



CENTRAL & EASTERN EUROPEAN  
SUSTAINABLE ENERGY NETWORK  
CEESEN-BENDER

# Statistical Tool to Prioritise Buildings for Renovation: User Manual



Co-funded by  
the European Union

CEESEN-BENDER project has developed a statistical tool that uses digital socioeconomic data, technical characteristics of buildings and energy consumption data to generate rankings of buildings that are the least energy-efficient and having occupants most likely to suffer from high levels of energy poverty.

For the analysis, R gramme should be installed.



CENTRAL & EASTERN EUROPEAN  
SUSTAINABLE ENERGY NETWORK  
·  
CEESEN-BENDER

# R programme installation



Co-funded by  
the European Union

- For statistical modelling of digital tool, which ranks buildings for renovation, R programme has been applied. This software has a number of advantages:
- it is open-source software, meaning it is freely available for anyone to use.
  - R runs on various operating systems, including Windows, macOS, and Linux, making it accessible to a broad range of users.
  - it has a large and active community of users and developers, providing extensive documentation, forums, and resources for troubleshooting and learning.
  - it can easily integrate with other programming languages (like Python, C++, and Java) and tools (like SQL databases and web applications).
  - it supports reproducible research through R Markdown, allowing users to create dynamic reports that combine code, results, and narrative text in a single document.
  - it is widely used in academia, and there are numerous resources, courses, and textbooks available for learning statistical analysis and R programming.
  - R provides a vast array of statistical techniques, including linear and nonlinear modelling, time-series analysis, classification, clustering, and more. It is continuously updated with new statistical methods.

The latest R can be downloaded from R homepage <http://www.r-project.org/>:

1. Choose “Download” from the left menu
2. Then choose a mirror server from where to download R (preferably a server closest to you).
3. Specify your operating system (Windows / Mac / Linux) and choose base (base program).



[Home]

**Download**

CRAN

**R Project**

About R

Logo

Contributors

What's New?

Reporting Bugs

Conferences

Search

Get Involved: Mailing Lists

Get Involved: Contributing

Developer Pages

R Blog

**R Foundation**

Foundation

Board

Members

Donors

Donate

## The R Project for Statistics

### Getting Started

R is a free software environment for statistical computing and graphics. It runs on a wide variety of UNIX platforms, Windows and MacOS. To [download R](#), see the [CRAN](#) mirror.

If you have questions about R like how to download and install, please read our [answers to frequently asked questions](#).

### News

- **R version 4.5.0 (How About a Twenty-Six)** has been released.
- **R version 4.4.3 (Trophy Case)** (wrap-up of 4.4.x) was released.
- The **useR! 2025** conference will take place at Duke University.
- We are deeply sorry to announce that our friend and colleague [Fritz](#) has passed away. [Read our tribute to Fritz here.](#)
- You can support the R Foundation with a renewable source of energy.

### News via Mastodon



**R\_Foundation**

New [#RStats](#) blog entry by Tomas Kalibera: Ser

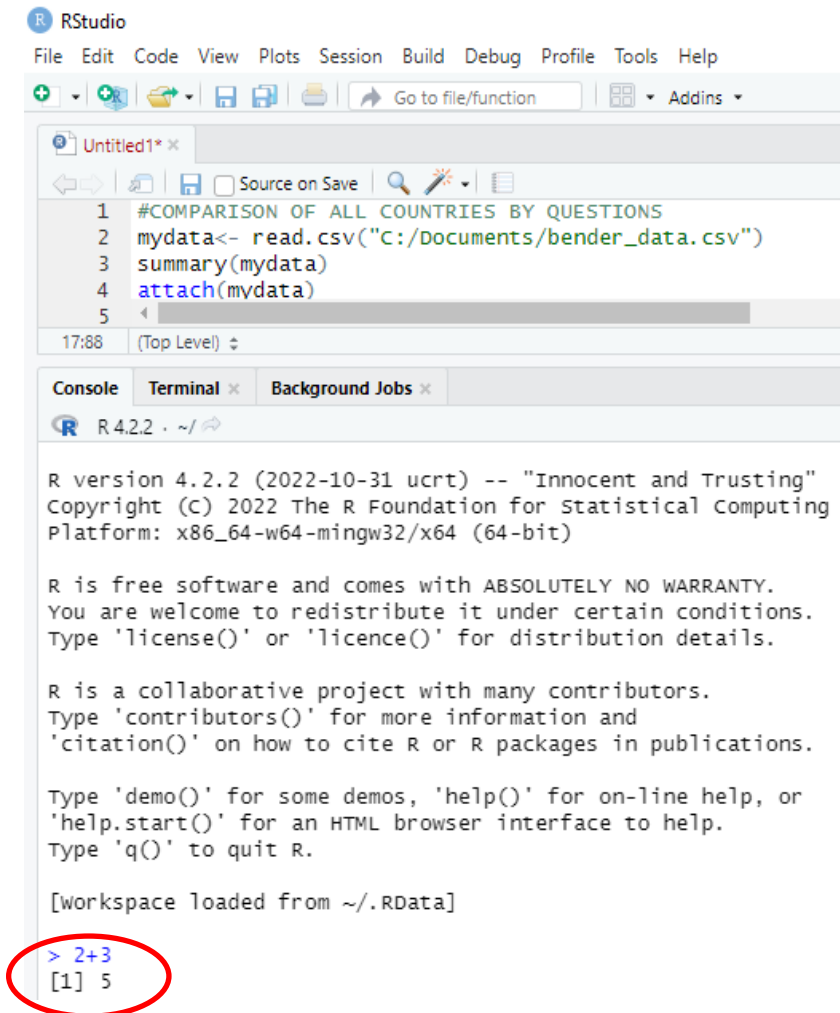
[blog.r-project.org/2025/04/24/...](http://blog.r-project.org/2025/04/24/...)

