

# HIGH-DIMENSIONAL METRICS IN R

VICTOR CHERNOZHUKOV, CHRISTIAN HANSEN, MARTIN SPINDLER

ABSTRACT. The package High-dimensional Metrics (`hdm`) is an evolving collection of statistical methods for estimation and quantification of uncertainty in high-dimensional approximately sparse models. It focuses on providing semi-parametrically efficient estimators, confidence intervals, and significance testing for low-dimensional subcomponents of the high-dimensional parameter vector. This vignette offers a brief introduction and a tutorial to the implemented methods. R and the package `hdm` are open-source software projects and can be freely downloaded from CRAN: <http://cran.r-project.org>.

## CONTENTS

1. Introduction	2
2. How to get started	3
3. Prediction using Approximate Sparsity	3
3.1. Prediction in Linear Models using Approximate Sparsity	3
3.2. A Joint Significance Test for Lasso Regression	6
R implementation	6
Example	7
4. Inference on Target Regression Coefficients	8
4.1. Intuition for the Orthogonality Principle in Linear Models via Partialling Out	9
4.2. Inference: Confidence Intervals and Significance Testing	11
4.3. Application: Estimation of the treatment effect in a linear model with many confounding factors	13
5. Instrumental Variable Estimation in a High-Dimensional Setting	15
5.1. Estimation and Inference	16
R Implementation	16
5.2. Application: Economic Development and Institutions	16
5.3. Application: Impact of Eminent Domain Decisions on Economic Outcomes	18
6. Inference on Treatment Effects in a High-Dimensional Setting	20
6.1. Treatment Effects Parameters – a short Introduction	21
6.2. Estimation and Inference of Treatment effects	22
R Implementation	22
6.3. Application: 401(k) plan participation	22
7. The Lasso Methods for Discovery of Significant Causes amongst Many Potential Causes, with Many Controls	26
8. Conclusion	28

Appendix A. Data Sets	29
A.1. Pension Data	29
A.2. Growth Data	29
A.3. Institutions and Economic Development – Data on Settler Mortality	30
A.4. Data on Eminent Domain	30
A.5. BLP data	30
A.6. CPS data	30
References	32

## 1. INTRODUCTION

Analysis of high-dimensional models, models in which the number of parameters to be estimated is large relative to the sample size, is becoming increasingly important. Such models arise naturally in modern data sets which have many measured characteristics available per individual observation as in, for example, population census data, scanner data, and text data. Such models also arise naturally even in data with a small number of measured characteristics in situations where the exact functional form with which the observed variables enter the model is unknown and we create many technical variables, a dictionary, from the raw characteristics. Examples covered by this scenario include semiparametric models with nonparametric nuisance functions. More generally, models with many parameters relative to the sample size often arise when attempting to model complex phenomena.

With increasing availability of such data sets in economics and other data science fields, new methods for analyzing those data have been developed. The R package `hdm` contains implementations of recently developed methods for high-dimensional approximately sparse models, mainly relying on forms of lasso and post-lasso as well as related estimation and inference methods. The methods are illustrated with econometric applications, but are also useful in other disciplines such as medicine, biology, sociology or psychology.

The methods which are implemented in this package are distinct from already available methods in other packages in the following three major ways:

- 1) First, we provide efficient estimators and uniformly valid confidence intervals for various low-dimensional causal/structural parameters appearing in high-dimensional approximately sparse models. For example, we provide efficient estimators and uniformly valid confidence intervals for a regression coefficient on a target variable (e.g., a treatment or policy variable) in a high-dimensional sparse regression model. We also provide estimates and confidence intervals for average treatment effect (ATE) and average treatment effect for the treated (ATET), as well extensions of these parameters to the endogenous setting.
- 2) Second, we provide theoretically grounded, data-driven choice of the penalty level  $\lambda$  in the lasso method. Because of this, we call the resulting method the “rigorous” lasso (`=rlasso`). The prefix `r` in function names should underscore this. In high-dimensions setting cross-validation is very popular, but it lacks a theoretical justification and some theoretical proposals for the choice of  $\lambda$  are often not feasible.

- 3) Third, we provide a version of lasso for linear regressions specifically that expressly handles non-Gaussian and heteroskedastic errors.

## 2. HOW TO GET STARTED

R is an open source software project and can be freely downloaded from the CRAN website along with its associated documentation. The R package `hdm` can be downloaded from [cran.r-project.org](http://cran.r-project.org). To install the `hdm` package from R we simply type,

```
install.packages("hdm")
```

The most current version of the package (development version) is maintained at R-Forge and can be installed by

```
install.packages("hdm", repos = "http://R-Forge.R-project.org")
```

Provided that your machine has a proper internet connection and you have write permission in the appropriate system directories, the installation of the package should proceed automatically. Once the `hdm` package is installed, it can be loaded to the current R session by the command,

```
library(hdm)
```

Online help is available in two ways. For example, you can type:

```
help(package = "hdm")
help(rlasso)
```

The former command gives an overview over the available commands in the package, and the latter gives detailed information about a specific command.

More generally one can initiate a web-browser help session with the command,

```
help.start()
```

and navigate as desired. The browser approach is better adapted to exploratory inquiries, while the command line approach is better suited to confirmatory ones.

A valuable feature of R help files is that the examples used to illustrate commands are executable, so they can be pasted into an R session or run as a group with a command like,

```
example(rlasso)
```

## 3. PREDICTION USING APPROXIMATE SPARSITY

**3.1. Prediction in Linear Models using Approximate Sparsity.** Consider high dimensional approximately sparse linear regression models. These models have a large number of regressors  $p$ , possibly much larger than the sample size  $n$ , but only a relatively small number  $s = o(n)$  of these regressors are important for capturing accurately the main features of the regression function. The latter assumption makes it possible to estimate these models effectively by searching for approximately the right set of regressors.

The model reads

$$y_i = x_i' \beta_0 + \varepsilon_i, \quad \mathbb{E}[\varepsilon_i x_i] = 0, \quad \beta_0 \in \mathbb{R}^p, \quad i = 1, \dots, n$$

where  $y_i$  are observations of the response variable,  $x_i = (x_{i,j}, \dots, x_{i,p})$ 's are observations of  $p$ -dimensional regressors, and  $\varepsilon_i$ 's are centered disturbances, where possibly  $p \gg n$ . Assume that the data sequence is i.i.d. for the sake of exposition, although the framework covered is considerably more general. An important point is that the errors  $\varepsilon_i$  may be non-Gaussian or heteroskedastic (Belloni, Chen, Chernozhukov, and Hansen, 2012).

The model can be exactly sparse, namely

$$\|\beta_0\|_0 \leq s = o(n),$$

or approximately sparse, namely that the values of coefficients, sorted in decreasing order,  $(|\beta_0|_{(j)})_{j=1}^p$  obey,

$$|\beta_0|_{(j)} \leq A j^{-a(\beta_0)}, \quad a(\beta_0) > 1/2, \quad j = 1, \dots, p.$$

An approximately sparse model can be well-approximated by an exactly sparse model with sparsity index

$$s \propto n^{1/(2a(\beta_0))}.$$

In order to get theoretically justified performance guarantees, we consider the Lasso estimator with data-driven penalty loadings:

$$\hat{\beta} = \arg \min_{\beta \in \mathbb{R}^p} \mathbb{E}_n[(y_i - x_i' \beta)^2] + \frac{\lambda}{n} \|\hat{\Psi} \beta\|_1$$

where  $\|\beta\|_1 = \sum_{j=1}^p |\beta_j|$ ,  $\hat{\Psi} = \text{diag}(\hat{\psi}_1, \dots, \hat{\psi}_p)$  is a diagonal matrix consisting of penalty loadings, and  $\mathbb{E}_n$  abbreviates the empirical average. The penalty loadings are chosen to insure basic equivariance of coefficient estimates to rescaling of  $x_{i,j}$  and can also be chosen to address heteroskedasticity in model errors. We discuss the choice of  $\lambda$  and  $\hat{\Psi}$  below.

Regularization by the  $\ell_1$ -norm naturally helps the Lasso estimator to avoid overfitting, but it also shrinks the fitted coefficients towards zero, causing a potentially significant bias. In order to remove some of this bias, consider the Post-Lasso estimator that applies ordinary least squares to the model  $\hat{T}$  selected by Lasso, formally,

$$\hat{T} = \text{support}(\hat{\beta}) = \{j \in \{1, \dots, p\} : |\hat{\beta}_j| > 0\}.$$

The Post-Lasso estimate is then defined as

$$\tilde{\beta} \in \arg \min_{\beta \in \mathbb{R}^p} \mathbb{E}_n \left( y_i - \sum_{j=1}^p x_{i,j} \beta_j \right)^2 : \beta_j = 0 \quad \text{if } \hat{\beta}_j = 0, \quad \forall j.$$

In words, the estimator is ordinary least squares applied to the data after removing the regressors that were not selected by Lasso. The Post-Lasso estimator was introduced and analysed in Belloni and Chernozhukov (2013).

