# A Unified Theory for Causal Inference: Direct Debiased Machine Learning via Bregman-Riesz Regression

Masahiro Kato[*]

The University of Tokyo

October 31, 2025

**Abstract**

This note introduces a unified theory for causal inference that integrates Riesz regression, covariate balancing, density-ratio estimation (DRE), targeted maximum likelihood estimation (TMLE), and the matching estimator in average treatment effect (ATE) estimation. In ATE estimation, the balancing weights and the regression functions of the outcome play important roles, where the balancing weights are referred to as the Riesz representer, bias-correction term, and clever covariates, depending on the context. Riesz regression, covariate balancing, DRE, and the matching estimator are methods for estimating the balancing weights, where Riesz regression is essentially equivalent to DRE in the ATE context, the matching estimator is a special case of DRE, and DRE is in a dual relationship with covariate balancing. TMLE is a method for constructing regression function estimators such that the leading bias term becomes zero. Nearest Neighbor Matching is equivalent to Least Squares Density Ratio Estimation and Riesz Regression.

## 1 Introduction

This note is written to convey and summarize the main ideas of Kato (2025a,b,c). These works propose the direct debiased machine learning (DDML) framework, which unifies existing treatment effect estimation methods such as Riesz regression, covariate balancing, density-ratio estimation (DRE), targeted maximum likelihood estimation (TMLE), and the matching estimator. For simplicity, we consider the standard setting of average treatment effect (ATE) estimation (Imbens & Rubin, 2015). Note that the arguments in this note can also be applied to other settings, such as estimation of the ATE on the treated (ATT). For details, see Kato (2025b). Throughout this study, we explain how the existing methods mentioned above can be unified from the viewpoint of targeted Neyman estimation via generalized Riesz regression, also called Bregman-Riesz regression.

---

[*]Email: `mkato-csecon@g.ecc.u-tokyo.ac.jp`

Specifically, these existing methods aim to estimate the nuisance parameters that minimize the estimation error between the oracle Neyman orthogonal score and an estimated Neyman orthogonal score. From this point of view, we can interpret Riesz regression, DRE, and the matching estimator as methods for estimating the Riesz representer, also called the bias-correction term or the clever covariates. Covariate balancing is in a dual relationship with these methods. TMLE is a method for regression function estimation to minimize the estimation error.

This generalization not only provides an integrated view of various methods proposed in different fields but also offers a practical guideline for choosing an ATE estimation algorithm. For example, for specific choices of basis functions and loss functions, we can automatically attain the covariate balancing property.

## 2 Setup

Let $(X, D, Y)$ be a triple of $k$-dimensional covariates $X \in \mathcal{X} \subseteq \mathbb{R}^k$, treatment indicator $D \in \{1, 0\}$, and outcome $Y \in \mathcal{Y} \subseteq \mathbb{R}$, where $\mathcal{X}$ and $\mathcal{Y}$ are the corresponding spaces, and $D = 1$ denotes treated while $D = 0$ denotes control. Following the Neyman-Rubin framework, let $Y(1) \in \mathcal{Y}$ and $Y(0) \in \mathcal{Y}$ be the potential outcomes for treated and control units. Let us define the ATE as

$$\tau_0 := \mathbb{E}\big[Y(1) - Y(0)\big].$$

We observe $n$ units with $\{(X_i, D_i, Y_i)\}_{i=1}^n$, where $(X_i, D_i, Y_i)$ is an i.i.d. copy of the predefined triple $(X, D, Y)$. Our goal is to estimate $\tau_0$ using the observations.

**Notations and assumptions** For simplicity, we assume that the covariate distributions of the treated and control groups have probability densities. We denote the probability density of covariates in the treated group by $p(x \mid D = 1)$, and that of the control group by $p(x \mid D = 0)$. We also denote the marginal probability density by $p(x)$ and the joint probability density of $(X, D)$ by $p(x, d)$. We introduce the propensity score, the probability of receiving treatment, by

$$e_0(X) := \frac{p(x, 1)}{p(x)}.$$

For $d \in \{1, 0\}$, we denote the expected outcome of $Y(d)$ conditional on $X$ by $\mu_0(d, X) = \mathbb{E}\big[Y(d) \mid X\big]$.

