

Online Appendix to Machine Labor

Joshua D. Angrist
MIT and NBER

Brigham Frandsen
BYU

We experimented with lasso commands in Stata 16, documented in [Stata \(2019\)](#), the Lassopack routines for Stata 15 documented in [Ahrens et al. \(2019\)](#), and the Elasticregress routines for Stata 15 documented in [Townsend \(2017\)](#). Different routines default to or allow the user to select different sorts of penalties. As in [Belloni et al. \(2013\)](#) and related work, our estimates use two types of penalty terms, one a [Bickel et al. \(2009\)](#)-type (BRT) plug-in penalty, described in [Belloni et al. \(2014\)](#), another using cross-validation. Each lasso estimation procedure has 3 parts: penalty selection; lasso estimation (to select controls or instruments); final estimation. Stata 16 lasso commands applied to problems other than estimation of elite college effects proved slow and, in some cases, numerically unstable, as documented in Table 3. This led us to Lassopack and Elasticregress for lasso estimation with large samples and/or high dimensional controls and instruments. Random forest estimates reported in Table A2 were computed in Stata; those in Table 5 use Stata and R.

Table 2: Post-Lasso Estimates of Elite College Effects

The sample size is 14,238. Estimates in columns 1-3 are from a post-double-selection (PDS) lasso procedure. Results in columns 4-6 are from a procedure applying lasso to a reduced-form regression of the outcome on the dictionary of controls. Controls include those used for the “self-revelation” model in Table 1 plus the following: indicators for being accepted to two colleges, three colleges, and four or more colleges; indicators for being rejected from one college, two colleges, three colleges, and four or more colleges; the number of schools applied to; the average SAT score among schools at which the applicant was accepted; the average SAT score among schools from which the applicant was rejected; the highest average SAT score across schools at which the applicant was accepted; the highest average SAT score across schools from which the applicant was rejected; the lowest average SAT score across schools at which the applicant was accepted; the lowest average SAT score across schools from which the applicant was rejected, and all two-way interactions and squared terms associated with the underlying list of possible controls. The dictionary of controls includes 384 variables, of which 303 are linearly independent. Penalties are computed once, using the original AK91 sample. Computational details are as follows:

- Columns 1 and 4 (Plug-in penalties)
 - Penalty: Uses Stata 16 `lasso linear` command to select controls, specifying a plug-in penalty, specifically,

```
lasso linear log-income 'dictionary-of-controls' [iw=weight],selection(plugin,het)
lasso linear elitetreatment 'dictionary-of-controls' [iw=weight],selection(plugin,het)
```

- Lasso (control selection): Uses `lasso linear` as specified above.
- Final estimates: Least squares regression of the outcome on the elite school variable, controlling for the variables selected by lasso (as needed for PDS and single selection), using College and Beyond sampling weights, with standard errors clustered at the institution level.
- Columns 2 and 5 (Cross-validated penalties)
 - Penalty: Computed using `lasso linear` to select controls using 10-fold cross-validation, which is the default in this case,


```
lasso linear log-income 'dictionary-of-controls' [iw=weight]
lasso linear elitetreatment 'dictionary-of-controls' [iw=weight]
```
 - Lasso (control selection): Uses `lasso linear` as specified above.
 - Final estimates: Least squares regression of the outcome on elite school variables, controlling for the variables selected by lasso (as needed for PDS and single selection), using College and Beyond sampling weights, with standard errors clustered at the institution level.
- Columns 3 and 6 (Cross-validated penalties using `cvlasso`)
 - Penalty: Computed using the Lasso pack command `cvlasso`. This yields a cross-validated MSE-minimizing penalty level, λ^{CV} . Lasso pack `rlasso` was also used to compute the default plug-in penalty, $\lambda^{default}$, specifying institutional weights and clustering (clustering induces robust, covariate-specific penalty loadings following Belloni et al. 2014). A cross-validated penalty scaling factor is then computed as $c^{CV} = 1.1\lambda^{CV}/\lambda^{default}$. The factor 1.1 arises because the default scaling factor in `rlasso` is 1.1.
 - Lasso (control selection): Uses `rlasso`, with c^{CV} replacing $c = 1.1$, specifying weights and clustering at the institution level.
 - Final estimates: Least squares regression of the outcome on elite school variables, controlling for the variables selected by lasso (as needed for PDS and single selection), using College and Beyond sampling weights, with standard errors clustered at the institution level.

