



Applied Econometrics
R Project - Case 1 - Part 02

Kelian IDRISOU and Ibrahim FOUSFOS

Table of Contents :

<u>Section</u>	<u>Page</u>
1. Overview	1-7
2. Sources	6-7
3. Data Manipulation	7-8

Overview

1) Data Import And Summary by Treatment Group :

We choose Savings and retirement as our two dependent variables (Y) because they directly capture the client's saving behaviour which is the core focus of our research question.

The purpose of the study is to determine whether the meeting with a salesperson influences how much clients save in general and more specially how much they contribute to their retirement savings.

- **Savings** → Reflects the overall saving behaviour of clients. It shows whether the meeting encourages people to save more money in general.
- **retirement** → Focuses on long term financial planning and measures whether the meeting motivates clients to invest more in their retirement plans.

Together, these two variables allow us to evaluate both short term and long term saving decisions giving a more complete view of the meeting impact.

Table 1: Summary statistics by treatment group

meeting	n	Retirement			Savings		
		Mean	SD	Median	Mean	SD	Median
Control	29530	214.89	388.77	75.38	9261.31	7024.26	7780.96
Treated	20470	1264.18	1467.31	807.12	11310.81	7814.36	9641.03

2) Creation Of New Variables And Data Transformations :

The variables `log_income`, `log_savings` and `log_retirement` help linearize relationships and reduce the effect of extreme values, allowing percentage-based interpretations.

We created:

- **log_saving_rate** → measures the proportion of income saved on a logarithmic scale. It helps understand saving intensity relative to income and compare clients with different income levels.
- **log_retirement_rate** → measures how much of a client's income goes to retirement savings on a logarithmic scale. It captures long-term saving effort and is less sensitive to income differences.

3) Distribution Of Outcome Variables with Histograms and interpretations :

We focus on the distribution of post-treatment savings (periods $t \geq 2$) at the individual level. The first histogram shows that savings are positively skewed, with most clients saving relatively small amounts and a few saving much more.

When we split the data by the treatment group (meeting), both distributions have a similar shape, but the clients who attended the meeting (in blue) tend to have slightly higher savings on average.

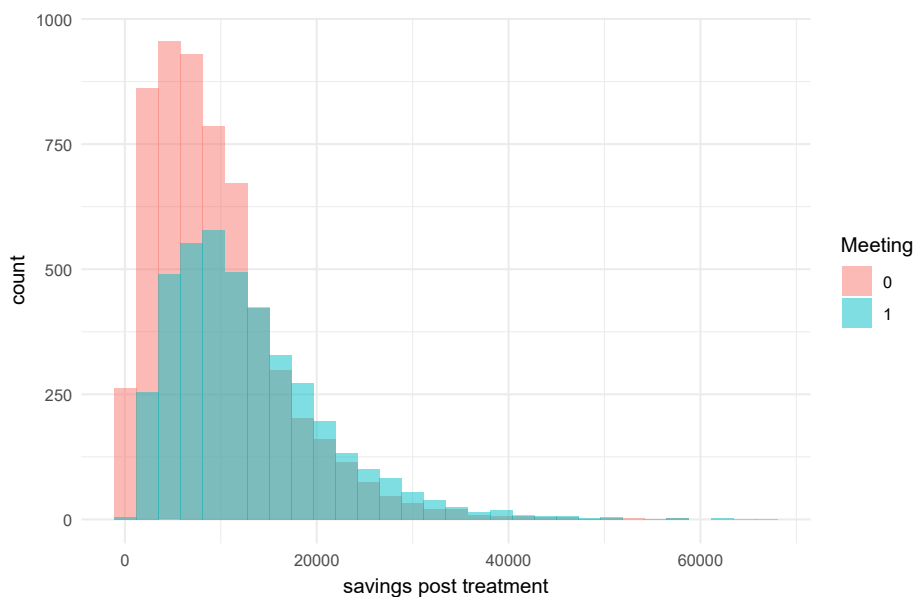


Figure 1: Distributions by treatment

When weighting the data to give each group the same total mass, the shape of the distribution becomes clearer and comparable across groups. The weighted histogram shows that although the two groups have similar savings patterns, the treated group maintains a small upward shift in saving levels.

