# Learning $L_2$-Continuous Regression Functionals via Regularized Riesz Representers

Victor Chernozhukov[*]     Whitney K. Newey[†]     Rahul Singh[‡]

*MIT*                      *MIT*                      *MIT*

September 12, 2018

### Abstract

Many objects of interest can be expressed as an $L_2$ continuous functional of a regression, including average treatment effects, economic average consumer surplus, expected conditional covariances, and discrete choice parameters that depend on expectations. Debiased machine learning (DML) of these objects requires a learning a Riesz representer (RR). We provide here Lasso and Dantzig learners of the RR and corresponding learners of affine and other nonlinear functionals. We give an asymptotic variance estimator for DML. We allow for a wide variety of regression learners that can converge at relatively slow rates. We give conditions for root-n consistency and asymptotic normality of the functional learner. We give results for non affine functionals in addition to affine functionals.

---

[*]Department of Economics, MIT, Cambridge, MA 02139, U.S.A E-mail: vchern@mit.edu.

[†]Department of Economics, MIT, Cambridge, MA 02139, U.S.A E-mail: wnewey@mit.edu.

[‡]Department of Economics, MIT, Cambridge, MA 02139, U.S.A E-mail: rahul.singh@mit.edu.

# 1   Introduction

Many statistical objects of interest can be expressed as an $L_2$ (mean square) continuous functional of a conditional expectation (regression). Examples of affine regression functionals include average treatment effects, policy effects, economic average consumer surplus, and the expected conditional covariance of two random variables. Nonlinear functionals include discrete choice models that depend on regressions. Often the regression may be high dimensional, depending on many variables. There may be many covariates for a treatment effect when treatment was assigned in a complicated way. There are often many prices and covariates in the economic demand for some commodity. This variety of important examples motivates the learning of $L_2$ continuous regression functionals.

Plugging a machine learner into a functional of interest can be badly biased; e.g. see Chernozhukov et al. (2018). We use debiased/double machine learning (DML, Chernozhukov et al. 2018), based on estimating equations that have zero derivative with respect to each nonparametric component. Such debiased estimating equations are sometimes referred to as Neyman orthogonal. They can be constructed by adding the influence function of a functional of the regression learner limit. We also debias using sample splitting (Bickel, 1982, Schick, 1986), where we average over data observations different than those used by the nonparametric learners. The resulting estimators of regression functionals have second order remainders which leads to root-n consistency under regularity conditions we give.

The influence function of an $L_2$ continuous functional of a regression limit is the product of the regression residual with the Riesz representer (RR) of the functional derivative, as shown in Newey (1994). Therefore, DML of regression functionals requires a machine learner of the RR. We provide here $\ell_1$ regularized RR learners: Lasso and Dantzig selector. These automatically learn the RR from the empirical analog of equations that implicitly characterize it, without needing to know its form. We derive convergence rates for these regularized RR's and give conditions sufficient for root-n consistency and asymptotic normality of the DML estimator. DML also requires a regression learner for its construction. We allow for a variety of regression learners, requiring only a sufficiently fast $L_2$ convergence rate for the regression. We give a consistent estimator of the asymptotic variance. Results are given for nonlinear functionals as well as for affine ones. We impose only $L_2$ convergence conditions on the RR and regression learners, so that our results apply to many possible machine learners.

Debiasing via DML is based on the zero derivative of the estimating equation with respect to each nonparametric component, as in Belloni, Chernozhukov, and Hansen (2014), Farrell (2015), and Robins et al. (2013). This kind of debiasing is different than bias correcting the regression learner, as in Zhang and Zhang (2014), Belloni Chernozhukov, and Wang (2014), Belloni, Chernozhukov, and Kato (2015), Javanmard and Montanari (2014a,b; 2015), van de Geer et al. (2014), Neykov et al. (2015), Ren et al. (2015), Jankova and van de Geer (2015,

2016a,b), Bradic and Kolar (2017); Zhu and Bradic (2018). These two debiasing approaches bear some resemblance when the functional of interest is a coefficient of a partially linear model (as discussed in Chernozhukov et al., 2018), but are quite different for other functionals. The differences between these methods seem analogous to the difference between nonparametric estimation and root-n consistent functional estimation in the semiparametric literature (see Bickel, Klassen, Ritov, and Wellner, 1993 and Van der Vaart, 1991). Inference for a nonparametric regression requires bias correcting or undersmoothing the regression estimator while root-n consistent functional estimation can be based on learners that are not debiased or undersmoothed (see Newey 1994 for series regression). Similarly, DML based inference does not require the use of debiased learners. As we show, any regression learner having a fast enough convergence rate will suffice when combined with the RR learners given here.

The functionals we consider are different than those analyzed in Cai and Guo (2017). We consider nonlinear functionals as well as linear functionals where the linear combination coefficients are estimated, neither of which is allowed for in Cai and Guo (2017). Also the $L_2$ continuity of the linear functionals provides additional structure that we exploit, involving the RR, which is not exploited in Cai and Guo (2017).

Targeted maximum likelihood (van der Laan and Rubin, 2006) based on machine learners has been considered by van der Laan and Rose (2011) and large sample theory given by Luedtke and van der Laan (2016), Toth and van der Laan (2016), and Zheng et al. (2016). Here we provide DML learners via regularized RR, which are relatively simple to implement and analyze and directly target functionals of interest.

$L_2$ continuity does place us squarely in a semiparametric setting where root-n consistent efficient semiparametric estimation of the object of interest is possible under sufficient regularity conditions; see Jankova and Van De Geer (2016a). Our results apply to different objects than considered by Ning and Liu (2017), who considered machine learning of the efficient score for a parameter of an explicit semiparametric form for the distribution of the data. Unlike Ning and Liu (2017), we do not work with an explicit semiparametric form for the distribution of the data. Instead we focus on learning functionals of a nonparametric regression. Our estimators can be thought of as being based on DML of a functional of interest rather than the efficient score for a parameter of interest in an explicit form of a semiparametric model. There are many interesting examples, including those we have given, where learning via DML is more convenient and natural than embedding the functional of interest in a large, explicit semiparametric form.

We build on previous work on debiased estimating equations constructed by adding an influence function. Hasminskii and Ibragimov (1979) and Bickel and Ritov (1988) suggested such estimators for functionals of a density. Doubly robust estimating equations as in Robins, Rotnitzky, and Zhao (1995) and Robins and Rotnitzky (1995) have this structure. Newey, Hsieh, and Robins (1998, 2004) and Robins et al. (2008) further developed theory. For an affine

functional the doubly robust learner we consider is given in Chernozhukov et al. (2016). We make use of simple and general regularity conditions in Chernozhukov et al. (2016) that only require $L_2$ convergence of nonparametric learners.

The RR learners we consider are linear in a dictionary of functions. Such RR learners were previously used in Newey (1994) for asymptotic variance estimation and in Robins et al. (2007) for estimation of the inverse of the propensity score with missing data. Recently Newey and Robins (2017) considered such RR learning in efficient semiparametric estimation of linear regression functionals with low dimensional regressors. Hirshberg and Wager (2018) gave different RR estimators when the regression is restricted to a Donsker class. None of these works are about machine learning.

The Athey, Imbens, and Wager (2018) learner of the average treatment effect is based on a specific regression learner and on approximate balancing weights when the regression is linear and sparse. Our estimator allows for a wide variety of regression learners and does not restrict the regression to be sparse or linear. We do this via regularized RR learning that can also be interpreted as learning of balancing weights or inverse propensity scores, as further discussed in Section 4.

Zhu and Bradic (2017) showed that it is possible to attain root-n consistency for the coefficients of a partially linear model when the regression function is dense. Our results apply to a wide class of affine and nonlinear functionals and similarly allow the regression learner to converge at relatively slow rates.

Chernozhukov, Newey, and Robins (2018) have previously given the Dantzig learner of the RR. We innovate here by allowing the functional to depend on data other than the regressors, by giving a Lasso learner of the RR, by deriving convergence rates for both Lasso and Dantzig as learners of the true RR rather than a sparse approximation to it, by allowing for a general regression learner rather than just Dantzig, and by providing learners for nonlinear functionals. These results are innovative relative to other previous work in the ways described in the previous paragraphs.

In Section 2 we describe the objects we are interested in, their DML estimators, give a Lasso learner of the RR, and an estimator of the asymptotic variance for DML. Section 3 derives $L_2$ convergence rates of Lasso and Dantzig RR learners. Section 4 gives conditions for root-n consistency and asymptotic normality of DML and consistency of the asymptotic variance, in general and for the examples. Section 5 shows how to construct Lasso and Dantzig RR learners for nonlinear functionals and gives large sample inference results for the DML estimator and its asymptotic variance estimator.

# 2 Learning Affine Functionals

For expositional purposes we first consider objects of interest that are $L_2$ continuous affine functionals of a conditional expectation. To describe such an object let $W$ denote a data observation and consider a subvector $(Y, X')'$ where $Y$ is a scalar outcome with finite second moment and the covariate vector $X$ that takes values $x \in \mathcal{X}$, a Borel subset of $\mathbb{R}^d$. Denote the conditional expectation of $Y$ given $X$ as

$$\gamma_0(x) = \mathrm{E}[Y \mid X = x].$$

Let $m(w, \gamma)$ denote an affine functional of a possible conditional expectation function $\gamma : X \longrightarrow \mathbb{R}$ that depends on the data observation $W$. The object of interest is

$$\theta_0 = \mathrm{E}[m(W, \gamma_0)]. \tag{2.1}$$

We focus on functionals where $E[m(W, \gamma) - m(W, 0)]$ is a mean square continuous linear functional of $\gamma$. This continuity property is equivalent to the semiparametric variance bound for $\theta_0$ being finite, as discussed in Newey (1994). In this case the Riesz representation theorem implies existence of $\alpha_0(x)$ with $E[\alpha_0(X)^2]$ finite and

$$E[m(W, \gamma) - m(W, 0)] = E[\alpha_0(X)\gamma(X)] \tag{2.2}$$

for all $\gamma(x)$ with $E[\gamma(X)^2]$ finite. We refer to $\alpha_0(x)$ as the RR.

There are many important examples of this type of object. One is the average treatment effect. Here $X = (D, Z)$ and $\gamma_0(x) = \gamma_0(d, z)$, where $D \in \{0, 1\}$ is the indicator of the receipt of the treatment and $Z$ are covariates. The object of interest is

$$\theta_0 = E[\gamma_0(1, Z) - \gamma_0(0, Z)].$$

When the treatment effect is mean independent of the treatment $D$ conditional on covariates $Z$ then $\theta_0$ is the average treatment effect, Rosenbaum and Rubin (1983). Here $m(w, \gamma) = \gamma(1, z) - \gamma(0, z)$ and the RR is $\alpha_0(x) = d/\pi_0(z) - (1 - d)/[1 - \pi_0(z)]$ where $\pi_0(z)$ is the propensity score $\pi_0(z) = \Pr(D = 1 | Z = z)$. Thus $E[m(W, \gamma)]$ is mean square continuous when $E[1/\pi_0(Z)] < \infty$ and $E[1/\{1 - \pi_0(Z)\}] < \infty$.

Another interesting example is the average effect of changing the conditioning variables according to the map $x \longrightarrow t(x)$. The object of interest is

$$\theta_0 = E[\gamma_0(t(X)) - \gamma_0(X)] = \int \gamma_0(x) dF_t(dx) - E[Y],$$

where $F_t$ denotes the CDF of $t(X)$. The object $\theta_0$ is the average policy effect of a counterfactual change of covariate values similar to Stock (1989). Here $m(w, \gamma) = \gamma(t(x)) - y$ and the RR is

$\alpha_0(x) = f_t(x)/f_0(x)$ where $f_0(x)$ is the pdf of $X$ and $f_t(x)$ is the pdf of $t(X)$. $E[m(W, \gamma)]$ is mean square continuous if $E[\alpha_0(X)^2] = \int f_0(x)^{-1} f_t(x)^2 dx < \infty$.

A third object of interest is a bound on average consumer surplus for economic demand functions. Here $Y$ is the share of income spent on a commodity and $X = (P_1, Z)$, where $P_1$ is the price of the commodity and $Z$ includes income $Z_1$, prices of other goods, and other observable variables affecting utility. Let $\check{p}_1 < \bar{p}_1$ be lower and upper prices over which the price of the commodity can change, $\kappa$ a bound on the income effect, and $\omega(z)$ some weight function. The object of interest is

$$\theta_0 = E[\omega(Z) \int_{\check{p}_1}^{\bar{p}_1} (\frac{Z_1}{u}) \gamma_0(u, Z) \exp(-\kappa[u - \check{p}_1]) du],$$

where $Z_1$ is income and $u$ is a variable of integration. When individual heterogeneity in consumer preferences is independent of $X$ and $\kappa$ is a lower (upper) bound on the derivative of consumption with respect to income across all individuals then $\theta_0$ is an upper (lower) bound on the weighted average over consumers of exact consumer surplus (equivalent variation) for a change in the price of the first good from $\check{p}_1$ to $\bar{p}_1$; see Hausman and Newey (2016). Here $m(w, \gamma) = \omega(z) \int_{\check{p}_1}^{\bar{p}_1} (z_1/u) \gamma_0(u, z) \exp(-\kappa[u - \check{p}_1]) du$ and the RR is

$$\alpha_0(x) = f_0(p_1|z)^{-1} \omega(z) 1(\check{p}_1 < p_1 < \bar{p}_1)(z_1/p_1) \exp(-\kappa[p_1 - \check{p}_1]),$$

where $f_0(p_1|z)$ is the conditional pdf of $P_1$ given $Z$.

