

Regression Project – Group 3 (Intermediate Report)

Reema and Muzammil

Table of contents

| | | |
|----------|--|-----------|
| 1 | 1. Introduction | 2 |
| 2 | 2. Data collection and description | 2 |
| 3 | 3. Load packages | 3 |
| 4 | 4. Load data and select variables | 5 |
| 5 | 5. Data preparation | 8 |
| 6 | 6. Descriptive statistics | 11 |
| 6.1 | 6.1 Numeric summaries | 11 |
| 6.2 | 6.2 Frequency tables | 11 |
| 7 | 7. Univariate visualizations | 12 |
| 7.1 | 7.1 Employment income (py010n) | 12 |
| 7.2 | 7.2 Age | 13 |
| 7.3 | 7.3 Household size (hsize) | 15 |
| 7.4 | 7.4 Gender | 17 |
| 7.5 | 7.5 Citizenship | 18 |
| 8 | 8. Bivariate plots (predictors vs response) | 19 |
| 8.1 | 8.1 Gender vs income | 19 |
| 8.2 | 8.2 Citizenship vs income | 20 |
| 8.3 | 8.3 Age vs income | 21 |
| 8.4 | 8.4 Household size vs income | 22 |
| 9 | 9. Interaction plots | 23 |
| 9.1 | 9.1 Gender \times Citizenship | 23 |

| | | |
|-----------|---|-----------|
| 9.2 | 9.2 Age × Gender | 24 |
| 9.3 | 9.3 Age × Citizenship | 25 |
| 9.4 | 9.4 Household Size × Gender | 26 |
| 9.5 | 9.5 Household Size × Citizenship | 27 |
| 10 | 10. Contingency tables (categorical × categorical) | 29 |
| 10.1 | 10.1 Extra: 2D Density Contours | 29 |
| 11 | 11. Summary | 30 |
| 12 | END OF DOCUMENT | 30 |

1 1. Introduction

- The aim of this intermediate report is to explore determinants of income in South Austria.
 - We focus on data management and descriptive statistics before conducting regression modelling.
 - The response variable is net employment income (`py010n`).
 - Explanatory variables include gender, citizenship, household size, and age.
 - For the intermediate report we restrict to data management and descriptive statistics.
-

2 2. Data collection and description

- Source: EU-SILC (European Union Statistics on Income and Living Conditions), Austria.
- Type of data: survey data, representative sample of private households.
- Data format: cross-sectional microdata with social, demographic, and income information.
- Variables used:
 - `py010n` — employment income (numeric)
 - `age` — age in years (numeric)
 - `hsize` — household size (categorical converted to numeric)
 - `gender` — male / female (categorical)
 - `citizenship` — grouped nationality categories (categorical)
 - `region` — Austrian federal region
- Missing value handling:
 - Keep only positive values of income (`py010n > 0`)

- Convert `hsize` to numeric
 - Remove observations with missing values
 - Subsetting:
 - Only individuals living in **Styria** and **Carinthia** (South Austria, NUTS-1 region)
-

3 3. Load packages

```
options(repos = c(CRAN = "https://cloud.r-project.org"))
install.packages("simFrame")
```

The downloaded binary packages are in
/var/folders/15/4lrzl3dn20j0wzcv5vd_8b4c0000gn/T//RtmpirnSGM/downloaded_packages

```
install.packages("dplyr")
```

The downloaded binary packages are in
/var/folders/15/4lrzl3dn20j0wzcv5vd_8b4c0000gn/T//RtmpirnSGM/downloaded_packages

```
install.packages("ggplot2")
```

The downloaded binary packages are in
/var/folders/15/4lrzl3dn20j0wzcv5vd_8b4c0000gn/T//RtmpirnSGM/downloaded_packages

```
install.packages("tidyr")
```

The downloaded binary packages are in
/var/folders/15/4lrzl3dn20j0wzcv5vd_8b4c0000gn/T//RtmpirnSGM/downloaded_packages

```
install.packages("forcats")
```