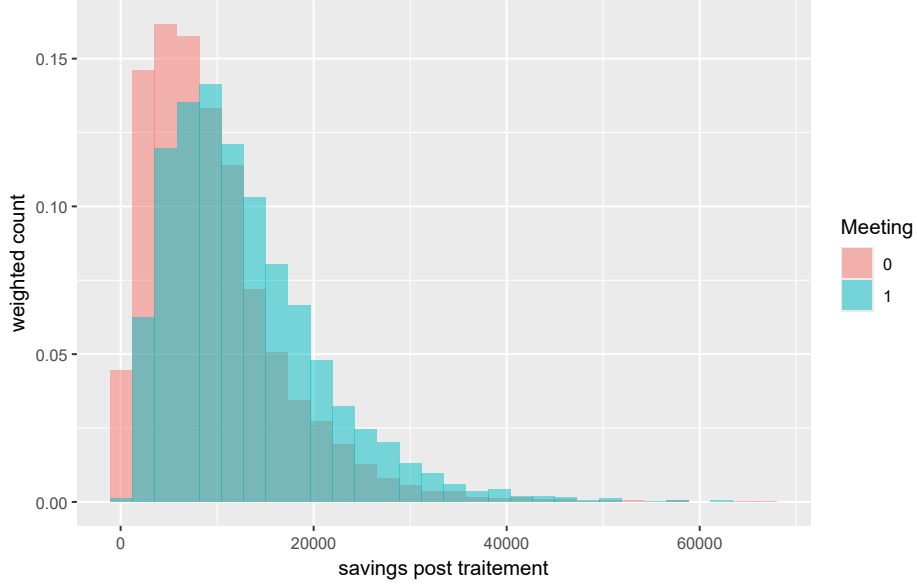


Figure 2: Distributions by treatment with weights

4) *scatter plot Analysis Between Income And Savings :*

The scatter plot of $\log \text{ savings}$ against $\log \text{ income}$ shows a clear positive relationship, meaning that clients with higher income tend to save more. The relationship is economically consistent with the life-cycle theory (Modigliani and Brumberg, 1957) and the permanent income theory (Friedman, 1957). According to these theories, individuals aim to smooth consumption over time and therefore save part of their income during high-earning periods to prepare for future needs or retirement.

In the plot, the fitted trend shows that as income increases, savings also rise proportionally. The log-log specification highlights a nearly linear relationship, suggesting that the elasticity of savings with respect to income is relatively stable.

When comparing groups, clients who attended the meeting (in blue) generally lie slightly above the untreated group, showing higher savings for a given income level. This suggests that the meeting may have reinforced financial awareness or encouraged better saving habits. The difference between groups is more visible among lower-income clients, who appear to benefit more from the intervention.

Overall, the data support the theoretical prediction that income positively influences savings, while also indicating a small behavioral effect of the meeting.

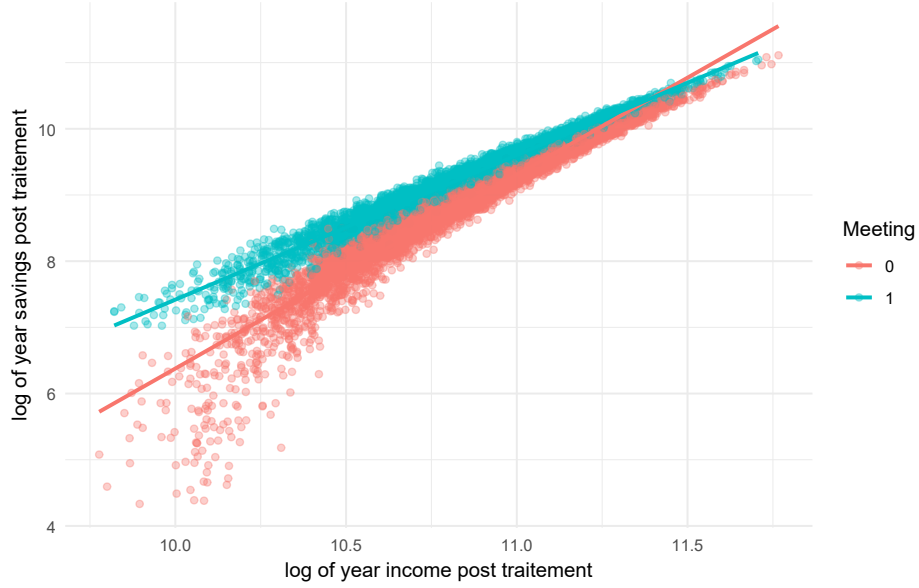


Figure 3: Income–Savings Relationship (log–log), post-treatment period

5) Creation Of A Dummy Variable And Probit Regression Analysis :