Appendix Table A1: Alternative Post-Lasso Estimates of Elite College Effects

This table reports estimates using Lasso pack command `rlasso` and Elasticregress command `lassoregress`.

- Columns 1, 5, and 9 use `rlasso` and a plug-in penalty
 - Penalty: Computed using Lasso pack `rlasso` with default plug-in penalty, specifying weights and clustering by institution.
 - Lasso (control selection): Uses `rlasso` with regressor-specific penalty loadings.
 - Final estimates: Least squares regression of the outcome on elite school variables, controlling for the variables selected by Lasso, using College and Beyond sampling weights, with standard errors clustered at the institution level.

- Remaining columns use the Elasticregress command `lassoregress` and cross-validated penalties done three ways
 - Penalty: Computed via 10-fold cross validation as implemented in `lassoregress`, specifying institutional weights. Columns 2, 6, and 10 use the default cross-validated penalty, which minimizes cross-validated MSE. Columns 3, 7, and 11 specify the option `, lambda1se`, which uses the largest penalty such that the cross-validated MSE is within one standard deviation of the minimum. Columns 4, 8, and 12 calculate the penalty as 10 times the default.
 - Lasso (control selection): Estimated by the `lassoregress` command call that computes penalties.
 - Final estimates: Least squares regression of the outcome on elite school variables, controlling for the variables selected by Lasso, using College and Beyond sampling weights, with standard errors clustered at the institution level.

Elite college effects computed using the R-based package `glmnet` (Friedman et al., 2010) are similar to those reported in Tables A1 and 2. When applied to estimate the propensity score, however, the number of controls retained under `glmnet`-determined cross-validated penalties is generally much larger than the number of controls retained by Lassopack and Elasticregress. This is the result of a smaller penalty chosen for equation (15); `glmnet` lasso with `cvlasso`-determined penalties behaves like the Stata lasso routines, as does `glmnet` lasso on the PDS reduced form, equation (16).

Table 3: Angrist and Krueger (1991) Simulation Results

For estimates using the large AK91 sample with many fixed effects, Lassopack was faster and appeared to be more stable than Stata 16's `poivregrss` (used for Table 4 and described below)

- Penalty: Lasso estimates computed with a plug-in penalty use Lassopack `ivlasso` with default parameters. Estimates using cross-validated penalties were computed using Lassopack `cvlasso` as applied to a first-stage equation, specifying the option `, fe` to control for a full set of state-of-birth and year-of-birth interactions. A scaling factor for `ivlasso` is then computed as described for Table 2, above.
- Lasso (instrument selection): Computed using Lassopack routine `ivlasso` controlling for state-of-birth and year-of-birth interactions via the `, fe` option. Estimates using plug-in penalties use Lassopack `ivlasso` defaults. Estimates using cross-validated penalties employ the scaling adjustment described for the `cvlasso` estimates reported in Table 2. Note that `ivlasso` computes first stage estimates by calling `rlasso`.
- Final estimates: Computed via post-lasso 2SLS using `ivlasso`.

LIML estimates were computed using Stata `ivregress liml`. Lassoed versions of LIML use the instrument lists chosen for post-lasso 2SLS. Pretested LIML estimates use the instrument list described in the text. SSIV estimates split the sample in equal-sized halves randomly. One half-sample is used to estimate first stage parameters by OLS; these are carried over to the second half to compute cross-sample fitted values. Cross-sample fitted values and covariates are used to compute second-stage parameters using `ivregress`. The sample halves are then swapped, and the two resulting estimates averaged.

Post-lasso LIML estimates use the instruments chosen for post-lasso 2SLS. Post-lasso SSIV recomputes the lasso first stage in each half sample.

IJIVE estimates use the improved jackknifed instrumental variables procedure described in [Akerberg and Devereux \(2009\)](#). We implement this by first partialing the covariates (state-of-birth and year-of-birth interactions) out of the outcome, endogenous regressor, and instruments, and then applying the `jive` Stata command to the residualized variables.