To identify the ATE, we assume unconfoundedness, positivity, and boundedness of the random variables; that is, $Y(1)$ and $Y(0)$ are independent of $D$ given $X$, there exists a universal constant $\epsilon \in (0, 1/2)$ such that $\epsilon < e_0(X) < 1 - \epsilon$, and $X$, $Y(1)$, and $Y(0)$ are bounded.

## 3 Riesz Representer and ATE Estimators

In ATE estimation, the following quantity plays an important role:

$$\alpha_0(D, X) := \frac{D}{e_0(X)} - \frac{1 - D}{1 - e_0(X)}.$$

This term is referred to by various names in different methods. In the classical semiparametric inference literature, it is called the bias-correction term (Schuler & van der Laan, 2024). In TMLE, it is called the clever covariates (van der Laan, 2006). In the debiased machine learning (DML) literature, it is called the Riesz representer (Chernozhukov et al., 2022). It may also be referred to as balancing weights (Imai & Ratkovic, 2013; Hainmueller, 2012), inverse propensity score (Horvitz & Thompson, 1952), or density ratio (Sugiyama et al., 2012).

This term has several uses. First, if we know the function $\alpha_0$, we can construct an inverse probability weighting (IPW) estimator as

$$\widehat{\tau}^{\mathrm{IPW}} := \frac{1}{n} \sum_{i=1}^{n} \alpha_0(D_i, X_i) Y_i.$$

This is known as one of the simplest unbiased estimators for the ATE $\tau_0$. Another usage is bias correction. Given an estimate $\widehat{\mu}(d, X)$ of $\mu_0(d, X)$, we can construct a naive plug-in estimator as

$$\widehat{\tau}^{\mathrm{PI}} := \frac{1}{n} \sum_{i=1}^{n} \big(\widehat{\mu}(1, X_i) - \widehat{\mu}(0, X_i)\big).$$

Such a naive estimator often includes bias caused by the estimation of $\mu_0$ that does not vanish at the $\sqrt{n}$ rate. Therefore, to obtain an estimator of $\tau_0$ with $\sqrt{n}$ convergence, we debias the estimator as

$$\widehat{\tau}^{\mathrm{OS}} := \frac{1}{n} \sum_{i=1}^{n} \Big( \alpha_0(D_i, X_i)\big(Y_i - \widehat{\mu}(D_i, X_i)\big) + \widehat{\mu}(1, X_i) - \widehat{\mu}(0, X_i) \Big).$$

This estimator is called the one-step estimator. There exists another direction of bias correction, called TMLE. In TMLE, we update the initial regression function estimates $\widehat{\mu}$ as

$$\widetilde{\mu}(d, X_i) = \widehat{\mu}(d, X_i) + \frac{\sum_{i=1}^{n} \alpha_0(D_i, X_i)\big(Y_i - \widehat{\mu}(D_i, X_i)\big)}{\sum_{i=1}^{n} \alpha_0(D_i, X_i)^2} \alpha_0(d, X_i).$$

Then, we redefine the ATE estimator as

$$\widehat{\tau}^{\mathrm{TMLE}} := \frac{1}{n} \sum_{i=1}^{n} \big(\widetilde{\mu}(1, X_i) - \widetilde{\mu}(0, X_i)\big).$$

Thus, the term $\alpha_0$ plays an important role. When $\alpha_0$ is unknown, its estimation becomes a core task in causal inference, along with the usually unknown regression function $\mu_0$. We can view Riesz regression, DRE, covariate balancing, and the matching estimator as methods for estimating $\alpha_0$ with different loss functions. In addition, TMLE has a close relationship with these estimation methods from the targeted Neyman estimation perspective, explained below.

