

# MATH 230A/STATS 310A NOTES

ARUN DEBRAY  
OCTOBER 1, 2014

These notes were taken in Stanford's Math 230A class (crosslisted as Stats 310A) in Fall 2014, taught by Perseus Diaconis. I live-T<sub>E</sub>Xed them using vim, and as such there may be typos; please send questions, comments, complaints, and corrections to [adebray@stanford.edu](mailto:adebray@stanford.edu).

## CONTENTS

1. Introduction and Examples of Probability Spaces: 9/22/14	1
2. The Strong Law of Coin Tossing: 9/24/14	4
3. The $\pi$ - $\lambda$ Theorem: 9/28/14	6
4. Independence: 10/1/14	8

## 1. INTRODUCTION AND EXAMPLES OF PROBABILITY SPACES: 9/22/14

*"I like inequalities. I hope you learn to like them too."*

In this class, we're going to develop measure theory; probability is one of the subjects in which it comes alive. So there will be proofs and constructions of many standard measure-theory ideas, but the explicit goal is for probability.

The course will cover roughly the first twenty sections of the book; there will be weekly homework of about six problems, due in class on Mondays. However, about 25% of the class material is not in the textbook, so come to class and take notes!

Before discussing the structure of the course, let's do some probability. The nice thing about probability is that it's very easy to say what it's about. In probability, there is a (for today, finite or countable) set  $\Omega$ , weights  $P(\omega) \geq 0$  such that

$$\sum_{\omega \in \Omega} P(\omega) = 1,$$

and a subset  $A \subseteq \Omega$ . These things are all told to us; the question of probability is to compute or approximate the probability

$$P(A) = \sum_{\omega \in A} P(\omega).$$

We'll refine this, develop it, etc., but it looks really trivial from this viewpoint, doesn't it?

**Example 1.1.** The Birthday Problem is a well-known problem in probability. It asks, "how many people do we need in a group to have a 50% or greater chance that two or more have the same birthday?" The answer is about 23, which is viscerally surprising. Even in this class, it's likely that two people share a birthday.

There are many variations; this is one of the building blocks of combinatorial probability. For example, you could do the same thing with the second hands on a watch (which, for the sake of the argument, are i.i.d.); this would be known as a *birthday problem with sixty categories*.

For a birthday problem with  $c$  categories and  $n$  people,  $\Omega$  is the set of  $n$ -tuples of numbers from 1 to  $c$ :  $\Omega = \{(\omega_1, \dots, \omega_n) : \omega_i \in \{1, \dots, c\}\}$ . Then,  $P(\omega) = 1/c^n$ , and for  $A$ , we'll calculate the probability of failure (since this is easier to deal with):  $A = \{\omega : \omega_i \neq \omega_j \text{ for all } i, j\}$ .

Then,

$$P(A) = \sum_{\omega \in A} P(\omega) = \frac{|A|}{c^n}.$$

Awesome. So what's the size of  $A$ ? The first person has  $c$  choices, then the next person needs to choose something different, so there are  $c - 1$  choices, then  $c - 2$  (different from the first two), and so on, so  $|A| = c(c - 1)(c - 2) \cdots (c - n + 1)$ . Thus,

$$P(A) = \frac{c(c - 1) \cdots (c - n + 1)}{c^n},$$

so the probability that two people share a birthday (the complement of  $A$ ) is

$$\left(1 - \frac{1}{c}\right)\left(1 - \frac{2}{c}\right) \cdots \left(1 - \frac{n-1}{c}\right).$$

This is a fine answer; it's 2014, so we can plug it into a computer and have fun. But it's not at all understandable to people. Recall that asymptotically,  $\log(1-x) \sim -x$  (since it's the first term in the power series).<sup>1</sup> So rewrite

$$\begin{aligned} P(A) &= e^{\sum_{i=1}^n \log(1-1/c)} \\ &\sim e^{-\sum_{i=1}^n i/c} \\ &= e^{-\binom{n}{2}/c}. \end{aligned}$$

Thus, set this equal to  $1/2$ , so

$$\frac{\binom{n}{2}}{c} = \log 2 \doteq 0.69$$

(in this class, all logs will be base  $e$ ), and therefore  $n^2/2c = 0.69$ , or  $n = 1.2\sqrt{c}$ .

This is human-understandable:  $n$  should be about 1.2 times the square root of  $c$ . For example, if  $c = 365$ , then  $1.2\sqrt{365} = 22.9$ , so the goal is 23 students. If we wanted it with 95% probability, one can calculate that  $n = 2.5\sqrt{c}$ .

In probability, these kinds of heuristics are commonplace: reasoning intuitively, but then backing it up with mathematics.

Let's refine this approximation:  $\log(1-x) = -x + O(x^2)$ , since  $-x - x^2 \leq \log(1-x) \leq -x$  when  $0 \leq x \leq 1/2$ . Thus,

$$\begin{aligned} P(A) &= e^{\sum_{i=1}^n \log(1-1/c)} \\ &= e^{-\sum_i i/c + O((1/c)^2)} \\ &= e^{-\binom{n}{2}/c + O(n^3/c^2)}. \end{aligned}$$

This is a rigorous mathematical result:

**Theorem 1.2.** If  $n, c \rightarrow \infty$  such that  $n^3/c^2 \rightarrow 0$  and  $\binom{n}{2}/c \rightarrow 1$ , then  $P(A) \rightarrow e^{-1}$ .

