# Metric embeddings and geometric inequalities

Lectures by Professor Assaf Naor
Scribes: Holden Lee and Kiran Vodrahalli

February 22, 2016

# Contents

# Introduction

The topic of this course is geometric inequalities with applications to metric embeddings; and we are actually going to do things more general than metric embeddings: Metric geometry. Today I will roughly explain what I want to cover and hopefully start proving a first major theorem. The strategy for this course is to teach novel work. Sometimes topics will be covered in textbooks, but a lot of these things will be a few weeks old. There are also some things which may not have been even written yet. I want to give you a taste for what's going on in the field. Notes from the course I taught last spring are also available.

One of my main guidelines in choosing topics will be topics that have many accessible open questions. I will mention open questions as we go along. I'm going to really choose topics that have proofs which are entirely self-contained. I'm trying to assume nothing. My intention is to make it completely clear and there should be no steps you don't understand.

Now this is a huge area. I will explain some of the background today. I'm going to use proofs of some major theorems as excuses to see the lemmas that go into the proofs. Some of the theorems are very famous and major, and you'er going to see some improvements, but along the way, we will see some lemmas which are **immensely powerful**. So we will always be proving a concrete theorem. But actually somewhere along the way, there are lemmas which have wide applicability to many many areas. These are excuses to discuss methods, though the theorems are important.

The course can go in many directions: If some of you have some interests, we can always change the direction of the course, so express your interests as we go along.

# Chapter 1

# The Ribe Program

## 1   The Ribe Program

The main motivation for most of what we will discuss is called the Ribe Program, which is a research program many hundreds of papers large. We will see some snapshots of it, and it all comes from a theorem from 1975, **Ribe's rigidity theorem** 1.1.9, which we will state now and prove later in a modern way. This theorem was Martin Ribe's dissertation, which started a whole direction of mathematics, but after he wrote his dissertation he left mathematics. He's apparently a government official in Sweden. The theorem is in the context of Banach spaces; a relation between their linear structure and their structure as metric spaces. Now for some terminology.

**Definition 1.1.1:** Banach space.
A Banach space is a complete, normed vector space. Therefore, a Banach space is equipped with a metric which defines vector length and distances between vectors. It is complete, so every Cauchy sequence of converges to a limit defined inside the space.

**Definition 1.1.2:** Let $(X, \|\cdot\|_X), (Y, \|\cdot\|_Y)$ be Banach spaces. We say that $X$ is (crudely) **finitely representable** in $Y$ if there exists some constant $K > 0$ such that for every finite-dimensional linear subspace $F \subseteq X$, there is a linear operator $S : F \to Y$ such that for every $x \in F$,

$$\|x\|_X \leq \|Sx\|_Y \leq K \|x\|_X.$$

   Note $K$ is decided once and for all, before the subspace $F$ is chosen.
   (Some authors use "finitely representable" to mean that this is true for any $K = 1 + \varepsilon$. We will not follow this terminology.)
   Finite representability is important because it allows us to conclude that $X$ has the same finite dimensional linear properties (**local properties**) as $Y$. That is, it preserves any invariant involves finitely many vectors, their lengths, etc.

Let's introduce some local properties like type. To motivate the definition, consider the triangle inequality, which says

$$\|y_1 + \cdots + y_n\|_Y \leq \|y_1\|_Y + \cdots + \|y_n\|_Y.$$

In what sense can we improve the triangle inequality? In $L^1$ this is the best you can say. In many spaces there are ways to improve it if you think about it correctly.

For any choice $\varepsilon_1, \ldots, \varepsilon_n \in \{\pm 1\}$,

$$\left\|\sum_{i=1}^n \varepsilon_i y_i\right\|_Y \leq \sum_{i=1}^n \|y_i\|_Y.$$

**Definition 1.1.3:** <span style="color:red">df:type</span> Say that $X$ has **type** $p$ if there exists $T > 0$ such that for every $n, y_1, \ldots, y_n \in Y$,

$$\underset{\varepsilon \in \{\pm 1\}^n}{\mathbb{E}} \left\|\sum_{i=1}^n \varepsilon_i y_i\right\|_Y \leq T \left[\sum_{i=1}^n \|y_j\|_Y^p\right]^{\frac{1}{p}}.$$

The $L^p$ norm is always at most the $L^1$ norm; if the lengths are spread out, this is asymptotically much better. Say $Y$ has **nontrivial type** if $p > 1$.

For example, $L_p(\mu)$ has type $\min(p, 2)$.

Later we'll see a version of "type" for metric spaces. How far is the triangle inequality from being an equality is a common theme in many questions. In the case of normed spaces, this controls a lot of the geometry. Proving a result for $p > 1$ is hugely important.

**Proposition 1.1.4:** <span style="color:red">pr:finrep-type</span> If $X$ is finitely representable and $Y$ has type $p$ then also $X$ has type $p$.

*Proof.* Let $x_1, \ldots, x_n \in X$. Let $F = \mathrm{span}\{x_1, \ldots, x_n\}$. Finite representability gives me $S : F \to Y$. Let $y_i = Sx_i$. What can we say about $\sum \varepsilon_i y_i$?

$$\underset{\varepsilon}{\mathbb{E}} \left\|\sum_{i=1}^n \varepsilon_i y_i\right\|_Y = \underset{\varepsilon}{\mathbb{E}} \left\|S(\sum_{i=1}^n \varepsilon_i x_i)\right\|_Y$$

$$\geq \underset{\varepsilon}{\mathbb{E}} \left\|\sum_{i=1}^n \varepsilon_i X_i\right\|_X$$

$$\underset{\varepsilon}{\mathbb{E}} \left\|\sum_{i=1}^n \varepsilon_i y_i\right\|_Y \leq T \left(\sum_{i=1}^n \|Sx_i\|^p\right)^{\frac{1}{p}}$$

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

Putting these two inequalities together gives the result.      $\square$

**Theorem 1.1.5** (Kahane's inequality). *For any normed space $Y$ and $q \geq 1$, for all $n$, $y_1, \ldots, y_n \in Y$,*

$$\mathbb{E} \left\| \sum_{i=1}^{n} \varepsilon_i y_i \right\| \gtrsim_q \left( \mathbb{E} \left[ \left\| \sum_{i=1}^{n} \varepsilon_i y_i \right\|_Y^q \right] \right)^{\frac{1}{q}}.$$

*Here $\gtrsim_q$ means "up to a constant"; subscripts say what the constant depends on. The constant here does not depend on the norm $Y$.*

Kahane's Theorem tells us that the LHS of Definition 1.1.3 can be replaced by any norm, if we change $\leq$ to $\lesssim$. We get that having type $p$ is equivalent to

$$\mathbb{E} \left\| \sum_{i=1}^{n} \varepsilon_i y_i \right\|_Y^p \lesssim T^p \sum_{i=1}^{n} \|y_i\|_Y^p.$$

Recall the **parallelogram identity** in a Hilbert space $H$:

$$\mathbb{E} \left\| \sum_{i=1}^{n} \varepsilon_i y_i \right\|^2 = \sum_{i=1}^{n} \|y_i\|_H^2.$$

A different way to understand the inequality in the definition of "type" is: how far is a given norm from being an Euclidean norm? The **Jordan-von Neumann Theorem** says that if parallelogram identity holds then it's a Euclidean space. What happes if we turn it in an inequality?

$$\mathbb{E} \left\| \sum_{i=1}^{n} \varepsilon_i y_i \right\|_H^2 \overset{\geq}{\underset{\leq}{}} T \sum_{i=1}^{n} \|y_i\|_H^2.$$

Either inequality *still* characterizes a Euclidean space.

What happens if we add constants or change the power? We recover the definition for type and cotype (which has the inequality going the other way):

$$\mathbb{E} \left\| \sum_{i=1}^{n} \varepsilon_i y_i \right\|_H^q \overset{\gtrsim}{\underset{\sim}{\lesssim}} \sum_{i=1}^{n} \|y_i\|_H^q.$$

**Definition 1.1.6:** Say it has **cotype** $q$ if

$$\sum_{i=1}^{n} \|y_i\|_Y^q \lesssim C^q \mathbb{E} \left\| \sum_{i=1}^{n} \varepsilon_i y_i \right\|_Y^q$$

R. C. James invented the local theory of Banach spaces, the study of geometry that involves properties involving finitely many vectors ($\forall x_1, \ldots, x_n, P(x_1, \ldots, x_n)$ holds). As a counterexample, reflexivity cannot be characterized using finitely many vectors (this is a theorem).

Ribe discovered link between metric and linear spaces.

First, terminology.

**Definition 1.1.7:** Two Banach spaces are **uniformly homeomorphic** if there exists $f : X \to Y$ that is 1-1 and onto and $f, f^{-1}$ are uniformly continuous.

Without the word "uniformly", if you think of the spaces as topological spaces, all of them are equivalent. Things become interesting when you quantify! "Uniformly" means you're controlling the quantity.

**Theorem 1.1.8** (Kadec). *Any two infinite-dimensional separable Banach spaces are homeomorphic.*

This is a amazing fact: these spaces are all topologically equivalent to Hilbert spaces!

Over time people people found more examples of Banach spaces that are homeomorphic but not uniformly homeomorphic. Ribe's rigidity theorem clarified a big chunk of what was happening.

**Theorem 1.1.9** (Rigidity Theorem, Martin Ribe (1975)). <span style="color:red">*thm:ribe*</span> *Suppose that $X, Y$ are uniformly homeomorphic Banach spaces. Then $X$ is finitely representable in $Y$ and $Y$ is finitely representable in $X$.*

For example, for $L^p$ and $L^q$, for $p \neq q$ it's always that case that one is not finitely representable in the other, and hence by Ribe's Theorem, $L^p, L^q$ are not uniformly homeomorphic. (When I write $L_p$, I mean $L_p(\mathbb{R})$.)

**Theorem 1.1.10.** *For every $p \geq 1, p \neq 2$, $L_p$ and $\ell_p$ are finitely representable in each other, yet not uniformly homeomorphic.*

(Here $\ell_p$ is the sequence space.)

**Exercise 1.1.11:** Prove the first part of this theorem: $L_p$ is finitely representable in $\ell_p$.

Hint: approximate using step functions. You'll need to remember some measure theory. When $p = 2$, $L_p, \ell_p$ are separable and isometric.

The theorem in various cases was proved by:

1. $p = 1$: Enflo

2. $1 < p < 2$: Bourgain

3. $p > 2$: Gorelik, applying the Brouwer fixed point theorem (topology)

Every linear property of a Banach signs which is local (type, cotype, etc.; involving summing, powers, etc.) is preserved under a general nonlinear deformation.

After Ribe's rigidity theorem, people wondered: can we reformulate the local theory of Banach spaces without mentioning anything about the linear structure? Ribe's rigidity theorem is more of an existence statement, we can't see anything about an explicit dictionary which maps statements about linear sums into statements about metric spaces. So people started to wonder whether we could reformulate the local theory of Banach spaces, but only

looks at distances between pairs instead of summing things up. Local theory is one of the hugest subjects in analysis. If you could actually find a dictionary which takes one linear theorem at a time, and restate it with only distances, there is a huge potential here! Because the definition of type only involves distances between points, we can talk about a metric space's type or cotype. So maybe we can use the intution given by linear arguments, and then state things for metric spaces which often for very different reasons remain true from the linear domain. And then now maybe you can apply these arguments to graphs, or groups! We could be able to prove things about the much more general metric spaces. Thus, we end up applying theorems on linear spaces in situations with *a priori nothing* to do with linear spaces. This is massively powerful.

There are very crucial entries that are missing in the dictionary. We don't even now how to define many of the properties! This program has many interesting proofs. Some of the most interesting conjectures are how to define things!

**Corollary 1.1.12.** *cor:uh-type If $X, Y$ are uniformly homeomorphic and if one of them is of type $p$, then the other does.*

This follows from Ribe's Theorem 1.1.9 and Proposition 1.1.4. Can we prove something like this theorem without using Ribe's Theorem 1.1.9? We want to reformulate the definition of type using only the distance, so this becomes self-evident.

Enflo had an amazing idea. Suppose $X$ is a Banach space, $x_1, \ldots, x_n \in X$. The type $p$ inequality says

$$\text{eq:type-p} \quad \mathbb{E}\left[ \left\| \sum_{i=1}^{n} \varepsilon_i x_i \right\|^p \right] \lesssim_X \sum_{i=1}^{n} \|x_i\|^p. \tag{1.1}$$

Let's rewrite this in a silly way. Define $f : \{\pm 1\}^n \to X$ by

$$f(\varepsilon_1, \ldots, \varepsilon_n) = \sum_{i=1}^{n} \varepsilon_i x_i.$$

Write $\varepsilon = (\varepsilon_1, \ldots, \varepsilon_n)$. Multiplying by $2^n$, we can write the inequality (1.1) as

$$\text{eq:type-gen} \quad \mathbb{E}\left[ \|f(\varepsilon) - f(-\varepsilon)\|^p \right] \lesssim_X \sum_{i=1}^{n} \mathbb{E}\left[ \|f(\varepsilon) - f(\varepsilon_1, \ldots, \varepsilon_{i-1}, -\varepsilon_i, \varepsilon_{i+1}, \ldots, \varepsilon_n)\|^p \right]. \tag{1.2}$$

This inequality just involves distances between points $f(\varepsilon)$, so it is the reformulation we seek.

**Definition 1.1.13:** *df:enflo A metric space* $(X, d_X)$ *has* **Enflo type** $p$ *if there exists* $T > 0$ *such that for every* $n$ *and every* $f : \{\pm 1\}^n \to X$,
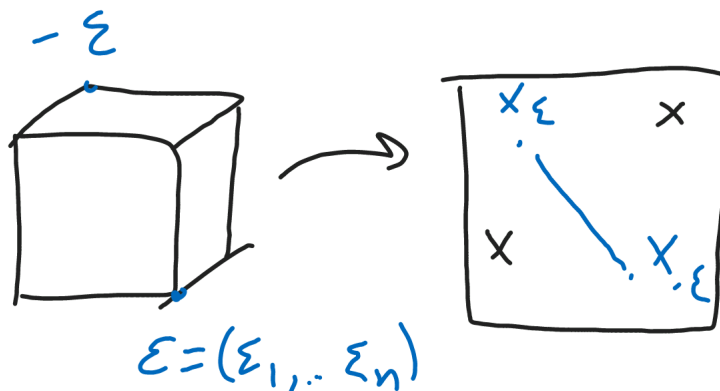
$$\mathbb{E}[d_X(f(\varepsilon), f(-\varepsilon))^p] \leq T^p \sum_{i=1}^{n} \mathbb{E}[d_X(f(\varepsilon), f(\varepsilon_1, \ldots, \varepsilon_{i-1}, -\varepsilon_i, \varepsilon_{i+1}, \ldots, \varepsilon_n))^p].$$

This is bold. It wasn't true before for a general function! The discrete cube (Boolean hypercube) $\{\pm 1\}^n$ is all the $\epsilon$ vectors, of which there are $2^n$. Our function just assigns $2^n$

points arbitrarily. No structure whatsoever. As they are indexed this way, you see nothing. But you're free to label them by vertices of the cube however you want. But there are many labelings! In (1.2), the points had to be vertices of a cube, but in Definition 1.1.13, they are arbitrary. The moment you choose the labelings, you impose a cube structure between the points. Some of them are diagonals of the cube, some of them are edges. $\epsilon$ and $-\epsilon$ are antipodal points. But it's not really a diagonal. They are points on a manifold, and are a function of how you decided to label them. What this sum says is that the sum over all diagonals, the length of the diagonals to the power $p$ is less than the sum over edges to the $p^{th}$ powers (these are the points where one $\epsilon_i$ is different). Thus we can see

$$\sum \text{diag}^p \lesssim_X \sum \text{edge}^p.$$

The total $p$th power of lengths of diagonals is up to a constant, at most the same thing over all edges.