To test whether the treatment (`meeting` = 1) was randomly assigned, we compared the pre-treatment characteristics of treated and control clients. We estimate a probit model to assess whether treatment assignment was random with respect to pre-treatment characteristics such as income, savings, retirement savings, and gender. From the probit estimation, we compute the *average marginal effects* (AME), which measure how a one-unit change in each explanatory variable affects the probability of being treated on average. The AME are all close to zero and statistically insignificant (all $p > 0.10$), indicating that the pre-treatment variables have no meaningful effect on the probability of being treated.

Variable	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	-0.22	0.10	-2.06	0.04
inc_pre	-0.00	0.00	-0.19	0.85
sav_pre	0.00	0.00	0.53	0.59
ret_pre	-0.00	0.00	-1.06	0.29
femaleTRUE	-0.02	0.03	-0.78	0.43

Table 2: Probit coefficients

The average income before the meeting is 49,864 for the control group and 50,053 for the treated group, with a difference of only 189 €; the associated p -value equals 0.55, which indicates that this difference is not statistically significant, meaning that the two groups had similar income levels before the intervention.

Mean savings are 8,943 € and 9,025 €, with a difference of 82 € ($p = 0.55$), and mean retirement

Table 3: Balance table: pre-treatment characteristics by treatment status

Variable	Control mean	Treated mean	Diff (1-0)
Pre-income (€)	49,863.97	50,052.83	188.86
Pre-savings (€)	8,942.85	9,025.15	82.29
Pre-retirement (€)	207.83	205.23	-2.61
Female (=1)	52.7%	51.9%	-0.8%

savings are 208 € and 205 € with a difference of -2.6 € ($p = 0.73$). Gender composition is also balanced with 52.7% women in the control group compared to 51.9% in the treated group (χ^2 test, $p = 0.44$).

Variable	Mean..0.	Mean..1.	Diff..1.0.	t	p.value
Pre-income (€)	49863.97	50052.83	188.86	-0.60	0.55
Pre-savings (€)	8942.85	9025.15	82.29	-0.59	0.55
Pre-retirement (€)	207.83	205.23	-2.61	0.35	0.73

Table 4: Equality-of-means tests (pre-treatment variables by treatment status)

Variable	Prop (0)	Prop (1)	Chi-square	df	p-value
Female (=1)	52.7%	51.9%	0.59	1	0.44

Table 5: Chi-square test for equality of proportions (Female)

None of these differences is statistically significant (all $p > 0.10$), meaning that before the policy intervention both groups had comparable income, saving behaviour, and gender composition. These results are consistent with random assignment of the meeting, so any differences observed after the meeting can plausibly be attributed to the effect of the meeting itself.

6) OLS Estimation And Correlation Analysis :

The graph of the mean savings (post) by `meeting` displays the mean *post-treatment* savings for clients who had a meeting (1) and those who did not (0), with 95% confidence intervals. It shows a clear upward shift for the treated group: on average, clients who attended the meeting saved substantially more after the intervention.