A fourth example is the average conditional covariance between $Y$ and some other variable, say $W_1$. In this case the object of interest is

$$\theta_0 = E[Cov(Y, W_1|X)] = E[W_1\{Y - \gamma_0(X)\}].$$

This object is useful in the analysis of covariance while controlling for regressors $X$ and is an important component in the coefficient $\beta_0$ of $W_1$ for a partially linear regression of $Y$ on $W_1$ and unknown functions of $x$. This object differs from the previous three examples in $m(w, \gamma)$ depending on $w$ other than the regressors $x$. Here $m(w, \gamma) = w_1\{y - \gamma(x)\}$ and the RR is $\alpha_0(x) = -E[W_1|X = x]$.

DML of $\theta_0$ can be carried out using the doubly robust moment function

$$\psi(w, \theta, \gamma, \alpha) = m(w, \gamma) - \theta + \alpha(x)[y - \gamma(x)],$$

given in Chernozhukov et al. (2016). This function has the doubly robust property that

$$0 = E[\psi(W, \theta_0, \gamma_0, \alpha)] = E[\psi(W, \theta_0, \gamma, \alpha_0)],$$

for all $\gamma$ and $\alpha$. Consequently, $\psi(w, \theta, \gamma, \alpha)$ is debiased in that any functional derivative of $E[\psi(W, \theta_0, \gamma_0, \alpha)]$ with respect to $\alpha$ and of $E[\psi(W, \theta_0, \gamma, \alpha_0)]$ with respect to $\gamma$ is zero. Therefore

a DML learner $\hat{\theta}$ can be constructed from machine learning estimators $\hat{\gamma}$ and $\hat{\alpha}$ by plugging these into the moment function $\psi(w, \theta, \gamma, \alpha)$ in place of $\gamma$ and $\alpha$ and solving for $\hat{\theta}$ from setting the sample moment of $\psi(w, \theta, \hat{\gamma}, \hat{\alpha})$ to zero.

To help avoid potentially severe finite sample bias and to avoid regularity conditions based on $\hat{\gamma}$ and $\hat{\alpha}$ being in a Donsker class, which machine learning estimators are usually not, we also use sample splitting. We construct $\hat{\gamma}$ and $\hat{\alpha}$ from observations that are not being averaged over. Let the data be $W_i$, $(i = 1, ..., n)$, assumed to be i.i.d.. Let $I_\ell$, $(\ell = 1, ..., L)$ be a partition of the observation index set $\{1, ..., n\}$ into $L$ distinct subsets of about equal size. Let $\hat{\gamma}_\ell$ and $\hat{\alpha}_\ell$ be estimators constructed from the observations that are *not* in $I_\ell$. We construct the estimator $\hat{\theta}$ by setting the sample average of $\psi(W_i, \theta, \hat{\gamma}_\ell, \hat{\alpha}_\ell)$ to zero and solving for $\theta$. This estimator has the explicit form

$$\hat{\theta} = \frac{1}{n} \sum_{\ell=1}^{L} \sum_{i \in I_\ell} \{m(W_i, \hat{\gamma}_\ell) + \hat{\alpha}_\ell(X_i)[Y_i - \hat{\gamma}_\ell(X_i)]\}. \tag{2.3}$$

A variety of regression learners $\hat{\gamma}_\ell$ of the nonparametric regression $E[Y|X]$ could be used here, as discussed in the Introduction. We also need an estimator $\hat{\alpha}_\ell$ to construct $\hat{\theta}$. We give here Lasso and Dantzig learners $\hat{\alpha}_\ell$. These learners make use of a $p \times 1$ dictionary of functions $b(x)$ where $p$ can be much bigger than $n$. The learners take the form

$$\hat{\alpha}(x) = b(x)'\hat{\rho}, \tag{2.4}$$

where $\hat{\rho}$ is a vector of estimated coefficients. For notational convenience we drop the $\ell$ subscript, with the understanding that the description which follows should be applied only to the observations not in $I_\ell$ for each $\ell$. The learners for $\alpha_0$ are based on the fact that the Riesz representation implies that for $m(w, b) = (m(w, b_1), ..., m(w, b_p))'$,

$$M = E[m(W, b) - m(W, 0)] = E[\alpha_0(X)b(X)].$$

Here we see that the cross moments $M$ between the true, unknown RR $\alpha_0(x)$ and the dictionary $b(x)$ are equal to the expectation of a known vector of functions $m(w, b) - m(w, 0)$. Consequently an unbiased estimator of $M = E[\alpha_0(X)b(X)]$ can be constructed as

$$\hat{M} = \frac{1}{n} \sum_{i=1}^{n} \{m(W_i, b) - m(W_i, 0)\}.$$

Likewise an unbiased estimator of $G = E[b(X)b(X)']$ can be constructed as

$$\hat{G} = \frac{1}{n} \sum_{i=1}^{n} \{b(X_i)b(X_i)'\}.$$

The estimator $\hat{M}$ is analogous to $\sum_{i=1}^{n} Y_i b(X_i)/n$ in Lasso and Dantzig regression. Just as $\sum_{i=1}^{n} Y_i b(X_i)/n$ is an unbiased estimator of $E[\gamma_0(X)b(X)]$ so is $\hat{M}$ an unbiased estimator of $M$.

Minimum distance versions of Lasso and Dantzig can be constructed by replacing $\sum_{i=1}^{n} Y_i b(X_i)/n$ in the Lasso objective function and Dantzig constraints by $\hat{M}$. Doing this for Lasso, while dropping $\sum_{i=1}^{n} Y_i^2/n$ term in the Lasso objective, gives an estimator

$$\hat{\rho}_L = \arg\min_{\rho}\{-2\hat{M}'\rho + \rho'\hat{G}\rho + 2r_L|\rho|_1\}. \tag{2.5}$$

The objective function here is a $\ell_1$ penalized approximation to the least squares regression of $\alpha_0(x)$ on $b(x)$, where $2r_L$ is the penalty. Making the analogous replacement in the constraints of the Dantzig selector gives a Dantzig estimator

$$\hat{\rho}_D = \arg\min_{\rho} |\rho|_1 \; s.t. |\hat{M} - \hat{G}\rho|_\infty \leq \lambda_D, \tag{2.6}$$

where $\lambda_D > 0$ is the slackness size. These two minimization problems can be thought of as minimum distance versions of Lasso and Dantzig respectively.

Either of these $\hat{\rho}$ may be used in equation (2.4) to form an estimator of the RR. This estimator of the RR may then be substituted in equation (2.3) along with a machine learning regression estimator to construct an estimator of the object of interest. We derive the properties of $\hat{\theta}$ under weak conditions that only require a relatively slow $L_2$ convergence rate for $\hat{\gamma}$. Our results on Lasso and Dantzig minimum distance can be applied to show that these produce fast enough convergence rates without assuming sparseness of the $\ell_1$ regularized approximation to the true regression.

It is interesting to note that the estimator $b(x)'\hat{\rho}$ of the RR does not require any knowledge of the form of $\alpha_0(x)$. In particular it does not depend on plugging in nonparametric estimates of components of $\alpha_0(x)$. Instead it is a linear in $b(x)$ estimator that uses $\hat{M}$ as an estimator of $M$ in an $\ell_1$ regularized least squares approximation of the least squares projection of $\alpha_0(x)$ on $b(x)$.

In the next Section we will derive convergence rates for the Lasso and Dantzig estimators of the RR and in Section 4 formulate sufficient conditions for root-n consistency and asymptotic normality of $\hat{\theta}$ from equation (2.3). For asymptotic inference we also need a consistent estimator of the asymptotic variance of $\sqrt{n}(\hat{\theta} - \theta_0)$. In order to rely only on $L_2$ rates for our asymptotic inference results we construct such an estimator by trimming the RR estimator. Let $T_\Delta(\alpha) = \alpha 1(|\alpha| \leq \Delta) + \Delta[1(\alpha > \Delta) - 1(\alpha < -\Delta)]$ for $\Delta > 0$ and

$$\hat{\alpha}_{\ell\Delta}(x) = T_\Delta(\hat{\alpha}_\ell(x)).$$

An estimator of the influence function of $\hat{\theta}$ for the ith observation is

$$\hat{\psi}_i = m(W_i, \hat{\gamma}_\ell) - \hat{\theta} + \hat{\alpha}_{\ell\Delta}(X_i)[Y_i - \hat{\gamma}_\ell(X_i)], \; i \in I_\ell, \; (\ell = 1, ..., L).$$

An estimator of the asymptotic variance is then the sample variance $\hat{V}$ of $\hat{\psi}_i$ given by

$$\hat{V} = \frac{1}{(n-1)} \sum_{i=1}^{n} (\hat{\psi}_i - \bar{\psi})^2, \; \bar{\psi} = \frac{1}{n} \sum_{i=1}^{n} \hat{\psi}_i. \tag{2.7}$$

To summarize, based on an estimated RR we have given a doubly robust machine learning estimator of a linear functional of a nonparametric regression. We have given Lasso and Dantzig estimators of the RR that are linear in approximating functions. We have also given an estimator of the asymptotic variance of $\sqrt{n}(\hat{\theta} - \theta_0)$.

# 3 Properties of Lasso and Dantzig Minimum Distance

In this Section we derive $L_2$ convergence rates for Lasso and Dantzig minimum distance estimators. We apply these result to obtain rates for regularized estimators of RRs. We also give $L_2$ rates for Lasso and Dantzig nonparametric regression estimators.

A key component of these results are regularized population approximations to unknown true functions via the Lasso and Dantzig selector. To describe these population approximations let $\alpha_0$ denote the true value of some unknown function, $\|\cdot\|$ denote the mean-square norm with $\|A\| = \sqrt{E[A(W)^2]}$ for a random variable $A(W)$, $|\rho|_1 = \sum_{j=1}^{p} |\rho_j|$ the $\ell_1$ norm of a $p \times 1$ vector $\rho$, and $|\rho|_\infty = \max_j |\rho_j|$. The coefficients of the Lasso $\ell_1$ regularized approximation are

$$\rho_L = \arg\min_\rho \|\alpha_0 - b'\rho\|^2 + 2r_0|\rho|_1,$$

where $r_0$ is a penalty that may be different than $r_L$ that appears in the estimator and will be further specified below. For $M = E[\alpha_0(X)b(X)]$ and $G = E[b(X)b(X)']$ the coefficients of the Dantzig $\ell_1$ regularized approximation are

$$\rho_D = \arg\min_\rho |\rho|_1 \text{ s.t. } |M - G\rho|_\infty \leq \lambda_0,$$

where $\lambda_0$ is a slackness constraint that may be different than the $\lambda_L$ that appears in the Dantzig estimation and will be further specified below. Because of the $\ell_1$ penalization for Lasso and the $\ell_1$ constraint for Dantzig these coefficient vectors will often be sparse.

The first condition we specify is a fundamental approximation rate hypothesis for the mean square approximation of $\alpha_0(x)$ by a linear combination of $b(x)$.

ASSUMPTION 1: *There is $\tilde{\rho}$, $\varepsilon_n$, and $B_n$ such that*

$$\|\alpha_0 - b'\tilde{\rho}\| \leq \varepsilon_n, \ |\tilde{\rho}|_1 \leq B_n.$$

For example, suppose that

$$\alpha(x) = \sum_{j=1}^{\infty} \rho_j b_j(x), \ \sum_{j=1}^{\infty} |\rho_j| < \infty,$$

where the first infinite sum is a mean square limit. Then Assumption 1 will be satisfied for

$$\tilde{\rho} = (\rho_1, ..., \rho_p)', \quad \varepsilon_n = \left\| \sum_{j=p+1}^{\infty} \rho_j b_j \right\|^2, \quad B_n = B = \sum_{j=1}^{\infty} |\rho_j|.$$

LEMMA 1: *If Assumption 1 is satisfied then for any $r_0 \geq \varepsilon_n/2$ and $\lambda_0 \geq \varepsilon_n/2$ both $\rho_L$ and $\rho_D$ exist and*

$$\|\alpha_0 - b'\rho_L\|^2 \leq 2(1 + B_n)r_0, |\rho_L|_1 \leq 1 + B_n,$$
$$\|\alpha_0 - b'\rho_D\|^2 \leq 12(1 + B_n)\lambda_0, |\rho_D|_1 \leq 1 + B_n.$$

The number $\mathcal{M}_L$ of nonzero elements of $\rho_L$ will be useful for our inference results. For $\rho = (\rho_1, ..., \rho_p)$ let $A_L$ and $A_D$ denote the index set for the components of $\rho_L$ and $\rho_D$ that are nonzero respectively and $A_L^c$ and $A_L^c$ denote the complement of those index sets in $\{1, ..., p\}$ respectively.