A crucial matter is the choice of the penalization parameter  $\lambda$ . With the right choice of the penalty level, Lasso and Post-Lasso estimators possess excellent performance guarantees: They both achieve the near-oracle rate for estimating the regression function, namely with probability  $1 - \gamma - o(1)$ ,

$$\sqrt{\mathbb{E}_n[(x_i'(\hat{\beta} - \beta_0))^2]} \lesssim \sqrt{(s/n) \log p}.$$

In high-dimensions setting, cross-validation is very popular in practice but lacks theoretical justification and so may not provide such a performance guarantee. In sharp contrast, the choice of the penalization parameter  $\lambda$  in the Lasso and Post-Lasso methods in this package is theoretical grounded and feasible. Therefore we call the resulting method the “rigorous”Lasso method and hence add a prefix **r** to the function names.

In the case of homoscedasticity, we set the penalty loadings  $\hat{\psi}_j = \sqrt{\mathbb{E}_n x_{i,j}^2}$ , which insures basic equivariance properties. There are two choices for penalty level  $\lambda$ : the  $X$ -independent choice and  $X$ -dependent choice. In the  $X$ -independent choice we set the penalty level to

$$\lambda = 2c\sqrt{n}\hat{\sigma}\Phi^{-1}(1 - \gamma/(2p)),$$

where  $\Phi$  denotes the cumulative standard normal distribution,  $\hat{\sigma}$  is a preliminary estimate of  $\sigma = \sqrt{\mathbb{E}\varepsilon^2}$ , and  $c$  is a theoretical constant, which is set to  $c = 1.1$  by default for the Post-Lasso method and  $c = .5$  for the Lasso method, and  $\gamma$  is the probability level, which is set to  $\gamma = .1$  by default. The parameter  $\gamma$  can be interpreted as the probability of mistakenly not removing  $X$ 's when all of them have zero coefficients. In the  $X$ -dependent case the penalty level is calculated as

$$\lambda = 2c\hat{\sigma}\Lambda(1 - \gamma|X),$$

where

$$\Lambda(1 - \gamma|X) = (1 - \gamma) - \text{quantile of } n\|\mathbb{E}_n[x_i e_i]\|_\infty|X,$$

where  $X = [x_1, \dots, x_n]'$  and  $e_i$  are iid  $N(0, 1)$ , generated independently from  $X$ ; this quantity is approximated by simulation. The  $X$ -independent penalty is more conservative than the  $X$ -dependent penalty. In particular the  $X$ -dependent penalty automatically adapts to highly correlated designs, using less aggressive penalization in this case Belloni, Chernozhukov, and Hansen (2010).

In the case of heteroskedasticity, the loadings are set to  $\hat{\psi}_j = \sqrt{\mathbb{E}_n[x_{ij}^2 \hat{\varepsilon}_i^2]}$ , where  $\hat{\varepsilon}_i$  are preliminary estimates of the errors. The penalty level can be  $X$ -independent (Belloni, Chen, Chernozhukov, and Hansen, 2012):

$$\lambda = 2c\sqrt{n}\Phi^{-1}(1 - \gamma/(2p)),$$

or it can be  $X$ -dependent and estimated by a multiplier bootstrap procedure (Chernozhukov, Chetverikov, and Kato, 2013):

$$\lambda = c \times c_W(1 - \gamma),$$

where  $c_W(1 - \gamma)$  is the  $1 - \gamma$ -quantile of the random variable  $W$ , conditional on the data, where

$$W := n \max_{1 \leq j \leq p} |2\mathbb{E}_n[x_{ij}\hat{\varepsilon}_i e_i]|,$$

where  $e_i$  are iid standard normal variables distributed independently from the data, and  $\hat{\varepsilon}_i$  denotes an estimate of the residuals.

Estimation proceeds by iteration. The estimates of residuals  $\hat{\varepsilon}_i$  are initialized by running least squares of  $y_i$  on five regressors that are most correlated to  $y_i$ . This implies conservative starting values for  $\lambda$  and the penalty loadings, and leads to the initial Lasso and Post-Lasso estimates, which are then further updated by iteration. The resulting iterative procedure is fully justified in the theoretical literature.

**3.2. A Joint Significance Test for Lasso Regression.** A basic question frequently arising in empirical work is whether the Lasso regression has explanatory power, comparable to a F-test for the classical linear regression model. The construction of a joint significance test follows (Chernozhukov, Chetverikov, and Kato, 2013) (Appendix M), and can be described as:

Based on the model  $y_i = a_0 + x_i' b_0 + \varepsilon_i$ , the null hypothesis of joint statistical in-significance is  $b_0 = 0$ . The alternative is that of the joint statistical significance:  $b_0 \neq 0$ . The null hypothesis implies that

$$\mathbb{E}[(y_i - a_0)x_i] = 0,$$

and restriction can be tested using the sup-score statistic:

$$S = \|\sqrt{n}\mathbb{E}_n[(y_i - \hat{a}_0)x_i]\|_\infty,$$

where  $\hat{a}_i = \mathbb{E}_n[y_i]$ . The critical value for this statistic can be approximated by the multiplier bootstrap procedure, which simulates the statistic:

$$S^* = \|\sqrt{n}\mathbb{E}_n[(y_i - \hat{a}_0)x_i g_i]\|_\infty,$$

where  $g_i$ 's are iid  $N(0, 1)$ , conditional on the data. The  $(1 - \alpha)$ -quantile of  $S^*$  serves as the critical value,  $c(1 - \alpha)$ . We reject the null if  $S > c(1 - \alpha)$  in favor of statistical significant, and we keep the null of non-significance otherwise. This test procedure is implemented in the package when calling the `summary-method` of `rlasso`-objects.

**R implementation.** The function `rlasso` implements Lasso and post-Lasso, where the prefix “r” signifies that these are theoretically rigorous versions of Lasso and post-Lasso. The default option is post-Lasso, `post=TRUE`. This function returns an object of S3 class `rlasso` for which methods like `predict`, `print`, `summary` are provided.

`lassoShooting.fit` is the computational algorithm that underlies the estimation procedure, which implements a version of the Shooting Lasso Algorithm (Fu, 1998). The user has several options for choosing the non-default options. Specifically, the user can decide if an unpenalized `intercept` should be included (`TRUE` by default). The option `penalty` of the function `rlasso` allows different choices for the penalization parameter and loadings. It allows for homoskedastic or heteroskedastic errors with default `homoscedastic = FALSE`. Moreover, the dependence structure of the design matrix might be taken into consideration for calculation of the penalization parameter with `X.dependent.lambda = TRUE`. In combination with these options, the option `lambda.start` allows the user to set a starting value for  $\lambda$  for the different algorithms. Moreover, the user can provide her own fixed value for the penalty level – instead of the data-driven methods discussed above – by setting `homoscedastic = "none"` and supplying the value via `lambda.start`.

The constants  $c$  and  $\gamma$  from above can be set in the option `penalty`. The quantities  $\hat{\varepsilon}$ ,  $\hat{\Psi}$ ,  $\hat{\sigma}$  are calculated in a iterative manner. The maximum number of iterations and the tolerance when the algorithms should stop can be set with `control`.

The method `summary` of `rlasso`-objects displays additionally for model diagnosis the  $R^2$  value, the adjusted  $R^2$  with degrees of freedom equal to the number of selected parameters, and the sup-score statistic for joint significance – described above – with corresponding p-value.

**Example.** (Prediction Using Lasso and Post-Lasso) Consider generated data from a sparse linear model:

```
set.seed(1)
n = 100 #sample size
p = 100 # number of variables
s = 3 # nubmer of variables with non-zero coefficients
X = matrix(rnorm(n * p), ncol = p)
beta = c(rep(5, s), rep(0, p - s))
Y = X %*% beta + rnorm(n)
```

Next we estimate the model, print the results, and make in-sample and out-of sample predictions. We can use methods `print` and `summarize` to print the results, where the option `all` can be set to `FALSE` to limit the print only to the non-zero coefficients.

```
lasso.reg = rlasso(Y ~ X, post = FALSE) # use lasso, not-Post-lasso
# lasso.reg = rlasso(X, Y, post=FALSE)
sum.lasso <- summary(lasso.reg, all = FALSE) # can also do print(lasso.reg, all=FALSE)
##
## Call:
## rlasso.formula(formula = Y ~ X, post = FALSE)
##
## Post-Lasso Estimation: FALSE
##
## Total number of variables: 100
## Number of selected variables: 9
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -2.07855 -0.50062 -0.04591  0.55878  1.86776
##
##           Estimate
## (Intercept)  0.043
## 1           4.738
## 2           4.863
## 3           4.829
## 24          0.016
## 83          -0.065
## 84           0.004
## 85           0.034
## 87          -0.031
## 89          -0.021
##
```

```

## Residual standard error: 0.8017
## Multiple R-squared:  0.9906
## Adjusted R-squared:  0.9897
## Joint significance test:
## the sup score statistic for joint significance test is 51.16 with a p-value of 0
yhat.lasso = predict(lasso.reg) #in-sample prediction
Xnew = matrix(rnorm(n * p), ncol = p) # new X
Ynew = Xnew %*% beta + rnorm(n) #new Y
yhat.lasso.new = predict(lasso.reg, newdata = Xnew) #out-of-sample prediction

post.lasso.reg = rlasso(Y ~ X, post = TRUE) #now use post-lasso
print(post.lasso.reg, all = FALSE) # or use summary(post.lasso.reg, all=FALSE)
##
## Call:
## rlasso.formula(formula = Y ~ X, post = TRUE)
##
## (Intercept)          1          2          3
## 0.01695      4.94547      5.04876      4.98373
yhat.postlasso = predict(post.lasso.reg) #in-sample prediction
yhat.postlasso.new = predict(post.lasso.reg, newdata = Xnew) #out-of-sample prediction
MAE <- apply(cbind(abs(Ynew - yhat.lasso.new), abs(Ynew - yhat.postlasso.new)),
             2, mean)
names(MAE) <- c("lasso MAE", "Post-lasso MAE")
print(MAE, digits = 2) # MAE for Lasso and Post-Lasso
##      lasso MAE Post-lasso MAE
##          0.88          0.82

```

In the example above the sup-score statistic for overall significance is 51.16 with a pvalue of 0. This means that the null hypothesis is rejected on level  $\alpha = 0.05$  and the model seems to have explanatory power.

