

Lista Prática 3

Exercício 1 (*swirl*). Usando o pacote *swirl*, instale o curso “Exploratory Data Analysis” usando a seguinte função `swirl::install_course("Exploratory Data Analysis")`¹ e faça as seguintes lições:

- a) 5: Base Plotting System
- b) 9: GGPlot2 Part²

Exercício 2. Neste exercício, usaremos a base de dados do General Social Survey (GSS), que coleta dados demográficos nos EUA e pode ser carregada no R usando o comando³:

```
1 load(url("http://bit.ly/dasi-gss-data"))
```

Queremos saber a relação entre anos de estudo (*educ*) e a renda das pessoas (*coninc*) que **trabalham em tempo integral** para homens e mulheres.

- a) Gere duas bases a partir da GSS: uma apenas com mulheres (*gss_mulher*) e outra apenas com homens (*gss_homem*), e estime os seguintes modelos:

$$\begin{aligned} \text{coninc} &= \beta_M \cdot \text{educ} + \varepsilon && \text{(base com mulheres)} \\ \text{coninc} &= \beta_H \cdot \text{educ} + \varepsilon && \text{(base com homens)} \end{aligned}$$

Note que o modelo não possui uma constante (insira um regressor 0 para não incluí-la na regressão via `lm()`). Quais são as estimativas para β_M e β_H ?

- b) Usando a base completa (*gss*), plote um gráfico de dispersão (*scatterplot*) entre anos de estudo \times renda, colorindo os pontos de acordo com o sexo da pessoa. Também, adicione as retas das regressões feitas no item (a) com cores distintas.
- c) Na base completa (*gss*), regrida um único modelo em que é possível calcular os mesmos β_M e β_H encontrados no item (a). A diferença entre β_M e β_H é estatisticamente significativa?

Resposta:

- a) $\beta_M = 3576.56$ e $\beta_H = 4229.01$

¹Caso tenha problemas, instalar manualmente: <http://swirlstats.com/scn/eda.html>

²O 1º exercício irá pedir para fazer um gráfico usando `qplot()`, você pode dar `skip()` para avançar para a parte com gráficos usando a função `ggplot()`

³Disponibilizado por Bryan Wheeler (2014)

```

1 library(dplyr)
2
3 # Gerando 2 bases com trab. em tempo integral de mulheres e de homens
4 gss_homem = gss %>% filter(sex=="Male", wrkstat=="Working Fulltime")
5 gss_mulher = gss %>% filter(sex=="Female", wrkstat=="Working Fulltime")
6
7 # Estimando os modelos sem constante
8 fit_homem = lm(coninc ~ 0 + educ, gss_homem)
9 fit_mulher = lm(coninc ~ 0 + educ, gss_mulher)
10
11 # Mostrando os betas estimados
12 fit_homem$coef
13 fit_mulher$coef

