “frequency of occurrence” of an event is one example of a particular interpretation of probability.
- Another possible interpretation is a subjective one, where we can think of the probability as a belief in the chance of an event occurring.

## Lecture 4: Aug 28

### Last time

- Set theory (1.1)

### Today

- Remember to go over DeMorgan's laws
- Axiomatic Foundations (1.2)
- Calculus of Probabilities (1.2)
- Conditional Probability (1.3)

**Theorem** For any three events,  $A$ ,  $B$ , and  $C$ , defined on a sample space  $S$ ,

1. Commutativity

$$\begin{aligned}A \cup B &= B \cup A, \\ A \cap B &= B \cap A;\end{aligned}$$

2. Associativity

$$\begin{aligned}A \cup (B \cup C) &= (A \cup B) \cup C, \\ A \cap (B \cap C) &= (A \cap B) \cap C;\end{aligned}$$

3. Distributive Laws

$$\begin{aligned}A \cap (B \cup C) &= (A \cap B) \cup (A \cap C), \\ A \cup (B \cap C) &= (A \cup B) \cap (A \cup C);\end{aligned}$$

4. DeMorgan's Laws

$$\begin{aligned}(A \cup B)^c &= A^c \cap B^c, \\ (A \cap B)^c &= A^c \cup B^c;\end{aligned}$$

We show the proof of  $A \cap (B \cup C) = (A \cap B) \cup (A \cap C)$  in the distributive laws. Caution: Venn diagrams are helpful in visualization, but they do not constitute a formal proof. To prove that two sets are equal, we need to show that each set contains the other.

*proof:*

- $A \cap (B \cup C) \subset (A \cap B) \cup (A \cap C)$ :  
Let  $x \in (A \cap (B \cup C))$ . By definition of intersection,  $x \in (B \cup C)$  that is, either  $x \in B$  or  $x \in C$ . Since  $x$  also must be in  $A$ , we have that either  $x \in (A \cap B)$  or  $x \in (A \cap C)$ ; therefore,  $x \in ((A \cap B) \cup (A \cap C))$ .
- $(A \cap B) \cup (A \cap C) \subset A \cap (B \cup C)$ :  
Let  $x \in ((A \cap B) \cup (A \cap C))$ . This implies that  $x \in (A \cap B)$  or  $x \in (A \cap C)$ . If  $x \in (A \cap B)$ , then  $x$  is in both  $A$  and  $B$ . Since  $x \in B$ , then  $x \in (B \cup C)$  and thus

$x \in (A \cap (B \cup C))$ . It follows the same argument when  $x \in (A \cap C)$ , we still have  $x \in (A \cap (B \cup C))$ .

**Definition** Two events  $A$  and  $B$  are *disjoint* (or *mutually exclusive*) if  $A \cap B = \emptyset$ . The events  $A_1, A_2, \dots$  are *pairwise disjoint* (or *mutually exclusive*) if  $A_i \cap A_j = \emptyset$  for all  $i \neq j$ .

**Definition** If  $A_1, A_2, \dots$  are pairwise disjoint and  $\cup_{i=1}^{\infty} A_i = A_1 \cup A_2 \cup \dots = S$ , then the collection of  $A_1, A_2, \dots$  forms a *partition* of  $S$ .

**Example** The sets  $A_i = [i, i + 1), i = 0, 1, 2, \dots$  form a partition of  $[0, \infty)$ .

## Basics of Probability Theory

When an experiment is performed, the realization of the experiment is an outcome in the sample space. If the experiment is performed a number of times, then

- different outcomes may occur each time
- some outcomes may repeat
- the “frequency of occurrence” of an outcome can be thought of as a probability

However, we **do not** define probabilities in terms of frequencies but instead take the mathematically simpler axiomatic approach. The axiomatic approach is not concerned with the interpretations of probabilities, but is concerned only that the probabilities are defined by a function satisfying the axioms. Interpretations of the probabilities are quite another matter:

- The “frequency of occurrence” of an event is one example of a particular interpretation of probability.
- Another possible interpretation is a subjective one, where we can think of the probability as a belief in the chance of an event occurring.

## Axiomatic Foundations

For each event  $A$  in the sample space  $S$ , we want to associate with  $A$  a number between zero and one that will be called the probability of  $A$ , denoted by  $\Pr(A)$ . The domain of  $\Pr$  is the set where the arguments of the function  $\Pr(\cdot)$  are defined. It is natural to define the domain of  $\Pr$  as all subsets of  $S$ , that is for each  $A \subset S$ , we define  $\Pr(A)$  as the probability that  $A$  occurs. However, there are some technical difficulties to overcome which requires us to familiarize with the following.

**Definition** A collection of subsets of  $S$  is called a *sigma algebra* (or *Borel field*), denoted by  $\mathcal{B}$ , if it satisfies the following three properties:

1.  $\emptyset \in \mathcal{B}$  (the empty set is an element of  $\mathcal{B}$ ).
2. If  $A \in \mathcal{B}$ , then  $A^c \in \mathcal{B}$  ( $\mathcal{B}$  is closed under complementation).

3. If  $A_1, A_2, \dots \in \mathcal{B}$ , then  $\cup_{i=1}^{\infty} A_i \in \mathcal{B}$  ( $\mathcal{B}$  is closed under countable unions).

From Property (1) and (2), we see that the empty set and its complement  $S$  (since  $S = \emptyset^c$ ) are always in a sigma algebra. In fact, they construct the *trivial* algebra  $\{\emptyset, S\}$  which is the smallest sigma algebra.

By DeMorgan's Law, (3) can be replaced by:

$$3'. \text{ if } A_1, A_2, \dots \in \mathcal{B}, \text{ then } \cap_{i=1}^{\infty} A_i \in \mathcal{B}.$$

This is because:

$$(\cup_{i=1}^{\infty} A_i^c)^c = \cap_{i=1}^{\infty} A_i.$$

**Example** If  $S$  is finite or countable (where the elements of  $S$  can be put into 1 – 1 correspondence with a subset of the integers), then these technicalities really do not arise, for we define for a given sample space  $S$ ,

$$\mathcal{B} = \{\text{all subsets of } S, \text{ including } S \text{ itself}\}.$$

If  $S$  has  $n$  elements, there are  $2^n$  sets in  $\mathcal{B}$  (why?). [hint: for each element, it is either in or out of a subset, so 2 choices].

**Example** Let  $S = (-\infty, \infty)$ , the real line. Then  $\mathcal{B}$  is chosen to contain all sets of the form

$$[a, b], (a, b], (a, b), \text{ and } [a, b)$$

for all real numbers  $a$  and  $b$ . Also, from the properties of  $\mathcal{B}$ , it follows that  $\mathcal{B}$  contains all sets that can be formed by taking (possibly countably infinite) unions and intersections of sets of the above varieties.

We now define a probability function.

**Definition** Given a sample space  $S$  and an associated sigma algebra  $\mathcal{B}$ , a *probability function* is a function  $\Pr$  with domain  $\mathcal{B}$  that satisfies

1.  $\Pr(A) \geq 0$  for all  $A \in \mathcal{B}$ .
2.  $\Pr(S) = 1$ .
3. If  $A_1, A_2, \dots \in \mathcal{B}$  are pairwise disjoint, then  $\Pr(\cup_{i=1}^{\infty} A_i) = \sum_{i=1}^{\infty} \Pr(A_i)$ .

The above three properties are usually referred to as the Axioms of Probability (or the Kolmogorov Axioms, after A. Kolmogorov, one of the fathers of probability theory). Any function that satisfies the Axioms of Probability is called a probability function.

**Example** Consider the simple experiment of tossing a fair coin (just once), so  $S = \{H, T\}$ . A reasonable probability function is the one that assigns equal probabilities to heads and tails, that is,

$$\Pr(\{H\}) = \Pr(\{T\}).$$

Since  $S = \{H\} \cup \{T\}$ , we have, from Axiom 2,  $\Pr(\{H\} \cup \{T\}) = 1$ . Also,  $\{H\}$  and  $\{T\}$  are disjoint, so  $\Pr(\{H\} \cup \{T\}) = \Pr(\{H\}) + \Pr(\{T\})$ . Collectively, we have

$$\begin{aligned}\Pr(\{H\}) &= \Pr(\{T\}) \\ \Pr(\{H\} \cup \{T\}) &= 1 \\ \Pr(\{H\} \cup \{T\}) &= \Pr(\{H\}) + \Pr(\{T\})\end{aligned}$$

Therefore,  $\Pr(\{H\}) = \Pr(\{T\}) = \frac{1}{2}$ .

**Example** If  $S$  is finite or countable (where the elements of  $S$  can be put into 1 – 1 correspondence with a subset of the integers), then these technicalities really do not arise, for we define for a given sample space  $S$ ,

$$\mathcal{B} = \{\text{all subsets of } S, \text{ including } S \text{ itself}\}.$$

If  $S$  has  $n$  elements, there are  $2^n$  sets in  $\mathcal{B}$  (why?). [hint: for each element, it is either in or out of a subset, so 2 choices].

**Example** Let  $S = (-\infty, \infty)$ , the real line. Then  $\mathcal{B}$  is chosen to contain all sets of the form

$$[a, b], (a, b], (a, b), \text{ and } [a, b)$$

for all real numbers  $a$  and  $b$ . Also, from the properties of  $\mathcal{B}$ , it follows that  $\mathcal{B}$  contains all sets that can be formed by taking (possibly countably infinite) unions and intersections of sets of the above varieties.

We now define a probability function.

**Definition** Given a sample space  $S$  and an associated sigma algebra  $\mathcal{B}$ , a *probability function* is a function  $\Pr$  with domain  $\mathcal{B}$  that satisfies

1.  $\Pr(A) \geq 0$  for all  $A \in \mathcal{B}$ .
2.  $\Pr(S) = 1$ .
3. If  $A_1, A_2, \dots \in \mathcal{B}$  are pairwise disjoint, then  $\Pr(\cup_{i=1}^{\infty} A_i) = \sum_{i=1}^{\infty} \Pr(A_i)$ .

The above three properties are usually referred to as the Axioms of Probability (or the Kolmogorov Axioms, after A. Kolmogorov, one of the fathers of probability theory). Any function that satisfies the Axioms of Probability is called a probability function.

**Example** Consider the simple experiment of tossing a fair coin (just once), so  $S = \{H, T\}$ . A reasonable probability function is the one that assigns equal probabilities to heads and tails, that is,

$$\Pr(\{H\}) = \Pr(\{T\}).$$

Since  $S = \{H\} \cup \{T\}$ , we have, from Axiom 1,  $\Pr(\{H\} \cup \{T\}) = 1$ . Also,  $\{H\}$  and  $\{T\}$  are disjoint, so  $\Pr(\{H\} \cup \{T\}) = \Pr(\{H\}) + \Pr(\{T\})$ . Collectively, we have

$$\begin{aligned}\Pr(\{H\}) &= \Pr(\{T\}) \\ \Pr(\{H\} \cup \{T\}) &= 1 \\ \Pr(\{H\} \cup \{T\}) &= \Pr(\{H\}) + \Pr(\{T\})\end{aligned}$$

Therefore,  $\Pr(\{H\}) = \Pr(\{T\}) = \frac{1}{2}$ .

## Caculus of Probabilities

We start with some fairly self-evident properties of the probability function when applied to a single event.

**Theorem** If  $\Pr$  is a probability function and  $A$  is any set in  $\mathcal{B}$ , then

1.  $\Pr(\emptyset) = 0$ , where  $\emptyset$  is the empty set;
2.  $\Pr(A) \leq 1$ ;
3.  $\Pr(A^c) = 1 - \Pr(A)$ .

*proof:*

- It's easy to prove (3) first. Since
  - $\Pr(A \cup A^c) = \Pr(S) = 1$ ,
  - $A$  and  $A^c$  are disjoint, by axiom (3),  $\Pr(A \cup A^c) = \Pr(A) + \Pr(A^c)$ .
 so that  $\Pr(A) + \Pr(A^c) = \Pr(S) = 1$
- with (3) proved, (1) is simple. because we know that
  - $S \cup \emptyset = S$ ,
  - $S \cap \emptyset = \emptyset$ , they are disjoint,
 so that  $\Pr(\emptyset) + \Pr(S) = \Pr(\emptyset \cup S) = \Pr(S)$ .
- now for (2),  $\Pr(A) = 1 - \Pr(A^c) \leq 1$ , by axiom (1).

## Lecture 5: Aug 30

Last time

- Set theory (1.1)

Today

- Calculus of Probabilities (1.2)
- Conditional Probability (1.3)

### Calculus of Probabilities

**Theorem** If  $\Pr$  is a probability function and  $A$  and  $B$  are any sets in  $\mathcal{B}$ , then

1.  $\Pr(B \cap A^c) = \Pr(B) - \Pr(A \cap B)$ ;
2.  $\Pr(A \cup B) = \Pr(A) + \Pr(B) - \Pr(A \cap B)$ ;
3. If  $A \subset B$ , then  $\Pr(A) \leq \Pr(B)$ .