# 4 Targeted Neyman Estimation

Following the debiased machine learning literature, we refer to $\alpha_0$ as the Riesz representer. We also focus on the Neyman orthogonal score, defined as

$$\psi(X, D, Y; \mu, \alpha, \tau) := \alpha(D, X)\big(Y - \mu(D, X)\big) + \mu(1, X) - \mu(0, X) - \tau.$$

From efficiency theory, we know that an estimator $\widehat{\tau}^{\mathrm{oracle}}$ is efficient if it satisfies

$$\frac{1}{n} \sum_{i=1}^{n} \psi(X, D, Y; \mu_0, \alpha_0, \widehat{\tau}^{\mathrm{oracle}}) = 0.$$

However, if we plug in estimators of $\mu_0$ and $\alpha_0$, we might incur bias caused by estimation. Even though the Neyman orthogonal score has the property that such bias can be asymptotically removed at a fast rate, it is desirable to construct estimators $\widehat{\mu}$ and $\widehat{\alpha}$ that behave well in finite samples.

Based on this motivation, Kato (2025b) proposes the targeted Neyman estimation procedure, which aims to estimate $\mu_0$, $\alpha_0$, and $\tau_0$ such that the following Neyman error becomes zero:

$$L(\mu, \alpha, \tau) := \frac{1}{n} \sum_{i=1}^{n} \psi(X, D, Y; \mu, \alpha, \tau).$$

This Neyman error can be decomposed as follows:

$$
\begin{aligned}
L(\mu, \alpha, \tau) &= \frac{1}{n} \sum_{i=1}^{n} \psi(X_i, D_i, Y_i; \mu, \alpha, \tau) \\
&= \frac{1}{n} \sum_{i=1}^{n} \Big( \big(\alpha_0(D_i, X_i) - \alpha(D_i, X_i)\big)\big(Y_i - \mu_0(D_i, X_i)\big) - \alpha(D_i, X_i)\big(Y_i - \mu(D_i, X_i)\big) \\
&\quad - \big(\mu(1, X_i) - \mu(0, X_i)\big) + \tau \Big).
\end{aligned}
$$

Therefore, in expectation, we have

$$
\begin{aligned}
&\mathbb{E}\left[L(\mu, \alpha, \tau)\right] \\
&= \mathbb{E}\left[\frac{1}{n} \sum_{i=1}^{n} \Big(\alpha_0(D_i, X_i) - \alpha(D_i, X_i)\Big)\Big(Y_i - \mu_0(D_i, X_i)\Big)\right] \\
&\quad + \mathbb{E}\left[\frac{1}{n} \sum_{i=1}^{n} \Big(\tau - \big(\mu(1, X_i) - \mu(0, X_i)\big)\Big)\right].
\end{aligned}
$$

Thus, the core error terms are the following two:

$$\frac{1}{n} \sum_{i=1}^{n} \Big(\alpha_0(D_i, X_i) - \alpha(D_i, X_i)\Big)\Big(Y_i - \mu_0(D_i, X_i)\Big), \tag{1}$$

$$\frac{1}{n} \sum_{i=1}^{n} \Big(\tau - \big(\mu(1, X_i) - \mu(0, X_i)\big)\Big). \tag{2}$$

We can interpret Riesz regression, covariate balancing, and nearest neighbor matching as methods for minimizing the error in (1) by estimating $\alpha_0$ well, while TMLE is a method that automatically sets (2) to zero by estimating $\mu_0$ well. We further point out that these existing methods for estimating the Riesz representer $\alpha_0$ can be generalized using the Bregman divergence and the duality of loss functions.

# 5  Bregman-Riesz Regression

This section reviews Bregman-Riesz regression, proposed in Kato (2025a) and Kato (2025b), which is also called generalized Riesz regression. Bregman-divergence regression generalizes Riesz regression in Chernozhukov et al. (2024) from the viewpoint of DRE via Bregman divergence minimization. As pointed out in Kato (2025b), we can derive covariate balancing methods as the dual of the Bregman divergence loss by extending the results in Bruns-Smith et al. (2025) and Zhao (2019). Note that this duality depends on the choice of models: when using Riesz regression, we need to use linear models for $\alpha_0$; when using Kullback-Leibler (KL) divergence, we need to use logistic models for $\alpha_0$.