ASSUMPTION 2: *There are constants $k_3$ and $k_1$ such that*

$$k_3 = \inf_{\{\delta : \delta \neq 0, |\delta_{A_L^c}|_1 \leq 3|\delta_{A_L}|_1\}} \frac{\delta'G\delta}{|\delta_{A_L}|_2^2} > 0$$

$$k_1 = \inf_{\{\delta : \delta \neq 0, |\delta_{A_D^c}|_1 \leq |\delta_{A_D}|_1\}} \frac{\delta'G\delta}{|\delta_{A_D}|_2^2} > 0.$$

This Assumption allows us to give results for Lasso in terms of the number $\mathcal{M}_L$ of nonzero elements of $\rho_L$. For the Dantzig selector we will state results in terms of the effective dimension of $\rho_D$ which is given by

$$s_D = \sup_{\delta \neq 0, |\rho_D + \delta|_1 \leq |\rho_D|_1} \frac{|\delta|_1^2}{\delta'G\delta}.$$

The effective dimension is the reciprocal of the identifiability factors that were introduced in Chernozhukov et al. (2013) as a generalization of the restricted eigenvalue of Bickel et al. (2009). Let $\mathcal{M}_D$ denote the number of nonzero elements of $\rho_D$. Note that $|\rho_D + \delta|_1 \leq |\rho_D|_1$ implies $|\delta_{A_D^c}|_1 \leq |\delta_{A_D}|_1$ which then implies that $|\delta|_1^2 \leq 2|\delta_{A_D}|_1^2 \leq 2\mathcal{M}_D|\delta_M|_2^2$ where the last inequality follows by Cauchy-Schwartz. Then

$$s_D \leq \sup_{\{\delta : \delta \neq 0, |\rho_D + \delta|_1 \leq |\rho_D|_1\}} \frac{2\mathcal{M}_D|\delta_{A_D}|_2^2}{\delta'G\delta} \leq 2\mathcal{M}_D \sup_{\{\delta : \delta \neq 0, |\delta_{A_D^c}|_1 \leq |\delta_{A_D}|_1\}} \frac{|\delta_{A_D}|_2^2}{\delta'G\delta} = \frac{2}{k_1}\mathcal{M}_D.$$

Hence we can view $s_D$ as a measure of effective dimension of $\rho_D$.

The usefulness of the eigenvalue condition of Assumption 2 depends on the sparseness of the $\ell_1$ regularized approximating coefficients $\rho_L$. Many rich function classes admit sparse linear

approximations with respect to conventional dictionaries $b$. For instance, Tsybakov (2009) and Belloni, Chernozhukov, and Wang (2014) give examples of Sobolev and rearranged Sobolev balls, respectively, as function classes and elements of the Fourier basis as the dictionary $b$, in which regularized approximations have small errors.

ASSUMPTION 3: *There is $B_n^b$ such that with probability one,*

$$max_{1 \leq j \leq p}|b_j(X)| \leq B_n^b.$$

As shown in Lemma A1 in the Appendix, this condition implies that

$$|\hat{G} - G|_\infty = O_p(\varepsilon_n^G), \ \varepsilon_n^G = (B_n^b)^2 \sqrt{\frac{\ln(p)}{n}}.$$

The rates of convergence of the estimators of the RR will also depend on the convergence rate for $|\hat{M} - M|_\infty$. Here we impose a general condition to obtain rates for a variety of cases.

ASSUMPTION 4: *There is $\varepsilon_n^M$ such that*

$$|\hat{M} - M|_\infty \leq O_p(\varepsilon_n^M),$$

In what follows we will give the form of $\varepsilon_n^M$ in specific settings.

The following result gives an $L_2$ convergence rate for the Lasso estimator of the RR.

THEOREM 2: *If Assumptions 1, 3, and 4 are satisfied and $\varepsilon_n^M + \varepsilon_n^G(1 + B_n) + r_0 = o(r_L)$ then $\hat{\rho}_L$ exists with probability approaching one and*

$$\|\hat{\alpha}_L - \alpha_0\|^2 = O_p((1 + B_n)r_L).$$

*If in addition Assumption 2 is satisfied then*

$$\|\hat{\alpha}_L - \alpha_0\|^2 = O_p(\varepsilon_n^L), \ \varepsilon_n^L = \min\{\mathcal{M}_L r_L^2 + (1 + B_n)r_0, (1 + B_n)r_L\}.$$

The first conclusion gives a relatively slow rate that does not depend on Assumption 2 (i.e. without sparseness) and the second result includes the faster rate. The choice of $r_L$ that yields these results is one that converges to zero slightly faster than the variance component $\varepsilon_n^M + \varepsilon_n^G(1 + B_n)$ and the bias component $(1 + B_n)r_0$. The convergence rates then depend on $r_L$ as in the conclusion of Theorem 2.

The next result gives an $L_2$ convergence rate for the Dantzig selector.

THEOREM 3: *If Assumptions 1, 3, and 4 are satisfied and $\varepsilon_n^M + \varepsilon_n^G(1 + B_n) + \lambda_0 = o(\lambda_D)$ then*

$$\|\hat{\alpha}_D - \alpha_0\|^2 = O_p(\varepsilon_n^D), \varepsilon_n^D = \min\{s_D \lambda_D^2 + (1 + B_n)\lambda_0, (1 + B_n)\lambda_D\}.$$

These rate results are useful in specifying conditions for root-n consistency and asymptotic normality of $\hat{\theta}$ and consistency of the asymptotic variance estimator, to which we now turn.

# 4 Large Sample Inference

In this Section we give conditions for root-n consistency and asymptotic normality of the estimator $\hat{\theta}$. We also show that the asymptotic variance estimator is consistent. These results allow us to carry out large sample inference about the object of interest in the usual way. We also apply the general results to each of the examples. Recall that the estimator is

$$\hat{\theta} = \frac{1}{n} \sum_{\ell=1}^{L} \sum_{i \in I_\ell} \{m(W_i, \hat{\gamma}_\ell) + \hat{\alpha}_\ell(X_i)[Y_i - \hat{\gamma}_\ell(X_i)]\}. \tag{4.1}$$

where $\hat{\alpha}_\ell(x) = b(x)' \hat{\rho}_\ell$.

We impose the following conditions.

ASSUMPTION 5: $Var(Y|X)$ is bounded, $\alpha_0(x)$ is bounded, $E[m(W, \gamma_0)^2] < \infty$, and $E[\{m(W, \gamma) - m(W, \gamma_0)\}^2]$ is continuous at $\gamma_0$ in $\|\gamma - \gamma_0\|$.

Boundedness of $Var(Y|X)$ is standard in the regression literature. It may be possible to weaken the second and third conditions but it is beyond the scope of the paper to do so. All of these conditions are imposed to make sure that only $L_2$ rates are needed for $\hat{\gamma}$ and for $\hat{\alpha}$. This helps the results apply to machine learning estimators where only $L_2$ convergence rates are available.

ASSUMPTION 6: There are $B_n^m$, $B_n^b$, and $A(W)$ such that $A(W)$ is sub-Gaussian and

$$max_{1 \leq j \leq p} |m(W, b_j)| \leq B_n^m A(W).$$

This is a primitive condition that leads to a convergence rate for $\hat{M}$.

LEMMA 4: If Assumption 6 is satisfied then

$$|\hat{M} - M|_\infty = O_p(B_n^m \sqrt{\frac{\ln(p)}{n}}).$$

Note that for $m(w, b_j) = y b_j(x)$ the minimization problems in equations (2.5) and (2.6) are those for the Lasso and Dantzig regression respectively. Thus the convergence rates of Theorems 2 and 3 apply to obtain population $L_2$ rates for Lasso and Dantzig learners for $\gamma_0$.

Our results for $\hat{\theta}$ will rely on a convergence rate for $\hat{\gamma}$. In order to allow these results to apply to as wide a variety of machine learning estimators $\hat{\gamma}$ as possible we just hypothesize such a rate.

ASSUMPTION 7: $\|\hat{\gamma} - \gamma_0\| = O_p(\varepsilon_n^\gamma)$, $\varepsilon_n^\gamma \longrightarrow 0$.

The results of Section 3 imply such a rate for Lasso or Dantzig selector. The next condition imposes rates that will be sufficient for root-n consistency of $\hat{\theta}$. Let

$$\varepsilon_n^\alpha = [B_n^m + \left(B_n^b\right)^2 (1 + B_n)]\sqrt{\frac{\ln(p)}{n}}$$

ASSUMPTION 8: *Either i) $\hat{\alpha} = \hat{\alpha}_L$ for $r_L$ such that $\varepsilon_n^\alpha + r_0 = o(r_L)$, $\mathcal{M}_L r_L^2 \longrightarrow 0$, $(1 + B_n)r_0 \longrightarrow 0$, and $\sqrt{n}\sqrt{\mathcal{M}_L r_L^2 + (1 + B_n)r_0}\varepsilon_n^\gamma \longrightarrow 0$; or ii) $\hat{\alpha} = \hat{\alpha}_D$ for $\lambda_D$ such that $\varepsilon_n^\alpha + \lambda_0 = o(\lambda_D)$, $s_D \lambda_D^2 \longrightarrow 0$, $(1 + B_n)\lambda_0 \longrightarrow 0$, and $\sqrt{n}\sqrt{s_D \lambda_D^2 + (1 + B_n)\lambda_0}\varepsilon_n^\gamma \longrightarrow 0$.*

The next results gives large sample inference based on $\hat{\theta}$.

THEOREM 5: *If Assumptions 1-3 and 5-8 are satisfied then for $\psi_0(w) = m(w, \gamma_0) - \theta_0 + \alpha_0(x)[y - \gamma_0(x)]$,*

$$\sqrt{n}(\hat{\theta} - \theta_0) = \frac{1}{\sqrt{n}} \sum_{i=1}^n \psi_0(W_i) + o_p(1).$$

*If in addition $\Delta \longrightarrow \infty$ and $\Delta \varepsilon_n^\gamma \longrightarrow 0$ then $\hat{V} \xrightarrow{p} V = E[\psi_0(W)^2]$.*

This result allows $\gamma_0$ to be "dense" and estimated at relatively slow rates if $\hat{\alpha}$ converges at a sufficiently fast $L_2$ rate. We now give more specific regularity conditions for the examples.

## 4.1 Average Treatment Effect

For the average treatment effect we consider a dictionary of the form $b(x) = [dq(z)', (1-d)q(z)']'$ where $q(z)$ is a $(p/2) \times 1$ dictionary of functions of the covariates $z$. Note that $m(w, b) = [q(z)', -q(z)']'$ so that

$$\hat{M}_\ell = \begin{pmatrix} \bar{q}_\ell \\ -\bar{q}_\ell \end{pmatrix}, \bar{q}_\ell = \frac{1}{n - n_\ell} \sum_{i \notin I_\ell} q(Z_i).$$

Let $\hat{\rho}_\ell^d$ be the estimated coefficients of $dq(z)$ and $\hat{\rho}_\ell^{1-d}$ the estimated coefficients of $(1 - d)q(z)$. Then the RR learner is given by

$$\hat{\alpha}_\ell(X_i) = D_i \hat{\omega}_{\ell i}^d + (1 - D_i)\hat{\omega}_{\ell i}^{1-d}, \ \hat{\omega}_{\ell i}^d = q(Z_i)'\hat{\rho}_\ell^d, \ \hat{\omega}_{\ell i}^{1-d} = q(Z_i)'\hat{\rho}_\ell^{1-d},$$

where $\hat{\omega}_{\ell i}^d$ and $\hat{\omega}_{\ell i}^{1-d}$ might be thought of as "weights." These weights sum to one if $q(z)$ includes a constant but need not be nonnegative. The first order conditions for Lasso and the constraints for Dantzig are that for each $j$,

$$\left| \frac{1}{n - n_\ell} \sum_{i \notin I_\ell} q_j(Z_i)[1 - D_i \hat{\omega}_{\ell i}^d] \right| \le r, \ \left| \frac{1}{n - n_\ell} \sum_{i \notin I_\ell} q_j(Z_i)[1 + (1 - D_i)\omega_{\ell i}^{1-d}] \right| \le r, \qquad (4.2)$$

where $r = r_L$ for Lasso and $r = \lambda_D$ for Dantzig. Here we see that RR learner sets the weights $\hat{\omega}^d_{\ell i}$ and $\hat{\omega}^{1-d}_{\ell i}$ to approximately "balance" the overall sample average with the treated and untreated averages for each element of the dictionary $q(z)$. The resulting learner of the ATE is

$$\hat{\theta} = \frac{1}{n} \sum_{\ell=1}^{L} \sum_{i \in I_\ell} \{\hat{\gamma}_\ell(1, Z_i) - \hat{\gamma}_\ell(0, Z_i) + \hat{\alpha}_\ell(X_i)[Y_i - \hat{\gamma}_\ell(X_i)]\}. \tag{4.3}$$

The conditions we give are sufficient for this estimator to be root-n consistent when $\hat{\gamma}_\ell$ has a sufficiently fast $L_2$ convergence rate. The constraints of equation (4.2) are similar to those of Zubizarreta (2015) and Athey, Imbens, and Wager (2017) though the source of these constraints is $\ell_1$ regularized best $L_2$ approximation of the RR $\alpha(x) = \pi_0(z)^{-1}d - [1 - \pi_0(z)]^{-1}(1 - d)$ by a linear combination of the dictionary $b(x)$. We show here that this type of balancing is sufficient to debias any regression learner under sufficient regularity conditions.

