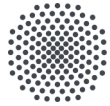


Robust Weighting and Matching Techniques for Causal Inference in Observational Studies with Continuous Treatment

Universität Stuttgart



Universität Stuttgart

Ioan Scheffel

December 21, 2022

Contents

1	Introduction	3
2	Balancing Weights	4
2.1	Introduction	4
2.2	Estimating the Population Mean of Potential Outcomes	4
2.3	Application of Convex Optimization	4
2.4	Application of Matrix Concentration Inequalities	5
3	Convex Analysis	6
3.1	Basic Notions	6
3.2	Relative Interior	6
3.3	Conjugate Calculus	6
3.4	Tseng Bertsekas	6
4	Random Matrix Inequalities	9
4.1	Exchangeable Pairs	9
4.1.1	Matrix Stein pairs	9
4.1.2	The method of exchangeable pairs	10
4.1.3	The conditional variance	10
4.2	Matrix Khintchin Inequality	11
4.3	Matrix Moment Inequality	13
4.4	Intrinsic Dimension	13
5	Empirical Processes	15
6	Simple yet useful Calculations	16

1 Introduction

2 Balancing Weights

2.1 Introduction

2.2 Estimating the Population Mean of Potential Outcomes

2.3 Application of Convex Optimization

Assumption 2.1. Assume that the map $f : \mathbb{R} \rightarrow \overline{\mathbb{R}}$ has the following properties.

- (i) f is strictly convex.
- (ii) f is lower-semicontinuous and continuously differentiable on $\text{int}(\text{dom}(f))$.
- (iii) The derivative of f on $\text{int}(\text{dom}(f))$ is a diffeomorphism.
- (iv) The Legendre transformation f^* of f is finite.
- (v) The function $x \mapsto xt - f(x)$ takes its supremum on $\text{int}(\text{dom}(f))$ for all $t \in \mathbb{R}$.

We consider the following optimization problem.

Problem 2.1.

$$\underset{w_1, \dots, w_n \in \mathbb{R}}{\text{minimize}} \quad \sum_{i=1}^n T_i f(w_i)$$

subject to the constraints

$$\begin{aligned} w_i T_i &\geq 0, & i &= 1, \dots, n, \\ \sum_{i=1}^n w_i T_i &= 1 \\ \left| \sum_{i=1}^n w_i T_i B_k(X_i) - \frac{1}{n} \sum_{i=1}^n B_k(X_i) \right| &\leq \delta_k, & k &= 1, \dots, K \end{aligned}$$

Theorem 2.1. *Under Assumption, the dual of the above Problem is the unconstrained optimization problem*

$$\underset{\lambda \in \mathbb{R}^K}{\text{minimize}} \quad \frac{1}{n} \sum_{i=1}^n nT_i f^*(\langle B(X_i), \lambda \rangle) - \langle B(X_i), \lambda \rangle + \langle \delta, |\lambda| \rangle,$$

where $t \mapsto f^*(t) = t(f')^{-1}(t) - f((f')^{-1}(t))$ is the Legendre transformation of f , $B(X_i) = [B_1(X_i), \dots, B_K(X_i)]^\top$ denotes the K basis functions of the covariates of unit $i \in \{1, \dots, n\}$ and $|\lambda| = [|\lambda_1|, \dots, |\lambda_K|]^\top$, where $|\cdot|$ is the absolute value of a real-valued scalar. Moreover, if λ^\dagger is an optimal solution then

$$w_i^* = (f')^{-1}(\langle B(X_i), \lambda^\dagger \rangle), \quad i \in \{1, \dots, n\} \quad (2.1)$$

are the unique optimal solutions to (P) .

Proof. We prove the following Lemma at the end of the section.

Lemma 2.1. *The dual of the optimization problem is*

$$\underset{\lambda \in \mathbb{R}^{2K}}{\text{minimize}} \quad \frac{1}{n} \sum_{i=1}^n nT_i f^*(\langle Q_{\bullet i}, \lambda \rangle) - \langle Q_{\bullet i}, \lambda \rangle + \langle d, \lambda \rangle$$

subject to

$$\lambda_k \geq 0 \quad \text{for all } k \in \{1, \dots, K\}, \quad (2.2)$$

where

$$\mathbf{Q} := \begin{bmatrix} \mathbf{I}_n \\ \mathbf{B}(\mathbf{X}) \\ -\mathbf{B}(\mathbf{X}) \end{bmatrix}, \quad \mathbf{B}(\mathbf{X}) := [B(X_1), \dots, B(X_n)], \quad \text{and} \quad d := \begin{bmatrix} 0_n \\ \delta \\ \delta \end{bmatrix}. \quad (2.3)$$

□

2.4 Application of Matrix Concentration Inequalities

3 Convex Analysis

3.1 Basic Notions

3.2 Relative Interior

3.3 Conjugate Calculus

3.4 Tseng Bertsekas

We present the relevant parts of the paper [BT03].