This is a vast generalization of type; we don't even know a Banach space satisfies this. The following is one of my favorite conjectures.

**Conjecture 1.1.14** (Enflo). *If a Banach space has type $p$ then it also has Enflo type $p$.*

This has been open for 40 years. We will prove the following.

**Theorem 1.1.15** (Bourgain-Milman-Wolfson, Pisier). *If $X$ is a Banach space of type $p > 1$ then $X$ also has type $p - \varepsilon$ for every $\varepsilon > 0$.*

If you know the type inequality for parallelograms, you get it for arbitrary sets of points, up to $\varepsilon$. Basically, you're getting arbitrarily close to $p$ instead of getting the exact result. We also know that the conjecture stated before is true for a lot of specific Banach spaces, though we do not yet have the general result. For instance, this is true for the $L_4$ norm. Index functions by vertices; some pairs are edges, some are diagonals; then the $L^4$ norm of the diagonals is at most that of the edges.

How do you go from knowing this for a linear function to deducing this for an arbitrary function?

Once you do this, you have a new entry in the hidden Ribe dictionary. If $X$ and $Y$ are uniformly homeomorphic Banach spaces and $Y$ has Enflo type $p$, then so is $X$. The minute you throw away the linear structure, Corollary 1.1.12 becomes easy. It requires a tiny argument. Now you can take a completely arbitrary function $f : \{\pm 1\}^n \to X$. There exists a homeomorphism $\psi : X \to Y$ such that $\psi, \psi^{-1}$ are uniformly continuous. Now we want to deduce that the same inequality in $Y$ gives the same inequality in $X$.

**Proposition 1.1.16:** <span style="color:red">pr:uh-enflo</span> If $X, Y$ are uniformly homeomorphic Banach spaces and $Y$ has Enflo type $p$, then so does $X$.

This is an example of making the abstract Ribe theorem explicit.

**Lemma 1.1.17** (Corson-Klee). <span style="color:red">lem:corson-klee</span> *If $X, Y$ are Banach spaces and $\psi : X \to Y$ are uniformly continuous, then for every $a > 0$ there exists $L(a)$ such that*

$$\|x_1 - x_2\|_X \geq a \implies \|\psi(x_1) - \psi(x_2)\| \leq L \|x_1 - x_2\|.$$

*Proof sketch of 1.1.16 given Lemma 1.1.17.* By definition of uniformly homeomorphic, there exists a homeomorphism $\psi : X \to Y$ such that $\psi, \psi^{-1}$ are uniformly continuous. Lemma 1.1.17 tells us that $\psi$ perserves distance up to a constant. Dividing so that the smallest nonzero distance you see is at least 1, we get the same inequality in the image and the preimage. $\square$

*Proof details.* Let $\varepsilon^i$ denote $\varepsilon$ with the $i$th coordinate flipped. We need to prove

$$\mathbb{E}(d_X(f(\varepsilon), f(-\varepsilon))^p) \leq T_f^p \sum_{i=1}^{n} \mathbb{E}(d_X(f(\varepsilon), f(\varepsilon^i))^p)$$
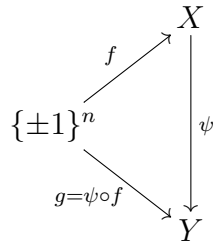
Without loss of generality, by scaling $f$ we may assume that all the points $f(\varepsilon)$ are distance at least 1 apart. ($X$ is a Banach space, so distance scales linearly; this doesn't affect whether the inequality holds.)

Let $\psi : X \to Y$ be such that $\psi, \psi^{-1}$ are uniform homeomorphisms. Because $\psi^{-1}$ is uniformly homeomorphic, there is $C$ such that $d_Y(y_1, y_2) \leq 1$ implies $d_X(\psi^{-1}(y_1), \psi^{-1}(y_2)) < C$. WLOG, by scaling $f$ we may assume that all the points $f(\varepsilon)$ are $\max(1, C)$ apart, so that the points $\psi \circ f(\varepsilon)$ are at least 1 apart.

We know that for any $g : \{\pm 1\}^n \to Y$ that

$$\mathbb{E}(d_X(g(\varepsilon), g(-\varepsilon))^p) \leq T_g^p \sum_{i=1}^{n} \mathbb{E}(d_X(g(\varepsilon), g(\varepsilon^i))^p).$$

We apply this to $g = \psi \circ f$,

to get

$$\begin{aligned}
\mathbb{E}(d_X(f(\varepsilon), f(-\varepsilon))^p) &= \mathbb{E}(d_X(\psi^{-1} \circ g(\varepsilon), \psi^{-1} \circ g(-\varepsilon))^p) \\
&\leq L_{\psi^{-1}}(1)\mathbb{E}(d_Y(g(\varepsilon), g(-\varepsilon))^p) \\
&\leq L_{\psi^{-1}}(1)T_g^p \sum_{i=1}^n \mathbb{E}(d_Y(g(\varepsilon), g(\varepsilon^i))) \\
&\leq L_{\psi^{-1}}(1)L_\psi(1)T_g^p \sum_{i=1}^n \mathbb{E}(d_X(g(\varepsilon), g(\varepsilon^i))^p)
\end{aligned}$$

as needed.                                                                                      $\square$

The parallelogram inequality for exponent 1 instead of 2 follows from using the triangle inequality on all possible paths for all paths of diagonals. Type $p > 1$ is a strengthening of the triangle inequality. For which metric spaces does it hold?

What's an example of a metric space where the inequality doesn't hold with $p > 1$? The cube itself (with $L^1$ distance).

$$n^p \nleq n.$$

I will prove to you that this is the only obstruction: given a metric space that doesn't contain bi-Lipschitz embeddings of arbitrary large cubes, the inequality holds.

We know an alternative inequality involving distance equivalent to type; I can prove it. It is, however, not a satisfactory solution to the Ribe program. There are other situations where we have complete success.

We will prove some things, then switch gears, slow down and discuss Grothendieck's inequality and applications. They will come up in the nonlinear theory later.

# 2   Bourgain's Theorem implies Ribe's Theorem

2-3-16

We will use the Corson-Klee Lemma 1.1.17.

*Proof of Lemma 1.1.17.* Suppose $x, y \in X$, $\|x - y\| \geq a$. Break up the line segment from $x, y$ into intervals of length $a$; let $x = x_0, x_1, \ldots, x_k = y$ be the endpoints of those intervals, with

$$\|x_{i+1} - x_i\| \leq a.$$

The **modulus of continuity** is defined as

$$W_f(t) = \sup \{\|f(u) - f(v)\| : u, v \in X, \|u - v\| \leq t\}.$$

Uniform continuity says $\lim_{t \to 0} W_f(t) = 0$. The number of intervals is

$$k \leq \frac{\|x - y\|}{a} + 1 \leq \frac{2\|x - y\|}{a}.$$

14

Then

$$\|f(x) - f(y)\| \leq \sum_{i=1}^{k} \|f(x_i) - f(x_{i-1})\|$$

$$\leq K W_f(a) \leq \frac{2W_f(a)}{a} \|x - y\|,$$

so we can let $L(a) = \frac{2W_f(a)}{a}$. $\qquad \square$

## 2.1 Bourgain's discretization theorem

There are 3 known proofs of Ribe's Theorem.

1. Ribe's original proof, 1987.

2. HK, 1990, a genuinely different proof.

3. Bourgain's proof, a Fourier analytic proof which gives a quantitative version. This is the version we'll prove.

Bourgain uses the Discretization Theorem 1.2.4. There is an amazing open problem in this context.

Saying $\delta$ is big says there is a not-too-fine net, which is enough. Therefore we are interested in lower bounds on $\delta$.

**Definition 1.2.1:** Discretization modulus.
Let $X$ be a finite-dimensional normed space $\dim(X) = n < \infty$. Let target space $Y$ be an arbitrary Banach space. Consider the unit ball $B_X$ in $X$. Take an $\delta$-net (the distance between points is at most $\delta$) $\mathcal{N}_\delta$ in $B_X$ (a maximal $\delta$-separated subset). Suppose we can embed $\mathcal{N}_\delta$ into $Y$ via $f : \mathcal{N}_\delta \to Y$. Suppose we know in $Y$ that

$$\|x - y\| \leq \|f(x) - f(y)\| \leq D \|x - y\|.$$

for all $x, y \in N_\delta$. (We say that $\mathcal{N}_\delta$ embeds with distortion $D$ into $Y$.)

You can prove using a nice compactness argument that if this holds for $\delta$ is small enough, then the entire space $X$ embeds into $Y$ with rough the same distortion. Bourgain's discretization theorem 1.2.4 says that you can choose $\delta = \delta_n$ to be independent of the geometry of $X$ and $Y$ such that if you give a $\delta$-approximation of the unit-ball in the $n$-dimensional norm, you succeed in embedding the whole space.

I often use this theorem in this way: I use continuous methods to show embedding $X$ into $Y$ requires big distortion; immediately I get an example with a finite object. Let us now make the notion of distortion more precise.

**Definition 1.2.2:** Distortion.

Suppose $(X, d_X), (Y, d_Y)$ are metric spaces $D \geq 1$. We say that $X$ embeds into $Y$ with **distortion** $D$ if there exists $f : X \to Y$ and $s > 0$ such that for all $x, y \in X$,

$$Sd_X(x, y) \leq d_Y(f(x), f(y)) \leq DSd_X(x, y).$$

The infimum over those $D \geq 1$ such that $X$ embeds into $Y$ with distortion is denoted $C_Y(X)$. This is a measure of how far $X$ is being from a subgeometry of $Y$.

**Definition 1.2.3:** Let be a $n$-dimensional normed space and $Y$ be any Banach space, $\varepsilon \in (0, 1)$. Let $\delta_{X \hookrightarrow Y}(\varepsilon)$ be the supremum over all those $\delta > 0$ such that for every $\delta$-net $\mathcal{N}_\delta$ in $B_X$,
$$C_Y(\mathcal{N}_\delta) \geq (1 - \varepsilon)C_Y(X).$$

Here $B_X := \{x \in X : \|x\| \leq 1\}$.

In other words, the distortion of the $\delta$-net is not much larger than the distortion of the whole space. That is, the discrete $\delta$-ball encodes almost all information about the space when it comes to embedding into $Y$: If you got $C_Y(\mathcal{N}_\delta)$ to be small, then the distortion of the entire object is not much larger.

**Theorem 1.2.4** (Bourgain's discretization theorem). *For every $n$, $\varepsilon \in (0, 1)$, for every $X, Y$, $\dim X = n$,*

$$\delta_{X \hookrightarrow Y}(\varepsilon) \geq e^{-\left(\frac{n}{\varepsilon}\right)^{Cn}}.$$

*Thus there is a delta which is just dependent on the dimension such that in any n-dim norm space if you look at the unit ball it encodes all the information of embedding $X$ into **anything else**. It's only a function of the dimension, not of any of the relevant geometry. Moreover for $\delta = e^{-(2n)^{Cn}}$, we have $C_Y(X) \leq 2C_Y(\mathcal{N}_\delta)$ via a linear operator.*

The theorem says that if you look at a $\delta$-net in the unit ball, it encodes all the information about $X$ when it comes to embedding into everything else. The amount you have to discretize is just a function of the dimension, and not of any of the other relevant geometry.

**Remark 1.2.5:** The proof is via a linear operator. All the inequality says is that you can find a *function* with the given distortion. The proof will actually give a *linear operator*.

The best known upper bound is

$$\delta_{X \hookrightarrow Y}\left(\frac{1}{2}\right) \lesssim \frac{1}{n}.$$

The latest progress was 1987, there isn't a better bound yet. You have a month to think about it before you get corrupted by Bourgain's proof.

There is a better bound when the target space is a $L^p$ space.

**Theorem 1.2.6** (Gladi, Naor, Shechtman)**.** *For every $p \geq 1$, if $\dim X = n$,*

$$\delta_{X \hookrightarrow L_p}(\varepsilon) \gtrsim \frac{\varepsilon^2}{n^{\frac{5}{2}}}$$

(We still don't know what the right power is.) The case $p = 1$ is important for applications. There are larger classes for spaces where we can write down axioms for where this holds. There are crazy Banach spaces which don't belong to this class, so we're not done. We need more tools to show this: Lipschitz extension theorems, etc.

## 2.2    Bourgain's Theorem implies Ribe's Theorem

With the "moreover," Bourgain's theorem implies Ribe's Theorem 1.1.9.

*Proof of Ribe's Theorem 1.1.9 from Bourgain's Theorem 1.2.4.* Let $X, Y$ be Banach spaces that are uniformly homeomorphic. By Corson-Klee 1.1.17, there exists $f : X \to Y$ such that

$$\|x - y\| \geq 1 \implies \|x - y\| \leq \|f(x) - f(y)\| \leq K \|x - y\| .$$

(Apply the Corson-Klee lemma for both $f$ and the inverse.)

In particular, if $R > 1$ and $\mathcal{N}$ is a 1-net in

$$RB_X = \{ x \in X : \|x\| \leq R \} ,$$

then $C_Y(\mathcal{N}) \leq K$. Equivalently, for every $\delta > 0$ every $\delta$-net in $B_X$ satisfies $C_Y(\mathcal{N}) \leq K$. If $F \subseteq X$ is a finite dimension subspace and $\delta = e^{-(2 \dim F)^{C \dim F}}$, then by the "moreover" part of Bourgain's Theorem 1.2.4, there exists a linear operator $T : F \to Y$ such that

$$\|x - y\| \leq \|Tx - Ty\| \leq 2K \|x - y\|$$

for all $x, y \in F$. This means that $X$ is finitely representable. $\qquad\square$

The motivation for this program comes in passing from continuous to discrete. The theory has many applications, e.g. to computer science whcih cares about finite things. I would like an improvement in Bourgain's Theorem 1.2.4.

First we'll prove a theorem that has nothing to do with Ribe's Theorem. There are lemmas we will be using later. It's an easier theorem. It looks unrelated to metric theory, but the lemmas are relevant.

