# Leave-Out Estimation of Variance Components in Two-Way Fixed Effects Models Using MATLAB

# January 19, 2021

This notebook describes the MATLAB package that implements the leave-out correction of Kline, Saggio, and Sølvsten (2020) (henceforth KSS) for two-way fixed effects models.

### **Contents**

1	Introduction	2
	1.1 The KSS Correction	
	1.1.2 The Bias in the Plug-in Estimator	
	1.1.3 The Problem with "Standard" Standard Errors in High- Dimensional Models	3
	1.1.4 The Leave-Out Correction	4
2	Computing the KSS Correction	5
	2.1 Setup	
	2.2 Importing the Data	
	2.3 Calling the Main Function	6
3	Interpreting the Output	7
4	What Does the Code Save?	7
5	Scaling to Large Datasets	8
	5.1 Computational Bottleneck	8
	5.2 Approximating the Statistical Leverages	8
6	Adding Controls	14
7	Leaving Out a Person-Year Observation vs. Leaving Out a Match	15
8	Variance of Person Effects When Leaving Out a Match	17
9	Regressing Firm Fixed Effects on Observables	18

#### 1 Introduction

Economists often study settings where units possess two or more group memberships, some of which can change over time. A prominent example comes from Abowd et al. (1999) (henceforth AKM) who proposed a panel model of log wage determination that is additive in worker and firm fixed effects.

This so-called "two-way" fixed effects or "AKM" model takes the form

$$y_{gt} = \alpha_g + \psi_{i(g,t)} + w'_{gt}\delta + \varepsilon_{gt} \qquad (g = 1, \dots, N, \ t = 1, \dots, T_g \ge 2), \tag{1}$$

where the function  $j(\cdot, \cdot): \{1, ..., N\} \times \{1, ..., \max_i T_g\} \to \{1, ..., J\}$  assigns a worker g and year t observation to one of J firms. Here  $\alpha_g$  is a person effect,  $\psi_{j(g,t)}$  is a firm effect,  $w_{gt}$  is a time-varying covariate, and  $\varepsilon_{gt}$  is a mean-independent time-varying error.

We can rewrite the original AKM model as:

$$y_i = x_i'\beta + \varepsilon_i \qquad i = 1, ..., n, \tag{2}$$

where *i* indexes a particular person-year observation (g, t),  $x_i$  is a vector that collects all the worker and firm dummies as well as the time-varying covariates  $w_{gt}$  so that  $\beta = (\alpha, \psi, \delta)'$  is a  $k \times 1$  vector that collects all the worker and firm fixed effects along with  $\delta$ .

Interest in AKM models often centers on understanding how much of the variability in log wages is attributable to firms. It is common to summarize the firm contribution to wage inequality using the following two parameters:

$$\sigma_{\psi}^{2} = \frac{1}{n} \sum_{g=1}^{N} \sum_{t=1}^{T_{g}} \left( \psi_{j(g,t)} - \bar{\psi} \right)^{2} \quad \text{and } \sigma_{\alpha,\psi} = \frac{1}{n} \sum_{g=1}^{N} \sum_{t=1}^{T_{g}} \left( \psi_{j(g,t)} - \bar{\psi} \right) \alpha_{g}$$
 (3)

where  $\bar{\psi} = \frac{1}{n} \sum_{g=1}^{N} \sum_{t=1}^{T_g} \psi_{j(g,t)}$ . The variance component  $\sigma_{\psi}^2$  measures the contribution of firm wage variability to inequality, while the covariance component  $\sigma_{\alpha,\psi}$  measures the additional contribution of systematic sorting of high-wage workers to high-wage firms.

The function leave\_out\_KSS provides unbiased estimates of  $\sigma_{\psi}^2$  and  $\sigma_{\alpha,\psi}$  as well as an estimate of  $\sigma_{\alpha}^2 = \frac{1}{n} \sum_{g=1}^{N} \sum_{t=1}^{T_g} \left(\alpha_g - \bar{\alpha}\right)^2$  using the leave-out bias-correction approach proposed by KSS.

#### 1.1 The KSS Correction

We now provide some general intuition about the KSS leave-out methodology. A more formal discussion can be found in KSS.

#### 1.1.1 The Plug-in Estimator

Suppose that the researcher is interested in the variance of firm effects,  $\sigma_{\psi}^2$ . To simplify the exposition, we normalize the firm effects so that their firm-size-weighted mean is equal to zero, i.e,

 $\bar{\psi}=0$ , and rewrite  $\sigma_{\psi}^2$  as

$$\sigma_{\psi}^2 = \sum_{j=1}^J s_j \psi_j^2 \tag{4}$$

where  $s_j$  gives the employment share of firm j, i.e.,  $s_j = \frac{1}{n} \sum_{g=1}^{N} \sum_{t=1}^{T_g} \mathbf{1} \{ j(g,t) = j \}.$ 

It is customary to report "plug-in" estimates of a given variance component using the corresponding OLS estimate. For instance, the plug-in estimate of the variance of firm effects  $\sigma_{\psi}^2$  is given by

$$\tilde{\sigma}_{\psi}^2 = \sum_{j=1}^{J} s_j \hat{\psi}_j^2 \tag{5}$$

where  $\hat{\psi}_i$  is the OLS estimate obtained after estimating equation (1) via OLS.

#### 1.1.2 The Bias in the Plug-in Estimator

The estimated firm effect,  $\hat{\psi}_j$ , represents a noisy estimate of the true firm effect,  $\psi_j$ . The presence of noise in  $\hat{\psi}_j$  is not an issue when one is interested in  $\psi_j$  as the OLS estimator  $\hat{\psi}_j$  is assumed to be unbiased in this context, i.e.,  $E[\hat{\psi}_j] = \psi_j$ .

