

MAT 314 LECTURE NOTES

1. ANALYSIS ON METRIC SPACES

1.1. Definitions, and open sets. A metric space is, essentially, a set of points together with a rule for saying how far apart two such points are:

Definition 1.1. A *metric space* consists of a set X together with a function $d: X \times X \rightarrow \mathbb{R}$ such that:

- (1) For each $x, y \in X$, $d(x, y) \geq 0$, and $d(x, y) = 0$ if and only if $x = y$.
- (2) For each $x, y \in X$, $d(x, y) = d(y, x)$.
- (3) For each $x, y, z \in X$, $d(x, z) \leq d(x, y) + d(y, z)$.

The last condition is known as the *triangle inequality*. If we think of $d(x, y)$ as representing, say, the smallest possible amount of time that it takes to get from x to y , the triangle inequality should make sense, since one way of getting from x to z is to first go from x to y and then from y to z , and doing this can take as little time as $d(x, y) + d(y, z)$. Of course, it will often be the case that the quickest path from x to z doesn't pass through y , in which case the inequality will be strict.

The first part of the course is concerned with developing notions like convergence and continuity in the context of a general metric space (X, d) . You've probably already had some exposure to these concepts at least in the context of the real numbers, which indeed are the first example:

Example 1.2. Take $X = \mathbb{R}$, and define $d(x, y) = |x - y|$. It should be fairly obvious that the three axioms for a metric space are satisfied in this case.

If you've had a good real analysis course, then a lot (though not all) of the proofs below should look somewhat familiar, essentially with absolute value signs replaced by 'd's.

Example 1.3. If $X = \mathbb{R}^n$, there are actually many metrics d that we can use which generalize the absolute value metric on \mathbb{R} . The most famous of these is surely

$$d_2((x_1, \dots, x_n), (y_1, \dots, y_n)) = \left(\sum_{i=1}^n (x_i - y_i)^2 \right)^{1/2},$$

which gives the distance from (x_1, \dots, x_n) to (y_1, \dots, y_n) that is dictated by the Pythagorean theorem.

Although (in view of the Pythagorean theorem) d_2 seems like a very natural sort of distance, the triangle inequality for it isn't completely obvious, but rather depends on the *Cauchy-Schwarz inequality*, which states that, where for $\vec{v} = (v_1, \dots, v_n), \vec{w} = (w_1, \dots, w_n) \in \mathbb{R}^n$ we write

$$\langle \vec{v}, \vec{w} \rangle = \sum_{i=1}^n v_i w_i,$$

one has

$$\langle \vec{v}, \vec{w} \rangle^2 \leq \langle \vec{v}, \vec{v} \rangle \langle \vec{w}, \vec{w} \rangle.$$

Proof of the triangle inequality for d_2 , assuming the Cauchy-Schwarz inequality. We have

$$\begin{aligned} d_2(\vec{x}, \vec{z})^2 &= \langle \vec{x} - \vec{z}, \vec{x} - \vec{z} \rangle = \langle (\vec{x} - \vec{y}) + (\vec{y} - \vec{z}), (\vec{x} - \vec{y}) + (\vec{y} - \vec{z}) \rangle \\ &= \langle \vec{x} - \vec{y}, \vec{x} - \vec{y} \rangle + 2\langle \vec{x} - \vec{y}, \vec{y} - \vec{z} \rangle + \langle \vec{y} - \vec{z}, \vec{y} - \vec{z} \rangle \\ &\leq \langle \vec{x} - \vec{y}, \vec{x} - \vec{y} \rangle + 2(\langle \vec{x} - \vec{y}, \vec{x} - \vec{y} \rangle \langle \vec{y} - \vec{z}, \vec{y} - \vec{z} \rangle)^{1/2} + \langle \vec{y} - \vec{z}, \vec{y} - \vec{z} \rangle \\ &= d_2(\vec{x}, \vec{y})^2 + 2d_2(\vec{x}, \vec{y})d_2(\vec{y}, \vec{z}) + d_2(\vec{y}, \vec{z})^2 = (d_2(\vec{x}, \vec{y}) + d_2(\vec{y}, \vec{z}))^2, \end{aligned}$$

and then taking the square root of both sides implies the triangle inequality. \square

Proof of the Cauchy-Schwarz inequality. If either \vec{v} or \vec{w} is the zero vector, then both sides of the Cauchy-Schwarz inequality ($\langle \vec{v}, \vec{w} \rangle^2 \leq \langle \vec{v}, \vec{v} \rangle \langle \vec{w}, \vec{w} \rangle$) are zero, so it holds in that case. So let us assume that \vec{v} and \vec{w} are both nonzero, so we can form the vectors

$$\vec{v}' = \frac{1}{\langle \vec{v}, \vec{v} \rangle^{1/2}} \vec{v}, \quad \vec{w}' = \frac{1}{\langle \vec{w}, \vec{w} \rangle^{1/2}} \vec{w},$$

which satisfy

$$\langle \vec{v}', \vec{v}' \rangle = \langle \vec{w}', \vec{w}' \rangle = 1.$$

Notice that

$$0 \leq \langle \vec{v}' \pm \vec{w}', \vec{v}' \pm \vec{w}' \rangle = \langle \vec{v}', \vec{v}' \rangle \pm 2\langle \vec{v}', \vec{w}' \rangle + \langle \vec{w}', \vec{w}' \rangle = 2 \pm 2\langle \vec{v}', \vec{w}' \rangle,$$

i.e.,

$$\pm \langle \vec{v}', \vec{w}' \rangle \leq 1 = \langle \vec{v}', \vec{v}' \rangle \langle \vec{w}', \vec{w}' \rangle.$$

Now square this last equation and then multiply it by $\langle \vec{v}, \vec{v} \rangle \langle \vec{w}, \vec{w} \rangle$ to get the desired inequality

$$\langle \vec{v}, \vec{w} \rangle^2 \leq \langle \vec{v}, \vec{v} \rangle \langle \vec{w}, \vec{w} \rangle.$$

\square

Example 1.4. Here are two examples of metrics on \mathbb{R}^n for which the triangle inequality is more obvious than for the standard Pythagorean distance d_2 :

$$d_1(\vec{x}, \vec{y}) = \sum_{i=1}^n |x_i - y_i|,$$

$$d_\infty(\vec{x}, \vec{y}) = \max_{1 \leq i \leq n} |x_i - y_i|.$$

The first of these, d_1 , is colloquially known as the “taxicab metric” (why?). Note that each of d_1, d_2, d_∞ specializes to the usual absolute value metric when $n = 1$. In fact, there’s an infinite family of metrics

$$d_p(\vec{x}, \vec{y}) = \left(\sum_{i=1}^n |x_i - y_i|^p \right)^{1/p},$$

defined for any real number $p \geq 1$. The triangle inequality for these follows from a generalization of the Cauchy-Schwarz inequality called the Minkowski inequality, which we’ll see later on in the course. Also later in the course (perhaps in a problem set) we’ll prove that

$$\lim_{p \rightarrow \infty} d_p(\vec{x}, \vec{y}) = d_\infty(\vec{x}, \vec{y}),$$

which explains the notation for d_∞ .

Of course, for $0 < p < 1$ one could try to define a metric d_p by the same formula above, but it turns out that for those values of p the “triangle inequality” would point in the wrong direction.

Although there are many metrics on \mathbb{R}^n , d_2 is generally the one that is used unless explicit mention otherwise is made.

Example 1.5. Let $[a, b]$ be any closed interval in \mathbb{R} . I hope that you’re familiar with the fact that if $h: [a, b] \rightarrow \mathbb{R}$ is a continuous function then there is some $x_0 \in [a, b]$ such that

$$h(x_0) = \max_{x \in [a, b]} h(x)$$

(if you’re not, we’ll be proving a generalization of this fact later on). In light of this, where

$$X = C([a, b]) := \{f: [a, b] \rightarrow \mathbb{R} \mid f \text{ is continuous}\},$$

we can define a metric on X by

$$d(f, g) = \max_{x \in [a, b]} |f(x) - g(x)|.$$

(Check for yourself that this is indeed a metric.)

As this example illustrates, metric space concepts apply not just to spaces whose elements are thought of as geometric points, but also sometimes to spaces of functions. Indeed, one of the major tasks later in the course, when we discuss Lebesgue integration theory, will be to understand convergence in various metric spaces of functions.

In calculus on \mathbb{R} , a fundamental role is played by those subsets of \mathbb{R} which are intervals. The analogues of open intervals in general metric spaces are the following:

Definition 1.6. If (X, d) is a metric space, $p \in X$, and $r > 0$, the **open ball of radius r around p** is

$$B_r(p) = \{q \in X \mid d(p, q) < r\}.$$

Exercise 1.7. In \mathbb{R}^2 , draw a picture of the open ball of radius 1 around the origin in the metrics d_2 , d_1 , and d_∞ .

One of the biggest themes of the whole unit on metric spaces in this course is that a major role is played by the following kinds of sets:

Definition 1.8. If (X, d) is a metric space and $U \subset X$, U is called an **open** subset of X if, for every $p \in U$, there is some $\epsilon > 0$ such that

$$B_\epsilon(p) \subset U.$$

(Of course, ϵ will typically depend on p .)

Fortunately, the terminology in the previous two definitions is consistent, because:

Proposition 1.9. *If (X, d) is a metric space, every open ball $B_r(p)$ is an open subset of X .*

Proof. Let $p \in X$ and $r > 0$. We need to show that if $q \in B_r(p)$ then some open ball $B_\epsilon(q)$ ($\epsilon > 0$) around q is contained in $B_r(p)$. Now by the definition of $B_r(p)$,

the fact that $q \in B_r(p)$ means that $d(p, q) < r$, so if we set $\epsilon = r - d(p, q)$ we have $\epsilon > 0$. If $x \in B_\epsilon(q)$, then the triangle inequality shows

$$d(p, x) \leq d(p, q) + d(q, x) < d(p, q) + \epsilon = r,$$

and so $x \in B_r(p)$. Thus $B_\epsilon(q) \subset B_r(p)$, as desired. \square

Open balls are conceptually fairly simple, but they don't behave very well when you apply set-theoretic operations to them (for instance the union of two open balls is only rarely an open ball). The more general notion of an open set is better in this regard.

Lemma 1.10. *Let (X, d) be a metric space.*

(i) *If $\{U_\alpha \mid \alpha \in A\}$ is any collection of open sets in X , then*

$$\bigcup_{\alpha \in A} U_\alpha$$

is open.

(ii) *If $\{U_1, \dots, U_n\}$ is a finite collection of open sets in X , then*

$$\bigcap_{i=1}^n U_i$$

is an open set.

Proof. For (i), if $x \in \bigcup_{\alpha \in A} U_\alpha$, then for some $\beta \in A$ we have $x \in U_\beta$ (this is just the definition of the union of a collection of sets). Since U_β is open, we have $B_\epsilon(x) \subset U_\beta$ for some $\epsilon > 0$. But of course $U_\beta \subset \bigcup_{\alpha \in A} U_\alpha$, so this shows that

$$B_\epsilon(x) \subset \bigcup_{\alpha \in A} U_\alpha.$$

Since x was an arbitrary element of $\bigcup_{\alpha \in A} U_\alpha$, this completes the proof of (i).

As for (ii), if $x \in \bigcap_{i=1}^n U_i$, then for each i there is $\epsilon_i > 0$ such that $B_{\epsilon_i}(x) \subset U_i$. Let

$$\epsilon = \min_{1 \leq i \leq n} \epsilon_i.$$

Then (since there are only finitely many i !) $\epsilon > 0$, and we have, for each i ,

$$B_\epsilon(x) \subset B_{\epsilon_i}(x) \subset U_i,$$

so that

$$B_\epsilon(x) \subset \bigcap_{i=1}^n U_i.$$

\square

Note the asymmetry between unions and intersections here; *arbitrary unions* of open sets are open, but we can only say that *finite intersections* of open sets are open. It should be easy for you to come up with examples of countable collections of open subsets of \mathbb{R} whose intersections fail to be open.

The following description of open sets in terms of open balls may be conceptually helpful:

Corollary 1.11. *If (X, d) is a metric space and $U \subset X$, U is open if and only if either $U = \emptyset$ or, for some open balls $B_{r_\alpha}(p_\alpha)$ ($\alpha \in A$, where A is some index set depending on U) we have*

$$U = \bigcup_{\alpha \in A} B_{r_\alpha}(p_\alpha).$$

Proof. For the forward implication, if U is open then for each $x \in U$ there is $r_x > 0$ such that $B_{r_x}(x) \subset U$. Thus (assuming $U \neq \emptyset$, so that the unions make sense)

$$U = \cup_{x \in U} \{x\} \subset \cup_{x \in U} B_{r_x}(x) \subset U,$$

in view of which we must have

$$U = \cup_{x \in U} B_{r_x},$$

which is precisely the sort of union of open balls promised in the statement of the corollary.

For the backward implication, note that if $U = \emptyset$ then the criterion for U to be open (“for all x in U a certain condition holds”) is satisfied vacuously since there are no $x \in U$ to test. On the other hand, if U is a union of open balls, we’ve shown that open balls are open and that unions of open sets are open, so U is open. \square

We’ll see as the unit progresses that many ideas in analysis can be expressed in terms of open sets. Nonetheless, we’ll presently see that knowing what the open sets are doesn’t tell you everything you might want to know about a metric; in particular, open sets fail to distinguish between metrics that are *equivalent* in the following sense.

Definition 1.12. Let X be a set and let d, d' be two metrics on X . We say that d is *equivalent* to d' provided that there is some $C > 0$ such that, for every $p, q \in X$, we have

$$d(p, q) \leq C d'(p, q)$$

and

$$d'(p, q) \leq C d(p, q).$$

It’s not hard to see that “equivalence” is an *equivalence relation*, i.e., d is equivalent to itself, if d is equivalent to d' then d' is equivalent to d , and if d and d' are both equivalent to d'' then they are equivalent to each other.

Example 1.13. Let $X = \mathbb{R}^n$ and consider the metrics d_1, d_2, d_∞ . We have, directly from the definitions,

$$d_\infty(\vec{x}, \vec{y}) \leq d_1(\vec{x}, \vec{y}),$$

$$d_\infty(\vec{x}, \vec{y}) \leq d_2(\vec{x}, \vec{y}),$$

$$d_1(\vec{x}, \vec{y}) \leq n d_\infty(\vec{x}, \vec{y}),$$

and

$$d_2(\vec{x}, \vec{y}) \leq \sqrt{n} d_\infty(\vec{x}, \vec{y}).$$

Thus d_1 is equivalent to d_∞ (use $C = n$) and d_2 is equivalent to d_∞ (use $C = \sqrt{n}$). Thus these three metrics (and indeed all of the d_p for $1 \leq p \leq \infty$) are equivalent to each other.

Note that while we use the term equivalent, it is of course not the case that all of these metrics are really the same (for instance their open balls look rather different). However, equivalence does guarantee that they share some properties, largely as a result of the following:

Proposition 1.14. *Suppose that X is a set with two equivalent metrics d and d' . If $U \subset X$ is open in the metric space (X, d) , then it is also open in the metric space (X, d') .*

Proof. Assume that U is open in (X, d) and that, in view of the definition of equivalence, that we have $d(p, q) \leq Cd'(p, q)$. Now if $p \in U$, we have, for some $\epsilon > 0$,

$$\{q \in X \mid d(p, q) < \epsilon\} \subset U.$$

Now if $d'(p, q) < \epsilon/C$, then $d(p, q) < \epsilon$, so we have

$$\{q \in X \mid d'(p, q) < \epsilon/C\} \subset U.$$

So since $\epsilon/C > 0$ and p was an arbitrary point in U , this shows that U is open in the metric space (X, d') . \square

As we study metric spaces, we will consistently find that certain of their fundamental properties (continuity, compactness, etc.) can be expressed purely in terms of the language of open sets. The above proposition then shows that, for any such property, as soon as we prove it for (X, d) we've also proven it for (X, d') for any d' equivalent to d . So, for example, to show that a subset of \mathbb{R}^n is compact for all of the infinitely many metrics d_p , we'll only have to check it for one of them.

1.2. Convergence and closed sets. If we have a sequence of points $\{x_n\}_{n=1}^{\infty}$ in a metric space X , we can ask whether the points x_n have a limit x ; roughly, this should mean that the x_n eventually all get so close to x as to be indistinguishable from it (by any device with finite resolution). The formal definition is just like that in the real numbers:

Definition 1.15. Let (X, d) be a metric space, $\{x_n\}_{n=1}^{\infty}$ a sequence in X , and $x \in X$. We say that the sequence $\{x_n\}_{n=1}^{\infty}$ converges to x (and write $x_n \rightarrow x$ or $\lim_{n \rightarrow \infty} x_n = x$) if the following holds: For every $\epsilon > 0$ there is some natural number N such that for all $n \geq N$ we have $d(x_n, x) < \epsilon$.

Another way of describing the condition without explicitly invoking distances is to say that $x_n \rightarrow x$ provided that whenever U is an open set such that $x \in U$, there is some N such that $x_n \in U$ for all $n \geq N$. (You should convince yourself that this is in fact an equivalent condition.) To connect with what we did above, this implies that if (X, d) and (X, d') are metric spaces such that the two metrics d and d' are equivalent, since the open sets in the two metric space are precisely the same it must be true that $x_n \rightarrow x$ in the metric space (X, d) if and only if $x_n \rightarrow x$ in the metric space (X, d') .

Definition 1.16. If (X, d) is a metric space and $F \subset X$, F is called *closed* if the following holds: whenever $\{x_n\}_{n=1}^{\infty}$ is a sequence of points in F such that $x_n \rightarrow x$ for some $x \in X$, we in fact have $x \in F$.

In other words, closed sets are “closed under taking limits,” which is the origin of the term. By contrast, it's not clear from the definition of an open set why the word “open” is an appropriate adjective. But now we arrive at an explanation for this odd choice of terms; an open set should by rights be in some sense the opposite of a closed set, and we have:

Theorem 1.17. *If (X, d) is a metric space and $F \subset X$, then F is closed if and only if its complement $X \setminus F = \{x \in X \mid x \notin F\}$ is open.*

Proof. “ \Rightarrow ”: Suppose that F is closed and that $x \in X \setminus F$. We need to show that, for some $\epsilon > 0$, we have $B_\epsilon(x) \subset X \setminus F$. Well, suppose that this were not the case. Then for any n we could find an $x_n \in X$ such that $x_n \notin X \setminus F$ (i.e., $x_n \in F$)

and $d(x, x_n) < 1/n$. But then the sequence $\{x_n\}_{n=1}^{\infty}$ is a sequence of points in F which converges to x . The fact that F is closed then forces x to be in F , which is a contradiction. This contradiction implies that it must be possible to find some $\epsilon > 0$ such that $B_{\epsilon}(x) \subset X \setminus F$, proving that $X \setminus F$ is open.

“ \Leftarrow ”: Suppose that $X \setminus F$ is open. Let $\{x_n\}_{n=1}^{\infty}$ be a sequence of points of F such that $x_n \rightarrow x$. If we had $x \in X \setminus F$, then there would be $\epsilon > 0$ such that $B_{\epsilon}(x) \subset X \setminus F$. So since each $x_n \in F$ we have $d(x_n, x) \geq \epsilon$. But this makes it impossible for $x_n \rightarrow x$, so we in fact have $x \in F$. Thus any limit of a sequence of points in F is also contained in F , so F is closed. \square

Unfortunately for linguistic purists, it is possible for a subset of a metric space (X, d) to be simultaneously open and closed; in fact the subsets \emptyset and X always have this property. For that matter, if X is a finite set (with any metric), then every subset of X is both open and closed (why?). However, it’s fairly common for a metric space (X, d) to have the property that its only simultaneously-open-and-closed subsets are \emptyset and X . Such an (X, d) is called *connected*, for reasons that we’ll probably get to later.

Corollary 1.18. *Let (X, d) be a metric space.*

(i) *If $\{F_{\alpha} | \alpha \in A\}$ is any collection of closed sets then*

$$\bigcap_{\alpha \in A} F_{\alpha}$$

is closed.

(ii) *If $\{F_1, \dots, F_n\}$ is a finite collection of closed sets then*

$$\bigcup_{i=1}^n F_i$$

is closed.

Proof. This follows directly from Theorem 1.17, Lemma 1.10 and the *deMorgan laws*

$$X \setminus \bigcup_{\alpha \in A} F_{\alpha} = \bigcap_{\alpha \in A} (X \setminus F_{\alpha}) \quad X \setminus \bigcap_{\alpha \in A} F_{\alpha} = \bigcup_{\alpha \in A} (X \setminus F_{\alpha}).$$

\square

Just as with Lemma 1.10, we really do need the collection in (ii) to be finite, since for instance if $X = \mathbb{R}$ (with the absolute value metric) and $F_n = [1/n, 1 - 1/n]$ then

$$\bigcup_{n=1}^{\infty} F_n = (0, 1),$$

which is certainly not closed even though all of the F_n are.

1.3. Limit points and closures.

Definition 1.19. Let (X, d) be a metric space and $A \subset X$. A point $x \in X$ is called a *limit point* of A if for every $\epsilon > 0$ there is some $a \in A$ such that

$$0 < d(x, a) < \epsilon.$$

In particular, if x is a limit point of A , then there is a sequence $\{a_n\}_{n=1}^{\infty}$ of points *distinct from x* such that $a_n \rightarrow x$. (We can take a_n to be an a in the definition with $\epsilon = 1/n$.) A point of A may or may not be a limit point of A and vice versa; for example if $X = \mathbb{R}$ with its usual metric and $A = (0, 1) \cup \{2\}$ then $1/2$ is both a limit point of A and an element of A , 0 is a limit point of A but not an element of A , and 2 is an element of A but not a limit point of A .

Definition 1.20. If $A \subset X$, the *closure* of A is the set

$$\bar{A} = cl_X(A) = \bigcap_{F \text{ closed}, A \subset F} F.$$

Of course, since X is closed, there is always at least one closed set containing A , so the intersection in the definition makes sense. Notice that \bar{A} is closed, since it's an intersection of closed sets. Also, if $A \subset F$ and F is closed, then F appears in the intersection defining \bar{A} , so $\bar{A} \subset F$. So \bar{A} is the *smallest closed set containing A* .

It's nice to know that a smallest closed set containing A exists, but the definition above doesn't give much of a feeling for what \bar{A} is. The following theorem remedies that:

Theorem 1.21. Let $A \subset X$ where (X, d) is a metric space. Let

$$B = \{x \in X | (\exists \{a_n\}_{n=1}^{\infty} \in A)(a_n \rightarrow x)\}.$$

Then

$$\bar{A} = B.$$

Remark 1.22. Equivalently,

$$B = \{x \in X | x \in A \text{ or } x \text{ is a limit point of } A\}$$

(convince yourself of this). So what this theorem says is that in order to turn an arbitrary set A into a closed set, we just have to add the limit points of A , and that moreover the limit points are the least we can add to A in order to get a closed set.

Proof. To show that $\bar{A} \subset B$ it's enough to show that B is closed. To see this, suppose that $\{x^{(n)}\}_{n=1}^{\infty}$ is a sequence in B and $x \in X$ with $x^{(n)} \rightarrow x$; we need to show that $x \in B$. In other words, we need to show that for any $\epsilon > 0$ there is some $a \in A$ such that $d(a, x) < \epsilon$. But for a sufficiently large N we will have $d(x^{(N)}, x) < \epsilon/2$, and since there is a sequence in A converging to $x^{(N)}$ we can find $a \in A$ such that $d(a, x^{(N)}) < \epsilon/2$. So the triangle inequality gives $d(a, x) < \epsilon$, as needed.

So it just remains to show that $B \subset \bar{A}$. Now if $b \in B$, there is a sequence $\{a_n\}_{n=1}^{\infty}$ of elements of A such that $a_n \rightarrow b$. But $A \subset \bar{A}$, so $a_n \in \bar{A}$. Since \bar{A} is closed, we then immediately get that $b = \lim_{n \rightarrow \infty} a_n \in \bar{A}$. \square

1.4. Completeness. Roughly speaking, a metric space (X, d) is complete if every sequence that should converge does converge. Here is what I mean by a sequence which "should converge."

Definition 1.23. If (X, d) is a metric space and $\{x_n\}_{n=1}^{\infty}$ is a sequence in X , we call $\{x_n\}_{n=1}^{\infty}$ a *Cauchy sequence* if the following holds: For every $\epsilon > 0$ there is some N such that whenever $n, m \geq N$ we have $d(x_m, x_n) < \epsilon$.

Notice that a convergent sequence is Cauchy, since if $x_n \rightarrow x$ and $\epsilon > 0$ then there is N such that when $n \geq N$ we have $d(x_n, x) < \epsilon/2$, and then the triangle inequality shows that for $m, n \geq N$ $d(x_n, x_m) < \epsilon$.

Definition 1.24. A metric space (X, d) is *complete* if for every Cauchy sequence $\{x_n\}_{n=1}^{\infty}$ in X , there is $x \in X$ such that $x_n \rightarrow x$.

If a sequence is Cauchy, then the locations of its points can be specified to arbitrarily good precision by looking far enough along in the sequence, so completeness means that if we specify with arbitrary precision the location of a point, then there is actually a point there.

Example 1.25. The set \mathbb{Q} of rational numbers, equipped with the absolute value metric, is not complete, since if $x \in \mathbb{R}$ is an irrational number and if x_n is the number gotten by taking the first n digits in the decimal expansion of x , if $n, m \geq N$ we have $d(x_n, x_m) < 10^{-N}$. Thus the x_n form a Cauchy sequence in \mathbb{Q} , but there's nothing in \mathbb{Q} to which they converge. \mathbb{R} , on the other hand, is complete; in fact if you've ever seen the set \mathbb{R} rigorously defined it's somewhat likely that you saw it defined precisely so that it would have this property. Using the fact that \mathbb{R} is complete, it's not too hard to show that \mathbb{R}^n is complete (with the metric d_2 , or indeed with any of the metrics d_p).

Here is a subtler example: Recall the metric d on the set $C([a, b])$ of continuous functions $f: [a, b] \rightarrow \mathbb{R}$, defined by

$$d(f, g) = \max_{x \in [a, b]} |f(x) - g(x)|.$$

Theorem 1.26. $(C([a, b]), d)$ is a complete metric space.

Proof. Let $\{f_n\}_{n=1}^\infty$ be a Cauchy sequence in $C([a, b])$. By the definitions, this means that if $\epsilon > 0$ there is N such that whenever $n, m \geq N$ we have

$$\max_{x \in [a, b]} |f_n(x) - f_m(x)| < \epsilon.$$

In particular, for any x , one has $|f_n(x) - f_m(x)| < \epsilon$ for each $n, m \geq N$. But then this shows that, for any x the sequence $\{f_n(x)\}_{n=1}^\infty$ of real numbers is Cauchy. Now \mathbb{R} is complete, so for every x the Cauchy sequence $\{f_n(x)\}_{n=1}^\infty$ has a limit $f(x) \in \mathbb{R}$.

This defines a function $f: [a, b] \rightarrow \mathbb{R}$; the claim now is that f is a continuous function and that $f_n \rightarrow f$ in the metric space $(C([a, b]), d)$. Now if $\epsilon > 0$ and N is such that $d(f_n, f_m) < \epsilon/2$ for every $n, m \geq N$ (i.e., N has the property that for each x we have $|f_n(x) - f_m(x)| < \epsilon/2$ whenever $n, m \geq N$), if $n \geq N$ we see that, for every x ,

$$|f(x) - f_n(x)| = \lim_{m \rightarrow \infty} |f_m(x) - f_n(x)| \leq \epsilon/2 < \epsilon,$$

i.e.

$$\max_{x \in [a, b]} |f(x) - f_n(x)| < \epsilon.$$

(In a different language, what we've shown is that if a sequence of continuous functions is "uniformly Cauchy" then it converges uniformly to some function f). Hence if we show the first part of the claim, namely that $f \in C([a, b])$, then by the definition of the metric d it will in fact be the case that $d(f_n, f) \rightarrow 0$, i.e., that $f_n \rightarrow f$.

Thus all that remains is to show that the function f is continuous. In fact, as you may already be aware from a previous course, whenever a sequence of continuous functions f_n converges uniformly to some other function f , that function is continuous. We recall the proof of this fact. Fix $x_0 \in [a, b]$. Let $\epsilon > 0$; there is then some N such that, for every $n \geq N$ and $x \in [a, b]$, we have $|f_n(x) - f(x)| < \epsilon/3$. (The "uniformity" of the convergence allows us to take N independent of x). Now f_N is continuous, so for some $\delta > 0$ it holds that if $|x_0 - x| < \delta$ then $|f_N(x_0) - f_N(x)| < \epsilon/3$.

But then for that same δ , we have, if $|x - x_0| < \delta$, then

$$|f(x_0) - f(x)| \leq |f(x_0) - f_N(x_0)| + |f_N(x_0) - f(x)| + |f(x) - f_N(x)| < \epsilon/3 + \epsilon/3 + \epsilon/3 = \epsilon.$$

Thus f is continuous at x_0 ; x_0 was an arbitrary point in $[a, b]$ so this shows that $f : [a, b] \rightarrow \mathbb{R}$ is continuous, completing the proof. \square

In the proof above, an important role was played by the fact that the limit of a uniformly converging sequence of continuous functions is continuous. Note that by contrast a pointwise limit of continuous functions might not be continuous, where by definition we say that $f_n \rightarrow f$ pointwise if for every x it is true that $f_n(x) \rightarrow f(x)$. (So in terms of ϵ 's and N 's, in uniform convergence the N can be taken independent of x , whereas for pointwise convergence it may depend on x .) For example, define functions $f_n : [0, 2] \rightarrow \mathbb{R}$ by

$$f_n(x) = \begin{cases} x^n & 0 \leq x \leq 1 \\ 1 & 1 \leq x \leq 2 \end{cases}$$

It's easy to see that where $f(x) = 0$ if $0 \leq x < 1$ and $f(x) = 1$ where $1 \leq x \leq 2$ we have $f_n(x) \rightarrow f(x)$ for every x even though f fails to be continuous.

In the homework, you'll show that another metric d' on $C([a, b])$, given by $d'(f, g) = \int_a^b |f(x) - g(x)| dx$, fails to be complete.

In general, if (X, d) is an incomplete metric space, there's a procedure for extending it to a smallest-possible complete metric space (called its completion), roughly speaking by adding in all of the limits that ought to be there. Call two Cauchy sequences $\{x_n\}_{n=1}^\infty, \{y_n\}_{n=1}^\infty$ *equivalent* if one has $d(x_n, y_n) \rightarrow 0$ as $n \rightarrow \infty$, and let \hat{X} be the collection of equivalence classes of Cauchy sequences. Given two such equivalence classes \hat{x}, \hat{y} , let $\{x_n\}_{n=1}^\infty, \{y_n\}_{n=1}^\infty$ be Cauchy sequences representing the classes x and y respectively, and let $\hat{d}(\hat{x}, \hat{y}) = \lim_{n \rightarrow \infty} d(x_n, y_n)$. One can show that this is definition of \hat{d} independent of the choices of representatives of \hat{x} and \hat{y} , and in fact defines a complete metric on \hat{X} . (\hat{X}, \hat{d}) contains a copy of the original metric space (X, d) inside it, consisting of equivalence classes of Cauchy sequences of the special form $x_n = x$.

While it's nice to know that any metric space can be "completed" as in the previous paragraph, the above construction is obviously a bit complicated and doesn't give much of a feeling for what the elements of the completion are. (For instance, \mathbb{R} can be viewed as the completion of \mathbb{Q} , but in practice you presumably don't think of real numbers as equivalence classes of Cauchy sequences.) In some cases, the completion can be given a more "hands-on" construction. In the case of the incomplete metric space $(C([a, b]), d')$ we'll learn about an explicit construction later on in the unit on measure and integration theory.

1.5. Continuity. Most fields of mathematics deal with sets carrying some additional structure (in our case, this additional structure is given by a metric); to understand the relationships between two different such objects one then looks at *functions* from the first set to the second which in some sense respect the additional structure. In the theory of metric spaces, the requirement that is imposed on such functions is *continuity*. We'll give three equivalent conditions for a map from one metric space to another to be continuous; one of these involves ϵ 's, another involves sequences, and the last involves open sets. When we get to another fundamental

concept, compactness, we'll similarly see that there are three equivalent conditions along these lines for that as well.

Theorem 1.27. *Let $(X, d_X), (Y, d_Y)$ be two metric spaces, and let $f: X \rightarrow Y$ be any map. Then the following are equivalent (i.e., any one of them implies the others):*

- (i) *For every $x \in X$ and $\epsilon > 0$ there is $\delta > 0$ such that whenever $x' \in X$ satisfies $d_X(x, x') < \delta$, we have $d_Y(f(x), f(x')) < \epsilon$.*
- (ii) *If $\{x_n\}_{n=1}^\infty$ is a sequence in X such that $x_n \rightarrow x \in X$, then $f(x_n) \rightarrow f(x)$.*
- (iii) *If V is any open set in the metric space (Y, d_Y) then $f^{-1}(V) = \{x \in X \mid f(x) \in V\}$ is an open set in (X, d_X) .*

Definition 1.28. If $(X, d_X), (Y, d_Y)$ and $f: X \rightarrow Y$, f is called *continuous* if it satisfies one (and hence, by Theorem 1.27, all three) of the above conditions (i)-(iii).

Proof of Theorem 1.27. (i) \Rightarrow (ii): If (i) holds, let $\{x_n\}_{n=1}^\infty$ be a sequence in X such that $x_n \rightarrow x \in X$. Given $\epsilon > 0$, let $\delta > 0$ be a number (whose existence is implied by (i)) such that $d_X(x', x) < \delta \Rightarrow d_Y(f(x'), f(x)) < \epsilon$. Since $x_n \rightarrow x$, there is some N such that whenever $n \geq N$ we have $d_X(x_n, x) < \delta$, so the property that we've assumed for δ implies that, when $n \geq N$, $d_Y(f(x_n), f(x)) < \epsilon$. ϵ was an arbitrary positive number, so this shows that $f(x_n) \rightarrow f(x)$, proving (ii).

