

Assorted notes

Sanjit Dandapanthula¹

sanjitd@cmu.edu

April 5, 2025

¹Carnegie Mellon University, Department of Statistics

Table of contents

| | | |
|----------|---|-----------|
| 1 | Probability theory | 3 |
| 1.1 | Basic theorems | 3 |
| 1.2 | Lévy's continuity theorem and inversion formula | 5 |
| 1.3 | Kolmogorov's three-series lemma and the SLLN | 7 |
| 1.4 | Disintegration of measure and regular conditional probability | 11 |
| 2 | Functional analysis | 13 |
| 2.1 | Key algebraic structures | 13 |
| 2.1.1 | Normed spaces and Banach spaces | 13 |
| 2.1.2 | Linear operators | 15 |
| 2.1.3 | Hilbert spaces | 16 |
| 2.1.4 | Orthogonality and Fourier series | 18 |
| 2.2 | The dual space and the Riesz-Fréchet representation theorem | 20 |
| 2.3 | The four pillars of functional analysis | 20 |
| 2.3.1 | The Hahn-Banach theorem | 21 |
| 2.3.2 | The uniform boundedness principle (Banach-Steinhaus) | 22 |
| 2.3.3 | The open mapping theorem | 24 |
| 2.3.4 | The closed graph theorem | 25 |
| 2.4 | Hilbert adjoint operators | 26 |
| 2.5 | The adjoint operator | 28 |
| 2.6 | Reflexive spaces and separability | 28 |
| 2.7 | Weak convergence and weak-* convergence | 29 |
| 3 | Convex analysis | 32 |
| 3.1 | Convex sets and functions | 32 |
| 3.2 | Lower semi-continuous functions | 33 |

| | | |
|----------|---|-----------|
| 3.3 | Separation theorems | 34 |
| 3.4 | Subgradients and subdifferentials | 35 |
| 3.5 | The Legendre-Fenchel transform | 36 |
| 3.6 | Cyclical monotonicity | 38 |
| 4 | Optimal transport | 39 |
| 4.1 | The Monge-Kantorovich problem | 39 |
| 4.2 | Transport maps between empirical averages | 41 |
| 4.3 | Wasserstein distances | 43 |
| 4.4 | The Kantorovich duality | 46 |
| 4.4.1 | Lower semi-continuous cost functions | 46 |
| 4.4.2 | Metric cost functions | 48 |
| 4.5 | Brenier's theorem | 49 |
| A | Supplementary results | 52 |
| A.1 | Zorn's lemma | 52 |
| A.2 | The Baire category theorem | 53 |
| A.3 | Tychonoff's theorem | 53 |

Chapter 1

Probability theory

In this chapter, we review some fundamental theorems from measure-theoretic probability theory. The sources in this chapter are widely varied, and depend mostly on my style preference for each topic. I'll mostly assume the basics of measure-theoretic probability and only cover theorems that I find interesting or useful. I'll also state most of the theorems for real-valued random variables, but many of them can be easily generalized to \mathbb{R}^d -valued random variables (or even more general spaces).

1.1 Basic theorems

We start with Skorokhod's representation theorem, which says that random variables converging in distribution can be coupled in such a way that they converge almost surely.

Definition 1.1.1 (Convergence in distribution). We say that $X_n \xrightarrow{d} X$ (X_n converges in distribution to X) if $F_{X_n}(x) \rightarrow F_X(x)$ for all continuity points x of F . Probabilists also call this *weak convergence*, but this actually corresponds to weak-* convergence in the sense of functional analysis.

Theorem 1.1.1 (Skorokhod representation). *Let $X_n \xrightarrow{d} X$. Then, there exists a probability space $(\tilde{\Omega}, \tilde{\mathcal{F}}, \tilde{\mathbb{P}})$ and random variables Y_n and Y on this space such that $Y_n \stackrel{d}{=} X_n$, $Y \stackrel{d}{=} X$, and $Y_n \xrightarrow{a.s.} Y$.*

Proof. The proof is easy; define the quantile transformation $F^{-1}(x) = \inf\{x \in \mathbb{R} : F(x) \geq x\}$, draw a uniformly random variable $U \sim \text{Unif}(0, 1)$, and define $Y_n = F_{X_n}^{-1}(U)$ and $Y = F_X^{-1}(U)$. \square

Next, we state the portmanteau lemma, which characterizes convergence in distribution in a variety of ways.

Theorem 1.1.2 (Portmanteau lemma). *The following are equivalent:*

1. $X_n \xrightarrow{d} X$.

2. $\mathbb{E}[f(X_n)] \rightarrow \mathbb{E}[f(X)]$ for all $f \in C_b(\mathbb{R})$.
3. $\mathbb{E}[f(X_n)] \rightarrow \mathbb{E}[f(X)]$ for all Lipschitz $f \in C_b(\mathbb{R})$.
4. $\liminf_{n \rightarrow \infty} \mathbb{E}[f(X_n)] \geq \mathbb{E}[f(X)]$ for all lower semi-continuous f taking values in $[0, \infty]$.
5. $\liminf_{n \rightarrow \infty} \mathbb{P}(X_n \in G) \geq \mathbb{P}(X \in G)$ for all open G .
6. $\mathbb{P}(X_n \in A) \rightarrow \mathbb{P}(X \in A)$ for all $A \in \mathcal{B}(\mathbb{R})$ with $\mathbb{P}(X \in \partial A) = 0$.

Proof. To show (1) implies (2), use Skorokhod's representation theorem (Theorem 1.1.1) and apply the dominated convergence theorem since f is bounded. Obviously, (2) implies (3). To show (3) implies (4), note that we can find Lipschitz functions $f_m \uparrow f$ pointwise since f is l.s.c. (Proposition 3.2.2); take the \liminf in n and then the limit in m using the monotone convergence theorem. It's clear that (4) implies (5) since $\mathbf{1}_G$ is l.s.c. Then, (5) implies (6) by taking complements to get a similar statement for closed set and applying these results to $\text{int}(A)$ and \bar{A} . (6) implies (1) is obvious. \square

Continuous mappings preserve convergence in distribution, in probability, and almost surely.

Theorem 1.1.3 (Continuous mapping). *If $X_n \rightarrow X$ in distribution, in probability, or almost surely, then $g(X_n) \rightarrow g(X)$ in the same sense for any continuous function g .*

Proof. The proof is obvious for almost sure convergence and convergence in probability, and for convergence in distribution, use (2) in the portmanteau lemma (Theorem 1.1.2). \square

Now, we state Prokhorov's theorem, which is a version of Bolzano-Weierstrass for probability measures.

Definition 1.1.2 (Tightness). A set $\{X_\alpha\}_{\alpha \in A}$ of random variables is *tight* if for all $\epsilon > 0$ there exists a compact set K such that $\sup_{\alpha \in A} \mathbb{P}(X_\alpha \notin K) \leq \epsilon$.

Theorem 1.1.4 (Prokhorov). *A set $\{X_\alpha\}_{\alpha \in A}$ of random variables is tight if and only if every sequence has a weakly convergent subsequence.*

Note that Prokhorov's theorem generalizes to any *Polish space* (complete separable metric space).

Proof. For the forward direction, pick a sequence $(X_n)_{n=1}^\infty \subseteq \{X_\alpha\}_{\alpha \in A}$ and enumerate \mathbb{Q} . By Cantor's diagonalization argument, extract a subsequence $(F_{X_{n_k}})_{k=1}^\infty$ such that $F_{X_{n_k}} \rightarrow F$ on all rationals. Define $F(x) = \inf_{q > x} F(q)$; it is easy to show that F is a cdf.

For the reverse direction, suppose every sequence in $\{X_\alpha\}_{\alpha \in A}$ has a weakly convergent subsequence but $\{X_\alpha\}_{\alpha \in A}$ isn't tight; in particular, there exists $\epsilon > 0$ such that for all compact K we have $\sup_{\alpha \in A} \mathbb{P}(X_\alpha \notin K) > \epsilon$. Then, pick $K_n = [-n, n]$ and choose X_n such that $\mathbb{P}(X_n \notin K_n) > \epsilon$. Extract a weakly convergent

subsequence $X_{n_k} \xrightarrow{d} X$ and pick $M > 0$ so that $\mathbb{P}(X \in (-M, M)) > 1 - \epsilon$. But we know by (5) in the portmanteau lemma ([Theorem 1.1.2](#)) that

$$1 - \epsilon < \mathbb{P}(X \in (-M, M)) \leq \liminf_{n \rightarrow \infty} \mathbb{P}(X_{n_k} \in (-M, M)) \leq 1 - \epsilon.$$

giving a contradiction. \square

In the context of the Banach-Alaoglu theorem ([Theorem 2.7.6](#)) and the Riesz-Markov representation theorem, one can view Prokhorov's theorem as saying that a family of probability measures is relatively closed in the unit ball of the dual space of $C_0(\mathbb{R})$ (the space of finite Radon measures) if and only if it is tight and closed. Furthermore, any probability measure in a Polish space is tight.

Theorem 1.1.5 (Ulam's lemma). *Any probability measure μ on a Polish space \mathcal{X} is σ -finite and therefore tight.*

Proof. Let $\{x_n\}_{n=1}^\infty$ be a countable dense set and define the sets $K_m = \bigcup_{n=1}^m \overline{B(x_n, 1/m)}$ for each $m \in \mathbb{N}$, which are compact because they are each totally bounded and complete. Then $\mathcal{X} = \bigcup_{m=1}^\infty K_m$, so \mathcal{X} is σ -finite. \square

Note that a Borel probability measure on a Polish space is also automatically *regular*, meaning that the measure of any set can be approximated from above by open sets and from below by closed sets; this is a consequence of the Riesz-Markov representation theorem. Finally, we show that the L^p norms are ordered in probability spaces.

Proposition 1.1.6 (Ordering of L^p norms). *If $1 \leq p < q \leq \infty$, then $\|X\|_p \leq \|X\|_q$ for all random variables X .*

Proof. Hölder's inequality with the conjugate exponents q/p and $q/(q-p)$ gives

$$\|X\|_p^p = \mathbb{E}[|X|^p] \leq \| |X|^p \|_{q/p} \|1\|_{q/(q-p)} = \|X\|_q^p.$$

This only works in a probability space because we needed that $\|1\|_{q/(q-p)} = 1$. \square

1.2 Lévy's continuity theorem and inversion formula

In this section, we give fundamental results about characteristic functions.

Definition 1.2.1 (Characteristic function). The *characteristic function* of a random variable X is $\varphi_X(t) = \mathbb{E}[e^{itX}]$.

We start by proving the Lévy continuity theorem, which characterizes convergence in distribution in terms of the convergence of characteristic functions.

Theorem 1.2.1 (Lévy continuity). *Suppose that $\varphi_{X_n} \rightarrow \varphi$ converges pointwise for all $t \in \mathbb{R}$. Then, there exists a random variable X such that $X_n \xrightarrow{d} X$ if and only if φ is continuous at 0.*

Proof. The forward direction is immediate from the dominated convergence theorem, so we'll focus on the reverse direction. If $(X_n)_{n=1}^\infty$ was tight, then every subsequence would have a further subsequence converging in distribution to some random variable X . But then by the portmanteau lemma, the characteristic functions would converge along all subsequences to φ_X and the result would immediately follow. So Lévy's continuity theorem is really about how the continuity of φ at 0 implies tightness.

The key idea in this proof is to use the Lebesgue differentiation theorem. Note that $\varphi_{X_n}(0) = 1$. Therefore, we study the following integral using the Lebesgue differentiation theorem and Fubini's theorem:

$$\begin{aligned} 0 &= \lim_{\delta \downarrow 0} \frac{1}{2\delta} \int_{-\delta}^{\delta} \Re(1 - \varphi_{X_n}(t)) dt \\ &= \lim_{\delta \downarrow 0} \mathbb{E} \left[\frac{1}{2\delta} \int_{-\delta}^{\delta} (1 - \cos(tX_n)) dt \right] \\ &= \lim_{\delta \downarrow 0} \mathbb{E} \left[1 - \frac{\sin(\delta X_n)}{\delta X_n} \right] \\ &\geq \lim_{\delta \downarrow 0} \mathbb{E} \left[\frac{\mathbf{1}_{|\delta X_n| \geq \pi}}{2} \right] \\ &= \lim_{\delta \downarrow 0} \frac{1}{2} \mathbb{P}(|X_n| \geq \pi/\delta). \end{aligned}$$

The inequality follows because when $x \geq \pi$, we have

$$\left| \frac{\sin(x)}{x} \right| \leq \frac{1}{\pi} \implies 1 - \frac{\sin(x)}{x} \geq 1 - \frac{1}{\pi} \geq \frac{1}{2}.$$

Picking δ small, we deduce that $(X_n)_{n=1}^\infty$ is tight as desired. \square

In fact, the proof shows that the continuity of $\Re(\varphi)$ at 0 is the only necessary condition for convergence in Lévy's continuity theorem. We can also use the Lévy continuity theorem to immediately get a converse to the weak law of large numbers.

Corollary 1.2.1.1 (Weak law of large numbers). *If X_1, \dots, X_n are i.i.d. with characteristic function φ and mean μ , then $\frac{1}{n} \sum_{i=1}^n X_i \xrightarrow{p} \mu$ if and only if φ is differentiable at 0 and $\varphi'(0) = i\mu$.*

Another corollary of the Lévy continuity theorem is the Lindeberg-Lévy central limit theorem.

Theorem 1.2.2 (Lindeberg-Lévy CLT). *If X_1, \dots, X_n are i.i.d. with mean μ and variance σ^2 , then $\frac{1}{\sqrt{n}} \sum_{i=1}^n X_i \xrightarrow{d} \mathcal{N}(0, \sigma^2)$.*

An alternative proof of the Lindeberg-Lévy CLT is to use the idea of a *Lindeberg exchange*; we start with Rademacher random variables in the sum and exchange them one-by-one for the X_i . Next, we state Lévy's inversion formula, which shows that random variables are uniquely determined by their characteristic functions.

Theorem 1.2.3 (Lévy inversion formula). *For continuity points $a < b$ of F_X , we have*

$$F_X(b) - F_X(a) = \frac{1}{2\pi} \int_{-\infty}^{\infty} \frac{e^{-ita} - e^{-itb}}{it} \varphi_X(t) dt.$$

Proof. The characteristic function φ_X is essentially the Fourier-Stieltjes transform of the distribution of X (up to sign changes), so the result follows from the general inversion formula from Fourier analysis. \square

1.3 Kolmogorov's three-series lemma and the SLLN

In this section, we give a proof of the strong law of large numbers using Kolmogorov's three-series lemma, which fully characterizes almost sure convergence of a series of random variables. We start by proving Kolmogorov's maximal inequality.

Theorem 1.3.1 (Kolmogorov's inequality). *Suppose X_1, \dots, X_n are independent with $\mathbb{E}[X_i] = 0$ and $\mathbb{E}[X_i^2] = \sigma_i^2$ and let $S_k = X_1 + \dots + X_k$. Then, for all $t > 0$ we have*

$$\mathbb{P}\left(\max_{1 \leq k \leq n} |S_k| \geq t\right) \leq \frac{1}{t^2} \sum_{i=1}^n \sigma_i^2.$$

Proof. Define the stopping time $\tau = \inf\{1 \leq k \leq n : |S_k| \geq t\} \wedge n$. Then, we compute the following L^2 estimate, since $\mathbf{1}_{i \leq \tau}$ is independent of X_i :

$$\mathbb{E}[S_\tau^2] = \sum_{i=1}^n \mathbb{E}[X_i^2 \mathbf{1}_{i \leq \tau}] = \sum_{i=1}^n \mathbb{E}[X_i^2] \mathbb{P}(i \leq \tau) \leq \sum_{i=1}^n \sigma_i^2.$$

Now the result follows from Chebyshev's inequality applied to S_τ . \square

From Kolmogorov's inequality we obtain the Khintchine-Kolmogorov two-series lemma.

Theorem 1.3.2 (Khintchine-Kolmogorov two-series lemma). *Let $(X_n)_{n=1}^\infty$ be independent with $\mathbb{E}[X_n] = \mu_n$ and $\mathbb{E}[X_n^2] = \sigma_n^2$ and let $S_k = X_1 + \dots + X_k$. Then, if $\sum_{n=1}^\infty \mu_n < \infty$ and $\sum_{n=1}^\infty \sigma_n^2 < \infty$, then S_k converges almost surely and in L^2 .*

Proof. Suppose without loss of generality that $\mu_n = 0$ for all $n \in \mathbb{N}$ (consider $X_n \mapsto X_n - \mu_n$). We know that

$$\mathbb{P}(S_k \text{ converges}) = \mathbb{P}\left(\bigcap_{k=1}^\infty \bigcup_{m=1}^\infty \bigcap_{n=m}^\infty \{|S_n - S_m| \leq 1/k\}\right).$$

By Kolmogorov's inequality ([Theorem 1.3.1](#)), we have

$$\mathbb{P}\left(\max_{m \leq i \leq n} |S_i - S_m| \geq 1/k\right) \leq k^2 \sum_{i=m}^n \sigma_i^2 \leq k^2 \sum_{i=m}^{\infty} \sigma_i^2.$$

Carefully letting $n \rightarrow \infty$ and then $m \rightarrow \infty$, it follows that for any $k \in \mathbb{N}$ we have

$$\mathbb{P}\left(\bigcap_{m=1}^{\infty} \bigcup_{n=m}^{\infty} \{|S_n - S_m| \geq 1/k\}\right) = 0 \implies \mathbb{P}\left(\bigcup_{m=1}^{\infty} \bigcap_{n=m}^{\infty} \{|S_n - S_m| \leq 1/k\}\right) = 1.$$

Almost sure convergence now follows by taking an intersection over all $k \in \mathbb{N}$; denote the limit as S . By Fatou's lemma, we have

$$\mathbb{E}[(S_n - S)^2] = \mathbb{E}\left[\liminf_{k \rightarrow \infty} (S_n - S_k)^2\right] \leq \liminf_{k \rightarrow \infty} \mathbb{E}[(S_n - S_k)^2] = \sum_{i=n}^{\infty} \sigma_i^2,$$

which tends to zero as $n \rightarrow \infty$; this gives the L^2 convergence. \square

Now, we use the two-series lemma to prove the three-series lemma.

