

# Online Appendix to Simulating Collusion: Challenging Conventional Estimation Methods

Nicole Bellert<sup>1,2\*</sup> and Andrea Günster<sup>3\*</sup>

<sup>1</sup>Institute for Wealth & Asset Management, Zurich University of Applied Sciences (ZHAW), Gertrudstrasse 8, Winterthur, 8401, Zurich, Switzerland.

<sup>2</sup>Department of Informatics, University of Zurich, Binzmühlestrasse 14, Zurich, 8050, Zurich, Switzerland.

<sup>3</sup>Institute of Business Information Technology, Zurich University of Applied Sciences (ZHAW), Theaterstrasse 17, Winterthur, 8400, Zurich, Switzerland.

\*Corresponding author(s). E-mail(s): [bell@zhaw.ch](mailto:bell@zhaw.ch); [gues@zhaw.ch](mailto:gues@zhaw.ch);

This online appendix contains summary statistics and result tables for the individual models we simulated. Linear regression, hazard rate (HR) estimation, Lasso cross-validation (CV) Regression ([Tibshirani \(1996\)](#)), and regressions corrected for Heckman Sample Selection ([Heckman \(1979\)](#)) are applied on data simulated based on Model I ([Stigler \(1964\)](#)), II ([Stigler \(1964\)](#)) and [Harrington and Wei \(2017\)](#), and III ([Stigler \(1964\)](#) and [Bos et al \(2018\)](#)).

## Appendix A Summary Statistics

**Table A1** Summary Statistics Population (Detected and Undetected Cartels) of Model I

	Mean	Median	SD	Min	Max	Skew	N
Number of Firms $n_f$	4.74	5	1.52	2	8	-0.59	77'822
Detection Probability $\sigma$	0.24	0.25	0.08	0.10	0.35	-0.22	77'822
Start	442.35	427	308.34	1	1000	0.12	77'822
End	545.66	539	306.41	1	1000	0	77'822
Duration	104.31	35	159.45	1	1000	2.52	77'822
Ln(Duration+1)	3.30	3.58	1.91	0.69	6.91	-0.08	77'822
Detected	0.47	0	0.50	0	1	0.12	77'822
Times Caught	2.12	1	2.54	0	18	1.30	77'822
Repeat Offender	0.47	0	0.50	0	1	0.14	77'822

For Model I, this table presents the following summary statistics of the 77'822 simulated detected and undetected cartels within a total period of 1'000 time units: the exogenous industry and enforcement characteristics, the start and end date of the cartel, the duration, if it got detected, detected repeatedly, and how often it got detected.

**Table A2** Summary Statistics Sample (Detected Cartels) of Model I

	Mean	Median	SD	Min	Max	Skew	N
Number of Firms $n_f$	3.79	4	1.33	2	8	0.15	36'615
Detection Probability $\sigma$	0.26	0.25	0.08	0.10	0.35	-0.46	36'615
Start	332.28	289	289.08	1	999	0.43	36'615
End	486.22	479	291.45	1	1000	0.06	36'615
Duration	154.94	103	158.53	1	1000	1.80	36'615
Ln(Duration+1)	4.47	4.64	1.24	0.69	6.91	-0.70	36'615
Detected	1	1	0	1	1		36'615
Times Caught	3.45	3	2.31	1	18	1.13	36'615
Repeat Offender	0.77	1	0.42	0	1	-1.29	36'615

For Model I, this table presents the following summary statistics of the 36'615 simulated detected cartels within a total period of 1'000 time units: the exogenous industry and enforcement characteristics, the start and end date of the cartel, the duration, if it got detected repeatedly, and how often it got detected.

**Table A3** Summary Statistics Population (Detected and Undetected Cartels) of Model II

	Mean	Median	SD	Min	Max	Skew	N
Number of Firms $n_f$	5.19	6	1.38	2	8	-1.16	53'178
Detection Probability $\sigma$	0.23	0.25	0.09	0.10	0.35	-0.10	53'178
Start	422.15	401	314.10	1	1000	0.18	53'178
End	574.27	571	309.25	1	1000	-0.05	53'178
Duration	153.12	11	264.18	1	1000	2	53'178
Ln(Duration+1)	3.06	2.48	2.23	0.69	6.91	0.35	53'178
Detected	0.22	0	0.41	0	1	1.35	53'178
Times Caught	0.62	0	1.05	0	8	2.01	53'178
Repeat Offender	0.17	0	0.37	0	1	1.77	53'178

For Model II, this table presents the following summary statistics of the 53'178 simulated detected and undetected cartels within a total period of 1'000 time units: the exogenous industry and enforcement characteristics, the start and end date of the cartel, the duration, if it got detected, detected repeatedly, and how often it got detected.

**Table A4** Summary Statistics Sample (Detected Cartels) of Model II

	Mean	Median	SD	Min	Max	Skew	N
Number of Firms $n_f$	3.82	4	1.33	2	7	0.12	11'733
Detection Probability $\sigma$	0.26	0.25	0.08	0.10	0.35	-0.43	11'733
Start	199.50	24	263.67	1	999	1.14	11'733
End	487.99	480	291.23	1	1000	0.06	11'733
Duration	289.49	219	246.85	1	1000	0.93	11'733
Ln(Duration+1)	5.14	5.39	1.24	0.69	6.91	-1.04	11'733
Detected	1	1	0	1	1		11'733
Times Caught	1.77	1	1.02	1	8	1.55	11'733
Repeat Offender	0.48	0	0.50	0	1	0.10	11'733

For Model II, this table presents the following summary statistics of the 11'733 simulated detected cartels within a total period of 1'000 time units: the exogenous industry and enforcement characteristics, the start and end date of the cartel, the duration, if it got detected repeatedly, and how often it got detected.