(ii) \Rightarrow (i): Assume that (ii) holds, and let $x \in X$ and $\epsilon > 0$. Suppose (to get a contradiction) that (i) failed to hold for this choice of x, ϵ . Then for every natural number n it would be possible to find x_n such that $d_X(x_n, x) < 1/n$ but $d_Y(f(x_n), f(x)) \geq \epsilon$. But then $\{x_n\}_{n=1}^\infty$ would be a sequence in X such that $x_n \rightarrow x$ but $f(x_n)$ does not converge to $f(x)$, and this contradicts (ii).

(i) \Rightarrow (iii): Assume (i), and suppose that $V \subset Y$ is open. Let $x \in f^{-1}(V)$. Since $f(x)$ belongs to the open set V , there is $\epsilon > 0$ such that whenever $d_Y(y, f(x)) < \epsilon$ we have $y \in V$. Let $\delta > 0$ be as given by (i). Then if $d_X(x', x) < \delta$, we have $d_Y(f(x'), f(x)) < \epsilon$, so that $f(x') \in V$, i.e., $x' \in f^{-1}(V)$. So we've shown that if $x \in f^{-1}(V)$ there is $\delta > 0$ such that $B_\delta(x) \subset f^{-1}(V)$, which is precisely what it means for $f^{-1}(V)$ to be open.

(iii) \Rightarrow (i): If (iii) holds, let $x \in X$ and $\epsilon > 0$. Now the open ball $B_\epsilon(f(x))$ is an open subset of Y , so $f^{-1}(B_\epsilon(f(x)))$ is open in X by (iii). Since $x \in f^{-1}(B_\epsilon(f(x)))$, there is then $\delta > 0$ such that $B_\delta(x) \subset f^{-1}(B_\epsilon(f(x)))$. So if $d_X(x', x) < \delta$, we have $x' \in f^{-1}(B_\epsilon(f(x)))$, i.e., $d_Y(f(x'), f(x)) < \epsilon$. \square

Although the above proof was not especially hard, the equivalence of the open set condition (iii) to the others might seem a bit odd if you've never seen it before. It can be a useful way to streamline some arguments, as the proof we give of the following shows (though admittedly this wouldn't be difficult to prove if we used (i) or (ii) instead).

Corollary 1.29. *If $(X, d_X), (Y, d_Y), (Z, d_Z)$ are metric spaces and $f: X \rightarrow Y$ and $g: Y \rightarrow Z$ are continuous functions, then $g \circ f: X \rightarrow Z$ is continuous (where as usual $g \circ f$ is defined by $(g \circ f)(x) = g(f(x))$).*

Proof. If $V \subset Z$ is open, just note that

$$(g \circ f)^{-1}(V) = \{x \in X \mid g(f(x)) \in V\} = \{x \in X \mid f(x) \in g^{-1}(V)\} = f^{-1}(g^{-1}(V)).$$

But $g^{-1}(V)$ is open in Y by the continuity of g , in light of which $f^{-1}(g^{-1}(V))$ is open in X by the continuity of f . \square

So, for instance, if $f: X \rightarrow \mathbb{R}$ is continuous, then so is, say, $e^f: X \rightarrow \mathbb{R}$ (since $x \rightarrow e^x$ is a continuous function from \mathbb{R} to itself). A bit more generally, suppose that $\Phi: \mathbb{R}^n \rightarrow \mathbb{R}$ is any continuous function. Now it's not hard to show that if $f_1, \dots, f_n: X \rightarrow \mathbb{R}$ are continuous then the function $F: X \rightarrow \mathbb{R}^n$ defined by $F(x) = (f_1(x), \dots, f_n(x))$ is continuous as well (I'll leave the proof to the reader). So the corollary shows that the function $X \rightarrow \mathbb{R}$ given by $x \mapsto \Phi(f_1(x), \dots, f_n(x))$ is continuous as soon as f_1, \dots, f_n and Φ are. As a simple example, taking $\Phi(x_1, \dots, x_n) = \sum_{i=1}^n x_i$ shows that if the f_i are continuous so is their sum. This latter fact isn't hard to prove directly, but the point here is that there are infinitely many different similar facts (one for each choice of Φ) which are all subsumed under this.

1.6. Compactness. Compactness is a fundamental property that distinguishes some metric spaces from others. The formal definition will appear below; remark first that the subsets of \mathbb{R} which are compact (with the metric restricted from \mathbb{R}) are precisely the sets which are both closed and bounded. For instance, finite sets are compact, as are closed, finite-length intervals $[a, b]$, while sets such as $(0, 1)$ and $[a, \infty)$ are not. We'll see later that compact metric spaces (X, d) share a number of properties with finite sets and closed intervals in \mathbb{R} , such as the fact that any continuous real-valued function on X attains its maximum.

Similarly to the continuity section, we'll define a metric space to be compact if it satisfies any one of three equivalent properties; that these three are equivalent will be an important theorem (with the proof a good deal harder than in the continuity case). First we need a background definition:

Definition 1.30. If (X, d) is a metric space, we call (X, d) *totally bounded* if, for every $\epsilon > 0$ there are finitely many points x_1, \dots, x_N such that

$$X = \cup_{i=1}^N B_\epsilon(x_i).$$

Incidentally, X is called bounded if there is $M \in \mathbb{R}$ such that every $x, y \in X$ we have $d(x, y) \leq M$. Any totally bounded set is bounded; taking $\epsilon = 1$ in the definition and letting x_1, \dots, x_N be such that $X = \cup_{i=1}^N B^1(x_i)$, we can take

$$M = 2 + \max_{1 \leq i, j \leq N} d(x_i, x_j).$$

It's true that a subset of \mathbb{R} (or more generally of \mathbb{R}^n) which is bounded is also totally bounded (why?); however it is possible to find bounded metric spaces which are not totally bounded.

Theorem 1.31. *Let (X, d) be a metric space. The following are equivalent:*

- (i) (X, d) satisfies the Heine-Borel Property: If $\{U_\alpha\}_{\alpha \in A}$ is a collection of open subsets of X such that

$$X = \cup_{\alpha \in A} U_\alpha,$$

then there are (finitely many) $\alpha_1, \dots, \alpha_n \in A$ such that

$$X = U_{\alpha_1} \cup \dots \cup U_{\alpha_n}.$$

- (ii) (X, d) is sequentially compact: If $\{x_n\}_{n=1}^\infty$ is any sequence in X , then there is a subsequence $\{x_{n_k}\}_{k=1}^\infty$ ($n_1 < \dots < n_k < \dots$) and $x \in X$ such that $x_{n_k} \rightarrow x$.
- (iii) (X, d) is complete and totally bounded.

Definition 1.32. A metric space (X, d) is called compact if it satisfies any (and hence all) of the above properties (i),(ii),(iii).

Proof. In outline, we'll show that (i) implies (ii), then that (ii) and (iii) are equivalent, and finally that, taken together, (ii) and (iii) imply (i).

(i) \Rightarrow (ii): Let $\{x_n\}_{n=1}^\infty$ be a sequence in X ; we need to extract a convergent subsequence assuming the Heine-Borel property. Define

$$A_N = \{x_n | n \geq N\}.$$

We claim first that

$$\bigcap_{N=1}^\infty \overline{A_N} \neq \emptyset.$$

Indeed, if this were not the case we would then have

$$X = \bigcup_{N=1}^\infty (X \setminus \overline{A_N}).$$

So the Heine-Borel property implies that just finitely many of the sets $X \setminus \overline{A_N}$ suffice to cover X ; thus for some (finite) k

$$X = \bigcup_{N=1}^k (X \setminus \overline{A_N}).$$

Now one has $A_{N+1} \subset A_N$, so $X \setminus \overline{A_N} \subset X \setminus \overline{A_{N+1}}$ for all N . Thus, for all $N \leq k$, $X \setminus \overline{A_N} \subset X \setminus \overline{A_k}$, so we in fact have $X = X \setminus \overline{A_k}$. But $\overline{A_k}$ is a nonempty subset of X (for instance, it contains x_k), so this is nonsense. This contradiction proves that, indeed,

$$\bigcap_{N=1}^\infty \overline{A_N} \neq \emptyset.$$

Accordingly, let

$$x \in \bigcap_{N=1}^\infty \overline{A_N};$$

we claim that x is the limit of some subsequence of $\{x_n\}_{n=1}^\infty$. We form the subsequence inductively; for the base step, the fact that $x \in \overline{A_1}$ shows that there is some $n_1 \geq 1$ such that $d(x_{n_1}, x) < 1$. Now assume that we've chosen $n_1 < \dots < n_k$ such that, for $1 \leq j \leq k$ we have $d(x_{n_j}, x) < 1/j$. Then since $x \in \overline{A_{n_k+1}}$, there is some $n_{k+1} > n_k$ such that $d(x_{n_{k+1}}, x) < 1/(k+1)$. By induction on k , this gives $\{x_{n_k}\}_{k=1}^\infty$ such that $n_1 < \dots < n_k < \dots$ and $d(x_{n_k}, x) < 1/k$ for every k ; hence $x_{n_k} \rightarrow x$, as desired.

(ii) \Rightarrow (iii): First we show that if (X, d) is sequentially compact then it is complete. Indeed, if $\{x_n\}_{n=1}^\infty$ is a Cauchy sequence, sequential compactness lets us extract a subsequence $\{x_{n_k}\}_{k=1}^\infty$ such that $x_{n_k} \rightarrow x \in X$. Let $\epsilon > 0$. There is then some N_1 with the property that if $n, m \geq N_1$ then $d(x_n, x_m) < \epsilon/2$, and some N_2 with the property that if k is such that $n_k \geq N_2$ then $d(x_{n_k}, x) < \epsilon/2$. Let $N = \max\{N_1, N_2\}$ and let $n \geq N$. Then if k is such that $n_k \geq N$ (of course since $n_k \nearrow \infty$ such k exists) we have

$$d(x_n, x) \leq d(x_n, x_{n_k}) + d(x_{n_k}, x) < \epsilon/2 + \epsilon/2 = \epsilon,$$

proving that $x_n \rightarrow x$. So sequential compactness implies that every Cauchy sequence converges, *i.e.* it implies completeness.

Now we prove that if (X, d) is sequentially compact then it is totally bounded. Rather, we'll prove the contrapositive; so suppose that (X, d) is not totally bounded, in other words, that, for some $\epsilon > 0$, X cannot be written as a finite union of open balls of radius ϵ . We will then construct a sequence $\{x_n\}_{n=1}^\infty$ which has no convergent subsequence. Namely, choose $x_1 \in X$ arbitrarily, and then assume inductively that we have chosen x_1, \dots, x_n such that $d(x_i, x_j) \geq \epsilon$ for all $i, j \in$

$\{1, \dots, n\}$. Our choice of ϵ shows that $X \neq \cup_{i=1}^n B_\epsilon(x_i)$, so we can then choose x_{n+1} such that, for each $i \leq n$, $x_{n+1} \notin B_\epsilon(x_i)$. We then have $d(x_i, x_j) \geq \epsilon$ whenever $1 \leq i, j \leq n+1$. Carrying this out successively for every n , we then get a sequence $\{x_n\}_{n=1}^\infty$ such that $d(x_i, x_j) \geq \epsilon$ for every i, j . If $\{x_{n_k}\}_{k=1}^\infty$ is any subsequence of this sequence, any two of the elements of this subsequence are a distance at least ϵ away from each other, so $\{x_{n_k}\}_{k=1}^\infty$ cannot be Cauchy and hence cannot converge. This proves that if (X, d) is not totally bounded then it is not sequentially compact, completing the proof that (ii) implies (iii).

(iii) \Rightarrow (ii): Assume that (X, d) is complete and totally bounded, and let $\{x_n\}_{n=1}^\infty$ be any sequence in X . Because X is totally bounded, for any integer k there are finitely many points $y_1^k, \dots, y_{m_k}^k$ with the property that

$$X = \cup_{i=1}^{m_k} B_{1/k}(y_i^k).$$

We claim that we can choose a sequence of integers $\{i_k\}_{k=1}^\infty$ with $1 \leq i_k \leq m_k$ such that, for every $l \geq 1$, the set

$$\cap_{k=1}^l B_{1/k}(y_{i_k}^k)$$

contains the point x_n for infinitely many different choices of n . Indeed, this holds for $k=1$ since there are just finitely many $B_1(y_i^1)$ and each $x_n \in X = \cup_{i=1}^{m_1} B_1(y_i^1)$. Assuming inductively that we've chosen i_1, \dots, i_l with the desired property, we know that there are infinitely many values of n such that

$$x_n \in \cap_{k=1}^l B_{1/k}(y_{i_k}^k) = \bigcup_{i=1}^{m_{l+1}} ((\cap_{k=1}^l B_{1/k}(y_{i_k}^k)) \cap B_{1/(l+1)}(y_i^{l+1})),$$

so since the union on the right is a finite union, we can choose i_{l+1} such that $\cap_{k=1}^{l+1} B_{1/k}(y_{i_k}^k)$ contains x_n for infinitely many values of n . This completes the inductive definition of $\{i_k\}_{k=1}^\infty$ with the desired property.

With this choice of the i_k , we can then form a strictly increasing sequence $\{n_l\}_{l=1}^\infty$ of integers n_l with the property that, for each l , $x_{n_l} \in \cap_{k=1}^l B_{1/k}(y_{i_k}^k)$. Notice that if $k \leq l$ then $x_{n_k}, x_{n_l} \in B_{1/k}(y_{i_k}^k)$, so

$$d(x_{n_k}, x_{n_l}) \leq d(x_{n_k}, y_{i_k}^k) + d(x_{n_l}, y_{i_k}^k) < 2/k,$$

which implies that $\{x_{n_l}\}_{l=1}^\infty$ is a Cauchy sequence. Hence since (X, d) is complete we have $x_{n_l} \rightarrow x$ for some $x \in X$, proving sequential compactness.

(ii) and (iii) \Rightarrow (i): We begin with the following lemma:

Lemma 1.33. *Let (X, d) be a sequentially compact metric space and let $\mathcal{U} = \{U_\alpha | \alpha \in A\}$ be a collection of open subsets of X such that $X = \cup_{\alpha \in A} U_\alpha$. Then there is a number $r > 0$ (called the Lebesgue number of \mathcal{U}) such that, for every $x \in X$ there is some $\beta \in A$ such that $B_r(x) \subset U_\beta$.*

The point of this lemma is that the number r can be taken independently of x .

Proof. Given $x \in X$, define the following subset of $(0, \infty)$

$$R_x = \{r > 0 | (\exists \beta)(B_r(x) \subset U_\beta)\}.$$

If there is a number r with the property that $r \in R_x$ for every x , then this r will satisfy the requirements of the lemma.

Since the U_β are each open, and since any x belongs to some U_β , each R_x is nonempty. Also, if $r \in R_x$ and $0 < s < r$ then $s \in R_x$, which implies that R_x is an interval, of the form $(0, r_x)$, $(0, r_x]$, or $(0, \infty)$.

Now let $x, y \in X$ and suppose that $M \in R_x$ with $M > d(x, y)$. One then has $B_M(x) \subset U_\beta$ for some β , so $B_{M-d(x,y)}(y) \subset U_\beta$ also. Thus

$$(1) \quad \text{if } d(x, y) < M \in R_x \text{ then } M - d(x, y) \in R_y.$$

Suppose now that there is *any* $x \in X$ with the property that $R_x = (0, \infty)$. Then if $y \in X$ and N is any positive number, we have $N + d(x, y) \in R_x$, so by (1) $N \in R_y$. Thus R_y contains all positive numbers, which is to say that $R_y = (0, \infty)$. So we've shown that if any of the R_x is $(0, \infty)$ then all of them are. In particular, for every $x \in X$ we have $1 \in R_x$, in view of which the lemma holds with $r = 1$.

So for the rest of the proof we can assume that none of the R_x is equal to $(0, \infty)$, so they are all equal to either $(0, r_x)$ or $(0, r_x]$ for some finite positive r_x . Now consider the function

$$\begin{aligned} \rho: X &\rightarrow \mathbb{R} \\ x &\mapsto r_x. \end{aligned}$$

I claim that ρ is continuous. Indeed, let $x, y \in \mathbb{R}$ and suppose that t is any number such that $d(x, y) \leq t < \rho(x) = r_x$. Then $t \in R_x$, so by (1) $t - d(x, y) \in R_y$, so $t - d(x, y) \leq r_y$. Since this holds for *every* $t < r_x$ we get that, provided $d(x, y) < r_x$,

$$r_x - d(x, y) \leq r_y, \text{ i.e., } \rho(x) - \rho(y) \leq d(x, y).$$

But the same argument applies equally well with the roles of x and y reversed as long as $d(x, y) < r_y$, giving that in this case $\rho(y) - \rho(x) \leq d(x, y)$. Now if $r_y \leq r_x$, the equation $\rho(y) - \rho(x) \leq d(x, y)$ is trivial, so in fact if $d(x, y) < r_x$ then we still have $\rho(y) - \rho(x) \leq d(x, y)$ (since either $d(x, y) < r_y$ or else $r_y \leq r_x$). Taken together, this shows that

$$|\rho(x) - \rho(y)| < d(x, y) \text{ whenever } d(x, y) < r_x.$$

So since r_x is always positive, we immediately get from the definition (i) of continuity that ρ is continuous.

Up to now we haven't used the sequential compactness assumption; now we finally do. What we need to show is that there is $r > 0$ such that $r \in R_x$ for every x ; in other words we need to show that $\rho(x) \geq r$ for some $r > 0$ and all x . If this were not the case, we could find a sequence $\{x_n\}_{n=1}^\infty$ in X such that for each n $\rho(x) < 1/n$. By sequential compactness, there would be $x \in X$ and a subsequence $\{x_{n_k}\}_{k=1}^\infty$ such that $x_{n_k} \rightarrow x$. But then since ρ is continuous, we would have $\rho(x) = \lim_{k \rightarrow \infty} \rho(x_{n_k}) = 0$. But this is a contradiction since ρ was constructed to have the property that $\rho(x) > 0$ for every x . This contradiction shows that the desired r must exist. \square

Now we show that if (X, d) is sequentially compact and totally bounded then it satisfies the Heine-Borel property. Let $\mathcal{U} = \{U_\alpha | \alpha \in A\}$ be an open cover of X . By sequential compactness, \mathcal{U} has a Lebesgue number $r > 0$, so that for each $x \in X$ there is $\beta(x) \in A$ such that

$$B_r(x) \subset U_{\beta(x)}.$$

But because (X, d) is totally bounded, there are $x_1, \dots, x_N \in X$ such that

$$X = \cup_{i=1}^N B_r(x_i).$$

But then

$$X = \cup_{i=1}^N U_{\beta(x_i)},$$

and we have found the desired finite subcover. \square

As we'll soon see, all this work makes it possible to quickly dispense with some classic results from the theory of functions of a real variable; first, though, we should make a digression concerning subspaces of metric spaces.

1.7. Subspaces.

Definition 1.34. Let (X, d) be a metric space, and let $S \subset X$ be a subset. The *subspace metric* on S is the metric $d_S: S \times S \rightarrow \mathbb{R}$ defined by $d_S(s, t) = d(s, t)$ for all $s, t \in S$. A *subspace* of the metric space X is defined to be a metric space of the form (S, d_S) where $S \subset X$.

We'll need a lemma about how open sets in a subspace are related to open sets in the larger space. The proof will need to talk about open balls both in the larger space and in the subspace, so if $S \subset X$ and $x \in S$ we'll write

$$B_r^X(p) = \{y \in X \mid d(x, y) < r\} \quad B_r^S(p) = \{y \in S \mid d_S(x, y) < r\}.$$

Thus $B_r^S(p)$ is an open set in (S, d_S) , but it's probably not an open set in X . $B_r^X(p)$, meanwhile, is an open set in X but is probably not contained in S . Notice that, since if $x, y \in S$ then $d(x, y) = d_S(x, y)$, we have

$$B_r^S(p) = \{y \in S \mid d(x, y) < r\} = S \cap B_r^X(p).$$

Lemma 1.35. *Let (S, d_S) be a subspace of the metric space (X, d) . Let $V \subset S$. Then V is open (considered as a subset of S) if and only if there is a subset $U \subset X$ such that U is open (considered as a subset of X) and*

$$V = S \cap U.$$

Proof. First we'll show that if $V \subset S$ has the form $V = S \cap U$ with $U \subset X$ open, then V is open. Indeed, let $x \in V$. Then $x \in U$, so there is $r > 0$ such that if $y \in X$ and $d(x, y) < r$ we have $y \in U$. But then if $y \in S$ and $d_S(x, y) < r$ then $y \in U$, while $y \in S$ by assumption, so $y \in V$. So we've shown that $B_r^S(x) \subset V$, proving that V is open when considered as a subset of S .

Conversely, suppose that $V \subset S$ is open (in the metric space (S, d_S)). By one of our first results (Corollary 1.11), one has $V = \cup_{\alpha \in A} B_{r_\alpha}^S(p_\alpha)$ for some $p_\alpha \in V$ and $r_\alpha > 0$. But by our remarks before the statement of the lemma, for each α it holds that $B_{r_\alpha}^S(p_\alpha) = S \cap B_{r_\alpha}^X(p_\alpha)$. Hence where

$$U = \cup_{\alpha \in A} B_{r_\alpha}^X(p_\alpha),$$

it follows that $V = S \cap U$. \square

Corollary 1.36. *Let (X, d) be a metric space and let $S \subset X$. Then the subspace (S, d_S) is a compact metric space if and only if the following holds: Whenever $\mathcal{U} = \{U_\alpha\}_{\alpha \in A}$ is a collection of open sets in X such that*

$$S \subset \cup_{\alpha \in A} U_\alpha,$$

there are $\alpha_1, \dots, \alpha_n$ such that

$$S \subset U_{\alpha_1} \cup \dots \cup U_{\alpha_n}.$$

Proof. Suppose (S, d_S) is compact, and let $S \subset \cup_{\alpha \in A} U_\alpha$ where each U_α is open in X . Then $S = \cup_{\alpha \in A} (S \cap U_\alpha)$, and by the lemma the $S \cap U_\alpha$ are all open subsets of S , so there are $\alpha_1, \dots, \alpha_n$ such that

$$S = (S \cap U_{\alpha_1}) \cup \dots \cup (S \cap U_{\alpha_n}), \quad \text{i.e., } S \subset U_{\alpha_1} \cup \dots \cup U_{\alpha_n}.$$

Conversely, suppose that S satisfies the property that we're trying to show is equivalent to compactness, and let $\{V_\alpha\}_{\alpha \in A}$ be a collection of open subsets of (S, d_S) with $S = \cup_{\alpha \in A} V_\alpha$. Then for each α there is U_α such that $V_\alpha \subset S \cap U_\alpha$, so in particular we have $S \subset \cup_{\alpha \in A} U_\alpha$. So by our assumption $S \subset U_{\alpha_1} \cup \dots \cup U_{\alpha_n}$, and so since $V_{\alpha_i} = S \cap U_{\alpha_i}$ we have $S = V_{\alpha_1} \cup \dots \cup V_{\alpha_n}$. This proves that (S, d_S) satisfies the Heine-Borel property, i.e., it is compact. \square

Remark 1.37. *In practice, if we are working in a given metric space (X, d) , we will say, e.g., “a compact subset S of X ” when we mean “a subset S of X with the property that the subspace (S, d_S) is compact.”*

1.8. Continuous functions on compact sets. Now we're ready to prove two basic results about the behavior of continuous functions on compact sets.

Theorem 1.38. *Let (X, d_X) and (Y, d_Y) be two metric spaces, and let $f: X \rightarrow Y$ be continuous. Suppose that (X, d_X) is compact. Then*

$$f(X) = \{y \in Y \mid (\exists x \in X)(f(x) = y)\}$$

is also compact (with respect to the subspace metric that it inherits from Y).

Proof. Suppose that $f(X) \subset \cup_{\alpha \in A} V_\alpha$ where the V_α are open subsets of Y . Then, for each $x \in X$, there is some α such that $f(x) \in V_\alpha$; thus

$$X \subset \cup_{\alpha \in A} f^{-1}(V_\alpha).$$

Note that each $f^{-1}(V_\alpha)$ is open because f is continuous. But since X is compact (and so satisfies the Heine-Borel property) this implies that there are $\alpha_1, \dots, \alpha_n$ such that

$$X \subset \cup_{i=1}^n f^{-1}(V_{\alpha_i}),$$

so that

$$f(X) \subset \cup_{i=1}^n V_{\alpha_i}.$$

Thus $f(X)$ satisfies the Heine-Borel property, i.e., it is compact. \square

Corollary 1.39. *Let (X, d_X) be a compact metric space and let $f: X \rightarrow \mathbb{R}$ be a continuous function. Then there are $x_{min}, x_{max} \in X$ such that*

$$f(x_{min}) \leq f(x) \leq f(x_{max}) \quad \text{for all } x \in X.$$

Proof. By the theorem, and the fact that a compact subset of \mathbb{R} is necessarily bounded, we have, for some $N > 0$, $f(X) \subset [-N, N]$. Let $m = \inf f(X)$, $M = \sup f(X)$ (i.e., m is the greatest lower bound of $f(X)$, and M is the least upper bound). Thus $m \leq f(x) \leq M$ for all $x \in X$, so we need to show that there are x_{min}, x_{max} such that $f(x_{min}) = m$, $f(x_{max}) = M$. Since $m = \inf f(X)$, there is a sequence $\{x_n\}_{n=1}^\infty$ such that $f(x_n) \rightarrow m$ (for instance, x_n could be chosen to be such that $f(x_n) \leq m + 1/n$). By the compactness of X , there is a subsequence $\{x_{n_k}\}_{k=1}^\infty$ of $\{x_n\}_{n=1}^\infty$ such that, for some $x_{min} \in X$, $x_{n_k} \rightarrow x_{min}$. But then by the continuity of f we have $f(x_{min}) = \lim_{k \rightarrow \infty} f(x_{n_k}) = m$, as desired. Similarly, there is a sequence $\{y_n\}_{n=1}^\infty$ such that $f(y_n) \rightarrow M$, and letting x_{max} be the limit

of a convergent subsequence of this sequence (which exists by compactness) we will have $f(x_{max}) = M$. \square

This in particular justifies the statement we made earlier in our discussion of the metric space $(C([a, b]), d)$ that any continuous function $f: [a, b] \rightarrow \mathbb{R}$ attains its maximum.

We now recall and generalize a stronger version of continuity from the theory of functions of a real variable.

Definition 1.40. Let (X, d_X) and (Y, d_Y) be metric spaces and let $f: X \rightarrow Y$ be a function. We say that f is *uniformly continuous* if the following holds: For every $\epsilon > 0$ there is $\delta > 0$ such that whenever $x, x' \in X$ are such that $d_X(x, x') < \delta$, we have $d_Y(f(x), f(x')) < \epsilon$.

If you don't recognize immediately why this is different than ordinary continuity, then you should review the ϵ - δ definition of continuity. The distinction lies in the fact that here we need to be able to choose δ independently of x .

Any uniformly continuous function is certainly continuous, but there are plenty of counterexamples to the reverse implication. For example, define $f: \mathbb{R} \rightarrow \mathbb{R}$ by $f(x) = x^2$. f is of course continuous. But since $f(x + \delta) - f(x) = 2x\delta + \delta^2$, even taking $\epsilon = 1$ we can see that there is no δ such that $f(x + \delta)$ differs from $f(x)$ by less than 1 for all x . Thus this f is not uniformly continuous.

However, we have the following.

Theorem 1.41. Let (X, d_X) and (Y, d_Y) be two metric spaces, and let $f: X \rightarrow Y$ be continuous. Suppose that (X, d_X) is compact. Then f is uniformly continuous.

Proof. Let $\epsilon > 0$. The continuity of f means that for each $x \in X$ there is $\delta_x > 0$ such that if $d_X(x, x') < \delta_x$ then $d_Y(f(x), f(x')) < \epsilon/2$. Of course we have

$$X = \cup_{x \in X} B_{\delta_x/2}(x).$$

So by the Heine-Borel property for X , there are $x_1, \dots, x_n \in X$ such that

$$X = \cup_{i=1}^n B_{\delta_{x_i}/2}(x_i).$$

Let

$$\delta = \min\{\delta_{x_1}/2, \dots, \delta_{x_n}/2\}.$$

Then if $x, x' \in X$ are such that $d_X(x, x') < \delta$, we know that there is x_i such that $d_X(x, x_i) < \delta_{x_i}/2$, in view of which

$$d_X(x', x_i) \leq d_X(x', x) + d_X(x, x_i) < \delta + \delta_{x_i}/2 < \delta_{x_i}.$$

Hence $d_Y(f(x'), f(x_i)) < \epsilon/2$ and $d_Y(f(x), f(x_i)) < \epsilon/2$, so by the triangle inequality $d_Y(f(x), f(x')) < \epsilon$. Thus f is uniformly continuous. \square

1.9. Connectedness. The last general property of metric spaces that we're going to discuss is *connectedness*. To try to motivate this, consider the two metric spaces $X_1 = [-1, 1]$, $X_2 = [0, 1] \cup [2, 3]$, both equipped with the subspace metric coming from the standard metric on \mathbb{R} . It should be intuitively clear that X_1 has "one piece" whereas X_2 has two pieces. A more careful way of saying what that means is that, in X_2 , the subset $[0, 1]$ is both open and closed (when considered as a subset of X_2 , not of \mathbb{R}), as is the subset $[2, 3]$; on the other hand, one can show that there are no subsets of X_1 other than the empty set and all of X_1 which are both open

and closed. Another manifestation of the fact that X_2 has more than one piece is that it's not possible to move continuously from 0 to 3 while staying in X_2 , whereas it is possible to move from any point in X_1 to any other point in X_1 without leaving X_1 . These ideas are formalized in the following:

Definition 1.42. Let (X, d) be a metric space.

- (1) (X, d) is called *connected* if, whenever $A \subset X$ is a subset which is both open and closed in (X, d) , we have either $A = \emptyset$ or $A = X$.
- (2) (X, d) is called *pathwise connected* if, for every $x, y \in X$, there is a *continuous* map $c: [0, 1] \rightarrow X$ such that $c(0) = x$ and $c(1) = y$.

Remark 1.43. *Since the complement of a closed set is open, one may equivalently say that X is connected if there are no nonempty disjoint, open sets $A, B \subset X$ such that $X = A \cup B$.*

Remark 1.44. *Sometimes instead of “pathwise connected” you’ll see the term “path-connected” or “arcwise connected.” Incidentally, if $c: [0, 1] \rightarrow X$ is a continuous function with $c(0) = x$ and $c(1) = y$, c is called a path from x to y .*

In particular, $[-1, 1]$ is connected and pathwise connected, whereas $[0, 1] \cup [2, 3]$ is neither.

Proposition 1.45. *If $A \subset \mathbb{R}$ is a connected subset and $a, b \in A$ then $[a, b] \subset A$.*

Proof. If the closed interval $[a, b]$ were not a subset of A , then there would be c such that $a < c < b$ and $c \notin A$. Now $A \cap (-\infty, c)$ (which contains a) is open in A by Lemma 1.35, as is $A \cap (c, \infty)$ (which contains b). But since $c \notin A$ we have $A = (A \cap (-\infty, c)) \cup (A \cap (c, \infty))$, so A is a disjoint union of two nonempty open sets, implying that A is not connected. \square

Proposition 1.46. *Any interval $[a, b] \subset \mathbb{R}$ is connected*

Proof. Suppose that $[a, b] = A \cup B$ where A and B are disjoint open subsets of $[a, b]$. Renaming A and B if necessary, assume that $b \in B$. Let $x = \sup A$. Thus there are $\{a_n\}_{n=1}^{\infty}$ in A such that $a_n \rightarrow x$. Now since B is open, $A = X \setminus B$ is closed, so it follows that $x \in A$. Hence since $b \in B$ $x < b$. So since A is open there is $\epsilon > 0$ (which may be taken smaller than $b - x$) such that $x + \epsilon \in A$. But this contradicts the definition of x as the supremum of A . \square

Theorem 1.47. *If a metric space (X, d) is pathwise connected, then it is connected.*

Proof. Suppose to the contrary that X were not connected, so that there is a set $A \subset X$, equal to neither X nor \emptyset , which is both open and closed. In particular, $X \setminus A$ is nonempty and open. So we can choose $a \in A$, $b \in X \setminus A$. Since X is assumed pathwise connected, there is a continuous function $c: [0, 1] \rightarrow X$ such that $c(0) = a$ and $c(1) = b$. So $0 \in c^{-1}(A)$ and $1 \in c^{-1}(X \setminus A)$. But $c^{-1}(A)$ and $c^{-1}(X \setminus A)$ are both open sets since c is continuous, and their union is all of $[0, 1]$. But we have just shown that $[0, 1]$ is connected, so this is a contradiction. \square

There exist connected metric spaces which are not pathwise connected, but the conditions do coincide for suitably “nice” metric spaces. You’ll explore the relationship between these two concepts further in the homework.

The following is a generalization of the intermediate value theorem (why?).

Theorem 1.48. *Let (X, d_X) , (Y, d_Y) be metric spaces such that X is connected, and let $f: X \rightarrow Y$ be a continuous function. Then $f(X)$ is connected (using the subspace metric that it inherits from Y).*

Proof. If $f(X) = A \cup B$ where A, B are nonempty, disjoint, and open, then $X = f^{-1}(A) \cup f^{-1}(B)$ where $f^{-1}(A)$ and $f^{-1}(B)$ are nonempty, disjoint, and (by the continuity of f) open. But X is connected so this is impossible. \square

2. SOME ODE THEORY

So far, we've just focused on some abstract properties of metric spaces; we'll now see how some of the general theory can be useful in studying a particular problem, namely the existence and uniqueness of solutions to ordinary differential equations.

We'll consider general (in particular, typically neither linear nor autonomous) first-order ODE's for a function

$$y: [a, b] \rightarrow \mathbb{R}^n$$

where $[a, b]$ is some interval in \mathbb{R} . Such an ODE has the form

$$(2) \quad y'(t) = F(t, y(t))$$

where $F: [a, b] \times \mathbb{R}^n \rightarrow \mathbb{R}^n$ is some continuous function. So we might think of y as the trajectory of a particle in n -dimensional space; the ODE then specifies the velocity of the particle at a given time, which is allowed to depend on either or both of the current time and the current position of the particle.

