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

**Universität Stuttgart**

**Universität Stuttgart**

Ioan Scheffel

December 16, 2022

# Contents

# 1  Introduction

Researchers are often left with observational studies to answer questions about causality. When confounders are present the task of infering causality can become arbitrarily complex. Propensity score methods [RR83], e.g. inverse probability weighting or matching, are popular methods to adjust for confounders. Usually these methods rely heavily on estimates of the true propensity score, which are known to suffer from model dependencies and misspecification [KS07]. This issue becomes more pressing when moving from binary to continuous treatment [HI05]. Therefore methods have been developed to directly target imbalances in the data [FHI18] [Hai12] [Zub15]. We take a closer look at [WZ19] and extend the analysis to settings with continuous treatment [VGC$^+$20] [Tüb20].

# 2 Balancing Weights

## 2.1 Introduction

We work in the Rubin Causal Model.

We assume a sample of $n$ units which is drawn from a population distribution.

In i.i.d. fashion.

We observe $(\mathbf{X}_i, T_i, Y_i)$, where $\mathbf{X}$ are covariates, $T$ is the indicator if treatment has been received and $Y$ is the observed outcome.

In the Rubin Causal Model we assume that for each unit the potential outcome exist, i.e. $(Y_i^0, Y_i^1)$ where $Y^1$ stands for the potential outcome had the unit received treatment and $Y^0$ for the potential outcome had the unit received **no** treatment.

It is clear that $Y_i = Y_i^{T_i}$ i.e. we can observe only one of the potential outcomes.

Thus there is a connection to missing data problems.

This is the dilemma of causal inference.

On the population level it is possible to estimate both.

Usually the means of the potential outcomes are compared against each other.

In randomized trials this is a valid approach to causal inference.

In observational studies however the treatment assignment is not known and direct comparison can lead to systematically wrong results.

This phenomenon is called **confounding**.

To address the issue of confounding many methods have been proposed.

An intuitive way to think about potential outcomes is to think of a stochastic process $Y(\cdot)$ indexed over $\{0, 1\}$. By observing $Y_i$ we in fact sample from this process at random index $T$, i.e. from $Y(T)$. We have

$$\mathbf{E}[Y(T)] = \mathbf{E}[Y(1)|T=1]\mathbf{P}[T=1] + \mathbf{E}[Y(0)|T=0]\mathbf{P}[T=0]. \tag{2.1}$$

Suppose we observe $T = 1$. Clearly we have

$$\mathbf{E}[Y(T)|T=1] = \mathbf{E}[Y(1)|T=1] \tag{2.2}$$

## 2.2 Estimating the Population Mean of Potential Outcomes

## 2.3 Application of Matrix Concentration Inequalities

**Analysis of $\mathbf{E}[\max_{i \leq r} \|\mathbf{A}_i\|^2]$**

We have

$$\mathbf{A}_i := \frac{1}{r}\left(\frac{1-\pi_i}{\pi_i}\right)\mathbf{B}(X_i) \qquad \text{for } i \in \{1, \ldots, r\}. \tag{2.3}$$

Since we take the maximum over a finite set it is attained for some $i^* \in \{1, \dots, r\}$:

$$\mathbf{E}[\max_{i \leq r} \|\mathbf{A}_i\|^2] = \mathbf{E}[\|\mathbf{A}_{i^*}\|^2]$$

$$= \frac{1}{r^2} \mathbf{E}\left[\left(\frac{1 - \pi_{i^*}}{\pi_{i^*}}\right)^2 \|\mathbf{B}(X_{i^*})\|^2\right] \leq \frac{1}{r^2} \mathbf{E}\left[\left(\frac{1 - \pi_{i^*}}{\pi_{i^*}}\right)^4\right]^{\frac{1}{2}} \mathbf{E}[\|\mathbf{B}(X_{i^*})\|^4]^{\frac{1}{2}} \quad (2.4)$$

$$\leq \frac{K}{r^2} \sqrt{C_\pi C_\mathbf{B}}$$

In the last two steps we applied the Cauchy-Schwarz inequality and Assumption. Note that

$$\sum_{i=1}^r \mathbf{E}[\|\mathbf{A}_i\|^2] \leq \frac{K}{r} \sqrt{C_\pi C_\mathbf{B}} \tag{2.5}$$

**Assumption 2.1.** *There exists $C_\pi \geq 1$ such that $\mathbf{E}\left[\left(\frac{1-\pi_i}{\pi_i}\right)^4\right] \leq C_\pi$ for all $i \in \{1, \dots, r\}$ .*

**Remark 2.1.** *If we assume a logistic regression model for the propensity score it holds for some $\theta \in \mathbb{R}^N$ (N is the number of covariates)*

$$\frac{1 - \pi(X)}{\pi(X)} = \exp(-\theta X) \qquad and \qquad \mathbf{E}\left[\left(\frac{1 - \pi(X)}{\pi(X)}\right)^4\right] = \mathbf{E}[\exp(-4\theta X)] = M_X(-4\theta), \quad (2.6)$$

*where $M_X$ is the momement-generating function of $X$. While the first quantity in (2.6) may be unbounded when $X$ has unbounded support, the latter quantity in (2.6) is still bounded for reasonable choices of $X$.* ◇