```

```

1 4229.012
2 3576.562

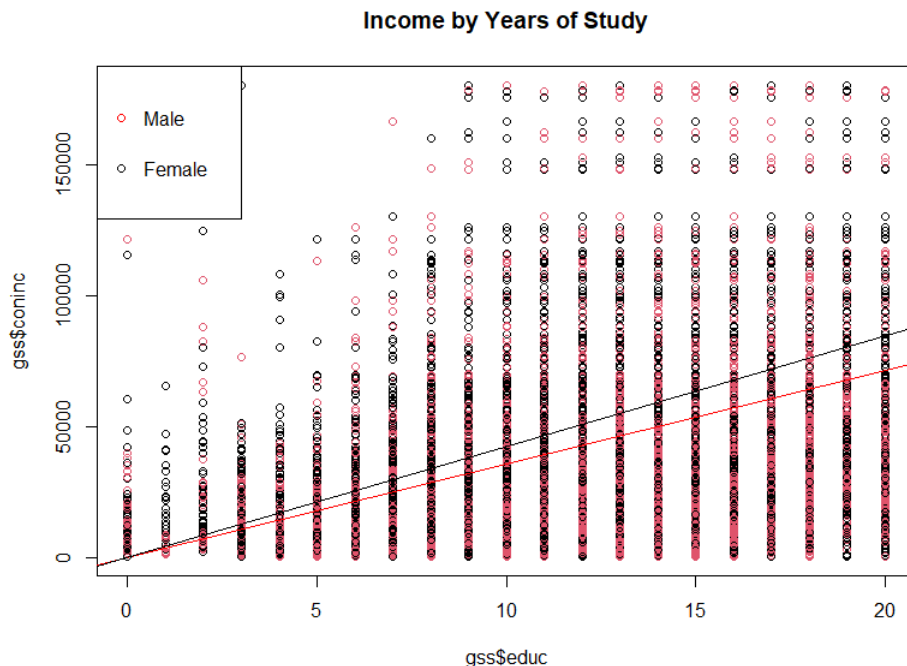
```

b) # Plotando anos de educação X renda

```

2 plot(gss$educ, gss$coninc, col=gss$sex, main="Income by Years of Study")
3 # Plotando as retas das regressões estimadas
4 abline(fit_mulher, col="red")
5 abline(fit_homem, col="black")
6 legend("topleft", pch=1, col=c("black", "red"), legend=c("Male", "Female"))

```

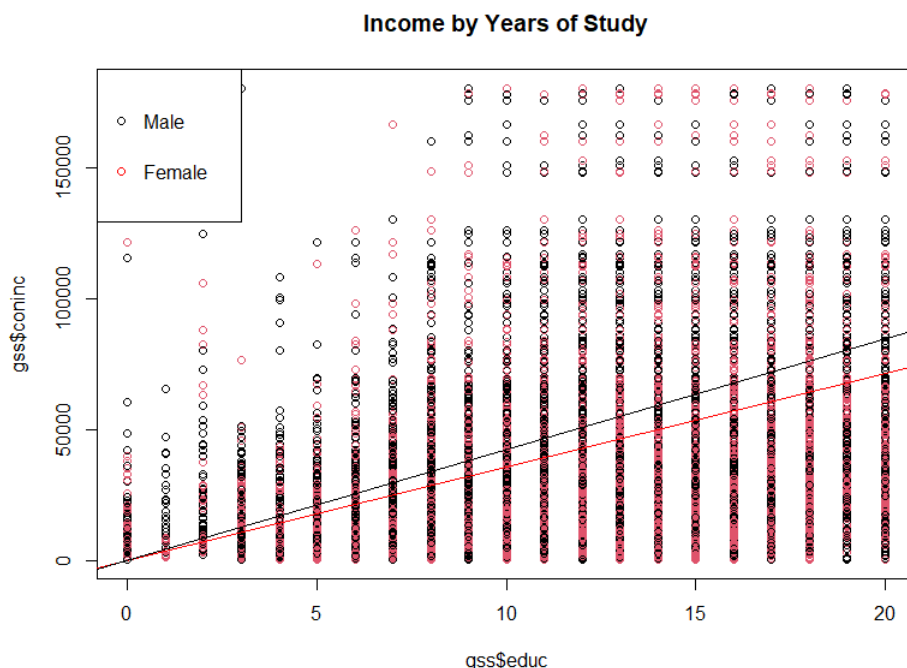


ou também por ggplot2:

```

1 library(ggplot2)
2 g = ggplot(gss, aes(educ, coninc))
3 g + geom_point(aes(color=sex), size=.9, alpha=0.3) +
4   geom_abline(intercept=c(0,0),
5               slope=c(fit_homem$coef, fit_mulher$coef),
6               color=c("darksalmon", "cyan3"))

```



- c) Vamos estimar o modelo que, além da covariada de anos de estudo, vamos incluir a interação entre anos de estudo e a *dummy* mulher:

$$\text{coninc} = \beta_1 \cdot \text{educ} + \beta_2 \cdot \text{educ} \cdot \text{mulher} + \varepsilon$$

```
1 # Na base completa, filtro trab tempo integral e criando Dummy mulher
2 gss = gss %>% filter(wrkstat=="Working Fulltime") %>%
3   mutate(mulher = ifelse(sex=="Female", 1, 0))
4
5 # Estimando o modelo com interação entre anos de educ. e dummy mulher
6 lm(coninc ~ 0 + educ + educ:mulher, gss) %>% summary()
```

```
1           Estimate Std. Error t value Pr(>|t|)
2 educ           4229.01      19.75   214.08  <2e-16 ***
3 educ:mulher    -652.45      29.84   -21.86  <2e-16 ***
4 ---
5 Signif. codes:  0  ***  0.001  **  0.01  *  0.05  .  0.1
```

Note que $\beta_H = \beta_1$ e $\beta_M = \beta_1 + \beta_2$. Além disso, a diferença entre as estimativas, $\beta_M - \beta_H = \beta_2$, é significativa a 0,1%.