However, the estimation error in  $\hat{\psi}_j$  is going to lead to biases if one is interested in estimating  $\psi_j^2$  using its "plug-in" analogue  $\hat{\psi}_j^2$  since

$$E[\hat{\psi}_{j}^{2}] = E[(\hat{\psi}_{j} - \psi_{j} + \psi_{j})^{2}] = \psi^{2} + \underbrace{\mathbb{V}[\hat{\psi}_{j}]}_{\text{Bias}},$$
(6)

where  $\mathbb{V}[\hat{\psi}_j]$  is the squared standard error of  $\hat{\psi}_j$ . Intuitively, when we take the square of  $\hat{\psi}_j$  we are not only squaring its signal,  $\psi_j$ , but also the estimation error in each  $\hat{\psi}_j$ . The latter is going to introduce a bias when estimating  $\psi_i^2$ .

The same logic applies when analyzing the bias of the plug-in estimator of the variance of firm effects since

$$E[\tilde{\sigma}_{\psi}^{2}] = \sigma_{\psi}^{2} + \underbrace{\sum_{j=1}^{J} s_{j} \mathbf{V}[\hat{\psi}_{j}]}_{\text{Bias}}$$

$$(7)$$

#### 1.1.3 The Problem with "Standard" Standard Errors in High-Dimensional Models

The above formula shows that the bias of the plug-in estimator of the variance of firm effects is

$$E[\tilde{\sigma}_{\psi}^2] - \sigma_{\psi}^2 = \sum_{j=1}^J s_j \mathbb{V}[\hat{\psi}_j]. \tag{8}$$

Therefore, all that is required for a bias correction is an estimate of the (squared) standard error of each firm effect,  $V[\hat{\psi}_j]$ . Similarly, if we are interested on the variance of person effects, then we would need the standard error on each of the person effects,  $V[\hat{\alpha}_i]$ . If we are interested in the covariance of worker and firm effects, then we would need the covariances in sampling errors between each  $\hat{\alpha}_i$  and  $\hat{\psi}_{i(i,t)}$ .

The above discussion highlights that an estimate of the sampling variability of the OLS coefficient vector  $\hat{\beta} = (\hat{\alpha}, \hat{\psi}, \hat{\delta})$  is required in order to derive an unbiased estimate of the variance components of the AKM model displayed in equation (2).

Recall that the sampling variability in  $\hat{\beta}$ , assuming independence across observations, is given by

$$\mathbb{V}[\hat{\beta}] = \left(\sum_{i=1}^{n} x_i x_i'\right)^{-1} \sum_{i=1}^{n} \sigma_i^2 x_i x_i' \left(\sum_{i=1}^{n} x_i x_i'\right)^{-1}, \tag{9}$$

where  $\sigma_i^2 = Var(\varepsilon_i)$ .

One might be tempted to provide an estimate of  $\mathbb{V}[\hat{\beta}]$  using heteroskedasticity *consistent* ("HC") or robust standard errors. Standard White (1980) HC standard-errors are calculated using a plug-in estimate of  $\sigma_i^2$  based on

$$\tilde{\sigma}_i^2 = (y_i - x_i'\hat{\beta})^2,\tag{10}$$

where the HC-based estimate of  $\mathbb{V}[\hat{\beta}]$  is given by

$$\tilde{\mathbb{V}}[\hat{\beta}] = \left(\sum_{i=1}^{n} x_i x_i'\right)^{-1} \sum_{i=1}^{n} \tilde{\sigma}_i^2 x_i x_i' \left(\sum_{i=1}^{n} x_i x_i'\right)^{-1}.$$
 (11)

However, HC standard errors based on  $\tilde{\sigma}_i^2$  are downward biased (MacKinnon and White, 1985). From an asymptotic perspective, HC standard errors are inconsistent in any high-dimensional model where the number of parameters grows in proportion to the sample size (Cattaneo et al., 2018). Such "many regressor" asymptotics are natural in the worker-firm fixed effects model as we often have fewer than 5 worker moves on average per firm.

#### 1.1.4 The Leave-Out Correction

KSS provides a heteroskedasticity-unbiased (HU) estimate of the standard error of any coefficient obtained from a linear regression model.

The KSS HU standard error estimate is based on a leave-out estimate of  $\sigma_i^2$ :

$$\hat{\sigma}_i^2 = y_i (y_i - x_i' \hat{\beta}_{-i}), \tag{12}$$

where  $\hat{\beta}_{-i}$  is the OLS estimate of  $\beta$  from equation (2) when observation i is left out.

KSS then replaces  $\sigma_i^2$  in  $\mathbb{V}[\hat{\beta}]$  with its unbiased estimate  $\hat{\sigma}_i^2$  to derive an HU estimate of  $\mathbb{V}[\hat{\beta}]$ :

$$\hat{\mathbb{V}}[\hat{\beta}] = \left(\sum_{i=1}^{n} x_i x_i'\right)^{-1} \sum_{i=1}^{n} \hat{\sigma}_i^2 x_i x_i' \left(\sum_{i=1}^{n} x_i x_i'\right)^{-1}.$$
 (13)

Going back to the example of the variance of the firm effects, we can extract from  $\hat{\mathbb{V}}[\hat{\beta}]$  the corresponding squared standard error of each firm effect,  $\hat{\mathbb{V}}[\hat{\psi}_j]$ . We can then use it to bias-correct the corresponding estimate of the variance of the firm effects as follows:

$$\hat{\sigma}_{\psi}^2 = \tilde{\sigma}_{\psi}^2 - \sum_{j=1}^{J} s_j \mathbb{V}[\hat{\psi}_j]. \tag{14}$$

The MATLAB function leave\_out\_KSS is going to print the bias-corrected variance of firm effects,  $\hat{\sigma}_{\psi}^2$ , the bias-corrected covariance of worker and firm effects and variance of person effects. leave\_out\_KSS also provides the correct standard errors — based on  $\hat{\mathbb{V}}[\hat{\beta}]$  as opposed to  $\tilde{\mathbb{V}}[\hat{\beta}]$  — when one regresses the firm effects on some observable characteristics; see Section 9.