The downloaded binary packages are in
/var/folders/15/4lrz13dn20j0wzcv5vd_8b4c0000gn/T//RtmpirnSGM/downloaded_packages

```
install.packages("effects")
```

The downloaded binary packages are in
/var/folders/15/4lrz13dn20j0wzcv5vd_8b4c0000gn/T//RtmpirnSGM/downloaded_packages

```
library(simFrame)
```

Loading required package: Rcpp

Loading required package: lattice

Loading required package: parallel

```
library(dplyr)      # data manipulation
```

Attaching package: 'dplyr'

The following objects are masked from 'package:stats':

filter, lag

The following objects are masked from 'package:base':

intersect, setdiff, setequal, union

```
library(ggplot2)    # visualization
library(tidyr)      # data tidying
library(forcats)    # factor management
library(effects)    # effect plots
```

Loading required package: carData

Use the command

```
lattice::trellis.par.set(effectsTheme())
```

to customize lattice options for effects plots.

See ?effectsTheme for details.

4 4. Load data and select variables

```
data(eusilcP)
dat = eusilcP
str(eusilcP)
```

```
'data.frame':  58654 obs. of  28 variables:
 $ hid      : int  1 1 2 2 3 4 4 4 5 6 ...
 $ region   : Ord.factor w/ 9 levels "Burgenland"<"Lower Austria"<...: 6 6 5 5 5 6 6 6 3 2
 $ hsize    : Factor w/ 9 levels "1","2","3","4",...: 2 2 2 2 1 3 3 3 1 5 ...
 $ eqsize   : num  1.5 1.5 1.5 1.5 1 1.8 1.8 1.8 1 2.6 ...
 $ eqIncome : num [1:58654(1d)] 11128 11128 19695 19695 5066 ...
 ..- attr(*, "dimnames")=List of 1
 .. ..$ : chr [1:58654] "2592313" "2592313" "2045000" "2045000" ...
 $ pid      : int  1 2 1 2 1 1 2 3 1 1 ...
 $ id       : chr  "0000101" "0000102" "0000201" "0000202" ...
 $ age      : num  25 24 57 53 30 32 33 8 77 34 ...
 $ gender    : Factor w/ 2 levels "male","female": 1 2 2 1 2 1 2 1 2 2 ...
 $ ecoStat   : Factor w/ 7 levels "1","2","3","4",...: 1 4 1 1 6 1 1 NA 5 2 ...
 $ citizenship: Factor w/ 3 levels "AT","EU","Other": 3 1 1 1 1 1 1 NA 1 1 ...
 $ py010n    : num  16693 0 0 16884 0 ...
 $ py050n    : num  0 0 12565 0 0 ...
 $ py090n    : num  0 0 0 0 5066 ...
 $ py100n    : num  0 0 0 0 0 ...
 $ py110n    : num  0 0 0 0 0 0 0 NA 0 0 ...
 $ py120n    : num  0 0 0 0 0 0 0 NA 0 0 ...
 $ py130n    : num  0 0 0 0 0 0 0 NA 0 0 ...
 $ py140n    : num  0 0 0 0 0 0 0 NA 0 0 ...
 $ hy040n    : num  0 0 0 0 0 0 0 0 0 0 ...
 $ hy050n    : num  0 0 0 0 0 ...
```

```

$ hy070n      : num  0 0 0 0 0 0 0 0 0 0 ...
$ hy080n      : num  0 0 0 0 0 0 0 0 0 0 ...
$ hy090n      : num  0 0 0 0 0 ...
$ hy110n      : num  0 0 0 0 0 ...
$ hy130n      : num  0 0 93.6 93.6 0 ...
$ hy145n      : num  0 0 -187 -187 0 ...
$ main        : logi  TRUE FALSE FALSE TRUE TRUE TRUE ...