Apr 28, 2025

4. After installing R you can install RStudio which is an IDE (integrated development environment), that makes the R experience much more easier and efficient. For a free version of RStudio go to the following website:

<https://www.rstudio.com/products/rstudio/download/>

5. Choose the free version for your operating system. Check if all of the installations worked by running RStudio and test if you will get the right answer for 2+3. Combination of Ctrl+Enter executes the command.



The screenshot shows the RStudio interface. The script editor contains the following code:

```
1 #COMPARISON OF ALL COUNTRIES BY QUESTIONS
2 mydata<- read.csv("C:/Documents/bender_data.csv")
3 summary(mydata)
4 attach(mydata)
5
```

The console shows the R version and license information:

```
R version 4.2.2 (2022-10-31 ucrt) -- "Innocent and Trusting"
Copyright (C) 2022 The R Foundation for Statistical Computing
Platform: x86_64-w64-mingw32/x64 (64-bit)

R is free software and comes with ABSOLUTELY NO WARRANTY.
You are welcome to redistribute it under certain conditions.
Type 'license()' or 'licence()' for distribution details.

R is a collaborative project with many contributors.
Type 'contributors()' for more information and
'citation()' on how to cite R or R packages in publications.

Type 'demo()' for some demos, 'help()' for on-line help, or
'help.start()' for an HTML browser interface to help.
Type 'q()' to quit R.

[workspace loaded from ~/.RData]
```

The console also shows the result of the command `> 2+3`, which is `[1] 5`. This result is circled in red in the image.

This is the **script window**, the commands that make the analysis are written here. The script can be saved in order to repeat the analysis later.

This is the **environment window**. Here we can see datasets, fitted models and other defined (saved) variables/objects that we can use.

The screenshot displays the RStudio interface with three main windows:

- Script Window:** Contains R code for data analysis, including loading a dataset, creating a factor, and generating a barplot. The code is as follows:

```
col=c("skyblue4", "darkslategray1", "lightskyblue", "deepskyblue", "lightskyblue3", "lightblue1"), ylab="%",
x1m=c(0,8), xlab="Countries", names=c("Greece", "Cyprus", "Estonia", "Spain"))
legend("topright", fill=c("skyblue4", "darkslategray1", "lightskyblue", "deepskyblue", "lightskyblue3", "lightblue1"), c("Str
prop.table(table(educ, country),2)*100
#I believe that I can recognize signs of physical domestic violence
mydata$physigns-factor(country, levels=1:4, labels=c("Greece", "Cyprus", "Estonia", "Spain"))
windows(width=6, height=4)
barplot(prop.table(table(physigns, country),2)*100,
col=c("skyblue4", "darkslategray1", "lightskyblue", "deepskyblue", "lightskyblue3", "lightblue1"), ylab="%",
x1m=c(0,8), xlab="Countries", names=c("Greece", "Cyprus", "Estonia", "Spain"),
main="I believe that I can recognize signs of physical domestic violence")
legend("topright", fill=c("skyblue4", "darkslategray1", "lightskyblue", "deepskyblue", "lightskyblue3", "lightblue1"), c("Str
prop.table(table(physigns, country),2)*100
```
- Console Window:** Shows the execution of the script, including error messages for failed file operations and the output of the barplot.

```
> 2+3
[1] 5
> #COMPARISON OF ALL COUNTRIES BY QUESTIONS
> mydata<- read.csv("C:/Users/marifa91/Documents/Projects/opep_data.csv")
Error in file(file, "rt") : cannot open the connection
In addition: warning message:
In file(file, "rt") :
cannot open file 'C:/Users/marifa91/Documents/Projects/opep_data.csv': No such file or directory
> #COMPARISON OF ALL COUNTRIES BY QUESTIONS
> mydata<- read.csv("C:/Users/marifa91/Documents/Projects/opep_data.csv")
Error in file(file, "rt") : cannot open the connection
In addition: warning message:
In file(file, "rt") :
cannot open file 'C:/Users/marifa91/Documents/Projects/opep_data.csv': No such file or directory
> #COMPARISON OF ALL COUNTRIES BY QUESTIONS
> mydata<- read.csv("C:/Users/marifa91/Documents/Projects/opep_data.csv")
> #COMPARISON OF ALL COUNTRIES BY QUESTIONS
> mydata<- read.csv("C:/Users/marifa91/Documents/Projects/opep_data.csv")
> attach(mydata)
> prop.table(table(role, country),2)*100
country
role      1      2      3      4
1 0.000000 5.882353 0.000000 2.000000
2 0.000000 5.882353 1.666667 0.000000
3 0.000000 5.882353 16.666667 0.000000
4 15.000000 6.47059 31.666667 6.000000
5 35.000000 35.294118 33.333333 38.000000
6 50.000000 29.411765 16.666667 54.000000
> barplot(prop.table(table(physigns, country),2)*100,
+ col=c("skyblue4", "darkslategray1", "lightskyblue", "deepskyblue", "lightskyblue3", "lightblue1"), ylab="%",
+ x1m=c(0,8), xlab="Countries", names=c("Greece", "Cyprus", "Estonia", "Spain"))
```
- Environment Window:** Shows the loaded dataset 'mydata' with 167 observations and 36 variables.
- Plot Window:** Displays a barplot titled "I believe that I can recognize signs of physical domestic violence". The x-axis is labeled "Countries" and includes Greece, Cyprus, Estonia, and Spain. The y-axis is labeled "%" and ranges from 0 to 100. The plot shows the distribution of responses (role) for each country.

**Console window.** This is where results, error messages and warnings are shown. You can copy commands straight into this window (if you wish to run a quick calculation not related to the current analysis).

This is where the **plot** will appear. Also, there are tabs on top, one of which is the help tab, containing helpful information on all the functions.

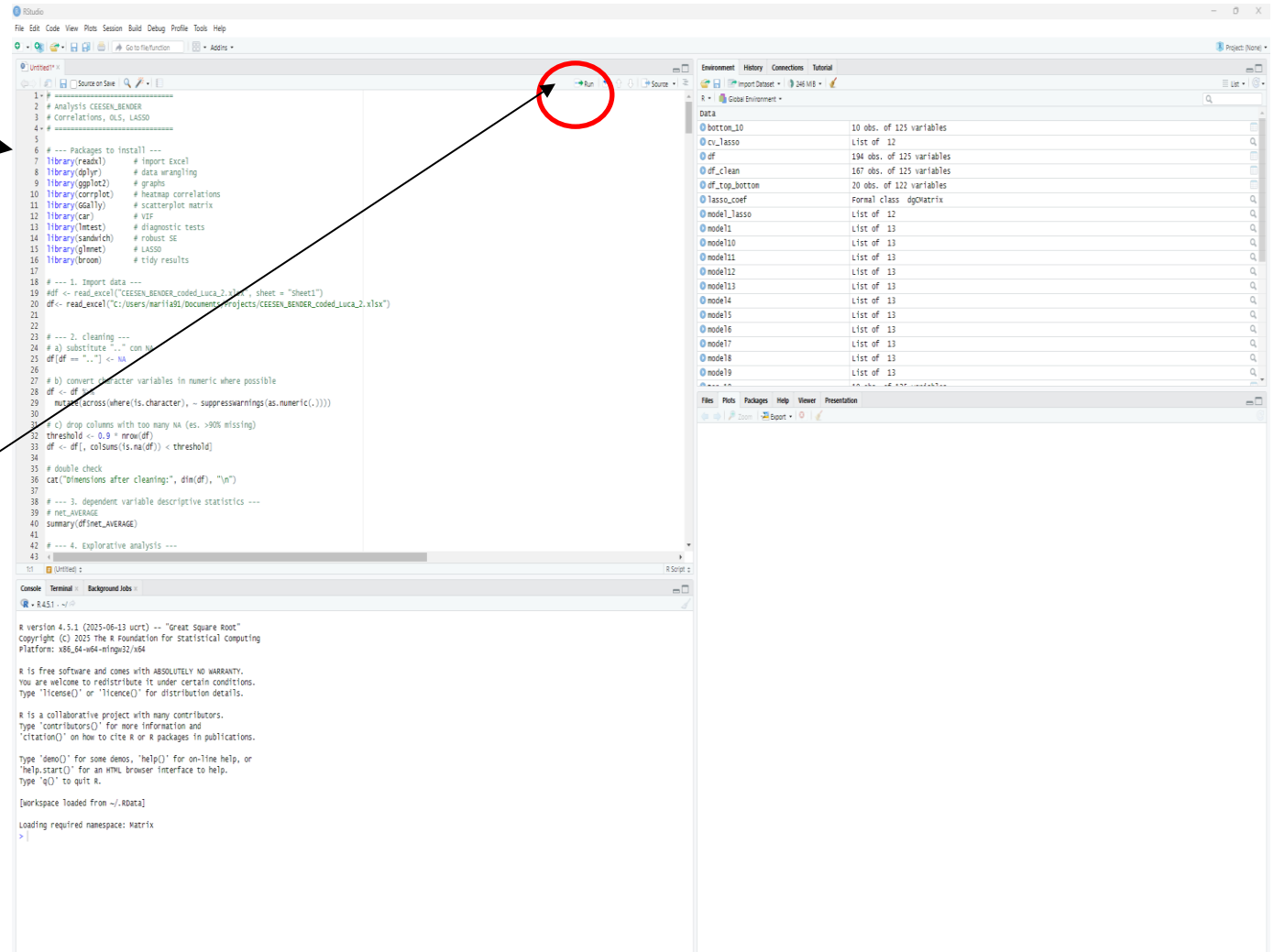
# Installation of additional R packages for the analysis (1)