**Assumption 2.2.** *There exists $C_\mathbf{B} \geq 1$ such that $\mathbf{E}[\mathbf{B}_k(X_i)^4] \leq C_\mathbf{B}$ for all $(k, i) \in \{1, \dots, K\} \times \{1, \dots, r\}$ .*

**Remark 2.2.** *With Assumption we also get a bound on the fourth moment of $\|\mathbf{B}(X_i)\|$. Indeed, by the convexity of $x \mapsto x^2$, the monotonicity and linearity of the expectation it holds*

$$\mathbf{E}[\|\mathbf{B}(X_i)\|^4] = \mathbf{E}\left[\left(\sum_{k=1}^K \mathbf{B}_k^2(X_i)\right)^2\right] = K^2 \mathbf{E}\left[\left(\sum_{k=1}^K \frac{1}{K} \mathbf{B}_k^2(X_i)\right)^2\right] \leq K^2 \mathbf{E}\left[\sum_{k=1}^K \frac{1}{K} \mathbf{B}_k^4(X_i)\right]$$

$$= K \sum_{k=1}^K \mathbf{E}\left[\mathbf{B}_k^4(X_i)\right] \leq K^2 C_\mathbf{B}$$

$$(2.7)$$

◇

**Analysis of $v(\mathbf{S})$**

We use the fact that $\|\mathbf{A}\|_2 \leq \|\mathbf{A}\|_F = \sqrt{\sum_{i,j} a_{ij}^2}$ It holds

$$\sum_{i=1}^r \mathbf{E}[\mathbf{A}_i \mathbf{A}_i^\top] = \frac{1}{r^2} \sum_{i=1}^r \mathbf{E}\left[\left(\frac{1 - \pi_i}{\pi_i}\right)^2 \mathbf{B}(X_i)\mathbf{B}(X_i)^\top\right] = \frac{1}{r^2}\left(\sum_{i=1}^r \mathbf{E}\left[\left(\frac{1 - \pi_i}{\pi_i}\right)^2 B_k(X_i)B_l(X_i)\right]\right)_{1 \leq k,l \leq K}.$$

$$(2.8)$$

Thus

$$\left\|\sum_{i=1}^{r}\mathbf{E}[\mathbf{A}_i\mathbf{A}_i^\top]\right\|_2^2$$

$$\leq \left\|\sum_{i=1}^{r}\mathbf{E}[\mathbf{A}_i\mathbf{A}_i^\top]\right\|_F^2 = \frac{1}{r^4}\sum_{k,l=1}^{K}\left(\sum_{i=1}^{r}\mathbf{E}\left[\left(\frac{1-\pi_i}{\pi_i}\right)^2 B_k(X_i)B_l(X_i)\right]\right)^2 \tag{2.9}$$

$$\leq \frac{1}{r^4}\sum_{k,l=1}^{K}\left(\sum_{i=1}^{r}\mathbf{E}\left[\left(\frac{1-\pi_i}{\pi_i}\right)^4\right]^{\frac{1}{2}}\mathbf{E}[B_k(X_i)^4]^{\frac{1}{4}}\mathbf{E}[B_l(X_i)^4]^{\frac{1}{4}}\right)^2 \leq \left(\frac{K}{r}\right)^2 C_\pi C_B$$

On the other hand

$$\left\|\sum_{i=1}^{r}\mathbf{E}[\mathbf{A}_i^\top\mathbf{A}_i]\right\|_2 = \sum_{i=1}^{r}\mathbf{E}[\mathbf{A}_i^\top\mathbf{A}_i] = \frac{1}{r^2}\sum_{i=1}^{r}\mathbf{E}\left[\left(\frac{1-\pi_i}{\pi_i}\right)^2 \|\mathbf{B}(X_i)\|_2^2\right]$$

$$\leq \frac{1}{r^2}\sum_{i=1}^{r}\mathbf{E}\left[\left(\frac{1-\pi_i}{\pi_i}\right)^4\right]^{\frac{1}{2}}\mathbf{E}[\|\mathbf{B}(X_i)\|_2^4]^{\frac{1}{2}} \leq \frac{K}{r}\sqrt{C_\pi C_B} \tag{2.10}$$

It follows

$$v(\mathbf{S}) \leq \frac{K}{r}\sqrt{C_\pi C_B} \tag{2.11}$$

Thus we can apply Theorem **??** to get

$$\mathbf{E}[\|\mathbf{S}\|_2] \leq \sqrt{2e\frac{K}{r}\sqrt{C_\pi C_B}\log(K+1)} + 4e\frac{\sqrt{K}}{r}\sqrt[4]{C_\pi C_B}\log(K+1) \leq 14C_\pi C_B\sqrt{\frac{K\log(K+1)}{r}} \tag{2.12}$$

# 3 Convex Analysis

## 3.1 Basic Notions

Excursively, we present some well known definitions and facts from convex analysis. For details, see, e.g., [MMN22].

A subset $C \subseteq \mathbb{R}^n$ is called **convex set**, if for all $x, y \in C$ and all $\lambda \in [0, 1]$, we have $\lambda x + (1 - \lambda)y \in C$. The Cartesian product of convex sets is convex. The intersection of a collection of convex sets is also convex. Given (not necessary convex) sets $\Omega, \Omega_1, \Omega_2 \subseteq \mathbb{R}^n$ and $\lambda \in \mathbb{R}$, define the **set addition** and **multiplication** by a real scalar as $\Omega_1 + \Omega_2 := \{x_1 + x_2 \colon x_1 \in \Omega_1, x_2 \in \Omega_2\}$ and $\lambda \Omega := \{\lambda x \colon x \in \Omega\}$. For convex sets the addition and multiplication by a real scalar are convex.