# Chapter 2

# Restricted invertibility principle

## 1   Restricted invertibility principle

### 1.1   The first restricted invertibility principles

We take basic facts in linear algebra and make things quantitative. This is the lesson of the course: when you make things quantitative, new mathematics appears.

**Proposition 2.1.1:** If $A : \mathbb{R}^m \to \mathbb{R}^n$ is a linear operator, then there exists a linear subspace $V \subseteq \mathbb{R}^n$ with $\dim(V) = \operatorname{rank}(A)$ such that $A : V \to A(V)$ is invertible.

What's the quantitative question we want to ask about this? Invertibility just says that an inverse exists. Can we find a large subspace where not only is $A$ invertible, but the inverse has small norm?

We insist that the subspace is a coordinate subspace. Let $e_1, \ldots, e_m$ be the standard basis of $\mathbb{R}^m$, $e_j = (0, \ldots, \underbrace{1}_{j}, 0, \ldots)$. The goal is to find a "large" subset $\sigma \subseteq \{1, \ldots, m\}$ such that $A$ is invertible on $\mathbb{R}^\sigma$ where

$$\mathbb{R}^\sigma := \{(x_1, \ldots, x_n) \in \mathbb{R}^m : x_i = 0 \text{ if } i \notin \sigma\}$$

and the norm of $A^{-1} : A(\mathbb{R}^\sigma) \to \mathbb{R}^\sigma$ is small.

A priori this seems a crazy thing to do; take a small multiple of the identity. But we can find conditions that allow us to achieve this goal.

**Theorem 2.1.2** (Bourgain-Tzafriri restricted invertibility principle, 1987). *thm:btrip Let $A : \mathbb{R}^m \to \mathbb{R}^m$ be a linear operator such that*

$$\|Ae_j\|_2 = 1$$

*for every $j \in \{1, \ldots, m\}$. Then there exist $\sigma \subseteq \{1, \ldots, m\}$ such that*

   *1. $|\sigma| \geq \frac{cm}{\|A\|^2}$, where $\|A\|$ is the operator norm of $A$.*

2. *A is invertible on $\mathbb{R}^\sigma$ and the norm of $A^{-1} : A(\mathbb{R}^\sigma) \to \mathbb{R}^\sigma$ is at most $C'$ (i.e., $\|AJ_\sigma\|_{S^\infty} \leq C'$, to use the notation introduced below).*

*Here $c, C'$ are universal constants.*

Suppose the biggest eigenvalue is at most 100. Then you can always find a coordinate subset of proportional size such that on this subset, $A$ is invertible and the inverse has norm bounded by a universal constant.

All of the proofs use something very amazing.

This proof is from 3 weeks ago. This has been reproved many times. I'll state a theorem that gives better bound than the entire history.

This was extended to rectangular matrices. (The extension is nontrivial.)

Given $V \subseteq \mathbb{R}^m$ a linear subspace with $\dim V = k$ and $A : V \to \mathbb{R}^m$ a linear operator, the singular values of $A$

$$s_1(A) \geq s_2(A) \geq \cdots \geq s_k(A)$$

are the eigenvalues of $(A^*A)^{\frac{1}{2}}$. We can decompose

$$A = UDV$$

where $D$ is a matrix with $s_i(A)$'s on the diagonal, and $U, V$ are unitary.

**Definition 2.1.3:** For $p \geq 1$ the **Schatten-von Neumann $p$-norm** of $A$ is

$$\|A\|_{S_p} := \left( \sum_{i=1}^{k} s_i(A)^p \right)^{\frac{1}{p}}$$
$$= (\mathrm{Tr}((A^*A)^{\frac{p}{2}}))^{\frac{1}{p}}$$
$$= (\mathrm{Tr}((AA^*)^{\frac{p}{2}}))^{\frac{1}{p}}$$

The cases $p = \infty, 2$ give the operator and Frobenius norm,

$$\|A\|_{S_\infty} = \text{operator norm}$$
$$\|A\|_{S_2} = \sqrt{\mathrm{Tr}(A^*A)} = \left( \sum a_{ij}^2 \right)^{\frac{1}{2}}.$$

**Exercise 2.1.4:** $\|\cdot\|_{S_p}$ is a norm on $\mathcal{M}_{n \times m}(\mathbb{R})$. You have to prove that given $A, B$,

$$(\mathrm{Tr}([(A+B)^*(A+B)]^{\frac{p}{2}}))^{\frac{1}{p}} \leq (\mathrm{Tr}((A^*A)^{\frac{p}{2}}))^{\frac{1}{p}} + (\mathrm{Tr}((B^*B)^{\frac{p}{2}}))^{\frac{1}{p}}.$$

This requires an idea. Note if $A, B$ commute this is trivial. Apparently von Neumann wrote a paper called "Metric Spaces" in the 1930s in which he just proves this inequality and doesn't know what to do with it, so it got forgotten for a while until the 1950s, when Schatten wrote books on applications. When I was a student in grad school, I was taking a class on random matrices. There was two weeks break, I was certain that it was trivial because the professor had not said it was not, and it completely ruined my break though I came up with a different proof of it. It's short, but not trivial: It's not typical linear algebra!. This is like another triangle inequality, which we may need later on.

Spielman and Srivastava have a beautiful theorem.

**Definition 2.1.5:** Stable rank.

Let $A : \mathbb{R}^m \to \mathbb{R}^n$. The **stable rank** is defined as

$$\mathrm{srank}(A) = \left( \frac{\|A\|_{S_2}}{\|A\|_{S_\infty}} \right)^2.$$

The numerator is the sum of squares of the singular values, and the denominator is the maximal value. Large stable rank means that many singular values are nonzero, and in fact large on average. Many people wanted to get the size of the subset in the Restricted Invertibility Principle to be close to the stable rank.

**Theorem 2.1.6** (Spielman-Srivastava). *thm:ss For every linear operator $A : \mathbb{R}^m \to \mathbb{R}^n$, $\varepsilon \in (0, 1)$, there exists $\sigma \subseteq \{1, \ldots, m\}$ with $|\sigma| \geq (1 - \varepsilon)\,\mathrm{srank}(A)$ such that*

$$\left\| (AJ_\sigma)^{-1} \right\|_{S_\infty} \lesssim \frac{\sqrt{m}}{\varepsilon \, \|A\|_{S_2}}.$$

*Here, $J_\sigma$ is the identity function restricted to $\mathbb{R}^\sigma$, $J : \mathbb{R}^\sigma \hookrightarrow \mathbb{R}^m$.*

This is stronger than Bourgain-Tzafriri. In Bourgain-Tzafriri the columns were unit vectors.

*Proof of Theorem 2.1.2 from Theorem 2.1.6.* Let $A$ be as in Theorem 2.1.2. Then $\|A\|_{S_2} = \sqrt{\mathrm{Tr}(A^*A)} = \sqrt{m}$ and $\mathrm{srank}(A) = \frac{m}{\|A\|_{S_\infty}^2}$. We obtain the existence of

$$|\sigma| \geq (1 - \varepsilon) \frac{m}{\|A\|_{S_\infty}^2}$$

with $\|(AJ_\sigma)^{-1}\|_{S_\infty} \lesssim \frac{\sqrt{m}}{=} \frac{1}{\varepsilon}$. $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad \square$

This is a sharp dependence on $\varepsilon$.

The proof introduces algebraic rather than analytic methods; it was eye-opening. Marcus even got sets bigger than the stable rank and looked at $pf\|A\|_{S_2}\|A\|_{S_4}^2$, which is much stronger.

## 1.2   A general restricted invertibility principle

I'll show a theorem that implies all these intermediate theorems. We use (classical) analysis and geometry instead of algebra. What matters is not the ratio of the norms, but the tail of the distribution of $s_1(A)^2, \ldots, s_m(A)^2$.

**Theorem 2.1.7.** *thm:gen-srank Let $A : \mathbb{R}^m \to \mathbb{R}^n$ be a linear operator. If $k < \mathrm{rank}(A)$ then there exist $\sigma \subseteq \{1, \ldots, m\}$ with $|\sigma| = k$ such that*

$$\left\| (AJ_\sigma)^{-1} \right\|_{S_\infty} \lesssim \min_{k < r \leq \mathrm{rank}(A)} \sqrt{\frac{mr}{(r - k) \sum_{i=r}^m s_i(A)^2}}.$$

You have to optimize over $r$. You can get the ratio of $L_p$ norms from the tail bounds. This implies all the known theorems in restricted invertibility. The subset can be as big as you want up to the rank, and we have sharp control in the entire range. This theorem generalizes Spielman-Srivasta (Theorem 2.1.6), which had generalized Bourgain-Tzafriri (Theorem 2.1.2).   2-8-16

Now we will go backwards a bit, and talk about a less general result. After Theorem 2.1.6, a subsequent theorem gave the same theorem but instead of the stable rank, used something better.

**Theorem 2.1.8** (Marcus, Spielman, Srivastava). *thm:mss4 If*

$$k < \frac{1}{4}\left(\frac{\|A\|_{S_2}}{\|A\|_{S_4}}\right)^4,$$

*there exists* $\sigma \subseteq \{1,\ldots,m\}$, $|\sigma| = k$ *such that*

$$\left\|(AJ_\sigma)^{-1}\right\|_{S_\infty} \lesssim \frac{\sqrt{m}}{\|A\|_{S_2}}.$$

A lot of these quotients of norms started popping up in people's results. The correct generalization is the following notion.

**Definition 2.1.9:** For $p > 2$, define the **stable $p$th rank** by

$$\mathrm{srank}_p(A) = \left(\frac{\|A\|_{S_2}}{\|A\|_{S_p}}\right)^{\frac{2p}{p-2}}.$$

**Exercise 2.1.10:** Show that if $p \geq q > 2$, then

$$\mathrm{srank}_p(A) \leq \mathrm{srank}_q(A).$$

(Hint: Use Hölder's inequality.)

Now we would like to prove how Theorem 2.1.7 generalizes the previously listed results:

*Proof of Generalizability of Theorem 2.1.7.* Using Hölder's inequality with $\frac{p}{2}$,

$$\begin{aligned}
\|A\|_{S_2}^2 &= \sum_{j=1}^m s_j(A)^2 \\
&= \sum_{j=1}^{r-1} s_j(A)^2 + \sum_{j=r}^m s_j(A)^2 \\
&\leq (r-1)^{1-\frac{2}{p}}\left(\sum_{j=1}^{r-1} s_j(A)^p\right)^{\frac{2}{p}} + \sum_{j=r}^m s_j(A)^2
\end{aligned}$$

$$\leq (r-1)^{1-\frac{2}{p}} \|A\|_{S_p}^2 + \sum_{j=r}^{m} s_j(A)^2$$

$$\sum_{j=r}^{m} s_j(A)^2 \geq \|A\|_{S_2}^2 \left(1 - (r-1)^{-\frac{2}{p}} \frac{\|A\|_{S_p}^2}{\|A\|_{S_2}^2}\right)$$

$$= \|A\|_{S_2}^2 \left(1 - \left(\frac{r-1}{\mathrm{srank}_p(A)}\right)^{1-\frac{2}{p}}\right)$$

Now we can plug the previous calculation into Theorem 2.1.7 to demonstrate the way the new theorem generalizes the previous results:

$$\left\|(AJ_\sigma)^{-1}\right\| \lesssim \min_{k+1 \leq r \leq \mathrm{rank}(A)} \sqrt{\frac{mr}{(r-k)\|A\|_{S_2}^2 \left(1 - \left(\frac{r-1}{\mathrm{srank}_p(A)}\right)^{1-\frac{2}{p}}\right)}}$$

$$= \frac{\sqrt{m}}{\|A\|_\infty} \min_{k+1 \leq r \leq \mathrm{rank}(A)} \sqrt{\frac{r}{(r-k)\left(1 - \left(\frac{r-1}{\mathrm{srank}_p(A)}\right)^{1-\frac{2}{p}}\right)}}$$

This equation implies all the earlier theorems. □

To optimize, fix the stable rank, differentiate in $r$, and set to 0. All theorems in the literature follow from this theorem; in particular, we get all the bounds we got before. There was nothing special about the number 4 in Theorem 2.1.8; this is about the distribution of the eigenvalues.

# 2   Ky Fan maximum principle

sec:kf We'll be doing linear algebra. It's mostly mechanical, except we'll need this lemma.

**Lemma 2.2.1** (Ky Fan maximum principle). lem:kf *Suppose that $P : \mathbb{R}^m \to \mathbb{R}^m$ is a rank $k$ orthogonal projection. Then*

$$Tr(A^*AP) \leq \sum_{i=1}^{k} s_i(A)^2$$

*where $s_i(A)$ are the singular values, i.e., $s_i(A)^2$ are the eigenvalues of $B := A^*A$.*

This material was given in class on 2-15. This lemma follows from the following general theorem.

**Theorem 2.2.2** (Convex function of dot products acheives maximum at eigenvectors). thm:eignmax *Let $B : \mathbb{R}^n \to \mathbb{R}^n$ be symmetric, and let $f : \mathbb{R}^n \to \mathbb{R}$ be a convex function. Then for every orthonormal basis $u_1, \ldots, u_n \in \mathbb{R}^n$, there exists a permutation $\pi$ such that*

$$\text{eq:kf0} \quad f(\langle Bu_1, u_1\rangle, \ldots, \langle Bu_n, u_n\rangle) \leq f(\lambda_{\pi(1)}, \ldots, \lambda_{\pi(n)}) \tag{2.1}$$

Essentially, we're saying using the eigenbasis maximizes the convex function.

**Remark 2.2.3:** We can weaken the condition on $f$ to the following: for every $i < j$, $t \mapsto f(x_1, \ldots, x_i + t, x_{i+1}, \ldots, x_{j-1}, x_j - t, x_{j+1}, \ldots, x_n)$ is convex as a function of $t$. If $f$ is smooth, this is equivalent to the second derivative in $t$ being $\geq 0$:

$$\text{eq:kf1}\quad \frac{\partial^2 f}{\partial x_i^2} + \frac{\partial^2 f}{\partial x_j^2} - 2\frac{\partial^2 f}{\partial x_i \partial x_j} \geq 0 \tag{2.2}$$

In other words, you just need that the Hessian is positive semidefinite. (Above, we wrote the determinant of the Hessian on the $x_i - x_j$ plane. This being true for all pairs $i, j$ is the same as the Hessian being positive definite.)

*Proof.* We may assume without loss of generality that

1. $f$ is smooth. If not, convolute with a good kernel.

2. Strict inequality holds:

$$\frac{\partial^2 f}{\partial x_i^2} + \frac{\partial^2 f}{\partial x_j^2} - 2\frac{\partial^2 f}{\partial x_i \partial x_j} > 0.$$

   To see this, note we can take any $\varepsilon > 0$ and perturb $f(x)$ into $f(x) + \varepsilon \|x\|_2^2$. The inequality to prove (2.1) is also perturbed by this slight change, and taking $\varepsilon \to 0$ gives the desired inequality.