## 5.1  Bregman Divergence

Our goal is to estimate $\alpha_0$ so that we can minimize

$$\frac{1}{n}\sum_{i=1}^{n}\Big(\alpha_0(D_i, X_i) - \alpha(D_i, X_i)\Big)\Big(Y_i - \mu_0(D_i, X_i)\Big).$$

For simplicity, let us ignore the term $\Big(Y_i - \mu_0(D_i, X_i)\Big)$. Then, our goal is merely to minimize the discrepancy between $\alpha_0(D_i, X_i)$ and $\alpha(D_i, X_i)$.

We first recap the Bregman divergence. Bregman divergence is defined via a differentiable and strictly convex function $g\colon \mathbb{R} \to \mathbb{R}$. Given $d \in \{1, 0\}$ and $x \in \mathcal{X}$, let us define the following pointwise Bregman divergence between $\alpha_0(d, x)$ and $\alpha(d, x)$:

$$\mathrm{BR}_g\big(\alpha_0(d, x) \mid \alpha(d, x)\big) := g(\alpha_0(d, x)) - g(\alpha(d, x)) - \partial g(\alpha(d, x))\big(\alpha_0(d, x) - \alpha(d, x)\big),$$

where $\partial g$ denotes the derivative of $g$. Taking the average over the distribution of $X$, we define the following average Bregman divergence:

$$\mathrm{BR}_g\big(\alpha_0 \mid \alpha\big) := \mathbb{E}\Big[g(\alpha_0(D, X)) - g(\alpha(D, X)) - \partial g(\alpha(D, X))\big(\alpha_0(D, X) - \alpha(D, X)\big)\Big].$$

Ideally, we want to estimate $\alpha_0$ by minimizing this average Bregman divergence, which is represented as

$$\alpha^* = \underset{\alpha \in \mathcal{A}}{\arg\min}\, \mathrm{BR}_g^{\dagger}\big(\alpha_0 \mid \alpha\big),$$

where $\mathcal{A}$ is a hypothesis class of $\alpha_0$. If $\alpha_0 \in \mathcal{A}$, then $\alpha^* = \alpha_0$ holds.

However, this formulation is infeasible because it includes the unknown $\alpha_0$. Surprisingly, by a simple computation, we can drop the unknown $\alpha_0$ and define an equivalent optimization problem as

$$\alpha^* = \underset{\alpha \in \mathcal{A}}{\arg\min}\, \mathrm{B}_g\big(\alpha\big),$$

where

$$\mathrm{B}_g(\alpha) := \mathbb{E}\Big[ -g(\alpha(D,X)) + \partial g(\alpha(D,X))\alpha(D,X) - \Big(\partial g(\alpha(1,X)) - \partial g(\alpha(0,X))\Big)\Big].$$

Finally, by replacing the expectation with sample approximations, we obtain the following feasible optimization problem for estimating the Riesz representer $\alpha_0$:

$$\widehat{\alpha} := \arg\min_{\alpha \in \mathcal{A}} \widehat{\mathrm{B}}_g(\alpha) + \lambda J(\alpha),$$

where $J(\alpha)$ is some regularization function, and

$$\widehat{\mathrm{B}}_g(\alpha) := \frac{1}{n}\sum_{i=1}^{n}\Big( -g(\alpha(D_i,X_i)) + \partial g(\alpha(D_i,X_i))\alpha(D_i,X_i) - \Big(\partial g(\alpha(1,X_i)) - \partial g(\alpha(0,X_i))\Big)\Big).$$