*proof:*

1. For (1), we have  $B = \{B \cap A\} \cup \{B \cap A^c\}$  and  $\{B \cap A\} \cap \{B \cap A^c\} = \emptyset$ , therefore

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

2. For (2), we plug in (1) first such that we only need to show  $\Pr(A \cup B) = \Pr(A) + \Pr(B \cap A^c)$ . Since  $A \cap \{B \cap A^c\} = \emptyset$  and  $A \cup B = A \cup \{B \cap A^c\}$  (use a Venn diagram, or see Exercise 1.2), we have  $\Pr(A \cup B) = \Pr(A) + \Pr(B \cap A^c)$ .
3. For (3), if  $A \subset B$ , then  $A \cap B = A$ . Then using (1), we have

$$0 \leq \Pr(B \cap A^c) = \Pr(B) - \Pr(A)$$

Formula (2) in the above theorem gives a useful inequality for the probability of an intersection (Bonferroni's Inequality):

$$\Pr(A \cap B) \geq \Pr(A) + \Pr(B) - 1.$$



## Lecture 6: Sept 1

Last time

- Calculus of Probabilities (1.2)

Today

- no class next Monday (Labor day)
- Conditional Probability (1.3)
- Independence (1.3)

**Theorem** If  $\Pr$  is a probability function, then

1.  $\Pr(A) = \sum_{i=1}^{\infty} \Pr(A \cap C_i)$  for any partition  $C_1, C_2, \dots$ ;
2.  $\Pr(\cup_{i=1}^{\infty} A_i) \leq \sum_{i=1}^{\infty} \Pr(A_i)$  for any sets  $A_1, A_2, \dots$

where (1) is also referred to as “Total probability” and (2) is Boole’s inequality.

*proof:*

By definition, since  $C_1, C_2, \dots$  form a partition, we have  $C_i \cap C_j = \emptyset$  for all  $i \neq j$ , and  $S = \cup_{i=1}^{\infty} C_i$ . Therefore,

$$A = A \cap S = A \cap (\cup_{i=1}^{\infty} C_i) = \cup_{i=1}^{\infty} (A \cap C_i),$$

where the last equality follows from the Distributive Law. Since  $\{A \cap C_i\} \cap \{A \cap C_j\} = \emptyset$  (i.e.  $A \cap C_i$  and  $A \cap C_j$  are disjoint), we have

$$\Pr(A) = \Pr(\cup_{i=1}^{\infty} (A \cap C_i)) = \sum_{i=1}^{\infty} \Pr(A \cap C_i).$$

To establish Boole’s Inequality, we first construct a disjoint collection  $A_1^*, A_2^*, \dots$ , with the property that  $\cup_{i=1}^{\infty} A_i^* = \cup_{i=1}^{\infty} A_i$ . We define  $A_i^*$  by

$$A_1^* = A_1, \quad A_i^* = A_i \setminus (\cup_{j=1}^{i-1} A_j), \quad i = 2, 3, \dots,$$

where the notation  $A \setminus B$  denotes the part of  $A$  that does not intersect with  $B$ . In other words,  $A \setminus B = A \cap B^c$ . It’s easy to see that  $\cup_{i=1}^{\infty} A_i^* = \cup_{i=1}^{\infty} A_i$ , and we have

$$\Pr(\cup_{i=1}^{\infty} A_i) = \Pr(\cup_{i=1}^{\infty} A_i^*) = \sum_{i=1}^{\infty} \Pr(A_i^*)$$

where the last equality holds because  $A_i^*$  are disjoint. To see this, consider any pair of  $A_i^* \cap A_k^*, i > k$ , then

$$\begin{aligned} A_i^* \cap A_k^* &= \{A_i \setminus (\cup_{j=1}^{i-1} A_j)\} \cap \{A_k \setminus (\cup_{j=1}^{k-1} A_j)\} \\ &= \{A_i \cap (\cup_{j=1}^{i-1} A_j)^c\} \cap \{A_k \cap (\cup_{j=1}^{k-1} A_j)^c\} \\ &= \{A_i \cap (\cap_{j=1}^{i-1} A_j^c)\} \cap \{A_k \cap (\cap_{j=1}^{k-1} A_j^c)\} \\ &= \emptyset. \end{aligned}$$

Lastly, we have  $\Pr(A_i^*) \leq \Pr(A_i)$ .

## Conditional Probability

All of the probabilities that we have dealt with thus far have been unconditional probabilities. A sample space was defined and all probabilities were calculated with respect to that sample space. In many instances, however, we are in a position to update the sample space based on new information. In such cases we want to be able to update probability calculations or to calculate *conditional probabilities*.

**Definition** If  $A$  and  $B$  are events in  $S$ , and  $\Pr(B) > 0$ , then the *conditional probability* of  $A$  given  $B$ , written  $\Pr(A|B)$ , is

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

Note that  $B$  becomes the sample space now:  $\Pr(B|B) = 1$ .

**Example** Four cards are dealt from the top of a well-shuffled deck. What is the probability that they are the four aces? What is the probability of getting four aces at the top if knowing the first card is an ace? (there are in total 52 cards)

*solution:*

We define two events first. Let  $A$  be the event {4 aces on top}, and  $B$  be the event {the first card on top is an ace}. For a well-shuffled deck, all groups of 4 cards are equally likely. For the 4 aces on top, we have  $4!$  ways of ordering (i.e., permutations of 4 distinct elements where order matters). For the rest of  $52 - 4 = 48$  cards, there are  $48!$  permutations (where again order matters). Therefore, the probability of event  $A$  is  $\Pr(A) = \frac{4!48!}{52!} = \frac{4!}{52 \cdot 51 \cdot 50 \cdot 49} = \frac{1}{270,725}$ .

Note,  $\binom{n}{m}$  reads “from  $n$  choose  $m$ ” (for  $m \leq n$ ) and calculates by  $\binom{n}{m} = \frac{n!}{(n-m)!m!}$  that gives the number of distinct combinations of choosing  $m$  elements from  $n$  total elements.

Now, let's calculate  $\Pr(A|B)$ . First of all,  $A \subset B$ , so that we have  $\Pr(A \cap B) = \Pr(A)$ . For  $\Pr(B)$ , having an ace on top instead of the other 12 kinds,  $\Pr(B) = \frac{1}{13}$ . Then  $\Pr(A|B) = \frac{\Pr(A \cap B)}{\Pr(B)} = \frac{\Pr(A)}{\Pr(B)} = \frac{1}{20,825}$ .

**Theorem** (Bayes' Rule) Let  $A_1, A_2, \dots$  be a partition of the sample space, and let  $B$  be any set. Then, for each  $i = 1, 2, \dots$ ,

$$\Pr(A_i|B) = \frac{\Pr(B|A_i) \Pr(A_i)}{\sum_{j=1}^{\infty} \Pr(B|A_j) \Pr(A_j)}.$$

*proof:*

By “Total probability”, we have  $\Pr(B) = \sum_{j=1}^{\infty} \Pr(B \cap A_j)$  which is the denominator. Therefore,  $\Pr(A_i|B) = \frac{\Pr(A_i \cap B)}{\Pr(B)} = \frac{\Pr(B|A_i) \Pr(A_i)}{\sum_{j=1}^{\infty} \Pr(B \cap A_j)}$ .

## Lecture 7: Sept 6

### Last time

- Conditional Probability (1.3)
- Independence (1.3)

### Today

- Independence (1.3)
- Random variables
- Distribution Functions

## Independence

**Definition** Two events,  $A$  and  $B$ , are *statistically independent* if

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

Note that independence could have been defined using Bayes' rule by  $\Pr(A|B) = \Pr(A)$  or  $\Pr(B|A) = \Pr(B)$  as long as  $\Pr(A) > 0$  or  $\Pr(B) > 0$ . More notation, often statisticians omit  $\cap$  when writing intersection in a probability function which means  $\Pr(AB) = \Pr(A \cap B)$ . Sometime, statisticians use comma  $(,)$  to replace  $\cap$  inside a probability function too,  $\Pr(A, B) = \Pr(A \cap B)$ .

**Theorem** If  $A$  and  $B$  are independent events, then the following pairs are also independent.

1.  $A$  and  $B^c$ ,
2.  $A^c$  and  $B$ ,
3.  $A^c$  and  $B^c$ .

*proof:*

For (1),

$$\begin{aligned}\Pr(A, B^c) &= \Pr(A) - \Pr(A, B) \\ &= \Pr(A) - \Pr(A) \Pr(B) \\ &= \Pr(A)(1 - \Pr(B)) \\ &= \Pr(A) \Pr(B^c)\end{aligned}$$

For (2), we just need to switch  $A$  and  $B$ .

For (3), we have  $A^c$  and  $B$  are independent, then we can treat  $A^c$  as  $A'$  and  $B$  as  $B'$ , then  $A'$  and  $B'^c$  are independent which is  $A^c$  and  $B^c$  are independent.

Alternatively, for (2),

$$\begin{aligned}\Pr(A^c, B) &= \Pr(A^c|B) \Pr(B) \\ &= [1 - \Pr(A|B)] \Pr(B) \\ &= [1 - \Pr(A)] \Pr(B) \\ &= \Pr(A^c) \Pr(B).\end{aligned}$$

And for (3),

$$\begin{aligned}\Pr(A^c, B^c) &= \Pr(A^c) - \Pr(A^c, B) \\ &= \Pr(A^c) - \Pr(A^c) \Pr(B) \\ &= \Pr(A^c) \Pr(B^c).\end{aligned}$$

**Example** Let the sample space  $S$  consist of the  $3!$  permutations of the letters  $a$ ,  $b$ , and  $c$  along with the three triples of each letter. Thus,

$$S = \left\{ \begin{array}{ccc} aaa & bbb & ccc \\ abc & bca & cba \\ acb & bac & cab \end{array} \right\}.$$

Furthermore, let each element of  $S$  have probability  $\frac{1}{9}$ . Define

$$A_i = \{i^{th} \text{ place in the triple is occupied by } a\}.$$

What are the values for  $\Pr(A_i), i = 1, 2, 3$ ? Are they pairwise independent?

*solution*

It is easy to count that

$$\Pr(A_i) = \frac{1}{3}, i = 1, 2, 3,$$

and

$$\Pr(A_1, A_2) = \Pr(A_1, A_3) = \Pr(A_2, A_3) = \frac{1}{9}$$

so that  $A_i$ s are pairwise independent.

**Definition\*** A collection of events  $A_1, \dots, A_n$  are *mutually independent* if for any subcollection  $A_{i_1}, \dots, A_{i_k}$ , we have

$$\Pr(\cap_{j=1}^k A_{i_j}) = \prod_{j=1}^k \Pr(A_{i_j}).$$

## Random Variables

In many experiments, it is easier to deal with a summary variable than with the original probability structure.

**Example** consider an opinion poll, we might decide to ask 50 people whether they agree or disagree with a certain issue. If we record a “1” for agree and “0” for disagree, the sample space for this experiment has  $2^{50}$  elements (all length 50 strings consist of 1s and 0s). However, if we are only interested in the number of people who agree, we may define a variable  $X =$  number of 1s recorded out of 50. Then, the sample space for  $X$  is the set of integers  $\{0, 1, 2, \dots, 50\}$ .

**Definition** A *random variable* (r.v.) is a function from a sample space  $S$  into the real numbers.

**Example** In some experiments random variables are implicitly used

#### Examples of random variables

Experiment	Random variable
Toss two dice	$X =$ sum of numbers
Toss a coin 25 times	$X =$ number of heads in 25 tosses
Apply different amounts of fertilizer to corn plants	$X =$ yield / acre

In defining a random variable, we have also defined a new sample space (the range of the random variable).

## Lecture 8: Sept 8

### Last time

- Conditional Probability (1.3)
- Independence (1.3)
- Random variables

### Today

- HW2 posted
- Random variables
- Distribution Functions

## Random Variables

In many experiments, it is easier to deal with a summary variable than with the original probability structure.

**Definition** A *random variable* (r.v.) is a function from a sample space  $S$  into the real numbers.

**Example** In some experiments random variables are implicitly used

Examples of random variables

Experiment	Random variable
Toss two dice	$X = \text{sum of numbers}$
Toss a coin 25 times	$X = \text{number of heads in 25 tosses}$
Apply different amounts of fertilizer to corn plants	$X = \text{yield / acre}$

In defining a random variable, we have also defined a new sample space (the range of the random variable).

**Induced probability function** Suppose we have a sample space  $S = \{s_1, s_2, \dots, s_n\}$  with a probability function  $\Pr$  defined on the original sample space. We define a random variable  $X$  with range  $\mathcal{X} = \{x_1, \dots, x_m\}$ . We can define a probability function  $\Pr_X$  on  $\mathcal{X}$  in the following way. We will observe  $X = x_i$  if and only if the outcome of the random experiment is an  $s_j \in S$  such that  $X(s_j) = x_i$ . Therefore,

$$\Pr_X(X = x_i) = \Pr(\{s_j \in S : X(s_j) = x_i\}),$$

defines an *induced* probability function on  $\mathcal{X}$ , defined in terms of the original function  $\Pr$ .