Now let $u_1, \ldots, u_n$ be an orthonormal basis at which $f(\langle Bu_1, u_1 \rangle, \ldots, \langle Bu_n, u_n \rangle)$ attains its maximum. Then for $u_i, u_j$, we want to rotate in the $i - j$ plane by angle $\theta$. Since $u_i, u_j$ span a two dimensional subspace, recall the 2-dimensional rotation matrix. Let

$$R_\theta = \begin{bmatrix} \cos(\theta) & \sin(\theta) \\ \sin(\theta) & -\cos(\theta) \end{bmatrix}; u_{i;j} = \begin{bmatrix} u_i \\ u_j \end{bmatrix}$$

Multiplying, we get

$$R_\theta u_{i;j} = \begin{bmatrix} \cos(\theta) & \sin(\theta) \\ \sin(\theta) & -\cos(\theta) \end{bmatrix} \begin{bmatrix} u_i \\ u_j \end{bmatrix} = \begin{bmatrix} \cos(\theta)u_i + \sin(\theta)u_j \\ \sin(\theta)u_i - \cos(\theta)u_j \end{bmatrix} = \begin{bmatrix} (R_\theta u_{i;j})_1 \\ (R_\theta u_{i;j})_2 \end{bmatrix}$$

Then, we replace $f$ with $g(\theta) =$

$$f\left( \langle Bu_1, u_1 \rangle, \ldots, \langle B\left(R_\theta u_{i;j}\right)_1, (R_\theta u_{i;j})_1 \rangle, \langle B\left(R_\theta u_{i;j}\right)_2, (R_\theta u_{i;j})_2 \rangle, \ldots, \langle Bu_n, u_n \rangle \right)$$

where we keep all other dot products the same. By assumption, $g$ attains its maximum at $\theta = 0$, so $g'(0) = 0, g''(0) \leq 0$. Expanding out the rotated dot products explicitly in $g(\theta)$, we get that the $i$th argument is

$$\cos^2(\theta)\langle Bu_i, u_i \rangle + \sin^2(\theta)\langle Bu_j, u_j \rangle + \sin(2\theta)\langle Bu_i, u_j \rangle$$

and the $j$th argument is

$$\sin^2(\theta)\langle Bu_i, u_i\rangle + \cos^2(\theta)\langle Bu_j, u_j\rangle - \sin(2\theta)\langle Bu_i, u_j\rangle$$

Then we can mechanically take the derivatives at $\theta = 0$ to get

$$0 = g'(0) = 2\langle Bu_i, u_j\rangle(f_{x_i} - f_{x_j})$$
$$0 \geq g''(0) = 2\left(\langle Bu_j, u_j\rangle - \langle Bu_i, u_i\rangle\right)\underbrace{(f_{x_i} - f_{x_j})}_{=0 \text{ if } \langle Bu_i, u_j\rangle \neq 0} + 4\langle Bu_i, u_j\rangle^2\underbrace{\left(f_{x_i x_i} + f_{x_j x_j} - 2f_{x_i x_j}\right)}_{\geq 0}.$$

This implies that for all $i \neq j$ $\langle Bu_i, u_j\rangle = 0$, which implies that for all $i$, $Bu_i = \mu_i u_i$ for some $\mu_i$. Thus any function applied to a vector of dot products is maximized at eigenvalues. $\square$

**Exercise 2.2.4:** <sub style="font-size:60%">exr:kf</sub> If $f : \mathbb{R}^n \to \mathbb{R}$ satisfies the conditions in Theorem 2.2.2 and $(u_1, \ldots, u_n)$, $(v_1, \ldots, v_n)$ are two orthonormal bases of $\mathbb{R}^n$, then for every $A : \mathbb{R}^n \to \mathbb{R}^n$, there exists $\pi \in S_n$, $(\varepsilon_1, \ldots, \varepsilon_n) \in \{\pm 1\}^n$ such that

$$f(\langle Au_1, v_1\rangle, \langle Au_2, v_2\rangle, \ldots, \langle Au_n, v_n\rangle) \leq f(\varepsilon_1 s_{\pi(1)}(A), \ldots, \varepsilon_n s_{\pi(n)}(A))$$

Show that choosing $u, v$ as the singular vectors maximizes $f$ (over all pairs of orthonormal bases).

To solve this problem, you can rotate both vectors in the same direction and take derivatives, and also rotate them in opposite directions and take derivatives to get enough information to prove that the singular values are the maximum.

Essentially, a lot of the inequalities you find in books follow from this. For instance, if you want to prove that the Schatten $p$-norm is a norm, it follows directly from this fact.

**Corollary 2.2.5.** *Let $\|\cdot\|$ be a norm on $\mathbb{R}^n$ that is invariant under premutations and sign:*

$$\|(x_1, \ldots, x_n)\| = \|(\varepsilon_1 x_{\pi(1)}, \ldots, \varepsilon_n x_{\pi(n)})\|$$

*for all $\varepsilon \in \{\pm 1\}^n$ and $\pi \in S_n$ (In the literature, we call this a **symmetric norm**). This induces a norm on matrices $M_{m \times n}(\mathbb{R})$ with*

$$\|A\| = \|(s_{\pi(1)}(A), \ldots, s_{\pi(n)}(A)\|)$$

*Then the triangle inequality holds for matrices $A, B$:*

$$\|A + B\| \leq \|A\| + \|B\|.$$

*Proof.* We have by Exercise 2.2.4

$$\|A + B\| = \max_{(u_i)\perp,(v_i)\perp} \|(\langle(A + B)u_i, v_i\rangle)_{i=1}^n\|$$
$$\leq \max_{(u_i)\perp,(v_i)\perp} \|(\langle Au_i, v_i\rangle)_{i=1}^n\| + \max_{(u_i)\perp,(v_i)\perp} \|(\langle Bu_i, v_i\rangle)_{i=1}^n\|$$
$$\leq \|A\| + \|B\|.$$

$\square$

Remember Theorem 2.2.2! For many results, you simply need to apply the right convex function to get the result.

Our lemma follows from setting $f(x) = \sum_{i=1}^{k} x_i$.

*Proof of Ky Fan Maximum Principle (Lemma 2.2.1).* Take an orthonormal basis $u_1, \ldots, u_n$ of $P$ such that $u_1, \ldots, u_k$ is a basis of the range of $P$. Then

$$\text{Tr}(BP) = \sum_{j=1}^{k} \langle Be_j, e_j \rangle \leq \sum_{i=1}^{k} s_i(B) = \sum_{i=1}^{k} s_i(A)^2$$

$\square$

# 3 Finding big subsets

We'll present 4 lemmas for finding big subsets with certain properties. We'll put them together at the end.

## 3.1 Little Grothendieck inequality

**Theorem 2.3.1** (Little Grothendieck inequality). *thm:lgi Fix $k, m, n \in \mathbb{N}$. Suppose that $T : \mathbb{R}^m \to \mathbb{R}^n$ is a linear operator. Then for every $x_1, \ldots, x_k \in \mathbb{R}^m$,*

$$\sum_{r=1}^{k} \|Tx_r\|_2^2 \leq \frac{\pi}{2} \|T\|_{\ell_\infty^m \to \ell_2^n}^2 \sum_{r=1}^{k} x_{ri}^2$$

*for some $i \in \{1, \ldots, m\}$ where $x_r = (x_{r1}, \ldots, x_{rm})$.*

Later we will show $\frac{\pi}{2}$ is sharp.

If we had only 1 vector, what does this say?

$$\|Tx_1\|_2 \leq \sqrt{\frac{\pi}{2}} \|T\|_{\ell_\infty^m \to \ell_2^n} \|x_1\|_\infty$$

We know the inequality is true for $k = 1$ with constant 1, by definition of the operator norm. The theorem is true for arbitrary many vectors, losing an universal constant $(\frac{\pi}{2})$. After we see the proof, the example where $\frac{\pi}{2}$ is attained will be natural.

We give Grothendieck's original proof.

The key claim is the following.

**Claim 2.3.2.** *clm:lgi*

$$eq:lgi1 \quad \sum_{j=1}^{m} \left( \sum_{r=1}^{k} (T^*Tx_r)_j^2 \right)^{\frac{1}{2}} \leq \sqrt{\frac{\pi}{2}} \|T\|_{\ell_\infty^m \to \ell_2^n} \left( \sum_{r=1}^{k} \|Tx_r\|^2 \right)^{\frac{1}{2}}. \tag{2.3}$$

*Proof of Theorem 2.3.1.* Assuming Claim 2.3.2, we prove the theorem.

$$\sum_{r=1}^{k} \|Tx_r\|_2^2 = \sum_{r=1}^{k} \langle Tx_r, Tx_r \rangle$$

$$= \sum_{r=1}^{k} \langle x_r, T^*Tx_r \rangle$$

$$= \sum_{r=1}^{k} \sum_{j=1}^{m} x_{rj}(T^*Tx_r)_j$$

$$\leq \sum_{j=1}^{m} \left( \sum_{r=1}^{k} x_{rj}^2 \right)^{\frac{1}{2}} \left( \sum_{r=1}^{k} (T^*Tx_r)_j^2 \right)^{\frac{1}{2}} \qquad \text{by Cauchy-Schwarz}$$

$$\leq \left( \max_{1 \leq j \leq m} \left( \sum_{r=1}^{k} x_{rj}^2 \right)^{\frac{1}{2}} \right) \left( \sum_{j=1}^{m} \sum_{r=1}^{k} (T^*Tx_r)_j^2 \right)^{\frac{1}{2}}$$

$$\leq \max_{1 \leq j \leq m} \left( \sum_{i=1}^{k} x_{ij}^2 \right)^{\frac{1}{2}} \sqrt{\frac{\pi}{2}} \|T\|_{\ell_\infty^m \to \ell_2^n} \left( \sum_{r=1}^{k} \|Tx_r\|_2^2 \right)^{\frac{1}{2}}$$

$$\sum_{r=1}^{k} \|Tx_r\|_2^2 \leq \frac{\pi}{2} \|T\|_{\ell_\infty^m \to \ell_2^n}^2 \max_j \sum_{r=1}^{k} x_{ij}^2.$$

We bounded by a square root of the multiple of the same term, a bootstrapping argument. In the last step, divide and square. $\qquad\qquad\square$

*Proof of Claim 2.3.2.* Let $g_1, \ldots, g_k$ be iid standard Gaussian random variables. For every fixed $j \in \{1, \ldots, m\}$,

$$\sum_{r=1}^{k} g_r (T^*Tx_r)_j.$$

This is a Gaussian random variable with mean 0 and variance $\sum_{r=1}^{k} (T^*Tx_r)_j^2$. Taking the expectation,[1]

$$\mathbb{E} \left| \sum_{r=1}^{k} g_r (T^*Tx_r)_j \right| = \left( \sum_{r=1}^{k} (T^*Tx_r)_j^2 \right)^{\frac{1}{2}} \sqrt{\frac{2}{\pi}}.$$

Sum these over $j$:

$$\mathbb{E} \left[ \sum_{j=1}^{m} \left| T^* (\sum_{r=1}^{k} g_r Tx_r)_j \right| \right] = \sqrt{\frac{2}{\pi}} \sum_{j=1}^{m} \left( \sum_{r=1}^{k} (T^*Tx_r)_j^2 \right)^{\frac{1}{2}}$$

$$\sum_{j=1}^{m} \left( \sum_{r=1}^{k} (T^*Tx_r)_j^2 \right)^{\frac{1}{2}} = \sqrt{\frac{\pi}{2}} \mathbb{E} \left[ \sum_{j=1}^{m} \left| T^* \sum_{r=1}^{k} g_r (Tx_r)_j \right| \right] \quad \text{.eq:lgi2} \qquad (2.4)$$

---

[1] $\sqrt{\frac{1}{2\pi}} \int_{-\infty}^{\infty} |x| e^{-\frac{x^2}{2}} = -2\sqrt{\frac{1}{2\pi}} [e^{-\frac{x^2}{2}}]_0^\infty = \sqrt{\frac{2}{\pi}}$

Define a random sign vector $z \in \{\pm 1\}^m$ by

$$z_j = \text{sign}\left( \left( T^* \sum_{r=1}^{k} g_r T x_r \right)_j \right)$$

Then

$$\sum_{j=1}^{m} \left| (T^* \sum_{r=1}^{k} g T x_r)_j \right| = \left\langle z, T^* \sum_{r=1}^{k} g_r T x_r \right\rangle$$

$$= \left\langle Tz, \sum_{r=1}^{k} g_r T x_r \right\rangle$$

$$\leq \|Tz\|_2 \left\| \sum_{r=1}^{k} g_r T x_r \right\|_2$$

$$\leq \|T\|_{\ell_\infty^m \to \ell_2^n} \left\| \sum_{r=1}^{k} g_r T x_r \right\|_2$$

This is a pointwise inequality. Taking expectations and using Cauchy-Schwarz,

$$\text{eq:lgi3}\quad \mathbb{E}\left[ \sum_{j=1}^{m} \left| \left( T^* \sum_{r=1}^{k} g_r T x_r \right)_j \right| \right] \leq \|T\|_{\ell_\infty^m \to \ell_2^n} \left( \mathbb{E}\left\| \sum_{r=1}^{k} g_r T x_r \right\|_2^2 \right)^{\frac{1}{2}}. \tag{2.5}$$

What is the second moment? Expand:

$$\text{eq:lgi4}\quad \mathbb{E}\left\| \sum_{r=1}^{k} g_i T x_r \right\|_2^2 = \mathbb{E}\left[ \sum_{ij} g_i g_j \langle T x_i, T x_j \rangle \right] = \sum_{r=1}^{k} \|T x_r\|_2^2. \tag{2.6}$$

Chaining together (2.4), (2.5), (2.6) gives the result. $\qquad\qquad\square$

Why use the Gaussians? The identity characterizes the Gaussians using rotation invariance. Using other random variables gives other constants that are not sharp.

There will be lots of geometric lemmas:

- A fact about restricting matrices.

- Another geometric argument to give a different method for selecting subsets.

- A combinatorial lemma for selecting subsets.

Finally we'll put them together in a crazy induction.

2-10-16: We were in the process of proving three or four subset selection principles, which we will somehow use to prove the RIP.

The little Grothendieck inequality (Theorem 2.3.1) is part of an amazing area of mathematics with many applications. It's little, but very useful. The proof is really Grothendieck's original proof, re-organized. For completeness, we'll show the fact that the inequality is sharp (cannot be improved).

### 3.1.1   Tightness of Grothendieck's inequality

**Corollary 2.3.3.** $\sqrt{\pi/2}$ *is the best constant in Theorem 2.3.1.*

From the proof, we reverse engineer vectors that make the inequality sharp. They are given in the following example.