**Table A5** Summary Statistics Population (Detected and Undetected Cartels) of Model III

	Mean	Median	SD	Min	Max	Skew	N
Number of Firms $n_f$	3.62	4	1.28	2	7	0.16	809'858
Fines $\gamma$ (% of Profit)	0.80	0.80	0.08	0.70	0.90	0.02	809'858
Leniency (% of Fine) $\theta$	0.40	0.50	0.41	0	1	0.37	809'858
Detection Probability $\sigma$	0.22	0.20	0.08	0.10	0.35	0.08	809'858
Structured	0.48	0	0.50	0	1	0.07	809'858
Start	416.30	397	317.51	1	1000	0.19	809'858
End	524.04	516	315.19	1	1000	0.02	809'858
Duration	108.75	43	162.41	1	1000	2.60	809'858
Ln(Duration+1)	3.50	3.78	1.79	0.69	6.91	-0.20	809'858
Detected	0.53	1	0.50	0	1	-0.11	809'858
Times Caught	2.38	1	2.82	0	25	1.64	809'858
Repeat Offender	0.49	0	0.50	0	1	0.02	809'858

For Model III, this table presents the following summary statistics of the 809'858 simulated detected and undetected cartels within a total period of 1'000 time units: the exogenous industry and enforcement characteristics, the start and end date of the cartel, the duration, if it got detected, detected repeatedly, and how often it got detected.

**Table A6** Summary Statistics Sample (Detected Cartels) of Model III

	Mean	Median	SD	Min	Max	Skew	N
Number of Firms $n_f$	2.96	3	0.98	2	7	0.70	427'108
Fines $\gamma$ (% of Profit)	0.80	0.80	0.08	0.70	0.90	0.03	427'108
Leniency (% of Fine) $\theta$	0.52	0.50	0.41	0	1	-0.06	427'108
Detection Probability $\sigma$	0.23	0.25	0.08	0.10	0.35	-0.07	427'108
Structured	0.49	0	0.50	0	1	0.04	427'108
Start	331.62	286	295.59	1	1000	0.43	427'108
End	476.10	464	291.09	1	1000	0.10	427'108
Duration	145.48	91	158.39	1	1000	1.97	427'108
Ln(Duration+1)	4.36	4.52	1.27	0.69	6.91	-0.60	427'108
Detected	1	1	0	1	1		427'108
Times Caught	3.52	3	2.71	1	25	1.59	427'108
Repeat Offender	0.74	1	0.44	0	1	-1.10	427'108

For Model III, this table presents the following summary statistics of the 427'108 simulated detected cartels within a total period of 1'000 time units: the exogenous industry and enforcement characteristics, the start and end date of the cartel, the duration, if it got detected repeatedly, and how often it got detected.

## Appendix B Linear Models Results

**Table B7** Sample Selection Bias Linear Regression - Model I

Coefficients	$\beta^s$	$\beta^l$	$\alpha_{IMR}$	$bias_{IMR}$
Number of Firms $n_f$	-0.12	-0.02	0.15	-0.10
Start	0	0	0	0
Detection Probability $\sigma$	-2.93	-3.11	-0.27	0.18
Times Caught	0.04	-0.03	-0.10	0.07
Repeat Offender	0.21	-0.18	-0.57	0.39
IMR		-0.68		
Constant	5.79	5.87	0.11	-0.07

This table shows, for the Linear Regression of Model I, the sample selection bias that we correct with the IMR following (Heckman (1979)).  $\beta^s$  is the estimated coefficients in the short model without IMR.  $\beta^l$  is the estimated coefficients in the corrected long model including IMR.  $\alpha_{IMR}$  is the coefficient in the auxiliary regression between each variable and IMR. The last column shows the sample selection bias:  $bias_{IMR} = \beta^l(IMR) * \alpha_{IMR}$ .

**Table B8** Sample Selection Bias Linear Regression - Model II

Coefficients	$\beta^s$	$\beta^l$	$\alpha_{IMR}$	$bias_{IMR}$
Number of Firms $n_f$	-0.16	-0.09	0.14	-0.07
Start	0	0	0	0
Detection Probability $\sigma$	-1.43	-1.67	-0.46	0.25
Times Caught	0.07	-0.17	-0.44	0.24
Repeat Offender	0.27	0.18	-0.16	0.09
IMR		-0.54		
Constant	6.26	6.56	0.56	-0.30

This table shows, for the Linear Regression of Model II, the sample selection bias that we correct with the IMR following (Heckman (1979)).  $\beta^s$  is the estimated coefficients in the short model without IMR.  $\beta^l$  is the estimated coefficients in the corrected long model including IMR.  $\alpha_{IMR}$  is the coefficient in the auxiliary regression between each variable and IMR. The last column shows the sample selection bias:  $bias_{IMR} = \beta^l(IMR) * \alpha_{IMR}$ .

**Table B9** Sample Selection Bias Linear Regression - Model III

Coefficients	$\beta^s$	$\beta^l$	$\alpha_{IMR}$	$bias_{IMR}$
Number of Firms $n_f$	-0.18	-0.01	0.25	-0.17
Fines $\gamma$ (% of Profit)	0.06	0.07	0.01	-0.01
Leniency (% of Fine) $\theta$	0.10	0.01	-0.13	0.09
Start	0	0	0	0
Structured	-0.40	-0.36	0.06	-0.04
Detection Probability $\sigma$	-3.85	-3.31	0.79	-0.54
Times Caught	-0.01	-0.05	-0.06	0.04
Repeat Offender	0.01	-0.37	-0.56	0.38
IMR		-0.68		
Constant	6.25	6.02	-0.33	0.23

This table shows, for the Linear Regression of Model II, the sample selection bias that we correct with the IMR following (Heckman (1979)).  $\beta^s$  is the estimated coefficients in the short model without IMR.  $\beta^l$  is the estimated coefficients in the corrected long model including IMR.  $\alpha_{IMR}$  is the coefficient in the auxiliary regression between each variable and IMR. The last column shows the sample selection bias:  $bias_{IMR} = \beta^l(IMR) * \alpha_{IMR}$ .