THEOREM 6: *If i) there is $C > 0$ with $C < \pi_0(z) = \Pr(D = 1|z) < 1 - C$, $Var(Y|X)$ is bounded; ii) there is $B^q_n$ with $\max_{j \leq p/2} \sup_z |q_j(Z)| \leq B^q_n$ and Assumption 8 is satisfied for*

$$\varepsilon^\alpha_n = [B^q_n + (B^q_n)^2 (1 + B_n)]\sqrt{\frac{\ln(p)}{n}};$$

*and iii) Assumptions 1, 2, and 7 are satisfied, then for $\alpha_0(x) = \pi_0(z)^{-1}d - [1 - \pi_0(z)]^{-1}(1 - d)$ and $\psi_0(w) = \gamma_0(1, z) - \gamma_0(0, z) - \theta_0 + \alpha_0(x)[y - \gamma_0(x)]$,*

$$\sqrt{n}(\hat{\theta} - \theta_0) = \frac{1}{\sqrt{n}} \sum_{i=1}^{n} \psi_0(W_i) + o_p(1).$$

*If in addition $\Delta \longrightarrow \infty$ and $\Delta \varepsilon^\gamma_n \longrightarrow 0$ then $\hat{V} \xrightarrow{p} V = E[\psi_0(W)^2]$.*

In comparison with Athey, Imbens, and Wager (2018) this result depends on relatively fast estimation of the RR, or equivalently the dictionary balancing weights, while allowing for relatively slow estimation of the regression. This result can be applied to any regression estimator $\hat{\gamma}$. Another advantage of the estimator here is its DML form allows us to trade-off rates at which the mean and the inverse propensity score are estimated while maintaining root-n consistency, when a Lasso or Dantzig regression estimator is used. Also, we do not require that the regression be linear and sparse.

## 4.2 Average Policy Effect

For the average policy effect let $b(x)$ be a dictionary satisfying Assumption 3. Note that $m(w, b) = b(t(x)) - y$, so that

$$\hat{M}_\ell = \frac{1}{n - n_\ell} \sum_{i \notin I_\ell} b(t(X_i)).$$

For $\hat{\rho}_\ell$ equal to the Lasso or Dantzig coefficients, the learner of the RR is given by $\hat{\alpha}_\ell(x) = b(x)'\hat{\rho}_\ell$. The first order conditions for Lasso and the Dantzig constraints are that for each $j$

$$\left| \frac{1}{n - n_\ell} \sum_{i \notin I_\ell} [b_j(t(X_i)) - b_j(X_i)\hat{\alpha}_\ell(X_i)] \right| \leq r.$$

Here $\hat{\alpha}_\ell(X_i)$ acts approximately as a reweighting scheme in making the sample average of the dictionary after transformation $b(t(X_i))$ be approximately equal to the sample average of the reweighted dictionary $b(X_i)\hat{\alpha}_\ell(X_i)$. The resulting learner of the average policy effect is

$$\hat{\theta} = \frac{1}{n} \sum_{\ell=1}^{L} \sum_{i \in I_\ell} \{\hat{\gamma}_\ell(t(X_i)) + \hat{\alpha}_\ell(X_i)[Y_i - \hat{\gamma}_\ell(X_i)] - Y_i\}. \tag{4.4}$$

THEOREM 7: *If i) there is $C > 0$ with $1/C \leq \alpha_0(x) = f_t(x)/f_0(x) \leq C$, $Var(Y|X)$ is bounded; ii) Assumptions 3 and 8 are satisfied for*

$$\varepsilon_n^\alpha = [B_n^b + (B_n^b)^2 (1 + B_n)]\sqrt{\frac{\ln(p)}{n}};$$

*iii) Assumptions 1, 2, and 7 are satisfied then for $\psi_0(w) = \gamma_0(t(x)) - y - \theta_0 + \alpha_0(x)[y - \gamma_0(x)]$,*

$$\sqrt{n}(\hat{\theta} - \theta_0) = \frac{1}{\sqrt{n}} \sum_{i=1}^{n} \psi_0(W_i) + o_p(1).$$

*If in addition $\Delta \longrightarrow \infty$ and $\Delta \varepsilon_n^\gamma \longrightarrow 0$ then $\hat{V} \xrightarrow{p} V = E[\psi_0(W)^2]$.*

The third example, estimation of a bound for average equivalent variation, is treated in detail in Chernozhukov, Hausman, and Newey (2018). We consider here the fourth example.

## 4.3 Expected Conditional Covariance

For the expected conditional covariance let $b(x)$ be a dictionary satisfying Assumption 3. Note that $m(w, b) - m(w, 0) = -w_1 b(x)$ so that

$$\hat{M}_\ell = \frac{-1}{n - n_\ell} \sum_{i \notin I_\ell} b(X_i) W_{1i}.$$

Here the Lasso or Dantzig are those obtained from Lasso or Dantzig regression where the dependent variable is $-W_{1i}$. The resulting learner of the expected conditional covariance is

$$\hat{\theta} = \frac{1}{n} \sum_{\ell=1}^{L} \sum_{i \in I_\ell} \{W_{1i} + \hat{\alpha}_\ell(X_i)\}[Y_i - \hat{\gamma}_\ell(X_i)]\}. \tag{4.5}$$

This estimator

THEOREM 8: *If i) $E[W_1^2|X]$ and $E[Y^2|X]$ are bounded, $E[W_1^2Y^2] < \infty$; ii) $W_1$ is sub-Gaussian and Assumptions 3 and 8 are satisfied for*

$$\varepsilon_n^\alpha = [B_n^b + \left(B_n^b\right)^2 (1 + B_n)]\sqrt{\frac{\ln(p)}{n}};$$

*and iii) Assumptions 1, 2, and 7 are satisfied then for $\psi_0(w) = [w + \alpha_0(x)][y - \gamma_0(x)] - \theta_0$,*

$$\sqrt{n}(\hat{\theta} - \theta_0) = \frac{1}{\sqrt{n}} \sum_{i=1}^n \psi_0(W_i) + o_p(1).$$

*If in addition $\Delta \longrightarrow \infty$ and $\Delta \varepsilon_n^\gamma \longrightarrow 0$ then $\hat{V} \xrightarrow{p} V = E[\psi_0(W)^2].$*

This result gives root-n consistency and asymptotic normality of the expected conditional covariance estimator when the regression estimator converges fast enough in $L_2$ and when $E[W_1|X]$ is estimated by Lasso or Dantzig. This asymmetric treatment may be useful in settings where one wants to allow one of the conditional expectation functions to be estimated at a slower rate.

For further bias reduction estimation of $E[Y|X]$ and $E[W_1|X]$ from different samples may be warranted, as in Newey and Robins (2018). It is beyond the scope of this paper to analyze such estimators.

# 5　Nonlinear Functionals

Debiased machine learning estimators of $\theta_0 = E[m(W, \gamma_0)]$ for nonlinear $m(w, \gamma)$ can also be constructed. The estimator is similar to the linear functional case except that the RR is that of a linearization and a different $\hat{M}$ is needed. In this Section we show how to construct $\hat{M}$ that can be used to machine learn the RR and give conditions that are sufficient for valid large sample inference for nonlinear functionals.

As before a RR is important in the construction of the estimator. Here the RR is that for a linearization of the functional. Suppose that $m(w, \gamma)$ has a Gateaux derivative $D(w, \zeta, \gamma)$ where $\zeta$ represents a deviation from $\gamma$ and $D(w, \zeta, \gamma)$ is linear in $\zeta$. That is suppose that

$$\left.\frac{d}{d\tau}m(w, \gamma + \tau\zeta)\right|_{\tau=0} = D(w, \zeta, \gamma),$$

where $\tau$ is a scalar. We will assume that $E[D(W, \gamma, \gamma_0)]$ is a linear mean square continuous functional of $\gamma$ so that there is a RR $\alpha_0(x)$ satisfying

$$E[D(W, \gamma, \gamma_0)] = E[\alpha_0(X)\gamma(X)],$$

16

for all $\gamma(x)$ with finite second moment. This Riesz representation is analogous to equation (2.2) with the functional $m(w, \gamma) - m(w, 0)$ replaced by the first order approximation $D(w, \gamma, \gamma_0)$. The Riesz representation implies that for $D(w, b, \gamma_0) = (D(w, b_1, \gamma_0), ..., D(w, b_p, \gamma_0))'$,

$$M = E[D(W, b, \gamma_0)] = E[\alpha_0(X)b(X)].$$

A learner $\hat{\theta}$ can be constructed from an estimator $\hat{\alpha}_\ell(x)$ of the RR $\alpha_0(x)$ and a learner $\hat{\gamma}_\ell(x)$ of $E[Y|X = x]$ exactly as in equation (2.3). This estimator may not be doubly robust due to the nonlinearity of $m(w, \gamma)$ in $\gamma$. Nevertheless it will have zero first order bias and so be root-n consistent and asymptotically normal under sufficient regularity conditions. It has zero first order bias because $\alpha_0(x)[y - \gamma_0(x)]$ is the influence function for $E[m(W, \gamma)]$, as shown in Newey (1994), and because a sample average plus an average of an estimate of that influence function has zero order bias; see Chernozhukov et al. (2016).

An estimator $\hat{\alpha}_\ell(x)$ is needed to construct $\hat{\theta}$. We continue to consider estimators $\hat{\alpha}_\ell(x)$ described in Section 2, but based on a different $\hat{M}_\ell$, where it is now convenient to include an $\ell$ subscript. For a machine learning estimator $\hat{\gamma}_{\ell,\ell'}$ of $E[Y|X]$ obtained from observations not in either $I_\ell$ or $I_{\ell'}$ the estimator $\hat{M}_\ell$ is given by

$$\hat{M}_\ell = (\hat{M}_{\ell 1}, ..., \hat{M}_{\ell p})',$$
$$\hat{M}_{\ell j} = \frac{d}{d\tau} \left( \frac{1}{n - n_\ell} \right) \sum_{\ell' \neq \ell} \sum_{i \in I_{\ell'}} m(W_i, \hat{\gamma}_{\ell,\ell'} + \tau b_j) = \left( \frac{1}{n - n_\ell} \right) \sum_{\ell' \neq \ell} \sum_{i \in I_{\ell'}} D(W_i, b_j, \hat{\gamma}_{\ell,\ell'}).$$

This estimator uses further sample splitting where $\hat{M}$ is constructed by averaging over observations that are not used in $\hat{\gamma}_{\ell,\ell'}$. For convenience we have used the same partitioning of the observations as before. This additional sample splitting helps us allow for $p$ to still be large in this setting where we are plugging in a nonparametric estimator into many sample moments.

Next we obtain a convergence rate for $\hat{M}$.

ASSUMPTION 9: *There is $\varepsilon > 0$, $B_n^D$, $B_n^\Delta$ and sub-Gaussian $A(W)$ such that for all $\gamma$ with $\|\gamma - \gamma_0\| \leq \varepsilon$, i)*

$$\max_j |D(W, b_j, \gamma)| \leq B_n^D A(W),$$

*ii)* $\max_j |E[D(W, b_j, \gamma) - D(W, b_j, \gamma_0)]| \leq B_n^\Delta \|\gamma - \gamma_0\|.$

LEMMA 9: *If Assumptions 7 and 8 are satisfied then*

$$|\hat{M} - M|_\infty = O_p(\varepsilon_n^D), \quad \varepsilon_n^D = \left( B_n^D \sqrt{\frac{\ln(p)}{n}} + B_n^\Delta \varepsilon_n^\gamma \right).$$

To allow for nonlinearity of $m(w, \gamma)$ in $\gamma$ we impose the following condition

17

ASSUMPTION 10: *There are $\varepsilon$, $C > 0$ such that for all $\gamma$ with $\|\gamma - \gamma_0\| \leq \varepsilon$,*

$$|E[m(W, \gamma) - m(W, \gamma_0) - D(W, \gamma - \gamma_0, \gamma_0)]| \leq C\|\gamma - \gamma_0\|^2.$$

This condition implies that $E[m(W, \gamma)]$ is Frechet differentiable in $\|\gamma - \gamma_0\|$ with derivative $E[D(W, \gamma - \gamma_0, \gamma_0)]$. It is a specific condition that corresponds to $E[m(W, \gamma)]$ being an $L_2$ differentiable function.

Let

$$\varepsilon_n^\alpha = [B_n^D + (B_n^b)^2(1 + B_n)]\sqrt{\frac{\ln(p)}{n}} + B_n^\Delta \varepsilon_n^\gamma$$

THEOREM 10: *If Assumptions 1-3, 5, and 7-10 are satisfied with $\varepsilon_n^\gamma = o(n^{-1/4})$ and $E[m(W, \gamma_0)^2] < \infty$, then for $\psi_0(w) = m(w, \gamma_0) - \theta_0 + \alpha_0(x)[y - \gamma_0(x)]$,*

$$\sqrt{n}(\hat{\theta} - \theta_0) = \frac{1}{\sqrt{n}} \sum_{i=1}^{n} \psi_0(W_i) + o_p(1).$$

*If in addition $\Delta_n \longrightarrow \infty$ and $\Delta_n \varepsilon_n^\gamma \longrightarrow 0$ then $\hat{V} \xrightarrow{p} V = E[\psi_0(W)^2]$.*

# 6   Appendix: Proofs of Results

In this Appendix we give the proofs of the results of the paper, partly based on useful Lemmas that are stated and proved in this Appendix. The first Lemma states a well known necessary condition for minimizing the Lasso objective function.