# 2 Computing the KSS Correction

We now demonstrate how one can implement the KSS correction in two-way models using MAT-LAB and the function <code>leave\_out\_KSS</code>. We continue to work with a simple example based on an AKM model. In what follows, we use the words "workers" and "firms" when describing the procedure but the function <code>leave\_out\_KSS</code> can in fact be applied to any two-way fixed effects model (e.g., patients and doctors, students and teachers, strata and treatment arms).

#### 2.1 Setup

We begin with some auxiliary lines of code that define the relevant paths, call the CMG package package developed by Yiannis Koutis and set-up the parallel environment within MATLAB.

```
[1]: %Setup Paths and Install CMG
clc
clear
cd '/Users/raffaelesaggio/Dropbox/LeaveOutTwoWay'
path(path,'codes'); %this contains the main LeaveOut Routines.
path(path,'CMG'); % CMG package http://www.cs.cmu.edu/~jkoutis/cmg.html
[result,output] = evalc('installCMG(1)'); %installs CMG routine (silently)
delete(gcp("nocreate")) %clear parallel envir.
c = parcluster('local'); %tell me # of available cores
nw = c.NumWorkers; %tell me # of available cores
pool=parpool(nw,'IdleTimeout', Inf); %all cores will be assigned to Matlab
```

Starting parallel pool (parpool) using the 'local' profile ... Connected to the parallel pool (number of workers: 6).

#### 2.2 Importing the Data

The GitHub Repo contains a matched employer-employee testing data where we observe the identity of the worker, the identity of the firm employing a given worker, the year in which the match is observed (either 1999 or 2001), and the associated log wage.

*Important!* The original data must be sorted by individual identifiers (id). For instance, one can see that the testing data is sorted by individual identifiers (and by year, using xtset id year in Stata)

```
[2]: %% Import Data
namesrc='data/test.csv'; %path to original testing data
data=importdata(namesrc); %import data
id=data(:,1); %worker identifiers
firmid=data(:,2); %firm identifiers
y=data(:,4); % outcome variable
clear data
```

#### 2.3 Calling the Main Function

The function leave\_out\_KSS relies on three mandatory inputs: (y,id,firmid). We can obtain an unbiased variance decomposition of the associated AKM model by simply calling

```
[3]: \%\% Run \KSS!
    [sigma2_psi,sigma_psi_alpha,sigma2_alpha] = leave_out_KSS(y,id,firmid);
   _*_*_*_*_*
   Running KSS Correction with the following options
   Leave Out Strategy: Leave match out
   Algorithm for Computation of Statistical Leverages: JLA with 200 simulations.
   _*-*-*-*-*-*-*-*-*-*-*-*
   _*-*-*-*-*-*-*-*-*-*-*-*-*-*
   SECTION 1
   _*-*-*-*-*-*-*-*-*-*-*-*
   _*-*-*-*-*-*-*-*-*-*-*-*
   Info on the leave one out connected set:
   _*_*_*_*_*
   mean wage: 4.7636
   variance of wage: 0.1245
   # of Movers: 6414
   # of Firms: 1684
   # of Person Year Observations: 56044
   _*_*_*_*_*_*
   _*_*_*_*_*_*
   SECTION 2
   _*-*-*-*-*-*-*-*-*-*-*-*-*
   _*-*-*-*-*-*-*-*-*-*-*
   Calculating the statistical leverages of the AKM model...
   Running JLA Algorithm...
```

Done!

Elapsed time is 5.602229 seconds.

SECTION 3

PLUG-IN ESTIMATES (BIASED)

Variance of Firm Effects: 0.019821

Covariance of Firm, Person Effects: -0.0039091

Variance of Person Effects: 0.10354

Correlation of Firm, Person Effects: -0.08629

BIAS CORRECTED ESTIMATES

Variance of Firm Effects: 0.010218

Covariance of Firm and Person Effects: 0.0047795

Variance of Person Effects: 0.085005

Correlation of Firm, Person Effects: 0.16217

### 3 Interpreting the Output

The code starts by printing its two key inputs: the algorithm used to compute the statistical leverages (exact vs. JLA) — we explain this distinction in Section 5 — and the level at which the leave-out correction is carried (observation vs. match) — we explain this in more detail in Section 7.

The output printed by leave\_out\_KSS is composed of three sections.

**Section 1**: Here we provide info on the leave-out connected set. This is the largest connected set of firms that remains connected after any worker from the associated graph is removed, see Lemma 1 and the Computational Appendix of KSS for details. The code provides some summary statistics (e.g. number of movers, number of firms, mean and variance of the outcome, etc.) of the leave-out connected set.

**Section 2**: After printing the summary statistics, the code computes the statistical leverages of the design, denoted by  $P_{ii}$ . Computation of  $\{P_{ii}\}_{i=1}^n$  represents the main computational bottleneck of the routine.

**Section 3**: The code then enters its third, and final stage, where the main results are printed. The code starts by reporting the — biased — estimates of the variance components that result from the "plug-in" approach of treating OLS estimates as measured without error. Finally, the code prints the bias-corrected variance of firm effects and the covariance of worker and firm effects.

#### 4 What Does the Code Save?

leave\_out\_KSS saves three scalars: the variance of firm effects (sigma2\_psi in [4]), the covariance of worker and firm effects (sigma\_psi\_alpha), and the variance of person effects (sigma2\_alpha).

leave\_out\_KSS also saves on disk one .csv file. This .csv contains information on the leave-out connected set. This file has 4 columns. First column reports the outcome variable, second and third columns the worker and the firm identifiers (as originally inputted by the user) and the fourth column reports the statistical leverages of the regression design. If the code is reporting a leave-out correction at the match-level, the .csv will be collapsed at the match level. By default, the .csv file is going to be saved in the main directory under the name leave\_out\_estimates. The user can specify an alternative path using the option filename when calling leave\_out\_KSS.