**Example 2.3.4:** Let $g_1, g_2, \ldots, g_k$ be iid Gaussians on the probability space $(\Omega, P)$. Let $T : L_\infty(\Omega, P) \to \ell_2^k$ be

$$Tf = (\mathbb{E}[fg_1], \ldots, \mathbb{E}[fg_k]).$$

Let $x_r \in L_\infty(\Omega, P)$,

$$x_r = \frac{g_r}{\left(\sum_{i=1}^k g_i^2\right)^{\frac{1}{2}}}.$$

*Proof.* Le $g_1, \ldots, g_k; T; x_1, \ldots x_k$ be as in the example. Note the $x_r$ are nothing more than vectors on the $k$-dimensional unit sphere, so they are bounded functions on the measure space $\Omega$. We can also write

$$\sum_{r=1}^k x_r(\omega)^2 = \sum_{r=1}^k \frac{g_r(\omega)^2}{\sum_{i=1}^r g_i(\omega)^2} = 1 \tag{2.7}$$

We use the Central Limit Theorem in order to replace the $\Omega$ by a discrete space. Let $\varepsilon_{r,i}$ be $\pm 1$ random variables. Then letting

$$g_r = \frac{\varepsilon_{r,1} + \ldots + \varepsilon_{r,N}}{\sqrt{N}}$$

instead, we have that $g_r$ approaches a standard Gaussian in distribution, so the statements we make will be asymptotically true. With this discretization, the random variables $\{g_r\}$ live in $\Omega = \{\pm 1\}^{NK}$. So $L_\infty(\Omega) = l_\infty^{2^{NK}}$, which is in a large but finite dimension. So $\omega$ will really be a coordinate in $\Omega$.

Now we show two things; they are nothing more than computations.

1. $\|T\|_{L_\infty(\Omega, \mathbb{P}) \to l_2^k} = \sqrt{2/\pi}$,

2. We also show $\sum_{r=1}^k \|Tx_r\|_2^2 \xrightarrow{k \to \infty} 1$.

From (2.7) and the 2 items, the little Grothendieck inequality is sharp in the limit.

For (1), we have

$$
\begin{aligned}
\|T\|_{\ell^\infty \to \ell^2} &= \sup_{\|f\|_\infty \leq 1} \left( \sum_{r=1}^{k} \mathbb{E}\left[ fg_r \right]^2 \right)^{1/2} \\
&= \sup_{\|f\|_\infty \leq 1} \sup_{\sum_{r=1}^{k} \alpha_r^2 = 1} \sum_{r=1} \alpha_r \mathbb{E}\left[ fg_r \right] \\
&= \sup_{\sum_{r=1}^{k} \alpha_r^2 = 1} \sup_{\|f\|_\infty \leq 1} \mathbb{E}\left[ f \sum_{i=1}^{k} \alpha_r g_r \right] \\
&= \sup_{\sum_{r=1}^{k}} \mathbb{E}\left| \sum_{r=1}^{k} \alpha_r g_r \right| = \mathbb{E}|g_1| = \sqrt{\frac{2}{\pi}}
\end{aligned}
\tag{2.8}
$$

as we claimed, since $\|\alpha\|_2 = 1$ implies $\sum_{r=1}^{k} \alpha_r g_r$ is also a gaussian.

Now for (2),

$$
\begin{aligned}
\sum_{r=1}^{k} \|Tx_r\|_2^2 &= \sum_{r=1}^{k} \left( \mathbb{E}\left[ \frac{g_r^2}{\left( \sum_{i=1}^{k} g_i^2 \right)^{1/2}} \right] \right)^2 \\
&= K \left( \mathbb{E}\left[ \frac{g_1^2}{\left( \sum_{i=1}^{k} g_i^2 \right)^{1/2}} \right] \right)^2 \\
&= K \left( \frac{1}{K} \mathbb{E}\left[ \sum_{r=1}^{k} \frac{g_r^2}{\left( \sum_{i=1}^{k} g_i^2 \right)^{1/2}} \right] \right)^2 \\
&= \frac{1}{K} \left( \mathbb{E}\left[ \left( \sum_{i=1}^{k} g_i^2 \right)^{1/2} \right] \right)^2
\end{aligned}
\tag{2.9}
$$

and you can use Stirling to finish. This is just a $\chi^2$-distribution.

In this case $\mathbb{E}\frac{g_1 g_2}{\left( \sum_i g_i^2 \right)^{1/2}} = \mathbb{E}\frac{g_1(-g_2)}{\left( \sum_i g_i^2 \right)^{1/2}}$. Also note that if $(g_1, \ldots, g_k) \in \mathbb{R}^k$ is a standard Gaussian, then $\frac{(g_1, \ldots, g_k)}{\left( \sum_{i=1}^{k} g_i^2 \right)^{1/2}}$ and $\left( \sum_{i=1}^{k} g_i^2 \right)^{1/2}$ are independent. In other words, the length and angle are independent: This is just polar coordinates, you can check this. $\qquad\square$

Now, how does this relate to the Restricted Invertibility Problem?

## 3.2    Pietsch Domination Theorem

**Theorem 2.3.5** (Pietsch Domination Theorem). *thm:pdt Fix $m, n \in \mathbb{N}$ and $M > 0$. Suppose that $T : \mathbb{R}^m \to \mathbb{R}^n$ is a linear operator such that for every $x_1, \ldots, x_k \in \mathbb{R}^m$ have*

$$
{\scriptstyle eq:pdt1} \left( \sum_{r=1}^{k} \|Tx_r\|_2^2 \right)^{1/2} \leq M \max_{1 \leq j \leq m} \left( \sum_{r=1}^{k} x_{rj}^2 \right)^{1/2}
\tag{2.10}
$$

*Then there exist $\mu = (\mu_1, \ldots, \mu_m) \in \mathbb{R}^m$ with $\mu_1 \geq 0$ and $\sum_{i=1}^m \mu_i = 1$ such that for every $x \in \mathbb{R}^m$*

$$\textcolor{red}{\scriptstyle eq:pdt2}\|Tx\|_2 \leq M \left( \sum_{i=1}^M \mu_i \|x_i\|^2 \right)^{1/2} \tag{2.11}$$

The theorem says that you can come up with a probability measure such that the norm of T as an operator as a standard norm from $l_\infty$ to $l_2$ (?), is bounded by $M$.

**Remark 2.3.6:** The theorem really an iff: (2.11) is a stronger statement than (2.10), and in fact they are equivalent.

*Proof.* Define $K \subseteq \mathbb{R}^m$ with

$$K = \left\{ y \in \mathbb{R}^m : y_i = \sum_{r=1}^k \|Tx_r\|_2^2 - M^2 \sum_{r=1}^m x_{ri}^2 \text{ for some } k, x_1, \ldots, x_k \in \mathbb{R}^m \right\}$$

Basically we cleverly select a convex set. Every $n$-tuple of vectors in $\mathbb{R}^m$ gives you a new vector in $\mathbb{R}^m$. Let's check that $K$ is convex. We have to check if two vectors $y, z \in K$, then all points on the line between them are in $K$. $y, z \in K$ means that there exist $(x_i)_{i=1}^k$, $(w_i)_{i=1}^l$,

$$y_i = \sum_{r=1}^k \|Tx_r\|_2^2 - M^2 \sum_{r=1}^m x_{ri}^2$$

$$z_i = \sum_{r=1}^l \|Tw_r\|_2^2 - M^2 \sum_{r=1}^l w_{ri}^2$$

for all $i$. Then $\alpha y_i + (1-\alpha)z_i$ comes from $(\sqrt{\alpha}x_1, \ldots, \sqrt{\alpha}x_k, \sqrt{1-\alpha}w_1, \ldots \sqrt{1-\alpha}w_k)$. So by design, $K$ is a convex set.

Now, the assumption of the theorem says that

$$\left( \sum_{r=1}^k \|Tx_r\|_2^2 \right)^{1/2} \leq M\max_{1 \leq j \leq m} \left( \sum_{r=1}^k x_{rj}^2 \right)^{1/2}$$

which implies

$$\|Tx_r\|_2^2 - M^2\max_{1 \leq j \leq m} \sum_{r=1}^m x_{rj}^2 \leq 0$$

which implies $K \cap (0, \infty)^m = \emptyset$. By the hyperplane separation theorem (for two disjoint convex sets in $\mathbb{R}^m$ with at least one compact, there is a hyperplane between them), there exists $0 \neq \mu = (\mu_1, \ldots, \mu_m) \in \mathbb{R}^m$ with

$$\langle \mu, y \rangle \leq \langle \mu, z \rangle$$

for all $y \in K$ and $z \in (0, \infty)^m$. By renormalizing, we may assume $\sum_{i=1}^m \mu_i = 1$. Moreover $\mu$ cannot have any strictly negative coordinate: Otherwise you could take $z$ to have arbitrarily

large value at a strictly negative coordinate with zeros everywhere else, implying $\langle u, z \rangle$ is no longer bounded from below, a contradiction. Therefore, $\mu$ is a probability vector and $\langle \mu, z \rangle$ can be arbitrarily small. So for every $y \in K$, $\sum_{i=1}^{m} \mu_i y_i \leq 0$. Write

$$y_i = \|Tx\|_2^2 - M^2 \|x_i\|^2 \in K.$$

Expanding this out,

$$\|Tx\|_2^2 - M^2 \sum_{i=1}^{n} \mu_i \|x_i\|^2 \leq 0,$$

which is exactly what we wanted. $\hfill\square$

## 3.3    A projection bound

**Lemma 2.3.7.** <span style="font-size:smaller">lem:projbound</span> $m, n \in \mathbb{N}$, $\varepsilon \in (0,1)$, $T : \mathbb{R}^n \to \mathbb{R}^m$ a linear operator. Then $\exists \sigma \subset \{1, \ldots, m\}$ with $|\sigma| \geq (1 - \varepsilon)m$ such that

$$\|Proj_{\mathbb{R}^\sigma} T\|_{S_\infty} \leq \sqrt{\frac{\pi}{2\varepsilon m}} \|T\|_{l_2^n \to l_1^m}$$

We will find ways to restrict a matrix to a big submatrix. We won't be able to control its operator norm, but we will be able to control the norm from $l_2^n$ to $l_1^m$. Then we pass to a further subset, which this becomes an operator norm on, which is an improvement which Grothendieck gave us. This is the first very useful tool to start finding big submatrices.

*Proof.* We have $T : l_2^n \to l_1^m$, $T^* : l_\infty^m \to l_2^n$. Now some abstract nonsense gives us that for Banach spaces, the norm of an operator and its adjoint are equal, i.e. $\|T\|_{l_2^n \to l_1^m} = \|T^*\|_{l_\infty^m \to l_2^n}$. This statement follows from the Hahn-Banach theorem (come see me if you haven't seen this before, I'll tell you what book to read). From the Little Grothendieck inequality (Theorem 2.3.1), $T^*$ satisfies the assumption of the Pietsch domination theorem 2.3.5 with $M = \sqrt{\frac{\pi}{2}}\|T\|_{l_2^n \to l_1^m}$ (we're applying it to $T^*$). By the theorem, there exists a probability vector $(\mu_1, \ldots, \mu_m)$ such that for every $y \in \mathbb{R}^m$,

$$\|T^* y\|_2 = M \left( \sum_{i=1}^{m} \mu_i y_i^2 \right)^{1/2}$$

with $M = \sqrt{\frac{\pi}{2}}\|T\|_{l_2^n \to l_1^m}$. Define $\sigma = \left\{ i \in \{1, \ldots, m\} : \mu_i \leq \frac{1}{m\varepsilon} \right\}$; then $|\sigma| \geq (1 - \varepsilon)m$ by Markov's inequality. We can also see this by writing

$$1 = \sum_{i=1}^{m} \mu_i = \sum_{i \in \sigma} \mu_i + \sum_{i \notin \sigma} \mu_i > \sum_{i \in \sigma} \mu_i + \frac{m - |\sigma|}{m\varepsilon}$$

which follows since for $j \notin \sigma$, $\mu_j > \frac{1}{m\varepsilon}$. Continuing,

$$\frac{m\varepsilon - m + |\sigma|}{m\varepsilon} \geq \sum_{i \in \sigma} \mu_i$$

$$|\sigma| \geq (m\varepsilon) \sum_{i \in \sigma} \mu_i + m(1 - \varepsilon)$$

Then, because $\mu$ is a probability distribution, $(m\varepsilon)\sum_{i\in\sigma}\mu_i \geq 0$ and we have

$$|\sigma| \geq m(1-\varepsilon)$$

Now take $x \in \mathbb{R}^n$ and choose $y \in \mathbb{R}^m$ with $\|y\|_2 = 1$. Then

$$\langle y, \mathrm{Proj}_{\mathbb{R}^\sigma} Tx\rangle^2 = \langle T^*\mathrm{Proj}_{\mathbb{R}^\sigma} y, x\rangle^2 \leq \|T^*\mathrm{Proj}_{\mathbb{R}^\sigma} y\|_2^2 \cdot \|x\|_2^2$$

$$\leq \frac{\pi}{2}\|T\|_{l_2^n \to l_1^m}\left(\sum_{i\in\sigma}\mu_i y_i^2\right)\|x\|_2^2 \leq \frac{\pi}{2}\|T\|_{l_2^n \to l_1^m}^2 \frac{1}{m\varepsilon}\|x\|_2^2$$

by Cauchy-Schwarz. Taking square roots gives the desired result.                    □

In the previous proof, we used a lot of duality to get an interesting subset.

**Remark 2.3.8:** In Lemma 2.3.7, I think that either the constant $\pi/2$ is sharp (no subset are bigger; it could come from the Gaussians), or there is a different constant here. If the constant is 1, I think you can optimize the previous argument and get the constant to be arbitrarily close to 1, which would have some nice applications: In other words, getting $\sqrt{\frac{\pi}{2\varepsilon m}}$ as close to 1 as possible would be good. I didn't check before class, but you might want to check if you can carry out this argument using the Gaussian argument we made for the sharpness of $\frac{\pi}{2}$ in Grothendieck's inequality (Theorem 2.3.1). It's also possible that there is a different universal constant.

## 3.4   Sauer-Shelah Lemma

Now we will give another lemma which is very easy and which we will use a lot.

**Lemma 2.3.9** (Sauer-Shelah). *lem:saushel Take integers $m, n \in \mathbb{N}$ and suppose that we have a large set $\Omega \subseteq \{\pm 1\}^n$ with*

$$|\Omega| > \sum_{k=0}^{m-1}\binom{n}{k}$$

*Then $\exists \sigma \subseteq \{1,\ldots,n\}$ such that with $|\sigma| = m$, if you project onto $\mathbb{R}^\sigma$ the set of vectors, you get the entire cube: $\mathrm{Proj}_{\mathbb{R}^\sigma}(\Omega) = \{\pm 1\}^\sigma$.[2] For every $\varepsilon \in \{\pm 1\}^\sigma$, there exists $\delta = (\delta_1,\ldots,\delta_n) \in \Omega$ such that $\delta_j = \varepsilon_j$ for $j \in \sigma$.*

Note that Lemma 2.3.9 is used in the proof of the Fundamental Theorem of Statistical Learning Theory.

