

Seminar 10: Machine Learning for Microeconometrics

A Monte Carlo exercise: OLS vs Naive Lasso vs Double Selection LASSO

We want to check how the different methods work.

Setup of the Monte Carlo Simulation

The following simulation is a slightly revised version of the simulation from S. Kranz "Lasso and the Methods of Causality" (<https://skrznz.github.io/2020/09/14/LassoCausality.html>)

We study different approaches for variable selection using a Monte-Carlo simulation with the following data generating process for

$$y = 1 + \alpha d + \sum_{k=1}^{K_C} \beta_k^* x_k^* + \sum_{k=1}^{K_C} \beta_k^* x_k^* + \varepsilon^y$$
$$d = 1 + \sum_{k=1}^{K_C} \beta_k^* x_k^* + \sum_{k=1}^{K_C} \beta_k^* x_k^* + \varepsilon^d$$

We have the following potential control variables that are all independently normally distributed from each other:

- d the endogenous variable
- x_k^* is one of $K_C=10$ confounders that directly affect both d and y . For a consistent OLS estimation of α , we need to control for all confounders.
- x_k^* is one of $K_y=15$ variables that only affect the dependent variable y but not the explanatory variable d . Whether we would add it or not in an OLS regression should not affect the bias of our estimator $\hat{\alpha}$.
- x_k^* is one of $K_C=15$ variables that only affect d but not through another channel the dependent variable y . It constitutes a source of exogenous variation in d . We can estimate in an OLS regression α more precisely if we don't add any x_k^* to the regression. Also, we will see that adding fewer x_k^* can reduce the bias of an OLS estimator that arises if we have not perfectly controlled for all confounders.
- We also observe $K_u=15$ variables x_k^* that neither affect n nor d . They are just uncorrelated noise variables that we ideally leave out of our regressions.

The following code simulates a data set with $n = 200$ observations, $K_C = 15$ confounders, $K_y = 15$ variables that affect only y , $K_u = 2$ that provides a source of exogenous variation and $K_u = 20$ explanatory variables that are uncorrelated with everything else. The causal effect of interest α of our regression coefficients and standard deviations are equal to 1.

```
In [1]: import stata_setup
stata_setup.config("C:/Program Files/Stata17", "se")

In [2]: %*data
clear all
set seed 20211109
set lsize 255
set obs 200

forvalves k=1/232(
    g v k'=normal(0,1)
)

egen sum_u=rowtotal(v1-v15)
egen sum_y=rowtotal(v16-v30)
egen sum_d=rowtotal(v31-v32)
g d=1+sum_u+sum_y+normal(1)
g v=d+sum_u+sum_y+normal(1)

OLS
```

We look at five different cases.

- kitchen sink specification (throw in all potential variables)
- Omitted Variable Bias
- Data Generating Process
- controlling variables that correlated with D but not impact y via other channels (valid IV for D)
- IV to correct OMB

```
In [4]: %*data
quietly reg y d v*
estimates store r1, title(kitchen)

quietly reg y d v1-v30
estimates store r2, title(OVB)

quietly reg y d v1-v30
estimates store r3, title(DGP)

quietly reg y d v1-v32
estimate store r4, title(D_e)

quietly lvarregss beta y v16-v30 (d=v31-v32)
estimates store r5, title(IV)

estout r*, keep(d) cells(b se(par fml(2))) ///
    legend label variabels=(cons Constant)

. quietly reg y d v*

. estimates store r1, title(kitchen)

.

. quietly

. estout r*, keep(d) cells(b se(par fml(2))) ///
> legend label variabels=(cons Constant)
```

	kitchen	OVB	DGP	D_e	IV
d	1.97646	1.494306	1.9467909	.9702836	1.8929811
	(.)	(0.07)	(0.05)	(0.07)	(0.44)

LASSO, Post-Lasso and Post Double Selection

In the following we estimate the relationship between y and predictors (d and V^*). The goal is to recover the causal parameter α when there are high numbers of potential controls. (NOTE: we assume that CIA holds.)

As we discussed in the lecture, the lasso regression selects a subset of explanatory variables whose estimated coefficients are non-zero. But also the coefficients of the selected variables will be typically attenuated towards 0 because of the penalty term. The post-lasso estimator avoids this attenuation by simply performing an OLS estimation using all the selected variables from the lasso estimation.

In the following, we use the stata package "lassopack" and "pdslasso".

LASSO and Post-Lasso

We use three different procedures to determine the penalty parameter λ .

1. cross validation (cvlasso)
2. using information criterion (lasso2)
3. Theoretical value (lasso)

```
In [5]: %*cvlasso
cvlasso y d v*, lopt postres plotcv

. cvlasso y d v*, lopt postres plotcv

K-fold cross-validation with 10 folds. Elastic net with alpha=.
Fold 1 2 3 4 5 6 7 8 9 10
-----
Lambdas MSFE st. dev.
1) 2875.711 70.645657 11.235848
2) 2620.2408 64.216446 11.392026
3) 2387.4658 56.658243 10.001639
4) 2175.3699 50.379822 8.8255175
5) 207.62640 45.164359 8.506878
6) 1806.0303 40.831648 6.9890587
7) 181.78741 37.223651 6.9890587
8) 1499.3981 34.241341 5.6749487
9) 1366.1957 31.756322 5.1659934
10) 1244.8266 29.691325 3.7345486
11) 1134.2397 27.975208 3.126884
12) 1033.477 26.548899 4.0674591
13) 941.46571 25.363304 3.808867
14) 858.01072 24.377709 3.5930362
15) 781.78741 23.563274 3.4388628
16) 712.33557 22.878668 3.257425
17) 649.03364 22.310172 3.126884
18) 591.93955 21.838448 2.7958673
19) 538.85574 21.449187 2.2896973
20) 490.98528 20.978089 2.8658107
21) 447.3675 20.001976 2.151707
22) 407.6246 18.85776 2.602543
23) 371.41236 17.759329 2.339319
24) 338.41711 16.762517 2.2863742
25) 280.33307 15.607381 2.1379607
26) 240.95983 16.008302 1.9386623
27) 256.00013 13.347922 1.7532399
28) 233.22779 12.165617 1.3734637
29) 212.53581 10.839476 1.1059256
30) 193.65471 9.723153 1.2727875
31) 176.47868 8.673214 1.4093862
32) 167.7554 7.6728737 95771969
33) 146.49269 6.8423937 82195163
34) 130.47868 6.154876 70718663
35) 121.6208 5.576866 61325231
36) 110.81635 5.099705 53452071
37) 100.97173 4.703839 47054332
38) 92.001677 4.3705292 41351308
39) 83.828502 4.094482 35467339
40) 76.381409 3.8508632 32494697
41) 69.593895 3.6448411 29419368
42) 63.43187 3.4620331 27295823
43) 57.779735 3.3231031 25195832
44) 52.646742 3.1990387 23413112
45) 47.367931 3.0823292 22083399
46) 43.70825 2.9914554 21537341
47) 40.143253 2.904462 20616781
48) 36.287357 2.8407593 20120308
49) 33.063688 2.7776465 19429731
50) 30.124601 2.6888654 18703396
51) 27.450055 2.6228881 18303298
52) 25.011468 2.6734949 18288842
53) 22.78318 2.680161 18288842
54) 20.76496 2.6842499 18826595
55) 18.902538 2.6949529 18698781
56) 17.239435 2.7429811 19776096
57) 15.707931 2.8041509 20779702
58) 14.312893 2.8596429 21616314
59) 13.041 2.9667291 23929786
60) 11.882474 3.062136 2561369
61) 10.828867 3.189131 2759459
62) 9.860383 3.260437 2979459
63) 8.988655 3.429839 32233398
64) 8.1301281 3.4294875 33857761
65) 7.4623397 3.5160464 35482791
66) 6.7939882 3.618356 37936623
67) 6.195315 3.7180339 38265628
68) 5.6451376 3.8202545 40139705
69) 5.1435391 3.9400839 42233398
70) 4.6866924 4.0623193 44629588
71) 4.2703396 4.1980939 46516769
72) 3.8909744 4.2891027 48256248
73) 3.5453109 4.4222282 50055491
74) 3.230353 4.539579 52067821
75) 2.9433794 4.7444922 55319404
76) 2.6818977 4.9204113 58626731
77) 2.4325533 5.0942629 61616314
78) 2.2265886 5.254475 64620122
```

* lopt = the lambda that minimizes MSFE.
Run model: cvlasso, lopt
* lse = largest lambda for which MSFE is within one standard error of the minimal MSFE.

Estimate Lasso with lambda=25.011 (lopt).