#### 4. INFERENCE ON TARGET REGRESSION COEFFICIENTS

Here we consider inference on the target coefficient  $\alpha$  in the model:

$$y_i = d_i \alpha + x_i' \beta_0 + \epsilon_i, \quad \mathbb{E} \epsilon_i (x_i', d_i')' = 0.$$

Here  $d_i$  is a target regressor such as treatment, policy or other variable whose regression coefficient  $\alpha_0$  we would like to learn (Belloni, Chernozhukov, and Hansen, 2014). If we are interested in coefficients of several or even many variables, we can simply write the model in the above form treating each variable of interest as  $d_i$  in turn and then applying the estimation and inference procedures described below.

We assume approximate sparsity for  $x_i' \beta_0$  with sufficient speed of decay of the sorted components of  $\beta_0$ , namely  $\mathbf{a}(\beta_0) > 1$ . This condition translates into having a sparsity index  $s \ll \sqrt{n}$ . In general  $d_i$  is



correlated to  $x_i$ , so  $\alpha_0$  cannot be consistently estimated by the regression of  $y_i$  on  $d_i$ . To keep track of the relationship of  $d_i$  to  $x_i$ , write

$$d_i = x_i' \pi_0^d + \rho_i^d, \quad \mathbb{E} \rho_i^d x_i = 0.$$

To estimate  $\alpha_0$ , we also impose approximate sparsity on the regression function  $x_i' \pi_0^d$  with sufficient speed of decay of sorted components of  $\pi_0^d$ , namely  $\mathfrak{a}(\pi_0^d) > 1$ .

**The Orthogonality Principle.** Note that we can not use naive estimates of  $\alpha_0$  based simply on applying Lasso and Post-Lasso to the first equation. Such a strategy in general does not produce root- $n$  consistent and asymptotically normal estimators of  $\alpha$ , due to the possibility of large omitted variable bias resulting from estimating the nuisance function  $x_i' \beta_0$  in high-dimensional setting. In order to overcome the omitted variable bias, we need to use orthogonalized estimating equations for  $\alpha_0$ . Specifically we seek to find a score  $\psi(w_i, \alpha, \eta)$ , where  $w_i = (y_i, x_i)'$  and  $\eta$  is the nuisance parameter, such that

$$\mathbb{E} \psi(w_i, \alpha_0, \eta_0) = 0, \quad \frac{\partial}{\partial \eta} \mathbb{E} \psi(w_i, \alpha_0, \eta_0) = 0.$$

The second equation is the orthogonality condition, which states that the equations are not sensitive to the first-order perturbations of the nuisance parameter  $\eta$  near the true value. The latter property allows estimation of these nuisance parameters  $\eta_0$  by regularized estimators  $\hat{\eta}$ , where regularization is done via penalization or selection. Without this property, regularization may have too much effect on the estimator of  $\alpha_0$  for regular inference to proceed.

The estimators  $\hat{\alpha}$  of  $\alpha_0$  solve the empirical analog of the equation above,

$$\mathbb{E}_n \psi(w_i, \hat{\alpha}, \hat{\eta}) = 0,$$

where we have plugged in the estimator  $\hat{\eta}$  for the nuisance parameter. Due to the orthogonality property the estimator is first-order equivalent to the infeasible estimator  $\tilde{\alpha}$  solving

$$\mathbb{E}_n \psi(w_i, \tilde{\alpha}, \eta_0) = 0,$$

where we use the true value of the nuisance parameter. The equivalence holds in a variety of models under plausible conditions. The systematic development of the orthogonality condition for inference on low-dimensional parameters in modern high-dimensional settings is given in Chernozhukov, Hansen, and Spindler (2015a).

It turns out that in the linear model the orthogonal equations are closely connected to the classical ideas of partialling out.

**4.1. Intuition for the Orthogonality Principle in Linear Models via Partialling Out.** One way to think about estimation of  $\alpha_0$  is to think of the regression model:

$$\rho_i^y = \alpha_0 \rho_i^d + \epsilon_i,$$

where  $\rho_i^y$  is the residual that is left after partialling out the linear effect of  $x_i$  from  $y_i$  and  $\rho_i^d$  is the residual that is left after partialling out the linear effect of  $x_i$  from  $d_i$ , both done in the population. Note that we have  $\mathbb{E} \rho_i^y x_i = 0$ , i.e.  $\rho_i^y = y_i - x_i' \pi_0^y$  where  $x_i' \pi_0^y$  is the linear projection of  $y_i$  on  $x_i$ . After partialling out,  $\alpha_0$  is the population regression coefficient in the univariate regression of  $\rho_i^y$  on  $\rho_i^d$ . This is the Frisch-Waugh-Lovell theorem. Thus,  $\alpha = \alpha_0$  solves the population equation:

$$\mathbb{E} (\rho_i^y - \alpha \rho_i^d) \rho_i^d = 0.$$

The score associated to this equation is:

$$\psi(w_i, \alpha, \eta) = (y_i - x_i' \pi^y) - \alpha(d_i - x_i' \pi^d)(d_i - x_i' \pi^d), \quad \eta = (\pi^y, \pi^d)',$$

$$\psi(w_i, \alpha_0, \eta_0) = (\rho_i^y - \alpha \rho_i^d) \rho_i^d, \quad \eta_0 = (\pi_0^y, \pi_0^d)'$$

It is straightforward to check that this score obeys the orthogonality principle; moreover, this score is the semi-parametrically efficient score for estimating the regression coefficient  $\alpha_0$ .

In low-dimensional settings, the empirical version of the partialling out approach is simply another way to do the least squares. Let's verify this in an example. First, we generate some data

```
set.seed(1)
n = 5000
p = 20
X = matrix(rnorm(n * p), ncol = p)
colnames(X) = c("d", paste("x", 1:19, sep = ""))
xnames = colnames(X)[-1]
beta = rep(1, 20)
y = X %*% beta + rnorm(n)
dat = data.frame(y = y, X)
```

We can estimate  $\alpha_0$  by running full least squares:

```
# full fit
fmla = as.formula(paste("y ~ ", paste(colnames(X), collapse = "+")))
full.fit = lm(fmla, data = dat)
summary(full.fit)$coef["d", 1:2]
## Estimate Std. Error
## 0.97807455 0.01371225
```

Another way to estimate  $\alpha_0$  is to first partial out the  $x$ -variables from  $y_i$  and  $d_i$ , and run least squares on the residuals:

```
fmla.y = as.formula(paste("y ~ ", paste(xnames, collapse = "+")))
fmla.d = as.formula(paste("d ~ ", paste(xnames, collapse = "+")))
# partial fit via ols
rY = lm(fmla.y, data = dat)$res
rD = lm(fmla.d, data = dat)$res
partial.fit.ls = lm(rY ~ rD)
summary(partial.fit.ls)$coef["rD", 1:2]
## Estimate Std. Error
## 0.97807455 0.01368616
```

One can see that the estimates are identical, while standard errors are nearly identical. In fact, the standard errors are asymptotically equivalent due to the orthogonality property of the estimating equations associated with the partialling out approach.

In high-dimensional settings, we can no longer rely on the full least-squares and instead may rely on Lasso or Post-Lasso for partialling out. This amounts to using orthogonal estimating equations, where we estimate the nuisance parameters by Lasso or Post-Lasso. Let's try this in the above example, using Post-Lasso for partialling out:

```
# partial fit via post-lasso
rY = rlasso(fmla.y, data = dat)$res
rD = rlasso(fmla.d, data = dat)$res
partial.fit.postlasso = lm(rY ~ rD)
summary(partial.fit.postlasso)$coef["rD", 1:2]
## Estimate Std. Error
## 0.97273870 0.01368677
```

We see that this estimate and standard errors are nearly identical to the previous estimates given above. In fact they are asymptotically equivalent to the previous estimates in the low-dimensional settings, with the advantage that, unlike the previous estimates, they will continue to be valid in the high-dimensional settings.