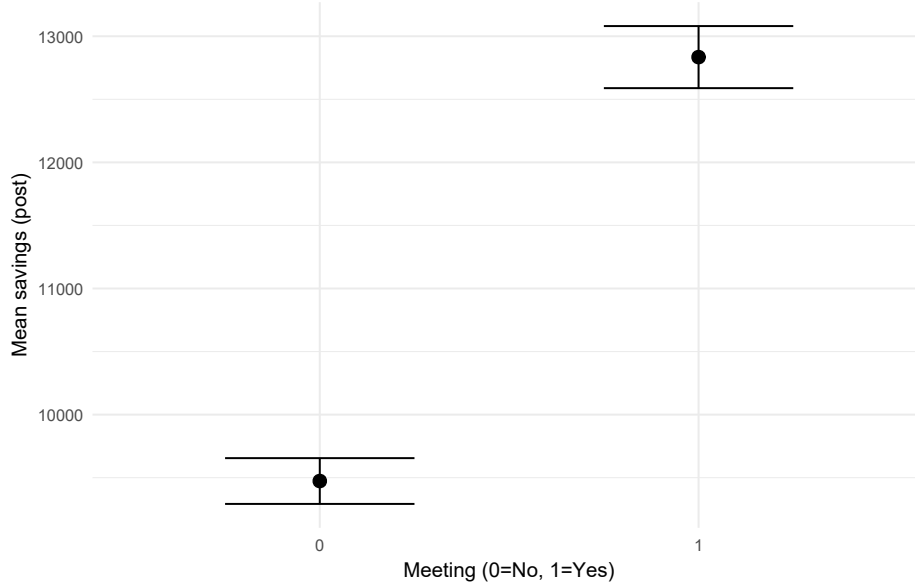


Figure 4: Mean post-treatment savings by meeting participation with 95% confidence interval

To assess the correlation between the policy intervention and saving behaviour, we estimate a simple OLS regression of post-treatment savings (Y_{post}) on the dummy variable (**meeting**).

Table 6: OLS: Y_{post} *meeting*

<i>Dependent variable:</i>	
	Y_{post}
Meeting	3,360.963*** (152.629)
Constant	9,473.620*** (97.659)
Observations	10,000
R^2	0.046
Adjusted R^2	0.046
F Statistic	484.903*** (df = 1; 9998)
<i>Note:</i>	*p<0.1; **p<0.05; ***p<0.01

The regression results show a positive and highly significant coefficient for **meeting** ($\hat{\beta} = 3,361$, $p < 0.001$). This means that, on average, clients who attended the meeting saved about 3,361 € more after the intervention than those who did not. The coefficient is statistically significant at the 1% level, indicating strong evidence of a positive association. The intercept (9,474) represents the average post-treatment savings for the control group (clients without a meeting).

The R^2 of 0.046 suggests that the meeting explains about 4.6% of the variation in savings,

meaning that while the meeting has a significant effect, most of the variation is still due to other factors.

Overall, the results indicate that the meeting has a positive and statistically significant correlation with post-treatment savings, suggesting that the intervention likely encouraged higher saving behaviour.

7) Extended OLS Regression With Covariates :

To understand the relationship between the meeting intervention and post-treatment savings, we estimate an extended OLS regression including several relevant covariates that were selected as:

- **Pre-income** controls for clients' financial capacity, since higher-income individuals typically save more.
- **Pre-savings** capture baseline saving behaviour, ensuring we compare clients with similar prior habits.

In our R code, we first defined the outcome (*savings*) and its log version to allow both level and percentage interpretations. The post-treatment (*Y_post*) was computed for each client ($t \geq 2$), while pre-treatment averages (*inc_pre*, *Y_pre*) served as covariates to control for income capacity and prior saving habits.

The results shown in the table indicate a positive and highly significant coefficient for *meeting* (treated) across both specifications. In the level model, treated clients save 3,275 € more on average than those without a meeting. In the log specification, the coefficient of 0.419 suggests an increase of approximately 42% in post-savings for treated clients, controlling for income, past savings and gender.

Both pre-income and pre-savings are also strong predictors of post-treatment savings, as expected. The models show excellent fit with $R^2 = 0.985$ for levels and $R^2 = 0.944$ for logs, suggesting that most variation in post-savings is explained by these covariates.

8) Economic Interpretation And Discussion Of Results:

From an economic perspective, our results are consistent with our initial expectations : clients who attended a meeting with a salesperson have significantly higher saving levels, both in total and in retirement savings. The positive and robust coefficient of the meeting variable across all models supports the hypothesis that follow-up meetings encourage stronger behaviour.

This finding aligns with the life cycle theory of Modigliani and Brumberg (1954) which suggests that individuals seek to smooth consumption and increase savings when their income or financial awareness rises. The meeting likely acted as a behavioral nudge improving financial planning and promoting long-term saving decisions.

Our extended OLS regressions confirm that this effect remains significant even after controlling for pre-income and pre-savings, meaning that the relationship is not driven by initial income differences or prior saving habits. The earlier randomness checks (T-tests and probit model) also show

that the treatment was assigned independently of observable characteristics, reducing selection bias and reinforcing the internal validity of the analysis.

However, while the estimated relationship is statistically strong and economically meaningful, it is incorrect to say that it is considered fully causal without accounting for potential unobserved heterogeneity, for instance differences in motivation or financial literacy that might also influence saving behaviour.