*Proof.* We want to prove by induction on $n$. First denote the **shattering set**

$$\mathrm{sh}(\Omega) = \{\sigma \subseteq \{1,\ldots,n\} : \mathrm{Proj}_{\mathbb{R}^\sigma}\Omega = \{\pm 1\}^\sigma\}$$

---

[2]I.e., the **VC dimension** of $\Omega$ is $\geq m$.

The claim is that the number of sets shattered by a given set is $|\text{sh}(\Omega)| \geq |\Omega|$. The empty set case is trivial. What happens when $n = 1$? $\Omega \subset \{-1, 1\}$, and thus the set is shattered. Assume that our claim holds for $n$, and now set $\Omega \subseteq \{\pm 1\}^{n+1} = \{\pm 1\}^n \times \{\pm 1\}$. Define

$$\Omega_+ = \{\omega \in \{\pm 1\}^n : (\omega, 1) \in \Omega\}$$

$$\Omega_- = \{\omega \in \{\pm 1\}^n : (\omega, -1) \in \Omega\}$$

Then, letting $\tilde{\Omega}_+ = \{(\omega, 1) \in \{\pm 1\}^{n+1} : \omega \in \Omega_+\}$ and $\tilde{\Omega}_-$ similarly, we have $|\Omega| = |\tilde{\Omega}_+| + |\tilde{\Omega}_-| = |\Omega_+| + |\Omega_-|$. By our inductive step, we have $\text{sh}(\Omega_+) \geq |\Omega_+|$ and $\text{sh}(\Omega_-) \geq |\Omega_-|$. Note that any subset that shatters $\Omega_+$ also shatters $\Omega$, and likewise for $\Omega_-$. Note that if a set $\Omega'$ shatters both of them, we are allowed to add on an extra coordinate to get $\Omega' \times \{\pm 1\}$ which shatters $\Omega$. Therefore,

$$\text{sh}(\Omega_+) \cup \text{sh}(\Omega_-) \cup \{\sigma \cup \{n+1\} : \sigma \in \text{sh}(\Omega_+) \cap \text{sh}(\Omega_-)\} \subseteq \text{sh}(\Omega)$$

where the last union is disjoint since the dimensions are different. Therefore, we can now use this set inclusion to complete the induction using the principle of inclusion-exclusion:

$$\begin{aligned}
|\text{sh}(\Omega)| &\geq |\text{sh}(\Omega_+) \cup \text{sh}(\Omega_-)| + |\text{sh}(\Omega_+) \cap \text{sh}(\Omega_-)| && \text{(disjoint sets)} \\
&= |\text{sh}(\Omega_+)| + |\text{sh}(\Omega_-)| - |\text{sh}(\Omega_+) \cap \text{sh}(\Omega_-)| + |\text{sh}(\Omega_+) \cap \text{sh}(\Omega_-)| \\
&= |\text{sh}(\Omega_+)| + |\text{sh}(\Omega_-)| \\
&\geq |\Omega_+| + |\Omega_-| = |\Omega|
\end{aligned}$$

which completes the induction as desired. $\qquad\qquad\square$

We will primarily use the theorem as the following corollary, which says that if you have half of the points in terms of cardinality, you get half of the dimension.

**Corollary 2.3.10.** *If* $|\Omega| \geq 2^{n-1}$ *then there exists* $\sigma \subseteq \{1, \dots, n\}$ *with* $|\sigma| \geq \lceil \frac{n+1}{2} \rceil \geq \frac{n}{2}$ *such that* $\text{Proj}_{\mathbb{R}^\sigma}\Omega = \{\pm 1\}^\sigma$.

2-15-16: Last time we left off with the proof of the Sauer-Shelah lemma. To remind you, we were finding ways to find interesting subsets where matrices behave well. Now recall we had a linear algebraic fact which I owe you; I will prove it in an analytic way. The proof has been moved to Section 2.

# 4   Proof of RIP

## 4.1   Step 1

Now we need another geometric lemma for the proof of Theorem 2.1.7, the restricted invertibility principle.

**Lemma 2.4.1** (Step 1). *lem:step1 Fix $m, n, r \in \mathbb{N}$. Let $A : \mathbb{R}^m \to \mathbb{R}^n$ be a linear operator with rank$(A) \geq r$. For every $\tau \subseteq \{1, \ldots, m\}$, denote*

$$E_\tau = (\text{span}((Ae_j)_{j \in \tau}))^\perp.$$

*Then there exists $\tau \subseteq \{1, \ldots, m\}$ with $|\tau| = r$ such that for all $j \in \tau$,*

$$\|\text{Proj}_{E_{\tau \setminus \{j\}}} Ae_j\|_2 \geq \frac{1}{\sqrt{m}} \left( \sum_{i=1}^m s_i(A)^2 \right)^{1/2}.$$

Basically we're taking the projection of the $j^{th}$ column onto the orthogonal completement of the span of the subspace of all columns in the set except for the $j^{th}$ one, and bounding the norm of that by a dimension term and the square root of the sum of the eigenvalues. (This is sharp asymptotically, and may in fact even be sharp as written too—I need to check. Check this?)
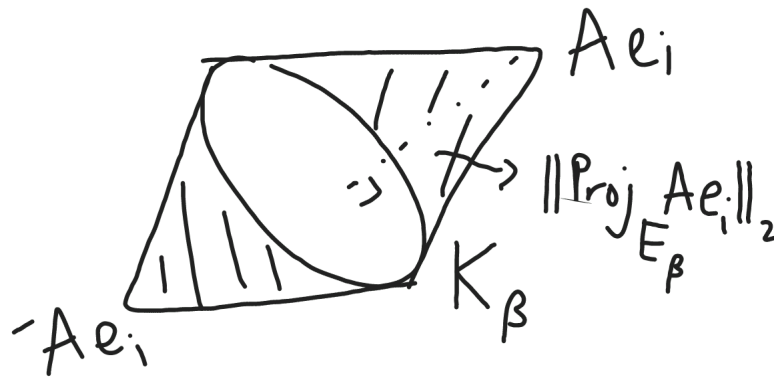
*Proof.* For every $\tau \subseteq \{1, \ldots, m\}$, denote

$$K_\tau = \text{conv}\left( \{\pm Ae_j\}_{j \in \tau} \right)$$

The idea is to make the convex hull have big volume. Once we do that, wewill get all these inequalities for free. Let $\tau \subseteq \{1, \ldots, m\}$ be the subset of size $r$ that maximizes $\text{vol}_r(K_\tau)$. We know that $\text{vol}_r(K_\tau) > 0$. Observe that for any $\beta \subseteq \{1, \ldots, m\}$ of size $r-1$ and $i \notin \beta$, we have

$$K_{\beta \cup \{i\}} = \text{conv}\left( K_\beta \cup \{\pm Ae_i\} \right),$$

which is a double cone.



What is the height of this cone? It is $\|\text{Proj}_{E_\beta} Ae_i\|_2$, as $E_\beta$ is the orthogonal complement of the space spanned by $\beta$. Therefore, the $r$-dimensional volume is given by

$$\text{vol}_r(K_{\beta \cup \{i\}}) = 2 \cdot \frac{\text{vol}_{r-1}(K_\beta) \cdot \|\text{Proj}_{E_\beta} Ae_i\|_2}{r}$$

Because $|\tau| = r$ is the maximizing subset of $F_\Omega$, for any $j \in \tau$ and $i \in \{1, \ldots, m\}$, choosing $\beta = \tau \setminus \{j\}$, we get

$$\mathrm{vol}_r(K_{\beta \cup \{j\}}) \geq \mathrm{vol}_r(K_{\beta \cup \{i\}})$$
$$\implies \|\mathrm{Proj}_{E_{\tau \setminus \{j\}}} A e_j\|_2 \geq \|\mathrm{Proj}_{E_{\tau \setminus \{j\}}} A e_i\|_2.$$

for every $j \in \tau$ and $i \in \{1, \ldots, m\}$. Summing,

$$m \|\mathrm{Proj}_{E_{\tau \setminus \{j\}}} A e_j\|_2^2 \geq \sum_{i=1}^{m} \|\mathrm{Proj}_{E_{\tau \setminus \{j\}}} A e_i\|_2^2 = \|\mathrm{Proj}_{E_{\tau \setminus \{j\}}} A\|_{S_2}^2.$$

Then, for all $j \in \tau$,

$$\text{\small eq:rip-s1-1} \|\mathrm{Proj}_{E_{\tau \setminus \{j\}}} A e_j\|_2 \geq \frac{1}{\sqrt{m}} \|\mathrm{Proj}_{E_{\tau \setminus \{j\}}} A\|_{S_2} \tag{2.12}$$

Let's denote $P = \mathrm{Proj}_{E_{\tau \setminus \{j\}}}$. Note $P$ is an orthogonal projection of rank $r - 1$. Then,

$$\|PA\|_{S_2}^2 = \mathrm{Tr}((PA)^*(PA)) = \mathrm{Tr}(A^* P^* P A) = \mathrm{Tr}(A^* P A) = \mathrm{Tr}(AA^* P)$$

$$\text{\small eq:rip-s1-2} \qquad = \mathrm{Tr}(AA^*) - \mathrm{Tr}(AA^*(I - P)) \geq \sum_{i=1}^{m} s_i(A)^2 - \sum_{i=1}^{r-1} s_i(A)^2 = \sum_{i=r}^{m} s_i(A)^2 \tag{2.13}$$

using the Ky Fan maximal principle 2.2.1, since $I - P$ is a projection of rank $m - r + 1$.

Putting (2.12) and (2.13) together gives the result. $\qquad\qquad\qquad\qquad \square$

## 4.2   Step 2

<span style="color:blue">In our proof of the restricted invertibility principle, this is the first step. Before proving it, let me just tell you what the second step looks like.</span>

**Lemma 2.4.2** (Step 2). _lem:step2 Let $k, m, n \in \mathbb{N}$, $A : \mathbb{R}^m \to \mathbb{R}^n$, $\mathrm{rank}(A) > k$. Let $\omega \subseteq \{1, \ldots, m\}$ with $|\omega| = \mathrm{rank}(A)$ such that $\{A e_j\}_{j \in \omega}$ are linearly independent. Denote for every $j \in \Omega$_

$$F_j = E_{\omega \setminus \{j\}} = \big(\mathrm{span}(A e_i)_{i \in \omega \setminus \{j\}}\big).$$

_Then there exists $\sigma \subseteq \omega$ with $|\sigma| = k$ such that_

$$\|(A J_\sigma)^{-1}\|_{S_\infty} \lesssim \frac{\sqrt{\mathrm{rank}(A)}}{\sqrt{\mathrm{rank}(A) - k}} \cdot \max_{j \in \omega} \frac{1}{\|\mathrm{Proj}_{F_j} A e_j\|}$$

Most of the work is in the second step. First we pass to a subset where we have some information about the shortest possible orthogonal project. But Step 1 saves us by bounding what this can be. Here we use the Grothendieck inequality, Sauer-Shelah, etc. Everything: It's simple, but it kills the restricted invertibility principle.

*Proof of Theorem 2.1.7 given Step 1 and 2.* Take $A : \mathbb{R}^m \to \mathbb{R}^n$. By Step 1 (Lemma 2.4.1), we can find subset $\tau \subseteq \{1, \ldots, m\}$ with $|\tau| = r$ such that for all $j \in \tau$,

$$\text{eq:step1} \|\text{Proj}_{E_{\tau \setminus \{j\}}} A e_j\|_2 \geq \frac{1}{\sqrt{m}} \left( \sum_{i=r}^{m} s_i(A)^2 \right)^{1/2}. \tag{2.14}$$

Now we apply Step 2 (Lemma 2.4.2) to $AJ_\tau$, using $\omega = \tau$, and find a further subset $\sigma \subseteq \tau$ such that

$$\| (AJ_\sigma)^{-1} \|_{S_\infty} \leq \min_{k < r < \text{rank}(A)} \sqrt{\frac{\text{rank}(A)}{\text{rank}(A) - r}} \max_{j \in \omega} \frac{1}{\left\| \text{Proj}_{F_j} A e_j \right\|}$$

$$\leq \min_{k < r < \text{rank}(A)} \sqrt{\frac{mr}{(r-k) \sum_{i=r}^{m} s_i(A)^2}}$$

which we get by plugging directly in $r$ for the rank and using Step 1 (2.14) to get the denominator. $\qquad \square$

2-17

Now we prove Step 2 (Lemma 2.4.2). Note we can assume $\omega = \{1, \ldots, m\}$ and that the rank is $m$.

First we need some lemmas.

**Lemma 2.4.3.** *lem:rip-step2-1 Let $A : \mathbb{R}^m \to \mathbb{R}^n$ be such that $\{Ae_j\}_{j=1}^m$ are linearly independent, and $\sigma \subseteq \{1, \ldots, m\}$, $t \in \mathbb{N}$. Then there exists $\tau \subseteq \sigma$ with*

$$|\tau| \geq \left( 1 - \frac{1}{2^t} \right) |\sigma|,$$

*such that denoting $\theta = \tau \cup (\{1, \ldots, m\} \setminus \sigma)$, $M = \max_{j \in \omega} \frac{1}{\|\text{Proj}_{F_j} A e_j\|}$, for all $\sigma \in \mathbb{R}^\theta$,*

$$\sum_{i \in \tau} |a_i| \leq 2^{\frac{t}{2}} M \sqrt{|\sigma|} \left\| \sum_{i \in \theta} a_i A e_i \right\|_2.$$
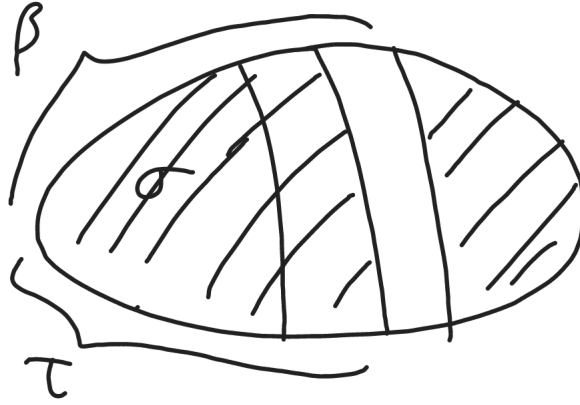
This is proved by a nice inductive argument.