The orthogonal estimating equations method – based on partialling out via Lasso or post-Lasso – is implemented by the function `rlassoEffect`, using `method= "partialling out"`:

```
Eff = rlassoEffect(X[, -1], y, X[, 1], method = "partialling out")
summary(Eff)$coef[, 1:2]
## Estimate. Std. Error
## 0.97273870 0.01368677
```

Another orthogonal estimating equations method – based on the double selection of covariates – is implemented by the the function `rlassoEffect`, using `method= "double selection"`. Algorithmically the method is as follows:

- (1) Select controls  $x_{ij}$ 's that predict  $y_i$  by Lasso.
- (2) Select controls  $x_{ij}$ 's that predict  $d_i$  by Lasso.
- (3) Run OLS of  $y_i$  on  $d_i$  and the union of controls selected in steps 1 and 2.

```
Eff = rlassoEffect(X[, -1], y, X[, 1], method = "double selection")
summary(Eff)$coef[, 1:2]
## Estimate. Std. Error
## 0.97807455 0.01415624
```

**4.2. Inference: Confidence Intervals and Significance Testing.** The function `rlassoEffects` does inference – namely construction of confidence intervals and significance testing – for target variables. Those can be specified either by the variable names, an integer valued vector giving their position in `x` or by a logical vector indicating the variables for which inference should be conducted. It returns an object of S3 class `rlassoEffects` for which the methods `summary`, `print`, `confint`, and `plot` are provided. `rlassoEffects` is a wrap function for the function `rlassoEffect` which does inference for a single target regressor. The theoretical underpinning is given in Belloni, Chernozhukov,

and Hansen (2014) and for a more general class of models in Belloni, Chernozhukov, and Kato (2014). The function `rlassoEffects` might either be used in the form `rlassoEffects(x, y, index)` where `x` is a matrix, `y` denotes the outcome variable and `index` specifies the variables of `x` for which inference is conducted. This can be done by an integer vector (position of the variables), a logical vector or the name of the variables. An alternative usage is as `rlassoEffects(formula, data, I)` where `I` is a one-sided formula which specifies the variables for which inference is conducted. For further details we refer to the help page of the function and the following examples where both methods for usage are shown.

Here is an example of the use.

```
set.seed(1)
n = 100 #sample size
p = 100 # number of variables
s = 3 # number of non-zero variables
X = matrix(rnorm(n * p), ncol = p)
colnames(X) <- paste("X", 1:p, sep = "")
beta = c(rep(3, s), rep(0, p - s))
y = 1 + X %*% beta + rnorm(n)
data = data.frame(cbind(y, X))
colnames(data)[1] <- "y"
fm = paste("y ~", paste(colnames(X), collapse = "+"))
fm = as.formula(fm)
```

We can do inference on a set of variables of interest, e.g. the first, second, third, and the fiftieth:

```
# lasso.effect = rlassoEffects(X, y, index=c(1,2,3,50))
lasso.effect = rlassoEffects(fm, I = ~X1 + X2 + X3 + X50, data = data)
print(lasso.effect)
##
## Call:
## rlassoEffects.formula(formula = fm, data = data, I = ~X1 + X2 +
##   X3 + X50)
##
## Coefficients:
##      X1      X2      X3      X50
## 2.94448 3.04127 2.97540 0.07196
summary(lasso.effect)
## [1] "Estimates and significance testing of the effect of target variables"
##   Estimate. Std. Error t value Pr(>|t|)
## X1    2.94448    0.08815  33.404 <2e-16 ***
## X2    3.04127    0.08389  36.253 <2e-16 ***
## X3    2.97540    0.07804  38.127 <2e-16 ***
## X50    0.07196    0.07765   0.927  0.354
```

```
## ---
## Signif. codes:
## 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
confint(lasso.effect)
##           2.5 %    97.5 %
## X1    2.77171308 3.1172421
## X2    2.87685121 3.2056979
## X3    2.82244962 3.1283583
## X50  -0.08022708 0.2241377
```

The two methods are first-order equivalent in both low-dimensional and high-dimensional settings under regularity conditions. Not surprisingly we see that in the numerical example of this section, the estimates and standard errors produced by the two methods are very close to each other.

It is also possible to estimate joint confidence intervals. The method relies on a multiplier bootstrap method as described in Belloni, Chernozhukov, and Kato (2014). Joint confidence intervals can be invoked by setting the option `joint` to `TRUE` in the method `confint` for objects of class `rlassoEffects`.

```
confint(lasso.effect, level = 0.95, joint = TRUE)
##           2.5 %    97.5 %
## X1    2.6470175 3.2419377
## X2    2.7608170 3.3217321
## X3    2.7115233 3.2392846
## X50  -0.1854959 0.3294065
```

Finally, we can also plot the estimated effects with their confidence intervals:

```
plot(lasso.effect, main = "Confidence Intervals")
```

For logistic regression we provide the functions `rlassologit` and `rlassologitEffects`. Further information can be found in the help.

#### 4.3. Application: Estimation of the treatment effect in a linear model with many confounding factors.

A part of empirical growth literature has focused on estimating the effect of an initial (lagged) level of GDP (Gross Domestic Product) per capita on the growth rates of GDP per capita. In particular, a key prediction from the classical Solow-Swan-Ramsey growth model is the hypothesis of convergence, which states that poorer countries should typically grow faster and therefore should tend to catch up with the richer countries, conditional on a set of institutional and societal characteristics. Covariates that describe such characteristics include variables measuring education and science policies, strength of market institutions, trade openness, savings rates and others.

Thus, we are interested in a specification of the form:

$$y_i = \alpha_0 d_i + \sum_{j=1}^p \beta_j x_{ij} + \varepsilon_i,$$

where  $y_i$  is the growth rate of GDP over a specified decade in country  $i$ ,  $d_i$  is the log of the initial level of GDP at the beginning of the specified period, and the  $x_{ij}$ 's form a long list of country  $i$ 's characteristics

at the beginning of the specified period. We are interested in testing the hypothesis of convergence, namely that  $\alpha_1 < 0$ . Given that in the Barro and Lee (1994) data, the number of covariates  $p$  is large relative to the sample size  $n$ , covariate selection becomes a crucial issue in this analysis. We employ here the partialling out approach (as well as the double selection version) introduced in the previous section.

First, we load and prepare the data

```
data(GrowthData)
dim(GrowthData)
## [1] 90 63
y = GrowthData[, 1, drop = F]
d = GrowthData[, 3, drop = F]
X = as.matrix(GrowthData)[, -c(1, 2, 3)]
varnames = colnames(GrowthData)
```

Now we can estimate the effect of the initial GDP level. First, we estimate by OLS:

```
xnames = varnames[-c(1, 2, 3)] # names of X variables
dandxnames = varnames[-c(1, 2)] # names of D and X variables
# create formulas by pasting names (this saves typing times)
fmla = as.formula(paste("Outcome ~ ", paste(dandxnames, collapse = "+")))
ls.effect = lm(fmla, data = GrowthData)
```

Second, we estimate the effect by the partialling out by Post-Lasso:

```
dX = as.matrix(cbind(d, X))
lasso.effect = rlassoEffect(x = X, y = y, d = d, method = "partialling out")
summary(lasso.effect)
## [1] "Estimates and significance testing of the effect of target variables"
##      Estimate. Std. Error t value Pr(>|t|)
## [1,] -0.04432    0.01532  -2.893  0.00381 **
## ---
## Signif. codes:
## 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Third, we estimate the effect by the double selection method:

```
dX = as.matrix(cbind(d, X))
doubleseleffect = rlassoEffect(x = X, y = y, d = d, method = "double selection")
summary(doubleseleffect)
## [1] "Estimates and significance testing of the effect of target variables"
##      Estimate. Std. Error t value Pr(>|t|)
## gdps465 -0.04432    0.01685  -2.631  0.00851 **
## ---
## Signif. codes:
```

```
## 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

We then collect results in a nice latex table:

```
library(xtable)
table = rbind(summary(ls.effect)$coef["gdpsh465", 1:2], summary(lasso.effect)$coef[,
  1:2], summary(doublesele.effect)$coef[, 1:2])
colnames(table) = c("Estimate", "Std. Error") #names(summary(full.fit)$coef)[1:2]
rownames(table) = c("full reg via ols", "partial reg
via post-lasso ",
  "partial reg via double selection")
tab = xtable(table, digits = c(2, 2, 5))
```

```
tab
```

	Estimate	Std. Error
full reg via ols	-0.01	0.02989
partial reg via post-lasso	-0.04	0.01532
partial reg via double selection	-0.04	0.01685

We see that the OLS estimates are noisy, which is not surprising given that  $p$  is comparable to  $n$ . The partial regression approaches, based on Lasso selection of covariates in the two projection equations, in contrast yields much more precise estimates, which does support the hypothesis of conditional convergence. We note that the partial regression approaches represent a special case of the orthogonal estimating equations approach.

## 5. INSTRUMENTAL VARIABLE ESTIMATION IN A HIGH-DIMENSIONAL SETTING

In many applied settings the researcher is interested in estimating the (structural) effect of a variable (treatment variable), but this variable is endogenous, i.e. correlated with the error term. In this case, instrumental variables (IV) methods are used for identification of the causal parameters.