Appendix Table A2: Angrist and Krueger (1991) Simulation Results for Random Forest

Random forest IV estimates in this table use Stata's `rforest` command (documented in [Schonlau \(2019\)](#)) to fit the first stage, with predictors YOB, POB, and QOB. Random forest estimates were computed with a minimum leaf size of 1 and 800, averaging results from 100 trees with no maximum depth (these are `rforest` defaults; [Oshiro et al. 2012](#) finds little payoff to more trees). The number of variables randomly investigated is equal to the square root of the number of right-hand-side variables (also a default setting). Random forest 2SLS uses random forest fitted values as excluded instruments, in a model with saturated control for year of birth and state of birth. Random forest SSIV fits the first stage with Stata command `rforest` in half the sample, assigning cross-sample fitted values to the relevant cells in the other half. Second-stage estimates are then obtained using these cross-fitted fitted values as instruments with saturated YOB-by-POB controls. As with the other SSIV estimates in the table, random forest SSIV swaps half samples and averages the resulting second-stage estimates from each.

Table 4: Simulation Results for Opening Weekend Effects

These simulations link a film's opening weekend ticket sales to subsequent ticket sales. OLS estimates are from a regression of second-weekend ticket sales on opening weekend ticket sales, controlling for year, week-of-year, day-of-week, and holiday dummies, as well as a set of second-weekend weather controls that includes indicators for the maximum temperature in 10-degree increments and indicators for rain, snow, and indicators for average precipitation in quarter inches per hour. 2SLS estimates (computed by `ivregress 2sls`) in columns 1-5 include the same exogenous controls used for OLS, with a set of 52 opening-weekend weather indicators used as excluded instruments. Estimates in columns 6-10 these instruments plus an additional 52 uniform noise instruments. Lasso estimates were computed as follows:

- Penalty: Lasso estimates computed with a plug-in penalty use Lassopack `ivlasso` with default parameters, partialing controls using the Lassopack, `partial()` option. Estimates using cross-validated penalties were computed using Lassopack `cvlasso` with opening-weekend ticket sales as the dependent variable and the controls and instruments as explanatory variables, partialing controls using the, `partial()` option. This yields a cross-validated MSE-minimizing penalty level, λ^{CV} . A scaling factor, c^{CV} , is then computed as described for Table 2, above. Penalties are recomputed for each simulation draw.
- Lasso (instrument selection): Computed using Lassopack `ivlasso`, with controls partialled via the, `partial()` option, using the scaling factor computed as described for Table 2, above.
- Final estimates: Computed via post-lasso 2SLS using `ivlasso`.

LIML estimates were computed using `ivregress liml`. SSIV estimates split the sample in equal-sized halves randomly. One half-sample is used to estimate first stage parameters by OLS; these are carried over to the second half to compute cross-sample fitted values. Cross-sample fitted values and covariates are used to compute second-stage parameters using `ivregress`. The sample halves are then swapped, and the two resulting estimates are averaged.

Appendix Table A3: Additional Simulation Results for Opening Weekend Effects

We briefly explored post-lasso IV estimators using the Stata 16 `poivregr` command (documented in [Stata 2019](#)). Motivated by [Chernozhukov et al. \(2015\)](#), `poivregr` allows the list of instruments and the list of exogenous covariates to be modeled as high-dimensional, applying lasso to the selection of variables in both. As with the estimators described in Tables 3 and 4, we focus first on the consequences of lasso for instrument selection.

The `poivregr` estimates discussed here use a plug-in penalty. Appendix Table A3 shows results generated using the `poivregr` command with the default plug-in penalty. This command is described on page 5 of [Stata \(2019\)](#) as

....partialing-out lasso instrumental-variables linear regression. This command estimates coefficients, standard errors, and confidence intervals and performs tests for variables of interest, both exogenous and endogenous, while using lassos to select from among potential control variables and instruments.