□

Exercício 3. Neste exercício, utilizaremos a base de dados de Card e Krueger (1994), a qual foi usada para analisar o efeito do aumento do salário mínimo estadual na quantidade de empregos. Uma versão modificada dela pode ser carregada no R usando:

```
1 fastfood = read.csv("https://fhnishida.github.io/fearp/eco1/CardKrueger_1994_
  modificado.csv")
```

- *id*: identificador da loja de fast food
- *state*: estado dos EUA (New Jersey ou Pennsylvania)
- *emptot_feb92*: total de empregados em fev/92
- *emptot_nov92*: total de empregados em nov/92
- *chain*: rede de fast food
- *pct_fte*: percentual de trabalhadores em tempo integral
- *wage_st*: salário por hora
- *hrsopen*: horas de funcionamento diário

Idealmente (para avaliação), haveria um experimento em que estados seriam aleatoriamente selecionados para aumentarem o salário mínimo e, logo, o efeito da intervenção poderia ser estimado por meio da diferença das médias de variação de emprego entre os estados que elevaram seus salários mínimos (tratados) e os que não alteraram (não-tratados/controle). No entanto, em políticas públicas, a aleatorização pode ser custosa (financeira e politicamente) ou operacionalmente inviável. Neste caso, para identificar esse efeito, Card e Krueger (1994) assumiram o estado vizinho de Pennsylvania como o controle, por não ter elevado o seu salário mínimo e ser considerado, ao menos no mercado de trabalho de fast food, similar a New Jersey. Sob essa premissa, os estados manteriam o mesmo comportamento se ambos não alterassem seu nível de salário mínimo e, portanto, o efeito do tratamento pode ser calculado pela diferença entre as variações de empregados dos dois estados. Por exemplo, se a variação do n^o de empregados em New Jersey foi consideravelmente maior em relação a de Pennsylvania, então haveria indício de um efeito positivo e significativo da mudança na política salarial no nível de emprego.

- a) *Calcule a diferença de médias das variações de empregos, entre fev/92 e nov/92, das lojas tratadas e das lojas de controle. A diferença foi significativa?*
- b) *Para avaliar o quão razoável é supor que ambos estados são similares no mercado de trabalho de fast food, compararemos características de New Jersey e de Pennsylvania por meio de diferenças de médias para as variáveis: (i-iv) dummies de cada rede de fast food, (v) n^o de empregados em fev/92, (vi) proporção de empregados em tempo integral, (vii) salário por hora e (viii) tempo de funcionamento. Construa uma tabela com as médias de cada estado, a diferença entre eles e o p valor:*

	Pennsylvania	New Jersey	dif.	p valor
dmy_bk	-	-	-	-
dmy_kfc	-	-	-	-
dmy_roys	-	-	-	-
dmy_wendys	-	-	-	-
emptot_feb92	-	-	-	-
pct_fte	-	-	-	-
wage_st	-	-	-	-
hrsopen	-	-	-	-

A partir desta análise, é razoável supor que os mercados de trabalho das redes de fast food eram similares em New Jersey e Pennsylvania?

Resposta:

```
a) library(dplyr)
2 fastfood = fastfood %>%
3   mutate(treated = ifelse(state=="New Jersey", 1, 0),
4          var_emp = emptot_nov92 - emptot_feb92)
5
6 lm(var_emp ~ treated, data = fastfood) %>% summary()
```

```
1          Estimate Std. Error t value Pr(>|t|)
2 (Intercept) -2.2833      0.7313  -3.122  0.00186 **
3 treated      2.7500      0.8152   3.373  0.00078 ***
4 ---
5 Signif. codes:  0 *** 0.001 ** 0.01 * 0.05 . 0.1
```