# 5 Scaling to Large Datasets

leave\_out\_KSS can be used on extremely large datasets. The code uses a variant of the random projection method, known as the Johnson–Lindenstrauss approximation (JLA) algorithm in KSS for its connection to the work of Johnson and Lindenstrauss (1984); see also Achlioptas (2003). We now discuss briefly the main computational bottleneck of the procedure and the JLA algorithm.

#### 5.1 Computational Bottleneck

Recall from the discussion in Section 1 that the KSS leave-out bias correction procedure relies on leave-out estimates of  $\sigma_i^2$ ,

$$\hat{\sigma}_i^2 = y_i (y_i - x_i' \hat{\beta}_{-i}), \tag{15}$$

where  $\hat{\beta}_{-i}$  is the OLS estimate of  $\beta$  from the AKM model in equation (2) after leaving observation i out.

Clearly, reestimating  $\hat{\beta}_{-i}$  by leaving a particular observation i for n times, is infeasible computationally. Fortunately, one can rewrite  $\hat{\sigma}_i^2$  as

$$\hat{\sigma}_i^2 = y_i \frac{(y_i - x_i' \hat{\beta})}{1 - P_{ii}},\tag{16}$$

where  $P_{ii}$  measures the influence or leverage of observation i, i.e.,  $P_{ii} = x_i' S_{xx}^{-1} x_i$ . The above expression highlights that all that is needed for computation of  $\hat{\sigma}_i^2$  are the n statistical leverages  $\{P_{ii}\}_{i=1}^n$ . However, exact computation of  $P_{ii}$  may remain prohibitive when n is in the order of tens of millions or higher.

### 5.2 Approximating the Statistical Leverages

The JLA algorithm introduced by KSS provides a stochastic approximation to  $\{P_{ii}\}_{i=1}^n$  using the random projection ideas developed by Johnson and Lindenstrauss (1984). We refer the reader to the Computational Appendix of KSS for further details.

The number of simulations underlying the JLA algorithm is governed by the input simulations\_JLA (which is denoted by p in the computational appendix). Intuitively, more simulations imply a higher accuracy – but also higher computation time — when estimating  $\{P_{ii}, B_{ii}\}_{i=1}^n$ .

**Note:** The user might want to prespecify a random-number generator seed to ensure replicability when calling the function leave\_out\_KSS.

We now demonstrate the performance of the code on a large dataset.

```
[4]: \"\" Running KSS on a large dataset
    websave('large_fake.csv', 'https://www.dropbox.com/s/ny5tef29ij7ran2/
     →large_fake_data.csv?dl=1'); %downloads and saves to disk a fake, large matched_
     →employer employee data
    namesrc='large_fake.csv'; %path to the large data
    data=importdata(namesrc); %import data
    id=data(:,1); %worker identifiers
    firmid=data(:,2); %firm identifiers
    y=data(:,4); % outcome variable
    clear data
    delete('large_fake.csv'); %delete original .csv data from disk
    %Run Leave Out Correction (50 simulations)
    type_of_algorithm='JLA'; %run random projection algorithm
    simulations_JLA=50;
    [sigma2_psi,sigma_psi_alpha,sigma2_alpha] =_
     →leave_out_KSS(y,id,firmid,[],[],type_of_algorithm,simulations_JLA);
   _*-*-*-*-*-*-*-*-*-*-*-*
   Running KSS Correction with the following options
   Leave Out Strategy: Leave match out
   Algorithm for Computation of Statistical Leverages: JLA with 50 simulations.
   _*_*_*_*_*
   _*-*-*-*-*-*-*-*-*-*-*-*-*-*
   SECTION 1
   _*-*-*-*-*-*-*
   _*_*_*_*_*
   Info on the leave one out connected set:
   _*_*_*_*_*
   mean wage: 4.7304
   variance of wage: 0.16248
   # of Movers: 916632
   # of Firms: 165360
   # of Person Year Observations: 13860616
   _*_*_*_*_*
   _*_*_*_*_*
   SECTION 2
   _*-*-*-*-*-*-*-*-*-*-*-*
   _*_*_*_*_*
   Calculating the statistical leverages of the AKM model...
   Running JLA Algorithm...
   Done!
   Elapsed time is 251.168051 seconds.
```

```
_*-*-*-*-*-*-*-*-*-*-*-*-*-*
_*-*-*-*-*-*-*-*-*-*-*-*-*
SECTION 3
_*_*_*_*_*
_*_*_*_*_*_*
PLUG-IN ESTIMATES (BIASED)
Variance of Firm Effects: 0.039448
Covariance of Firm, Person Effects: 0.0084313
Variance of Person Effects: 0.080329
Correlation of Firm, Person Effects: 0.14978
_*-*-*-*-*-*-*-*-*-*-*-*
_*-*-*-*-*-*-*-*-*-*-*-*
BIAS CORRECTED ESTIMATES
Variance of Firm Effects: 0.030371
Covariance of Firm and Person Effects: 0.014639
Variance of Person Effects: 0.048803
Correlation of Firm, Person Effects: 0.38024
```

We can see from the output that the leave-out connected set has almost 14 million person-year observations. The code is able to complete in around 4 minutes (on a 2020 Macbook Pro with 6 cores and 16GB of RAM).

The computational appendix in KSS shows that the JLA algorithm can cut computation time by a factor of 100 while introducing an approximation error of roughly  $10^{-4}$ .