We will write  $\Pr(X = x_i)$  rather than  $\Pr_X(X = x_i)$  for simplicity. Note on notation: random variables will always be denoted with uppercase letters and the realized values of the variable (or its range) will be denoted by the corresponding lowercase letters.

**Example** Consider the experiment of tossing a fair coin three times. Define the random variable  $X$  to be the number of heads obtained in the three tosses. A complete enumeration of the value of  $X$  for each point in the sample space is

$s$	HHH	HHT	HTH	THH	TTH	THT	HTT	TTT
$X(s)$	3	2	2	2	1	1	1	0

What is the range of  $X$ ? What is the induced probability function  $\Pr_X$ ?

*solution:*

The range for the random variable  $X$  is  $\mathcal{X} = \{0, 1, 2, 3\}$ . Assuming all 8 points in  $S$  has probability  $\frac{1}{8}$ . By simply counting, we see that the induced probability function on  $\mathcal{X}$  is

$x$	0	1	2	3
$\Pr_X(X = x)$	$\frac{1}{8}$	$\frac{3}{8}$	$\frac{3}{8}$	$\frac{1}{8}$

So far, we have seen finite  $S$  and finite  $\mathcal{X}$ , and the definition of  $\Pr_X$  is straightforward. If  $\mathcal{X}$  is uncountable, we define the induced probability function,  $\Pr_X$  for any set  $A \subset \mathcal{X}$ ,

$$\Pr_X(X \in A) = \Pr(\{s \in S : X(s) \in A\}).$$

This defines a legitimate probability function for which the Kolmogorov Axioms can be verified.

## Distribution Functions

Distribution Functions are used to describe the behavior of a r.v.

### Cumulative distribution function

**Definition** The *cumulative distribution function* or *cdf* of a random variable  $X$ , denoted by  $F_X(x)$ , is defined by

$$F_X(x) = \Pr_X(X \leq x), \text{ for all } x.$$

**Definition** The *survival function* of a random variable  $X$ , is defined by

$$S_X(x) = 1 - F_X(x) = \Pr_X(X > x).$$

**Example** Consider the experiment of tossing three fair coins, and let  $X$  = number of heads observed. The cdf of  $X$  is

$$F_X(x) = \begin{cases} 0 & \text{if } -\infty < x < 0 \\ \frac{1}{8} & \text{if } 0 \leq x < 1 \\ \frac{1}{2} & \text{if } 1 \leq x < 2 \\ \frac{7}{8} & \text{if } 2 \leq x < 3 \\ 1 & \text{if } 3 \leq x < \infty \end{cases}$$

Some properties of the cdf:

Let  $F(x)$  be a cdf. Then

1.  $0 \leq F(x) \leq 1$
2.  $\lim_{x \rightarrow -\infty} F(x) = 0$
3.  $\lim_{x \rightarrow \infty} F(x) = 1$
4.  $F$  is nondecreasing: if  $a < b$ , then  $F(a) \leq F(b)$
5.  $F$  is right-continuous:  $\lim_{x \downarrow b} F(x) = F(b)$ , or  $\lim_{x \rightarrow b^+} F(x) = F(b)$
6.  $\Pr(a < X \leq B) = F(b) - F(a)$

**Theorem** The function  $F(x)$  is a cdf if and only if the following three conditions hold:

1.  $\lim_{x \rightarrow -\infty} F(x) = 0$  and  $\lim_{x \rightarrow \infty} F(x) = 1$
2.  $F$  is nondecreasing: if  $a < b$ , then  $F(a) \leq F(b)$
3.  $F$  is right-continuous:  $\lim_{x \downarrow b} F(x) = F(b)$ , or  $\lim_{x \rightarrow b^+} F(x) = F(b)$

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

**Definition** Two random variables  $X$  and  $Y$  are identically distributed if, for every Borel set  $A \subset \mathbb{R}$ ,  $\Pr(X \in A) = \Pr(Y \in A)$ .

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

**Theorem** The following two statements are equivalent:

1. The random variables  $X$  and  $Y$  are *identically distributed*.
2.  $F_X(x) = F_Y(x)$  for every  $x$ .



## Lecture 9: Sept 11

Last time

- Random variables

Today

- Distribution Functions
- Types of Random Variables

### Distribution Functions

Distribution Functions are used to describe the behavior of a r.v.

Cumulative distribution function

**Definition** The *cumulative distribution function* or *cdf* of a random variable  $X$ , denoted by  $F_X(x)$ , is defined by

$$F_X(x) = \Pr_X(X \leq x), \text{ for all } x.$$

**Definition** The *survival function* of a random variable  $X$ , is defined by

$$S_X(x) = 1 - F_X(x) = \Pr_X(X > x).$$

**Example** Consider the experiment of tossing three fair coins, and let  $X$  = number of heads observed. The cdf of  $X$  is

$$F_X(x) = \begin{cases} 0 & \text{if } -\infty < x < 0 \\ \frac{1}{8} & \text{if } 0 \leq x < 1 \\ \frac{1}{2} & \text{if } 1 \leq x < 2 \\ \frac{7}{8} & \text{if } 2 \leq x < 3 \\ 1 & \text{if } 3 \leq x < \infty \end{cases}$$

Some properties of the cdf:

Let  $F(x)$  be a cdf. Then

1.  $0 \leq F(x) \leq 1$
2.  $\lim_{x \rightarrow -\infty} F(x) = 0$
3.  $\lim_{x \rightarrow \infty} F(x) = 1$
4.  $F$  is nondecreasing: if  $a < b$ , then  $F(a) \leq F(b)$
5.  $F$  is right-continuous:  $\lim_{x \downarrow b} F(x) = F(b)$ , or  $\lim_{x \rightarrow b^+} F(x) = F(b)$
6.  $\Pr(a < X \leq B) = F(b) - F(a)$

**Theorem** The function  $F(x)$  is a cdf if and only if the following three conditions hold:

1.  $\lim_{x \rightarrow -\infty} F(x) = 0$  and  $\lim_{x \rightarrow \infty} F(x) = 1$
2.  $F$  is nondecreasing: if  $a < b$ , then  $F(a) \leq F(b)$
3.  $F$  is right-continuous:  $\lim_{x \downarrow b} F(x) = F(b)$ , or  $\lim_{x \rightarrow b^+} F(x) = F(b)$

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

**Definition** Two random variables  $X$  and  $Y$  are identically distributed if, for every Borel set  $A \subset \mathbb{R}$ ,  $\Pr(X \in A) = \Pr(Y \in A)$ .

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

**Theorem** The following two statements are equivalent:

1. The random variables  $X$  and  $Y$  are *identically distributed*.
2.  $F_X(x) = F_Y(x)$  for every  $x$ .

## Types of Random Variables

**Definition** A random variable  $X$  can be

- *discrete*:
  - $X$  takes on a finite or countably infinite number of values
  - $F_X(x)$  is step-wise constant
- *continuous*:
  - the range of  $X$  consists of subsets of the real line
  - $F_X(x)$  is continuous.
- *mixed*:  $F_X(x)$  is piecewise continuous.

**Example** A random variable has cdf

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

Is this a valid cdf? Is it a discrete random variable or continuous random variable or mixed?  
*solution:*

$F(x)$  satisfies the three properties of a cdf that

1.  $\lim_{x \rightarrow -\infty} F(x) = 0$  and  $\lim_{x \rightarrow \infty} F(x) = 1$
2.  $F$  is nondecreasing: if  $a < b$ , then  $F(a) \leq F(b)$
3.  $F$  is right-continuous:  $\lim_{x \downarrow b} F(x) = F(b)$ , or  $\lim_{x \rightarrow b^+} F(x) = F(b)$ .

Therefore,  $F(x)$  is a valid cdf. The random variable  $X$  is a mixed type.

## Lecture 10: Sept 13

### Last time

- Distribution Functions
- Types of Random Variables
- Discrete Random Variables

### Today

- Discrete Random Variables
- Continuous Random Variables
- Counting Techniques

### Discrete Random Variables

Suppose a random variable  $X$  takes only a finite or countable number of values. Let the sample space of  $X$  be  $S = \{x_1, x_2, \dots\}$ . Then the cdf can be expressed as:

$$F(x) = \sum_{x_i \leq x} \Pr(X = x_i).$$

**Definition** The *probability mass function* (pmf) of a discrete random variable  $X$  is given by

$$f_X(x) = \Pr(X = x) \text{ for all } x.$$

If the sample space of  $X$  is  $X = \{x_1, x_2, \dots\}$ , then

$$f(x_i) = \Pr(X = x_i) = \Pr(x_{i-1} < X \leq x_i) = F(x_i) - F(x_{i-1}).$$

**Example** (Geometric probabilities) Suppose we do an experiment that consists of tossing a coin until a head appears. Let  $p$  = probability of a head on any given toss, and define a random variable  $X$  = number of tosses required to get a head. Then for any  $x = 1, 2, \dots$ ,

$$\Pr(X = x) = (1 - p)^{x-1}p,$$

since we must get  $x - 1$  tails followed by a head for the event to occur and all trials are independent. What is the pmf of the above Geometric distribution? What is the cdf?

*solution:*

We have the pmf

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

For cdf, we have

$$\begin{aligned}
 F(x) &= \Pr(X \leq x) = \sum_{i=1}^{\lfloor x \rfloor} f(i) \\
 &= \begin{cases} f(1) + f(2) + \cdots + f(\lfloor x \rfloor) & \text{for } x \geq 1 \\ 0 & \text{otherwise} \end{cases} \\
 &= \begin{cases} 1 - (1-p)^{\lfloor x \rfloor} & \text{for } x \geq 1 \\ 0 & \text{for } x < 1 \end{cases}
 \end{aligned}$$

where  $\lfloor x \rfloor$  denote the floor function that returns the largest integer smaller or equal to  $x$  and we used the summation of a geometric sequence.

**Definition** The *domain* of a random variable  $X$  is the set of all values of  $x$  for which  $f(x) > 0$ . This is also called *range*, *sample space* or *support*.

Properties of the pmf:

1.  $f(x) > 0$  for at most a countable number of values  $x$ . For all other values  $x$ ,  $f(x) = 0$ .
2. Let  $\{x_1, x_2, \dots\}$  denote the domain of  $X$ . Then

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

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

**Example** What is the pmf of a deterministic random variable (a constant)?  
*solution:*

$$f(x) = \Pr(X = x) = \begin{cases} 1 & \text{for } x = c \\ 0 & \text{otherwise.} \end{cases}$$

This is equivalent as a constant of value  $c$ .

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

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

How would we determine if this is an allowable pmf?

*solution:*

We show that  $f(x)$  satisfies the properties of pmf.

1.  $f(x) > 0$  for a countable number of values  $x$ . For all other values  $x$ ,  $f(x) = 0$ .

2. Let  $\{x_1, x_2, \dots\}$  denote the domain of  $X$ . Then

$$\sum_{i=1}^{\infty} f(x_i) = \sum_{i=1}^{\infty} f(i) = \sum_{i=1}^{\infty} \left( \frac{1}{i} - \frac{1}{i+1} \right) = 1.$$

## Lecture 11: Sept 15

### Last time

- Discrete Random Variables

### Today

- Continuous Random Variables
- Transformations of Random Variables

### Continuous Random Variables

**Definition** A random variable  $X$  is *continuous* if  $F_X(x)$  is a continuous function of  $x$ .

**Definition** A random variable  $X$  is *absolutely continuous* if  $F_X(x)$  is an absolutely continuous function of  $x$ .

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

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

Absolute continuity is stronger than continuity but weaker than differentiability. An example of an absolutely continuous function is one that is:

- continuous everywhere
- differentiable everywhere, except possibly for a countable number of points.

**Definition** The *probability density function* or pdf,  $f_X(x)$ , of a continuous random variable  $X$  is the function that satisfies

$$F_X(x) = \int_{-\infty}^x f_X(t)dt \quad \text{for all } x.$$

**Notation:** We write  $X \sim F_X(x)$  for the expression “ $X$  has a distribution given by  $F_X(x)$ ” where we read the symbol “ $\sim$ ” as “is distributed as”. Similarly, we can write  $X \sim f_X(x)$  or  $X \sim F_X(x)$ , if  $X$  and  $Y$  have the same distribution,  $X \sim Y$ .

**Theorem** A function  $f_X(x)$  is a pdf (or pmf) of a random variable  $X$  if and only if

1.  $f_X(x) \geq 0$  for all  $x$ .
2.  $\int_{-\infty}^{\infty} f_X(x)dx = 1$  (pdf)    or     $\sum_x f_X(x) = 1$  (pmf).

**Example** Suppose  $\lambda > 0$ ,  $F(x) = 1 - e^{-\lambda x}$  for  $x > 0$  and  $F(x) = 0$  otherwise. Is  $F(x)$  a cdf? What is the associated pdf?