Houve um aumento significativo na variação do n^o de empregados nas redes de fast food de New Jersey (tratado) em relação ao de Pennsylvania (controle).

```
b) # Diferenças de médias - variáveis pré-tratamento
2 fastfood = fastfood %>%
3   mutate(dmy_bk = ifelse(chain=="bk", 1, 0),
4          dmy_kfc = ifelse(chain=="kfc", 1, 0),
5          dmy_roys = ifelse(chain=="roys", 1, 0),
6          dmy_wendys = ifelse(chain=="wendys", 1, 0))
7
8 lista_var = c("dmy_bk", "dmy_kfc", "dmy_roys", "dmy_wendys",
9              "emptot_feb92", "pct_fte", "wage_st", "hrsopen")
10
11 resultados = matrix(0, length(lista_var), 4) %>%
12 as.data.frame()
13 colnames(resultados) = c("Pennsylvania", "New Jersey", "dif", "p valor")
14 rownames(resultados) = lista_var
15
16 for (var in lista_var) {
17   reg = lm(eval(parse(text=var)) ~ treated, fastfood) %>% summary() %>% coef()
18   resultados[var, 1] = reg["(Intercept)", "Estimate"]
19   resultados[var, 2] = reg["(Intercept)", "Estimate"] + reg["treated", "
20     Estimate"]
21   resultados[var, 3] = reg["treated", "Estimate"]
22   resultados[var, 4] = reg["treated", "Pr(>|t|)"]
23 }
24 resultados %>% round(4)
```

	Pennsylvania	New Jersey	dif	p valor
dmy_bk	0.4430	0.4109	-0.0322	0.4619
dmy_kfc	0.1519	0.2054	0.0535	0.1274
dmy_roys	0.2152	0.2477	0.0325	0.3910
dmy_wendys	0.1899	0.1360	-0.0539	0.0851
emptot_feb92	23.3312	20.4394	-2.8918	0.0009
pct_fte	32.7265	34.3406	1.6141	0.4582
wage_st	4.6240	4.8480	0.2240	0.0000
hrsopen	14.5892	14.4191	-0.1701	0.4917

`*eval(parse(text="nome_variavel"))` foi utilizado para transformar um texto em um nome de objeto. Assim, podemos “chamar” um objeto e, neste caso, usá-lo dentro da fórmula do modelo na função `lm()`.

Note que, para a maioria das variáveis analisadas, as diferenças não são estatisticamente significantes a 5%. Há diferença apenas para o salário por hora (4,8% superior em New Jersey) e na quantidade de empregados por loja (Pennsylvania possui 3 funcionários a mais em média), o que pode indicar alguma fragilidade na hipótese de que as tendências de comportamento dos dois estados seguiriam semelhantes na ausência de tratamento. \square

Exercício 4. *Nos EUA, discute-se bastante a questão da discriminação dos povos hispânicos (mexicanos, porto riquenhos, cubanos, etc.). Analisando **apenas trabalhadores em tempo integral e parcial**, verificaremos a possível discriminação na renda usando a base de dados General Social Survey (GSS), a mesma utilizada no Exercício 2 e que pode ser carregada no R usando o comando:*

```
1 load(url("http://bit.ly/dasi_gss_data"))
```

- a) *Regrida a renda (`coninc`) em relação às dummies dos povos hispânicos e liste os que possuem diferença significativa de renda em relação aos não-hispânicos. Utilize como variáveis de controle: o sexo, a idade, a idade², a raça, os anos de estudo (`educ`) e o status de trabalho (`wrkstat`).*
- b) *Você acha razoável considerar que o resultado do item (a) como representativo de toda população trabalhadora nos EUA? Justifique resumindo/analizando os dados.*
- Dica:** *Verifique se todas observações de gss foram utilizadas na regressão no item (a).*

Resposta:

- a) Mexicanos a 0,1%, Filipinos a 5%, Argentinos e Porto Riquenhos a 10%:

```
1 library(dplyr)
2 load(url("http://bit.ly/dasi_gss_data"))
3
4 gss = gss %>%
5   # filtrando apenas trabalhadores em tempo integral e parcial
6   filter(wrkstat %in% c("Working Fulltime", "Working Parttime")) %>%
7
8   # definindo Not Hispanic como referencia
9   mutate(hispanic = relevel(hispanic, ref="Not Hispanic"))
10
11 # Regressão
12 fit = lm(coninc ~ hispanic + sex + age + I(age^2) + race + educ + wrkstat, gss)
13 summary(fit)
```

```
1 Call:
2 lm(formula = coninc ~ hispanic + sex + age + I(age^2) + race +
3 educ + wrkstat, data = gss)
4
5 Residuals:
6 Min      1Q  Median      3Q      Max
7 -102538  -25548   -8266   14571  163415
8
9 Coefficients:
```



```
2 gss %>% select(coninc, hispanic, sex, age, race, educ, wrkstat) %>%
3   is.na() %>% apply(2, sum)
```

```
1   coninc  hispanic  sex   age   race  educ  wrkstat
2     2600    22434    0   109     0    75      0
```

e, por isso, as proporções dos povos hispânicos alterou muito pouco entre as bases:

```
1 # Comparando hispanic na base cheia e na usada em regressão
2 rbind(prop.table(table(gss$hispanic)), prop.table(table(gss_reg$hispanic))) %>%
   round(4) %>% t()
```