**Exercise 1.** Extend this to three people: how many people are needed in a room so there are even odds that three of them share a birthday?

This is one of a million versions of the birthday problem. For example, the birthday distribution is nonuniform (for example, people are more likely to be born on weekdays, because one birth in five is induced, and doctors don't like working on weekends). By the end of the course, we'll have tools (e.g. Poisson approximation and Stein's process) to analyze this where  $P(\omega)$  is nonuniform.

**Example 1.3.** Put  $N$  points at random in the unit square  $[0, 1]^2$ , and an  $\varepsilon$ -ball around each one. Then, what are the odds that the square is covered? What's the size of the largest disc covered? And so on. This is an example of a continuous distribution; it is common, for example, to put probabilities on  $C([0, 1])$ , the space of continuous functions on the interval. One may even put probability distributions on manifolds, such as

$$\Omega = \{x_1, \dots, x_{35} \in \mathbb{R}_+^{35} : \sum x_n = s, \prod x_n = p\},$$

which is a 33-dimensional manifold that appears in the study of the gamma distribution; or on the set of probability measures  $M_1(\mathbb{R})$  with certain properties.

Let's take a halfway house,  $\Omega = [0, 1]$ , and define probabilities on a certain class of sets: left-open, right-closed sets  $(a_i, b_i]$ . Then, if  $A = \bigcup_{i=1}^n I_n$ , where the  $I_n$  are disjoint, then define

$$P(A) \stackrel{\text{def}}{=} \sum_{i=1}^n (b_i - a_i).$$

This can be used as a model for lots of different things, e.g. fair coin tossing. Right now, we're only considering finite sums of disjoint sets. For fair coin tossing, write each number in  $[0, 1]$  as its nonterminal binary expansion,  $\omega = 0.\omega_1\omega_2\omega_3\dots$ , e.g.  $1/2 = 0.01111\dots$ . Let  $d_n : [0, 1] \rightarrow \{0, 1\}$  denote the function which retrieves the  $n^{\text{th}}$  binary digit from a number, so that  $\omega = \sum_{n=1}^{\infty} d_n(\omega)/2^n$ .

Each  $d_i$  defines intervals where  $d_i(\omega) = 1$ ; the set  $A = \{\omega : d_i(\omega) = 1\}$  has probability  $1/2$  for all  $i$ , and similarly  $P(d_1 = d_2 = 1) = 1/4$ , and so forth. Then (as will be explicated on the homework),  $P(d_1 = \varepsilon_1, \dots, d_n = \varepsilon_n) = 1/2^n$ .

This is a nice model for flipping a coin, and is even the standard model, but... coins have physics, air resistance, coefficients of restitution, and so on. The difference is that it's a model, so we should calculate with it and then compare to the real world. But it does seem to model coins very well.<sup>2</sup>

<sup>1</sup>The notation  $f(n) \sim g(n)$ , said " $f(n)$  is asymptotic to  $g(n)$ ," means that  $\lim_{n \rightarrow \infty} f(n)/g(n) = 1$ .

<sup>2</sup>"I am possibly the world's most accomplished expert on flipping coins... don't miss a chance to ask me about that."

**Theorem 1.4** (Bernoulli's Law of Large Numbers (1713)). *For this model of coin tossing and for all  $\varepsilon > 0$ , the probability*

$$P\left(\left|\frac{d_1 + \cdots + d_n}{n} - \frac{1}{2}\right| > \varepsilon\right) \rightarrow 0$$

as  $n \rightarrow \infty$ .

This implies that the proportion of heads that happens as one flips more and more coins should be fair, tending to 1/2.

*Proof of Theorem 1.4.* To prove the theorem, we know that

$$A = \left\{ \omega : \left| \frac{1}{n} \sum d_i(\omega) - \frac{1}{2} \right| > \varepsilon \right\}$$

has a probability; we need to find it.

Define

$$r_i(\omega) = 2d_i(\omega) - 1 = \begin{cases} 1, & d_i = 1 \\ -1, & d_i = 0. \end{cases}$$

This is a simple renormalization, but makes life a bit easier. In particular, it's enough to prove that for all  $\varepsilon > 0$ ,

$$P\left(\left|\frac{1}{n} \sum r_i\right| > \varepsilon\right) \rightarrow 0.$$

We can use the Riemann integral, and in particular know that

$$\int_0^1 r_i(\omega) d\omega = 0,$$

and the different  $r_i$  are orthogonal:

$$\int_0^1 r_i(\omega) r_j(\omega) d\omega = \delta_{ij} = \begin{cases} 1, & i = j \\ 0, & i \neq j. \end{cases}$$

Thus, we can find the first and second moments of the sum.

$$\int_0^1 (\sum r_i) d\omega = \int_0^1 (\sum r_i)^2 d\omega = n.$$

Thus,

$$\begin{aligned} P\left(\left|\frac{1}{n} \sum r_i\right| > 2\varepsilon\right) &= P\left((\sum r_i)^2 > 4\varepsilon^2 n\right) \\ &\leq \frac{1}{4\varepsilon^2 n} \int (\sum r_i)^2 d\omega = \frac{1}{4\varepsilon^2 n}, \end{aligned} \tag{1}$$

which goes to 0 as  $n \rightarrow \infty$ . □

(1) comes from *Markov's inequality*:

**Lemma 1.5** (Markov's inequality). *If  $f : [0, 1] \rightarrow \mathbb{R}$  is a nonnegative step function, then for all  $a > 0$ ,*

$$P(\omega : f(\omega) \geq a) \leq \frac{1}{a} \int_0^1 f(\omega) d\omega.$$

*Proof.* Let  $A = \{\omega : f(\omega) \geq a\}$ ; then,

$$\int_0^1 f(\omega) d\omega \geq \int_A f(\omega) d\omega \geq a \cdot P(A). \quad \square$$

We'll also want another theorem. We want to say that

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n r_i(\omega) = 0,$$

but this isn't true: for example, if  $\omega = 0.111111\dots$ , then the limit is 1, and if  $\omega = 0.01100001111111\dots$  (each time, twice as many as before), the limit simply doesn't exist. So we need to make this theorem a little more precise.