One might more generally want to consider higher-order differential equations, say of the form

$$(3) \quad y^{(m)}(t) = F(t, y(t), \dots, y^{(m-1)}(t))$$

(for a path $y: [a, b] \rightarrow \mathbb{R}^n$). But in fact such an equation can be converted into a first order equation, by changing the dimension n . Namely, consider a general function $x: [a, b] \rightarrow \mathbb{R}^{nm}$ which we write as an m -tuple of functions $x_i: [a, b] \rightarrow \mathbb{R}^n$, $x = (x_1, \dots, x_m)$. Then it's easy to see that y satisfies (3) if and only if the function $x = (y, y', \dots, y^{(m-1)})$ satisfies the first order ODE (system) defined by $x'_i = x_{i+1}$ for $1 \leq i \leq m-1$ and $x'_m = F(t, x_1, \dots, x_{m-1})$; this system is a first-order ODE of the type mentioned earlier.

With this said, if you've had much physics you'll realize that a great many problems in physics (indeed, probably basically every problem in classical mechanics) can be expressed in these terms. The general hope for an equation like (2) is that if we specify an *initial condition* $y(a) = y_0$ it should then be possible to find precisely one solution y which satisfies both (2) and the initial condition. If n were, say, one-sixth of the number of particles in the universe, with the vector-valued function y describing their positions and velocities (and so having 6 entries per particle), and the function F somehow encoded their interactions with each other in such a way that their evolution is prescribed by (2), this would then in some sense justify the following famous assertion of Laplace:

An intellect which at any given moment knew all of the forces that animate nature and the mutual positions of the beings that compose it, if this intellect were vast enough to submit the data to analysis, could condense into a single formula the movement of the greatest bodies of the universe and that of the lightest atom; for such an

intellect nothing could be uncertain and the future just like the past would be present before its eyes.

This isn't the proper forum to address either the physical or philosophical issues with such statements, but we can address the key mathematical motivation for it, namely the expectation that (2) should possess one and only one solution satisfying a given initial condition. Of course, there are two issues here: existence (is there at least one solution?) and uniqueness (can there be two different solutions?).

Here are some examples showing that this can be a subtle problem. We'll specialize to the case where $n = 1$ and $a = 0$, and where, for some $\alpha > 0$ we have $F(t, y) = |y|^\alpha$. For different values of α and different initial conditions, we'll find cases where either existence or uniqueness fails.

First suppose $\alpha > 1$; and consider the initial value problem (for a function $y: [0, b] \rightarrow \mathbb{R}$)

$$(4) \quad y'(t) = |y(t)|^\alpha \quad y(0) = 1.$$

Note that it's easy to see that any solution of this equation will necessarily remain positive, so $|y(t)|^\alpha = y(t)^\alpha$, so this is equivalent to solving $y' = y^\alpha$, $y(0) = 1$ (or if you don't want to deal with absolute values at all, you could just consider the case where α is an even integer). It's likely that you've learned how to solve this equation at some point in your past, using separation of variables:

$$\int_{y(0)}^{y(t)} y^{-\alpha} dy = \int_0^t dt,$$

so

$$t = \frac{1}{1-\alpha}(y(t)^{1-\alpha} - 1),$$

i.e.,

$$y(t) = ((1-\alpha)t + 1)^{1/(1-\alpha)}.$$

So for instance if $\alpha = 2$ then $y(t) = (1-t)^{-1}$. If you think about the logic of the argument, what we've shown is that any solution to (4) is necessarily given by the formula $y(t) = ((1-\alpha)t + 1)^{1/(1-\alpha)}$; thus uniqueness, at least, does hold in this case. If you put this formula for y back into the equation you see that it does indeed satisfy it. This might seem to take care of existence as well, but the problem is that we wanted y to be a function defined on the interval $[0, b]$, and the formula that we derived blows up as t approaches $1/(\alpha-1)$, so if $b > 1/(\alpha-1)$ then we've failed in our search. Moreover, since we reasoned that y was the only possible solution that could work (on any interval), it follows that there is *no* solution to (4) defined on an interval $[0, b]$ if $b > 1/(\alpha-1)$. (However, if $b < 1/(\alpha-1)$ we've seen that there is a solution, and in fact it is the unique one). Summing up, in this case we've found that uniqueness holds, and that existence holds *for sufficiently short time intervals*, but long-time existence may fail.

Now suppose $\alpha < 1$, and consider the initial value problem (for $y: [0, b] \rightarrow \mathbb{R}$)

$$(5) \quad y'(t) = |y(t)|^\alpha \quad y(0) = 0.$$

Now the same calculation as in the previous paragraph (but subbing in $y(0) = 0$ instead of $y(0) = 1$ and being somewhat cavalier about dividing by zero in the separation of variables argument) suggests the formula $y(t) = ((1-\alpha)t)^{1/(1-\alpha)}$. Since we assumed $\alpha < 1$, so that $1/(1-\alpha) > 0$, this is a continuous function, and one verifies that it does indeed solve (5). But there's another solution to (5),

namely $y(0) = 0$. Both of these solutions are defined for all time, so in this case existence does hold regardless of which b we choose in the domain interval $[0, b]$; however since there are two different solutions uniqueness fails. Thus if the universe evolved according to an ODE of this form, Laplace's hypothetical intellect would be unable to predict the future, since there would be two equally good possible solutions to the equations.

On the other hand, if $\alpha = 1$ and we take any initial condition $y(0) = C$, it's straightforward to show that on any interval $[0, b]$ the only solution to $y'(t) = y$, $y(0) = C$ is $y(t) = Ce^t$. So in this case both (long-time) existence and uniqueness do hold.

We now turn to metric space theory to get some positive results rather than the negative ones we found above. We'll see that what led to the nonuniqueness in the case that we looked at was that the function $y \mapsto y^\alpha$ isn't differentiable at $y = 0$; when differentiability holds we will get uniqueness. We'll also find "short-time" existence quite generally, and we'll see that if the function $F(t, y)$ doesn't grow too fast as $|y|$ gets large then long-time existence will hold as well.

The basic observation is the following:

Lemma 2.1. *Given a continuous function $F: [a, b] \times \mathbb{R}^n \rightarrow \mathbb{R}^n$, $y_0 \in \mathbb{R}^n$, and $y: [a, b] \rightarrow \mathbb{R}^n$, the following are equivalent:*

- (i) *y is continuously differentiable (C^1) and is a solution to the initial value problem*

$$y'(t) = F(t, y(t)) \quad (t \in [a, b]) \quad y(a) = y_0.$$

- (ii) *y is continuous, and for each $t \in [a, b]$ we have*

$$y(t) = y_0 + \int_a^t F(s, y(s)) ds.$$

Proof. Assuming (i), we have (by the Fundamental Theorem of Calculus)

$$y(t) - y(a) = y(t) - y(0) = \int_a^t y'(s) ds = \int_a^t F(s, y(s)) ds,$$

which proves (ii).

Conversely, assuming (ii), we obviously have $y(a) = y_0$, and for $h \neq 0$ we have

$$\frac{y(t+h) - y(t)}{h} = \frac{1}{h} \int_t^{t+h} F(s, y(s)) ds;$$

the right hand side tends to $F(t, y(t))$ as $h \rightarrow 0$ by (the other part of) the Fundamental Theorem of Calculus, so we deduce that the derivative $y'(t)$ exists and is equal to $F(t, y(t))$. Since $F(t, y(t))$ is a continuous function of t , this shows that y is continuously differentiable, proving (i). \square

Let

$$C([a, b]; \mathbb{R}^n) = \{y: [a, b] \rightarrow \mathbb{R}^n \mid y \text{ is continuous}\},$$

and for $y, z \in C([a, b]; \mathbb{R}^n)$ define

$$d(y, z) = \max_{t \in [a, b]} d_2(y(t), z(t)).$$

When $n = 1$, we showed earlier (Theorem 1.26) that $(C([a, b]; \mathbb{R}), d)$ is a *complete* metric space, and it's straightforward to use the same argument to show that $(C([a, b]; \mathbb{R}^n), d)$ is a complete metric space for any n .

Define the function

$$\Phi: C([a, b]; \mathbb{R}^n) \rightarrow C([a, b]; \mathbb{R}^n)$$

by

$$(\Phi(y))(t) = y_0 + \int_a^t F(s, y(s)) ds.$$

What Lemma 2.1 shows is that solving our initial value problem for y is exactly the same as finding a “fixed point” of Φ , *i.e.*, an element $y \in C([a, b]; \mathbb{R}^n)$ with the property that $\Phi(y) = y$.

Fortunately, there’s a general criterion for guaranteeing that a function on a complete metric space has a unique fixed point:

Theorem 2.2 (Contractive Mapping Principle). *Let (X, d) be a complete metric space and let $\phi: X \rightarrow X$ be a function satisfying the following property: For some $r < 1$, we have, for each $x, x' \in X$,*

$$d(\phi(x), \phi(x')) \leq rd(x, x').$$

Then there is one and only one point $x_0 \in X$ such that

$$\phi(x_0) = x_0.$$

Proof. Certainly there is no more than one fixed point of ϕ , since if $\phi(x_0) = x_0$ and $\phi(x_1) = x_1$ then our assumption shows $d(x_0, x_1) \leq rd(x_0, x_1)$, which since $r < 1$ forces $x_0 = x_1$. So we just need to prove existence.

To do so, let $x_1 \in X$ be any point, and for $n \geq 1$ define $x_{n+1} = \phi(x_n)$. We will show that the sequence $\{x_n\}_{n=1}^{\infty}$ converges, and that its limit x_0 is our desired fixed point.

To see this, set $A = d(x_1, x_2)$. For any n we then have $d(x_{n+1}, x_{n+2}) = d(\phi(x_n), \phi(x_{n+1})) \leq rd(x_n, x_{n+1})$; by induction on n this implies that

$$d(x_{n+1}, x_{n+2}) \leq Ar^n.$$

So by the triangle inequality, we find

$$d(x_n, x_{n+k}) \leq \sum_{j=0}^{k-1} d(x_{n+j}, x_{n+j+1}) \leq \sum_{j=0}^{k-1} Ar^{n-1+j} = Ar^{n-1} \sum_{j=0}^{k-1} r^j \leq \frac{Ar^{n-1}}{1-r}.$$

Hence if $n, m \geq N$, we get

$$d(x_n, x_m) \leq \frac{Ar^N}{1-r}.$$

So since $r^N \rightarrow 0$ as $N \rightarrow \infty$, this shows that $\{x_n\}_{n=1}^{\infty}$ is a Cauchy sequence. We assumed (X, d) to be complete, so for some $x_0 \in X$ we have $x_n \rightarrow x_0$. Now the assumption on ϕ obviously implies that ϕ is continuous (given ϵ we can just take $\delta = \epsilon$), so it follows that

$$\phi(x_0) = \lim_{n \rightarrow \infty} \phi(x_n) = \lim_{n \rightarrow \infty} x_{n+1} = x_0.$$

□

Note that the proof gives an explicit method of constructing the fixed point x_0 : just pick any point in X , apply ϕ repeatedly to it, and x_0 will be the limit of the resulting sequence. There are various other theorems in mathematics asserting that a map satisfying various hypotheses has a fixed point (the Brouwer fixed point theorem for example), but most of these have non-constructive proofs, in which one

assumes that no fixed point exists and derives a contradiction; such proofs give no hint as to how the fixed point might be found.

Anyway, in light of the Contractive Mapping Principle, to prove existence and uniqueness for our ODE it's enough to show that the map $\Phi: C([a, b]) \rightarrow C([a, b])$ is contractive (*i.e.*, satisfies the hypotheses of the Contractive Mapping Principle). Of course, in light of our examples from earlier, this evidently isn't always true, but we'll see next that it often does hold.

Definition 2.3. $F: [a, b] \times \mathbb{R}^n \rightarrow \mathbb{R}^n$ is called *uniformly Lipschitz* with Lipschitz constant M if, for every $t \in \mathbb{R}$ and $x, y \in \mathbb{R}^n$ we have

$$|F(t, x) - F(t, y)| \leq M|x - y|.$$

(Here $|x|$ denotes the standard Pythagorean distance of $x \in \mathbb{R}^n$ from the origin.)

In particular, if $F(t, y) = A(t)y$ for some continuous $n \times n$ matrix-valued function $A: [a, b] \rightarrow M_{n \times n}(\mathbb{R})$ then F is uniformly Lipschitz. In fact, you can check that for a Lipschitz constant M in this case one could take

$$\max_{t \in [a, b]} \sqrt{\sum_{i, j} A(t)_{i, j}^2}.$$

On the other hand, $F(t, y) = |y|^\alpha$ is uniformly Lipschitz only if $\alpha = 1$; if $\alpha < 1$ the estimate that we need fails for $x - y$ small, while for $\alpha > 1$ it fails for $x - y$ large.

Proposition 2.4. *Suppose that $F: [a, b] \times \mathbb{R}^n \rightarrow \mathbb{R}^n$ is continuous and uniformly Lipschitz with Lipschitz constant M and $y_0 \in \mathbb{R}^n$, and define $\Phi: C([a, b]; \mathbb{R}^n) \rightarrow C([a, b]; \mathbb{R}^n)$ by*

$$(\Phi(y))(t) = y_0 + \int_a^t F(s, y(s)) ds.$$

Then, for each $y, z \in C([a, b]; \mathbb{R}^n)$ we have

$$d(\Phi(y), \Phi(z)) \leq M(b - a)d(y, z).$$

Proof. By definition, we see that

$$d(\Phi(y), \Phi(z)) = \max_{t \in [a, b]} \left| \int_a^t (F(s, y(s)) - F(s, z(s))) ds \right|.$$

So by the uniform Lipschitz assumption, we get

$$d(\Phi(y), \Phi(z)) \leq \max_{t \in [a, b]} \int_a^t M|y(s) - z(s)| ds \leq \max_{t \in [a, b]} \int_a^t M d(f, g) ds \leq M(b - a)d(f, g).$$

□

Corollary 2.5. *Let $F: [a, b] \times \mathbb{R}^n \rightarrow \mathbb{R}^n$ be continuous and uniformly Lipschitz with Lipschitz constant M satisfying $b - a < 1/M$. Then for any $y_0 \in \mathbb{R}^n$, the initial value problem*

$$y'(t) = F(t, y(t)) \quad y(a) = y_0$$

has exactly one solution.

Proof. Proposition 2.4 gives $d(\Phi(y), \Phi(z)) \leq M(b - a)d(y, z)$, and our assumption is that $M(b - a) < 1$. Hence we can apply the contractive mapping principle to Φ , which shows that there is one and only one $y \in C([a, b]; \mathbb{R}^n)$ such that $\Phi(y) = y$.

But we observed earlier that $y \in C([a, b]; \mathbb{R}^n)$ satisfies $\Phi(y) = y$ precisely when y is a solution to the initial value problem. \square

We've now established "short-time existence and uniqueness" for the case when F is uniformly Lipschitz. There are two ways in which we could hope to improve this result: show that in fact we have long-time existence when F is uniformly Lipschitz (thus preserving the hypothesis and strengthening the conclusion); or show that we still have short-time existence even when F satisfies some weaker assumption (thus preserving the conclusion while weakening the hypothesis). We'll do both. Recall again the bad behavior of solutions to $y' = y^2$, $y(0) = 1$, in view of which we can't hope to get long-time existence for non-uniformly-Lipschitz F .

Theorem 2.6. *Suppose that $F: [a, b] \times \mathbb{R}^n \rightarrow \mathbb{R}^n$ is continuous and uniformly Lipschitz. Then there is one and only one solution $y: [a, b] \rightarrow \mathbb{R}^n$ to the initial value problem*

$$y'(t) = F(t, y(t)) \quad y(a) = y_0.$$

Thus we get existence and uniqueness on any finite-length time interval on which F stays uniformly Lipschitz, regardless of the Lipschitz constant; in fact it's not hard to extend the argument we give below to show that if $F: [a, \infty) \times \mathbb{R}^n \rightarrow \mathbb{R}^n$ is uniformly Lipschitz then there's one and only one solution $y: [a, \infty) \rightarrow \mathbb{R}^n$ to the initial value problem (you should check this).

Proof. Where M is a Lipschitz constant for F , let $c = 1/(2M)$, and let N be the unique integer with the property that $a + Nc < b \leq a + (N + 1)c$. For $n = 0, \dots, N - 1$, let $I_n = [a + nc, a + (n + 1)c]$, and let $I_N = [a + Nc, b]$; thus for each $i \in \{0, \dots, N\}$ I_i has length strictly less than $1/M$, and the left endpoint of I_i is the right endpoint of I_{i-1} ; also the union of the I_i 's is precisely $[a, b]$.

For $i \in \{0, \dots, N\}$ define $x_i: I_i \rightarrow \mathbb{R}^n$ inductively on i as follows. First, by the Corollary, there is a unique continuous function $x_0: I_0 \rightarrow \mathbb{R}^n$ such that $x_0'(t) = F(t, x_0(t))$ and $x_0(a) = y_0$. Next, assuming $i \leq N$ and x_0, \dots, x_{i-1} have already been defined (in particular, we have $x_{i-1}: [a + (i-1)c, a + ic] \rightarrow \mathbb{R}^n$), let $x_i: I_i \rightarrow \mathbb{R}^n$ be the unique function (using the corollary, since the length of I_i is less than $1/M$) satisfying $x_i'(t) = F(t, x_i(t))$ and $x_i(a + ic) = x_{i-1}(a + ic)$.

Now define $y: [a, b] \rightarrow \mathbb{R}^n$ by setting $y(t) = x_i(t)$ where i is such that $t \in I_i$. Since the only points in $[a, b]$ belonging to more than one I_i are those of form $a + ic$ (which belong only to I_i and I_{i-1}), and since we've arranged for x_{i-1} and x_i to coincide at $a + ic$, this definition makes sense, and y is continuous since the x_i are. Moreover, if $t \in [a, b]$, letting i be such that $t \in I_i$ we have

$$y'(t) = x_i'(t) = F(t, x_i(t)) = F(t, y(t)).$$

So since $y(a) = x_0(a) = y_0$, y is indeed a solution to the initial value problem.

The fact that y is the unique such solution follows fairly straightforwardly from the uniqueness of the x_i . To give a careful proof, we shall prove by induction on j that any solution $z: \cup_{i=0}^j I_i \rightarrow \mathbb{R}^n$ to the initial value problem necessarily coincides with y ; the $j = N$ version of this statement suffices for what we need. For $j = 0$, this is just the statement that $x_0: I_0 \rightarrow \mathbb{R}^n$ is the unique solution on I_0 , noted earlier. Assuming that the claim holds for $j - 1$, it in particular follows that any such solution z satisfies $z(a + jc) = y(a + jc) = x_{j-1}(a + jc)$. But then $z|_{I_j}$ is a solution to the initial value problem $z'(t) = F(t, z(t))$, $z(a + jc) = x_{j-1}(a + jc)$. But x_j is the unique solution to this initial value problem, so z coincides with x_j

on I_j . So since $y|_{I_j} = x_j$, z coincides with y on I_j . This completes the induction, and hence also the proof of uniqueness. \square

The situation for uniformly Lipschitz F is thus quite satisfactory. However, there are many fairly straightforward and well-behaved functions (such as $F(t, y) = y^2$) which are not uniformly Lipschitz. To deal with some such functions, we make some small modifications to our earlier approach. For any $y_0 \in \mathbb{R}^n$ and $N > 0$, define

$$C([a, b]; \mathbb{R}^n, y_0, N) = \{f \in C([a, b]; \mathbb{R}^n) \mid (\forall t \in [a, b]) (|f(t) - y_0| \leq N)\}.$$

Now if $\{f_n\}_{n=1}^\infty$ is a sequence of continuous functions from $[a, b] \rightarrow \mathbb{R}^n$ such that for all t and n we have $|f(t) - y_0| \leq N$, and if $f \in C([a, b]; \mathbb{R}^n)$ is a continuous function such that $f_n \rightarrow f$ uniformly (or indeed even pointwise, but the uniform case is what's relevant here), then it is also the case that $|f(t) - y_0| \leq N$ for all t . This is another way of saying that $C([a, b]; \mathbb{R}^n, y_0, N)$ is a *closed* subset of $C([a, b]; \mathbb{R}^n)$. Hence, by a homework problem, $C([a, b]; \mathbb{R}^n, y_0, N)$ is a complete metric space (with respect to the subspace metric $d(y, z) = \max_t |y(t) - z(t)|$). So we can apply the contractive mapping principle to this smaller metric space as well.

Definition 2.7. Let $y_0 \in \mathbb{R}^n$. A function $F: [a, b] \times \mathbb{R}^n \rightarrow \mathbb{R}^n$ is called *locally Lipschitz near y_0* if there are $M, N > 0$ such that, for all $t \in [a, b]$, we have

$$|F(t, y) - F(t, z)| \leq M|y - z| \text{ whenever } |y - y_0| \leq N \text{ and } |z - y_0| \leq N.$$

It's not difficult to verify (using the mean value theorem) that if F is continuously differentiable on some neighborhood of y_0 then F is locally Lipschitz near y_0 . Thus $F(t, y) = y^\alpha$ is locally Lipschitz near any point for any $\alpha \geq 1$. However, if $\alpha < 1$ this function fails to be locally Lipschitz near 0.

We'll use a function just like our earlier Φ ; however, we need to vary the interval that we work over. So if $0 < \delta < b - a$ define

$$\Phi^\delta: C([a, a + \delta]; \mathbb{R}^n) \rightarrow C([a, a + \delta]; \mathbb{R}^n) \text{ by } (\Phi^\delta(y))(t) = y_0 + \int_a^t F(s, y(s)) ds.$$

Theorem 2.8. Assume that F is continuous and locally Lipschitz near y_0 , with constants M, N as in the definition above. Then there is $\epsilon > 0$ such that whenever $\delta \leq \epsilon$ we have

$$\Phi^\delta(C([a, a + \delta]; \mathbb{R}^n, y_0, N)) \subset C([a, a + \delta]; \mathbb{R}^n, y_0, N),$$

and, for some $r < 1$, if $y, z \in C([a, a + \delta]; \mathbb{R}^n, y_0, N)$ then

$$d(\Phi^\delta(y), \Phi^\delta(z)) \leq rd(y, z).$$

Proof. Since F is continuous, and since $[a, b] \times \{x : |x - y_0| \leq N\}$ is compact, there is M_0 such that $|F(t, x)| \leq M_0$ whenever $a \leq t \leq b$, $|x - y_0| \leq N$. Then if $y \in C([a, a + \delta]; \mathbb{R}^n, y_0, N)$, for each $t \in [a, a + \delta]$ we have

$$|(\Phi^\delta(y))(t) - y_0| \leq \left| \int_a^t F(s, y(s)) ds \right| \leq \delta M_0.$$

So if $\epsilon \leq N/M_0$ the first condition in the theorem will be satisfied whenever $\delta \leq \epsilon$.

Further, for $y, z \in C([a, a + \delta]; \mathbb{R}^n, y_0, N)$ we have

$$\begin{aligned} d(\Phi^\delta(y), \Phi^\delta(z)) &\leq \max_{a \leq t \leq a + \delta} \int_a^t |F(s, y(s)) - F(s, z(s))| ds \\ &\leq \int_a^{a + \delta} M |y(s) - z(s)| ds \leq \delta M d(y, z) \end{aligned}$$

Thus if ϵ is also chosen smaller than $1/M$ the second condition is also satisfied whenever $\delta \leq \epsilon$. \square

Theorem 2.9. *Let $F: [a, b] \times \mathbb{R}^n \rightarrow \mathbb{R}^n$ be continuous and let $y_0 \in \mathbb{R}^n$. Assume that F is uniformly Lipschitz near y_0 . Then if $a \in \mathbb{R}$ there is $\epsilon > 0$ such that, for each $\delta \leq \epsilon$, there is one and only one solution $y: [a, a + \delta] \rightarrow \mathbb{R}^n$ to the initial value problem*

$$y'(t) = F(t, y(t)) \quad y(a) = y_0.$$

Proof. If ϵ and N are as in Theorem 2.8 and $\delta \leq \epsilon$, then the unique fixed point of Φ^δ on $C([a, a + \delta]; \mathbb{R}^n, y_0, N)$, as given by the contractive mapping principle, is a solution to the initial value problem on $[a, a + \delta]$. This will be the only solution in $C([a, a + \delta]; \mathbb{R}^n, y_0, N)$, but it doesn't quite prove uniqueness because we don't yet know that any solution to the initial value problem in fact belongs to $C([a, a + \delta]; \mathbb{R}^n, y_0, N)$. To deal with this, first note that since our solution y is continuous and satisfies $y(a) = y_0$, after decreasing ϵ we can assume that $|y(t) - y_0| < N$ for all $t \in [a, a + \epsilon]$ (before this the inequality would only have been $|y(t) - y_0| \leq N$). y is still the unique solution in $C([a, a + \delta]; \mathbb{R}^n, y_0, N)$ for any $\delta \leq \epsilon$. Suppose $z: [a, a + \delta] \rightarrow \mathbb{R}^n$ were some other solution (only assumed to be in $C([a, b]; \mathbb{R}^n)$) to the problem, with $\delta \leq \epsilon$. Let

$$\gamma_0 = \sup\{\gamma | z(a + t) = y(a + t) \text{ for all } t \in [0, \gamma]\}.$$

So $\gamma_0 \geq 0$. Note since there are $t_n \in [0, \gamma_0)$ such that $t_n \rightarrow \gamma_0$ and $z(a + t_n) = y(a + t_n)$, the continuity of y and z shows that $y(a + \gamma_0) = z(a + \gamma_0)$. We'll be done if we show that $\gamma_0 = \delta$. If not, then since $|z(a + \gamma_0) - y_0| < N$, the continuity of z shows that there is some $\eta \in (0, \delta - \gamma_0)$ such that $|z(a + \gamma_0 + t) - y_0| < N$ for all $t \leq \eta$. But then $z \in C([a, a + \gamma_0 + \eta]; \mathbb{R}^n, y_0, N)$. But (since $\gamma_0 + \eta \leq \epsilon$) we've established that y is the unique solution to the initial value problem among functions in $C([a, a + \gamma_0 + \eta]; \mathbb{R}^n, y_0, N)$, so we must have $y(t) = z(t)$ for all $t \leq a + \gamma_0 + \eta$. But this contradicts the definition of γ_0 . This contradiction shows that $\gamma_0 = \delta$, which completes the proof. \square

As seen earlier, ϵ may be quite small, as in the case of the initial value problem $y' = y^{m+1}$, $y(0) = 1$ for large m (which has solution $y(t) = (1 - mt)^{-1/m}$, so ϵ in this case must be taken smaller than $1/m$). Of course, the problem in this case is that the solution shoots off to infinity in finite time, and it's possible to show that, if a solution y exists on the interval $[a, a + T)$ but not on $[a, a + T]$, then it necessarily holds that for some sequence $t_i \nearrow T$ we have $|y(t_i)| \rightarrow \infty$. Thus these sorts of divergences are the only issues that prevent us from going from short-time existence to long-time existence.

Incidentally, in case the function F is only continuous, but not locally Lipschitz, it turns out that short-time existence still holds for the initial value problem; however, as we saw in the case of the equation $y' = y^\alpha$, $y(0) = 0$ with

$0 < \alpha < 1$, uniqueness may fail to hold. Proving existence requires in this context involves quite different methods, which relate to the compactness of certain subsets of $C([a, b]; \mathbb{R}^n)$.

3. MEASURE THEORY

Our goal at this point will be to develop a theory of integration for functions on \mathbb{R} (and for that matter on other spaces) which will be somewhat more general than the Riemann theory of integration that you've learned about in calculus. Recall that to define the Riemann integral, one starts from the basic notion that the area of a rectangle is its width times its height, and then, to integrate any given function, one chops the domain into very small subintervals and approximates the region under the graph of the function by a collection of rectangles whose bases are the subintervals of the domain. The new, Lebesgue, theory of integration, will instead operate by chopping up the *range* into small intervals, and for any value c in the range, looking at the (potentially rather complicated) subset of the domain consisting of points which are mapped to a value near c . Whereas the Riemann theory starts from the formula for the area of a rectangle, the Lebesgue theory starts from the idea that if, for some set A , we define the *characteristic function* of A by the formula

$$\chi_A(x) = \begin{cases} 1 & x \in A \\ 0 & x \notin A \end{cases},$$

the integral of the characteristic function should be given by

$$\int \chi_A = m(A)$$

where the $m(A)$ on the right hand side is to be thought of as the “size of A .” In particular, if A is an interval then $m(A)$ should be its length. Moreover, just as in the Riemann theory, the integral should be linear (in other words, $\int (af + bg) = a \int f + b \int g$). These two prescriptions specify what the integral of any linear combination of characteristic functions should be, and to define the integral of a general function we (try to) approximate the given function by a linear combination of characteristic functions. (By contrast, the Riemann integral approximates the function by a step function, which is a special kind of linear combination of characteristic functions.)

It's reasonable to ask why anyone would want to bother doing this, given that one can do a great deal with the Riemann theory that you already know about. Here are three answers to that question:

- (i) Although any continuous function on a compact interval can be integrated using Riemann's theory, there are many functions for which his theory fails to give a meaningful answer; for instance this is true for the Dirichlet “salt-and-pepper” function defined by $d(x) = 1$ if x is rational and $d(x) = 0$ otherwise. The Lebesgue theory works fine for the salt-and-pepper function, and for basically any function that you could foreseeably encounter or construct (there are examples of “non-measurable” functions for which Lebesgue's theory fails, but to prove their existence one needs to use the axiom of choice; they can't be explicitly constructed).
- (ii) Even if you only care about continuous functions, which can be handled in the Riemann theory, using Lebesgue's theory one can do a much better

job of answering the following sort of question: If $\{f_n\}_{n=1}^{\infty}$ is a sequence of continuous functions which converge (in some sense) to some function f , is it the case that $\int f_n \rightarrow \int f$? This is a question that comes up very naturally in many contexts, and Lebesgue used his integration theory to discover results (such as the Dominated Convergence Theorem) in this direction that were not previously known, even though around half a century had elapsed between Riemann's work and his.

- (iii) The modern, rigorous study of the theory of probability, whose development started with Kolmogorov's work around 1930, is based in a very fundamental way on Lebesgue's theory of measure and integration (unlike the first two, this is a development that Lebesgue himself didn't anticipate).

As the description of the strategy above should suggest, before we can introduce the Lebesgue integral we need to figure out how to make sense of the size $m(A)$ of a (possibly rather complicated) subset $A \subset \mathbb{R}$. This will take some work, which we now commence.

3.1. Lebesgue (outer) measure on \mathbb{R} .

Notation 3.1. *Intervals will frequently be appearing in this subsection, and it will often be the case that I'll want to say something about an interval with certain endpoints, irrespective of whether it is closed, open, or half-open. Accordingly, I will use the notation $\langle a, b \rangle$ to denote any of the intervals (a, b) , $[a, b]$, $(a, b]$, or $[a, b)$.*

Furthermore, the length of such an interval will be denoted by

$$l(\langle a, b \rangle) = b - a.$$

We'll allow $b = \infty$ and/or $a = -\infty$, in either of which case we have $l(\langle a, b \rangle) = \infty$. We also allow $a = b$, in which case the interval is the singleton $\{a\}$ if the interval is closed and the empty set otherwise, and in our convention has length zero.

We wish to assign to each set $A \subset \mathbb{R}$ a "measure" $m(A)$ with the following properties; I'll leave it to you to convince yourself that these are natural properties that the size of a set should satisfy:

- (i) If $A = \langle a, b \rangle$ is an interval (and so has a length $l(A) = b - a$), we have $m(A) = l(A)$.
- (ii) If $A \subset B$ then $m(A) \leq m(B)$ ("monotonicity").
- (iii) If $A \cap B = \emptyset$, then $m(A \cup B) = m(A) + m(B)$ ("finite additivity").
- (iv) For *any* A, B , $m(A \cup B) \leq m(A) + m(B)$ ("finite subadditivity").

I should caution that we won't completely succeed in doing this; what we'll do is arrange for (i), (ii), (iv) to hold, and for (iii) to hold as long as A and B are taken from a certain collection of subsets of \mathbb{R} (the "measurable" subsets); this collection is quite large and contains every set that you're likely to encounter in practice, but the axiom of choice can be used to "construct" sets which aren't in the collection.

(iii) and (iv) can immediately be extended (by induction on n) to the statement that $m(\cup_{i=1}^n A_i) \leq \sum_{i=1}^n m(A_i)$, with equality holding if the A_i are disjoint. Can this additivity property be extended to infinite unions? Certainly not all, since we have $[0, 1] = \cup_{t \in [0, 1]} [t, t]$, and property (i) tells us that $l([0, 1]) = 1$ whereas each $l([t, t]) = 0$, so that

$$\sum_{t \in [0, 1]} l([t, t]) = 0 < l([0, 1]),$$

violating any hypothetical “infinite additivity” (or even “infinite subadditivity”) property. However, it turns out that requiring that additivity hold for *countably infinite* unions is not unreasonable; justification will be provided for this in the next theorem. So we provisionally add the following desired conditions for our not-yet-constructed m :

- (v) For any sequence $\{A_n\}_{n=1}^{\infty}$ of subsets of \mathbb{R} , we have $m(\cup_{n=1}^{\infty} A_n) \leq \sum_{n=1}^{\infty} m(A_n)$ (“ σ -subadditivity”)
- (vi) For any sequence $\{A_n\}_{n=1}^{\infty}$ of *disjoint* subsets of \mathbb{R} , we have $m(\cup_{n=1}^{\infty} A_n) = \sum_{n=1}^{\infty} m(A_n)$ (“ σ -additivity”)

Again, we’ll end up getting (v) to hold for all subsets of \mathbb{R} , and (vi) to hold for subsets that are chosen from the collection of measurable sets alluded to earlier.

The following result both shows that (v) and (vi) are not inconsistent with (i), and suggests a formula for m .

Theorem 3.2. *Let $\{I_n\}_{n=1}^{\infty}$ be any sequence of intervals in \mathbb{R} , and I any other interval.*

- (i) *If $I \subset \cup_{n=1}^{\infty} I_n$, then*

$$l(I) \leq \sum_{n=1}^{\infty} l(I_n).$$

- (ii) *If the I_n are disjoint, and if $I = \cup_{n=1}^{\infty} I_n$, then*

$$l(I) = \sum_{n=1}^{\infty} l(I_n).$$

Proof. Let us assume first that $l(I)$ is finite.