	[,1]	[,2]
2 Not Hispanic	0.8830	0.8844
3 Mexican, Mexican American, Chicano/A	0.0710	0.0712
4 Puerto Rican	0.0145	0.0145
5 Cuban	0.0046	0.0043
6 Salvadorian	0.0040	0.0040
7 Guatemalan	0.0024	0.0023
8 Panamanian	0.0009	0.0009
9 Nicaraguan	0.0007	0.0008
10 Costa Rican	0.0003	0.0003
11 Central American	0.0012	0.0011
12 Honduran	0.0016	0.0014
13 Dominican	0.0022	0.0020
14 West Indian	0.0001	0.0001
15 Peruvian	0.0011	0.0012
16 Ecuadorian	0.0016	0.0012
17 Columbian	0.0016	0.0013
18 Venezuelan	0.0007	0.0008
19 Argentinian	0.0002	0.0002
20 Chilean	0.0000	0.0000
21 Spanish	0.0047	0.0049
22 Basque	0.0001	0.0001
23 Filipino/A	0.0003	0.0004
24 Latin American	0.0003	0.0003
25 South American	0.0009	0.0009
26 Latin	0.0002	0.0000
27 Latino/A	0.0007	0.0006
28 Hispanic	0.0004	0.0003
29 Other, Not Specified	0.0007	0.0007

Em gss, criaremos a variável binária que indica se a observação (idcase) está presente em gss_reg, ou seja, se possui algum missing value nas variáveis utilizadas na regressão. Também, incluiremos dummies de mulher, raças e trabalho em tempo integral.

```
1 # criando variáveis dummies
2 gss = gss %>% mutate(
3   dmy_reg = ifelse(caseid %in% gss_reg$caseid, 1, 0),
4   dmy_female = ifelse(sex=="Female", 1, 0),
5   dmy_race_white = ifelse(race=="White", 1, 0),
6   dmy_race_black = ifelse(race=="Black", 1, 0),
7   dmy_race_other = ifelse(race=="Other", 1, 0),
8   dmy_fulltime = ifelse(wrkstat=="Working Fulltime", 1, 0)
9 )
```

Agora, calcularemos as diferenças de média das variáveis de renda, sexo, idade, raça, anos de educação e status de trabalho:


```

1 # diferenças de médias
2 lista_var = c("coninc", "dmy_female", "age", "dmy_race_white",
3               "dmy_race_black", "educ", "dmy_fulltime")
4
5 resultados = matrix(0, length(lista_var), 4) %>%
6 as.data.frame()
7 colnames(resultados) = c("base completa", "base reg", "dif", "p valor")
8 rownames(resultados) = lista_var
9
10 for (var in lista_var) {
11   reg = lm(eval(parse(text=var)) ~ dmy_reg, gss) %>% summary() %>% coef()
12   resultados[var, 1] = reg["(Intercept)", "Estimate"]
13   resultados[var, 2] = reg["(Intercept)", "Estimate"] + reg["dmy_reg", "
14     Estimate"]
15   resultados[var, 3] = reg["dmy_reg", "Estimate"]
16   resultados[var, 4] = reg["dmy_reg", "Pr(>|t|)"]
17 }
18 resultados %>% round(4)

```

	base completa	base reg	dif	p valor
coninc	48686.4779	56103.1536	7416.6757	0.0000
dmy_female	0.4752	0.4938	0.0186	0.0015
age	39.8771	41.7068	1.8297	0.0000
dmy_race_white	0.8342	0.7646	-0.0696	0.0000
dmy_race_black	0.1299	0.1373	0.0073	0.0663
educ	13.2103	13.8810	0.6707	0.0000
dmy_fulltime	0.8279	0.8297	0.0018	0.6893

Note que apenas as proporções de negros e de trabalhadores em tempo integral são estatisticamente iguais a 5%. Já as diferenças de idade, de proporção de mulheres e de anos de educação, embora sejam significativas, não são tão diferentes em números absolutos.

□