Theorem 1.3.3 (Kolmogorov's three-series lemma). *Let $\lambda > 0$ be any constant (for instance, one could pick $\lambda = 1$). Suppose $(X_n)_{n=1}^{\infty}$ are independent and define $X_n^{(\lambda)} := X_n \mathbf{1}_{|X_n| \leq \lambda}$. Then, $\sum_{n=1}^{\infty} X_n$ converges almost surely if and only if the following three series converge:*

1. $\sum_{n=1}^{\infty} \mathbb{P}(|X_n| \geq \lambda)$
2. $\sum_{n=1}^{\infty} \mathbb{E}[X_n^{(\lambda)}]$
3. $\sum_{n=1}^{\infty} \text{Var}(X_n^{(\lambda)})$.

Proof. First, we'll show the reverse direction. By the first Borel-Cantelli lemma and convergence of series (1), we have that the probability that $\{|X_n| \geq \lambda\}$ infinitely often is zero; in particular, $X_n = Y_n$ eventually. By the Khintchine-Kolmogorov two-series lemma ([Theorem 1.3.2](#)) and convergence of series (2) and (3), we have that Y_n (and therefore X_n) converges almost surely.

Next, we show the forward direction. If series (1) didn't converge, then by the second Borel-Cantelli lemma (and independence of the X_n), we would have $|X_n| \geq \lambda$ infinitely often with probability 1, and $\sum_{n=1}^{\infty} X_n$ would diverge almost surely. Note that convergence of series (3) implies convergence of series (2): $\sum_{n=1}^{\infty} (X_n^{(\lambda)} - \mathbb{E}[X_n^{(\lambda)}])$ converges almost surely by the two-series lemma ([Theorem 1.3.2](#)) and for $\sum_{n=1}^{\infty} X_n$ to converge almost surely, we need $\sum_{n=1}^{\infty} \mathbb{E}[X_n^{(\lambda)}]$ to converge. If series (3) didn't converge, then by a slight generalization of the Lindeberg-Lévy CLT, we would find that

$$\frac{1}{\sqrt{\sum_{i=1}^n \text{Var}(X_i^{(\lambda)})}} \sum_{i=1}^n (X_i^{(\lambda)} - \mathbb{E}[X_i^{(\lambda)}]) \xrightarrow{d} \mathcal{N}(0, 1).$$

This simple generalization follows from the Lévy continuity theorem ([Theorem 1.2.1](#)) and holds whenever the sum of the variances of i.i.d. random variables diverges. Note that $\sum_{i=1}^n X_i^{(\lambda)}$ converges almost surely whenever $\sum_{i=1}^n X_i$ converges almost surely since the summands are almost surely eventually equal. But then this means that

$$\frac{1}{\sqrt{\sum_{i=1}^n \text{Var}(X_i^{(\lambda)})}} \sum_{i=1}^n X_i^{(\lambda)} \xrightarrow{p} 0$$

since the denominator converges to 0. This is a contradiction (since $\mathcal{N}(0, 1)$ is nondegenerate), so series (3) must converge. \square

Now, we need one last lemma before we can prove the strong law of large numbers.

Lemma 1.3.4 (Kronecker's lemma). *If $a_n \uparrow \infty$ and $\sum_{n=1}^{\infty} b_n/a_n$ converges then $\frac{1}{a_n} \sum_{m=1}^n b_m \rightarrow 0$.*

Proof. The proof follows from *summation by parts*. \square

Kronecker's lemma is useful because it changes questions about averages into questions about sums, and we can use the three-series lemma to handle sums.

Theorem 1.3.5 (Strong law of large numbers). *Let $(X_n)_{n=1}^{\infty}$ be i.i.d. with $\mathbb{E}[X_n] = \mu$. Then, we have $\frac{1}{n} \sum_{i=1}^n X_i \xrightarrow{a.s.} \mu$.*

Proof. Apply Kolmogorov's three series lemma ([Theorem 1.3.3](#)) with $\tilde{X}_n = X_n/n$. Then, use Kronecker's lemma to upgrade almost sure convergence of the series to almost sure convergence of the average. \square

We chose to prove the SLLN here using the Kolmogorov three series lemma in order to learn how to deal with almost sure convergence of series. However, there is a slightly simpler proof of the SLLN using Riesz's lemma of the rising sun, which generalizes to prove the strong ergodic theorem.

Alternate proof of Theorem 1.3.5. The intuition for this proof comes from thinking of S_n as a stochastic process which is asymptotically linear with slope μ . Let $E_\alpha = \{\limsup_{n \rightarrow \infty} S_n/n > \alpha\}$ and $F_\beta = \{\liminf_{n \rightarrow \infty} S_n/n < \beta\}$, with $G = \{\lim_{n \rightarrow \infty} S_n/n \text{ exists and is equal to } \mu\}$. In particular, we have

$$G^c = \bigcup_{m=1}^{\infty} (E_{\mu+1/m} \cup F_{\mu-1/m}),$$

so it suffices to show that $\mathbb{P}(E_\alpha) = 0$ for $\alpha > \mu$. Since E_α is tail-measurable, its probability is either 0 or 1 by Kolmogorov's 0-1 law. Suppose that there is $\alpha \in \mathbb{R}$ such that $\mathbb{P}(E_\alpha) = 1$; we want to show that this forces $\alpha \leq \mu$. By stationarity of $(X_n)_{n=1}^{\infty}$, we have

$$\mathbb{P}\left(\sup_{n > k} \frac{S_n - S_k}{n - k} > \alpha\right) = 1.$$

Now, we need a lemma before we continue the proof of the SLLN.

Lemma 1.3.6 (Riesz's lemma of the rising sun). *Fix $\{X_1, \dots, X_M\}$ and $n \in [M]$. Then, we say $n \in L$ if $\max_{n < t \leq M} \frac{S_t - S_n}{t - n} \leq \alpha$ and $n \in D$ otherwise. Intuitively, imagine we're plotting S_n over time and there is light shining from the right with slope α . Here, L is the set of points which are in the light and D is the set of points in the dark. Then, we have*

$$\frac{1}{|D|} \sum_{n-1 \in D} X_n \geq \alpha.$$

This means that jumps after dark points contribute an average of at least α to the sum.

Proof. The trick of this proof is mostly visual. We can decompose D into disjoint intervals and it suffices to obtain a bound on each dark interval $I = [i, j]$. Then, it is easy to show that S_{j+1} keeps all of I in the dark; in particular, this means that

$$S_{j+1} - S_i = \sum_{n=i+1}^{j+1} X_n \geq \alpha(j - i) = \alpha|I|. \quad \square$$

We now complete the proof of the SLLN using Riesz's lemma of the rising sun. Fix $\epsilon > 0$ and choose N large enough that $\mathbb{P}(\max_{1 \leq n \leq N} S_n/n > \alpha) > 1 - \epsilon$ and

$$\mathbb{E} \left[|X_1| \mathbf{1}_{\max_{1 \leq n \leq N} S_n/n \leq \alpha} \right] < \epsilon,$$

by the dominated convergence theorem. Then if $M \geq N$, we know that

$$\mu M = \mathbb{E}[S_M] = \mathbb{E}[X_1] + \mathbb{E} \left[\sum_{n-1 \in D} X_n \right] + \mathbb{E} \left[\sum_{\substack{n-1 \in L \\ n \leq M-N}} X_n \right] + \mathbb{E} \left[\sum_{\substack{n-1 \in L \\ n > M-N}} X_n \right].$$

To deal with the first term, note that $\mathbb{P}(n \in D) > 1 - \epsilon$ for $n \leq M - N$. By Riesz's rising sun lemma, we have the estimate

$$\mathbb{E} \left[\sum_{n-1 \in D} X_n \right] \geq \alpha \mathbb{E}[|D|] = \alpha \sum_{n=1}^M \mathbb{P}(n \in D) \geq \alpha(1 - \epsilon)(M - N).$$

We also obtain

$$\mathbb{E} \left[\sum_{n=1}^{M-N} X_n \mathbf{1}_{n-1 \in L} \right] \geq - \sum_{n=1}^{M-N} \mathbb{E}[|X_n| \mathbf{1}_{n-1 \in L}] \geq -\epsilon(M - N)$$

and

$$\mathbb{E} \left[\sum_{\substack{n-1 \in L \\ n > M-N}} X_n \right] \geq - \sum_{n > M-N} \mathbb{E}[|X_n|] \geq -E[|X_1|] N.$$

Putting all the bounds together and taking $M \rightarrow \infty$, we obtain $\mu \geq \alpha(1 - \epsilon) - \epsilon$. Letting $\epsilon \downarrow 0$ gives the result. \square

1.4 Disintegration of measure and regular conditional probability

In this section, we formalize the notion of conditioning on a set of measure zero using *regular conditional probabilities*. In the following, denote by $\mathcal{P}(\mathcal{X})$ the set of probability measures over \mathcal{X} . Intuitively, the disintegration of measure theorem says that we can write any measure as an average of measures for each point $x \in \mathcal{X}$, and this disintegration is essentially unique.

Theorem 1.4.1 (Disintegration of measure). *Let \mathcal{X} and \mathcal{Y} be Polish spaces (complete and separable metric spaces), let $\mu \in \mathcal{P}(\mathcal{X} \times \mathcal{Y})$, and denote by $\mu_{\mathcal{X}}$ the marginal of μ on \mathcal{X} . Then, there exists a measurable map $x \mapsto \pi_x$ from \mathcal{X} into $\mathcal{P}(\mathcal{Y})$ (uniquely determined $\mu_{\mathcal{X}}$ -almost everywhere) such that*

$$\mu = \int_{\mathcal{X}} (\delta_x \times \pi_x) d\mu_{\mathcal{X}}(x)$$

Here, the topology on $\mathcal{P}(\mathcal{Y})$ is the weak-* topology induced by weak convergence.

The proof is lengthy so we outline it here but leave some technical details to the reader.

Proof. In a Polish space, there is always a countable family of bounded continuous functions $\{f_n\}_{n=1}^{\infty}$ that separates point; the distance function is continuous, so we can consider the sequence of distance functions from elements of a countable dense subset by separability (potentially capping them off at rationals). For each $f \in C_b(\mathcal{Y})$, we can define a measure ν_f on \mathcal{X} by

$$\nu_f(A) = \int_{A \times \mathcal{Y}} f(y) d\mu(x, y).$$

Since $\nu_f \ll \mu_{\mathcal{X}}$, we can define the Radon-Nikodym derivative $h_f : \mathcal{X} \rightarrow \mathbb{R}$ such that

$$\nu_f(A) = \int_A h_f(x) d\mu_{\mathcal{X}}(x).$$

Since the h_f are defined $\mu_{\mathcal{X}}$ -almost everywhere, we can define the linear functional $L_x(f_n) = h_{f_n}(x)$ for each $n \in \mathbb{N}$, for $\mu_{\mathcal{X}}$ -almost every $x \in \mathcal{X}$. Since L_x is linear and continuous, we can extend the definition to all of $C_b(\mathcal{Y})$ by the Stone-Weierstrass theorem and the fact that $\{f_n\}_{n=1}^{\infty}$ separates points. In particular, linear combinations of the f_n are dense in $C_b(\mathcal{Y})$ with respect to uniform convergence on compact sets and we can easily define L_x on these linear combinations. Note that $L_x(f) \geq 0$ when $f \geq 0$ and $L_x(1) = 1$, so by the Riesz-Markov representation theorem, there exists a unique probability measure π_x on \mathcal{Y} such that

$$h_f(x) = L_x(f) = \int_{\mathcal{Y}} f(y) d\pi_x(y).$$

It can be verified that any bounded continuous function $f : \mathcal{X} \rightarrow \mathcal{P}(\mathcal{Y})$, the map $x \mapsto h_f(x)$ is measurable so $x \mapsto \pi_x$ is measurable as well. It's now clear that $\mu_{\mathcal{X}}$ -almost every $x \in \mathcal{X}$, the measure π_x concentrates

on the *fiber* $\pi_{\mathcal{X}}^{-1}(x)$ (where $\pi_{\mathcal{X}}$ is the projection from $\mathcal{X} \times \mathcal{Y}$ into the first coordinate). For bounded measurable functions of the form $g(x)f(y)$, one can verify that

$$\int_{\mathcal{X} \times \mathcal{Y}} g(x)f(y) d\mu(x, y) = \int_{\mathcal{X}} g(x) \left(\int_{\mathcal{Y}} f(y) d\pi_x(y) \right) d\mu_{\mathcal{X}}(x),$$

so we can extend this to all bounded measurable functions by Dynkin's π - λ theorem, showing that

$$\mu = \int_{\mathcal{X}} (\delta_x \times \pi_x) d\mu_{\mathcal{X}}(x).$$

Essential uniqueness follows immediately because for any bounded measurable function f ,

$$\int_{\mathcal{X}} \int_{\mathcal{Y}} f(x, y) d\pi_x(y) d\mu_{\mathcal{X}}(x) = \int_{\mathcal{X} \times \mathcal{Y}} f(x, y) d\mu(x, y)$$

forces π_x to be unique $\mu_{\mathcal{X}}$ -almost everywhere. □

The function $(x, A) \mapsto \pi_x(A)$ is called a *regular conditional probability* and is denoted by $\mathbb{P}(A \mid X = x)$. The disintegration of measure theorem shows that we can think of π_x as a conditional probability given x , but that such a thing is only uniquely defined for $\mu_{\mathcal{X}}$ -almost every $x \in \mathcal{X}$. In particular, if we let $(X, Y) \sim \mu$ we can think of π_x as a conditional probability given x in the sense that

$$\int_A \int_{\mathcal{Y}} f(y) d\pi_x(y) d\mu_{\mathcal{X}}(x) = \int_A \mathbb{E}[f(Y) \mid X = x] d\mu_{\mathcal{X}}(x) = \int_{A \times \mathcal{Y}} f(y) d\mu(x, y)$$

for any bounded measurable function $f : \mathcal{Y} \rightarrow \mathbb{R}$ and any measurable set $A \in \mathcal{X}$.

Chapter 2

Functional analysis

The study of functional analysis is a generalization of linear algebra to infinite-dimensional spaces. The key algebraic structures are normed spaces, inner product spaces, Banach spaces, and Hilbert spaces. In this chapter, we cover several important theorems of functional analysis to provide intuition for these spaces. We assume an undergraduate background in linear algebra and real analysis. A lot of this chapter is taken from *Introductory Functional Analysis with Applications* by Erwin Kreyszig, although some proofs come from other sources and the content has been reorganized so that it makes the most sense to me.

2.1 Key algebraic structures

First, we begin with a few basic results about normed spaces. Note that all theorems for arbitrary normed spaces or inner product spaces can be specialized to Banach spaces or Hilbert spaces respectively.

2.1.1 Normed spaces and Banach spaces

Definition 2.1.1 (Banach space). A *Banach space* is a complete normed space.

Here are a few examples of normed spaces and Banach spaces to keep in mind.

- \mathbb{Q} with the absolute value is not a Banach space.
- $C([0, 1])$ with the 2-norm is not a Banach space.
- \mathbb{R}^n with the Euclidean norm is a Banach space.
- $C([0, 1])$ with the sup-norm is a Banach space.
- $L^p(\mathbb{R})$ with the L^p norm is a Banach space.
- $L^p(\mathbb{N})$ (also called ℓ^p) with the p -norm is a Banach space.

We have two notions of basis in infinite dimensions.

Definition 2.1.2 (Hamel basis). A *Hamel basis* for a vector space X is a set of linearly independent vectors that span X .

Note that every vector space has a Hamel basis due to Zorn's lemma ([Axiom A.1.1](#)). However, Hamel bases are not very useful in infinite dimensions because they are not unique and may not be countable.

Definition 2.1.3 (Schauder basis). A *Schauder basis* for a normed space X is a sequence of vectors $(e_n)_{n \in \mathbb{N}}$ such that every $x \in X$ can be written as a unique series $\sum_{n=1}^{\infty} \alpha_n e_n$ where the series converges in the norm of X .

A Hamel basis corresponds to a basis in the sense of linear algebra. A Schauder basis allows us to write any vector as an infinite series and not just a finite linear combination. As an example, note that the standard basis is a Schauder basis for $L^p(\mathbb{N})$ for $1 \leq p < \infty$ but not for $L^\infty(\mathbb{N})$ (e.g., take $x = (1, 1, 1, \dots)$).

Proposition 2.1.1. *If a normed space X has a Schauder basis, then it is separable.*

Proof. Pick rational coefficients in the series representation. □

Definition 2.1.4 (Absolute convergence). A series $\sum_{n=1}^{\infty} x_n$ (where the x_n are in a normed space) is said to *converge absolutely* if the series $\sum_{n=1}^{\infty} \|x_n\|$ converges.

The key lemma that we use to get basic results in finite-dimensional normed spaces is the following.

Lemma 2.1.2 (Quantitative independence bound). *If $(x_i)_{i=1}^n$ is independent, then there exists $c > 0$ such that for all scalars $(\alpha_i)_{i=1}^n$, we have*

$$\left\| \sum_{i=1}^n \alpha_i x_i \right\| \geq c \sum_{i=1}^n |\alpha_i|.$$

Proof. Suppose w.l.o.g. that $\sum_{i=1}^n |\alpha_i| = 1$ and proceed by contradiction. Now pick a sequence of vectors $\alpha^{(k)}$ such that

$$\left\| \sum_{i=1}^n \alpha_i^{(k)} x_i \right\| \leq \frac{1}{k}.$$

The $\alpha^{(k)}$ are bounded so extract a convergent subsequence by Bolzano-Weierstrass and deduce a contradiction by independence and continuity of the norm. □

There are a few immediate consequences of the quantitative independence bound for finite-dimensional normed spaces. The proofs are easy so we omit them.

Corollary 2.1.2.1. *Finite-dimensional normed spaces are complete (and therefore closed).*

Corollary 2.1.2.2. *All norms on finite-dimensional spaces are equivalent (i.e., there are constants $0 < c_1 < c_2$ such that $c_1\|x\|_0 \leq \|x\|_1 \leq c_2\|x\|_0$ for all x).*

Using only the key lemma, we can in fact show the Riesz's lemma.

Lemma 2.1.3 (Riesz). *If Y is a strict closed subspace of a normed space X , then for all $\theta \in (0, 1)$ there exists $u \in X$ with $\|u\| = 1$ such that*

$$\inf_{y \in Y} \|u - y\| \geq \theta.$$

Riesz's lemma says that you can find unit vectors far from a closed space in any normed space. The closedness assumption is important because c_0 (the set of sequences with finitely many nonzero terms) is a non-closed subspace of ℓ^2 but there is no unit vector in ℓ^2 far from c_0 .