Consider the following optimization problem

$$\underset{x \in \mathbb{R}^m}{\text{minimize}} \quad f(x)$$

subject to the constraints

$$\mathbf{A}x \geq b, \tag{3.1}$$

Where $f : \mathbb{R}^m \rightarrow \overline{\mathbb{R}}$, \mathbf{A} is a given $n \times m$ matrix, and b is a vector in \mathbb{R}^n .

Assumption 3.1. Assume that the map $f : \mathbb{R}^m \rightarrow \overline{\mathbb{R}}$ has the following properties.

- (i) f is strictly convex.
- (ii) f is lower-semicontinuous and continuous $\text{dom}(f)$.
- (iii) The convex conjugate f^* of f is finite.

The dual optimization problem associated with (P) is

$$\underset{p \in \mathbb{R}^n}{\text{maximize}} \quad q(p)$$

subject to the constraints

$$p \geq 0, \tag{3.2}$$

where $q : \mathbb{R}^n \rightarrow \overline{\mathbb{R}}$ is the concave function given by

$$q(p) := \min_{x \in \mathbb{R}^m} f(x) + \langle p, b - \mathbf{A}x \rangle = \langle p, b \rangle - f^*(\mathbf{A}^\top p). \tag{3.3}$$

The dual problem (D) is a concave program with simple nonnegativity constraints. Furthermore, strong duality holds for (P) and (D), i.e., the optimal value of (P) equals the optimal value of (D).

Since f^* is real-valued and f is strictly convex, f^* and q are continuously differentiable.

Theorem 3.1. [Roc70, Theorem 26.3] *A closed proper convex function is (essentially) strictly convex if and only if its conjugate is essentially smooth.*

We will denote the gradient of q at p by $d(p)$ and its i th coordinate by $d_i(p)$. Since q is continuously differentiable, $d_i(p)$ is continuous, and since q is concave, $d_i(p)$ is nonincreasing in p_i .

By differentiating and by using the chain rule, we obtain the dual cost gradient

$$d(p) = b - \mathbf{A}x, \quad \text{where} \quad x := \nabla f^*(\mathbf{A}^\top p) = \operatorname{argsup}_{\xi \in \mathbb{R}^m} \langle p, \mathbf{A}\xi \rangle - f(\xi). \quad (3.4)$$

The last equality follows from Danskin's Theorem and [Roc70, Theorem 23.5]

Proposition 3.1. (Danskin's Theorem [BT03, page 649]) *Let $Z \subseteq \mathbb{R}^m$ be a non-empty set, and let $\phi : \mathbb{R}^n \times Z \rightarrow \mathbb{R}$ be a continuous function such that $\phi(\cdot, z) : \mathbb{R}^n \rightarrow \mathbb{R}$, viewed as a function of its first argument, is convex for each $z \in Z$. Then the function*

$$f : \mathbb{R}^n \rightarrow \mathbb{R}, \quad x \mapsto \sup_{z \in Z} \phi(x, z) \quad (3.5)$$

is convex and has directional derivative given by

$$f'(x; y) = \sup_{z \in Z(x)} \phi'(x, z; y), \quad (3.6)$$

where $\phi'(x, z; y)$ is the directional derivative of the function $\phi(\cdot, z)$ at x in the direction y , and

$$Z(x) := \left\{ \bar{z} \in \mathbb{R}^m : \phi(x, \bar{z}) = \sup_{z \in Z} \phi(x, z) \right\}. \quad (3.7)$$

In particular, if $Z(x)$ consists of a unique point \bar{z} and $\phi(\cdot, \bar{z})$ is differentiable at x , and $\nabla f(x) = \nabla_x \phi(x, \bar{z})$, where $\nabla_x \phi(x, \bar{z})$ is the vector with coordinates $(\partial \phi / \partial x_i)(x, \bar{z})$

Note that x is the unique vector satisfying

$$\mathbf{A}p \in \partial f(x). \quad (3.8)$$

From the optimality conditions for (D) it follows that a dual vector is an optimal solution of (D) if and only if

$$p = [p + d(p)]^+, \quad (3.9)$$

where $[\cdot]^+$ is the projection onto the positive orthant, i.e., $[y]^+ = [0 \vee y_1, \dots, 0 \vee y_n]^\top$.

Given an optimal dual solution p , we may obtain an optimal primal solution from the equation $x = \nabla f^*(\mathbf{A}^\top p)$. To see this, note that

$$\mathbf{A}x \geq b \quad \text{and} \quad p_i = 0 \quad \text{for all } i \text{ such that} \quad \sum_{j=1}^m a_{ij}x_j > b_i. \quad (3.10)$$

We can show that p and x satisfy the KKT conditions and thus x is an optimal solution to (P) .

Definition 3.1. [Roc70, §28] By an **ordinary convex program** (P) we mean an optimization problem of the following form

$$\underset{x \in C}{\text{minimize}} \quad f_0(x)$$

subject to the constraints

$$f_1(x) \leq 0, \dots, f_r(x) \leq 0, \quad f_{r+1}(x) = 0, \dots, f_m(x) = 0, \quad (3.11)$$

where $C \subseteq \mathbb{R}^n$ is a non-empty convex set, f_i is a finite convex function on C for $i \in \{1, \dots, r\}$ and f_i is an affine function on C for $i \in \{r+1, \dots, m\}$.