The proof in this case splits into two parts: Proving (i), and proving that the inequality “ \geq ” holds in (ii). The latter of these is easier, so we’ll do it first. So we assume that $I = \cup_{n=1}^{\infty} I_n$, with the intervals I_n disjoint. Now to show that $l(I) \geq \sum_{n=1}^{\infty} l(I_n)$, it suffices to show that, for any integer N , we have $l(I) \geq \sum_{n=1}^N l(I_n)$ (since, by definition, $\sum_{n=1}^{\infty} l(I_n) = \lim_{N \rightarrow \infty} \sum_{n=1}^N l(I_n)$). So we just need to consider finitely many intervals I_1, \dots, I_N which are *disjoint* and which satisfy $I_1 \cup \dots \cup I_N \subset I$. The advantage of working with just finitely many intervals is that now they can be put in order: if (recalling our notational convention) $I_n = \langle a_n, b_n \rangle$, we can choose the distinct numbers $k_1, \dots, k_N \in \{1, \dots, N\}$ so that $a_{k_1} \leq a_{k_2} \leq \dots \leq a_{k_N}$. Now the I_n are disjoint, which implies that, for any j , we have

$$b_{k_j} \leq a_{k_{j+1}}$$

(since otherwise I_{k_j} and $I_{k_{j+1}}$ would intersect). In particular, the largest of the b_n must be b_{k_N} (since it’s the only b_n which can be larger than a_{k_N}). In light of this, the fact that $\cup_{n=1}^N I_n \subset I$ forces $b \geq b_{k_N}$, and $a \leq a_{k_1}$. So we have

$$a \leq a_{k_1} \leq b_{k_1} \leq a_{k_2} \leq b_{k_2} \leq \dots \leq a_{k_N} \leq b_{k_N} \leq b.$$

Hence

$$\begin{aligned} \sum_{n=1}^N l(I_n) &= \sum_{j=1}^N (b_{k_j} - a_{k_j}) = (b_{k_N} - a_{k_1}) + \sum_{j=1}^{N-1} (b_{k_j} - a_{k_{j+1}}) \\ &\leq b_{k_N} - a_{k_1} \leq b - a = l(I), \end{aligned}$$

as desired.

We now prove (i). To do so, it suffices to show that if $\epsilon > 0$ is given then

$$\sum_{n=1}^{\infty} l(I_n) \leq l(I) + 4\epsilon$$

(for we can then take the limit as $\epsilon \rightarrow 0$ to deduce the result, as both sides of (i) are independent of ϵ). Let $\epsilon > 0$. If $I = \langle a, b \rangle$ and $I_n = \langle a_n, b_n \rangle$, denote

$$I' = [a + \epsilon, b - \epsilon], \quad I'_n = (a_n - \epsilon 2^{-n}, b_n + \epsilon 2^{-n});$$

we then have

$$I' \subset I \subset \cup_{n=1}^{\infty} I_n \subset I'_n,$$

with the I'_n now *open* and I' *compact* (since at the outset we assumed I had finite length). I' thus satisfies the Heine-Borel property, so for some N it must be the case that

$$I' \subset I'_1 \cup \dots \cup I'_N.$$

It should be intuitively clear that if one interval J is covered by finitely many other intervals J_1, \dots, J_M then $l(J) \leq \sum_{m=1}^M l(J_m)$; to prove this carefully we can use induction on M . If $M = 1$, then the left endpoint of J_1 is no larger than that of J , and the right endpoint of J_1 is no smaller than that of J , so it immediately follows that $l(J_1) \geq l(J)$. Assume the result proven for all covers by at most $M - 1$ intervals, and suppose that $J \subset J_1 \cup \dots \cup J_M$. If we in fact have $J \subset J_1 \cup \dots \cup J_{M-1}$ then by induction $l(J) \leq \sum_{m=1}^{M-1} l(J_m) \leq \sum_{m=1}^M l(J_m)$ and the result is proven. So assume that J is not contained in $J_1 \cup \dots \cup J_{M-1}$. Write $J_M = \langle c, d \rangle$. Then since $J \subset J_1 \cup \dots \cup J_M$, the interval $J_M = \langle c, d \rangle$ is not contained in any of the other J_m . Write $I_- = \{t \in J | t < c\}$, $I_0 = \{t \in J | c < t < d\}$, and $I_+ = \{t \in J | t > d\}$. I_- , I_0 , and I_+ are each intervals, and we have

$$l(J) = l(I_-) + l(I_0) + l(I_+) \leq l(I_-) + l(I_+) + l(J_M).$$

Write $S_- = \{m | 1 \leq m \leq M - 1, J_m \cap I_- \neq \emptyset\}$ and $S_+ = \{m | 1 \leq m \leq M - 1, J_m \cap I_+ \neq \emptyset\}$. If $m \in S_- \cap S_+$ then J_m would be an interval containing points both less than c and greater than d , so J_m would contain J_M , which we know not to be the case. Hence S_- and S_+ are *disjoint* sets, each with at most $M - 1$ elements; by construction, we have

$$I_- \subset \cup_{m \in S_-} J_m \quad I_+ \subset \cup_{m \in S_+} J_m.$$

So by the inductive hypothesis

$$l(I_-) \leq \sum_{m \in S_-} l(J_m), \quad l(I_+) \leq \sum_{m \in S_+} l(J_m),$$

and so

$$l(J) \leq l(I_-) + l(I_+) + l(J_M) \leq \sum_{m=1}^M l(J_m).$$

Applying this general fact to our intervals I', I'_1, \dots, I'_N , we get

$$l(I') \leq \sum_{n=1}^N l(I'_n).$$

Now, by construction,

$$l(I) = l(I') + 2\epsilon, \quad l(I_n) = l(I'_n) - 2\epsilon 2^{-n}.$$

So

$$l(I) \leq 2\epsilon + \sum_{n=1}^N (l(I_n) + 2\epsilon 2^{-n}) \leq \sum_{n=1}^N l(I'_n) + 4\epsilon \leq \sum_{n=1}^{\infty} l(I_n) + 4\epsilon,$$

from which (i) follows by taking the limit as ϵ tends to zero. This completes the proof of the theorem in case $l(I) < \infty$.

For the case that $l(I) = \infty$, we just need to show that if some collection of intervals $\{I_n\}_{n=1}^{\infty}$ covers I then $\sum_{n=1}^{\infty} l(I_n) = \infty$. Now that $l(I) = \infty$ implies that, for any $N > 0$, there is M such that $l(I \cap [-M, M]) > N$. In this case, the sets $\{I_n \cap [-M, M]\}$ (which are intervals) cover $I \cap [-M, M]$, and since $I \cap [-M, M]$ has finite length we can apply what we've already done to deduce that

$$\sum_{n=1}^{\infty} l(I_n) \geq \sum_{n=1}^{\infty} l(I_n \cap [-M, M]) \geq l(I \cap [-M, M]) \geq N.$$

That this holds for every real N then implies that

$$\sum_{n=1}^{\infty} l(I_n) = \infty,$$

as desired. □

The penultimate paragraph of the proof employed a useful trick, which exploits the fact that $\sum_{n=1}^{\infty} 2^{-n} = 1$ (and hence $\sum_{n=1}^{\infty} \epsilon 2^{-n} = \epsilon$) to show that, essentially, if one has a countable collection of errors that can be made arbitrarily small, then the infinite sum of these errors can be made arbitrarily small. This will come in handy often.

Corollary 3.3. *If I is an interval, we have*

$$l(I) = \inf\left\{\sum_{n=1}^{\infty} l(I_n) \mid I \subset \cup_{n=1}^{\infty} I_n \text{ and each } I_n \text{ is an interval}\right\}.$$

Proof. By choosing $I_1 = I$ and $I_n = \emptyset$ for $n \geq 2$ (recall $\emptyset = (0, 0)$ is an interval in our convention) we see that “ \geq ” holds, while “ \leq ” follows immediately from part (i) of Theorem 3.2. □

This serves as partial justification for the following definition:

Definition 3.4. If $A \subset \mathbb{R}$ is any subset, the *Lebesgue outer measure* of A is

$$m^*(A) = \inf\left\{\sum_{n=1}^{\infty} l(I_n) \mid A \subset \cup_{n=1}^{\infty} I_n \text{ and each } I_n \text{ is an interval}\right\}.$$

Here are some properties of m^* :

- (1) If I is an interval, then $m^*(I) = l(I)$ (this is Corollary 3.3).
- (2) If $A \subset B$, then $m^*(A) \leq m^*(B)$. (Indeed, if $\{I_n\}$ is a collection of intervals covering B , then the collection $\{I_n\}$ also covers A , so by definition we have $m^*(A) = \inf S_A$, $m^*(B) = \inf S_B$ for certain sets S_A, S_B with $S_B \subset S_A$. But enlarging a set can only make its infimum smaller).

(3) $m^*(\cup_{n=1}^{\infty} A_n) \leq \sum_{n=1}^{\infty} m^*(A_n)$. *Proof:* Let $\epsilon > 0$. Then for each n there is a sequence $\{I_{n,k}\}_{k=1}^{\infty}$ such that $A_n \subset \cup_{k=1}^{\infty} I_{n,k}$ and $\sum_{k=1}^{\infty} l(I_{n,k}) \leq m^*(A_n) + \epsilon 2^{-n}$. But then $\cup_{n=1}^{\infty} A_n \subset \cup_{n=1}^{\infty} \cup_{k=1}^{\infty} I_{n,k}$, this being a countable union of intervals. Hence

$$\sum_{n=1}^{\infty} m^*(A_n) \leq \sum_{n=1}^{\infty} \sum_{k=1}^{\infty} l(I_{n,k}) \leq \sum_{n=1}^{\infty} (m^*(A_n) + \epsilon 2^{-n}) = \sum_{n=1}^{\infty} m^*(A_n) + \epsilon.$$

For this to hold for every $\epsilon > 0$ it must indeed be the case that

$$m^*(\cup_{n=1}^{\infty} A_n) \leq \sum_{n=1}^{\infty} m^*(A_n).$$

Now is as good a time as any to introduce the following definition:

Definition 3.5. If X is any set, a function $\mu^*: \{\text{subsets of } X\} \rightarrow [0, \infty]$ is called an *outer measure* on X if

- $\mu^*(\emptyset) = 0$;
- If $A \subset B \subset X$ then $\mu^*(A) \leq \mu^*(B)$; and
- For any sequence $\{A_n\}_{n=1}^{\infty}$ of subsets of X we have

$$\mu^*(\cup_{n=1}^{\infty} A_n) \leq \sum_{n=1}^{\infty} \mu^*(A_n).$$

Note that we allow μ^* to take the value ∞ . Addition in the “interval” $[0, \infty]$ is defined in the obvious way: if $a, b \in \mathbb{R}$ then $a + b$ is its usual value in \mathbb{R} , while $a + \infty = \infty + a = \infty$. Note also that, for any sequence $\{a_n\}_{n=1}^{\infty}$ of elements of $[0, \infty]$, the sum $\sum_{n=1}^{\infty} a_n$ is a well-defined value in $[0, \infty]$: it’s ∞ if either some $a_n = \infty$ or if the a_n are real and the sum $\sum_{n=1}^{\infty} a_n$ diverges (in which case it diverges to ∞ (*i.e.*, for any real M the partial sums $\sum_{n=1}^N a_n$ are all eventually larger than M) since the a_n are all nonnegative), while in any other case $\sum_{n=1}^{\infty} a_n$ converges to some real number. This would no longer be the case if we allowed some of the a_n to be negative. Of course $[0, \infty]$ lacks the metric space structure of $[0, \infty)$, since distances in a metric space can’t be infinite.

At any rate, an immediate consequence of the properties above is:

Corollary 3.6. *Lebesgue outer measure is an outer measure on \mathbb{R} .*

There are other outer measures; the simplest example is the “counting outer measure,” which can be defined on any set X and whose value on a set A is the number of elements of A if this number is finite and ∞ otherwise.

Let’s consider two examples of sets of Lebesgue outer measure zero:

Example 3.7. Any countable set $S \subset \mathbb{R}$ has Lebesgue outer measure zero. Indeed, that S is countable means that $S = \{r_n | n \in \mathbb{N}\}$ (where \mathbb{N} is the set of natural numbers) for some $r_n \in \mathbb{R}$, so we have

$$m^*(S) = m^*(\cup_{n=0}^{\infty} \{r_n\}) \leq \sum_{n=0}^{\infty} m^*(\{r_n\}) = \sum_{n=0}^{\infty} m^*([r_n, r_n]) = \sum_{n=0}^{\infty} 0 = 0.$$

This gives a new proof of the following familiar fact from set theory.

Corollary 3.8. *Any interval in \mathbb{R} with positive length is uncountable.*

Proof. Indeed, any such interval has positive Lebesgue outer measure, so since a countable set has zero Lebesgue outer measure the interval must not be countable. \square

Example 3.9 (The Cantor set). The Cantor set C is given by $C = \bigcap_{n=0}^{\infty} C_n$ where the C_n are constructed as follows. First, $C_0 = [0, 1]$. Next, inductively assume that we have constructed C_n as a union $C_n = \bigcup_{k=1}^{2^n} I_{n,k}$ of 2^n closed intervals, each having length 3^{-n} (say $I_{n,k} = [a_{n,k}, a_{n,k} + 3^{-n}]$) (of course this is true for $n = 0$, with $a_{0,0} = 0$). Then set

$$I_{n+1,2k-1} = [a_{n,k}, a_{n,k} + 3^{-n-1}], \quad I_{n+1,2k} = [a_{n,k} + 2 \cdot 3^{-n}, a_{n,k} + 3^{-n}]$$

and

$$C_{n+1} = \bigcup_{k=1}^{2^{n+1}} I_{n+1,k}.$$

Thus the intervals of C_{n+1} are obtained by deleting the open middle thirds of the intervals of C_n . In particular $C_{n+1} \subset C_n$.

Since the C_n are all compact and nonempty and satisfy $C_{n+1} \subset C_n$, $C = \bigcap_{n=0}^{\infty} C_n$ is nonempty by one of your homework problems. In fact, C is uncountable. To sketch a proof of this, recall that the set of sequences $\{\epsilon_n\}_{n=1}^{\infty}$ where each $\epsilon_n \in \{0, 1\}$ is uncountable (in fact, it can be put into one-to-one correspondence with $[0, 1]$). To any such sequence we can associate a sequence of intervals $I_{n,k_n} \subset C_n$ such that $I_{n+1,k_{n+1}} \subset I_{n,k_n}$ as follows: $k_0 = 0$ (so $I_{0,k_0} = I_{0,0} = [0, 1]$), and assuming inductively that we've chosen I_{n,k_n} , there are two intervals $I_{n+1,2k_n-1}, I_{n+1,2k_n}$ of C_{n+1} contained in I_{n,k_n} , and we choose $I_{n+1,2k_n-1}$ if $\epsilon_{n+1} = 0$ and $I_{n+1,2k_n}$ if ϵ_{n+1} is odd. The resulting intervals I_{n,k_n} are a decreasing sequence of nonempty compact sets, so have nonempty intersection, which is contained in C . In fact, since the lengths of the I_{n,k_n} tend to zero the intersection $\bigcap_{n=1}^{\infty} I_{n,k_n}$ is a single point, and it is this point which we associate to the sequence $\{\epsilon_n\}_{n=1}^{\infty}$ of zeros and ones. Different sequences of ϵ_n 's result in different sequences of intervals, and so since the $I_{n,k}$ are disjoint as k varies and n is fixed they result in different points of C . This gives a one-to-one correspondence between the set of countable sequences of zeros and ones and a subset of C (in fact, it's not too hard to see that the subset is all of C), proving that C is uncountable.

Now the C_n are each unions of 2^n intervals of length 3^{-n} ; hence by the subadditivity property of m^* we have $m^*(C_n) = (2/3)^n$ (in fact this is an equality). So since $C \subset C_n$, $m^*(C) \leq (2/3)^n$. But then for this to hold for all n it must be that $m^*(C) = 0$.

Thus C is an uncountable set, which nonetheless has Lebesgue outer measure zero. So even though C is as large as $[0, 1]$ from the point of view of set theory, from a measure-theoretic point of view it is essentially negligible.

If $A \subset \mathbb{R}$, $m^*(A)$ is defined by looking at all countable covers of A by intervals; it can sometimes be useful to consider more specific kinds of covers.

Lemma 3.10. *If $\delta > 0$ and $A \subset \mathbb{R}$, we have*

$$m^*(A) = \inf \left\{ \sum_{n=1}^{\infty} l(I_n) \mid A \subset \bigcup_{n=1}^{\infty} I_n \text{ and each } I_n \text{ is an open interval of length at most } \delta \right\}.$$

Proof. Since any countable cover of A by open intervals of length at most δ is a countable cover of A by intervals, the inequality " \leq " follows directly from the definition of $m^*(A)$ (and the fact that enlarging a set causes its infimum to either

stay the same or decrease). Also, the inequality “ \geq ” is trivial when $m^*(A) = \infty$, so we’ll be done if we show “ \geq ” in case $m^*(A) < \infty$. To do so, it’s enough to show that, for any $\epsilon > 0$, if $A \subset \cup_{n=1}^{\infty} J_n$ where the J_n are intervals of finite length, then there are open intervals I_n each having length at most δ such that $A \subset \cup_{n=1}^{\infty} I_n$ and

$$\sum_{n=1}^{\infty} l(I_n) \leq \sum_{n=1}^{\infty} l(J_n) + \epsilon.$$

In turn, to do this, it suffices to, for each n , find intervals $I_{n,1}, \dots, I_{n,m_n}$, each open and of length at most δ , such that $J_n \subset \cup_{k=1}^{m_n} I_{n,k}$ and

$$\sum_{k=1}^{m_n} l(I_{n,k}) \leq l(J_n) + \epsilon 2^{-n}.$$

(Indeed, if we do this, then the countable collection $\{I_{n,k} | n \geq 1, 1 \leq k \leq m_n\}$ will be the desired countable collection of open intervals of length at most δ). We now construct $I_{n,1}, \dots, I_{n,m_n}$. Assume that $J_n = \langle a_n, b_n \rangle$. Choose m_n to be the smallest integer with the property that $m_n \delta / 2 \geq b_n - a_n$. For $k = 1, \dots, m_n - 1$, set

$$I_{n,k} = \left(a_n + (i-1) \frac{\delta}{2} - \min\{\delta/4, \epsilon/(2m_n 2^n)\}, a_n + i\delta/2 + \min\{\delta/4, \epsilon/(2m_n 2^n)\} \right),$$

and set

$$I_{n,m_n} = (a_n - (m_n - 1)\delta/2, b_n + \min\{\delta/4, \epsilon/(2m_n 2^n)\}).$$

These $I_{n,k}$ have been arranged to each have length $\delta/2 + 2 \min\{\delta/4, \epsilon/(2m_n 2^n)\} \leq \delta$, and to have the property that their union contains J_n and the sum of their lengths exceeds the length of J_n by no more than $m_n (2\epsilon/2m_n 2^{-n}) = \epsilon 2^{-n}$. As explained above, this suffices to complete the proof. \square

Definition 3.11. If (X, d) is a metric space, and outer measure μ^* on X is called a *metric outer measure* if the following holds: whenever $\epsilon > 0$ and $A, B \subset X$ are such that, if $a \in A$ and $b \in B$ we have $d(a, b) \geq \epsilon$, we have

$$\mu^*(A \cup B) = \mu^*(A) + \mu^*(B).$$

Theorem 3.12. *Lebesgue outer measure m^* is a metric outer measure on \mathbb{R} with its standard metric.*

Proof. Of course, we already know that m^* is an outer measure on \mathbb{R} . Let ϵ, A, B be as in the definition of a metric outer measure. Let $\eta > 0$. By Lemma 3.10, we can write $A \cup B \subset \cup_{n=1}^{\infty} I_n$ with each I_n an open interval of length at most ϵ and

$$\sum_{n=1}^{\infty} l(I_n) \leq m^*(A \cup B) + \eta.$$

Write $N_A = \{n | I_n \cap A \neq \emptyset\}$ and $N_B = \{n | I_n \cap B \neq \emptyset\}$. Then $A \subset \cup_{n \in N_A} I_n$ and $B \subset \cup_{n \in N_B} I_n$, so $m^*(A) \leq \sum_{n \in N_A} l(I_n)$ and $m^*(B) \leq \sum_{n \in N_B} l(I_n)$. Further, if $n \in N_A$, then I_n contains a point, say a , of A , but then since I_n is an open interval of length at most ϵ all of its points are a distance less than ϵ from a ,

and so $I_n \cap B = \emptyset$ since all points of B has distance at least ϵ from a . Thus $n \in N_A \Rightarrow n \notin N_B$, proving that N_A and N_B are disjoint. So

$$m^*(A) + m^*(B) \leq \sum_{n \in N_A} l(I_n) + \sum_{n \in N_B} l(I_n) \leq \sum_{n=1}^{\infty} l(I_n) \leq m^*(A \cup B) + \eta.$$

Since this holds for all $\eta > 0$ we have $m^*(A \cup B) \geq m^*(A) + m^*(B)$. The reverse inequality is a consequence of the subadditivity property of m^* , so in fact $m^*(A \cup B) = m^*(A) + m^*(B)$. □

The goal at this point is to improve the σ -subadditivity property of m^* to additivity or σ -additivity, as in properties (iii) or (vi) at the start of this subsection. That m^* is a metric outer measure evidently amounts to the statement that it is additive on certain pairs of disjoint subsets. Later, we'll improve this to get σ -additivity on a rather large family of sets. Before this, we'll address the question of whether it's reasonable to hope for m^* (or anything like it) to be σ -additive on *all* subsets of \mathbb{R} .

Given $E \subset \mathbb{R}$ and $t \in \mathbb{R}$, we define the t -translate of E as

$$E + t = \{e + t | e \in E\}.$$

So for instance $[a, b] + t = [a + t, b + t]$. Note that I is an interval if and only if every $I + t$ is an interval, in which case we have $l(I + t) = l(I)$. As such

$$\begin{aligned} m^*(E + t) &= \inf \left\{ \sum_{n=1}^{\infty} l(I_n + t) \mid E + t \subset \cup_{n=1}^{\infty} (I_n + t), I_n \text{ are intervals} \right\} \\ &= \inf \left\{ \sum_{n=1}^{\infty} l(I_n) \mid E \subset \cup_{n=1}^{\infty} I_n, I_n \text{ are intervals} \right\} = m^*(E). \end{aligned}$$

(Thus we say that Lebesgue outer measure is *translation invariant*)

For the purposes of the theorem which follows, it's more useful to work with subsets of the half-open interval $[0, 1)$. Accordingly if $E \subset [0, 1)$ and $t \in [0, 1)$, define

$$E \oplus t = \{e + t | e \in E, 0 \leq e + t < 1\} \cup \{e + t - 1 | 1 \leq e + t < 2\}.$$

Thus $E \oplus t$ consists of E translated by t , but with the part of $E + t$ that leaves the interval $[0, 1)$ on the right reappearing on the left part of the interval $[0, 1)$ (like in Pac-Man). If $I \subset [0, 1)$ is an interval, then $I \oplus t$ is either a subinterval of $[0, 1)$ or the disjoint union of two subintervals of $[0, 1)$, with the sum of the lengths of these one or two intervals equal to $l(I)$. Given a cover of E by intervals, performing this \oplus version of translation on the intervals gives a cover of $E \oplus t$ by intervals, with the sum of the lengths unchanged, so by the same argument as in the $E + t$ case we have

$$m^*(E \oplus t) = m^*(E).$$

The following theorem is therefore bad news with respect to our goal of having m^* be σ -additive.

Theorem 3.13. *It is impossible to assign to each subset $E \subset \mathbb{R}$ a "size" $m(E)$ satisfying each of the following conditions:*

- (i) $m([0, 1)) = 1$.
- (ii) If $E \subset [0, 1)$ and $t \in [0, 1)$ then $m(E \oplus t) = m(E)$.

(iii) If $\{E_n\}_{n=1}^\infty$ is a sequence of disjoint subsets of \mathbb{R} then

$$m(\cup_{n=1}^\infty E_n) = \sum_{n=1}^\infty m(E_n).$$

Proof. Let us assume that conditions (i) and (ii) hold, and produce a family $\{E_n\}_{n=1}^\infty$ of disjoint sets for which condition (iii) fails. More specifically, we will arrange that $\cup_{n=1}^\infty E_n = [0, 1)$ and, for some set $A \subset [0, 1)$ and numbers $r_n \in [0, 1)$, we have $E_n = A \oplus r_n$. Condition (i) then gives $m(\sum_{n=1}^\infty E_n) = 1$, while by condition (ii) we have $m(E_n) = m(A)$, so

$$\sum_{n=1}^\infty m(E_n) = \sum_{n=1}^\infty m(A) = \begin{cases} \infty & m(A) \neq 0 \\ 0 & m(A) = 0 \end{cases}$$

Thus if we can produce a countably infinite family of disjoint sets E_n whose union is $[0, 1)$ and each of which is a \oplus -translate of some fixed set A , then it will follow that (i), (ii), and (iii) are mutually incompatible and the theorem will be proven.

We can obtain such E_n by appealing to the axiom of choice (there is no explicit construction, and if we worked in a model of set theory in which the axiom of choice failed, then this theorem would likely be false; however, we have generally tacitly assumed the axiom throughout this course, so it would be inconsistent for us to change our minds about this set-theoretic matter now). Given $s, t \in [0, 1)$, let us say that $s \sim t$ if and only if $s - t \in \mathbb{Q}$ (here \mathbb{Q} is the set of rational numbers). \sim is then an equivalence relation ($s \sim s$, $s \sim t \Rightarrow t \sim s$, and $(s \sim t, t \sim u) \Rightarrow s \sim u$), so \mathbb{R} is partitioned into equivalence classes (i.e., $\mathbb{R} = \cup_{\alpha \in A} S_\alpha$ where the S_α are disjoint and if $s \in S_\alpha$ then $t \in S_\alpha$ if and only if $s \sim t$.) Using the axiom of choice, we can form a set A which consists of one element from each of the equivalence classes S_α . (If you prefer the language of group theory, $[0, 1)$ is a group with respect to the operation of addition modulo 1, with normal subgroup given by $\mathbb{Q} \cap [0, 1)$; the equivalence classes are the cosets of $\mathbb{Q} \cap [0, 1)$ in $[0, 1)$, and our set A consists of precisely one element from each coset).

Since $\mathbb{Q} \cap [0, 1)$ is countable, we can enumerate it as $\mathbb{Q} \cap [0, 1) = \{r_n | n \geq 1\}$. Let $E_n = A \oplus r_n$. If $s \in [0, 1)$, let α be such that $s \in S_\alpha$, and let a_α be the unique element of $A \cap S_\alpha$. Since s and a_α belong to the same equivalence class, they differ by a rational number. If $s \geq a_\alpha$, we have $s = a_\alpha + r_n \in [0, 1)$ for some n , while if $s < a_\alpha$ we have $s = a_\alpha + r_n - 1$ for some r_n with $a_\alpha + r_n \in [1, 2)$. So in any event $s \in A \oplus r_n$. By definition $A \oplus r_n \subset [0, 1)$, so we've shown that

$$[0, 1) = \cup_{n=1}^\infty (A \oplus r_n).$$

It remains only to show that the $E_n = A \oplus r_n$ are pairwise disjoint. Indeed, if $x \in (A \oplus r_m) \cap (A \oplus r_n)$, then we have $a_1 \oplus r_m = a_2 \oplus r_n$ for some $a_1, a_2 \in A$, so a_1, a_2 differ by a rational number, and so belong to the same equivalence class (here for $s, t \in [0, 1)$ $s \oplus t$ is whichever of $s + t$ or $s + t - 1$ belongs to $[0, 1)$). But A only contains one element from each equivalence class, so $a_1 = a_2$ and we have $a_1 \oplus r_m = a_1 \oplus r_n$. But from the definition of the modular addition operation \oplus it immediately follows that $r_m = r_n$. Thus for $m \neq n$ the sets $A \oplus r_m$ and $A \oplus r_n$ are disjoint.

So we've written $[0, 1)$ as a countable disjoint union of sets $E_n = A \oplus r_n$ with equal measure; as explained earlier this completes the proof, since the sum of the

measures of an infinite collection of sets with equal measure is either zero (if their common measure is zero) or infinity (if their common measure is nonzero). \square

In particular, since Lebesgue outer measure m^* satisfies (i) and (ii), it can't satisfy (iii). This is moderately disappointing, but since any reasonable definition of the size of a set really should satisfy (i) and (ii), this also suggests that the failure of (iii) isn't really m^* 's fault; any other candidate that we might have would have the same problem. What we'll ultimately do is show that m^* does satisfy (iii) if we restrict attention to a suitable collection of subsets, which is quite large and contains any set that you're likely to encounter in practice (indeed, basically any set that can be explicitly constructed; note that in the proof above we needed the axiom of choice to make A , so that there's no way of figuring out whether a given number belongs to A).

3.2. Measures from outer measures. The collection of subsets that we'll use is a certain example of the following type of collection:

Definition 3.14. If X is a set, a collection \mathcal{A} of subsets of X is called a σ -algebra in X provided that the following conditions hold:

- (i) $\emptyset \in \mathcal{A}$;
- (ii) If $E \in \mathcal{A}$, then $E^c \in \mathcal{A}$ (here and elsewhere E^c denotes the complement $X \setminus E$ of E in X).
- (iii) If $\{E_n\}_{n=1}^\infty$ is any sequence of sets belonging to \mathcal{A} , then $\cup_{n=1}^\infty E_n \in \mathcal{A}$.

Some immediate consequences of the definition are that $X \in \mathcal{A}$ (since $X = \emptyset^c$ and $\emptyset \in \mathcal{A}$); that if $A, B \in \mathcal{A}$ then $A \cup B \in \mathcal{A}$ (use condition (iii) with $E_1 = A, E_2 = B$, and $E_n = \emptyset$ for $n \geq 3$), $A \cap B \in \mathcal{A}$ (use condition (ii) and the fact that $A \cap B = (A^c \cup B^c)^c$), and $A \setminus B \in \mathcal{A}$ (since $A \setminus B = A \cap B^c$); and that if $\{E_n\}_{n=1}^\infty$ is a sequence of sets in \mathcal{A} then $\cap_{n=1}^\infty E_n \in \mathcal{A}$ (since $\cap_{n=1}^\infty E_n = (\cup_{n=1}^\infty E_n^c)^c$). Thus a σ -algebra is closed under each of the standard set-theoretic operations of union, intersection, and complement, provided that these are carried out on finite or countable collections of sets.

Two obvious examples of σ -algebras are $\mathcal{A} = \{\emptyset, X\}$ and $\mathcal{A} = \{\text{all subsets of } X\}$. While the first of these examples shows that σ -algebras can be quite small, it is typically the case that, thanks to the fact that σ -algebras are closed under the set theoretic operations mentioned in the previous paragraph, as soon as a σ -algebra contains a few subsets it automatically contains quite a lot.

A more subtle example of a σ -algebra comes from the theory of outer measures.

Definition 3.15. If μ^* is an outer measure on X , a subset $E \subset X$ is called μ^* -measurable if, for every $A \subset X$, we have

$$\mu^*(A) = \mu^*(A \cap E) + \mu^*(A \cap E^c).$$

Note that the inequality " \leq " in the above follows directly from the subadditivity property in the definition of an outer measure, so if we want to show that a set is measurable we only need to check " \geq ". Also, the inequality " \geq " is trivial in the case that $\mu^*(A) = \infty$, so in fact we only need to check that " \geq " holds for those sets A with $\mu^*(A) < \infty$.

Theorem 3.16. Let X be a set, let μ^* be an outer measure on X , and let

$$\mathcal{M} = \{E \subset X \mid E \text{ is } \mu^*\text{-measurable}\}.$$

Then

- (i) \mathcal{M} is a σ -algebra.
- (ii) If $\{E_n\}_{n=1}^{\infty}$ is a sequence of disjoint sets which belong to \mathcal{M} , then

$$\mu^*(\cup_{n=1}^{\infty} E_n) = \sum_{n=1}^{\infty} \mu^*(E_n).$$

Proof. \mathcal{M} contains the empty set, since for any $A \subset X$ we have

$$\mu^*(A \cap \emptyset) + \mu^*(A \cap \emptyset^c) = \mu^*(\emptyset) + \mu^*(A) = \mu^*(A).$$

If $E \in \mathcal{M}$, then since $(E^c)^c = E$ we have, for any $A \subset X$

$$\mu^*(A \cap E^c) + \mu^*(A \cap (E^c)^c) = \mu^*(A \cap E^c) + \mu^*(A \cap E) = \mu^*(A)$$

by the measurability of E , and so $E^c \in \mathcal{M}$. So to prove part (i) of the theorem we just need to show that \mathcal{M} is closed under countable unions.

As a first step in this direction, we'll prove that if $E, F \in \mathcal{M}$ then $E \cup F \in \mathcal{M}$. So let $A \subset X$; as mentioned earlier it's enough to restrict to the case that $\mu^*(A) < \infty$. Note that, since F is measurable, we have

$$\mu^*(A \cap E^c) = \mu^*(A \cap E^c \cap F) + \mu^*(A \cap E^c \cap F^c).$$

Now since $E \cup F = E \cup (E^c \cap F)$, the subadditivity of μ^* gives that

$$\begin{aligned} \mu^*(A \cap (E \cup F)) &\leq \mu^*(A \cap E) + \mu^*(A \cap E^c \cap F) \\ &= \mu^*(A \cap E) + (\mu^*(A \cap E^c) - \mu^*(A \cap E^c \cap F^c)). \end{aligned}$$

Now $E^c \cap F^c = (E \cup F)^c$, so this gives

$$\mu^*(A \cap (E \cup F)) + \mu^*(A \cap (E \cup F)^c) \leq \mu^*(A \cap E) + \mu^*(A \cap E^c) = \mu^*(A),$$

where the last equality comes from the measurability of E . As mentioned earlier, this inequality suffices to show that $E \cup F \in \mathcal{M}$.

Now that we know that the union of any two elements of \mathcal{M} belongs to \mathcal{M} , it follows by induction that, for any finite N , the union of any N elements of \mathcal{M} belongs to \mathcal{M} .