A mapping $A : \mathbb{R}^n \to \mathbb{R}^m$ is called **affine mapping** if there exist a linear mapping $L : \mathbb{R}^n \to \mathbb{R}^m$ and a vector $b \in \mathbb{R}^m$ such that $A(x) = L(x) + b$ for all $x \in \mathbb{R}^n$. The image and inverse image/preimage of convex sets under affine mappings are also convex.

## 3.2 Relative interiors

**Definition 3.1.** (Affine set and hull) *A set $A \subseteq \mathbb{R}^n$ is called **affine**, if*

$$\alpha x + (1 - \alpha)y \in A \quad \text{for all } x, y \in A \text{ and } \alpha \in \mathbb{R}. \tag{3.1}$$

*The **affine hull** $\mathrm{aff}(\Omega)$ of a set $\Omega \subseteq \mathbb{R}^n$ is the smallest affine set that includes $\Omega$.*

**Definition 3.2.** (Relative interior) *Let $\Omega \subseteq \mathbb{R}^n$. We define the **relative interior** of $\Omega$ by*

$$\mathrm{ri}(\Omega) := \left\{x \in \Omega \colon \text{there exists } \gamma > 0 \text{ such that } \mathrm{B}_\gamma(x) \cap \mathrm{aff}(\Omega) \in \Omega\right\}. \tag{3.2}$$

**Proposition 3.1.** *Let $C$ be a non-empty convex set in $\mathbb{R}^n$. Then we get the representation*

$$\mathrm{ri}(C) = \{z \in C \colon \text{for all } x \in C \text{ there exists } t > 0 \text{ such that } z + t(z - x) \in C\}. \tag{3.3}$$

**Proof.** [Roc70, Theorem 6.4] □

**Proposition 3.2.** *Let $C_1 \subseteq \mathbb{R}^{n_1}$ and $C_2 \subseteq \mathbb{R}^{n_2}$ be two non-empty convex sets. Then it holds*

$$\mathrm{ri}(C_1 \times C_2) = \mathrm{ri}(C_1) \times \mathrm{ri}(C_2). \tag{3.4}$$

**Proof.** Let $(z_1, z_2) \in \mathrm{ri}(C_1 \times C_2)$. Then for all $(x_1, x_2) \in C_1 \times C_2$ there exists $t > 0$ such that

$$z_i + t(z_i - x_i) \in C_i \qquad \text{for } i \in \{1, 2\}. \tag{3.5}$$

This proves $\subseteq$. Suppose $z_1 \in \mathrm{ri}(C_1)$ and $z_2 \in \mathrm{ri}(C_2)$. Let $(x_1, x_2) \in C_1 \times C_2$ with corresponding $t_1, t_2 > 0$. If $t_1 = t_2$ everything is clear. W.l.o.g. assume $t_1 < t_2$. Define $\theta := \frac{t_1}{t_2} \in (0, 1)$. By the convexity of $C_2$ it follows

$$z_2 + t_1(z_2 - x_2) = \theta(z_2 + t_2(z_2 - x_2)) + (1 - \theta)z_2 \in C_2. \tag{3.6}$$

Thus $(z_1, z_2) \in \mathrm{ri}(C_1 \times C_2)$. This proves $\supseteq$ and equality. □

## 3.3 Convex Separation

**Definition 3.3.** *Let $C_1$ and $C_2$ be two non-empty convex sets in $\mathbb{R}^n$. A hyperplane $H$ is said to **separate** $C_1$ and $C_2$ if $C_1$ is contained in one of the closed half-spaces associated with $H$ and $C_2$ lies in the opposite closed half-space. It is said to separate $C_1$ and $C_2$ **properly** if $C_1$ and $C_2$ are not both actually contained in $H$ itself.*

**Theorem 3.1.** *Let $C_1$ and $C_2$ be two non-empty convex sets in $\mathbb{R}^n$. There exists a hyperplane separating $C_1$ and $C_2$ properly if and only if there exists a vector $b \in \mathbb{R}^n$ such that*

$$\sup_{x \in C_2} \langle x, b \rangle \leq \inf_{x \in C_1} \langle x, b \rangle \quad and \quad \inf_{x \in C_2} \langle x, b \rangle < \sup_{x \in C_1} \langle x, b \rangle. \tag{3.7}$$

**Proof.** [Roc70, Theorem 11.1] □

**Theorem 3.2.** (Convex separation in finite dimension) *Let $C_1$ and $C_2$ be two non-empty convex sets in $\mathbb{R}^n$. Then $C_1$ and $C_2$ can be properly separated if and only if $\mathrm{ri}(C_1) \cap \mathrm{ri}(C_2) = \emptyset$.*

**Proof.** [Roc70, Theorem 11.3] □

**Definition 3.4.** (Support function intersection rule) (Support function) *Given a nonempty subset $\Omega \subseteq \mathbb{R}^n$ the **support function** $\sigma_\Omega : \mathbb{R}^n \to \overline{\mathbb{R}}$ of $\Omega$ is defined by*

$$\sigma_\Omega(x^*) := \sup_{x \in \Omega} \langle x^*, x \rangle \qquad for\ x^* \in \mathbb{R}^n. \tag{3.8}$$

