

# Lecture Notes on Random Matrix Theory in Data Science and Statistics

Dmitriy (Tim) Kunisky

Fall 2025 (Last Updated: September 12, 2025)

# Contents

|          |   |           |
|----------|---|-----------|
| <b>1</b> | <b>Random Vector Theory</b>   | <b>2</b>  |
| 1.1      | Natural Random Models and Orthogonal Invariance . . . . .                   | 2         |
| 1.2      | Gaussian Random Vectors . . . . .   | 3         |
| 1.3      | Concentration of Gaussian Vector Norms . . . . .                            | 4         |
| 1.4      | Consequences for High-Dimensional Geometry . . . . .                        | 7         |
| <b>2</b> | <b>First Steps With Rectangular Matrices</b>                                | <b>9</b>  |
| 2.1      | Gaussian Random Field Viewpoint . . . . .                                   | 9         |
| 2.2      | Application: Dimensionality Reduction and Johnson-Lindenstrauss Lemma . . . | 12        |
| 2.2.1    | Simple Extensions . . . . .   | 13        |
| 2.2.2    | Lower Bounds . . . . .  | 14        |
| 2.2.3    | Sparse and Fast JLTs . . . . .  | 14        |
| 2.2.4    | Proof Technique: First Moment Method . . . . .                              | 14        |
| 2.3      | Rectangular Gaussian Matrices: Singular Values . . . . .                    | 15        |
| 2.3.1    | The <b>Geometric Method</b> of Random Matrix Theory . . . . .               | 16        |
| 2.3.2    | Nets, Coverings, and Packings . . . . .                                     | 18        |
| 2.3.3    | Coarse Non-Asymptotic Bound on Singular Values . . . . .                    | 20        |
| 2.4      | Rectangular Gaussian Matrices: Singular Vectors . . . . .                   | 23        |
| 2.4.1    | Distinctness of Singular Values . . . . .                                   | 24        |
| 2.4.2    | Orthogonal Invariance for Matrices . . . . .                                | 26        |
| 2.5      | Application: Compressed Sensing . . . . .                                   | 27        |
| 2.5.1    | Null Space and Restricted Isometry Properties . . . . .                     | 28        |
| 2.5.2    | Random Sensing Matrices . . . . .   | 29        |
|          | <b>Bibliography</b>   | <b>31</b> |

# 1 | RANDOM VECTOR THEORY

Before we continue to random matrix theory, it will be useful to establish some preliminaries about random *vectors*, which will also lead us to a few important intuitions about high-dimensional probability and geometry.

## 1.1 NATURAL RANDOM MODELS AND ORTHOGONAL INVARIANCE

As in the case of matrices, vectors have at least two natural interpretations: they are a **list of numbers**, as well as a geometric object with a **magnitude and direction** (related notions about vector and Banach spaces are abstractions of this latter viewpoint) with respect to a basis of the space to which they belong. Each interpretation gives rise to a natural model<sup>1</sup> of random vector.

In the list-of-numbers interpretation, it is most logical to model a random vector where the entries are as statistically decoupled as possible. Thus, we consider

$$\mathbf{x} \sim \mu^{\otimes d}, \tag{1.1.1}$$

a vector with entries i.i.d. from some probability measure  $\mu$  on  $\mathbb{R}$ .

In the magnitude-and-direction interpretation, it is most logical to model a random vector where the magnitude and direction are as statistically decoupled as possible. Thus, we take these two properties to be independent (that is, a random  $\mathbf{x}$  where  $\|\mathbf{x}\|$  and  $\mathbf{x}/\|\mathbf{x}\|$  are independent). Since the magnitude can be any random non-negative scalar, the more interesting new information of such a vector distribution is in its direction. Thus, we consider

$$\mathbf{x} \sim \text{Unif}(\mathbb{S}^{d-1}), \tag{1.1.2}$$

a vector drawn from the uniform surface measure on the unit sphere  $\mathbb{S}^{d-1} := \{\mathbf{x} \in \mathbb{R}^d : \|\mathbf{x}\| = 1\}$  (so denoted because it is a  $(d-1)$ -dimensional manifold).

More generally, the following property of a vector distribution is natural to consider in the context of the magnitude-and-direction interpretation.

**Definition 1.1.1** (Orthogonal invariance). *We say that a random vector  $\mathbf{x} \in \mathbb{R}^d$  (or its law  $\text{Law}(\mathbf{x})$ ) is orthogonally invariant if, for each (deterministic) orthogonal matrix  $\mathbf{Q} \in \mathcal{O}(d)$ , we have  $\text{Law}(\mathbf{Q}\mathbf{x}) = \text{Law}(\mathbf{x})$ .*

---

<sup>1</sup>We will sometimes use this word “model” instead of “law” or “distribution” to emphasize that such a choice includes some judgement of what distributional assumption is natural or reasonable—we make such a choice to try to *model* typical situations we might encounter in applications.

Without going into more technical details of actually defining this probability measure explicitly, the following result should clarify the meaning of the measure  $\text{Unif}(\mathbb{S}^{d-1})$ .

**Proposition 1.1.2.** *There exists a unique probability measure supported on  $\mathbb{S}^{d-1}$  that is orthogonally invariant, which we denote  $\text{Unif}(\mathbb{S}^{d-1})$ .*

The following also shows that such distributions are nothing more than vectors distributed as  $\text{Unif}(\mathbb{S}^{d-1})$  rescaled by some independent random scalar, and thus are precisely the more general class of models natural to the magnitude-and-direction interpretation that we indicated above.

**Proposition 1.1.3.** *Suppose that  $\mathbf{x}$  is orthogonally invariant. Then,  $\|\mathbf{x}\|$  is independent of  $\mathbf{x}/\|\mathbf{x}\|$ , and  $\text{Law}(\mathbf{x}/\|\mathbf{x}\|) = \text{Unif}(\mathbb{S}^{d-1})$ .*

Each of these distributions makes the random vector easy to reason about in its “native” interpretation: the entries of an i.i.d. vector and the direction and magnitude of an orthogonally invariant vector are both described directly in the definitions. On the other hand, the information “across” interpretations is harder to access. We will build up some tools to address these questions now: what do the entries of  $\mathbf{x} \sim \text{Unif}(\mathbb{S}^{d-1})$  look like? And, what are the magnitude and direction of an i.i.d. random vector?

## 1.2 GAUSSIAN RANDOM VECTORS

As we will see, in many ways the most natural of the i.i.d. random vectors is the standard multivariate Gaussian  $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ . We also denote this law in the more conventional multivariate Gaussian notation  $\mathcal{N}(\mathbf{0}, \mathbf{I}_d) = \mathcal{N}(\mathbf{0}, \Sigma)$ . The general meaning of this notation, which you should be familiar with, is as follows.

**Definition 1.2.1.**  $\mathcal{N}(\boldsymbol{\mu}, \Sigma)$  for parameters  $\boldsymbol{\mu} \in \mathbb{R}^d$  and  $\Sigma \in \mathbb{R}_{\geq 0}^{d \times d}$  is the probability measure with density

$$\frac{1}{\sqrt{\det^*(2\pi\Sigma)}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \Sigma^+ (\mathbf{x} - \boldsymbol{\mu})\right) \quad (1.2.1)$$

relative to the Lebesgue measure on the row space of  $\Sigma$ . Here,  $\Sigma^+$  is the Moore-Penrose pseudoinverse, and  $\det^*$  is the product of all non-zero eigenvalues. When  $\Sigma$  is invertible, then  $\det^*$  above is the ordinary determinant,  $\Sigma^+ = \Sigma^{-1}$ , and the above density is with respect to the Lebesgue measure on all of  $\mathbb{R}^d$ . We call  $\mathbf{x}$  a Gaussian random vector if  $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \Sigma)$  for some  $\boldsymbol{\mu}$  and  $\Sigma$ , and a standard Gaussian random vector if  $\boldsymbol{\mu} = \mathbf{0}$  and  $\Sigma = \mathbf{I}_d$ .

The following are the two fundamental properties of Gaussian random vectors.

**Proposition 1.2.2.** *The parameters  $\boldsymbol{\mu}$  and  $\Sigma$  are the mean vector and covariance matrix of the random vector  $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \Sigma)$ :*

$$\boldsymbol{\mu} = \mathbb{E}\mathbf{x}, \quad (1.2.2)$$

$$\Sigma = \mathbb{E}(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^\top = \mathbb{E}\mathbf{x}\mathbf{x}^\top - \boldsymbol{\mu}\boldsymbol{\mu}^\top \quad (1.2.3)$$

Thus, the law of a Gaussian random vector is determined by its mean and covariance, or equivalently by its linear and quadratic moments,  $\mathbb{E}\mathbf{x}$  and  $\mathbb{E}\mathbf{x}\mathbf{x}^\top$ .

**Proposition 1.2.3.** *If  $\mathbf{x} \in \mathbb{R}^d$  is a Gaussian random vector and  $\mathbf{A} \in \mathbb{R}^{d' \times d}$  and  $\mathbf{b} \in \mathbb{R}^{d'}$ , then  $\mathbf{Ax} + \mathbf{b}$  is also a Gaussian random vector. If  $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ , then the corresponding parameters after this linear transformation are*

$$\text{Law}(\mathbf{Ax} + \mathbf{b}) = \mathcal{N}(\mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^\top). \quad (1.2.4)$$

These two facts taken together make it easy to do certain computations with Gaussian random vectors: provided we only make linear transformations, we both remain within the class of Gaussian random vectors and can keep track of all distributions involved merely by linear algebra.<sup>2</sup>

The following remarkable facts about standard Gaussian random vectors then follow by simple calculations using the above. These show that, in fact, a standard Gaussian random vector is *both* an i.i.d. random vector *and* an orthogonally invariant one, thus bridging the two natural classes of models of vectors we were discussing before.

**Proposition 1.2.4.** *Suppose  $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ . For any  $\mathbf{a} \in \mathbb{S}^{d-1}$ , we have  $\text{Law}(\langle \mathbf{a}, \mathbf{x} \rangle) = \mathcal{N}(0, 1)$ . In particular, this law does not depend on  $\mathbf{a}$ ; the distribution of  $\mathbf{x}$  is “the same in any direction.” Further,  $\mathbf{x}$  is orthogonally invariant.*

*Proof.* Using the previous results, for  $\mathbf{Q} \in \mathcal{O}(d)$ ,

$$\text{Law}(\mathbf{Qx}) = \mathcal{N}(\mathbf{0}, \mathbf{Q}^\top \mathbf{I}_d \mathbf{Q}) = \mathcal{N}(\mathbf{0}, \mathbf{I}_d) = \text{Law}(\mathbf{x}), \quad (1.2.5)$$

giving the second result. The first follows from the second by considering a  $\mathbf{Q}$  whose first row is  $\mathbf{a}$ .  $\square$

**Corollary 1.2.5.** *If  $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ , then  $\text{Law}(\mathbf{x}/\|\mathbf{x}\|) = \text{Unif}(\mathbb{S}^{d-1})$ .*

*Proof.* The result follows directly from the above together with Proposition 1.1.3.  $\square$

In fact, the following elegant result, which we will not prove here, shows that the Gaussian scalar measure is the *only* one having this property.

**Theorem 1.2.6 (Maxwell).** *Suppose that  $\mathbf{x} \sim \mu^{\otimes d}$  for some  $\mu$  a probability measure on  $\mathbb{R}$ . If  $d \geq 2$  and  $\mathbf{x}$  is orthogonally invariant, then  $\mu = \mathcal{N}(0, \sigma^2)$  for some  $\sigma^2 \geq 0$ .*

## 1.3 CONCENTRATION OF GAUSSIAN VECTOR NORMS

The above tells us that the direction of a Gaussian random vector is uniformly distributed; in the exercises, you will see some similar more quantitative statements for more general classes of i.i.d. random vectors. What about the magnitudes of such vectors? Again, let us focus here on the Gaussian case.

A simple calculation gives the expectation of the norm (squared): if  $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ , then

$$\mathbb{E}\|\mathbf{x}\|^2 = d \cdot \mathbb{E}x_1^2 = d. \quad (1.3.1)$$

---

<sup>2</sup>We will not use it in this class, but another important fact to be aware of is that the more complicated operation of *conditioning* a Gaussian random vector  $\mathbf{x}$  on the values of some linear forms  $\langle \mathbf{a}_i, \mathbf{x} \rangle$  also yields another Gaussian random vector.

And, since  $\|\mathbf{x}\|^2 = \sum_{i=1}^d x_i^2$  is a sum of i.i.d. random variables, the basic theory of such sums from scalar probability (law of large numbers, central limit theorem, etc.) tells us to expect that we should have  $\|\mathbf{x}\|^2 \approx d$  with high probability; more precisely, we should have  $\|\mathbf{x}\|^2 = d + O(\sqrt{d})$ .

Let us first try to establish this using some basic concentration inequalities. Markov's inequality gives

$$\mathbb{P}[\|\mathbf{x}\|^2 \geq d + t] \leq \frac{d}{d + t} = \frac{1}{1 + \frac{t}{d}}, \quad (1.3.2)$$

which only tells us that  $\|\mathbf{x}\|^2 = O(d)$  with high probability (say,  $\mathbb{P}[\|\mathbf{x}\|^2 \leq 100d] \geq 0.99$  and similar statements); the inequality “kicks in” and gives useful information only when  $t \gtrsim d$ . Chebyshev's inequality gives

$$\mathbb{P}[\|\mathbf{x}\|^2 \geq d + t] = \mathbb{P}[\|\mathbf{x}\|^2 \geq \mathbb{E}\|\mathbf{x}\|^2 + t] \leq \frac{\text{Var}[\|\mathbf{x}\|^2]}{t^2} = \frac{2d}{t^2}, \quad (1.3.3)$$

where we calculate  $\text{Var}[\|\mathbf{x}\|^2] = d \text{Var}[x_1^2] = 2d$ . This gives the correct scaling, kicking in when  $t \gtrsim \sqrt{d}$ . However, as the following result shows, the actual dependence on  $t$  above is far from optimal.

**Lemma 1.3.1.** *For  $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$  and any  $t \geq 0$ ,*

$$\begin{aligned} \mathbb{P}\left[\left|\|\mathbf{x}\|^2 - d\right| \geq t\right] &= \mathbb{P}\left[\left|\sum_{i=1}^d x_i^2 - d\right| \geq t\right] \\ &\leq 2 \begin{cases} \exp\left(-\frac{t^2}{8d}\right) & \text{if } t \leq d, \\ \exp\left(-\frac{t}{8}\right) & \text{if } t \geq d \end{cases} \end{aligned} \quad (1.3.4)$$

$$= 2 \exp\left(-\frac{1}{8} \min\left\{\frac{t^2}{d}, t\right\}\right). \quad (1.3.5)$$

*Proof.* The proof uses the “Chernoff method,” a proof technique behind concentration inequalities like the Hoeffding, Bernstein, Azuma-Hoeffding, and McDiarmid inequalities. The method is a general instance of something you might call a “nonlinear Markov” inequality that also appears in the proofs of other bounds like Chebyshev's inequality. This whole family of approaches is important to know about.