Overall, the results are in line with theoretical expectations and indicate that the meeting intervention likely had a positive and substantial impact on the client's saving behaviour. Yet, the relationship should be interpreted as a robust correlation with plausible causal direction, rather than definitive proof of causality.

Sources

lien régression scatter

- **Applied Econometrics lecture slides :**
 - <https://cours.univ-paris1.fr/course/view.php?id=44026>
- **Life-Cycle Hypothesis developed by Modigliani in 1954 :**
 - <https://www.lafinancepourtous.com/decryptages/finance-perso/epargne-et-placement/epargne/la-theorie-du-cycle-de-vie/>
- **The Table function for making a table with R studio :**
 - <https://bookdown.org/yihui/rmarkdown-cookbook/kable.html>
- **To know how to use LaTeX :**
 - https://www.overleaf.com/learn/latex/Learn_LaTeX_in_30_minutes
 - <https://latex-tutorial.com/tutorials/pgfplotstable/>
 - <https://www.youtube.com/watch?v=lgICpA4zzGU&t=1s>
- **To know how to make a Histogram with ggplot :**
 - <https://www.datacamp.com/tutorial/make-histogram-basic-r>
 - https://r-graph-gallery.com/histogram_several_group.html
 - <https://www.youtube.com/watch?v=onEumD5xUOE&t=1s>
- **To know how to scatter :**
 - <https://www.datacamp.com/doc/r/scatterplot-in-r>
 - <https://www.r-bloggers.com/2020/07/create-a-scatter-plot-with-ggplot/>

- To know how to make a regression :
 - <https://www.datacamp.com/tutorial/linear-regression-R>
 - <https://www.youtube.com/watch?v=MEPP5oJ4rWc>
- To know how to make an OLS regression :
 - <https://www.r-bloggers.com/2017/07/ordinary-least-squares-ols-linear-regression-in-r/>
 - <https://www.geeksforgeeks.org/r-machine-learning/ordinary-least-squares-ols-regression/>
 - <https://cran.r-project.org/web/packages/olsrr/vignettes/intro.html>

