

Measure and Probability on General Spaces

Michael Betancourt

July 2023

Table of contents

1	Allocation Over Elements	2
2	Allocation Over All Subsets	4
2.1	Consistent Allocations	4
2.2	Sub-Additivity and Super-Additivity	6
3	σ-Algebras	7
3.1	Filtering Subsets	7
3.2	Generating σ -Algebras	8
4	Measures and Probability Distributions	9
4.1	Formal Definitions	9
4.2	Derived Properties	12
4.3	Null Subsets	15
4.4	Measures and Probability Distributions In Practice	15
5	Uniform Measures	16
5.1	Counting Measures	16
5.2	Lebesgue Measures	18
5.3	Uniformity, Ignorance, and Information	21
6	Interpretations of Measure And Probability	22
6.1	Modeling Physical Distributions	22
6.2	Modeling Populations	23
6.3	Modeling Frequencies	23
6.4	Modeling Uncertainties	24
6.5	Everyone Play Nicely	25
7	Conclusion	26

8 Acknowledgements	27
References	28
License	28

Previously in [Chapter One](#) we introduced measure and probability theory over sets with only a finite number of elements. We saw in [Chapter Two](#), however, that many of the most mathematical spaces we encounter in practical applications, like the integers and the real line, feature not a finite number of elements but rather countably infinite and even uncountably infinite numbers of elements. Unfortunately extending measure and probability theory to more general spaces like these is not always straightforward.

In this chapter we will investigate the difficulties in defining measure and probability theory on general mathematical spaces, with a focus on concepts instead of technical details. We will first discuss why measures allocated to individual elements does not, in general, provide enough information to define a consistent allocation for all subsets. Then we will consider how certain pathological subsets on some spaces can obstruct consistent allocations over the full power set, and how we can systematically remove these obstructions in practice. Finally we will present the most general form of measure and probability theory that can be applied to any mathematical space and then discuss some common applications.

1 Allocation Over Elements

Recall that in [Chapter One](#) we first defined measures and probability distributions as allocations over the individual elements in a finite set. More formally we were able to define a measure as a function that mapped each element to its allocation of the total measure,

$$\begin{aligned}\mu : X &\rightarrow [0, \infty] \\ x &\mapsto \mu(x) \ .\end{aligned}$$

This element-wise allocation then allowed us to define the measure allocated to subsets. In particular the measure allocated to a subset $x \subset X$ was unambiguously determined by summing up the measures allocated to the included elements,

$$\mu(x) = \sum_{x \in x} \mu(x).$$

On finite spaces this construction gives us a *consistent* allocation in the sense that the total measure is always preserved no matter how we might decompose the ambient set into subsets.

Conveniently this construction does extend to spaces with countably infinite numbers of elements, such as the integers. In these spaces every subset contains at most a countably infinite number of elements and sums of measures will always converge to well-defined values.

Element-wise measure allocations on finite and countably infinite spaces are also known as **mass functions**, with element-wise probability allocations also known as **probability mass functions**.

Mass functions are particularly straightforward to visualize when X is not only countable but also ordered, such as the integers or a subset of the integers. In this case we can visualize the element-wise allocations with a sequence of vertical bars stacked next to each other (Figure 1).

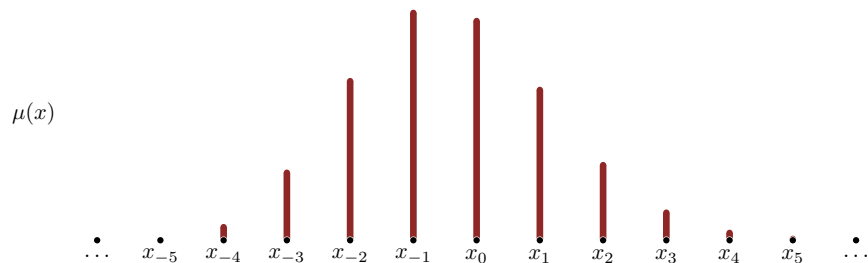


Figure 1: On countable, ordered spaces we can visualize mass functions with a sequence of vertical bars representing the amount of measure allocated to each individual element.

Unfortunately the element-wise construction does not extend any further. Once we consider spaces with uncountably infinite numbers of elements, such as the real numbers, we have to confront subsets with uncountably infinite numbers of elements where sums start to misbehave.

Consider for example a subset x where each of the included elements has been allocated exactly zero measure. If x contains only a finite or countably infinite number of elements then the sum of these zero measures *always* yields zero.

When x contains an uncountably infinite number of elements, however, the sum of the individual element measures is not necessarily zero. In fact it can give *any* value between zero and infinity; uncountably infinite spaces have so many elements that we can very much get something from nothing!

Ultimately this means that on general spaces the allocation of measure to individual elements *does not provide enough information* to uniquely determine what measure should be allocated to every combination of those elements. In order to completely define a measure we need to specify what those subset allocations are ourselves.

2 Allocation Over All Subsets

In [Chapter One](#) we also considered defining a measure by specifying allocations to each subset in the power set,

$$\begin{aligned}\mu : 2^X &\rightarrow [0, \infty] \\ \mathbf{x} &\mapsto \mu(\mathbf{x}) \quad .\end{aligned}$$

Importantly these subsets allocations needed to be consistent with each other to match the behavior of those derived from individual element allocations. For any finite collection of *disjoint* subsets we should have

$$\mu(\cup_{i=1}^I \mathbf{x}_i) = \sum_{i=1}^I \mu(\mathbf{x}_i).$$

For finite spaces this construction is excessive; the subset allocations contain an abundance of redundant information. Because we also can derive subset allocations from element-wise allocations on countably infinite spaces, this construction is unnecessary there as well.

On the other hand at least some subset allocation is strictly necessary for fully defining measures on uncountably infinite spaces, and hence mathematical spaces in general. The only question is whether or not *consistent* subset allocations are even possible on these more sophisticated spaces.

2.1 Consistent Allocations

Before answering this question let's take a second to define exactly what kind of consistency we need. Because finite spaces feature only a finite number of subsets we only ever have to consider the consistency of a finite collection of subsets at a time. More formally if

$$\{\mathbf{x}_1, \dots, \mathbf{x}_i, \dots, \mathbf{x}_I\}$$

is any finite collection of disjoint subsets,

$$\mathbf{x}_i \cap \mathbf{x}_{i' \neq i} = \emptyset,$$

then a consistent measure should give

$$\mu(\cup_{i=1}^I \mathbf{x}_i) = \sum_{i=1}^I \mu(\mathbf{x}_i).$$

Regardless of how many elements the ambient space contains, consistency of a measure over any finite collection of subsets is known as **finite additivity**.

More general spaces can feature infinitely many subsets, and hence different possible notions of additive consistency. For example on a countably infinite space the subset allocations derived from a mass function are consistent across countably infinite collections of subsets. If

$$\{x_1, \dots, x_i, \dots\}$$

is any countably infinite collection of disjoint subsets with with

$$x_i \cap x_{i' \neq i} = \emptyset$$

then

$$\mu(\cup_i x) = \sum_i \mu(x_i).$$

This is known as **countable additivity**.

The question is then whether measures with finite additivity are sufficiently useful for practical application or if we need to consider countably additive measures, let alone measures that might be additive over even larger collections of subsets.