Proof. Start by picking a nonzero $x_0 \in X \setminus Y$. Then $a = \inf_{y \in Y} \|x_0 - y\| > 0$ since Y is closed. Find $y_0 \in Y$ which is pretty close to x_0 :

$$a \leq \|x_0 - y_0\| \leq \frac{a}{\theta}.$$

Letting $u = (x_0 - y_0)/\|x_0 - y_0\|$, it's easy to see that u satisfies the conditions in Riesz's lemma. \square

Using Riesz's lemma, we can prove the following fundamental result.

Theorem 2.1.4 (Compactness of the unit ball). *The closed unit ball in a normed space is compact if and only if the space is finite-dimensional.*

Proof. The reverse direction is obvious by Heine-Borel. For the forward direction, one can inductively extract a sequence of points in the unit ball without any convergent subsequence using Riesz's lemma (Lemma 2.1.3). \square

2.1.2 Linear operators

Definition 2.1.5 (Bounded operator). A linear operator $T : X \rightarrow Y$ between normed spaces is *bounded* if

$$\|T\| := \sup_{x \neq 0} \frac{\|T(x)\|}{\|x\|} < \infty.$$

Proposition 2.1.5 (Boundedness is equivalent to continuity). *A linear operator $T : X \rightarrow Y$ between normed spaces is bounded if and only if it is continuous (or even continuous at a point).*

Proof. Assuming boundedness, $\|T(x - y)\| \leq c\|x - y\|$ implies continuity. For the converse, suppose T is continuous at $x_0 \in X$. Then, for all $\epsilon > 0$, there exists $\delta > 0$ such that $\|x - x_0\| \leq \delta$ implies $\|T(x) - T(x_0)\| \leq \epsilon$. The rest of the proof proceeds by picking a small vector pointing from x_0 in the direction of x , which is a common technique in functional analysis. For any $x \neq 0 \in X$ associate with it the vector $\tilde{x} = x + \delta x / \|x\|$ so that $\|\tilde{x} - x_0\| = \delta \implies (\delta / \|x\|) T(x) = \|T(\tilde{x}) - T(x_0)\| \leq \epsilon$. Rearrange to obtain the result. \square

The set of bounded linear operators between two normed spaces X and Y (denoted $B(X, Y)$) is itself a normed space with the operator norm.

Proposition 2.1.6. *$B(X, Y)$ is a normed space with the operator norm and is complete if Y is complete.*

Note that only the output space needs to be complete for $B(X, Y)$ to be complete.

Proof. It's obvious that $B(X, Y)$ is a normed space. If $(T_n)_{n=1}^\infty$ is Cauchy then $(T_n x)_{n=1}^\infty$ is Cauchy for all x and converges in Y to y_x . Define $Tx = y_x$ and show it's in $B(X, Y)$. \square

2.1.3 Hilbert spaces

Definition 2.1.6 (Hilbert space). A *Hilbert space* is a complete inner product space.

Easy algebra gives the following lemma.

Lemma 2.1.7 (Parallelogram equation). *For all x, y in any inner product space X , we have*

$$\|x + y\|^2 + \|x - y\|^2 = 2(\|x\|^2 + \|y\|^2).$$

Notably, not all norms come from inner products (e.g., the sup-norm on $C([0, 1])$ or the p -norm on ℓ^p spaces). The proofs all rely on the fact that given a norm in an inner product space, we can recover the inner product.

Lemma 2.1.8 (Polarization identity). *If X is an inner product space over \mathbb{R} , then*

$$\langle x, y \rangle = \frac{1}{4} (\|x + y\|^2 - \|x - y\|^2).$$

If X is an inner product space over \mathbb{C} , then

$$\langle x, y \rangle = \frac{1}{4} (\|x + y\|^2 - \|x - y\|^2 + i(\|x + iy\|^2 - \|x - iy\|^2)).$$

Interestingly, “rotations” don't exist in complex inner product spaces.

Proposition 2.1.9. *Let $Q : X \rightarrow X$ be a bounded linear operator.*

- First, $\langle Qx, y \rangle = 0$ for all $x \in X$ and $y \in Y$ if and only if $Q = 0$.
- In fact, if X is complex and $Q : X \rightarrow X$ is a bounded linear operator then $\langle Qx, x \rangle = 0$ for all $x \in X$ suffices to force $Q = 0$.

Proof. The first statement is obvious by setting $y = Qx$. For the second statement, notice that for all $\alpha \in \mathbb{C}$ and all $x, y \in X$, we have

$$\begin{aligned} 0 &= \langle Q(x + \alpha y), x + \alpha y \rangle \\ &= \langle Qx, x \rangle + \alpha \langle Qy, x \rangle + \bar{\alpha} \langle Qx, y \rangle + |\alpha|^2 \langle Qy, y \rangle \\ &= \alpha \langle Qy, x \rangle + \bar{\alpha} \langle Qx, y \rangle. \end{aligned}$$

We conclude by setting $\alpha = i$ and $\alpha = 1$ and using the first statement. \square

Now we prove the Cauchy-Schwarz inequality.

Theorem 2.1.10 (Cauchy-Schwarz). *For all x, y in an inner product space X , we have*

$$|\langle x, y \rangle| \leq \|x\| \|y\|,$$

with equality if and only if x and y are linearly dependent.

Proof. We study $0 \leq \|x - \alpha y\|^2$ for $\alpha \in \mathbb{C}$. Standard geometry in \mathbb{R}^n endowed with the usual inner product suggests that putting $\alpha = \frac{\langle x, y \rangle}{\|y\|^2}$ should intuitively minimize the expression on the RHS. Now just expand and rearrange. The derivation then shows that equality holds if and only if $x - \alpha y = 0$ or $y = 0$. \square

The Cauchy-Schwarz inequality implies the triangle inequality and continuity of the inner product. One of the most important properties of a Hilbert space is the existence of orthogonal projections, which we now show.

Theorem 2.1.11 (Projection onto a convex set). *If $M \subseteq X$ is convex and complete for any inner product space X , then for all x there exists a unique $y \in M$ such that*

$$\delta = \inf_{\tilde{y} \in M} \|x - \tilde{y}\| = \|x - y\|.$$

Proof. Pick $y_n \in M$ such that $\|x - y_n\| \downarrow \delta$. Then if $v_n := y_n - x$, the parallelogram equality gives

$$\|y_n - y_m\|^2 = 2\|v_n\|^2 + 2\|v_m\|^2 - \|v_n + v_m\|^2.$$

The first two terms go to δ , and convexity of M gives

$$\|v_n + v_m\| = \|y_n + y_m - 2x\| = 2 \left\| \frac{y_n + y_m}{2} - x \right\| \geq 2\delta.$$

So $(y_n)_{n=1}^\infty$ is Cauchy, and completeness of M gives the result. Another application of the parallelogram equality shows that the projection is unique. \square

Theorem 2.1.11 leads to the following theorem.

Theorem 2.1.12 (Orthogonal decomposition). *If Y is a closed subspace of a Hilbert space X then $X = Y \oplus Y^\perp$ (everything in X is a unique sum of something in Y and something in Y^\perp).*

2.1.4 Orthogonality and Fourier series

Suppose we have an orthonormal set $(e_n)_{n=1}^\infty$ in an inner product space X .

Theorem 2.1.13 (Bessel's inequality). *For all $x \in X$, we have*

$$\sum_{n=1}^{\infty} |\langle x, e_n \rangle|^2 \leq \|x\|^2.$$

Proof. Let $y = \sum_{k=1}^n \langle x, e_k \rangle e_k$ and let $z = x - y$ so that $\langle z, e_k \rangle = 0$ for all $1 \leq k \leq n$. So $\|x\|^2 = \|y\|^2 + \|z\|^2$ and the result follows from $\|z\| \geq 0$. \square

The terms $\langle x, e_n \rangle$ are called the *Fourier coefficients* of x .

Theorem 2.1.14 (Convergence of orthonormal series). *The series*

$$\sum_{n=1}^{\infty} \alpha_n e_n$$

converges if and only if

$$\sum_{n=1}^{\infty} |\alpha_n|^2 < \infty.$$

Also if the limit is x then $\alpha_n = \langle x, e_n \rangle$ are the Fourier coefficients.

Proof. Use the Pythagorean identity $\left\| \sum_{n=1}^N \alpha_n e_n \right\|^2 = \sum_{n=1}^N |\alpha_n|^2$ to show the first statement. The second statement follows from continuity of the inner product, since the partial sums are assumed to converge to x . \square

Remarkably, there can only be countably many non-zero Fourier coefficients, no matter how large the orthonormal set is.

Theorem 2.1.15 (Riesz-Fischer). *If $(e_i)_{i \in I}$ is an orthonormal set in a Hilbert space X , then for all $x \in X$ the set $\{i : \langle x, e_i \rangle \neq 0\}$ is countable.*

Proof. By Bessel's inequality, the set $\{i : |\langle x, e_i \rangle| \geq 1/k\}$ is finite for all $k \in \mathbb{N}$. \square

A Schauder basis may not be spanning, so we define the following generalization of a spanning set.

Definition 2.1.7 (Total set). A *total set* $M \subseteq X$ has $\overline{\text{span}(M)} = X$.

All nontrivial Hilbert spaces have total orthonormal sets by Gram-Schmidt and Zorn's lemma (Axiom A.1.1). All total orthonormal sets have the same cardinality, which is called the *Hilbert dimension* of X . The proof is set-theoretic and not difficult, so we omit it. We can also equivalently characterize total sets in Hilbert space.

Proposition 2.1.16. *If X is a Hilbert space, then M is total if and only if $x \perp M$ implies $x = 0$.*

Proof. The forward direction is true in any inner product space by continuity of the inner product. The reverse direction follows from the orthogonal decomposition $X = \overline{\text{span}(M)} \oplus \overline{\text{span}(M)}^\perp$. \square

We have another characterization of totality for orthonormal sets.

Theorem 2.1.17 (Parseval). *An orthonormal set M in a Hilbert space X is total if and only if for all $x \in X$, we have*

$$\|x\|^2 = \sum_{m \in M} |\langle x, m \rangle|^2.$$

This is called the Parseval relation.

Proof. For the reverse direction, if there existed a nonzero $x \perp M$ (applying Proposition 2.1.16), then x cannot satisfy the Parseval relation. For the forward direction, we can define $y = \sum_{m \in M} |\langle x, m \rangle|^2 m$ and show that $x - y \perp M$. Then $x - y = 0$ by Proposition 2.1.16. \square

Orthonormal sets in a separable Hilbert space behave nicely.

Theorem 2.1.18. *If X is a separable Hilbert space then every orthonormal set is countable. Also, if there is a total orthonormal set then X is separable.*

Proof. For the first statement, let $(e_n)_{n=1}^\infty$ be an orthonormal set. Since $\|e_i - e_j\|^2 = 2$ for all $i \neq j$, neighborhoods of size $\sqrt{2}/2$ around the e_i are disjoint, but we need a countable dense set. The second statement follows from picking rational coefficients. \square

Finally, we have the following result about isomorphism, which is analogous to the usual one for finite-dimensional spaces.

Theorem 2.1.19. *Hilbert spaces are isomorphic if and only if they have the same Hilbert dimension.*

Proof. The forward direction is obvious and the reverse direction follows from letting $T(x) = \sum_{k \in K} \langle x, e_k \rangle f_k$ where $(e_k)_{k \in K}$ and $(f_k)_{k \in K}$ are orthonormal bases for the two Hilbert spaces respectively. \square

2.2 The dual space and the Riesz-Fréchet representation theorem

We call the space $X^* := B(X, \mathbb{F})$ the *dual space* of X where \mathbb{F} is \mathbb{R} or \mathbb{C} ; [Proposition 2.1.6](#) gives us some intuition about the dual space.

Proposition 2.2.1. *The dual space X^* is always a Banach space.*

Proof. \mathbb{R} and \mathbb{C} are complete, so use [Proposition 2.1.6](#) (the space of bounded operators between X and Y is complete whenever Y is complete). \square

Now, we have the Riesz-Fréchet representation theorem.

Theorem 2.2.2 (Riesz-Fréchet representation). *If X is a Hilbert space, then for all $f \in X^*$ there exists a unique $y \in X$ such that*

$$f(x) = \langle x, y \rangle$$

for all $x \in X$ and $\|y\| = \|f\|$.

The Riesz-Fréchet representation theorem isn't surprising; \mathbb{R} or \mathbb{C} are one-dimensional Hilbert spaces over themselves, so f has to squash everything into one dimension. The only way to do that is by taking the inner product with a fixed vector.

Proof. The proof is trivial when $f = 0$ so assume $f \neq 0$. Now we want to show $\ker(f)^\perp$ is one-dimensional. For any nonzero $z_1, z_2 \in \ker(f)^\perp$, we have

$$f\left(z_1 - \frac{f(z_1)}{f(z_2)} z_2\right) = 0$$

So $z_1 - \frac{f(z_1)}{f(z_2)} z_2$ is in $\ker(f) \cap \ker(f)^\perp$, which shows that z_1 and z_2 are dependent. Now pick any $z_0 \in \ker(f)^\perp$ with norm 1 and define $y = \overline{f(z_0)} z_0$ so that $f(x) = \langle x, y \rangle$ for all x . Essentially, we are projecting onto $\text{span}(z_0)$ and applying f . Uniqueness is immediate since $0 = \langle y_1 - y_2, y_1 \rangle - \langle y_1 - y_2, y_2 \rangle = \|y_1 - y_2\|^2$ implies $y_1 = y_2$ for alternatives y_1 and y_2 . Finally, we have $\|y\|^2 = f(y) \leq \|f\| \|y\|$ so $\|y\| \leq \|f\|$, and the reverse inequality follows from Cauchy-Schwarz. \square

2.3 The four pillars of functional analysis

In this section, we prove several important theorems in functional analysis, which are often called the “four pillars” of the subject. The Hahn-Banach theorem works in general normed spaces, but the other three theorems are specific to Banach spaces. The latter three will follow mostly from the Baire category theorem ([Theorem A.2.1](#)).

2.3.1 The Hahn-Banach theorem

The Hahn-Banach theorem allows us to extend bounded linear functionals from subspaces of any vector space to the entire space, and shows that the dual contains a lot of things. We begin with a related definition.

Definition 2.3.1 (Sublinear functional). A *sublinear functional* on a vector space X is a function $p : X \rightarrow \mathbb{R}$ such that

$$p(\alpha x) = \alpha p(x)$$

for $\alpha \geq 0$ and

$$p(x + y) \leq p(x) + p(y)$$

for all $x, y \in X$.

An important example of a sublinear functional is the norm on a normed space.

Theorem 2.3.1 (Hahn-Banach). *Suppose Z is a subspace of a vector space X and $p : X \rightarrow \mathbb{R}$ is a sublinear functional. If $f : Z \rightarrow \mathbb{R}$ is a linear functional such that $f(x) \leq p(x)$ for all $x \in Z$, then there exists a linear functional $\tilde{f} : X \rightarrow \mathbb{R}$ such that $\tilde{f}|_Z = f$ and $\tilde{f}(x) \leq p(x)$ for all $x \in X$.*

Proof. Apply Zorn's lemma (Axiom A.1.1) with $f_1 \leq f_2$ if f_2 extends f_1 on the poset $M = \{g : g \text{ linear, } g \text{ extends } f, g(x) \leq p(x)\}$. Now we have the existence of a maximal element $\tilde{f} : Y \rightarrow \mathbb{R}$. Suppose for a contradiction that \tilde{f} isn't defined at some $z \in X$. Then, define $g(y + \alpha z) = \tilde{f}(y) + \alpha \tilde{f}(z)$ for all $y \in Y$ and $\alpha \in \mathbb{R}$; we need to pick $\tilde{f}(z)$ such that $g(x) \leq p(x)$. Rephrasing: for all $y \in Y$ and $\alpha > 0$, we need

$$\begin{aligned} \tilde{f}(y) + \alpha \tilde{f}(z) &\leq p(y + \alpha z), \\ \tilde{f}(y) - \alpha \tilde{f}(z) &\leq p(y - \alpha z). \end{aligned}$$

Solving for the constraint on $\tilde{f}(z)$ and setting $y \mapsto y/\alpha$, we find:

$$\sup_{y \in Y} \{f(y) - p(y - z)\} \leq \tilde{f}(z) \leq \inf_{y \in Y} \{p(y + z) - f(y)\}.$$

As long as the left-hand side is not larger than the right-hand side, we will have a valid choice for $\tilde{f}(z)$. But we have

$$f(y_2) - p(y_2 - z) \leq p(y_1 + z) - f(y_1) \iff f(y_1 + y_2) \leq p(y_1 + z) + p(y_2 - z),$$

and the latter condition is implied by sublinearity of p . □

We can slightly generalize the Hahn-Banach theorem to complex vector spaces.

Theorem 2.3.2 (Generalized Hahn-Banach). *Suppose Z is a subspace of a vector space X and $p : X \rightarrow \mathbb{R}$ is a subadditive functional; for all $x, y \in X$ and $\alpha \in \mathbb{C}$, we have*

$$p(\alpha x) = |\alpha|p(x)$$

and

$$p(x + y) \leq p(x) + p(y).$$

If $f : Z \rightarrow \mathbb{C}$ is a linear functional such that $|f(x)| \leq p(x)$ for all $x \in Z$, then there exists a linear functional $\tilde{f} : X \rightarrow \mathbb{C}$ such that $\tilde{f}|_Z = f$ and $|\tilde{f}(x)| \leq p(x)$ for all $x \in X$.

Proof. Use Hahn-Banach on the real and imaginary parts of f separately and set $\tilde{f}(x) = \tilde{f}_{\Re}(x) - i\tilde{f}_{\Im}(x)$. \square

Note that Hahn-Banach is a generalization of the Riesz-Fréchet representation theorem ([Theorem 2.2.2](#)) to arbitrary vector spaces; if Z is a closed subspace of a Hilbert space X then the Riesz-Fréchet representation gives a linear extension $\tilde{f}(x) = \langle x, z \rangle$ for some $z \in X$. We now state the most important corollary of the Hahn-Banach theorem, which shows that there are lots of bounded linear functionals on a normed space.