*solution:*

$F(x)$  satisfies the three properties of cdf

1.  $\lim_{x \rightarrow -\infty} F(x) = 0$  and  $\lim_{x \rightarrow \infty} F(x) = 1$
2.  $F$  is nondecreasing: if  $a < b$ , then  $F(a) \leq F(b)$
3.  $F$  is right-continuous:  $\lim_{x \downarrow b} F(x) = F(b)$ , or  $\lim_{x \rightarrow b^+} F(x) = F(b)$ .

$F(x)$  is a cdf. Actually,  $F(x)$  is the cdf of exponential distribution.

To get the pdf, we only need to differentiate the cdf.

$$f(x) = \frac{dF(x)}{dx} = \begin{cases} \lambda e^{-\lambda x} & \text{for } x > 0 \\ 0 & \text{otherwise.} \end{cases}$$

Note

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

*proof*

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

- Because  $\Pr(X = a) = 0$ , all the following are equivalent:

$$\Pr(a \leq X \leq b), \quad \Pr(a \leq X < b) \quad , \quad \Pr(a < X \leq b) \quad \text{and} \quad \Pr(a < X < b)$$

- $f(x)$  can exceed one!



## Lecture 12: Sept 18

Last time

- Continuous Random Variables

Today

- Transformations of Random Variables

### Transformations of Random Variables

**Theorem** If  $X$  is a r.v. 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  $Y$ .

**Probability mapping:** For any set  $A \subset \mathcal{Y}$ :

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

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 necessarily bijective.

**Example** (Binomial transformation) A discrete random variable  $X$  has a *binomial distribution* if its pmf is of the form

$$f_X(x) = \Pr(X = x) = \binom{n}{x} p^x (1-p)^{n-x}, \quad x = 0, 1, \dots, n,$$

where  $n$  is a positive integer and  $0 \leq p \leq 1$ . Values such as  $n$  and  $p$  that can be set to different values, producing different probability distributions, are called *parameters*. Consider a random variable  $Y = g(X)$ , where  $g(x) = n - x$ ; that is,  $Y = n - X$ . Here  $\mathcal{X} = \{0, 1, \dots, n\}$  and  $\mathcal{Y} = \{y : y = g(x), x \in \mathcal{X}\} = \{0, 1, \dots, n\}$ . For any  $y \in \mathcal{Y}$ ,  $n - x = g(x) = y$  if and only if  $x = n - y$ . Therefore,  $g^{-1}(y) = n - y$  and

$$\begin{aligned}f_Y(y) &= \sum_{x \in g^{-1}(y)} f_X(x) \\ &= f_X(n - y) \\ &= \binom{n}{n-y} p^{n-y} (1-p)^{n-(n-y)} \\ &= \binom{n}{y} (1-p)^y p^{n-y}.\end{aligned}$$

Therefore,  $Y$  also has a binomial distribution, but with parameters  $n$  and  $1 - p$ .

**Example** (exercise 2.3) Suppose  $X$  has the geometric pmf  $f_X(x) = \frac{1}{3}(\frac{2}{3})^x, x = 0, 1, 2, \dots$ . Determine the probability distribution of  $Y = X/(X + 1)$ . Note that here both  $X$  and  $Y$  are discrete random variables. To specify the probability distribution of  $Y$ , specify its pmf.  
*Solution:*

$$\Pr(Y = y) = \Pr\left(\frac{X}{X + 1} = y\right) = \Pr\left(X = \frac{y}{1 - y}\right) = \frac{1}{3}\left(\frac{2}{3}\right)^{y/(1-y)}, y = 0, \frac{1}{2}, \frac{2}{3}, \dots, \frac{x}{x + 1}, \dots$$

## Lecture 13: Sept 20

Last time

- Transformations of Random Variables

Today

- Transformations of Random Variables

### Transformations of Random Variables

**Theorem** Suppose a continuous random variable  $X$  has cdf  $F_X(x)$ , let  $Y = g(X)$ , and let  $\mathcal{X}$  and  $\mathcal{Y}$  be defined as

$$\mathcal{X} = \{x : f(x) > 0\} \quad \text{and} \quad \mathcal{Y} = \{y : y = g(x) \text{ for some } x \in \mathcal{X}\}.$$

Then,

1. If  $g$  is an increasing function on  $\mathcal{X}$ ,  $F_Y(y) = F_X(g^{-1}(y))$  for  $y \in \mathcal{Y}$ .
2. If  $g$  is a decreasing function on  $\mathcal{X}$ ,  $F_Y(y) = 1 - F_X(g^{-1}(y))$  for  $y \in \mathcal{Y}$ .

*Proof:* We start with

$$\begin{aligned} F_Y(y) &= \Pr(Y \leq y) \\ &= \Pr(g(X) \leq y) \end{aligned}$$

1. If  $g$  is an increasing function, then  $g(X) \leq y$  if and only if  $X \leq g^{-1}(y)$ . Therefore,  $F_Y(y) = \Pr(g(X) \leq y) = \Pr(X \leq g^{-1}(y)) = F_X(g^{-1}(y))$ .
2. Similarly, if  $g$  is a decreasing function, then  $g(X) \leq y$  if and only if  $X \geq g^{-1}(y)$ . And  $F_Y(y) = \Pr(g(X) \leq y) = \Pr(X \geq g^{-1}(y)) = 1 - F_X(g^{-1}(y))$ .

**Theorem** Let  $X$  have pdf  $f_X(x)$  and let  $Y = g(X)$ , where  $g$  is a monotone function. Let  $\mathcal{X}$  and  $\mathcal{Y}$  be defined as

$$\mathcal{X} = \{x : f(x) > 0\} \quad \text{and} \quad \mathcal{Y} = \{y : y = g(x) \text{ for some } x \in \mathcal{X}\}.$$

Suppose that  $f_X(x)$  is continuous on  $\mathcal{X}$  and that  $g^{-1}(y)$  has a continuous derivative on  $\mathcal{Y}$ . Then the pdf of  $Y$  is given by

$$f_Y(y) = \begin{cases} f_X(g^{-1}(y)) \left| \frac{d}{dy} g^{-1}(y) \right| & y \in \mathcal{Y} \\ 0 & \text{otherwise.} \end{cases}$$

*Proof:*

From last theorem, we have the cdf forms  $F_Y(y)$ . Then  $f_Y(y) = \frac{d}{dy} F_Y(y)$ . (finish the proof)  
From last theorem, we have

$$F_Y(y) = \begin{cases} F_X(g^{-1}(y)) & \text{if } g \text{ is increasing} \\ 1 - F_X(g^{-1}(y)) & \text{if } g \text{ is decreasing.} \end{cases}$$

We have, by the chain rule,

$$f_Y(y) = \frac{d}{dy}F_Y(y) = \begin{cases} f_X(g^{-1}(y))\frac{d}{dy}g^{-1}(y) & \text{if } g \text{ is increasing} \\ -f_X(g^{-1}(y))\frac{d}{dy}g^{-1}(y) & \text{if } g \text{ is decreasing,} \end{cases}$$

where  $\frac{d}{dy}g^{-1}(y) < 0$  when  $g$  is decreasing such that  $-\frac{d}{dy}g^{-1}(y) = |\frac{d}{dy}g^{-1}(y)|$ .

**Example** (Square transformation) Suppose  $X$  is a continuous random variable. For  $y > 0$ , the cdf of  $Y = X^2$  is

$$F_Y(y) = \Pr(Y \leq y) = \Pr(X^2 \leq y) = \Pr(-\sqrt{y} \leq X \leq \sqrt{y}).$$

Because  $x$  is continuous, we can drop the equality from the left endpoint and obtain

$$\begin{aligned} F_Y(y) &= \Pr(-\sqrt{y} < X \leq \sqrt{y}) \\ &= \Pr(X \leq \sqrt{y}) - \Pr(X \leq -\sqrt{y}) = F_X(\sqrt{y}) - F_X(-\sqrt{y}). \end{aligned}$$

The pdf of  $Y$  can now be obtained from the cdf by differentiation:

$$\begin{aligned} f_Y(y) &= \frac{d}{dy}F_Y(y) \\ &= \frac{d}{dy} [F_X(\sqrt{y}) - F_X(-\sqrt{y})] \\ &= \frac{1}{2\sqrt{y}}f_X(\sqrt{y}) + \frac{1}{2\sqrt{y}}f_X(-\sqrt{y}), \end{aligned}$$

where we use the chain rule to differentiate  $F_X(\sqrt{y})$  and  $F_X(-\sqrt{y})$ .

## Lecture 14: Sept 22

Last time

- Transformations of Random Variables

Today

- 09/18 attendance void (Jewish Holiday)
- Transformations of Random Variables

### Transformations of Random Variables

**Example** (Linear transformation) Suppose  $X$  is a continuous random variable 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 random variable  $Y$  with a distribution that has the same shape as that of  $X$  but has mean  $a$  and variance  $b^2$ .

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” transformation.

**Example** (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.

**Example** (Square root of an exponential RV) Suppose  $X \sim \text{exp}(\lambda)$ , so that

$$f_X(x) = \begin{cases} \lambda e^{-\lambda x} & x \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

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 Weibull distribution.

## Lecture 15: Sept 25

Last time

- Transformations of Random Variables

Today

- 09/18 attendance still valid (actually today is a Jewish Holiday)
- One-page one-sided letter-size cheat sheet for midterm 1
- Transformations of Random Variables
- Expected Values

### Transformations of Random Variables

**Theorem** (Probability integral transformation) Let  $X$  have continuous cdf  $F_X(x)$  and define the random variable  $Y$  as  $Y = F_X(X)$ . Then  $Y$  is uniformly distributed on  $(0, 1)$ , that is,  $\Pr(Y \leq y) = y, 0 < y < 1$ .

Before we prove this theorem, we will digress for a moment and look at  $F_X^{-1}$ , the inverse of the cdf  $F_X$ , in some detail. If  $F_X$  is strictly increasing, then  $F_X^{-1}$  is well defined by

$$F_X^{-1}(y) = x \iff F_X(x) = y.$$

However, if  $F_X$  is constant on some interval, then  $F_X^{-1}$  is not well defined as Figure 14.1 illustrates. Any  $x_1 \leq x \leq x_2$  satisfies  $F_X(x) = y$



Figure 14.1: Figure 2.1.2. (a)  $F_X(x)$  strictly increasing; (b)  $F_X(x)$  nondecreasing

This problem is avoided by defining  $F_X^{-1}$  for  $0 < y < 1$  by

$$F_X^{-1}(y) = \inf\{x : F_X(x) \geq y\}.$$

With this definition, for Figure 14.1(b), we have  $F_X^{-1}(y) = x_1$ .

*Proof:*

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

$$\begin{aligned} \Pr(Y \leq y) &= \Pr(F_X(X) \leq y) \\ &= \Pr(F_X^{-1}[F_X(X)] \leq F_X^{-1}(y)) \quad (F_X^{-1} \text{ is increasing}) \\ &= \Pr(X \leq F_X^{-1}(y)) \\ &= F_X(F_X^{-1}(y)) \quad (\text{definition of } F_X) \\ &= y. \end{aligned}$$

One application of the probability integral transformation is in the generation of random samples from a particular distribution. If it is required to generate an observation  $X$  from a population with cdf  $F_X$ , we need only generate a uniform random number  $U$ , between 0 and 1, and solve for  $x$  in the equation  $F_X(x) = u$ .

## Expected Values

**Definition** The *expected value* or *mean* of a random variable  $g(X)$ , denoted by  $Eg(X)$ , is

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

Provided the integral or summation exists.

If we let  $g(X) = X$ , then we get

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

**Example** (Exponential mean) Suppose  $X$  has an *exponential* ( $\lambda$ ) *distribution*,  $X \sim \text{Exp}(\lambda)$ , that is, it has pdf given by

$$f_X(x) = \frac{1}{\lambda} e^{-x/\lambda}, \quad 0 \leq x < \infty, \lambda > 0.$$

Find out  $EX$ .

*Solution:*



$$\begin{aligned}
EX &= \int_0^{\infty} \frac{1}{\lambda} x e^{-x/\lambda} dx \\
&= -x e^{-x/\lambda} \Big|_0^{\infty} + \int_0^{\infty} e^{-x/\lambda} dx \\
&= \int_0^{\infty} e^{-x/\lambda} dx \\
&= \lambda
\end{aligned}$$

## Lecture 16: Sept 27

### Last time

- Transformations of Random Variables
- Expected Values
- Moments

### Today

- One-page one-sided letter-size hand-written cheat sheet for midterm 1
- Expected Values

### Expected Values

**Definition** The *expected value* or *mean* of a random variable  $g(X)$ , denoted by  $Eg(X)$ , is

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

Provided the integral or summation exists.

