

Notes for High-Dimensional Probability Second Edition by
Roman Vershynin

Gallant Tsao

July 30, 2025

Contents

0	Appetizer: Using Probability to Cover a Set	4
0.1	Covering Geometric Sets	5
1	A Quick Refresher on Analysis and Probability	7
1.1	Convex Sets and Functions	7
1.2	Norms and Inner Products	7
1.3	Random Variables and Random Vectors	7
1.4	Union Bound	8
1.5	Conditioning	9
1.6	Probabilistic Inequalities	9
1.7	Limit Theorems	11
2	Concentration of Sums of Independent Random Variables	13
2.1	Why Concentration Inequalities?	13
2.2	Hoeffding Inequality	14
2.3	Chernoff Inequality	16
2.4	Application: Median-of-means Estimator	18
2.5	Application: Degrees of Random Graphs	20
2.6	Subgaussian Distributions	21
2.6.1	The Subgaussian Norm	23
2.7	Subgaussian Hoeffding and Khintchine Inequalities	23
2.7.1	Subgaussian Hoeffding Inequality	24
2.7.2	Subgaussian Khintchine Inequality	24
2.7.3	Maximum of Subgaussians	25
2.7.4	Centering	26
2.8	Subexponential Distributions	26
2.8.1	Subexponential Properties	26
2.8.2	The Subexponential Norm	28
2.9	Bernstein Inequality	29
3	Random Vectors in High Dimensions	32
3.1	Concentration of the Norm	32
3.2	Covariance Matrices and PCA	33
3.2.1	Learning from the Covariance Matrix	33
3.2.2	Principle Component Analysis	34
3.2.3	Isotropic Distributions	35
3.3	Examples of High-dimensional Distributions	35
3.3.1	Standard Normal	35
3.3.2	General Normal	36
3.3.3	Uniform on the Sphere	37
3.3.4	Uniform on a Convex Set	38
3.3.5	Frames	38
3.4	Subgaussian Distributions in High Dimensions	40
3.4.1	Gaussian, Rademacher, and More	40
3.4.2	Uniform on the Sphere	40
3.4.3	Non-examples	41
3.5	Application: Grothendieck Inequality and Semidefinite Programming	42
3.5.1	Semidefinite Programming	44
3.6	Application: Maximum Cut for Graphs	46
3.6.1	A Simple 0.5-approximation Algorithm	46
3.6.2	Semidefinite Relaxation	47
3.7	Kernel Trick and Tightening of Grothendieck Inequality	48
3.7.1	Tensors	49
3.7.2	Proof of Theorem 3.5.1	51
3.7.3	Kernels and Feature Maps	51

4	Random Matrices	52
4.1	A Quick Refresher on Linear Algebra	52
4.1.1	Singular Value Decomposition	52
4.1.2	Min-max Theorem	53
4.1.3	Frobenius and Operator Norms	54
4.1.4	The Matrix Norms and the Spectrum	54
4.1.5	Low-rank Approximation	55
4.1.6	Perturbation Theory	55
4.1.7	Isometries	57
4.2	Nets, Covering, and Packing	57
4.2.1	Covering Numbers and Volume	59
4.3	Application: Error Correcting Codes	60
4.3.1	Metric Entropy and Complexity	61
4.3.2	Error Correcting Codes	61
4.4	Upper Bounds on Subgaussian Random Matrices	63
4.4.1	Computing the Norm on an ε net	63
4.4.2	The Norms of Subgaussian Random Matrices	63
4.4.3	Symmetric Matrices	65
4.5	Application: Community Detection in Networks	65
4.5.1	Stochastic Block Model	65
4.5.2	The Expected Adjacency Matrix Holds the Key	66
4.5.3	The Actual Adjacency Matrix is a Good Approximation	66
4.5.4	Perturbation Theory	67
4.5.5	Spectral Clustering	67
4.6	Two-sided Bounds on Subgaussian Matrices	68
4.7	Application: Covariance Estimation and Clustering	69
4.7.1	Application: Clustering of Point Sets	71
5	Concentration Without Independence	73
5.1	Concentration of Lipschitz Functions on the Sphere	73
5.1.1	Lipschitz Functions	73
5.1.2	Concentration via Isoperimetric Inequalities	73
5.1.3	Blow-up of Sets on the Sphere	74
5.1.4	Proof of Theorem 5.1.3	75
5.2	Concentration on Other Metric Measure Spaces	76
5.2.1	Gaussian Concentration	76
5.2.2	Hamming Cube	76
5.2.3	Symmetric Group	77
5.2.4	Riemannian Manifolds with Strictly Positive Curvature	77
5.2.5	Special Orthogonal Group	77
5.2.6	Grassmannian	78
5.2.7	Continuous Cube and Euclidean Ball	78
5.2.8	Densities of the Form $e^{-U(x)}$	78
5.2.9	Random Vectors with Independent Bounded Coordinates	79
5.3	Application: Johnson-Lindenstrauss Lemma	79
5.4	Matrix Bernstein Inequality	81
5.4.1	Matrix Calculus	81
5.4.2	Trace Inequalities	83
5.4.3	Proof of Matrix Bernstein Inequality	83
5.4.4	Matrix Hoeffding and Khintchine Inequalities	85
5.5	Application: Community Detection in Sparse Networks	86
5.6	Application: Covariance Estimation for General Distributions	88
5.7	Extra notes	90