**Theorem 3.3.** *Let $C_1$ and $C_2$ be two non-empty convex sets in $\mathbb{R}^n$ with $\mathrm{ri}(C_1) \cap \mathrm{ri}(C_2) \neq \emptyset$. Then the support function of the intersection $C_1 \cap C_2$ is represented as*

$$(\sigma_{C_1 \cap C_2})(x^*) = (\sigma_{C_1} \square \sigma_{C_2})(x^*) \qquad for\ all\ x^* \in \mathbb{R}^n. \tag{3.9}$$

*Furthermore, for any $x^* \in \mathrm{dom}(\sigma_{C_1 \cap C_2})$ there exist dual elements $x_1^*, x_2^* \in \mathbb{R}^n$ such that $x^* = x_1^* + x_2^*$. and*

$$(\sigma_{C_1 \cap C_2})(x^*) = \sigma_{C_1}(x_1^*) + \sigma_{C_2}(x_2^*). \tag{3.10}$$

**Proof.** *[MMN22, Theorem 4.23]* We define

$$\Theta_1 := C_1 \times [0, \infty) \quad and \quad \Theta_2 := \{(x, \lambda) \in \mathbb{R}^n : x \in C_2 \text{ and } \lambda \leq \langle x^*, x \rangle - \alpha\}. \tag{3.11}$$

Clearly $\Theta_1$ is convex by the convexity of $C_1$. Consider the affine function

$$\varphi : \mathbb{R}^{n+1} \to \mathbb{R}, \quad (x, \lambda) \mapsto \alpha - \langle x^*, x \rangle - \lambda. \tag{3.12}$$

It holds $\Theta_2 = \varphi^{-1}((-\infty, 0]) \cap (C_2 \times \mathbb{R})$. Thus, by the convexity of the sets $\varphi^{-1}((-\infty, 0])$ and $C_2$ it follows that $\Theta_2$ is convex. We want to apply convex separation to $\Theta_1$ and $\Theta_2$. To this end we show $\mathrm{ri}(\Theta_1) \cap \mathrm{ri}(\Theta_2) = \emptyset$. First note that

$$\mathrm{ri}(\Theta_1) = \mathrm{ri}(C_1) \times \mathrm{ri}([0, \infty)) \subseteq \mathrm{ri}(C_1) \times (0, \infty). \tag{3.13}$$

Indeed, if $0 \in \mathrm{ri}([0, \infty))$ then there exists $t > 0$ such that $-tx \geq 0$ for some $x > 0$. A contradition. Furthermore

$$\mathrm{ri}(\Theta_2) \subseteq \{(x, \lambda) \in \mathbb{R}^n \colon x \in \mathrm{ri}(C_2) \text{ and } \lambda < \langle x^*, x \rangle - \alpha\}. \tag{3.14}$$

To see this, assume there is $(x, \lambda) \in \mathrm{ri}(\Theta_2)$ with $\lambda = \langle x^*, x \rangle - \alpha$. Then for some $(y, \mu) \in \Theta_2$ with $\mu < \langle x^*, y \rangle - \alpha$ there exists $t > 0$ such that $(x, \lambda) + t((x, \lambda) - (y, \mu)) \in \Theta_2$. It follows

$$0 \leq (1 + t)(\langle x^*, x \rangle - \alpha - \lambda) + t(\mu - \langle x^*, y \rangle + \alpha) < 0, \tag{3.15}$$

a contradiction. The first inequality is due to $(x, \lambda) + t((x, \lambda) - (y, \mu)) \in \Theta_2$ and the second inequality due to $\mu < \langle x^*, y \rangle - \alpha$ and $\lambda = \langle x^*, x \rangle - \alpha$. But then $\mathrm{ri}(\Theta_1) \cap \mathrm{ri}(\Theta_2) = \emptyset$. Indeed, suppose that there exists $(x, \lambda) \in \mathrm{ri}(\Theta_1) \cap \mathrm{ri}(\Theta_2)$. Then it holds $\langle x^*, x \rangle - \alpha \leq 0$ and $\lambda > 0$ since $x \in \mathrm{ri}(C_1) \cap \mathrm{ri}(C_2) \subseteq C_1 \cap C_2$. On the other hand

$$0 < \lambda < \langle x^*, x \rangle - \alpha \leq 0, \tag{3.16}$$

a contradiction. $\qquad\square$

> **Takeaways** Lorem ipsum dolor sit amet, consectetuer adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris. Nam arcu libero, nonummy eget, consectetuer id, vulputate a, magna. Donec vehicula augue eu neque. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Mauris ut leo. Cras viverra metus rhoncus sem. Nulla et lectus vestibulum urna fringilla ultrices. Phasellus eu tellus sit amet tortor gravida placerat. Integer sapien est, iaculis in, pretium quis, viverra ac, nunc. Praesent eget sem vel leo ultrices bibendum. Aenean faucibus. Morbi dolor nulla, malesuada eu, pulvinar at, mollis ac, nulla. Curabitur auctor semper nulla. Donec varius orci eget risus. Duis nibh mi, congue eu, accumsan eleifend, sagittis quis, diam. Duis eget orci sit amet orci dignissim rutrum.