LEMMA A0: *For any $p \times 1$ vector $\hat{M}$, $p \times p$ positive semi-definite $\hat{G}$, and $r > 0$, if $\rho^* = \arg\min_\rho\{-2\hat{M}'\rho + \rho'\hat{G}\rho + 2r|\rho|_1\}$ then*

$$|\hat{M} - \hat{G}\rho^*|_\infty \leq r.$$

Proof: Because the objective function is convex in $\rho$, a necessary condition for minimization is that 0 belongs to the sub-differential of the objective, i.e.

$$0 \in -2\hat{M} + 2\hat{G}\rho^* + 2r([-1, 1] \times ... \times [-1, 1])'.$$

Therefore for each $j$ we have

$$0 \leq -2\hat{M}_j + 2e_j'\hat{G}\rho^* + 2r, \ 0 \geq -2\hat{M}_j + 2e_j'\hat{G}\rho^* - 2r,$$

where $e_j$ is the $j^{th}$ unit vector. Dividing through by 2 and adding $\hat{M}_j - e_j'\hat{G}\rho^*$ both sides of each inequality gives

$$-r \leq \hat{M}_j - e_j'\hat{G}\rho^* \leq r,$$

that is,

$$|\hat{M}_j - e_j'\hat{G}\rho^*| \leq r.$$

The conclusion follows because this inequality holds for each $j$. $Q.E.D.$

**Proof of Lemma 1:** For $r_0 \geq \varepsilon_n/2$,

$$\|\alpha_0 - b'\rho_L\|^2 + 2r_0|\rho_L|_1 \leq \|\alpha_0 - b'\tilde{\rho}\|^2 + 2r_0|\tilde{\rho}|_1 \leq \varepsilon_n + 2r_0B_n \leq 2r_0(1 + B_n).$$

The first inequality in the conclusion follows from

$$\|\alpha_0 - b'\rho_L\|^2 \leq \|\alpha_0 - b'\rho_L\|^2 + 2r_0|\rho_L|_1 \leq 2r_0(1 + B_n).$$

The second inequality in the conclusion follows from

$$2r_0|\rho_L|_1 \leq \|\alpha_0 - b'\rho_L\|^2 + 2r_0|\rho_L|_1 \leq 2r_0(1 + B_n),$$

and dividing through by $2r_0$.

Next, note that the first order conditions for $\rho_L$ with $r_0 = \lambda_0$ imply that

$$|M - G\rho_L|_\infty \leq \lambda_0.$$

Therefore $\rho_L$ for $r_0 = \lambda_0$ satisfies the constraints of the Dantzig minimization problem, and hence

$$|\rho_D|_1 \leq |\rho_L|_1 \leq 1 + B_n,$$
$$|M - G\rho_D|_\infty \leq \lambda_0, |M - G\rho_L|_\infty \leq \lambda_0.$$

Furthermore by the triangle inequality we have

$$\|\alpha_0 - b'\rho_D\| \leq \|\alpha_0 - b'\rho_L\| + \|b'(\rho_L - \rho_D)\| \leq \sqrt{2\lambda_0(1 + B_n)} + \sqrt{(\rho_L - \rho_D)'G(\rho_L - \rho_D)}$$
$$\leq \sqrt{2\lambda_0(1 + B_n)} + \sqrt{|\rho_L - \rho_D|_1|M - G\rho_D - (M - G\rho_L)|_\infty}$$
$$\leq \sqrt{2\lambda_0(1 + B_n)} + \sqrt{2(1 + B_n)2\lambda_0}.$$

Squaring both sides then gives

$$\|\alpha_0 - b'\rho_D\|^2 \leq 12(1 + B_n)\lambda_0. \ Q.E.D.$$

We next give some Lemmas that will be useful in the proof of Theorem 2.

19

The following result gives the rate of convergence of $|\hat{G} - G|_\infty$. Let $\|A(W)\|_{\Psi_2}$ be the sub-Gaussian norm of a random variable $A(W)$ as in Vershynin (2018).

LEMMA A1: *If Assumption 3 is satisfied then*

$$|\hat{G} - G|_\infty = O_p(\varepsilon_n^G), \ \varepsilon_n^G = (B_n^b)^2 \sqrt{\frac{\ln(p)}{n}}.$$

Proof: Define

$$T_{ijk} = b_j(X_i)b_k(X_i) - E[b_j(X_i)b_k(X_i)], \ U_{jk} = \frac{1}{n}\sum_{i=1}^n T_{ijk}.$$

For any constant $C$,

$$\Pr(|\hat{G} - G|_\infty \geq C\varepsilon_n^G) \leq \sum_{j,k=1}^p \mathbb{P}(|U_{jk}| > C\varepsilon_n^G) \leq p^2 \max_{j,k} \mathbb{P}(|U_{jk}| > C\varepsilon_n^G)$$

Note that $E[T_{ijk}] = 0$ and

$$|T_{ijk}| \leq |b_j(X_i)| \cdot |b_k(X_i)| + E[|b_j(X_i)| \cdot |b_k(X_i)|] \leq 2(B_n^b)^2.$$

Define $K = \|T_{ijk}\|_{\Psi_2} \leq 2(B_n^b)^2$. By Hoeffding's inequality there is a constant $c$ such that

$$p^2 \max_{j,k} \mathbb{P}(|U_{jk}| > C\varepsilon_n^G) \leq 2p^2 \exp\left(-\frac{cn(C\varepsilon_n^G)^2}{K^2}\right) \leq 2p^2 \exp\left(-\frac{cn(C\varepsilon_n^G)^2}{4(B_n^b)^4}\right)$$

$$\leq 2\exp\left(\ln(p)[2 - \frac{cC^2}{4}]\right) \longrightarrow 0.$$

For any $C > \sqrt{8/c}$. Thus for large enough $C$, $\Pr(|\hat{G} - G|_\infty \geq C\varepsilon_n^G) \longrightarrow 0$, implying the conclusion. *Q.E.D.*

LEMMA A2: *For any $\alpha(x)$ and $\tilde{M}_\alpha = \sum_{i=1}^n b(X_i)\alpha(X_i)/n$,*

$$-2\hat{M}'\rho + \rho'\hat{G}\rho = \|\alpha - b'\rho\|_n^2 - 2(\hat{M} - \tilde{M}_\alpha)'\rho - \frac{1}{n}\sum_{i=1}^n \alpha(X_i)^2.$$

*where $\|A\|_n = \sqrt{\sum_{i=1}^n [A(W_i)^2]/n}$.*

Proof: Follows by adding and subtracting $\tilde{M}_\alpha'\rho + \sum_{i=1}^n \alpha(X_i)^2/n$. *Q.E.D.*

LEMMA A3: *For any $\rho$ and $\alpha(x)$,*

$$\|\alpha - b'\hat{\rho}_L\|_n^2 + 2r_L|\hat{\rho}_L|_1 \leq \|\alpha - b'\rho\|_n^2 + 2r_L|\rho|_1 + 2|\hat{M} - \tilde{M}_\alpha|_\infty|\hat{\rho}_L - \rho|_1.$$

Proof: $\hat{\rho}_L$ minimizes $\hat{S}(\rho) = -2\hat{M}'\rho + \rho'\hat{G}\rho + 2r_L|\rho|_1$. Beginning with the inequality $\hat{S}(\hat{\rho}_L) \leq \hat{S}(\rho)$, Lemma A2 and cancellation of $-\frac{1}{n}\sum_{i=1}^{n}\alpha(X_i)^2$ yields

$$\|\alpha - b'\hat{\rho}_L\|_n^2 - 2(\hat{M} - \tilde{M}_\alpha)'\hat{\rho}_L + 2r_L|\hat{\rho}_L|_1 \leq \|\alpha - b'\rho\|_n^2 - 2(\hat{M} - \tilde{M}_\alpha)'\rho + 2r_L|\rho|_1.$$

Adding $2(\hat{M} - \tilde{M}_\alpha)'\hat{\rho}_L$ to both sides then gives

$$\|\alpha - b'\hat{\rho}_L\|_n^2 + 2r_L|\hat{\rho}_L|_1 \leq \|\alpha - b'\rho\|_n^2 + 2r_L|\rho|_1 + 2(\hat{M} - \tilde{M}_\alpha)'(\hat{\rho}_L - \rho)$$
$$\leq \|\alpha - b'\rho\|_n^2 + 2r_L|\rho|_1 + 2|\hat{M} - \tilde{M}_\alpha|_\infty|\hat{\rho}_L - \rho|_1. \ Q.E.D.$$

LEMMA A4: *If Assumptions 1, 3, and 4 are satisfied then for $\alpha_L(x) = b(x)'\rho_L$,*

$$|\hat{M} - \tilde{M}_{\alpha_L}|_\infty \leq O_p(\varepsilon_n^M + \varepsilon_n^G(1 + B_n) + r_0).$$

.

Proof: The first order conditions for $\rho_L$ imply

$$|M - G\rho_L|_\infty \leq r_0.$$

Then by the triangle and Holder inequality, Lemmas 1 and A1, and $\tilde{M}_{\alpha_L} = \hat{G}\rho_L$

$$|\hat{M} - \tilde{M}_{\alpha_L}|_\infty \leq |\hat{M} - M|_\infty + |M - G\rho_L|_\infty + |(G - \hat{G})\rho_L|_\infty$$
$$\leq |\hat{M} - M|_\infty + |M - G\rho_L|_\infty + |G - \hat{G}|_\infty|\rho_L|_1$$
$$= O_p(\varepsilon_n^M + r_0 + \varepsilon_n^G(1 + B_n)). \ Q.E.D.$$

LEMMA A5: *If Assumptions 1, 3, and 4 are satisfied and $\varepsilon_n^M + \varepsilon_n^G(1 + B_n) + r_0 = o(r_L)$ then $|\hat{\rho}_L|_1 = O_p(1 + B_n)$.*

Proof: By the definition of $\alpha_L(x) = b(x)'\rho_L$, Lemma A3 for $\alpha(x) = \alpha_L(x)$ and $\rho = \rho_L$, and Lemma A4,

$$2r_L|\hat{\rho}_L|_1 \leq \|\alpha_L - b'\hat{\rho}_L\|_n^2 + 2r_L|\hat{\rho}_L|_1 \leq 2r_L|\rho_L|_1 + 2|\hat{M} - \tilde{M}_{\alpha_L}|_\infty|\hat{\rho}_L - \rho_L|_1$$
$$\leq 2r_L|\rho_L|_1 + [O_p(\varepsilon_n^M + \varepsilon_n^G(1 + B_n) + r_0)]|\hat{\rho}_L - \rho_L|_1.$$

Dividing through by $2r_L$ and using $[\varepsilon_n^M + \varepsilon_n^G(1 + B_n) + r_0]/r_L \longrightarrow 0$ it follows that

$$|\hat{\rho}_L|_1 \leq |\rho_L|_1 + o_p(1)|\hat{\rho}_L - \rho_L|_1 \leq |\rho_L|_1 + o_p(1)(|\hat{\rho}_L|_1 + |\rho_L|_1). \tag{6.1}$$

Subtracting $o_p(1)|\hat{\rho}_L|_1$ from both sides then gives

$$[1 - o_p(1)]|\hat{\rho}_L|_1 \leq [1 + o_p(1)]|\rho_L|_1 \leq [1 + o_p(1)](1 + B_n) = O_p(1 + B_n),$$

implying the conclusion. *Q.E.D.*

LEMMA A6: *If Assumptions 1, 3, and 4 are satisfied and $\varepsilon_n^M + \varepsilon_n^G(1 + B_n) + r_0 = o(r_L)$ then with probability approaching one $\sum_{j \in A_L^c} |\hat{\rho}_{Lj} - \rho_{Lj}| \le 3 \sum_{j \in A_L} |\hat{\rho}_{Lj} - \rho_{Lj}|$.*

Proof: It follows as in equation (6.1) of the proof of Lemma A5 that

$$|\hat{\rho}_L|_1 \le |\rho_L|_1 + o_p(1)|\hat{\rho}_L - \rho_L|_1.$$

Therefore with probability approaching one,

$$|\hat{\rho}_L|_1 \le |\rho_L|_1 + \frac{1}{2}|\hat{\rho}_L - \rho_L|_1.$$

Note that $|\rho_{Lj}| + |\hat{\rho}_{Lj} - \rho_{Lj}| - |\hat{\rho}_{Lj}| = 0$ when $\rho_{Lj} = 0$ and that $|\rho_{Lj}| - |\hat{\rho}_{Lj}| \le |\hat{\rho}_{Lj} - \rho_{Lj}|$ by the triangle inequality. Then adding $|\hat{\rho}_L - \rho_L|_1/2$ to and subtracting $|\hat{\rho}_L|_1$ from both sides gives

$$\frac{1}{2}|\hat{\rho}_L - \rho_L|_1 \le |\rho_L|_1 + |\hat{\rho}_L - \rho_L|_1 - |\hat{\rho}_L|_1 = \sum_{j=1}^p (|\rho_{Lj}| + |\hat{\rho}_{Lj} - \rho_{Lj}| - |\hat{\rho}_{Lj}|)$$

$$= \sum_{j \in A_L} (|\rho_{Lj}| + |\hat{\rho}_{Lj} - \rho_{Lj}| - |\hat{\rho}_{Lj}|) \le 2 \sum_{j \in A_L} |\hat{\rho}_{Lj} - \rho_{Lj}|.$$

Note that $|\hat{\rho}_L - \rho_L|_1 = \sum_{j \in A_L^c} |\hat{\rho}_{Lj} - \rho_{Lj}| + \sum_{j \in A_L} |\hat{\rho}_{Lj} - \rho_{Lj}|$, so multiplying both sides by 2 and subtracting $\sum_{j \in A_L} |\hat{\rho}_{Lj} - \rho_{Lj}|$ from both sides gives the result. *Q.E.D.*

LEMMA A7: *If Assumptions 1, 3, and 4 are satisfied and $\varepsilon_n^M + \varepsilon_n^G(1 + B_n) + r_0 = o(r_L)$ then $|G(\hat{\rho}_L - \rho_L)|_\infty = O_p(r_L)$.*

Proof: The population and sample Lasso first order conditions give

$$|M - G\rho_L|_\infty \le r_0, |\hat{M} - \hat{G}\hat{\rho}_L|_\infty \le r_L.$$

Then by the triangle and Holder inequalities we have

$$|G(\hat{\rho}_L - \rho_L)|_\infty \le |(G - \hat{G})\hat{\rho}_L|_\infty + |\hat{G}\hat{\rho}_L - \hat{M}|_\infty + |\hat{M} - M|_\infty + |M - G\rho_L|_\infty$$

$$\le |G - \hat{G}|_\infty |\hat{\rho}_L|_1 + |\hat{G}\hat{\rho}_L - \hat{M}|_\infty + |\hat{M} - M|_\infty + |M - G\rho_L|_\infty$$

$$= O_p(\varepsilon_n^G(1 + B_n) + r_L + \varepsilon_n^M + r_0)$$

$$= O_p(r_L). \ Q.E.D.$$

**Proof of Theorem 2:** Note that by Lemma 1,

$$\|\hat{\alpha}_L - \alpha_0\|^2 \le 2\|\hat{\alpha}_L - \alpha_L\|^2 + 2\|\alpha_L - \alpha_0\|^2 \tag{6.2}$$

$$\le 2(\hat{\rho}_L - \rho_L)'G(\hat{\rho}_L - \rho_L) + 4r_0(1 + B_n)$$

By Lemmas A5 and A7,

$$(\hat{\rho}_L - \rho_L)'G(\hat{\rho}_L - \rho_L) \leq |(\hat{\rho}_L - \rho_L)'G|_\infty(|\hat{\rho}_L|_1 + |\rho_L|_1)$$
$$= O_p(r_L(1 + B_n)).$$

The first conclusion then follows by eq. (6.2) and $r_0 = o(r_L)$.