The current code uses an improved estimator of both  $P_{ii}$  and  $M_{ii} = 1 - P_{ii}$ , which are both guaranteed to lie in [0,1]. These improved estimators are then combined to derive an asymptotically unbiased JLA estimator of a given variance component provided that  $\frac{n}{p^4} = o(1)$ ; see this document for further details..

We can check the stability of the estimates for different values of simulations\_JLA. For instance, if we double simulations\_JLA from 50 to 100 and run the code again on the same data:

[5]: | %% Compute estimates while doubling number of simulations

```
mean wage: 4.7304
variance of wage: 0.16248
# of Movers: 916632
# of Firms: 165360
# of Person Year Observations: 13860616
_*_*_*_*_*
_*_*_*_*_*_*
SECTION 2
_*_*_*_*_*_*
_*_*_*_*_*_*
Calculating the statistical leverages of the AKM model...
Running JLA Algorithm...
Done!
Elapsed time is 458.093692 seconds.
_*_*_*_*_*
_*-*-*-*-*-*-*-*-*-*-*-*-*
SECTION 3
_*-*-*-*-*-*-*-*-*-*-*-*-*
_*_*_*_*_*_*
PLUG-IN ESTIMATES (BIASED)
Variance of Firm Effects: 0.039448
Covariance of Firm, Person Effects: 0.0084313
Variance of Person Effects: 0.080329
Correlation of Firm, Person Effects: 0.14978
_*_*_*_*_*_*
_*_*_*_*_*_*
BIAS CORRECTED ESTIMATES
Variance of Firm Effects: 0.030382
Covariance of Firm and Person Effects: 0.014598
Variance of Person Effects: 0.048738
Correlation of Firm, Person Effects: 0.37937
```

We obtain virtually the same variance components as when simulations\_JLA=50 while significantly increasing the computational time! If the user does not specify a value for simulations\_JLA, the code defaults to simulations\_JLA=200.

We conclude this section by noting that the user can also calculate an exact version of  $\{P_{ii}\}_{i=1}^n$ . This can be done by setting the option type\_of\_algorithm to exact.

**Warning!** Calling the option exact in large datasets can be very time-consuming! We now load again the original, smaller, testing data and then compare the exact and JLA-based estimates of the variance components,

```
[6]: %% Compare Exact vs. JLA Estimates
namesrc='data/test.csv'; %path to original testing data
data=importdata(namesrc); %import data
id=data(:,1); %worker identifiers
firmid=data(:,2); %firm identifiers
y=data(:,4); % outcome variable
```

```
clear data
%Run Leave Out Correction with exact
type_of_algorithm='exact'; %run random projection algorithm;
[sigma2_psi,sigma_psi_alpha,sigma2_alpha] = __
 →leave_out_KSS(y,id,firmid,[],[],type_of_algorithm);
%Run Leave Out Correction with JLA
simulations_JLA=100;
type_of_algorithm='JLA'; %run random projection algorithm;
[sigma2_psi,sigma_psi_alpha,sigma2_alpha] =__
 →leave_out_KSS(y,id,firmid,[],[],type_of_algorithm,simulations_JLA);
_*_*_*_*_*
Running KSS Correction with the following options
Leave Out Strategy: Leave match out
Algorithm for Computation of Statistical Leverages: Exact
_*_*_*_*_*_*
_*_*_*_*_*_*
SECTION 1
_*-*-*-*-*-*-*-*-*-*-*-*-*
_*-*-*-*-*-*-*-*-*-*-*-*
Info on the leave one out connected set:
_*-*-*-*-*-*-*-*-*-*-*-*
mean wage: 4.7636
variance of wage: 0.1245
# of Movers: 6414
# of Firms: 1684
# of Person Year Observations: 56044
_*_*_*_*_*
_*_*_*_*_*
SECTION 2
_*_*_*_*_*_*
_*_*_*_*_*
Calculating the statistical leverages of the AKM model...
Running Exact Algorithm...
Done!
Elapsed time is 162.242224 seconds.
_*-*-*-*-*-*-*-*-*-*-*
_*_*_*_*_*
SECTION 3
_*-*-*-*-*-*-*-*-*-*-*-*
_*-*-*-*-*-*-*-*-*-*-*
PLUG-IN ESTIMATES (BIASED)
```

Variance of Firm Effects: 0.019821

Covariance of Firm, Person Effects: -0.0039091

```
Variance of Person Effects: 0.10354
Correlation of Firm, Person Effects: -0.08629
_*-*-*-*-*-*-*-*-*-*-*-*
_*_*_*_*_*_*
BIAS CORRECTED ESTIMATES
Variance of Firm Effects: 0.010289
Covariance of Firm and Person Effects: 0.0046293
Variance of Person Effects: 0.085204
Correlation of Firm, Person Effects: 0.15635
_*_*_*_*_*_*
Running KSS Correction with the following options
Leave Out Strategy: Leave match out
Algorithm for Computation of Statistical Leverages: JLA with 100 simulations.
_*-*-*-*-*-*-*-*-*-*-*-*
_*_*_*_*_*
SECTION 1
_*-*-*-*-*-*-*-*-*-*-*-*
_*-*-*-*-*-*-*-*-*-*-*-*
Info on the leave one out connected set:
_*-*-*-*-*-*-*-*-*-*-*-*
mean wage: 4.7636
variance of wage: 0.1245
# of Movers: 6414
# of Firms: 1684
# of Person Year Observations: 56044
_*-*-*-*-*-*-*-*-*-*-*
_*_*_*_*_*_*
SECTION 2
_*-*-*-*-*-*-*-*-*-*-*-*
_*-*-*-*-*-*-*-*-*-*-*-*
Calculating the statistical leverages of the AKM model...
Running JLA Algorithm...
Done!
Elapsed time is 2.434287 seconds.
_*_*_*_*_*_*
_*_*_*_*_*
SECTION 3
_*-*-*-*-*-*-*-*-*-*-*-*
_*_*_*_*_*
PLUG-IN ESTIMATES (BIASED)
Variance of Firm Effects: 0.019821
Covariance of Firm, Person Effects: -0.0039091
Variance of Person Effects: 0.10354
Correlation of Firm, Person Effects: -0.08629
_*_*_*_*_*
_*-*-*-*-*-*-*-*-*-*-*-*
BIAS CORRECTED ESTIMATES
Variance of Firm Effects: 0.01044
```

```
Covariance of Firm and Person Effects: 0.0044957
Variance of Person Effects: 0.085326
Correlation of Firm, Person Effects: 0.15063
```

The variance components estimated using JLA are extremely close to the exact estimates but only take a fraction of the time to compute. If the input data has more than 10,000 observations, the code defaults to using the JLA algorithm unless the user specifies type\_of\_algorithm as "exact".