**Table B10** Sample Selection Bias Linear Regression - Models I, II and III

Coefficients	$\beta^s$	$\beta^l$	$\alpha_{IMR}$	$bias_{IMR}$
Number of Firms $n_f$	-0.16	0	0.23	-0.16
Fines $\gamma$ (% of Profit)	0.06	0.07	0.01	-0.01
Leniency (% of Fine) $\theta$	0.09	0	-0.13	0.09
Detection Probability $\sigma$	-3.66	-3.24	0.61	-0.43
Start	-0.40	-0.36	0	0
Structured	0.53	0.67	0.05	-0.04
Model II	-0.24	-0.07	0.19	-0.14
Model III	0	0	0.24	-0.17
Times Caught	0	-0.05	-0.07	0.05
Repeat Offender	0.04	-0.37	-0.58	0.41
IMR		-0.71		
Constant	6.38	6.06	-0.45	0.32

This table shows, for the Linear Regression of the combined Models I, II, and III, the sample selection bias that we correct with the inverse Mill's ratio (IMR) following (Heckman (1979)).  $\beta^s$  is the estimated coefficients in the short model without IMR.  $\beta^l$  is the estimated coefficients in the corrected long model including IMR.  $\alpha_{IMR}$  is the coefficient in the auxiliary regression between each variable and IMR. The last column shows the sample selection bias:  $bias_{IMR} = \beta^l(IMR) * \alpha_{IMR}$ .

**Table B11** Linear Regression and HR for Cartel Duration on Model I - ICC on Stigler - Detection independent of Collusion

	mlrSample	Ln(Duration+1)		mlrHeck	HRSample	Cartel Death		HRHeck
		mlrUndetect	mlrCartels			HRUndetect	HRCartels	
N Firms $n_f$	-0.12*** (0.005)	-0.57*** (0.01)	-0.39*** (0.004)	-0.02*** (0.01)	0.09*** (0.004)	0.79*** (0.01)	0.30*** (0.003)	0.01* (0.01)
Start	-0.001*** (0.0000)	-0.001*** (0.0000)	-0.002*** (0.0000)	-0.0002*** (0.0001)	0.001*** (0.0000)	0.001*** (0.0000)	0.001*** (0.0000)	0.001*** (0.0001)
Detection Prob. $\sigma$	-2.93*** (0.08)	-1.60*** (0.08)	-2.61*** (0.06)	-3.11*** (0.08)	2.78*** (0.07)	1.42*** (0.06)	2.73*** (0.05)	2.94*** (0.07)
Times Caught	0.04*** (0.004)	-0.01 (0.01)	0.13*** (0.003)	-0.03*** (0.01)	-0.03*** (0.004)	-0.56*** (0.02)	-0.13*** (0.003)	0.02*** (0.005)
Repeat Offender	0.21*** (0.02)	2.06*** (0.03)	1.20*** (0.02)	-0.18*** (0.03)	-0.22*** (0.02)	-0.19*** (0.05)	-0.59*** (0.01)	0.08*** (0.02)
IMR				-0.68*** (0.03)				0.53*** (0.03)
Constant	5.79*** (0.03)	6.13*** (0.05)	5.81*** (0.03)	5.87*** (0.03)				
Observations	36'615	41'207	7'822	36'615	36'615	41'207	7'822	36'615
R <sup>2</sup>	0.10	0.47	0.50	0.11				
Adjusted R <sup>2</sup>	0.10	0.47	0.50	0.11				
Log Likelihood					-218'572.00	-114'465.60	-350'666.30	-218'383.30

Note: This table shows the estimation results of linear cross-sectional regressions to explain cartel duration (ln(duration+1)) and the estimation results of a Weibull Hazard Model to explain cartel death, both at the industry level, for data simulated for Model I. Columns 2 - 5 estimate linear regression coefficients, while columns 6 - 9 estimate HR coefficients, both on the sample of detected cartels, the group of undetected cartels, the population of all cartels, and the sample corrected for Heckman Sample Selection, respectively. The estimated coefficients show standard errors in the sample, but do not test for the real population. The estimated HR coefficients show the change of risk for cartel breakdown if the covariate increases by 1 unit, keeping all others fixed.

Standard errors are in parentheses. Significance at the 1%, 5%, and 10% level is indicated by \*\*\*, \*\*, and \*, respectively.

**Table B12** Linear Regression and HR for Cartel Duration on Model II - ICC on Stigler - Detection depends on number of Firms

