# PISA Data Analysis leveraging R:

## pros and cons

Mariia Mazorchuk, Anna Bychko

Ukrainian Center for Educational Quality Assessment, Kyiv, Ukraine

*Dortmund, 27-28 February, 2020*

# What is PISA?

**PISA** is the OECD's Programme for International Student Assessment

**Goals:**

- PISA measures 15-year-olds' ability to use their reading, mathematics and science knowledge and skills to meet real-life challenges.

- PISA makes it possible to compare Educational systems of different countries.

- PISA gives us opportunity to define factors, which influences students achievements.

## 2000-2018:

PISA has involved more than 90 countries and economies and about 3 000 000 students worldwide
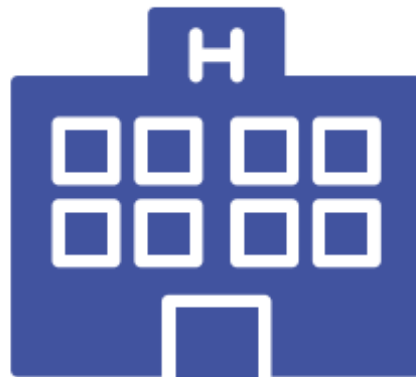
## PISA-2018:

**600 000** students representing about **32** million 15-year-olds in schools of **79** participating countries and economies sat the 2-hour PISA test in 2018
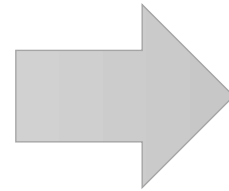
# **Ukraine in PISA-2018**

- Ukraine participated in PISA for the first time
- Over 6 000 students from 250 Ukrainian schools representing about 315 000 Ukrainian students sat the 2-hour PISA test and filled in questionnaires
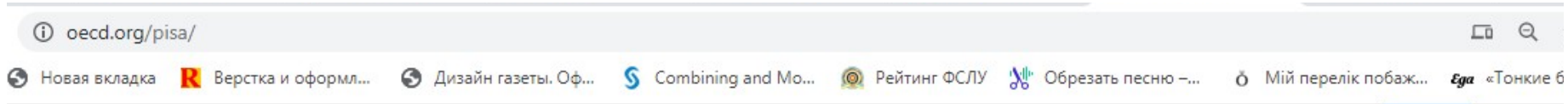
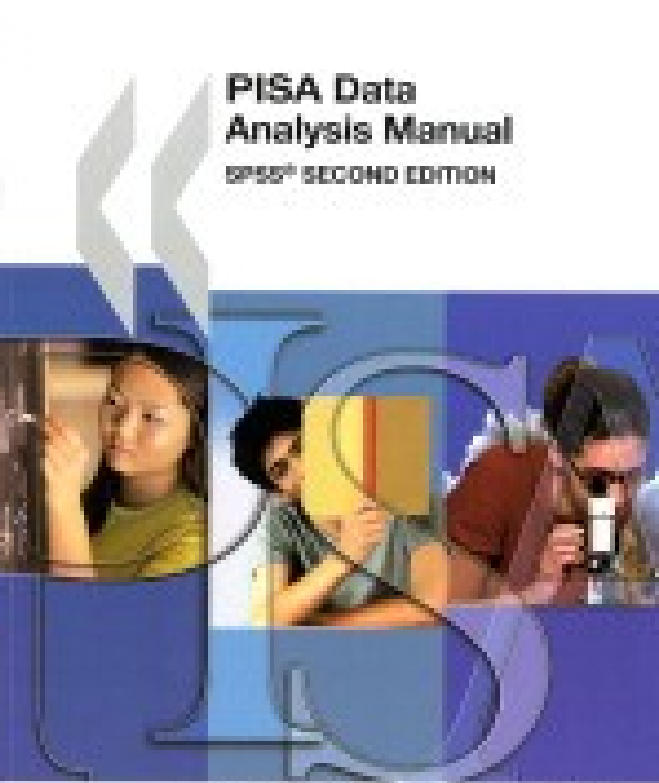# Our goal: Ukraine PISA National report



Data Processing

Національний звіт за результатами міжнародного дослідження якості освіти PISA-2018

https://www.oecd.org/pisa/

# Content of technical report

- Usefulness of PISA Data for Policy Makers, Researchers and Experts on Methodology