## 3.4 Relative Interior

**Definition 3.5.** *Let $\Omega \subseteq \mathbb{R}^n$. We define the **relative interior** of $\Omega$ by*

$$\mathrm{ri}(\Omega) := \{x \in \Omega \colon \text{there exists } \gamma > 0 \text{ such that } \mathrm{B}_\gamma(x) \cap \mathrm{aff}(\Omega) \in \Omega\}. \tag{3.17}$$

Next we collect some useful properties of relative interiors.

**Theorem 3.4.**

> **Theorem 3.5.** *Let $C$ be a non-empty convex set in $\mathbb{R}^n$. Then we get the representation*
>
> *(i)* $\mathrm{ri}(C) = \{z \in C \colon \text{for all } x \in C \text{ there exists } t > 0 \text{ such that } z + t(z - x) \in C\}.$
>
> *(ii)* $\mathrm{cl}(C)$ *and* $\mathrm{ri}(C)$ *are convex sets.*
>
> *(iii)* $\mathrm{cl}(\mathrm{ri}(C)) = \mathrm{cl}(C)$ *and* $\mathrm{ri}(\mathrm{cl}(C)) = \mathrm{ri}(C)$.
>
> *(iv)* *Suppose* $\bigcap_{i \in I} C_i \neq \emptyset$ *for a finite index set* $I$. *Then* $\mathrm{ri}\left(\bigcap_{i \in I} C_i\right) = \bigcap_{i \in I} \mathrm{ri}(C_i)$.
>
> *(v)* *Let* $L \colon \mathbb{R}^n \to \mathbb{R}^m$ *be a linear mapping. Then* $\mathrm{ri}(L(C)) = L(\mathrm{ri}(C))$. *If additionaly it holds* $L^{-1}(\mathrm{ri}(C)) \neq \emptyset$ *we have* $\mathrm{ri}(L^{-1}(C)) = L^{-1}(\mathrm{ri}(C))$.
>
> *(vi)* $\mathrm{ri}(C_1 \times C_2) = \mathrm{ri}(C_1) \times \mathrm{ri}(C_2)$.

*(vii)* $\operatorname{ri}(C_1) \cap \operatorname{ri}(C_2) = \emptyset$ *if and only if* $0 \notin \operatorname{ri}(C_1 - C_2)$.

## 3.5 Conjugate Calculus

When studying different primal problems such as (**??**) we often turn to the dual instead. Therefore we need some reliable tools. Begin able to compute specific convex conjugates is one tool required.

**Definition 3.6.** (Convex conjugate) *Given a function* $f : \mathbb{R}^n \to \overline{\mathbb{R}}$ *, the* **convex conjugate** $f^* : \mathbb{R}^n \to \overline{\mathbb{R}}$ *of* $f$ *is defined as*

$$f^*(x^*) := \sup_{x \in \mathbb{R}^n} (x^*)^T x - f(x) \tag{3.18}$$

Note that $f$ in Definition 3.6 does not have to be convex. On the other hand, the convex conjugate is always convex:

**Proposition 3.3.** *Let* $f : \mathbb{R}^n \to (-\infty, \infty]$ *be a proper function. Then its convex conjugate* $f^* : \mathbb{R}^n \to (-\infty, \infty]$ *is convex.*

**Definition 3.7.** *Given a nonempty subset* $\Omega \subseteq \mathbb{R}^n$ *the* **support function** $\sigma_\Omega : \mathbb{R}^n \to \overline{\mathbb{R}}$ *of* $\Omega$ *is defined by*

$$\sigma_\Omega(x^*) := \sup_{x \in \Omega} \langle x^*, x \rangle \qquad \text{for } x^* \in \mathbb{R}^n. \tag{3.19}$$

**Lemma 3.1.** *For any proper function* $f : \mathbb{R}^n \to \overline{\mathbb{R}}$ *we have*

$$f^*(x^*) = \sigma_{\operatorname{epi}(f)}(x^*, -1) \qquad \text{for } x^* \in \mathbb{R}^n. \tag{3.20}$$

***Proof.*** Let $x^* \in \mathbb{R}^n$ and $(x, \lambda) \in \operatorname{epi}(f)$. Then $x \in \operatorname{dom}(f)$ and $f(x) \leq \lambda$. Thus

$$\langle x^*, x \rangle - f(x) \geq \langle x^*, x \rangle - \lambda \qquad \text{for all } (x, \lambda) \in \operatorname{epi}(f). \tag{3.21}$$

On the other hand $(x, f(x)) \in \operatorname{epi}(f)$ for all $x \in \operatorname{dom}(f)$. It follows

$$\langle x^*, x \rangle - f(x) \leq \sup_{(x,\lambda) \in \operatorname{epi}(f)} \langle x^*, x \rangle - \lambda \qquad \text{for all } x \in \operatorname{dom}(f). \tag{3.22}$$

Taking the supremum in the last two displays yields

$$f^*(x^*) = \sup_{x \in \operatorname{dom}(f)} \langle x^*, x \rangle - f(x) = \sup_{(x,\lambda) \in \operatorname{epi}(f)} \langle x^*, x \rangle - \lambda \tag{3.23}$$

$$= \sup_{(x,\lambda) \in \operatorname{epi}(f)} \langle (x^*, -1), (x, \lambda) \rangle = \sigma_{\operatorname{epi}(f)}(x^*, -1). \tag{3.24}$$

$\square$

**Proposition 3.4.**

**Theorem 3.6.** (Conjugate Chain Rule) *Let* $A : \mathbb{R}^m \to \mathbb{R}^n$ *be a linear map (matrix) and* $g : \mathbb{R}^n \to (-\infty, \infty]$ *a proper convex function. If* $\operatorname{Im}(A) \cap \operatorname{ri}(\operatorname{dom}(g)) \neq \emptyset$ *it follows*

$$(g \circ A)^*(x^*) = \inf_{y^* \in (A^*)^{-1}(x^*)} g^*(y^*). \tag{3.25}$$

*Furthermore, for any* $x^* \in \operatorname{dom}(g \circ A)^*$ *there exists* $y^* \in (A^*)^{-1}(x^*)$ *such that* $(g \circ A)^*(x^*) = g^*(y^*)$.