We consider the linear instrumental variables model:

$$(1) \quad y_i = \alpha_0 d_i + \gamma_0 x_i' + \varepsilon_i,$$

$$(2) \quad d_i = z_i' \Pi + \beta_0 x_i' + v_i,$$

where  $\mathbb{E}[\varepsilon_i(x_i', z_i')] = 0$ ,  $\mathbb{E}[v_i(x_i', z_i')] = 0$ , but  $\mathbb{E}[\varepsilon_i v_i] \neq 0$  leading to endogeneity. In this setting  $d_i$  is a scalar endogenous variable of interest,  $z_i$  is a  $p_z$ -dimensional vector of instruments and  $x_i$  is a  $p_x$ -dimensional vector of control variables.

In this section we present methods to estimate the effect  $\alpha_0$  in a setting where either  $x$  is high-dimensional or  $z$  is high-dimensional. Instrumental variables estimation with very many instruments was analysed in Belloni, Chen, Chernozhukov, and Hansen (2012), the extension with many instruments and many controls in Chernozhukov, Hansen, and Spindler (2015b).

**5.1. Estimation and Inference.** To get efficient estimators and uniformly valid confidence intervals for the structural parameters there are different strategies which are asymptotically equivalent where again orthogonalization (via partialling out) is a key concept.

In the case of the high-dimensional instrument  $z_i$  and low-dimensional  $x_i$ . We predict the endogenous variable  $d_i$  using (Post-)Lasso regression of  $d_i$  on the instruments  $z_i$  and  $x_i$ , forcing the inclusion of  $x_i$ . Then we compute the IV estimator (2SLS)  $\hat{\alpha}$  of the parameter  $\alpha_0$  using the predicted value  $\hat{d}_i$  as instrument and using  $x_i$ 's as controls. We then perform inference on  $\alpha_0$  using  $\hat{\alpha}$  and conventional heteroskedasticity robust standard errors.

In the case of the low-dimensional instrument  $z_i$  and high-dimensional  $x_i$ , we first partial out the effect of  $x_i$  from  $d_i$ ,  $y_i$ , and  $z_i$  by (Post-)Lasso. Second, we then use the residuals to compute the IV estimator (2SLS)  $\hat{\alpha}$  of the parameter  $\alpha_0$ . We then perform inference on  $\alpha_0$  using  $\hat{\alpha}$  and conventional heteroskedasticity robust standard errors.

In the case of the high-dimensional instrument  $z_i$  and high-dimensional  $x_i$  the algorithm described in Chernozhukov, Hansen, and Spindler (2015b) is adopted.

**R Implementation.** The wrap function `rlassoIV` handles all of the previous cases. It has the options `select.X` and `select.Z` which implement selection of either covariates or instruments, both with default values set to `TRUE`. The class of the return object depends on the chosen options, but the methods `summary`, `print` and `confint` are available for all. The functions `rlassoSelectX` and `rlassoSelectZ` do selection on  $x$ -variables only and  $z$ -variables only. Selection on both is done in `rlassoIV`. All functions employ the option `post=TRUE` as default, which uses post-Lasso for partialling out. By setting `post=FALSE` we can employ Lasso instead of Post-Lasso. Finally, the package provides the standard function `tsls`, which implements the “classical” two-stage least squares estimation.

**Function usage** Both the family of `rlassoIV`-functions and the family of the functions for treatment effects, which are introduced in the next section, allow use with both formula-interface and by handing over the prepared model matrices. Hence the general pattern for use with formula is `function(formula, data, ...)` where formula consists of two-parts and is a member of the class `Formula`. These formulas are of the pattern `y ~ d + x | x + z` where `y` is the outcome variable, `x` are exogenous variables, `d` endogenous variables (if several ones are allowed depends on the concrete function), and `z` denote the instrumental variables. A more primitive use of the functions is by simply hand over the required group of variables as matrices: `function(x=x, d=d, y=y, z=z)`. In some of the following examples both alternatives are displayed.

**5.2. Application: Economic Development and Institutions.** Estimating the causal effect of institutions on output is complicated by the simultaneity between institutions and output: specifically, better institutions may lead to higher incomes, but higher incomes may also lead to the development of better institutions. To help overcome this simultaneity, Acemoglu, Johnson, and Robinson (2001) use mortality rates for early European settlers as an instrument for institution quality. The validity of this instrument hinges on the argument that settlers set up better institutions in places where they are more likely to establish long-term settlements, that where they are likely to settle for the long term is related to settler mortality at the time of initial colonization, and that institutions are highly persistent. The exclusion restriction for the instrumental variable is then motivated by the argument



that GDP, while persistent, is unlikely to be strongly influenced by mortality in the previous century, or earlier, except through institutions.

In this application, we consider the problem of selecting controls. The input raw controls are the Latitude and the continental dummies. The technical controls can include various polynomial transformations of the Latitude, possibly interacted with the continental dummies. Such flexible specifications of raw controls results in the high-dimensional  $x$ , with dimension comparable to the sample size.

First, we process the data

```
data(AJR)
y = AJR$GDP
d = AJR$Exprop
z = AJR$logMort
x = model.matrix(~-1 + (Latitude + Latitude2 + Africa + Asia +
  Namer + Samer)^2, data = AJR)
dim(x)
## [1] 64 21
```

Then we estimate an IV model with selection on the  $X$

```
# AJR.Xselect = rlassoIV(x=x, d=d, y=y, z=z, select.X=TRUE,
# select.Z=FALSE)
AJR.Xselect = rlassoIV(GDP ~ Exprop + (Latitude + Latitude2 +
  Africa + Asia + Namer + Samer)^2 | logMort + (Latitude +
  Latitude2 + Africa + Asia + Namer + Samer)^2, data = AJR,
  select.X = TRUE, select.Z = FALSE)
summary(AJR.Xselect)
## [1] "Estimation and significance testing of the effect of target variables in the IV regression mo
##      coeff.      se. t-value p-value
## Exprop 0.8792 0.2975   2.956 0.00312 **
## ---
## Signif. codes:
## 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
confint(AJR.Xselect)
##              2.5 %   97.5 %
## Exprop 0.2961694 1.462225
```

It is interesting to understand what the procedure above is doing. In essence, it partials out  $x_i$  from  $y_i$ ,  $d_i$  and  $z_i$  using Post-Lasso and applies the 2SLS to the residual quantities.

Let us investigate partialling out in more detail in this example. We can first try to use OLS for partialling out:

```
# partialling out by linear model
fmla.y = GDP ~ (Latitude + Latitude2 + Africa + Asia + Namer +
  Samer)^2
```

```
fmla.d = Exprop ~ (Latitude + Latitude2 + Africa + Asia + Namer +
  Samer)^2
fmla.z = logMort ~ (Latitude + Latitude2 + Africa + Asia + Namer +
  Samer)^2
rY = lm(fmla.y, data = AJR)$res
rD = lm(fmla.d, data = AJR)$res
rZ = lm(fmla.z, data = AJR)$res
# ivfit.lm = tsls(y=rY,d=rD, x=NULL, z=rZ, intercept=FALSE)
ivfit.lm = tsls(rY ~ rD | rZ, intercept = FALSE)
print(cbind(ivfit.lm$coef, ivfit.lm$se), digits = 3)
##      [,1] [,2]
## rD 1.27 1.73
```

We see that the estimates exhibit large standard errors. The imprecision is expected because dimension of  $x$  is quite large, comparable to the sample size.

Next, we replace the OLS operator by post-Lasso for partialling out

```
# partialling out by lasso
rY = rlasso(fmla.y, data = AJR)$res
rD = rlasso(fmla.d, data = AJR)$res
rZ = rlasso(fmla.z, data = AJR)$res
# ivfit.lasso = tsls(y=rY,d=rD, x=NULL, z=rZ,
# intercept=FALSE)
ivfit.lasso = tsls(rY ~ rD | rZ, intercept = FALSE)
print(cbind(ivfit.lasso$coef, ivfit.lasso$se), digits = 3)
##      [,1] [,2]
## rD 0.879 0.297
```

This is exactly what command `rlassoIV` is doing by calling the command `rlassoSelectX`, so the numbers we see are identical to those reported first. In comparison to OLS results, we see precision is improved by doing selection on the exogenous variables.

**5.3. Application: Impact of Eminent Domain Decisions on Economic Outcomes.** Here we investigate the effect of pro-plaintiff decisions in cases of eminent domain (government's takings of private property) on economic outcomes. The analysis of the effects of such decisions is complicated by the possible endogeneity between judicial decisions and potential economic outcomes. To address the potential endogeneity, we employ an instrumental variables strategy based on the random assignment of judges to the federal appellate panels that make the decisions. Because judges are randomly assigned to three-judge panels, the exact identity of the judges and their demographics are randomly assigned conditional on the distribution of characteristics of federal circuit court judges in a given circuit-year. Under this random assignment, the characteristics of judges serving on federal appellate panels can only be related to property prices through the judges' decisions; thus the judge's characteristics will plausibly satisfy the instrumental variable exclusion restriction. For further information on this

application and the data set we refer to Chen and Yeh (2010) and Belloni, Chen, Chernozhukov, and Hansen (2012).