- Exploratory Analysis Procedures

- Sample Weights

- **Replicate Weights**

- **Computation of Standard Errors** (take into account the complex sample design)

- **Plausible Values**

- **Analyses with Plausible Values** (take into account rotated test forms )

- Use of Proficiency Levels

- The Rasch Model

# Content of technical report

- Analyses with School-Level Variables

- Standard Error on a Difference

- **OECD Total and OECD Average**

- Trends

- **Studying the Relationship between Student Performance and Indices Derived from Contextual Questionnaires**

- Multilevel Analyses

- PISA and Policy Relevance – Three Examples of Analyses

- SPSS® Macro; SAS® Macro

- SAS Macro for 10 Plausible Values

# Statistical software

- SPSS
- STATA
- SAS
- R

| | SAS | SPSS | R |
|---|---|---|---|
| **Advantages** | 1. High adoption rate in major industries<br>2. Flow based interface with drag and drop<br>3. Official support<br>4. Handling large datasets<br>5. 'PROC SQL' | 1. Used a lot in universities<br>2. Good user interface with extensive documentation<br>3. Click & Play functionality<br>4. Writing code made easy using the 'paste' button.<br>5. Official support | 1. Big community who creates libraries<br>2. Free<br>3. Early adopter in explanatory and predictive modeling.<br>4. Easy to connect to data sources, including NoSQL and webscraping. |
| **Disadvantages** | 1. Relatively high cost<br>2. For not-standard options not in interface, you'll need to write the code<br>3. Slow adapting to new techniques<br>4. Different programs for visualization or Data Mining | 1. Relatively high cost<br>2. different licenses for different functionalities.<br>3. Syntax limited<br>4. Slow adapting to new techniques<br>5. Slow in handling large datasets | 1. Can be slow with big datasets<br>2. Steep learning curve<br>3. No official support<br>4. No user interface |

- IDB Analyzer (https://www.iea.nl/data-tools/tools#section-308)
- PISA Data Explorer (https://pisadataexplorer.oecd.org/ide/idepisa/)

# intsvy package

R 'intsvy'

**intsvy: International Assessment Data Manager**

intsvy is an R package for working with international assessment data from PISA

**Consulting services**

Daniel Caro, Data Scientist

PISA
TIMSS
PIRLS
PIAAC
ICILS

https://cran.r-project.org/web/packages/intsvy/intsvy.pdf

http://danielcaro.net/r-intsvy/

# Preparation of the PISA data files

- **Importing my data**
- **Merge the PISA data files**
- **Recode variables**
- **Creating new variables**

library(intsvy)

pisa.var.label

pisa.select.merge

library(survey)

svydesign

svyquantile

# Calculating main estimates

- **Average students performance**

pisa2015.mean.pv(pvlabel = "READ", data = pisa)

| | Freq | Mean | s.e. | SD | s.e |
|---|---|---|---|---|---|
| 1 | 5998 | 465.95 | 3.5 | 93.34 | 1.7 |

- **Frequency tables**

pisa2015.table(variable="TFGender", data = pisa)

| | TFGender | Freq | Percentage | Std.err. |
|---|---|---|---|---|
| 1 | Female | 2857 | 47.37 | 1.02 |
| 2 | Male | 3141 | 52.63 | 1.02 |

library(intsvy)

pisa2015.mean.pv

pisa2015.mean

pisa2015.table

# Calculating main estimates

- **Proficiency levels**

pisa2015.ben.pv(pvlabel="READ", cutoff = c(189.33, 262.04, 334.75, 407.47, 480.18, 552.89, 625.61, 698.32), data=pisa)

|   | Benchmarks | Percentage | Std. err. |
|---|------------|------------|-----------|
| 1 | <= 189.33 | 0.17 | 0.08 |
| 2 | (189.33, 262.04] | 1.80 | 0.29 |
| 3 | (262.04, 334.75] | 7.21 | 0.69 |
| 4 | (334.75, 407.47] | 16.73 | 0.87 |
| 5 | (407.47, 480.18] | 27.73 | 0.81 |
| 6 | (480.18, 552.89] | 28.48 | 0.97 |
| 7 | (552.89, 625.61] | 14.47 | 0.82 |
| 8 | (625.61, 698.32] | 3.24 | 0.44 |
| 9 | > 698.32 | 0.17 | 0.11 |