**Definition.** Define a set  $A \subset \Omega$  to be *negligible* if for any  $\varepsilon > 0$ ,  $A$  can be covered by countably many intervals with total length less than  $\varepsilon$ .

For example, any one-point set is negligible, as is  $\mathbb{Q}$ .

Now, we can state the theorem.

**Theorem 1.6** (Strong Law of Large Numbers (Borel)). *Except on a negligible  $N \subseteq \Omega$ , for all  $\omega \in N^c$ ,*

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n r_i(\omega) = 0.$$

*Remark.* Look at the set  $A = N^c$ , the set of places where this statement is true. It can be characterized as

$$A = \bigcap_{h=1}^{\infty} \bigcup_{m=1}^{\infty} \bigcap_{n=m}^{\infty} \left\{ \omega : \left| \frac{1}{n} \sum_{i=1}^n r_i(\omega) \right| < \frac{1}{h} \right\}.$$

This is very complicated, but we'll develop general methods for dealing with these sets.

This class will intersperse measure theory and probability, since the latter is more interesting and the former is very useful.

## 2. THE STRONG LAW OF COIN TOSSING: 9/24/14

*"The chalk knows how to do this. I need to stop thinking."*

Recall the setup we had yesterday: that  $\Omega = [0, 1]$ , an interval  $I$  is  $I = (a, b]$ , with  $P(I) = b - a$ , and the base-2 expansion of an  $\omega \in \Omega$  is  $\omega = \sum_{n=1}^{\infty} d_n(\omega)/2^n$ . Then, rescale these to the functions  $r_n(\omega) = 2d_n(\omega) - 1$ . Finally, we define

$$S_n(\omega) = \sum_{n=1}^{\infty} r_n(\omega).$$

**Theorem 2.1.**

$$\lim_{n \rightarrow \infty} \frac{s_n(\omega)}{n} = 0,$$

except on a negligible set.

Recall that a negligible set is a set that can be covered by a countable collection of intervals with total measure less than  $\varepsilon$  for any  $\varepsilon > 0$ .

*Proof.* Let  $B = \{\omega : \lim_n s_n/n^2 = 0\}$ , since  $s_n/n \rightarrow 0$  iff  $s_{n^2}/n^2 \rightarrow 0$ ; we want to show that  $B^c$  is negligible.

Choose a sequence  $d_n$  approaching 0 from above; then,

$$\left\{ \omega : \frac{|s_{n^2}|}{n^2} < d_n \text{ for large } n \right\} \subseteq B.$$

Thus,  $B^c$  is the set of things such that  $|s_{n^2}|/n^2 \geq d_n$  infinitely often, and thus is contained in  $\bigcup_{n=1}^{\infty} B_n$ , where

$$B_n = \left\{ \omega : \frac{|s_{n^2}|}{n^2} < d_n \right\}.$$

Since  $B_n$  is a disjoint union of integers, then by Markov's inequality,  $P(B_n) \leq 1/(n^2 \delta_n^2)$ . Thus, we're done if we choose  $\delta_n$  such that this converges, e.g.  $\delta_n = 1/n^{1/4}$ .

We still have to prove the result on subsequences: that  $s_n/n \rightarrow 0$  iff  $s_{n^2}/n^2 \rightarrow 0$ . Unfortunately, I'm really far from the board and can't read the professor's handwriting, so this is going to remain an exercise... and full of tricky typos.

It looks like we are assuming that  $y_{n^2}/n^2 \rightarrow 0$  and

$$\lim_{n \rightarrow \infty} \frac{(y_n - y_{k_n+1})}{k_n^2} \rightarrow 0,$$

where  $k_n = \lfloor \sqrt{n} \rfloor$ . Then, we want to show that  $y_n/n \rightarrow 0$ .

Since  $k_n^2 \leq n \leq (k_n + 1)^2$ , then

$$\begin{aligned} \left| \frac{y_n}{n} \right| &\leq \frac{|y_n| k_n^2}{n} \frac{|y_{k_n} - (y_{k_n} - y_n)|}{k_n^2} \\ &\leq \frac{|y_{k_n}| k_n^2}{n} \frac{|y_{k_n} - y_n|}{k_n^2} + \frac{|y_n - y_{k_n}|}{n}. \end{aligned}$$

Thus, if  $y_n/n$  is bounded, then we can expand over the bound and get the second term to go to zero.  $\square$

For the strong law, we used something more general than just Chebyshev's inequality; either you need better bounds (e.g. the method of fourth moments, as in the book), or subsequences as we used. The two statements are different, and it's important to see how: if  $P(j) = c/(j^2 \log |j|)$  for  $|j| \geq 2$ , then pick  $X_1, \dots, X_n$  independently from  $P_j$ . These are symmetric, so the mean tends to 0 as  $n \rightarrow \infty$ , but the strong law fails: we can't know how many to try until it's below a certain  $\varepsilon$ . It has no finite content.