# 6 Adding Controls

We have demonstrated the functioning of leave\_out\_KSS using a simple AKM model with no controls ( $w_{gt} = 0$ ). It is easy to add a matrix of controls to the routine. Suppose for instance that we want to add year fixed effects to the original AKM model. This can be done as follows.

```
[7]: \%% How to add controls
    namesrc='data/test.csv'; %path to original testing data
    data=importdata(namesrc); %import data
    id=data(:,1); %worker identifiers
    firmid=data(:,2); %firm identifiers
    year=data(:,3); %year identifier
    y=data(:,4); % outcome variable
    clear data
    %Specify year fixed effects as controls
    [~,~,controls] = unique(year);
    controls
                      = sparse((1:size(y,1))',controls',1,size(y,1),max(controls));
    controls
                  = controls(:,1:end-1); %to avoid collinearity issues, omit last_
     \rightarrow year fixed effects.
    %Call KSS with matrix of controls
    [sigma2_psi,sigma_psi_alpha,sigma2_alpha] = leave_out_KSS(y,id,firmid,controls);
    _*_*_*_*_*_*
   Running KSS Correction with the following options
   Leave Out Strategy: Leave match out
   Algorithm for Computation of Statistical Leverages: JLA with 200 simulations.
    _*_*_*_*_*
    _*_*_*_*_*_*
   SECTION 1
    _*_*_*_*_*_*
    _*-*-*-*-*-*-*-*-*-*-*-*
   Info on the leave one out connected set:
   _*-*-*-*-*-*-*-*-*-*-*
   mean wage: 4.7636
   variance of wage: 0.1245
   # of Movers: 6414
   # of Firms: 1684
   # of Person Year Observations: 56044
```

```
_*-*-*-*-*-*-*-*-*-*-*-*-*-*
_*-*-*-*-*-*-*-*-*-*-*-*-*
SECTION 2
_*_*_*_*_*_*
_*_*_*_*_*_*
pcg converged at iteration 58 to a solution with relative residual 8.7e-11.
Calculating the statistical leverages of the AKM model...
Running JLA Algorithm...
Done!
Elapsed time is 4.271175 seconds.
_*-*-*-*-*-*-*-*-*-*-*-*-*
_*-*-*-*-*-*-*-*-*-*-*-*-*
SECTION 3
_*-*-*-*-*-*-*-*-*-*-*-*-*-*
_*-*-*-*-*-*-*-*-*-*-*-*-*
PLUG-IN ESTIMATES (BIASED)
Variance of Firm Effects: 0.019479
Covariance of Firm, Person Effects: -0.004008
Variance of Person Effects: 0.10404
Correlation of Firm, Person Effects: -0.089031
_*_*_*_*_*_*
_*_*_*_*_*_*
BIAS CORRECTED ESTIMATES
Variance of Firm Effects: 0.0097812
Covariance of Firm and Person Effects: 0.0046239
Variance of Person Effects: 0.08578
Correlation of Firm, Person Effects: 0.15963
```

When controls are specified, the code proceeds by partialling them out. That is, it first estimates by OLS the AKM model in the leave-out connected set

$$y_{gt} = \alpha_g + \psi_{j(g,t)} + w'_{gt}\delta + \varepsilon_{gt}$$
(17)

from which we obtain  $\hat{\delta}$ . We then work with a residualized model where the outcome variable is now defined as  $y_{gt}^{new} = y_{gt} - w_{gt}' \hat{\delta}$  and project this residualized outcome on worker and firm indicators and report the associated (bias-corrected) variance components.

# 7 Leaving Out a Person-Year Observation vs. Leaving Out a Match

By default, the code reports leave-out corrections for the variance of firm effects and the covariance of firm and worker effects that are robust to unrestricted heteroskedasticity and serial correlation of the error term within a given match (defined as the unique combination of the worker and firm identifier); see Remark 3 of KSS. Intuitively, leaving out matches is analogous to "clustering" the standard error estimates at the match level. Section 8 discusses the interpretation of the leave-out bias-corrected variance of person effects when leaving a match out

The user can specify the function to run the KSS correction when leaving only an observation out using the option leave\_out\_level. When the user leaves a person-year observation out, the