## 5.2 Squared Loss

We consider the following convex function:

$$g^{\mathrm{LS}}(\alpha) = (\alpha - 1)^2.$$

Under this choice of $g$, the estimation problem is written as

$$\widehat{\alpha} := \arg\min_{\alpha \in \mathcal{A}} \widehat{\mathrm{BR}}_{g^{\mathrm{LS}}}(\alpha) + \lambda J(\alpha), \tag{3}$$

where

$$\widehat{\mathrm{BR}}_{g^{\mathrm{LS}}}(\alpha) := \frac{1}{n}\sum_{i=1}^{n}\Big( -2\big(\alpha(1,X_i) + \alpha(0,X_i)\big) + \mathbb{1}[D_i = 1]\alpha(1,X_i)^2 + \mathbb{1}[D_i = 0]\alpha(0,X_i)^2 \Big).$$

This estimation method corresponds to Riesz regression in debiased machine learning (Chernozhukov et al., 2024) and least-squares importance fitting (LSIF) in DRE (Kanamori et al., 2009). Moreover, if we define $\mathcal{A}$ appropriately, we can yield nearest neighbor matching, as pointed out in Kato (2025c), which extends the argument in Lin et al. (2023).

**Stable balancing weights**  We can use various models for $\mathcal{A}$. For example, we can use neural networks, though it is known to cause serious overfitting problems for this kind of DRE objective (Rhodes et al., 2020; Kato & Teshima, 2021).

This study focuses on linear-in-parameter models for squared loss, defined as

$$\alpha(D,X) = \beta^{\top}\Phi(D,X),$$

where $\Phi\colon \{1,0\} \times \mathcal{X} \to \mathbb{R}^p$ is some basis function that maps $(D,X)$ to a $p$-dimensional feature space, and $\beta$ is a $p$-dimensional parameter. For such a choice of basis function, the dual of the problem (3) can be written as

$$\min_{\alpha \in \mathbb{R}^n} \|\alpha\|_2^2$$

$$\text{s.t.} \quad \sum_{i=1}^{n} \alpha_i \Phi(D_i, X_i) - \left( \sum_{i=1}^{n} \Big( \Phi(1, X_i) - \Phi(0, X_i) \Big) \right) = \mathbf{0}_p,$$

where $\mathbf{0}_p$ is the $p$-dimensional zero vector. Here, for simplicity, we let $\lambda = 0$ in this argument.

This optimization problem is the same as the one in stable balancing weights (Zubizarreta, 2015). This result is shown in Bruns-Smith et al. (2025), and Kato (2025b) calls it automatic covariate balancing since we can attain the covariate balancing property without explicitly solving the covariate balancing problem.

## 5.3 KL Divergence Loss

Next, we consider the following KL-divergence-motivated convex function:

$$g^{\mathrm{KL}}(\alpha) = (|\alpha| - 1) \log (|\alpha| - 1) - |\alpha|.$$

Then, we estimate $\alpha_0$ by minimizing the empirical objective:

$$\widehat{\alpha} := \arg\min_{\alpha \in \mathcal{A}} \widehat{\mathrm{BR}}_{g^{\mathrm{E}}}(\alpha) + \lambda J(\alpha),$$

where

$$\widehat{\mathrm{BR}}_{g^{\mathrm{E}}}(\alpha) := \frac{1}{n} \sum_{i=1}^{n} \Big( \log (|\alpha(D_i, X_i)| - 1) + |\alpha(D_i, X_i)| - \log (\alpha(1, X_i) - 1) - \log (-\alpha(0, X_i) - 1) \Big).$$

For the derivation of this loss, see Kato (2025b). If we use $g(\alpha) = |\alpha| \log |\alpha| - |\alpha|$ instead of $g^{\mathrm{KL}}$, the optimization problem aligns with LSIF in DRE (Sugiyama et al., 2007). On the other hand, if we use $g^{\mathrm{KL}}$, the optimization problem aligns with the tailored loss in covariate balancing (Zhao, 2019). Under this choice, we obtain the following duality result for entropy balancing weights (Hainmueller, 2012).