First, we load the data and construct the matrices with the controls ( $x$ ), instruments ( $z$ ), outcome ( $y$ ), and treatment variables ( $d$ ). Here we consider regional GDP as the outcome variable.

```
data(EminentDomain)
z <- EminentDomain$logGDP$z
x <- EminentDomain$logGDP$x
y <- EminentDomain$logGDP$y
d <- EminentDomain$logGDP$d
```

As mentioned above,  $y$  is the economic outcome, the logarithm of the GDP,  $d$  the number of plaintiff appellate takings decisions in federal circuit court  $c$  and year  $t$ ,  $x$  is a matrix with control variables, and  $z$  is the matrix with instruments. Here we consider socio-economic and demographic characteristics of the judges as instruments.

First, we estimate the effect of the treatment variable by simple OLS and 2SLS using two instruments:

```
ED.ols = lm(y ~ cbind(d, x))
ED.2sls = tsls(y = y, d = d, x = x, z = z[, 1:2], intercept = FALSE)
```

Next, we estimate the model with selection on the instruments.

```
lasso.IV.Z = rlassoIV(x = x, d = d, y = y, z = z, select.X = FALSE,
  select.Z = TRUE)
# or lasso.IV.Z = rlassoIVselectZt(x=X, d=d, y=y, z=z)
summary(lasso.IV.Z)
## [1] "Estimates and significance testing of the effect of target variables in the IV regression model"
##      coeff.      se. t-value p-value
## d1 0.01543 0.01926  0.801  0.423
confint(lasso.IV.Z)
##           2.5 %    97.5 %
## d1 -0.0223192 0.0531804
```

Finally, we do selection on both the  $x$  and  $z$  variables.

```
lasso.IV.XZ = rlassoIV(x = x, d = d, y = y, z = z, select.X = TRUE,
  select.Z = TRUE)
summary(lasso.IV.XZ)
## Estimates and Significance Testing of the effect of target variables in the IV regression model
##      coeff.      se. t-value p-value
## d1 -0.03488 0.15881  -0.22  0.826
confint(lasso.IV.XZ)
##           2.5 %    97.5 %
```

```
## d1 -0.3461475 0.2763868
```

Comparing the results we see, that the OLS estimates indicate that the influence of pro plaintiff appellate takings decisions in federal circuit court is significantly positive, but the 2SLS estimates which account for the potential endogeneity render the results insignificant. Employing selection on the instruments yields similar results. Doing selection on both the  $x$ - and  $z$ -variables results in extremely imprecise estimates.

Finally, we compare all results

```
library(xtable)
table = matrix(0, 4, 2)
table[1, ] = summary(ED.ols)$coef[2, 1:2]
table[2, ] = cbind(ED.2sls$coef[1], ED.2sls$se[1])
table[3, ] = summary(lasso.IV.Z)[, 1:2]
## [1] "Estimates and significance testing of the effect of target variables in the IV regression model"
##      coeff.      se. t-value p-value
## d1 0.01543 0.01926  0.801  0.423
table[4, ] = summary(lasso.IV.XZ)[, 1:2]
## Estimates and Significance Testing of the effect of target variables in the IV regression model
##      coeff.      se. t-value p-value
## d1 -0.03488 0.15881  -0.22  0.826
colnames(table) = c("Estimate", "Std. Error")
rownames(table) = c("ols regression", "IV estimation ", "selection on Z",
                    "selection on X and Z")
tab = xtable(table, digits = c(2, 2, 7))
```

```
tab
```

	Estimate	Std. Error
ols regression	0.01	0.0052848
IV estimation	0.02	0.0183376
selection on Z	0.02	0.0192605
selection on X and Z	-0.03	0.1588127

## 6. INFERENCE ON TREATMENT EFFECTS IN A HIGH-DIMENSIONAL SETTING

In this section, we consider estimation and inference on treatment effects when the treatment variable  $d$  enters non-separably in determination of the outcomes. This case is more complicated than the additive case, which is covered as a special case of Section 3. However, the same underlying principle – the orthogonality principle – applies for the estimation and inference on the treatment effect parameters. Estimation and inference of treatment effects in a high-dimensional setting is analysed in Belloni, Chernozhukov, Fernández-Val, and Hansen (2013).

**6.1. Treatment Effects Parameters – a short Introduction.** In many situations researchers are asked to evaluate the effect of a policy intervention. Examples are the effectiveness of a job-related training program or the effect of a newly developed drug. We consider  $n$  units or individuals,  $i = 1, \dots, n$ . For each individual we observe the treatment status. The treatment variable  $D_i$  takes the value 1, if the unit received (active) treatment, and 0, if it received the control treatment. For each individual we observe the outcome for only one of the two potential treatment states. Hence, the observed outcome depends on the treatment status and is denoted by  $Y_i(D_i)$ . One important parameter of interest is the average treatment effect (ATE):

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

This quantity can be interpreted as the average effect of the policy intervention.

Researchers might also be interested in the average treatment effect on the treated (ATET) given by

$$\mathbb{E}[Y(1) - Y(0)|D = 1] = \mathbb{E}[Y(1)|D = 1] - \mathbb{E}[Y(0)|D = 1].$$

This is the average treatment effect restricted to the population the treated individuals.

When treatment  $D$  is randomly assigned conditional on confounding factors  $X$ , the ATE and ATET can be identified by regression or propensity score weighting methods. Our identification and estimation method rely on the combination of two methods to create orthogonal estimating equations for these parameters.<sup>1</sup>

In observational studies, the potential treatments are endogenous, i.e. jointly determined with the outcome variable. In such cases, causal effects may be identified with the use of a binary instrument  $Z$  that affects the treatment but is independent of the potential outcomes. An important parameter in this case is the local average treatment effect (LATE):

$$\mathbb{E}[Y(1) - Y(0)|D(1) > D(0)].$$

The random variables  $D(1)$  and  $D(0)$  indicate the potential participation decisions under the instrument states 1 and 0. LATE is the average treatment effect for the subpopulation of compliers – those units that respond to the change in the instrument. Another parameter of interest is the local average treatment effect of the treated (LATET):

$$\mathbb{E}[Y(1) - Y(0)|D(1) > D(0), D = 1],$$

which is the average effect for the subpopulation of the treated compliers.

When the instrument  $Z$  is randomly assigned conditional on confounding factors  $X$ , the LATE and LATET can be identified by instrumental variables regression or propensity score weighting methods. Our identification and estimation method rely on the combination of two methods to create orthogonal estimating equations for these parameters.

---

<sup>1</sup>It turns out that the orthogonal estimating equations are the same as doubly robust estimating equations, but emphasizing the name "doubly robust" could be confusing in the present context.

**6.2. Estimation and Inference of Treatment effects.** We consider the estimation of the effect of an endogenous binary treatment,  $D$ , on an outcome variable,  $Y$ , in a setting with very many potential control variables. In the case of endogeneity, the presence of a binary instrumental variable,  $Z$ , is required for the estimation of the LATE and LATET.

When trying to estimate treatment effects, the researcher has to decide what conditioning variables to include. In the case of a non-randomly assigned treatment or instrumental variable, the researcher must select the conditioning variables so that the instrument or treatment is plausibly exogenous. Even in the case of random assignment, for a precise estimation of the policy variable selection of control variables is necessary to absorb residual variation, but overfitting should be avoided. For uniformly valid post-selection inference, “orthogonal” estimating equations as described above are they key to the efficient estimation and valid inference. We refer to Belloni, Chernozhukov, Fernández-Val, and Hansen (2013) for details.

**R Implementation.** The package contains the functions `rlassoATE`, `rlassoATET`, `rlassoLATE`, and `rlassoLATE` to estimate the corresponding treatment effects. All functions have as arguments the outcome variable  $y$ , the treatment variable  $d$ , and the conditioning variables  $x$ . The functions `rlassoLATE`, and `rlassoLATE` have additionally the argument  $z$  for the binary instrumental variable. For the calculation of the standard errors bootstrap methods are implemented, with options to use `Bayes`, `normal`, or `wild` bootstrap. The number of repetitions can be specified with the argument `nRep` and the default is set to 500. By default no bootstrap standard errors are provided (`bootstrap="none"`). For the functions the logicals `intercept` and `post` can be specified to include an intercept and to do Post-Lasso at the selection steps. The family of treatment functions returns an object of class `rlassoTE` for which the methods `print`, `summary`, and `confint` are available.

**6.3. Application: 401(k) plan participation.** Though it is clear that 401(k) plans are widely used as vehicles for retirement saving, their effect on assets is less clear. The key problem in determining the effect of participation in 401(k) plans on accumulated assets is saver heterogeneity coupled with nonrandom selection into participation states. In particular, it is generally recognized that some people have a higher preference for saving than others. Thus, it seems likely that those individuals with the highest unobserved preference for saving would be most likely to choose to participate in tax-advantaged retirement savings plans and would also have higher savings in other assets than individuals with lower unobserved saving propensity. This implies that conventional estimates that do not allow for saver heterogeneity and selection of the participation state will be biased upward, tending to overstate the actual savings effects of 401(k) and IRA participation.