Corollary 2.3.2.1. *If X is a normed space and $x_0 \in X$ is nonzero, then there exists a bounded linear functional $\tilde{f} : X \rightarrow \mathbb{F}$ such that $\tilde{f}(x_0) = \|x_0\|$ and $\|\tilde{f}\| = 1$.*

Proof. Apply the Hahn-Banach theorem with $Z = \text{span}(x_0)$ and $f(z) = \|z\|$ defined from Z to \mathbb{F} . \square

This implies that the set of bounded linear functionals separates points in a normed space, and X^* therefore has a rich structure. This becomes useful later in the characterization of weak convergence and the adjoint operator. For applications of Hahn-Banach, see [Sections 2.5](#) and [2.6](#).

2.3.2 The uniform boundedness principle (Banach-Steinhaus)

The uniform boundedness principle allows us to upgrade pointwise convergence of operators to uniform convergence, and follows from the Baire category theorem ([Theorem A.2.1](#)).

Theorem 2.3.3 (Banach-Steinhaus). *If X is a Banach space and Y is a normed space, then for all $T_n : X \rightarrow Y$ such that*

$$\sup_{n \in \mathbb{N}} \|T_n x\| < \infty$$

for all $x \in X$, we have

$$\sup_{n \in \mathbb{N}} \|T_n\| < \infty.$$

This is sometimes called the uniform boundedness principle.

Proof. Define $A_k = \{x \in X : \|Tx\| \leq k\}$ so that A_k is closed and $X = \bigcup_{k=1}^{\infty} A_k$. By the Baire category theorem (since X is complete), there exists A_{k_0} containing a ball $B(x_0, r)$. We are almost done, since $B(x_0, r) - x_0$ contains vectors at a fixed length in all directions and $\|Tx\|$ is uniformly bounded in the ball.

To formalize this argument, for any $x \in X$ we can set $z = x_0 + \gamma x$ where $\gamma = r/(2\|x\|)$ so that $z \in B(x_0, r)$. So then we have the inequality

$$\|T_n x\| = \left\| \frac{1}{\gamma} T_n(z - x_0) \right\| \leq \frac{1}{\gamma} (\|T_n z\| + \|T_n x_0\|) \leq \frac{2k_0}{\gamma} = \frac{4k_0}{r} \|x\|.$$

This estimate shows that $\|T_n\|$ is uniformly bounded by $4k_0/r$, and we are done. \square

A common use of the uniform boundedness principle is to construct operators which are pointwise bounded but not uniformly bounded, thereby showing that a space X is not complete. The uniform boundedness principle can also reveal interesting structure in many spaces, as we show in the following example.

Example 2.3.1 (Most continuous functions aren't locally differentiable). There are lots of famous examples of continuous and nowhere differentiable functions (sample paths of Brownian motion, the Weierstrass function, etc.) but the uniform boundedness principle can be used to show that *most* continuous functions are actually nowhere locally differentiable. Consider $X = C([0, 1])$ endowed with the supremum norm and for each $x \in [0, 1]$ and $n \in \mathbb{N}$ we define the linear functional:

$$T_{n,x}(f) = \begin{cases} n(f(x + \frac{1}{n}) - f(x)) & x + \frac{1}{n} \leq 1 \\ n(f(x) - f(x - \frac{1}{n})) & x + \frac{1}{n} > 1. \end{cases}$$

Essentially, $T_{n,x}$ are the difference quotients with $\Delta x = 1/n$; if f is differentiable at x , then $\sup_n |T_{n,x}(f)| < \infty$. Formally, define $D_x = \{f \in X : f \text{ is differentiable at } x\}$ so that $T_{n,x}$ is pointwise bounded on D_x . It is easy to see that $\sup_n \|T_{n,x}\| = \infty$ by constructing a sequence of functions that rise more and more sharply. In particular, by a slight generalization of the uniform boundedness principle to any nonmeager subset of a Banach space, we deduce that D_x must be meager. The set of anywhere locally differentiable functions is exactly the set of functions which are differentiable at any rational, so we deduce that the set of anywhere locally differentiable functions is meager in $C([0, 1])$.

For further example applications of the uniform boundedness principle, see [Section 2.7](#).

2.3.3 The open mapping theorem

The open mapping theorem says that a surjective bounded linear operator between Banach spaces sends open sets to open sets. In particular, a bijective bounded linear operator between Banach spaces has a bounded inverse.

Definition 2.3.2 (Open mapping). $T : X \rightarrow Y$ between metric spaces is an *open mapping* if $T(U)$ is open in Y whenever U is open in X .

Theorem 2.3.4 (Open mapping). *If X and Y are Banach spaces and $T : X \rightarrow Y$ is a surjective bounded linear operator, then T is an open mapping.*

Proof. Let $B_r = B(0, r) \subseteq X$. It will suffice to show that $T(B_1)$ contains an open ball around 0 in Y . We know that

$$Y = \bigcup_{n=1}^{\infty} \overline{T(B_n)},$$

so the Baire category theorem implies that we can fit a ball in $\overline{T(B_n)}$ for some n . Shrinking this ball by a factor of n , we can find a ball $B(y_0, r_0) \subseteq \overline{T(B_1)}$. It is easy to show that $\overline{T(B_1)} - y_0 \subseteq \overline{T(B_2)}$, so that we can center the ball. Our goal is now to show that there is a ball in $T(B_1)$, thereby removing the closure.

Dilating the previous result, there exists $r > 0$ such that as long as $\|y\| < r/2^n$, then $y \in \overline{T(B_{2^{-n}})}$. We're going to show that $B_{r/2} \subseteq T(B_1)$, thereby removing the closure. As long as $\|y\| < r/2$, we can find $x_1 \in B_{1/2}$ such that $\|y - Tx_1\| < r/4$, and inductively, $x_n \in B_{2^{-n}}$ such that $\|y - \sum_{k=1}^n Tx_k\| < r/2^{n+1}$. Since X is complete, $\sum_{k=1}^{\infty} x_k$ converges to some $x \in X$. Now $\|x\| < 1$ and $y = Tx$, so the result follows. \square

The open mapping theorem immediately has a useful corollary, called the bounded inverse theorem.

Corollary 2.3.4.1 (Bounded inverse). *If X and Y are Banach spaces and $T : X \rightarrow Y$ is a bijective bounded linear operator, then T^{-1} is bounded.*

The bounded inverse theorem is used to show that solutions to linear equations in Banach spaces are continuous with respect to the desired output. For instance, if we are trying to solve $Tx = y$ for x and T is a bijective bounded linear operator between Banach spaces, then the solution $x = T^{-1}y$ is continuous in y ; as an example, T might be a differential operator. Even if the map is not bijective, we have the following corollary.

Corollary 2.3.4.2. *Suppose X , Y , and T are as in the open mapping theorem. Then there exists a constant $c > 0$ such that for any $y \in Y$ there exists $x \in X$ with $Tx = y$ and $\|x\| \leq c\|y\|$.*

Proof. By the open mapping theorem (Theorem 2.3.4), $T(B_1)$ contains an open ball $B(0, r)$. For any $y \in Y$, there exists $x \in B_1$ with $Tx = ry/(2\|y\|)$. In particular, we have $T(x(2\|y\|)/r) = y$ with $\|x(2\|y\|)/r\| \leq 2\|y\|/r$. Let $c = 2/r$. \square

So the solution to a linear equation between Banach spaces is always bounded in norm by the desired output, as long as the linear map is surjective, which has important implications for numerical stability. If y has a small measurement error δ , then the solution x will have an error of at most $c\delta$, where c didn't depend on y .

2.3.4 The closed graph theorem

The closed graph theorem gives a sufficient condition for a linear operator between Banach spaces to be bounded, and follows from the open mapping theorem.

Definition 2.3.3 (Closed operator). A linear operator $T : \text{dom}(T) \rightarrow Y$ between normed spaces is *closed* if the graph of T is closed in $X \times Y$ under the norm $\|(x, y)\| = \|x\| + \|y\|$. Here, the graph of T is the set $\mathcal{G}(T) = \{(x, Tx) : x \in \text{dom}(T)\}$.

We can characterize closed operators in another equivalent way.

Proposition 2.3.5. A linear operator $T : \text{dom}(T) \rightarrow Y$ between Banach spaces is closed if and only if for all sequences $(x_n)_{n=1}^\infty$ in $\text{dom}(T)$ such that $x_n \rightarrow x$ and $Tx_n \rightarrow y$, we have $x \in \text{dom}(T)$ and $Tx = y$.

Example 2.3.2 (Differentiation operator). The differentiation operator from $C^1([0, 1])$ to $C([0, 1])$ is closed but not bounded.

Now, we state the closed graph theorem, which gives sufficient conditions for a linear operator between Banach spaces to be bounded.

Theorem 2.3.6 (Closed graph). If X and Y are Banach spaces and $T : \text{dom}(T) \rightarrow Y$ is a closed linear operator, then T is bounded.

Proof. $X \times Y$ is complete, so consider the mapping $p(x, Tx) = x$ from $\mathcal{G}(T)$ to $\text{dom}(T)$. This mapping is bounded and bijective, so by the open mapping theorem, p^{-1} is bounded (since $\mathcal{G}(T)$ and $\text{dom}(T)$ are complete). In particular, this shows that T is bounded. \square

We give sufficient conditions for the converse of the closed graph theorem to hold.

Proposition 2.3.7. If $T : \text{dom}(T) \rightarrow Y$ is a linear operator between any normed spaces then if $\text{dom}(T)$ is complete and T is bounded, then T is closed. Also, if T is closed and Y is complete, then $\text{dom}(T)$ is closed.

Proof. The first statement is immediate from continuity of T . For the second statement, take a sequence $x_n \rightarrow x$ with $x_n \in \text{dom}(T)$. Then, we have

$$\|T(x_n - x_m)\| \leq \|T\| \|x_n - x_m\| \rightarrow 0,$$

so $(Tx_n)_{n=1}^\infty$ is Cauchy. Now $Tx_n \rightarrow y$ for some $y \in Y$ by completeness of Y and $y = Tx$ by closedness of T . \square

The closed graph theorem is used to show that an operator is bounded, since in many cases closedness is easier to verify directly than boundedness.

2.4 Hilbert adjoint operators

First, we slightly generalize the Riesz-Fréchet representation theorem to sesquilinear forms.

Definition 2.4.1 (Sesquilinear form). A *sesquilinear form* is $h : X \times Y \rightarrow \mathbb{F}$ which is linear in the first coordinate and conjugate linear in the second coordinate (where X and Y are vector fields over \mathbb{F}). Define the norm of h as

$$\|h\| = \sup_{\substack{x \neq 0 \\ y \neq 0}} \frac{|h(x, y)|}{\|x\| \|y\|}.$$

Theorem 2.4.1 (Generalized Riesz-Fréchet representation). *If X and Y are Hilbert spaces and h is a bounded sesquilinear form on $X \times Y$ then we can write*

$$h(x, y) = \langle Sx, y \rangle,$$

where $S : X \rightarrow X$ is a unique linear operator with $\|S\| = \|h\|$.

Proof. $\overline{h(x, y)}$ is a bounded linear functional in y for fixed x , so the Riesz-Fréchet representation theorem gives $\overline{h(x, y)} = \langle y, z \rangle$ for some $z \in Y$; this means that $h(x, y) = \langle z, y \rangle$. Define $Sx = z$ for all x ; it's easy to show S is linear and unique. By Cauchy-Schwarz, we have

$$\|h\| \leq \sup_{\|x\|=1} \frac{\|Sx\|}{\|x\|} = \|S\|.$$

Similarly, we get the reverse inequality by

$$\|h\| \geq \sup_{\substack{x \neq 0 \\ y \neq 0}} \frac{|h(x, y)|}{\|x\| \|y\|} = \sup_{\substack{x \neq 0 \\ Sx \neq 0}} \frac{|\langle Sx, Sx \rangle|}{\|x\| \|Sx\|} = \|S\|.$$

\square

One main point of this generalization is to prove that the Hilbert adjoint operator exists.

Definition 2.4.2 (Hilbert adjoint). If $T : X \rightarrow Y$ is bounded linear where X and Y are Hilbert spaces, then the *Hilbert adjoint* of T is the unique operator $T^* : Y \rightarrow X$ such that

$$\langle Tx, y \rangle = \langle x, T^*y \rangle$$

for all $x \in X$ and $y \in Y$.

We now verify that the Hilbert adjoint is well-defined.

Theorem 2.4.2. T^* exists, is unique, and $\|T^*\| = \|T\|$.

Proof. The function $h(x, y) = \langle y, Tx \rangle$ is a bounded sesquilinear form on $X \times Y$, so use the generalized Riesz-Fréchet representation theorem ([Theorem 2.4.1](#)). \square

The adjoint operator is a generalization of the conjugate transpose for finite-dimensional spaces. If $\langle x, y \rangle = x^\top \bar{y}$ then $\langle Bx, y \rangle = (Bx)^\top \bar{y} = x^\top \overline{B^*y}$, so $B^* = \overline{B}^\top$. The properties of the adjoint operator all follow immediately from the definition, except perhaps for the following.

Proposition 2.4.3. The adjoint satisfies $\|T^*T\| = \|T\|^2$.

Proof. We immediately have $\|T^*T\| \leq \|T^*\|\|T\| = \|T\|^2$. To show the reverse inequality, we have for all $x \in X$ that

$$\|Tx\|^2 = \langle Tx, Tx \rangle = \langle x, T^*Tx \rangle \leq \|x\|\|T^*Tx\|.$$

Taking the supremum over all x with $\|x\| = 1$ gives the result. \square

We now have the following important definitions related to the adjoint operator.

Definition 2.4.3 (Self-adjoint operator). An operator $T : X \rightarrow X$ is *self-adjoint* or *Hermitian* if $T = T^*$.

Definition 2.4.4 (Unitary operator). An operator $T : X \rightarrow X$ is *unitary* if $T^* = T^{-1}$.

Self-adjoint operators are a generalization of conjugate symmetric matrices, and unitary operators are a generalization of orthogonal matrices. The adjoint operator has a few nice properties.

Proposition 2.4.4. If T is self-adjoint then $\langle Tx, x \rangle$ is real for all x . Also, the reverse implication is true in complex Hilbert spaces.

Proof. The first statement is obvious and the second statement follows from [Proposition 2.1.9](#). \square

Proposition 2.4.5. The limit of bounded self-adjoint linear operators is bounded and self-adjoint.

Proof. We have by the triangle inequality

$$\|T - T^*\| \leq \|T - T_n\| + \|T_n - T_n^*\| + \|T_n^* - T^*\| = 2\|T - T_n\| \rightarrow 0. \quad \square$$

Unitary operators behave nicely in Hilbert spaces too.

Proposition 2.4.6. *Unitary operators are isometries ($\|Ux\| = \|x\|$), and a bounded linear operator on a complex Hilbert space is unitary if and only if it is isometric and surjective.*

Proof. The first statement and the forward implication of the second statement are easy to show. For the reverse implication, isometries are injective so T is bijective. Now $\langle x, x \rangle = \langle T^*Tx, x \rangle$ implies that $\langle (T^*T - I)x, x \rangle = 0$ for all x ; conclude by [Proposition 2.1.9](#) that $T^*T = I$. \square

2.5 The adjoint operator

Suppose we have a linear functional g on Y . Then, if $T \in B(X, Y)$ then define $f(x) = g(Tx)$ for all $x \in X$.

Definition 2.5.1 (Adjoint operator). The *adjoint operator* $T^\times : Y^* \rightarrow X^*$ is defined by $T^\times(g) = g \circ T$ for all $g \in Y^*$.

Here, $T : X \rightarrow Y$ but $T^\times : Y^* \rightarrow X^*$.

Proposition 2.5.1. *The adjoint operator is linear with $\|T^\times\| = \|T\|$.*

Proof. Linearity is clear and $\|T^\times g\| \leq \|T\| \|g\|$ implies that $\|T^\times\| \leq \|T\|$. Now, for all $x_0 \in X$, [Corollary 2.3.2.1](#) gives $g_0 \in Y^*$ such that $g_0(Tx_0) = \|Tx_0\|$ and $\|g_0\| = 1$. So, we obtain

$$\|Tx_0\| = g_0(Tx_0) = T^\times(g_0)(x_0) \leq \|T^\times\| \|g_0\| \|x_0\| = \|T^\times\| \|x_0\|.$$

This implies that $\|T\| \leq \|T^\times\|$ and the result follows. \square

Notice that for the adjoint operator to have the a large enough norm, we needed the dual space to be large enough. Also, if T has matrix T_E , then the adjoint operator has matrix T_E^\top .

2.6 Reflexive spaces and separability

Proposition 2.6.1. *The canonical mapping $x \mapsto g_x$ where $g_x(f) = f(x)$ is an isomorphism from X to a subset of X^{**} with $\|g_x\| = \|x\|$.*

Proof. The proof is immediate from definitions. \square

Definition 2.6.1 (Reflexive space). A normed space X is *reflexive* if the canonical mapping is surjective.

Note that if X is reflexive then X is complete because the double dual is complete. Furthermore, every Hilbert space is reflexive due to the Riesz-Fréchet representation theorem. Now, we would like to show that if X^* is separable in a normed space then X is separable.

Lemma 2.6.2. *Suppose that Y is a proper closed subspace of a normed space X and fix $x_0 \in X \setminus Y$ with $\delta := \text{dist}(x_0, Y)$. Then there exists $\tilde{f} \in X^*$ such that $\tilde{f}(x_0) = \delta$, and $\tilde{f}(y) = 0$ for all $y \in Y$, and $\|\tilde{f}\| = 1$.*

Proof. Apply the Hahn-Banach theorem (Theorem 2.3.2). \square

Theorem 2.6.3. *If X^* is separable then X is separable for any normed space X .*

Proof. The unit sphere $U^* = \{f : \|f\| = 1\} \subseteq X^*$ is closed so it contains a countable dense subset $(f_n)_{n=1}^\infty$. Find $(x_n)_{n=1}^\infty$ such that $|f_n(x_n)| \geq 1/2$ since $\|f_n\| = 1$. Let $Y = \overline{\text{span}(\{x_n\}_{n=1}^\infty)}$. If $Y \neq X$, then by Lemma 2.6.2 there exists $\tilde{f} \in X^*$ such that $\tilde{f}(y) = 0$ for all $y \in Y$ and $\|\tilde{f}\| = 1$. In particular, we have (since $\tilde{f}(x_n) = 0$ for all n) that

$$\frac{1}{2} \leq |f_n(x_n)| \leq |f_n(x_n) - \tilde{f}(x_n)| \leq \|f_n - \tilde{f}\| \|x_n\| = \|f_n - \tilde{f}\|.$$

This contradicts the density of $(f_n)_{n=1}^\infty$ in U^* , so $Y = X$. Now pick rational coefficients in the span defining Y and conclude. \square

2.7 Weak convergence and weak-* convergence

Definition 2.7.1 (Weak convergence). A sequence $(x_n)_{n=1}^\infty$ converges weakly to x (written $x_n \xrightarrow{w} x$) if $f(x_n) \rightarrow f(x)$ for all $f \in X^*$.

