

Data Sharing Under Confidentiality: CBRT Case

Timur Hülagü, Ph. D.

AUGUST 30, 2018 | BASEL



This is a joint project between CBRT and METU. All views expressed here are those of authors and do not necessarily reflect those of the two institutions.





► Main Goal: Address the growing need of accessing micro data for academic research

► G-20 Data Gaps Initiative 2, Recommendation II.20: Promotion of Data Sharing by G-20 Economies

Share information and ideas on ways to apply confidential rules/arrangements in a manner that would allow sharing of more granular data

Eurostat Peer review report on the compliance with the Code of Practice and the coordination role of the National Statistical Institute in Turkey

Recommendation 22: TurkStat should introduce remote access facilities for researchers, who are permitted to use its anonymized microdata for research purposes (European Statistics Code of Practice, indicator 15.4)



Data Sharing Trade-off









Main Aspects



Accuracy

- ► Descriptive Analysis
- ► Univariate Regression Analysis
- ► Multivariate Regression Analysis 🛛 🗹
- ► Logistic Regression 🗹

 $\mathbf{\nabla}$

Logarithmic Regression



Accuracy

Descriptive analysis

Variable	B30	B50	B15	G600
Seed	123000	234000	345000	456000
Mean (O/M)	0.9756	0.9756	0.9756	0.9756
Standart Deviation (O/M)	0.9997	0.9995	0.9999	0.9998
Skewness (O/M)	1.001	1.003	1.002	1.002
Kurtosis (O/M)	1.002	1.008	1.004	1.004

Multiple Linear Regression Analysis

Original			Masked									
Coefficients (Intercept) B2L B10L B32L Signif. code Residual sta Multiple R-S F-statistic:	: Estimate ST 8.08888 0.04716 0.21006 0.32332 s: 0 '*** ndard error quared: 0. 2.56e+03 (td. Error t 0.11373 0.00616 0.00486 0.00583 ' 0.001 '** r: 0.856 on 511, A on 3 and 73	value 71.13 < (7.66 43.22 < (55.44 < (' 0.01 '* 7342 degr djusted R 42 DF, p	Pr(> t) 0.00000000000000022 0.0000000000000000	*** *** *** 1 000000002	Coefficients (Intercept) B2L_M B10L_M B32L_M Signif. code Residual sta Multiple R-S F-statistics	Estimate S 8.28657 0.04777 0.20994 0.32314 es: 0 '*** endard erro equared: 0 2.55e+03	td. Error 1 0.11670 0.00617 0.00486 0.00584 ' 0.001 '*' r: 0.857 or .51, / on 3 and 7	value 71.01 7.75 43.16 55.36 ' 0.01 n 7342 de Adjusted 342 DF,	< 0.000000 0.00000 0.000000 0.000000 **' 0.05 ' egrees of R-squared p-value:	<pre>Pr(> t) 000000002 0000000011 0000000002 00000000</pre>	*** *** *** 1 00000002

annt







Privacy

Spectral Filtering

	Sinusoidal	Triangular
λ_{\min}	0,01121	0,007627
λ _{max}	0,06104	0,04152
m x n	250x40	250x40
k	2	1
d (O , M)	1410	1163
d(0, E)	319,6	205, 2
m	0,227	0, 176

Sinusoidal



Triangular



	G69
λ_{min}	45870034800
λ _{max}	194658444176
m x n	250x30
k	30
d(0, M)	1964402616
d(0, E)	1964402616
m	1



TÜRKİYE CUMHURİYET MERKEZ BANKASI

Privacy

Random Number Generation In Our System

We can use 2³⁰ different seed for generating random numbers. So that, brute force attacks may be a threat. To solve this problem we offer the following algorithm.

Original Data	Seed	Noise	Masked Data
<i>x</i> ₁	$s_1 = f(IV + x_1)mod \ 2^{30}$	$\varepsilon_1 = \text{RNG}(s_1)$	$x'_1 = x_1 + \varepsilon_1$
<i>x</i> ₂	$s_2 = f(x_2 + x'_1)mod \ 2^{30}$	$\varepsilon_2 = \text{RNG}(s_2)$	$x'_2 = x_2 + \varepsilon_2$
:	:	:	:
x _n	$s_n = f(x_n + x'_{n-1}) \mod 2^{30}$	$\varepsilon_n = \operatorname{RNG}(s_n)$	$x'_n = x_n + \varepsilon_n$

We choose nonlinear f(x), such that : $f(x) = \mu x^3 + \sigma$





What is done

- ▶ We achieve, measurable accuracy on masked data.
- ▶ We observed that our system is secure for the attacks we mentioned.

Future Works

- ► We will study two new attacks called map estimation and distribution analysis.
- ▶ In the masked data set, we will check the accuracy of other statistical functions.
- ► Finally, we will produce a user friendly software product developed by using Java and R.