Definition 3.2. We define $[\lambda_1, \dots, \lambda_m] \in \mathbb{R}^m$ to be a **Karush-Kuhn-Tucker (KKT) vector** for (P) , if

- (i) $\lambda_i \geq 0$ for all $i \in \{1, \dots, r\}$.
- (ii) The infimum of the proper convex function $f_0 + \sum_{i=1}^m \lambda_i f_i$ is finite and equal to the optimal value in (P) .

Theorem 3.2. (Karush-Kuhn-Tucker conditions) Let (P) be an ordinary convex program, $\bar{\alpha} \in \mathbb{R}^m$, and $\bar{z} \in \mathbb{R}^n$. Then $\bar{\alpha}$ is a KKT vector for (P) and \bar{z} is an optimal solution to (P) if and only if \bar{z} and the components α_i of $\bar{\alpha}$ satisfy the following conditions.

- (i) $\alpha_i \geq 0$, $f_i(\bar{z}) \leq 0$, and $\alpha_i f_i(\bar{z}) = 0$ for all $i \in \{1, \dots, r\}$.
- (ii) $f_i(\bar{z}) = 0$ for $i \in \{r+1, \dots, m\}$.
- (iii) $0_n \in [\partial f_0(\bar{z}) + \sum_{\alpha_i \neq 0} \alpha_i \partial f_i(\bar{z})]$.

Proof. [Roc70, Theorem 28.3] □

Takeaways For strictly convex functions we can derive duality in terms of the optimal solutions.

4 Random Matrix Inequalities

In our application we want to bound moments of vector-valued random variables. For this we choose the theory of random matrix inequalities which lately received a lot of attention. In particular an approach via the method of exchangeable pairs [MJC⁺14] has been fruitful in simplifying the proofs of long standing results such as the matrix Khintchin inequality. We base our exposition on [MJC⁺14]. A lot will be exact copy of this paper, so no originality is claimed. Where it seemed fit, we conducted some calculations in more detail than presented in the paper.

We will first introduce the method of exchangeable pairs and derive auxiliary theorems to establish the matrix Khintchin inequality. Then we will derive inequalities for moments of matrices, first for psd matrices and then via the Hermitian dilataition for general rectangular matrices. In a last step we will introduce the notion of intrinsic dimension to improve the bounds.

4.1 Exchangeable Pairs

4.1.1 Matrix Stein pairs

We first define an exchangeable pair.

Definition 4.1. *Let Z and Z' random variables taking values in a Polish space \mathcal{Z} . We say that (Z, Z') is an **exchangeable pair** if it has the same distribution as (Z', Z) . In particular, Z and Z' must share the same distribution.*

The following approach originates in the work of Charles Stein [Ste72] on normal approximation for a sum of dependent random variable. We will explain how some central ideas of this theory extends to matrices.

We can obtain a lot of information about the fluctuation of a random matrix \mathbf{X} if we can construct a good exchangeable pair $(\mathbf{X}, \mathbf{X}')$. With this motivation in mind, let us introduce a special class of exchangeable pairs.

Definition 4.2. *Let (Z, Z') be an exchangeable pair of random variables taking values in a Polish space \mathcal{Z} , and let $\Psi : \mathcal{Z} \rightarrow \mathbb{H}_d$ be a measurable function. Define the random Hermitian matrices*

$$\mathbf{X} := \Psi(Z) \quad \text{and} \quad \mathbf{X}' := \Psi(Z'). \quad (4.1)$$

*We say that $(\mathbf{X}, \mathbf{X}')$ is a **matrix Stein pair** if there is a constant $\alpha \in (0, 1]$ for which*

$$\mathbf{E}[\mathbf{X} - \mathbf{X}' | Z] = \alpha \mathbf{X} \quad \text{almost surely.} \quad (4.2)$$

*The constant α is called the **scale factor** of the pair. We always assume $\mathbf{E}[\|\mathbf{X}\|^2] < \infty$.*

A matrix Stein pair $(\mathbf{X}, \mathbf{X}')$ has several useful properties. First, $(\mathbf{X}, \mathbf{X}')$ always forms an exchangeable pair. Second, it must be the case that $\mathbf{E}[\mathbf{X}] = \mathbf{0}$. Indeed,

$$\mathbf{E}[\mathbf{X}] = \frac{1}{\alpha} \mathbf{E}[\mathbf{E}[\mathbf{X} - \mathbf{X}' | Z]] = \frac{1}{\alpha} \mathbf{E}[\mathbf{X} - \mathbf{X}'] = \mathbf{0}.$$

4.1.2 The method of exchangeable pairs

A well-chosen matrix Stein pair $(\mathbf{X}, \mathbf{X}')$ provides a surprisingly powerful tool for studying the random matrix \mathbf{X} . The technique depends on a fundamental technical lemma.