If we let  $g(X) = X$ , then we get

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

**Example** (Exponential mean) Suppose  $X$  has an *exponential* ( $\lambda$ ) *distribution*,  $X \sim \text{Exp}(\lambda)$ , that is, it has pdf given by

$$f_X(x) = \frac{1}{\lambda}e^{-x/\lambda}, \quad 0 \leq x < \infty, \lambda > 0.$$

Find out  $EX$ .

*Solution:*

$$\begin{aligned}
EX &= \int_0^{\infty} \frac{1}{\lambda} x e^{-x/\lambda} dx \\
&= -x e^{-x/\lambda} \Big|_0^{\infty} + \int_0^{\infty} e^{-x/\lambda} dx \\
&= \int_0^{\infty} e^{-x/\lambda} dx \\
&= \lambda
\end{aligned}$$

**Example** (Binomial mean) if  $X$  has a *binomial distribution*,  $X \sim \text{Binomial}(n, p)$ , its pmf is given by

$$\Pr(X = x) = \binom{n}{x} p^x (1-p)^{n-x}, \quad x = 0, 1, \dots, n,$$

where  $n$  is a positive integer,  $0 \leq p \leq 1$ , and for every fixed pair  $n$  and  $p$  the pmf sums to 1. Find out  $EX$ .

*Solution:*

$$EX = \sum_{x=0}^n x \binom{n}{x} p^x (1-p)^{n-x} = \sum_{x=1}^n x \binom{n}{x} p^x (1-p)^{n-x}.$$

Using the identity  $x \binom{n}{x} = n \binom{n-1}{x-1}$ , we have

$$\begin{aligned}
EX &= \sum_{x=1}^n n \binom{n-1}{x-1} p^x (1-p)^{n-x} \\
&= \sum_{y=0}^{n-1} n \binom{n-1}{y} p^{y+1} (1-p)^{n-(y+1)} \\
&= np \sum_{y=0}^{n-1} \binom{n-1}{y} p^y (1-p)^{n-1-y} \\
&= np,
\end{aligned}$$

since the last summation must be 1, being the sum over all possible values of a binomial( $n-1, p$ ) pmf.

The process of taking expectations is a linear operation, which means that the expectation of a linear function of  $X$  can be easily evaluated by noting that for any constants  $a$  and  $b$ , such that

$$E(aX + b) = aEX + b$$

**Theorem** Let  $X$  be a random variable and let  $a$ ,  $b$ , and  $c$  be constants. Then for any functions  $g_1(x)$  and  $g_2(x)$  whose expectations exist,

1.  $E(ag_1(X) + bg_2(X) + c) = aEg_1(X) + bEg_2(X) + c$ .
2. If  $g_1(x) \geq 0$  for all  $x$ , then  $Eg_1(X) \geq 0$ .
3. If  $g_1(x) \geq g_2(x)$  for all  $x$ , then  $Eg_1(X) \geq Eg_2(X)$ .
4. If  $a \leq g_1(x) \leq b$  for all  $x$ , then  $a \leq Eg_1(X) \leq b$ .

*Proof:*

We will give details for only the continuous case, the discrete case being similar. By definition

$$\begin{aligned}
 E(ag_1(X) + bg_2(X) + c) &= \int_{-\infty}^{\infty} [ag_1(x) + bg_2(x) + c] f_X(x) dx \\
 &= \int_{-\infty}^{\infty} ag_1(x) f_X(x) dx + \int_{-\infty}^{\infty} bg_2(x) f_X(x) dx + \int_{-\infty}^{\infty} cf_X(x) dx \\
 &= aEg_1(X) + bEg_2(X) + c
 \end{aligned}$$

The other three properties are proved in a similar manner (shown in class).

## Lecture 17: Sept 29

### Last time

- Expected Values

### Today

- One-page one-sided letter-size hand-written cheat sheet for midterm 1
- Expected Values
- Moments

### Expected Values

**Example** (Method of indicators) An example of how the above properties are useful. Let  $X \sim \text{Binomial}(n, p)$  for  $n$  positive integer and  $0 \leq p \leq 1$  ( $n$  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  $i^{\text{th}}$  trial is a success (i.e.  $I_i \stackrel{\text{i.i.d.}}{\sim} \text{Bernoulli}(p)$ ). We have

$$EI_i = 1 \cdot p + 0 \cdot (1 - p) = p.$$

Therefore,

$$EX = \sum_{i=1}^n EI_i = \sum_{i=1}^n p = np.$$

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

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

*Proof:*

We prove the continuous case first,

$$\begin{aligned}\int_0^{\infty} [1 - F(x)] dx &= \int_0^{\infty} [1 - \Pr(X \leq x)] dx \\&= \int_0^{\infty} \Pr(X > x) dx \\&= \int_0^{\infty} \int_{y=x}^{\infty} f_X(y) dy dx \\&= \int_0^{\infty} \int_{x=0}^y f_X(y) dx dy \\&= \int_0^{\infty} y f_X(y) dy \\&= EX.\end{aligned}$$

Then, for discrete case, we have

$$\begin{aligned}\sum_{x=0}^{\infty} (1 - F(x)) &= \sum_{x=0}^{\infty} \Pr(X > x) \\&= \sum_{x=0}^{\infty} \sum_{y=x+1}^{\infty} \Pr(X = y) \\&= \sum_{y=1}^{\infty} \sum_{x=0}^{y-1} \Pr(X = y) \\&= \sum_{y=1}^{\infty} y \Pr(X = y) \\&= EX\end{aligned}$$

## Lecture 18: Oct 9

### Last time

- Midterm exam 1

### Today

- Distribute Midterm Exam 1
- Vote for alternative weighting on Wednesday:
  - Old one still works
  - New: if Midterm exam 2 score is higher than Midterm exam 1 score, then use Midterm exam 2 score for both
  - Final grade will use the higher one
- Moments

### Moments

**Example** (Minimizing distance) The expected value of a random variable has another property, one that we can think of as relating to the interpretation of  $EX$  as a good guess at a value of  $X$ .

Suppose we measure the distance between a random variable  $X$  and a constant  $b$  by  $(X - b)^2$ . The closer  $b$  is to  $X$ , the smaller this quantity is. We can now determine the value of  $b$  that minimizes  $E[(X - b)^2]$  and, hence, will provide us with a good predictor of  $X$ . (Note that it does no good to look for a value of  $b$  that minimizes  $(X - b)^2$ , since the answer would depend on  $X$ , making it a useless predictor of  $X$ .)

We could proceed with the minimization of  $E(X - b)^2$  by using calculus, but there is a simpler method:

$$\begin{aligned} E(X - b)^2 &= E(X - EX + EX - b)^2 \\ &= E[(X - EX) + (EX - b)]^2 \\ &= E(X - EX)^2 + (EX - b)^2 + 2E[(X - EX)(EX - b)], \end{aligned}$$

where we have expanded the square. Note that  $E[(X - EX)(EX - b)] = (EX - b)E(X - EX) = 0$ , since  $EX - b$  is constant and comes out of the expectation,  $E(X - EX) = EX - EX = 0$ . This means

$$E(X - b)^2 = E(X - EX)^2 + (EX - b)^2.$$

Such that  $E(X - b)^2$  is minimized at  $b = EX$ . And  $E(X - EX)^2$  is actually the variance of  $X$  (i.e.,  $Var(X) = E[(X - EX)^2]$ ).

## Lecture 19: Oct 11

### Last time

- Distribute Midterm Exam 1

### Today

- Vote for alternative weighting (remind me 10 min before class ends):
  - Old one still works
  - New: if Midterm exam 2 score is higher than Midterm exam 1 score, then use Midterm exam 2 score for both
  - Final grade will use the higher one
- Moments
- Moment Generating Function

### Moments

The various moments of a distribution are an important class of expectations.

**Definition** For each integer  $n$ , the  $n$ th *moment* of  $X$  (or  $F_X(x)$ ),  $\mu'_n$ , is

$$\mu'_n = EX^n.$$

The  $n$ th *central moment* of  $X$ ,  $\mu_n$ , is

$$\mu_n = E(X - \mu)^n,$$

where  $\mu = \mu'_1 = EX$ .

Notes:

- $\mu'_0 = EX^0 = 1$
- $\mu'_1$  is the *mean*, usually denoted by  $\mu$ .
- $\mu_0 = E(X - \mu)^0 = 1$
- $\mu_1 = 0$
- $\mu_2 = E(X - EX)^2$  is the *variance*
- $\mu_3 = E(X - EX)^3$  is related to the *skewness*.
- $\mu_4 = E(X - EX)^4$  is related to the *kurtosis*.



**Definition** The *variance* of a random variable  $X$  is its second central moment,  $\text{Var}(X) = E[(X - EX)^2]$ . The positive square root of  $\text{Var}(X)$  is the *standard deviation* of  $X$ .

The variance gives a measure of the degree of spread of a distribution around its mean. Figure 30.4 shows a plot of two samples, one sample draws 100 numbers from a normal distribution with mean 0 and variance 1,  $N(0, 1)$ . The other sample draws 100 numbers from a normal distribution with mean 0 and variance 100,  $N(0, 100)$ .



Figure 18.2: Figure 2.1.2. Two samples of 100 numbers drawn from  $N(0, 1)$  and  $N(0, 100)$ .

**Example** (Exponential variance) Let  $X$  have the exponential( $\lambda$ ) distribution. We can calculate the variance of  $X$  now.

*Solution:*

$$\begin{aligned}\text{Var}(X) &= E(X - \lambda)^2 \\ &= \int_0^{\infty} (x - \lambda)^2 \frac{1}{\lambda} e^{-x/\lambda} dx \\ \text{Var}(X) &= \int_0^{\infty} (x^2 - 2x\lambda + \lambda^2) \frac{1}{\lambda} e^{-x/\lambda} dx \\ &= \int_0^{\infty} x^2 \frac{1}{\lambda} e^{-x/\lambda} dx - 2 \int_0^{\infty} x\lambda \frac{1}{\lambda} e^{-x/\lambda} dx + \lambda^2 \\ &= EX^2 - \lambda^2 \\ &= \lambda^2\end{aligned}$$

**Theorem** If  $X$  is a random variable with finite variance, then for any constants  $a$  and  $b$ ,

$$\text{Var}(aX + b) = a^2 \text{Var}(X).$$

*Proof:*

From the definition, we have

$$\begin{aligned}\mathrm{Var}(aX + b) &= E[(aX + b) - E(aX + b)]^2 \\ &= E(aX - aEX)^2 \\ &= a^2 E(X - EX)^2 \\ &= a^2 \mathrm{Var}(X).\end{aligned}$$

It is sometimes to use an alternative formula for the variance, given by

$$\mathrm{Var}(X) = E(X^2) - (EX)^2,$$

which is easily established by

$$\begin{aligned}\mathrm{Var}(X) &= E(X - EX)^2 = E[X^2 - 2XEX + (EX)^2] \\ &= EX^2 - 2(EX)^2 + (EX)^2 \\ &= EX^2 - (EX)^2.\end{aligned}$$

## Lecture 20: Oct 13

### Last time

- Distribute Midterm Exam 1
- Alternative weighting passed
- Moments

### Today

- Internal “evaluation” open where you can make anonymous suggestions
- Moments
- Moment Generating Function

### Moments

The various moments of a distribution are an important class of expectations.

**Definition** For each integer  $n$ , the  $n$ th *moment* of  $X$  (or  $F_X(x)$ ),  $\mu'_n$ , is

$$\mu'_n = EX^n.$$

The  $n$ th *central moment* of  $X$ ,  $\mu_n$ , is

$$\mu_n = E(X - \mu)^n,$$

where  $\mu = \mu'_1 = EX$ .

**Theorem** If  $X$  is a random variable with finite variance, then for any constants  $a$  and  $b$ ,

$$\text{Var}(aX + b) = a^2 \text{Var}(X).$$

It is sometimes to use an alternative formula for the variance, given by

$$\text{Var}(X) = E(X^2) - (EX)^2,$$

which is easily established by

$$\begin{aligned} \text{Var}(X) &= E(X - EX)^2 = E[X^2 - 2XEX + (EX)^2] \\ &= EX^2 - 2(EX)^2 + (EX)^2 \\ &= EX^2 - (EX)^2. \end{aligned}$$

**Example** (Binomial variance) Let  $X \sim \text{Binomial}(n, p)$ , that is ,

$$\Pr(X = x) = \binom{n}{x} p^x (1 - p)^{n-x}.$$

What is the variance of  $X$ ?