If Assumption 2 is satisfied let $\hat{\delta} = \hat{\rho}_L - \rho_L$. Then by Lemma A6 with probability approaching one

$$|\hat{\delta}|_1^2 = (\sum_{j \in A_L^c} |\hat{\delta}_j| + \sum_{j \in A_L} |\hat{\delta}_j|)^2 \leq (4 \sum_{j \in A_L} |\hat{\delta}_j|)^2 \leq 16\mathcal{M}_L \sum_{j \in A_L} |\hat{\delta}_j|^2 \leq \frac{16\mathcal{M}_L}{k_3}\hat{\delta}'G\hat{\delta}$$
$$\leq \frac{16\mathcal{M}_L}{k_3}|G\hat{\delta}|_\infty|\hat{\delta}|_1 = O_p(\mathcal{M}_L r_L)|\hat{\delta}|_1.$$

Dividing through by $|\hat{\delta}|_1$ then gives

$$|\hat{\delta}|_1 = O_p(\mathcal{M}_L r_L).$$

It follows that

$$\hat{\delta}'G\hat{\delta} \leq |G\hat{\delta}|_\infty|\hat{\delta}|_1 = O_p(\mathcal{M}_L r_L^2).$$

The conclusion then follows from eq. (6.2) $Q.E.D.$

The following result is useful for the proof of Theorem 3:

LEMMA A8: *If* $\varepsilon_n^M + \varepsilon_n^G(1 + B_n) + \lambda_0 = o(\lambda_D)$ *then with probability approaching one* $|\hat{\rho}_D|_1 \leq |\rho_D|_1$ *and* $|G(\hat{\rho}_D - \rho_D)|_\infty = O_p(\lambda_D)$.

Proof: By $\rho_D$ satisfying the constraint given in the definition of $\rho_D$, Lemma 1, and the triangle and Holder inequalities

$$|\hat{M} - \hat{G}\rho_D|_\infty \leq |\hat{M} - M|_\infty + |M - G\rho_D|_\infty + |(G - \hat{G})\rho_D|_\infty$$
$$\leq |\hat{M} - M|_\infty + |M - G\rho_D|_\infty + |G - \hat{G}|_\infty|\rho_D|_1$$
$$= O_p(\varepsilon_n^M + \lambda_0 + \varepsilon_n^G(1 + B_n))$$
$$= o_p(\lambda_D).$$

It follows that with probability approaching one $|\hat{M} - \hat{G}\rho_D|_\infty \leq \lambda_D$, so that $\rho_D$ is feasible for the sample Dantzig constrained minimization problem, implying the first conclusion.

To show the second conclusion note that by the first conclusion, with probability approaching one,

$$|G(\hat{\rho}_D - \rho_D)|_\infty \leq |(G - \hat{G})\hat{\rho}_D|_\infty + |\hat{G}\hat{\rho}_D - \hat{M}|_\infty + |\hat{M} - M|_\infty + |M - G\rho_D|_\infty$$
$$\leq |G - \hat{G}|_\infty|\hat{\rho}_D|_1 + |\hat{G}\hat{\rho}_D - \hat{M}|_\infty + |\hat{M} - M|_\infty + |M - G\rho_D|_\infty$$
$$= O_p(\varepsilon_n^G(1 + B_n) + \lambda_D + \varepsilon_n^M + \lambda_0)$$
$$= O_p(\lambda_D). \ Q.E.D.$$

**Proof of Theorem 3:** For this proof let $\hat{\delta} = \hat{\rho}_D - \rho_D$. By Lemma 1, for $\alpha_D = b'\rho_D$ we have

$$\|\hat{\alpha}_D - \alpha_0\|^2 \leq 2\|\hat{\alpha}_D - \alpha_D\|^2 + 2\|\alpha_D - \alpha_0\|^2 \tag{6.3}$$
$$\leq 2(\hat{\rho}_D - \rho_D)'G(\hat{\rho}_D - \rho_D) + 24\lambda_0(1 + B_n)$$

By Lemma A8

$$(\hat{\rho}_D - \rho_D)'G(\hat{\rho}_D - \rho_D) \leq |(\hat{\rho}_D - \rho_D)'G|_\infty(|\hat{\rho}_D|_1 + |\rho_D|_1)$$
$$= O_p(\lambda_D(1 + B_n)),$$

giving the slow rate by $\lambda_0 = o(\lambda_D)$. Also by the first conclusion of Lemma A8 $|\rho_D + \hat{\delta}|_1 = |\hat{\rho}_D|_1 \leq |\rho_D|_1$ with probability approaching one, so that by the definition of $s_D$ and the second conclusion of Lemma A8,

$$|\hat{\delta}|_1^2 \leq s_D\hat{\delta}'G\hat{\delta} \leq s_D|G\hat{\delta}|_\infty|\hat{\delta}|_1 = O_p(s_D\lambda_D)|\hat{\delta}|_1.$$

Dividing through by $|\hat{\delta}|_1$ then gives

$$|\hat{\delta}|_1 = O_p(s_D\lambda_D).$$

It then follows from the second conclusion of Lemma A8 that

$$\hat{\delta}'G\hat{\delta} \leq |G\hat{\delta}|_\infty|\hat{\delta}|_1 = O_p(s_D\lambda_D^2).$$

The fast rate then follows from eq. (6.3). *Q.E.D.*

**Proof of Lemma 4:** Define

$$T_{ij} = m(W_i, b_j) - E[m(W_i, b_j)], \ U_j = \frac{1}{n}\sum_{i=1}^n T_{ij}.$$

For any constant $C$,

$$\Pr(|\hat{M} - M|_\infty \geq C\varepsilon_n^M) \leq \sum_{j=1}^p \mathbb{P}(|U_j| > C\varepsilon_n^M) \leq p \cdot \max_j \mathbb{P}(|U_j| > C\varepsilon_n^M).$$

Note that $E[T_{ij}] = 0$ and

$$|T_{ij}| \leq |m(W_i, b_j)| + E[|m(W_i, b_j)|] \leq B_n^M\{A(W_i) + E[|A(W_i)|]\}.$$

Define $C_A = \|A(W_i)\|_{\Psi_2} + E[|A(W_i)|]$ and let $K = \|T_{ij}\|_{\Psi_2} \leq C_A B_n^M$. By Hoeffding's inequality there is a constant $c$ such that

$$p \cdot \max_j \mathbb{P}(|U_j| > C\varepsilon_n^M) \leq 2p\exp\left(-\frac{cn(C\varepsilon_n^M)^2}{K^2}\right) \leq 2p\exp\left(-\frac{cn(C\varepsilon_n^M)^2}{C_A^2(B_n^M)^2}\right)$$
$$\leq 2\exp\left(\ln(p)[1 - \frac{cC^2}{C_A^2}]\right) \longrightarrow 0,$$

24

for any $C > C_A/\sqrt{c}$. Thus for large enough $C$, $\Pr(|\hat{M} - M|_\infty \geq C\varepsilon_n^M) \longrightarrow 0$, implying the conclusion. $Q.E.D.$

**Proof of Theorem 5:** $\sqrt{n}\mathcal{M}_L r_L \geq C > 0$ since $\sqrt{n}\varepsilon_n^\alpha \to \infty$ and $\varepsilon_n^\alpha = o(r_L)$. By hypothesis, $\sqrt{n}\sqrt{\mathcal{M}_L r_L^2 + (1 + B_n)r_0}\varepsilon_n^\gamma \longrightarrow 0$ so that $\sqrt{n}\mathcal{M}_L r_L \varepsilon_n^\gamma \longrightarrow 0$ and hence $\varepsilon_n^\gamma \longrightarrow 0$. Then Assumption 5 implies that

$$\int [m(W, \hat{\gamma}) - m(W, \gamma_0)]^2 F_0(dW) \overset{p}{\longrightarrow} 0.$$

By $\sqrt{\mathcal{M}_L} r_L \longrightarrow 0$ and $(1 + B_n)r_0 \to 0$ and the fast bound in Theorem 2, we have $\|\hat{\alpha}_L - \alpha_0\| \overset{p}{\longrightarrow} 0$. The same argument appealing to the fast bound in Theorem 3 shows $\|\hat{\alpha}_D - \alpha_0\| \overset{p}{\longrightarrow} 0$.

Furthermore, by Assumption 8 $\sqrt{n}\|\hat{\alpha}_L - \alpha_0\|\|\hat{\gamma} - \gamma_0\| \overset{p}{\longrightarrow} 0$ or $\sqrt{n}\|\hat{\alpha}_D - \alpha_0\|\|\hat{\gamma} - \gamma_0\| \overset{p}{\longrightarrow} 0$. Then the first conclusion follows by Theorem 13 of Chernozhukov et al. (2016).

To prove the second conclusion let $\psi_i = \psi_0(W_i)$ and $\varepsilon_i = Y_i - \gamma_0(X_i)$. Then

$$(\hat{\psi}_i - \psi_i)^2 \leq 4(R_{i1} + R_{i2} + R_{i3} + R_{i4})$$
$$R_{i1} = [m(W_i, \hat{\gamma}) - m(W_i, \gamma_0)]^2, R_{i2} = \alpha_0(X_i)^2\{\hat{\gamma}(X_i) - \gamma_0(X_i)\}^2,$$
$$R_{i3} = \{\hat{\alpha}_\Delta(X_i) - \alpha_0(X_i)\}^2\{Y_i - \gamma_0(X_i)\}^2, R_{i4} = \{\hat{\alpha}_\Delta(X_i) - \alpha_0(X_i)\}^2\{\hat{\gamma}(X_i) - \gamma_0(X_i)\}^2,$$

where for notational convenience we drop the $\ell$ subscript on $\hat{\gamma}_\ell$ and $\hat{\alpha}_{\ell\Delta}$. Let $E_{-\ell}[\cdot]$ denote the expectation conditional on the subvector of the data where $i \notin I_\ell$. Then for $i \in I_\ell$, by Assumption 5

$$E_{-\ell}[R_{i1}] = \int [m(W, \hat{\gamma}) - m(W, \gamma_0)]^2 F_0(dW) \overset{p}{\longrightarrow} 0, \ E_{-\ell}[R_{i2}] \leq C\|\hat{\gamma} - \gamma_0\|^2 \overset{p}{\longrightarrow} 0,$$

$$E_{-\ell}[R_{i3}] = E_{-\ell}[\{\hat{\alpha}_\Delta(X) - \alpha_0(X)\}^2 Var(Y|X)] \leq CE_{-\ell}[\{\hat{\alpha}_\Delta(X) - \alpha_0(X)\}^2]$$
$$\leq C\|\hat{\alpha} - \alpha_0\|^2 \overset{p}{\longrightarrow} 0,$$

where the last inequality follows by $\alpha_0(x)$ bounded which implies $\|\hat{\alpha}_\Delta - \alpha_0\| \leq \|\hat{\alpha} - \alpha_0\|$ for large enough $\Delta$. Also by $\alpha_0(x)$ bounded, for large enough $n$ we have $|\alpha_0(x)| \leq \Delta$, so that

$$R_{i4} \leq 4\Delta^2\{\hat{\gamma}(X) - \gamma_0(X)\}^2.$$

Then

$$E_{-\ell}[R_{i4}] = 4\Delta^2 E_{-\ell}[\{\hat{\gamma}(X) - \gamma_0(X)\}^2] = 4\Delta^2\|\hat{\gamma} - \gamma_0\|^2 \overset{p}{\longrightarrow} 0.$$

It then follows that

$$E_{-\ell}[\frac{1}{n}\sum_{i \in I_\ell}(\hat{\psi}_i - \psi_i)^2] \leq 4\sum_{j=1}^4 \frac{1}{n}\sum_{i \in I_\ell}E_{-\ell}[R_{ij}] \overset{p}{\longrightarrow} 0.$$

It then follows by the triangle and conditional Markov inequalities and summing over $\ell$ that

$$\frac{1}{n}\sum_{i=1}^{n}(\hat{\psi}_i - \psi_i)^2 \xrightarrow{p} 0.$$

Then $\hat{V} \xrightarrow{p} 0$ follows by the law of large numbers in the usual way. $Q.E.D.$

**Proof of Theorem 6:** Note that

$$m(w, b_j) = 1(j \le p/2)q_j(z) - 1(j > p/2)q_{j-p/2}(z).$$

Therefore

$$\max_{1\le j\le p}|m(W, b_j)| \le \max_{1\le j\le p/2}|q_j(Z)| \le B_n^q$$

It then follows by hypothesis ii) of the statement of Theorem 6 that Assumption 3 is satisfied with $B_n^b = B_n^q$, Assumption 6 is satisfied with $A(W) = 1$ and $B_n^m = B_n^q$, and Assumption 8 is satisfied.