We're now ready to show that \mathcal{M} is closed under countable unions. Let $\{E_n\}_{n=1}^{\infty}$ be any sequence of sets in \mathcal{M} . For any n write $F_n = \cup_{k=1}^n E_k$; by what we've already shown we have $F_n \in \mathcal{M}$. Also $G_n := F_n \setminus F_{n-1} = F_n \cap F_{n-1}^c \in \mathcal{M}$. Thus the F_n form an increasing sequence of sets in \mathcal{M} , while the G_n form a sequence of disjoint sets in \mathcal{M} , and we have

$$\cup_{n=1}^{\infty} E_n = \cup_{n=1}^{\infty} F_n = \cup_{n=1}^{\infty} G_n.$$

Let $A \subset X$ be any subset. Now for any n , the fact that F_{n-1} is measurable implies that

$$\mu^*(A \cap F_n) = \mu^*(A \cap F_n \cap F_{n-1}) + \mu^*(A \cap F_n \cap F_{n-1}^c) = \mu^*(A \cap F_{n-1}) + \mu^*(A \cap G_n).$$

(The second equality follows from the definitions of the F_k and G_k .) By induction on n , it then follows that

$$\mu^*(A \cap F_n) = \sum_{k=1}^n \mu^*(A \cap G_k).$$

Meanwhile since $F_n \subset \cup_{k=1}^{\infty} E_k$, monotonicity gives that $\mu^*(A \cap F_n^c) \geq \mu^*(A \cap (\cup_{k=1}^{\infty} E_k)^c)$, and that $\mu^*(A \cap F_n) \leq \mu^*(A \cap (\cup_{k=1}^{\infty} E_k))$. Sending $n \rightarrow \infty$ in the latter relation, and using our formula for $\mu^*(A \cap F_n)$, then gives

$$(6) \quad \mu^*(A \cap (\cup_{k=1}^{\infty} E_k)) \geq \sum_{k=1}^{\infty} \mu^*(A \cap G_k).$$

So, using the measurability of F_n , we obtain

$$\mu^*(A) = \mu^*(A \cap F_n) + \mu^*(A \cap F_n^c) \geq \sum_{k=1}^n \mu^*(A \cap G_k) + \mu^*(A \cap (\cup_{k=1}^{\infty} E_k)^c)$$

for all n , and hence

$$\mu^*(A) \geq \sum_{k=1}^{\infty} \mu^*(A \cap G_k) + \mu^*(A \cap (\cup_{k=1}^{\infty} E_k)^c).$$

But the subadditivity of μ^* shows that $\mu^*(A \cap (\cup_{n=1}^{\infty} E_n)) \leq \sum_{k=1}^{\infty} \mu^*(A \cap G_k)$, (since the union of the G_k is the same as that of the E_k), so we have

$$\mu^*(A) \geq \mu^*(A \cap (\cup_{n=1}^{\infty} E_n)) + \mu^*(A \cap (\cup_{n=1}^{\infty} E_n)^c),$$

which proves that $\cup_{n=1}^{\infty} E_n \in \mathcal{M}$. This completes the proof that \mathcal{M} is a σ -algebra.

For the σ -additivity property (ii), if the E_n form a disjoint sequence from \mathcal{M} , then with the notation as above we will have $G_k = E_k$, and applying (6) with $A = X$, we get

$$\mu^*(\cup_{k=1}^{\infty} E_k) \geq \sum_{k=1}^{\infty} \mu^*(E_k).$$

The reverse inequality follows by σ -subadditivity of the outer measure μ^* , so this proves (ii). □

Thus our outer measure has the σ -additivity property that we wanted, provided that we restrict attention to its measurable subsets. For this to be useful, we need to know that there are a lot of measurable subsets; at this point we know that the measurable subsets form a σ -algebra, but in principle the σ -algebra might just be the trivial one $\{\emptyset, X\}$ (and in fact one can construct outer measures for which this is true). Recall, though, that Lebesgue outer measure satisfies the additional property of being a metric outer measure (*i.e.*, it is additive on pairs of sets which are separated from each other by a positive distance), and the following theorem implies the existence of many measurable sets for any metric outer measure:

Theorem 3.17. *Let (X, d) be a metric space, and suppose that μ^* is a metric outer measure on X . Then every closed subset of X is μ^* -measurable.*

Proof. Let $F \subset X$ be a closed set, and let A be any subset of X satisfying $\mu^*(A) < \infty$. For any positive integer n , let

$$A_n = \{x \in A \cap F^c \mid (\forall f \in F)(d(x, f) \geq 1/n)\}.$$

Now since F is closed, if $x \in F^c$ there is some n such that $B_{1/n}(x) \subset F^c$, and for this n we will have $d(x, f) \geq 1/n$ for all $f \in F$. Thus

$$A \cap F^c = \cup_{n=1}^{\infty} A_n.$$

Lemma 3.18. *We have*

$$\mu^*(A \cap F^c) = \lim_{n \rightarrow \infty} \mu^*(A_n).$$

Proof of Theorem 3.17, assuming Lemma 3.18. The fact that μ^* is a metric outer measure implies that, for any n , we have $\mu^*((A \cap F) \cup A_n) = \mu^*(A \cap F) + \mu^*(A_n)$ (since by the definition of the A_n , any point of $A \cap F$ is a distance at least $1/n$ from A_n). Now by monotonicity we have $\mu^*(A) \geq \mu^*((A \cap F) \cup A_n)$, so

$$\mu^*(A) \geq \mu^*(A \cap F) + \mu^*(A_n).$$

Then sending $n \rightarrow \infty$ and applying Lemma 3.18 shows that $\mu^*(A) \geq \mu^*(A \cap F) + \mu^*(A \cap F^c)$. A was an arbitrary subset of X with finite outer measure, so this proves that F is measurable. \square

Proof of Lemma 3.18. Let $B_n = A_{n+1} \setminus A_n$. Now if $a \in A_{n-1}$ and $b \in B_n$, since $b \notin A_n$ there must be some $f \in F$ such that $d(b, f) < 1/n$. Since $a \in A_{n-1}$, $d(a, f) \geq 1/(n-1)$. Hence the triangle inequality gives $d(a, b) \geq d(a, f) - d(b, f) \geq 1/(n-1) - 1/n$. Thus A_{n-1} and B_n are separated from each other by a positive distance (namely $1/(n-1) - 1/n$), so the fact that μ^* is a metric outer measure shows that

$$\mu^*(A_{n-1} \cup B_n) = \mu^*(A_{n-1}) + \mu^*(B_n).$$

So since $A_{n-1} \cup B_n \subset A_{n+1}$,

$$\mu^*(A_{n+1}) \geq \mu^*(A_{n-1}) + \mu^*(B_n).$$

But now we can apply the same argument with n replaced by $n-2$, and repeating this inductively (both for odd and even choices of n) shows, for any m

$$\mu^*(A_{2m+1}) \geq \sum_{k=1}^m \mu^*(B_{2k}), \quad \mu^*(A_{2m}) \geq \sum_{k=1}^m \mu^*(B_{2k-1}).$$

Now $\mu^*(A)$ is at least as large as any $\mu^*(A_n)$ since $A_n \subset A$, so it follows that, for any n ,

$$2\mu^*(A) \geq \sum_{k=1}^n \mu^*(B_k).$$

In particular, the infinite sum $\sum_{k=1}^{\infty} \mu^*(B_k)$ converges (recall we're assuming $\mu^*(A)$ is finite). So given $\epsilon > 0$, there is N_0 such that, for $N \geq N_0$, $\sum_{k=N}^{\infty} \mu^*(B_k) < \epsilon$. Now, by construction, we have

$$A \cap F^c = A_N \cup (\cup_{k=N}^{\infty} B_k),$$

so by the σ -subadditivity of μ^* we get, if $N \geq N_0$,

$$\mu^*(A \cap F^c) - \mu^*(A_N) \leq \sum_{k=N}^{\infty} \mu^*(B_k) < \epsilon.$$

Now $A_N \subset A \cap F^c$, so we in fact have $0 \leq \mu^*(A \cap F^c) - \mu^*(A_N) < \epsilon$ for $N \geq N_0$. Since $\epsilon > 0$ was arbitrary, this proves that $\mu^*(A \cap F^c) = \lim_{n \rightarrow \infty} \mu^*(A_n)$, completing the proof. \square

\square

\square

We now make the following definition.

Definition 3.19. If X is a set and \mathcal{M} is a σ -algebra in X , a function $\mu: \mathcal{M} \rightarrow [0, \infty]$ is called a *measure on \mathcal{M}* provided that we have:

- (i) $\mu(\emptyset) = 0$,
- (ii) If $\{E_n\}_{n=1}^{\infty}$ is a countable collection of disjoint sets in \mathcal{M} , then

$$\mu(\cup_{n=1}^{\infty} E_n) = \sum_{n=1}^{\infty} \mu(E_n).$$

Note that measures automatically satisfy the monotonicity axiom, since if $A \subset B$ with $A, B \in \mathcal{M}$, then also $B \setminus A \in \mathcal{M}$, and by (ii) $\mu(B) = \mu(A) + \mu(B \setminus A) \geq \mu(A)$. Measures are also σ -subadditive, since if $\{F_n\}_{n=1}^{\infty}$ is an arbitrary sequence of sets in \mathcal{M} then we can let $E_n = F_n \setminus \cup_{k=1}^{n-1} F_k$ and apply (ii) to the sets E_n and use the fact that $\mu(E_n) \leq \mu(F_n)$ to deduce that $\mu(\cup_{n=1}^{\infty} F_n) \leq \sum_{n=1}^{\infty} \mu(F_n)$.

What we've proven above gives a large family of examples of measures:

Corollary 3.20. *If μ^* is an outer measure on X , then the collection \mathcal{M} of μ^* -measurable subsets is a σ -algebra, and the function $\mu: \mathcal{M} \rightarrow [0, \infty]$ given by $\mu = \mu^*|_{\mathcal{M}}$ defines a measure on \mathcal{M} . Moreover, if (X, d) is a metric space and μ^* is a metric outer measure on X , then the σ -algebra \mathcal{M} contains all closed subsets of X .*

Note that as soon as a σ -algebra contains all closed sets, it contains all open sets since it's closed under complements. Moreover, it contains all countable unions of closed sets (F_{σ} 's), and all countable intersections of open sets (G_{δ} 's), and then, respectively, countable intersections and countable unions of these ($F_{\sigma\delta}$'s and $G_{\delta\sigma}$'s) and so on ad infinitum. It follows from a homework problem that it makes sense to talk about the smallest σ -algebra containing the closed sets; this smallest such σ -algebra is called the *Borel σ -algebra*, and its elements are called Borel sets. Thus any open or closed set is a Borel set, as are the G_{δ} 's, $F_{\sigma\delta}$'s, etc.

The special case that we're most interested in is where μ^* is the Lebesgue outer measure m^* on \mathbb{R} ; in that case we obtain a σ -algebra \mathcal{M} (called the Lebesgue σ -algebra) which contains the Borel σ -algebra (in fact the containment is strict, but we won't prove this), and a measure m on \mathcal{M} , called Lebesgue measure. By Theorem 3.13, it must not be the case that the Lebesgue σ -algebra contains every subset of \mathbb{R} ; however it contains essentially any set that you're likely to be able to explicitly construct, and will be quite adequate for our goal of constructing a theory of integration.

The general setting that we'll work in from now on is that of a *measure space*, which is formally defined as a triple (X, \mathcal{M}, μ) where X is some set, \mathcal{M} is a σ -algebra in X , and μ is a measure on \mathcal{M} . Of course this subsumes the Lebesgue case, but also is a useful notion in probability theory, where X is essentially viewed as the set of all possible states of the universe; where the measure μ satisfies $\mu(X) = 1$, and for a set $A \in \mathcal{M}$ one interprets $\mu(A)$ as the probability that the state of the universe will be represented among the elements of A . (The σ -algebra \mathcal{M} which is used should reflect the information about the universe that one has and/or is interested in. In particular, in contrast to the Lebesgue case, \mathcal{M} might be rather small; if one just cares about the outcome of a coin toss, one could use a σ -algebra with just four elements: the empty set; all of X ; all states of the universe in which the coin comes up heads; and all states of the universe in which the coin comes up tails.)

Here is a basic fact about the behavior of measures.

Proposition 3.21. *Let (X, \mathcal{M}, μ) be a measure space, and let $\{E_n\}_{n=1}^\infty$ be a sequence of sets in \mathcal{M} . Then:*

(i) *If for all n we have $E_n \subset E_{n+1}$, then*

$$\mu(\cup_{n=1}^\infty E_n) = \lim_{n \rightarrow \infty} \mu(E_n).$$

(ii) *If for all n we have $E_{n+1} \subset E_n$, and if $\mu(E_1) < \infty$, then*

$$\mu(\cap_{n=1}^\infty E_n) = \lim_{n \rightarrow \infty} \mu(E_n).$$

Note the restriction to the case $\mu(E_1) < \infty$ in the second part; you should have shown in your homework that this hypothesis really is necessary, as is illustrated by the example that μ is Lebesgue measure on $X = \mathbb{R}$ and $E_n = (n, \infty)$. (Of course, if your main interest is in probability theory, then you'll be working with measures which only take finite values, so this won't be so much of a concern to you.)

Proof. (i) Let $G_n = E_n \setminus E_{n-1}$. Since if $n > m$ we have $E_m \subset E_{n-1}$, we see that $G_n \cap G_m = \emptyset$ for $n > m$. Also, $E_n = \cup_{k=1}^n G_k$, so since the G_k are disjoint we have $\mu(E_n) = \sum_{k=1}^n \mu(G_k)$. Also, $\cup_{n=1}^\infty E_n = \cup_{n=1}^\infty G_n$, with the union on the right a disjoint one. Hence

$$\mu(\cup_{n=1}^\infty E_n) = \sum_{n=1}^\infty \mu(G_n) = \lim_{N \rightarrow \infty} \sum_{n=1}^N \mu(G_n) = \lim_{N \rightarrow \infty} \mu(E_N),$$

as desired.

(ii) This second part is proven simply by reducing to part (i). Let $F_n = E_1 \setminus E_n$. Then

$$\cup_{n=1}^\infty F_n = \cup_{n=1}^\infty (E_1 \setminus E_n) = E_1 \setminus (\cap_{n=1}^\infty E_n).$$

So part (i) gives

$$(7) \quad \mu(E_1 \setminus (\cap_{n=1}^\infty E_n)) = \lim_{n \rightarrow \infty} \mu(E_1 \setminus E_n).$$

Now as noted above for any two measurable sets A and B , if $A \subset B$ we have $\mu(A) + \mu(B \setminus A) = \mu(B)$. So if $\mu(B) < \infty$ (in which case both terms on the left are (finite) real numbers) we can subtract $\mu(A)$ from both sides to get $\mu(B \setminus A) = \mu(B) - \mu(A)$. Applying this to both sides of (7) then gives

$$\mu(E_1) - \mu(\cap_{n=1}^\infty E_n) = \mu(E_1) - \lim_{n \rightarrow \infty} \mu(E_n).$$

Since $\mu(E_1)$ is finite, we can then cancel it to obtain the desired result. \square

3.3. Measurable functions. Throughout this section we will work in a fixed measure space (X, \mathcal{M}, μ) . A subset $A \subset X$ will be called *measurable* if it belongs to \mathcal{M} .

As mentioned earlier, the original purpose of the development of measure theory was to set up a theory of integration for certain functions; here is the class of functions that we will try to integrate:

Definition 3.22. A function $f: X \rightarrow [-\infty, \infty]$ is called *measurable* if, for every $\alpha \in \mathbb{R}$, the set

$$f^{-1}((\alpha, \infty]) = \{x \in X \mid f(x) > \alpha\}$$

belongs to the σ -algebra \mathcal{M} .

Obviously, whether a function is measurable or not depends on the σ -algebra \mathcal{M} that we're working with. However, given (X, \mathcal{M}) the notion of the measurability of a function doesn't depend on the measure μ .

Note that we allow our functions to take the values ∞ and $-\infty$; this will be convenient when we take suprema, infima, and limits later. One disadvantage of doing this is that it is no longer always possible to add two numbers in the range $[-\infty, \infty]$ of our functions, since one can't make much sense of $\infty + (-\infty)$. So if $f, g: X \rightarrow [-\infty, \infty]$ are two functions, their sum $f + g$ might not be well-defined (specifically, $f + g$ is only well-defined on the complement of the set $(f^{-1}(\{-\infty\}) \cap g^{-1}(\{\infty\})) \cup (f^{-1}(\{\infty\}) \cap g^{-1}(\{-\infty\}))$). I'll tend to ignore that issue here, leaving it to the reader to mentally add caveats along the lines of "assuming that $f + g$ is defined" when they are appropriate.

Perhaps the simplest family of measurable functions are given as follows. If $A \subset X$, define the *characteristic function* χ_A of A by

$$\chi_A(x) = \begin{cases} 1 & x \in A \\ 0 & x \notin A \end{cases}$$

Then χ_A is a measurable function if and only if $A \in \mathcal{M}$. Indeed, the possible preimages $\chi_A^{-1}((\alpha, \infty])$ as α ranges over \mathbb{R} are just \emptyset , A , and X , of which the first and last automatically belong to \mathcal{M} , so the only set whose measurability needs to be checked is A .

If X happens to be the set underlying a metric space (X, d) , and if the σ -algebra \mathcal{M} includes all of the open subsets (as is true in the case of Lebesgue measure on \mathbb{R} , or more generally in the case of any measure obtained from a metric outer measure, as we showed last section), then if $f: X \rightarrow \mathbb{R}$ is continuous each $f^{-1}((\alpha, \infty])$ is open and therefore is measurable (*i.e.*, belongs to \mathcal{M}). Thus for Lebesgue measure and its cousins any continuous function is measurable. However, we'll find that measurability is a much more flexible notion than continuity. Of course, the characteristic functions of the previous paragraph already give a substantial supply of measurable functions which aren't continuous.

Given a countable collection of functions, there are various ways of making new functions out of them; it tends to be the case that if these functions that you start with are measurable then so will be the new functions that you create from them. A first example of this (for the case of starting with just two functions) is given by the following:

Proposition 3.23. *Let $f, g: X \rightarrow [-\infty, \infty]$ be two measurable functions. Then*

- (i) $\max\{f, g\}$ is measurable.
- (ii) $\min\{f, g\}$ is measurable.
- (iii) $f + g$ is measurable.

Proof. (i) For $\alpha \in \mathbb{R}$, we have

$$\max\{f, g\}^{-1}((\alpha, \infty]) = \{x \mid \max\{f(x), g(x)\} > \alpha\} = f^{-1}((\alpha, \infty]) \cup g^{-1}((\alpha, \infty]).$$

Since f and g are measurable, both $f^{-1}((\alpha, \infty])$ and $g^{-1}((\alpha, \infty])$ belong to \mathcal{M} , so since \mathcal{M} , being a σ -algebra, is closed under finite unions, we see that $\max\{f, g\}^{-1}((\alpha, \infty])$ belongs to \mathcal{M} , so that $\max\{f, g\}$ is measurable.

(ii) Similarly, for $\alpha \in \mathbb{R}$, we have

$$\min\{f, g\}^{-1}((\alpha, \infty]) = f^{-1}((\alpha, \infty]) \cap g^{-1}((\alpha, \infty]),$$

so since \mathcal{M} is closed under finite intersections we deduce that $\min\{f, g\}$ is measurable.

(iii) (Of course, as noted shortly after the definition of measurability, we should assume here that $f + g$ is well defined.) If $\alpha \in \mathbb{R}$, notice that $(f + g)(x) > \alpha$ if and only if there is $\beta \in \mathbb{Q}$ such that $g(x) > \beta$ and $f(x) > \alpha - \beta$ (here \mathbb{Q} denotes the set of rational numbers). Indeed, the backward implication is obvious by adding the two inequalities, and for the forward implication we can choose any rational β such that $0 < g(x) - \beta < f(x) + g(x) - \alpha$; such a β exists since the rationals are dense in \mathbb{R} (*i.e.*, by the fact that any real number can be arbitrarily-well approximated by rationals). Hence

$$(f + g)^{-1}((\alpha, \infty]) = \cup_{\beta \in \mathbb{Q}} (g^{-1}((\beta, \infty]) \cup f^{-1}((\alpha - \beta, \infty])).$$

By the measurability of f and of g , this is a countable union of sets in \mathcal{M} and so belongs to \mathcal{M} . □

Proposition 3.24. *If $I \subset [-\infty, \infty]$ is any interval and $f: X \rightarrow [-\infty, \infty]$ is measurable then $f^{-1}(I) \in \mathcal{M}$.*

Proof. Of course the definition of measurability is that this proposition holds for intervals of the special form $(\alpha, \infty]$; we just need to go through the other kinds of intervals. α, β will denote arbitrary real numbers.

If $I = [\alpha, \infty]$, we see that $f^{-1}([\alpha, \infty]) = \cap_{n=1}^{\infty} f^{-1}((\alpha - 1/n, \infty])$ is a countable intersection of measurable sets and hence is measurable.

If $I = [-\infty, \beta]$, $f^{-1}(I) = (f^{-1}((\beta, \infty]))^c$ is the complement of a measurable set and hence is measurable.

If $I = [-\infty, \beta)$, $f^{-1}(I) = (f^{-1}([\beta, \infty]))^c$ is the complement of a measurable set and hence is measurable.

If $I = \langle \alpha, \beta \rangle$ (where as earlier angle-brackets mean either (or), $f^{-1}(I) = f^{-1}(\langle \alpha, \infty \rangle) \cap f^{-1}([-\infty, \beta))$ is an intersection of two measurable sets and hence is measurable.

If $I = \langle \alpha, \infty \rangle$, then $f^{-1}(I) = \cup_{n=1}^{\infty} f^{-1}(\langle \alpha, n \rangle)$ is a countable union of measurable sets and hence is measurable.

If $I = (-\infty, \beta)$, then $f^{-1}(I) = \cup_{n=1}^{\infty} f^{-1}((-n, \beta))$ is a countable union of measurable sets and hence is measurable.

If $I = (-\infty, \infty)$, $f^{-1}(I) = \cup_{n=1}^{\infty} f^{-1}((-n, n))$ is measurable.

If $I = [\infty, \infty] = \{\infty\}$ $f^{-1}(I) = \cap_{n=1}^{\infty} f^{-1}((n, \infty])$ is a countable intersection of measurable sets, so is measurable.

If $I = [-\infty, -\infty] = \{-\infty\}$ $f^{-1}(I) = \cap_{n=1}^{\infty} f^{-1}([-\infty, -n])$ is a countable intersection of measurable sets, so is measurable.

If $I = \langle -\infty, \infty \rangle$, $f^{-1}(I)$ is the union of $f^{-1}((-\infty, \infty))$ and zero, one, or two of the measurable sets $f^{-1}(\{\infty\})$, $f^{-1}(\{-\infty\})$, and so is measurable. □

Consider any sequence $\{a_n\}_{n=1}^{\infty}$ of extended real numbers (*i.e.*, elements of $[-\infty, \infty]$). Since a set of real numbers either has a real supremum (least upper bound) or is unbounded above (in which case we say its supremum is ∞), and likewise either has a real infimum (greatest lower bound) or is unbounded below (in which case we say its infimum is $-\infty$), our sequence $\{a_n\}_{n=1}^{\infty}$ will have both a well-defined supremum and a well-defined infimum, which will likewise be elements of $[-\infty, \infty]$.

We may also define

$$\overline{\lim} a_n = \limsup a_n := \inf_{n \geq 1} \sup\{a_k | k \geq n\}$$

and

$$\underline{\lim} a_n = \liminf a_n := \sup_{n \geq 1} \inf\{a_k | k \geq n\}.$$

By our remarks above, both $\overline{\lim} a_n$ and $\underline{\lim} a_n$ are well-defined elements of $[-\infty, \infty]$, for *any* sequence $\{a_n\}_{n=1}^{\infty}$ of extended real numbers.

If you haven't seen the concepts of $\overline{\lim}$, $\underline{\lim}$ before, it might be instructive to check that $\overline{\lim} a_n$ is the largest limit of any subsequence of $\{a_n\}_{n=1}^{\infty}$ while $\underline{\lim} a_n$ is the smallest limit of any subsequence, and that moreover $\overline{\lim} a_n = \underline{\lim} a_n$ if and only if $\lim_{n \rightarrow \infty} a_n$ exists (in which case the limit is equal to the common value of the \liminf and the \limsup).

With this said, given a sequence of functions $f_n: X \rightarrow [-\infty, \infty]$, we may define new functions

$$\overline{\lim} f_n \quad \underline{\lim} f_n$$

by saying that $(\overline{\lim} f_n)(x) = \overline{\lim}(f_n(x))$ and $(\underline{\lim} f_n)(x) = \underline{\lim}(f_n(x))$. More simply, we can also define functions $\inf f_n$, $\sup f_n$ by $(\inf f_n)(x) = \inf(f_n(x))$, and $(\sup f_n)(x) = \sup(f_n(x))$. A nice feature of our notion of measurability is that it behaves well with respect to these operations:

Theorem 3.25. *If $\{f_n\}_{n=1}^{\infty}$ is a sequence of measurable functions $f_n: X \rightarrow [-\infty, \infty]$, then each of*

$$\inf f_n, \sup f_n, \overline{\lim} f_n, \underline{\lim} f_n$$

is also measurable.

Proof. We have

$$x \in (\sup f_n)^{-1}((\alpha, \infty]) \Leftrightarrow \sup(f_n(x)) > \alpha \Leftrightarrow x \in \cup_{n=1}^{\infty} f_n^{-1}((\alpha, \infty]),$$

and this last set is a countable union of measurable sets and hence is measurable. Thus $\sup f_n$ is measurable. Somewhat similarly,

$$x \in (\inf f_n)^{-1}([\alpha, \infty]) \Leftrightarrow x \in \cap_{n=1}^{\infty} f_n^{-1}([\alpha, \infty]),$$

so each $(\inf f_n)^{-1}([\alpha, \infty])$ is measurable, in view of which $(\inf f_n)^{-1}((\alpha, \infty]) = \cup_{n=1}^{\infty} (\inf f_n)^{-1}([\alpha + 1/n, \infty])$ is measurable. Thus $\inf f_n$ is measurable.

The statements about $\overline{\lim}$ and $\underline{\lim}$ then follow directly from this, since $\overline{\lim} f_n = \inf_{n \geq 1} \sup\{f_k | k \geq n\}$ and $\underline{\lim} f_n = \sup_{n \geq 1} \inf\{f_k | k \geq n\}$. \square

Corollary 3.26. *If $f_n: X \rightarrow [-\infty, \infty]$ and $f: X \rightarrow [-\infty, \infty]$ are functions such that each f_n is measurable and, for all $x \in X$, $f_n(x) \rightarrow f(x)$, then f is also measurable.*

Proof. Indeed, in this case $f = \overline{\lim} f_n$. \square

Thus measurability behaves well with respect to taking limits, which is one of the nicer features of the whole theory (recall, by contrast, that the pointwise limit of a sequence of continuous functions is very often not continuous).

3.4. Simple functions, and the integration of nonnegative measurable functions. We're now ready to start defining the Lebesgue integral. The first class of functions that we will learn how to integrate consists of the following:

Definition 3.27. Let (X, \mathcal{M}, μ) be a measure space. A function $s: X \rightarrow \mathbb{R}$ is called *simple* if

- (i) s is measurable; and
- (ii) the range $s(X)$ of s is a finite set.

(A few authors don't include the measurability assumption in the definition of a simple function.) Note that, whereas a general measurable function is allowed to take the values $-\infty$ or ∞ , all the values of a simple function are required to be real numbers.

Let s be a simple function; condition (ii) above tells us that, for certain real numbers a_1, \dots, a_n , we have $s(X) = \{a_1, \dots, a_n\}$. Write $A_i = s^{-1}(\{a_i\})$. Then the sets A_i are pairwise disjoint, and their union is all of X . So since $s(x) = a_i$ if $x \in A_i$, we see that, recalling from above the definition of the characteristic function χ_{A_i} ,

$$s = \sum_{i=1}^n a_i \chi_{A_i}.$$

Furthermore, each A_i is the preimage under s of an interval, so by Proposition 3.24 each set A_i is measurable. So any simple function is a *linear combination of characteristic functions of measurable subsets*. Conversely, it's not hard to see that any linear combination of characteristic functions of measurable subsets is a simple function; such a linear combination is measurable because characteristic functions of measurable subsets are measurable and linear combinations of measurable functions are measurable, and a moment's thought shows that such a linear combination only takes finitely many values (this is still true if the subsets aren't disjoint; generally if there are n subsets then there will be no more than 2^n values).

So suppose that $s: X \rightarrow [0, \infty)$ is a simple function, which takes only non-negative values. As mentioned earlier, if $s(X) = \{a_1, \dots, a_n\}$ we have a *canonical representation* of s as

$$s = \sum_{i=1}^n a_i \chi_{A_i},$$

where the $A_i = s^{-1}(\{a_i\})$ are pairwise disjoint.

Convention 3.28. To multiply arbitrary elements of $[0, \infty]$, we set $a \cdot \infty = \infty \cdot a = \infty \cdot \infty = \infty$ for all $a \in (0, \infty)$, but

$$0 \cdot \infty = \infty \cdot 0 = 0.$$

We now define the integral of s , over any measurable subset $E \in \mathcal{M}$:

$$(8) \quad \int_E s d\mu = \sum_{i=1}^n a_i \mu(A_i \cap E).$$

(So in particular, for instance, $\int_X \chi_A d\mu = \mu(A)$, in accordance with what we suggested at the start of the unit). Note that the above convention implies, for instance, that the integral of the function which is identically zero is always zero (even if it's being integrated over a set of infinite measure), as probably seems appropriate.

Now although there is a canonical way, described above, of writing a simple function as a linear combination of characteristic functions, it may happen that we are given an expression for a simple function as a linear combination of characteristic functions which doesn't have this special canonical form (for instance the subsets A_i might not be disjoint). Accordingly, the following fact is often helpful.

Lemma 3.29. *If for some $a_i, b_j \geq 0$ and $A_i, B_j \in \mathcal{M}$ we have*

$$\sum_{i=1}^n a_i \chi_{A_i} = \sum_{j=1}^m b_j \chi_{B_j},$$

then

$$\sum_{i=1}^n a_i \mu(A_i) = \sum_{j=1}^m b_j \mu(B_j).$$

Proof. This is an exercise on Problem Set 7. □

Corollary 3.30. *For any $a_1, \dots, a_n \geq 0$ and $A_1, \dots, A_n, E \in \mathcal{M}$ (not necessarily disjoint), we have*

$$\int_E \left(\sum_{i=1}^n a_i \chi_{A_i} \right) d\mu = \sum_{i=1}^n a_i \mu(A_i \cap E).$$

Proof. Express the simple function $\sum_{i=1}^n a_i \chi_{A_i}$ in its canonical form $\sum_{j=1}^m b_j \chi_{B_j}$, and apply (8) (with the a 's and A 's replaced by b 's and B 's) and the previous lemma. □

Note that the sum of two simple functions is simple.

Proposition 3.31. *If s and t are nonnegative simple functions and $E \in \mathcal{M}$ then*

$$\int_E (s + t) d\mu = \int_E s d\mu + \int_E t d\mu.$$

Proof. Write $s = \sum_{i=1}^n a_i \chi_{A_i}$ and $t = \sum_{j=1}^m b_j \chi_{B_j}$, with $A_i \cap A_{i'} = B_j \cap B_{j'} = \emptyset$ for $i \neq i'$ and $j \neq j'$, and $\cup_{i=1}^n A_i = \cup_{j=1}^m B_j = X$. Since the A_i (respectively, the B_j) are disjoint and have union equal to all of X , we have $\sum_{i=1}^n \chi_{A_i} = 1$ (respectively, $\sum_{j=1}^m \chi_{B_j} = 1$). Further, for any sets C and D inspection of the definition shows that $\chi_C \chi_D = \chi_{C \cap D}$. So

$$\chi_{A_i} = \sum_{j=1}^m \chi_{A_i \cap B_j} \quad \chi_{B_j} = \sum_{i=1}^n \chi_{A_i \cap B_j}.$$

As such, we obtain

$$s + t = \sum_{i=1}^n \sum_{j=1}^m (a_i + b_j) \chi_{A_i \cap B_j}.$$

Hence, by the previous corollary,

$$\begin{aligned} \int_E (s+t)d\mu &= \sum_{i=1}^n \sum_{j=1}^m (a_i + b_j) \mu(A_i \cap B_j \cap E) \\ &= \sum_{i=1}^n a_i \sum_{j=1}^m \mu(A_i \cap B_j \cap E) + \sum_{j=1}^m b_j \sum_{i=1}^n \mu(A_i \cap B_j \cap E) \\ &= \sum_{i=1}^n a_i \mu(A_i \cap E) + \sum_{j=1}^m b_j \mu(B_j \cap E) = \int_E s d\mu + \int_E t d\mu, \end{aligned}$$

where in the second to last inequality we've used the additivity of the measure μ and the fact that the sets $A_i \cap B_j \cap E$ are pairwise disjoint (with union $A_i \cap E$ when i is fixed and j varies, and union $B_j \cap E$ when j is fixed and i varies). \square

Corollary 3.32. *If s and t are simple functions with $0 \leq s(x) \leq t(x)$ for all $x \in X$, then, for any $E \in \mathcal{M}$,*

$$\int_E s d\mu \leq \int_E t d\mu.$$

Proof. $t - s$ is a simple function, and since $t(x) \geq s(x)$ for all x $t - s$ is in fact a nonnegative simple function. So, using Proposition 3.31,

$$\int_E t d\mu = \int_E s d\mu + \int_E (t - s) d\mu \geq \int_E s d\mu,$$

since the integral of a nonnegative simple function is nonnegative. \square

We can now define the integral of a general nonnegative measurable function:

Definition 3.33. Let $f: X \rightarrow [0, \infty]$ be any nonnegative measurable function, and $E \in \mathcal{M}$. Then we define

$$\int_E f d\mu = \sup \left\{ \int_E s d\mu \mid s \text{ is simple, } (\forall x \in X)(0 \leq s(x) \leq f(x)) \right\}.$$

Note that Corollary 3.32 implies that, in the special case that f happens to be simple, this definition is consistent with the previous one.