Weak convergence behaves like strong convergence in the following ways.

Proposition 2.7.1. *If $x_n \xrightarrow{w} x$ then all subsequences converge weakly and the limit is unique. Also, $(\|x_n\|)_{n=1}^\infty$ is bounded.*

Proof. The first two statements are easy to show. For the third statement, consider the canonical mapping $g_n(f) = f(x_n)$. Since $(f(x_n))_{n=1}^\infty$ converges, it is bounded for all $f \in X^*$. Because X^* is complete, the uniform boundedness principle (Theorem 2.3.3) implies that $(\|g_n\|)_{n=1}^\infty$ is bounded. Conclude using $\|g_n\| = \|x_n\|$. \square

Weak convergence is a less strict form of strong convergence when $\dim(X) = \infty$.

Proposition 2.7.2. *Strong convergence implies weak convergence, and the converse is true when $\dim(X) < \infty$.*

Proof. To show strong convergence implies weak convergence, we know that $|f(x_n) - f(x)| \leq \|f\| \|x_n - x\| \rightarrow 0$. To show the converse, suppose that $(e_i)_{i=1}^n$ is a basis for X and use the canonical basis for the dual (the one composed of functionals that send e_i to 1 and all other vectors to 0). \square

We can easily write down an equivalent characterization of weak convergence using a total subset of the dual (a set $M \subseteq X^*$ such that $\overline{\text{span}(M)} = X^*$).

Proposition 2.7.3. *Weak convergence $x_n \xrightarrow{w} x$ is equivalent to $(\|x_n\|)_{n=1}^\infty$ being bounded and $f(x_n) \rightarrow f(x)$ for all f in a total set $M \subseteq X^*$.*

Proof. The forward implication is obvious. For the reverse implication, pick any $f \in X^*$ and $(f_j)_{j=1}^\infty$ in $\text{span}(M)$ such that $f_j \rightarrow f$. Then use a 3ϵ argument, since $f(x_n)$ is close to $f_j(x_n)$, $f_j(x_n)$ is close to $f_j(x)$, and $f_j(x)$ is close to $f(x)$. We needed boundedness of $(\|x_n\|)_{n=1}^\infty$ to ensure that $f_j(x_n)$ is close to $f_j(x)$. \square

Weak convergence intuitively captures many forms of convergence not captured by strong convergence; for instance, the sequence $(e_n)_{n=1}^\infty$ in ℓ^2 converges weakly to 0 but not strongly. We can also define notions of convergence for operators.

Definition 2.7.2 (Uniform operator convergence). A sequence of operators $(T_n)_{n=1}^\infty$ converges uniformly to T if $\|T_n - T\| \rightarrow 0$.

Definition 2.7.3 (Strong operator convergence). A sequence of operators $(T_n)_{n=1}^\infty$ is strongly operator convergent to T if $(T_n x)_{n=1}^\infty$ converges for all $x \in X$.

Definition 2.7.4 (Weak operator convergence). A sequence of operators $(T_n)_{n=1}^\infty$ converges weakly to T if $(T_n x)_{n=1}^\infty$ converges weakly for all $x \in X$.

For example, $T_n = (\text{set first } n \text{ coordinates to } 0)$ in ℓ^2 converges strongly but not uniformly to 0. Similarly, $T_n = (\text{shift forward by } n)$ in ℓ^2 converges weakly but not strongly to 0.

Note that for bounded linear functionals, strong operator convergence is equivalent to weak operator convergence. In this case, we call uniform operator convergence *strong convergence* and we call weak operator convergence *weak-* convergence*.

Example 2.7.1 (Convergence in distribution). Convergence of random variables in distribution is an example of weak-* convergence. For a compact set X , the dual of $C(X)$ is the space of regular signed Borel measures on X (sometimes denoted $\mathcal{M}(X)$); this is called the Riesz-Markov representation theorem. The set of probability measures $\mathcal{P}(X)$ is a subset of $\mathcal{M}(X)$, so weak convergence of probability measures exactly means that $\mathbb{E}[f(X_n)] \rightarrow \mathbb{E}[f(X)]$ for all $f \in C(X) = C_b(X)$. Of course, this is the characterization of convergence in distribution given by the portmanteau lemma. Similar analysis works over unbounded

sets X by considering the space of continuous functions vanishing at infinity (called $C_0(X)$), whose dual is the space of Radon measures on \mathbb{R} .

Now, strong operator convergence implies that the limit operator is bounded, assuming that the input comes from a Banach space.

Lemma 2.7.4. *If $T_n \rightarrow T$ is strongly operator convergent then $T \in B(X, Y)$ as long as X is a Banach space.*

Proof. Linearity is immediate, and boundedness follows from the uniform boundedness principle ([Theorem 2.3.3](#)). \square

We have the following characterization of strong operator convergence.

Proposition 2.7.5. *A sequence of operators $(T_n)_{n=1}^\infty$ converges strongly to T if and only if $(\|T_n\|)_{n=1}^\infty$ is bounded and $(T_n x)_{n=1}^\infty$ is Cauchy for all $x \in M$, where $M \subseteq X$ is total.*

Proof. The proof is more or less identical to the proof of [Proposition 2.7.3](#). \square

Since strong operator convergence is equivalent to weak-* convergence for bounded operators, the same characterization holds for weak-* convergence. We now prove the Banach-Alaoglu theorem, which gives a way to show weak-* convergence.

Theorem 2.7.6 (Banach-Alaoglu). *If X is a normed space and $(f_n)_{n=1}^\infty$ is a sequence in X^* with $\|f_n\| \leq 1$ for all $n \in \mathbb{N}$, then there exists a subsequence such that $(f_{n_k})_{k=1}^\infty$ converges in weak-* to some $f \in X^*$.*

This theorem can be stated more generally in the context of topological vector spaces: the closed unit ball of the dual of a normed space is compact in the weak-* topology. However, we state the result here (equivalently) in terms of subsequential weak-* convergence.

Proof. The trick of this proof is to embed our space into a compact product space, such that convergence in the product space automatically gives our desired convergence. We associate with each point $x \in X$ the compact set $K_x = \{ |z| \leq \|x\| \} \subseteq \mathbb{R}$. By Tychonoff's theorem (which states that the product of compact sets is always compact), the product $\prod_{x \in X} K_x$ is compact. Define the map Φ on X^* by $\Phi(f) = (f(x))_{x \in X}$ so that $\Phi(f) \in \prod_{x \in X} K_x$ when $\|f\| \leq 1$. Extract a convergent subsequence $\Phi(f_{n_k}) \rightarrow \Phi(f)$ in the product topology on $\prod_{x \in X} K_x$. Now, for any $x \in X$, we have $f_{n_k}(x) \rightarrow f(x)$ so that $f_{n_k} \xrightarrow{w} f$ as desired. \square

Chapter 3

Convex analysis

In this chapter, we discuss many classical concepts in convex analysis which appear in many other areas of applied math. We will assume that the reader knows some basic definitions, and most of the material will come from *Convex Analysis* by R. Tyrrell Rockafellar.

3.1 Convex sets and functions

Most of this section is intended as review.

Proposition 3.1.1. *Intersections of convex sets are convex.*

Definition 3.1.1 (Epigraph). If $f : \mathbb{R}^n \rightarrow [-\infty, \infty]$ is a function, then the *epigraph* of f is the set

$$\text{epi}(f) = \{(x, \mu) \in \mathbb{R}^n \times \mathbb{R} : f(x) \leq \mu\}.$$

Essentially, if you draw the graph of f , the epigraph is the set of points above the graph (including the graph itself).

Definition 3.1.2 (Proper function). A function f is *proper* if it maps into $(-\infty, \infty]$ and is not identically ∞ .

Definition 3.1.3 (Convex function). A function f is *convex* if its epigraph is a convex set. If f is proper then this is equivalent to Jensen's inequality.

Proposition 3.1.2. *If f is convex then the sublevel sets $\{x : f(x) \leq \alpha\}$ are convex for all $\alpha \in \mathbb{R}$. So are the strict sublevel sets $\{x : f(x) < \alpha\}$.*

Proof. The α -sublevel set is the intersection of $\text{epi}(f)$ with the hyperplane $\mu = \alpha$ (restricted to the first coordinate), and the convexity of the strict sublevel sets follows from Jensen's inequality. \square

Proposition 3.1.3. *If f and φ are proper and convex with φ nondecreasing and $\varphi(\infty) = \infty$, then $\varphi \circ f$ is convex.*

The proof is easy so we omit it. For example, this implies that $e^{f(x)}$ is convex if $f > 0$ is convex.

Proposition 3.1.4. *$f(x) = \inf\{\mu : (x, \mu) \in K\}$ is convex if K is convex.*

Again, we omit the proof. We know that we can associate a convex set to every convex function by $\text{epi}(f)$, but it now follows that we can almost treat any convex set like $\text{epi}(g)$ for some convex function g .

3.2 Lower semi-continuous functions

To motivate this section, notice that convex functions need not have consistent limiting behavior. For example, take the function $f(x) = x^2 + \mathbf{1}_{x \neq 0}$ with $\text{dom}(f) = [0, \infty)$, which is convex. One could try to computationally optimize this function and obtain $x^* = 0$ as an output, but this is totally meaningless. In some sense, this is the only problem that convex functions can have. They can have points where they are unexpectedly high, but they can't have points where they are unexpectedly low. We need to impose some sort of continuity condition to ensure that the function behaves as expected, and the right condition is lower semi-continuity.

Definition 3.2.1 (Lower semi-continuity). A function f is *lower semi-continuous* (l.s.c.) if $f(x) \leq \liminf_{y \rightarrow x} f(y)$ for all $x \in \text{dom}(f)$.

In particular, if we tried to minimize the function only using surrounding values, there would not be an unwelcome surprise at the limit no matter how we approach. We can equivalently characterize lower semi-continuity as follows, highlighting the deep connection between lower semi-continuity and convexity.

Proposition 3.2.1. *A function f is l.s.c. if and only if the sublevel sets $\{x : f(x) \leq \alpha\}$ is closed for all $\alpha \in \mathbb{R}$. Also, f is l.s.c. if and only if $\text{epi}(f)$ is closed.*

The proof is easy so we omit it. We will use this characterization later.

Proposition 3.2.2. *If f takes values in $[0, \infty]$ and is l.s.c. then it can be written as the pointwise supremum of Lipschitz functions.*

Note that the proof can be adapted so that this holds in any metric space.

Proof. For $\lambda \geq 0$, define $f_\lambda(x) = \inf\{f(y) + \lambda\|x - y\|_2 : y \in \mathbb{R}^n\}$; this is called the *Moreau-Yosida trick*. Here, f_λ is the infimal convolution of f with $\lambda\|\cdot\|_2$. It is easy to check that f_λ is λ -Lipschitz and $f_\lambda \uparrow f$ pointwise as $\lambda \rightarrow \infty$. □

3.3 Separation theorems

In this section, we cover the supporting hyperplane theorem and the separating hyperplane theorem, which are fundamental results in convex analysis.

Definition 3.3.1 (Gauge function). The *gauge function* of a set C is the function $\gamma_C(x) = \inf\{\lambda > 0 : x \in \lambda C\}$, and is convex if C is convex. The gauge function is sometimes called the *Minkowski functional* of C .

Proposition 3.3.1. *The gauge function γ_C is a sublinear functional (as defined in Definition 2.3.1) when C is convex.*

Proof. Positive homogeneity is clear. For subadditivity, fix any $\alpha > \gamma_C(x)$ and $\beta > \gamma_C(y)$ so that $x \in \alpha C$ and $y \in \beta C$. Then $x + y \in (\alpha + \beta)C$ so that $\gamma_C(x + y) \leq \alpha + \beta$. Take an infimum over all such α, β to get the result. \square

Lemma 3.3.2 (Separating a point from a set). *If C is a nonempty convex set then for all $x_0 \notin C$, there exists a nonzero $w \in \mathbb{R}^n$ such that $\langle w, x \rangle \leq \langle w, x_0 \rangle$ for all $x \in C$.*

Proof. Suppose without loss of generality that $0 \in C$. Define $L = \text{span}(\{x_0\})$ and define a linear functional on L by $f(\alpha x_0) = \alpha$ for $\alpha \in \mathbb{R}$. Now, we see that $f(v) \leq \gamma_C(v)$ for all $v \in L$ since $\gamma_C(\alpha x_0) \geq \alpha = f(\alpha x_0)$ for $\alpha > 0$ and $f(\alpha x_0) \leq 0 \leq \gamma_C(\alpha x_0)$ for $\alpha \leq 0$.

By the Hahn-Banach theorem (Theorem 2.3.2), we can extend f to a linear functional \tilde{f} on all of \mathbb{R}^n satisfying $\tilde{f}(v) \leq \gamma_C(v)$ for all $v \in \mathbb{R}^n$. By the Riesz-Fréchet representation theorem (Theorem 2.2.2), there exists a nonzero vector $w \in \mathbb{R}^n$ such that $\tilde{f}(x) = \langle w, x \rangle$ for all $x \in \mathbb{R}^n$. For any $x \in C$, we have $\langle w, x \rangle = \tilde{f}(x) \leq \gamma_C(x) \leq 1 = \tilde{f}(x_0)$. \square

Note that the Hahn-Banach theorem wasn't strictly needed in \mathbb{R}^n , but the point is that the same proof works in arbitrary normed spaces. From this lemma, we derive the separating hyperplane theorem.

Theorem 3.3.3 (Separating hyperplane). *If A and B are nonempty disjoint convex sets, then there exists a nonzero $w \in \mathbb{R}^n$ and $\alpha \in \mathbb{R}$ such that $\sup_{x \in A} \langle w, x \rangle \leq \alpha \leq \inf_{x \in B} \langle w, x \rangle$.*

Proof. Let $C = A - B = \{x - y : x \in A, y \in B\}$. Then C is convex and contains 0, so by Lemma 3.3.2 there exists a nonzero w such that $\langle w, x \rangle \leq \langle w, 0 \rangle = 0$ for all $x \in C$. This implies that $\langle w, a - b \rangle \leq 0$ for all $a \in A$ and $b \in B$, so $\langle w, a \rangle \leq \langle w, b \rangle$ for all $a \in A$ and $b \in B$. \square

Sometimes, the separating hyperplane theorem is also called the *Hahn-Banach separation theorem*. Now, we are ready to show the supporting hyperplane theorem.

Theorem 3.3.4 (Supporting hyperplane). *If C is a nonempty convex set and x_0 is a boundary point of C , then there exists a nonzero $w \in \mathbb{R}^n$ such that $\langle w, x \rangle \leq \langle w, x_0 \rangle$ for all $x \in C$.*

Proof. Pick a sequence $(y_n)_{n=1}^\infty$ in C° converging to x_0 . By the separating hyperplane theorem (Theorem 3.3.3), there exists a nonzero continuous linear functional f_n such that $f_n(y) \leq f_n(y_n)$ for all $y \in C$ for all $n \in \mathbb{N}$. Renormalize the f_n to have norm 1 (preserving the separation); by the Banach-Alaoglu theorem (Theorem 2.7.6), there exists a subsequence $(f_{n_k})_{k=1}^\infty$ converging in weak-* to some continuous linear functional f with $\|f\| = 1$. For all $k \in \mathbb{N}$, we know that $f_{n_k}(x) \leq f_{n_k}(y_{n_k})$ for all $x \in C$, so $f_{n_k}(x - x_0) \leq f_{n_k}(y_{n_k} - x_0)$ for all $x \in C$. As $k \rightarrow \infty$, the left-hand side converges to $f(x - x_0)$ by weak-* convergence and the right-hand side converges to zero since $\|f_{n_k}\| = 1$ and $y_{n_k} \rightarrow x_0$. Hence, we deduce that $f(x) \leq f(x_0)$ for all $x \in C$, and we conclude by the Riesz-Fréchet representation theorem (Theorem 2.2.2). \square

We could have given a purely elementary proof here as well, but the proof above is more general and gives a simple example of the usefulness of the Banach-Alaoglu theorem. The previous theorems have an immediate corollary.

Corollary 3.3.4.1. *A closed and convex set C is the intersection of all closed half-spaces containing it.*

In fact, we can give a useful property of convex l.s.c. functions using the supporting hyperplane theorem.

Theorem 3.3.5. *If f is convex and l.s.c. then f is the supremum of affine functions lying below f .*

Proof. Since f is l.s.c., we know that $\text{epi}(f)$ is closed and convex. By Corollary 3.3.4.1, $\text{epi}(f)$ is the intersection of all closed half-spaces containing it, each represented as $\{x \in \mathbb{R}^n : \langle a, x \rangle + b\mu \leq c\}$ for some $a \in \mathbb{R}^n$ and $b, c \in \mathbb{R}$. Not all of these half-spaces can have $b = 0$, since otherwise f is trivial and the theorem follows immediately. We can't have $b > 0$ at all, since the epigraph of f extends upwards infinitely. When $b < 0$, we can normalize to $b = -1$ so that $\langle a, x \rangle - c \leq \mu$; this defines an affine function $\langle a, x \rangle - c$ lying below f . The constraints induced by half-spaces with $b < 0$ are the only nontrivial ones, so the result follows. \square

3.4 Subgradients and subdifferentials

The idea of subdifferentials is to use a supporting hyperplane to generalize the notion of a derivative for non-differentiable functions.

Definition 3.4.1 (Subgradient). A vector x^* is a *subgradient* of f at x if $f(z) \geq f(x) + \langle x^*, z - x \rangle$ for all $z \in \mathbb{R}^n$.

Definition 3.4.2 (Subdifferential). The subdifferential of f at x is the set of all subgradients at x , denoted $\partial f(x)$.

Although one can theoretically define a subdifferential for any function, we almost always restrict ourselves to convex functions so that the subdifferential enjoys several nice properties.