library(intsvy)

pisa2015.ben.pv

# Regression models

- **Linear regression analysis**

pisa2015.reg.pv(pvlabel = "READ", x="TFGender", data = pisa)

|  | Estimate | Std. Error | t value |
|---|---|---|---|
| (Intercept) | 483.56 | 3.63 | 133.30 |
| TFGenderMale | -33.46 | 3.86 | -8.67 |
| R-squared | 0.03 | 0.01 | 4.54 |

library(intsvy)

pisa2015.reg.pv

pisa2015.reg

# Regression models

- **Logistic regression analysis**

fit1<-pisa2015.log.pv(pvlabel = "MATH", x="TFGender",cutoff=420, data=pisa)

|  | Coef. | Std. Error | t value | OR | CI95low | CI95up |
|---|---|---|---|---|---|---|
| (Intercept) | 1.47 | 0.09 | 16.33 | 4.34 | 3.64 | 5.18 |
| TFGenderMale | -0.72 | 0.09 | -7.94 | 0.49 | 0.41 | 0.58 |

library(intsvy)

pisa2015.log.pv

pisa2015.log

odds_female=exp(fit1$Coef.[1])

[1] 4.349235

odds_male=exp(fit1$Coef.[2])*exp(fit1$Coef.[1])

[1] 2.117

**Plots**

plot

# **Plots**

**PISA Plots**

library(ggplot2)

ggplot

# Additional PISA data analysis

- **Multilevel regression models**

fm2<-lmer(stud484$PV1READ~1+stud484$ESCS +(1+stud484$Mean|

stud484$CNTSCHID),weights=stud484$norm_weight, stud484)

summary(fm2)

```
        Random effects:
Groups                  Name            Variance Std.Dev. Corr
 stud484$CNTSCHID (Intercept)      631.9    25.14
                 stud484$Mean      8047.4    89.71    0.44
                    Residual       5799.9   76.16
Number of obs: 5998, groups:  stud484$CNTSCHID, 250
        Fixed effects:
                  Estimate Std. Error t value
(Intercept)        473.445    2.423   195.43
stud484$ESCS        24.111    1.481    16.28
Correlation of Fixed Effects:
        (Intr) std484$ESCS 0.110
```

library(lme4)

lmer

library(survey)

svydesign

**Only for one PV !**

# Additional PISA data analysis



*Null model. Only random intercept PV1READ ~ ESCS*

*Fixed effect and random effects: Reading ~ ESCS and Mean of ESCS by Schools*

library(ggplot2)

ggplot

# Pros and cons

- R is a free open source package.
- You can create a flexible script and repeat calculation process.
- You can use effective functions for calculating estimates taking into account the complex sample design and rotated test form of PISA data (using **intsvy** package).

- A major drawback of R is that most of its functions have to load all the data into memory before execution, which sets a **limit** to the volumes that can be handled.
- R requires some programming skills.
- The process of graph creating sometimes is very complex.

# Results and discussion

We have published a Ukrainian National report, where the results of our work have been shown.
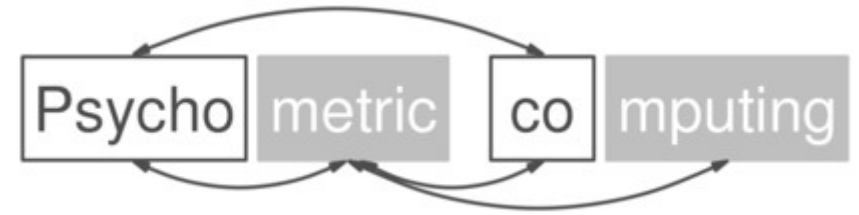You can look them up in the reference:
http://pisa.testportal.gov.ua/wp-content/uploads/2019/12/PISA_2018_Report_UKR.pdf

We don't have a unified system for all calculations and for forming reports.
There aren't a lot of functions for calculating different indicators. We want to extend intsvy package.

Національний звіт за результатами міжнародного дослідження якості освіти PISA-2018

# Thank you!

mazorchuk.mary@gmail.com     bychko.anya@gmail.com