resulting KSS variance components are robust to unrestricted heteroskedasticity but not to serial correlation within a match. Below we demonstrate how to compute KSS- adjusted variance components when leaving a single (person-year) observation out.

```
[8]: | %% Leaving out a Person-Year Observation vs. Leaving Out a Match
    leave_out_level='obs'; %leave a single person-year observation out
    [sigma2_psi,sigma_psi_alpha,sigma2_alpha] =_
     →leave_out_KSS(y,id,firmid,[],leave_out_level);
   _*-*-*-*-*-*-*-*-*-*-*-*
   Running KSS Correction with the following options
   Leave Out Strategy: Leave person-year observation out
   Algorithm for Computation of Statistical Leverages: JLA with 200 simulations.
   _*-*-*-*-*-*-*-*-*-*-*-*
   _*_*_*_*_*_*
   SECTION 1
   _*_*_*_*_*
   _*-*-*-*-*-*-*-*-*-*-*-*-*-*
   Info on the leave one out connected set:
   _*_*_*_*_*
   mean wage: 4.7636
   variance of wage: 0.1245
   # of Movers: 6414
   # of Firms: 1684
   # of Person Year Observations: 56044
   _*_*_*_*_*_*
   _*_*_*_*_*_*
   SECTION 2
   _*-*-*-*-*-*-*-*-*-*-*-*
   _*-*-*-*-*-*-*-*-*-*-*-*
   Calculating the statistical leverages of the AKM model...
   Running JLA Algorithm...
   Done!
   Elapsed time is 4.700531 seconds.
   _*_*_*_*_*
   _*-*-*-*-*-*-*-*-*-*-*-*
   SECTION 3
   _*-*-*-*-*-*-*-*-*-*-*-*
   _*_*_*_*_*
   PLUG-IN ESTIMATES (BIASED)
   Variance of Firm Effects: 0.019821
   Covariance of Firm, Person Effects: -0.0039091
   Variance of Person Effects: 0.10354
   Correlation of Firm, Person Effects: -0.08629
   _*-*-*-*-*-*-*-*-*-*-*-*
   _*_*_*_*_*_*
   BIAS CORRECTED ESTIMATES
```

Variance of Firm Effects: 0.010306

Covariance of Firm and Person Effects: 0.0046087

Variance of Person Effects: 0.085253

Correlation of Firm, Person Effects: 0.15548

When T=2 (i.e., the underlying matched employer-employee data spans only two years), as in this example, it turns out that the KSS-adjusted variance of firm effects and covariance of firm and worker effects are robust to any arbitrary correlation between  $\varepsilon_{g2}$  and  $\varepsilon_{g1}$ .

### 8 Variance of Person Effects When Leaving Out a Match

By leaving a match-out, we can bias-correct the variance of firm effects and the covariance of worker and firm effects while allowing for unrestricted hetoreskedasticity and serial correlation of the error term  $\varepsilon_{gt}$  within each worker-firm match.

However, the person effects,  $\alpha_g$ , of "stayers" — workers that never leave a particular firm — are not leave-match-out estimable. This implies that we cannot compute an unbiased estimate of  $\Omega_g = Var(\varepsilon_{g1},...,\varepsilon_{gT_g})$  for stayers. An estimate of  $\Omega_g$  for both stayers and movers is required in order to provide a bias-correction for the variance of person effects; see Section 1 and Remark 3 in KSS.

The current implementation of the code estimates  $\Omega_g$  for stayers by leaving only a single observation out, that is, by assuming  $\Omega_g$  is diagonal. This approach yields an upper bound estimate on the variance of person effects (computed across both stayers and movers).

There are several alternatives that the user can explore:

- 1. Estimate a variance decomposition in a sample of movers only: For movers, it is possible to estimate a leave-out bias-corrected variance of person effects that is robust to both unrestricted heteroskedasticity and serial correlation in the error term of the AKM model within a given match. Therefore, one can provide an unbiased variance decomposition of all the three components of the two-way fixed effects model by simply feeding to the function leave\_out\_KSS a movers-only sample.
- 2. Drop adjacent wage observations for stayers: Under the assumption that the errors are serially independent after m periods, it suffices to keep every mth stayer observation and apply the estimator after leaving a person-year observation out. For example, if m=2 and we have a balanced panel with T=5, we can restore independence of the errors in the stayer sample by keeping any of the following pairs of stayer time periods: (1,4), (2,5), (1,5). One can choose randomly from the available pairs for each stayer with equal probability.
- 3. Drop interior wage observations for stayers: To minimize concerns regarding serial correlation, the user can drop all but the first and last wage observations of each stayer. Note that dropping stayer wage observations reduces their weight in the variance components. Future versions of the code will allow the variance components to be defined in terms of weights other than the number of micro-observations.

<sup>&</sup>lt;sup>1</sup>This is because leaving a match-out means leaving *all* the observations associated with a stayer and therefore we cannot estimate her  $\alpha_g$ .

# 9 Regressing Firm Fixed Effects on Observables

It is common in empirical applications to regress the fixed effects estimated from the two-way model on some observable characteristics. Using the AKM model again as our leading example, suppose that we are interested in the linear projection of the firm effects  $\psi_{gt}$  on some observables  $Z_{gt}$ 

$$\psi_{i(g,t)} = Z'_{gt} \gamma + e_{gt}. \tag{18}$$

The standard practice is to estimate  $\gamma$  using a simple regression where the estimated firm effects,  $\hat{\psi}_{i(g,t)}$ , are regressed on  $Z_{gt}$ 

$$\hat{\gamma} = \left(\sum_{g,t} Z_{gt} Z_{gt}'\right)^{-1} Z_{gt} \hat{\psi}_{gt}. \tag{19}$$

KSS show that inference on  $\hat{\gamma}$  needs to be adjusted because the estimated firm fixed effects  $\{\hat{\psi}_j\}_{j=1}^J$  are correlated with one another.

To see this, suppose that we have a simple AKM model with only two time periods, set  $w_{gt}=0$ , and take first differences  $\Delta y_g \equiv y_{g2}-y_{g1}$  to eliminate the worker fixed effects so that the AKM model becomes

$$\Delta y_{g} = \Delta f_{g}' \psi + \varepsilon_{g}, \tag{20}$$

where  $\Delta f_g = f_{g,2} - f_{g,1}$  and  $f_{gt} = \{\mathbf{1}_{j(g,t)=1},...,\mathbf{1}_{j(g,t)=J}\}$  is the vector containing the firm dummies. In this model,

$$\hat{\psi} = \psi + \underbrace{\sum_{g=1}^{N} (\Delta f_g \Delta f_g')^{-1} \Delta f_g \varepsilon_g}_{\text{Correlated Noise}}.$$
(21)

Note how the dependence in the vector of estimated firm fixed effects,  $\hat{\psi}$ , is induced by the regressor design  $\sum_{g=1}^{N} (\Delta f_g \Delta f_g')^{-1}$ . As shown in Table 3 of KSS, ignoring this correlation can easily lead to underestimating standard errors by an order of magnitude in practice.