We will only deal with one of the tails; the other follows similarly and the factor of 2 appears from taking a union bound. Note that  $\mathbb{E}x_i^2 = 1$ , so, defining  $s_i := x_i^2 - 1$ , we have

$$\begin{aligned} \mathbb{P}\left[\sum_{i=1}^d x_i^2 - d \geq t\right] &= \mathbb{P}\left[\sum_{i=1}^d s_i \geq t\right] \\ &= \mathbb{P}\left[\exp\left(\lambda \sum_{i=1}^d s_i\right) \geq \exp(\lambda t)\right] \\ &\leq \frac{\mathbb{E} \exp\left(\lambda \sum_{i=1}^d s_i\right)}{\exp(\lambda t)} && \text{(Markov inequality)} \\ &= \frac{(\mathbb{E} \exp(\lambda s_1))^d}{\exp(\lambda t)} && \text{(independence)} \\ &= \exp(-\lambda t + d \log \mathbb{E} \exp(\lambda s_1)). \end{aligned} \quad (1.3.6)$$

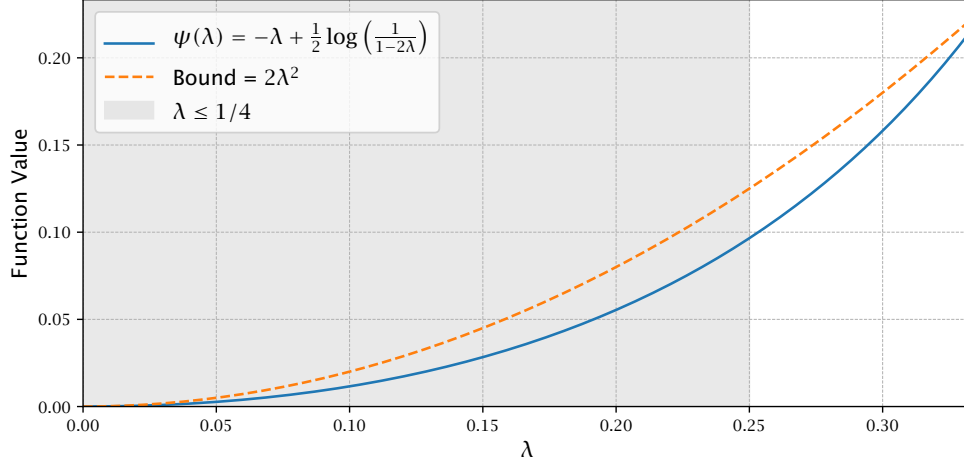


Figure 1.1: The cumulant generating function  $\psi(\lambda)$  appearing in Lemma 1.3.1, our bound on it, and the region where we apply this bound.

Let us write

$$\psi(\lambda) := \log \mathbb{E} \exp(\lambda s_1) = \log \mathbb{E} \exp(\lambda(x_1^2 - 1)), \quad (1.3.7)$$

the *cumulant generating function* of the random variable  $x_1^2 - 1$  (where we retain the  $-1$  so that the expectation is zero. You may check as a calculus exercise that the remaining expectation is finite if and only if  $\lambda < \frac{1}{2}$ , and in this case

$$\mathbb{E} \exp(\lambda x_1^2) = \frac{1}{\sqrt{1 - 2\lambda}}, \quad (1.3.8)$$

whereby

$$\psi(\lambda) = -\lambda + \frac{1}{2} \log \left( \frac{1}{1 - 2\lambda} \right). \quad (1.3.9)$$

Taylor expansion shows that  $\psi(\lambda) = \lambda^2 + O(\lambda)$ . Turning this into a concrete bound, you may convince yourself by plotting (see Figure 1.1) and check by a slightly tedious algebra exercise that  $\psi(\lambda) \leq 2\lambda^2$  for all  $\lambda \leq 1/4$ .

Thus we have that, provided  $\lambda \leq 1/4$ ,

$$\begin{aligned} \mathbb{P} \left[ \sum_{i=1}^d x_i^2 - d \geq t \right] &\leq \exp(-\lambda t + d\psi(\lambda)) \\ &\leq \exp(-\lambda t + 2d\lambda^2). \end{aligned}$$

Inside the exponential now remains a convex quadratic function of  $\lambda$ , which has a unique minimizer that we compute to be at  $\lambda^* := \frac{t}{4d}$ . We would like to set  $\lambda := \lambda^*$ , but we must be careful to respect that we have used a bound only holding for  $\lambda \leq \frac{1}{4}$ . So, we consider two cases:

*Case 1:*  $t \leq d$ . In this case, we may set  $\lambda := \lambda^* \leq \frac{1}{4}$ . This gives

$$\mathbb{P} \left[ \sum_{i=1}^d x_i^2 - d \geq t \right] \leq \exp \left( -\frac{t^2}{8d} \right),$$

as claimed.

*Case 2:  $t \geq d$ .* In this case, we set  $\lambda := \frac{1}{4}$ . This gives

$$\mathbb{P} \left[ \sum_{i=1}^d x_i^2 - d \geq t \right] \leq \exp \left( -\frac{t}{4} + \frac{d}{8} \right) \leq \exp \left( -\frac{t}{8} \right),$$

using that  $t \geq d$ . □

The heart of the calculation, you can convince yourself, is that  $\psi(\lambda) = O(\lambda^2)$  for  $\lambda$  smaller than some constant. Note that the bound on  $\lambda$  is required, since  $\psi(\lambda)$  becomes infinite for  $\lambda \geq \frac{1}{2}$  in our case. This property, with the bound on  $\lambda$ , is called a random variable's being *subexponential*, a weaker version of the property of being *subgaussian* that you might have encountered before.

The result of the Lemma is characteristic of sums of i.i.d. subexponential random variables: they have “Gaussian tails” scaling as  $\exp(-ct^2/d)$  up to a certain cutoff of  $t \sim d$ , beyond which they only have “exponential tails” scaling as  $\exp(-ct)$ . *Bernstein's inequality* is the general tool expressing this behavior; see Chapter 1 of [RH17] for more on that. A rough intuitive explanation for this is that, when  $x_i$  themselves have exponential tails, then large deviations of  $\sum_{i=1}^d x_i$  of order  $t \ll d$  are driven by the  $x_i$  each being slightly unusually large, while large deviations of order  $t \gg d$  are driven by the largest of  $x_i$  being unusually large, whereby the tail behavior becomes the same as that of an individual  $x_i$ .

## 1.4 CONSEQUENCES FOR HIGH-DIMENSIONAL GEOMETRY

Lemma 1.3.1 is one of the fundamental facts of high-dimensional geometry. Let us unpack its geometric meaning: informally, the result says that  $\|\mathbf{x}\|^2 = d + O(\sqrt{d})$  with high probability. Taking square roots and a Taylor approximation, we see that

$$\|\mathbf{x}\| = \sqrt{d + O(\sqrt{d})} = \sqrt{d} \cdot \sqrt{1 + O\left(\frac{1}{\sqrt{d}}\right)} = \sqrt{d} \left(1 + O\left(\frac{1}{\sqrt{d}}\right)\right) = \sqrt{d} + O(1).$$

That is, a random standard Gaussian vector usually falls close to the spherical shell of width  $O(1)$  around the sphere of radius  $\sqrt{d}$ .

This is quite counterintuitive if you have not seen it before: we think of a one- or two-dimensional Gaussian as having its “typical set” being a solid blob around the origin.<sup>3</sup> But, a high-dimensional Gaussian actually has a *non-convex* typical set, consisting of a hollow spherical shell!<sup>4</sup>

Some other “paradoxes” of high-dimensional geometry that you might be aware of may be explained in the same way. For instance, consider  $\mathbf{x} \sim \text{Unif}([-1, +1])^{\otimes d}$ , a uniformly random point in the solid box  $[-1, +1]^d$ . An integral calculation gives  $\mathbb{E}\|\mathbf{x}\|^2 = d\mathbb{E}x_1^2 = d/3$ . Similar concentration inequalities to the Gaussian case show that  $\|\mathbf{x}\| = \sqrt{d/3} + O(1)$  with high probability. In particular, then, unlike in low dimension, almost all of the mass of this

<sup>3</sup>Let us not try to define a precise notion of typical set; think of it as meaning the smallest set in which a random variable falls with high probability.

<sup>4</sup>Sometimes such random vectors are called *thin-shell* in convex geometry and probability.



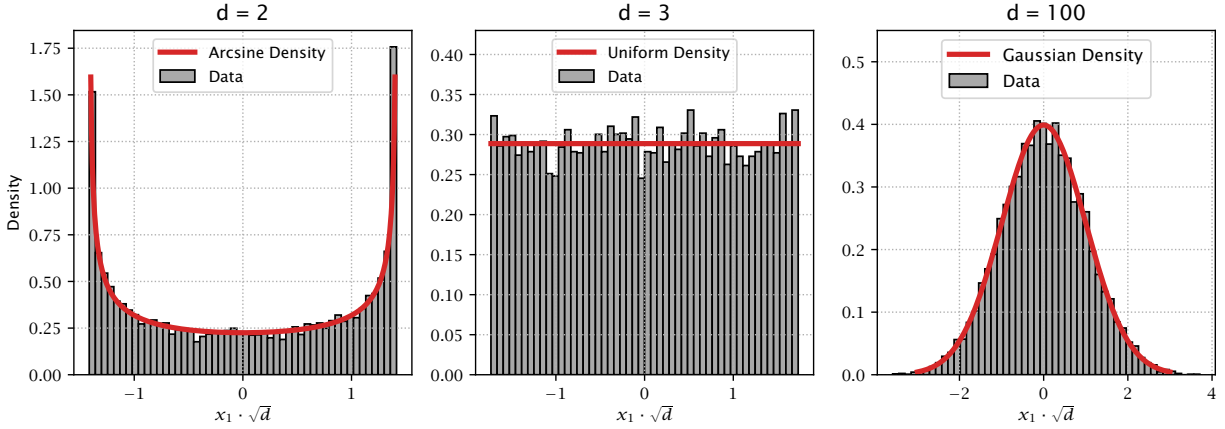


Figure 1.2: An illustration of Borel's theorem, Theorem 1.4.1, showing two low-dimensional examples  $d = 2, 3$  where the distribution of  $x_1 \cdot \sqrt{d}$  is known exactly (and is not Gaussian) as well as the Gaussian limit for large  $d$ .

box lies very close to its corners, far from the inscribed solid unit ball  $\mathbb{B}^d := \{\mathbf{x} \in \mathbb{R}^d : \|\mathbf{x}\| \leq 1\}$ . Using the analogous concentration inequality also implies that the ball  $\mathbb{B}^d$  relative to its circumscribing box  $[-1, +1]^d$  is bounded by

$$\begin{aligned}
 \frac{\text{vol}(\mathbb{B}^d)}{\text{vol}([-1, +1]^d)} &= \mathbb{P}_{\mathbf{x} \sim \text{Unif}([-1, +1]^d)} [\|\mathbf{x}\| \leq 1] \\
 &\leq \mathbb{P}_{\mathbf{x} \sim \text{Unif}([-1, +1]^d)} [|\|\mathbf{x}\|^2 - d| \geq d - 1] \\
 &= \exp(-\Omega(d)).
 \end{aligned} \tag{1.4.1}$$

Thus, these same probabilistic tools show that the volume of the unit sphere is exponentially smaller than the volume of its circumscribing box.<sup>5</sup>

Another way to look at our result is as implying the following approximation of laws, which we emphasize is an informal one:

$$“ \mathcal{N}(\mathbf{0}, \mathbf{I}_d) \approx \text{Unif}(\mathbb{S}^{d-1}(\sqrt{d})) ”$$

(1.4.2)

We arrive at this heuristic by combining our observations that, when  $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ , then  $\text{Law}(\mathbf{x}/\|\mathbf{x}\|) = \text{Unif}(\mathbb{S}^{d-1})$  (Corollary 1.2.5) and  $\|\mathbf{x}\| \approx \sqrt{d}$  (Lemma 1.3.1). This idea can be made rigorous, for instance in the following result.

**Theorem 1.4.1** (Borel's central limit theorem). *Let  $\mathbf{x} \sim \text{Unif}(\mathbb{S}^{d-1})$ . Then,  $x_1 \cdot \sqrt{d}$  converges in distribution as  $d \rightarrow \infty$  to  $\mathcal{N}(0, 1)$ . In fact, for any fixed  $k \geq 1$ ,  $(x_1 \cdot \sqrt{d}, \dots, x_k \cdot \sqrt{d})$  converges in distribution as  $d \rightarrow \infty$  to  $\mathcal{N}(\mathbf{0}, \mathbf{I}_k)$ .*

<sup>5</sup>Actually, as you may look up, the true volume is even smaller than that, but to show that requires being more careful than simply mimicking the proof of the bound for the Gaussian case in Lemma 1.3.1.

## 2 | FIRST STEPS WITH RECTANGULAR MATRICES

We now take our first steps into random matrix theory, using the results from the previous chapter to study the behavior of asymmetric or rectangular matrices with i.i.d. Gaussian entries. In particular, we will focus on the random matrix model

$$\mathbf{G} \sim \mathcal{N}(0, 1)^{\otimes d \times m}, \quad (2.0.1)$$

a  $d \times m$  random matrix with i.i.d. Gaussian entries.

### 2.1 GAUSSIAN RANDOM FIELD VIEWPOINT

Let us start by trying to gain some intuition about what  $\mathbf{G}$  does to a single  $\mathbf{y} \in \mathbb{R}^m$  as well as several given  $\mathbf{y}_1, \dots, \mathbf{y}_n \in \mathbb{R}^m$ . Immediately our observations about Gaussian random vectors become very useful, since we can easily completely characterize the laws of these images.