Lemma 4.1. *Suppose that $(\mathbf{X}, \mathbf{X}')$ is a matrix Stein pair with scale factor α . Let $\mathbf{F} : \mathbb{H}_d \rightarrow \mathbb{H}_d$ be a measurable function that satisfies the regularity condition $\mathbf{E} [\|(\mathbf{X} - \mathbf{X}')\mathbf{F}(\mathbf{X})\|] < \infty$. Then*

$$\mathbf{E}[\mathbf{X} \cdot \mathbf{F}(\mathbf{X})] = \frac{1}{2\alpha} \mathbf{E}[(\mathbf{X} - \mathbf{X}')(\mathbf{F}(\mathbf{X}) - \mathbf{F}(\mathbf{X}'))]. \quad (4.3)$$

In short, the randomness in the Stein pair furnishes an alternative expression for the expected product of \mathbf{X} and a function \mathbf{F} . It allows us to estimate the expectation using the smoothness properties of the function \mathbf{F} and the discrepancy between \mathbf{X} and \mathbf{X}' .

Proof. [MJC⁺14, Lemma 2.4] Suppose that $(\mathbf{X}, \mathbf{X}')$ constructed from an auxiliary exchangeable pair (Z, Z') . The defining property implies

$$\alpha \cdot \mathbf{E}[\mathbf{X} \cdot \mathbf{F}(\mathbf{X})] = \mathbf{E}[\mathbf{E}[\mathbf{X} - \mathbf{X}' | Z] \cdot \mathbf{F}(\mathbf{X})] = \mathbf{E}[(\mathbf{X} - \mathbf{X}')\mathbf{F}(\mathbf{X})] \quad (4.4)$$

□

4.1.3 The conditional variance

To each matrix Stein pair $(\mathbf{X}, \mathbf{X}')$, we may associate a random matrix called the conditional variance of \mathbf{X} . The purpose of this section is to argue that the spectral norm of \mathbf{X} is unlikely to be large, when the conditional variance is small.

Definition 4.3. *Suppose that $(\mathbf{X}, \mathbf{X}')$, is a matrix Stein pair, constructed from an auxiliary exchangeable pair (Z, Z') . The **conditional variance** is the random matrix*

$$\Delta_{\mathbf{X}} := \Delta_{\mathbf{X}}(Z) := \frac{1}{2\alpha} \mathbf{E}[(\mathbf{X} - \mathbf{X}')^2 | Z], \quad (4.5)$$

where α is the scale factor of the pair. We may take any version of the conditional expectation in this definition.

The conditional variance $\Delta_{\mathbf{X}}$ can be regarded as a stochastic estimate for the variance of the random matrix \mathbf{X} . To see this, assume the second moment of \mathbf{X} exists. Then it follows from Lemma with $\mathbf{F}(\mathbf{X}) = \mathbf{X}$

$$\mathbf{E}[\Delta_{\mathbf{X}}] = \mathbf{E}[\mathbf{X}^2]. \quad (4.6)$$

To verify the regularity condition, note that

$$\mathbf{E}[\|(\mathbf{X} - \mathbf{X}')\mathbf{X}\|] \leq \mathbf{E}[\|\mathbf{X}\|^2] + \mathbf{E}[\|\mathbf{X}\| \cdot \|\mathbf{X}'\|] \leq 2\mathbf{E}[\|\mathbf{X}\|^2] < \infty. \quad (4.7)$$

Example 4.1. [MJC⁺14, Example 2.4] ◇

nrt
is not necessary, as, for example, the normal distribution does not fulfill it. For

4.2 Matrix Khintchin Inequality

The goal of this section is to derive the matrix Khintchin inequality and show some important applications. For this we need an auxiliary theorem which is an extension of the *Burkholder-Davis-Gundy (BDG) inequality* from classical martingale theory [Bur73]. We prepare for the proof of this theorem by assembling some analytic tools.

Proposition 4.1. (Generalized Klein inequality) *Let u_1, \dots, u_n and v_1, \dots, v_n be real-valued functions on an interval I of the real line. Suppose*

$$\sum_{k=1}^n u_k(a)v_k(b) \geq 0 \quad \text{for all } a, b \in I. \quad (4.8)$$

Then

$$\overline{\text{tr}} \left(\sum_{k=1}^n u_k(\mathbf{A})v_k(\mathbf{B}) \right) \geq 0 \quad \text{for all } \mathbf{A}, \mathbf{B} \in \mathbb{H}_d(I). \quad (4.9)$$