Selected	Lasso	Post-est OLS
d	1.7671647	1.6330105
V2	0.2107179	0.4018278
V3	0.1498505	0.2391398
V4	0.161942	0.2587519
V5	0.1262491	0.3053922
V6	0.089593	0.2522352
V7	0.089593	0.2522352
V8	0.089593	0.2522352
V9	0.089593	0.2522352
V10	0.089593	0.2522352
V11	0.089593	0.2522352
V12	0.089593	0.2522352
V13	0.089593	0.2522352
V14	0.089593	0.2522352
V15	0.089593	0.2522352
V16	0.089593	0.2522352
V17	0.089593	0.2522352
V18	0.089593	0.2522352
V19	0.089593	0.2522352
V20	0.089593	0.2522352
V21	0.089593	0.2522352
V22	0.089593	0.2522352
V23	0.089593	0.2522352
V24	0.089593	0.2522352
V25	0.089593	0.2522352
V26	0.089593	0.2522352
V27	0.089593	0.2522352
V28	0.089593	0.2522352
V29	0.089593	0.2522352
V30	0.089593	0.2522352
V31	0.089593	0.2522352
V32	0.089593	0.2522352
V33	0.089593	0.2522352
V34	0.089593	0.2522352
V35	0.089593	0.2522352
V36	0.089593	0.2522352
V37	0.089593	0.2522352
V38	0.089593	0.2522352
V39	0.089593	0.2522352
V40	0.089593	0.2522352
V41	0.089593	0.2522352
V42	0.089593	0.2522352
V43	0.089593	0.2522352
V44	0.089593	0.2522352
V45	0.089593	0.2522352
V46	0.089593	0.2522352
V47	0.089593	0.2522352
V48	0.089593	0.2522352
V49	0.089593	0.2522352
V50	0.089593	0.2522352
V51	0.089593	0.2522352
V52	0.089593	0.2522352
V53	0.089593	0.2522352
V54	0.089593	0.2522352
V55	0.089593	0.2522352
V56	0.089593	0.2522352
V57	0.089593	0.2522352
V58	0.089593	0.2522352
V59	0.089593	0.2522352
V60	0.089593	0.2522352
V61	0.089593	0.2522352
V62	0.089593	0.2522352
V63	0.089593	0.2522352
V64	0.089593	0.2522352
V65	0.089593	0.2522352
V66	0.089593	0.2522352
V67	0.089593	0.2522352
V68	0.089593	0.2522352
V69	0.089593	0.2522352
V70	0.089593	0.2522352
V71	0.089593	0.2522352
V72	0.089593	0.2522352
V73	0.089593	0.2522352
V74	0.089593	0.2522352
V75	0.089593	0.2522352
V76	0.089593	0.2522352
V77	0.089593	0.2522352
V78	0.089593	0.2522352
V79	0.089593	0.2522352
V80	0.089593	0.2522352
V81	0.089593	0.2522352
V82	0.089593	0.2522352
V83	0.089593	0.2522352
V84	0.089593	0.2522352
V85	0.089593	0.2522352
V86	0.089593	0.2522352
V87	0.089593	0.2522352
V88	0.089593	0.2522352
V89	0.089593	0.2522352
V90	0.089593	0.2522352
V91	0.089593	0.2522352
V92	0.089593	0.2522352
V93	0.089593	0.2522352
V94	0.089593	0.2522352
V95	0.089593	0.2522352
V96	0.089593	0.2522352
V97	0.089593	0.2522352
V98	0.089593	0.2522352
V99	0.089593	0.2522352
V100	0.089593	0.2522352
V101	0.089593	0.2522352
V102	0.089593	0.2522352
V103	0.089593	0.2522352
V104	0.089593	0.2522352
V105	0.089593	0.2522352
V106	0.089593	0.2522352
V107	0.089593	0.2522352
V108	0.089593	0.2522352
V109	0.089593	0.2522352
V110	0.089593	0.2522352
V111	0.089593	0.2522352
V112	0.089593	0.2522352
V113	0.089593	0.2522352
V114	0.089593	0.2522352
V115	0.089593	0.2522352
V116	0.089593	0.2522352
V117	0.089593	0.2522352
V118	0.089593	0.2522352
V119	0.089593	0.2522352
V120	0.089593	0.2522352
V121	0.089593	0.2522352
V122	0.089593	0.2522352
V123	0.089593	0.2522352
V124	0.089593	0.2522352
V125	0.089593	0.2522352
V126	0.089593	0.2522352
V127	0.089593	0.2522352
V128	0.089593	0.2522352
V129	0.089593	0.2522352
V130	0.089593	0.2522352
V131	0.089593	0.2522352
V132	0.089593	0.2522352
V133	0.089593	0.2522352
V134	0.089593	0.2522352
V135	0.089593	0.2522352
V136	0.089593	0.2522352
V137	0.089593	0.2522352
V138	0.089593	0.2522352
V139	0.089593	0.2522352
V140	0.089593	0.2522352
V141	0.089593	0.2522352
V142	0.089593	0.2522352
V143	0.089593	0.2522352
V144	0.089593	0.2522352
V145	0.089593	0.2522352
V146	0.089593	0.2522352
V147	0.089593	0.2522352
V148	0.089593	0.2522352
V149	0.089593	0.2522352
V150	0.089593	0.2522352
V151	0.089593	0.2522352
V152	0.089593	0.2522352
V153	0.089593	0.2522352
V154	0.089593	0.2522352
V155	0.089593	0.2522352
V156	0.089593	0.2522352
V157	0.089593	0.2522352
V158	0.089593	0.2522352
V159	0.089593	0.2522352
V160	0.089593	0.2522352
V161	0.089593	0.2522352
V162	0.089593	0.2522352
V163	0.089593	0.2522352
V164	0.089593	0.2522352
V165	0.089593	0.2522352
V166	0.089593	0.2522352
V167	0.089593	0.2522352
V168	0.089593	0.2522352
V169	0.089593	0.2522352
V170	0.089593	0.2522352
V171	0.089593	0.2522352
V172	0.089593	0.2522352
V173	0.089593	0.2522352
V174	0.089593	0.2522352
V175	0.089593	0.2522352
V176	0.089593	0.2522352
V177	0.089593	0.2522352
V178	0.089593	0.2522352
V179	0.089593	0.2522352
V180	0.089593	0.2522352
V181	0.089593	0.2522352
V182	0.089593	0.2522352
V183	0.089593	0.2522352
V184	0.089593	0.2522352
V185	0.089593	0.2522352
V186	0.089593	0.2522352
V187	0.089593	0.2522352
V188	0.089593	0.2522352


```

capture drop d_hat

. quietly lasso2 d V1-V32, l1c(aic) postres

. predict d_hat, xb
Use e(b) from previous lasso2 estimation (lambda=152.4505926358)

. reg y V1-V30 d_hat

-----
Source |      SS      df       MS      Number of obs
-----+-----+-----+-----+-----
Model | 6921.1918   16  432.57482   F(15, 184)
Residual | 7242.2108   183   39.57492   Prob > F
-----+-----+-----+-----+-----
Total | 14163.4092   199 71.1729105   R-squared
                                           Root MSE
-----+-----+-----+-----+-----

y | Coefficient      Std. err.      z    P>|z|    [95%
-----+-----+-----+-----+-----
V16 | 1.68601         .660915   2.56   0.000   .7766
V17 | 635924         .4570133   1.83   0.069   .505
V18 | .915455         .489502   1.04   0.042   .303
V19 | .283935         .4534881   0.53   0.600   .6347
V20 | 1.02539         .4751977   2.16   0.032   .087
V21 | 1.14743         .6609464   2.49   0.014   .2379
V22 | .4507072        .60582   0.9   0.322   .4516
V23 | .468789         .4448258   1.05   0.295   .4108
V24 | -.878771        .4510332   -1.95   0.053   .014
V25 | .4082248        .451781   0.89   0.373   .4937
V26 | .923806         .466177   1.98   0.049   .020
V27 | 1.066985        .457428   2.33   0.021   .1644
V28 | .7401802        .441445   1.68   0.095   .320
V29 | .931296         .4436147   2.10   0.037   .0560
V30 | .1744806        .4518461   0.39   0.700   .7170
d_hat | 1.001819        .8804784   1.12   0.262   .2311
-----+-----+-----+-----+-----
_cons | -1.79978        .8527855   -2.06   0.003   -.2908

. ivregress 2sls y V16-V30 (d=V1-V32)

Instrumental variables 2SLS regression

Number of obs = 185
Hald chi2(16) = 1.21
Prob > chi2 = .928
R-squared = 0.331
Root MSE = .4229

-----
y | Coefficient      Std. err.      z    P>|z|    [95%
-----+-----+-----+-----+-----
V1 | .889291         .3963259   2.24   0.027   .029
V16 | .9647399        .3369955   2.83   0.015   .1889
V17 | .6057348        .406781   1.48   0.139   .206
V18 | .7173051        .3352889   2.14   0.034   .037
V19 | .5365854        .4021282   1.28   0.200   .2884
V20 | 1.13304         .4079193   2.78   0.007   .0385
V21 | 1.527868        .434471   3.52   0.000   .6763
V22 | .8292419        .3993792   2.08   0.038   .0646
V23 | .8808632        .3897153   2.25   0.027   .029
V24 | .9584635        .3940326   2.43   0.015   .1881
V25 | -.2622715       .4046056   -0.65   0.515   .4645
V26 | .4607674        .442242   1.04   0.309   .4589
V27 | 1.426428        .4095524   3.48   0.001   .4229
V28 | .6518562        .4078193   1.60   0.110   .3945
V29 | .549518         .4058813   1.35   0.176   .2491
V30 | .4970616        .439328   1.13   0.258   .3640
d_hat | 1.612411        .8845886   1.75   0.006   .4229

Instrumented: d_hat
Instrumentals: V16 V17 V18 V19 V20 V21 V22 V23 V24 V25 V26 V27 V28 V29 V30

. drop d_hat

. qui lasso2 d V1-V32, l1c(aic) postres

. predict d_hat, xb

. reg y V1-V30 d_hat
note: d_hat omitted because of collinearity.

-----
Source |      SS      df       MS      Number of obs
-----+-----+-----+-----+-----
Model | 2366.1758     15  157.78019   F(15, 184)
Residual | 11796.6938    184  64.1124669   Prob > F
-----+-----+-----+-----+-----
Total | 14163.4092    199 71.1729105   R-squared
                                           Root MSE
-----+-----+-----+-----+-----

y | Coefficient      Std. err.      z    P>|z|    [95%
-----+-----+-----+-----+-----
V16 | .8903259        .5790085   1.54   0.126   .225