*Proof.* TODO next time

$\square$

**Lemma 2.4.4.** <sub>lem:rip-step2-2</sub> *Let* $m, n, t \in \mathbb{N}$ *and* $\beta \subseteq \{1, \ldots, m\}$. *Let* $A : \mathbb{R}^m \to \mathbb{R}^n$ *be a linear operator such that* $\{Ae_j\}_{j=1}^m$ *are linearly independent. Then there exist two subsets* $\sigma \subseteq \tau \subseteq \beta$ *such that* $|\tau| \geq \left(1 - \frac{1}{2^t}\right)|\beta|$, $|\tau \backslash \sigma| \leq \frac{|\beta|}{4}$, *and if we denote* $\theta = \tau \cup (\{1, \ldots, m\} \backslash \beta)$, $M = \max_{j \in \omega} \frac{1}{\|\mathrm{Proj}_{F_j} Ae_j\|}$, *then*

$$\left\|\mathrm{Proj}_{\mathbb{R}^\sigma}(AJ_\theta)^{-1}\right\|_{S_\infty} \lesssim 2^{\frac{t}{2}} M.$$



*Proof of Lemma 2.4.4 from Lemma 2.4.3.* Apply Lemma 2.4.3 with $\sigma = \beta$. Basically, we're going to inductively construct a subset $\sigma = \beta$ of $\{1, \cdots, m\}$ onto which to project, so that we can control the $l_1$ norm of the subset $\tau$ by the $l_2$ norm on a slightly greater set $\theta$ via Lemma 2.4.3.

We find $\tau \subseteq \beta$ with $|\tau| \geq \left(1 - \frac{1}{2^t}\right)|\beta|$ such that

$$\sum_{i \in \tau} |a_i| \leq 2^{\frac{t}{2}} M \sqrt{|\beta|} \left\|\sum_{i \in \theta} a_i Ae_i\right\|_2.$$

Rewriting gives that

$$\forall a \in \mathbb{R}^\sigma, \qquad \|\mathrm{Proj}_{\mathbb{R}^\tau} a\| \leq 2^{\frac{t}{2}} M \sqrt{|\beta|} \|AJ_\theta a\|_2$$
$$\implies \left\|\mathrm{Proj}_{\mathbb{R}^\tau}(AJ_\theta)^{-1}\right\|_{\ell_2^\theta \to \ell_1^\tau} \lesssim 2^{\frac{t}{2}} M \sqrt{|\beta|}.$$

Denote $\varepsilon = \frac{|\beta|}{4|\tau|}$. By Lemma 2.3.7, there exists $\sigma \subseteq \tau$, $|\sigma| \geq (1 - \varepsilon)|\tau|$ such that

$$\left\|\mathrm{Proj}_{\mathbb{R}^\sigma}(AJ_\theta)^{-1}\right\|_{S^\infty} = \left\|\mathrm{Proj}_{\mathbb{R}^\sigma} \mathrm{Proj}_{\mathbb{R}^\tau}(AJ_\theta)^{-1}\right\|_{S^\infty} \lesssim \sqrt{\frac{\pi}{2\varepsilon|\tau|}} 2^{\frac{t}{2}} M \sqrt{|\beta|}$$
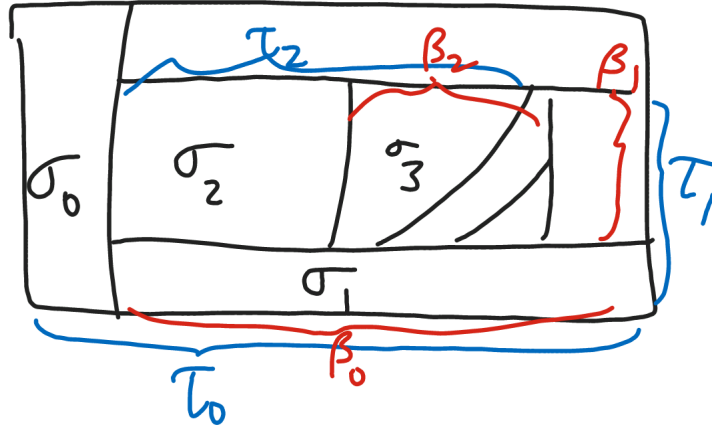$$\lesssim 2^{\frac{t}{2}} M.$$

$\square$

Now, we have basically finished making use of Lemma 2.4.3 in the proof of Lemma 2.4.4. It remains to use Lemma 2.4.4 to prove the full second step in our proof of the Restricted Invertibility Principle.

*Proof of Lemma 2.4.2 (Step 2).* Fix an integer $r$ such that $\frac{1}{2^{2r+1}} \leq 1 - \frac{k}{m} \leq \frac{1}{2^{2r}}$. Proceed inductively as follows. First set

$$\tau_0 = \{1, \ldots, m\}$$
$$\sigma_0 = \phi.$$

Suppose $u \in \{0, \ldots, r+1\}$ and we constructed $\sigma_k, \tau_k \subseteq \{1, \ldots, m\}$ such that if we denote $\beta_u = \tau_u \backslash \sigma_u$, $\theta_u = \tau_u \cup (\{1, \ldots, m\} \backslash \beta_{u-1})$, then

1. $\sigma_u \subseteq \tau_u \subseteq \beta_{u-1}$

2. $|\tau_u| \geq \left(1 - \frac{1}{2^{2r-u+4}}\right) |\beta_{u-1}|$

3. $|\beta_u| \leq \frac{1}{4} |\beta_{u-1}|$

4. $\left\| \text{Proj}_{\mathbb{R}^{\sigma_u}} (A J_{\theta_u})^{-1} \right\|_{S_\infty} \lesssim 2^{r - \frac{u}{2}} M.$



Let $H = 2r - u + 4$. For instance, $|\tau_1| \geq \left(1 - \frac{1}{2^{2r+3}}\right) |\beta_0|$. What is the new $\beta$?

For the inductive step, apply Lemma 2.4.4 on $\beta_{u-1}$ with $t = 2r - u + 4$ to get $\sigma_u \subseteq \tau_u \subseteq \beta_{u-1}$ such that $|\tau_u| \geq \left(1 - \frac{1}{2^{2r-u+4}}\right) |\beta_{u-1}|$, $|\tau_u \backslash \sigma_u| \leq \frac{|\beta_{u-1}|}{4}$ As we induct, we are essentially building up more $\sigma_u$ to eventually produce the invertible subset over which we will project. Note that the size of the $\theta$ set is decreasing as we proceed inductively.

$$\text{eq:rip-s2-1} \quad \left\| \text{Proj}_{\mathbb{R}^{\sigma_u}} (A J_{\theta_u})^{-1} \right\| \lesssim 2^{r-u/2} M. \tag{2.15}$$

We know $|\beta_{u-1}| \leq \frac{m}{4^{u-1}}$,

$$\beta_{u-1} = \beta_u \sqcup \sigma_u \sqcup (\beta_{u-1} \backslash \tau_u)$$
$$|\beta_{u-1}| = |\beta_u| + |\sigma_u| + (|\beta_{u-1}| - |\tau_u|)$$
$$|\sigma_u| = |\beta_{u-1}| - |\beta_u| - (|\beta_{u-1}| - |\tau_u|)$$
$$\geq |\beta_{u-1}| - |\beta_u| - \frac{|\beta_{u-1}|}{2^{2r-u+4}}$$
$$\geq |\beta_{u-1}| - |\beta_u| - \frac{m}{2^{2r+u+2}}.$$

Our choice for the invertible subset is

$$\sigma = \bigsqcup_{u=1}^{r+1} \sigma_u.$$

Telescoping gives

$$|\sigma| = \sum_{u=1}^{r+1} |\sigma_u| \geq |\beta_0| - |\beta_{r+1}| - \frac{m}{2^{2r+2}} \sum_{u=1}^{\infty} \frac{1}{2^u}$$
$$\geq m - \frac{m}{4^{r+1}} - \frac{m}{2^{2r+2}}$$
$$= m\left(1 - \frac{1}{2^{2r+1}}\right) \geq m\frac{k}{m} = k.$$

Observe that $\sigma \subseteq \bigcap_{u=1}^{r+1} \theta_u$ and for every $u$,

$$\sigma_u, \ldots, \sigma_{r+1} \subseteq \tau_u$$
$$\sigma_1, \ldots, \sigma_{u-1} \subseteq \{1, \ldots, m\} \backslash \beta_{u-1}$$

This allows us to use the conclusion for all the $\sigma_u$'s at once.
For $a \in \mathbb{R}^\sigma$, $J_\sigma a \subseteq J_{\theta_u} \mathbb{R}^{\theta_u}$,

$$\text{Proj}_{\mathbb{R}^{\sigma_u}} (AJ_{\theta_u})^{-1}(AJ_\sigma)a = \text{Proj}_{\mathbb{R}^{\sigma_u}} J_\sigma a.$$

since $A^{-1}A = I$ and projecting a subset onto its containing set is just the subset itself. Then, breaking $J_\sigma a$ into orthogonal components,

$$\|J_\sigma a\|_2^2 = \sum_{u=1}^{r+1} \|\text{Proj}_{\mathbb{R}^{\sigma_u}} J_\sigma a\|_2^2$$
$$= \sum_{u=1}^{r+1} \left\|\text{Proj}_{\mathbb{R}^{\sigma_u}} (AJ_{\sigma_u})^{-1}(AJ_\sigma)a\right\|_2^2$$
$$\lesssim \sum_{u=1}^{r+1} 2^{2r-u} M^2 \|AJ_\sigma a\|_2^2 \qquad \text{by (2.15)}$$
$$\lesssim 2^{2r} M^2 \|AJ_\sigma a\|_2^2$$

40

$$\leq \frac{M^2}{1 - \frac{k}{m}} \left\| AJ_\sigma a \right\|_2^2$$

$$\left\| J_\sigma a \right\|_2 \leq \sqrt{\frac{m}{m-k}} M \left\| AJ_\sigma a \right\|_2$$

Since this is true for all $a \in \mathbb{R}^\sigma$,

$$\implies \left\| (AJ_\sigma)^{-1} \right\|_{S_\infty} \lesssim \sqrt{\frac{m}{m-k}} M.$$

$\square$

An important aspect of this whole proof to realize is the way we took a first subset to give one bound, and then took another subset in order to apply the Little Grothendieck inequality. In other words, we first took a subset where we could control $l_1$ to $l_2$, and then took a further subset to get the operator norm bounds.

Also, notice that in this approach, we construct our "optimal subset" of the space inductively: However, doing this in general is inefficient. It would be nice to construct the set greedily instead for algorithmic reasons, if we wanted to come up with a constructive proof. I have a feeling that if we were to avoid this kind of induction, it would involve not using Sauer-Shelah at all.

How can we make this theorem algorithmic?

The way the Pietsch Domination Theorem 2.3.5 worked was by duality. We look at a certain explicitly defined convex set. We found a separating hyperplane which must be a probability measure. Then we had a probabilistic construction. This part is fine.

The bottleneck for making this an algorithm (I do believe this will become an algorithm) consists of 2 parts:

1. Sauer-Shelah lemma 2.3.9: We have some cloud of points in the boolean cube, and we know there is some large subset of coordinates (half of them) such that when you project to it you see the full cube. I'm quite certain that it's NP-hard in general. (Work is necessary to formulate the correct algorithmic Sauer-Shelah lemma. How is the set given to you?). In fact, formulating the right question for an algorithmi Sauer-Shelah is the biggest difficulty.

   We only need to answer the algorithmic question tailored to our sets, which have a simple description: the intersection of an ellipsoid with a cube. There is a good chance that there is a polynomial time algorithm in this case. This question has other applications as well (perhaps one could generalize to the intersection of the cube with other convex bodies). [3]

2. The second bottleneck is finding a subset with maximal volume.

   It's another place where we chose a subset, the subset that maximizes the volume out of all subsets of a given size (Lemma 2.4.1). Specifically, we want the set of columns

---

[3]There could be an algorithm out there, but I couldn't find one in the literature.

of a matrix that maximizes the volume of the convex hull. Computing the volume of the convex hull requires some thought. Also, there are $\binom{n}{r}$ subsets of size $r$; if $r = n/2$ there are exponentially many. We need a way to find subsets with maximum volume fast. There might be a replacement algorithm which approximately maximizes this volume.

2-22

Let me just remind you were we were. We were at the final lemma in the Restricted Invertiblity Theorem. We have a linear map $A : \mathbb{R}^m \to \mathbb{R}^n$, and have $\{Ae_j\}_{j=1}^m$ linearly independent, and denote

$$F_j = (\text{span}(\{Ae_j\}_{i \neq j}))^\perp$$

and then denote $M = \max_{i \leq j \leq m} \frac{1}{\|\text{Proj}_{F_1} Ae_j\|}$.

We have reduced everything to the following lemma:

**Lemma 2.4.5.** *lem:SS-induct (Final lemma).*
*Suppose $\sigma \subseteq \{1, \cdots, m\}$ with $t \geq 0$ an integer. Then there exists $\tau \subseteq \sigma$ with $|\tau| \geq (1 - \frac{1}{2^t})|\sigma|$ such that if we denote $\theta = \tau \cup (\{1, \cdots, m\} \setminus \sigma)$, then*

$$\sum_{i \in \tau} |a_i| \leq 2^{t/2} M \sqrt{|\sigma|} \| \sum_{i \in \theta} a_i Ae_i \|_2$$

The proof of this will be an inductive application of the Sauer-Shelah lemma. A very important idea comes from Giannopoulos. If you naively try to use Sauer-Shelah, it won't work out. We will give a stronger statement of the previous lemma which we can prove by induction.

**Lemma 2.4.6.** *lem:SS-induct-stronger (Stronger version of Final lemma).*
*Take $m, n \in \mathbb{N}$, $A : \mathbb{R}^m \to \mathbb{R}^n$ a linear operator such that $\{Ae_j\}_{j=1}^m$ are linearly independent. Suppose that $k \geq 0$ is an integer and $\sigma \subseteq \{1, \cdots, m\}$. Then there exists $\tau \subseteq \sigma$ with $|\tau| \geq (1 - \frac{1}{2^k})|\sigma|$ such that for every $\theta \supseteq \tau$ for all $a \in \mathbb{R}^m$ we have*

$$\sum_{i \in \tau} |a_i| \leq M \sqrt{|\sigma|} \left( \sum_{r=1}^k 2^{r/2} \right) \left\| \sum_{i \in \theta} a_i Ae_i \right\|_2 + (2^k - 1) \sum_{i \in \theta \cap (\sigma \setminus \tau)} |a_i| (*)$$

Our first lemma (all we need to complete Restricted Invertibility Principle) is the case where $\theta = \tau \cup \{1, \cdots, m\} \setminus \sigma$, $t = k$.

We prove the stronger version via induction on $k$.

*Proof.* As $k$ becomes bigger, we're eating more and more out from the set $\sigma$. So we're going to use Sauer-Shelah, taking half of $\sigma$, and then a bit more and a bit more. For $k = 0$, this is vacuous, since we take $\tau$ to be an empty set. Now via induction assume for $k$ that we found already $\tau \subseteq \sigma$, with $|\tau| \geq (1 - \frac{1}{2^k})|\sigma|$ and satisfies (*) for every $\tau \subseteq \theta$. If $\sigma = \tau$ already, then $\tau$ satisfies for $k + 1$ as well, since WLOG $|\sigma \setminus \tau| > 0$. Now define $v_j$ is the projection

$$v_j = \frac{\text{Proj}_{F_j} Ae_j}{\left\| \text{Proj}_{F_j} Ae_j \right\|_2^2}$$

42

Then $\langle v_i, Ae_j \rangle = \delta_{ij}$, by definition since we're looking at a dual basis for the $Ae_j$s.