Proposition 3.4.1. *If f is convex then $\partial f(x)$ is a closed convex set for all $x \in \text{dom}(f)$. Also, the subdifferential is nonempty for all $x \in \text{int}(\text{dom}(f))$, and for all $x \in \text{dom}(f)$ if f is l.s.c..*

Proof. It follows from the supporting hyperplane theorem (Theorem 3.3.4) that the subdifferential is nonempty under the given conditions. Also, closedness and convexity are immediate from the definition. \square

Example 3.4.1 (Subdifferential of the norm). The subdifferential of the Euclidean norm $\|\cdot\|_2$ at x is $\{x/\|x\|_2\}$ at every $x \neq 0$ and $B(0, 1)$ at $x = 0$, by Cauchy-Schwarz.

Note that if the gradient exists at a point x , then the subdifferential at that point is a singleton containing only the gradient; this can be seen by using Jensen's inequality when taking the directional derivative at x in the direction of $z - x$.

3.5 The Legendre-Fenchel transform

The Legendre-Fenchel transform gives a rich duality theory for convex functions.

Definition 3.5.1 (Legendre-Fenchel transform). The *Legendre-Fenchel transform* of f is the function f^* defined by

$$f^*(y) = \sup_{x \in \mathbb{R}^n} \{\langle y, x \rangle - f(x)\}.$$

The function f^* is also called the *convex conjugate* of f .

In a general normed space X , the Legendre-Fenchel transform can be thought of as a transformation from X^* to \mathbb{R} , defined by

$$f^*(g) = \sup_{x \in X} \{g(x) - f(x)\}.$$

The Legendre-Fenchel transform represents the maximal difference between a linear function $\langle y, x \rangle$ and $f(x)$. If f is convex, the Legendre-Fenchel transform is immediately convex as it is the supremum of affine functions. Intuitively, if $f(x)$ is the cost to buy a portfolio x of products, then $f^*(g)$ can be thought of as the maximum profit that you can make by setting linear prices $g(x)$ (under any portfolio). The following is a simple result, but is fundamental in the theory.

Theorem 3.5.1 (Fenchel-Young inequality). *If f is convex then $f(x) + f^*(y) \geq \langle y, x \rangle$ for all $x, y \in \mathbb{R}^n$. Equality holds if and only if $y \in \partial f(x)$.*

Proof. The inequality is obvious from the previous interpretation of the Legendre-Fenchel transform, and the equality condition follows from the definition of the subdifferential. \square

Theorem 3.5.2 (Fenchel-Moreau). *If f is l.s.c. and convex, then $f^{**} = f$.*

Proof. The inequality $f \geq f^{**}$ is immediate from definitions. Now, note that for every affine function $a \leq f$, we have $a^{**} = a \leq f^{**}$; in particular, the affine functions lying below f are the same ones below f^{**} . Taking a supremum over all affine functions $a \leq f$ and applying [Theorem 3.3.5](#), it follows that $f^{**} \geq f$. \square

Further, we have the following useful strong duality theorem, which gives a dual problem for the minimization of a function subject to constraints.

Theorem 3.5.3 (Fenchel-Rockafellar). *Suppose f_1 and f_2 are convex functions on a normed space X taking values in $[0, \infty]$. Also, assume that there exists $x_0 \in X$ such that $f_1(x_0) + f_2(x_0) < \infty$ and f_1 is continuous at x_0 . Then we have the minimax principle*

$$\inf_{x \in X} \{f_1(x) + f_2(x)\} = \sup_{g \in X^*} \{-f_1^*(-g) - f_2^*(g)\} = \sup_{g \in X^*} \inf_{x, y \in X} \{f_1(x) + f_2(y) + g(x - y)\}.$$

Intuitively, suppose x denotes a portfolio of different products. Suppose $f_1(x)$ is the cost of producing x , and $f_2(x)$ is a convex penalty describing whether producing x is feasible or not. Then, the primal problem is to minimize the cost of producing x subject to feasibility constraints. On the other hand, suppose g denotes a vector of market prices for the products. Then, $-f_1^*(-g) = \inf_{x \in X} \{g(x) + f_1(x)\}$ is the minimum net cost under prices g and $-f_2^*(g) = \inf_{x \in X} \{f_2(x) - g(x)\}$ is the minimum gross profit under prices g and feasibility constraints. The dual problem is to price in a way that maximizes the worst-case profit for the producer. By the Fenchel-Rockafellar duality, minimizing the cost of producing x subject to feasibility constraints is equivalent to pricing in a way that maximizes the worst-case profit for the producer.

Proof. By choosing $x = y$, we see that the right hand side is less than or equal to the left hand side. The reverse inequality will now follow from the separating hyperplane theorem. Define $m = \inf_{x \in X} \{f_1(x) + f_2(x)\}$, which is finite because we assumed $f_1(x_0) + f_2(x_0) < \infty$. Let $C = \text{epi}(f_1)$ and $C' = \{(x, \mu) \in X \times \mathbb{R} : \mu < m - f_2(x)\}$; these are convex sets and disjoint by definition of m . Furthermore, C has nonempty interior since $(x_0, f_1(x_0) + 1) \in \text{int}(C)$ by continuity of f_1 at x_0 , so C has nonempty interior. Note that in an infinite-dimensional space, we need one of the sets to have nonempty interior in order to apply the separating hyperplane theorem.

By the separating hyperplane theorem (Theorem 3.3.3), there exists a nonzero $g \in X^*$ such that $g(x) + \alpha\lambda \geq g(y) + \alpha\mu$ for all $(x, \lambda) \in C$ and $(y, \mu) \in C'$. It is easy to see that $\alpha > 0$, so assume w.l.o.g. that $\alpha = 1$. Then, we can rearrange

$$g(x) + f_1(x) \geq g(x) + \lambda \geq g(y) + \mu \geq g(y) + m - f_2(y)$$

to obtain the reverse inequality. □

3.6 Cyclical monotonicity

Cyclical monotonicity is a characterizing property of the subdifferential of convex functions, and is used in optimal transport theory to prove Brenier's theorem.

Definition 3.6.1 (Cyclical monotonicity). A subset $\Gamma \subseteq \mathbb{R}^n \times \mathbb{R}^n$ is *cyclically monotone* if for all finite collections $(x_1, y_1), \dots, (x_m, y_m) \in \Gamma$, we have

$$\sum_{i=1}^m \langle y_i, x_{i+1} - x_i \rangle \leq 0,$$

under the convention $x_{m+1} = x_1$.

In some sense, cyclical monotonicity is the generalization of monotonicity to \mathbb{R}^n , since monotone functions on \mathbb{R} are the derivatives of convex functions. The motivation for this definition is Rockafellar's theorem.

Theorem 3.6.1 (Rockafellar). A nonempty $\Gamma \subseteq \mathbb{R}^n \times \mathbb{R}^n$ is cyclically monotone if and only if it is included in the graph of $\partial\varphi$ for some proper l.s.c. convex function $\varphi : \mathbb{R}^n \rightarrow (-\infty, \infty]$.

Proof. For the reverse direction, suppose that Γ is included in the graph of $\partial\varphi$ for a l.s.c. convex function φ . If $(x_1, y_1), \dots, (x_m, y_m) \in \Gamma$, then for $1 \leq i \leq m$, we have

$$\varphi(x_{i+1}) \geq \varphi(x_i) + \langle y_i, x_{i+1} - x_i \rangle$$

by definition of the subdifferential. Summing over i shows that Γ is cyclically monotone. For the forward direction, suppose that Γ is cyclically monotone and nonempty. Pick $(x_0, y_0) \in \Gamma$ and define the function

$$\varphi(x) = \sup\{\langle y_m, x - x_m \rangle + \langle y_{m-1}, x_m - x_{m-1} \rangle + \dots + \langle y_0, x_1 - x_0 \rangle : (x_1, y_1), \dots, (x_m, y_m) \in \Gamma\}.$$

Since φ is the supremum of affine functions it is l.s.c. and convex. Also, φ is proper since $\varphi(x_0) \leq 0$ by cyclical monotonicity. It is now easy to verify by the definition of φ that Γ is contained in the graph of $\partial\varphi$. □

Chapter 4

Optimal transport

In this chapter, we write about optimal transport theory, primarily using material from *Topics in Optimal Transportation* by Cédric Villani and *Statistical Optimal Transport* by Sinho Chewi, Jonathan Niles-Weed, and Philippe Rigollet. Some material may also come from *Optimal Transport for Applied Mathematicians* by Filippo Santambrogio or *Gradient Flows in Metric Spaces and in the Space of Probability Measures* by Luigi Ambrosio, Nicola Gigli, and Giuseppe Savaré.

4.1 The Monge-Kantorovich problem

The Monge problem is that of finding a deterministic transport map between measures. We let $\mathcal{P}(\mathcal{X})$ denote the set of probability measures over \mathcal{X} and $T_{\#}$ denote the pushforward operator. In everything that follows, we will assume \mathcal{X} and \mathcal{Y} are Polish spaces (complete and separable metric spaces). The main reason to assume a Polish space setting is so that we can use Prokhorov's theorem ([Theorem 1.1.4](#)) to extract weakly convergent subsequences from tight sets, as well as for disintegration of measure ([Theorem 1.4.1](#)).

Definition 4.1.1 (Monge problem). If $\mu \in \mathcal{P}(X)$ and $\nu \in \mathcal{P}(Y)$ and $c : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$ is a cost function, the *Monge problem* is

$$\inf_{T_{\#}\mu=\nu} \int c(x, T(x)) d\mu(x).$$

The Kantorovich problem is the problem of finding a stochastic transport plan; one motivation for this is that there is literally no deterministic transport map from δ_0 to $\frac{1}{2}\delta_{-1} + \frac{1}{2}\delta_1$ and another is that the Monge problem is highly nonconvex and hard to solve. Let $\Pi(\mu, \nu)$ denote the set of couplings between μ and ν (joint distributions with marginals μ and ν).

Definition 4.1.2 (Kantorovich problem). If $\mu \in \mathcal{P}(\mathcal{X})$ and $\nu \in \mathcal{P}(\mathcal{Y})$ and $c : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$ is a cost function, the *Kantorovich problem* is

$$\inf_{\pi \in \Pi(\mu, \nu)} \int c(x, y) d\pi(x, y).$$

Equivalently, if X and Y are random variables, the Kantorovich problem is to minimize $\mathbb{E}_\pi[c(X, Y)]$ over all joint distributions π of X and Y . Note that the Kantorovich problem is a linear program since the objective and constraints are linear. Here are a few simple facts about the set of couplings.

Proposition 4.1.1. *The set of couplings $\Pi(\mu, \nu)$ is a nonempty and convex. Also, if $(\pi_n)_{n=1}^\infty$ is a sequence in $\Pi(\mu, \nu)$, then there exists a weakly convergent subsequence with a limit also in $\Pi(\mu, \nu)$.*

Proof. Nonemptiness follows from $\mu \times \nu \in \Pi(\mu, \nu)$ and convexity is obvious. By Prokhorov's theorem (Theorem 1.1.4), it suffices to show that $\Pi(\mu, \nu)$ is tight and closed. Tightness follows because μ and ν are simultaneously concentrated on a compact set K , and $\pi((K \times K)^c) \leq \mu(K^c) + \nu(K^c)$. Note that π is a coupling if and only if $\int f(x) d\pi(x, y) = \int f(x) d\mu(x)$ for all $f \in C_b(X)$ and $\int f(y) d\pi(x, y) = \int f(y) d\nu(y)$ for all $f \in C_b(Y)$, by the Riesz-Markov representation theorem. Therefore, $\Pi(\mu, \nu)$ is closed by the portmanteau lemma (Theorem 1.1.2). \square

Proposition 4.1.2. *If the cost function c is l.s.c., then the Kantorovich problem has a solution.*

Proof. Let $\pi_n \in \Pi(\mu, \nu)$ be such that

$$\int c(x, y) d\pi_n(x, y) \rightarrow \inf_{\pi \in \Pi(\mu, \nu)} \int c(x, y) d\pi(x, y).$$

and extract a convergent subsequence $\pi_{n_k} \rightarrow \pi \in \Pi(\mu, \nu)$ by Proposition 4.1.1. By (4) in the portmanteau lemma (Theorem 1.1.2) and because c is l.s.c., we have

$$\int c(x, y) d\pi(x, y) \leq \liminf_{k \rightarrow \infty} \int c(x, y) d\pi_{n_k}(x, y) = \inf_{\pi \in \Pi(\mu, \nu)} \int c(x, y) d\pi(x, y).$$

In particular, π solves the Kantorovich problem. \square

Note that this statement is a special case of the *extreme value theorem*, which states that l.s.c. functions must attain their infimum on compact sets. Even though the Kantorovich problem often has a solution, it is not always unique; for example, consider the case $\mu = \frac{1}{2}\delta_{(-1,0)} + \frac{1}{2}\delta_{(1,0)}$ and $\nu = \frac{1}{2}\delta_{(0,-1)} + \frac{1}{2}\delta_{(0,1)}$ with the quadratic cost function.

4.2 Transport maps between empirical averages

In this section, we show that we can fully characterize optimal transport maps between measures associated to empirical averages. We start by proving the Krein-Milman theorem, which states that points in a compact convex set in a Banach space can be written as limits of convex combinations of extreme points.

Theorem 4.2.1 (Krein-Milman). *Let K be a nonempty, compact, and convex subset of a Banach space and let $\mathcal{E}(K)$ denote the set of extremal points of K (points that cannot be written as the convex combination of any two other points). Then for all $x \in K$ there exists a probability measure ρ_x on $\mathcal{E}(K)$ such that $x = \int_{\mathcal{E}(K)} y d\rho_x(y)$.*

Proof. We would first like to show that $C := \overline{\text{conv}(\mathcal{E}(K))} = K$. Suppose not for a contradiction; then there exists $x \in K \setminus C$. By the separating hyperplane theorem (Theorem 3.3.3), there exists a continuous linear functional f and a constant c such that $f(x) > \max_{y \in C} f(y) \geq \max_{y \in \mathcal{E}(K)} f(y)$ (using compactness of C and closedness of $\mathcal{E}(K)$). Let $M = f^{-1}(\max_{x \in C} f(x))$, which is a nonempty, closed, and compact; this fact is sometimes called the *Krein-Milman lemma* and is easy to show. So M has an extreme point z (which must also be in $\mathcal{E}(K)$) and we have the inequality

$$f(z) \geq f(x) > \max_{y \in C} f(y) \geq \max_{y \in \mathcal{E}(K)} f(y) \geq f(z),$$

giving a contradiction. Then for each $x \in K$, we can extract a sequence of discrete probability measures $\rho_x^{(n)}$ such that $\int y d\rho_x^{(n)} \rightarrow x$ as $n \rightarrow \infty$. By the Banach-Alaoglu theorem (Theorem 2.7.6), we can extract a weakly convergent subsequence $\rho_x^{(n_k)}$ with limit ρ_x , and the result follows since $\int y d\rho_x = x$. \square

Using the Krein-Milman theorem, we can show Choquet's theorem, which states that continuous linear functionals on a compact convex set in a Banach space achieve their minimum at an extreme point.

Theorem 4.2.2 (Choquet). *Let K be a nonempty, compact, and convex subset of a Banach space and let $f : K \rightarrow \mathbb{R}$ be the restriction of a continuous linear functional. Then the minimum of f over K is achieved at an extreme point of K .*

Proof. Note that f has a minimizer $x \in K$ since K is compact. By the Krein-Milman theorem (Theorem 4.2.1), there exists a probability measure ρ_x on $\mathcal{E}(K)$ such that $x = \int_{\mathcal{E}(K)} y d\rho_x(y)$. Then, by continuity and linearity of f , we have

$$f(x) = f\left(\int_{\mathcal{E}(K)} y d\rho_x(y)\right) = \int_{\mathcal{E}(K)} f(y) d\rho_x(y).$$

It is easy to formalize this exchange by approximating $y \mapsto y$ uniformly on K by simple functions and applying the dominated convergence theorem. Clearly, this means that ρ_x cannot give positive mass to

any point $y \in \mathcal{E}(K)$ such that $f(y) > f(x)$, so there must exist $y \in \mathcal{E}(K)$ such that $f(y) = f(x)$, which is the result we wanted. \square

Using Choquet's theorem, we can show that the solution to the Kantorovich problem between discrete measures must be a deterministic permutation; in particular, the Kantorovich problem is equivalent to the Monge problem in this case.

Theorem 4.2.3 (Birkhoff). *Let \mathcal{B}_n denote the set of bistochastic $n \times n$ matrices (matrices with nonnegative entries such that each row and column sums to 1). Then the set of extremal points of \mathcal{B}_n is exactly the set of permutation matrices.*

Proof. It's clear that all permutation matrices are extremal points, so we show the converse. It suffices to show that whenever M is an extremal point of \mathcal{B}_n , all entries of M are either 0 or 1. Suppose that $M \in \mathcal{B}_n$ with $M_{ij} \in (0, 1)$ for some $i, j \in [n]$. Then make a list of indices in the matrix, first finding another entry in the same column in $(0, 1)$ and then an entry in the same row from that new point in $(0, 1)$, and so on. Continue this process until we either end up at the same row or the same column of M . The list will look something like $((i, j), \dots, (i_f, j_f))$; if we ended up at the same column, then we delete (i, j) from the beginning of our list.

The list is now of even length and only contains entries in $(0, 1)$, so pick $\epsilon > 0$ smaller than the gap $\min\{\min\{M_{ij}\}_{ij}, \min\{1 - M_{ij}\}_{ij}\}$. Then we can create a matrix $M^{(1)}$ by bumping the entry of M at the first index in the list up by ϵ , the entry at the second index down by ϵ , and so on. Similarly, we create a matrix $M^{(2)}$ by bumping the entry of M at the first index in the list down by ϵ , the entry at the second index up by ϵ , and so on. Then $M = \frac{1}{2}(M^{(1)} + M^{(2)})$, and because $M^{(1)}, M^{(2)} \in \mathcal{B}_n$, M cannot be an extremal point of \mathcal{B}_n . \square

Proposition 4.2.4. *The solution to the Kantorovich problem between discrete spaces \mathcal{X} and \mathcal{Y} where $\mu = \frac{1}{n} \sum_{i=1}^n \delta_{x_i}$ and $\nu = \frac{1}{n} \sum_{i=1}^n \delta_{y_i}$ is a deterministic transport map $T : \mathcal{X} \rightarrow \mathcal{Y}$ such that $T(x_i) = y_{\sigma(i)}$ for some permutation $\sigma \in S_n$.*

Proof. Let \mathcal{B}_n denote the set of bistochastic $n \times n$ matrices, so that the Kantorovich problem is

$$\inf_{\pi \in \mathcal{B}_n} \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^n \pi_{ij} c(x_i, y_j).$$

By Choquet's theorem (Theorem 4.2.2), the minimum is achieved at an extremal point of \mathcal{B}_n , which by Birkhoff's theorem (Theorem 4.2.3) is a permutation matrix. \square

4.3 Wasserstein distances

Suppose (\mathcal{X}, d) is a Polish space (a complete and separable metric space). Let $\mathcal{P}_p(\mathcal{X})$ denote the set of probability measures over \mathcal{X} with finite p th moment, meaning that

$$\int d(x_0, x)^p d\mu(x) < \infty$$

for all $x_0 \in \mathcal{X}$. Let $\mathcal{T}_p(\mu, \nu)$ denote the optimal transportation cost between measures μ and ν if $c(x, y) = d(x, y)^p$; we can now define the Wasserstein distances on $\mathcal{P}_p(\mathcal{X})$ as follows.