```

```
head(eusilcP)
```

| | hid | region | hsize | eqsize | eqIncome | pid | id | age | gender | ecoStat |
|-------|-------------|---------------|----------|---------|----------|---------|---------|---------|--------|---------|
| 39993 | 1 | Upper Austria | 2 | 1.5 | 11128.45 | 1 | 0000101 | 25 | male | 1 |
| 39994 | 1 | Upper Austria | 2 | 1.5 | 11128.45 | 2 | 0000102 | 24 | female | 4 |
| 31004 | 2 | Styria | 2 | 1.5 | 19694.85 | 1 | 0000201 | 57 | female | 1 |
| 31005 | 2 | Styria | 2 | 1.5 | 19694.85 | 2 | 0000202 | 53 | male | 1 |
| 29071 | 3 | Styria | 1 | 1.0 | 5066.24 | 1 | 0000301 | 30 | female | 6 |
| 41322 | 4 | Upper Austria | 3 | 1.8 | 31480.01 | 1 | 0000401 | 32 | male | 1 |
| | citizenship | py010n | py050n | py090n | py100n | py110n | py120n | py130n | py140n | |
| 39993 | Other | 16692.67 | 0.00 | 0.00 | 0 | 0 | 0 | 0 | 0 | |
| 39994 | AT | 0.00 | 0.00 | 0.00 | 0 | 0 | 0 | 0 | 0 | |
| 31004 | AT | 0.00 | 12564.59 | 0.00 | 0 | 0 | 0 | 0 | 0 | |
| 31005 | AT | 16884.06 | 0.00 | 0.00 | 0 | 0 | 0 | 0 | 0 | |
| 29071 | AT | 0.00 | 0.00 | 5066.24 | 0 | 0 | 0 | 0 | 0 | |
| 41322 | AT | 25047.39 | 0.00 | 0.00 | 0 | 0 | 0 | 0 | 0 | |
| | hy040n | hy050n | hy070n | hy080n | hy090n | hy110n | hy130n | hy145n | main | |
| 39993 | 0 | 0.00 | 0 | 0 | 0.00 | 0.00 | 0.00 | 0.00 | TRUE | |
| 39994 | 0 | 0.00 | 0 | 0 | 0.00 | 0.00 | 0.00 | 0.00 | FALSE | |
| 31004 | 0 | 0.00 | 0 | 0 | 0.00 | 0.00 | 93.63 | -187.26 | FALSE | |
| 31005 | 0 | 0.00 | 0 | 0 | 0.00 | 0.00 | 93.63 | -187.26 | TRUE | |
| 29071 | 0 | 0.00 | 0 | 0 | 0.00 | 0.00 | 0.00 | 0.00 | TRUE | |
| 41322 | 0 | 7167.39 | 0 | 0 | 31.15 | 1349.91 | 0.00 | 0.00 | TRUE | |

```
summary(eusilcP)
```

| | hid | region | hsize | eqsize |
|----------|--------|----------------|--------|--------|
| Min. | : 1 | Vienna | :11657 | 2 |
| 1st Qu.: | 6262 | Lower Austria: | 11127 | 4 |
| Median | :12465 | Upper Austria: | 10310 | 3 |
| Mean | :12488 | Styria | : 8142 | 1 |
| 3rd Qu.: | 18719 | Tyrol | : 4796 | 5 |
| Max. | :25000 | Carinthia | : 4111 | 6 |

| | hsize | eqsize |
|----------|--------|---------------|
| Min. | :14128 | Min. :1.000 |
| 1st Qu.: | :13180 | 1st Qu.:1.500 |
| Median | :12429 | Median :2.000 |
| Mean | : 8602 | Mean :1.943 |
| 3rd Qu.: | : 6745 | 3rd Qu.:2.400 |
| Max. | : 2094 | Max. :4.500 |