```

```

V17      4529499      579905    0.78   0.436   - .691
V18      6068135      615201    0.99   0.326   - .676
V19      334905      570884    0.58   0.563   - .8047
V20      1.256609      604201    0.20   0.039   - .3643
V21      1.770703      580515    0.34   0.126   - .2301
V22      7699611      5850536    1.32   0.190   - .0843
V23      392334      5661102    0.09   0.489   - .7248
V24      4846222      575926   0.14   0.003   - .822
V25      -8305067      5630318   -1.48   0.142   - .1341
V26      3559354      5895141   0.60   0.547   - .00717
V27      1.583662      5789728   0.24   0.007   - .4416
V28      4561087      5686807    0.81   0.417   - .6504
V29      5934885      5405436   0.64   0.522   - .7464
V30      2264722      575078   0.39   0.694   - .1951
d_hat      0 (omitted)
Cons      1.98787      578513   3.80   0.000   1.0061

.
. drop d_hat

. qui clvarso d v31-v32, lopt: postures
Warning: the value is at the limit of the lambda range.

. predict d_hat, xb
Use eby from previous lasso2 estimation (lambda=152.450562938)

. reg y v16-v30 d_hat

Source      SS          df           MS      Number of obs
-----
Model      6921.1818      16   432.574882   R-squared
Residual   2422.1108      183   13.254784   F Prob > F
Total     14143.4092      199   71.172905   Adj R-squared
               Root MSE

y | Coefficient      Std. err.      t P>|z|      [95%
-----+-----
V16      1.68601      4609155    3.66   0.000   .7766
V17      8339242      4570313    1.83   0.069   -.0657
V18      9931405      4570313    2.18   0.032   -.0522
V19      2389375      4534881    5.33   0.000   -.6583
V20      1.025299      4751977    2.16   0.032   -.0877
V21      1.14743      4603464    2.49   0.014   -.2375
V22      4570792      4605982   9.90   0.322   -.4516
V23      4667879      4448285    1.05   0.295   -.4398
V24      8787      4538023    0.02   0.983   -.0271
V25      4008248      4571781   0.89   0.373   -.4307
V26      923036      4646777    1.99   0.049   -.0040
V27      1.066985      4574228    2.33   0.021   -.1644
V28      7401802      4451425    1.68   0.095   -.1331
V29      9322917      4538023    2.06   0.041   -.0560
V30      1740886      4518461    3.90   0.700   -.7370
d_hat      0.461889      3804794    10.73   0.000   .31301
Cons      -1.789789      4878758   -1.06   0.003   -2.9388

The Vlasso and the kitchen sink version of IV

In [2]: %>>> vlasso2

vlasso2 zeta y v16-v30 (d=v31-v232)
* matrix list e(b)

vlasso2 y v16-v20 (d=v31-v232)
* matrix list e(b)

. ivregress 2sls y v16-v30 (d=v31-V232)

note: V230 omitted because of collinearity.
note: V231 omitted because of collinearity.
note: V232 omitted because of collinearity.
note: V233 omitted because of collinearity.
note: V234 omitted because of collinearity.
note: V235 omitted because of collinearity.
note: V236 omitted because of collinearity.
note: V237 omitted because of collinearity.
note: V238 omitted because of collinearity.
note: V239 omitted because of collinearity.
note: V240 omitted because of collinearity.
note: V241 omitted because of collinearity.
note: V242 omitted because of collinearity.
note: V243 omitted because of collinearity.
note: V244 omitted because of collinearity.
note: V245 omitted because of collinearity.
note: V246 omitted because of collinearity.
note: V247 omitted because of collinearity.
note: V248 omitted because of collinearity.
note: V249 omitted because of collinearity.
note: V250 omitted because of collinearity.
note: V251 omitted because of collinearity.
note: V252 omitted because of collinearity.
note: V253 omitted because of collinearity.
note: V254 omitted because of collinearity.
note: V255 omitted because of collinearity.
note: V256 omitted because of collinearity.
note: V257 omitted because of collinearity.
note: V258 omitted because of collinearity.
note: V259 omitted because of collinearity.
note: V260 omitted because of collinearity.
note: V261 omitted because of collinearity.
note: V262 omitted because of collinearity.
note: V263 omitted because of collinearity.
note: V264 omitted because of collinearity.
note: V265 omitted because of collinearity.
note: V266 omitted because of collinearity.
note: V267 omitted because of collinearity.
note: V268 omitted because of collinearity.
note: V269 omitted because of collinearity.
note: V270 omitted because of collinearity.
note: V271 omitted because of collinearity.
note: V272 omitted because of collinearity.
note: V273 omitted because of collinearity.
note: V274 omitted because of collinearity.
note: V275 omitted because of collinearity.
note: V276 omitted because of collinearity.
note: V277 omitted because of collinearity.
note: V278 omitted because of collinearity.
note: V279 omitted because of collinearity.
note: V280 omitted because of collinearity.
note: V281 omitted because of collinearity.
note: V282 omitted because of collinearity.
note: V283 omitted because of collinearity.
note: V284 omitted because of collinearity.
note: V285 omitted because of collinearity.
note: V286 omitted because of collinearity.
note: V287 omitted because of collinearity.
note: V288 omitted because of collinearity.
note: V289 omitted because of collinearity.
note: V290 omitted because of collinearity.
note: V291 omitted because of collinearity.
note: V292 omitted because of collinearity.
note: V293 omitted because of collinearity.
note: V294 omitted because of collinearity.
note: V295 omitted because of collinearity.
note: V296 omitted because of collinearity.
note: V297 omitted because of collinearity.
note: V298 omitted because of collinearity.
note: V299 omitted because of collinearity.
note: V300 omitted because of collinearity.
note: V301 omitted because of collinearity.
note: V302 omitted because of collinearity.
note: V303 omitted because of collinearity.
note: V304 omitted because of collinearity.
note: V305 omitted because of collinearity.
note: V306 omitted because of collinearity.
note: V307 omitted because of collinearity.
note: V308 omitted because of collinearity.
note: V309 omitted because of collinearity.
note: V310 omitted because of collinearity.
note: V311 omitted because of collinearity.
note: V312 omitted because of collinearity.
note: V313 omitted because of collinearity.
note: V314 omitted because of collinearity.
note: V315 omitted because of collinearity.
note: V316 omitted because of collinearity.
note: V317 omitted because of collinearity.
note: V318 omitted because of collinearity.
note: V319 omitted because of collinearity.
note: V320 omitted because of collinearity.
note: V321 omitted because of collinearity.
note: V322 omitted because of collinearity.
note: V323 omitted because of collinearity.
note: V324 omitted because of collinearity.
note: V325 omitted because of collinearity.
note: V326 omitted because of collinearity.
note: V327 omitted because of collinearity.
note: V328 omitted because of collinearity.
note: V329 omitted because of collinearity.
note: V330 omitted because of collinearity.
note: V331 omitted because of collinearity.
note: V332 omitted because of collinearity.
note: V333 omitted because of collinearity.
note: V334 omitted because of collinearity.
note: V335 omitted because of collinearity.
note: V336 omitted because of collinearity.
note: V337 omitted because of collinearity.
note: V338 omitted because of collinearity.
note: V339 omitted because of collinearity.
note: V340 omitted because of collinearity.
note: V341 omitted because of collinearity.
note: V342 omitted because of collinearity.
note: V343 omitted because of collinearity.
note: V344 omitted because of collinearity.
note: V345 omitted because of collinearity.
note: V346 omitted because of collinearity.
note: V347 omitted because of collinearity.
note: V348 omitted because of collinearity.
note: V349 omitted because of collinearity.
note: V350 omitted because of collinearity.
note: V351 omitted because of collinearity.
note: V352 omitted because of collinearity.
note: V353 omitted because of collinearity.
note: V354 omitted because of collinearity.
note: V355 omitted because of collinearity.
note: V356 omitted because of collinearity.
note: V357 omitted because of collinearity.
note: V358 omitted because of collinearity.
note: V359 omitted because of collinearity.
note: V360 omitted because of collinearity.
note: V361 omitted because of collinearity.
note: V362 omitted because of collinearity.
note: V363 omitted because of collinearity.
note: V364 omitted because of collinearity.
note: V365 omitted because of collinearity.
note: V366 omitted because of collinearity.
note: V367 omitted because of collinearity.
note: V368 omitted because of collinearity.
note: V369 omitted because of collinearity.
note: V370 omitted because of collinearity.
note: V371 omitted because of collinearity.
note: V372 omitted because of collinearity.
note: V373 omitted because of collinearity.
note: V374 omitted because of collinearity.
note: V375 omitted because of collinearity.
note: V376 omitted because of collinearity.
note: V377 omitted because of collinearity.
note: V378 omitted because of collinearity.
note: V379 omitted because of collinearity.
note: V380 omitted because of collinearity.
note: V381 omitted because of collinearity.
note: V382 omitted because of collinearity.
note: V383 omitted because of collinearity.
note: V384 omitted because of collinearity.
note: V385 omitted because of collinearity.
note: V386 omitted because of collinearity.
note: V387 omitted because of collinearity.
note: V388 omitted because of collinearity.
note: V389 omitted because of collinearity.
note: V390 omitted because of collinearity.
note: V391 omitted because of collinearity.
note: V392 omitted because of collinearity.
note: V393 omitted because of collinearity.
note: V394 omitted because of collinearity.
note: V395 omitted because of collinearity.
note: V396 omitted because of collinearity.
note: V397 omitted because of collinearity.
note: V398 omitted because of collinearity.
note: V399 omitted because of collinearity.
note: V400 omitted
```

```

V28 | 1.191967 | .1629848 | 7.31 | 0.000 | .8728 |
V29 | 1.073892 | .166811 | 6.44 | 0.000 | .7459 |
V30 | 1.243678 | .1738472 | 7.15 | 0.000 | .9029 |
IV using ChR post-lasso-orthogonalized vars
-----
y | Coefficient | Std. err. | z | P>|z| | [95% |
-----+-----+-----+-----+-----+-----
V16 | 1.189087 | .1635883 | 7.15 | 0.000 | .8484 |
V17 | 1.028709 | .1663215 | 6.19 | 0.000 | .7037 |
V18 | 1.225169 | .1778933 | 6.88 | 0.000 | .8763 |
V19 | 1.102003 | .1688166 | 6.53 | 0.000 | .7711 |
V20 | .920998 | .171461 | 4.62 | 0.000 | .4560 |
V21 | .837832 | .174585 | 5.00 | 0.000 | .5037 |
V22 | .9928107 | .1629563 | 6.09 | 0.000 | .6734 |
V23 | .725135 | .1608674 | 4.51 | 0.000 | .4098 |
V24 | 1.146276 | .1587704 | 7.32 | 0.000 | .8515 |
V25 | 1.302616 | .1713506 | 6.01 | 0.000 | .8799 |
V26 | 1.306311 | .1784872 | 7.32 | 0.000 | .9564 |
V27 | .992584 | .1678479 | 5.91 | 0.000 | .6636 |
V28 | 1.191967 | .1629848 | 7.31 | 0.000 | .8728 |
V29 | 1.073892 | .166811 | 6.44 | 0.000 | .7459 |
V30 | 1.243678 | .1738472 | 7.15 | 0.000 | .9029 |
IV with PDS-selected variables and full regressor set
-----
y | Coefficient | Std. err. | z | P>|z| | [95% |
-----+-----+-----+-----+-----+-----
d | 0 | (omitted) |
V16 | 0 | (omitted) |
V17 | 0 | (omitted) |
V18 | 0 | (omitted) |
V19 | 0 | (omitted) |
V20 | 0 | (omitted) |
V21 | 0 | (omitted) |
V22 | 0 | (omitted) |
V23 | 0 | (omitted) |
V24 | 0 | (omitted) |
V25 | 0 | (omitted) |
V26 | 0 | (omitted) |
V27 | 0 | (omitted) |
V28 | 0 | (omitted) |
V29 | 0 | (omitted) |
V30 | 0 | (omitted) |
V31 | 0 | (omitted) |
-----
Standard errors and test statistics valid for the following vars:
d V16 V17 V18 V19 V20 V21 V22 V23 V24 V25 V26 V27 V28 V29 V30
Warning: not enough instruments selected: model unidentified

. * matrix list e(b)
Unknown command

Distribution of estimators

Run the simulation for 100 times and look at the distribution of the different estimators

In [1]: %data
clear all
set seed 20211
capture postout1 clear
tempfile simresults
postfile sim_mem OVB OLS LASSO2 PDS DS P_out IV IVlasso using
forvalues i=1/100 {
    quietly set obs 5000
    di `i'
    forvalues k=1/232{
        g V`k'=rnormal(0,1)
    }
    egen sum_c=rowtotal(V1-V15)
    egen sum_y=rowtotal(V16-V30)
    egen c=rowtotal(V31-V32)
    g d=1+sum_c*sum_e+normal()
    g y=1+d*sum_c*sum_y+normal()
    * OVB
    quietly reg y d
    mat R= e(b)
    * DGR OLS
    quietly reg y d V1-V30
    mat A= e(b)
    *lasso2
    quietly lasso2 y d V*, l1c(aicc) postres
    mat B=e(beta)
    ** own post double selection using lasso2
    quietly lasso2 y V*, l1c(aicc) postres
    local string1 = e(selected)
    quietly lasso2 d V*, l1c(aicc) postres
    local string2 = e(selected)
    local newlist: list string1 | string2
    * di `newlist'
    quietly reg y d `newlist'
    mat F= e(b)
    ** packaged
    quietly pldlasso2 y d (V*)
    mat D=e(beta_pds)
    ** Partialling out
    quietly reg y `string1'
    predict y_dot, res
    quietly reg d `string2'
    predict d_dot, res
    quietly reg y_dot d_dot
    mat G=e(b)
    ** IV
    quietly ivregress 2sls y V16-V30 (d=V31-V232)
    mat IV=e(b)
    ** IVlasso
    quietly ivlasso2 y V16-V30 (d=V31-V232)
    mat IV2=e(b)
    post sim_mem (A[1,1]) (A[1,1]) (B[1,1]) (D[1,1]) (F[1,1])
    capture drop y d V* sum* d_dot
}
postclose sim_mem
use `simresults', clear
save sim1_15000.dta, replace

case I: 200 obs vs 5000 obs, all parameters equal to 1.


$$y = ad + \sum_{k=1}^K \beta_k^O x_k^O + \sum_{k=1}^K \beta_k^I x_k^I$$


$$d = \sum_{k=1}^K \beta_k^O x_k^O + \sum_{k=1}^K \beta_k^I x_k^I$$


In [13]: %data
clear all
use sim1_1_200.dta, clear
su
qui {
    twoway hist OLS, xline(1.0) name(g)
    twoway hist PDS, xline(1.0) name(gb)
    twoway hist LASSO2, xline(1.0) name(gc)
    twoway hist DS, xline(1.0) name(gd)
}
graph combine gb qa qc qd, xcommon
. clear all
.
. use sim1_1_200.dta, clear