6	Quadratic Forms, Symmetrization, and Contraction	91
6.1	Decoupling	91
6.2	Hanson-Wright Inequality	93
6.3	Symmetrization	96
6.4	Random Matrices with non-i.i.d. Entries	97
6.5	Application: Matrix Completion	98
6.6	Contraction Principle	100
7	Random Processes	92
7.1	Basic Concepts and Examples	92
7.1.1	Covariance and Increments	93
7.1.2	Gaussian Processes	93
7.2	Slepian, Sudakov-Fernique, and Gordon Inequalities	94
7.2.1	Gaussian Interpolation	95
7.2.2	Proof of Slepian Inequality	97
7.2.3	Sudakov-Fernique and Gordon Inequalities	98
7.3	Application: Sharp Bounds for Gaussian Matrices	99
7.4	Sudakov Inequality	99
7.4.1	Application for covering numbers in \mathbb{R}^n	100
7.5	Gaussian Width	101
7.5.1	Geometric Meaning of Width	102
7.5.2	Examples	103
7.5.3	Gaussian Complexity and Effective Dimension	104
7.6	Application: Random Projection of Sets	105
8	Chaining	106
8.1	Dudley Inequality	106
8.1.1	Variations and Examples	109
8.2	Application: Empirical Processes	110
8.3	VC Dimension	110
8.3.1	Definition and Examples	110
8.3.2	Pajor's Lemma	112
8.3.3	Sauer-Shelah Lemma	113
8.3.4	Growth Function	113
8.3.5	Covering Numbers via VC Dimension	115
8.3.6	VC Law of Large Numbers	116
8.4	Application: Statistical Learning Theory	119
8.5	Generic Chaining	119
8.5.1	A Makeover of Dudley's Inequality	119
8.5.2	The γ_2 Functional and Generic Chaining	119
8.5.3	Majorizing Measure and Comparison Theorems	121
8.6	Chevet Inequality	123
9	Deviations of Random Matrices on Sets	125
9.1	Matrix Deviation Inequality	125
9.2	Random Matrices, Covariance Estimation, and Johnson-Lindenstrauss	128
9.2.1	Singular Values of Random Matrices	128
9.2.2	Random Projections of Sets	129
9.2.3	Covariance Estimation for Low-dimensional Distributions	129
9.2.4	Johnson-Lindenstrauss Lemma for Infinite Sets	129
9.3	Random Sections: The M^* Bound and Escape Theorem	130
9.3.1	The M^* Bound	130
9.3.2	The Escape Theorem	131
9.4	Application: High-dimensional Linear Models	132
9.5	Application: Exact Sparse Recovery	132
9.6	Deviations of Random Matrices for General Norms	132
9.7	Two-sided Chevet Inequality and Dvoretzky-Milman Theorem	134
9.7.1	Two-sided Chevet's Inequality	134
9.7.2	Dvoretzky-Milman Theorem	135

6 Quadratic Forms, Symmetrization, and Contraction

This section concerns mostly with decoupling, concentration of quadratic forms, symmetrization, and contraction, which are a number of basic tools of high-dimensional probability.

6.1 Decoupling

We'll look at quadratic forms of the form

$$\sum_{i,j=1}^n a_{ij} X_i X_j = X^T A X = \langle X, A X \rangle$$

where $A = (a_{ij})$ is an $n \times n$ coefficient matrix and $X = (X_1, \dots, X_n)$ is a random vector with independent coordinates. Such quadratic forms are known as chaos.

We can compute the expectation of a chaos. If X_i have zero means and unit variances, then

$$\mathbb{E}[X^T A X] = \sum_{i,j=1}^n a_{ij} \mathbb{E}[X_i X_j] = \sum_{i=1}^n a_{ii} = \text{tr}(A).$$

However, establishing concentration on a chaos is harder, because the terms of the sum above are not independent. However, we can overcome this difficulty via decoupling. We'll replace the quadratic form above with the bilinear form

$$\sum_{i,j=1}^n a_{ij} X_i X'_j = X^T A X' = \langle X, A X' \rangle,$$

where $X' = (X'_1, \dots, X'_n)$ is an independent copy of X . Bilinear forms are easier to analyze than quadratic forms as they are linear in X . Therefore if we condition on X' , we may treat the bilinear form as a sum of independent random variables

$$\sum_{i=1}^n \left(\sum_{j=1}^n a_{ij} X'_j \right) X_i = \sum_{i=1}^n b_i X_i$$

with fixed coefficients b_i .

Theorem 6.1.1 (Decoupling). Let A be an $n \times n$ diagonal free matrix, i.e. all diagonal entries are zero. Let X be a random vector in \mathbb{R}^n with independent mean zero coordinates, and let X' be an independent copy. Then for every convex function $F : \mathbb{R} \rightarrow \mathbb{R}$,

$$\mathbb{E}[F(X^T A X)] \leq \mathbb{E}[F(4X^T A X')].$$

Proof. We'll replace the chaos by a partial chaos, which we extend back to the original chaos later via Jensen's inequality. The partial chaos is defined by

$$\sum_{(i,j) \in I \times I^c} a_{ij} X_i X_j$$

where $I \subset \{1, \dots, n\}$ is a randomly chosen subset of indices.

(Step 1: Randomly selecting a partial sum) To specify a random subset of indices I , we'll use selectors - independent Bernoulli random variables $\delta_1, \dots, \delta_n \sim_{iid} \text{Ber}(1/2)$. We define the index set

$$I := \{i : \delta_i = 1\}.$$

Condition on X . Since by assumption $a_{ii} = 0$ and

$$\mathbb{E}[\delta_i(1 - \delta_j)] = \frac{1}{2} \cdot \frac{1}{2} = \frac{1}{4} \text{ for all } i \neq j$$

we may express the chaos as

$$X^T A X = \sum_{i \neq j} a_{ij} X_i X_j = 4\mathbb{E}_\delta \left[\sum_{i \neq j} \delta_i (1 - \delta_j) a_{ij} X_i X_j \right] = 4\mathbb{E}_I \left[\sum_{(i,j) \in I \times I^c} a_{ij} X_i X_j \right].$$

(In the expression above, the subscripts δ and I indicate the source of randomness in the conditional expectations. Since X is fixed, the expectations are taken over the random selection of $\delta = (\delta_1, \dots, \delta_n)$, or equivalently, the random index set I).