Proof. [Pet94, Proposition 3] Let $\mathbf{A} = \sum_{i=1}^d \lambda_i \mathbf{P}_{\bullet i} \mathbf{P}_{\bullet i}^*$ and $\mathbf{B} = \sum_{j=1}^d \mu_j \mathbf{Q}_{\bullet j} \mathbf{Q}_{\bullet j}^*$ be the orthonormal decompositions of \mathbf{A} and \mathbf{B} . Then

$$\overline{\text{tr}} \left(\sum_{k=1}^n u_k(\mathbf{A})v_k(\mathbf{B}) \right) = \sum_{k=1}^n \sum_{i,j=1}^d \overline{\text{tr}} (u_k(\lambda_i) \mathbf{P}_{\bullet i} \mathbf{P}_{\bullet i}^* v_k(\mu_j) \mathbf{Q}_{\bullet j} \mathbf{Q}_{\bullet j}^*) \quad (4.10)$$

$$= \sum_{i,j=1}^d \overline{\text{tr}} (\mathbf{P}_{\bullet i} \mathbf{P}_{\bullet i}^* \mathbf{Q}_{\bullet j} \mathbf{Q}_{\bullet j}^*) \sum_{k=1}^n u_k(\lambda_i) v_k(\mu_j) \geq 0 \quad (4.11)$$

by the hypothesis. To see that $\overline{\text{tr}} (\mathbf{P}_{\bullet i} \mathbf{P}_{\bullet i}^* \mathbf{Q}_{\bullet j} \mathbf{Q}_{\bullet j}^*)$ is non-negative for all $i, j \in \{1, \dots, d\}$, we apply a well known extension of von Neumann's trace inequality [Ruh70, Lemma 1], namely

$$\text{tr}(\mathbf{P}\mathbf{Q}) \geq \sum_{i=1}^d p_i q_{d-i+1} \geq 0 \quad \text{for all } \mathbf{P}, \mathbf{Q} \in \mathbb{H}_d([0, \infty)), \quad (4.12)$$

where the eigenvalues $p_1 \geq \dots \geq p_d$ and $q_1 \geq \dots \geq q_d$ are sorted decreasingly. \square

Lemma 4.2. (Mean value trace inequality) *Let I be an interval of the real line. Suppose that $g : I \rightarrow \mathbb{R}$ is a weakly increasing function and that $h : I \rightarrow \mathbb{R}$ is a function whose derivative h' is convex. Then for all matrices $\mathbf{A}, \mathbf{B} \in \mathbb{H}_d(I)$ it holds*

$$\overline{\text{tr}}[(g(\mathbf{A}) - g(\mathbf{B})) \cdot (h(\mathbf{A}) - h(\mathbf{B}))] \leq \frac{1}{2} \overline{\text{tr}}[(g(\mathbf{A}) - g(\mathbf{B})) \cdot (\mathbf{A} - \mathbf{B}) \cdot (h'(\mathbf{A}) + h'(\mathbf{B}))]. \quad (4.13)$$

When h' is concave, the inequality is reversed. The same result holds for the standard trace.

Proof. [MJC⁺14, Lemma 3.4] Fix $a, b \in I$. Since g is weakly increasing, $(g(a) - g(b)) \cdot (a - b) \geq 0$. The fundamental theorem of calculus and the convexity of h' yield the estimate

$$(g(a) - g(b)) \cdot (h(a) - h(b)) = (g(a) - g(b)) \cdot (a - b) \int_0^1 h'(\tau a + (1 - \tau)b) d\tau \quad (4.14)$$

$$\leq (g(a) - g(b)) \cdot (a - b) \int_0^1 [\tau h'(a) + (1 - \tau)h'(b)] d\tau \quad (4.15)$$

$$= \frac{1}{2} [(g(a) - g(b)) \cdot (a - b) \cdot (h'(a) + h'(b))]. \quad (4.16)$$

The inequality is reversed, if h' is concave. To apply the Kleins inequality we expand the terms. The RHS is

$$\begin{aligned} & (g(a) - g(b)) \cdot (a - b) \cdot (h'(a) + h'(b)) \\ &= [g(a) \cdot a \cdot h'(a)] + [g(a) \cdot a] \cdot h'(b) - b \cdot [h'(a) \cdot g(a)] - [b \cdot h'(b)] \cdot g(a) \\ &+ [\text{the same as above with } a \text{ and } b \text{ reversed}](a \rightleftharpoons b) \end{aligned} \quad (4.17)$$

Taking the trace yields

$$\begin{aligned} & \text{tr}[g(\mathbf{A}) \cdot \mathbf{A} \cdot (h'(\mathbf{A}) + h'(\mathbf{B}))] - \text{tr}[\mathbf{B} \cdot (h'(\mathbf{A}) + h'(\mathbf{B})) \cdot g(\mathbf{A})] + (\mathbf{A} \rightleftharpoons \mathbf{B}) \\ &= \text{tr}[g(\mathbf{A}) \cdot \mathbf{A} \cdot (h'(\mathbf{A}) + h'(\mathbf{B}))] - \text{tr}[g(\mathbf{A}) \cdot \mathbf{B} \cdot (h'(\mathbf{A}) + h'(\mathbf{B}))] + (\mathbf{A} \rightleftharpoons \mathbf{B}) \\ &= \text{tr}[g(\mathbf{A}) \cdot (\mathbf{A} - \mathbf{B}) \cdot (h'(\mathbf{A}) + h'(\mathbf{B}))] + (\mathbf{A} \rightleftharpoons \mathbf{B}) \\ &= \text{tr}[(g(\mathbf{A}) - g(\mathbf{B})) \cdot (\mathbf{A} - \mathbf{B}) \cdot (h'(\mathbf{A}) + h'(\mathbf{B}))]. \end{aligned} \quad (4.18)$$