Script R du Projet

```
1  ### library
2  library(margins)
3  library(dplyr)
4  library(ggplot2)
5
6  #### we import the dataset
7  groupe_41_pt2 = read.csv("C:/Users/ibrah/OneDrive/Documents/ cole /S7/
   econometrie appliquee/projet/group41.csv")
8
9  ### Q1) choosing 2 variable of interests which are saving and retirements
10
11 groupe_41_pt2 = groupe_41_pt2 %>% mutate(meeting = as.integer(meeting))
12
13 sum_by_group <- groupe_41_pt2 %>%
14   group_by(meeting) %>%
15   summarise(
16     n = n(),
17     mean_ret = mean(retirement, na.rm=TRUE), # na.rm = True icase we have NAAN
18     sd_ret = sd(retirement, na.rm=TRUE),
19     p50_ret = median(retirement, na.rm=TRUE),
20     mean_sav = mean(savings, na.rm=TRUE),
21     sd_sav = sd(savings, na.rm=TRUE),
22     p50_sav = median(savings, na.rm=TRUE)
23   )
24 sum_by_group
25
26
27 ### Q2) data transformation
28
29 ## logarithmic transformation on the data if yincome = 0 then the log should =
   0
30
31 groupe_41_pt2$log_yincome = ifelse(groupe_41_pt2$yincome == 0, log(groupe_41_pt2
   $yincome + 1), log(groupe_41_pt2$yincome))
32 groupe_41_pt2$log_savings = ifelse(groupe_41_pt2$savings == 0, log(groupe_41_pt2
   $savings + 1), log(groupe_41_pt2$savings))
33 groupe_41_pt2$log_retirement = ifelse(groupe_41_pt2$retirement == 0, log(groupe_
   41_pt2$retirement + 1), log(groupe_41_pt2$retirement))
34
35
36 ## simples rates of 2 main variables
37 groupe_41_pt2$saving_rate = groupe_41_pt2$savings / groupe_41_pt2$yincome
38 groupe_41_pt2$retirement_rate = groupe_41_pt2$retirement / groupe_41_pt2$
   yincome
39
40 head(groupe_41_pt2)
41
42 ### same transformation but after traitement
43 groupe_post_D = subset(groupe_41_pt2, time>1)
44
45 ## logarithmic transformation on the data
46 groupe_post_D$log_yincome = ifelse(groupe_post_D$yincome == 0, log(groupe_post_D
   $yincome + 1), log(groupe_post_D$yincome))
```

```

47 groupe_post_D$log_savings = ifelse(groupe_post_D$savings == 0, log(groupe_post_D
  $savings + 1), log(groupe_post_D$savings))
48 groupe_post_D$log_retirement = ifelse(groupe_post_D$retirement == 0, log(groupe_
  post_D$retirement + 1), log(groupe_post_D$retirement))
49
50 head(groupe_post_D)
51
52 ### Q3) hist of yearly savings in periods t>1 (we can do the same with
  retirement variable)
53 ## we do a common unit, a mean savings of the 3 periods (2;3;4)
54
55 ### common unit
56 unit_post_savings <- groupe_41_pt2 %>%
57   filter(time >= 2) %>%
58   group_by(id, meeting) %>%
59   summarise(sav_post = mean(savings, na.rm = TRUE), .groups = "drop")
60 head(unit_post_savings)
61
62
63 ## A/ histogramme per group
64 ggplot(unit_post_savings, aes(x = sav_post, fill = factor(meeting))) +
65   geom_histogram(alpha = .5, position = "identity", bins = 30) +
66   labs(title="Distribution by treatment", x="savings post treatment", fill="
    Meeting") +
67   theme_minimal()
68
69 ## B/ histogramme per group but with weights
70 #B.1/ new comon unit
71 weight_tab <- unit_post_savings %>%
72   count(meeting, name = "n") %>%
73   mutate(p = n / sum(n),
74          w_equal = 1 / n,
75          w_prop = p / n)
76 unit_weighted <- unit_post_savings %>% left_join(weight_tab, by = "meeting")
77
78 # B.2/ histograme with weight
79 ggplot(unit_weighted, aes(x = sav_post, fill = factor(meeting))) +
80   geom_histogram(aes(weight = w_equal), position = "identity", alpha = .5, bins
    = 30) +
81   labs(title="Distributions by treatment with weights", x = "savings post
    traitement", y = "weighted count", fill = "Meeting")
82
83 ## Q4) Scatter
84
85 # 1) we have 3 periods so we use the mean and use the log of the means and we
  put it in unit_post
86
87 unit_post <- groupe_41_pt2 %>%
88   filter(time >= 2) %>%
89   group_by(id, meeting) %>%
90   summarise(sav_post = mean(savings, na.rm=TRUE),
91             inc_post = mean(yincome, na.rm=TRUE), .groups="drop")
92
93 unit_post$log_inc_post <- log1p(unit_post$inc_post)
94 unit_post$log_sav_post <- log1p(unit_post$sav_post)
95 # B/ we draw a scatter
96 ggplot(unit_post, aes(x = log_inc_post, y = log_sav_post, color = factor(
  meeting))) +
97   geom_point(alpha = 0.35) +
98   geom_smooth(method = "lm", se = FALSE) +
99   labs(title = "Income Savings Relationship (log log), post-treatment
    period",
100        x = "log of year incom post traitement", y = "log of year savings post

```

```

    traitement", color = "Meeting") +
101   theme_minimal()
102
103
104   ### Q5)
105   # A/ Construire les covariables PRE au niveau unit (t < 2)
106   pre_unit <- groupe_41_pt2 %>%
107     filter(time < 2) %>%
108     group_by(id, meeting, female) %>%
109     summarise(
110       inc_pre = mean(yincome, na.rm = TRUE),
111       sav_pre = mean(savings, na.rm = TRUE),
112       ret_pre = mean(retirement, na.rm = TRUE),
113       .groups = "drop"
114     )
115   # we do a regression with probit
116   m_probit <- glm(meeting ~ inc_pre + sav_pre + ret_pre + female,
117     data = pre_unit,
118     family = binomial(link = "probit"))
119
120   summary(m_probit) # signs & z-stats (quick look)
121
122   #Average marginal effects (ame)
123   ame <- margins(m_probit)
124   summary(ame)
125   #####
126
127
128   # B/ difference of means table
129   bal_table <- pre_unit %>%
130     summarise(
131       mean_inc_0 = mean(inc_pre [meeting==0], na.rm=TRUE),
132       mean_inc_1 = mean(inc_pre [meeting==1], na.rm=TRUE),
133       diff_inc   = mean_inc_1 - mean_inc_0,
134
135       mean_sav_0 = mean(sav_pre [meeting==0], na.rm=TRUE),
136       mean_sav_1 = mean(sav_pre [meeting==1], na.rm=TRUE),
137       diff_sav   = mean_sav_1 - mean_sav_0,
138
139       mean_ret_0 = mean(ret_pre [meeting==0], na.rm=TRUE),
140       mean_ret_1 = mean(ret_pre [meeting==1], na.rm=TRUE),
141       diff_ret   = mean_ret_1 - mean_ret_0,
142
143       prop_fem_0 = mean(female [meeting==0], na.rm=TRUE),
144       prop_fem_1 = mean(female [meeting==1], na.rm=TRUE),
145       diff_fem   = prop_fem_1 - prop_fem_0
146     )
147   bal_table
148
149   ## C/Simple Tests of equality between groups
150   #for continous variables
151   t_inc <- t.test(inc_pre ~ meeting, data = pre_unit)
152   t_sav <- t.test(sav_pre ~ meeting, data = pre_unit)
153   t_ret <- t.test(ret_pre ~ meeting, data = pre_unit)
154
155   # for limited variable (binaire)
156   tab_f <- table(pre_unit$female, pre_unit$meeting)
157   chi_f <- chisq.test(tab_f)
158
159   # print of results
160   t_inc; t_sav; t_ret
161   chi_f
162