```

```

. su
+-----+-----+
| Variable | Obs   | Mean   | Std. dev. | Min   |
+-----+-----+
| OVB | 100 | 1.850334 | .0715134 | 1.689236 |
| OLS | 100 | 1.003209 | .0376376 | .8977358 |
| LASSO2 | 100 | 1.824524 | .0382321 | 1.735735 |
| PDS | 100 | 1.659872 | .218394 | .0774625 |
| DS | 100 | .8869708 | .1114914 | .6057009 |
+-----+-----+
| P_out | 100 | .7294227 | .1274495 | .4364092 |
| IV | 100 | 1.834283 | .0352793 | 1.754402 |
| IVlasso | 100 | .5022424 | .6092893 | 0 |
+-----+-----+

. qui

. graph combine gb ga gc gd, xcommon

Density
15
10
5
0
.9 .95 OLS 1.05

Density
2
1
0
.5 1 PDS 1.5 2

Density
10
5
0
1.75 1.8 LASSO2 1.85 1.9

Density
4
3
2
1
0
.6 .8 DS 1 1.2

Density
2
1.5
1
0.5
0
.5 1 PDS 1.5 2

Density
15
10
5
0
.5 1 LASSO2 1.5 2

Density
4
3
2
1
0
.5 1 DS 1.5 2

In [14]: %stata
clear all

use sim1_1_5000.dta, clear

su

qui {
    twoway hist OLS, xline(1.0) name(ga)
    twoway hist PDS, xline(1.0) name(gb)
    twoway hist LASSO2, xline(1.0) name(gc)
    twoway hist DS, xline(1.0) name(gd)
}

graph combine gb ga gc gd, xcommon

. clear all

. use sim1_1_5000.dta, clear

. su
+-----+-----+
| Variable | Obs   | Mean   | Std. dev. | Min   |
+-----+-----+
| OVB | 100 | 1.834414 | .0156417 | 1.785204 |
| OLS | 100 | .9959777 | .0076402 | .9836014 |
| LASSO2 | 100 | 1.074105 | .0227339 | 1.026595 |
| PDS | 100 | 1.00037 | .0148173 | .9543703 |
| DS | 100 | .9952539 | .0152882 | .9499198 |
+-----+-----+
| P_out | 100 | .9895131 | .0160422 | .9427422 |
| IV | 100 | 1.228091 | .0270834 | 1.1654 |
| IVlasso | 100 | .9954752 | .687869 | .9373788 |
+-----+-----+

```

`graph combine gb ga gc gd, xcommon`

case II: 200 obs vs 5000 obs, $\beta_k^{od} = 10$, all other parameters equal

$$y = 1 + \alpha d + \sum_{k=1}^{K_1} \beta_k^{mu} x_k + \sum_{k=1}^{K_2} \beta_k^{\sigma} x_k$$

$$d = 1 + \sum_{k=1}^{K_1} 10x_k^+ + \sum_{k=1}^{K_2} \beta_k^{\sigma} x_k$$

```

In [15]: %!stata
clear all
use sim1_10_200.dta, clear
su
twoway hist OLS, xline(1.0) name(ga)
twoway hist PDS, xline(1.0) name(gb)
twoway hist LASSO2, xline(1.0) name(gc)
twoway hist DS, xline(1.0) name(gd)
graph combine gb ga gc gd, xcommon