We have a bit of work to do before it will be fully clear why this is a good definition. We've started with a natural notion of what the integral of a nonnegative simple function should be, and the integral of a general nonnegative measurable function f is then obtained by looking at simple functions sandwiched between 0 and f . For this to be appropriate, we'd want to know that there are simple functions between 0 and f that do a good job of reflecting the properties of the general measurable function f . Here is a result in that direction:

Theorem 3.34. *Let $f: X \rightarrow [0, \infty]$ be a nonnegative, measurable function. Then there is a sequence $\{s_n\}_{n=1}^{\infty}$ of simple functions such that for all $x \in X$ and $n \in \mathbb{N}$ we have $0 \leq s_n(x) \leq s_{n+1}(x) \leq f(x)$, and for all $x \in X$ we have $s_n(x) \rightarrow f(x)$ as $n \rightarrow \infty$.*

The conclusion is typically phrased as saying that the sequence s_n "converges monotonically" to f , and is abbreviated $s_n \nearrow f$.

Proof. The idea is to chop the range of f into finer and finer pieces, and define s_n to be constant on the preimages of each of these pieces (also, since the s_n all have finite range contained in \mathbb{R} , whereas the range of f is unbounded and may even contain ∞ , we need to have the ranges of the s_n expand upward as n grows). Here is one way of achieving this. Given n , consider integer values of k ranging from 0 to $n2^n$. Define

$$A_{n,k} = f^{-1}([k2^{-n}, (k+1)2^{-n})) \quad (0 \leq k \leq n2^n - 1), \quad A_{n,n2^n} = f^{-1}([n, \infty)).$$

Since f is measurable, the sets $A_{n,k}$ are all measurable, and by definition we have $f(x) \geq k2^{-n}$ whenever $x \in A_{n,k}$, while if $f(x) \leq n$ we have $f(x) < (k+1)2^{-n}$ when $x \in A_{n,k}$. Also the $A_{n,k}$ are pairwise disjoint as k varies with n fixed, with union equal to all of X . Let

$$s_n = \sum_{k=0}^{n2^n} \frac{k}{2^n} \chi_{A_{n,k}}.$$

Since $s_n(x) = k2^{-n}$ for $x \in A_{n,k}$, we have $s_n(x) \leq f(x)$ for all n and x . If $f(x) = \infty$, then $s_n(x) = n$ for all x , while if $f(x) < \infty$, then there is N such that $f(x) \leq N$, in which case we have, for each $n \geq N$, $0 \leq f(x) - s_n(x) < 2^{-n}$. In view of this we have $s_n(x) \rightarrow f(x)$ for all x . It remains only to show that $s_n(x) \leq s_{n+1}(x)$ for all x . But this too follows quickly from the definitions: if $s_n(x) < n$, so $x \in A_{n,k}$ for some $k < n2^n$ and $s_n(x) = k2^{-n}$, we have either $f(x) \in [(2k)2^{-n-1}, (2k+1)2^{-n-1})$, in which case $s_{n+1}(x) = (2k)2^{-n-1} = k2^{-n}$, or else $f(x) \in [(2k+1)2^{-n-1}, (2k+2)2^{-n-1})$, in which case $s_{n+1}(x) = (2k+1)2^{-n-1}$. Similarly, if $s_n(x) = n$, then $x \in A_{k,n+1}$ for some $k \geq n2^{n+1}$, in view of which $s_{n+1}(x)$ is some number in $[n, n+1]$. So (since by definition $s_n(x) \leq n$ for all x), we indeed have $s_{n+1}(x) \geq s_n(x)$ for all n and x . \square

If $f: X \rightarrow [0, \infty]$ is nonnegative and measurable, we now know that there are nonnegative simple functions s_n such that $s_n \nearrow f$; meanwhile by the definition of the integral we know that there are nonnegative simple functions t_n such that $\int_X t_n d\mu \nearrow \int_X f d\mu$. It is reasonable to expect that the t_n can be taken equal to s_n ; this doesn't yet follow from what we've proven, but rather is a special case of the Monotone Convergence Theorem, to be proven below.

Let us now record some properties of the integral; each of these follows quite quickly from the definition, so I'll leave the proofs to the reader. f and g denote arbitrary nonnegative measurable functions, and E and F denote measurable subsets of X (i.e., elements of the σ -algebra \mathcal{M}).

1. If $0 \leq f \leq g$ then $\int_E f d\mu \leq \int_E g d\mu$.
2. If $E \subset F$ then $\int_E f d\mu \leq \int_F f d\mu$.
3. If c is a nonnegative real number then $\int_E (cf) d\mu = c \int_E f d\mu$.
4. If $f(x) = 0$ for all $x \in E$ then $\int_E f d\mu = 0$ (even if $\mu(E) = \infty$).
5. If $\mu(E) = 0$ then $\int_E f d\mu = 0$ (even if $f(x) = \infty$ for all $x \in E$).
6. $\int_E f d\mu = \int_X f \chi_E d\mu$.

Noticeably absent from the above list is a statement to the effect that $\int (f+g) = \int f + \int g$; this too will have to wait for the Monotone Convergence Theorem.

Here is one of the most fundamental results about the Lebesgue integral:

Theorem 3.35 (Fatou's Lemma). *If $f_n: X \rightarrow [0, \infty]$ is a sequence of nonnegative, measurable functions on the measure space (X, \mathcal{M}, μ) , and $E \in \mathcal{M}$, then*

$$\int_E (\liminf f_n) d\mu \leq \liminf \int_E f_n d\mu.$$

Note that the \liminf 's in the theorem always exist without any hypothesis on the f_n ; as such this is a very widely applicable result (more so than the Monotone and Dominated Convergence Theorems, which we'll deduce as corollaries of Fatou's Lemma later; these other results give an equality rather than an inequality, but depend on fairly strong assumptions about the f_n).

Example 3.36. Where $X = \mathbb{R}$ and μ is Lebesgue measure, consider the sequence of functions $f_n: \mathbb{R} \rightarrow [0, \infty]$ given by $f_n(x) = n\chi_{[0, 1/n]}$. So for all x $f_n(x) \rightarrow f(x)$, where $f(x) = 0$ if $x \neq 0$, and $f(0) = \infty$. So $\int_{\mathbb{R}} \liminf f_n d\mu = \int_{\mathbb{R}} f d\mu = 0$, while $\int_{\mathbb{R}} f_n d\mu = 1$ for all n . This shows that equality sometimes doesn't hold in Fatou's Lemma, and is a good example to remember if ever you can't recall which way the inequality goes.

Proof of Fatou's Lemma. By the definition of the integral, it is enough to show that whenever s is a simple function with $0 \leq s \leq \liminf f_n$, we have $\liminf \int_E f_n d\mu \geq \int_E s d\mu$. The proof naturally splits into two cases, though the idea in both is rather similar.

Case 1: $\int_E s d\mu = \infty$. Now that s is a simple function means that it takes only finitely many values, each of which is finite. So if $A = \{x \in E | s(x) \neq 0\}$, it must be that $\mu(A) = \infty$ (since $\int_E s d\mu$, which is assumed infinite, is at most the (finite) maximum of s times $\mu(A)$). Also, for some $a > 0$, we have $s(x) > a$ whenever $x \in A$ (for instance a could be taken equal to half the smallest nonzero value taken by s). So if $x \in A$, we have $\liminf f_n(x) > a$, which implies that, for some N (depending on x), it is the case that whenever $n \geq N$ we have $f_n(x) > a$. Thus where

$$A_N = \{x \in A | f_n(x) > a \text{ whenever } n \geq N\},$$

we have $A = \cup_{N=1}^{\infty} A_N$. From the definition, it's clear that $A_N \subset A_{N+1}$ for each N , so by proposition 3.21 we have $\lim_{N \rightarrow \infty} \mu(A_N) = \mu(A) = \infty$.

So if $M > 0$ there is N such that $\mu(A_N) \geq M/a$, so since $f_n(x) > a$ for all $n \geq N$ and $x \in A_N$ we have, for $n \geq N$,

$$\int_E f_n d\mu \geq \int_{A_N} f_n d\mu \geq \int_{A_N} a d\mu \geq M.$$

So indeed

$$\int_E f_n d\mu \rightarrow \infty = \int_E s d\mu,$$

as desired.

Case 2: $\int_E s d\mu < \infty$. In this case we need to show that for all $\delta > 0$ there is N such that if $n \geq N$ we have $\int_E f_n d\mu \geq \int_E s d\mu - \delta$.

Let $\epsilon > 0$ be arbitrary. As before, let $A = \{x \in E | s(x) > 0\}$. Now since s is simple (so its restriction to A is bounded below by a positive number, namely its minimal nonzero value), the fact that $\int_E s d\mu < \infty$ implies that

$$\mu(A) < \infty.$$

If $x \in A$, we have $(1 - \epsilon)s(x) < s(x)$, so there is N (depending on x) such that $f_n(x) \geq (1 - \epsilon)s(x)$ for all $n \geq N$. So if

$$A_N = \{x \in A | f_n(x) \geq (1 - \epsilon)s(x) \text{ whenever } n \geq N\},$$

the A_N form an increasing sequence of sets whose union is A . Hence the $A \setminus A_N$ form a decreasing sequence of sets (each with finite measure since $\mu(A) < \infty$) whose intersection is empty. Hence by part (ii) of Proposition 3.21, we have $\mu(A \setminus A_N) \rightarrow 0$ as $N \rightarrow \infty$. So, for some N , we have $\mu(A \setminus A_N) < \epsilon$. Then if $n \geq N$, we have

$$\int_E f_n d\mu \geq \int_{A_N} f_n d\mu \geq (1 - \epsilon) \int_{A_N} s d\mu,$$

while, writing M for the largest value attained by the simple function s ,

$$\int_E s d\mu = \int_{A_N} s d\mu + \int_{A \setminus A_N} s d\mu \leq \int_{A_N} s d\mu + M\epsilon.$$

So we've shown that, for any $\epsilon > 0$, there is N such that when $n \geq N$ such that

$$\int_E f_n d\mu \geq (1 - \epsilon) \left(\int_E s d\mu - M\epsilon \right).$$

Of course, $\int_E s d\mu$ and M are finite numbers which depend only on the simple function s (not on ϵ or N). So for any given $\delta > 0$ we can make $\epsilon > 0$ so small that the right hand side above differs from $\int_E s d\mu$ by at most δ , and then the N corresponding to this ϵ will indeed have the property that for $n \geq N$ it holds that $\int_E f_n d\mu \geq \int_E s d\mu - \delta$, as desired. \square

Corollary 3.37 (Monotone Convergence Theorem). *If $\{f_n\}_{n=1}^\infty$ is a sequence of nonnegative measurable functions from X to $[0, \infty]$ such that $f_n \nearrow f$, and if $E \in \mathcal{M}$, then*

$$\int_E f d\mu = \lim_{n \rightarrow \infty} \int_E f_n d\mu.$$

(Again, $f_n \nearrow f$ means that $f_n(x) \rightarrow f(x)$ for all x , and $f_n(x) \leq f_{n+1}(x)$ for all n and x .)

Proof. The fact that $f(x) \geq f_n(x)$ for all n and x implies that $\int_E f d\mu \geq \int_E f_n d\mu$ for all n , and therefore that $\int_E f d\mu \geq \overline{\lim} \int_E f_n d\mu$. On the other hand, since $f = \underline{\lim} f_n$, Fatou's Lemma asserts that $\int_E f d\mu \leq \underline{\lim} \int_E f_n d\mu$. So we've shown that $\underline{\lim} \int_E f_n d\mu \geq \overline{\lim} \int_E f_n d\mu$. But the reverse inequality $\overline{\lim} \geq \underline{\lim}$ always holds for any sequence of extended real numbers. Hence the liminf and the limsup must be equal (and so the limit $\lim_{n \rightarrow \infty} \int_E f_n d\mu$ exists and is equal to their common value), so since we've shown that $\int_E f d\mu$ is sandwiched between the liminf and the limsup it must in fact be equal to both (and thus to the limit). \square

In particular, the simple functions from Theorem 3.34 have integrals converging to the integral of f . So $\int_E f d\mu$ is equal to the limit of the integrals of some (and indeed every) sequence of simple functions converging monotonically to f . The existence of these simple functions, combined with the Monotone Convergence theorem, can come in handy, as in the following:

Corollary 3.38. *If $f, g: X \rightarrow [0, \infty]$ are nonnegative, measurable functions and $E \in \mathcal{M}$ we have*

$$\int_E (f + g) d\mu = \int_E f d\mu + \int_E g d\mu.$$

Proof. Choose sequences of simple functions $s_n \nearrow f$ and $t_n \nearrow g$. Then $s_n + t_n \nearrow f + g$, so the Monotone Convergence Theorem gives

$$\int_E (f+g)d\mu = \lim_{n \rightarrow \infty} \int_E (s_n+t_n)d\mu = \lim_{n \rightarrow \infty} \left(\int_E s_n d\mu + \int_E t_n d\mu \right) = \int_E f d\mu + \int_E g d\mu,$$

where in the second equality we use Proposition 3.31. \square

Corollary 3.39. *If $\{f_n\}_{n=1}^\infty$ is a sequence of nonnegative measurable functions on X and $E \in \mathcal{M}$ then*

$$\int_E \left(\sum_{n=1}^\infty f_n \right) d\mu = \sum_{n=1}^\infty \left(\int_E f_n d\mu \right).$$

Proof. This follows directly from the previous corollary applied inductively to give $\sum_{n=1}^N \int_E f_n d\mu = \int_E \sum_{n=1}^N f_n d\mu$, and then the Monotone Convergence Theorem (recalling that the notation $\sum_{n=1}^\infty$ precisely means $\lim_{N \rightarrow \infty} \sum_{n=1}^N$). \square

Now might be a good time to record the following:

Proposition 3.40. *Let $f: \mathbb{R} \rightarrow [0, \infty)$ be a nonnegative continuous function. Then $\int_{[a,b]} f d\mu$ is equal to the Riemann integral $\int_a^b f(x)dx$.*

Proof. Recall that the Riemann integral of f is defined as the limit of a sequence of “lower sums” L_n (which, in the continuous case, is equal to the limit of a sequence of “upper sums”), which may for example be chosen by dividing the interval $[a, b]$ into equal-length subintervals $I_{n,k}$ ($1 \leq k \leq 2^n$) of length $(b-a)2^{-n}$, and adding up the areas of the rectangles of width $(b-a)2^{-n}$ and height equal to the minimum value of f on the interval $I_{n,k}$. But this area sum L_n is precisely equal to the (Lebesgue) integral of a step function $s_n = \sum_{k=1}^{2^n} c_k \chi_{I_{n,k}}$, with the c_k equal to the minimal value of f on $I_{n,k}$. Now the continuity of f implies that $s_n \nearrow f$, and so by the Monotone Convergence Theorem

$$\int_{[a,b]} f d\mu = \lim_{n \rightarrow \infty} \int_{[a,b]} s_n d\mu = \lim_{n \rightarrow \infty} L_n.$$

But, as mentioned earlier, the Riemann integral of f is precisely defined to be $\lim_{n \rightarrow \infty} L_n$. \square

In fact, any function which is Riemann integrable on an interval $[a, b]$ has its Riemann integral equal to its Lebesgue integral; to actually prove this would require somewhat more of a digression into the properties of the Riemann integral than would be appropriate given time constraints. But, of course, there are many more functions that one can integrate in the Lebesgue theory than there are in the Riemann theory (in fact, Lebesgue proved that a function is Riemann integrable if and only if it is continuous except on a set of Lebesgue measure zero, which is of course a much stronger requirement than measurability).

3.5. Integrable functions. We’ve now seen how to integrate an arbitrary *nonnegative* measurable function on a measure space (X, \mathcal{M}, μ) . Now of course we may occasionally want to integrate a function that takes negative values. We (try to) achieve this by just reducing to the nonnegative case, as follows. Let

$f: X \rightarrow [-\infty, \infty]$ be an arbitrary measurable function. Define the *positive part* f^+ and the *negative part* f^- of f by

$$f^+(x) = \max\{f(x), 0\} \quad f^-(x) = \max\{-f(x), 0\}.$$

Both f^+ and f^- are nonnegative functions, and by Proposition 3.23 both are measurable. Hence for any $E \in \mathcal{M}$ we have well-defined integrals $\int_E f^+ d\mu$ and $\int_E f^- d\mu$. Further, one has

$$f = f^+ - f^- \quad |f| = f^+ + f^-.$$

In view of this, the obvious definition for the integral of f over E would be $\int_E f d\mu = \int_E f^+ d\mu - \int_E f^- d\mu$. Now there's one potential problem with this, namely that $\int_E f^+ d\mu$ and $\int_E f^- d\mu$ might both be infinite, and then this definition would give $\infty - \infty$, which has no natural interpretation. So we assume this problem away by requiring f to be as in the following definition.

Definition 3.41. A function $f: X \rightarrow [-\infty, \infty]$ is called *integrable* if f is measurable and $\int_X f^+ d\mu$ and $\int_X f^- d\mu$ are both finite. Further, we define

$$L^1(\mu) = \{f: X \rightarrow [-\infty, \infty] \mid f \text{ is integrable}\}.$$

Note that f is integrable if and only if it is measurable and $\int_X |f| d\mu$ is finite.

Definition 3.42. If $f \in L^1(\mu)$ and $E \in \mathcal{M}$ then

$$\int_E f d\mu = \int_E f^+ d\mu - \int_E f^- d\mu.$$

(Slightly more generally, we could require only that one of $\int_E f^+ d\mu$ and $\int_E f^- d\mu$ be finite, since $\int_E f^+ d\mu - \int_E f^- d\mu$ would then still make sense even if one of the terms is infinite.)

Proposition 3.43. If $f, g \in L^1(\mu)$ and $E \in \mathcal{M}$ then

$$\int_E (f + g) d\mu = \int_E f d\mu + \int_E g d\mu.$$

Proof. We have

$$(f + g)^+ - (f + g)^- = f + g = f^+ - f^- + g^+ - g^-,$$

so

$$(f + g)^+ + f^- + g^- = f^+ + g^+ + (f + g)^-.$$

All the terms above are nonnegative measurable functions, so by Corollary 3.38 we have

$$\int_E (f + g)^+ d\mu + \int_E f^- d\mu + \int_E g^- d\mu = \int_E f^+ d\mu + \int_E g^+ d\mu + \int_E (f + g)^- d\mu.$$

Rearranging these terms (and using that all of them are finite) then proves the result. \square

Here is the fundamental convergence theorem for sequences of (not necessarily nonnegative) integrable functions:

Theorem 3.44. Let $f_n, f: X \rightarrow [-\infty, \infty]$ be a sequence of measurable functions, such that $f_n(x) \rightarrow f(x)$ for all $x \in X$, and such that there exists an integrable function $g: X \rightarrow [-\infty, \infty]$ satisfying $|f_n(x)| \leq g(x)$ for all $x \in X$. Then

$$\int_X |f_n - f| d\mu \rightarrow 0 \text{ as } n \rightarrow \infty.$$

Note that the hypothesis implies that $|f(x)| \leq g(x)$ for all x , and that $f_n, f \in L^1(\mu)$. The conclusion implies that (and is stronger than the statement that), for all $E \in \mathcal{M}$, $\int_E f_n d\mu \rightarrow \int_E f d\mu$.

Proof. Since $g \in L^1(\mu)$, $2g \in L^1(\mu)$, and we have, for each $x \in X$,

$$|f_n(x) - f(x)| \leq |f_n(x)| + |f(x)| \leq 2g(x).$$

Hence, for each n , $2g - |f_n - f|$ is a nonnegative measurable function. By hypothesis, for each $x \in X$, $\lim_{n \rightarrow \infty} (2g(x) - |f_n(x) - f(x)|) = 2g(x)$. Hence Fatou's Lemma gives

$$\int_X 2g d\mu \leq \liminf \int_X (2g - |f_n - f|) d\mu.$$

But the right hand side is $\int_X 2g d\mu - \overline{\lim} \int_X |f_n - f| d\mu$, so we in fact have

$$\overline{\lim} \int_X |f_n - f| d\mu \leq 0.$$

So since each $\int_X |f_n - f| d\mu$ is nonnegative, it follows that

$$\lim_{n \rightarrow \infty} \int_X |f_n - f| d\mu = 0.$$

□

3.6. Sets of measure zero. If (X, \mathcal{M}, μ) is a measure space, the sets E in \mathcal{M} such that $\mu(E) = 0$ are thought of as being negligible. Accordingly if $f, g: X \rightarrow [-\infty, \infty]$ are two functions, we say “ $f = g$ almost everywhere” (abbreviated “a.e.”) if $\mu(\{x \in X | f(x) \neq g(x)\}) = 0$. More generally, if we have some statement $P(x)$ for each $x \in X$, we say that P holds a.e. if $\{x \in X | P(x) \text{ is false}\}$ is a set of measure zero.

Recall from Problem Set 7 that (X, \mathcal{M}, μ) is called complete if whenever $A \in \mathcal{M}$ with $\mu(A) = 0$ and $B \subset A$ we have $B \in \mathcal{M}$ (which implies that $\mu(B) = 0$). Intuition would suggest that, if a set is known to be “negligible,” then we could conclude that any subset of that set is also negligible; the measure space is complete if and only if this intuition is valid. Now a problem on Problem Set 7 shows that \mathcal{M} can be enlarged and μ can be extended in a fairly straightforward way so as to make the measure space complete, as a result of which there's typically no harm in restricting attention to complete measure spaces, though some naturally occurring measure spaces aren't complete, notably the case where $X = \mathbb{R}$, \mathcal{M} is the Borel σ -algebra, and μ is Lebesgue measure (we won't prove this; it's not obvious). If instead \mathcal{M} is the σ -algebra of Lebesgue measurable subsets, or more generally if (X, \mathcal{M}, μ) is obtained as in Section 3.2 from an outer measure, then (X, \mathcal{M}, μ) is complete (why?).

Proposition 3.45. If (X, \mathcal{M}, μ) is any measure space and $h: X \rightarrow [-\infty, \infty]$ is measurable, then $h = 0$ a.e. if and only if for all $E \in \mathcal{M}$, $\int_E h d\mu = 0$.

Proof. \Rightarrow : If $h = 0$ a.e., then since $h^+ = \max\{h, 0\}$ and $h^- = \max\{-h, 0\}$, we have $h^+ = 0$ a.e. and $h^- = 0$ a.e. If s is a simple function such that $0 \leq s \leq h^+$ or $0 \leq s \leq h^-$, the set on which s is nonzero then has measure zero, so from the definition of the integral of a simple function we get that $\int_E s d\mu = 0$. So the definition of the integral of a nonnegative measurable function shows that $\int_E h^+ d\mu = \int_E h^- d\mu = 0$. So $\int_E h d\mu = 0$.

\Leftarrow : Let $A_n = \{x \in X | h(x) \geq 1/n\}$ and $B_n = \{x \in X | h(x) \leq -1/n\}$. Since h is measurable, we have $A_n, B_n \in \mathcal{M}$. So

$$0 = \int_{A_n} h d\mu \geq \frac{\mu(A_n)}{n}$$

and

$$0 = \int_{B_n} (-h) d\mu \geq \frac{\mu(B_n)}{n}.$$

Thus $\mu(A_n) = \mu(B_n) = 0$ for each n . So

$$\{x \in X | h(x) \neq 0\} = \cup_{n=1}^{\infty} (A_n \cup B_n)$$

is a countable union of sets of measure zero and so has measure zero. \square

Corollary 3.46. *If $f, g: X \rightarrow [-\infty, \infty]$ are measurable then $f = g$ a.e. if and only if, for all $E \in \mathcal{M}$, we have $\int_E f d\mu = \int_E g d\mu$.*

Proof. Apply the previous proposition with $h = f - g$. \square

Remark 3.47. *A quick examination of the proofs of the Monotone Convergence and Dominated Convergence Theorems shows that their conclusions still hold if instead of assuming that $f_n(x) \rightarrow f(x)$ for all $x \in X$ we only assume that $f_n(x) \rightarrow f(x)$ for almost every $x \in X$.*

3.7. Notions of convergence. For an arbitrary measure space (X, \mathcal{M}, μ) , here are four different criteria for the convergence of a sequence of measurable functions $\{f_n\}_{n=1}^{\infty}$ to some other measurable function f :

- (1) Convergence almost everywhere (a.e.): For almost every $x \in X$, we have $f_n(x) \rightarrow f(x)$.
- (2) Almost uniform (a.u.) convergence: For every $\delta > 0$ there is a set $A \in \mathcal{M}$ with $\mu(A) \leq \delta$ such that $f_n \rightarrow f$ uniformly on $X \setminus A$ (i.e., for each $\epsilon > 0$ there is N such that if $n \geq N$ then $|f_n(x) - f(x)| < \epsilon$ for all $x \in X \setminus A$).
- (3) Convergence in the mean of order p ($p \geq 1$ is a real number; this is also called “convergence in L^p ,” or, when $p = 1$, just “convergence in the mean”):

$$\lim_{n \rightarrow \infty} \int_X |f_n - f|^p d\mu = 0.$$

- (4) Convergence in measure (known in probability theory (in the case that μ is a probability measure, i.e., $\mu(X) = 1$) as “convergence in probability”): For all $\epsilon > 0$, we have

$$\lim_{n \rightarrow \infty} \mu(\{x \in X | |f_n(x) - f(x)| \geq \epsilon\}) = 0.$$

(The criterion is sometimes expressed as follows, which you should convince yourself is equivalent: For any $\epsilon > 0$ there is N such that if $n \geq N$ then

$$\mu(\{x \in X | |f_n(x) - f(x)| \geq \epsilon\}) < \epsilon.)$$

As we'll see, none of these four criteria is equivalent to any of the others in full generality; this subsection will be concerned with circumstances under which one of them implies another. Note that there are 12 potential implications to check here. The following three examples (all for the case where $X = \mathbb{R}$, μ is Lebesgue measure, and \mathcal{M} is the collection of Lebesgue-measurable subsets) are useful ones to keep in mind when considering different notions of convergence; in particular, between them, they rule out 9 of the 12 possible implications.

Example 3.48. Let $f_n = n\chi_{[0,1/n]}$ (we saw this example earlier in the discussion of Fatou's Lemma). Then $f_n \rightarrow 0$ almost everywhere, almost uniformly, and in measure, as can be verified directly from the definitions. But $\int_{\mathbb{R}} |f_n - 0|^p dm = n^{p-1}$ does not converge to zero (as we're assuming $p \geq 1$), so f_n does not converge to f in L^p for any $p \geq 1$. Thus criterion 3 above is not implied by any of the others.

Example 3.49. Let $f_n = \chi_{[n,n+1]}$. Then $f_n(x) \rightarrow 0$ for all x , so $f_n \rightarrow 0$ a.e. But since large values of x require passing to correspondingly large values of n before $f_n(x)$ becomes zero, it is not the case that $f_n \rightarrow 0$ a.u., and since $m(\{x \in X \mid |f_n(x) - f(x)| \geq 1\}) = 1$ for all n , it is not the case that $f_n \rightarrow f$ in measure. Also, $\int_{\mathbb{R}} |f_n - 0|^p dm = 1$ for all n . Thus criterion 1 above does not imply any of the others.

Example 3.50. Let $f_1 = \chi_{[0,1]}$, $f_2 = \chi_{[0,1/2]}$, $f_3 = \chi_{[1/2,1]}$, and in general, expressing an arbitrary natural number in the form $2^k + m$ where $0 \leq m < 2^k$, $f_{2^k+m} = \chi_{[m2^{-k},(m+1)2^{-k}]}$. Now if $0 < \epsilon < 1$ and $0 \leq m < 2^{-k}$, $\{x \mid |f_{2^k+m}(x) - 0| \geq \epsilon\}$ is a single interval of length 2^{-k} , so for any n we have $\mu(\{x \mid |f_n(x) - 0| \geq \epsilon\}) < 2/n$. Thus $f_n \rightarrow 0$ in measure. Further $\int_{\mathbb{R}} |f_n - 0|^p dm < 2/n \rightarrow 0$, so $f_n \rightarrow 0$ in L^p for any $p \geq 1$. But for every $x \in [0, 1]$ there are arbitrarily large n (depending on x) such that $f_n(x) = 1$, so the set on which $f_n(x)$ fails to converge to 0 has measure 1; thus f_n does not converge to 0 either a.e. or a.u. So neither of conditions 3 or 4 implies either of conditions 1 or 2.

Note that in examples 3.48 or 3.50 we could have used $X = [0, 1]$ instead of $X = \mathbb{R}$ to get an example of a probability measure for which the relevant implications fail; example 3.49, on the other hand, depended on using the whole positive real line. Egoroff's Theorem, proven below, will show that this isn't a coincidence.

Having shown that 9 of the 12 possible implications sometimes don't hold, we'll now prove that the remaining three always do.

Proposition 3.51. *If $f_n \rightarrow f$ almost uniformly then $f_n \rightarrow f$ in measure.*

Proof. This is really just a matter of unpacking the definitions. Let $\epsilon > 0$, and $\delta > 0$. To prove convergence in measure, we just need to show that there is N such that if $n \geq N$ then $\mu(\{x \in X \mid |f_n(x) - f(x)| \geq \epsilon\}) < \delta$. But almost uniform convergence shows that there is $A \in \mathcal{M}$ such that $f_n \rightarrow f$ uniformly on $X \setminus A$ and $\mu(A) < \delta$; hence there is N such that if $n \geq N$ then $|f_n(x) - f(x)| < \epsilon$ for all $x \in X \setminus A$. Thus for $n \geq N$ we have

$$\{x \in X \mid |f_n(x) - f(x)| \geq \epsilon\} \subset A,$$

so since $\mu(A) < \delta$ the conclusion follows. (Note the sense in which a.u. convergence is stronger than convergence in measure: in a.u. convergence the set A can be taken independent of the choice of $n \geq N$, whereas convergence in measure allows

the sets $\{x \in X \mid |f_n(x) - f(x)| \geq \epsilon\}$ to vary as long as they remain small, as is illustrated by Example 3.50). \square

Proposition 3.52. *If $f_n \rightarrow f$ in L^p for some $p \geq 1$ then $f_n \rightarrow f$ in measure.*

Proof. If $\epsilon > 0$, let $A_n = \{x \in X \mid |f_n(x) - f(x)| \geq \epsilon\}$; we need to show that $\mu(A_n) \rightarrow 0$. But we have

$$\epsilon^p \mu(A_n) \leq \int_{A_n} |f_n - f|^p d\mu \leq \int_X |f_n - f|^p d\mu \rightarrow 0 \text{ as } n \rightarrow \infty$$

by hypothesis, so since $\epsilon^p > 0$ is independent of n we must indeed have $\mu(A_n) \rightarrow 0$. \square

Proposition 3.53. *If $f_n \rightarrow f$ almost uniformly then $f_n \rightarrow f$ almost everywhere.*

Proof. We know that for any m there is $A_m \in \mathcal{M}$ such that $\mu(A_m) < 1/m$ and $f_n \rightarrow f$ uniformly on $X \setminus A_m$. So if $x \in X$ is such that $f_n(x)$ fails to converge to $f(x)$, it must be that $x \notin X \setminus A_m$ for all m , i.e., $x \in \bigcap_{m=1}^{\infty} A_m$. But $\mu(\bigcap_{m=1}^{\infty} A_m) \leq \mu(A_m) < 1/m$ for all m , so $\mu(\bigcap_{m=1}^{\infty} A_m) = 0$. Thus the set of points $x \in X$ such that $f_n(x)$ fails to converge to $f(x)$ has measure zero, i.e., $f_n \rightarrow f$ a.e. \square

We've now considered all twelve implications; the three that we proved each had fairly easy proofs at least once we understood the definitions reasonably well. Somewhat more subtly, some of the implications that fail in general do hold under weaker hypotheses. One example of this is:

Theorem 3.54 (Egoroff's Theorem). *If $\mu(X) < \infty$ and $f_n \rightarrow f$ almost everywhere then $f_n \rightarrow f$ almost uniformly.*

Thus to a probabilist (who only deals with the case $\mu(X) = 1$) almost everywhere convergence and almost uniform convergence are completely equivalent.

Proof. For any $\epsilon > 0$, define

$$E_N(\epsilon) = \{x \in X \mid (\exists n \geq N)(|f_n(x) - f(x)| \geq \epsilon)\}.$$

If $x \in X$ is such that $f_n(x) \rightarrow f(x)$, then we have $x \notin E_N(\epsilon)$ for some N . Thus $\bigcap_{N=1}^{\infty} E_N(\epsilon)$ is contained in the set of x at which the f_n do not converge to f , and so has measure zero by hypothesis. Further, it's obviously the case that $E_{N+1}(\epsilon) \subset E_N(\epsilon)$ for each N . Thus, for each $\epsilon > 0$, by Proposition 3.21 and the fact that $\mu(X) < \infty$, we have

$$\lim_{N \rightarrow \infty} \mu(E_N(\epsilon)) = 0.$$

Accordingly, given a real number $\delta > 0$ and a natural number $k \geq 1$, we can find N_k such that

$$\mu(E_{N_k}(1/k)) < \delta 2^{-k}.$$

Let

$$A = \bigcup_{k=1}^{\infty} E_{N_k}(1/k),$$

so that $\mu(A) < \delta$. The theorem will be proven when we show that $f_n \rightarrow f$ uniformly on $X \setminus A$. Indeed, if $\epsilon > 0$, we can find k such that $1/k < \epsilon$, and then if $n \geq N_k$ and $x \in X \setminus A$, we in particular have $x \notin E_{N_k}(1/k)$ and so $|f_n(x) - f(x)| < 1/k < \epsilon$. So since N_k is chosen independently of $x \in X \setminus A$, we indeed have $f_n \rightarrow f$ uniformly on $X \setminus A$. \square

The Dominated Convergence Theorem gives a situation in which almost everywhere convergence implies convergence in the mean (of order 1); only slightly changing the proof we get a statement for general L^p convergence.