	Ln(Duration+1)				Cartel Death		HRHeck
	mlrSample	mlrUndetect	mlrCartels	mlrHeck	HRUndetect	HRCartels	
N Firms $n_f$	-0.16*** (0.01)	-0.97*** (0.01)	-0.72*** (0.01)	-0.09*** (0.01)	0.96*** (0.01)	0.56*** (0.01)	0.07*** (0.01)
Start	-0.002*** (0.0001)	-0.002*** (0.0000)	-0.002*** (0.0000)	-0.001*** (0.0001)	0.001*** (0.0000)	0.001*** (0.0000)	0.002*** (0.0001)
Detection Prob. $\sigma$	-1.43*** (0.14)	-0.86*** (0.09)	-1.17*** (0.08)	-1.67*** (0.14)	0.67*** (0.06)	1.12*** (0.06)	1.35*** (0.12)
Times Caught	0.07*** (0.02)	0.75*** (0.02)	0.68*** (0.01)	-0.17*** (0.03)	-1.52*** (0.03)	-0.49*** (0.01)	0.10*** (0.03)
Repeat Offender	0.27*** (0.04)	-0.54*** (0.06)	-0.31*** (0.03)	0.18*** (0.04)	1.26*** (0.08)	0.34*** (0.03)	-0.21*** (0.03)
IMR				-0.54*** (0.06)			0.33*** (0.05)
Constant	6.26*** (0.05)	8.82*** (0.05)	7.60*** (0.04)	6.56*** (0.06)			
Observations	11'733	41'445	53'178	11'733	41'445	53'178	11'733
R <sup>2</sup>	0.17	0.51	0.56	0.17			
Adjusted R <sup>2</sup>	0.16	0.51	0.56				
Log Likelihood				-77'014.56	-120'764.40	-209'666.50	-76'995.30

Note: This table shows the estimation results of linear cross-sectional regressions to explain cartel duration (ln(duration+1)) and the estimation results of a Weibull Hazard Model to explain cartel death, both at the industry level, for data simulated for Model II. Columns 2 - 5 estimate linear regression coefficients, while columns 6 - 9 estimate HR coefficients, both on the sample of detected cartels, the group of undetected cartels, the population of all cartels, and the sample corrected for Heckman Sample Selection, respectively. The estimated coefficients show standard errors in the sample, but do not test for the real population. The estimated HR coefficients show the change of risk for cartel breakdown if the covariate increases by 1 unit, keeping all others fixed.

Standard errors are in parentheses. Significance at the 1%, 5%, and 10% level is indicated by \*\*\*, \*\*, and \*, respectively.



**Table B13** Linear Regression and HR for Cartel Duration on Model III - ICC on Harrington et al.

	Ln(Duration+1)		Cartel Death		HRHeck
	mlrSample	mlrUndetect	mlrCartels	mlrHeck	
N Firms $n_f$	-0.18*** (0.002)	-0.56*** (0.003)	-0.49*** (0.002)	-0.01** (0.003)	-0.001 (0.003)
Fines $\gamma$ (% of Profit)	0.06*** (0.02)	0.02 (0.03)	0.02 (0.02)	0.07*** (0.02)	-0.11*** (0.02)
Leniency (% Fine) $\theta$	0.10*** (0.004)	1.59*** (0.01)	0.87*** (0.004)	0.01 (0.005)	0.02*** (0.004)
Start	-0.001*** (0.0000)	-0.001*** (0.0000)	-0.002*** (0.0000)	-0.0002*** (0.0000)	0.001*** (0.0000)
Structured	-0.40*** (0.004)	-0.52*** (0.005)	-0.56*** (0.003)	-0.36*** (0.004)	0.34*** (0.003)
Detection Prob. $\sigma$	-3.85*** (0.02)	-4.83*** (0.03)	-5.16*** (0.02)	-3.31*** (0.03)	3.10*** (0.02)
Times Caught	-0.01*** (0.001)	-0.004** (0.002)	0.04*** (0.001)	-0.05*** (0.001)	0.04*** (0.001)
Repeat Offender	0.01** (0.01)	1.16*** (0.01)	0.76*** (0.005)	-0.37*** (0.01)	0.29*** (0.01)
IMR				-0.68*** (0.01)	0.54*** (0.01)
Constant	6.25*** (0.02)	6.04*** (0.03)	6.53*** (0.02)	6.02*** (0.02)	
Observations	427'108	382'750	809'858	427'108	427'108
R <sup>2</sup>	0.14	0.43	0.41	0.14	
Adjusted R <sup>2</sup>	0.14	0.43	0.41	0.14	
Log Likelihood				-2'509'783.00	-2'507'956.00

Note: This table shows the estimation results of linear cross-sectional regressions to explain cartel duration (ln(duration+1)) and the estimation results of a Weibull Hazard Model to explain cartel death, both at the industry level, for data simulated for Models IIIa and IIIb. Columns 2 - 5 estimate linear regression coefficients, while columns 6 - 9 estimate HR coefficients, both on the sample of detected cartels, the group of undetected cartels, the population of all cartels, and the sample corrected for Heckman Sample Selection, respectively. The estimated coefficients show standard errors in the sample, but do not test for the real population. The estimated HR coefficients show the change of risk for cartel breakdown if the covariate increases by 1 unit, keeping all others fixed. Standard errors are in parentheses. Significance at the 1%, 5%, and 10% level is indicated by \*\*\*, \*\*, and \*, respectively.

## Appendix C Lasso Results

**Table C14** Sample Selection Bias Lasso Regression - Model I

Coefficients	$\beta^s$	$\beta^l$	$\alpha_{IMR}$	$bias_{IMR}$
Number of Firms $n_f$	0.26	0.16	-0.22	0.11
$n_f^3$	-0.01	0	0.01	0
Detection Probability $\sigma$	-4.01	-4.25	-0.50	0.24
$\sigma^3$	3.93	5.97	4.30	-2.04
$n_f\sigma$	0.12	0.04	-0.17	0.08
Start	0	0	0	0
Times Caught	0.02	-0.02	-0.09	0.04
Repeat Offender	0.18	-0.05	-0.48	0.23
IMR		-0.47		
Constant	5.09	5.54	0.95	-0.45

This table shows, for the Lasso CV Linear Regression of Model I, the sample selection bias that we correct with the IMR following (Heckman (1979)).  $\beta^s$  is the estimated coefficients in the short model without IMR.  $\beta^l$  is the estimated coefficients in the corrected long model including IMR.  $\alpha_{IMR}$  is the coefficient in the auxiliary regression between each variable and IMR. The last column shows the sample selection bias:  $bias_{IMR} = \beta^l(IMR) * \alpha_{IMR}$ .