```

```

163
164 ##### Q6) simple OLS regression
165 # A/ mean of savings post traitement at t=2;3;4
166 post_unit_6 <- groupe_41_pt2 %>%
167   filter(time >= 2) %>%
168   group_by(id, meeting) %>%
169   summarise(
170     Y_post = mean(.data[["savings"]], na.rm = TRUE),
171     .groups = "drop"
172   )
173
174 # B/ simple OLS : Y_post ~ meeting
175 m6 <- lm(Y_post ~ meeting, data = post_unit_6)
176 summary(m6)
177
178 # C/ Petit plot moyennes par groupe avec IC 95%
179 m6_means <- post_unit_6 %>%
180   group_by(meeting) %>%
181   summarise(
182     n = n(),
183     mu = mean(Y_post, na.rm = TRUE),
184     se = sd(Y_post, na.rm = TRUE) / sqrt(n),
185     ci_low = mu - 1.96 * se,
186     ci_high = mu + 1.96 * se,
187     .groups = "drop"
188   )
189
190 ggplot(m6_means, aes(x = factor(meeting), y = mu)) +
191   geom_point(size = 3) +
192   geom_errorbar(aes(ymin = ci_low, ymax = ci_high), width = .5) +
193   labs(title = paste("Mean", "savings", "(post) by Meeting"),
194        x = "Meeting (0=No, 1=Yes)", y = paste("Mean", "savings", "(post)")) +
195   theme_minimal()
196
197
198 #####7 extended OLS regression
199
200 #A/ mean of variable post traitement
201 post_unit_7 <- groupe_41_pt2 %>%
202   filter(time >= 2) %>%
203   group_by(id, meeting, female) %>%
204   summarise(
205     Y_post = mean(.data[["savings"]], na.rm = TRUE),
206     logY_post = mean(.data[["log_savings"]], na.rm = TRUE),
207     .groups = "drop"
208   )
209
210 # 2) mean of variables pre traitement
211 pre_unit_7 <- groupe_41_pt2 %>%
212   filter(time < 2) %>%
213   group_by(id) %>%
214   summarise(
215     inc_pre = mean(yincome, na.rm = TRUE),
216     log_inc_pre = mean(log_yincome, na.rm = TRUE),
217     Y_pre = mean(.data[["savings"]], na.rm = TRUE),
218     logY_pre = mean(.data[["log_savings"]], na.rm = TRUE),
219     .groups = "drop"
220   )
221
222 # C/ we do an intercection of the 2 tables to do the regression easely
223 dat7 <- left_join(post_unit_7, pre_unit_7, by = "id")
224
225 # D/ 1) level OLS

```

```

226 m7_levels <- lm(Y_post ~ meeting + inc_pre + Y_pre + female, data = dat7)
227 summary(m7_levels)
228 dat7
229 # D/ 2) log OLS (option : lecture      % ; garde ta fa on de loguer)
230 m7_logs <- lm(logY_post ~ meeting + log_inc_pre + logY_pre + female, data =
      dat7)
231 summary(m7_logs)

```