In models with a fixed number of controls, we computed `poivregr` estimates using the command

```
poivregr week2tickets 'low-dim exogenous vars' (week1tickets = 'high-dim instruments'), vce(robust).
```

For `poivregr` estimates treating controls as high dimensional, we used

```
poivregr week2tickets (week1tickets = 'high-dim instruments'), controls('high-dim exogenous vars') vce(robust).
```

The simulation results using `poivregr` to choose instruments and controls look much like those using `Lassopack` commands with a plug-in penalty. The resulting estimates exhibit bias on the order of 0.13 for models estimated with 52 instruments and 0.14 for models estimated with 104 instruments, only slightly below the bias of estimates computed using `Lassopack` commands. This in spite of the fact that `poivregr` retains an average of 1.3 instruments conditional on retaining any. It's noteworthy that `poivregr` fails to select *any* instruments in about two-thirds of the runs when starting with a dictionary of 52, while no instruments are selected in about three-fourths of the runs starting with a dictionary of 104. Surprisingly, `poivregr` reports second stage estimates even for these runs. Unsurprisingly, IV estimates generated without excluded instruments are biased and imprecise, with an MAE approaching that of OLS.

The bottom two rows of Table A3 show `poivregr` estimates computed for a procedure in which control variable coefficients are penalized along with first stage coefficients. That is, controls are treated as high-dimensional. Of 142 possible controls, 15-16 are retained in models that also retain instruments. Again, the average number of instruments retained is close to 1, conditional on having any retained. These estimates are less biased than post-lasso estimates computed with a plug-in penalty. Compare, for example, a bias of -0.096 using `poivregr` with a bias for plug-in post-lasso of -0.132 in the 52-instrument model. The better performance of `poivregr` with high-dimensional controls appears to reflect a higher first stage F when redundant controls are dropped. Even so, SSIV, IJIVE, and LIML estimates of this model are better

on MAE grounds. For the model with 104 instruments, which perhaps hews closest to the idea of sparsity, `poivregr` MAE beats that of SSIV (but not LIML or IJIVE) in runs for which instruments are selected.

Application of `poivregr` to models that select controls as well as instruments also yields second-stage estimates in runs with no instruments retained. In this case, instruments are retained in about half of 999 simulations, while no instruments are selected and an estimate of zero reported for 372 and 409 runs in the 52- and 104-instrument models, respectively. Remaining runs generate non-zero IV estimates even with no instruments chosen (though instrument-free IV estimates are widely dispersed).

Pages 267-8 of [Stata \(2019\)](#) describe the multi-step sequence of regression and post-lasso partialing implemented by this command. The fact that `poivregr` reports IV estimates with no instruments retained appears to be an artifact of numerical imprecision in the construction of first-stage residuals computed in the final partialing step.

Figure 6 and Table 5: IV Estimates After Random Forest Partialing

These exhibits use the [Angrist and Evans \(1998\)](#) sample of married women from the 1980 Census. Random forest partialing for Figure 6 uses the Stata `rforest` command discussed in the context of Table 3. Estimates in Table 5 use `rforest` and the `regression_forest` command contained in the Generalized Random Forest (GRF) software package referenced by [Athey et al. \(2019\)](#). Residuals plotted in the figure and used as instruments were computed using leaf sizes indicated in legends and column headings. GRF parameter settings mostly equal to those use by [Athey et al. \(2019\)](#). The number of variables randomly investigated is equal to the square root of the number of right-hand-side variables plus 20, with a subsample rate of 5 percent. Our implementation computes 100 trees. We obtained similar estimates using much larger numbers of trees. GRF `regression_forest` reports leave-out fitted values, as suggested by [Athey et al. \(2019\)](#).

References

- Akerberg, Daniel A and Paul J Devereux. 2009. Improved JIVE estimators for overidentified linear models with and without heteroskedasticity. *The Review of Economics and Statistics* 91 (2):351–362. URL <https://doi.org/10.1162/rest.91.2.351>.
- Ahrens, Achim, Christian B Hansen, and Mark E Schaffer. 2019. Lassopack: Model selection and prediction with regularized regression in Stata. *arXiv preprint arXiv:1901.05397*.
- Angrist, Joshua and William Evans. 1998. Children and their parents' labor supply: Evidence from exogenous variation in family size. *American Economic Review* 88 (3):450–77.
- Angrist, Joshua D and Alan B Krueger. 1991. Does compulsory school attendance affect schooling and earnings? *The Quarterly Journal of Economics* 106 (4):979–1014.
- Athey, Susan, Julie Tibshirani, and Stefan Wager. 2019. Generalized random forests. *The Annals of Statistics* 47 (2):1148–1178. URL <https://doi.org/10.1214/18-AOS1709>.