On the LHS we have only products of two factors which commute under the trace operation. Thus we may use the same expression as in the scalar case without further calculations. The result follows immediately from the Klein inequality. \square

Proposition 4.2. (Hölder inequality for trace) *Let p and q be Hölder conjugate indices. Then*

$$\text{tr}(\mathbf{BC}) \leq \|\mathbf{B}\|_p \|\mathbf{C}\|_q \quad \text{for all } \mathbf{B}, \mathbf{C} \in \mathbb{M}_d. \quad (4.19)$$

Proof. [Bha97, Corollary IV.2.6] \square

We are now ready to prove the auxiliary theorem.

Theorem 4.1. (Matrix BDG inequality) *Let $p = 1$ or $p \geq 3/2$. Suppose that $(\mathbf{X}, \mathbf{X}')$ is a matrix Stein pair where $\mathbf{E}[\|\mathbf{X}\|_{2p}^{2p}] < \infty$. Then*

$$\mathbf{E}[\|\mathbf{X}\|_{2p}^{2p}]^{1/(2p)} \leq \sqrt{2p-1} \mathbf{E}[\|\Delta_{\mathbf{X}}\|_p^p]^{1/(2p)}, \quad (4.20)$$

where $\Delta_{\mathbf{X}}$ is the conditional variance .

Proof. [MJC⁺14, §7.3] Suppose that $(\mathbf{X}, \mathbf{X}')$ is a matrix Stein pair with scale factor α . First, observe that the result for $p = 1$ already follows from $\mathbf{E}[\Delta_{\mathbf{X}}] = \mathbf{E}[\mathbf{X}^2]$. Therefore we may assume that $p \geq 3/2$. We introduce the notation for the quantity of interest,

$$E := \mathbf{E}[\|\mathbf{X}\|_{2p}^{2p}] = \mathbf{E}[\text{tr}(|\mathbf{X}|^{2p})]. \quad (4.21)$$

We rewrite the expression for E by peeling off a copy of $|\mathbf{X}|$. This yields

$$E = \mathbf{E}[\text{tr}(|\mathbf{X}| \cdot |\mathbf{X}|^{2p-1})] = \mathbf{E}[\text{tr}(\mathbf{X} \cdot \text{sgn}(\mathbf{X}) \cdot |\mathbf{X}|^{2p-1})]. \quad (4.22)$$

Apply the method of exchangeable pairs with $\mathbf{F}(\mathbf{X}) = \text{sgn}(\mathbf{X}) \cdot |\mathbf{X}|^{2p-1}$ to reach

$$E = \frac{1}{2\alpha} \mathbf{E}[\text{tr}((\mathbf{X} - \mathbf{X}') \cdot (\text{sgn}(\mathbf{X}) \cdot |\mathbf{X}|^{2p-1} - \text{sgn}(\mathbf{X}') \cdot |\mathbf{X}'|^{2p-1}))] \quad (4.23)$$

Apply method of exchangeable pairs, generalized Klein inequality, trace Hölder \square

Theorem 4.2. [MJC⁺14, Corollary 7.3] Suppose that $p = 1$ or $p \geq 3/2$. Consider a finite sequence $(\mathbf{Y}_k)_{k \geq 1}$ of independent, random, Hermitian matrices and a deterministic sequence $(\mathbf{A}_k)_{k \geq 1}$ for which

$$\mathbf{E}[\mathbf{Y}_k] = 0 \quad \text{and} \quad \mathbf{Y}_k^2 \preceq \mathbf{A}_k^2 \quad \text{almost surely for all } k \geq 1. \quad (4.24)$$

Then

$$\mathbf{E} \left[\left\| \sum_{k \geq 1} \mathbf{Y}_k \right\|_{2p}^{2p} \right]^{1/(2p)} \leq \sqrt{p - \frac{1}{2}} \left\| \left(\sum_{k \geq 1} (\mathbf{A}_k^2 + \mathbf{E}[\mathbf{Y}_k^2]) \right)^{1/2} \right\|_{2p}. \quad (4.25)$$

In particular, when $(\xi_k)_{k \geq 1}$ is an independent sequence of Rademacher random variables,

$$\mathbf{E} \left[\left\| \sum_{k \geq 1} \xi_k \mathbf{A}_k \right\|_{2p}^{2p} \right]^{1/(2p)} \leq \sqrt{2p - 1} \left\| \left(\sum_{k \geq 1} \mathbf{A}_k^2 \right)^{1/2} \right\|_{2p}. \quad (4.26)$$