**Table C15** Sample Selection Bias Lasso Regression - Model II

Coefficients	$\beta^s$	$\beta^l$	$\alpha_{IMR}$	$bias_{IMR}$
Number of Firms $n_f$	0.33	0.27	-0.20	0.05
$n_f^3$	-0.01	-0.01	0.01	0
Detection Probability $\sigma$	-2.34	-2.64	-1.09	0.30
$\sigma^3$	3.19	4.81	5.96	-1.62
$n_f\sigma$	0.11	0.08	-0.13	0.04
Start	0	0	0	0
Times Caught	0.04	-0.08	-0.42	0.11
Repeat Offender	0.21	0.17	-0.15	0.04
IMR		-0.27		
Constant	5.32	5.71	1.44	-0.39

This table shows, for the Lasso CV Linear Regression of Model II, the sample selection bias that we correct with the IMR following (Heckman (1979)).  $\beta^s$  is the estimated coefficients in the short model without IMR.  $\beta^l$  is the estimated coefficients in the corrected long model including IMR.  $\alpha_{IMR}$  is the coefficient in the auxiliary regression between each variable and IMR. The last column shows the sample selection bias:  $bias_{IMR} = \beta^l(IMR) * \alpha_{IMR}$ .

**Table C16** Sample Selection Bias Lasso Regression - Model III

Coefficients	$\beta^s$	$\beta^l$	$\alpha_{IMR}$	$bias_{IMR}$
Number of Firms $n_f$	0.40	0.19	-0.46	0.21
$n_f^3$	-0.01	-0.01	0.01	-0.01
Detection Probability $\sigma$	-5.44	-9.09	-8.03	3.65
$\sigma^2$	11.03	22.02	24.20	-10.99
$\sigma^3$	-6.74	-20.94	-31.27	14.20
$n_f\sigma$	-0.83	-0.36	1.02	-0.46
$\gamma^3$	0.03	0.03	0.01	0
Leniency (% of Fine) $\theta$	-0.22	-0.14	0.17	-0.08
$n_f\theta$	0.11	0.06	-0.11	0.05
Structured	-0.39	-0.37	0.04	-0.02
Start	0	0	0	0
Times Caught	-0.01	-0.04	-0.07	0.03
Repeat Offender	-0.04	-0.26	-0.48	0.22
IMR		-0.45		
Constant	5.41	6.18	1.68	-0.76

This table shows, for the Lasso CV Linear Regression of Model III, the sample selection bias that we correct with the IMR following (Heckman (1979)).  $\beta^s$  is the estimated coefficients in the short model without IMR.  $\beta^l$  is the estimated coefficients in the corrected long model including IMR.  $\alpha_{IMR}$  is the coefficient in the auxiliary regression between each variable and IMR. The last column shows the sample selection bias:  $bias_{IMR} = \beta^l(IMR) * \alpha_{IMR}$ .

**Table C17** Sample Selection Bias Lasso Regression - Models I, II and III

Coefficients	$\beta^s$	$\beta^l$	$\alpha_{IMR}$	$bias_{IMR}$
Number of Firms $n_f$	0.16	0.02	-0.23	0.14
$n_f^3$	-0.01	0	0.01	-0.01
Detection Probability $\sigma$	-8.44	-11.34	-4.81	2.90
$\sigma^2$	14.02	26.38	20.47	-12.36
$\sigma^3$	-8.44	-25.83	-28.81	17.40
$n_f\sigma$	-0.07	0.09	0.28	-0.17
Fines $\gamma$ (% of Profit)	1.77	1.66	-0.17	0.10
$\gamma^3$	-1.43	-1.34	0.14	-0.09
Leniency (% of Fine) $\theta$	0.03	-0.05	-0.14	0.08
$n_f\theta$	0.02	0.02	0	0
Structured	-0.37	-0.36	0.02	-0.01
Model II	0.54	0.65	0.18	-0.11
Model III	-0.67	-0.44	0.37	-0.23
Start	0	0	0	0
Times Caught	0	-0.04	-0.07	0.04
Repeat Offender	0.02	-0.30	-0.53	0.32
IMR		-0.60		
Constant	6.25	6.75	0.82	-0.49

This table shows, for the Lasso CV Linear Regression of the combined Models I, II, and III, the sample selection bias that we correct with the IMR following (Heckman (1979)).  $\beta^s$  is the estimated coefficients in the short model without IMR.  $\beta^l$  is the estimated coefficients in the corrected long model including IMR.  $\alpha_{IMR}$  is the coefficient in the auxiliary regression between each variable and IMR. The last column shows the sample selection bias:  $bias_{IMR} = \beta^l(IMR) * \alpha_{IMR}$ .

**Table C18** Lasso CV Regression and HR for Cartel Duration on Model I - ICC on Stigler - Detection independent of Collusion