**Definition 3.8.** (Infimal convolution) *Given functions $f_i : \mathbb{R}^n \to (-\infty, \infty]$ for $i = 1, \ldots, n$ the* ***infimal convolution*** *of these functions as defined as*

$$(f_1 \square \ldots \square f_m)(x) := \inf_{\substack{x_i \in \mathbb{R}^n \\ \sum_{i=1}^{m} x_i = x}} \sum_{i=1}^{m} f_i(x_i) \tag{3.26}$$

**Theorem 3.7.** *Let $f, g : \mathbb{R}^n \to (-\infty, \infty]$ be proper convex functions and $ri(dom(f)) \cap ri(dom(g)) \neq \emptyset$. Then we have the conjugate sum rule*

$$(f + g)^*(x^*) = (f^* \square g^*)(x^*) \tag{3.27}$$

*for all $x^* \in \mathbb{R}^n$. Moreover, the infimum in $(f^* \square g^*)(x^*)$ is attained, i.e., for any $x^* \in dom(f+g)^*$ there exists vectors $x_1^*, x_2^*$ for which*

$$(f + g)^*(x^*) = f^*(x_1^*) + g^*(x_2^*), \quad x^* = x_1^* + x_2^*. \tag{3.28}$$

***Proof.*** Let $x^* \in \mathbb{R}^n$ and fix $x_1^*, x_2^* \in \mathbb{R}^n$ such that $x^* = x_1^* + x_2^*$. We get

$$\begin{aligned} f^*(x_1^*) + g^*(x_2^*) &= \sup_{x \in \mathbb{R}^n} \langle x_1^*, x \rangle - f(x) + \sup_{x \in \mathbb{R}^n} \langle x_2^*, x \rangle - g(x) \\ &\geq \sup_{x \in \mathbb{R}^n} \langle x_1^*, x \rangle - f(x) + \langle x_2^*, x \rangle - g(x) = \sup_{x \in \mathbb{R}^n} \langle x_1^* + x_2^*, x \rangle - (f(x) + g(x)) \\ &= \sup_{x \in \mathbb{R}^n} \langle x^*, x \rangle - (f + g)(x) = (f + g)^*(x^*) \end{aligned}$$

Taking the infimum over $x_1^*, x_2^* \in \mathbb{R}^n$ in the above display gives $(f^* \square g^*)(x^*) \geq (f + g)^*(x^*)$. Let us prove now $\leq$ under the condition $\mathrm{ri}\,(\mathrm{dom}(f)) \cap \mathrm{ri}\,(\mathrm{dom}(g)) \neq \emptyset$. The only case we need to consider is $(f + g)^*(x^*) < \infty$. Define two convex sets by

$$\Omega_1 := \left\{ (x, \alpha, \beta) \in \mathbb{R}^{n+2} \colon \alpha \geq f(x) \right\} = \mathrm{epi}(f) \times \mathbb{R}, \tag{3.29}$$
$$\Omega_2 := \left\{ (x, \alpha, \beta) \in \mathbb{R}^{n+2} \colon \beta \geq g(x) \right\}. \tag{3.30}$$

Similar to Lemma we get the representation

$$(f + g)^*(x^*) = \sigma_{\Omega_1 \cap \Omega_2}(x^*, -1, -1). \tag{3.31}$$

Indeed, the only thing we need to verify is $\mathrm{dom}(f) \cap \mathrm{dom}(g) = \mathrm{dom}(f + g)$. The inclusion $\subseteq$ is clear. Assume towards a contradiction that $(f + g)(x) < \infty$ and $f(x) = \infty$. Since $g(x) > -\infty$ it holds

$$\infty = \infty + g(x) = f(x) + g(x) = (f + g)(x) < \infty. \tag{3.32}$$

This is a contradiction. The same holds for $f$ and $g$ reversed. It follows the inclusion $\supseteq$ and equality. By the support function intersection rule there exist triples

$$(x_1^*, -\alpha_1, -\beta_1), (x_2^*, -\alpha_2, -\beta_2) \in \mathbb{R}^{n+2} \quad \text{such that} \quad (x^*, -1, -1) = (x_1^* + x_2^*, -(\alpha_1 + \alpha_2), -(\beta_1 + \beta_2)) \tag{3.33}$$

and

$$(f + g)^*(x^*) = \sigma_{\Omega_1 \cap \Omega_2}(x^*, -1, -1) = \sigma_{\Omega_1}(x_1^*, -\alpha_1, -\beta_1) + \sigma_{\Omega_2}(x_2^*, -\alpha_2, -\beta_2). \tag{3.34}$$