Definition 4.3.1 (Wasserstein distances). For $p \geq 0$, the p -Wasserstein distance between measures $\mu, \nu \in \mathcal{P}_p$ is defined as $W_p(\mu, \nu) = \mathcal{T}_p(\mu, \nu)^{\min\{1, 1/p\}}$.

Theorem 4.3.1. *The p -Wasserstein distance is actually a metric on $\mathcal{P}_p(\mathcal{X})$.*

Proof. We may assume $p \geq 1$, since otherwise d^p is topologically equivalent to d (although it may not be a metric) and is subadditive. It's obvious that W_p is finite, symmetric, and nonnegative. If $\mu = \nu$, then the measure $d\mu(x) \times d\delta_x(y)$ is a coupling, showing that $W_p(\mu, \nu) = 0$. On the other hand, if $W_p(\mu, \nu) = 0$, there exists a coupling π such that $\mathbb{E}_{(X,Y) \sim \pi}[d(X,Y)^p] = 0$ by [Proposition 4.1.2](#). But now this forces $X = Y$ almost surely so $\mu = \nu$.

Finally, we prove the triangle inequality. First, we need the following lemma, which allows us to put together measures with a common marginal; essentially, this is a technical lemma that will allow us to put together optimal transport plans to make a (possibly suboptimal) transport plan.

Lemma 4.3.2 (Gluing). *Let μ_1, μ_2, μ_3 be measures on Polish spaces $\mathcal{X}_1, \mathcal{X}_2, \mathcal{X}_3$ with $\pi_{12} \in \Pi(\mu_1, \mu_2)$ and $\pi_{23} \in \Pi(\mu_2, \mu_3)$. Then there exists a measure $\pi \in \mathcal{P}(\mathcal{X}_1 \times \mathcal{X}_2 \times \mathcal{X}_3)$ such that π has marginals π_{12} and π_{23} in its restriction to the first two or last two coordinates respectively.*

Proof. Applying the disintegration theorem ([Theorem 1.4.1](#)) along the second coordinate of π_{12} and the first coordinate of π_{23} , we have

$$\pi_{12} = \int_{\mathcal{X}_2} (\pi_{12})_y \otimes \delta_y d\mu_2(y)$$

and

$$\pi_{23} = \int_{\mathcal{X}_2} \delta_y \otimes (\pi_{23})_y d\mu_2(y),$$

such that $y \mapsto (\pi_{12})_y$ and $y \mapsto (\pi_{23})_y$ are measurable. Then we can define π by

$$\pi = \int_{\mathcal{X}_2} (\pi_{12})_y \otimes \delta_y \otimes (\pi_{23})_y d\mu_2(y),$$

which satisfies the properties of the gluing lemma. □

Now, we are ready to complete the proof of the triangle inequality. Consider $\mu_1, \mu_2, \mu_3 \in \mathcal{P}_p(\mathcal{X})$ and let $\pi_{12} \in \Pi(\mu_1, \mu_2)$ and $\pi_{23} \in \Pi(\mu_2, \mu_3)$ be optimal transport plans. Glue these together using the gluing lemma (Lemma 4.3.2) to get $\pi \in \mathcal{P}(\mathcal{X}_1 \times \mathcal{X}_2 \times \mathcal{X}_3)$ and let $(X, Y, Z) \sim \pi$. Then, we have by Minkowski's inequality that

$$\begin{aligned} W_p(\mu_1, \mu_3) &\leq \|d(X, Z)\|_p \\ &\leq \|d(X, Y) + d(Y, Z)\|_p \\ &\leq \|d(X, Y)\|_p + \|d(Y, Z)\|_p \\ &= W_p(\mu_1, \mu_2) + W_p(\mu_2, \mu_3). \end{aligned} \quad \square$$

Example 4.3.1 (Point masses). Note that $W_p(\delta_x, \delta_y) = d(x, y)$, so (\mathcal{X}, d) is isometrically embedded in $(\mathcal{P}_p(\mathcal{X}), W_p)$ by the map $x \mapsto \delta_x$.

By the ordering of the L^p norms in a probability space (Proposition 1.1.6), we have $W_p(\mu, \nu) \leq W_q(\mu, \nu)$ for $1 \leq p \leq q$. Note that the Wasserstein distances metrize weak convergence.

Proposition 4.3.3. For $p \geq 1$, $W_p(\mu_n, \mu) \rightarrow 0$ if and only if $\mu_n \xrightarrow{d} \mu$.

Proof. For the forward direction, by ordering of the W_p distances, we have $W_1(\mu, \nu) \rightarrow 0$. \square

Now, we give a few techniques to bound the Wasserstein distance. First, one can bound the Wasserstein distance from below by coming up with an inequality that holds uniformly over couplings, pulling out an $X - Y$ using Hölder's inequality, and then choosing the optimal coupling to get a tight bound.

Proposition 4.3.4. If $X \sim \mu$ and $Y \sim \nu$ are sub-exponential in the sense that $\mathbb{E}[e^{|X|}] \leq 2$ and $\mathbb{E}[e^{|Y|}] \leq 2$, then for any $p > 1$ there exists a universal constant $c_p > 0$ so that for all $k \in \mathbb{N}$, we have

$$\mathbb{E}[||X|^k - |Y|^k|] \leq (c_p k)^k W_p(\mu, \nu).$$

Hence, a bound on the Wasserstein distance between μ and ν simultaneously implies that all of their moments are close, as long as the tails decay sufficiently fast.

Proof. Note that if $f(x) = |x|^k$, then $k|x|^{k-1} \in \partial f(x)$ for all $x \in \mathbb{R}$ and $k \in \mathbb{N}$. In particular, we know that:

$$|X|^k - |Y|^k \leq |X - Y| \cdot k |X|^{k-1}.$$

A symmetric bound holds if we swap the roles of X and Y , so we can take the expectation on both sides to get that for any coupling of X and Y ,

$$\mathbb{E}[||X|^k - |Y|^k|] \leq k \mathbb{E}[|X - Y| (|X| \vee |Y|)^{k-1}].$$

Applying Hölder's inequality with the conjugate exponents p and $q = p/(p-1)$, we have

$$k \mathbb{E}[|X - Y| (|X| \vee |Y|)^{k-1}] \leq k \|X - Y\|_p \|(|X| \vee |Y|)^{k-1}\|_q.$$

In particular, this inequality holds for the optimal coupling, so we can replace the right-hand side with

$$k W_p(\mu, \nu) \|(|X| \vee |Y|)^{k-1}\|_q.$$

But now it is easy to show that $|X| \vee |Y|$ is sub-exponential and there exists a constant c_r such that $\|Z\|_r \leq c_r r$ for sub-exponential Z and $r \geq 1$. Therefore, we have that

$$\|(|X| \vee |Y|)^{k-1}\|_q \leq (c_p(k-1))^{k-1} \leq (c_p k)^k,$$

which concludes the proof. \square

Bounding Wasserstein distances from above is actually easier, since we only need to come up with a (potentially suboptimal) coupling and find the resulting cost.

Proposition 4.3.5. *Suppose $\mu, \nu \in \mathcal{P}_p(\mathcal{X})$ and fix $x_0 \in \mathcal{X}$. Then, for all $p \geq 0$, we have*

$$W_p(\mu, \nu)^{(p \vee 1)} \leq (2^{p-1} \vee 1) \int d(x_0, x)^p d|\mu - \nu|(x) = (2^{p-1} \vee 1) \|d(x_0, \cdot)^p(\mu - \nu)\|_{\text{TV}}.$$

This proposition says that the Wasserstein distance is loosely bounded above by a weighted total variation distance.

Proof. The proof is to consider the coupling where we keep all the shared mass between μ and ν fixed in place and then distribute the rest uniformly. More formally, let

$$\pi = (\text{Id} \times \text{Id})_{\#}(\mu \wedge \nu) + \frac{1}{\alpha}(\mu - \nu)_+ \times (\mu - \nu)_-,$$

where $\alpha = (\mu - \nu)_+(\mathcal{X}) = (\mu - \nu)_-(\mathcal{X})$ is the total excess mass. Hence, we have the bound

$$W_p(\mu, \nu)^{(p \vee 1)} \leq \int d(x, y)^p d\pi(x, y) = \frac{1}{\alpha} \int d(x, y)^p d(\mu - \nu)_+(x) d(\mu - \nu)_-(y).$$

When $p \geq 1$, we can use Jensen's inequality on the convex function $x \mapsto |x|^p$ to get

$$\begin{aligned} d(x, y)^p &\leq (d(x_0, x) + d(x_0, y))^p \\ &= \left(\frac{1}{2}(2d(x_0, x)) + \frac{1}{2}(2d(x_0, y)) \right)^p \\ &\leq \frac{1}{2}(2d(x_0, x))^p + \frac{1}{2}(2d(x_0, y))^p \\ &= 2^{p-1}(d(x_0, x)^p + d(x_0, y)^p). \end{aligned}$$

If $0 \leq p < 1$, we immediately have $d(x, y)^p \leq d(x_0, x)^p + d(x_0, y)^p$ by subadditivity of $x \mapsto |x|^p$. Hence, we deduce

$$\begin{aligned}
& \frac{1}{\alpha} \int d(x, y)^p d(\mu - \nu)_+(x) d(\mu - \nu)_-(y) \\
& \leq \frac{(2^{p-1} \vee 1)}{\alpha} \left(\int d(x_0, x)^p d(\mu - \nu)_+(x) d(\mu - \nu)_-(y) + \int d(x_0, y)^p d(\mu - \nu)_+(x) d(\mu - \nu)_-(y) \right) \\
& = (2^{p-1} \vee 1) \left(\int d(x_0, x)^p d(\mu - \nu)_+(x) + \int d(x_0, y)^p d(\mu - \nu)_-(y) \right) \\
& = (2^{p-1} \vee 1) \int d(x_0, x)^p d|\mu - \nu|(x). \quad \square
\end{aligned}$$

4.4 The Kantorovich duality

In this section, we develop the dual formulation of the Kantorovich problem, and show that strong duality holds under very general conditions.

4.4.1 Lower semi-continuous cost functions

Theorem 4.4.1 (Kantorovich duality). *Suppose \mathcal{X} and \mathcal{Y} are Polish spaces with $\mu \in \mathcal{P}(\mathcal{X})$ and $\nu \in \mathcal{P}(\mathcal{Y})$, and suppose c is an l.s.c. cost function. We define the set of dual variables by*

$$\Phi_c = \{(\phi, \psi) \in L^1(\mu) \times L^1(\nu) : \phi(x) + \psi(y) \leq c(x, y) \text{ for } \mu\text{-a.e. } x \in \mathcal{X} \text{ and } \nu\text{-a.e. } y \in \mathcal{Y}\}.$$

Then the Kantorovich problem enjoys the following strong duality:

$$\inf_{\pi \in \Pi(\mu, \nu)} \int c(x, y) d\pi(x, y) = \sup_{(\phi, \psi) \in \Phi_c} \left\{ \int \phi(x) d\mu(x) + \int \psi(y) d\nu(y) \right\}.$$

Intuitively, suppose that a shipper offers to pick up goods from location x at cost $\phi(x)$ and deliver them to location y at cost $\psi(y)$ in a way which is beneficial to you; namely, so that your total cost through this service ($\phi(x) + \psi(y)$) is at most your cost ($c(x, y)$) of carrying the goods directly from x to y . Then, strong duality says that a clever shipper can price their service such that you pay them almost as much as you would have paid anyways.

Proof. We first show the result in the case where \mathcal{X} and \mathcal{Y} are compact, and c is continuous. Define the functional f_1 on $C_b(\mathcal{X}, \mathcal{Y})$ by

$$f_1(u) = \begin{cases} 0 & u(x, y) \geq -c(x, y) \\ \infty & \text{otherwise.} \end{cases}$$

Similarly, we define the functional f_2 on $C_b(\mathcal{X}, \mathcal{Y})$ by

$$f_2(u) = \begin{cases} \int \phi(x) d\mu(x) + \int \psi(y) d\nu(y) & u(x, y) = \phi(x) + \psi(y) \\ \infty & \text{otherwise.} \end{cases}$$

Let $\mathcal{M}(\mathcal{X}, \mathcal{Y})$ denote the set of Radon measures on $\mathcal{X} \times \mathcal{Y}$, which is the dual of $C_b(\mathcal{X}, \mathcal{Y})$ by the Riesz-Markov representation theorem; similarly, let $\mathcal{M}_+(\mathcal{X}, \mathcal{Y})$ denote the set of nonnegative Radon measures on $\mathcal{X} \times \mathcal{Y}$. Then, f_1 and f_2 are both convex with $f_1(1) + f_2(1) < \infty$ and f_1 is continuous at 1. By the Fenchel-Rockafellar duality ([Theorem 3.5.3](#)), we have

$$\inf_{u \in C_b(\mathcal{X}, \mathcal{Y})} \{f_1(u) + f_2(u)\} = \sup_{\pi \in \mathcal{M}(\mathcal{X}, \mathcal{Y})} \{-f_1^*(-\pi) - f_2^*(\pi)\}.$$

The left-hand side is

$$\begin{aligned} & \inf \left\{ \int \phi(x) d\mu(x) + \int \psi(y) d\nu(y) : \phi(x) + \psi(y) \geq -c(x, y) \right\} \\ &= - \sup_{(\phi, \psi) \in \Phi_c} \left\{ \int \phi(x) d\mu(x) + \int \psi(y) d\nu(y) \right\}. \end{aligned}$$

Next, we compute the Legendre-Fenchel transform of f_1 :

$$f_1^*(-\pi) = \sup_{u \in C_b(\mathcal{X}, \mathcal{Y})} \left\{ \int u(x, y) d\pi(x, y) : u(x, y) \leq c(x, y) \right\} = \begin{cases} \int c(x, y) d\pi(x, y) & \pi \in \mathcal{M}_+(\mathcal{X}, \mathcal{Y}) \\ \infty & \text{otherwise.} \end{cases}$$

Similarly, the Legendre-Fenchel transform of f_2 is

$$f_2^*(\pi) = \begin{cases} 0 & \pi \in \Pi(\mu, \nu) \\ \infty & \text{otherwise.} \end{cases}$$

Since we can approximate functions in $L^1(\mu)$ and $L^1(\nu)$ by continuous functions, the proof is complete in the case when \mathcal{X} and \mathcal{Y} are compact and c is continuous. The rest of the proof is technical, since we have to carefully relax the assumptions of compactness and continuity.

First, we relax the assumption of compactness but keep the assumption that c is uniformly continuous and bounded. Suppose π_* is an optimal coupling of μ and ν with respect to c , which exists by [Proposition 4.1.2](#). By tightness of π_* in the Polish space $\mathcal{X} \times \mathcal{Y}$ due to Ulam's lemma ([Theorem 1.1.5](#)), there is a compact set $\mathcal{X}_0 \times \mathcal{Y}_0$ such that $\pi_*(\mathcal{X}_0 \times \mathcal{Y}_0) = 1 - \delta$. Pick $\tilde{\pi}_0$ which is optimal in $\mathcal{X}_0 \times \mathcal{Y}_0$ and construct $\tilde{\pi} = \pi_*(\mathcal{X}_0 \times \mathcal{Y}_0) \tilde{\pi}_0 + \mathbf{1}_{(\mathcal{X}_0 \times \mathcal{Y}_0)^c} \tilde{\pi}_* \in \Pi(\mu, \nu)$, which is close to optimal. Also, by our previous result we may pick $\tilde{\phi}_0$ and $\tilde{\psi}_0$ which are dual-optimal for $\tilde{\pi}_0$. The goal is to upgrade these to functions in Φ_c which are close to optimal for the primal problem.

Note that there exists $(x_0, y_0) \in \mathcal{X} \times \mathcal{Y}$ such that $\tilde{\phi}_0(x_0) + \tilde{\psi}_0(y_0) \geq -1$ since the choice $\phi = \psi = 0$ is always feasible. Then, since replacing $(\tilde{\phi}_0, \tilde{\psi}_0) \mapsto (\tilde{\phi}_0 + \epsilon, \tilde{\psi}_0 - \epsilon)$ for $\epsilon \in \mathbb{R}$ maintains feasibility, we may assume $\tilde{\phi}_0(x_0) \geq -1/2$ and $\tilde{\psi}_0(y_0) \geq -1/2$. So for all $(x, y) \in \mathcal{X} \times \mathcal{Y}$, we have

$$\begin{aligned} \tilde{\phi}_0(x) &\leq c(x, y_0) - \tilde{\psi}_0(y_0) \leq c(x, y_0) + 1/2, \\ \tilde{\psi}_0(y) &\leq c(x_0, y) - \tilde{\phi}_0(x_0) \leq c(x_0, y) + 1/2. \end{aligned}$$

We improve the admissible pair by using *Rüschendorf's trick*, and define

$$\bar{\phi}_0(x) = \inf_{y \in \mathcal{Y}_0} \{c(x, y) - \tilde{\psi}_0(y)\}.$$

Essentially, this is the best possible choice of ϕ if we hold $\tilde{\psi}_0$ fixed. Since $\bar{\phi}_0 \geq \tilde{\phi}_0$, we can bound it above and below by the cost function:

$$\inf_{y \in \mathcal{Y}_0} \{c(x, y) - c(x_0, y)\} - 1/2 \leq \bar{\phi}_0(x) \leq c(x, y_0) + 1/2.$$

Then, if we define $\bar{\psi}_0(y) = \inf_{x \in \mathcal{X}} \{c(x, y) - \bar{\phi}_0(x)\}$, it's easy to bound $\bar{\psi}_0$ from below and above as before and to show that $(\bar{\phi}_0, \bar{\psi}_0) \in \Phi_c$. This shows that $\bar{\phi}_0(x) \geq -\|c\|_\infty - 1/2$ and $\bar{\psi}_0(y) \geq -\|c\|_\infty - 1/2$, and putting together our estimates gives the result; namely, it follows that $(\bar{\phi}_0, \bar{\psi}_0)$ is close to being a dual pair. We then finish the proof by approximating a general l.s.c. cost function from below by bounded Lipschitz functions (as in [Proposition 3.2.2](#)) and using the monotone convergence theorem. \square

Definition 4.4.1 (*c-concave functions*). A function $f : \mathcal{X} \rightarrow \mathbb{R}$ is called *c-concave* if it can be written as

$$f(x) = \inf_{y \in \mathcal{Y}} \{c(x, y) - g(y)\}$$

for some function $g : \mathcal{Y} \rightarrow \mathbb{R}$. If this is the case, we say that $f = g^c$ is the *c-concave conjugate* of g .