1. Copy these codes and paste into the script window:

```
library(readxl)      # import Excel  
library(dplyr)       # data wrangling  
library(ggplot2)     # graphs  
library(corrplot)    # heatmap correlations  
library(GGally)      # scatterplot matrix  
library(car)         # VIF  
library(lmtest)      # diagnostic tests  
library(sandwich)    # robust SE  
library(glmnet)      # LASSO  
library(broom)       # tidy results
```

2. Select them all at once and run by either clicking “Run” command in the right hand part of the script window or as a combination of “Ctrl+Enter”



# Installation of additional R packages for the analysis (2)



If these packages have not been previously installed, they should be installed one by one manually by clicking on

## Tools – Install Packages

The screenshot displays the RStudio interface. The 'Tools' menu is open, and the 'Install Packages...' option is highlighted. The console window shows the following R code and output:

```
R version 4.5.1 (2025-06-13 ucrt) -- "Great Square Root"
Copyright (C) 2025 The R Foundation for Statistical Computing
Platform: x86_64-w64-mingw32/x64

R is free software and comes with ABSOLUTELY NO WARRANTY.
You are welcome to redistribute it under certain conditions.
Type 'license()' or 'licence()' for distribution details.

R is a collaborative project with many contributors.
Type 'contributors()' for more information and
'citation()' on how to cite R or R packages in publications.

Type 'demo()' for some demos, 'help()' for on-line help, or
'help.start()' for an HTML browser interface to help.
Type 'q()' to quit R.

[workspace loaded from ~/.Rdata]

Loading required namespace: Matrix
Failed with error: 'reached elapsed time limit'
```

The script being executed in the editor includes the following code:

```
1 # -----
2 # Analysis CEESEN_BENDER - all count
3 # correlations, OLS, LASSO
4 # -----
5
6 # --- Packages to install ---
7 library(readxl) # import Excel
8 library(dplyr) # data wranglin
9 library(ggplot2) # graphs
10 library(corrplot) # heatmap corr
11 library(ggally) # scatterplot
12 library(car) # VIF
13 library(lmtest) # diagnostic te
14 library(sandwich) # robust SE
15 library(glmnet) # LASSO
16 library(broom) # tidy results
17
18 # --- 1. Import data ---
19 df <- read_excel("C:/Users/mariia91/D
20
21
22 # --- 2. cleaning ---
23 # a) substitute " " with NA
24 df[df == " "] <- NA
25
26 # b) convert character variables in numeric where possible
27 df <- df %>%
28 mutate(across(where(is.character), ~ suppresswarnings(as.numeric(.))))
29
30 # c) drop columns with too many NA (es. >90% missing)
31 threshold <- 0.9 * nrow(df)
32 df <- df[, colSums(is.na(df)) < threshold]
33
34 # double check
35 cat("Dimensions after cleaning:", dim(df), "\n")
36
37 # --- 3. dependent variable descriptive statistics ---
38 # avheat
39 summary(df$avheat)
40
41 # --- 4. Explorative analysis ---
42 # avheat distribution
43
```