So this probabilistic argument may be slightly confusing, but it stands as a motivating example of dealing with complicated sets probabilistically. The next topic will be the language that allows us to define probability spaces: algebras,  $\sigma$ -algebras, and so forth.

**Definition.** Let  $\Omega$  be a set; then, a *field of subsets* of  $\Omega$  is a collection  $\mathcal{F} \subseteq \Omega$  such that

- (1)  $\emptyset \in \mathcal{F}$ .

- (2) If  $F \in \mathcal{F}$ , then  $F^c \in \mathcal{F}$ .
- (3) If  $A, B \in \mathcal{F}$ , then  $A \cup B \in \mathcal{F}$ .

**Example 2.2.**

- (1) The simplest example is  $\mathcal{F} = \{\emptyset, \Omega\}$ .
- (2) If  $\Omega = (0, 1]$ , then we can take  $\mathcal{F}$  to be the set of all finite unions of disjoint intervals of half-open sets; the reason we use (and used above) half-open sets is so that complements work out correctly.

**Definition.** A  $\sigma$ -field is a field such that (3) is generalized to countable collections: if  $F_n \in \mathcal{F}$  for  $n = 1, 2, \dots$ , then  $\bigcup_{n=1}^{\infty} F_n \in \mathcal{F}$ .

Thus, all subsets of  $\Omega$  is a  $\sigma$ -field, and if  $\{\mathcal{F}_n\}_{n \in I}$  is a collection of  $\sigma$ -fields, then their intersection  $\bigcap_{n \in I} \mathcal{F}_n$  is also a  $\sigma$ -field.

Thus, we can take any collection  $\mathcal{O}$  of subsets of  $\Omega$  and define  $\mathcal{F}(\mathcal{O})$  to be the smallest  $\sigma$ -field containing  $\mathcal{O}$ , which is the intersection of all  $\sigma$ -fields containing  $\mathcal{O}$  (since that intersection itself is a  $\sigma$ -field containing  $\mathcal{O}$ , and is certainly the smallest).

This is right about where things get less trivial, but we also lost contact with physical intuition.

**Definition.** The *Borel sets* are the elements of the  $\sigma$ -field generated by  $\mathcal{O} = \{(a, b] : 0 < a < b \leq 1\}$ .

It's easy to work with Borel sets, but trying to describe them is very difficult; some are very complicated. This can be unsettling, but we'll see how it's possible to work with them anyways.

**Definition.** A *probability space*  $(\Omega, \mathcal{F}, P)$  is a set  $\Omega$ ,  $\sigma$ -algebra  $\mathcal{F}$  of subsets of  $\Omega$ , and a function  $P : \mathcal{F} \rightarrow [0, 1]$  such that

- (1)  $P(\emptyset) = 0$ ,
- (2)  $P(A^c) = 1 - P(A)$  for all  $A \in \mathcal{F}$ , and
- (3) If  $A_1, A_2, \dots \in \mathcal{F}$  are all pairwise disjoint, then

$$P\left(\bigcup_{n=1}^{\infty} A_n\right) = \sum_{n=1}^{\infty} P(A_n).$$

Apparently one of Perseus's friends has the license plate F OMEGA P.

We've learned some things since the Greeks; for example:

- (1) It's not possible to define a probability on the  $\sigma$ -algebra of all subsets of some sets  $\Omega$ . The professor's Halloween lecture will explain this further — this is when the monsters come out.
- (2) In some sense, measure-theoretic probability is a wonderful approximate theory.

Suppose  $\Omega$  is a set and  $\mathcal{F}_0$  is a field (not necessarily a  $\sigma$ -field) and  $P : \mathcal{F}_0 \rightarrow [0, 1]$  is a probability on  $\mathcal{F}_0$ . Then, for any  $A \subseteq \Omega$ , define

$$P^*(A) = \inf \sum_{i=1}^{\infty} P(F_i),$$

where the infimum is taken over all coverings  $F_1, \dots$  of  $A$  within  $\mathcal{F}_0$ . Basically, take any horrible set, and cover it by sets that we understand how to assign probabilities to. This is called the *outer measure* of  $A$ .

Here's a good idea in this field, given by Carathéodory: define a collection of *measurable sets*  $\mathcal{M}$  whose probability we can understand. Specifically, we say that  $M \in \mathcal{M}$  if for all  $E \subseteq \Omega$ ,  $P^*(E) = P^*(E \cap M) + P^*(E \cap M^c)$ . That is, it splits every test set, no matter how horrible.

All right, but how did anyone think of *that*? It works and is useful... I guess sometimes people just have ideas sometimes.

**Theorem 2.3.** Given an algebra  $\mathcal{F}$  of subsets of  $\Omega$  and a probability  $P : \mathcal{F} \rightarrow [0, 1]$ ,

- (1)  $\mathcal{M}$  is a  $\sigma$ -algebra containing  $\mathcal{F}$ , and
- (2)  $P^*|_{\mathcal{M}}$  is a probability, and is equal to  $P$ .
- (3)  $P^*$  is unique.