The package provides the HU standard errors on  $\hat{\gamma}$  using the function lincom\_KSS, which is designed to emulate the Stata function lincom and therefore works as a post-estimation command. We demonstrate the functioning of lincom\_KSS with an example.

In this example, we are interested in testing whether the difference in person-year weighted mean firm effects between region 1 and region 2 is statistically different from zero. This amounts to running a regression where the dependent variable is the vector of estimated firm effects and the set of observables,  $Z_{gt}$ , is here represented by a constant and a dummy for whether the firm of worker g in year t belongs to region 2.

The resulting coefficient (and standard error) can be computed by calling the function  $leave_out_KSS$  specifying that we want to run the lincom option and using the region dummy as  $Z_{gt}$  (the constant is automatically added by the code).

```
[9]: \( \text{Regressing firm effects on observables} \)
    namesrc='data/lincom.csv'; %testing data for the lincom function
    data=importdata(namesrc);
    id=data(:,1);
    firmid=data(:,2);
    y=data(:,5);
    region=data(:,4); %Region indicator. Value -1 for region 1, Value 1 for region 2;
    region_dummy=region;
    region_dummy(region_dummy==-1)=0; %Make it a proper dummy variable
    %Run the KSS correction and "lincom"
    labels_lincom={'Region 2 Dummy'}; %give me the label of the columns of Z.
    lincom_do=1; "tell the function leave_out_KSS that we want to project the firm_
     \rightarroweffects on some Z.
    Z=region_dummy; %we're going to project the firm effects on a constant + the
     →region dummy. Constant automatically added by the code
    "Ready to call KSS with lincom option!
    [sigma2_psi,sigma_psi_alpha,sigma2_alpha] =_
     →leave_out_KSS(y,id,firmid,[],[],[],lincom_do,Z,labels_lincom);
    _*_*_*_*_*_*
   Running KSS Correction with the following options
   Leave Out Strategy: Leave match out
   Algorithm for Computation of Statistical Leverages: JLA with 200 simulations.
    _*_*_*_*_*
    _*_*_*_*_*
   SECTION 1
    _*-*-*-*-*-*-*-*-*-*-*-*
    _*-*-*-*-*-*-*-*-*-*-*-*
   Info on the leave one out connected set:
    _*-*-*-*-*-*-*-*-*-*-*-*
   mean wage: 4.7047
   variance of wage: 0.14653
   # of Movers: 9972
   # of Firms: 2974
   # of Person Year Observations: 89666
    _*_*_*_*_*
    _*_*_*_*_*
   SECTION 2
    _*-*-*-*-*-*-*-*-*-*-*-*
    _*_*_*_*_*
   Calculating the statistical leverages of the AKM model...
   Running JLA Algorithm...
   Done!
   Elapsed time is 9.324213 seconds.
    _*-*-*-*-*-*-*-*-*-*-*-*
```

```
_*-*-*-*-*-*-*-*-*-*-*-*-*-*
SECTION 3
_*-*-*-*-*-*-*-*-*-*-*-*
_*_*_*_*_*
PLUG-IN ESTIMATES (BIASED)
Variance of Firm Effects: 0.060695
Covariance of Firm, Person Effects: -0.012603
Variance of Person Effects: 0.10318
Correlation of Firm, Person Effects: -0.15926
_*_*_*_*_*_*
_*-*-*-*-*-*-*-*-*-*-*-*-*
BIAS CORRECTED ESTIMATES
Variance of Firm Effects: 0.044613
Covariance of Firm and Person Effects: 0.0025688
Variance of Person Effects: 0.079191
Correlation of Firm, Person Effects: 0.043218
Regressing the firm effects on observables...
pcg converged at iteration 115 to a solution with relative residual 8.6e-11.
************
************
RESULTS ON LINCOM
***********
***********
Coefficient on Region 2 Dummy: 0.25982
Robust "White" Standard Error: 0.050155
KSS Standard error: 0.088374
T-stat: 2.94
***********
```

We can see from the above output (make sure to scroll until the end) that the difference in person-year weighted mean firm effects between the two regions is equal to 0.26. The traditional HC or "robust" standard errors on this coefficient is around 0.05 while the HU standard error derived in KSS is roughly twice as large (0.09).

#### References

- Abowd, J. M., F. Kramarz, and D. N. Margolis (1999). High wage workers and high wage firms. *Econometrica* 67(2), 251–333.
- Achlioptas, D. (2003). Database-friendly random projections: Johnson-lindenstrauss with binary coins. *Journal of computer and System Sciences* 66(4), 671–687.
- Cattaneo, M. D., M. Jansson, and W. K. Newey (2018). Inference in linear regression models with many covariates and heteroscedasticity. *Journal of the American Statistical Association* 113(523), 1350–1361.
- Johnson, W. B. and J. Lindenstrauss (1984). Extensions of lipschitz mappings into a hilbert space. *Contemporary mathematics* 26(189-206), 1.

- Kline, P., R. Saggio, and M. Sølvsten (2020). Leave-out estimation of variance components. *Econometrica* 88(5), 1859–1898.
- MacKinnon, J. G. and H. White (1985). Some heteroskedasticity-consistent covariance matrix estimators with improved finite sample properties. *Journal of econometrics* 29(3), 305–325.
- White, H. (1980). A heteroskedasticity-consistent covariance matrix estimator and a direct test for heteroskedasticity. *Econometrica*, 817–838.