Next we show $\beta_1 = \alpha_2 = 0$. Suppose towards a contradiction that $\beta_1 \neq 0$. We fix $(\overline{x}, \overline{\alpha}) \in \mathrm{epi}(f)$. Then

$$\sigma_{\Omega_1}(x_1^*, -\alpha_1, -\beta_1) = \sup_{(x, \alpha, \beta) \in \mathrm{epi}(f) \times \mathbb{R}} \langle x^*, x \rangle - \alpha \alpha_1 - \beta \beta_1 \geq \sup_{\beta \in \mathbb{R}} \langle x^*, \overline{x} \rangle - \overline{\alpha} \alpha_1 - \beta \beta_1 = \infty. \tag{3.35}$$

This contradicts $(f+g)^*(x^*) < \infty$. In a similar fashion we can derive a contradiction for $\alpha_2 \neq 0$. Employing Lemma and taking into account the structures of the sets $\Omega_1$ and $\Omega_2$ this implies

$$(f+g)^*(x^*) = \sigma_{\Omega_1 \cap \Omega_2}(x^*, -1, -1) = \sigma_{\Omega_1}(x_1^*, -1, 0) + \sigma_{\Omega_2}(x_2^*, 0, -1) \tag{3.36}$$

$$= \sigma_{\mathrm{epi}(f)}(x_1^*, -1) + \sigma_{\mathrm{epi}(g)}(x_2^*, -1) = f^*(x_1^*) + g^*(x_2^*) \geq (f^* \square g^*)(x^*). \tag{3.37}$$

This finishes the proof. $\qquad\square$

# 4 Random Matrix Inequalities

## 4.1 Matrix Analysis

The **trace** of a square matrix, denoted by tr, is the sum of its diagonal entries, i.e. $\mathrm{tr}(\mathbf{B}) = \sum_{j=1}^{d} b_{jj}$ for $\mathbf{B} \in \mathbb{M}_d$. The trace is unitarily invariant, i.e. $\mathrm{tr}(\mathbf{B}) = \mathrm{tr}(\mathbf{Q}\mathbf{B}\mathbf{Q}^*)$ for all $\mathbf{B} \in \mathbb{M}_d$ for all unitary $\mathbf{Q} \in \mathbb{M}_d$. In particular, the existence of an eigenvalue value decomposition shows that the trace of a Hermitian matrix equals the sum of its eigenvalues. Let $f : I \to \mathbb{R}$ where $I \subseteq \mathbb{R}$ is an interval. Consider a matrix $\mathbf{A} \in \mathbb{H}_d$ whose eigenvalues are contained in $I$. We define the matrix $f(\mathbf{A}) \in \mathbb{H}_d$ using an eigenvalue decomposition of $\mathbf{A}$ :

$$
f(\mathbf{A}) = \mathbf{Q} \begin{bmatrix} f(\lambda_1) & & \\ & \ddots & \\ & & f(\lambda_d) \end{bmatrix} \mathbf{Q}^* \quad \text{where} \quad \mathbf{A} = \mathbf{Q} \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_d \end{bmatrix} \mathbf{Q}^*. \quad (4.1)
$$

The definition of $f(\mathbf{A})$ does not depend on which eigenvalue decomposition we choose. Any matrix function that arises in this fashion is called a **standard matrix function**.

**Proposition 4.1.** *Let $f, g : I \to \mathbb{R}$ be real-valued functions on an interval $I \subseteq \mathbb{R}$, and let $\mathbf{A} \in \mathbb{H}_d$ be a Hermitian matrix whose eigenvalues are contained in $I$.*

*(i) If $\lambda$ is an eigenvalue of of $\mathbf{A}$, then $f(\lambda)$ is an eigenvalue of $f(\mathbf{A})$.*

*(ii) $f(a) \le g(a)$ for all $a \in I$ implies $f(\mathbf{A}) \preccurlyeq g(\mathbf{A})$.*

## 4.2 Matrix Khintchin Inequality

**Theorem 4.1.** [MJC$^+$14, Corollary 7.3] *Suppose that $p = 1$ or $p \ge 3/2$. Consider a finite sequence $(\mathbf{Y}_k)_{k \ge 1}$ of independent, random, Hermitian matrices and a deterministic sequence $(\mathbf{A}_k)_{k \ge 1}$ for which*

$$
\mathbf{E}[\mathbf{Y}_k] = 0 \quad \text{and} \quad \mathbf{Y}_k^2 \preccurlyeq \mathbf{A}_k^2 \quad \text{almost surely for all } k \ge 1. \quad (4.2)
$$

*Then*

$$
\mathbf{E}\left[ \left\| \sum_{k \ge 1} \mathbf{Y}_k \right\|_{2p}^{2p} \right]^{1/(2p)} \le \sqrt{p - \frac{1}{2}} \left\| \left( \sum_{k \ge 1} (\mathbf{A}_k^2 + \mathbf{E}[\mathbf{Y}_k^2]) \right)^{1/2} \right\|_{2p}. \quad (4.3)
$$

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

$$\mathbf{E}\left[\left\|\sum_{k\geq 1}\xi_k\mathbf{A}_k\right\|_{2p}^{2p}\right]^{1/(2p)} \leq \sqrt{2p-1}\left\|\left(\sum_{k\geq 1}\mathbf{A}_k^2\right)^{1/2}\right\|_{2p}. \tag{4.4}$$