We'll do this in three steps: first showing that  $P^*(\emptyset) = 0$ , then that if  $A \subseteq B$ , then  $P^*(A) \leq P^*(B)$ , and finally that if  $\{A_n\}_{n=1}^{\infty}$  is any collection of sets, then

$$P^*\left(\bigcup_{n=1}^{\infty} A_n\right) \leq \sum_{n=1}^{\infty} P^*(A_n).$$

Let's look at the last one in more detail. For all  $\varepsilon > 0$ , there are coverings  $F_{i,k} \in \mathcal{F}$ , where  $\bigcup_k F_{i,k} \supseteq A_i$  and such that

$$\sum_{k=1}^{\infty} P(F_{i,k}) \leq P^*(A_i) + \frac{\varepsilon}{2}.$$

To be continued...

### 3. THE $\pi$ - $\lambda$ THEOREM: 9/28/14

*“By chalk, by math, by logic, it is defined!”*

Recall that if we have a ground set  $\Omega$  and a field  $\mathcal{F}$  of subsets, with a probability  $P$  defined on  $\mathcal{F}$ , then for all  $E \subseteq \Omega$ , we can define the outer measure as the infimum over all countable covers  $\{F_n\}$  with  $F_n \in \mathcal{F}$ :

$$P^*(E) = \inf \sum_1^\infty P(F_n).$$

Then, we defined

$$\mathcal{M} = \{A \subseteq \Omega : P^*(E) = P^*(A \cap E) + P^*(E \cap A^c) \text{ for all } E \subseteq \Omega\}.$$

Then, we stated Theorem 2.3, which provides that  $P^*$  is a probability on  $\mathcal{M}$ , but we have to prove it.

*Proof of Theorem 2.3.* First, we need to prove that  $\mathcal{M}$  is a field. For any  $A \subseteq \Omega$ ,  $P^*(E) = P^*((A \cap E) \cup (A \cap E^c)) \leq P^*(A \cap E) + P^*(A^c \cap E)$ , so  $A \in \mathcal{M}$  iff  $P^*(E) \geq P^*(A \cap E) + P^*(A^c \cap E)$ . Thus,  $\emptyset \in \mathcal{M}$  of course, and compliments are just as easy, since the expression is unchanged if  $A$  and  $A^c$  are switched. Thus, we need to show that if  $A, B \in \mathcal{M}$ , then  $A \cap B \in \mathcal{M}$ .

At this point, drawing a picture of a Venn diagram may be useful.

$$\begin{aligned} P^*(E) &= P^*(B \cap E) + P^*(B^c \cap E) \\ &= P^*(A \cap B \cap E) + P^*(A^c \cap B \cap E) + P^*(A \cap B^c \cap E) + P^*(A^c \cap B^c \cap E) \\ &= P^*(A \cap B \cap E) + P^*((A \cap B \cap E) \cup (A \cap B^c \cap E) \cup (B^c \cap A^c \cap E)) \\ &= P^*(A \cap B \cap E) + P^*((A \cap B)^c \cap E). \end{aligned}$$

The next thing we need to prove is countable additivity: if  $\{A_n\}_{n=1}^\infty \in \mathcal{M}$  and are pairwise disjoint, then for any  $E \subseteq \Omega$ ,

$$P\left(E \cap \bigcup_{n=1}^\infty A_n\right) = \sum_{n=1}^\infty P^*(E \cap A_n).$$

We'll show it for the finite case by induction, and then use limits to generalize to the countable case. The induction is obvious for  $n = 1$ , and for  $n = 2$ , if  $A_1 = A_2^c$ , then it's obvious and otherwise,

$$\begin{aligned} L^*(E \cap (A_1 \cup A_2)) &= P^*(E \cap (A_1 \cup A_2) \cap A_1) + P^*(E \cap (A_1 \cup A_2) \cap A_1^c) \\ &= P^*(E \cap A_1) + P^*(E \cap A_2). \end{aligned}$$

This argument also works for finite  $n$ : pick out one, and then use the above reasoning to reduce to the  $n - 1$  case.

For the countable case, we have

$$P^*\left(E \cap \bigcup_1^\infty A_n\right) \geq P^*\left(E \cap \bigcup_1^n A_n\right) = \sum_{j=1}^n P^*(A_j \cap E);$$

then, let  $n \rightarrow \infty$ .

The next thing we need to show is that  $\mathcal{M}$  is a  $\sigma$ -algebra, and  $P^*$  is countably additive on  $\mathcal{M}$ . Given any countable collection  $A_1, A_2, \dots \in \mathcal{M}$ , disjointify<sup>3</sup> them, letting  $A'_1 = A_1$ ,  $A'_2 = A_2 \cap (A'_1)^c$ , and so forth, in general letting

$$A'_n = A_n \cap \left(\bigcup_{i=1}^n A'_i\right)^c.$$

Thus, the unions of the  $A_n$  and the  $A'_n$  are the same; call this union  $A$ , and let  $F_n = \bigcup_{i=1}^n A'_i$ , which we know is in  $\mathcal{M}$  for each  $n$ . Thus,

$$\begin{aligned} P^*(E) &= P^*(E \cap F_n) + P^*(E \cap F_n^c) \geq P^*(E \cap F_n) + P^*(E \cap A^c) \\ &= \sum_{i=1}^n P^*(E \cap A'_i) + P^*(E \cap A^c). \end{aligned}$$

Thus,

$$\begin{aligned} P^*(E) &\geq \sum_1^\infty P^*(E \cap A'_n) + P^*(E \cap A^c) \\ &= P^*\left(E \cap \bigcup_1^\infty A'_n\right) + P^*(E \cap A^c) \\ &= P^*(E \cap A) + P^*(E \cap A^c). \end{aligned}$$

Moreover, we have already shown that  $P^*$  is additive on  $\mathcal{M}$ .