For example a common problem that arises in practice is reconstructing the measure allocated to a general subset from the measures allocated to particularly nice subsets that are easier with which to work. If we could always decompose a generic subset into the disjoint union of a finite number of nice subsets then finite additivity would be sufficient for this task. On the other hand if we could decompose a generic subset into the disjoint union of only a countably infinite number of nice subsets then countable additivity would be sufficient. Potentially some subsets might be decomposable only into an uncountably infinite number of subsets in which case we would need even stronger notions of additivity!

Fortunately for us we don't have to go to that last extreme. It turns out that on most spaces that we'll encounter in practice, and typical notions of "nice" subsets, countable additivity is sufficient for reconstructing the measure allocated to more general subsets.

To demonstrate let's consider the two-dimensional real plane \mathbb{R}^2 and a measure that is partially defined through its allocations to **rectangular** subsets (Figure 2). In general a non-rectangular subset, in this case a disk, can be crudely approximated by a single rectangular subset. The disk can be approximated more precisely as the disjoint union of many different rectangular subsets, but that will never exactly reconstruct the disk. Only when we incorporate a countably infinite number of rectangular subsets can we reconstruct the disk without any error.

Ultimately countable additive measures give us the mathematical flexibility we need for many practical applications.

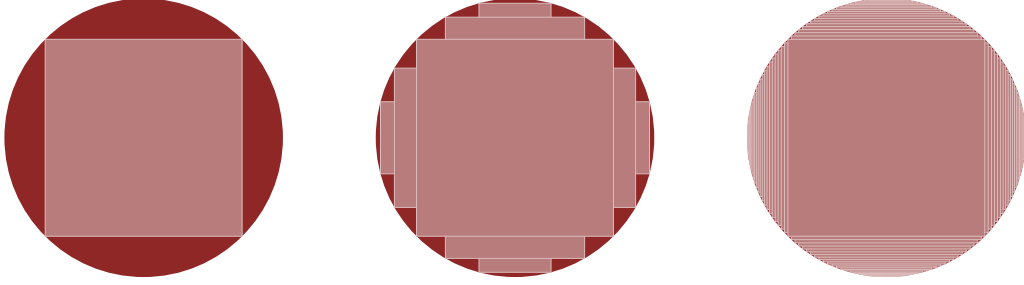


Figure 2: On a two-dimensional real plane \mathbb{R}^2 a non-rectangular disc can be approximated, but not exactly reconstructed, by the finite union of different rectangular subsets. In order to exactly reconstruct a non-rectangular subset we need to include a countably infinite number of rectangular subsets. If measures are countably additive then we can use this decomposition to reconstruct the measure allocated to the disk by adding up the measures allocated to the infinitely many rectangular subsets in the reconstruction.

2.2 Sub-Additivity and Super-Additivity

Ideally we would be able to define measures that are additive over any countably infinite collection of disjoint subsets on any space. Unfortunately mathematics is not always kind, and many seemingly well-behaved space feature pathological subsets that obstruct countable additivity.

Specifically many uncountably infinite spaces feature disjoint subsets that will always behave **sub-additivity**,

$$\begin{aligned} x_1 \cup x_2 &= \emptyset \\ \mu(x_1 \cup x_2) &< \mu(x_1) + \mu(x_2) \end{aligned}$$

no matter how we try to define the allocation! In other words the power set will always be infiltrated by certain subsets that are always less than the sum of their parts and obstruct a consistent definition of measure.

At the same time we can generally show that there exist disjoint subsets that are **super-additive**,

$$\begin{aligned} x_1 \cup x_2 &= \emptyset \\ \mu(x_1 \cup x_2) &> \mu(x_1) + \mu(x_2). \end{aligned}$$

In other words if the measure allocated to these subsets and the measure allocated to their union then we will always appear to end up with measure than what had been initially allocated.

What makes these pathological subsets even more awkward is that we can't actually construct them from explicit conditions. Given typical assumptions about infinity all we can do is prove that these subsets exist. These phantom subsets are known as **non-constructive** objects.

That said because the misbehaving subsets are non-constructive we don't really need to consider them in any practical application of measure theory. If we could consistently filter them out of the full power set then we would be able to define consistent measures over the remaining subsets, and that would be sufficient for any practical application.

3 σ -Algebras

Because the term " σ -algebra" is often thrown around in measure and probability theory without much explanation it is often seen as an impenetrable concept that defies explanation. In reality, however, σ -algebra are simply a way to consistently filter out subsets from the power set.

3.1 Filtering Subsets

We can always filter the power set by removing certain subsets. The difficulty is ensuring that no application of the three set operations would ever lead us back to the excised subsets and reveal a "hole" in the remaining collection of subsets. In other words we need our filtered collection of subsets to be *closed* under the three set operations so that there is no risk of accidentally recreating a subset outside of the collection.

In particular if the subset $x \subset X$ is in our filtered collection then so too should be the complement x^c . If this is true then anytime we apply the complement operator to a subset in our collection we are guaranteed to always see another subset in our collection.

Similarly for every pair of subsets $x_1 \subset X$ and $x_2 \subset X$ in a filtered collection the union $x_1 \cup x_2$ and intersection $x_1 \cap x_2$ should also be in the collection. In order to ensure closure under repeated applications of the union and intersection operators we need the union and intersection of any countably infinite sequences of subsets to also be in the filtered collection.

A **σ -algebra** is any collection of subsets that is closed under complements, countable unions, and countable intersections. In other words a σ -algebra is just any consistent filtering of the power set. I will use a calligraphic font to refer to σ -algebras so that if X is a space then $\mathcal{X} \subset 2^X$ will denote a σ -algebra defined on that space.

A set equipped with a σ -algebra, (X, \mathcal{X}) is known as a **measurable space**. I will refer to X as the **ambient set**, or the **ambient space** if it is also equipped with additional structure. Similarly the elements of a σ -algebra are known as **measurable subsets** while any subsets in the power set but not in the σ -algebra are referred to as **non-measurable** subsets.

When non-measurable subsets are misbehaving subsets they reveals the subtle, and often counterintuitive, pathologies inherent to that space. By working with σ -algebras directly we can avoid these awkward pathologies entirely.

3.2 Generating σ -Algebras

Now that we've defined how a consistent sub-collection of subsets behaves we need to consider how to construct these σ -algebras in practice. One particularly useful way to build up σ -algebras is to *generate* them by repeatedly applying the three set operations to an initial collection of subsets.

For example consider an initial collection of two subsets

$$\{x_1, x_2\}.$$

Applying the complement operator gives us two subsets that fall outside of the initial collection,

$$\{x_1^c, x_2^c\},$$

Similarly applying the union operator gives

$$\{x_1 \cup x_2\}$$

while applying the intersection operator gives

$$\{x_1 \cap x_2\}.$$

To ensure closure we have to add *all* of these subsets to our initial collection,

$$\{x_1, x_2, x_1^c, x_2^c, x_1 \cup x_2, x_1 \cap x_2\}.$$

At this point we iterate, applying the complement operator to every subset and the union and intersection operators to every finite and countably infinite sub-collection of subsets to generate an even larger collection of subsets. When the set operations no longer return new subsets the final collection of subsets defines a σ -algebra.