**Proposition 2.1.1.**  $\text{Law}(\mathbf{G}\mathbf{y}) = \mathcal{N}(0, \|\mathbf{y}\|^2 \mathbf{I}_d)$ .

*Proof.* Since  $\mathbf{G}\mathbf{y}$  is linear in  $\mathbf{G}$ , whose entries form a Gaussian random vector,  $\mathbf{G}\mathbf{y}$  is itself a Gaussian random vector. So, it suffices to compute its mean and covariance:

$$\begin{aligned} \mathbb{E}(\mathbf{G}\mathbf{y})_i &= \mathbb{E} \sum_{a=1}^m G_{ia} y_a \\ &= 0, \\ \text{Cov}((\mathbf{G}\mathbf{y})_i, (\mathbf{G}\mathbf{y})_j) &= \mathbb{E}(\mathbf{G}\mathbf{y})_i (\mathbf{G}\mathbf{y})_j \\ &= \mathbb{E} \left( \sum_{a=1}^m G_{ia} y_a \right) \left( \sum_{b=1}^m G_{jb} y_b \right) \\ &= \sum_{a,b=1}^m y_a y_b \cdot \mathbb{E} G_{ia} G_{jb} \\ &= \begin{cases} 0 & \text{if } i \neq j, \\ \|\mathbf{y}\|^2 & \text{if } i = j \end{cases}, \end{aligned}$$

where we have used that the  $G_{ia}$  are i.i.d. and have mean zero.

Another approach is to work with matrices; it is good to get used to using matrices inside of expectations in the way written below. For the mean, the following operation is valid in general when an expectation of a matrix product only has one random matrix involved:

$$\mathbb{E}[\mathbf{G}\mathbf{y}] = \mathbb{E}[\mathbf{G}]\mathbf{y} = \mathbf{0}. \quad (2.1.1)$$

You may check that this is just another form of the linearity of expectation. For the covariance, if we write  $\mathbf{g}_1, \dots, \mathbf{g}_m \in \mathbb{R}^d$  for the columns of  $\mathbf{G}$ , so that  $\mathbf{g}_a \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$  are i.i.d. standard Gaussian random vectors, we may treat the entire covariance using similar manipulations:

$$\begin{aligned} \text{Cov}(\mathbf{G}\mathbf{y}) &= \mathbb{E}(\mathbf{G}\mathbf{y})(\mathbf{G}\mathbf{y})^\top \\ &= \mathbb{E} \left( \sum_{a=1}^m \mathcal{Y}_a \mathbf{g}_a \right) \left( \sum_{b=1}^m \mathcal{Y}_b \mathbf{g}_b^\top \right) \\ &= \sum_{a,b=1}^m \mathcal{Y}_a \mathcal{Y}_b \cdot \mathbb{E} \mathbf{g}_a \mathbf{g}_b^\top \\ &= \sum_{a=1}^m \mathcal{Y}_a^2 \cdot \mathbf{I}_d \\ &= \|\mathbf{y}\|^2 \mathbf{I}_d, \end{aligned}$$

now using that the  $\mathbf{g}_a$  are i.i.d. and have mean zero and covariance identity.  $\square$

By the same kinds of calculations, we may also deduce the joint law of any collection of finitely many  $\mathbf{G}\mathbf{y}_1, \dots, \mathbf{G}\mathbf{y}_n$ .

**Proposition 2.1.2.** *Let  $\mathbf{y}_1, \dots, \mathbf{y}_n \in \mathbb{R}^m$ . Then,*

$$\text{Law} \left( \begin{bmatrix} \mathbf{G}\mathbf{y}_1 \\ \vdots \\ \mathbf{G}\mathbf{y}_n \end{bmatrix} \right) = \mathcal{N} \left( \mathbf{0}, \begin{bmatrix} \|\mathbf{y}_1\|^2 \mathbf{I}_d & \langle \mathbf{y}_1, \mathbf{y}_2 \rangle \mathbf{I}_d & \cdots & \langle \mathbf{y}_1, \mathbf{y}_n \rangle \mathbf{I}_d \\ \langle \mathbf{y}_2, \mathbf{y}_1 \rangle \mathbf{I}_d & \|\mathbf{y}_2\|^2 \mathbf{I}_d & \cdots & \langle \mathbf{y}_2, \mathbf{y}_n \rangle \mathbf{I}_d \\ \vdots & \vdots & \ddots & \vdots \\ \langle \mathbf{y}_n, \mathbf{y}_1 \rangle \mathbf{I}_d & \langle \mathbf{y}_n, \mathbf{y}_2 \rangle \mathbf{I}_d & \cdots & \|\mathbf{y}_n\|^2 \mathbf{I}_d \end{bmatrix} \right). \quad (2.1.2)$$

Written more concisely, if  $\mathbf{Y} \in \mathbb{R}^{m \times n}$  has the  $\mathbf{y}_i$  as its columns, then the covariance matrix above is  $(\mathbf{Y}^\top \mathbf{Y}) \otimes \mathbf{I}_d$ .

*Proof.* The result follows from identical calculations to the previous ones in each block of the covariance matrix.  $\square$

The previous result fully characterizes the *Gaussian random field* (a term for a higher-dimensional Gaussian process) that attaches the random vector  $\mathbf{G}\mathbf{y} \in \mathbb{R}^d$  to each point  $\mathbf{y} \in \mathbb{R}^m$ , since we have characterized the joint law of any finitely many points of this field.

From this characterization, we may start to probe more particular geometric properties of the action of  $\mathbf{G}$ . We may, for example, compute what multiplication by  $\mathbf{G}$  does to expected lengths:

$$\mathbb{E} \|\mathbf{G}\mathbf{y}\|^2 = \text{Tr}(\text{Cov}(\mathbf{G}\mathbf{y})) = d \|\mathbf{y}\|^2. \quad (2.1.3)$$

What about angles? By a similar calculation you can find

$$\mathbb{E}\langle \mathbf{G}\mathbf{y}_1, \mathbf{G}\mathbf{y}_2 \rangle = \text{Tr}(\text{Cov}(\mathbf{G}\mathbf{y}_1, \mathbf{G}\mathbf{y}_2)) = d\langle \mathbf{y}_1, \mathbf{y}_2 \rangle, \quad (2.1.4)$$

(Actually, this also follows directly from the calculation for distances by the polarization identity  $\langle \mathbf{y}_1, \mathbf{y}_2 \rangle = \frac{1}{4}\|\mathbf{y}_1 + \mathbf{y}_2\|^2 - \frac{1}{4}\|\mathbf{y}_1 - \mathbf{y}_2\|^2$ . Try working it out.) Indeed,  $\mathbf{G}$  acts in the same way on the entire *Gram matrix* of any finite collection of vectors in expectation: for  $\mathbf{y}_1, \dots, \mathbf{y}_n \in \mathbb{R}^m$  organized into the columns of  $\mathbf{Y}$ , we have

$$\mathbb{E}(\mathbf{G}\mathbf{Y})^\top (\mathbf{G}\mathbf{Y}) = d\mathbf{Y}^\top \mathbf{Y}. \quad (2.1.5)$$

**Remark 2.1.3** (Gram matrices). *The Gram matrix is a fundamental geometric object that does not always get its due in a linear algebra class. For a set of vectors  $\mathbf{y}_i$  as above,  $\mathbf{Y}^\top \mathbf{Y}$  describes the entire relative geometry of this set. Looking at the entries, we see that the Gram matrix contains the lengths of the  $\mathbf{y}_i$  (on the diagonal) and the angles between each pair (on the off-diagonal). Actually, this information fully specifies the geometry of this set of vectors: if any other  $\mathbf{y}'_1, \dots, \mathbf{y}'_n \in \mathbb{R}^m$  have the same Gram matrix, then there must be an orthogonal  $\mathbf{Q} \in \mathcal{O}(m)$  such that  $\mathbf{Q}\mathbf{y}_i = \mathbf{y}'_i$ ; that is, it is possible to get from one point cloud to the other by rotations and reflections.*

The above calculations suggest that we should normalize away the factors of  $d$  throughout by considering instead

$$\widehat{\mathbf{G}} := \frac{1}{\sqrt{d}}\mathbf{G}, \quad \text{Law}(\widehat{\mathbf{G}}) = \mathcal{N}\left(0, \frac{1}{d}\right)^{\otimes d \times m}. \quad (2.1.6)$$

This matrix will satisfy the appealing properties

$$\mathbb{E}\|\widehat{\mathbf{G}}\mathbf{y}\|^2 = \|\mathbf{y}\|^2, \quad (2.1.7)$$

$$\mathbb{E}\langle \widehat{\mathbf{G}}\mathbf{y}_1, \widehat{\mathbf{G}}\mathbf{y}_2 \rangle = \langle \mathbf{y}_1, \mathbf{y}_2 \rangle, \quad (2.1.8)$$

$$\mathbb{E}(\widehat{\mathbf{G}}\mathbf{Y})^\top (\widehat{\mathbf{G}}\mathbf{Y}) = \mathbf{Y}^\top \mathbf{Y}. \quad (2.1.9)$$

That is, in expectation,  $\widehat{\mathbf{G}}$  preserves lengths, angles, and Gram matrices, and thus looks as if it preserves all of the geometry of  $\mathbb{R}^m$ ! We will see in short order that this is true in a specific quantitative sense with high probability as well, making multiplication by  $\widehat{\mathbf{G}}$  a very useful algorithmic and computational tool.

First, though, let us clarify an important caveat. Again, the above makes it look like  $\widehat{\mathbf{G}}$  preserves the relative geometry of any given point cloud  $\mathbf{y}_1, \dots, \mathbf{y}_n \in \mathbb{R}^m$  while mapping it to  $\mathbb{R}^d$ . But, this holds for any  $d$ , including  $d \ll m$ , and even  $d = 1$ ! Thus, part of the above optimistic intuition is a kind of illusion created by taking expectations.

There are two important cautions about this intuition. First, once  $d \ll m$ , once we draw  $\widehat{\mathbf{G}}$  at random, of course we can find  $\mathbf{y} \neq \mathbf{0}$  such that  $\widehat{\mathbf{G}}\mathbf{y} = \mathbf{0}$  (any  $\mathbf{y}$  in the kernel of  $\widehat{\mathbf{G}}$ ), and for such  $\mathbf{y}$  we will have  $0 = \|\widehat{\mathbf{G}}\mathbf{y}\| \ll \|\mathbf{y}\|$ . Thus,  $\widehat{\mathbf{G}}$  cannot and will not preserve the relative geometry of *any* vectors we multiply it by, and in particular we can easily find vectors *depending* on  $\widehat{\mathbf{G}}$  whose relative geometry it badly distorts. On the other hand, if we *first* pick  $\mathbf{y}$  and then draw a random  $\widehat{\mathbf{G}}$ , then we have reason to hope that  $\|\widehat{\mathbf{G}}\mathbf{y}\|$  will concentrate

around the value  $\|\mathbf{y}\|$  suggested by the previous calculations; indeed, our calculations show that  $\text{Law}(\widehat{\mathbf{G}}\mathbf{y}) = \mathcal{N}(\mathbf{0}, \frac{1}{d}\|\mathbf{y}\|^2\mathbf{I}_d)$ , to which Lemma 1.3.1 should apply.

Second, even provided we make sure  $\widehat{\mathbf{G}}$  is independent of the  $\mathbf{y}$  that we multiply by it, the extent to which the above properties hold with high probability rather than merely in expectation will depend on how large  $d$  is: when  $d = 1$ , then  $\widehat{\mathbf{G}} = \mathbf{g}^\top$  for  $\mathbf{g} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ , and we still have  $\mathbb{E}\|\widehat{\mathbf{G}}\mathbf{y}_i\|^2 = \mathbb{E}\langle \mathbf{g}, \mathbf{y}_i \rangle^2 = \|\mathbf{y}_i\|^2$ . But, it should be intuitively clear that we cannot actually have  $\langle \mathbf{g}, \mathbf{y}_i \rangle^2 \approx \|\mathbf{y}_i\|^2$  simultaneously for very many  $\mathbf{y}_i$ , even if we choose them in advance of drawing  $\widehat{\mathbf{G}}$ . (Exercise: Try constructing concrete counterexamples.)

We will see next that what we have mentioned above are the only obstructions to  $\widehat{\mathbf{G}}$  behaving as we wish: provided that (1)  $d$  is not too small and (2) that we fix a point cloud  $\mathbf{y}_1, \dots, \mathbf{y}_n$  in advance of drawing  $\widehat{\mathbf{G}}$ , we will have that multiplication by  $\widehat{\mathbf{G}}$  approximately preserves the relative geometry of the  $\mathbf{y}_i$  (that is,  $\widehat{\mathbf{G}}$  acts as an *approximate isometry* on these points), even for  $d \ll m$ .

## 2.2 APPLICATION: DIMENSIONALITY REDUCTION AND JOHNSON-LINDENSTRAUSS LEMMA

Validating our heuristic reasoning above, the Johnson-Lindenstrauss lemma states that, in a certain sense, a deterministic point cloud that is typically deformed substantially by  $\widehat{\mathbf{G}}$  *must* be quite large relative to the dimension we are mapping it into, so the above obstruction to  $\widehat{\mathbf{G}}$  preserving geometric structure is the only one. However, and this is an important condition that we will return to later, this result only concerns whether  $\widehat{\mathbf{G}}$  deforms *pairwise* distances between the  $\mathbf{y}_i$ . This is not the same as preserving all aspects of the global geometry of the  $\mathbf{y}_i$ , but is an intuitive notion and is relevant to various applications.

**Definition 2.2.1.** Given  $\mathbf{y}_1, \dots, \mathbf{y}_n \in \mathbb{R}^m$ , a function  $f : \mathbb{R}^m \rightarrow \mathbb{R}^d$  is pairwise  $\epsilon$ -faithful<sup>1</sup> on the  $\mathbf{y}_i$  if, for all  $i, j \in [n]$ ,

$$(1 - \epsilon)\|\mathbf{y}_i - \mathbf{y}_j\|^2 \leq \|f(\mathbf{y}_i) - f(\mathbf{y}_j)\|^2 \leq (1 + \epsilon)\|\mathbf{y}_i - \mathbf{y}_j\|^2. \quad (2.2.1)$$

**Theorem 2.2.2** (Johnson-Lindenstrauss [JL82, IM98]). Let  $\mathbf{y}_1, \dots, \mathbf{y}_n \in \mathbb{R}^m$  be arbitrary and  $\epsilon \in (0, 1)$ . Suppose that

$$d \geq 24 \frac{\log n}{\epsilon^2}. \quad (2.2.2)$$

Note that the lower bound on  $d$  does not depend on the dimension  $m$  at all! Then, with probability at least  $1 - \frac{1}{n^k}$ , multiplication by  $\widehat{\mathbf{G}} \sim \mathcal{N}(\mathbf{0}, \frac{1}{d}\mathbf{I}_d)^{\otimes d \times m}$  is pairwise  $\epsilon$ -faithful on the  $\mathbf{y}_i$ .