(Step 2: Applying F) Applying the function F to both sides and take expectation over X . By Jensen inequality and Fubini theorem, we get

$$\mathbb{E}_X[F(X^T A X)] \leq \mathbb{E}_I \left[\mathbb{E}_X \left[F \left(4 \sum_{(i,j) \in I \times I^c} a_{ij} X_i X_j \right) \right] \right].$$

It follows that there exists a realization of a subset I such that

$$\mathbb{E}_X[F(X^T A X)] \leq \mathbb{E}_X \left[F \left(4 \sum_{(i,j) \in I \times I^c} a_{ij} X_i X_j \right) \right].$$

Fix such a realization I until the end of the proof, and drop the subscript X on the expectation for convenience. Since the random variables $(X_i)_{i \in I}$ are independent from $(X_j)_{j \in I^c}$, the distribution of the sum in the right side will not change if we replace X_j by X'_j hence

$$\mathbb{E}_X[F(X^T A X)] \leq \mathbb{E} \left[F \left(4 \sum_{(i,j) \in I \times I^c} a_{ij} X_i X'_j \right) \right].$$

(Step 3: Completing the partial sum) It remains to complete the sum on the RHS to the sum over all pairs of indices. We want to show that

$$\mathbb{E} \left[F \left(4 \sum_{(i,j) \in I \times I^c} a_{ij} X_i X'_j \right) \right] \leq \mathbb{E} \left[F \left(4 \sum_{(i,j) \in [n] \times [n]} a_{ij} X_i X'_j \right) \right]$$

where $[n] = \{1, \dots, n\}$. To do this, we can decompose the sum on the right side as

$$\sum_{(i,j) \in [n] \times [n]} a_{ij} X_i X'_j = \underbrace{\sum_{(i,j) \in I \times I^c} a_{ij} X_i X'_j}_Y + \underbrace{\sum_{(i,j) \in I \times I} a_{ij} X_i X'_j + \sum_{(i,j) \in I^c \times [n]} a_{ij} X_i X'_j}_Z$$

Condition on all $(X_i)_{i \in I}$ and $(X'_j)_{j \in I^c}$, and denote this expectation by \mathbb{E}' . This fixes Y , while Z has zero conditional expectation (check). Thus, by Jensen inequality, we get

$$F(4Y) = F(4Y + \mathbb{E}'[4Z]) = F(\mathbb{E}'[4Y + 4Z]) \leq \mathbb{E}'[F(4Y + 4Z)].$$

Finally, taking expectations over all remaining random variables, we get

$$\mathbb{E}[F(4Y)] \leq \mathbb{E}[F(4Y + 4Z)].$$

Hence the proof is complete. □

Remark 6.1.2 (Diagonal-free assumption). The assumption is essential in Theorem 6.1.1, since the conclusion fails for diagonal matrices when $F(x) = x$. But we can include the diagonal on the right

hand side: for any $n \times n$ matrix $A = (a_{ij})$, we get

$$\mathbb{E} \left[F \left(\sum_{i \neq j} a_{ij} X_i X_j \right) \right] \leq \mathbb{E} \left[F \left(4 \sum_{i,j} a_{ij} X_i X'_j \right) \right]$$

This is shown in Exercise 6.1, and there are other variants of decoupling (Exercises 6.2-6.4).

6.2 Hanson-Wright Inequality

If X is a subgaussian random vector in \mathbb{R}^n , what can we say about its norm? If X has independent entries, then it concentrated (Theorem 3.1.1). But in general, it does not have to - it can be too small with high probability (Exercise 3.37). However, it can't be too large:

Proposition 6.2.1 (Norm of subgaussian random vector). Let X be a mean zero subgaussian random vector in \mathbb{R}^n with $\|X\|_{\psi_2} \leq K$. Then for every $t \geq 0$,

$$P(\|X\|_2 \geq CK(\sqrt{n} + t)) \leq e^{-t^2}.$$

Proof. WLOG, we can assume that $K = 1$. Squaring and exponentiating both sides and using Markov's inequality, we get

$$P(c\|X\|_2 \geq \sqrt{n} + t) \leq e^{-(n+t^2)} \mathbb{E} [\exp(c^2\|X\|_2^2)].$$

Now we will use a **Gaussian replacement** trick: for some absolute constant $c > 0$, we claim that

$$\mathbb{E} [\exp(c^2\|X\|_2^2)] \leq \mathbb{E} [\exp(\|g\|_2^2/6)] \text{ where } g \sim N(0, I_n).$$

To see this, condition on X (treating it as a fixed vector); then $\langle g, X \rangle \sim N(0, \|X\|_2^2)$ by Corollary 3.3.2, so (using the MGF of a normal distribution)

$$\mathbb{E} [\exp(c^2\|X\|_2^2)] = \mathbb{E}_g [\exp(\sqrt{2}c \langle g, X \rangle)],$$

where \mathbb{E}_g denotes the conditional expectation over g . Now takes expectation over X over both sides and apply Fubini:

$$\mathbb{E}_X [\exp(c^2\|X\|_2^2)] = \mathbb{E}_X [\mathbb{E}_g [\exp(\sqrt{2}c \langle g, X \rangle)]] = \mathbb{E}_g [\mathbb{E}_X [\exp(\sqrt{2}c \langle X, g \rangle)]].$$

When we condition on g (treating g as a fixed vector), the subgaussian norm of $\langle X, g \rangle$ is at most 1 by assumption ($K = 1$), so Proposition 2.6.6 (iv) gives

$$\mathbb{E}_X [\exp(\sqrt{2}c \langle X, g \rangle)] \leq \exp(\|g\|_2^2/4)$$

for some absolute constant $c > 0$. Substitute this into the bound above, then our claim for the Gaussian replacement is proved. After that, substitute into the first equation, and the proof is complete. \square