	Ln(Duration+1)		Cartel Death					
	LasSample	LasUndetec	LasCartels	LasHeck	HRLasSample	HRLasUnd	HRLasCartels	HRLasHeck
Start	-0.001*** (0.0000)	-0.001*** (0.0000)	-0.002*** (0.0000)	-0.0004*** (0.0001)	0.001*** (0.0000)	0.001*** (0.0000)	0.001*** (0.0000)	0.001*** (0.0001)
N Firms $n_f$	-1.37*** (0.14)	4.87*** (0.17)	2.50*** (0.09)	0.16*** (0.02)	1.04*** (0.12)	12.61*** (1.42)	-1.49*** (0.07)	-0.14*** (0.02)
$n_f^2$	0.43*** (0.04)	-1.29*** (0.04)	-0.64*** (0.02)		-0.33*** (0.03)	-1.56*** (0.25)	0.38*** (0.02)	
$n_f^3$	-0.04*** (0.003)	0.09*** (0.003)	0.04*** (0.002)	-0.005*** (0.0004)	0.03*** (0.002)	0.07*** (0.01)	-0.02*** (0.001)	0.004*** (0.0004)
Detection Prob. $\sigma$	-5.04*** (2.38)	-17.14*** (2.25)	-12.64*** (1.71)	-4.25*** (0.36)	4.97*** (2.05)	18.48*** (2.12)	12.14*** (1.38)	3.72*** (0.31)
$\sigma^2$	4.74 (10.82)	26.19** (10.53)	13.83* (7.93)		-6.61 (9.34)	-15.51* (9.02)	-15.50** (6.41)	
$\sigma^3$	-2.64 (15.46)	-34.76*** (15.51)	-17.19 (11.53)	5.97*** (1.49)	5.93 (13.35)	22.24* (13.29)	18.06* (9.30)	-4.94*** (1.29)
$n_f \sigma$	0.12*** (0.06)	1.69*** (0.07)	1.39*** (0.04)	0.04 (0.06)	-0.06 (0.05)	-2.31*** (0.16)	-1.12*** (0.03)	0.005 (0.05)
Times Caught	0.02*** (0.004)	0.03*** (0.01)	0.12*** (0.003)	-0.02*** (0.01)	-0.02*** (0.004)	-0.49*** (0.02)	-0.12*** (0.003)	0.01 (0.004)
Repeat Offender	0.17*** (0.02)	1.75*** (0.03)	0.99*** (0.02)	-0.05* (0.02)	-0.19*** (0.02)	-0.23*** (0.05)	-0.51*** (0.01)	-0.03 (0.02)
IMR				-0.47*** (0.03)				0.33*** (0.03)
Constant	7.03*** (0.24)	1.06*** (0.28)	3.17*** (0.17)	5.54*** (0.09)				
Observations	36'615	41'207	7'822	36'615	36'615	41'207	7'822	36'615
R <sup>2</sup>	0.12	0.49	0.53	0.13				
Adjusted R <sup>2</sup>	0.12	0.49	0.53	0.13				
Log Likelihood					-218'190.40	-113'185.60	-348'310.50	-218'178.10

Note: This table shows the estimation results of Lasso CV linear cross-sectional regressions to explain cartel duration (ln(duration+1)) and the estimation results of a Lasso CV Weibull Hazard Model to explain cartel death, both at the industry level, for data simulated for Model I. Columns 2 - 5 estimate linear regression coefficients, while columns 6 - 9 estimate HR coefficients, both on the sample of detected cartels, the group of undetected cartels, the population of all cartels, and the sample corrected for Heckman Sample Selection, respectively. The estimated coefficients show standard errors in the sample, but do not test for the real population. The estimated HR coefficients show the change of risk for cartel breakdown if the covariate increases by 1 unit, keeping all others fixed. Standard errors are in parentheses. Significance at the 1%, 5%, and 10% level is indicated by \*\*\*, \*\*, and \*, respectively.

**Table C19** Lasso CV Regression and HR for Cartel Duration on Model II - ICC on Stigler - Detection depends on number of Firms

	Ln(Duration+1)		Cartel Death		Cartel Death			
	LasSample	LasUndetec	LasCartels	LasHeck	HRLasSample	HRLasUnd	HRLasCartels	HRLasHeck
Start	-0.002*** (0.0001)	-0.002*** (0.0000)	-0.002*** (0.0000)	-0.001*** (0.0001)	0.002*** (0.0001)	0.001*** (0.0000)	0.001*** (0.0000)	0.002*** (0.0001)
N Firms $n_f$	-1.98*** (0.24)	8.53*** (0.18)	5.88*** (0.13)	0.27*** (0.04)	1.37*** (0.20)	7.67*** (1.48)	-4.30*** (0.11)	-0.26*** (0.03)
$n_f^2$	0.61*** (0.06)	-2.16*** (0.04)	-1.48*** (0.03)		-0.43*** (0.05)	-0.55*** (0.26)	1.08*** (0.03)	
$n_f^3$	-0.06*** (0.01)	0.15*** (0.003)	0.10*** (0.002)	-0.01*** (0.001)	0.04*** (0.004)	-0.0004 (0.01)	-0.07*** (0.002)	0.01*** (0.001)
Detection Prob. $\sigma$	-2.36*** (0.61)	-5.81*** (0.51)	-7.30*** (0.40)	-2.64*** (0.61)	4.17 (3.61)	5.79*** (2.11)	8.33*** (1.68)	1.84*** (0.55)
$\sigma^2$					-11.18 (16.44)	0.25 (9.01)	-5.39 (7.83)	
$\sigma^3$	3.12 (2.52)	-2.29 (1.64)	-0.19 (1.46)	4.81* (2.55)	12.46 (23.52)	1.97 (13.27)	7.29 (11.46)	-4.34* (2.29)
$n_f \sigma$	0.13 (0.10)	0.97*** (0.08)	1.19*** (0.05)	0.08 (0.10)	0.01 (0.09)	-0.94*** (0.16)	-1.11*** (0.05)	0.04 (0.09)
Times Caught	0.02 (0.02)	0.62*** (0.02)	0.54*** (0.01)	-0.08** (0.03)	-0.01 (0.02)	-1.30*** (0.03)	-0.41*** (0.01)	0.03 (0.03)
Repeat Offender	0.20*** (0.03)	-0.44*** (0.05)	-0.23*** (0.03)	0.17*** (0.04)	-0.22*** (0.03)	1.07*** (0.08)	0.26*** (0.03)	-0.21*** (0.03)
IMR				-0.27*** (0.06)				0.12** (0.06)
Constant	7.99*** (0.30)	-2.81*** (0.26)	-0.25 (0.19)	5.71*** (0.16)				
Observations	11'733	41'445	53'178	11'733	11'733	41'445	53'178	11'733
R <sup>2</sup>	0.20	0.54	0.59	0.20				
Adjusted R <sup>2</sup>	0.20	0.54	0.59	0.19				
Log Likelihood					-76'836.46	-118'987.30	-207'949.80	-76'867.61