Theorem 3.55. *If $\{f_n\}_{n=1}^\infty$ is a sequence of functions such that $f_n \rightarrow f$ a.e. and, for some measurable $g: X \rightarrow [-\infty, \infty]$ we have $\int_X |g|^p d\mu < \infty$ and $|f_n(x)| \leq g(x)$ for all x , then $f_n \rightarrow f$ in L^p .*

Proof. First, redefine g to equal ∞ on the set of measure zero on which $f_n(x)$ fails to converge to $f(x)$; this certainly doesn't change the condition $|f_n| \leq g$, and doesn't change the integral of $|g|^p$ since we've only changed the function on a set of measure zero. With this change made, we then have $|f(x)| \leq g(x)$ for all x , and so $|f_n(x) - f(x)|^p \leq 2^p g^p$ for all x . We then apply Fatou's Lemma to the sequence of (nonnegative) functions $2^p g^p - |f_n - f|^p$ (which converge pointwise to $2^p g^p$); the exact same argument as is given in the proof of the Dominated Convergence Theorem yields the result. \square

We close this investigation of notions of convergence with the following, which serves as a counterpoint to Example 3.50.

Theorem 3.56. *If $f_n \rightarrow f$ in measure, then there is a subsequence $\{f_{n_k}\}_{k=1}^\infty$ such that $f_{n_k} \rightarrow f$ almost uniformly.*

Proof. We first "speed up" the convergence by choosing numbers n_k with the property that, if $n \geq n_k$, then

$$\mu(\{x \in X \mid |f_n(x) - f(x)| \geq 2^{-k}\}) < 2^{-k}.$$

Let

$$E_k = \{x \in X \mid |f_{n_k}(x) - f(x)| \geq 2^{-k}\}.$$

For any m , let $G_m = \cup_{k=m}^\infty E_k$. Then

$$\mu(G_m) \leq \sum_{k=m}^\infty \mu(E_k) < \sum_{k=m}^\infty 2^{-k} = 2^{-m+1}.$$

Given $\delta > 0$, choose m so large that $2^{-m+1} < \delta$; thus $\mu(G_m) < \delta$. Then if $\epsilon > 0$, choose K so large that both $K \geq m$ and $2^{-K} < \epsilon$. Then for any $x \in X \setminus G_m$, if $k \geq K$ we have $k \geq m$, so $x \notin E_k$, in view of which

$$|f_{n_k}(x) - f(x)| < 2^{-k} \leq 2^{-K} < \epsilon.$$

So since K was chosen independently of $x \in X \setminus G_m$, we've shown that $f_{n_k} \rightarrow f$ uniformly on $X \setminus G_m$. δ was arbitrary, so this proves almost uniform convergence. \square

Of course, since convergence in L^p implies convergence in measure, and since almost uniform convergence implies a.e. convergence, this shows that any L^p -convergent sequence has a subsequence that converges almost everywhere. We'll see an independent proof of a closely related fact shortly.

3.8. L^p spaces. The use of the terminology “convergence in L^p ” would seem to suggest that there is a *space* called L^p in which one can talk about convergence. This is indeed the case, and gives some interesting examples of some of the phenomena that were discussed at the very start of this course.

We continue to work in an arbitrary measure space (X, \mathcal{M}, μ) . We may tacitly assume that the measure space is complete (so that we can infer that a set has measure zero as soon as we show that it is a subset of a set of measure zero, without showing that the first of these sets is measurable); however, at least in most cases, with a little more work we would be able to show directly that the set is indeed measurable, making the completeness assumption superfluous.

For any real number $p \in [1, \infty)$, and any measurable $f: X \rightarrow [-\infty, \infty]$, we define

$$\|f\|_p = \left(\int_X |f|^p d\mu \right)^{1/p}.$$

Note that if $\{f_n\}_{n=1}^\infty$ is a sequence of measurable functions, then $f_n \rightarrow f$ in L^p if and only if $\|f_n - f\|_p \rightarrow 0$.

Perhaps a bit surprisingly, there is a natural generalization of this to the case that $p = \infty$. Namely, we define

$$\|f\|_\infty = \inf\{M \geq 0 \mid |f(x)| \leq M \text{ for almost every } x \in X\}.$$

The right hand side above is sometimes called the “essential supremum” of X . By definition, if n is any positive integer we have $|f(x)| \leq \|f\|_\infty + 1/n$ for a.e. x ; in other words, $\mu(|f|^{-1}(\|f\|_\infty + 1/n, \infty)) = 0$. But

$$|f|^{-1}(\|f\|_\infty, \infty) = \cup_{n=1}^\infty |f|^{-1}(\|f\|_\infty + 1/n, \infty),$$

so in fact $|f(x)| \leq \|f\|_\infty$ for almost every x . Thus $\|f\|_\infty$ is the smallest (extended real) number M such that $|f|$ is almost everywhere bounded by M (in particular, the above argument shows that such a smallest number exists).

Definition 3.57. If $p \in [1, \infty]$, we define

$$L^p(\mu) = \{f: X \rightarrow [-\infty, \infty] \mid f \text{ is measurable and } \|f\|_p < \infty\}.$$

Remark 3.58. (1) *In the case $p = 1$, this is consistent with our previous definition. (Namely, $L^1(\mu)$ is the space of integrable functions.)*

(2) *For any p , we have $\|f\|_p = 0$ if and only if $f = 0$ a.e.*

(3) *As mentioned earlier, if $1 \leq p < \infty$, $f_n \rightarrow f$ in L^p if and only if $\|f_n - f\|_p \rightarrow 0$. As for the case $p = \infty$, it's not too hard to see that $\|f_n - f\|_\infty \rightarrow 0$ if and only if there is a set E of measure zero such that $f_n \rightarrow f$ uniformly on $X \setminus E$. (Take for E the union over n of the sets $\{x \mid |f_n(x) - f(x)| > \|f_n - f\|_\infty\}$; details are left to the reader.)*

(4) *If $c \in \mathbb{R}$ and $f \in L^p(\mu)$, we clearly have $\|cf\|_p = |c|\|f\|_p$.*

Now we certainly could have made these definitions for $0 < p < 1$ as well; the main reason for sticking to the case $p \geq 1$ is that doing so is required for the following to hold:

Theorem 3.59 (Minkowski inequality). *If $f, g \in L^p(\mu)$ ($1 \leq p \leq \infty$), then*

$$\|f + g\|_p \leq \|f\|_p + \|g\|_p.$$

(Of course, the same inequality holds trivially when one of f or g is not in $L^p(\mu)$, since then the right hand side is infinite).

There are two easy cases of this theorem, namely when $p = 1$ and when $p = \infty$. Indeed, for $p = 1$, the inequality just says that $\int_X |f + g| d\mu \leq \int_X |f| d\mu + \int_X |g| d\mu$, a fact which follows immediately from the triangle inequality for \mathbb{R} , since for all $x \in X$ we have $|f(x) + g(x)| \leq |f(x)| + |g(x)|$. As for the case $p = \infty$, we know that there are sets $E_1, E_2 \in \mathcal{M}$ with $\mu(E_1) = \mu(E_2)$ such that if $x \notin E_1$ then $|f(x)| \leq \|f\|_\infty$ and $|g(x)| \leq \|g\|_\infty$. So if $E = E_1 \cup E_2$, then $\mu(E) = 0$ and if $x \notin E$ then $|f(x) + g(x)| \leq |f(x)| + |g(x)| \leq \|f\|_\infty + \|g\|_\infty$. So since $\|f + g\|_\infty$ is the *smallest* number which gives an almost everywhere bound for $|f + g|$, it must be that $\|f + g\|_\infty \leq \|f\|_\infty + \|g\|_\infty$.

It remains to prove the Minkowski inequality for $1 < p < \infty$. This is subtle, and along the way we'll prove another important inequality for integrals (namely the Hölder inequality). However, most of the subtleties involve tricky arithmetic facts about the real numbers, rather than deep theoretical issues. The first of these facts is:

Lemma 3.60. *If $x, y, \alpha \in \mathbb{R}$, $x, y \geq 0$, and $0 < \alpha < 1$, then*

$$x^\alpha y^{1-\alpha} \leq \alpha x + (1 - \alpha)y.$$

Proof. If either $x = 0$ or $y = 0$ the result is trivial, so let us assume that x and y are both positive. So we can set $w = \log x$, $z = \log y$, and we will be done as soon as we show that

$$(9) \quad e^{\alpha w + (1-\alpha)z} \leq \alpha e^w + (1 - \alpha)e^z,$$

for any $w, z \in \mathbb{R}$ (of course we can assume $w < z$ at the cost of replacing α by $1 - \alpha$). That this holds is an expression of the *convexity* of the exponential function. Let $f(t) = e^t$. We claim that the (obvious) fact that f' is an increasing function is enough to imply (9). Namely, let $u = \alpha w + (1 - \alpha)z$, so $w < u < z$. The mean value theorem gives $s \in (w, u)$, $t \in (u, z)$ such that

$$f'(s) = \frac{f(u) - f(w)}{u - w}, \quad f'(t) = \frac{f(z) - f(u)}{z - u}.$$

Since f' is increasing, $f'(s) \leq f'(t)$. Now $u - w = (1 - \alpha)(z - w)$ and $z - u = \alpha(z - w)$. So we obtain

$$\frac{f(u) - f(w)}{(1 - \alpha)(z - w)} \leq \frac{f(z) - f(u)}{\alpha(z - w)},$$

so

$$\alpha(f(u) - f(w)) \leq (1 - \alpha)(f(z) - f(u)),$$

i.e.,

$$f(u) \leq \alpha f(w) + (1 - \alpha)f(z),$$

which is precisely (9). □

Corollary 3.61 (Young's inequality). *If $a, b \in \mathbb{R}$ are nonnegative, if $1 < p < \infty$, and if $q \in (1, \infty)$ is the unique number such that $\frac{1}{p} + \frac{1}{q} = 1$, then*

$$ab \leq \frac{a^p}{p} + \frac{b^q}{q}.$$

Proof. This is just the previous lemma with $\alpha = 1/p$, $x = a^{1/\alpha}$, and $y = b^{1/(1-\alpha)}$. □

Remark 3.62. In general if $p \in [1, \infty]$, the dual exponent to p is by definition the number q satisfying $\frac{1}{p} + \frac{1}{q} = 1$, with the convention that this means that if $p = 1$ then $q = \infty$, and vice versa. Note that the only case here where $p = q$ is where $p = 2$; in that case, Young's inequality has a very simple proof (can you find it?).

Theorem 3.63 (Hölder inequality). If p and q are dual exponents ($\frac{1}{p} + \frac{1}{q} = 1$) and if $f, g: X \rightarrow [-\infty, \infty]$ are measurable functions then

$$\int_X |fg| d\mu = \|f\|_p \|g\|_q.$$

Proof. The case that either $p = 1$ and $q = \infty$ or vice versa is trivial, since (in the first case) we have $|fg| \leq |f| \|g\|_\infty$ a.e., and similarly for the second case. Accordingly let us assume that $1 < p < \infty$.

If either $\|f\|_p = 0$ or $\|g\|_q = 0$, then $fg = 0$ a.e., so that $\int_X |fg| d\mu = 0$ and the theorem holds. Meanwhile if $\|f\|_p$ and $\|g\|_q$ are both nonzero and if one or both of them is infinite then the right hand side is ∞ , so the inequality is trivial. So we can assume that $\|f\|_p$ and $\|g\|_q$ are both finite and nonzero. Let $F = f/\|f\|_p$ and $G = g/\|g\|_q$; thus $\|F\|_p = \|G\|_q = 1$. Then, using Young's inequality,

$$\int_X |FG| d\mu \leq \frac{1}{p} \int_X |F|^p d\mu + \frac{1}{q} \int_X |G|^q d\mu = \frac{1}{p} + \frac{1}{q} = 1.$$

Hence

$$\int_X |fg| d\mu = \|f\|_p \|g\|_q \int_X |FG| d\mu \leq \|f\|_p \|g\|_q,$$

as desired. □

Proof of the Minkowski inequality. We've already proven the result when $p = 1$ and when $p = \infty$, so assume that $1 < p < \infty$. Let q be the dual exponent to p . Note that the equation $\frac{1}{p} + \frac{1}{q} = 1$ is equivalent to $(p-1)q = p$.

The argument will involve dividing by a power of $\|f + g\|_p$, so first we need to reduce to the case where $\|f + g\|_p$ is finite and nonzero. Of course, when $\|f + g\|_p = 0$, the inequality is trivial. On the other hand, note that (using the fact that $p \geq 1$), we have $\left(\frac{|f(x)+g(x)|}{2}\right)^p \leq \frac{1}{2}|f(x)|^p + \frac{1}{2}|g(x)|^p$ (for instance, this follows from the convexity of the function $t \mapsto t^p$, as in the proof of Lemma 3.60). Hence $\int_X |f + g|^p d\mu \leq 2^{p-1} (\int_X |f|^p d\mu + \int_X |g|^p d\mu)$, so that if $\|f + g\|_p = \infty$ then one of $\|f\|_p$ or $\|g\|_p$ is infinite, and the Minkowski inequality holds (both sides are ∞). This reduces us to the case that $\|f + g\|_p \in (0, \infty)$.

Now the Hölder inequality gives (using $(p-1)q = p$)

$$\int_X |f| |f + g|^{p-1} d\mu \leq \|f\|_p \left(\int_X |f + g|^{(p-1)q} d\mu \right)^{1/q} = \|f\|_p \|f + g\|_p^{p/q}.$$

Likewise

$$\int_X |g| |f + g|^{p-1} d\mu \leq \|g\|_p \|f + g\|_p^{p/q}.$$

So

$$\|f + g\|_p^p \leq \int_X (|f| + |g|) |f + g|^{p-1} d\mu \leq (\|f\|_p + \|g\|_p) \|f + g\|_p^{p/q}.$$

Now divide by (the finite, positive) $\|f + g\|_p^{p/q}$ and notice that $p - p/q = 1$ to deduce the result. □

The fact that the sum of two functions in $L^p(\mu)$ is again in $L^p(\mu)$ (as follows from the Minkowski inequality, and indeed was proven en route to it), together with the obvious fact that $L^p(\mu)$ is closed under scalar multiplication, shows that $L^p(\mu)$ is a vector space. We thus have a function $\|\cdot\|_p: L^p(\mu) \rightarrow \mathbb{R}$ satisfying (i) $\|f\|_p \geq 0$ for all $f \in L^p(\mu)$, (ii) $\|cf\|_p = |c|\|f\|_p$ for $c \in \mathbb{R}$, $f \in L^p(\mu)$, and the triangle inequality (iii) $\|f + g\|_p \leq \|f\|_p + \|g\|_p$ for all $f, g \in L^p(\mu)$. These facts are summarized in saying that $\|\cdot\|_p$ is a *pseudonorm* on $L^p(\mu)$. A *norm* $\|\cdot\|$ on a vector space V is a function $\|\cdot\|: V \rightarrow \mathbb{R}$ satisfying the analogues of (i),(ii),(iii) for elements of V together with the additional condition that if $\|v\| = 0$ then $v = 0$. Since $\|f\|_p = 0$ if $f(x) = 0$ for a.e. x , it is *not* the case that $\|\cdot\|_p$ defines a norm on $L^p(\mu)$, except in the case every nonempty measurable subset of X has positive measure (since in this latter case the only function which is zero a.e. is the zero function).

Note that if $\|\cdot\|$ is a norm on V , then setting $d(v, w) = \|v - w\|$ defines a metric on V .

Example 3.64. Let (X, \mathcal{M}, μ) be given as follows. $X = \{1, \dots, n\}$, \mathcal{M} consists of all subsets of X , and for $A \subset X$ $\mu(A)$ is the number of elements of A . Then any function $f: X \rightarrow \mathbb{R}$ is measurable (since every subset of X is measurable). Given $f: X \rightarrow \mathbb{R}$, let $x_i = f(i)$ for $i = 1, \dots, n$. This gives an identification of the set of measurable functions $f: X \rightarrow \mathbb{R}$ with \mathbb{R}^n . We have $\|f\|_p = (\sum_{i=1}^n |x_i|^p)^{1/p}$ for $1 \leq p < \infty$, and $\|f\|_\infty = \max_{1 \leq i \leq n} |x_i|$. Every nonempty subset of X has positive measure, so the Minkowski inequality justifies the assertion, made at the very start of the course, that

$$d_p(\vec{x}, \vec{y}) = \left(\sum_{i=1}^n |x_i - y_i|^p \right)^{1/p}$$

gives a metric on \mathbb{R}^n .

Although $(L^p(\mu), \|\cdot\|_p)$ in most cases is not a normed vector space due to the existence of nonzero functions which are zero almost everywhere, we can get a normed vector space out of it by the following simple algebraic construction. Notice that where $Z = \{f \in L^p(\mu) | f = 0 \text{ a.e.}\}$, Z is a vector subspace of $L^p(\mu)$ (since $f + g$ is nonzero on a subset of the union of the set on which f is nonzero with that on which g is nonzero, and the union of two sets of measure zero has measure zero). So one can form the quotient vector space

$$\mathcal{L}^p(\mu) = \frac{L^p(\mu)}{Z}.$$

Thus the elements of $\mathcal{L}^p(\mu)$ are *equivalence classes* of functions in $L^p(\mu)$, where $f, g \in L^p(\mu)$ are deemed equivalent iff $f - g \in Z$, i.e., iff $f = g$ a.e. Write $[f]$ for the equivalence class of $f \in L^p(\mu)$ (so $[f] \in \mathcal{L}^p(\mu)$). Since if $f = g$ a.e. then $\|f\|_p = \|g\|_p$, we get a well-defined function $\|\cdot\|_p: \mathcal{L}^p(\mu) \rightarrow \mathbb{R}$ by setting $\|[f]\|_p = \|f\|_p$. What we've shown immediately implies that $\|\cdot\|_p$ is a norm on $\mathcal{L}^p(\mu)$ (if this isn't clear, you should convince yourself of it).

In particular, defining $d_p([f], [g]) = \|f - g\|_p$ gives a metric on $\mathcal{L}^p(\mu)$. An important property of this metric space is:

Theorem 3.65. For any $p \in [1, \infty]$, $(\mathcal{L}^p(\mu), d_p)$ is a complete metric space.

Proof. The case that $p = \infty$ will be left as an exercise on Problem Set 9. So we assume that $1 \leq p < \infty$.

First we should translate the statement of the theorem into the language of functions, rather than equivalence classes of functions. We are to show that if $\{[f_n]\}_{n=1}^\infty$ is a Cauchy sequence in $\mathcal{L}^p(\mu)$, then there is $[f] \in \mathcal{L}^p(\mu)$ such that $d_p([f_n], [f]) \rightarrow 0$ as $n \rightarrow \infty$. Now (as the notation suggests), the elements $[f_n], [f]$ of $\mathcal{L}^p(\mu)$ are equivalence classes of certain functions $f_n, f \in L^p(\mu)$; the statement that $\{[f_n]\}_{n=1}^\infty$ is a Cauchy sequence amounts to the statement that if $\epsilon > 0$ there is N such that whenever $n, m \geq N$ we have $\|f_n - f_m\|_p < \epsilon$; and the statement that $d_p([f_n], [f]) \rightarrow 0$ amounts to the statement that $\|f_n - f\|_p \rightarrow 0$ (i.e., $f_n \rightarrow f$ in L^p). Note also that since if a Cauchy sequence has a convergent subsequence, then the entire sequence converges, it's enough to find a subsequence $\{f_{n_k}\}_{k=1}^\infty$ of $\{f_n\}_{n=1}^\infty$ such that $f_{n_k} \rightarrow f$ in L^p (for some f yet to be determined).

So let $\{f_n\}_{n=1}^\infty$ be a sequence in $L^p(\mu)$ such that $\{[f_n]\}_{n=1}^\infty$ is a Cauchy sequence in $\mathcal{L}^p(\mu)$. For any integer k , we can then find n_k such that if $n, m \geq n_k$ then $\|f_n - f_m\|_p < 2^{-k}$, and such that $n_{k+1} > n_k$. For any k , define

$$g_k = |f_{n_1}| + \sum_{j=1}^{k-1} |f_{n_{j+1}} - f_{n_j}|,$$

and set

$$g(x) = \lim_{k \rightarrow \infty} g_k(x) = |f_{n_1}(x)| + \sum_{j=1}^{\infty} |f_{n_{j+1}}(x) - f_{n_j}(x)|.$$

Thus $|g_k|^p \nearrow |g|^p$. Hence the monotone convergence theorem gives

$$\int_X |g|^p d\mu = \lim_{k \rightarrow \infty} \int_X |g_k|^p d\mu = \lim_{k \rightarrow \infty} \|g_k\|_p^p.$$

But for each k we have, by the Minkowski inequality,

$$\|g_k\|_p = \|f_{n_1}\|_p + \sum_{j=1}^{k-1} \|f_{n_{j+1}} - f_{n_j}\|_p < \|f_{n_1}\|_p + \sum_{j=1}^{k-1} 2^{-j} < \|f_{n_1}\|_p + 1.$$

Hence

$$\int_X |g|^p d\mu \leq (\|f_{n_1}\|_p + 1)^p < \infty.$$

This in particular implies that, where

$$E = \{x \in X | g(x) = \infty\}, \text{ we have } \mu(E) = 0.$$

Thus for every $x \in X \setminus E$, the sum $|f_{n_1}(x)| + \sum_{j=1}^{\infty} |f_{n_{j+1}}(x) - f_{n_j}(x)|$ converges to a finite real number. So the sum $f_{n_1}(x) + \sum_{j=1}^{\infty} (f_{n_{j+1}}(x) - f_{n_j}(x))$ is (absolutely) convergent to some real number, say $f(x)$, whenever $x \in X \setminus E$. For $x \in E$, set $f(x) = 0$. But notice that

$$f_{n_k}(x) = f_{n_1}(x) + \sum_{j=1}^{k-1} (f_{n_{j+1}}(x) - f_{n_j}(x));$$

hence

$$f_{n_k}(x) \rightarrow f(x) \text{ for all } x \in X \setminus E.$$

In particular, we have shown that $f_{n_k} \rightarrow f$ a.e. It remains to show that $f_{n_k} \rightarrow f$ in L^p . For $x \in X \setminus E$, note that

$$|f(x) - f_{n_k}(x)| = \left| \sum_{j=k}^{\infty} (f_{n_{j+1}}(x) - f_{n_j}(x)) \right| \leq \sum_{j=k}^{\infty} |f_{n_{j+1}}(x) - f_{n_j}(x)| = g(x) - g_k(x).$$

Hence (using that E has measure zero and so does not affect integrals)

$$\int_X |f - f_{n_k}|^p d\mu = \int_{X \setminus E} |f - f_{n_k}|^p \leq \int_{X \setminus E} |g - g_k|^p \rightarrow 0$$

as $k \rightarrow \infty$, where the fact that $|g - g_k|^p \rightarrow 0$ follows from the Dominated Convergence Theorem (as $|g - g_k|^p \leq |g|^p \in L^1(\mu)$). Thus $f_{n_k} \rightarrow f$ in L^p , i.e., $[f_{n_k}] \rightarrow [f]$, and so since $\{[f_n]\}_{n=1}^{\infty}$ is Cauchy it must be that $[f_n] \rightarrow [f]$, completing the proof. \square

Now let us restrict to the case where $X = \mathbb{R}$ and μ is the Lebesgue measure m on the σ -algebra \mathcal{M} of Lebesgue measurable subsets. The case $p = 1$ of the following is an exercise on Problem Set 8; the case for general p is proven in precisely the same way, so we omit the proof. (The existence of the constant M isn't asked for in the statement of the homework problem, but if you followed the outline in the hints to the problem the g that you produced will certainly have the stated property).

Proposition 3.66. *Let $1 \leq p < \infty$, and let $f \in L^p(m)$. Then for any $\epsilon > 0$ there is a continuous function $g: \mathbb{R} \rightarrow \mathbb{R}$ such that*

$$\int_X |f - g|^p d\mu < \epsilon.$$

Further, g may be chosen so that there is $M > 0$ such that $g(x) = 0$ whenever $|x| > M$.

In general if (M, d) is a metric space, a subset $D \subset M$ is called *dense* in M if for each $m \in M$ and $\epsilon > 0$ there is $a \in D$ such that $d(a, m) < \epsilon$. (Equivalently $D \subset M$ is dense if $\overline{D} = M$). For instance, the rational numbers are dense in the real numbers.

Let $C_c(\mathbb{R})$ be the set of continuous functions $f: \mathbb{R} \rightarrow \mathbb{R}$ such that there is some M such that $f(x) = 0$ whenever $|x| > M$. The above proposition looks like the statement that $C_c(\mathbb{R})$ is dense in $L^p(m)$, except of course that $L^p(m)$ isn't a metric space; rather, $\mathcal{L}^p(m)$ is.

Notice that if $f, g \in C_c(\mathbb{R})$, then since $\{x | f(x) \neq g(x)\} = (f - g)^{-1}(\mathbb{R} \setminus \{0\})$ is open by the continuity of $f - g$, if $f \neq g$ then $\{x | f(x) \neq g(x)\}$ contains an open interval of positive length, and hence has positive measure. Thus if $f, g \in C_c(\mathbb{R})$ and $f = g$ a.e., then in fact $f = g$. Certainly any $f \in C_c(\mathbb{R})$ lies in $L^p(m)$, since f is bounded by continuity, and zero except on a set of finite measure. Thus if $f \in C_c(\mathbb{R})$, there is one and only one equivalence class $[f] \in \mathcal{L}^p(m)$ containing f , and further f is the only continuous function contained in the equivalence class $[f]$. Hence $C_c(\mathbb{R})$ may be thought of as a subset of the metric space $\mathcal{L}^p(m)$, and the above proposition shows that, for $1 \leq p < \infty$, $C_c(\mathbb{R})$ is dense in $\mathcal{L}^p(m)$.

Now if $D \subset M$ is dense, any element of M is the limit of some Cauchy sequence in D (with respect to the metric d). So this gives us another interpretation of the space $\mathcal{L}^p(m)$. Since $\mathcal{L}^p(m)$ is *complete* and contains $C_c(\mathbb{R})$ as a dense subspace, we can think of $\mathcal{L}^p(m)$ as consisting of precisely the things that need to be added

to $C_c(\mathbb{R})$ in order to make every Cauchy sequence in $C_c(\mathbb{R})$ (with respect to the metric $d_p(f, g) = (\int_X |f - g|^p dm)^{1/p}$) converge; that is, $\mathcal{L}^p(m)$ is the completion of $C_c(\mathbb{R})$ with respect to d_p (if one sets up an appropriate definition of isomorphism of metric spaces, it's not hard to show that any two complete metric space containing a given metric space (D, d) as a dense subspace are isomorphic, justifying the use of the word "the"). As was mentioned long ago, one can always form a completion of a given metric space (such as $(C_c(\mathbb{R}), d_p)$) via an abstract construction involving equivalence classes of Cauchy sequences; however, in this case we've shown that the completion has a somewhat more explicit interpretation, as $\mathcal{L}^p(m)$.

The situation with $L^\infty(m)$ is a bit different. It's still true that $C_c(\mathbb{R})$ can be identified with a subset of the complete metric space $\mathcal{L}^\infty(m)$, but this subset isn't dense. Indeed, a sequence of continuous functions which converges in $L^\infty(m)$ actually converges uniformly, so its limit is continuous. One can show that the closure of $C_c(\mathbb{R})$ in $(\mathcal{L}^\infty(m), d_\infty)$ is

$$C_0(\mathbb{R}) = \{f: \mathbb{R} \rightarrow \mathbb{R} \mid f \text{ is continuous and } (\forall \epsilon > 0)(\exists M)(|x| > M \Rightarrow |f(x)| < \epsilon)\}$$

(which is also known as the set of continuous functions which "vanish at infinity").

4. FOURIER ANALYSIS

4.1. Convergence of Fourier series. Fourier analysis has its roots in Fourier's attempts to solve the heat equation. Consider a one-dimensional wire, and choose units of length so that the wire has length 2π and extends from $x = -\pi$ to $x = \pi$. Assume the ends of the wire are in thermal equilibrium with each other (or suppose that the wire is actually circular and x is the angular coordinate, so that $-\pi$ and π represent the same point). Newton's Law of Cooling asserts that the temperature $u(x, t)$ of the wire at position x and time t evolves according to the heat equation

$$(10) \quad \frac{\partial u}{\partial t} = k \frac{\partial^2 u}{\partial x^2};$$

here $k > 0$ is some constant reflecting the thermal conductivity of the material of which the wire is composed. For our purposes here, it's convenient to think of u as being a *complex-valued* function $u: [-\pi, \pi] \rightarrow \mathbb{C}$; of course this seems at odds with our statement that u was supposed to represent temperature, but notice that if $u = f + ig$ where f and g are real-valued functions, then u satisfies (10) if and only if both f and g do, so we can always go back to the physical situation by taking the real part of u .

Fourier's goal was as follows: assume that we have measured the wire's initial temperature to be $f(x)$ for some function f (necessarily with $f(-\pi) = f(\pi)$); then find a solution $u(x, t)$ for all positive t satisfying both (10) and the initial condition $u(x, 0) = f(x)$. For certain special choices of f , namely $f(x) = e^{inx}$ where n is an integer, Fourier noticed that such a solution was given by

$$u_n(x, t) = e^{-n^2 kt} e^{inx}.$$

Further, (10) is a *linear* equation (*i.e.*, linear combinations of solutions are still solutions), so initial data of the form $f(x) = \sum_{n=-N}^N a_n e^{inx}$ would give a solution

$$u(x, t) = \sum_{n=-N}^N a_n e^{-n^2 kt} e^{inx}.$$

Fourier's contention now was that any function $f: [-\pi, \pi] \rightarrow \mathbb{C}$ should be expressible as an infinite sum $f(x) = \sum_{n=-\infty}^{\infty} a_n e^{inx}$, in which case (at least if one makes some assumptions about interchanging limits and derivatives), solutions of the above form (extended to infinite sums) would solve the heat equation for any initial data.

Fourier didn't prove this contention (which actually isn't true in full generality, as will be discussed below), and in fact the desire to put his work on a rigorous foundation was historically a major stimulus for the development of analysis. Fourier did give a simple way of finding the constants a_n . Namely, if one assumes that

$$f(x) = \sum_{n=-\infty}^{\infty} a_n e^{inx}$$

(and if one assumes that limits and integrals can be interchanged; this would be valid if the partial sums $S_N f(x) = \sum_{n=-N}^N a_n e^{inx}$ converged to f uniformly, for example), then from the observation that $\int_{-\pi}^{\pi} e^{i(n-m)x} dx$ is equal to 2π when $n = m$ and 0 otherwise, one finds that, for all m ,

$$\int_{-\pi}^{\pi} f(x) e^{-imx} dx = 2\pi a_m.$$

We'll now discuss some results which partially justify Fourier's intuitions. The question to be studied is:

Question 4.1. *If $f: [-\pi, \pi] \rightarrow \mathbb{C}$ is a continuous function with $f(-\pi) = f(\pi)$, and if we set*

$$a_n = \frac{1}{2\pi} \int_{-\pi}^{\pi} f(x) e^{-inx} dx, \quad (S_N f)(x) = \sum_{n=-N}^N a_n e^{inx},$$

is it the case that $S_N f \rightarrow f$ as $N \rightarrow \infty$? If so, in what sense?

First let us rewrite the partial sum $S_N f$. We have

$$\begin{aligned} S_N f(x) &= \frac{1}{2\pi} \sum_{n=-N}^N \left(\int_{-\pi}^{\pi} f(y) e^{-iny} dy \right) e^{inx} \\ &= \int_{-\pi}^{\pi} \left(\frac{1}{2\pi} \sum_{n=-N}^N e^{in(x-y)} \right) f(y) dy = \int_{-\pi}^{\pi} D_N(x-y) f(y) dy, \end{aligned}$$

where we define the *Dirichlet kernel* D_N by

$$D_N(u) = \frac{1}{2\pi} \sum_{n=-N}^N e^{inu}.$$

Note that D_N is a periodic function of period 2π , and that $\int_{-\pi}^{\pi} D_N(x) dx = 1$. Now f is a function defined on $[-\pi, \pi]$ with $f(-\pi) = f(\pi)$; let us extend f to be a periodic function defined on all of \mathbb{R} by setting $f(a + 2k\pi) = f(a)$ whenever $a \in [-\pi, \pi]$ and k is an integer. We then have, making the substitution $u = x - y$

$$\begin{aligned} S_N f(x) &= \int_{-\pi}^{\pi} D_N(x-y) f(y) dy = \int_{x+\pi}^{x-\pi} -D_N(u) f(x-u) du = \int_{x-\pi}^{x+\pi} D_N(u) f(x-u) du \\ &= \int_{-\pi}^{\pi} D_N(u) f(x-u) du, \end{aligned}$$

where the final equality follows from the fact that $u \mapsto D_N(x-u)f(u)$ is periodic of period 2π , and therefore has the same integral over any interval of length 2π .

Now we record a convenient alternative formula for D_N .

Proposition 4.2.

$$D_N(u) = \frac{1}{2\pi} \frac{\sin((N+1/2)u)}{\sin(u/2)}.$$

(Note that L'Hôpital's rule shows that the right hand side extends continuously to $u=0$ as $D_N(0) = (2N+1)/2\pi$, consistently with the definition of D_N)

Proof. We have

$$\begin{aligned} 2\pi D_N(u) &= \sum_{n=-N}^N e^{inu} = e^{-iNu} \sum_{n=0}^{2N} e^{inu} = e^{-iNu} \frac{1 - e^{(2N+1)u}}{1 - e^{iu}} \\ &= \frac{e^{-i(N+1/2)u} (1 - e^{(2N+1)u})}{e^{-iu/2} (1 - e^{iu})} = \frac{-2i \sin((N+1/2)u)}{-2i \sin(u/2)} = \frac{\sin(N+1/2)u}{\sin(u/2)}, \end{aligned}$$

as desired. \square

Definition 4.3. If $I \subset \mathbb{R}$ is an interval, $k \geq 0$ is an integer, and $f: I \rightarrow \mathbb{C}$, f is said to be of class C^k if the function f is continuous and each of its first k derivatives $f', \dots, f^{(k)}$ exist and are continuous. f is said to be of class C^∞ if it is of class C^k for every nonnegative integer k .

Here is the first positive result about the convergence of $S_N f$ to f .

Theorem 4.4. Assume that the extended periodic function $f: \mathbb{R} \rightarrow \mathbb{C}$ is of class C^2 . Then $S_N f \rightarrow f$ uniformly as $N \rightarrow \infty$.