4.3 Matrix Moment Inequality

Theorem 4.3. Assume $n \geq 3$

(i) Suppose that $p \geq 1$, and fix $r \geq p \vee 2 \log(n)$. Consider a finite sequence $(\mathbf{S}_k)_{k \geq 1}$ of independent, random, positive-semidefinite matrices with dimension $n \times n$. Then

$$\mathbf{E} \left[\left\| \sum_{k \geq 1} \mathbf{S}_k \right\|^p \right]^{1/p} \leq \left[\left\| \sum_{k \geq 1} \mathbf{E}[\mathbf{S}_k] \right\|^{1/2} + 2\sqrt{er} \mathbf{E}[\max_{k \geq 1} \|\mathbf{S}_k\|^p]^{1/(2p)} \right]^2. \quad (4.27)$$

(ii) Suppose that $p \geq 2$, and fix $r \geq p \vee 2 \log(n)$. Consider a finite sequence $(\mathbf{Y}_k)_{k \geq 1}$ of independent, symmetric, random, self-adjoint matrices with dimension $n \times n$. Then

$$\mathbf{E} \left[\left\| \sum_{k \geq 1} \mathbf{Y}_k \right\|^p \right]^{1/p} \leq \sqrt{er} \left\| \left(\sum_{k \geq 1} \mathbf{E}[\mathbf{Y}_k^2] \right)^{1/2} \right\| + 2er \mathbf{E}[\max_{k \geq 1} \|\mathbf{S}_k\|^p]^{1/p}. \quad (4.28)$$

4.4 Intrinsic Dimension

Definition 4.4. For a positive-semidefinite matrix \mathbf{S} , the *intrinsic dimension* is the quantity

$$\text{intdim}(\mathbf{A}) := \frac{\text{tr} \mathbf{A}}{\|\mathbf{A}\|}.$$

Lemma 4.3. (Intrinsic dimension) *Let $\varphi : [0, \infty) \rightarrow \mathbb{R}$ be a convex function with $\varphi(0) = 0$. For any positive-semidefinite matrix \mathbf{S} it holds that*

$$\mathrm{tr}(\varphi(\mathbf{S})) \leq \mathrm{intdim}(\mathbf{S}) \cdot \varphi(\|\mathbf{S}\|).$$

Proof. [Tro15, Lemma 7.5.1] Since φ is convex on any interval $[0, L]$ with $L > 0$ and $\varphi(0) = 0$, it holds

$$\varphi(a) \leq \left(1 - \frac{a}{L}\right) \varphi(0) + \frac{a}{L} \varphi(L) = \frac{a}{L} \varphi(L) \quad \text{for all } a \in [0, L]. \quad (4.29)$$

Since \mathbf{S} is positive-semidefinite, the eigenvalues of \mathbf{S} fall in the interval $[0, L]$, where $L = \|\mathbf{S}\|$.

$$\mathrm{tr}(\varphi(\mathbf{S})) = \sum_{i=1}^d \varphi(\lambda_i) \leq \frac{\sum_{i=1}^d \lambda_i}{\|\mathbf{S}\|} \varphi(\|\mathbf{S}\|) = \frac{\mathrm{tr}(\mathbf{S})}{\|\mathbf{S}\|} \varphi(\|\mathbf{S}\|) = \mathrm{intdim}(\mathbf{S}) \cdot \varphi(\|\mathbf{S}\|). \quad (4.30)$$

□

5 Empirical Processes

Let $(\Omega, \mathcal{A}, \mathbf{P})$ be a probability space and (\mathcal{X}, Σ) a measurable space. Let $X_j : (\Omega, \mathcal{A}, \mathbf{P}) \rightarrow (\mathcal{X}, \Sigma), j = 1, \dots, n$ be independent and identically-distributed (i.i.d.) random variables with probability distribution \mathbf{P}_X and \mathcal{F} a family of measurable functions $f : (\mathcal{X}, \Sigma) \rightarrow (\mathbb{R}, \mathcal{B}(\mathbb{R}))$. Consider the map

$$f \mapsto G_n f := \sqrt{n} \left(\frac{1}{n} \sum_{i=1}^n f(X_i) - \mathbf{P}_X f \right), \quad (5.1)$$

where $\mathbf{P}_X f := \int_{\mathcal{X}} f d\mathbf{P}_X$. We call $(G_n f)_{f \in \mathcal{F}}$ the empirical process indexed by \mathcal{F} . Furthermore

$$\|G_n f\|_{\mathcal{F}} := \sup_{f \in \mathcal{F}} |G_n f|. \quad (5.2)$$

Lemma 5.1. (Bernstein Inequality for Empirical Processes) *For any bounded, measurable function f it holds for all $t > 0$*

$$\mathbf{P}(|G_n f| > t) \leq 2 \exp \left(-\frac{1}{4} \frac{t^2}{\mathbf{P}_X(f^2) + t \|f\|_{\infty} / \sqrt{n}} \right) \quad (5.3)$$

Proof. By the Markov inequality it holds for all $\lambda > 0$

$$\mathbf{P}(G_n f > t) \leq e^{-\lambda t} \mathbf{E} \exp(\lambda G_n f) \quad (5.4)$$