The Gaussian replacement trick will also be useful when we are proving concentration regarding a chaos - namely, the Hanson-Wright inequality:

Theorem 6.2.2 (Hanson-Wright inequality). Let A be an $n \times n$ matrix. Let $X = (X_1, \dots, X_n) \in \mathbb{R}^n$ be a random vector with independent, mean-zero, subgaussian coordinates. Then, for every $t \geq 0$, we have

$$P(|X^T A X - \mathbb{E}[X^T A X]| \geq t) \leq 2 \exp \left[-c \min \left(\frac{t^2}{K^4 \|A\|_F^2}, \frac{t}{K^2 \|A\|} \right) \right],$$

where $K = \max_i \|X_i\|_{\psi_2}$.

The proof will be based on bounding the MGF of $X^T A X$. Here is the plan:

- (a) replace $X^T AX$ by $X^T AX'$ by decoupling;
- (b) replace $X^T AX'$ by $g^T Ag'$ using Gaussian replacement, for $g \sim N(0, I_n)$;
- (c) compute $g^T Ag'$ by diagonalizing A using the rotational invariance of $N(0, I_n)$.

We start with part (b):

Lemma 6.2.3 (Gaussian replacement). Let A be an $n \times n$ matrix. Let X be a mean zero, subgaussian random vector in \mathbb{R}^n with $\|X\|_{\psi_2} \leq K$, and X' be its independent copy. Let $g, g' \sim N(0, I_n)$ be independent. Then for any $\lambda \in \mathbb{R}$,

$$\mathbb{E} [\exp (\lambda X^T A X')] \leq \mathbb{E} [\exp (C K^2 \lambda g^T A g')].$$

Proof. Condition on X' and take expectation over X , which we denote \mathbb{E}_X . Then the random variable $\langle X, A X' \rangle$ is (conditionally) subgaussian, with subgaussian norm $\leq K \|A X'\|_2$. Then Proposition 2.6.6 (iv) gives

$$\mathbb{E}_X [\exp (\lambda X^T A X')] \leq \exp (C \lambda^2 K^2 \|A X'\|_2^2), \lambda \in \mathbb{R}.$$

Compare the above to the normal MGF formula. Applied to the normal random variable $g^T A X' = \langle g, A X' \rangle$ (still conditionally on X'), it gives

$$\mathbb{E}_g [\exp (\mu g^T A X')] = \exp (\mu^2 \|A X'\|_2^2 / 2), \mu \in \mathbb{R}.$$

Setting $\mu = \sqrt{2CK} \lambda$, we match the right hand sides of the two equations above and obtain

$$\mathbb{E}_X [\exp (\lambda X^T A X')] = \mathbb{E}_g [\exp (\sqrt{2CK} \lambda g^T A X')].$$

Then, taking expectation over X' on both sides, we see that we have replace X by g in the chaos, at a cost of the factor $\sqrt{2CK}$. Repeating the same thing for X' , we can replace X' with g' and get another factor of $\sqrt{2CK}$. \square

We now move to step (c):

Lemma 6.2.4 (MGF of a Gaussian quadratic form). Let $A = (a_{ij})$ be an $n \times n$ matrix, and let $g, g' \sim N(0, I_n)$ be independent. Then

$$\mathbb{E} [\exp (\lambda g^T A g')] \leq \exp (\lambda^2 \|A\|_F^2) \text{ whenever } |\lambda| \leq \frac{1}{2\|A\|}.$$

Proof. Let's use rotational invariance of the normal distribution to diagonalize A . With its singular value decomposition, $A = U \Sigma V^T$, we can write

$$g^T A g' = (U^T g)^T \Sigma (V^T g').$$

By the rotational invariance of the normal distribution (Proposition 3.3.1), $U^T g$ and $V^T g'$ are independent standard normal random vectors in \mathbb{R}^n . So,

$$g^T A g' \sim g^T \Sigma g' = \sum_{i=1}^n \sigma_i g_i g_i^T.$$

This is a sum of independent random variables, so

$$\mathbb{E} [\exp (\lambda g^T A g')] = \mathbb{E} \left[\prod_i \mathbb{E} [\exp (\lambda \sigma_i g_i g_i^T)] \right] = \prod_i \mathbb{E} [\exp (\lambda \sigma_i g_i g_i^T)].$$

Now, for each i and $t \in \mathbb{R}$, we have

$$\mathbb{E} [\exp (t g_i g_i^T)] = \mathbb{E} \left[\exp \left(\frac{t^2 g_i^2}{2} \right) \right] = \frac{1}{\sqrt{1-t^2}} \leq \exp (t^2) \text{ if } t^2 \leq \frac{1}{2}.$$

The first identity is done by conditioning on g_i and using the MGF formula for the normal random variable g' ; the other steps are just direct calculations. Substituting this bound with $t = \lambda\sigma_i$ into the product, we get

$$\mathbb{E} [\exp(\lambda g^T A g')] \leq \exp\left(\lambda^2 \sum_i \sigma_i^2\right) \text{ if } \lambda^2 \leq \frac{1}{2 \max_i \sigma_i^2}.$$

Since σ_i are the singular values of A , $\sum_i \sigma_i^2 = \|A\|_F^2$ and $\max_i \sigma_i = \|A\|$, hence the lemma is proved. \square

Now we move to the main proof!

Proof of Hanson-Wright inequality. Without loss of generality, assume $K = 1$. As usual, it is enough to bound the one-sided tail

$$p := P(X^T X - \mathbb{E}[X^T A X] \geq t).$$

This is because we can find the lower tail by just replacing A with $-A$. By combining the two tails, the proof would be complete.