The Environment pane on the right shows the following data objects:

Object	Details
bottom_10	10 obs. of 125 variables
cv_lasso	List of 12
df	194 obs. of 125 variables
df_clean	167 obs. of 125 variables
df_top_bottom	20 obs. of 122 variables
lasso_coef	Formal class dgcmatrix
model_lasso	List of 12
model1	List of 13
model10	List of 13
model11	List of 13
model12	List of 13
model13	List of 13
model14	List of 13
model15	List of 13
model16	List of 13
model17	List of 13
model18	List of 13
model19	List of 13



# Loading dataset into R environment (1)



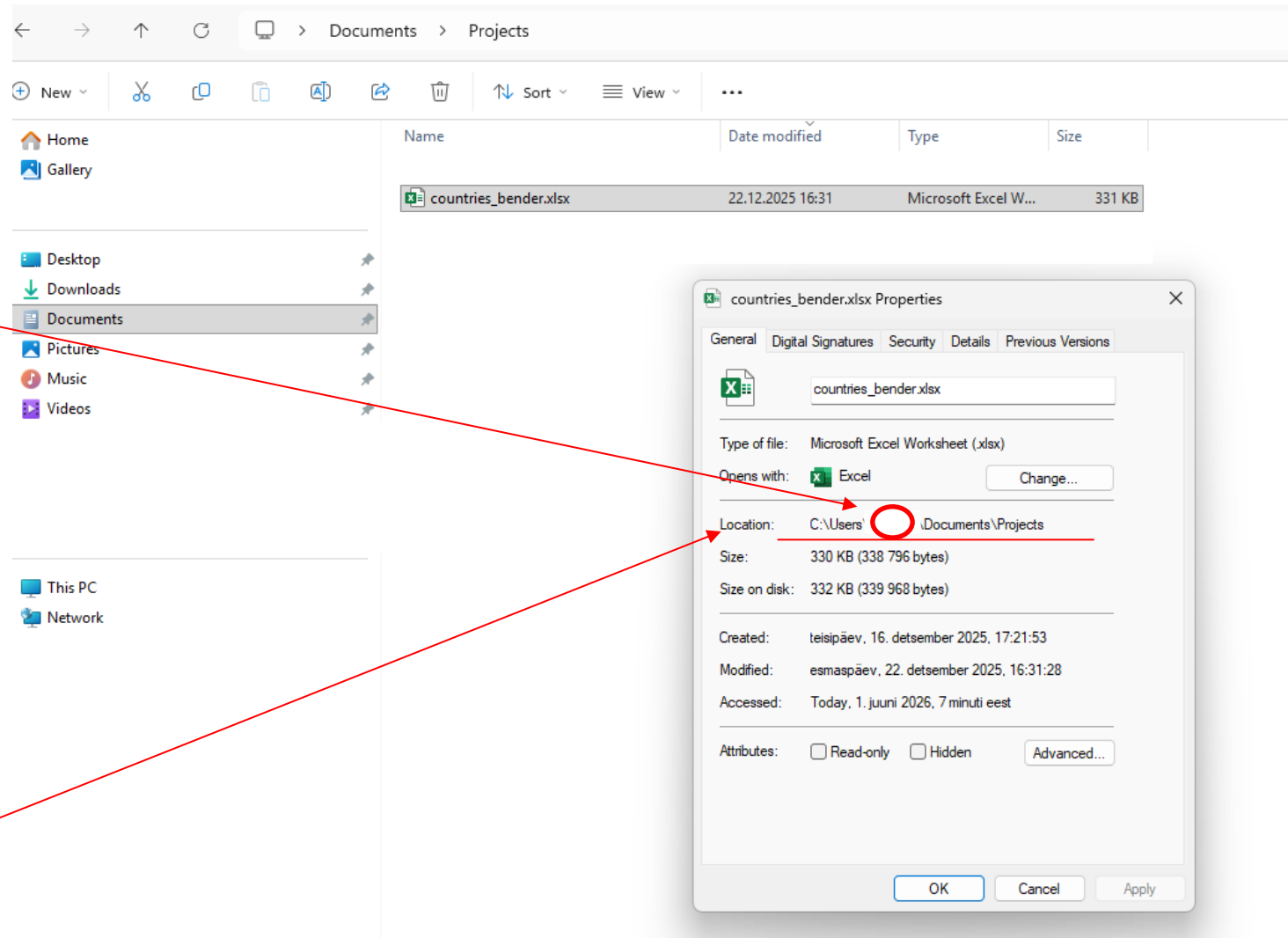
1. Copy all the codes from the file called "**R codes.txt**" into script window in R programme.