```



```
clear all
use siml_10_200.dta, clear

su

Variable | Obs      Mean      Std. dev.      Min      Max
-----+-----
OVB |      100      1.101343      .0070703      1.080156      1.119817
OLS |      100      1.001209      .0374376      .8977359      1.055103
LASSO2 |      100      1.03286      .0024848      1.088108      1.1101014
PDS |      100      1.059699      .1101767      .7438613      1.402718
DS |      100      .8432662      .0835669      .6320081      1.156364
-----+-----
P_out |      100      .6107089      .0921731      .4086039      .978497
IV |      100      1.100013      .001948      1.093504      1.103477
IVlasso |      100      0      0      0      0

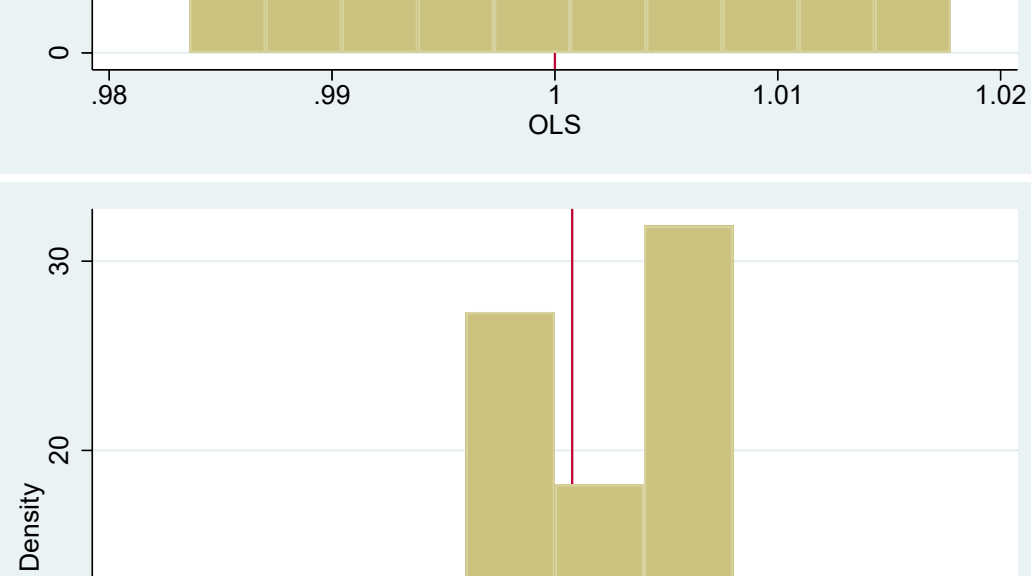
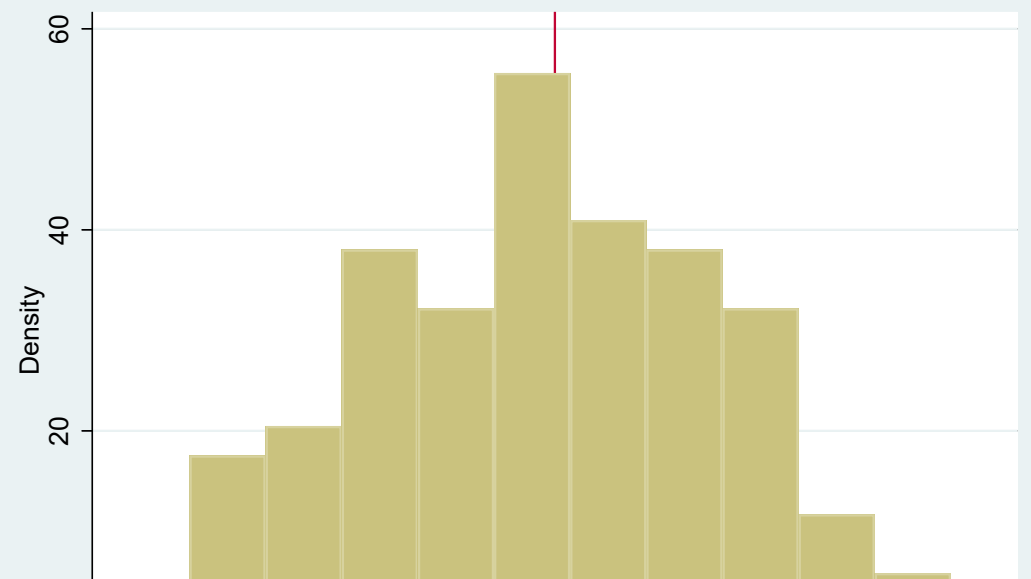
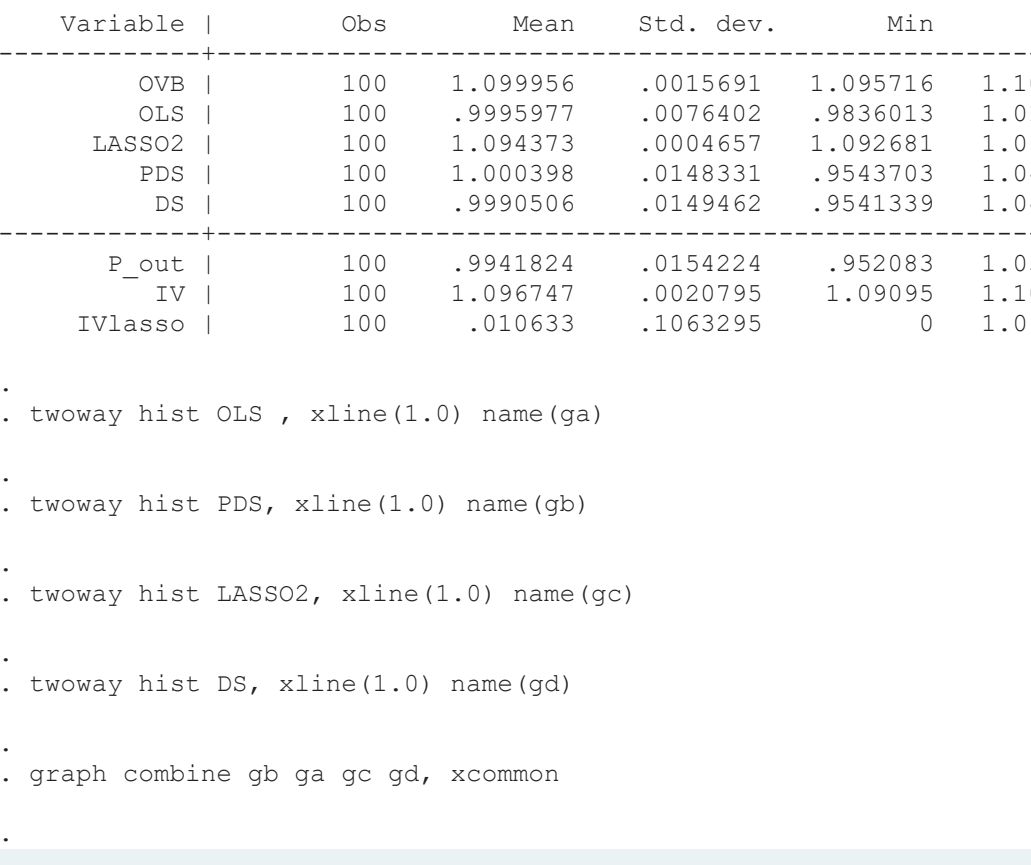
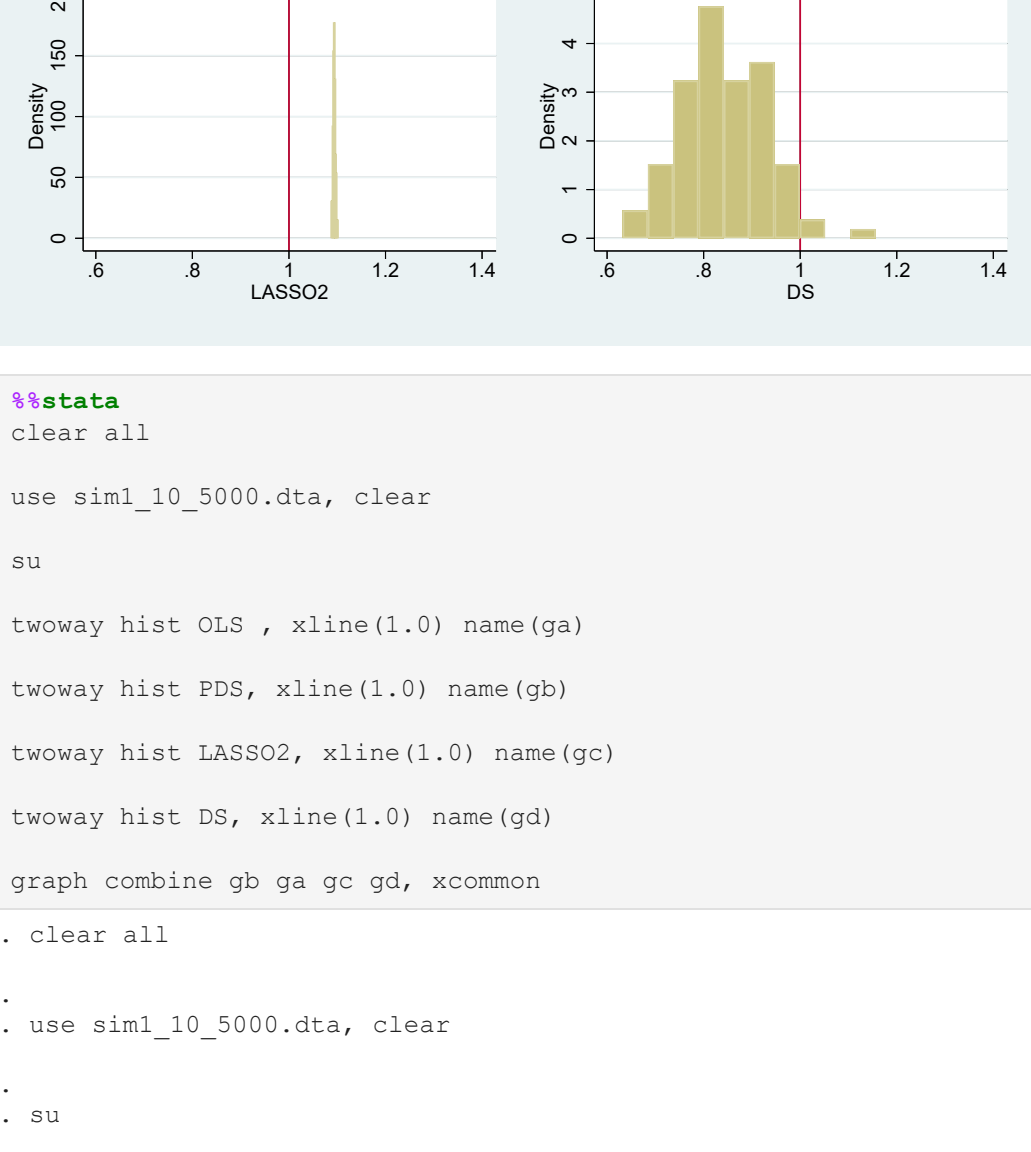
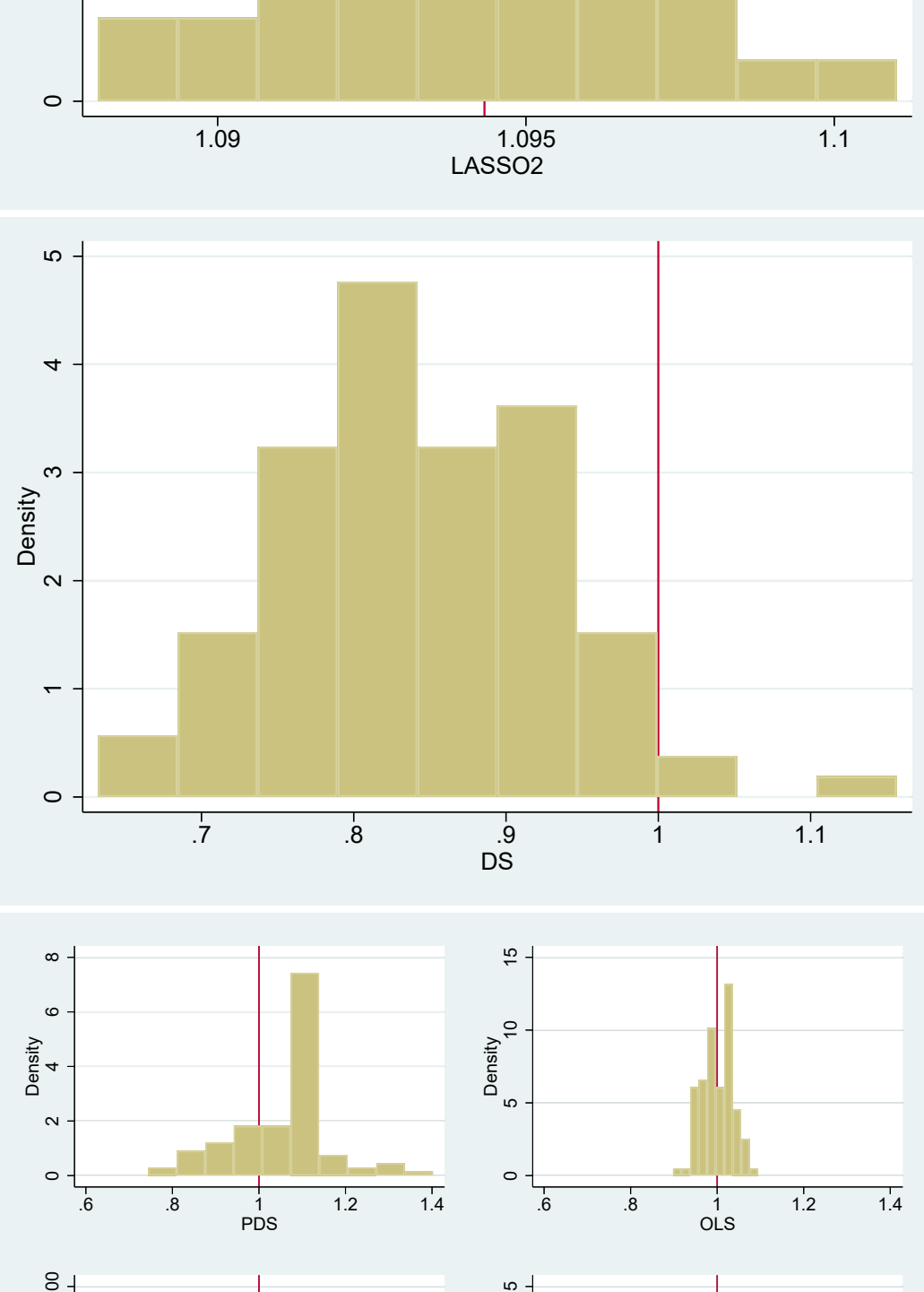
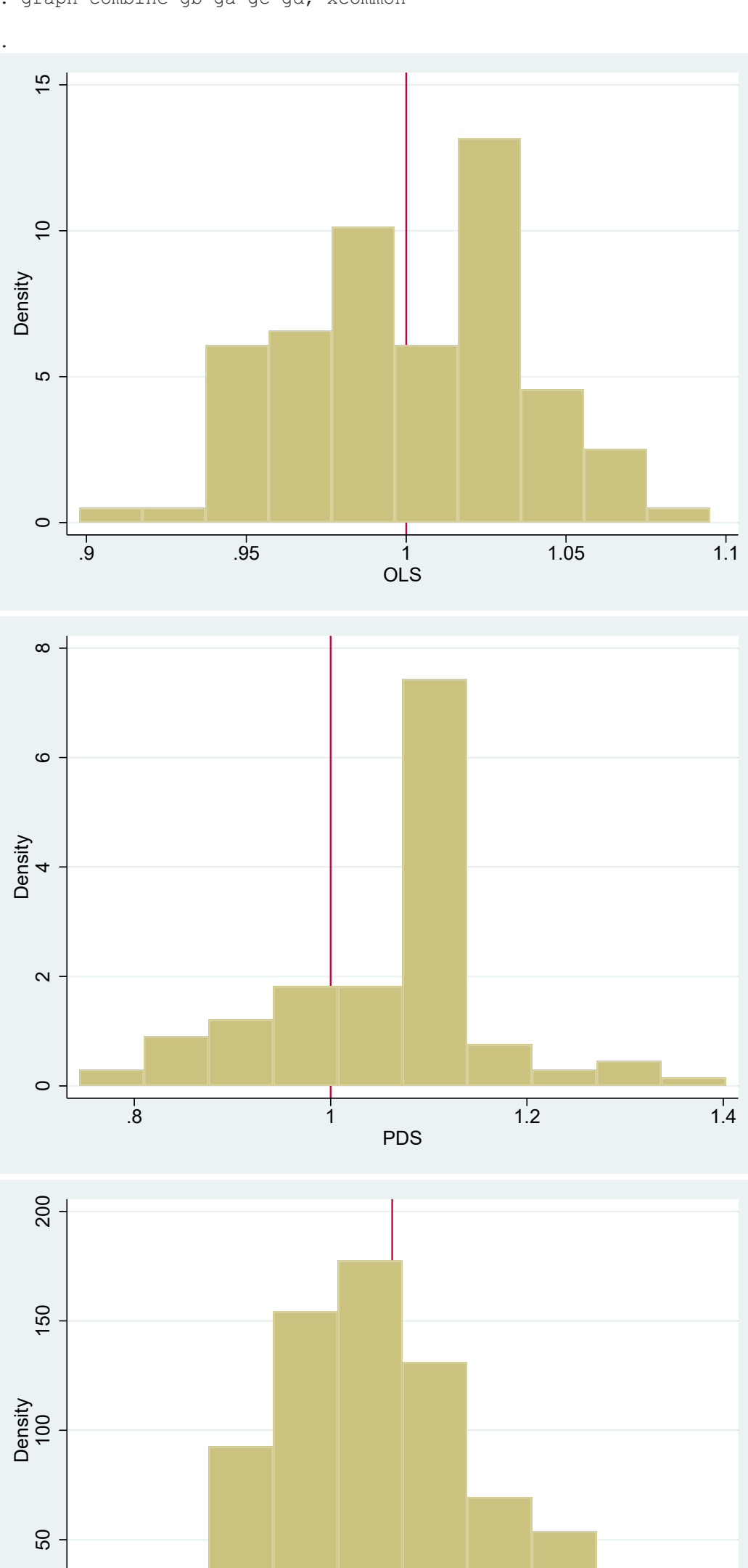
twoway hist OLS , xline(1.0) name(ga)

twoway hist PDS, xline(1.0) name(gb)

twoway hist LASSO2, xline(1.0) name(gc)

twoway hist DS, xline(1.0) name(gd)

graph combine gb ga gc gd, xcommon
```



```
In [16]: %*data
clear all

use siml_10_5000.dta, clear

su

twoway hist OLS , xline(1.0) name(ga)

twoway hist PDS, xline(1.0) name(gb)

twoway hist LASSO2, xline(1.0) name(gc)

twoway hist DS, xline(1.0) name(gd)

graph combine gb ga gc gd, xcommon

clear all
```

```
. use siml_10_5000.dta, clear

. su

Variable | Obs      Mean      Std. dev.      Min      Max
-----+-----
OVB |      100      1.099956      .0015691      1.095716      1.104368
OLS |      100      .9989577      .0076402      .9836013      1.017771
LASSO2 |      100      1.094373      .0004637      1.092681      1.095432
PDS |      100      1.000398      .0148331      .9543703      1.042218
DS |      100      .9990506      .0149462      .9541339      1.041394
-----+-----
P_out |      100      .9841824      .0154224      .952083      1.036947
IV |      100      1.096747      .0020795      1.090995      1.102232
IVlasso |      100      .910633      .1063295      0      1.063295
```