In terms of the entries of $A = (a_{ij})$, we have

$$X^T A X = \sum_{i,j} a_{ij} X_i X_j \text{ and } \mathbb{E}[X^T A X] = \sum_i a_{ii} \mathbb{E}[X_i^2],$$

where we used the mean zero assumption and independence. So

$$X^T A X - \mathbb{E}[X^T A X] = \sum_i a_{ii}(X_i^2 - \mathbb{E}[X_i^2]) + \sum_{i \neq j} a_{ij} X_i X_j.$$

The problem then reduces to estimating the diagonal and off-diagonal sums:

$$p \leq P\left(\sum_i a_{ii}(X_i^2 - \mathbb{E}[X_i^2]) \geq t/2\right) + P\left(\sum_{i \neq j} a_{ij} X_i X_j \geq t/2\right) := p_1 + p_2.$$

Let's bound these probabilities!

Step 1: Diagonal sum. Since X_i are independent and subgaussian, $X_i^2 - \mathbb{E}[X_i^2]$ are independent, mean-zero, and subexponential. Also,

$$\|X_i^2 - \mathbb{E}[X_i^2]\|_{\psi_2} \lesssim \|X_i^2\|_{\psi_1} \lesssim \|X_i\|_{\psi_2}^2 \lesssim 1.$$

(The above follows from centering and [cramer:2.8.5](#)). Then, Bernstein's inequality ([Corollary 2.9.2](#)) gives

$$p_1 \leq \exp\left[-c \min\left(\frac{t^2}{\sum_i a_{ii}^2}, \frac{t}{\max_i |a_{ii}|}\right)\right] \leq \exp\left[-c \min\left(\frac{t^2}{\|A\|_F^2}, \frac{t}{\|A\|}\right)\right].$$

Step 2: Off-diagonal sum. Now we bound the off-diagonal sum

$$S := \sum_{i \neq j} a_{ij} X_i X_j.$$

Let $\lambda > 0$ be a parameter to be determined later. By [Merkov's inequality](#), we have

$$p_2 = P(S \geq t/2) = p(\lambda S \geq \lambda t/2) \leq \exp(-\lambda t/2) \mathbb{E}[\exp(\lambda S)].$$

We get

$$\begin{aligned} \mathbb{E}[\exp(\lambda S)] &\leq \mathbb{E}[\exp(4\lambda X^T A X')] \quad (\text{By decoupling}) \\ &\leq \mathbb{E}[\exp(C_1 \lambda g^T A g')] \quad (\text{By Lemma 6.2.3}) \\ &\leq \exp(C \lambda^2 \|A\|_F^2) \quad (\text{By Lemma 6.2.4}) \end{aligned}$$

whenever $|\lambda| \leq \frac{1}{2\|A\|}$. Putting this bound into the exponential bound above, we obtain

$$p_2 \leq \exp(-\lambda t/2 + C \lambda^2 \|A\|_F^2).$$

Optimizing over $0 \leq \lambda \leq \frac{1}{2\|A\|}$, we conclude that

$$p_2 \leq \exp\left[-c \min\left(\frac{t^2}{\|A\|_F^2}, \frac{t}{\|A\|}\right)\right].$$

To summarize, we obtained the desired bounds for the probabilities of the diagonal deviation p_1 and the off-diagonal deviation p_2 . Putting them together, we complete the proof of [Theorem 6.2.2](#). \square

6.3 Symmetrization

A random variable X is called symmetric if it has the same distribution as $-X$. A basic example is the Rademacher random variable, which takes values -1 and 1 with equal probabilities. Mean-zero normal random variables are also symmetric, while the exponential and Poisson distributions are not.

This section introduces symmetrization, a useful trick for reducing problems to symmetric distributions - and sometimes even to the Rademacher distribution. It is based on the following:

Lemma 6.3.1 (Constructing symmetric distributions). Let X be a random variable and ξ be an independent Rademacher random variables. Then

- (a) ξX and $\xi|X|$ are identically distributed and symmetric.
- (b) If X is symmetric, both ξX and $\xi|X|$ have the same distribution as X .
- (c) If X' is an independent copy of X , then $X - X'$ is symmetric.

Proof. We'll check that ξX is symmetric. For any interval $A \subset \mathbb{R}$, the law of total probability gives

$$\begin{aligned} P(\xi X \in A) &= P(\xi X \in A \mid \xi = 1) \cdot \frac{1}{2} + P(\xi X \in A \mid \xi = -1) \cdot \frac{1}{2} \\ &= \frac{1}{2}(P(X \in A) + P(-X \in A)). \end{aligned}$$

Let's also do this for $-\xi X$:

$$\begin{aligned} P(-\xi X \in A) &= P(-\xi X \in A \mid \xi = 1) \cdot \frac{1}{2} + P(-\xi X \in A \mid \xi = -1) \cdot \frac{1}{2} \\ &= \frac{1}{2}(P(-X \in A) + P(X \in A)). \end{aligned}$$

Therefore ξX and $-\xi X$ have the same CDF, meaning they have the same distribution.