2. Execute the command (Ctrl+Enter):

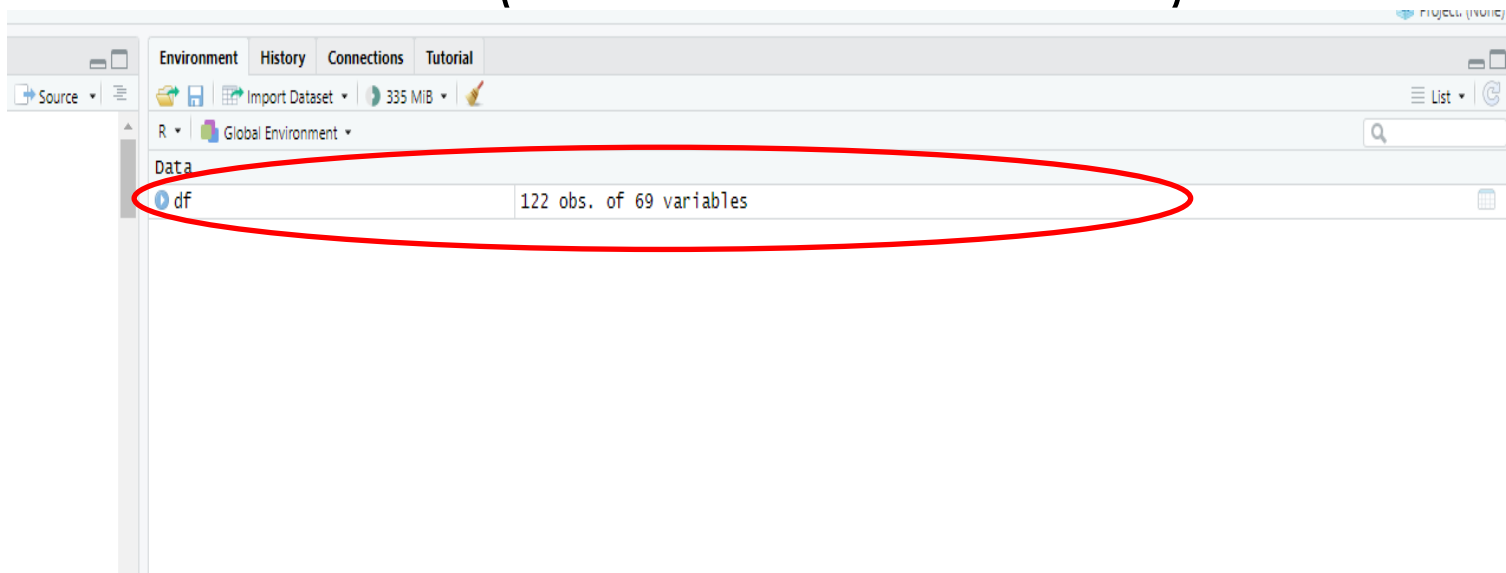
```
df<- read_excel("C:/Users/YOUR NAME/Documents/Projects/count  
ries_bender.xlsx")
```

**NB!** The path to your file in "" will vary, you have to check it under "**countries\_bender.xlsx**" file's "Properties"

For that: Go to Documents → folder "Projects" → right click on the file "**countries\_bender.xlsx**" → Properties → under "General", you will see the correct location of your file.



3. If there were no errors running these commands (due to missing parenthesis for example), then a dataset named “df” should appear in the top right corner window (environment window):



If you click on the dataset row, the dataset or its information is shown. For more details on how to use R programme, one can use a user-friendly online guidebook available [here](#) or manuals allocated on R website [here](#).

4. Run all codes from “**R codes.txt**” file



CENTRAL & EASTERN EUROPEAN  
SUSTAINABLE ENERGY NETWORK  
CEESEN-BENDER

# Description of the model and results



Co-funded by  
the European Union

# Dataset overview



Cross-sectional annual data has been collected by BENDER project partners in four countries – Croatia, Romania, Slovenia and Poland – to pilot the tool developed initially on Estonian data. Three groups of data are used:

Socio-economic data of MAB* residents	Technical characteristics of buildings	Energy data
<ul style="list-style-type: none"> <li>• Number of (mostly) empty flats with 0 residents</li> <li>• Number of flats with 1 resident</li> <li>• Number of flats with 2 residents</li> <li>• Number of flats with 3 residents or more</li> <li>• Average living area per flat (m<sup>2</sup>)</li> <li>• Number of residents living in the MAB</li> <li>• Average living area per person (m<sup>2</sup>)</li> <li>• Percentage of owners living in the MAB (%)</li> <li>• Percentage of tenants living in the MAB (%)</li> <li>• Children below 18 of age (%)</li> <li>• Working adults 18-60 of age (%)</li> <li>• Unemployed adults 18-60 of age (%)</li> <li>• Pensioners above 60 of age (%)</li> <li>• Average salary per dwelling (€ gross)</li> <li>• Dwellings with residents receiving social assistance (%)</li> </ul>	<ul style="list-style-type: none"> <li>• GFA - Gross Floor Area (m<sup>2</sup>)</li> <li>• Conditioned area (m<sup>2</sup>)</li> <li>• Conditioned area % from GFA</li> <li>• Year of construction</li> <li>• Year of last renovation (if available)</li> </ul>	<ul style="list-style-type: none"> <li>• MAB electricity consumption [MWh]</li> <li>• MAB natural gas consumption [MWh]</li> <li>• MAB heating costs [€]</li> </ul>

\*MAB – multi-apartment buildings

The sample consists of 122 observations from four countries.

## **Independent variables:**

- Country
- Household composition (number of dwellings with 0, 1, 2, 3 or more residents)
- Average living area per flat
- Average living area per person
- Ownership structure:
  - Share of owners (%)
  - Share of renters (%)
- Employment structure:
  - Children below 18 of age (%)
  - Working adults 18-60 of age (%)
  - Unemployed adults 18-60 of age (%)
  - Pensioners above 60 of age (%)
- Income (measured as average monthly salary per building (€ gross))
- Flats with residents receiving social assistance (%)
- GFA - Gross Floor Area (m<sup>2</sup>)
- Conditioned area (m<sup>2</sup>)
- Year of construction
- Average electricity consumption of building per year [MWh]
- Average natural gas consumption of building per year [MWh]

## **Dependent variable:**

Average annual heating costs, €

- *Correlation analysis* showed which independent variables correlate with each other → only one should be included
- 16 *Multiple Linear Regression (MLR)* models were tested to identify which factors are consistently significant across all the models
- *LASSO model* was applied to double check the results of OLS. When many predictors are involved LASSO identifies the most important from them.

- VIF (Variance Inflation Factor) is applied to detect a multicollinearity problem in tested models. VIF between 1 and 5 is an acceptable range.
- Residuals
  - Check for linearity and homoscedasticity
  - Check for normality of residuals
- Robust Standard Error (test for heteroscedasticity)

# LASSO model for Croatia, Poland and Romania



Significant factors	LASSO coefficient	Weights*
<b>Average monthly salary per building</b>	0.000002101748	0.05
<b>Country</b>	0.5015637	0.4
<b>Percent of unemployed adults age 18-60</b>	0.02386472	0.15
<b>Number of (mostly) empty dwellings with 0 residents</b>	0.06683785	0.3
<b>Year of construction</b>	0.006892765	0.1

\*Weights are selected based on the magnitude of LASSO coefficients

Calculation of a weighted score:

$$\text{country} * \text{weights}[0.4] + \text{resident0\_std} * \text{weights}[0.3] + \\ \text{unemployed\_std} * \text{weights}[0.15] + \text{year\_std} * \text{weights}[0.1] + \\ \text{bavsalary\_std} * \text{weights}[0.05]$$

# LASSO model for Slovenia



For Slovenia, factor “unemployed” has to be excluded from the analysis, since no observations for this factor are available.

Significant factors	LASSO coefficient	Weights*
<b>Average monthly salary per building</b>	0.000002101748	0.1
<b>Country</b>	0.5015637	0.4
<b>Number of (mostly) empty dwellings with 0 residents</b>	0.06683785	0.3
<b>Year of construction</b>	0.006892765	0.2