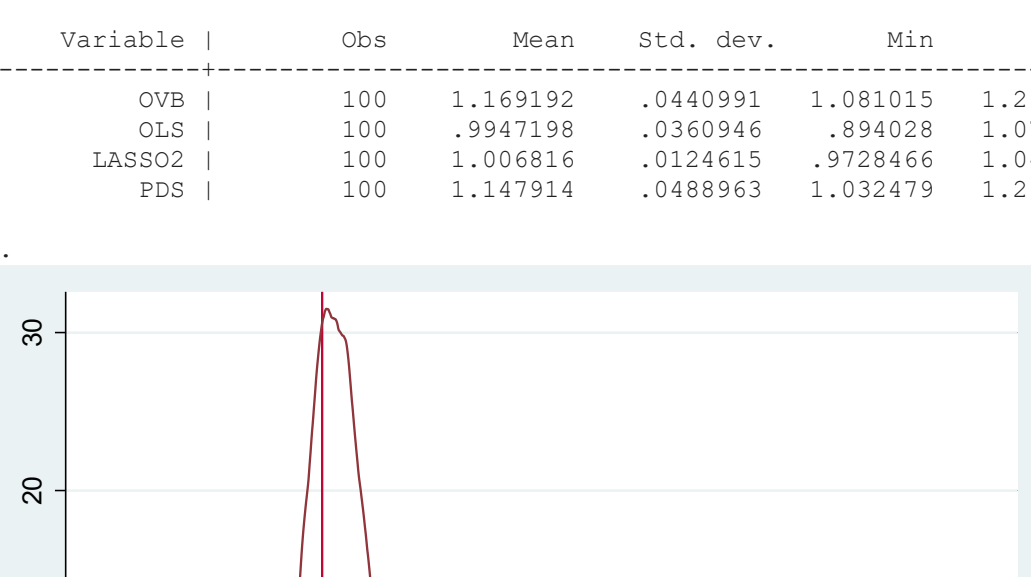
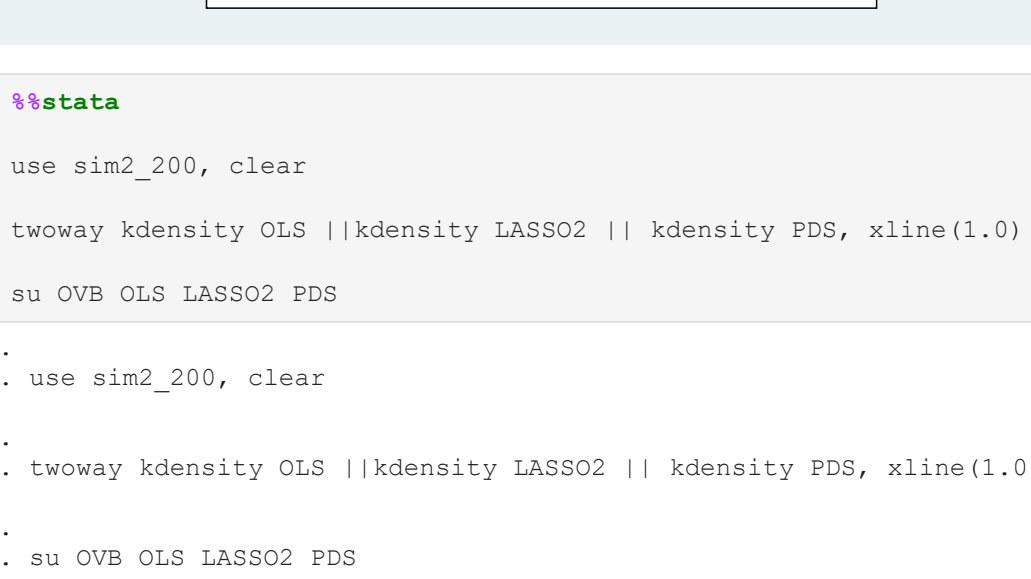
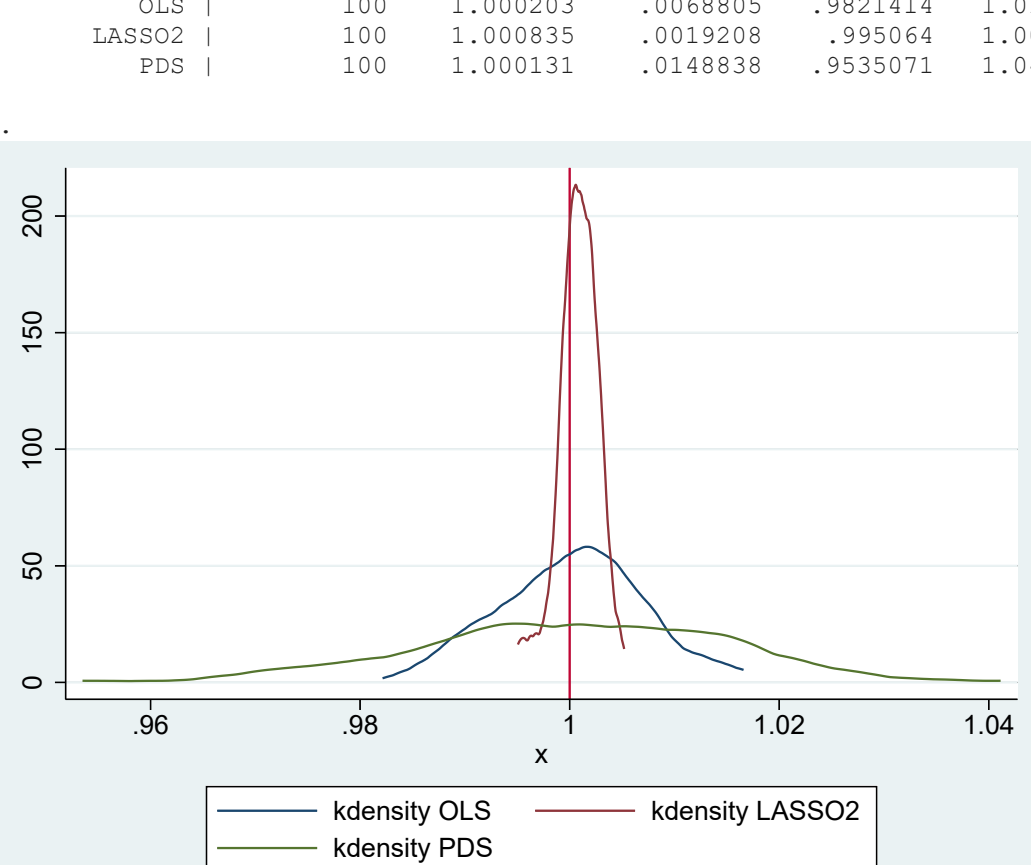
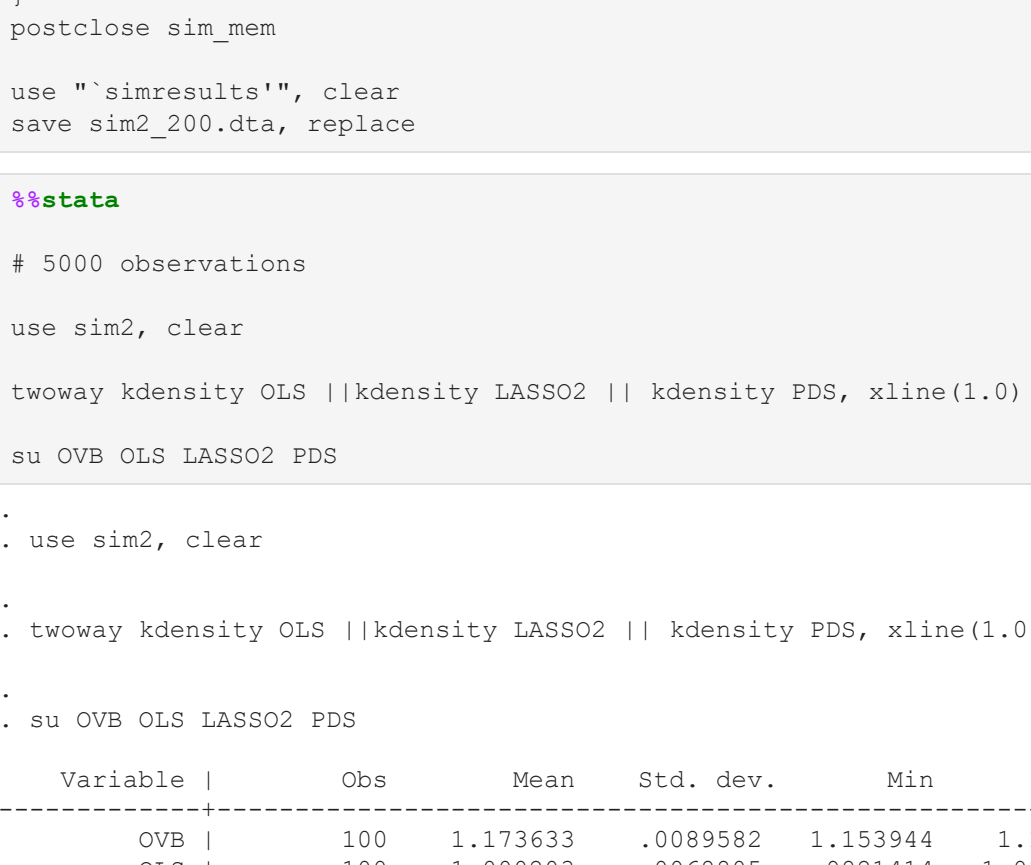
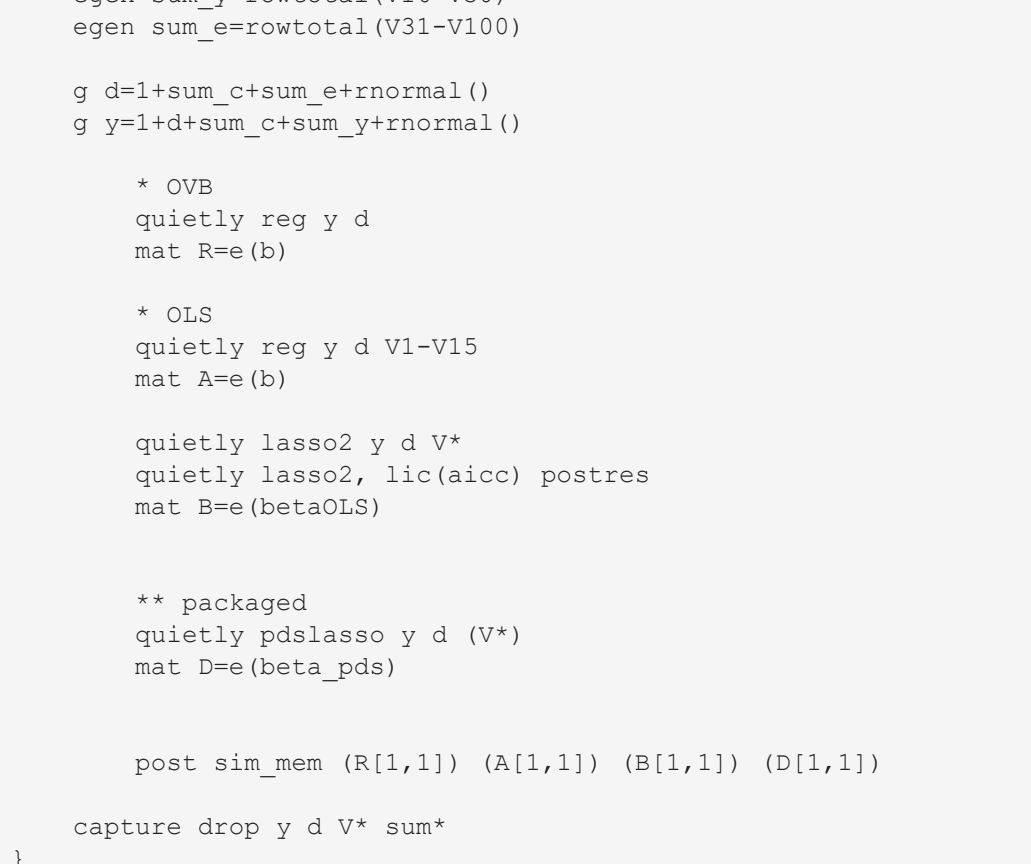
```
. twoway hist OLS , xline(1.0) name(ga)

. twoway hist PDS, xline(1.0) name(gb)

. twoway hist LASSO2, xline(1.0) name(gc)

. twoway hist DS, xline(1.0) name(gd)

graph combine gb ga gc gd, xcommon
```