The proof of [Theorem 4.4.1](#) shows that if c is bounded, the supremum can be taken over pairs of *c-concave conjugates* (ϕ^{cc}, ϕ^c) for bounded functions ϕ .

4.4.2 Metric cost functions

If the cost function $c(x, y) = d(x, y)$ is a metric, we have a stronger duality theorem.

Theorem 4.4.2 (Kantorovich-Rubinstein duality). *Suppose $\mathcal{X} = \mathcal{Y}$ is a Polish space with an l.s.c. metric cost d , and define $\text{Lip}_1(\mathcal{X})$ to be the set of 1-Lipschitz functions on \mathcal{X} . Then, we have strong duality:*

$$\inf_{\pi \in \Pi(\mu, \nu)} \int d(x, y) d\pi(x, y) = \sup_{\phi \in \text{Lip}_1(\mathcal{X})} \left\{ \int \phi(x) d(\mu - \nu)(x) \right\}.$$

Proof. We define $d_n = d/(1 + d/n)$, which is a distance bounded by 1 and whose 1-Lipschitz functions are a subset of $\text{Lip}_1(\mathcal{X})$. By a similar argument to the proof of the Kantorovich duality theorem ([Theorem 4.4.1](#)), it suffices to show the result when $d = d_n$ (since $d_n \uparrow d$); hence, we assume d is bounded so all 1-Lipschitz functions are bounded. From the Kantorovich duality theorem, we have

$$\sup_{(\phi, \psi) \in \Phi_d} \left\{ \int \phi(x) d\mu(x) + \int \psi(y) d\nu(y) \right\} = \sup_{\phi \in L^1(\mu)} \left\{ \int \phi^{dd}(x) d\mu(x) + \int \phi^d(y) d\nu(y) \right\}.$$

But since ϕ^d is defined as an infimum of 1-Lipschitz functions (bounded from below), it is also 1-Lipschitz. Furthermore, we see from the definition that $\phi^{dd} = -\phi^d$, and the result follows. \square

For instance, the Kantorovich-Rubinstein duality theorem holds for $W_1(\mu, \nu)$.

4.5 Brenier's theorem

Assume that c is the quadratic cost function on \mathbb{R}^d and let λ denote the Lebesgue measure on \mathbb{R}^d . In this section, we show the astonishing result (due to Brenier) that as long as $\mu \ll \lambda$, the optimal transport map is unique and is given by the graph of the gradient of a convex function. Furthermore, *any* valid transport coupling which is the gradient of a convex function is optimal.

Lemma 4.5.1. *Suppose π_* is an optimal coupling of μ and ν with respect to the quadratic cost function $c(x, y) = \|x - y\|_2^2$. Then, the support of π_* is cyclically monotone. In particular, $\text{supp}(\pi_*)$ is a subset of the graph of a proper l.s.c. convex function.*

Proof. Suppose that $\text{supp}(\pi_*)$ is not cyclically monotone, so that there exist points $(x_1, y_1), \dots, (x_m, y_m) \in \text{supp}(\pi_*)$ such that

$$0 < \sum_{i=1}^m \langle y_i, x_{i+1} - x_i \rangle.$$

Rearranging the terms, this is equivalent to

$$\sum_{i=1}^m \|x_i - y_i\|^2 > \sum_{i=1}^m \|x_{i+1} - y_i\|^2.$$

By continuity of the norm, we can find neighborhoods U_i of x_i and V_i of y_i such that the inequality holds for all $x'_i \in U_i$ and $y'_i \in V_i$. The point of this is that now we can contradict optimality of π_* by writing down a better coupling π . Define Borel measures π_i on $\mathbb{R}^d \times \mathbb{R}^d$ by $\pi_i(\cdot) = \pi_*(\cdot | U_i \times V_i)$, and let $\pi_i^{(1)}$ and $\pi_i^{(2)}$ denote the marginals of π_i . Now, for any $\epsilon > 0$ we define

$$\pi = \pi_* + \frac{\epsilon}{m} \sum_{i=1}^m (\pi_{i+1}^{(1)} \times \pi_i^{(2)} - \pi_i).$$

Note that this definition is motivated by the analogous structure of the inequality defining cyclical monotonicity. For all $A \in \mathcal{B}(\mathbb{R}^d \times \mathbb{R}^d)$, we have

$$\pi(A) \geq \pi_*(A) - \frac{\epsilon}{m} \sum_{i=1}^m \pi_i(A) \geq \pi_*(A) - \frac{\epsilon \pi_*(A)}{m} \sum_{i=1}^m \frac{1}{\pi_*(U_i \times V_i)}.$$

If ϵ is chosen small enough, then we can guarantee $\pi(A) \geq 0$. Now, it is easy to check that $\pi \in \Pi(\mu, \nu)$. But then this implies that

$$\begin{aligned} & \int \|x - y\|_2^2 d\pi(x, y) - \int \|x - y\|_2^2 d\pi_*(x, y) \\ &= \frac{\epsilon}{m} \sum_{i=1}^m \left(\int_{U_{i+1} \times V_i} \|x - y\|_2^2 d\pi_{i+1}^{(1)}(x) d\pi_i^{(2)}(y) - \int_{U_i \times V_i} \|x - y\|_2^2 d\pi_i(x, y) \right). \end{aligned}$$

This expression is strictly negative (since each term in the sum is strictly negative) which contradicts the optimality of π_* . The second statement follows from Rockafellar's theorem ([Theorem 3.6.1](#)). \square

Note. The converse of [Lemma 4.5.1](#) is also true: if $\pi_* \in \Pi(\mu, \nu)$ has cyclically monotone support (for $\mu, \nu \in \mathcal{P}_2(\mathbb{R}^d)$) then it is optimal for the quadratic cost.

In the following, let λ denote the Lebesgue measure on \mathbb{R}^d .

Theorem 4.5.2 (Brenier). *Let $\mu, \nu \in \mathcal{P}_2(\mathbb{R}^d)$ be such that $\mu \ll \lambda$. Then, the following are equivalent.*

1. $\pi_* \in \Pi(\mu, \nu)$ is optimal for the Kantorovich problem with $c(x, y) = \|x - y\|_2^2$.
2. There exists a proper l.s.c. convex function $\varphi : \mathbb{R}^d \rightarrow \mathbb{R}$ such that $\pi_* = (\text{Id} \times \nabla \varphi)_\# \mu$, where $\nabla \varphi$ is defined μ -almost everywhere.
3. The supremum in the Kantorovich duality ([Theorem 4.4.1](#))

$$\int \|x - y\|_2^2 d\pi_*(x, y) = \sup_{(\phi, \psi) \in \Phi_c} \left\{ \int \phi(x) d\mu(x) + \int \psi(y) d\nu(y) \right\}$$

is attained for $\phi(x) = \|x\|_2^2 - 2\varphi(x)$ and $\psi(y) = \|y\|_2^2 - 2\varphi^*(y)$. The dual variables ϕ and ψ are called Kantorovich potentials for (μ, ν) .

Notice that Brenier's theorem shows that the solution to the Kantorovich problem is essentially unique and coincides with the solution to the Monge problem.

Proof. First, we show (1) implies (2). We know by [Lemma 4.5.1](#) that $\text{supp}(\pi_*)$ is contained in the graph of $\partial\varphi$ for some proper l.s.c. convex function φ . Since proper convex functions are λ -a.e. differentiable on the interior of their domain, it follows that $\pi_* = (\text{Id} \times \nabla \varphi)_\# \mu$, where $\nabla \varphi$ is defined μ -almost everywhere by absolute continuity.

Next, we show that (2) implies (3). By the Fenchel-Young inequality ([Theorem 3.5.1](#)), we have

$$\phi(x) + \psi(y) = \|x\|_2^2 + \|y\|_2^2 - 2(\varphi(x) + \varphi^*(y)) \leq \|x\|_2^2 + \|y\|_2^2 - 2\langle x, y \rangle = \|x - y\|_2^2.$$

In fact, by the equality condition of the Fenchel-Young inequality, we find that

$$\begin{aligned} \int \phi(x) d\mu(x) + \int \psi(y) d\nu(y) &= \int (\phi(x) + \psi(y)) d\pi_*(x, y) \\ &= \int (\phi(x) + \psi(\nabla \varphi(x))) d\mu(x) \\ &= \int \|x - y\|_2^2 d\pi_*(x, y). \end{aligned}$$

so ϕ and ψ are dual-optimal. All that's left is to show that $\phi \in L^1(\mu)$ and $\psi \in L^1(\nu)$, which will imply that $(\phi, \psi) \in \Phi_c$. Note that $\varphi = \varphi^{**}$ by the Fenchel-Moreau theorem ([Theorem 3.5.2](#)), so φ^* and $\varphi = \varphi^{**}$ are bounded below by affine functions (meaning that ϕ_+ and ψ_+ are integrable). Since

$$\int \phi(x) d\mu(x) + \int \psi(y) d\nu(y) = \int \|x - y\|_2^2 d\pi_*(x, y) \geq 0,$$

it follows that ϕ and ψ are integrable with respect to μ and ν respectively.

Finally, (3) implies (1) is obvious from the Kantorovich duality ([Theorem 4.4.1](#)) and that for any $\pi \in \Pi(\mu, \nu)$, we have by the Fenchel-Young inequality ([Theorem 3.5.1](#)) that

$$\int \|x - y\|_2^2 d\pi_*(x, y) = \int \phi(x) d\mu(x) + \int \psi(y) d\nu(y) \leq \int \|x - y\|_2^2 d\pi(x, y). \quad \square$$

This theorem allows us to easily construct optimal transport maps for the quadratic cost.

Corollary 4.5.2.1. *If $\mu, \nu \in \mathcal{P}_2(\mathbb{R})$ are such that $\mu \ll \lambda$, then the optimal transport map is $F_\nu^{-1} \circ F_\mu$, where F^{-1} denotes the quantile transformation.*

Proof. It's clear that $(F_\nu^{-1} \circ F_\mu)_\# \mu = \nu$ by inverse transform sampling, and $F_\nu^{-1} \circ F_\mu$ is monotone so it is the derivative of a convex function; for instance, this follows from Rockafellar's theorem ([Theorem 3.6.1](#)). Now apply Brenier's theorem ([Theorem 4.5.2](#)) to see that it is optimal. \square

Corollary 4.5.2.2. *If $\mu = \mathcal{N}(m_1, \Sigma_1)$ and $\nu = \mathcal{N}(m_2, \Sigma_2)$ are two Gaussian measures on \mathbb{R}^d with $\Sigma_1 \succ 0$ and $\Sigma_2 \succ 0$, then the optimal transport map for the quadratic cost is given by*

$$T(x) = \Sigma_1^{-1/2} (\Sigma_1^{1/2} \Sigma_2 \Sigma_1^{1/2})^{1/2} \Sigma_1^{-1/2} (x - m_1) + m_2,$$

which induces the Wasserstein distance

$$W_2^2(\mu, \nu) = \|m_1 - m_2\|^2 + \text{tr} \left(\Sigma_1 + \Sigma_2 - 2 (\Sigma_1^{1/2} \Sigma_2 \Sigma_1^{1/2})^{1/2} \right).$$

Essentially, the motivation for this corollary is as follows. We make the ansatz that the transport map $T(x) = Ax + b$ is affine. Furthermore, we pick a symmetric $A \succeq 0$ so that T is the gradient of a convex function, which is required by Brenier's theorem ([Theorem 4.5.2](#)). There are no constraints on b , so we may assume w.l.o.g. that $m_1 = m_2$. Then, since we need

$$\Sigma_2 = \int (Ax)(Ax)^\top d\mu(x) = A \Sigma_1 A^\top,$$

this forces

$$\Sigma_1^{1/2} \Sigma_2 \Sigma_1^{1/2} = \Sigma_1^{1/2} (A \Sigma_1 A^\top) \Sigma_1^{1/2} = (\Sigma_1^{1/2} A \Sigma_1^{1/2})^2.$$

Solving for A yields the formula in the corollary.

Proof. Note that T is linear and $J_T = \Sigma_1^{-1/2} (\Sigma_1^{1/2} \Sigma_2 \Sigma_1^{1/2})^{1/2} \Sigma_1^{-1/2} \succeq 0$ is symmetric. Therefore $T = \nabla f$ is a conservative vector field and $H_f \succeq 0$ implies that f is convex. It is also easy to check that $T_\# \mu = \nu$. The optimality of T now follows from Brenier's theorem ([Theorem 4.5.2](#)), and the resulting Wasserstein distance is easy to compute. \square

Appendix A

Supplementary results

We use this appendix to discuss some results that are used in the main text (and which may be of interest in their own right).

A.1 Zorn's lemma

In this section, we discuss Zorn's lemma, which is a generalization of induction to infinite sets. It is equivalent to the axiom of choice and transfinite induction, so we omit the proof of equivalence here and assume it as an axiom.

Definition A.1.1 (Poset). A *partially ordered set* (or *poset*) is a set P with a binary relation \leq that is reflexive, antisymmetric, and transitive.

Note that there may be a, b in a poset where \leq isn't defined between them.

Definition A.1.2 (Chain). A *chain* (or *totally ordered set*) in a poset P is a subset $C \subseteq P$ such that \leq is defined between all elements.

Definition A.1.3 (Maximal element). An element $m \in P$ is *maximal* if $m \leq p$ implies $m = p$.

Axiom A.1.1 (Zorn's lemma). *If all chains have an upper bound in a nonempty poset P , then there is a maximal element $m \in P$.*

For instance, Zorn's lemma implies that all vector spaces have Hamel bases and all Hilbert spaces have a total orthonormal set.

A.2 The Baire category theorem

In this section, we discuss the Baire category theorem, which prevents points in a complete metric space from being too sparse.

Definition A.2.1 (Nowhere dense). A set $M \subseteq X$ is *nowhere dense* if \overline{M} has empty interior.

Definition A.2.2 (Meager). A set $M \subseteq X$ is *meager* if it is a countable union of nowhere dense sets and is *nonmeager* otherwise.

Theorem A.2.1 (Baire category theorem). *If X is a complete metric space, then X is nonmeager. Equivalently, the intersection of countably many open dense sets in X is also dense.*

Proof. The proof proceeds by contradiction; we construct a Cauchy sequence using a nested sequence of balls. Let $X = \bigcup_{n=1}^{\infty} M_n$ where each M_n is nowhere dense. Then, $\overline{M_1}^c$ is nonempty and open, so we can find a ball $B_1 \subseteq \overline{M_1}^c$. We can then find a ball $B_2 \subseteq B_1 \cap \overline{M_2}^c$ with radius at most half of B_1 's radius, and so on. This gives a nested sequence of balls $B_1 \supseteq B_2 \supseteq \cdots$. By the completeness of X , there is a unique point $x \in \bigcap_{n=1}^{\infty} B_n$. But now, by construction, $x \notin M_n$ for any n , which gives the desired contradiction. \square

A.3 Tychonoff's theorem

In this section, we discuss Tychonoff's theorem, which states that the product of compact spaces is compact.

Definition A.3.1 (Finite intersection property). A collection \mathcal{F} of sets has the *finite intersection property* if every finite subcollection of \mathcal{F} has nonempty intersection.

The following lemma is a useful characterization of compactness in topological spaces, and is easily equivalent to the usual definition.

Lemma A.3.1. *A topological space X is compact if and only if every collection of closed sets with the finite intersection property has nonempty intersection.*

Proof. The forward direction is elementary and follows from Cantor's intersection theorem. We'll now show the contrapositive of the reverse direction. Suppose that X is not compact and \mathcal{U} is a collection of open sets with no finite subcover. Then, the complements of these sets form a collection of closed sets with the finite intersection property, but their intersection is empty. \square

Theorem A.3.2 (Tychonoff). *If X_i is compact for all $i \in I$, then $\prod_{i \in I} X_i$ is compact.*

Proof. Let \mathcal{F} be a family of closed sets in $\prod_{i \in I} X_i$ with the finite intersection property; by [Lemma A.3.1](#), it suffices to show that \mathcal{F} has nonempty intersection. Now, let \mathcal{P} be the set of pairs (B, x_B) such that $B \subseteq I$ is finite, $x_B \in \prod_{i \in B} X_i$, and for any finite subfamily $\mathcal{G} \subseteq \mathcal{F}$ there exists $y \in \bigcap_{G \in \mathcal{G}} G$ such that $\pi_B(y) = x_B$. We define a partial order on \mathcal{P} by $(B, x_B) \leq (C, x_C)$ if $B \subseteq C$ and $\pi_B(x_C) = x_B$ (where π_B denotes the projection onto the coordinates in B). It's clear that any chain has an upper bound, so by Zorn's lemma ([Axiom A.1.1](#)), there exists a maximal element (M, x_M) . If $M \neq I$, pick $i_0 \in I \setminus M$ and define (for each finite $\mathcal{G} \subseteq \mathcal{F}$):

$$S_{\mathcal{G}} = \left\{ x_0 \in X_{i_0} : \exists y \in \bigcap_{G \in \mathcal{G}} G \text{ such that } \pi_M(y) = x_M \text{ and } \pi_{i_0}(y) = x_0 \right\}.$$

Now, the collection of possible sets $S_{\mathcal{G}}$ satisfies the finite intersection property since $(M, x_M) \in \mathcal{P}$. Hence, the intersection of the possible sets $S_{\mathcal{G}}$ is nonempty by compactness of X_{i_0} . But this gives an extension of (M, x_M) , which is a contradiction. \square