Now we want to user Sauer-Shelah so we're going to define a certain subset of the cube. Define

$$\Omega = \{\epsilon \in \{\pm 1\}^{\sigma \setminus \tau} : \left\| \sum_{i \in \sigma \setminus \tau} \epsilon_i v_i \right\|_2 \leq M\sqrt{2|\sigma \setminus \tau|}\}$$

So this is really an ellipsoid intersected with the cube, since the $v_i$s are not orthogonal. Then we have

$$M^2|\sigma \setminus \tau| \geq \sum_{i \in \sigma \setminus \tau} \frac{1}{\|\mathrm{Proj}_{F_j} Ae_j\|^2} = \sum_{j \in \sigma \setminus \tau} \|v_j\|_2^2$$

$$= \frac{1}{2^{|\sigma \setminus \tau|}} \sum_{\epsilon \in \{\pm 1\}^{\sigma \setminus \tau}} \| \sum_{j \in \sigma \setminus \tau} \epsilon_j v_j \|_2^2$$

where the last step is true for any vectors (sum the squares and the pairwise correlations disappear).

Now we're using Markov's inequality to get

$$\geq \frac{1}{2^{|\sigma \setminus \tau|}} \left( 2^{|\sigma \setminus \tau|} - |\Omega| \right) M^2 2|\sigma \setminus \tau|$$

which gives

$$|\Omega| > 2^{|\sigma \setminus \tau| - 1}$$

Then by Sauer-Shelah lemma, there exists $\beta \subseteq \sigma \setminus \tau$ such that

$$\mathrm{Proj}_{\mathbb{R}^\beta} \Omega = \{\pm 1\}^\beta$$

and

$$|\beta| \geq \frac{1}{2}|\sigma \setminus \tau|$$

Now define $\tau^* = \tau \cup \beta$. We will show that $\tau^*$ satisfies the inductive hypothesis with $k+1$. Each time we find a certain set of coordinates to add to what we have before. $|\tau^*|$ is the correct size because

$$|\tau^*| = |\tau| + |\beta| \geq |\tau| + \frac{|\sigma| - |\tau|}{2} = \frac{|\tau| + |\sigma|}{2} \geq \left( 1 - \frac{1}{2^{k+1}} \right) |\sigma|$$

where we used that $|\tau| \geq \left( 1 - \frac{1}{2^k} \right) |\sigma|$. So at least $\tau^*$ is the right size.

Now, suppose $\theta \supseteq \tau^*$. For every $a \in \mathbb{R}^m$, we claim there exists some $\epsilon \in \Omega$ such that $\forall j \in \beta$ such that $\epsilon_j = \mathrm{sign}(a_j)$. For any $\beta$, we can find some vector in the cube that has the sign pattern of our given vector $a$. What does being in $\Omega$ mean? It means that at least the dual basis is small there. $\epsilon \in \Omega$ says that

$$\| \sum_{i \in \sigma \setminus \tau} \epsilon_i v_i \|_2 \leq M\sqrt{2|\sigma \setminus \tau|} \leq \frac{M\sqrt{2|\sigma|}}{2^{k/2}}$$

That was how we chose our ellipsoid. So we know a bound for just $\tau$ already, now let's do it with the addition of $\beta$. Well,

$$\sum_{i\in\beta} |a_i| = \langle \sum_{i\in\beta} a_i Ae_i, \sum_{i\in\sigma\setminus\tau} \epsilon v_i \rangle$$

which is precisely because the $v_i$'s were a dual basis and the dot products will be one. We only know the $\epsilon_i$ are the signs when you're inside $\beta$. This equals

$$= \langle \sum_{i\in\theta} a_i Ae_i, \sum_{i\in\sigma\setminus\tau} \epsilon v_i \rangle - \sum_{i\in(\theta\setminus\beta)\cap(\sigma\setminus\tau)} \epsilon_i a_i$$

Note that $(\theta\setminus\beta)\cap(\sigma\setminus\tau) = \theta\cap(\sigma\setminus\tau^*)$. We can't control the signs $\epsilon_i$ any more. Then, we get applying Sauer-Shelah

$$\leq \left\| \sum_{i\in\theta} a_i Ae_i \right\|_2 \cdot \left\| \sum_{i\in\sigma\setminus\tau} \epsilon v_i \right\|_2 + \sum_{i\in\theta\cap(\sigma\setminus\tau^*)} |a_i|$$

$$\leq \left\| \sum_{i\in\theta} a_i Ae_i \right\|_2 \cdot \frac{M\sqrt{2|\sigma|}}{2^{k/2}} + \sum_{i\in\theta\cap(\sigma\setminus\tau^*)} |a_i|$$

since Sauer-Shelah told us nothing about the signs of $\epsilon_i$, so we just take the worst possible thing.

Then

$$\sum_{i\in\beta} |a_i| \leq \frac{M\sqrt{2|\sigma|}}{2^{k/2}} \left\| \sum_{i\in\theta} a_i Ae_i \right\|_2 + \sum_{i\in\theta\setminus(\sigma\setminus\tau^*)} |a_i|$$

Using the inductive step,

$$\sum_{i\in\tau^*} |a_i| = \sum_{i\in\tau} |a_i| + \sum_{i\in\beta} |a_i|$$

$$\leq M\sqrt{|\sigma|}\alpha_k \left\| \sum_{i\in\theta} a_i Ae_i \right\|_2 + (2^k - 1) \sum_{i\in\theta\cap(\sigma\setminus\tau)} |a_i| + \sum_{i\in\beta} |a_i|$$

$$= \alpha_k\sqrt{|\sigma|} \left\| \sum_{i\in\theta} a_i Ae_i \right\|_2 + (2^k - 1) \sum_{i\in\theta\cap(\sigma\setminus\tau^*)} |a_i| + 2^k \sum_{i\in\beta} |a_i|$$

In the last step, we throw $\beta$ away, but with weight $2^k - 1$. And now use what we got before for the bound on $\sum_{i\in\beta} |a_i|$ and plug it in to get

$$\leq \left( \alpha_k + 2^{(k+1)/2} \right) \sqrt{|\sigma|} \left\| \sum_{i\in\theta} a_i Ae_i \right\|_2 + \left( 2^{k+1} - 1 \right) \sum_{i\in\theta\cap(\sigma\setminus\tau^*)} |a_i|$$

which is exactly the inductive hypothesis. I looked through the original Gianpopoulous paper, and it was clear he tried out many many things to find which inductive hypothesis makes everything go through cleanly. You want to bound an $l_1$ sum from above, so you want to use duality, and then use Sauer-Shelah to get signs such that the norm of the dual-basis is small. $\qquad\square$

**Remark 2.4.7:** Algorithmic Sauer-Shelah.
Now regarding the use of Sauer-Shelah, we can see that we are only using it for intersecting cubes with ellipsoid. So regarding what I said last time, what we need is an algorithm for finding these intersections. The reason these ellipsoids are big is because we are actually multiplying by $\sqrt{2}$ times the expectation. So this algorithm is probably do-able. Maybe afterwards, you could ask for higher dimensional shapes. I've seen some references that worked for Sauer-Shelah when sets were of a special form, namely of size $o(n)$. This is something more geometric. I don't think there's literature about Sauer-Shelah for intersection of surfaces with small degree. This is a tiny motivation to do it, but it's still interesting independently.

# 5 Bourgain's Discretization Theorem

We will prove Bourgain's Discretization Theorem. This will take maybe two weeks, and has many interesting ingredients along the way. By doing this, we will also prove Ribe's theorem, which is what we stated at the beginning.

Let's remind ourselves of the definition:

**Definition 2.5.1:** Discretization Modulus.
$(X, \|\cdot\|_X), (Y, \|\cdot\|_Y)$ are Banach spaces. Let $\epsilon \in (0,1)$. Then $\delta_{X\hookrightarrow Y}(\epsilon)$ is the supremum over $\delta > 0$ such that for every $\delta$-net $N_\delta$ of the unit ball of $X$, the distortion

$$C_Y(X) \leq \frac{C_Y(N_\delta)}{1 - \epsilon}$$

$C_Y$ is smallest bi-Lipschitz distortion by which you can embed $X$ into $Y$. There are ideas required to get $1 - \epsilon$. I'll decide by next time if I want to do the full thing or just do a constant. Think of $\epsilon = 1/2$. What we're saying is that if we succeed in embedding a $\delta$-net into $Y$, then we succeeded in the full space with a distortion twice as much. A priori it's not even clear that there exists such a $\delta$. There's a nontrivial compactness argument to prove you can (Lesbegue density points). But we will just prove bounds on it assuming it exists.

Now, Bourgain's discretization theorem says

**Theorem 2.5.2.** *Bourgain's discretization theorem.*
*If $dim(X) = n$, $dim(Y) = \infty$, then*

$$\delta_{X\hookrightarrow Y}(\epsilon) \geq e^{-\left(\frac{n}{\epsilon}\right)^{C*n}}$$

*for $C$ a universal constant.*

The way to read this is: If $\delta$ is bigger than this number which is only dependent on $n$, then given any Banach spaces with dimension $n$, there are mappings with this granularity for any such spaces.

**Remark 2.5.3:** It doesn't matter what mapping you're actually using, the proof will give a linear mapping and we won't end up needing them. Assuming linear map in the definition is not necessary. Rademacher's theorem says that for any mapping from $\mathbb{R}^n \to \mathbb{R}^n$ is differentiable almost everywhere for bi-Lipschitz derivative, and you can extend this to $n = \infty$, but you need some additional properties: You need to be embedding into the dual space, and the limit needs to be in the weak$-*$ topology. That derivative is almost everywhere, but you can definitely have a sequence of norm$-1$ vectors that tend to 0. A weak $*$ limit of a function can degenerately become 0, the upper bound is not the issue. But you can prove that this doesn't happen almost everywhere. You can look at the Principle of Local Reflexivity, which says $Y^{**}$ and $Y$ are not the same for infinite dimensions. The double dual of all sequences which tend to 0 is $L_\infty$, a bigger space, but the difference between these is never appearing in finite dimensional phenomena.

From now on, $B_X = \{x \in X : \|X\|_X \leq 1\}$, the ball. $S_X = \partial B = \{x \in X : \|X\|_X = 1\}$, the boundary.

Later on we will be differentiating things without thinking about it, so I just want to prove to you first that $X \to \|X\|_X$ is smooth on $X \setminus \{0\}$.

**Lemma 2.5.4.** *For all $\delta \in (0,1)$ there exists some $\delta$-net of $S_X$ with $|N_\delta| \leq \left(1 + \frac{2}{\delta}\right)^n$.*

*Proof.* Let $N_\delta \subseteq S_X$ be maximal with respect to inclusion such that $\|x - y\|_X > \delta$ for every distinct $x, y \in N_\delta$. We want it to be both separated and $\delta$-dense. For every $z \in S_X$, if $z \in N_\delta$, $\{z\} \cup N_\delta$ implies there exists $x \in N_\delta$, $\|x - y\|_X \leq \delta$. The balls $\{x + \frac{\delta}{2}B_X\}_{x \in M_\delta}$ are pairwise disjoint. Moreover, the balls are all contained in $1 + \frac{\delta}{2}B_X$. And then we get volume $(\frac{\delta}{2})^n \text{vol}(B_X)$, and we can get

$$\text{vol}((1 + \frac{\delta}{2})B_X) \geq \sum_{X \in M_\delta} \text{vol}(x + \frac{\delta}{2}B_X)$$

$$\left(1 + \frac{\delta}{2}\right)^n \text{vol}(B_X) = |N_\delta|\left(\frac{\delta}{2}\right)^n \text{vol}(B_X)$$

If you ask what the smallest size of a $\delta$-net is, there are bounds, but they are not sharp. There is a lot of literature about the relations between these things, we just need an upper bound.

Now say you have your convex body, and you find your $\delta$-net $N_\delta$ of $S_X$ with $|N_\delta| = N \leq (1 + \frac{2}{\delta})^n$ which is finite. Then for every $x \in N_\delta$, choose any $z^* \in X^*$ unit vector on the sphere ($\langle z^*, z \rangle = 1$) and $\|z^*\|_{X^*} = 1$ by Hahn-Banach. It normalizes the net-point. What would be a good approximation? Let $k$ be an integer such that $N^{1/(2k)} \leq 1 + \delta$. Then define

$$\|x\| := \left(\sum_{z \in N_\delta} \langle z^*, x \rangle^{2k}\right)^{1/(2k)}$$

Each term, separately $|\langle z^*, x \rangle| \leq \|x\|_X$, thus we know $(1 - \delta)\|x\|_X \leq \|x\| \leq N^{1/2k}\|x\|_X \leq (1 + \delta)\|x\|_X$. If $x \in S_X$, then choose $z \in N_\delta$ such that $\|x - z\|_X \leq \delta$, and thus $1 - \langle z^*, x \rangle =$

$\langle z^*, z - x \rangle \leq \|z - x\| \leq \delta$, and $\langle z^*, x \rangle \geq 1 - \delta$. Thus any norm is up to $1 + \delta$ some really nice smooth norm. $\delta$ was arbitrary, if you prove for Bourgain, you prove it for any norm, and now without loss of generality I can differentiate. $\qquad\square$

Let me just explain the strategy of how we will prove Bourgain's discretization theorem. We're given a $\delta$-net $N_\delta \subseteq B_X$ with $\delta \leq e^{-(n/\epsilon)^{Cn}}$, and we know $\exists\ f : N_\delta \to Y$ such that $\frac{1}{D}\|x - y\|_X \leq \|f(x) - f(y)\|_Y \leq \|x - y\|_X$, which means you can embed with distortion $D$. Our goal: If $\delta < e^{-(D/\epsilon)^{Cn}}$, then there exists a linear operator $T : X \to Y$ such that $\|T\| \cdot \|T^{-1}\| \leq 1 + 20\epsilon$.

We will need a little background in convex geometry. We're going to find the correct coorindate system (John ellipsoid), which will give us a dot product structure and a natural Laplacian. A priori $f$ is defined on the net. We're going to find that it extends to the whole space in a nice way (Bourgain's extension theorem) which doesn't coincide with the function on the net, but is not too far away from it. Then we will solve the Laplace equation. We will start at initial condition, and then evolve $f$ according to the Poisson semigroup. This extended function is going to be smooth the minute you flow a little bit away from your discrete function. And you can differentiate it! We will prove that there's a point where the derivative satisfies what we want: The point cannot not exist (pigeonhole style argument), but we won't be able to pinpoint where the derivative behaves nicely. And this will come from estimates of the Poisson kernel, and we will jump into the Fourier analysis.