The rest of the proof is in Exercise 6.16. □

Lemma 6.3.2 (Symmetrization). Let X_1, \dots, X_N be independent, mean zero random vectors in a normed space, and let $\varepsilon_1, \dots, \varepsilon_N$ be independent Rademacher random variables. Then

$$\frac{1}{2} \mathbb{E} \left[\left\| \sum_{i=1}^N \varepsilon_i X_i \right\| \right] \leq \mathbb{E} \left[\left\| \sum_{i=1}^N X_i \right\| \right] \leq 2 \mathbb{E} \left[\left\| \sum_{i=1}^N \varepsilon_i X_i \right\| \right].$$

Proof. (**Upper bound**) Let (X'_i) be an independent copy of (X_i) . Since $\sum_{i=1}^N X'_i$ has mean zero, we have

$$p := \mathbb{E} \left[\left\| \sum_i X_i \right\| \right] \leq \mathbb{E} \left[\left\| \sum_i X_i - \sum_i X'_i \right\| \right] = \mathbb{E} \left[\left\| \sum_i (X_i - X'_i) \right\| \right].$$

The inequality above comes from the fact that for independent random vectors Y and Z ,

$$\mathbb{E}[Z] = 0 \implies \mathbb{E}[\|Y\|] \leq \mathbb{E}[\|Y + Z\|].$$

Since $X_i - X'_i$ are symmetric random vectors, they have the same distribution as $\varepsilon_i(X_i - X'_i)$ by Lemma 6.3.1 (b). Then

$$\begin{aligned} p &\leq \mathbb{E} \left[\left\| \sum_i \varepsilon_i (X_i - X'_i) \right\| \right] \\ &\leq \mathbb{E} \left[\left\| \sum_i \varepsilon_i X_i \right\| \right] + \mathbb{E} \left[\left\| \sum_i \varepsilon_i X'_i \right\| \right] \quad (\text{Triangle inequality}) \\ &= 2 \mathbb{E} \left[\left\| \sum_i \varepsilon_i X_i \right\| \right] \quad (\text{The two terms are identically distributed}). \end{aligned}$$

(**Lower bound**) The argument is similar as the proof for the upper bound:

$$\begin{aligned}
\mathbb{E} \left[\left\| \sum_i \varepsilon_i X_i \right\| \right] &\leq \mathbb{E} \left[\left\| \sum_i \varepsilon_i (X_i - X'_i) \right\| \right] \\
&= \mathbb{E} \left[\left\| \sum_i (X_i - X'_i) \right\| \right] \quad (\text{Same distribution}) \\
&\leq \mathbb{E} \left[\left\| \sum_i X_i \right\| \right] + \mathbb{E} \left[\left\| \sum_i X'_i \right\| \right] \quad (\text{Triangle inequality}) \\
&= 2\mathbb{E} \left[\left\| \sum_i X_i \right\| \right] \quad (\text{Identical distribution}).
\end{aligned}$$

Question: Where did we use X_i 's independence? Do we need mean zero for both upper and lower bounds? \square

There are also other versions of symmetrization lemmas (Exercises 6.19-6.21).

6.4 Random Matrices with non-i.i.d. Entries

A typical application of symmetrization consist of two steps: First, replace random variables X_i with symmetric ones $\varepsilon_i X_i$, then condition on X_i so that all randomness comes from the signs ε_i . Hence this reduces the problems to Rademacher random variables. To illustrate this technique, let's bound the operator norm of a random matrix with independent, non-identically distributed entries:

Theorem 6.4.1 (Norm of random matrices with non-i.i.d. entries). Let A be an $n \times n$ symmetric random matrix with independent, mean zero entries above and on the diagonal. Then

$$\mathbb{E} \left[\max_i \|A_i\|_2 \right] \leq \mathbb{E} [\|A\|] \leq C\sqrt{\log n} \cdot \mathbb{E} \left[\max_i \|A_i\|_2 \right],$$

where A_i denotes the rows of A .

Proof. The lower bound is already done in Exercise 4.7.

For the upper bound, we will use symmetrization and the matrix Khintchine inequality (Theorem 5.4.14). Let's decompose A entry-by-entry, keeping symmetry in mind, like the proof of Theorem 5.5.1. Denote the standard basis of \mathbb{R}^n by e_1, \dots, e_n , then A can be expressed as a sum of independent, mean zero random matrices:

$$A = \sum_{i \leq j} Z_{ij}, \text{ where } Z_{ij} = \begin{cases} A_{ij}(e_i e_j^T + e_j e_i^T) & \text{if } i < j, \\ A_{ii} e_i e_i^T & \text{if } i = j. \end{cases}$$

By applying symmetrization (Lemma 6.3.2), we get

$$\mathbb{E} [\|A\|] = \mathbb{E} \left[\left\| \sum_{i \leq j} Z_{ij} \right\| \right] \leq 2\mathbb{E} \left[\left\| \sum_{i \leq j} \varepsilon_{ij} Z_{ij} \right\| \right] \quad (*)$$

where ε_{ij} are independent Rademacher random variables.

Condition on (Z_{ij}) , apply the matrix Khintchine inequality (Theorem 5.4.14) for $p = 1$, and take expectation over (Z_{ij}) using the law of total expectation, which gives

$$\mathbb{E} \left[\left\| \sum_{i \leq j} \varepsilon_{ij} Z_{ij} \right\| \right] \leq C\sqrt{\log n} \mathbb{E} \left[\left\| \sum_{i \leq j} Z_{ij}^2 \right\|^{1/2} \right]. \quad (**)$$

Since (Z_{ij}) is a diagonal matrix,

$$Z_{ij}^2 = \begin{cases} A_{ij}^2 (e_i e_j^T + e_j e_i^T) & \text{if } i < j, \\ A_{ii}^2 e_i e_i^T & \text{if } i = j. \end{cases}$$

Therefore,

$$\sum_{i \leq j} Z_{ij}^2 = \sum_{i=1}^n \left(\sum_{j=1}^n A_{ij}^2 \right) e_i e_i^T = \sum_{i=1}^n \|A_i\|_2^2 e_i e_i^T.$$

In other words, this is a diagonal matrix with diagonal entries equal to $\|A_i\|_2^2$. Since the operator norm of a diagonal matrix is the maximal absolute value of its entries, we get

$$\left\| \sum_{i \leq j} Z_{ij}^2 \right\| = \max_i \|A_i\|_2^2.$$

Substitute the bound above into (**) then into (*) completes the proof. \square

There is more practice on symmetrization as well (Exercises 6.22-6.29).