<sup>3</sup>Is this a word? I hope so.

Finally, we have to show that  $\mathcal{F} \subseteq \mathcal{M}$ . For any  $\varepsilon > 0$  and  $E \subseteq \Omega$ , choose an  $A \in \mathcal{F}$  and  $A_1, A_2, \dots \in \mathcal{F}$  such that  $E$  is covered by  $\{A_n\}$  but

$$\sum_1^\infty P(A_n) \leq P^*(E) + \varepsilon.$$

Set  $B_n = A \cap A_n$  and  $C_n = A^c \cap A_n$ , so

$$E \cap A \subset \bigcup_1^\infty B_n \quad \text{and} \quad E \cap A^c \subset \bigcup_1^\infty C_n.$$

Then,

$$\begin{aligned} P^*(A \cap E) + P^*(A^c \cap E) &= \sum_n P(B_n) + P(C_n) \\ &= \sum_n P(A_n) \leq P^*(E) + \varepsilon. \end{aligned}$$

This is about the only place where we used the actual definition of the outer measure.

Finally, we can prove that  $P^*(F) = P(F)$ ; we already know it's at least as much, and we can cover  $F$  by  $F_1, F_2, \dots$  in  $\mathcal{F}$ . Then,

$$P(F) \leq \sum_n P(F \cap F_n) \leq \sum_1^\infty P(F_n),$$

so  $P^*(F) \leq P(F)$ . □

Though this is the right definition, in that this theorem is demonstrated, it's very abstract: nobody actually uses it to compute, and it can be specialized to Lebesgue measure, where it's a bit more tractable.

In this case,  $\Omega = (0, 1]$ , and  $\mathcal{F}$  is taken to be the set of finite unions of disjoint intervals, with

$$P\left(\sum_{i=1}^n (a_i, b_i]\right) = \sum_{i=1}^n b_i - a_i.$$

Unfortunately, it's a bit  $\odot$  to show this,<sup>4</sup> but we can break it down into showing that if  $\bigcup_1^\infty I_n \subseteq I$ , then  $\sum(b_n - a_n) \leq b - a$  (where  $I = (a, b]$  and  $I_n = (a_n, b_n]$ ; and that if the inclusion is switched, then the inequality is also switched.

The first one follows from induction on  $n$ , at which point it's easy from disjointness, and the second point actually requires choosing an  $\varepsilon > 0$ , such that  $\varepsilon < b - a$  and

$$[a + \varepsilon, b] \subseteq \bigcup_{n=1}^i \text{nty}\left(a_n, b_n + \frac{\varepsilon}{2}\right).$$

Then, the Heine-Borel theorem implies there's a finite subcover of  $(a_1, b_1], \dots, (a_N, b_N]$ . This is really analysis (where the rest was definition-chasing), and

$$(b - a - \varepsilon) \leq \sum_1^\infty (b_n - a_n) + \varepsilon,$$

so the result is shown when  $\varepsilon \rightarrow 0$ .

One of the most useful things to have happened in the last fifty years is something called the  $\pi$ - $\lambda$  theorem, which is a bit technical. We'll prove it, but first we'll use it to avoid getting too technical too soon.

**Definition.** A collection of sets  $\mathcal{L}$  in  $\Omega$  is a  $\pi$ -system if it is closed under finite intersections. It is a  $\lambda$ -system if  $\emptyset \in \mathcal{L}$ ,  $\mathcal{L}$  is closed under complements, and countable unions of disjoint sets in  $\mathcal{L}$  are in  $\mathcal{L}$ .

**Example 3.1.** For example, if  $\Omega = \{1, 2, 3, 4\}$ , then  $\mathcal{L} = \{\emptyset, \Omega, \{1, 2\}, \{3, 4\}, \{1, 3\}, \{2, 4\}\}$  is a  $\lambda$ -system but not a field.

**Theorem 3.2 ( $\pi$ - $\lambda$ ).** Let  $\mathcal{P}$  be a  $\pi$ -system on  $\Omega$  and  $\mathcal{L}$  be a  $\lambda$ -system on  $\Omega$ ; then, if  $\mathcal{P} \subseteq \mathcal{L}$ , then  $\sigma(\mathcal{P}) = \mathcal{L}$ .

This theorem has plenty of applications:

- (1) If two probability measures agree on a  $\pi$ -system  $\mathcal{P}$  (e.g. on the set of intervals!), then they agree on  $\sigma(\mathcal{P})$ .

*Proof.* Call these two probability measures  $P$  and  $Q$ , and let  $\mathcal{L} = \{A : P(A) = Q(A)\}$ ; this happens to be a  $\lambda$ -system (which is the bulk of the argument, but devolves to definition-checking), and it contains  $\mathcal{P}$ , so  $\sigma(\mathcal{P}) \subseteq \mathcal{L}$ . □

- (2) As a corollary, we get the uniqueness in Carathéodory's theorem.

---

<sup>4</sup>The professor actually drew the frowny face on the chalkboard.



#### 4. INDEPENDENCE: 10/1/14

*“They’re trying to make me not heard, BUT THEY WILL NOT SUCCEED!”*

Recall that a class  $\mathcal{P}$  is a  $\pi$ -system if whenever  $A, B \in \mathcal{P}$ , then  $A \cap B \in \mathcal{P}$ , and that a class  $\mathcal{L}$  is a  $\lambda$ -system on  $\Omega$  if  $\Omega \in \mathcal{L}$ ,  $\mathcal{L}$  is closed under complements, and  $\mathcal{L}$  is closed under countable disjoint unions.