*Solutions:*

Method #1:

We want to find  $EX^2$  first. We use the

$$EX^2 = \sum_{x=0}^n x^2 \binom{n}{x} p^x (1 - p)^{n-x}.$$

we use the same property  $x^2 \binom{n}{x} = xn \binom{n-1}{x-1}$ . We then have

$$\begin{aligned} EX^2 &= n \sum_{x=1}^n x \binom{n-1}{x-1} p^x (1 - p)^{n-x} \\ &= n \sum_{y=0}^{n-1} (y+1) \binom{n-1}{y} p^{y+1} (1 - p)^{n-1-y} \\ &= np \sum_{y=0}^{n-1} y \binom{n-1}{y} p^y (1 - p)^{n-1-y} + np \sum_{y=0}^{n-1} \binom{n-1}{y} p^y (1 - p)^{n-1-y} \\ &= np \cdot (n-1)p + np \\ &= n(n-1)p^2 + np. \end{aligned}$$

And now

$$\begin{aligned} \text{Var}(X) &= EX^2 - (EX)^2 \\ &= n(n-1)p^2 + np - (np)^2 \\ &= np - np^2 \\ &= np(1 - p). \end{aligned}$$

Method #2:

Recall that we could write  $X = \sum_{i=1}^n I_i$ , where  $I_i \stackrel{\text{i.i.d.}}{\sim} \text{Bernoulli}(p)$ . Then

$$\begin{aligned} \text{Var}(X) &= \text{Var}\left(\sum_{i=1}^n I_i\right) \\ &= \sum_{i=1}^n \text{Var}(I_i) \quad (I_i\text{'s are independent}) \\ &= n \text{Var}(I_i) \quad (I_i\text{'s are identically distributed}) \\ &= n [E(I_i^2) - (EI_i)^2] \\ &= n [p - p^2] \\ &= np(1 - p). \end{aligned}$$

## Moment Generating Function

**Definition** Let  $X$  be a random variable with cdf  $F_X$ . The *moment generating function (mgf)* of  $X$  (or  $F_X$ ), denoted by  $M_X(t)$ , is

$$M_X(t) = Ee^{tX},$$

provided that the expectation exists for  $t$  in some neighborhood of 0. That is, there is an  $h > 0$  such that, for all  $t$  in  $-h < t < h$ ,  $Ee^{tX}$  exists. If the expectation does not exist in a neighborhood of 0, we say that the moment generating function does not exist.

More explicitly, we can write the mgf of  $X$  as

$$M_X(t) = \int_{-\infty}^{\infty} e^{tx} f_X(x) dx, \quad \text{if } X \text{ is continuous,}$$

or

$$M_X(t) = \sum_x e^{tx} \Pr(X = x), \quad \text{if } X \text{ is discrete.}$$

It is easy to see how the mgf generates moments as in the following theorem.

## Lecture 21: Oct 16

### Last time

- Moments
- Moment Generating Function

### Today

- Internal “evaluation” open where you can make anonymous suggestions
- Moment Generating Function
- Common Discrete Distribution

### Moment Generating Function

**Definition** Let  $X$  be a random variable with cdf  $F_X$ . The *moment generating function (mgf)* of  $X$  (or  $F_X$ ), denoted by  $M_X(t)$ , is

$$M_X(t) = Ee^{tX},$$

provided that the expectation exists for  $t$  in some neighborhood of 0. That is, there is an  $h > 0$  such that, for all  $t$  in  $-h < t < h$ ,  $Ee^{tX}$  exists. If the expectation does not exist in a neighborhood of 0, we say that the moment generating function does not exist.

More explicitly, we can write the mgf of  $X$  as

$$M_X(t) = \int_{-\infty}^{\infty} e^{tx} f_X(x) dx, \quad \text{if } X \text{ is continuous,}$$

or

$$M_X(t) = \sum_x e^{tx} \Pr(X = x), \quad \text{if } X \text{ is discrete.}$$

It is easy to see how the mgf generates moments as in the following theorem.

**Theorem** If  $X$  has mgf  $M_X(t)$ , then

$$EX^n = M_X^{(n)}(0),$$

where we define

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

That is, the  $n^{th}$  moment is equal to the  $n^{th}$  derivative of  $M_X(t)$  evaluated at  $t = 0$ .

*Proof:*

$$\begin{aligned}
\frac{d}{dt}M_X(t) &= \frac{d}{dt} \int_{-\infty}^{\infty} e^{tx} f_X(x) dx \\
&= \int_{-\infty}^{\infty} \left( \frac{d}{dt} e^{tx} \right) f_X(x) dx \\
&= \int_{-\infty}^{\infty} (x e^{tx}) f_X(x) dx \\
&= E(X e^{tX}).
\end{aligned}$$

Therefore,

$$\left. \frac{d}{dt} M_X(t) \right|_{t=0} = E(X e^{tX}) \Big|_{t=0} = EX.$$

Proceeding in an analogous manner, we can establish that

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

**Example** (Binomial mgf) Let  $X \sim \text{Binomial}(n, p)$ , then its mgf is

$$\begin{aligned}
M_X(t) &= \sum_{x=0}^n e^{tx} \binom{n}{x} p^x (1-p)^{n-x} \\
&= \sum_{x=0}^n \binom{n}{x} (pe^t)^x (1-p)^{n-x} \\
&= [pe^t + (1-p)]^n.
\end{aligned}$$

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

1. If  $X$  and  $Y$  have **bounded support**, then  $F_X(u) = F_Y(u)$  for all  $u$  if and only if  $EX^r = EY^r$  for all integers  $r = 0, 1, 2, \dots$ .
2. If the moment generating functions exist and  $M_X(t) = M_Y(t)$  for all  $t$  in some neighborhood of 0, then  $F_X(u) = F_Y(u)$  for all  $u$ .

**Theorem** (Convergence of mgfs) Suppose  $\{X_i, i = 1, 2, \dots\}$  is a sequence of random variables, each with mgf  $M_{X_i}(t)$ . Furthermore, suppose that

$$\lim_{i \rightarrow \infty} M_{X_i}(t) = M_X(t), \quad \text{for all } t \text{ in a neighborhood of } 0,$$

and  $M_X(t)$  is an mgf. Then there is a unique cdf  $F_X$  whose moments are determined by  $M_X(t)$  and, for all  $x$  where  $F_X(x)$  is continuous, we have

$$\lim_{i \rightarrow \infty} F_{X_i}(x) = F_X(x).$$

That is, *convergence*, for  $|t| < h$ , of mgfs to an mgf implies *convergence* of cdfs.

**Poisson approximation** One approximation that is usually taught in elementary statistics courses is that binomial probabilities can be approximated by Poisson probabilities. It is taught that the Poisson approximation is valid “when  $n$  is large and  $np$  is small”, and rules of thumb are sometimes given.

The *Poisson*( $\lambda$ ) pmf is given by

$$\Pr(X = x) = \frac{e^{-\lambda} \lambda^x}{x!}, \quad x = 0, 1, 2, \dots,$$

where  $\lambda$  is a positive constant. The approximation states that if  $X \sim \text{Binomial}(n, p)$  and  $Y \sim \text{Poisson}(\lambda)$ , with  $\lambda = np$ , then

$$\Pr(X = x) \approx \Pr(Y = x)$$

for large  $n$  and small  $np$ . We now show that the mgf converge, lending credence to this approximation. Recall that

$$M_X(t) = [pe^t + (1 - p)]^n.$$

For the *Poisson*( $\lambda$ ) distribution, we can calculate (HW4, exercise 2.33)

$$M_Y(t) = e^{\lambda(e^t - 1)},$$

and if we define  $p = \lambda/n$ , then  $M_X(t) = [1 + (e^t - 1)\lambda/n]^n$  such that  $M_X(t) \rightarrow M_Y(t)$  as  $n \rightarrow \infty$ .

**Theorem** For any constant  $a$  and  $b$ , the mgf of the random variable  $aX + b$  is given by

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

*Proof:*

By definition,

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



## Lecture 22: Oct 18

### Last time

- Moment Generating Function

### Today

- Internal “evaluation” open where you can make anonymous suggestions
- Common Discrete Distribution

### Common Discrete Distribution

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.
- In this course, we emphasize parametric models.

**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:

$$F_X(x) = \Pr(X \leq x) = \begin{cases} 0 & x < 1 \\ [x]/N & 1 \leq x < N \\ 1 & N \leq x \end{cases}$$

- moments:

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

**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:  $\Pr(Y = 1) = p$  and  $\Pr(Y = 0) = 1 - p$ . We can write this as:

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

- what are the cdf, mean and variance?

**Binomial Distribution** A  $Binomial(n, p)$  random variable  $X$  is defined as the number of successes in  $n$  i.i.d. (independent, identically distributed) Bernoulli trials, each with probability  $p$  of success:

$$X = \sum_{i=1}^n Y_i, \quad Y_1, \dots, Y_n \stackrel{\text{i.i.d.}}{\sim} Bernoulli(p)$$

- Sample space:  $\{0, 1, \dots, n\}$
- pmf:

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

- cdf:

$$F_X(x) = \sum_{s=0}^x \binom{n}{s} p^s (1-p)^{n-s} \quad (\text{no closed form})$$

**Poisson Distribution** The Poisson distribution was derived by the French mathematician Poisson in 1837 as a limiting version of the binomial distribution. The Poisson distribution is often used to model the number of occurrences in a given time interval. One of the basic assumptions on which the Poisson distribution is built is that, for small time intervals, the probability of an arrival is proportional to the length of waiting time. This makes it a reasonable model for situations such as waiting for a bus, waiting for customers to arrive in a bank.

The Poisson distribution has a single parameter  $\lambda$ , sometimes called the intensity parameter. A Poisson random variable  $X$ , takes values in the nonnegative integers with pmf

$$\Pr(X = x|\lambda) = \frac{e^{-\lambda} \lambda^x}{x!}, \quad x = 0, 1, \dots$$

To see that  $\sum_{x=0}^{\infty} P(X = x|\lambda) = 1$ , recall the Taylor series expansion of  $e^\lambda = \sum_{i=0}^{\infty} \frac{\lambda^i}{i!}$ . Thus

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

What is the mean and variance of  $X$ ?

$$\begin{aligned} EX &= \sum_{x=0}^{\infty} x \frac{e^{-\lambda} \lambda^x}{x!} \\ &= \sum_{x=1}^{\infty} \frac{e^{-\lambda} \lambda^x}{(x-1)!} \\ &= \lambda \sum_{x=1}^{\infty} \frac{e^{-\lambda} \lambda^{x-1}}{(x-1)!} \\ &= \lambda \sum_{y=0}^{\infty} \frac{e^{-\lambda} \lambda^y}{y!} \\ &= \lambda \end{aligned}$$

Similarly

$$\begin{aligned} EX^2 &= \sum_{x=0}^{\infty} x^2 \frac{e^{-\lambda} \lambda^x}{x!} \\ &= \sum_{x=1}^{\infty} x \frac{e^{-\lambda} \lambda^x}{(x-1)!} \\ &= \sum_{x=1}^{\infty} \frac{e^{-\lambda} \lambda^x}{(x-1)!} + \sum_{x=2}^{\infty} \frac{e^{-\lambda} \lambda^x}{(x-2)!} \\ &= \lambda + \lambda^2 \end{aligned}$$

So that

$$\text{Var}(X) = EX^2 - (EX)^2 = \lambda$$

- Sample space:  $\{0, 1, \dots\}$
- pmf:  $\Pr(X = x) = \frac{e^{-\lambda} \lambda^x}{x!}$
- cdf:  $F_X(x) = \sum_{s=0}^x \frac{e^{-\lambda} \lambda^s}{s!}$

## Lecture 23: Oct 20

### Last time

- Common Discrete Distribution

### Today

- Internal “evaluation” open where you can make anonymous suggestions
- Common Discrete Distribution
- Common Continuous Distribution

### Common Discrete Distribution

**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 (one parameterization). Then,  $X$  has a geometric distribution.

- Sample space:  $\{1, 2, \dots\}$
- pmf:

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

- cdf:

$$F(x) = \Pr(X \leq x) = 1 - (1-p)^x$$

- Moments:

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

**Memoryless property.** Suppose  $k > i$ , then

$$\Pr(X > k | X > i) = \Pr(X > k - i)$$

*Proof:*

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

**Example** Suppose  $X$  is number of years you live, and  $X$  follows a geometric distribution, then

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

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

## Lecture 24: Oct 23

### Last time

- Common Discrete Distribution

### Today

- Internal “evaluation” open where you can make anonymous suggestions
- Common Discrete Distribution
- Common Continuous Distribution

### Common Discrete Distribution

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

### Notes

- Why the name? See Casella & Berger 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

**Other parameterizations** The negative binomial distribution is sometimes defined in terms of the random variable  $Y$  = number of failures before the  $r$ th success. Then

- Sample space:  $\{0, 1, 2, \dots\}$
- pmf

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

- cdf: no closed form
- Expectation:  $EY = r(1 - p)/p$
- Variance:  $Var(Y) = r(1 - p)/p^2$

**Negative binomial vs. Poisson** The negative binomial distribution is often good for modeling count data as an alternative to the Poisson. In the previous parameterization, define

$$\lambda = \frac{r(1 - p)}{p} \iff p = \frac{r}{r + \lambda}$$

Then we have

$$\begin{aligned} EX &= \lambda \\ Var(X) &= \frac{\lambda}{p} = \lambda(1 + \frac{\lambda}{r}) = \lambda + \frac{\lambda^2}{r} \end{aligned}$$

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.