Proof. By the calculation above, and the fact that $\int_{-\pi}^{\pi} D_N(u) du = 1$, we have

$$\begin{aligned} S_N f(x) - f(x) &= \int_{-\pi}^{\pi} D_N(u) f(x-u) du - f(x) = \int_{-\pi}^{\pi} D_N(u) (f(x-u) - f(x)) du \\ &= \frac{1}{2\pi} \int_{-\pi}^{\pi} \frac{\sin((N+1/2)u)}{\sin(u/2)} (f(x-u) - f(x)) du. \end{aligned}$$

The theorem will follow fairly easily from the following lemma:

Lemma 4.5. Let $g_x(u) = \frac{f(x-u) - f(x)}{\sin(u/2)}$. Then g_x is of class C^1 , and there is a constant M , independent of x , such that $|g'_x(u)| \leq M$ for all x and u .

We'll prove the lemma below; let us first show how it implies the theorem. We have, integrating by parts,

$$\begin{aligned} S_N f(x) - f(x) &= \int_{-\pi}^{\pi} g_x(u) \sin(N+1/2)u du \\ &= - \left. \frac{g_x(u) \cos((N+1/2)u)}{N+1/2} \right|_{-\pi}^{\pi} + \int_{-\pi}^{\pi} g'_x(u) \frac{\cos((N+1/2)u)}{N+1/2} du \\ &= \frac{1}{N+1/2} \int_{-\pi}^{\pi} g'_x(u) \cos((N+1/2)u) du. \end{aligned}$$

Hence

$$|S_N f(x) - f(x)| \leq \frac{1}{N+1/2} \int_{-\pi}^{\pi} |g'_x(u)| |\cos((N+1/2)u)| du \leq \frac{2\pi M}{N+1/2}.$$

So if $\epsilon > 0$ if we choose N_0 so large that $2\pi M/(N_0 + 1/2) < \epsilon$ (note this N_0 is chosen independently of x), if $N \geq N_0$ we have $|S_N f(x) - f(x)| < \epsilon$. Thus $S_N f \rightarrow f$ uniformly. \square

Proof of Lemma 4.5. Write $h_x(u) = \frac{f(x-u)-f(x)}{u}$ and $p(u) = \frac{\sin(u/2)}{u}$ (we only consider values of u belonging to $[-\pi, \pi]$ throughout the proof). Note that $p(u)$ extends continuously to $u = 0$ as $p(0) = 1/2$ by L'Hôpital's rule, and the definition of the derivative shows that $h_x(u)$ extends continuously to $u = 0$ as $h_x(0) = -f'(x)$. We have $g_x = h_x/p$, so since p is strictly positive and bounded below on $[-\pi, \pi]$ the quotient rule implies that the result will be proven if we show that each of h_x, h'_x, p, p' are continuous and are bounded above and below independently of x (since we will then have $g'_x(u) = (ph'_x - p'h_x)/p^2$).

Since f, f' , and f'' are assumed continuous, there is M_0 such that $|f(y)|, |f'(y)|, |f''(y)| \leq M_0$ for all $y \in [-\pi, \pi]$; since f (and hence also f', f'') is periodic of period 2π , this bound also holds for all $y \in \mathbb{R}$. Taylor's theorem shows that, for any u , there is $\xi \in (x - u, x)$ such that

$$f(x - u) = f(x) - uf'(x) + \frac{u^2}{2}f''(\xi),$$

in view of which $h_x(u) = -f'(x) + uf''(\xi)/2$. So if $|u| \leq 1$ we get $|h_x(u)| \leq \frac{3M_0}{2}$, while if $|u| \geq 1$ we have $|h_x(u)| \leq 2M_0$ directly from the definition of h_x . Now

$$\frac{h_x(u) - h_x(0)}{u} = \frac{1}{u}((f(x-u)-f(x))/u + f'(x)) = \frac{1}{u} \left(-f'(x) + \frac{u}{2}f''(\xi) + f'(x) \right) = \frac{f''(\xi)}{2}$$

where $\xi \in (x - u, x)$, so sending $u \rightarrow 0$ shows (since f'' is continuous) that $h'_x(0) = \frac{f''(x)}{2}$. For $u \neq 0$ the quotient rule gives

$$h'_x(u) = \frac{-f'(x-u)}{u} - \frac{f(x-u)-f(x)}{u^2} = \frac{f'(x) - f'(x-u)}{u} - \frac{1}{2}f''(\xi_1) = f''(\xi_2) - \frac{1}{2}f''(\xi_1)$$

for certain $\xi_1, \xi_2 \in (x - u, x)$. This proves that $|h'_x(u)| \leq \frac{3}{2}M_1$ for all x and u , and also proves that h'_x is continuous at $u = 0$.

As for $p(u) = \frac{\sin(u/2)}{u}$, we simply note that p has a Taylor expansion

$$p(u) = \sum_{n=0}^{\infty} \frac{(-1)^n}{2((2n+1)!)} (u/2)^{2n+1}$$

which has infinite radius of convergence; p is clearly continuous, and term-by-term differentiation implies the continuity of p' . So $|p|$ and $|p'|$ are both bounded above on the compact interval $[-\pi, \pi]$ by some universal constant.

Since p never vanishes and is bounded below on $[-\pi, \pi]$ the continuity and boundedness of each of h_x, h'_x, p, p' then implies the same for $g'_x = \frac{1}{p^2}(ph'_x - p'h_x)$. \square

This result is encouraging, but isn't as strong as might be hoped: the Fourier partial sums $S_N f$ are constructed without taking any derivatives, so *a priori* differentiability seems like it should have little to do with whether $S_N f \rightarrow f$. Here's a sampling of results indicating the extent to which the hypothesis of Theorem 4.4 can be weakened. (We won't prove these, and you're not responsible for them.)

- (i) If one works a bit harder, one can improve the proof of Theorem 4.4 so that it still applies when f is just C^1 . The idea here is to, for suitable $\delta > 0$,

divide the integral

$$S_N f(x) - f(x) = \int_{-\pi}^{\pi} g_x(u) \sin((N + 1/2)u) du$$

into a region $|u| \geq \delta$ (where g_x is still C^1 and so integration by parts can still be carried out) and a region $|u| < \delta$ (where the integral can be made small because the region of integration is small). This is a bit tricky; one needs to take δ on the order of $N^{-1/2}$. (Both this and Theorem 4.4 are due to Dirichlet in the mid-19th century.)

- (ii) Fejer showed around the turn of the 20th century that if one sets $\sigma_N f = \frac{1}{N+1} \sum_{n=0}^N S_N f$, then for any continuous f it is true that $\sigma_N f \rightarrow f$ uniformly. It turns out that where $K_N(x) = \frac{1}{N+1} \sum_{n=0}^N D_N(x)$, so that

$$(11) \quad \sigma_N f(x) = \int_{-\pi}^{\pi} K_N(u) f(x-u) du,$$

K_N is nonnegative everywhere (!) and one has, for any $\delta > 0$, $\int_{|x| > \delta} K_N(u) du \rightarrow 0$ as $N \rightarrow \infty$. As a result of this, for N large nearly all of the integral in (11) is concentrated very close to x , which (using that $\int_{-\pi}^{\pi} K_N(u) du = 1$) can be seen to imply that $\sigma_N f \rightarrow f$ uniformly.

- (iii) Notwithstanding the above positive results, there are continuous functions f such that $\sup_N S_N f(0) = \infty$, so that $S_N f$ cannot converge to f , even pointwise. Fejer constructed one such f ; his example appears in the chapter on Fourier series in Strichartz's book. More strikingly, using results about complete metric spaces that are only slightly beyond the scope of this course, it's possible to give a more indirect argument that there is a dense set $A \subset C([-\pi, \pi])$ such that for each $f \in A$ $S_N f(x)$ diverges for every rational number x . (In fact, for each $f \in A$, the set of x at which $S_N f$ diverges is uncountable). If you're curious, the argument proving this (including all the necessary theoretical background) may be found on pp. 95-103 of Rudin's *Real and Complex Analysis*.
- (iv) In 1966, Lennart Carleson proved that for any L^2 function f , $S_N f \rightarrow f$ almost everywhere. (This is a very hard theorem, and is a large part of the reason that Carleson won the Abel Prize in 2006.)

While we won't be proving the above results, we will explore another form of convergence, namely convergence in L^2 . Let μ be the measure on the σ -algebra of Lebesgue measurable subsets of $[-\pi, \pi]$ obtained by restricting Lebesgue measure to such subsets (so our measure space (X, \mathcal{M}, μ) has $X = [-\pi, \pi]$). In keeping with the periodicity requirement on our functions it might make more sense to consider the points π and $-\pi$ as being identified with each other in X (or to consider X as being the unit circle in the complex plane, with measure obtained from Lebesgue measure via the usual parametrization of the unit circle by $[-\pi, \pi]$), but since we'll be looking at L^2 functions on X and since any L^2 function on $[-\pi, \pi]$ agrees a.e. with a periodic function this doesn't really make a difference.

In any case, given $f \in L^2(\mu)$, we also have $f \in L^1(\mu)$ since $[-\pi, \pi]$ has finite measure, and so where we define

$$e_n(x) = e^{inx}$$

we can form the Fourier coefficients

$$a_n = \frac{1}{2\pi} \int_{[-\pi, \pi]} f e_{-n} d\mu$$

and then the Fourier partial sums

$$S_N f = \sum_{n=-N}^N a_n e_n.$$

Of course each $e_n \in L^2(\mu)$, so since the a_n are (finite) scalars we always have $S_N f \in L^2(\mu)$. Our goal is to show that $S_N f \rightarrow f$ in L^2 .

We'll be able to do this without too much difficulty by using the fact that $\mathcal{L}^2(\mu)$ carries a natural inner product, compatible with the L^2 norm. Namely, for $f, g \in L^2(\mu)$, define

$$\langle f, g \rangle = \int_{[-\pi, \pi]} f \bar{g} d\mu.$$

Note that the Hölder inequality gives

$$\int_{[-\pi, \pi]} |fg| d\mu \leq \|f\|_2 \|g\|_2,$$

so $\langle f, g \rangle$ is always a finite complex number. Also, we clearly have

$$\langle f, f \rangle = \|f\|_2^2.$$

In particular, $\langle f, f \rangle = 0$ only if $f = 0$ a.e., so if you know what an inner product space (over \mathbb{C}) is you can easily show that this induces an inner product on $L^2(\mu)$.

Note that since the n th Fourier coefficient of $f \in L^2(\mu)$ is $a_n = \frac{1}{2\pi} \int_0^{2\pi} f(x) e_n(x) dx$, we have

$$S_N f = \sum_{n=-N}^N \frac{1}{2\pi} \langle f, e_n \rangle e_n.$$

Also note that

$$\langle e_n, e_m \rangle = \begin{cases} 2\pi & n = m \\ 0 & n \neq m \end{cases}$$

(This follows from our earlier computation of $\int_{-\pi}^{\pi} e^{i(n-m)x} dx$.)

Define a sequence of subspaces V_N ($N \geq 1$) of $L^2(\mu)$ as follows:

$$V_N = \left\{ \sum_{n=-N}^N a_n e_n \mid a_n \in \mathbb{C} \right\},$$

so we evidently have $S_N f \in V_N$ for any $f \in L^2(\mu)$. In fact $S_N f$ is what is called the orthogonal projection of f into V_N , which has the following consequence:

Lemma 4.6. *If $f \in L^2(\mu)$ and $g \in V_N$, then*

$$\|f - S_N f\|_2 \leq \|f - g\|_2,$$

with equality only if $g = S_N f$.

Proof. For each n with $-N \leq n \leq N$ we have

$$\langle S_N f, e_n \rangle = \left\langle \sum_{m=-N}^N \frac{1}{2\pi} \langle f, e_m \rangle e_m, e_n \right\rangle = \langle f, e_n \rangle.$$

Thus,

$$\langle f - S_N f, e_n \rangle = 0 \text{ when } -N \leq n \leq N$$

(We've used a linearity property of the inner product $\langle \cdot, \cdot \rangle$ here, which follows immediately from its definition). So if $h \in V_N$, in which case we can write $h = \sum_{n=-N}^N a_n e_n$ for some $a_n \in \mathbb{C}$, it then follows (again from the linearity of the inner product) that

$$\langle f - S_N f, h \rangle = \langle h, f - S_N f \rangle = 0.$$

Hence (using $h = S_N f - g \in V_N$)

$$\begin{aligned} \|f - g\|_2^2 &= \langle f - g, f - g \rangle = \langle (f - S_N f) + (S_N f - g), (f - S_N f) + (S_N f - g) \rangle \\ &= \langle f - S_N f, f - S_N f \rangle + \langle f - S_N f, S_N f - g \rangle + \langle S_N f - g, f - S_N f \rangle + \langle S_N f - g, S_N f - g \rangle \\ &= \|f - S_N f\|_2^2 + \|S_N f - g\|_2^2 \geq \|f - S_N f\|_2^2, \end{aligned}$$

with equality only if $\|S_N f - g\|_2^2 = 0$, *i.e.*, only if $S_N f = g$ a.e.. But $S_N f$ and g , both being elements of V_N , are continuous functions, so if they are equal a.e. they are equal everywhere, completing the proof. \square

Lemma 4.6 shows that for each N $S_N f$ is the best approximation to f in V_N with respect to the L^2 norm. We want to argue that $S_N f$ approximates f arbitrarily well for suitably large N , and the lemma implies that we'll be able to conclude this as soon as we find any element g of V_N which approximates f to within the prescribed error. Now if f is L^2 -close to some other function h such that $S_N h \rightarrow h$ in L^2 , then for large enough N f will be close to $S_N h$. As a result of this it will be enough to show that any $f \in L^2(\mu)$ is L^2 -close to some other function whose Fourier series behaves in the desired way. This is indeed the case:

Lemma 4.7. *If $f \in L^2(\mu)$ and $\epsilon > 0$ there is a C^∞ function $h: [-\pi, \pi] \rightarrow \mathbb{C}$ such that $\|f - h\|_2 < \epsilon$.*

(In fact, it would be enough for our purposes to have h be C^2 , thanks to Theorem 4.4).

Proof. The proof is essentially the same as a problem from Problem Set 8 up until one final step. First one approximates f by a simple function $\sum_{i=1}^n a_i \chi_{A_i}$. Then one shows that each χ_{A_i} can be approximated (in L^2) by the characteristic function of a finite disjoint union of intervals, which is equal to $\sum_{j=1}^m \chi_{I_{ij}}$ for certain intervals I_{ij} . The last step then is to show that if $I = (a, b)$ is an interval and $\epsilon > 0$ then there is a C^∞ function h such that $\|\chi_I - h\|_2 < \epsilon$. In the homework you did this with h merely a continuous function (and with the L^2 norm replaced by the L^1 norm, but this doesn't affect the argument). Most likely, the h that you used was nondifferentiable at certain points. With these steps outlined, the proof of Lemma 4.7 will be complete when we prove the following, which is of interest in its own right:

Lemma 4.8. *If $I = (a, b)$ is any open interval in \mathbb{R} , and if $\delta > 0$, there is a C^∞ function h such that $0 \leq h(x) \leq 1$ for all x , $h(x) = 0$ when $x \notin (a - \delta, b + \delta)$ and $h(x) = 1$ when $x \in (a + \delta, b - \delta)$.*

This suffices to prove Lemma 4.7, since we then have $\int_{[-\pi, \pi]} |h - \chi_I|^2 d\mu \leq 4\delta$, which can be made arbitrarily small by taking δ small. \square

Proof of Lemma 4.8. We will build h out of the function $g: \mathbb{R} \rightarrow \mathbb{R}$ defined by $g(x) = 0$ for $x \leq 0$ and $g(x) = e^{-1/x}$ for $x > 0$. By induction on k , one can show that the k th derivative of g at a point $x > 0$ is given by $g^{(k)}(x) = R_k(x)e^{-1/x}$, where R_k is a rational function (*i.e.* a quotient of two polynomials; somewhat more specifically, $R_0(x) = 1$ and $R_{k+1}(x) = \frac{R_k(x)}{x^2} + R'_k(x)$ for $k \geq 0$). But then for all k we have $\lim_{x \rightarrow 0^+} g^{(k)}(x) = 0$. In light of this latter fact and the mean value theorem, we see that for every k we have $|g(x)| \leq |x|^{k+1}$ for all sufficiently small x , which can then be used to show (again by induction on k) that $g^{(k)}(0) = 0$ for all k . Thus g is C^∞ , with $g(x) = 0$ for $x \leq 0$ and $g(x) > 0$ for $x > 0$. So $g(\delta + x) = 0$ for $x \leq -\delta$ and $g(x) > 0$ for $x > -\delta$, and $g(\delta - x) = 0$ for $x \geq \delta$ and $g(\delta - x) > 0$ for $x < \delta$. In particular, $g(\delta + x) + g(\delta - x) > 0$ for all x . Hence

$$x \mapsto \frac{g(\delta + x)}{g(\delta + x) + g(\delta - x)}$$

defines a C^∞ function which is equal to 0 for $x \leq -\delta$ and to 1 for $x \geq \delta$. The function h defined by

$$h(x) = \left(\frac{g(\delta + x - a)}{g(\delta + x - a) + g(\delta - x + a)} \right) \left(\frac{g(\delta - x + b)}{g(\delta + x - b) + g(\delta - x + b)} \right)$$

will then satisfy the required properties. □

With this established, we can now prove L^2 -convergence of Fourier series.

Theorem 4.9. *If $f \in L^2(\mu)$ then $\|S_N f - f\|_2 \rightarrow 0$ as $N \rightarrow \infty$.*

Proof. Let $\epsilon > 0$, and by Lemma 4.7 choose a C^∞ function h such that $\|f - h\|_2 < \epsilon/2$. By Theorem 4.4, $S_N h \rightarrow h$ uniformly on $[-\pi, \pi]$. So there is N_0 such that when $N \geq N_0$ we have, for all $x \in [-\pi, \pi]$, $|S_N h(x) - h(x)|^2 < \frac{\epsilon^2}{8\pi}$. So

$$\int_{-\pi}^{\pi} |S_N h - h|^2 d\mu \leq \frac{\epsilon^2}{4},$$

i.e., $\|S_N h - h\|_2 < \epsilon/2$. So the Minkowski inequality gives

$$\|f - S_N h\|_2 \leq \|f - h\|_2 + \|h - S_N h\|_2 < \epsilon.$$

But $S_N h \in V_N$, so by Lemma 4.6 we have $\|f - S_N f\|_2 \leq \|f - S_N h\|_2 < \epsilon$ whenever $N \geq N_0$, as desired. □

Accordingly Theorem 3.56 implies that for some subsequence $\{S_{N_k} f\}_{k=1}^\infty$ of $\{S_N f\}_{N=1}^\infty$ we have $S_{N_k} f \rightarrow f$ almost everywhere. However the theorem of Carleson mentioned earlier actually implies that one doesn't need to pass to a subsequence.

Note that Theorem 4.9 immediately implies:

Corollary 4.10 (Parseval's Theorem). *If $f \in L^2(\mu)$, then where $a_n = \frac{1}{2\pi} \int_{-\pi}^{\pi} f(x)e^{-inx} dx$ is the n th Fourier coefficient of f , we have*

$$\int_{-\pi}^{\pi} |f|^2 d\mu = 2\pi \sum_{n=-\infty}^{\infty} |a_n|^2.$$

Proof. Indeed, since $S_N f \rightarrow f$ in L^2 , we have

$$\|f\|_2^2 = \lim_{N \rightarrow \infty} \|S_N f\|_2^2 = \lim_{N \rightarrow \infty} \left\langle \sum_{n=-N}^N a_n e_n, \sum_{m=-N}^N a_m e_m \right\rangle = 2\pi \sum_{n=-\infty}^{\infty} |a_n|^2,$$

since $\langle e_n, e_m \rangle$ is zero when $n \neq m$ and 2π when $n = m$. \square

4.2. The Fourier transform on \mathbb{R} . In the previous subsection, we tried to represent a periodic function of period 2π as a superposition of the standard periodic functions e^{inx} where $n \in \mathbb{Z}$. The Fourier transform does something similar for general integrable functions $f: \mathbb{R} \rightarrow \mathbb{R}$, but now there is a continuum of standard periodic functions, namely $x \mapsto e^{itx}$ where t can be any real number. Throughout, we let m denote Lebesgue measure on \mathbb{R} , and always use this measure for each integral that appears (though the notation will not reflect this; contrary to our earlier notation but consistently with what you've seen in calculus, we'll write things like " $\int_{-\infty}^{\infty} g(x) dx$ " when we mean " $\int_{\mathbb{R}} g dm$, where g is the function $x \mapsto g(x)$ ").

Definition 4.11. Let $f \in L^1(m)$. The *Fourier transform* of f is the function

$$\hat{f}: \mathbb{R} \rightarrow \mathbb{C}$$

defined by

$$\hat{f}(t) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} f(x) e^{-itx} dx.$$

This is formally similar to how we defined the Fourier coefficients of a 2π -periodic function. If instead of being integrable on \mathbb{R} f was a C^2 function on $[-\pi, \pi]$, we could denote $\hat{f}(n) = \frac{1}{\sqrt{2\pi}} \int_{-\pi}^{\pi} f(x) e^{-inx} dx$, and then Theorem 4.4 says that we have $f(x) = \frac{1}{\sqrt{2\pi}} \sum_{n \in \mathbb{Z}} \hat{f}(n) e^{inx}$. Our goal will be to establish an analogous result for the Fourier transform.

It will be useful to compute the Fourier transform of one particular function in order to do this; we'll use the following:

Example 4.12. Let $H(x) = e^{-|x|}$. We then find

$$\begin{aligned} \hat{H}(t) &= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{-|x|-itx} dx \\ &= \frac{1}{\sqrt{2\pi}} \left(\int_{-\infty}^0 e^{x-itx} dx + \int_0^{\infty} e^{-x-itx} dx \right) = \frac{1}{\sqrt{2\pi}} \left(\frac{1}{1-it} + \frac{1}{1+it} \right) \\ &= \sqrt{\frac{2}{\pi}} \frac{1}{1+t^2}. \end{aligned}$$

From the definition of the Fourier transform, it's quickly apparent that it turns certain natural operations on functions into certain other natural operations: for instance if $f \in L^1(m)$ and we put $g(x) = e^{iax} f(x)$ then one calculates $\hat{g}(t) = \hat{f}(t - a)$; if instead we put $g(x) = f(x - a)$ one calculates $\hat{g}(t) = e^{-iat} \hat{f}(t)$.

A similar fact that will be important for us later is the following: If $f \in L^1(m)$ and $\lambda > 0$, set

$$\sigma_{\lambda} f(x) = f(\lambda x).$$

Then

$$\begin{aligned}\widehat{\sigma_\lambda f}(t) &= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} f(\lambda x) e^{-itx} dx = \frac{1}{\sqrt{2\pi}} \frac{1}{\lambda} \int_{-\infty}^{\infty} f(u) e^{-itu/\lambda} du \\ &= \frac{1}{\lambda} \widehat{f}(t/\lambda) = \frac{1}{\lambda} \sigma_{1/\lambda} \widehat{f}(t).\end{aligned}$$

We'll actually only need this for the functions of Example 4.12: Write $h_1(t) = \frac{1}{\pi(1+t^2)}$, so that $\widehat{H} = \sqrt{2\pi}h_1$, and put

$$h_\lambda(t) = \frac{1}{\lambda} \sigma_{1/\lambda} h_1(t) = \frac{1}{\lambda\pi} \frac{1}{1 + (t/\lambda)^2} = \frac{1}{\pi} \frac{\lambda}{\lambda^2 + t^2}.$$

Then

$$(12) \quad \widehat{\sigma_\lambda H}(t) = \sqrt{2\pi}h_\lambda(t) \quad (H(x) = e^{-|x|}).$$

We'll make use here of what happens as $\lambda \rightarrow 0$: $\sigma_\lambda H(x) = e^{-\lambda|x|}$ “flattens out” and approaches one for any given x , whereas $h_\lambda(t) = \frac{\lambda}{\pi(\lambda^2+t^2)}$ becomes quite small outside a small neighborhood of zero but quite large at zero, while retaining the property that $\int_{-\infty}^{\infty} h_\lambda(t) dt = 1$ for all λ (as can be seen by the trigonometric substitution $t = \lambda \tan \theta$).

Remark 4.13. *In some of the arguments to follow, we will be making statements along the lines of*

$$(13) \quad \int_{-\infty}^{\infty} \left(\int_{-\infty}^{\infty} F(x, y) dx \right) dy = \int_{-\infty}^{\infty} \left(\int_{-\infty}^{\infty} F(x, y) dy \right) dx,$$

where $F: \mathbb{R}^2 \rightarrow \mathbb{C}$ is some function. If, for instance, F is continuous and vanishes outside some compact set, then this follows from facts about the Riemann integral, as is familiar from multivariable calculus. In the generality that we need, it follows from a result called the Fubini theorem, which would require a rather lengthy detour to prove and whose proof we therefore omit. (13) holds for any F satisfying a certain measurability condition and also having the property that each function $|F|_x: y \mapsto |F(x, y)|$ is integrable and also $x \mapsto \int |F|_x(y) dy$ is integrable, and these properties will hold in every situation we consider below. There are examples of functions $F: \mathbb{R}^2 \rightarrow \mathbb{C}$ (not satisfying the properties of the last sentence) such that both sides of (13) exist but they are unequal.

Proposition 4.14. *Let $f, g \in L^1(\mathbb{R})$, and define*

$$f * g(x) = \int_{-\infty}^{\infty} f(x - y)g(y)dy$$

for all x such that the integral on the right exists (and $(f * g)(x) = 0$ otherwise). Then $f * g \in L^1(\mathbb{R})$ (in fact $\|f * g\|_1 \leq \|f\|_1 \|g\|_1$), and

$$\widehat{f * g}(t) = \sqrt{2\pi} \widehat{f}(t) \widehat{g}(t).$$

Proof. We have

$$\begin{aligned}\int_{-\infty}^{\infty} \left(\int_{-\infty}^{\infty} |f(x - y)| |g(y)| dy \right) dx &= \int_{-\infty}^{\infty} |g(y)| \left(\int_{-\infty}^{\infty} |f(x - y)| dx \right) dy \\ &= \|f\|_1 \int_{-\infty}^{\infty} |g(y)| dy = \|f\|_1 \|g\|_1.\end{aligned}$$

Now $\|f\|_1\|g\|_1$ is finite, so it must be that for a.e. x the function $y \mapsto f(x-y)g(y)$ is integrable, and so we have $(f * g)(x) = \int_{-\infty}^{\infty} f(x-y)g(y)dy$ for a.e. x . Further, the calculation above (and the fact that $|\int h| \leq \int |h|$ for any h) shows that $\|f * g\|_1 \leq \|f\|_1\|g\|_1$, so $f * g \in L^1(m)$.

We now find

$$\begin{aligned} \widehat{f * g}(t) &= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} \left(\int_{-\infty}^{\infty} f(x-y)g(y)dy \right) e^{-it(x-y)} dx \\ &= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} g(y)e^{-ity} \left(\int_{-\infty}^{\infty} f(x-y)e^{-it(x-y)} dx \right) dy = \sqrt{2\pi} \hat{f}(t) \hat{g}(t). \end{aligned}$$

□

Lemma 4.15. *Suppose that $g: \mathbb{R} \rightarrow \mathbb{C}$ is continuous, and that, for some $M > 0$, we have $g(x) = 0$ whenever $|x| \geq M$. Then*

$$\lim_{\lambda \rightarrow 0} \|g * h_\lambda - g\|_1 = 0.$$

Proof. Since $\int_{-\infty}^{\infty} h_\lambda(y)dx = 1$, we have

$$(g * h_\lambda)(x) - g(x) = \int_{-\infty}^{\infty} (g(x-y) - g(x))h_\lambda(y)dy.$$

Let $\epsilon > 0$. g is continuous, and vanishes outside the compact set $[-M, M]$, so g is uniformly continuous. Thus there is $\delta > 0$ (and we may as well suppose $\delta < 1$) such that whenever $|y| \leq \delta$ and $x \in \mathbb{R}$ we have $|g(x-y) - g(x)| < \epsilon/(4M+4)$ (this follows from Theorem 1.41). Also, if $y \leq \delta < 1$, then $g(x-y)$ and $g(x)$ are both zero unless $x \in [-M-1, M+1]$, so that $\int_{-\infty}^{\infty} |g(x-y) - g(x)|dx \leq \frac{\epsilon}{4M+4}(2M+2) = \frac{\epsilon}{2}$. Thus, for any λ ,

$$(14) \quad \int_{|y| < \delta} \int_{-\infty}^{\infty} |h_\lambda(y)| \left(\int_{-\infty}^{\infty} |g(x-y) - g(x)|dx \right) dy < \epsilon/2.$$

Also g is bounded (again since it is continuous and vanishes outside some compact set), so there is $N > 0$ such that $|g(x)| < N$ for all $x \in \mathbb{R}$. Now

$$\int_{|y| \geq \delta} |h_\lambda(y)|dy = 2 \int_{\delta}^{\infty} \frac{1}{\lambda} h_1(y/\lambda)dy = 2 \int_{\delta/\lambda}^{\infty} \frac{du}{\pi(1+u^2)} \rightarrow 0$$

as $\lambda \rightarrow 0$, so there is $\lambda_0 > 0$ such that whenever $0 < \lambda < \lambda_0$ we have

$$\int_{|y| \geq \delta} |h_\lambda(y)|dy < \frac{\epsilon}{16NM}.$$

Given y , the set of x such that one or both of the terms in the expression $g(x-y) - g(x)$ is nonzero is contained in $[y-M, y+M] \cup [-M, M]$, and $|g(x-y) - g(x)| < 2N$ for all such x and y , so we obtain, for $\lambda < \lambda_0$,

$$\int_{|y| \geq \delta} \int_{-\infty}^{\infty} |h_\lambda(y)| \int_{-\infty}^{\infty} |g(x-y) - g(x)|dx dy \leq \frac{\epsilon}{16NM} 2N(4M) = \epsilon/2.$$

Combining this with (14) implies that for $\lambda < \lambda_0$ we have $\|g * h_\lambda - g\|_1 < \epsilon$, which completes the proof. □

Corollary 4.16. *If $f \in L^1(m)$ then*

$$\lim_{\lambda \rightarrow 0} \|f * h_\lambda - f\|_1 = 0.$$

Proof. If $\epsilon > 0$, choose a continuous function $g: \mathbb{R} \rightarrow \mathbb{C}$ which vanishes outside some interval $[-M, M]$ such that $\|f - g\|_1 \leq \epsilon/3$, and let $\lambda_0 > 0$ be so small that $\|g * h_{\lambda_0} - g\|_1 < \epsilon/3$, using the preceding lemma.

Then, for $\lambda < \lambda_0$

$$\begin{aligned} \|f * h_\lambda - f\|_1 &\leq \|f * h_\lambda - g * h_\lambda\|_1 + \|g * h_\lambda - g\|_1 + \|g - f\|_1 < \|(f - g) * h_\lambda\|_1 + \frac{2\epsilon}{3} \\ &\leq \|f - g\|_1 \|h_\lambda\|_1 + \frac{2\epsilon}{3} < \epsilon. \end{aligned}$$

□

We can now prove the Fourier Inversion Theorem:

Theorem 4.17. *Suppose that $f \in L^1(\mathbb{R})$, and also that $\hat{f} \in L^1(\mathbb{R})$. Then*

$$f(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} \hat{f}(t) e^{itx} dt$$

for almost every $x \in \mathbb{R}$.

Note the sign change in the exponent, and also the similarity to the situation with Fourier series. One immediate consequence of Theorem 4.17 that is worth noting is that if $f \in L^1(\mathbb{R})$ and $\hat{f} = 0$ then $f = 0$ a.e.

Proof. Recall that the functions $h_\lambda(t) = \frac{\lambda}{\pi(\lambda^2 + t^2)}$ satisfy (where $\lambda > 0$)

$$\sqrt{2\pi} h_\lambda(t) = \widehat{\sigma_\lambda H}(t) \quad (\sigma_\lambda H)(t) = H(\lambda t) = e^{-\lambda|t|}.$$

Now the h_λ are obviously even functions, so we have

$$h_\lambda(t) = h_\lambda(-t) = \frac{1}{2\pi} \int_{-\infty}^{\infty} H(\lambda u) e^{itu} du.$$

So

$$\begin{aligned} f * h_\lambda(x) &= \int_{-\infty}^{\infty} f(x-t) h_\lambda(t) dt = \frac{1}{2\pi} \int_{-\infty}^{\infty} f(x-t) \int_{-\infty}^{\infty} H(\lambda u) e^{itu} du dt \\ &= \frac{1}{2\pi} \int_{-\infty}^{\infty} H(\lambda u) \int_{-\infty}^{\infty} f(x-t) e^{itu} dt du \\ &= \frac{1}{2\pi} \int_{-\infty}^{\infty} H(\lambda u) \int_{-\infty}^{\infty} f(y) e^{i(x-y)u} dy du = \frac{1}{2\pi} \int_{-\infty}^{\infty} H(\lambda u) e^{ixu} \int_{-\infty}^{\infty} f(y) e^{-iyu} dy du \\ &= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} \hat{f}(u) H(\lambda u) e^{ixu} du. \end{aligned}$$

Now recalling that $H(\lambda u) = e^{-\lambda|u|}$, if $\{\lambda_n\}_{n=1}^{\infty}$ is any sequence of positive numbers such that $\lambda_n \rightarrow 0$, for $\lambda = \lambda_n$ the integrand above has absolute value at most $|\hat{f}(u)|$ for each λ_n and u , and converges pointwise as $n \rightarrow \infty$ to $\hat{f}(u) e^{ixu}$. So the Dominated Convergence Theorem (which applies since we assume $\hat{f} \in L^1(\mathbb{R})$) shows that, if $\lambda_n \rightarrow 0$, then for all x we have

$$(f * h_{\lambda_n})(x) \rightarrow \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} \hat{f}(u) e^{ixu} du.$$

But Corollary 4.16 shows that $f * h_{\lambda_n} \rightarrow f$ in L^1 , so for some subsequence we have $f * h_{\lambda_{n_k}} \rightarrow f$ a.e. The only way for these two statements to be compatible is if in fact we have, for a.e. x ,

$$f(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} \hat{f}(u) e^{ixu} du.$$

□