In fact, as the proof will make clear, by increasing the constant  $C$  in the assumption  $d \geq C \frac{\log n}{\epsilon^2}$  to some  $C = C(K)$ , you can achieve a success probability  $1 - \frac{1}{n^k}$  for any given  $K > 0$ .

*Proof.* The proof actually has little to do with the  $\mathbf{y}_i$ . Notice that being pairwise  $\epsilon$ -faithful is a matter of preserving  $\binom{n}{2}$  many vector norms to within a factor in  $[1 - \epsilon, 1 + \epsilon]$ . We will

---

<sup>1</sup>My own non-standard terminology.

show that this is true with high probability for *any* choice of  $\binom{n}{2}$  vectors, and therefore also for the  $\widehat{\mathbf{G}}(\mathbf{y}_i - \mathbf{y}_j) = \widehat{\mathbf{G}}\mathbf{y}_i - \widehat{\mathbf{G}}\mathbf{y}_j$ .

Consider an arbitrary  $\mathbf{x} \in \mathbb{R}^m$ . Recall that  $\text{Law}(\widehat{\mathbf{G}}\mathbf{x}) = \mathcal{N}(0, \frac{1}{d}\|\mathbf{x}\|^2 \mathbf{I}_d)$ . We then have

$$\begin{aligned} \mathbb{P}_{\widehat{\mathbf{G}} \sim \mathcal{N}(0, \frac{1}{d})^{\otimes d \times m}} \left[ \left| \|\widehat{\mathbf{G}}\mathbf{x}\|^2 - \|\mathbf{x}\|^2 \right| > \epsilon \|\mathbf{x}\|^2 \right] &= \mathbb{P}_{\mathbf{g} \sim \mathcal{N}(0, \frac{1}{d}\|\mathbf{x}\|^2 \mathbf{I}_d)} \left[ \left| \|\mathbf{g}\|^2 - \|\mathbf{x}\|^2 \right| > \epsilon \|\mathbf{x}\|^2 \right] \\ &= \mathbb{P}_{\mathbf{g} \sim \mathcal{N}(0, \mathbf{I}_d)} \left[ \left| \frac{1}{d}\|\mathbf{x}\|^2 \cdot \|\mathbf{g}\|^2 - \|\mathbf{x}\|^2 \right| > \epsilon \|\mathbf{x}\|^2 \right] \\ &= \mathbb{P}_{\mathbf{g} \sim \mathcal{N}(0, \mathbf{I}_d)} \left[ \left| \|\mathbf{g}\|^2 - d \right| > \epsilon d \right] \end{aligned}$$

and, using that  $\epsilon < 1$  whereby  $\epsilon d < d$ , by Lemma 1.3.1 we have

$$\leq 2 \exp \left( -\frac{1}{8} \epsilon^2 d \right). \quad (2.2.3)$$

This gives a bound on the probability of embedding a single vector with low distortion of the length:

$$\mathbb{P}_{\widehat{\mathbf{G}} \sim \mathcal{N}(0, \frac{1}{d})^{\otimes d \times m}} \left[ (1 - \epsilon) \|\mathbf{x}\|^2 \leq \|\widehat{\mathbf{G}}\mathbf{x}\|^2 \leq (1 + \epsilon) \|\mathbf{x}\|^2 \right] \geq 1 - 2 \exp \left( -\frac{1}{8} \epsilon^2 d \right). \quad (2.2.4)$$

Finally, let  $E_{ij}$  be the event that  $(1 - \epsilon) \|\mathbf{y}_i - \mathbf{y}_j\|^2 \leq \|\widehat{\mathbf{G}}\mathbf{y}_i - \widehat{\mathbf{G}}\mathbf{y}_j\|^2 \leq (1 + \epsilon) \|\mathbf{y}_i - \mathbf{y}_j\|^2$ . Then, by the union bound we have that

$$\begin{aligned} \mathbb{P}[\text{some } E_{ij} \text{ does not occur}] &\leq \binom{n}{2} \cdot 2 \exp \left( -\frac{1}{8} \epsilon^2 d \right) \\ &\leq n^2 \exp \left( -\frac{1}{8} \epsilon^2 d \right) \\ &= \frac{1}{n} \exp \left( 3 \log n - \frac{1}{8} \epsilon^2 d \right), \end{aligned} \quad (2.2.5)$$

and the result follows since, if  $d \geq 24 \frac{\log n}{\epsilon^2}$ , then this is at most  $1/n$ .  $\square$

Multiplication by a random matrix for the purposes of dimensionality reduction is sometimes called the *Johnson-Lindenstrauss transform (JLT)* after this result. The original JLT and its algorithmic applications sparked many directions of further research. Let us briefly go over a few of these.

### 2.2.1 SIMPLE EXTENSIONS

It is easy to extend the proof to show that we also have  $(1 - \epsilon) \|\mathbf{y}_i\| \leq \|f(\mathbf{y}_i)\| \leq (1 + \epsilon) \|\mathbf{y}_i\|$  with the same high probability up to constants (say, by increasing  $n$  by one and adding  $\mathbf{0}$  to the  $\mathbf{y}_i$ , or by including further appropriate “bad events” in the union bound argument above). As mentioned before, it is also possible to get an error bound of  $1/n^K$  for any  $K > 0$  by suitably increasing the constant 24 in the statement of the Theorem and adjusting the last step accordingly.

### 2.2.2 LOWER BOUNDS

It is natural to ask if the result is optimal: could it be possible to reduce  $d$  even further? In fact it is not. The original paper of Johnson and Lindenstrauss [JL82] showed that dimension  $\Omega(\log n)$  is needed for any given constant  $\epsilon$ , but with no dependence on  $\epsilon$ . Only relatively recently was it shown that the Johnson-Lindenstrauss result is actually tight up to constants. This was first shown for embeddings produced by applying a linear map to the  $\mathbf{y}_i$  (i.e., multiplying them by a matrix) [LN14], and finally extended to arbitrary embeddings [LN17].

### 2.2.3 SPARSE AND FAST JLTs

Another line of work has explored how to speed up the matrix-vector multiplication required to form the embeddings  $f(\mathbf{y}_i) = \widehat{\mathbf{G}}\mathbf{y}_i$ . All such work is basically predicated on finding special structured matrices  $\widehat{\mathbf{G}}$  that allow for faster matrix-vector multiplication. The simplest such structure is *sparsity* of the matrix involved; see [DKS10, KN14] for some results in this direction. Perhaps the most elegant idea studied in this direction is quite different: the algorithm of [AC09] uses that the fast Fourier transform algorithm may be used to multiply quickly by the discrete Fourier transform matrix. Their version of the JLT is based on performing this multiplication (which does not reduce the dimension of the vector, as the discrete Fourier transform matrix is square) and then subsampling a random subset of entries of the result. There is a further caveat that some vectors have sparse Fourier transforms, so subsampling their Fourier transform runs a risk of ending up with an unusually small vector. Thus, a further preprocessing step is required to avoid this particular set of “bad” vectors; see the paper for more details.

### 2.2.4 PROOF TECHNIQUE: FIRST MOMENT METHOD

As we have seen, the proof technique is merely an application of the union bound. However, it is perhaps worth emphasizing some aspects of how specifically we are using the union bound that also hold in many modern applications. While the union bound might seem trivial, and indeed is when applied to just a few events, the union bound when applied to many carefully chosen events is actually a quite sophisticated tool.<sup>2</sup> If one wishes to bound the probability of a “bad event”  $E$ , one may decompose it into  $E = E_1 \cup \dots \cup E_N$  and bound

$$\mathbb{P}[E] \leq \sum_{i=1}^N \mathbb{P}[E_i] \leq N \cdot \max_{i=1}^N \mathbb{P}[E_i]. \quad (2.2.6)$$

Often, the  $\mathbb{P}[E_i]$  are roughly the same. Then, the above becomes a clear “competition” between  $N$ , the number of ways that a bad event can happen, and  $\mathbb{P}[E_i]$ , the probability of a given way that it can happen.<sup>3</sup> In high-dimensional probability, often  $N$  is very large and

---

<sup>2</sup>See in particular the book [Tal05], whose results show that the values of a tremendous number of random optimization-like quantities including the norms of very general Gaussian random matrices can, in principle, be controlled accurately by nothing more than a sufficiently careful use of the union bound.

<sup>3</sup>Borrowing the language of statistical physics, these quantities are sometimes called “entropy” and “energy,” respectively.

$\mathbb{P}[E_i]$  is very small, leading to the kind of competition on an exponential scale that we saw in the proof.

As an aside, it is also worth being aware that this approach is also often called the “first moment method,” because one may phrase the bound as

$$\mathbb{P}[\text{some } E_i \text{ occurs}] = \mathbb{P}\left[\sum_{i=1}^N \mathbb{1}\{E_i \text{ occurs}\} \geq 1\right] \leq \mathbb{E}\left[\sum_{i=1}^N \mathbb{1}\{E_i \text{ occurs}\}\right] = \sum_{i=1}^N \mathbb{P}[E_i] \quad (2.2.7)$$

by Markov’s inequality. Thus one is computing the first moment (or expectation) of the random variable  $\#\{i : E_i \text{ occurs}\}$ . There is also a widely used “second moment method” which involves computing the second moment and thus variance of this random variable, using which you can show that some  $E_i$  *does* occur with high probability. For instance, you could use this to show that, for a suitable point cloud, if  $d$  is too small then multiplication by  $\widehat{\mathbf{G}}$  is *not* pairwise  $\epsilon$ -faithful. (Note that this does not answer the question about lower bounds posed above, since it only applies to this particular method of embedding by multiplication by a random matrix.)

## 2.3 RECTANGULAR GAUSSIAN MATRICES: SINGULAR VALUES

Let us next see how to extend our approach to the JLT to understand more intrinsic properties of  $\widehat{\mathbf{G}}$ , namely its singular values and vectors. This will both give us a first treatment of our first fundamental random matrix model—that of rectangular matrices with i.i.d. Gaussian entries—and will hopefully further clarify why the JLT works well. We focus on “short fat” matrices, that is, the setting  $d \leq m$  (often we will also mention how our results specialize to the asymptotic regime  $d \ll m$ ). We also return to working with  $\mathbf{G} \sim \mathcal{N}(0, 1)^{\otimes d \times m}$  again in order to simplify the notation.

As we have mentioned, the same argument as used in the proof of Theorem 2.2.2 actually implies the following more general result.

**Theorem 2.3.1.** *There is a constant  $C > 0$  such that the following holds. For any  $\mathbf{x}_1, \dots, \mathbf{x}_N \in \mathbb{R}^m$ , if  $d \geq C \frac{\log N}{\epsilon^2}$  and  $\mathbf{G} \sim \mathcal{N}(0, 1)^{\otimes d \times m}$ , then*

$$\mathbb{P}\left[\left|\|\mathbf{G}\mathbf{x}_i\|^2 - d\|\mathbf{x}_i\|^2\right| \leq \epsilon d\|\mathbf{x}_i\|^2 \text{ for all } i \in [N]\right] \geq 1 - \frac{1}{N}. \quad (2.3.1)$$

Before, we took  $N \asymp n^2$  and the  $\mathbf{x}_i$  to be all possible  $\mathbf{y}_j - \mathbf{y}_k$ , but we can apply the same tools to any  $\mathbf{x}_i$ .

We may reframe this by rewriting

$$\left|\|\mathbf{G}\mathbf{x}_i\|^2 - \|\mathbf{x}_i\|^2\right| = \left|\mathbf{x}_i^\top (\mathbf{G}^\top \mathbf{G}) \mathbf{x}_i - \mathbf{x}_i^\top (d\mathbf{I}_m) \mathbf{x}_i\right| = \left|\mathbf{x}_i^\top (\mathbf{G}^\top \mathbf{G} - d\mathbf{I}_m) \mathbf{x}_i\right|. \quad (2.3.2)$$

Thus it seems that, from the “point of view” of a small number of deterministic quadratic form evaluations, the random matrix  $\mathbf{G}^\top \mathbf{G}$  behaves like  $d\mathbf{I}_m$ , and indeed the JLT is based on this fact. Had we stuck with our previous normalization, we would see the different normalization that  $\widehat{\mathbf{G}}^\top \widehat{\mathbf{G}}$  behaves like  $\mathbf{I}_m$ . In any case, this of course cannot be true without



restricting the number of quadratic form evaluations, since the matrices  $\mathbf{G}^\top \mathbf{G}$  and  $\widehat{\mathbf{G}}^\top \widehat{\mathbf{G}}$  have rank at most  $d \ll m$ .

This is another “paradox” (that we will resolve in short order) of high-dimensional probability and geometry: somehow, a very low-rank random matrix can actually act like the identity against a large number of test quadratic forms. This is basically just a variation on the same paradox we saw in  $\widehat{\mathbf{G}}$  preserving lengths and angles in expectation, but note that the above statement does actually state that an event happens with high probability, so it cannot merely be accounted for by saying that the random variables involved are not close enough to their expectations.

We will now work towards a more sophisticated and “more spectral” (or more intrinsic and geometric) explanation of this phenomenon than the union bound based proof we saw, which is relatively simple but maybe not very enlightening. To do that, we will undertake an analysis of the eigendecomposition of the matrix  $\mathbf{G}^\top \mathbf{G} \in \mathbb{R}_{\text{sym}}^{m \times m}$ , seeking to make claims about the distribution of its eigenvectors and eigenvalues (the same information as the right singular vectors and singular values of  $\mathbf{G}$ ).

Note that  $\mathbf{G}^\top \mathbf{G} \succeq \mathbf{0}$ , so its eigenvalues are non-negative. Call these  $\lambda_1 \geq \dots \geq \lambda_m \geq 0$ . Since  $\text{rank}(\mathbf{G}^\top \mathbf{G}) \leq d$ , we have  $\lambda_{d+1} = \dots = \lambda_m = 0$ . The remaining eigenvalues are related to the singular values of  $\mathbf{G}$  by  $\lambda_i = \sigma_i(\mathbf{G})^2$ . We begin by studying these eigenvalues, or equivalently the singular values of  $\mathbf{G}$ .