Again, we start first with the data preparation:

```
data(pension)
y = pension$tw
d = pension$p401
z = pension$e401
X = pension[, c("i2", "i3", "i4", "i5", "i6", "i7", "a2", "a3",
               "a4", "a5", "fsize", "hs", "smcol", "col", "marr", "twoearn",
               "db", "pira", "hown")] # simple model
```

```
xvar = c("i2", "i3", "i4", "i5", "i6", "i7", "a2", "a3", "a4",
        "a5", "fsize", "hs", "smcol", "col", "marr", "twoearn", "db",
        "pira", "hown")
xpart = paste(xvar, collapse = "+")
form = as.formula(paste("tw ~ ", paste(c("p401", xvar), collapse = "+"),
                        "|", paste(xvar, collapse = "+")))
formZ = as.formula(paste("tw ~ ", paste(c("p401", xvar), collapse = "+"),
                        "|", paste(c("e401", xvar), collapse = "+")))
```

Now we can compute the estimates of the target treatment effect parameters. For ATE and ATET we report the the effect of eligibility for 401(k).

```
# pension.ate = rlassoATE(X,d,y)
pension.ate = rlassoATE(form, data = pension)
summary(pension.ate)

## Estimation and significance testing of the treatment effect
## Type: ATE
## Bootstrap: not applicable
##   coeff.   se. t-value p-value
## TE 10490 1920   5.464 4.67e-08 ***
## ---
## Signif. codes:
## 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

# pension.atet = rlassoATET(X,d,y)
pension.atet = rlassoATET(form, data = pension)
summary(pension.atet)

## Estimation and significance testing of the treatment effect
## Type: ATET
## Bootstrap: not applicable
##   coeff.   se. t-value p-value
## TE 11810 2844   4.152 3.29e-05 ***
## ---
## Signif. codes:
## 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

For LATE and LATET we estimate the effect of 401(k) participation (d) with plan eligibility (z) as instrument.

```
pension.late = rlassoLATE(X, d, y, z)
# pension.late = rlassoLATE(formZ, data=pension)
summary(pension.late)

## Estimation and significance testing of the treatment effect
## Type: LATE
```

```

## Bootstrap: not applicable
##   coeff.   se. t-value p-value
## TE 12189 2734 4.458 8.27e-06 ***
## ---
## Signif. codes:
## 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
pension.latet = rlassoLATET(X, d, y, z)
# pension.latet = rlassoLATET(formZ, data=pension)
summary(pension.latet)
## Estimation and significance testing of the treatment effect
## Type: LATET
## Bootstrap: not applicable
##   coeff.   se. t-value p-value
## TE 12687 3590 3.534 0.00041 ***
## ---
## Signif. codes:
## 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

The results are summarized into a table, which is then processed using `xtable` to produce the latex output:

```

library(xtable)
table = matrix(0, 4, 2)
table[1, ] = summary(pension.ate)[, 1:2]
## Estimation and significance testing of the treatment effect
## Type: ATE
## Bootstrap: not applicable
##   coeff.   se. t-value p-value
## TE 10490 1920 5.464 4.67e-08 ***
## ---
## Signif. codes:
## 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
table[2, ] = summary(pension.atet)[, 1:2]
## Estimation and significance testing of the treatment effect
## Type: ATET
## Bootstrap: not applicable
##   coeff.   se. t-value p-value
## TE 11810 2844 4.152 3.29e-05 ***
## ---
## Signif. codes:
## 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```



```

table[3, ] = summary(pension.late)[, 1:2]
## Estimation and significance testing of the treatment effect
## Type: LATE
## Bootstrap: not applicable
##   coeff.   se. t-value p-value
## TE 12189  2734   4.458 8.27e-06 ***
## ---
## Signif. codes:
## 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
table[4, ] = summary(pension.latet)[, 1:2]
## Estimation and significance testing of the treatment effect
## Type: LATET
## Bootstrap: not applicable
##   coeff.   se. t-value p-value
## TE 12687  3590   3.534 0.00041 ***
## ---
## Signif. codes:
## 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
colnames(table) = c("Estimate", "Std. Error")
rownames(table) = c("ATE", "ATET ", "LATE", "LATET")
tab = xtable(table, digits = c(2, 2, 2))

```

	Estimate	Std. Error
ATE	10490.07	1919.99
ATET	11810.45	2844.33
LATE	12188.66	2734.12
LATET	12686.87	3590.09

Finally, we estimate a model including all interaction effects:

```

# X = model.matrix(~ -1 + (i2 + i3 + i4 + i5 + i6 + i7 + a2 +
# a3 + a4 + a5 + fsize + hs + smcol + col + marr + twoearn +
# db + pira + hown)^2, data = pension) # model with
# interactions
xvar2 <- paste("(", xvar, ")^2", sep = "")
formExt = as.formula(paste("tw ~ ", paste(c("p401", xvar2), collapse = "+"),
  "|", paste(xvar2, collapse = "+")))
formZExt = as.formula(paste("tw ~ ", paste(c("p401", xvar2),
  collapse = "+"), "|", paste(c("e401", xvar2), collapse = "+")))

```

```

pension.ate = rlassoATE(X, z, y)
pension.atet = rlassoATET(X, z, y)
pension.late = rlassoLATE(X, d, y, z)
pension.latet = rlassoLATET(X, d, y, z)
# pension.ate = rlassoATE(formExt, data = pension)
# pension.atet = rlassoATET(formExt, data = pension)
# pension.late = rlassoLATE(formZExt, data = pension)
# pension.latet = rlassoLATET(formZExt, data = pension)
table = matrix(0, 4, 2)
table[1, ] = summary(pension.ate)[, 1:2]
table[2, ] = summary(pension.atet)[, 1:2]
table[3, ] = summary(pension.late)[, 1:2]
table[4, ] = summary(pension.latet)[, 1:2]
colnames(table) = c("Estimate", "Std. Error")
rownames(table) = c("ATE", "ATET ", "LATE", "LATET")
tab = xtable(table, digits = c(2, 2, 2))

```

This gives the following results:

	Estimate	Std. Error
ATE	8449.80	1895.43
ATET	8938.00	2529.25
LATE	12188.66	2734.12
LATET	12686.87	3590.09

## 7. THE LASSO METHODS FOR DISCOVERY OF SIGNIFICANT CAUSES AMONGST MANY POTENTIAL CAUSES, WITH MANY CONTROLS

Here we consider the model

$$\underbrace{Y_i}_{\text{Outcome}} = \underbrace{\sum_{l=1}^{p_1} D_{il}\alpha_l}_{\text{Causes}} + \underbrace{\sum_{j=1}^{p_2} W_{ij}\beta_j}_{\text{Controls}} + \underbrace{\epsilon_i}_{\text{Noise}}$$

where the number of potential causes  $p_1$  could be very large and the number of controls  $p_2$  could also be very large. The causes are randomly assigned conditional on controls.

Under approximate sparsity of  $\alpha = (\alpha_l)_{l=1}^{p_1}$  and  $\beta = (\beta_l)_{l=1}^{p_2}$ , we can use Lasso-based method of Belloni, Chernozhukov, and Kato (2014) for estimating  $(\alpha_l)_{l=1}^{p_1}$  and constructing a joint confidence band on  $(\alpha_l)_{l=1}^{p_1}$  and then checking which  $\alpha_l$ 's are significantly different from zero. The approach is based on building orthogonal estimating equations for each of  $(\alpha_l)_{l=1}^{p_1}$ , and can be interpreted as doing Frisch-Waugh procedure for each coefficient of interest, where we do partialling out via Lasso or OLS-after-Lasso.

This procedure is implemented in the R package `hdm`. Here is an example in which  $n = 100$ ,  $p_1 = 20$ , and  $p_2 = 20$ , so that total number of regressors is  $p = p_1 + p_2 = 40$ . In this example  $\alpha_1 = 5$  and

$\beta_1 = 5$ , i.e. there is only one true cause  $D_{i1}$ , among the large number of causes,  $D_{i1}, \dots, D_{i20}$ , and only one true control  $W_{i1}$ . This example is made super-simple for clarity sake. The Belloni, Chernozhukov, and Kato (2014) procedure, implemented by `rlasso.effects` command in R package `hdm`.