A convenient feature of this procedure is that if we start with a collection of constructive subsets then we will *always* end up with a σ -algebra that is free of any non-constructive subsets and their pathological behaviors. To ensure that we don't filter out any well-behaved subsets in the process we just have to make sure that our initial collection is sufficiently large.

Conveniently when working on a topological space we already have a natural collection of subsets that we can use to generate a σ -algebra – the defining topology itself! The σ -algebra generated by repeatedly applying all three set operations to the subsets in a topology is known as a **Borel** σ -algebra. In other words a Borel σ -algebra is the unique σ -algebra comprised of all

of the open and closed subsets. If X is a topological space then I will denote the corresponding Borel σ -algebra by $\mathcal{B}(X)$.

Every space that we will consider in this book will be a topological space. Consequently we can always use the corresponding Borel σ -algebra to remove any undesired subsets that might obstruct the definition of a consistent measures and probability distributions. Indeed Borel σ -algebras are so common that they are often take for granted, with any reference to a “measurable space” implicitly assuming a topological space and its corresponding Borel σ -algebras to filter out any inconsistent behavior.

For example finite and countably infinite spaces are almost always equipped with discrete topologies. Because discrete topologies contain all of the atomic sets the σ -algebras derived from them will always be the full power set. In these cases there are no pathological behaviors that we have to avoid at all!

On the the other hand the Borel σ -algebra derived from the topology that defines the real line filters out all of the non-constructive subsets and their undesired behaviors. This results in a σ -algebra that is strictly smaller than the full power set of the real line.

A Borel σ -algebra is sufficient for removing any counterintuitive behavior from a topological space, but in more technical mathematical work there are circumstances where slightly larger σ -algebras may be more convenient; see for example the end of [Section 5.2](#). In more applied practice we can always safely assume a Borel σ -algebra or any extension of that σ -algebra that might be needed to resolve any technical issues.

4 Measures and Probability Distributions

With all of that work we are finally ready to define a theory for allocating any conserved, but not necessarily finite, quantity across a general mathematical space.

4.1 Formal Definitions

A **measure** on any measurable space (X, \mathcal{X}) is a function from the σ -algebra \mathcal{X} to the extended positive real line,

$$\begin{aligned} \mu : \mathcal{X} &\rightarrow [0, \infty] \\ x &\mapsto \mu(x) \end{aligned} ,$$

that is **countably additive**,

$$\mu(\cup_i x_i) = \sum_i \mu(x_i)$$

for any countably infinite collection of subsets

$$\{x_1, \dots, x_i, \dots\}$$

that are mutually disjoint,

$$x_i \cap x_{i' \neq i} = \emptyset.$$

On finite and countably infinite spaces we can always take $\mathcal{X} = 2^X$ and ensure countable additivity by allocating measure to individual elements and then deriving the measure allocated to subsets by summing over the measure allocated to the included elements. When working with more sophisticated ambient spaces, however, the pair $(X, 2^X)$ may not admit *any* consistent measures. In these cases we have to consider smaller σ -algebras in order for measures to exist.

A set equipped with not only a σ -algebra but also a measure, in other words the triple (X, \mathcal{X}, μ) is known as a **measure space**. Again I will refer to X as the ambient set or ambient space as appropriate.

If the total measure is finite, $\mu(X) < \infty$, then μ is referred to as a **finite measure**. In this case we can always normalize the measure by $\mu(X)$ to define a proportional allocation.

A **probability distribution** (Figure 3) on any measurable space (X, \mathcal{X}) is a function from the σ -algebra \mathcal{X} to the closed unit interval,

$$\begin{aligned} \pi : \mathcal{X} &\rightarrow [0, 1] \\ x &\mapsto \pi(x), \end{aligned}$$

with

$$\pi(X) = 1$$

and

$$\pi(\cup_i x_i) = \sum_i \pi(x_i)$$

for any countably infinite collection of subsets

$$\{x_1, \dots, x_i, \dots\}$$

that are mutually disjoint,

$$x_i \cap x_{i' \neq i} = \emptyset.$$

These properties are also known collectively as the **Kolmogorov axioms**.

A set equipped with a σ -algebra and a probability distribution is known as a **probability space**. Sometimes the combination (X, \mathcal{X}, π) is also known as a **probability triple**.

Probability spaces are also sometimes denoted by

$$x \sim \pi,$$

where $x \in X$ indicates the ambient set and a σ -algebra is taken for granted. In words this reads “the variable x is distributed according to π ” or “the variable x follows the distribution π ”.

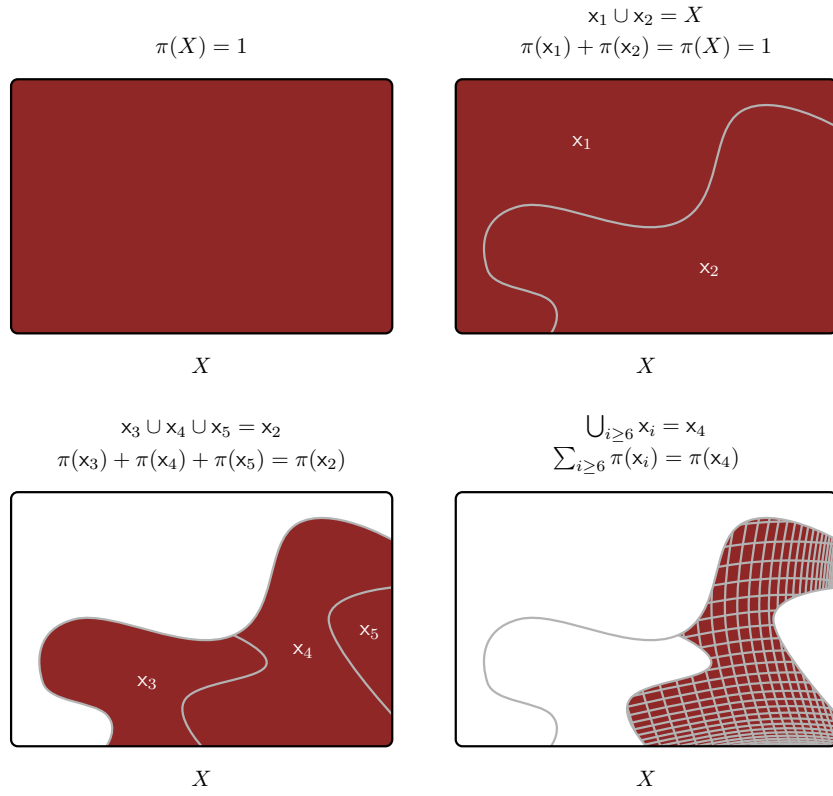


Figure 3: A probability distribution defines a proportional allocation across a measurable space such that no matter how we slice up the ambient set X into measurable subsets the total probability is always preserved.

That said because probability distributions are generally defined over measurable subsets and not individual elements a more precise description would be “the variable x takes values in a set X that is equipped with a probability distribution π ”. The emphasis on variables instead of spaces in this notation is related to the awkward notion of a “random variable” which we will discuss in more detail in Chapter Ten.