Note: This table shows the estimation results of Lasso CV linear cross-sectional regressions to explain cartel duration (ln(duration+1)) and the estimation results of a Lasso CV Weibull Hazard Model to explain cartel death, both at the industry level, for data simulated for Model II. Columns 2 - 5 estimate linear regression coefficients, while columns 6 - 9 estimate HR coefficients, both on the sample of detected cartels, the group of undetected cartels, the population of all cartels, and the sample corrected for Heckman Sample Selection, respectively. The estimated coefficients show standard errors in the sample, but do not test for the real population. The estimated HR coefficients show the change of risk for cartel breakdown if the covariate increases by 1 unit, keeping all others fixed. Standard errors are in parentheses. Significance at the 1%, 5%, and 10% level is indicated by \*\*\*, \*\*, and \*, respectively.

**Table C20** Cartel Duration with CV Lasso on Model III - ICC on Harrington et al.

	Ln(Duration+1)		Cartel Death		
	LasSample	LasUndetec	LasCartels	LasHeck	
Start	-0.001*** (0.0000)	-0.001*** (0.0000)	-0.002*** (0.0000)	-0.0004*** (0.0000)	0.001*** (0.0000)
N Firms $n_f$	0.41*** (0.01)	-0.40*** (0.06)	0.58*** (0.04)	0.19*** (0.01)	0.001*** (0.0000)
$n_f^2$		-0.13*** (0.02)	-0.22*** (0.01)		1.96*** (0.03)
$n_f^3$	-0.01*** (0.0002)	0.01*** (0.001)	0.01*** (0.001)	0.08*** (0.01)	0.001*** (0.0000)
Fines $\gamma$ (% of Profit)	1.40*** (0.31)		0.04*** (0.02)	-0.01*** (0.0003)	-0.73*** (0.01)
$\gamma^3$	-0.70*** (0.16)	0.02* (0.01)		-0.02*** (0.001)	0.22*** (0.01)
Leniency (% Fine) $\theta$	-0.22*** (0.01)	1.06*** (0.02)	-0.60*** (0.01)	-0.04*** (0.01)	0.01*** (0.001)
$n_f\theta$	0.11*** (0.004)	0.12*** (0.005)	0.42*** (0.003)	0.14*** (0.01)	0.50*** (0.01)
Structured	-0.39*** (0.004)	-0.49*** (0.004)	0.06*** (0.005)	0.06*** (0.004)	0.03*** (0.002)
Detection Prob. $\sigma$	-3.09*** (0.12)	-8.71*** (0.78)	-11.16*** (0.10)	0.37*** (0.003)	-0.31*** (0.002)
$\sigma^2$		-60.58*** (3.67)		0.52*** (0.01)	0.42*** (0.002)
$\sigma^3$	9.32*** (0.46)	121.65*** (5.49)	21.48*** (0.39)	27.77*** (0.68)	9.44*** (0.09)
$n_f\sigma$	-0.83*** (0.03)	2.69*** (0.03)	0.56*** (0.02)	47.40*** (3.06)	-14.55*** (2.59)
Times Caught	-0.01*** (0.001)	0.05*** (0.002)	0.07*** (0.001)	-122.98*** (3.31)	11.95*** (3.79)
Repeat Offender	-0.04*** (0.01)	1.00*** (0.01)	0.62*** (0.005)	-5.49*** (0.04)	0.29*** (0.03)
IMR				-0.31*** (0.003)	0.03*** (0.001)
Constant	4.52*** (0.16)	7.19*** (0.09)	6.05*** (0.05)	-0.18*** (0.01)	0.20*** (0.01)
Observations	427'108	382'750	809'858	427'108	809'858
R <sup>2</sup>	0.15	0.45	0.45	382'750	427'108
Adjusted R <sup>2</sup>	0.15	0.45	0.45		
Log Likelihood				-1'183'409.00	-2'506'846.00

Note: This table shows the estimation results of Lasso CV linear cross-sectional regressions to explain cartel duration (ln(duration+1)) and the estimation results of a Lasso CV Weibull Hazard Model to explain cartel death, both at the industry level, for data simulated for Models IIIa and IIIb combined. Columns 2 - 5 estimate linear regression coefficients, while columns 6 - 9 estimate HR coefficients, both on the sample of detected cartels, the group of undetected cartels, the population of all cartels, and the sample corrected for Heckman Sample Selection, respectively. The estimated coefficients show standard errors in parentheses. The estimated HR coefficients show the change of risk for cartel breakdown if the covariate increases by 1 unit, keeping all others fixed. Standard errors are in parentheses. Significance at the 1%, 5%, and 10% level is indicated by \*\*\*, \*\*, and \*, respectively.