## 4.3 Matrix Moment Inequality

**Theorem 4.2.** *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*

$$\mathbf{E}\left[\left\|\sum_{k\geq 1}\mathbf{S}_k\right\|^p\right]^{1/p} \leq \left[\left\|\sum_{k\geq 1}\mathbf{E}[\mathbf{S}_k]\right\|^{1/2} + 2\sqrt{er}\mathbf{E}[\max_{k\geq 1}\|\mathbf{S}_k\|^p]^{1/(2p)}\right]^2. \tag{4.5}$$

*(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*

$$\mathbf{E}\left[\left\|\sum_{k\geq 1}\mathbf{Y}_k\right\|^p\right]^{1/p} \leq \sqrt{er}\left\|\left(\sum_{k\geq 1}\mathbf{E}[\mathbf{Y}_k^2]\right)^{1/2}\right\| + 2er\mathbf{E}[\max_{k\geq 1}\|\mathbf{S}_k\|^p]^{1/p}. \tag{4.6}$$

## 4.4 Intrinsic Dimension

**Definition 4.1.** *For a positive-semidefinite matrix $\mathbf{S}$, the **intrinic dimension** is the quantity*

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

**Lemma 4.1.** *(Intrinsic dimenision) Let $\varphi : [0, \infty) \to \mathbb{R}$ 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 $\varphi$ is convex on any interval $[0, L]$ with $L > 0$ and $\varphi(0) = 0$, it holds

$$\varphi(a) \leq \left(1 - \frac{a}{L}\right)\varphi(0) + \frac{a}{L}\varphi(L) = \frac{a}{L}\varphi(L) \qquad \text{for all } a \in [0, L]. \tag{4.7}$$

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

$$\mathrm{tr}(\varphi(\mathbf{S})) = \sum_{i=1}^{d}\varphi(\lambda_i) \leq \frac{\sum_{i=1}^{d}\lambda_i}{\|\mathbf{S}\|}\varphi(\|\mathbf{S}\|) = \frac{\mathrm{tr}(\mathbf{S})}{\|\mathbf{S}\|}\varphi(\|\mathbf{S}\|) = \mathrm{intdim}(\mathbf{S}) \cdot \varphi(\|\mathbf{S}\|). \tag{4.8}$$

$\square$

# 5 Empirical Processes

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

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

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

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

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

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

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

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

$\square$

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

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

# 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*

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

**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*

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

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

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

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

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

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

Plugging (6.3) and (6.4) into (6.1) yields (6.2). $\qquad\square$

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

$$
|x + y| - |x| \geq -|y|
\tag{6.5}
$$

**Proof.** Checking all 6 combinations of $x + y, x, y$ being nonnegative or negative yields the result. $\qquad\square$

# 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 $X$

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

$\overset{\mathcal{D}}{\longrightarrow}$    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 $\mu$

# Bibliography

[FHI18]    Christian Fong, Chad Hazlett, and Kosuke Imai. Covariate balancing propensity score for a continuous treatment: Application to the efficacy of political advertisements. *The Annals of Applied Statistics*, 12(1):156–177, March 2018.

[Hai12]    Jens Hainmueller. Entropy Balancing for Causal Effects: A Multivariate Reweighting Method to Produce Balanced Samples in Observational Studies. *Political Analysis*, 20(1):25–46, 2012.

[HI05]     Keisuke Hirano and Guido W. Imbens. The Propensity Score with Continuous Treatments. In Andrew Gelman and Xiao-Li Meng, editors, *Wiley Series in Probability and Statistics*, pages 73–84. John Wiley & Sons, Ltd, Chichester, UK, July 2005.

[KS07]     Joseph D. Y. Kang and Joseph L. Schafer. Demystifying Double Robustness: A Comparison of Alternative Strategies for Estimating a Population Mean from Incomplete Data. *Statistical Science*, 22(4):523–539, November 2007.

[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.

[MMN22]    Boris S. Mordukhovich and Nguyen Mau Nam. ENHANCED CALCULUS AND FENCHEL DUALITY. In Boris S. Mordukhovich and Nguyen Mau Nam, editors, *Convex Analysis and Beyond: Volume I: Basic Theory*, Springer Series in Operations Research and Financial Engineering, pages 255–310. Springer International Publishing, Cham, 2022.

[Roc70]    R. Tyrrell Rockafellar. *Convex Analysis*. Princeton University Press, 1970.

[RR83]     Paul R. Rosenbaum and Donald B. Rubin. The Central Role of the Propensity Score in Observational Studies for Causal Effects. *Biometrika*, 70(1):41–55, 1983.

[Tro15]    Joel A. Tropp. An Introduction to Matrix Concentration Inequalities, January 2015.

[Tüb20]    Stefan Tübbicke. Entropy Balancing for Continuous Treatments, May 2020.

[VGC⁺20]   Brian G. Vegetabile, Beth Ann Griffin, Donna L. Coffman, Matthew Cefalu, and Daniel F. McCaffrey. Nonparametric Estimation of Population Average Dose-Response Curves using Entropy Balancing Weights for Continuous Exposures, March 2020.

[WZ19]     Yixin Wang and José R. Zubizarreta. Minimal Dispersion Approximately Balancing Weights: Asymptotic Properties and Practical Considerations. *Biometrika*, page asz050, October 2019.

[Zub15]    José R. Zubizarreta. Stable Weights that Balance Covariates for Estimation With Incomplete Outcome Data. *Journal of the American Statistical Association*, 110(511):910–922, July 2015.