4.2 Derived Properties

Although these definitions might appear to be a bit stark, we can derive all of the usual rules of measure and probability theory from them.

For example consider one measurable subset that is strictly smaller than another,

$$x_1 \subset x_2 \in \mathcal{X}.$$

In this case we can always write

$$x_2 = x_1 \cup x_3$$

for the non-empty, measurable subset of elements that are in x_2 but not in x_1 . Applying countable additivity then gives

$$\begin{aligned}\pi(x_2) &= \pi(x_1 \cup x_3) \\ &= \pi(x_1) + \pi(x_3) \\ &\geq \pi(x_1)\end{aligned}$$

because $\pi(x_3) \geq 0$. In other words larger measurable subsets are always allocated more or equal probability than smaller subsets.

Similarly because any subset and its complement are disjoint and combine to reconstruct the full set we always have

$$\begin{aligned}1 &= \pi(X) \\ &= \pi(x \cup x^c) \\ &= \pi(x) + \pi(x^c)\end{aligned}$$

or

$$\pi(x^c) = 1 - \pi(x).$$

In order to work with two measurable subsets $x_1, x_2 \in \mathcal{X}$ that might not be disjoint (Figure 4) we have to consider the elements that are unique to each,

$$x_{1\ 2} = \{x \in X \mid x \in x_1, x \notin x_2\}$$

and

$$x_{2\ 1} = \{x \in X \mid x \in x_2, x \notin x_1\},$$

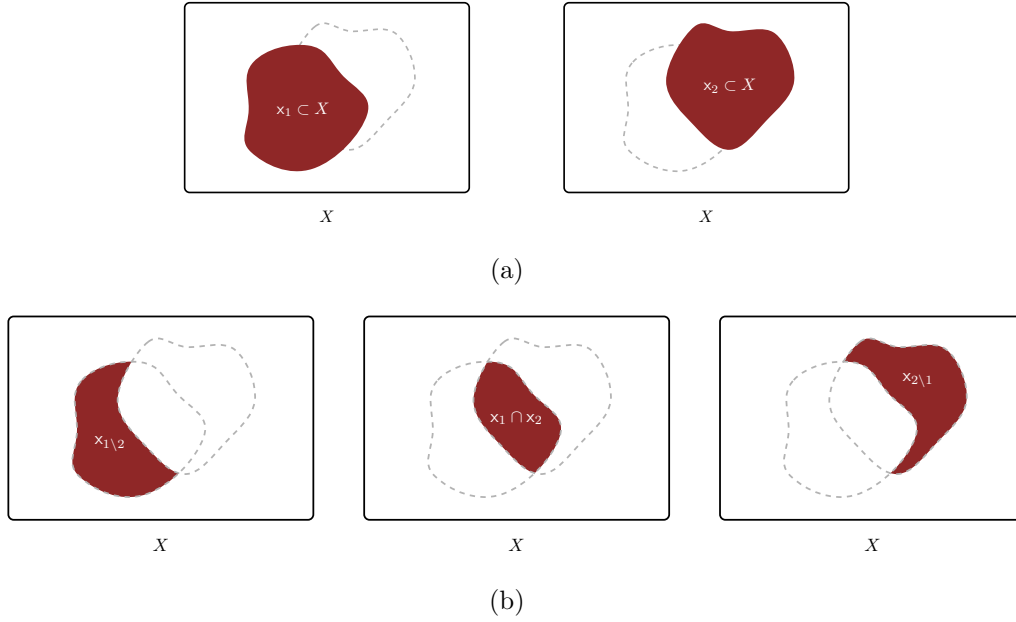


Figure 4: The union of (a) two overlapping subsets x_1 and x_2 can always be decomposed into the union of (b) three disjoint subsets. One disjoint subset $x_1 \setminus x_2$ encapsulates the elements unique to x_1 , another $x_2 \setminus x_1$ encapsulates the elements unique to x_2 , and finally the intersection $x_1 \cap x_2$ encapsulates the elements shared by the two input subsets.

and the elements that are shared,

$$x_1 \cap x_2 = \{x \in X \mid x \in x_1, x \in x_2\}.$$

This then allows us to decompose x_1 , x_2 , and their union into disjoint, measurable subsets (Figure 5,)

$$\begin{aligned} x_1 &= x_{1 \setminus 2} \cup (x_1 \cap x_2) \\ x_2 &= x_{2 \setminus 1} \cup (x_1 \cap x_2) \\ x_1 \cup x_2 &= x_{1 \setminus 2} \cup (x_1 \cap x_2) \cup x_{2 \setminus 1}. \end{aligned}$$

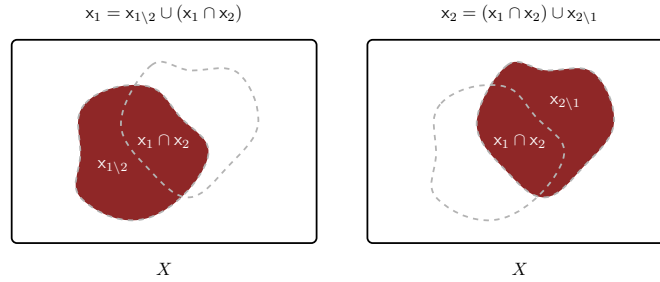


Figure 5: The disjoint subsets introduced in (Figure 4) can also be used to reconstruct the two input subsets individually.

Applying countable additivity to these three decompositions gives a system of equations

$$\begin{aligned} \pi(x_1) &= \pi(x_{1 \setminus 2}) + \pi(x_1 \cap x_2) \\ \pi(x_2) &= \pi(x_{2 \setminus 1}) + \pi(x_1 \cap x_2) \\ \pi(x_1 \cup x_2) &= \pi(x_{1 \setminus 2}) + \pi(x_1 \cap x_2) + \pi(x_{2 \setminus 1}). \end{aligned}$$

Adding the first two together gives

$$\pi(x_1) + \pi(x_2) = \pi(x_{1 \setminus 2}) + 2\pi(x_1 \cap x_2) + \pi(x_{2 \setminus 1})$$

or

$$\pi(x_{1 \setminus 2}) + \pi(x_1 \cap x_2) + \pi(x_{2 \setminus 1}) = \pi(x_1) + \pi(x_2) - \pi(x_1 \cap x_2).$$

Substituting this into the third equation finally gives

$$\pi(x_1 \cup x_2) = \pi(x_1) + \pi(x_2) - \pi(x_1 \cap x_2).$$

4.3 Null Subsets

The measure allocated to a measurable subset quantifies the weight of that subset relative to any other measurable subsets. Those measurable subsets that are allocated *zero* measure are the least important subsets in terms of the overall allocation. At the same time these negligible subsets can be useful for characterizing certain properties of a given measure.

Any measurable subset $x \in \mathcal{X}$ that is allocated zero measure

$$\mu(x) = 0$$

is referred to as a **null subset** of the measure space (X, \mathcal{X}, μ) or, more compactly, a μ -null subset. Similarly if

$$\pi(x) = 0$$

then the measurable subset x is denoted a null subset of the probability space (X, \mathcal{X}, π) , or simply a π -null subset.

Most properties of measures depend on the detailed allocation of the total across all measurable subsets. Some useful properties, however, are completely characterized by which measurable subsets receive a non-zero allocation and which measurable subsets receive a zero allocation.