□

Lemma 5.2. *For any finite class \mathcal{F} of bounded, measurable, square-integrable functions, with $|\mathcal{F}|$ elements, it holds*

$$\mathbf{E} \|G_n f\|_{\mathcal{F}} \lesssim \max_{f \in \mathcal{F}} \frac{\|f\|_{\infty}}{\sqrt{n}} \log(1 + |\mathcal{F}|) + \max_{f \in \mathcal{F}} \|f\|_{\mathbf{P}, 2} \sqrt{\log(1 + |\mathcal{F}|)}. \quad (5.5)$$

6 Simple yet useful Calculations

Theorem 6.1. (Multivariate Taylor Theorem) *Let $f \in C^2(\mathbb{R}^n, \mathbb{R})$. Then for all $x, \Delta \in \mathbb{R}^n$ there exists $\xi \in [0, 1]$ such that it holds*

$$\begin{aligned} f(x + \Delta) = f(x) &+ \sum_{i=1}^n \frac{\partial f(x)}{\partial x_i} \Delta_i + \sum_{\substack{i,j=1 \\ i \neq j}}^n \frac{\partial^2 f(x + \xi \Delta)}{\partial x_i \partial x_j} \Delta_i \Delta_j \\ &+ \frac{1}{2} \sum_{i=1}^n \frac{\partial^2 f(x + \xi \Delta)}{\partial x_i^2} \Delta_i^2 \end{aligned} \quad (6.1)$$

Corollary 6.1.1. *Let $f \in C^2(\mathbb{R})$. Then for all $a, x, \Delta \in \mathbb{R}^n$ there exist $\xi \in [0, 1]$ such that it holds*

$$f(a^T(x + \Delta)) - f(a^T x) = f'(a^T x) \Delta^T a + \frac{1}{2} f''(a^T(x + \xi \Delta)) \Delta^T A \Delta, \quad (6.2)$$

where $A := aa^T \in \mathbb{R}^{n \times n}$.

Proof. By the chain rule we have for all $a, x, \Delta \in \mathbb{R}^n$ and $\xi \in [0, 1]$

$$\frac{\partial^2 f(a^T(x + \xi \Delta))}{\partial x_i \partial x_j} = f''(a^T(x + \xi \Delta)) a_i a_j. \quad (6.3)$$

Since $A := aa^T$ is symmetric we have

$$\Delta^T A \Delta = 2 \sum_{\substack{i,j=1 \\ i \neq j}}^n a_i a_j \Delta_i \Delta_j + \sum_{i=1}^n a_i^2 \Delta_i^2. \quad (6.4)$$

Plugging (6.3) and (6.4) into (6.1) yields (6.2). □

Proposition 6.1. *For all $x, y \in \mathbb{R}$ it holds*

$$|x + y| - |x| \geq -|y| \quad (6.5)$$

Proof. Checking all 6 combinations of $x + y, x, y$ being nonnegative or negative yields the result. □

Notation Index

$\#A$ cardinality of the set A

$\mathbf{E}[X|Y]$ conditional expectation of the random variable X with respect to $\sigma(Y)$

$\mathbf{E}[X]$ expectation of the random variable X

$\mathbf{Var}[X]$ variance of the random variable X

$\overline{\mathbb{R}} = \mathbb{R} \cup \{+\infty\}$ extension of the real numbers

$\xrightarrow{\mathcal{D}}$ convergence of distributions

\mathbf{P} generic probability measure

$\mathbf{P}_X = \mathbf{P} \circ X^{-1}$ distribution of the random variable X

\mathbb{R} set of real numbers

$x \vee y, x \wedge y, x^+, x^-$ maximum, minimum, positive part, negative part of real numbers

$X \sim \mu$ the random variable has distribution μ

Bibliography

- [Bha97] Rajendra Bhatia. *Matrix Analysis*, volume 169 of *Graduate Texts in Mathematics*. Springer, New York, NY, 1997.
- [BT03] Dimitri P. Bertsekas and John N. Tsitsiklis. Parallel and Distributed Computation: Numerical Methods. November 2003.
- [Bur73] D. L. Burkholder. Distribution Function Inequalities for Martingales. *The Annals of Probability*, 1(1):19–42, 1973.
- [MJC⁺14] Lester Mackey, Michael I. Jordan, Richard Y. Chen, Brendan Farrell, and Joel A. Tropp. Matrix concentration inequalities via the method of exchangeable pairs. *The Annals of Probability*, 42(3), May 2014.
- [Pet94] Dénes Petz. A survey of certain trace inequalities. *Banach Center Publications*, 30(1):287–298, 1994.
- [Roc70] R. Tyrrell Rockafellar. *Convex Analysis*. Princeton University Press, 1970.
- [Ruh70] Axel Ruhe. Perturbation bounds for means of eigenvalues and invariant subspaces. *BIT Numerical Mathematics*, 10(3):343–354, September 1970.
- [Ste72] Charles Stein. A bound for the error in the normal approximation to the distribution of a sum of dependent random variables. *Proceedings of the Sixth Berkeley Symposium on Mathematical Statistics and Probability, Volume 2: Probability Theory*, 6.2:583–603, January 1972.
- [Tro15] Joel A. Tropp. An Introduction to Matrix Concentration Inequalities, January 2015.