### 2.3.1 THE GEOMETRIC METHOD OF RANDOM MATRIX THEORY

It will be more convenient to instead consider the eigenvalues of  $\mathbf{G}\mathbf{G}^\top \in \mathbb{R}_{\text{sym}}^{d \times d}$ . Since the eigenvalues of  $\mathbf{G}\mathbf{G}^\top$  and  $\mathbf{G}^\top \mathbf{G}$  are equal except that the latter has  $m - d$  extra eigenvalues equal to 0,  $\mathbf{G}\mathbf{G}^\top$  will provide us with the same information; at the end, we will return to  $\mathbf{G}^\top \mathbf{G}$  and deduce some geometric intuitions about it. The probabilistic reason for working with  $\mathbf{G}\mathbf{G}^\top$  instead is that, if  $\mathbf{g}_1, \dots, \mathbf{g}_m \in \mathbb{R}^d$  are the columns of  $\mathbf{G}$  (distributed i.i.d. as  $\mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ ), then we have

$$\mathbf{G}\mathbf{G}^\top = \sum_{i=1}^m \mathbf{g}_i \mathbf{g}_i^\top$$

which is a sum of  $m \gg d$  i.i.d. random matrices of dimension  $d$ . As we will see, in such a situation it is justified to hope for a matrix-valued law of large numbers to hold, giving us

$$\begin{aligned} &\approx m \cdot \mathbb{E} \mathbf{g}_1 \mathbf{g}_1^\top \\ &= m \mathbf{I}_d \\ &= \mathbb{E} \mathbf{G}\mathbf{G}^\top. \end{aligned} \tag{2.3.3}$$

Roughly speaking, this heuristic is more accurate than the same applied to  $\mathbf{G}^\top \mathbf{G}$  where the roles of  $d$  and  $m$  are reversed because we want in our appeal to the law of large numbers to work with an expression having “as much averaging as possible.”

**Remark 2.3.2.** *To emphasize the difference in how we are discussing these two matrices: above we observed that  $\mathbf{G}^\top \mathbf{G}$  looks like  $d\mathbf{I}_m$  from the point of view of a small enough number*

of quadratic forms, but we also observed that it cannot actually be close to this matrix (for instance by simple rank considerations). On the other hand, the above heuristic suggests that  $\mathbf{G}\mathbf{G}^\top$  looks like  $m\mathbf{I}_d$ , and we will show below that this is actually true, and use that to understand the trickier behavior of  $\mathbf{G}^\top\mathbf{G}$ .

Concretely, we expect the difference of operator norms  $\|\mathbf{G}\mathbf{G}^\top - m\mathbf{I}_d\|$  to be small. This kind of statement that a sum of independent or weakly dependent random matrices is close in operator norm to its expectation is often called a *matrix concentration inequality*. Let us write  $\Delta := \mathbf{G}\mathbf{G}^\top - m\mathbf{I}_d$ .

**Remark 2.3.3** (Statistical interpretation). *Expanding in a sum of rank one terms again, we may also view*

$$\frac{1}{m}\Delta = \frac{1}{m} \sum_{i=1}^m \mathbf{g}_i \mathbf{g}_i^\top - \mathbb{E} \mathbf{g}_1 \mathbf{g}_1^\top \quad (2.3.4)$$

as the difference between a sample covariance and the true covariance of several i.i.d. draws  $\mathbf{g}_1, \dots, \mathbf{g}_m$  from  $\mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ . Thus, our investigation will also address the natural statistical question of how good of a statistical estimator the sample covariance is (in this particular case).

How can we show that  $\|\Delta\|$  is relatively small? The chief trouble is that the operator norm is a quantity whose definition involves an uncountable quantification over test vectors. We adopt the following notation for closed balls in  $\mathbb{R}^d$ , which will be useful soon:

$$\mathbb{B}^d(\mathbf{x}, r) := \{\mathbf{y} \in \mathbb{R}^d : \|\mathbf{y} - \mathbf{x}\| \leq r\}, \quad (2.3.5)$$

$$\mathbb{B}^d := \mathbb{B}^d(\mathbf{0}, 1) \supset \mathbb{S}^{d-1}. \quad (2.3.6)$$

The operator norm can be formulated in either of the equivalent ways

$$\|\Delta\| = \sup_{\mathbf{x} \in \mathbb{S}^{d-1}} |\mathbf{x}^\top \Delta \mathbf{x}| = \sup_{\mathbf{x} \in \mathbb{B}^d} |\mathbf{x}^\top \Delta \mathbf{x}|. \quad (2.3.7)$$

In either case, if  $\mathbb{S}^{d-1}$  or  $\mathbb{B}^d$  were replaced by a finite set, we could use the first moment method, bounding  $\mathbb{P}[|\mathbf{x}^\top \Delta \mathbf{x}| > t]$  (which as we will see boils down to yet another application of Lemma 1.3.1 on the concentration of norms of Gaussian random vectors) and using a union bound to control  $\mathbb{P}[\|\Delta\| > t]$ , but this is not the case.

But, it will turn out to be useful to take this hypothetical situation more seriously, so let us define notation for these notions.

**Definition 2.3.4.** For any bounded set  $X \subset \mathbb{R}^d$  and  $\mathbf{M} \in \mathbb{R}_{\text{sym}}^{d \times d}$ , we define

$$\|\mathbf{M}\|_X := \sup_{\mathbf{x} \in X} |\mathbf{x}^\top \mathbf{M} \mathbf{x}|. \quad (2.3.8)$$

This is not always actually a norm; depending on  $X$ , it may only be a *seminorm*, where  $\|\mathbf{M}\|_X = 0$  can hold for  $\mathbf{M} \neq \mathbf{0}$ . The definition will be useful in any case, however. Applying the first moment method to a general random matrix gives the following:

**Proposition 2.3.5.** Let  $\mathbf{M} \in \mathbb{R}_{\text{sym}}^{d \times d}$  be a random matrix and  $X \subset \mathbb{R}^d$  be a finite set. Then,

$$\mathbb{P}[\|\mathbf{M}\|_X > t] \leq |X| \cdot \max_{\mathbf{x} \in X} \mathbb{P}[|\mathbf{x}^\top \mathbf{M} \mathbf{x}| > t]. \quad (2.3.9)$$

We will apply this to  $M = \Delta$ . As mentioned above, controlling the second factor above involving scalar probabilities will be possible using Lemma 1.3.1. Thus, the key point will be the tradeoff between the factor of  $|\mathcal{X}|$  in the bound above, which makes the bound deteriorate as  $|\mathcal{X}|$  grows, versus the error in approximating  $\|\Delta\|$  by  $\|\Delta\|_{\mathcal{X}}$ , which we will soon see decreases as  $|\mathcal{X}|$  increases. We will therefore reduce the matter of bounding  $\|\Delta\|$  to merely that of finding a sufficiently small set  $|\mathcal{X}|$  that behaves like a sufficiently dense “grid” of  $\mathbb{S}^{d-1}$  or  $\mathbb{B}^d$  to approximate the operator norm accurately (the sets we work with will not actually be grids but will behave similarly). We call this general approach the **Geometric Method** of random matrix theory, because it reduces the matter of understanding norms of random matrices to (possibly more complicated variants of) the above purely geometric question.

### 2.3.2 NETS, COVERINGS, AND PACKINGS

To do this, our first task will be to develop some tools for relating  $\|M\|_{\mathcal{X}}$  for finite  $\mathcal{X}$  to  $\|M\|$ , *discretizing* the continuous optimization appearing in the operator norm.

**Definition 2.3.6** ( $\epsilon$ -net). A set  $\mathcal{X} = \{x_1, \dots, x_N\} \subseteq \mathcal{Y} \subset \mathbb{R}^d$  is an  $\epsilon$ -net or  $\epsilon$ -covering of  $\mathcal{Y}$  if

$$\mathcal{Y} \subseteq \bigcup_{i=1}^N \mathbb{B}^d(x_i, \epsilon), \quad (2.3.10)$$

or equivalently if, for all  $y \in \mathcal{Y}$ , there is  $x_i \in \mathcal{X}$  such that  $\|y - x_i\| \leq \epsilon$ .

The following key result shows that we may indeed discretize the operator norm and bound it by a maximum over a suitable net.

**Lemma 2.3.7.** Suppose that  $M \in \mathbb{R}_{\text{sym}}^{d \times d}$  and  $\mathcal{X}$  is an  $\epsilon$ -net of  $\mathbb{S}^{d-1}$  for  $\epsilon < \frac{1}{2}$ . Then,

$$\|M\|_{\mathcal{X}} \leq \|M\| \leq \frac{1}{1 - 2\epsilon} \|M\|_{\mathcal{X}}. \quad (2.3.11)$$

The same is true if  $\mathcal{X}$  is instead an  $\epsilon$ -net of  $\mathbb{B}^d$ .

*Proof.* The first inequality is immediate. For the second, let  $x \in \mathbb{S}^{d-1}$  be such that  $|x^\top M x| = \|M\|$ , and let  $x_i \in \mathcal{X}$  be such that  $\|x - x_i\| \leq \epsilon$ . We then have

$$\begin{aligned} \|M\|_{\mathcal{X}} &\geq |x_i^\top M x_i| \\ &= |(x + x_i - x)^\top M (x + x_i - x)| \\ &= |x^\top M x + x^\top M (x_i - x) + (x_i - x)^\top M x_i| \\ &\geq |x^\top M x| - |x^\top M (x_i - x)| - |(x_i - x)^\top M x_i| \\ &\geq (1 - 2\epsilon) \|M\|, \end{aligned}$$

and rearranging gives the result. The same argument applies verbatim in the case of a net of  $\mathbb{B}^d$ .  $\square$

Next, we need to develop some further tools for constructing  $\epsilon$ -nets in a way such that we have some control over their size. As you can convince yourself, it is rather tricky to construct explicit  $\epsilon$ -nets by hand by completely describing the vectors of  $\mathcal{X}$ . Fortunately, there is a more abstract approach that lets us non-constructively show that small  $\epsilon$ -nets exist. This proceeds by considering a complementary notion:

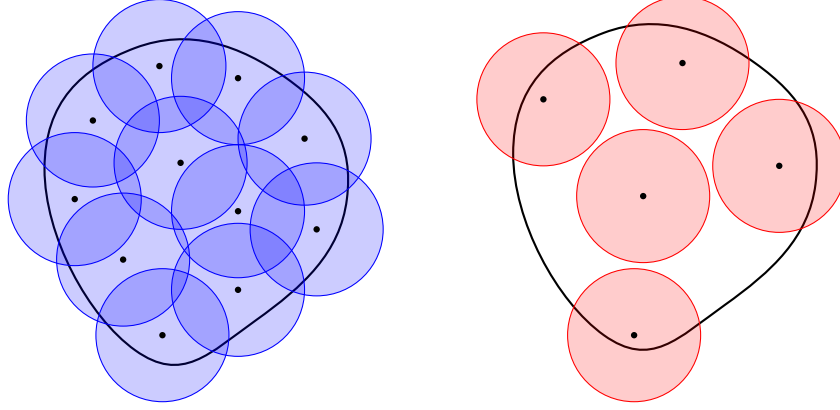


Figure 2.1: An illustration of a covering (left) and a packing (right) of a given set in  $\mathbb{R}^2$  by balls of the same radius.

**Definition 2.3.8** ( $\epsilon$ -packing). A set  $X = \{x_1, \dots, x_N\} \subseteq \mathcal{Y} \subset \mathbb{R}^d$  is an  $\epsilon$ -packing of  $\mathcal{Y}$  if the balls  $\mathbb{B}^d(x_i, \epsilon)$  are pairwise disjoint, or equivalently if  $\|x_i - x_j\| > 2\epsilon$  for all  $i, j \in [N]$  distinct.

**Proposition 2.3.9.** Let  $X$  be a maximal  $\epsilon$ -packing of  $\mathcal{Y}$ , one that cannot be extended to a larger packing by adding more points. Then,  $X$  is a  $2\epsilon$ -net of  $\mathcal{Y}$ .

*Proof.* Let  $y \in \mathcal{Y} \setminus X$ . Since adding  $y$  to  $X$  must yield a set that is not an  $\epsilon$ -packing, there must be some  $x \in X$  such that  $\mathbb{B}^d(x, \epsilon) \cap \mathbb{B}^d(y, \epsilon) \neq \emptyset$ . Letting  $z$  be in the intersection, we have  $\|x - y\| \leq \|x - z\| + \|y - z\| \leq 2\epsilon$ . Thus, either  $y \in X$  or  $y \in \mathbb{B}^d(x, 2\epsilon)$  for some  $x \in X$ , giving the result.  $\square$

Thus, to show that small nets exist, it suffices to show that small maximal packings exist. But, there is a fundamental *volumetric* bound on the size of arbitrary packings, including maximal ones: roughly speaking, the total volume of the union of the balls in a packing cannot be much bigger than the volume of the set being packed (at least for a nice full-dimensional convex set).

**Proposition 2.3.10.** Let  $X$  be any  $\epsilon$ -packing of  $\mathbb{B}^d$ . Then,  $|X| \leq (1 + \frac{1}{\epsilon})^d$ .

*Proof.* We have

$$\bigcup_{x \in X} \mathbb{B}^d(x, \epsilon) \subseteq \mathbb{B}^d(0, 1 + \epsilon), \quad (2.3.12)$$

and the union on the left-hand side is disjoint. Thus, taking volumes,

$$|X| \cdot \text{vol}(\mathbb{B}^d(0, \epsilon)) \leq \text{vol}(\mathbb{B}^d(0, 1 + \epsilon)). \quad (2.3.13)$$

Rearranging,

$$|X| \leq \frac{\text{vol}(\mathbb{B}^d(0, 1 + \epsilon))}{\text{vol}(\mathbb{B}^d(0, \epsilon))} = \left(\frac{1 + \epsilon}{\epsilon}\right)^d = \left(1 + \frac{1}{\epsilon}\right)^d, \quad (2.3.14)$$

by scaling considerations of the volumes of  $d$ -dimensional balls.  $\square$

**Corollary 2.3.11.** For any  $\epsilon > 0$ , there exists an  $\epsilon$ -net  $X$  of  $\mathbb{B}^d$  of size  $|X| \leq (1 + \frac{2}{\epsilon})^d$ .

*Proof.* Let  $X$  be a maximal  $\frac{\epsilon}{2}$ -packing of  $\mathbb{B}^d$ . By Proposition 2.3.10,  $|X| \leq (1 + \frac{2}{\epsilon})^d$ , while by Proposition 2.3.9,  $X$  is an  $\epsilon$ -net.  $\square$

### 2.3.3 COARSE NON-ASYMPTOTIC BOUND ON SINGULAR VALUES

Finally, we may put together the pieces to bound  $\|\Delta\|$ : we use the existence of a not-too-large  $\epsilon$ -net by the above non-constructive method, approximate the operator norm by discretizing to test vectors in this net, and use the first moment method.