**Claim.** If  $\mathcal{L}$  is a  $\lambda$ -system and  $A, B \in \mathcal{L}$  are such that  $A \subseteq B$ , then  $B \setminus A = B \cap A^c \in \mathcal{L}$ .

*Proof.* Since  $B \in \mathcal{L}$ , then  $B^c \in \mathcal{L}$  and it’s disjoint from  $A$ , so  $A \cup B^c \in \mathcal{L}$ , and thus so is its complement  $A^c \cap B$ .  $\square$

**Claim.** If  $\mathcal{L}$  is a  $\sigma$ -system and a  $\pi$ -system, then it is a  $\sigma$ -algebra.

*Proof.* We just need to show that non-disjoint unions are in  $\mathcal{L}$ , since the rest follow by definition. For any  $A, B \in \mathcal{L}$ ,  $A \cup B = (A \setminus (A \cap B)) \cup (A \cap B) \cup (B \setminus (A \cap B))$  as a disjoint union, so it’s in  $\mathcal{L}$ ; similarly, one can “disjointify” a countable collection of sets into a disjoint countable collection, so its union is in  $\mathcal{L}$ .  $\square$

Recall also the  $\pi$ - $\lambda$  theorem, Theorem 3.2.

*Proof of Theorem 3.2.* Let  $\mathcal{L}_0$  be the smallest  $\lambda$ -system containing  $\mathcal{P}$ ; we’ll show that  $\mathcal{L}_0$  is also a  $\pi$ -system, so  $\sigma(\mathcal{P}) \subseteq \mathcal{L}_0 \subseteq \mathcal{L}$ .

For any  $A \in \mathcal{P}$ , let  $\mathcal{L}_A = \{B \subseteq \Omega : A \cap B \in \mathcal{L}_0\}$ . Then,  $\mathcal{L}_A$  is a  $\lambda$ -system:  $\Omega \in \mathcal{L}_A$ , certainly, and if  $B \in \mathcal{L}_A$ , then  $A \cap B \in \mathcal{L}_0$ , so  $A \setminus (A \cap B) \in \mathcal{L}_0$  and  $A \cap (A^c \cup B^c) = A \cap B^c \in \mathcal{L}_0$ , so  $B^c \in \mathcal{L}_A$ . Finally, if  $B_1, \dots$  is a disjoint collection in  $\mathcal{L}_A$ , then each  $B_i \cap A \in \mathcal{L}_0$  are still disjoint, so

$$\left(\bigcup_i B_i\right) \cap A = \bigcup_i B_i \cap A \in \mathcal{L}_0.$$

Thus, the union of the  $B_i$  is in  $\mathcal{L}_A$ .

Moreover, if  $A \in \mathcal{P}$ , then  $\mathcal{L}_0 \subseteq \mathcal{L}_A$ , because if  $A, B \in \mathcal{P}$ , then  $\mathcal{P} \subseteq \mathcal{L}_A$ , so  $\mathcal{L}_0 \subseteq \mathcal{L}_A$  (there’s something to show here, but only a single line).

If  $B \in \mathcal{L}_0$ , then  $\mathcal{L}_0 \in \mathcal{L}_B$ , because if  $A \in \mathcal{P}$ , then  $B \in \mathcal{L}_A$ , so  $A \cap B \in \mathcal{L}_0$  and thus  $A \in \mathcal{L}_B$ . Thus,  $\mathcal{P} \subseteq \mathcal{L}_B$ , so  $\mathcal{L}_0 \in \mathcal{L}_B$ .

Finally, if  $B, C \in \mathcal{L}_0$  and  $C \in \mathcal{L}_B$ , then  $\mathcal{L}_0$  is a  $\pi$ -system; thus,  $\mathcal{L}_0$  is a  $\sigma$ -algebra in  $\mathcal{L}$ , and  $\sigma(\mathcal{P}) \subseteq \mathcal{L}_0 \subseteq \mathcal{L}$ .  $\square$

**Independence.** Kolmogorov, when asked to explain the difference between probability and analysis, responded that probability is analysis plus independence; that may have been truer in the 1930s than today, but independence is still important.

**Definition.** Let  $(\Omega, \mathcal{F}, P)$  be a probability space; then,  $A, B \in \mathcal{F}$  are independent if  $P(A \cap B) = P(A)P(B)$ . A collection  $\{A_n\}_{n \in I}$  is *independent* if for any finite subcollection  $A_{n_1}, \dots, A_{n_m}$  we have

$$P\left(\bigcap_{i=1}^m A_{n_i}\right) = \prod_{i=1}^m P(A_{n_i}).$$

**Example 4.1.** Remember  $\Omega = (0, 1]$  and  $w = \sum_{n=1}^{\infty} d_n(\omega)/2^n$ . If  $A_n = \{d_n = 0\}$ , then the collection  $\{A_n\}$  is independent.

Let  $\mathcal{A}_i = \sigma(A_i) = \{\emptyset, \Omega, \{d_i = 0\}, \{d_i = 1\}\}$  (in some sense, all you can know about the  $i^{\text{th}}$  digit). Then, let  $\mathcal{A}_{\text{ODD}} = \sigma(\mathcal{A}_{2n+1}, n \in \mathbb{N})$  and  $\mathcal{A}_{\text{EVEN}} = \sigma(\mathcal{A}_{2n}, n \in \mathbb{N})$ . Intuitively,  $\mathcal{A}_{\text{ODD}}$  and  $\mathcal{A}_{\text{EVEN}}$  are independent, but how do we show this?