Any two measures that share the same null subsets will share any properties that are derived from those null subsets. Consequently the overlap in null subsets, or the lack thereof, is often a useful way to determine how compatible two measures are with each other. This compatibility will be particularly important when we construct density functions in Chapter 5.

4.4 Measures and Probability Distributions In Practice

The formal definition of measures and probability distributions tell us what form the consistent allocation of any quantity on any measurable space has to take, but it does not necessarily provide a way for constructing explicit allocations in practice. Specifically in almost all circumstances it is infeasible, if not outright impossible, to exhaustively specify the measure or probability allocated to *every* subset in the ambient σ -algebra. Constructing and then storing infinitely large databases linking each measurable subset, or even just the non-null subsets, to their allocations is not particularly practical!

In some cases we can define useful measures and probability distributions by specifying the allocation to only *some* of the measurable subsets and then deriving the allocations to the rest with countable additivity. For example in finite and countable spaces we need to specify only the allocations to atomic subsets. Similarly measures over $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$ can be completely specified by allocations to interval subsets.

That said in most of these cases those reduced allocations are still impractical to specify one-by-one. Because of that our introduction to general measure and probability theory will have to remain a bit abstract, without many explicit examples, for the time being.

In applied problems measures and probability distributions are almost always defined *algorithmically*, with rules to evaluate the measure or probability allocated to a subset on the fly instead of storing and retrieving the allocation from an exhaustive specification. We will introduce two of these algorithmic representations, and use them to define many useful allocations, in Chapter Five and Chapter Eight.

5 Uniform Measures

Any given measurable space (X, \mathcal{X}) can be equipped with infinitely many measures and probability distributions. Some of these objects, however, are more useful in applied practice than others.

This section will introduce two measures that encode distinct notions of *uniformity* that are applicable to different types of ambient spaces. In the following chapters we will see how the properties of these **uniform measures** make them particularly useful in practical applications of probability theory.

5.1 Counting Measures

Intuitively a uniform measure should allocate the same measure to as many measurable subsets as possible. The consistency of measures, however, limits just how many subsets can receive the same allocation. For example if two disjoint subsets $x_1 \in \mathcal{X}$ and $x_2 \in \mathcal{X}$ are allocated the same measure,

$$\mu(x_1) = \mu(x_2) = \mu_0,$$

then their union will be allocated the measure,

$$\mu(x_1) \cup \mu(x_2) = \mu(x_1) + \mu(x_2) = 2\mu_0;$$

all three allocations will be equal only if $\mu_0 = \infty$.

In order to define a notion of uniformity that doesn't abuse infinity we need to restrict our consideration from the entire σ -algebra to collections of measurable subsets that we want to behave in the same way. If this collection is large enough then we can completely define a uniform measure by enforcing the same allocation to those distinguished subsets.

For example on any countable measure space $(X, 2^X)$ a natural collection of subsets to consider are the atomic subsets. Because in this case the measure allocated to these subsets completely specifies the allocations to every other subset we can fully define a uniform measure by enforcing the same allocation to each element. This results in a constant mass function (Figure 6).

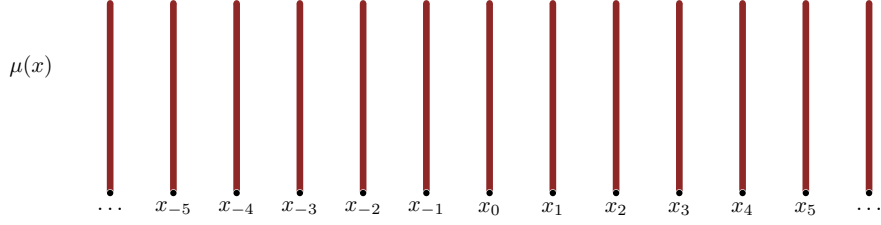


Figure 6: On a countable space any uniform measure allocates the same measure to every individual element resulting in a constant mass function.

In particular the **counting measure** on $(X, 2^X)$ is defined by a unit allocation to each atomic set,

$$\chi(\{x\}) = 1.$$

Equivalently we can define the counting measure with a uniform mass function that assigns each element to a unit allocation,

$$\chi(x) = 1.$$

Given these element-wise allocations we can derive the measure allocated to any other subset $x \subset X$ by countable additivity,

$$\begin{aligned} \chi(x) &= \sum_{x \in x} \chi(\{x\}) \\ &= \sum_{x \in x} 1, \end{aligned}$$

which always results in the total number of elements in x . The total counting measure,

$$\chi(X) = \sum_{x \in X} 1,$$

just counts the number of elements in the ambient set. In other words a counting measure formalizes our intuitive notion of counting discrete objects.

An immediate consequence of these derived allocations is that *any* subset with the same number of elements will receive the same allocation. Uniformity over the individual elements induces uniformity over other subsets as well.

When there are only a finite number of elements in X the total measure will also be finite. In this case we can normalize the counting measure into a uniform probability distribution,

$$\pi(x) = \frac{\sum_{x \in x} 1}{\sum_{x \in X} 1},$$

which quantifies the proportion of all of the elements in X that are contained in x . If there are an infinite number of elements in X then this normalization is no longer possible; for example there is no well-defined notion of a uniform probability distribution over the integers \mathbb{Z} .

Counting measures are not the only uniform measures that we could define over a countable ambient set. More generally we can define a uniform measure by allocating any positive real number $c \in \mathbb{R}^+$ to each element,

$$\kappa(\{x\}) = c.$$

That said these other uniform measures are somewhat redundant in the sense that their allocations can be recovered by scaling the corresponding counting measure allocations,

$$\kappa(x) = c \cdot \chi(x).$$

One important feature of a counting measure, indeed any uniform measure, is that every subset except the empty set receives a non-zero allocation. That is to say that the empty set is the *only* χ -null subset.

5.2 Lebesgue Measures

Unfortunately on uncountable spaces element-wise allocations do not completely define measures. In order to generalize any notion of uniform measure we have to specify a larger class of measurable subsets that should receive equal allocations.

An uncountable set alone, however, offers no criteria for preferring any collection of subsets to any other, and hence no criteria for defining a consistent notion of uniform measure. Additional structure on X , however, may be able to break this ambiguity.

Consider for example a real line \mathbb{R} equipped with an appropriate ordering, algebra, metric, and topology as discussed in [Chapter Two](#). Using the ordering we can construct closed interval subsets,

$$[x_1, x_2] = \{x \in \mathbb{R} \mid x_1 \leq x \leq x_2\}.$$

We can then use the metric to characterize these intervals by the distance between the end points,

$$l([x_1, x_2]) = d(x_1, x_2) = |x_2 - x_1|,$$

otherwise known as the interval **length**. Moreover if we use a Borel σ -algebra derived from the real topology then these closed intervals will all be measurable subsets.

Any notion of uniformity that is compatible with all of this additional structure should treat all interval subsets with the *same length* in the same way. In other words a uniform measure over \mathbb{R} should allocate the same measure to all equal-length intervals (Figure 7). For example because the intervals

$$l([-2, -1]) = l([5, 6]) = l([150, 151])$$

all have the same length any uniform measure should give

$$\mu([-2, -1]) = \mu([5, 6]) = \mu([150, 151]).$$

Likewise because

$$l([-350, -300]) = l([0, 50])$$

we should have

$$\mu([-350, -300]) = \mu([0, 50]),$$

and so on.