**Theorem 2.3.12.** *If  $d \leq m$ , then*

$$\mathbb{P}[\|\Delta\| > 64\sqrt{dm}] \leq 2 \exp\left(-\frac{d}{8}\right). \quad (2.3.15)$$

*Proof.* Let us fix  $\epsilon := 1/4$ . Let  $\mathcal{X}$  be an  $\epsilon$ -net of  $\mathbb{B}^d$  as in Corollary 2.3.11, by which we can assume  $|\mathcal{X}| \leq (1 + \frac{2}{\epsilon})^d = 9^d$ . By Lemma 2.3.7, we have  $\|\Delta\| \leq \frac{1}{1-2\epsilon}\|\Delta\|_{\mathcal{X}} = 2\|\Delta\|_{\mathcal{X}}$ . Thus we may bound

$$\begin{aligned} \mathbb{P}[\|\Delta\| \geq t] &\leq \mathbb{P}\left[\|\Delta\|_{\mathcal{X}} \geq \frac{t}{2}\right] && \text{(Lemma 2.3.7)} \\ &\leq 9^d \cdot \max_{x \in \mathcal{X}} \mathbb{P}\left[|x^\top \Delta x| \geq \frac{t}{2}\right] && \text{(Proposition 2.3.5)} \end{aligned}$$

and writing  $\hat{x} := x/\|x\|$ , since  $\|x\| \leq 1$  for all  $x \in \mathcal{X}$ , we may bound

$$\leq 9^d \cdot \max_{x \in \mathcal{X}} \mathbb{P}\left[|\hat{x}^\top \Delta \hat{x}| \geq \frac{t}{2}\right]. \quad (2.3.16)$$

Let us work with this inner probability for any given  $x$  (and its corresponding normalization  $\hat{x}$ ) for a moment: expanding the definition of  $\Delta$ ,

$$\begin{aligned} \mathbb{P}\left[|\hat{x}^\top \Delta \hat{x}| \geq \frac{t}{2}\right] &= \mathbb{P}\left[|\|G^\top \hat{x}\|^2 - m| > \frac{t}{2}\right] \\ &= \mathbb{P}_{g \sim \mathcal{N}(0, I_m)}\left[|\|g\|^2 - m| > \frac{t}{2}\right]. \end{aligned} \quad (2.3.17)$$

These probabilities are all equal, not depending on  $x$ , so we have

$$\mathbb{P}[\|\Delta\| \geq t] \leq 9^d \cdot \mathbb{P}_{g \sim \mathcal{N}(0, I_m)}\left[|\|g\|^2 - m| > \frac{t}{2}\right]$$

To the remaining probability we can apply Lemma 1.3.1, which it will be convenient to do specifically in the form (1.3.5). We can also bound  $9 \leq \exp(3)$ . Together, these observations give

$$\leq 2 \exp\left(3d - \frac{1}{8} \min\left\{\frac{t^2}{4m}, \frac{t}{2}\right\}\right) \quad (2.3.18)$$

Now, let us take  $t := C\sqrt{dm}$  for some  $C$  to be chosen momentarily. Noting that  $\sqrt{dm} \geq d$  since  $d \leq m$ , we may further bound

$$\begin{aligned} &\leq 2 \exp\left(3d - \frac{1}{8} \min\left\{\frac{C^2}{4}d, \frac{C}{2}d\right\}\right) \\ &= 2 \exp\left(d \cdot \left[3 - \min\left\{\frac{C^2}{32}, \frac{C}{16}\right\}\right]\right) \end{aligned} \quad (2.3.19)$$

and finally we see that taking  $C := 64$  gives the result.  $\square$

To give a corollary in a form that is maybe easier and more intuitive to understand concerning the singular values of  $\mathbf{G}$ , we will use the following basic eigenvalue perturbation inequalities.

**Proposition 2.3.13.** *Let  $\mathbf{A}, \Delta \in \mathbb{R}_{\text{sym}}^{d \times d}$ . Then,*

$$\lambda_d(\mathbf{A}) - \|\Delta\| \leq \lambda_d(\mathbf{A} + \Delta) \leq \lambda_1(\mathbf{A} + \Delta) \leq \lambda_1(\mathbf{A}) + \|\Delta\|. \quad (2.3.20)$$

*Proof.* For both bounds, we use the variational description of the eigenvalues. For the upper bound, we have

$$\begin{aligned} \lambda_1(\mathbf{A} + \Delta) &= \sup_{\mathbf{x} \in \mathbb{S}^{d-1}} \mathbf{x}^\top (\mathbf{A} + \Delta) \mathbf{x} \\ &\leq \sup_{\mathbf{x} \in \mathbb{S}^{d-1}} (\mathbf{x}^\top \mathbf{A} \mathbf{x} + \|\Delta\|) \\ &= \left( \sup_{\mathbf{x} \in \mathbb{S}^{d-1}} \mathbf{x}^\top \mathbf{A} \mathbf{x} \right) + \|\Delta\| \\ &= \lambda_1(\mathbf{A}) + \|\Delta\|. \end{aligned}$$

Similarly, for the lower bound,

$$\begin{aligned} \lambda_d(\mathbf{A} + \Delta) &= \inf_{\mathbf{x} \in \mathbb{S}^{d-1}} \mathbf{x}^\top (\mathbf{A} + \Delta) \mathbf{x} \\ &\geq \inf_{\mathbf{x} \in \mathbb{S}^{d-1}} (\mathbf{x}^\top \mathbf{A} \mathbf{x} - \|\Delta\|) \\ &= \left( \inf_{\mathbf{x} \in \mathbb{S}^{d-1}} \mathbf{x}^\top \mathbf{A} \mathbf{x} \right) - \|\Delta\| \\ &= \lambda_d(\mathbf{A}) - \|\Delta\|. \end{aligned} \quad \square$$

These are special cases of *Weyl's inequalities*, which are a good more general tool to be aware of. The proof of the more general inequalities similarly uses a more general variational form of arbitrary eigenvalues given by the *Courant-Fischer min-max theorem*.

**Corollary 2.3.14.** *For any  $d \leq m$ ,*

$$\mathbb{P} \left[ \sqrt{m} - 64\sqrt{d} \leq \sigma_d(\mathbf{G}) \leq \sigma_1(\mathbf{G}) \leq \sqrt{m} + 32\sqrt{d} \right] \geq 1 - 2 \exp(-d).$$

*Proof.* We will use the basic bounds

$$\sqrt{1+x} \leq 1 + \frac{x}{2} \quad \text{for all } x \geq 0, \quad (2.3.21)$$

$$\sqrt{1-x} \geq 1 - x \quad \text{for all } 0 \leq x \leq 1. \quad (2.3.22)$$

Both may be viewed as consequences of the concavity of the square root function: the first says that one affine transformation of the square root lies below a tangent line to it, while the second says that another affine transformation lies above a secant segment connecting

two points on its graph. Then, on the event that  $\|\Delta\| > 64\sqrt{dm}$  from Theorem 2.3.12, using that  $\mathbf{G}\mathbf{G}^\top = m\mathbf{I}_d + \Delta$ , we have

$$\begin{aligned}
\sigma_1(\mathbf{G}) &= \sqrt{\lambda_1(\mathbf{G}\mathbf{G}^\top)} \\
&\leq \sqrt{m + \|\Delta\|} && \text{(Proposition 2.3.13)} \\
&\leq \sqrt{m + 64\sqrt{dm}} \\
&= \sqrt{m} \cdot \sqrt{1 + 64\sqrt{\frac{d}{m}}} \\
&\leq \sqrt{m} \cdot \left(1 + 32\sqrt{\frac{d}{m}}\right) && (2.3.21) \\
&= \sqrt{m} + 32\sqrt{d}.
\end{aligned}$$

Likewise, for  $\sigma_d(\mathbf{G})$  we have a similar argument. First, though, note that if  $m - 64\sqrt{dm} < 0$ , then  $\sqrt{m} - 64\sqrt{d} < 0$  as well, so the result holds vacuously in this case. Otherwise,

$$\begin{aligned}
\sigma_d(\mathbf{G}) &= \sqrt{\lambda_d(\mathbf{G}\mathbf{G}^\top)} \\
&\geq \sqrt{m - \|\Delta\|} && \text{(Proposition 2.3.13)} \\
&\geq \sqrt{m - 64\sqrt{dm}} \\
&= \sqrt{m} \cdot \sqrt{1 - 64\sqrt{\frac{d}{m}}} \\
&\geq \sqrt{m} \cdot \left(1 - 64\sqrt{\frac{d}{m}}\right) && (2.3.22) \\
&= \sqrt{m} - 64\sqrt{d},
\end{aligned}$$

completing the proof.  $\square$

Of course, what you should remember is just the softer statement that the singular values of  $\mathbf{G}$  are, with high probability, all in the range

$\sqrt{m} - O(\sqrt{d}) \leq \sigma_1(\mathbf{G}) \leq \dots \leq \sigma_d(\mathbf{G}) \leq \sqrt{m} + O(\sqrt{d}) \quad (\text{when } d \leq m)$

(2.3.23)

We highlight this because it is a fundamental fact about the spectra of random matrices that you should try to commit to your intuition. We will see considerably more precise elaborations of this claim soon, too. Note that, indeed, once  $m \gg d$ , then all of these singular values are  $\sqrt{m}$  to leading order and  $\mathbf{G}\mathbf{G}^\top \approx m\mathbf{I}_d$ , as we proposed earlier.

The top  $d$  eigenvalues of the larger matrix  $\mathbf{G}^\top \mathbf{G} \in \mathbb{R}_{\text{sym}}^{m \times m}$  are also then approximately  $m$ ; however, for this matrix, there are  $m - d$  more eigenvalues that are identically zero. Thus, the eigenvalues of  $\frac{1}{m}\mathbf{G}^\top \mathbf{G}$  are all either zero or close to 1; that is, it is close to a *projection matrix*. It is then natural to ask: a projection to what subspace? To answer this question we must understand the singular vectors of  $\mathbf{G}$ , to which we now turn our attention.

## 2.4 RECTANGULAR GAUSSIAN MATRICES: SINGULAR VECTORS

While the ordered singular values  $\sigma_1(\mathbf{G}) \geq \dots \geq \sigma_d(\mathbf{G})$  are well-defined for any  $\mathbf{G}$ , the individual singular vectors are not always well-defined. Suppose we take the singular value decomposition (SVD),

$$\mathbf{G} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top = \sum_{i=1}^d \sigma_i \mathbf{u}_i \mathbf{v}_i^\top,$$

for  $\{\mathbf{u}_1, \dots, \mathbf{u}_d\} \subset \mathbb{R}^d$  and  $\{\mathbf{v}_1, \dots, \mathbf{v}_d\} \subset \mathbb{R}^m$  orthonormal sets. When are we justified in viewing  $\mathbf{u}_i(\mathbf{G})$  and  $\mathbf{v}_i(\mathbf{G})$  as well-defined functions?

Actually, the answer is *never*, because the SVD does not disambiguate between the pair of singular vectors  $(\mathbf{u}_i, \mathbf{v}_i)$  and the pair  $(-\mathbf{u}_i, -\mathbf{v}_i)$ ; either pair would give rise to the same decomposition. To get around this minor issue, we may instead consider the rank-one matrices associated to the pairs of singular vectors,  $\mathbf{Z}_i = \mathbf{u}_i \mathbf{v}_i^\top$ . These matrices are unchanged by a simultaneous change of sign, and from them we may also recover the more geometrically meaningful projection matrices

$$\mathbf{L}_i := \mathbf{Z}_i \mathbf{Z}_i^\top = \mathbf{u}_i \mathbf{u}_i^\top, \quad (2.4.1)$$

$$\mathbf{R}_i := \mathbf{Z}_i^\top \mathbf{Z}_i = \mathbf{v}_i \mathbf{v}_i^\top. \quad (2.4.2)$$

It is sometimes useful for the sake of intuition to identify a projection matrix with the subspace it projects to; thus,  $\mathbf{u}_i \mathbf{u}_i^\top$  and  $\mathbf{v}_i \mathbf{v}_i^\top$  describe not the singular *vectors* but the entire *lines* in those directions.  $\mathbf{Z}_i$  contains the information of this pair of lines, and its being one of the rank-one matrices in the SVD says that  $\mathbf{G}$  maps the line spanned by  $\mathbf{v}_i$  to that spanned by  $\mathbf{u}_i$ .  $\mathbf{Z}_i$  also contains one more bit of information, saying whether that mapping maps  $\mathbf{v}_i$  to a positive or negative multiple of  $\mathbf{u}_i$ .

However, there are more nuances. For instance, if  $\sigma_i = 0$ , then  $\mathbf{Z}_i$  is equivalent to  $-\mathbf{Z}_i$  in the SVD. And, more substantially, if  $\sigma_i = \sigma_j$ , even if neither are zero, then  $\mathbf{Z}_i$  and  $\mathbf{Z}_j$  are not uniquely determined. Indeed, writing

$$\tilde{\mathbf{U}} := \begin{bmatrix} | & | \\ \mathbf{u}_i & \mathbf{u}_j \\ | & | \end{bmatrix} \in \mathbb{R}^{2 \times d}, \quad \tilde{\mathbf{V}} := \begin{bmatrix} | & | \\ \mathbf{v}_i & \mathbf{v}_j \\ | & | \end{bmatrix} \in \mathbb{R}^{2 \times m},$$

we have that the sum of the two corresponding terms in the SVD is

$$\sigma_i(\mathbf{R}_i + \mathbf{R}_j) = \sigma_i(\mathbf{u}_i \mathbf{v}_i^\top + \mathbf{u}_j \mathbf{v}_j^\top) = \sigma_i \tilde{\mathbf{U}} \tilde{\mathbf{V}}^\top = \sigma_i (\tilde{\mathbf{U}} \mathbf{Q}) (\tilde{\mathbf{V}} \mathbf{Q})^\top$$

for any  $\mathbf{Q} \in \mathcal{O}(2)$ . This gives rise to equivalent decompositions of  $\mathbf{R}_i + \mathbf{R}_j$  as a sum of rank-two matrices. Since  $\sigma_i = \sigma_j$ , the expression  $\sigma_i(\mathbf{R}_i + \mathbf{R}_j)$  is what appears in the SVD for these two components, so these equivalent decompositions mean that  $\mathbf{R}_i$  and  $\mathbf{R}_j$  are not uniquely determined in this case.