**Proposition 4.2.** Let  $(\Omega, \mathcal{F}, P)$  be a probability space and  $\mathcal{A}_n$ ,  $1 \leq n \leq I$  be independent  $\pi$ -systems. Then,  $\sigma(\mathcal{A}_n)$  for  $1 \leq n \leq I$  are independent.

*Proof.* Let  $\mathcal{B}_i = \mathcal{A}_i \cup \Omega$ , which is still a  $\pi$ -system, and the  $\mathcal{B}_i$  are still independent.

Let

$$\mathcal{L}_1 = \left\{B_1 \in \Omega : \text{for all choices } B_2, \dots, B_I \in \Omega, P\left(\bigcup_{j=1}^I B_j\right) = \prod_{j=1}^I P(B_j)\right\}.$$

Then,  $\mathcal{B}_1 \subseteq \mathcal{L}_1$ , and  $\mathcal{L}_1$  is a  $\lambda$ -system (just check all of the axioms). Thus,  $\sigma(\mathcal{B}_1) \subseteq \mathcal{L}_1$ . Now, independence follows.  $\square$

We’re about to do two extremely useful, trivial theorems, and one hard, useless theorem.

**Definition.** The statement that  $A_n$  occurs *infinitely often* denotes the set

$$\bigcap_{n=1}^{\infty} \bigcup_{m=n}^{\infty} A_m,$$

which is the set of  $\omega$  in which  $A_n$  occurs infinitely often (e.g. in a decimal expansion).

**Theorem 4.3 (Borel-Cantelli).** Let  $(\Omega, \mathcal{F}, P)$  be a probability space and  $\{A_n\}_{n=1}^{\infty} \in \mathcal{F}$ . Then,

- (1) If  $\sum_{n=1}^{\infty} P(A_n)$  is finite, then  $P\{A_n \text{ infinitely often}\} = 0$ .
- (2) If the sum is infinite and the  $A_n$  are independent, then  $P\{A_n \text{ infinitely often}\} = 1$ .



*Proof.* For the first part, given an  $\varepsilon > 0$ ,

$$\begin{aligned} P(A_n \text{ infinitely often}) &\leq P\left(\bigcup_{i=N}^{\infty} A_i\right) \\ &\leq \sum_{N}^{\infty} P(A_i) < \varepsilon. \end{aligned}$$

For the second part, we'll use that  $1 - x \leq e^{-x}$  for all  $x > 0$ , and study

$$P((A_i \text{ infinitely often})^c) = P\left(\bigcup_{n=1}^{\infty} \bigcap_{i=n}^{\infty} A_i^c\right).$$

Now,

$$\begin{aligned} P\left(\bigcap_{i=n}^{\infty} A_i^c\right) &\leq P\left(\bigcap_{i=n}^N A_i^c\right) \\ &= \prod_{i=n}^N P(A_i^c) \\ &= \prod_{i=n}^N (1 - P(A_i)) \\ &\leq e^{-\sum_{i=n}^N P(A_i)}. \end{aligned}$$

Letting  $n \rightarrow \infty$ , the right-hand side goes to 0, so the probability of the complement of  $A_i$  happening infinitely often is 0, so  $P(A \text{ infinitely often}) = 1$ .  $\square$

Now for a useless, but beautiful, theorem.

Return to  $\Omega = (0, 1]$ ,  $\omega = \sum d_n/2^n$ , and  $r_n = 2d_n - 1 \in \{\pm 1\}$ . The strong law showed that  $s_n = \sum_{i=1}^n r_i$ , and  $s_n/n \rightarrow 0$ .

Recall the definitions of the  $\underline{\lim}$  and  $\overline{\lim}$ : if  $x_n \in \mathbb{R}$ , then

$$\begin{aligned} \overline{\lim} x_n &= \lim_{n \rightarrow \infty} \sup_{i \geq n} x_i, \text{ and} \\ \underline{\lim} x_n &= \lim_{n \rightarrow \infty} \inf_{i \geq n} x_i. \end{aligned}$$

The sequence  $\{x_n\}$  has a limit iff  $\overline{\lim} x_n = \underline{\lim} x_n$ .

Interestingly,  $\overline{\lim} s_n = +\infty$  and  $\underline{\lim} s_n = -\infty$ , and this is also true for  $s_n/\sqrt{n}$ , and Hardy showed  $\lim s_n/(\sqrt{n} \log n) = 0$ . But we can prove a tighter bound on the fluctuation.

**Theorem 4.4 (Law of the Iterated Logarithm).**

$$\overline{\lim} \frac{s_n}{\sqrt{2n \log \log n}} = 1 \quad \text{and} \quad \underline{\lim} \frac{s_n}{\sqrt{2n \log \log n}} = -1.$$

*Remark.*

- (1) In statistical testing (in the real world), people look at the event  $\{s_n/n > .96\}$ , which is usually taken as evidence for the hypothesis, but the theorem shows that if you wait long enough, this will eventually happen, even for a fair coin!
- (2) Well, but why is it useless? The issue is that it takes a long time:  $\log \log 10^{100} = 5.43$ , so don't sit around waiting for it to happen.

The proof needs tail bounds for the binomial distribution, which come from the homework (and aren't too bad with hints).