- Belloni, Alexandre, Victor Chernozhukov, and Christian Hansen. 2014. Inference on treatment effects after selection among high-dimensional controls. *The Review of Economic Studies* 81 (2):608–650.
- Belloni, Alexandre, Victor Chernozhukov, and Christian B Hansen. 2013. Inference for high-dimensional sparse econometric models. In *Advances in economics and econometrics: Tenth world congress of econometric society, volume iii*, eds. Daron Acemoglu, Manuel Arellano, and Eddie Dekel, chap. 7. Cambridge: Cambridge University Press, 245–295.
- Bickel, Peter J, Yaacov Ritov, Alexandre B Tsybakov et al. 2009. Simultaneous analysis of lasso and Dantzig selector. *The Annals of Statistics* 37 (4):1705–1732.
- Chernozhukov, Victor, Christian Hansen, and Martin Spindler. 2015. Post-selection and post-regularization inference in linear models with many controls and instruments. *American Economic Review* 105 (5):486–90.
- Friedman, Jerome, Trevor Hastie, and Robert Tibshirani. 2010. Regularization paths for generalized linear models via coordinate descent. *Journal of Statistical Software* 33 (1):1–22. URL <http://www.jstatsoft.org/v33/i01/>.
- Oshiro, Thais Mayumi, Pedro Santoro Perez, and José Augusto Baranauskas. 2012. How many trees in a random forest? In *International workshop on machine learning and data mining in pattern recognition*. New York: Springer, 154–168.
- Schonlau, Matthias. 2019. RFOREST: Stata module to implement Random Forest algorithm. Working Paper S458614, Boston College Department of Economics. URL <https://ideas.repec.org/c/boc/bocode/s458614.html>.
- Stata. 2019. *Stata lasso reference manual version 16*. College Station, Texas: Stata Press.
- Townsend, Wilbur. 2017. ELASTICREGRESS: Stata module to perform elastic net regression, lasso regression, ridge regression. Working Paper S458397, Boston College Department of Economics. URL <https://ideas.repec.org/c/boc/bocode/s458397.html>.

Table A1: Alternative Post-Lasso Estimates of Elite College Effects

	Double-selection (PDS)				Treatment (score) selection				Outcome selection			
	plug-in (15)	C.V. λ	1 S.E. λ	10xC.V. λ	plug-in (15)	C.V. λ	1 S.E. λ	10xC.V. λ	plug-in (15)	C.V. λ	1 S.E. λ	10xC.V. λ
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)	(10)	(11)	(12)
A. Private School Effects												
Estimate	0.055 (0.041)	0.021 (0.039)	0.023 (0.038)	0.037 (0.041)	0.058 (0.062)	0.021 (0.040)	0.021 (0.039)	0.048 (0.054)	0.205 (0.050)	0.039 (0.042)	0.040 (0.041)	0.035 (0.041)
Treatment residual s.d	0.339	0.304	0.308	0.330	0.339	0.305	0.308	0.332	0.492	0.333	0.335	0.356
No. of controls	6	107	91	27	5	93	86	25	1	34	10	2
B. Effects of School Average SAT/100												
Estimate	-0.014 (0.023)	-0.014 (0.018)	-0.015 (0.019)	-0.014 (0.020)	-0.028 (0.027)	-0.009 (0.020)	-0.009 (0.021)	-0.015 (0.021)	0.111 (0.018)	-0.013 (0.020)	-0.009 (0.021)	-0.018 (0.017)
Treatment residual s.d	0.414	0.390	0.391	0.409	0.415	0.393	0.393	0.409	0.943	0.411	0.414	0.528
No. of controls	15	90	64	15	14	71	61	14	1	34	6	2
C. Effects of Attending Schools Rated Highly Competitive +												
Estimate	0.079 (0.029)	0.057 (0.034)	0.050 (0.034)	0.073 (0.032)	0.076 (0.043)	0.054 (0.034)	0.052 (0.035)	0.075 (0.037)	0.220 (0.036)	0.076 (0.033)	0.079 (0.031)	0.063 (0.030)
Treatment residual s.d	0.342	0.289	0.288	0.330	0.342	0.289	0.288	0.332	0.470	0.335	0.337	0.357
No. of controls	8	106	82	30	7	88	79	28	1	34	6	2