In the previous parameterization, the pmf becomes

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

Letting  $r \rightarrow \infty$ , we get

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

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

## Lecture 25: Oct 25

Last time

- Common Discrete Distribution

Today

- Common Continuous Distribution

### Common continuous distributions

Uniform Distribution A random variable  $X$  having a pdf

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

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

The cdf is:

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

- Unifrom;  $Y \sim U[a, b]$
- sample space:  $[a, b]$
- pdf:

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

- cdf:

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

- moments:

$$E(Y) = (a + b)/2$$
$$Var(Y) = \frac{(b - a)^2}{12}$$

Notes

- The uniform extends to the continuous case the idea of equally likely outcomes.
- If  $Y \sim U[0, 1]$ , then  $a + (b - a)Y \sim U[a, b]$

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

- sample space:  $x \geq 0$
- pdf:

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

- cdf:

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

- moments:

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



## Lecture 26: Oct 27

Last time

- Common Continuous Distribution

Today

- Common Continuous Distribution

### Common continuous distributions

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

- sample space:  $x \geq 0$
- pdf:

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

- cdf:

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

- moments:

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

**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.

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

and the number of events in  $[0, t]$  follows a Poisson distribution with parameter  $\lambda t$ . Therefore,

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

The cdf of  $Y$  is

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

and hence the density is  $f(t) = \lambda e^{-\lambda t}$ .

**Alternative parameterization** Many books write the density as

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

so that  $E(Y) = \theta$  and  $Var(Y) = \theta^2$ . In this case  $\theta = 1/\lambda$  is called the *mean parameter*, while  $\lambda = 1/\theta$  is called the *rate parameter*.

**Memoryless property** The exponential has a memoryless property, just like the geometric.

$$\Pr(Y > s + t | Y > t) = \Pr(Y > s)$$

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  $\lambda$  is constant. (In the geometric case, the probability,  $p$  of getting an event in every discrete time unit is constant).

## Lecture 27: Oct 30

Last time

- Common Continuous Distribution

Today

- Show HW4 Q9
- Common Continuous Distribution

### Common continuous distributions

**Shifted exponential** Let  $X \sim \text{Exp}(\lambda)$  and  $Y = X + v, v \in \mathbb{R}$ . Then,  $Y$  has the *shifted exponential distribution* with pdf:

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

Interpretation:

- $v > 0$ : Event is delayed
- $v < 0$ : The news of the event is delayed

Does the shifted exponential maintain the memoryless property?

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

$$f(x) = \frac{1}{2} \lambda e^{-\lambda|x|}, \quad x \in \mathbb{R}$$

**Laplace distribution** 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 *Laplace distribution* with pdf:

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

with moments

$$EY = \mu, \quad \text{Var}(Y) = 2\sigma^2$$

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

## Lecture 28: Nov 1

### Last time

- Show HW4 Q9
- Common Continuous Distribution

### Today

- Common Continuous Distribution

### Common continuous distributions

**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:  $\mathbb{R} = (-\infty, \infty)$
- pdf: For  $Y \sim N(\mu, \sigma^2)$ ,

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

- cdf: There is no closed form.
- When  $\mu = 0$  and  $\sigma = 1$ , the distribution is called *standard normal*:

$$\Phi(y) = \Pr(Y \leq y), \quad \Phi(-y) = 1 - \Phi(y)$$

- Mean:

$$EY = \mu$$

- Variance:

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

- Higher central moments:

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

- In particular:

$$\begin{aligned} \mu_3 &= E(Y - \mu)^3 = 0 \text{ (Skewness)} \\ \mu_4 &= E(Y - \mu)^4 = 3\sigma^4 \end{aligned}$$

- Moment generating function:

$$M_Y(t) = \exp(\mu t + \sigma^2 t^2 / 2)$$

## Lecture 29: Nov 3

Last time

- Common Continuous Distribution

Today

- Common Continuous Distribution

### Common continuous distributions

**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:  $\mathbb{R} = (-\infty, \infty)$
- pdf: For  $Y \sim N(\mu, \sigma^2)$ ,

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

- cdf: There is no closed form.
- When  $\mu = 0$  and  $\sigma = 1$ , the distribution is called *standard normal*:

$$\Phi(y) = \Pr(Y \leq y), \quad \Phi(-y) = 1 - \Phi(y)$$

- Mean:

$$EY = \mu$$

- Variance:

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

- Higher central moments:

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

- In particular:

$$\begin{aligned} \mu_3 &= E(Y - \mu)^3 = 0 \text{ (Skewness)} \\ \mu_4 &= E(Y - \mu)^4 = 3\sigma^4 \end{aligned}$$

- Moment generating function:

$$M_Y(t) = \exp(\mu t + \sigma^2 t^2 / 2)$$

### Standardization

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

Shifting and scaling:

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

### Notes

- Normal distribution is useful in many practical settings. E.g. measurement error.
- Plays an important role in *sampling distributions* in *large samples*, since the Central Limit Theorem says that the 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.  $\chi^2$ ,  $t$ ,  $F$ ). We will briefly present the pdf's and sample spaces of these distributions.

**$\chi^2$  distribution** If  $Z \sim N(0, 1)$ , then  $X = Z^2$  has the  $\chi^2$  distribution with 1 degree of freedom. More generally, we have the  $\chi^2$  distribution with  $v$  degrees of freedom with pdf:

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

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

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

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

### Facts about the Gamma function

- $\Gamma(a+1) = a\Gamma(a), a > 0$
- $\Gamma(1) = 1$
- $\Gamma(n) = (n-1)!$
- $\Gamma(1/2) = \sqrt{\pi}$

**Student's  $t$  and  $F$  distributions**  $Y$  has a  $t_k$  distribution ( $t$  with  $k$  degrees of freedom) if its pdf can be written as:

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

$Y$  has an  $F(v_1, v_2)$  distribution if its pdf can be written as:

$$f(y) = \frac{(v_1/v_2)\Gamma[(v_1 + v_2)/2] (v_1 y/v_2)^{v_1/2-1}}{\Gamma(v_1/2)\Gamma(v_2/2)(1 + v_1 y/v_2)^{(v_1+v_2)/2}}, \quad 0 \leq y < \infty$$

There are many important properties and relationships between these three distributions (e.g.,  $\chi_k^2$  is the distribution of the sum of the squares of  $k$  independent standard normals).

**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 gamma function,

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

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

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

- MGF:

$$M_Y(t) = \left( \frac{1}{1 - t/\lambda} \right)^a, \quad t < \lambda$$

## Lecture 30: Nov 6

Last time

- Common Continuous Distribution

Today

- HW4 Q1.3
- HW4 Q6
- Common Continuous Distribution

### Common continuous distributions

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 gamma function,

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

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

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

- MGF:

$$M_Y(t) = \left( \frac{1}{1 - t/\lambda} \right)^a, \quad t < \lambda$$

Another parameterization   Same as the exponential distribution, we can let  $\beta = \frac{1}{\lambda}$ , then we have

- pdf:

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

- moments:

$$\begin{aligned} EX &= \alpha\beta \\ \text{Var}(X) &= \alpha\beta^2 \end{aligned}$$



- MGF:

$$M_Y(t) = \left( \frac{1}{1 - t\beta} \right)^a, \quad t < \frac{1}{\beta}$$

Notes:

- The special case  $a = 1$  corresponds to an *exponential*( $\lambda$ )
- The parameter  $a$  is known as the *shape parameter*, since it most influences the peakedness of the distribution.
- The parameter  $\beta$  is called the *scale parameter* since most of its influence is on the spread of the distribution.
- The special case  $\text{Gamma}(a = n/2, \lambda = 1/2)$ , for integer  $n$ , corresponds to the  $\chi_n^2$  distribution with  $n$  degrees of freedom.
- The gamma distribution can be derived as the sum of  $a$  independent *exponential*( $\lambda$ ) distributions.

**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 Beta function,

$$B(a, b) = \int_0^1 x^{a-1}(1-x)^{b-1} dx = \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)},$$

and  $\Gamma(a)$  is the gamma function. Note that if  $a$  and  $b$  are integers, then  $B(a, b)$  can be calculated in closed form.

- cdf: In general, there is no closed form, except if  $a$  and  $b$  are integers.
- moments:

$$EY = \frac{a}{a+b}$$

$$\text{Var}(Y) = \frac{ab}{(a+b)^2(a+b+1)}$$

The beta distribution is very flexible, and can take a wide variety of shapes by varying its parameters.

- Special case:  $\text{Beta}(1, 1) = U(0, 1)$ .

Omitted distributions: Weibull distribution, and Cauchy distribution.

## Lecture 31: Nov 8

Last time

- HW4 Q1.3
- HW4 Q6
- Common Continuous Distribution

Today

- Location Scale Families
- Exponential Family

### Location and Scale families

Let  $Z$  be a continuous random variable with pdf  $f(z)$ . Define the class of rvs

$$X_{\mu,\sigma} = \sigma Z + \mu, \quad \mu \in \mathbb{R}, \sigma > 0$$

Then

1.  $X_{\mu,\sigma}$  has pdf

$$f_{\mu,\sigma}(x) = \frac{1}{\sigma} f\left(\frac{x - \mu}{\sigma}\right)$$

- 2.

$$E(X) = \sigma E(Z) + \mu, \quad \text{Var}(X) = \sigma^2 \text{Var}(Z)$$

3. The variable  $Z = X_{0,1}$  is called the *generator* and is a member of the class.

### Location families and scale families

- The family of pdfs  $f_{\mu,\sigma}(x)$  is called a *location-scale* family where  $\mu$  is called the *location parameter*, and  $\sigma$  is called the *scale parameter*.
- The family of pdfs

$$f_{\mu,1}(x) = f(x - \mu)$$

with  $\sigma = 1$  is called a *location* family.

- The family of pdfs

$$f_{0,\sigma}(x) = \frac{1}{\sigma} f\left(\frac{x}{\sigma}\right)$$

with  $\mu = 0$  is called a *scale* family.

**Example (Exponential location family)** Let  $f(x) = e^{-x}, x \geq 0$ , and  $f(x) = 0, x < 0$ . To form a location family we replace  $x$  with  $x - \mu$  to obtain

$$f(x|\mu) = \begin{cases} e^{-(x-\mu)} & x - \mu \geq 0 \\ 0 & x - \mu < 0 \end{cases}$$

$$= \begin{cases} e^{-(x-\mu)} & x \geq \mu \\ 0 & x < \mu \end{cases}$$

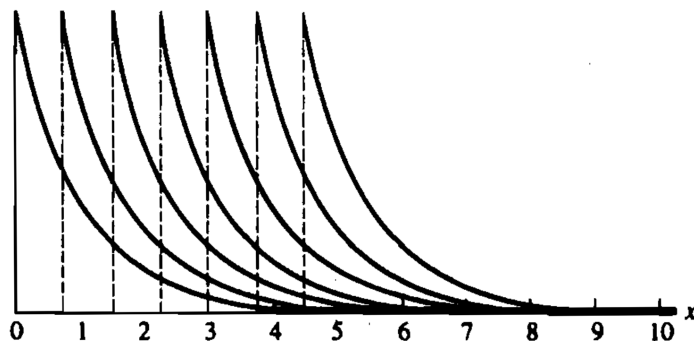


Figure 3.5.2. *Exponential location densities*

Figure 30.3: Figure 3.5.2. Exponential location densities.

As shown in the above graph, the densities are shifted. Now the positive part of the density starts at  $\mu$  rather than at 0. If  $X$  measures time, then  $\mu$  might be restricted to be nonnegative so that  $X$  will be positive with probability 1 for every value of  $\mu$ . In this type of model, where  $\mu$  denotes a bound on the range of  $X$ ,  $\mu$  is sometimes called a *threshold parameter*.

The effect of introducing the scale parameter  $\sigma$  is either to stretch ( $\sigma > 1$ ) or to contract ( $\sigma < 1$ ) the graph of  $f(x)$  while still maintaining the same basic shape of the graph. This is illustrated in the Figure below.

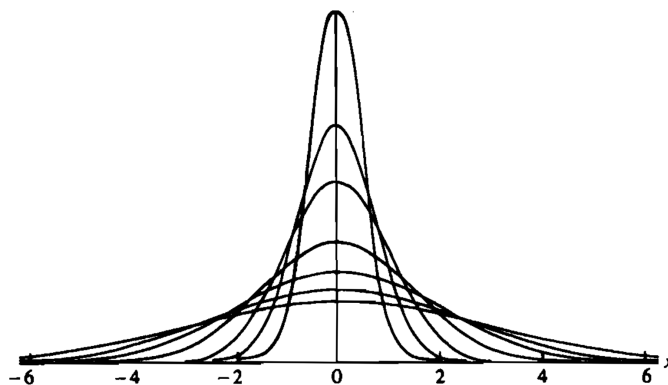


Figure 30.4: Figure 3.5.3. Members of the same scale family

**Exponential Families** A family of pdfs or pmfs with vector parameter  $\boldsymbol{\theta}$  is called an *exponential family* if it can be expressed as

$$f(x|\boldsymbol{\theta}) = h(x)c(\boldsymbol{\theta})\exp\left(\sum_{j=1}^k w_j(\boldsymbol{\theta})t_j(x)\right), \quad x \in S \subset \mathbb{R} \quad (1)$$

where  $S$  is not defined in terms of  $\boldsymbol{\theta}$ ,  $h(x)$ ,  $c(\boldsymbol{\theta}) \geq 0$  and the functions are just functions of the parameters specified; i.e.  $h$  is free of  $\boldsymbol{\theta}$ ,  $c(\boldsymbol{\theta})$  is free of  $x$ , etc...