Sometimes Double Selection might not be better than Naive Lasso

- more likely to mistakenly include valid IV variables into the model when using DS
- less likely to drop confounders in Naive Lasso
- sparsity condition! Now we have 70 valid IV (V31 – V100) for D.

```
In [17]: %*data

clear all

set seed 202111

capture postutil clear
templite simresults
postfile sim_mem OVB OLS LASSO2 PDS using "simresults", replace

forvalues i=1/100 {

    quietly set obs 200

    * di `i'

    su

    forvalues k=1/232{
        g V`k'=rnormal(0,1)
    }

    egen sum_c=rowtotal(V1-V15)
    egen sum_y=rowtotal(V16-V30)
    egen sum_e=rowtotal(V31-V100)

    g d=(sum_c+sum_e*rnormal())
    g y=d+sum_c+sum_y*rnormal()

    * OVB
    quietly reg y d
    mat R= e(b)

    * OLS
    quietly reg y d V1-V15
    mat R= e(b)

    quietly lasso2 y d V*
    quietly lasso2, l1c(aicc) postres
    mat B=(betaOLS)

    ** packaged
    quietly pdlasso y d (V*)
    mat B=(beta_pds)

    post sim_mem (R(1,1)) (A(1,1)) (B(1,1)) (D(1,1))

    capture drop y d V* sum*
}

postclose sim_mem

use "simresults", clear
save siml_200.dta, replace
```

```
In [18]: %*data

# 5000 observations

use siml_200, clear

twoway kdensity OLS || kdensity LASSO2 || kdensity PDS, xline(1.0)

su OVB OLS LASSO2 PDS
```

```
. use siml_200, clear

. twoway kdensity OLS || kdensity LASSO2 || kdensity PDS, xline(1.0)

. su OVB OLS LASSO2 PDS

Variable | Obs      Mean      Std. dev.      Min      Max
-----+-----
OVB |      100      1.173633      .0089392      1.153944      1.19357
OLS |      100      1.002003      .0068905      .9921414      1.016598
LASSO2 |      100      1.000635      .0019208      .995064      1.005203
PDS |      100      1.000131      .0148838      .9535071      1.041129
```



```
In [18]: %*data

use siml_200, clear

twoway kdensity OLS || kdensity LASSO2 || kdensity PDS, xline(1.0)

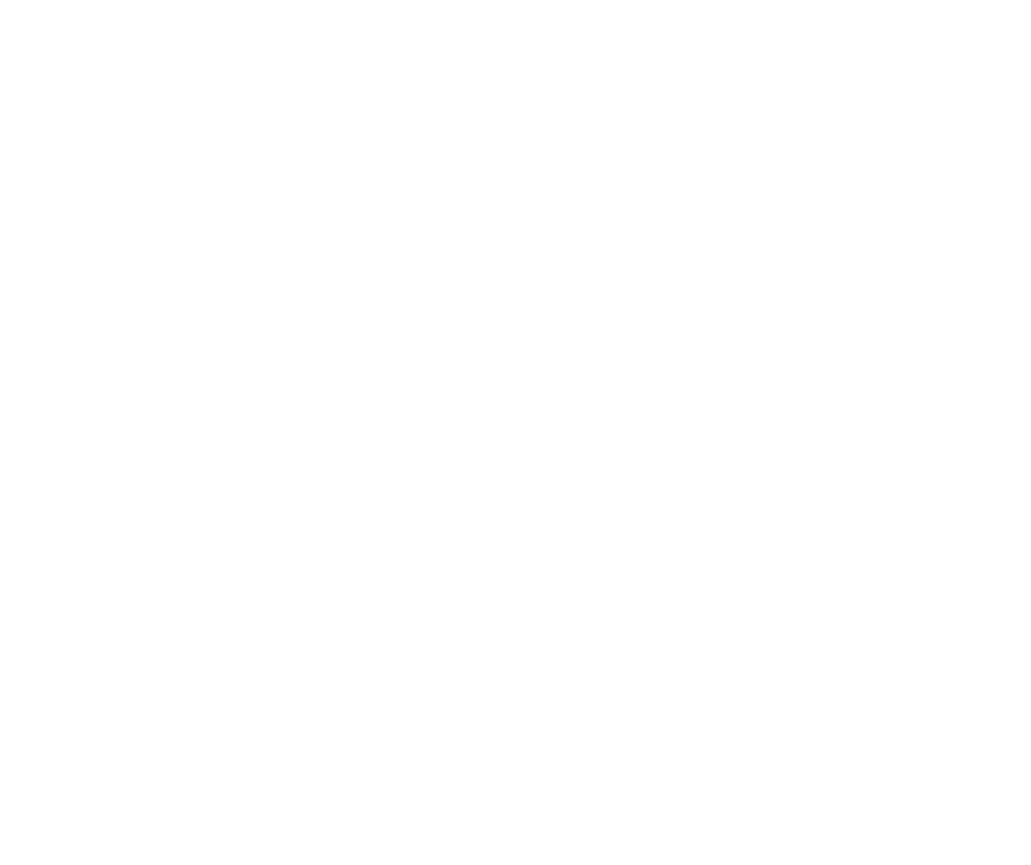
su OVB OLS LASSO2 PDS
```

```
. use siml_200, clear

. twoway kdensity OLS || kdensity LASSO2 || kdensity PDS, xline(1.0)

. su OVB OLS LASSO2 PDS

Variable | Obs      Mean      Std. dev.      Min      Max
-----+-----
OVB |      100      1.169192      .0468991      1.091015      1.260978
OLS |      100      .9947198      .0360946      .894028      1.071052
LASSO2 |      100      1.006816      .0124615      .9728466      1.042932
PDS |      100      1.147914      .0468963      1.032479      1.260978
```



```
In [19]: %*data

clear all

use siml_200, clear

twoway kdensity OLS || kdensity LASSO2 || kdensity PDS, xline(1.0)

su OVB OLS LASSO2 PDS
```