Notes: This table reports estimates computed using alternative lasso routines. Estimators are as described in the note to Table 2. Columns 1, 5, and 9 use the default plug-in penalty implemented in `Lassopack rlasso`; these estimates (like those in Table 2) use regressor-specific penalty loadings. Columns 2, 6, and 10 use a cross-validated penalty, as implemented in the `Elasticregress lassoregress` command (see Townsend, 2018); these estimates standardize regressors but omit regressor-specific penalty loadings. Columns 3, 7, and 11 use the largest penalty such that the cross-validated mean squared error is within one standard error of the minimum cross-validated mean squared error, a variation implemented in `lassoregress`. Hastie, Tibishrani, and Wainwright (2016) suggest this modification. Columns 4, 8, and 12 use a penalty equal to 10 times the cross-validated penalty used for columns 2, 6, and 10.

Table A2: Angrist and Krueger (1991) Random Forest Simulation Results

Estimator	Bias (1)	Standard deviation (2)	Median abs. dev. (3)	Median abs. error (4)
Random forest first stage, 2SLS using RF fits as instruments (min leaf size=1)	0.0611	0.0047	0.0030	0.0612
Random forest 2SLS, min leaf size = 800	0.0567	0.0065	0.0045	0.0567
Random forest first stage, SSIV using RF fits as instruments (min leaf size =1)	-0.0003	0.0158	0.0109	0.0108
Random forest SSIV, min leaf size = 800	-0.0005	0.0158	0.0104	0.0103

Notes: The table describes simulation results for 999 Monte Carlo estimates of the economic returns to schooling using simulated samples constructed from the Angrist and Krueger (1991) census sample of men born 1930-39 (N=329,509). The causal effect of schooling is calibrated to 0.1; the OLS estimand is 0.207. The instruments used to compute the estimates are quarter-of-birth-by-year-of-birth-by-state-of-birth interactions (average F-stat = 1.7, average concentration parameter = 1050). All models include saturated year of birth by state of birth controls. Random forest routines are described in the appendix.

Table A3: Additional Simulation Results for Opening Weekend Effects

Estimator	Original Instruments (F=2.85)					Original plus 52 noise instruments (F=2.06)				
	Avg. IVs retained (1)	Bias (2)	Standard deviation (3)	Median abs. dev. (4)	Median abs. error (5)	Avg. IVs retained (6)	Bias (7)	Standard deviation (8)	Median abs. dev. (9)	Median abs. error (10)
<u>poivregress</u> (fixed # of controls)										
Some IVs selected	1.32	-0.133	0.084	0.053	0.137	1.26	-0.145	0.077	0.050	0.143
No IVs selected*	0	-0.287	0.360	0.198	0.316	0	-0.291	0.352	0.189	0.325
<u>poivregress</u> (high-dim controls)										
Some IVs selected (15-16 ctls retained)	1.22	-0.096	0.108	0.069	0.107	1.22	-0.098	0.104	0.068	0.106
No IVs selected** (49-50 ctls retained)	0	-1.58	18.5	0.354	0.477	0	-1.90	19.5	0.375	0.477

Notes: The table reports simulation results for 999 Monte Carlo estimates of the effect of opening weekend ticket sales on second weekend ticket sales using simulated samples constructed from the data used by Gilchrist and Sands (2016) (N=1,671). The causal effect of interest is calibrated to 0.6. Columns 1-5 show results using the original instruments. Columns 6-10 report the results of adding 52 randomly generated (standard uniform) instruments to the original 52-instrument dictionary. Estimates produced using Stata's `poivregress` command which partials covariates out of the outcome, endogenous regressor, and instruments, and uses lasso to select instruments.

*`poivregress` reports estimates with zero instruments selected in 640 of 999 runs for the original 52-instrument set and in 741 of 999 runs using the 104 instrument set.

**`poivregress` with high-dimensional controls reports estimates with zero instruments selected in 141 of 999 runs for the original 52-instrument set and in 132 of 999 runs using the 104 instrument set. In 372 runs with 52 instruments and 409 runs with 104 instruments, this version of `poivregress` selects zero instruments and reports an estimate of zero.