Next, it also follows by hypothesis i) and the form of $\alpha_0(x)$ that $Var(Y|X)$ and $\alpha_0(x)$ are bounded. In addition, by iterated expectations,

$$E[\gamma_0(1, Z)^2] = E[\frac{D}{\pi_0(Z)}\gamma_0(1, Z)^2] = E[\frac{D}{\pi_0(Z)}\gamma_0(X)^2] \le CE[\gamma_0(X)^2] < \infty,$$

$$E[\{\gamma(1, Z) - \gamma_0(1, Z)\}^2] = E[\frac{D}{\pi_0(Z)}\{\gamma(1, Z) - \gamma_0(1, Z)\}^2] = E[\frac{D}{\pi_0(Z)}\{\gamma(X) - \gamma_0(X)\}^2]$$

$$\le C\|\gamma - \gamma_0\|^2.$$

Combining these inequalities with the analogous inequalities for $\gamma(0, z)$ it follows that Assumption 5 is satisfied. The conclusion then follows by Theorem 5. $Q.E.D.$

**Proof of Theorem 7:** Note that

$$m(w, b_j) - m(w, 0) = b_j(t(x)).$$

By $\alpha_0(x)$ bounded, the distribution of $t(X)$ is absolutely continuous with respect to the distribution of $X$ so that by Assumption 3,

$$\max_{1\le j\le p}|m(W, b_j)| \le \max_{1\le j\le p}|b_j(t(X))| \le B_n^b$$

It then follows by hypothesis ii) of the statement of Theorem 7 that Assumption 6 is satisfied with $A(W) = 1$ and $B_n^m = B_n^b$.

Next, it also follows by hypothesis i) that $Var(Y|X)$ and $\alpha_0(x)$ are bounded. In addition, by iterated expectations,

$$E[m(W, \gamma_0)^2] \leq CE[\gamma_0(t(X))^2] + C = C \int \frac{f_t(x)}{f_0(x)} \gamma_0(x)^2 f_0(x) dx + C$$

$$\leq CE[\gamma_0(X)^2] + C < \infty,$$

$$E[\{m(W, \gamma) - m(W, \gamma_0)\}^2] = E[\{\gamma(t(X)) - \gamma_0(t(X))\}^2] = \int \frac{f_t(x)}{f_0(x)} \{\gamma(x) - \gamma_0(x)\}^2 f_0(x) dx \leq C\|\gamma - \gamma_0\|^2.$$

Thus we see that Assumption 5 is satisfied. The conclusion then follows by Theorem 5. $Q.E.D.$

**Proof of Theorem 8:** We have $m(w, \gamma) = w_1[y - \gamma(x)]$ so that

$$m(w, b_j) - m(w, 0) = -w_1 b_j(x).$$

Therefore by Assumption 3,

$$\max_{1 \leq j \leq p} |m(W, b_j) - m(W, 0)| \leq |W_1| \max_{1 \leq j \leq p} |b_j(X)| \leq B_n^b |W_1|$$

It then follows by hypothesis ii) of the statement of Theorem 8 that Assumption 6 is satisfied with $A(W) = |W_1|$ and $B_n^m = B_n^b$.

Next, it also follows by hypothesis i) that $Var(Y|X)$ and $\alpha_0(x) = -E[W_1|x]$ are bounded. In addition, by $W_1$ and $Y$ having fourth moments,

$$E[m(W, \gamma_0)^2] \leq CE[W_1^2 \gamma_0(X)^2] + C < \infty,$$

$$E[\{m(W, \gamma) - m(W, \gamma_0)\}^2] = E[E[W_1^2|X]\{\gamma(X) - \gamma_0(X)\}^2] \leq C\|\gamma - \gamma_0\|^2.$$

Thus we see that Assumption 5 is satisfied. The conclusion then follows by Theorem 5. $Q.E.D.$

**Proof of Lemma 9:** Define

$$\hat{M}_\ell = (\hat{M}_{\ell 1}, ..., \hat{M}_{\ell p})', \ \hat{M}_{\ell j} = \left(\frac{1}{n - n_\ell}\right) \sum_{\tilde{\ell} \neq \ell} \sum_{i \in I_{\tilde{\ell}}} D(W_i, b_j, \hat{\gamma}_{\ell, \tilde{\ell}}),$$

$$\bar{M}(\gamma) = (\bar{M}_1(\gamma), ..., \bar{M}_p(\gamma))', \ \bar{M}_j(\gamma) = \int D(W, b_j, \gamma) F_0(dW).$$

Note that $M = \bar{M}(\gamma_0)$. Let $\Gamma_{\ell, \tilde{\ell}}$ be the event that $\|\hat{\gamma}_{\ell, \tilde{\ell}} - \gamma_0\| < \varepsilon$ and note that $\Pr(\Gamma_{\ell, \tilde{\ell}}) \longrightarrow 1$ for each $\ell$ and $\tilde{\ell}$. When $\Gamma_{\ell, \tilde{\ell}}$ occurs,

$$\max_j |D(W_i, b_j, \hat{\gamma}_{\ell, \tilde{\ell}})| \leq B_n^D A(W_i)$$

by Assumption 9. Define

$$T_{ij}(\gamma) = D(W_i, b_j, \gamma) - \bar{M}_j(\gamma), \ (i \in I_{\tilde{\ell}}), \ U_{\tilde{\ell} j}(\gamma) = \frac{1}{n_{\tilde{\ell}}} \sum_{i \in I_{\tilde{\ell}}} T_{ij}(\gamma).$$

27

Note that for any constant $C$ and the event $\mathcal{A} = \{\max_j |U_{\tilde{\ell}j}(\hat{\gamma}_{\ell,\tilde{\ell}})| \geq C\varepsilon_n^D\}$

$$\Pr(\mathcal{A}) = \Pr(\mathcal{A}|\Gamma_{\ell,\tilde{\ell}})\Pr(\Gamma_{\ell,\tilde{\ell}}) + \Pr(\mathcal{A}|\Gamma_{\ell,\tilde{\ell}}^c)\left[1 - \Pr(\Gamma_{\ell,\tilde{\ell}})\right]$$

$$\leq \Pr(\max_j |U_{\tilde{\ell}j}(\hat{\gamma}_{\ell,\tilde{\ell}})| \geq C\varepsilon_n^D|\Gamma_{\ell,\tilde{\ell}}) + 1 - \Pr(\Gamma_{\ell,\tilde{\ell}}).$$

Also

$$\Pr(\max_j |U_{\tilde{\ell}j}(\hat{\gamma}_{\ell,\tilde{\ell}})| \geq C\varepsilon_n^D|\Gamma_{\ell,\tilde{\ell}}) \leq p \cdot \max_j \Pr(|U_{\tilde{\ell}j}(\hat{\gamma}_{\ell,\tilde{\ell}})| > C\varepsilon_n^D|\Gamma_{\ell,\tilde{\ell}}).$$

Note that $E[T_{ij}(\hat{\gamma}_{\ell,\tilde{\ell}})|\hat{\gamma}_{\ell,\tilde{\ell}}] = 0$ for $i \in I_{\tilde{\ell}}$. Also, conditional on the event $\Gamma_{\ell,\tilde{\ell}}$,

$$|T_{ij}(\hat{\gamma}_{\ell,\tilde{\ell}})| \leq B_n^D\{A(W_i) + E[|A(W_i)|]\}, \ i \in I_{\tilde{\ell}}.$$

Define $C_A = \|A(W_i)\|_{\Psi_2} + E[|A(W_i)|]$ and let $K(\hat{\gamma}_{\ell,\tilde{\ell}}) = \|T_{ij}(\hat{\gamma}_{\ell,\tilde{\ell}})\|_{\Psi_2} \leq CB_n^D, \ i \in I_{\tilde{\ell}}$. By Hoeffding's inequality and the independence of $(W_i)_{i \in I_{\tilde{\ell}}}$ and $\hat{\gamma}_{\ell,\tilde{\ell}}$ there is a constant $c$ such that

$$p \cdot \max_j \Pr(|U_{\tilde{\ell}j}(\hat{\gamma}_{\ell,\tilde{\ell}})| > C\varepsilon_n^D|\Gamma_{\ell,\tilde{\ell}}) = p \cdot \max_j E[\Pr(|U_{\tilde{\ell}j}(\hat{\gamma}_{\ell,\tilde{\ell}})| > C\varepsilon_n^D|\hat{\gamma}_{\ell,\tilde{\ell}})|\Gamma_{\ell,\tilde{\ell}}]$$

$$\leq 2pE[\exp\left(-\frac{cn(C\varepsilon_n^D)^2}{K(\hat{\gamma}_{\ell,\tilde{\ell}})^2}\right)|\Gamma_{\ell,\tilde{\ell}}] \leq 2p\exp\left(-\frac{cn(C\varepsilon_n^D)^2}{C_A^2(B_n^D)^2}\right)$$

$$\leq 2\exp\left(\ln(p)[1 - \frac{cC^2}{C_A^2}]\right) \longrightarrow 0,$$

for any $C > C_A/\sqrt{c}$. Let $U_{\tilde{\ell}}(\gamma) = (U_{\tilde{\ell}1}(\gamma), ..., U_{\tilde{\ell}p}(\gamma))'$. It then follows from above that for large $C$, $\Pr(|U_{\tilde{\ell}}(\hat{\gamma}_{\ell,\tilde{\ell}})|_\infty \geq C\varepsilon_n^D) \longrightarrow 0$. Therefore $|U_{\tilde{\ell}}(\hat{\gamma}_{\ell,\tilde{\ell}})|_\infty = O_p(\varepsilon_n^D)$.

Next, for each $\ell$,

$$\left|\hat{M}_\ell - \sum_{\tilde{\ell} \neq \ell} \frac{n_{\tilde{\ell}}}{n - n_\ell} \bar{M}(\hat{\gamma}_{\ell,\tilde{\ell}})\right|_\infty = \left|\sum_{\tilde{\ell} \neq \ell} \frac{n_{\tilde{\ell}}}{n - n_\ell} U_{\tilde{\ell}}(\hat{\gamma}_{\ell,\tilde{\ell}})\right|_\infty \leq \sum_{\tilde{\ell} \neq \ell} \frac{n_{\tilde{\ell}}}{n - n_\ell} |U_{\tilde{\ell}}(\hat{\gamma}_{\ell,\tilde{\ell}})|_\infty = O_p(\varepsilon_n^D).$$

Also by Assumption 9 ii) and the fact that $\Pr(\Gamma_{\ell,\tilde{\ell}}) \longrightarrow 1$ for each $\ell$ and $\tilde{\ell}$

$$\left|\sum_{\tilde{\ell} \neq \ell} \frac{n_{\tilde{\ell}}}{n - n_\ell} \bar{M}(\hat{\gamma}_{\ell,\tilde{\ell}}) - M\right|_\infty = \left|\sum_{\tilde{\ell} \neq \ell} \frac{n_{\tilde{\ell}}}{n - n_\ell} [\bar{M}(\hat{\gamma}_{\ell,\tilde{\ell}}) - M]\right|_\infty \leq B_n^\Delta \sum_{\tilde{\ell} \neq \ell} \frac{n_{\tilde{\ell}}}{n - n_\ell} \|\hat{\gamma}_{\ell,\tilde{\ell}} - \gamma_0\|$$

$$= O_p(B_n^\Delta \varepsilon_n^\gamma).$$

The conclusion then follows by the triangle inequality. *Q.E.D.*

**Proof of Theorem 10:** The first conclusion follows exactly as in the proof of Theorem 5. We prove the second conclusion by verifying the conditions of Lemma 14 of Chernozhukov et

al. (2016). Let $\lambda$ in Chernozhukov et al. (2016) be $\alpha$ here and $\phi(w, \gamma, \lambda)$ in Chernozhukov et al. (2016) be $\lambda(x)[y - \gamma(x)]$. By Assumption 5, $\varepsilon_n^\gamma \longrightarrow 0$, and $\varepsilon_n^\alpha \longrightarrow 0$ it follows that

$$\int [\phi(W, \hat{\gamma}, \lambda_0) - \phi(W, \gamma_0, \lambda_0)]^2 F_0(dW) = \int \lambda_0(X)^2 [\hat{\gamma}(X) - \gamma_0(X)]^2 F_0(dW) \leq C\|\hat{\gamma} - \gamma_0\|^2 \xrightarrow{p} 0.$$

$$\int [\phi(W, \gamma_0, \hat{\lambda}) - \phi(W, \gamma_0, \lambda_0)]^2 F_0(dW) = \int [\hat{\lambda}(X) - \lambda_0(X)]^2 [Y - \gamma_0(X)]^2 F_0(dW)$$

$$= \int [\hat{\lambda}(X) - \lambda_0(X)]^2 Var(Y|X) F_0(dX) \leq C\|\hat{\lambda} - \lambda_0\|^2 \xrightarrow{p} 0.$$

Also by Assumption 5 $\int [m(W, \hat{\gamma}) - m(W, \gamma_0)]^2 F_0(dW) \xrightarrow{p} 0$, so all the conditions of Assumption 4 of Chernozhukov et al. (2018b) are satisfied.