6.5 Application: Matrix Completion

Matrix completion is the process of recovering missing entries from a partially observed matrix. Of course, this is not possible without knowing something extra about the matrix. Let's show that for low-rank matrices, we can recover the missing entries algorithmically.

To describe the problem mathematically, consider an $n \times n$ matrix X with

$$\text{rank}(X) = r$$

where $r \ll n$. Suppose we are shown a few *randomly chosen entries* of X . Each entry X_{ij} is revealed to us independently with some probability $p \in (0, 1)$ and is hidden from us with probability $1 - p$. In other words, assume that we observe the $n \times n$ matrix Y with entries

$$Y_{ij} = \delta_{ij} X_{ij} \text{ where } \delta_{ij} \sim_{i.i.d.} \text{Ber}(p).$$

These δ_{ij} are *selectors*. If

$$p = \frac{m}{n^2}$$

then we observe m entries on average.

The question is, how can we recover X from Y ? Although X has small rank r , Y may not have small rank. To fix this, we can pick the best rank r approximation to Y . Properly scaled, this gives a good estimate of the original matrix X :

Theorem 6.5.1 (Matrix completion). Let \hat{X} be a best rank r approximation to $p^{-1}Y$. Then

$$\mathbb{E} \left[\frac{1}{n} \|\hat{X} - X\|_F \right] \leq C \sqrt{\frac{rn \log n}{m}} \|X\|_\infty,$$

as long as $m \geq n \log n$. Here $\|X\|_\infty = \max_{i,j} |X_{ij}|$ is the largest entry, NOT the usual matrix infinity norm!

Before we prove this, note that the recovery error

$$\frac{1}{n} \|\hat{X} - X\|_F = \left(\frac{1}{n} \sum_{i,j=1}^{n^2} \sum_{i,j=1}^n |\hat{X}_{ij} - X_{ij}|^2 \right)^{1/2}$$

represents the average error per entry (in the L^2 sense). If we choose the average number of observed entries m so that

$$M \geq C' rn \log n$$

with large constant C' , then Theorem 6.5.1 guarantees that the average error is much smaller than $\|X\|_\infty$. So, matrix completion is possible if the number of observed entries exceeds rn by a logarithmic margin.

Proof. We first bound the recovery error in the operator norm, and then pass to the Frobenius norm using the low-rank assumption.

Step 1: Bounding the error in the operator norm. Using the triangle inequality, we can split the error as follows:

$$\|\hat{X} - X\| \leq \|\hat{X} - p^{-1}Y\| + \|p^{-1}Y - X\|.$$

Since we have chosen \hat{X} as a best rank r approximation to $p^{-1}Y$, the second summand dominates, i.e. $\|\hat{X} - p^{-1}Y\| \leq \|p^{-1}Y - X\|$, so we have

$$\|\hat{X} - X\| \leq 2\|p^{-1}Y - X\| = \frac{2}{p}\|Y - pX\|.$$

Note that the matrix \hat{X} , which is tricky to handle, is gone in the bound. Instead, we get $Y - pX$, which is easier to understand since its entries,

$$(Y - pX)_{ij} = (\delta_{ij} - p)X_{ij},$$

are independent, mean-zero random variables. Using Theorem 6.4.1 (more precisely, Exercise 6.28), we get

$$\mathbb{E}[\|Y - pX\|] \leq C\sqrt{\log n} \left(\mathbb{E} \left[\max_{i=1,\dots,n} \|(Y - pX)_{i:}\|_2 \right] + \mathbb{E} \left[\max_{j=1,\dots,n} \|(Y - pX)_{:j}\|_2 \right] \right). \quad (*)$$

To bound the norms of the rows of $Y - pX$, we write them as

$$\|(Y - pX)_{i:}\|_2^2 = \sum_{j=1}^n (\delta_{ij} - p)^2 X_{ij}^2 \leq \sum_{j=1}^n (\delta_{ij} - p)^2 \cdot \|X\|_\infty^2,$$

and similarly for columns. These sums of independent random variables can be easily bounded using Bernstein's (or Chernoff's) inequality, which yields (Exercise 6.30)

$$\mathbb{E} \left[\max_{i=1,\dots,n} \sum_{j=1}^n (\delta_{ij} - p)^2 \right] \lesssim pn.$$

Combining with a similar bound for the columns and substituting into (*), we obtain

$$\mathbb{E}[\|Y - pX\|] \lesssim \sqrt{pn \log n} \|X\|_\infty.$$

Then, by the bound for $\|\hat{X} - X\|$ from earlier, we get

$$\mathbb{E}[\|\hat{X} - X\|] \lesssim \sqrt{\frac{n \log n}{p}} \|X\|_\infty.$$

Passing to the Frobenius norm. We have not used the low rank assumption yet, so we'll do this now. Since $\text{rank}(X) \leq r$ by assumption and $\text{rank}(\hat{X}) \leq r$ by construction, we have (Exercise 4.4)

$$\text{rank}(\hat{X} - X) \leq 2r \implies \|\hat{X} - X\|_F \leq \sqrt{2r} \|\hat{X} - X\|.$$

Taking expectations and using the bound on the error in the operator norm from step 1, we get

$$\mathbb{E}[\|\hat{X} - X\|_F] \lesssim \sqrt{\frac{rn \log n}{p}} \|X\|_\infty.$$

Dividing both sides by n , we can rewrite this bound as

$$\mathbb{E} \left[\frac{1}{n} \|\hat{X} - X\|_F \right] \lesssim \sqrt{\frac{rn \log n}{pn^2}} \|X\|_\infty.$$

From the definition above, $pn^2 = m$ so plugging in finishes the proof. \square

Remark 6.5.2 (Extensions). Theorem 6.5.1 can be extended and improved in many ways, such as to rectangular matrices (Exercise 6.31) and matrices with noisy observations (Exercise 6.32). It is less trivial but possible to remove the logarithmic factor from the error bound, and achieve zero error for noiseless observations!