Examples:

- One-dimensional: Exponential, Poisson
- Two-dimensional: Gaussian

Exponential family parameterizations are unique except for multiplying constant factors.

**Example: Gaussian** Let  $f(x|\mu, \sigma^2)$  be the  $n(\mu, \sigma^2)$  family of pdfs, where  $\boldsymbol{\theta} = (\mu, \sigma)$ . Then

$$\begin{aligned} f(x|\mu, \sigma^2) &= \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right) \\ &= \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{\mu^2}{2\sigma^2}\right) \exp\left(-\frac{x^2}{2\sigma^2} + \frac{\mu x}{\sigma^2}\right) \end{aligned}$$

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  $(\mu, \sigma^2) \in \mathbb{R} \times (0, \infty)$ .

**Example: Binomial** Let  $f(x|p)$  be the *binomial*( $n, p$ ),  $0 < p < 1$  family of pmfs.

$$\begin{aligned} f(x|p) &= \binom{n}{x} p^x (1-p)^{n-x} = \binom{n}{x} (1-p)^n \left[\frac{p}{1-p}\right]^x \\ &= \binom{n}{x} (1-p)^n \exp\left[\log\left(\frac{p}{1-p}\right) x\right] \end{aligned}$$

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. Also,  $p$  cannot be 0 or 1. Otherwise, the range changes.

**More examples** The following distributions belong to Exponential families:

- Continuous: exponential, Gaussian, gamma, beta,  $\chi^2$
- Discrete: Poisson, geometric, binomial (fixed # trials), negative binomial (fixed # successes)

The following distributions not exponential families:

- Continuous:  $t$ ,  $F$ , uniform E.g.:  $X \sim U(0, \theta)$

$$f_X(x) = \theta^{-1} 1(0 < x < \theta)$$

- Discrete: uniform, hypergeometric

**Theorem** If  $X$  is a random variable with pdf or pmf of the form 3, then

$$\begin{aligned} E \left( \sum_{i=1}^k \frac{\partial w_i(\boldsymbol{\theta})}{\partial \theta_j} t_i(X) \right) &= -\frac{\partial}{\partial \theta_j} \log c(\boldsymbol{\theta}) \\ \text{Var} \left( \sum_{i=1}^k \frac{\partial w_i(\boldsymbol{\theta})}{\partial \theta_j} t_i(X) \right) &= -\frac{\partial^2}{\partial \theta_j^2} \log c(\boldsymbol{\theta}) - E \left( \sum_{i=1}^k \frac{\partial^2 w_i(\boldsymbol{\theta})}{\partial \theta_j^2} t_i(X) \right). \end{aligned}$$

Although these equations may look formidable, when applied to specific cases they can work out quite nicely. Their advantage is that we can replace integration or summation by differentiation, which is often more straightforward.

**Example (Normal exponential family)** Let  $f(x|\mu, \sigma^2)$  be the  $N(\mu, \sigma^2)$  family of pdfs, where  $\boldsymbol{\theta} = (\mu, \sigma)$ ,  $-\infty < \mu < \infty, \sigma > 0$ . Then

$$\begin{aligned} f(x|\mu, \sigma^2) &= \frac{1}{\sqrt{2\pi}\sigma} \exp \left( -\frac{(x - \mu)^2}{2\sigma^2} \right) \\ &= \frac{1}{\sqrt{2\pi}\sigma} \exp \left( -\frac{\mu^2}{2\sigma^2} \right) \exp \left( -\frac{x^2}{2\sigma^2} + \frac{\mu x}{\sigma^2} \right) \end{aligned}$$

Define

$$\theta_1 = \frac{1}{\sigma^2} > 0, \quad \theta_2 = \frac{\mu}{\sigma^2} \in \mathbb{R}$$

Then

$$f_X(x) = \frac{\sqrt{\theta_1}}{\sqrt{2\pi}} \exp \left( -\frac{\theta_2^2}{2\theta_1} \right) \exp \left( -\theta_1 \frac{x^2}{2} + \theta_2 x \right)$$

and

$$\begin{aligned} h(x) &= 1 \text{ for all } x; \\ c(\boldsymbol{\theta}) &= c(\theta_1, \theta_2) = \exp \left( -\frac{\theta_2^2}{2\theta_1} \right), \quad (\theta_1, \theta_2) \in (0, \infty) \times \mathbb{R} \\ w_1(\boldsymbol{\theta}) &= \theta_1 & t_1(x) &= -x^2/2 \\ w_2(\boldsymbol{\theta}) &= \theta_2 & t_2(x) &= x \end{aligned}$$

Therefore, by the above theorem

$$\begin{aligned} E(X) &= -\frac{\partial}{\partial \theta_2} \log c(\boldsymbol{\theta}) = \frac{\theta_2}{\theta_1} = \mu \\ Var(X) &= -\frac{\partial^2}{\partial \theta_2^2} \log c(\boldsymbol{\theta}) = -\frac{1}{\theta_1} = \sigma^2 \end{aligned} \tag{2}$$

## Lecture 32: Nov 10

### Last time

- Location Scale Families
- Exponential Family

### Today

- Exponential Family (C&B 3.4)
- Multiple Random Variables (Chapter 4)

**Exponential Families** A family of pdfs or pmfs with vector parameter  $\boldsymbol{\theta}$  is called an *exponential family* if it can be expressed as

$$f(x|\boldsymbol{\theta}) = h(x)c(\boldsymbol{\theta})\exp\left(\sum_{j=1}^k w_j(\boldsymbol{\theta})t_j(x)\right), \quad x \in S \subset \mathbb{R} \quad (3)$$

where  $S$  is not defined in terms of  $\boldsymbol{\theta}$ ,  $h(x)$ ,  $c(\boldsymbol{\theta}) \geq 0$  and the functions are just functions of the parameters specified; i.e.  $h$  is free of  $\boldsymbol{\theta}$ ,  $c(\boldsymbol{\theta})$  is free of  $x$ , etc...

Examples:

- One-dimensional: Exponential, Poisson
- Two-dimensional: Gaussian

Exponential family parameterizations are unique except for multiplying constant factors.

**Example: Binomial** Let  $f(x|p)$  be the *binomial*( $n, p$ ),  $0 < p < 1$  family of pmfs.

$$\begin{aligned} f(x|p) &= \binom{n}{x} p^x (1-p)^{n-x} = \binom{n}{x} (1-p)^n \left[ \frac{p}{1-p} \right]^x \\ &= \binom{n}{x} (1-p)^n \exp \left[ \log \left( \frac{p}{1-p} \right) x \right] \end{aligned}$$

Thus,

$$\begin{aligned} h(x) &= \binom{n}{x}, \quad x = 0, \dots, n \quad w_1(p) = \log \left( \frac{p}{1-p} \right) \\ c(p) &= (1-p)^n, \quad 0 < p < 1 \quad t_1(x) = x \end{aligned}$$

Note that this works when  $p$  is considered the parameter, while  $n$  is fixed. Also,  $p$  cannot be 0 or 1. Otherwise, the range changes.

**More examples** The following distributions belong to Exponential families:

- Continuous: exponential, Gaussian, gamma, beta,  $\chi^2$
- Discrete: Poisson, geometric, binomial (fixed # trials), negative binomial (fixed # successes)

The following distributions not exponential families:

- Continuous:  $t$ ,  $F$ , uniform E.g.:  $X \sim U(0, \theta)$

$$f_X(x) = \theta^{-1} 1(0 < x < \theta)$$

- Discrete: uniform, hypergeometric

**Theorem** If  $X$  is a random variable with pdf or pmf of the form 3, then

$$\begin{aligned} E \left( \sum_{i=1}^k \frac{\partial w_i(\boldsymbol{\theta})}{\partial \theta_j} t_i(X) \right) &= -\frac{\partial}{\partial \theta_j} \log c(\boldsymbol{\theta}) \\ \text{Var} \left( \sum_{i=1}^k \frac{\partial w_i(\boldsymbol{\theta})}{\partial \theta_j} t_i(X) \right) &= -\frac{\partial^2}{\partial \theta_j^2} \log c(\boldsymbol{\theta}) - E \left( \sum_{i=1}^k \frac{\partial^2 w_i(\boldsymbol{\theta})}{\partial \theta_j^2} t_i(X) \right). \end{aligned}$$

Although these equations may look formidable, when applied to specific cases they can work out quite nicely. Their advantage is that we can replace integration or summation by differentiation, which is often more straightforward.

**Example (Normal exponential family)** Let  $f(x|\mu, \sigma^2)$  be the  $N(\mu, \sigma^2)$  family of pdfs, where  $\boldsymbol{\theta} = (\mu, \sigma)$ ,  $-\infty < \mu < \infty, \sigma > 0$ . Then

$$\begin{aligned} f(x|\mu, \sigma^2) &= \frac{1}{\sqrt{2\pi}\sigma} \exp \left( -\frac{(x - \mu)^2}{2\sigma^2} \right) \\ &= \frac{1}{\sqrt{2\pi}\sigma} \exp \left( -\frac{\mu^2}{2\sigma^2} \right) \exp \left( -\frac{x^2}{2\sigma^2} + \frac{\mu x}{\sigma^2} \right) \end{aligned}$$

Define

$$\theta_1 = \frac{1}{\sigma^2} > 0, \quad \theta_2 = \frac{\mu}{\sigma^2} \in \mathbb{R}$$

Then

$$f_X(x) = \frac{\sqrt{\theta_1}}{\sqrt{2\pi}} \exp \left( -\frac{\theta_2^2}{2\theta_1} \right) \exp \left( -\theta_1 \frac{x^2}{2} + \theta_2 x \right)$$

and

$$\begin{aligned} h(x) &= 1 \text{ for all } x; \\ c(\boldsymbol{\theta}) &= c(\theta_1, \theta_2) = \exp \left( -\frac{\theta_2^2}{2\theta_1} \right), \quad (\theta_1, \theta_2) \in (0, \infty) \times \mathbb{R} \\ w_1(\boldsymbol{\theta}) &= \theta_1 & t_1(x) &= -x^2/2 \\ w_2(\boldsymbol{\theta}) &= \theta_2 & t_2(x) &= x \end{aligned}$$



Therefore, by the above theorem

$$\begin{aligned} E(X) &= -\frac{\partial}{\partial \theta_2} \log c(\boldsymbol{\theta}) = \frac{\theta_2}{\theta_1} = \mu \\ \text{Var}(X) &= -\frac{\partial^2}{\partial \theta_2^2} \log c(\boldsymbol{\theta}) = -\frac{1}{\theta_1} = \sigma^2 \end{aligned} \quad (4)$$

## Joint and Marginal Distributions

In previous lectures, we have discussed probability models and computation of probability for events involving only one random variable. These are called *univariate models*.

In an experimental situation, it would be very unusual to observe only the value of one random variable. For example, in an experiment designed to gain information about some health characteristics of a population of people, the body weights of several people in the population might be measured. These different weights would be observations on different random variables, one for each person measured. Multiple observations could also arise because several physical characteristics were measured on each person. Thus, we need to know how to describe and use probability models that deal with more than one random variable at a time.

**Definition:** An  $n$ -dimensional random vector  $\mathbf{X} = (X_1, \dots, X_n)$  is a function from a sample space  $S$  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 each Borel set  $B$ ,

$$\Pr\{\mathbf{X} \in B\} = \Pr\{\mathbf{X}^{-1}(B)\} \quad (5)$$

where

$$\mathbf{X}^{-1}(B) = \{w : \mathbf{X}(w) \in B\}$$

**Example (Bivariate random variable)** 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	$(x, y)$	$\Pr(\text{outcome})$
(H, H, H)	(1, 3)	1/8
(H, T, H), (T, H, H)	(1, 2)	2/8
(H, H, T)	(0, 2)	1/8
(T, T, H)	(1, 1)	1/8
(T, H, T), (H, T, T)	(0, 1)	2/8
(T, T, T)	(0, 0)	1/8

**Definition** Two random variables  $X$  and  $Y$  are said to be jointly *discrete* if there is an associated *joint probability mass function*,

$$f_{X,Y}(x, y) = \Pr\{X = x, Y = y\}$$

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:

$$f_X(x) = \Pr(X = x) = \sum_y f_{X,Y}(x, y)$$

$$f_Y(y) = \Pr(Y = y) = \sum_x f_{X,Y}(x, y)$$

**Example** Back to the fair coin example again. From the definition, we can construct the joint pmf 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.