As an extreme example, consider the case  $d = m$  and  $\mathbf{G} = \mathbf{I}_d$ , the identity matrix. All  $d$  of its singular values are 1, and it admits many SVDs of the form  $\mathbf{G} = \sum_{i=1}^d \mathbf{u}_i \mathbf{u}_i^\top$  for *any* orthonormal basis  $\{\mathbf{u}_1, \dots, \mathbf{u}_d\}$ . Thus, the more repeated singular values there are, the more dire this issue becomes. However, as the following fact of linear algebra shows, these are the only issues.



**Proposition 2.4.1.** *Suppose that  $\mathbf{G} \in \mathbb{R}^{d \times m}$  with  $d \leq m$ ,  $\sigma_d(\mathbf{G}) > 0$ , and the singular values  $\sigma_1(\mathbf{G}), \dots, \sigma_d(\mathbf{G})$  are distinct. Then, there exist unique  $\mathbf{Z}_i = \mathbf{Z}_i(\mathbf{G}) \in \mathbb{R}^{d \times m}$  such that each  $\mathbf{R}_i$  has rank one, has  $\sigma_1(\mathbf{Z}_i) = 1$ , and*

$$\mathbf{G} = \sum_{i=1}^d \sigma_i(\mathbf{G}) \mathbf{Z}_i(\mathbf{G}).$$

#### 2.4.1 DISTINCTNESS OF SINGULAR VALUES

Remarkably, the conditions of the above Proposition 2.4.1 hold not merely with high probability but *almost surely* for Gaussian random matrices (and, though we will not discuss it here, for other matrices whose entries are independent random variables whose laws are absolutely continuous with respect to the Lebesgue measure).

**Theorem 2.4.2.** *Let  $\mathbf{G} \sim \mathcal{N}(0, 1)^{\otimes d \times m}$  with  $d \leq m$ . Then, almost surely, the following hold:*

1.  $\sigma_d(\mathbf{G}) > 0$ .
2. *The singular values  $\sigma_1(\mathbf{G}), \dots, \sigma_d(\mathbf{G})$  are distinct, i.e., we have the strict inequalities  $\sigma_1(\mathbf{G}) > \dots > \sigma_d(\mathbf{G})$ .*

We will use the following result to establish both parts, which should be intuitive but is a little bit clunky to prove. You can prove it as an exercise—one approach is to show the more geometric fact that any hypersurface has Lebesgue measure zero.

**Lemma 2.4.3.** *Suppose  $p \in \mathbb{R}[x_1, \dots, x_N]$  is a non-constant polynomial, and  $\mathbf{g} \sim \mathcal{N}(0, \mathbf{I}_N)$ . Then, for any  $a \in \mathbb{R}$ ,  $\mathbb{P}[p(\mathbf{g}) = a] = 0$ .*

The first part of Theorem 2.4.2 can be proved quite directly using this result, as follows.

*Proof of Theorem 2.4.2, Claim 1.* Note that we have  $\sigma_d(\mathbf{G}) = 0$  if and only if  $\mathbf{G}$  is rank-deficient, which occurs if and only if  $P(\mathbf{G}) := \det(\mathbf{G}\mathbf{G}^\top) = 0$ . (Note that this is not the case if we take  $\det(\mathbf{G}^\top\mathbf{G})$  instead, which is zero whenever  $d < m$ .) It is easy to show that  $p$  is a non-constant polynomial of the  $dm$  entries of the matrix argument: it is non-constant because, for example, for rank-deficient matrix inputs it is zero while for full-rank inputs it is non-zero. Thus  $\mathbb{P}[\sigma_d(\mathbf{G}) = 0] = \mathbb{P}[P(\mathbf{G}) = 0] = 0$  by Lemma 2.4.3.  $\square$

For the second part we will need to establish some tools from the theory of symmetric polynomials as well as a relationship between polynomials of the *eigenvalues* of a matrix and polynomials of the *entries* of a matrix. First, the following are the basic notions of the general theory of symmetric polynomials.

**Definition 2.4.4.** *A polynomial  $f \in \mathbb{R}[t_1, \dots, t_d]$  is symmetric if, for all permutations  $\pi \in S_d$ , we have  $f(t_{\pi(1)}, \dots, t_{\pi(d)}) = f(t_1, \dots, t_d)$ .*

**Theorem 2.4.5** (Fundamental theorem of symmetric polynomials). *Define the power sum polynomials*

$$s_k(t_1, \dots, t_d) := \sum_{i=1}^d t_i^k.$$

Then, if  $f \in \mathbb{R}[t_1, \dots, t_d]$  is a symmetric polynomial, then there exists some polynomial  $g$  of  $d$  variables such that

$$f(t_1, \dots, t_d) = g(s_1(t_1, \dots, t_d), \dots, s_d(t_1, \dots, t_d)).$$

In words, this result says that any symmetric polynomial is some polynomial of the power sum polynomials, or that the power sum polynomials *algebraically generate* the symmetric polynomials (which form a subring of the ring of all polynomials).

The power sum polynomials have an important role in random matrix theory, because they serve as a bridge between the properties of the entries of a matrix and the properties of the eigenvalues, a connection that we will see much more sophisticated uses of later. The basic underlying fact is quite simple, however. Suppose here and below that  $\mathbf{X} \in \mathbb{R}_{\text{sym}}^{d \times d}$ . Let us define the *power trace* polynomials of  $\mathbf{X}$ ,

$$S_k(\mathbf{X}) := \text{Tr}(\mathbf{X}^k) = \sum_{a_1, \dots, a_k=1}^d X_{a_1 a_2} X_{a_2 a_3} \cdots X_{a_{k-1} a_k} X_{a_k a_1}.$$

The latter expression shows that  $S_k(\mathbf{X})$  is a polynomial of the entries of  $\mathbf{X}$ .

**Proposition 2.4.6.**  $s_k(\lambda_1(\mathbf{X}), \dots, \lambda_d(\mathbf{X})) = S_k(\mathbf{X})$ .

This is just a simple consequence of the fact that the eigenvalues of  $\mathbf{X}^k$  are the  $\lambda_i(\mathbf{X})^k$ . However, it leads to the following important corollary.

**Corollary 2.4.7.** Suppose that  $F(\mathbf{X}) = f(\lambda_1(\mathbf{X}), \dots, \lambda_d(\mathbf{X}))$  for  $f$  a symmetric polynomial (i.e.,  $F(\mathbf{X})$  is a symmetric polynomial of the eigenvalues of  $\mathbf{X}$ ). Then,  $F(\mathbf{X})$  is also a (not necessarily symmetric) polynomial of the entries of  $\mathbf{X}$ .

**Remark 2.4.8.** The assumption that  $f$  is symmetric is crucial. For instance, you may check as an exercise that  $F(\mathbf{X}) = \lambda_1(\mathbf{X})$  is not a polynomial of  $\mathbf{X}$ .

With these tools, we are now ready for the proof of the second part of the main Theorem.

*Proof of Theorem 2.4.2, Claim 2.* Consider first a function of  $\mathbf{X} \in \mathbb{R}_{\text{sym}}^{d \times d}$ , defined as:

$$\begin{aligned} f(t_1, \dots, t_d) &:= \prod_{1 \leq i < j \leq d} (t_i - t_j)^2, \\ F(\mathbf{X}) &:= f(\lambda_1(\mathbf{X}), \dots, \lambda_d(\mathbf{X})) \\ &= \prod_{1 \leq i < j \leq d} (\lambda_i(\mathbf{X}) - \lambda_j(\mathbf{X}))^2. \end{aligned}$$

Clearly we have that  $\mathbf{X}$  has a repeated eigenvalue if and only if  $F(\mathbf{X}) = 0$ . Further,  $F(\mathbf{X})$  satisfies the conditions of Corollary 2.4.7. Therefore,  $F(\mathbf{X})$  is a polynomial of the entries of  $\mathbf{X}$ . (Note that, remarkably, we can draw this conclusion non-constructively from the tools above.)

Now, returning to the setting of the Theorem, we see that  $\mathbf{G}$  has a repeated singular value if and only if  $P(\mathbf{G}) = 0$ , where  $P(\mathbf{G}) := F(\mathbf{G}\mathbf{G}^\top)$ . As before,  $P$  is non-constant merely

because some matrices do have repeated singular values and some do not. Thus, we may conclude as before by Lemma 2.4.3 that

$$\mathbb{P}[\mathbf{G} \text{ has a repeated singular value}] = \mathbb{P}[P(\mathbf{G}) = 0] = 0,$$

as claimed.  $\square$

**Remark 2.4.9.** *The proof above shows that the set of symmetric matrices with a repeated eigenvalue is a hypersurface in the space  $\mathbb{R}_{\text{sym}}^{d \times d}$ , which may be identified with  $\mathbb{R}^{d(d+1)/2}$ . A priori, since we described this set as  $\{\mathbf{X} : F(\mathbf{X}) = 0\}$  for a polynomial  $F$ , we only know that its codimension is at most 1 (i.e., the dimension of this hypersurface is at most  $\frac{d(d+1)}{2} - 1$ ). Actually, its codimension is exactly 2. For example, the dimension of  $\mathbb{R}_{\text{sym}}^{2 \times 2}$  is  $\frac{2(2+1)}{2} = 3$ , while the  $2 \times 2$  matrices with a repeated eigenvalue are merely multiples of the identity, whose dimension is clearly  $1 = 3 - 2$  (indeed, these matrices form a line in the space of all symmetric matrices). Note that one polynomial constraint can in fact encode several, for instance if  $F(\mathbf{X}) = F_1(\mathbf{X})^2 + F_2(\mathbf{X})^2 = 0$ , then  $F_1(\mathbf{X}) = F_2(\mathbf{X}) = 0$ . As discussed for instance by [BKL18], a more complicated version of this sum-of-squares structure is indeed what happens in the case of  $F(\mathbf{X})$ .*

We have established that, on an almost sure event, we are justified in speaking of the singular vectors of a Gaussian random matrix, at least via the rank-one matrices  $\mathbf{Z}_i(\mathbf{G})$  and the projections to the left and right singular vectors  $\mathbf{L}_i(\mathbf{G})$  and  $\mathbf{R}_i(\mathbf{G})$ . We are now ready to study the laws of these objects.

## 2.4.2 ORTHOGONAL INVARIANCE FOR MATRICES

We will see that orthogonal invariance, as we discussed for vectors earlier, is crucial to understanding these distributions. The appropriate definitions for random matrices are as follows:

**Definition 2.4.10** (Orthogonal invariance of matrices). *We say that a random matrix  $\mathbf{X} \in \mathbb{R}^{d \times m}$  (or its law  $\text{Law}(\mathbf{x})$ ) is:*

- Left-orthogonally invariant if, for each (deterministic) orthogonal matrix  $\mathbf{Q} \in \mathcal{O}(d)$ , we have  $\text{Law}(\mathbf{Q}\mathbf{X}) = \text{Law}(\mathbf{X})$ .
- Right-orthogonally invariant if, for each (deterministic) orthogonal matrix  $\mathbf{Q} \in \mathcal{O}(m)$ , we have  $\text{Law}(\mathbf{X}\mathbf{Q}) = \text{Law}(\mathbf{X})$ .
- Bi-orthogonally invariant if it is both left- and right-orthogonally invariant.

Further, we say that a random symmetric matrix  $\mathbf{X} \in \mathbb{R}_{\text{sym}}^{d \times d}$  is symmetric-orthogonally invariant if  $\text{Law}(\mathbf{Q}\mathbf{X}\mathbf{Q}^\top) = \text{Law}(\mathbf{X})$  for each  $\mathbf{Q} \in \mathcal{O}(d)$ .

The following property is then simple to check using either explicit mean and covariance calculations or by using that the rows and columns of an i.i.d. Gaussian random matrix are i.i.d. standard Gaussian vectors, which are themselves orthogonally invariant.

**Proposition 2.4.11.**  $\mathbf{G} \sim \mathcal{N}(0, 1)^{\otimes d \times m}$  is bi-orthogonally invariant.

**Corollary 2.4.12.** *If  $G \sim \mathcal{N}(0, 1)^{\otimes d \times m}$ , then  $L_i(G)$  and  $R_i(G)$  are both symmetric-orthogonally invariant, for each  $i = 1, \dots, d$ . Further, the stronger equalities of law hold that*

$$\begin{aligned} \text{Law}\left((QL_1(G)Q^\top, \dots, QL_d(G)Q^\top)\right) &= \text{Law}\left((L_1(G), \dots, L_d(G))\right) \quad \text{for all } Q \in \mathcal{O}(d), \\ \text{Law}\left((QR_1(G)Q^\top, \dots, QR_d(G)Q^\top)\right) &= \text{Law}\left((R_1(G), \dots, R_d(G))\right) \quad \text{for all } Q \in \mathcal{O}(m). \end{aligned}$$

The following basic linear algebra observation is a useful preliminary to the proof.

**Proposition 2.4.13** (Invariance and equivariance of SVD). *Let  $G \in \mathbb{R}^{d \times m}$  be a (deterministic) matrix such that  $\sigma_1(G) > \dots > \sigma_d(G) > 0$ . Let  $Q_L \in \mathcal{O}(d)$  and  $Q_R \in \mathcal{O}(m)$ . Then,*

$$\begin{aligned} \sigma_i(Q_L G Q_R^\top) &= \sigma_i(G), \\ Z_i(Q_L G Q_R^\top) &= Q_L Z_i(G) Q_R^\top, \\ L_i(Q_L G Q_R^\top) &= Q_L L_i(G) Q_L^\top, \\ R_i(Q_L G Q_R^\top) &= Q_R R_i(G) Q_R^\top. \end{aligned}$$

*Proof of Corollary 2.4.12.* Let us prove the simpler claim about the individual  $L_i$ : we have, by combining the above results,

$$\begin{aligned} \text{Law}(QL_i(G)Q^\top) &= \text{Law}(L_i(QG)) && \text{(Proposition 2.4.13)} \\ &= \text{Law}(L_i(G)). && \text{(Proposition 2.4.11)} \end{aligned}$$

The result for the individual  $R_i$  and the collections of all  $L_i$  and  $R_i$  follows similarly.  $\square$

Let us comment on the intuition behind what we have derived so far. Consider the  $L_i(G)$ . Recall that, for  $u_i$  some choice of the left singular vectors in the SVD (which for  $G$  are almost surely well-defined up to a sign flip),  $L_i(G) = u_i u_i^\top$ . As we have discussed, we may identify this rank-one projection matrix with the line formed by  $\text{span}(u_i)$ . Thus the collection  $(L_1(G), \dots, L_d(G))$  is a collection of  $d$  orthogonal lines, i.e., a coordinate system for  $\mathbb{R}^d$  consisting of  $d$  orthogonal axes, although without a “positive direction” chosen for any of the axes. The above result says that this object for  $G \sim \mathcal{N}(0, 1)^{\otimes d \times m}$  has an orthogonally invariant distribution. We will next try to formally explain in what sense this makes  $(L_1(G), \dots, L_d(G))$  have the law of the unique “uniformly random” choice of coordinate system.