Also let $\Delta_n^\alpha = \sqrt{\mathcal{M}_L r_L^2 + (1 + B_n) r_0}$ for Lasso and $\Delta_n^\alpha = \sqrt{s_D \lambda_D^2 + (1 + B_n) \lambda_0}$ for Dantzig. Then by the Cauchy-Schwartz inequality,

$$\sqrt{n} \int |\phi(W, \hat{\gamma}_\ell, \hat{\lambda}_\ell) - \phi(W, \gamma_0, \hat{\lambda}_\ell) - \phi(W, \hat{\gamma}_\ell, \lambda_0) + \phi(W, \gamma_0, \lambda_0)| F_0(dW)$$

$$= \sqrt{n} \int |\hat{\alpha}_\ell(X) - \alpha_0(X)| |\hat{\gamma}_\ell(X) - \gamma_0(X)| F_0(dW) \leq \sqrt{n} \|\hat{\alpha}_\ell - \alpha_0\| \|\hat{\gamma}_\ell - \gamma_0\|$$

$$= O_p(\sqrt{n} \Delta_n^\alpha \varepsilon_n^\gamma) \xrightarrow{p} 0.$$

Therefore Assumption 5 of Chernozhukov et al. (2016) is satisfied.

Also, we have by Assumption 10

$$\sqrt{n} \left| \int [m(W, \hat{\gamma}_\ell) - m(W, \gamma_0) + \alpha_0(X)\{Y - \hat{\gamma}_\ell(X)\}] F_0(dW) \right|$$

$$= \sqrt{n} \left| \int [m(W, \hat{\gamma}_\ell) - m(W, \gamma_0) + \alpha_0(X)\{\gamma_0(X) - \hat{\gamma}_\ell(X)\}] F_0(dW) \right|$$

$$= \sqrt{n} \left| \int [m(W, \hat{\gamma}_\ell) - m(W, \gamma_0) - D(W, \hat{\gamma}_\ell - \gamma_0, \gamma_0)] F_0(dW) \right|$$

$$\leq C\sqrt{n} \|\hat{\gamma}_\ell - \gamma_0\|^2 = C\sqrt{n} o_p(1/\sqrt{n}) \xrightarrow{p} 0.$$

Also,

$$\sqrt{n} \left| \int \hat{\alpha}(X)\{Y - \gamma_0(X)\}] F_0(dW) \right| = 0.$$

Therefore Assumption 6 of Chernozhukov et al. (2016) is satisfied, so the first conclusion follows by Lemma 14 of Chernozhukov et al. (2016). *Q.E.D.*

# 7    References

Athey, S., G. Imbens, and S. Wager (2018): "Approximate Residual Balancing: Debiased Inference of Average Treatment Effects in High Dimensions," *Journal of the Royal Statistical Society, Series B* 80, 597–623.

Belloni, A., V. Chernozhukov, and C. Hansen (2014): "Inference on Treatment Effects after Selection among High-Dimensional Controls," *Review of Economic Studies* 81, 608–650.

Belloni, A., V. Chernozhukov, and K. Kato (2015): "Uniform Post Selection Inference for Least Absolute Deviation Regression and Other $Z$-Estimation Problems," *Biometrika*, 102: 77–94. ArXiv, 2013.

Belloni, A., V. Chernozhukov, L. Wang (2014): "Pivotal Estimation via Square-Root Lasso in Nonparametric Regression," *Annals of Statistics* 42, 757–788.

Bickel, P.J. (1982): "On Adaptive Estimation," *Annals of Statistics* 10, 647–671.

Bickel, P.J. and Y. Ritov (1988): "Estimating Integrated Squared Density Derivatives: Sharp Best Order of Convergence Estimates," *Sankhyā: The Indian Journal of Statistics, Series A* 238, 381–393.

Bickel, P.J., C.A.J. Klaassen, Y. Ritov and J.A. Wellner (1993): *Efficient and Adaptive Estimation for Semiparametric Models*, Baltimore: Johns Hopkins University Press.

Bickel, P.J., Y.Ritov, and A.Tsybakov (2009): "Simultaneous Analysis of Lasso and Dantzig Selector," *Annals of Statistics* 37, 1705–1732.

Bradic, J. and M. Kolar (2017): "Uniform Inference for High-Dimensional Quantile Regression: Linear Functionals and Regression Rank Scores," *arXiv preprint arXiv:1702.06209*.

Cai, T.T. and Z. Guo (2017): "Confidence Intervals for High-Dimensional Linear Regression: Minimax Rates and Adaptivity," *Annals of Statistics* 45, 615-646.

Candes, E. and T. Tao (2007): "The Dantzig Selector: Statistical Estimation when $p$ is much Larger than $n$," *Annals of Statistics* 35, 2313–2351.

Chernozhukov, V., D. Chetverikov, and K. Kato (2013): "Gaussian Approximations and Multiplier Bootstrap for Maxima of Sums of High-Dimensional Random Vectors," *Annals of Statistics* 41, 2786–2819.

Chernozhkov, V., C. Hansen, and M. Spindler (2015): "Valid Post-Selection and Post-Regularization Inference: An Elementary, General Approach," *Annual Review of Economics* 7, 649–688.

Chernozhukov, V., D. Chetverikov, M. Demirer, E. Duflo, C. Hansen, W. Newey and J. Robins (2018): "Debiased/Double Machine Learning for Treatment and Structural Parameters," *Econometrics Journal* 21, C1-C68.

Chernozhukov, V., J. C. Escanciano, H. Ichimura, W.K. Newey, and J. Robins (2016): "Locally Robust Semiparametric Estimation," arXiv preprint arXiv:1608.00033.

Chernozhukov, V., W.K. Newey, and J. Robins (2018): "Double/De-Biased Machine Learning Using Regularized Riesz Representers," arXiv.

Chernozhukov, V., J.A. Hausman, and W.K. Newey (2018): "Demand Analysis with Many Prices," forthcoming.

Farrell, M. (2015): "Robust Inference on Average Treatment Effects with Possibly More Covariates than Observations," *Journal of Econometrics* 189, 1–23.

Hasminskii, R.Z. and I.A. Ibragimov (1979): "On the Nonparametric Estimation of Functionals," in P. Mandl and M. Huskova (eds.), *Proceedings of the 2nd Prague Symposium on Asymptotic Statistics, 21-25 August 1978*, Amsterdam: North-Holland, pp. 41-51.

Hausman, J.A. and W.K. Newey (2016): "Individual Heterogeneity and Average Welfare," *Econometrica* 84, 1225–1248.

Hirshberg, D.A. and S. Wager (2018): "Augmented Minimax Linear Estimation," arXiv.

Jankova, J. and S. Van De Geer (2015): "Confidence Intervals for High-Dimensional Inverse Covariance Estimation," *Electronic Journal of Statistics* 90, 1205–1229.

Jankova, J. and S. Van De Geer (2016a): "Semi-Parametric Efficiency Bounds and Efficient Estimation for High-Dimensional Models," arXiv preprint arXiv:1601.00815.

Jankova, J. and S. Van De Geer (2016b): "Confidence Regions for High-Dimensional Generalized Linear Models under Sparsity," arXiv preprint arXiv:1610.01353.

Javanmard, A. and A. Montanari (2014a): "Hypothesis Testing in High-Dimensional Regression under the Gaussian Random Design Model: Asymptotic Theory," *IEEE Transactions on Information Theory* 60, 6522–6554.

Javanmard, A. and A. Montanari (2014b): "Confidence Intervals and Hypothesis Testing for High-Dimensional Regression," *Journal of Machine Learning Research* 15: 2869–2909.

Javanmard, A. and A. Montanari (2015): "De-Biasing the Lasso: Optimal Sample Size for Gaussian Designs," arXiv preprint arXiv:1508.02757.

Jing, B.Y., Q.M. Shao, and Q. Wang (2003): "Self-Normalized Cramér-Type Large Deviations for Independent Random Variables," *Annals of Probability* 31, 2167–2215.

Luedtke, A. R. and M. J. van der Laan (2016): "Optimal Individualized Treatments in Resource-limited Settings," *The International Journal of Biostatistics* 12, 283-303.

Newey, W.K. (1994): "The Asymptotic Variance of Semiparametric Estimators," *Econometrica* 62, 1349–1382.

Newey, W.K., F. Hsieh, and J.M. Robins (1998): "Undersmoothing and Bias Corrected Functional Estimation," MIT Dept. of Economics working paper 98-17.

Newey, W.K., F. Hsieh, and J.M. Robins (2004): "Twicing Kernels and a Small Bias Property of Semiparametric Estimators," *Econometrica* 72, 947–962.

Newey, W.K. and J.M. Robins (2017): "Cross Fitting and Fast Remainder Rates for Semiparametric Estimation," arxiv.

Neykov, M., Y. Ning, J.S. Liu, and H. Liu (2015): "A Unified Theory of Confidence Regions and Testing for High Dimensional Estimating Equations," arXiv preprint arXiv:1510.08986.

Ning, Y. and H. Liu (2017): "A General Theory of Hypothesis Tests and Confidence Regions for Sparse High Dimensional Models," *Annals of Statistics* 45, 158-195.

Ren, Z., T. Sun, C.H. Zhang, and H. Zhou (2015): "Asymptotic Normality and Optimalities in Estimation of Large Gaussian Graphical Models," *Annals of Statistics* 43, 991–1026.

Robins, J.M. and A. Rotnitzky (1995): "Semiparametric Efficiency in Multivariate Regression Models with Missing Data," *Journal of the American Statistical Association* 90 (429): 122–129.

Robins, J.M., A. Rotnitzky, and L.P. Zhao (1995): "Analysis of Semiparametric Regression Models for Repeated Outcomes in the Presence of Missing Data," *Journal of the American Statistical Association* 90, 106–121.

Robins, J.M., M. Sued, Q. Lei-Gomez, and A. Rotnitzky (2007): "Comment: Performance of Double-Robust Estimators When 'Inverse Probability' Weights Are Highly Variable," *Statistical Science* 22, 544–559.

Robins, J.M., L. Li, E. Tchetgen, and A. van der Vaart (2008): "Higher Order Influence Functions and Minimax Estimation of Nonlinear Functionals," *IMS Collections Probability and Statistics: Essays in Honor of David A. Freedman, Vol 2,* 335-421.

Robins, J., P. Zhang, R. Ayyagari, R. Logan, E. Tchetgen, L. Li, A. Lumley, and A. van der Vaart (2013): "New Statistical Approaches to Semiparametric Regression with Application to Air Pollution Research," Research Report Health E Inst..

Rosenbaum, P.R. and D. B. Rubin (1983): "The Central Role of the Propensity Score in Observational Studies for Causal Effects," *Biometrika* 70: 41–55.

Schick, A. (1986): "On Asymptotically Efficient Estimation in Semiparametric Models," *Annals of Statistics* 14, 1139–1151.

Stock, J.H. (1989): "Nonparametric Policy Analysis," *Journal of the American Statistical Association* 84, 567–575.

Toth, B. and M. J. van der Laan (2016), "TMLE for Marginal Structural Models Based On An Instrument," U.C. Berkeley Division of Biostatistics Working Paper Series, Working Paper 350.

Tsybakov, A.B. (2009): *Introduction to Nonparametric Estimation.* New York: Springer.

Van De Geer, S., P. Bühlmann, Y. Ritov, and R. Dezeure (2014): "On Asymptotically Optimal Confidence Regions and Tests for High-Dimensional Models," *Annals of Statistics*, 42: 1166–1202.

Van der Laan, M. and D. Rubin (2006): "Targeted Maximum Likelihood Learning," *International Journal of Biostatistics* 2.

Van der Laan, M. J. and S. Rose (2011): *Targeted Learning: Causal Inference for Observa-*

*tional and Experimental Data,* Springer.

Van der Vaart, A.W. (1991): "On Differentiable Functionals," *Annals of Statistics*, 19: 178–204.

Van der Vaart, A.W. (1998): *Asymptotic Statistics.* New York: Cambridge University Press.

Van der Vaart, A.W. and J.A. Wellner (1996): *Weak Convergence and Empirical Processes*, New York: Springer.

Vershynin, R. (2018): *High-Dimensional Probability*, New York: Cambridge University Press.

Zhang, C. and S. Zhang (2014): "Confidence Intervals for Low-Dimensional Parameters in High-Dimensional Linear Models," *Journal of the Royal Statistical Society, Series B* 76, 217–242.

Zheng, W., Z. Luo, and M. J. van der Laan (2016), "Marginal Structural Models with Counterfactual Effect Modifiers," U.C. Berkeley Division of Biostatistics Working Paper Series, Working Paper 348.

Zhu, Y. and J. Bradic (2017): "Linear Hypothesis Testing in Dense High-Dimensional Linear Models," *Journal of the American Statistical Association* 112.

Zhu, Y. and J. Bradic (2018): "Breaking the Curse of Dimensionality in Regression," *Journal of Machine Learning Research*, forthcoming.

Zubizarreta, J.R. (2015): "Stable Weights that Balance Covariates for Estimation With Incomplete Outcome Data," *Journal of the American Statistical Association* 110, 910-922.