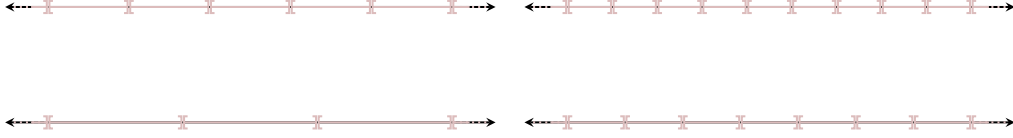


Figure 7: The structure that defines a real line also distinguishes interval subsets that we can characterize by their length. Any uniform measure over a real line should allocate the same measure to any interval with the same length.

The easiest way to accomplish this uniformity is to allocate to each interval a measure directly equal to the its length,

$$\lambda([x_1, x_2]) = l([x_1, x_2]) = |x_2 - x_1|.$$

Allocations to more general measurable subsets can then be derived from these interval allocations. The resulting uniform measure over a given real line is known as the **Lebesgue measure**.

Just as the counting measure formalizes intuition notions of counting on countable spaces, the Lebesgue measure formalizes intuitive notions of length on a real line. This formalization of length, however, is a bit more subtle. Counting behaves the same on all countable spaces, but length can behave differently across different real lines!

Two real lines with incompatible metrics will assign different lengths to the same intervals, resulting in different Lebesgue measures. Equivalently the Lebesgue measure can behave differently for two different parameterizations of a flexible real line. When there might be any chance of confusion we have to be careful to communicate which parameterization we're using in any given application.

Because the total Lebesgue measure is infinitely large,

$$\begin{aligned} \lambda(\mathbb{R}) &= \lim_{x \rightarrow \infty} \lambda([-x, x]) \\ &= \lim_{x \rightarrow \infty} 2|x| \\ &= \infty, \end{aligned}$$

it cannot be normalized into a probability distribution. As with the integers there is no well-defined notion of a uniform probability distribution over a real line.

Every other uniform measure over a real line is defined by allocating a measure to each interval *proportional* to its length,

$$\nu([a, b]) \propto l([a, b]).$$

Consequently every uniform measure over a real line reduces to a constant scaling of the Lebesgue measure,

$$\nu(x) \propto \lambda(x),$$

similar to how every uniform measure over countable spaces reduces to a scaling of the counting measure.

By definition the Lebesgue measure on any real line will allocate zero measure to individual points,

$$\begin{aligned}\lambda(\{x\}) &= \lambda([x, x]) \\ &= d(x, x) \\ &= 0.\end{aligned}$$

Indeed any measurable subset with only a countable number of elements will also be λ -null,

$$\begin{aligned}\lambda(x) &= \lambda(\cup_i \{x_i\}) \\ &= \sum_i \lambda(\{x_i\}) \\ &= \sum_i 0 \\ &= 0.\end{aligned}$$

Critically these properties follow from $d(x, x) = 0$ which is true for *any* well-behaved metric. Two real lines with incompatible metrics might feature different Lebesgue measures, but those Lebesgue measures will always share the same null subsets and hence share any properties derived from those null subsets.

Finally there is one formal detail that I want to mention to prepare any readers who may continue on to more technical literature. Allocations based on interval lengths can be used to derive consistent allocations over any subset constructed from open and closed subsets, and hence every measurable subset in the Borel σ -algebra $\mathcal{B}(\mathbb{R})$. They can also be used, however, to derive consistent allocations to some subsets that are *not* in $\mathcal{B}(\mathbb{R})$.

Accordingly Lebesgue measures are usually formally defined not over $\mathcal{B}(\mathbb{R})$ but rather over a slightly larger σ -algebra referred to as, unsurprisingly, the **Lebesgue σ -algebra**. That said these additional subsets are negligible in the sense that they are always allocated zero Lebesgue measure. Because the differences between these σ -algebras are limited to null subsets they can be effectively ignored in most practical applications.

5.3 Uniformity, Ignorance, and Information

The concepts of **ignorance** and **information** are related to uniformity; formalizing the relationships between these concepts, however, is subtle. In order to avoid confusing these concepts we have take care to recognize not only their similarities but also their differences.

When two elements of a countable space are allocated the same measure then the overall allocation will be the same even if we permute those two elements before allocating measures. In other words any measure that allocates the same measure to two elements is not able to distinguish between any permutations of those elements.

The more regular the individual allocations are the less sensitive the resulting measure will be to any rearrangement of the elements. Conversely the more the allocations vary from element to element the more the resulting measure will be able to discern one permutation from another. Informally we might say that the more uniform the measure the less information it encodes. Because a uniform measure allocates the same measure to every element it is ignorant to *any* bijective transformation of the elements, capturing the least information possible on a countable space.

On uncountable spaces these concepts become more delicate. The allocations defined by the Lebesgue measure, for example, are *not* invariant to arbitrary transformations of the real line. Any transformation that warps the metric will also warp lengths and hence the measures allocated by the Lebesgue measure. Instead the Lebesgue measure is ignorant only to transformations that preserve distances.

In order to formalize heuristic concepts like “ignorance” and “information” we have to embrace a bit more abstraction. Recall that in [Chapter Two](#) we discussed the notion of a structure-preserving transformation. More generally if $\phi : X \rightarrow X$ is a structure-preserving automorphism then we say that the structure is **symmetric** to ϕ while the transformation ϕ is a **symmetry** of the structure.

In other words if ϕ is a symmetry of a structure \mathfrak{x} then the behavior of \mathfrak{x} is the same before and after we apply the transformation. The structure cannot detect whether or not we apply the transformation.

Some structures admit multiple symmetries at the same time. For example the discrete topology on a countable set is invariant to any permutation of the elements while the metric on a real line is invariant to all translations of the elements. The more symmetric a structure is the less it can distinguish between arbitrary transformations to the ambient set. If we formalize information as the ability of a structure to distinguish between transformations of the ambient set then the more symmetric a structure is the less information it encodes.

From an abstract perspective measures and probability distributions are, like orderings, algebras, metrics, topologies, and σ -algebras, just structures that we can endow onto a set. The more invariant a measure is to transformations of that set the less information it will contain.

We will more formally consider how measures transform, and hence how to precisely define symmetries of a measure, in Chapter Six.

Uniform measures are built to be symmetric to some at least some transformations and hence encode less information than other measures. For example the counting measure is invariant to any permutation of a countable ambient set while the Lebesgue measure is invariant to any translation of a real line. We can also extend this construction to more elaborate spaces, for instance defining uniform measures on spheres that invariant to any rotation.

Critically, however, not every uniform measure is invariant to every possible transformation of the ambient set. Some uniform measures are more informative than others! Consequently notions of uniformity do not define any *universal* notions of ignorance, just ignorance to the particular transformations that are used to define uniformity in a given context.

In practice this means that we have to be careful not to make broad claims about uniform measures that are least informative or most ignorant but instead specify *with respect to which transformations* a measure might be least informative or most ignorant.