```
# library(hdm) library(stats)
set.seed(1)
n = 100
p1 = 20
p2 = 20
D = matrix(rnorm(n * p1), n, p1) # Causes
W = matrix(rnorm(n * p2), n, p2) # Controls
X = cbind(D, W) # Regressors
Y = D[, 1] * 5 + W[, 1] * 5 + rnorm(n) #Outcome
confint(rlassoEffects(X, Y, index = c(1:p1)), joint = TRUE)
##           2.5 %      97.5 %
## V1  4.45254352 5.27634920
## V2 -0.36919897 0.35985454
## V3 -0.40022202 0.23459987
## V4 -0.30226979 0.33541589
## V5 -0.32560180 0.32529338
## V6 -0.37606144 0.34882071
## V7 -0.27299706 0.34768804
## V8 -0.09355296 0.51986258
## V9 -0.23770854 0.44138018
## V10 -0.28169032 0.30856654
## V11 -0.36118733 0.25593690
## V12 -0.36016826 0.31669947
## V13 -0.22301062 0.42568025
## V14 -0.38644116 0.44829939
## V15 -0.37829104 0.36943955
## V16 -0.31779430 0.38385080
## V17 -0.23207059 0.46982430
## V18 -0.40623680 0.08387136
## V19 -0.15173470 0.43811152
## V20 -0.25749589 0.29721605
# BCK Joint Confidence Band for Reg Coefficients 1 to 20
```

As you can see the procedure correctly tells that only the first cause  $D_{i1}$ , among the large number of causes,  $D_{i1}, \dots, D_{i20}$ , is a statistically significant cause of  $Y$  (see the confidence interval for variable V1).

## 8. CONCLUSION

We have provided an introduction to some of the capabilities of the R package `hdm`. Inevitably, new applications will demand new features and, as the project is in its initial phase, unforeseen bugs will show up. In either case comments and suggestions of users are highly appreciated. We shall update the documentation and the package periodically. The most current version of the R package and its accompanying vignette will be made available at the homepage of the maintainer and `cran.r-project.org`. See the R command `vignette()` for details on how to find and view vignettes from within R .

## APPENDIX A. DATA SETS

In this section we describe briefly the data sets which are contained in the package and used afterwards. They might also be of general interest either for illustrating methods or for classroom presentation.

**A.1. Pension Data.** In the United States 401(k) plans were introduced to increase private individual saving for retirement. They allow the individual to deduct contributions from taxable income and allow tax-free accrual of interest on assets held within the plan (within an account). Employers provide 401(k) plans, and employers may also match a certain percentage of an employee's contribution. Because 401(k) plans are provided by employers, only workers in firms offering plans are eligible for participation. This data set contains individual level information about 401(k) participation and socio-economic characteristics.

The data set can be loaded with

```
data(pension)
```

A description of the variables and further references are given in Chernozhukov and Hansen (2004) and at the help page, for this example called by

```
help(pension)
```

The sample is drawn from the 1991 Survey of Income and Program Participation (SIPP) and consists of 9,915 observations. The observational units are household reference persons aged 25-64 and spouse if present. Households are included in the sample if at least one person is employed and no one is self-employed. All dollar amounts are in 1991 dollars. The 1991 SIPP reports household financial data across a range of asset categories. These data include a variable for whether a person works for a firm that offers a 401(k) plan. Households in which a member works for such a firm are classified as eligible for a 401(k). In addition, the survey also records the amount of 401(k) assets. Households with a positive 401(k) balance are classified as participants, and eligible households with a zero balance are considered nonparticipants. Available measures of wealth in the 1991 SIPP are total wealth, net financial assets, and net non-401(k) financial assets. Net non-401(k) assets are defined as the sum of checking accounts, U.S. saving bonds, other interest-earning accounts in banks and other financial institutions, other interest-earning assets (such as bonds held personally), stocks and mutual funds less non-mortgage debt, and IRA balances. Net financial assets are net non-401(k) financial assets plus 401(k) balances, and total wealth is net financial assets plus housing equity and the value of business, property, and motor vehicles.

**A.2. Growth Data.** Understanding what drives economic growth, measured in GDP, is a central question of macroeconomics. A well-known data set with information about GDP growth for many countries over a long period was collected by Barro and Lee (1994). This data set can be loaded by

```
data(GrowthData)
```

This data set contains the national growth rates in GDP per capita (Outcome) for many countries with additional covariates. A very important covariate is `gdpsh465`, which is the initial level of per-capita GDP. For further information we refer to the help page and the references herein, in particular the online descriptions of the data set.

**A.3. Institutions and Economic Development – Data on Settler Mortality.** This data set was collected by Acemoglu, Johnson, and Robinson (2001) to analyse the effect of institutions on economic development. The data can be loaded with

```
data(AJR)
```

The data set contains measurements of GDP, settler morality, an index measuring protection against expropriation risk and geographic information (latitude and continent dummies). There are 64 observations on 11 variables.

**A.4. Data on Eminent Domain.** Eminent domain refers to the government’s taking of private property. This data set was collected to analyse the effect of the number of pro-plaintiff appellate takings decisions on economic outcome variables such as house price indices.

The data set is loaded into R by

```
data(EminentDomain)
```

The data set consists of four “sub data sets” which have the following structure:

- *y*: outcome variable, a house price index or a local GDP index,
- *d*: the treatment variable, represents the number of pro-plaintiff appellate takings decisions in federal circuit court *c* and year *t*
- *x*: exogenous control variables that include a dummy variable for whether there were relevant cases in that circuit-year, the number of takings appellate decisions, and controls for the distribution of characteristics of federal circuit court judges in a given circuit-year
- *z*: instrumental variables, here characteristics of judges serving on federal appellate panels

The four data sets differ mainly in the dependent variable, which can be repeat-sales FHFA/OFHEO house price index for metro (FHFA) and non-metro (NM) areas, the Case-Shiller home price index (CS), and state-level GDP from the Bureau of Economic Analysis.

**A.5. BLP data.** This data set was analyzed in the seminal contribution of Berry, Levinsohn, and Pakes (1995) and stems from annual issues of the Automotive News Market Data Book. The data set includes information on all models marketed during the the period beginning 1971 and ending in 1990 containing 2217 model/years from 997 distinct models. A detailed description is given in Berry, Levinsohn, and Pakes (1995), p. 868–871. The function `constructIV` constructs instrumental variables along the lines described and used in Berry, Levinsohn, and Pakes (1995). The data set is loaded by

```
data(BLP)
```

It contains information on the price (in logarithm), the market share, and car characteristics like miles per gallon, miles per dollar, horse power per weight, space and air conditioning.

**A.6. CPS data.** The CPS is a monthly U.S. household survey conducted jointly by the U.S. Census Bureau and the Bureau of Labor Statistics. The data were collected for the year 2012. The sample comprises white non-hispanic, ages 25-54, working full time full year (35+ hours per week at least 50 weeks), exclude living in group quarters, self-employed, military, agricultural, and private household sector, allocated earning, inconsistent report on earnings and employment, missing data. It can be inspected with the command

```
data(cps2012)
```

## REFERENCES

- ACEMOGLU, D., S. JOHNSON, AND J. A. ROBINSON (2001): “The Colonial Origins of Comparative Development: An Empirical Investigation,” *American Economic Review*, 91(5), 1369–1401.
- BARRO, R. J., AND J.-W. LEE (1994): “Data set for a panel of 139 countries,” *NBER*, <http://www.nber.org/pub/barro.lee.html>.
- BELLONI, A., D. CHEN, V. CHERNOZHUKOV, AND C. HANSEN (2012): “Sparse Models and Methods for Optimal Instruments with an Application to Eminent Domain,” *Econometrica*, 80, 2369–2429, ArXiv, 2010.
- BELLONI, A., AND V. CHERNOZHUKOV (2013): “Least Squares After Model Selection in High-dimensional Sparse Models,” *Bernoulli*, 19(2), 521–547, ArXiv, 2009.
- BELLONI, A., V. CHERNOZHUKOV, I. FERNÁNDEZ-VAL, AND C. HANSEN (2013): “Program Evaluation with High-Dimensional Data,” *arXiv:1311.2645*, ArXiv, 2013.
- BELLONI, A., V. CHERNOZHUKOV, AND C. HANSEN (2010): “Inference for High-Dimensional Sparse Econometric Models,” *Advances in Economics and Econometrics. 10th World Congress of Econometric Society. August 2010*, III, 245–295, ArXiv, 2011.
- (2014): “Inference on Treatment Effects After Selection Amongst High-Dimensional Controls,” *Review of Economic Studies*, 81, 608–650, ArXiv, 2011.
- BELLONI, A., V. CHERNOZHUKOV, AND K. KATO (2014): “Uniform post-selection inference for least absolute deviation regression and other Z-estimation problems,” *Biometrika*.
- BERRY, S., J. LEVINSOHN, AND A. PAKES (1995): “Automobile Prices in Market Equilibrium,” *Econometrica*, 63, 841–890.
- CHEN, D. L., AND S. YEH (2010): “The Economic Impacts of Eminent Domain,” unpublished manuscript.
- 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.
- CHERNOZHUKOV, V., AND C. HANSEN (2004): “The impact of 401(k) participation on the wealth distribution: An instrumental quantile regression analysis,” *Review of Economics and Statistics*, 86(3), 735–751.
- CHERNOZHUKOV, V., C. HANSEN, AND M. SPINDLER (2015a): “Valid Post-Selection and Post-Regularization Inference: An Elementary, General Approach,” *Annual Review of Economics*, 7(1), 649–688.
- (2015b): “Valid Post-Selection and Post-Regularization Inference In Linear Models with Many Controls and Instruments,” *American Economic Review: Papers and Proceedings*.
- FU, W. J. (1998): “Penalized Regressions: The Bridge versus the Lasso,” *Journal of Computational and Graphical Statistics*, 7(3), 397–416.