| | | | |
|----------------|--------------|------------------------------|---------------|
| | | (Other) : 8511 (Other): 1476 | |
| eqIncome | pid | id | age |
| Min. : 0 | Min. :1.00 | Length:58654 | Min. : -1.00 |
| 1st Qu.: 13539 | 1st Qu.:1.00 | Class :character | 1st Qu.:22.00 |
| Median : 18322 | Median :2.00 | Mode :character | Median :40.00 |
| Mean : 20163 | Mean :2.07 | | Mean :39.75 |
| 3rd Qu.: 24277 | 3rd Qu.:3.00 | | 3rd Qu.:57.00 |
| Max. :179946 | Max. :9.00 | | Max. :97.00 |

| | | | | |
|--------------|---------------|-------------|----------------|--------------|
| gender | ecoStat | citizenship | py010n | py050n |
| male :28539 | 1 :20900 | AT :44066 | Min. : 0 | Min. : -6895 |
| female:30115 | 5 :12836 | EU : 1257 | 1st Qu.: 0 | 1st Qu.: 0 |
| | 7 : 4607 | Other: 3162 | Median : 2382 | Median : 0 |
| | 2 : 4362 | NA's :10169 | Mean : 9062 | Mean : 1288 |
| | 4 : 2921 | | 3rd Qu.: 16820 | 3rd Qu.: 0 |
| | (Other): 2859 | | Max. :199075 | Max. :129874 |
| | NA's :10169 | | NA's :10169 | NA's :10169 |

| | | | |
|---------------|--------------|---------------|----------------|
| py090n | py100n | py110n | py120n |
| Min. : 0.0 | Min. : 0 | Min. : 0.0 | Min. : 0.00 |
| 1st Qu.: 0.0 | 1st Qu.: 0 | 1st Qu.: 0.0 | 1st Qu.: 0.00 |
| Median : 0.0 | Median : 0 | Median : 0.0 | Median : 0.00 |
| Mean : 444.6 | Mean : 3713 | Mean : 72.9 | Mean : 51.22 |
| 3rd Qu.: 0.0 | 3rd Qu.: 0 | 3rd Qu.: 0.0 | 3rd Qu.: 0.00 |
| Max. :29887.1 | Max. :101777 | Max. :22546.8 | Max. :46398.44 |
| NA's :10169 | NA's :10169 | NA's :10169 | NA's :10169 |

| | | | |
|---------------|----------------|----------------|---------------|
| py130n | py140n | hy040n | hy050n |
| Min. : 0.0 | Min. : 0.00 | Min. : -2962.5 | Min. : -11857 |
| 1st Qu.: 0.0 | 1st Qu.: 0.00 | 1st Qu.: 0.0 | 1st Qu.: 0 |
| Median : 0.0 | Median : 0.00 | Median : 0.0 | Median : 0 |
| Mean : 393.7 | Mean : 41.73 | Mean : 879.9 | Mean : 2826 |
| 3rd Qu.: 0.0 | 3rd Qu.: 0.00 | 3rd Qu.: 0.0 | 3rd Qu.: 4558 |
| Max. :53183.6 | Max. :18643.46 | Max. :129586.6 | Max. :118309 |
| NA's :10169 | NA's :10169 | | |

| | | | |
|----------------|----------------|-----------------|----------------|
| hy070n | hy080n | hy090n | hy110n |
| Min. : 0.00 | Min. : 0.0 | Min. : -457.46 | Min. : 0.00 |
| 1st Qu.: 0.00 | 1st Qu.: 0.0 | 1st Qu.: 0.75 | 1st Qu.: 0.00 |
| Median : 0.00 | Median : 0.0 | Median : 58.45 | Median : 0.00 |
| Mean : 93.12 | Mean : 744.6 | Mean : 462.45 | Mean : 32.97 |
| 3rd Qu.: 0.00 | 3rd Qu.: 0.0 | 3rd Qu.: 234.78 | 3rd Qu.: 0.00 |
| Max. :17954.97 | Max. :124206.2 | Max. :112011.03 | Max. :14506.49 |

| | | |
|--------------|-----------------|---------------|
| hy130n | hy145n | main |
| Min. : -5490 | Min. : -29519.3 | Mode :logical |

| | | | | |
|----------|-------|----------|---------|-------------|
| 1st Qu.: | 0 | 1st Qu.: | -256.8 | FALSE:33654 |
| Median : | 0 | Median : | 0.0 | TRUE :25000 |
| Mean : | 339 | Mean : | -108.8 | |
| 3rd Qu.: | 0 | 3rd Qu.: | 0.0 | |
| Max. : | 40763 | Max. : | 49768.0 | |