**Entropy balancing weights**   This study focuses on logistic models for KL-divergence loss, defined as

$$\alpha(D, X) = \mathbb{1}[D = 1] r(1, Z) - \mathbb{1}[D = 0] r(0, Z),$$

where $r(1, Z) = \frac{1}{e(X)}$, $r(0, Z) = \frac{1}{1 - e(X)}$, and

$$e(X) := \frac{1}{1 + \exp\left( - \beta^\top \Phi(Z) \right)}.$$

Here, $\Phi \colon \mathcal{X} \to \mathbb{R}^p$ is a basis function that does not include $D$, unlike the basis function for squared loss. Under this choice, we can write the optimization problem as

$$\widehat{r} := \arg\min_{r \in \mathcal{R}} \frac{1}{n} \sum_{i=1}^{n} \Bigg( \mathbb{1}[D_i = 1] \left( - \log \left( \frac{1}{r(1, X_i) - 1} \right) + r(1, X_i) \right)$$

$$+ \mathbb{1}[D_i = 0] \left( - \log \left( \frac{1}{r(0, X_i) - 1} \right) + r(0, X_i) \right) \Bigg),$$

where $\mathcal{R}$ is the set of functions $r$ defined above. This objective function is called the tailored loss in Zhao (2019).

Then, as shown in Zhao (2019), from the duality, it is known that this problem is equivalent to solving

$$\min_{w \in (1,\infty)^n} \sum_{i=1}^{n} (w_i - 1) \log(w_i - 1)$$

$$\text{s.t.} \quad \sum_{i=1}^{n} \left( \mathbb{1}[D_i = 1] w_i \Phi(1, X_i) - \mathbb{1}[D_i = 0] w_i \Phi(0, X_i) \right) = \mathbf{0}_p.$$

This optimization problem aligns with that in entropy balancing (Hainmueller, 2012). Here, note that for estimated $\widehat{w}_i$, we can write $\widehat{\alpha}(D_i, X_i) = \mathbb{1}[D_i = 1]\widehat{w}_i - \mathbb{1}[D_i = 0]\widehat{w}_i$

# 6 Implementation Suggestion

In practice, one of our recommendations is the following procedure:

- Estimate the regression function $\mu_0$ in some way.

- Model the Riesz representer using the logistic model $e(X) = \frac{1}{1 + \exp\left(-\beta^\top \Phi(Z)\right)}$.

- Estimate $r_0(D, X)$ as

$$\widehat{r} := \arg\min_{r \in \mathcal{R}} \frac{1}{n} \sum_{i=1}^{n} \left( \mathbb{1}[D_i = 1] \left( -\log\left( \frac{1}{r(1, X_i) - 1} \right) + r(1, X_i) \right) \right.$$

$$\left. + \mathbb{1}[D_i = 0] \left( -\log\left( \frac{1}{r(0, X_i) - 1} \right) + r(0, X_i) \right) \right) (Y_i - \widehat{\mu}(D_i, X_i))^2,$$

where $r(1, Z) = \frac{1}{e(X)}$, $r(0, Z) = \frac{1}{1 - e(X)}$. Here, we used weights $(Y_i - \widehat{\mu}(D_i, X_i))^2$, motivated by targeted Neyman estimation.

- Apply TMLE to $\widehat{\mu}$ and update it to $\widetilde{\mu}$, as in Section 3.

That is, we recommend using entropy balancing to estimate the Riesz representer and applying TMLE to obtain the final ATE estimator. As shown above, both squared loss (Riesz regression) and KL divergence correspond to the same error minimization problem with different losses. On the other hand, KL divergence uses a basis function $\Phi(X)$ that depends only on $X$, while squared loss uses a basis function $\Phi(D, X)$ with an additional input. Although we can use logistic models for squared loss, we lose the covariate balancing property.

However, the combination of squared loss (Riesz regression) and linear models is also effective in important applications. As discussed in Kato (2025c), Riesz regression includes nearest neighbor matching as a special case. By changing the kernel (basis function), we can also derive various matching methods. Moreover, Bruns-Smith et al. (2025) finds that under Riesz regression, we can write the Neyman orthogonal score as linear in $Y$, similar to standard OLS or Ridge regression.