6 Interpretations of Measure And Probability

To this point our treatment of measure and probability theory has been purely mathematical. A measure defines the allocations of some abstract conserved quantity across some abstract measurable space; a probability distribution defines a proportional allocations. This mathematical construction cannot be endowed with any particular interpretation until we use it to model something.

In this section we'll review some of the most common applications of measure and probability theory and the particular interpretations those applications create.

6.1 Modeling Physical Distributions

One immediate application of measure theory is to model the behavior of a physical quantity, such as mass or electric charge. For example physical mass can be distributed across a solid object in a variety of different ways, with the exact distribution affecting how that object interacts with the surrounding environment. Similarly the distribution of charge across the surface of a conducting object defines its electrostatic properties.

In some physical systems the distribution can also change with time and influence the dynamics of the system. Time-dependent measures that quantify how the distribution of a physical quantity evolves are a common feature of many mathematical physical theories.

6.2 Modeling Populations

A similar application is modeling the selection of individuals, or the properties of individuals, from a larger population. Each time we *sample* a subset of individuals from the population we will observe a different ensemble of behaviors, such as political or consumer preferences, heights, or ages. The heterogeneity of these characteristics across the population can often be quantified with measures, and their relative occurrences modeled with probability distributions.

For example is 30% of the individuals in a population have a height between 0 feet and 5 feet then a probability distribution modeling the variation in heights would give

$$\pi([0, 5]) = 0.3.$$

6.3 Modeling Frequencies

An application particular to probability theory concerns the frequencies of repeated events.

Consider an abstract event whose outcomes take *unpredictable* values in some space X . Perfectly replicating the circumstances of this event N times defines a sequence of values in Y ,

$$\{x_1, \dots, x_n, \dots x_N\}.$$

While we cannot predict what values the individual events in this sequence will take, we may be able to characterize how often certain outcomes appear relative to others. In particular we can define the **frequency** of a subset $x \subset X$ by the number of events that take values in x ,

$$f_N(x) = \frac{\sum_{n=1}^N \mathbb{I}_x(x_n)}{N},$$

where

$$\mathbb{I}_x(x) = \begin{cases} 1, & x \in x \\ 0, & x \notin x \end{cases}.$$

Replicating the event a countably infinite number of times defines the **asymptotic** or **long-run** frequency of a subset,

$$\begin{aligned} f(x) &= \lim_{N \rightarrow \infty} f_N(x) \\ &= \lim_{N \rightarrow \infty} \frac{\sum_{n=1}^N \mathbb{I}_x(x_n)}{N}. \end{aligned}$$

In other words the more frequent subsets contain more common event outcomes.

If the frequencies are the same for *any* sequence of events then we can model them with probability theory. Specifically we can interpret the allocated probabilities as the proportion of the total event outcomes that fall into each subset of outcome values.

In this case the particular ordering of the event sequences doesn't matter and we can also interpret them as defining a population of possible events. From this perspective the application of probability theory is equivalent to the application in the previous section. At the same time if we interpret repeated samples from a population as events then the population probabilities can be interpreted as frequencies.

6.4 Modeling Uncertainties

Probability theory can also be used to consistently quantify uncertain information.

Consider a space of possible statements X . Under perfect knowledge we would be able to specify a particular statement $x \in X$ as true with all other statements in X being false. In other words certainty is quantified with binary true/false assignments. When our knowledge is not quite so certain, however, we have to soften those claims.

To quantify uncertain information we have to generalize beyond binary true/false assignments to continuous values that *interpolate* between absolute truth and falsity. The larger the value we assign to a subset of statements the more our uncertain information supports one of those statements being true. Conversely the smaller the value we assign to a subset the more our uncertain information supports all of the included statements being false.

Applying probability theory allows us to enforce consistent uncertainty assignments across all of the possible statements. The individual probability allocations can then be interpreted as quantifying how strongly our information supports that one of the statements within a measurable set is true. In this setting the allocated probabilities are sometimes referred to as “plausibilities”, “credibilities”, and “beliefs”.

For example the property that $\pi(X) = 1$ corresponds to the fact that at least one of the statements in X always has to be true. A probability distribution that concentrates around the statement x encodes confidence that one of the statements near x is true. The singular limit where all of the available probability collapses onto a single statement, $\pi(\{x\}) = 1$, communicates certainty that x is true.

This kind of probabilistic uncertainty quantification can be interpreted in many ways. For example we can use it to model the personal, subjective beliefs that an individual holds about the behavior of a system. In particular we can use it to model our own specific beliefs. At the same time we can use it to model the collective understanding of entire communities. We can also use probability theory to model only certain aspects of individual or community knowledge and not attempt to quantify the entirety of that knowledge at once.

More formally this application defines one way to generalize classical propositional logic to a many-valued logic. Using probability theory to generalize other logical systems can sometimes also be possible, although the technical details quickly become more complicated.

6.5 Everyone Play Nicely

A key point of confusion in probability theory is the confounding of its abstract mathematical structure with the interpretations that arise in particular applications. This confusion is made all the worse by the long history of attempts to *derive* probability theory from these particular applications.

For example many have tried to derive probabilities as asymptotic frequencies of physical events. The key motivations of this approach is that the resulting probabilities would be objective in the sense that everyone who could implement those infinite trials would attain the same probabilities. Even if we ignore the impracticality of perfectly repeating an event an infinite number of times within a finite lifetime it turns out that there are also some subtle mathematical complications with this approach. For a comprehensive discussion see Diaconis and Skyrms (2017).

Similarly many have tried to derive probability theory from uncertainty quantification. For example the Cox postulates (Van Horn, Kevin S. (2003)) define basic intuitions about uncertainty quantification. On simpler spaces these rules are equivalent to probability theory, but that equivalence doesn't persist to more general spaces. Because of that this approach is not able to recover the full generality of probability theory.

A common reaction to these technical difficulties is to resort to a sort of philosophical bait and switch. When one cannot derive probability theory from a particular application one might define probability theory abstractly, as we have done above, but then impose an *arbitrary* restriction that it can only ever be applied to that one application. For example those trying to derive probability theory from frequencies might argue that probability theory can be applied to model only frequencies, in which case all probabilities are frequencies. Others trying to derive probability theory from the Cox axioms might argue that any application of probability theory always models uncertain information.

These interpretational restrictions then force some awkward philosophical contortions when trying to apply probability theory in practice. For example after imposing that all probabilities are frequencies the only way to model uncertainty in the value of some quantity is to treat it as the outcome of some hypothetical, and completely non-existent, event. The introduction of these hypothetical events to real events makes the entire system more difficult to understand.

In this book we will avoid these restrictions, respect the full generality of probability theory, and take advantage of any consistent applications that might be useful in a given problem. Indeed we will often take advantage of multiple applications *at the same time*.

Consider, for example, a binary space $X = 0, 1$ that corresponds to the two sides of a coin. In particular let 0 denote tails and 1 denote heads. Any probability distribution over X can be quantified with the probability $p \in [0, 1]$ allocated to the point 1, which gives the consistent

probability allocations

$$\begin{aligned}\pi(\emptyset; p) &= 0 \\ \pi(0; p) &= 1 - p \\ \pi(1; p) &= p \\ \pi(X; p) &= 1.\end{aligned}$$

There are many ways to flip a coin, but let's say that we flip our coin in a way that results in an unpredictable sequence of heads and tails. The asymptotic frequencies of these outcomes can then be modeled with an application of probability theory. In other words we can use a probability distribution $\pi(\cdot; p)$ to model the physical outcomes of the flips.