## 2.5 APPLICATION: COMPRESSED SENSING

We will now see one last application of the ideas developed so far. This is for the problem of *compressed sensing*, recovering  $x \in \mathbb{R}^m$  from  $y = Gx \in \mathbb{R}^d$  with  $G \in \mathbb{R}^{d \times m}$  (the same dimensions as before).

For general  $x \in \mathbb{R}^m$ , we need  $G$  to be injective, which requires  $d \geq m$ . Compressed sensing concerns making  $d$  much smaller, provided we are promised that  $x$  is *sparse*. In this setting, while we will discuss taking  $G$  random, we view ourselves as having control over  $G$  (the so-called *sensing matrix*) in general, so the sparsity assumption should be seen as in a basis of our choosing. Thus to take advantage of compressed sensing we should encode any prior knowledge of  $x$  in a basis that makes  $x$  sparse. (For example, if  $x$  encodes an image, we might choose a wavelet basis.)

### 2.5.1 NULL SPACE AND RESTRICTED ISOMETRY PROPERTIES

Let us denote sparsity by

$$\|\mathbf{x}\|_0 := \#\{i \in [m] : x_i \neq 0\}. \quad (2.5.1)$$

Note that this “ $\ell^0$ -norm” is not actually a norm, because it is not homogeneous in  $\mathbf{x}$ .

Assuming that  $\|\mathbf{x}\|_0 \leq k$  makes the task of recovering  $\mathbf{x}$  easier. At the extreme, if  $k = 1$ , then  $\mathbf{y} = \mathbf{G}\mathbf{x}$  is just (up to scaling) one of the columns of  $\mathbf{G}$ . Thus provided the columns of  $\mathbf{G}$  are well-separated, it will be easy to recover  $\mathbf{x}$ . You can show that it is possible to construct  $\exp(\Omega(d))$  unit vectors in  $\mathbb{R}^d$  that have, say, pairwise inner products of magnitude each at most  $1/2$ , which we may use as the columns of  $\mathbf{G}$  (and such a choice will be robust to a small amount of noise). Thus when  $k = 1$  then we may take  $d$  as small as  $O(\log m)$  while still being able to recover any 1-sparse  $\mathbf{x}$ .

How does the situation change for larger  $k$ ? We will show below that it actually does not change that much.

First, let us specify what we mean by  $\mathbf{x}$  being “recoverable.” The following definition is not standard but useful.

**Definition 2.5.1.** *We say that  $\mathbf{G}$  distinguishes  $k$ -sparse vectors if  $\mathbf{G}\mathbf{x} \neq \mathbf{G}\mathbf{x}'$  for any  $\mathbf{x} \neq \mathbf{x}'$  with  $\|\mathbf{x}\|_0, \|\mathbf{x}'\|_0 \leq k$ .*

If  $\mathbf{G}$  distinguishes  $k$ -sparse vectors, then compressed sensing is *information-theoretically* possible: by, say, brute force search over an  $\epsilon$ -net of sparse vectors followed by a rounding procedure, we may exactly recover any  $k$ -sparse  $\mathbf{x}$  from  $\mathbf{y} = \mathbf{G}\mathbf{x}$ . We will not go into the *computational* feasibility of compressed sensing here, which is a deep area in its own right. You may look up the role of  $\ell^1$ -norm minimization for such algorithms to get started with computational approaches.

The following more linear-algebraic definition is actually equivalent to distinguishing  $k$ -sparse vectors.

**Definition 2.5.2** (Null space property). *We say that  $\mathbf{G}$  has the  $k$ -null space property ( $k$ -NSP) if, for all  $\mathbf{x} \neq \mathbf{0}$  with  $\|\mathbf{x}\|_0 \leq k$ ,  $\mathbf{G}\mathbf{x} \neq \mathbf{0}$ .*

**Proposition 2.5.3.**  *$\mathbf{G}$  distinguishes  $k$ -sparse vectors if and only if  $\mathbf{G}$  has the  $2k$ -NSP.*

*Proof.* If  $\mathbf{G}$  does not distinguish  $k$ -sparse vectors, then there exist  $\mathbf{x} \neq \mathbf{x}'$  with  $\|\mathbf{x}\|_0, \|\mathbf{x}'\|_0 \leq k$  such that  $\mathbf{G}\mathbf{x} = \mathbf{G}\mathbf{x}'$ . In particular,  $\mathbf{G}(\mathbf{x} - \mathbf{x}') = \mathbf{0}$ , and  $\|\mathbf{x} - \mathbf{x}'\|_0 \leq 2k$ , so  $\mathbf{G}$  does not have the  $2k$ -NSP. Conversely, if  $\mathbf{G}$  does distinguish  $k$ -sparse vectors and  $\mathbf{x}'' \neq \mathbf{0}$  has  $\|\mathbf{x}''\|_0 \leq 2k$ , then we may write  $\mathbf{x}'' = \mathbf{x} - \mathbf{x}'$  for  $\mathbf{x} \neq \mathbf{x}'$  and  $\|\mathbf{x}\|_0, \|\mathbf{x}'\|_0 \leq k$  (by partitioning the  $2k$  indices on which  $\mathbf{x}''$  is non-zero in some arbitrary way). Reversing the above argument then shows that  $\mathbf{G}$  has the  $2k$ -NSP.  $\square$

We will in fact be able to show the following, more quantitative property.

**Definition 2.5.4** (Restricted isometry property). *We say that  $\mathbf{G}$  has the  $(k, \delta)$ -restricted isometry property ( $(k, \delta)$ -RIP) if, for all  $\mathbf{x} \neq \mathbf{0}$  with  $\|\mathbf{x}\|_0 \leq k$ ,*

$$(1 - \delta)\|\mathbf{x}\|^2 \leq \|\mathbf{G}\mathbf{x}\|^2 \leq (1 + \delta)\|\mathbf{x}\|^2. \quad (2.5.2)$$

The following implication is immediate.

**Proposition 2.5.5.** *If  $G$  has the  $(k, \delta)$ -RIP for any  $\delta < 1$ , then  $G$  has the  $k$ -NSP.*

The RIP is more useful than the NSP both for analyzing compressed sensing when some noise is further added to  $y = Gx$ , and for analyzing algorithms for recovering  $x$ , but here we will just view it as an analytic tool for establishing the NSP.

## 2.5.2 RANDOM SENSING MATRICES

The following is the main result that we will show.

**Theorem 2.5.6.** *For any  $1 \leq k \leq m$  and  $\delta \in (0, 1)$ , if*

$$d \geq \frac{160}{\delta^2} k \left( 1 + \log \left( \frac{m}{k} \right) \right)$$

*and  $G \sim \mathcal{N}(0, 1)^{\otimes d \times m}$ , then*

$$\mathbb{P}[G \text{ has the } k\text{-NSP}] \geq \mathbb{P}[G \text{ has the } (k, \delta)\text{-RIP}] \geq 1 - \frac{2}{\binom{m}{k}} \geq 1 - 2 \left( \frac{k}{m} \right)^k. \quad (2.5.3)$$

Note that this scaling of  $d$  is consistent with our earlier observation that  $d \gtrsim \log m$  was the right condition for  $k = 1$ . Remarkably, the result implies that, while

$$d \geq m$$

is needed for *arbitrary*  $x$  to be recoverable from  $Gx$ , for the specific case of  $k$ -sparse vectors it suffices to have only the far smaller

$$d \gtrsim k \log \left( \frac{m}{k} \right) \quad (2.5.4)$$

*Proof.* For  $S \subseteq [m]$  with  $|S| = k$  and  $x \in \mathbb{R}^m$ , let  $x^{(S)} \in \mathbb{R}^k$  be the restriction of  $x$  to the indices in  $S$ . Likewise, for  $G \in \mathbb{R}^{d \times m}$ , let  $G^{(S)} \in \mathbb{R}^{d \times k}$  be the restriction of  $G$  to the columns whose indices are in  $S$ . If  $\|x\|_0 \leq k$  and the non-zero indices of  $x$  are contained in such  $S$ , then

$$\|x\|^2 = \|x^{(S)}\|^2, \quad (2.5.5)$$

$$Gx = G^{(S)} x^{(S)}. \quad (2.5.6)$$

We may then view the  $(k, \delta)$ -RIP as requiring that, for all  $S \subseteq [m]$  with  $|S| = k$  and all  $x^{(S)} \in \mathbb{R}^k$ , we have

$$(1 - \delta) \|x^{(S)}\|^2 \leq \|G^{(S)} x^{(S)}\|^2 \leq (1 + \delta) \|x^{(S)}\|^2. \quad (2.5.7)$$

But this just amounts to asking that

$$1 - \delta \leq \lambda_k(G^{(S)\top} G^{(S)}) \leq \lambda_1(G^{(S)\top} G^{(S)}) \leq 1 + \delta, \quad (2.5.8)$$

or again equivalently that

$$\|\mathbf{G}^{(S)\top} \mathbf{G}^{(S)} - \mathbf{I}_k\| \leq \delta. \quad (2.5.9)$$

We will then proceed by union bounding,

$$\mathbb{P}[\mathbf{G} \text{ does not have the } (k, \delta)\text{-RIP}] \leq \sum_{S \in \binom{[m]}{k}} \mathbb{P}[\|\mathbf{G}^{(S)\top} \mathbf{G}^{(S)} - \mathbf{I}_k\| > \delta]. \quad (2.5.10)$$

This is almost exactly the kind of deviation that we bounded before in Theorem 2.3.12. First, let's decode our statement into that form. Introduce  $\mathbf{H} \sim \mathcal{N}(0, 1)^{\otimes k \times d}$ , the scaling and aspect ratio that Theorem 2.3.12 addresses. We have  $\text{Law}(\mathbf{G}^{(S)}) = \text{Law}(d^{-1/2} \mathbf{H}^\top)$ . Thus,

$$\begin{aligned} \mathbb{P}[\|\mathbf{G}^{(S)\top} \mathbf{G}^{(S)} - \mathbf{I}_k\| > \delta] &= \mathbb{P}\left[\left\|\frac{1}{d} \mathbf{H} \mathbf{H}^\top - \mathbf{I}_k\right\| > \delta\right] \\ &= \mathbb{P}[\|\mathbf{H} \mathbf{H}^\top - d \mathbf{I}_k\| > \delta d]. \end{aligned}$$

The only difference between this and what Theorem 2.3.12 covered is that there we were concerned with deviations of the order  $O(\sqrt{kd})$ , while here we are concerned with  $O(d)$ , which is much larger since  $d \gg k$  under our assumption. But, we may still apply the intermediate result (2.3.18) from that proof, which gives

$$\begin{aligned} &\leq 2 \exp\left(3k - \frac{1}{8} \min\left\{\frac{\delta^2 d^2}{4d}, \frac{\delta d}{2}\right\}\right) \\ &= 2 \exp\left(3k - \frac{\delta^2}{32} d\right), \end{aligned} \quad (2.5.11)$$

using that  $\delta < 1$  so  $\delta^2/4 = (\delta/2)^2 < \delta/2$ . Returning to the union bound, we first recall the non-asymptotic bounds on binomial coefficients

$$\left(\frac{m}{k}\right)^k \leq \binom{m}{k} \leq \left(\frac{em}{k}\right)^k,$$

which we will use twice below. We may bound the terms in the union bound uniformly as

$$\begin{aligned} \mathbb{P}[\mathbf{G} \text{ does not have the } (k, \delta)\text{-RIP}] &\leq \binom{m}{k} \cdot 2 \exp\left(3k - \frac{\delta^2}{32} d\right) \\ &\leq \left(\frac{em}{k}\right)^k \cdot 2 \exp\left(3k - \frac{\delta^2}{32} d\right) \\ &\leq 2 \exp\left(k \left(1 + \log\left(\frac{m}{k}\right)\right) + 3k - \frac{\delta^2}{32} d\right) \\ &\leq 2 \exp\left(4k \left(1 + \log\left(\frac{m}{k}\right)\right) - \frac{\delta^2}{32} d\right) \\ &\leq 2 \exp\left(-k \left(1 + \log\left(\frac{m}{k}\right)\right)\right) \\ &\leq 2 \left(\frac{em}{k}\right)^{-k} \\ &\leq \frac{2}{\binom{m}{k}} \end{aligned} \quad (2.5.12)$$

completing the proof. The alternate form of the bound follows from the lower bound on binomial coefficients above.  $\square$

# BIBLIOGRAPHY

- [AC09] Nir Ailon and Bernard Chazelle. The fast johnson-lindenstrauss transform and approximate nearest neighbors. *SIAM Journal on Computing*, 39(1):302–322, 2009.
- [BKL18] Paul Breiding, Khazhgali Kozhasov, and Antonio Lerario. On the geometry of the set of symmetric matrices with repeated eigenvalues. *Arnold Mathematical Journal*, 4(3):423–443, 2018.
- [DKS10] Anirban Dasgupta, Ravi Kumar, and Tamás Sarlós. A sparse Johnson-Lindenstrauss transform. In *Proceedings of the forty-second ACM symposium on Theory of computing*, pages 341–350, 2010.
- [IM98] Piotr Indyk and Rajeev Motwani. Approximate nearest neighbors: towards removing the curse of dimensionality. In *Proceedings of the thirtieth annual ACM symposium on Theory of computing*, pages 604–613, 1998.
- [JL82] William Johnson and J. Lindenstrauss. Extensions of Lipschitz mappings into a Hilbert space. *Conference in Modern Analysis and Probability*, 26:189–206, 01 1982.
- [KN14] Daniel M Kane and Jelani Nelson. Sparser Johnson-Lindenstrauss transforms. *Journal of the ACM (JACM)*, 61(1):1–23, 2014.
- [LN14] Kasper Green Larsen and Jelani Nelson. The johnson-lindenstrauss lemma is optimal for linear dimensionality reduction. *arXiv preprint arXiv:1411.2404*, 2014.
- [LN17] Kasper Green Larsen and Jelani Nelson. Optimality of the johnson-lindenstrauss lemma. In *2017 IEEE 58th Annual Symposium on Foundations of Computer Science (FOCS)*, pages 633–638. IEEE, 2017.
- [RH17] Phillippe Rigollet and Jan-Christian Hütter. Lecture notes on high dimensional statistics. 2017.
- [Tal05] Michel Talagrand. *The generic chaining: upper and lower bounds of stochastic processes*. Springer, 2005.