\*Weights are selected based on the magnitude of LASSO coefficients

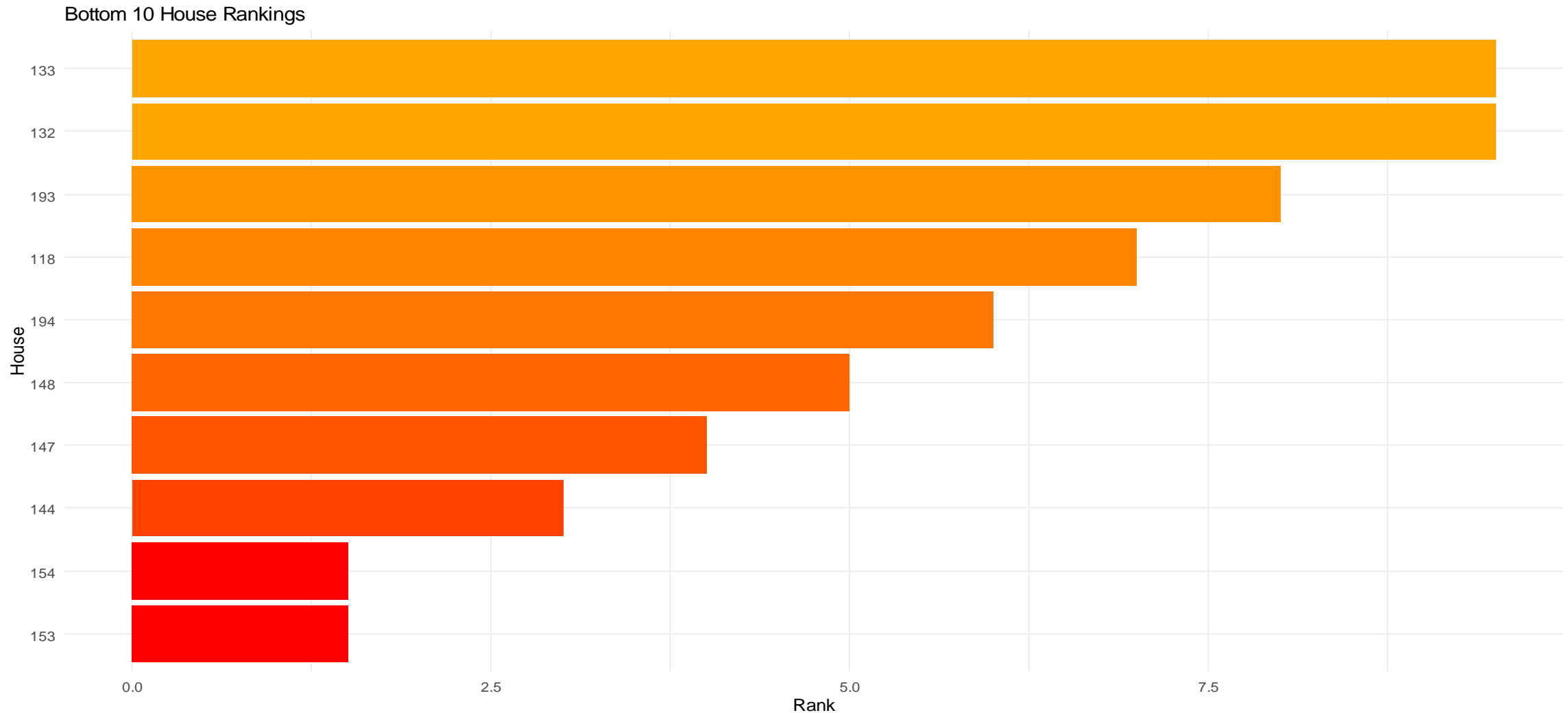
Calculation of a weighted score:

$$\text{country\_std} * \text{weights}[0.4] + \text{resident0\_std} * \text{weights}[0.3] + \text{year\_std} * \text{weights}[0.2] + \text{bavsalary\_std} * \text{weights}[0.1]$$

# Houses which require renovation first – lowest rank score



The outcome of the model is ranking of the buildings as it is exemplified below. Where the houses that collected less points (dark red) require renovation first.

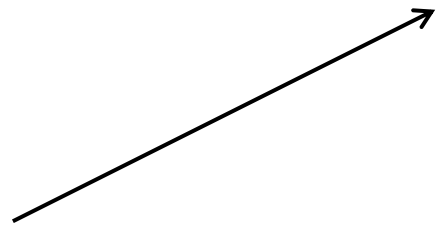


- Small dataset!
  - Longitudinal data would be preferable (e.g. 5 years)
  - More variables to include (e.g. some other technical characteristics of the buildings are omitted)
- Seasonality is not captured (cross-sectional data)
- Gas data was not available for all the buildings.
- Given a richer dataset, a different statistical model can be applied.
- Slovenian sample was not full.

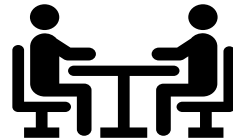
# Who can use the model?



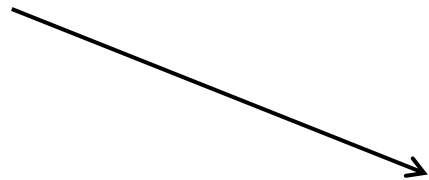
**Data owners!**



Local authorities (e.g., for renovation of municipal buildings)



Financial intermediaries (e.g., KredEx in Estonia) for deciding on the grants



Energy agencies for data driven advise





**CEE  
SEN**