At the same time we use probability theory to model any uncertainty in which of the possible frequency models best matches the true behavior of the coin. In particular we can construct a probability distribution over the unit interval to quantify how compatible each probability allocation $p \in [0, 1]$ is with our knowledge of the coin.

If we have a bag of I coins then could also model the variation in the probability parameters

$$\{p_1, \dots, p_i, \dots, p_I\}$$

for each coin. In this case we can apply probability theory once again, this time to model the population of coin behaviors.

To be clear the interpretations inherent to particular applications of probability theory are important for ensuring that we implement those applications correctly in practice. Elevating one interpretation to the exclusion of others, however, excludes the corresponding applications and limits the full potential of probability theory. To take full advantage of the practical utility of probability theory we have to respect all of consistent applications!

7 Conclusion

Conceptually measure and probability theory are straightforward. Measure theory quantifies how we can consistently allocate a conserved quantify across a general mathematical space and probability theory considers the special case of proportional allocations. In order to quantify that conceptual simplicity, however, we need to resort to some careful mathematics. In particular we need to incorporate σ -algebras to surgically remove any pathological behavior that can arise, even on seemingly well-behaved spaces such as the real line, and obstruct consistent allocations.

Once we've safely constructed these theories in full generality we can use the apply abstract mathematics to model particular systems. Within these applications the math inherits particular interpretations, but we have to be careful to not take these circumstantial interpretations too seriously lest we abandon the full utility of the abstract mathematics.

The technical exploration of measures and probability distributions goes far beyond the introduction in this chapter. Unfortunately many textbooks that cover this material can be difficult to parse without extensive mathematical experience. My personal favorite is Folland (1999) which, while technically rigorous, provides more exposition and motivation than I have found in other treatments.

8 Acknowledgements

I thank Jeff Helzner and Simon Duane for helpful discussion.

A very special thanks to everyone supporting me on Patreon: Adam Fleischhacker, Adriano Yoshino, Alan Chang, Alessandro Varacca, Alexander Bartik, Alexander Noll, Alexander Petrov, Alexander Rosteck, Anders Valind, Andrea Serafino, Andrew Mascioli, Andrew Rouillard, Andrew Vigotsky, Angie_Hyunji Moon, Ara Winter, Austin Rochford, Austin Rochford, Avraham Adler, Ben Matthews, Ben Swallow, Benjamin Glemain, Bradley Kolb, Brandon Liu, Brynjolfur Gauti Jónsson, Cameron Smith, Canaan Breiss, Cat Shark, Charles Naylor, Chase Dwelle, Chris Jones, Chris Zawora, Christopher Mehrvarzi, Colin Carroll, Colin McAuliffe, Damien Mannion, Damon Bayer, dan mackinlay, Dan Muck, Dan W Joyce, Dan Waxman, Dan Weitzenfeld, Daniel Edward Marthaler, Darshan Pandit, Darthmaluus, David Burdelski, David Galley, David Wurtz, Doug Rivers, Dr. Jobo, Dr. Omri Har Shemesh, Ed Cashin, Edgar Merkle, Eric LaMotte, Erik Banek, Ero Carrera, Eugene O’Friel, Felipe González, Fergus Chadwick, Finn Lindgren, Florian Wellmann, Francesco Corona, Geoff Rollins, Granville Matheson, Greg Sutcliffe, Guido Biele, Hamed Bastan-Hagh, Haonan Zhu, Hector Munoz, Henri Wallen, hs, Hugo Botha, Håkan Johansson, Ian Costley, Ian Koller, idontgetoutmuch, Ignacio Vera, Iliaria Prosdocimi, Isaac Vock, J, J Michael Burgess, Jair Andrade, James Hodgson, James McInerney, James Wade, Janek Berger, Jason Martin, Jason Pekos, Jason Wong, Jeff Burnett, Jeff Dotson, Jeff Helzner, Jeffrey Erlich, Jesse Wolfhagen, Jessica Graves, Joe Wagner, John Flournoy, Jonathan H. Morgan, Jonathon Vallejo, Joran Jongerling, Joseph Despres, Josh Weinstock, Joshua Duncan, Joshua Griffith, JU, Justin Bois, Karim Naguib, Karim Osman, Kejia Shi, Kevin Foley, Kristian Gårdhus Wichmann, Kádár András, Lars Barquist, lizzie, LOU ODETTE, Marc Dotson, Marcel Lüthi, Marek Kwiatkowski, Mark Donoghoe, Markus P., Martin Modrák, Matt Moores, Matthew, Matthew Kay, Matthieu LEROY, Maurits van der Meer, Merlin Noel Heidemanns, Michael DeWitt, Michael Dillon, Michael Lerner, Mick Cooney, Márton Vaitkus, N Sanders, Name, Nathaniel Burbank, Nic Fishman, Nicholas Clark, Nicholas Cowie, Nick S, Nicolas Frisby, Octavio Medina, Ole Rogeberg, Oliver Crook, Olivier Ma, Pablo León Villagrà, Patrick Kelley, Patrick Boehnke, Pau Pereira Batlle, Peter Smits, Pieter van den Berg, ptr, Putra Manggala, Ramiro Barrantes Reynolds, Ravin Kumar, Raúl Peralta Lozada, Riccardo Fusaroli, Richard Nerland, Robert Frost, Robert Goldman, Robert kohn, Robert Mitchell V, Robin Taylor, Ross McCullough, Ryan Grossman, Rémi, S Hong, Scott Block, Sean Pinkney, Sean Wilson, Seth Axen, shira, Simon Duane, Simon Lilburn, sssz, Stan_user, Stefan, Stephanie Fitzgerald, Stephen Lienhard, Steve Bertolani, Stew Watts, Stone Chen, Susan Holmes, Svilup, Sören Berg, Tao Ye, Tate Tunstall, Tatsuo Okubo, Teresa

Ortiz, Thomas Lees, Thomas Vladeck, Tiago Cabaço, Tim Radtke, Tobbychev , Tom McEwen, Tony Wuersch, Utku Turk, Virginia Fisher, Vitaly Druker, Vladimir Markov, Wil Yegelwel, Will Farr, Will Tudor-Evans, woejozney, yolhaj , yureq , Zach A, Zad Rafi, and Zhengchen Cai.

References

- Diaconis, Persi, and Brian Skyrms. 2017. *Ten Great Ideas about Chance*. Princeton University Press.
- Folland, G. B. 1999. *Real Analysis: Modern Techniques and Their Applications*. New York: John Wiley; Sons, Inc.
- Van Horn, Kevin S. 2003. “Constructing a logic of plausible inference: a guide to Cox’s theorem.” *International Journal of Approximate Reasoning* 34 (1): 3–24.

License

A repository containing all of the files used to generate this chapter is available on [GitHub](#).

The text and figures in this chapter are copyrighted by Michael Betancourt and licensed under the CC BY-NC 4.0 license:

<https://creativecommons.org/licenses/by-nc/4.0/>