

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.

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

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).

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):

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

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? (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.

In total, there are $\binom{52}{4} = \frac{52!(52-4)!}{4!} = 270,725$ distinct groups. Therefore, the probability of event A is $\Pr(A) = \frac{1}{270,725}$.

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

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 .

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

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).

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):

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

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? (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.

In total, there are $\binom{52}{4} = \frac{52!(52-4)!}{4!} = 270,725$ distinct groups. Therefore, the probability of event A is $\Pr(A) = \frac{1}{270,725}$.

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

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 .