# 7  Conclusion

This note presents a unified framework for causal inference by connecting Riesz regression, covariate balancing, density-ratio estimation, TMLE, and matching estimators under the lens of targeted Neyman estimation. Central to this framework is the estimation of the Riesz representer, which plays a crucial role in constructing efficient ATE estimators. We demonstrate that several existing methods can be interpreted as minimizing a common error term with different loss functions, and we propose a practical implementation that combines entropy balancing and TMLE. This unified view not only clarifies the relationships among these diverse methods but also provides guidance for applied researchers in choosing robust and better estimation strategies. For theoretical details and simulation studies, see Kato (2025a,b,c).

# References

David Bruns-Smith, Oliver Dukes, Avi Feller, and Elizabeth L Ogburn. Augmented balancing weights as linear regression. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 04 2025. 5, 7, 8

Victor Chernozhukov, Whitney K. Newey, and Rahul Singh. Automatic debiased machine learning of causal and structural effects. *Econometrica*, 90(3):967–1027, 2022. 3

Victor Chernozhukov, Whitney K. Newey, Victor Quintas-Martinez, and Vasilis Syrgkanis. Automatic debiased machine learning via riesz regression, 2024. arXiv:2104.14737. 5, 6

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. 3, 7, 8

Daniel G. Horvitz and Donovan J. Thompson. A generalization of sampling without replacement from a finite universe. *Journal of the American Statistical Association*, 47(260): 663–685, 1952. 3

Kosuke Imai and Marc Ratkovic. Estimating treatment effect heterogeneity in randomized program evaluation. *The Annals of Applied Statistics*, 7(1):443 – 470, 2013. 3

Guido W. Imbens and Donald B. Rubin. *Causal Inference for Statistics, Social, and Biomedical Sciences: An Introduction.* Cambridge University Press, 2015. 1

Takafumi Kanamori, Shohei Hido, and Masashi Sugiyama. A least-squares approach to direct importance estimation. *Journal of Machine Learning Research*, 10(Jul.):1391–1445, 2009. 6

Masahiro Kato. Direct bias-correction term estimation for propensity scores and average treatment effect estimation, 2025a. arXiv: 2509.22122. 1, 5, 9

Masahiro Kato. Direct debiased machine learning via bregman divergence minimization, 2025b. aXiv: 2510.23534. 1, 4, 5, 7, 9

Masahiro Kato. Nearest neighbor matching as least squares density ratio estimation and riesz regression, 2025c. arXiv: 2510.24433. 1, 6, 8, 9

Masahiro Kato and Takeshi Teshima. Non-negative bregman divergence minimization for deep direct density ratio estimation. In *International Conference on Machine Learning (ICML)*, 2021. 6

Zhexiao Lin, Peng Ding, and Fang Han. Estimation based on nearest neighbor matching: from density ratio to average treatment effect. *Econometrica*, 91(6):2187–2217, 2023. 6

Benjamin Rhodes, Kai Xu, and Michael U. Gutmann. Telescoping density-ratio estimation. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2020. 6

Alejandro Schuler and Mark van der Laan. Introduction to modern causal inference, 2024. URL https://alejandroschuler.github.io/mci/introduction-to-modern-causal-inference.html. 3

Masashi Sugiyama, Matthias Krauledat, and Klaus-Robert Müller. Covariate shift adaptation by importance weighted cross validation. *Journal of Machine Learning Research*, 8(35): 985–1005, 2007. URL http://jmlr.org/papers/v8/sugiyama07a.html. 7

Masashi Sugiyama, Taiji Suzuki, and Takafumi Kanamori. *Density Ratio Estimation in Machine Learning*. Cambridge University Press, 2012. 3

van der Laan. Targeted maximum likelihood learning, 2006. U.C. Berkeley Division of Biostatistics Working Paper Series. Working Paper 213. https://biostats.bepress.com/ucbbiostat/paper213/. 3

Qingyuan Zhao. Covariate balancing propensity score by tailored loss functions. *The Annals of Statistics*, 47(2):965 – 993, 2019. 5, 7, 8

José R. Zubizarreta. Stable weights that balance covariates for estimation with incomplete outcome data. *Journal of the American Statistical Association*, 110(511):910–922, 2015. 7