6.6 Contraction Principle

There is one more useful inequality the text covers in the chapter:

Theorem 6.6.1 (Contraction principle). Let x_1, \dots, x_N be any vectors in a normed space, $(a_1, \dots, a_N) \in \mathbb{R}^N$, and $\varepsilon_1, \dots, \varepsilon_N$ be independent Rademacher random variables. Then

$$\mathbb{E} \left[\left\| \sum_{i=1}^N a_i \varepsilon_i x_i \right\| \right] \leq \|a\|_\infty \cdot \mathbb{E} \left[\left\| \sum_{i=1}^N \varepsilon_i x_i \right\| \right].$$

Proof. WLOG, assume that $\|a\|_\infty \leq 1$. Define the function

$$f(a) := \mathbb{E} \left[\left\| \sum_{i=1}^N a_i \varepsilon_i x_i \right\| \right].$$

Then $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is convex (Exercise 6.35).

We want to bound for f the set of points a satisfying $\|a\|_\infty \leq 1$, i.e. on the unit cube $[-1, 1]^N$. By the maximum principle (Exercises 1.4 & 1.5), the maximum of a convex function on the cube is attained at a vertex, where all $a_i = \pm 1$. For such a , the random variables $(\varepsilon_i a_i)$ have the same distribution as ε_i by symmetry. Thus

$$\mathbb{E} \left[\left\| \sum_{i=1}^N a_i \varepsilon_i x_i \right\| \right] = \mathbb{E} \left[\left\| \sum_{i=1}^N \varepsilon_i x_i \right\| \right],$$

thus

$$f(a) \leq \mathbb{E} \left[\left\| \sum_{i=1}^N \varepsilon_i x_i \right\| \right] \text{ whenever } \|a\|_\infty \leq 1,$$

which completes the proof. \square

As an application, we can prove a version of symmetrization but with Gaussian random variables $g_i \sim N(0, 1)$ instead of Rademachers.

Lemma 6.6.2 (Symmetrization with Gaussians). Let X_1, \dots, X_N be independent, mean zero random vectors in a normed space. Let $g_1, \dots, g_N \sim N(0, 1)$ be independent Gaussian random variables, which are also independent of X_i . Then

$$\frac{c}{\sqrt{\log N}} \mathbb{E} \left[\left\| \sum_{i=1}^N g_i X_i \right\| \right] \leq \mathbb{E} \left[\left\| \sum_{i=1}^N X_i \right\| \right] \leq 3 \mathbb{E} \left[\left\| \sum_{i=1}^N g_i X_i \right\| \right].$$

Proof. (**Upper bound**) By symmetrization (Lemma 6.3.2), we have

$$E := \mathbb{E} \left[\left\| \sum_{i=1}^N X_i \right\| \right] \leq 2 \mathbb{E} \left[\left\| \sum_{i=1}^N \varepsilon_i X_i \right\| \right].$$

To interject Gaussian random variables, recall that $\mathbb{E}[|g_i|] = \sqrt{2/\pi}$. Then we can continue the bound as follows:

$$\begin{aligned}
E &\leq 2\sqrt{\frac{\pi}{2}} \mathbb{E}_X \left[\left\| \sum_{i=1}^N \varepsilon_i \mathbb{E}_g[|g_i|] X_i \right\| \right] \\
&\leq 2\sqrt{\frac{\pi}{2}} \mathbb{E} \left[\left\| \sum_{i=1}^N \varepsilon_i |g_i| X_i \right\| \right] \quad (\text{Jensen inequality}) \\
&= 2\sqrt{\frac{\pi}{2}} \mathbb{E} \left[\left\| \sum_{i=1}^N g_i X_i \right\| \right].
\end{aligned}$$

The last equality holds since the random variables $(\varepsilon_i |g_i|)$ have the same joint distribution as (g_i) (Lemma 6.3.1 (b)).

(Lower bound) We have

$$\begin{aligned}
\mathbb{E} \left[\left\| \sum_{i=1}^N g_i X_i \right\| \right] &= \mathbb{E} \left[\left\| \sum_{i=1}^N \varepsilon_i g_i X_i \right\| \right] \quad (\text{Symmetry of } g_i) \\
&\leq \mathbb{E}_g \left[\mathbb{E}_X \left[\|g\|_\infty \mathbb{E}_\varepsilon \left[\left\| \sum_{i=1}^N \varepsilon_i X_i \right\| \right] \right] \right] \quad (\text{Theorem 6.6.1}) \\
&= \mathbb{E}_g \left[\|g\|_\infty \mathbb{E}_\varepsilon \left[\mathbb{E}_X \left[\left\| \sum_{i=1}^N \varepsilon_i X_i \right\| \right] \right] \right] \quad (\text{Independence}) \\
&\leq 2\mathbb{E}_g \left[\|g\|_\infty \mathbb{E}_X \left[\left\| \sum_{i=1}^N X_i \right\| \right] \right] \quad (\text{Lemma 6.3.2}) \\
&= 2\mathbb{E}[\|g\|_\infty] \cdot \mathbb{E} \left[\left\| \sum_{i=1}^N X_i \right\| \right] \quad (\text{Independence}).
\end{aligned}$$

Moreover, by Proposition 2.7.6,

$$\mathbb{E}[\|g\|_\infty] \leq C\sqrt{\log N}.$$

Plugging back gives the result. □

Remark 6.6.3 (Log factor is unavoidable). The logarithmic factor in Lemma 6.6.2 is necessary and optimal in general (Exercise 6.37), making Gaussian symmetrization weaker than Rademacher's.