**Table C21** Cartel Duration with CV Lasso on Model I, II and III

	Ln(Duration+1)			LasHeck	Cartel Death			
	LasSample	LasUndetec	LasCartels		HRLasSample	HRLasUnd	HRLasCartels	
N Firms $n_f$	0.26*** (0.05)	0.45*** (0.05)		0.02** (0.01)	-0.28*** (0.04)	1.15*** (0.05)		-0.01 (0.01)
$n_f^2$	-0.03** (0.01)	-0.30*** (0.01)	-0.11*** (0.001)		0.04*** (0.01)	0.02** (0.01)	0.07*** (0.001)	
$n_f^3$	-0.01*** (0.001)	0.02*** (0.001)	-0.001*** (0.0002)	-0.003*** (0.0002)	0.003** (0.001)	-0.002*** (0.001)	0.003*** (0.0001)	0.002*** (0.0001)
$\gamma^2$				1.66*** (0.38)				-0.64* (0.33)
$\gamma^3$	0.03*** (0.01)	0.02 (0.01)	0.03*** (0.01)	-1.34*** (0.32)	-0.05*** (0.01)	-0.03** (0.01)	-0.03*** (0.01)	0.48* (0.27)
Leniency (% Fine) $\theta$	0.03*** (0.01)	0.66*** (0.02)	-0.55*** (0.01)	-0.05*** (0.01)	-0.05*** (0.01)	-2.69*** (0.03)	0.44*** (0.01)	0.03** (0.01)
$n_f\theta$	0.02*** (0.004)	0.21*** (0.005)	0.40*** (0.003)	0.02*** (0.004)	-0.001 (0.004)	0.29*** (0.01)	-0.28*** (0.002)	-0.002 (0.003)
Detection Prob. $\sigma$	-5.43*** (0.11)	-18.35*** (0.14)	-13.46*** (0.08)	-11.34*** (0.62)	4.77*** (0.09)	23.45*** (0.15)	11.52*** (0.07)	9.07*** (0.53)
$\sigma^2$				26.38*** (2.88)				-18.72*** (2.46)
$\sigma^3$	11.91*** (0.42)	24.02*** (0.50)	23.64*** (0.35)	-25.83*** (4.21)	-10.07*** (0.36)	-29.63*** (0.43)	-22.95*** (0.27)	16.65*** (3.60)
$n_f\sigma$	-0.09*** (0.02)	2.19*** (0.02)	1.26*** (0.01)	0.09*** (0.02)	0.10*** (0.02)	-2.97*** (0.02)	-0.93*** (0.01)	-0.07*** (0.02)
Structured	-0.37*** (0.004)	-0.51*** (0.005)	-0.51*** (0.003)	-0.36*** (0.004)	0.35*** (0.003)	0.35*** (0.004)	0.40*** (0.002)	0.34*** (0.003)
Model II	0.54*** (0.01)	0.30*** (0.01)	0.36*** (0.01)	0.65*** (0.01)	-0.48*** (0.01)	-0.39*** (0.01)	-0.40*** (0.01)	-0.57*** (0.01)
Model III	-0.30*** (0.01)	-0.90*** (0.01)	-0.90*** (0.01)	-0.44*** (0.08)	0.24*** (0.01)	0.62*** (0.01)	0.73*** (0.01)	0.21*** (0.07)
Start	-0.001*** (0.0000)	-0.001*** (0.0000)	-0.002*** (0.0000)	-0.0003*** (0.0000)	0.001*** (0.0000)	0.001*** (0.0000)	0.001*** (0.0000)	0.001*** (0.0000)
Times Caught	-0.001 (0.001)	0.05*** (0.002)	0.07*** (0.001)	-0.04*** (0.001)	0.005*** (0.001)	-0.35*** (0.003)	-0.07*** (0.001)	0.04*** (0.001)
Repeat Offender	0.02*** (0.01)	1.06*** (0.01)	0.69*** (0.004)	-0.30*** (0.01)	-0.02*** (0.004)	-0.21*** (0.01)	-0.32*** (0.003)	0.24*** (0.01)
IMR				-0.60*** (0.01)				0.48*** (0.01)
Constant	5.95*** (0.06)	7.57*** (0.07)	7.83*** (0.02)	6.75*** (0.05)				
Observations	475'456	465'402	940'858	475'456	475'456	465'402	940'858	475'456
R <sup>2</sup>	0.15	0.45	0.46	0.15				
Adjusted R <sup>2</sup>	0.15	0.45	0.46					
Log Likelihood					-2'804'453.00	-1'432'521.00	-4'409'298.00	-2'803'121.00

Note: This table shows the estimation results of Lasso CV linear cross-sectional regressions to explain cartel duration (ln(duration+1)) and the estimation results of a Lasso CV Weibull Hazard Model to explain cartel death, both at the industry level, for data simulated for Models I, II, IIIa, and IIIb combined. Columns 2 - 5 estimate linear regression coefficients, while columns 6 - 9 estimate HR coefficients, both on the sample of detected cartels, the group of undetected cartels, the population of all cartels, and the sample corrected for Heckman Sample Selection, respectively. The estimated coefficients show standard errors in the sample, but do not test for the real population. The estimated HR coefficients show the change of risk for cartel breakdown if the covariate increases by 1 unit, keeping all others fixed. Standard errors are in parentheses. Significance at the 1%, 5%, and 10% level is indicated by \*\*\*, \*\*, and \*, respectively.



## References

- Bos I, Davies SW, Harrington JE, et al (2018) Does Enforcement Deter Cartels? A Tale of Two Tails. *International Journal of Industrial Organization* 59:372–405
- Harrington JE, Wei Y (2017) What Can the Duration of Discovered Cartels Tell Us About the Duration of All Cartels? *The Economic Journal* 127(604):1977–2005
- Heckman JJ (1979) Sample Selection Bias as a Specification Error. *Econometrica: Journal of the econometric society* pp 153–161
- Stigler GJ (1964) A Theory of Oligopoly. *Journal of Political Economy* 72(1):44–61
- Tibshirani R (1996) Regression Shrinkage and Selection via the Lasso. *Journal of the Royal Statistical Society Series B (Methodological)* 58(1):267–288