5 5. Data preparation

- Filter for South Austria regions (Styria and Carinthia)
- Keep only positive income observations
- Convert household size to numeric
- Remove missing values
- Group citizenship categories if needed

```
dat <- eusilcP %>%
  select(py010n, gender, citizenship, hsize, age, region) %>%
  filter(region %in% c("Carinthia", "Styria")) %>%
  filter(py010n > 0) %>%
  na.omit()

dat$gender <- as.factor(dat$gender)
dat$citizenship <- as.factor(dat$citizenship)
dat$hsize <- as.factor(dat$hsize)

model_int <- lm(py010n ~ gender * citizenship + hsize + age, data = dat)
anova(model_int)
```

Analysis of Variance Table

Response: py010n

| | Df | Sum Sq | Mean Sq | F value | Pr(>F) |
|--------------------|------|------------|------------|----------|---------------|
| gender | 1 | 7.5470e+10 | 7.5470e+10 | 833.1618 | < 2.2e-16 *** |
| citizenship | 2 | 5.0537e+09 | 2.5268e+09 | 27.8954 | 8.894e-13 *** |
| hsize | 8 | 6.7445e+09 | 8.4306e+08 | 9.3071 | 7.907e-13 *** |
| age | 1 | 1.4708e+10 | 1.4708e+10 | 162.3697 | < 2.2e-16 *** |
| gender:citizenship | 2 | 3.0402e+07 | 1.5201e+07 | 0.1678 | 0.8455 |
| Residuals | 5252 | 4.7574e+11 | 9.0583e+07 | | |

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```
summary(model_int)
```

Call:

```
lm(formula = py010n ~ gender * citizenship + hsize + age, data = dat)
```

Residuals:

| Min | 1Q | Median | 3Q | Max |
|--------|-------|--------|------|-------|
| -23632 | -5930 | -1065 | 4652 | 95265 |

Coefficients:

| | Estimate | Std. Error | t value | Pr(> t) | |
|-------------------------------|----------|------------|---------|----------|-----|
| (Intercept) | 15907.30 | 589.59 | 26.980 | < 2e-16 | *** |
| genderfemale | -7711.77 | 271.54 | -28.400 | < 2e-16 | *** |
| citizenshipEU | 1487.22 | 1520.96 | 0.978 | 0.3282 | |
| citizenshipOther | -5343.40 | 987.60 | -5.410 | 6.56e-08 | *** |
| hsize2 | -220.35 | 461.65 | -0.477 | 0.6332 | |
| hsize3 | -606.41 | 445.42 | -1.361 | 0.1734 | |
| hsize4 | -1072.25 | 468.18 | -2.290 | 0.0220 | * |
| hsize5 | -809.71 | 567.70 | -1.426 | 0.1538 | |
| hsize6 | -3380.42 | 686.65 | -4.923 | 8.78e-07 | *** |
| hsize7 | -2815.78 | 1425.09 | -1.976 | 0.0482 | * |
| hsize8 | -4156.36 | 1633.33 | -2.545 | 0.0110 | * |
| hsize9 | -1552.22 | 1650.84 | -0.940 | 0.3471 | |
| age | 132.13 | 10.37 | 12.741 | < 2e-16 | *** |
| genderfemale:citizenshipEU | 161.21 | 2161.44 | 0.075 | 0.9405 | |
| genderfemale:citizenshipOther | 874.02 | 1517.33 | 0.576 | 0.5646 | |

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

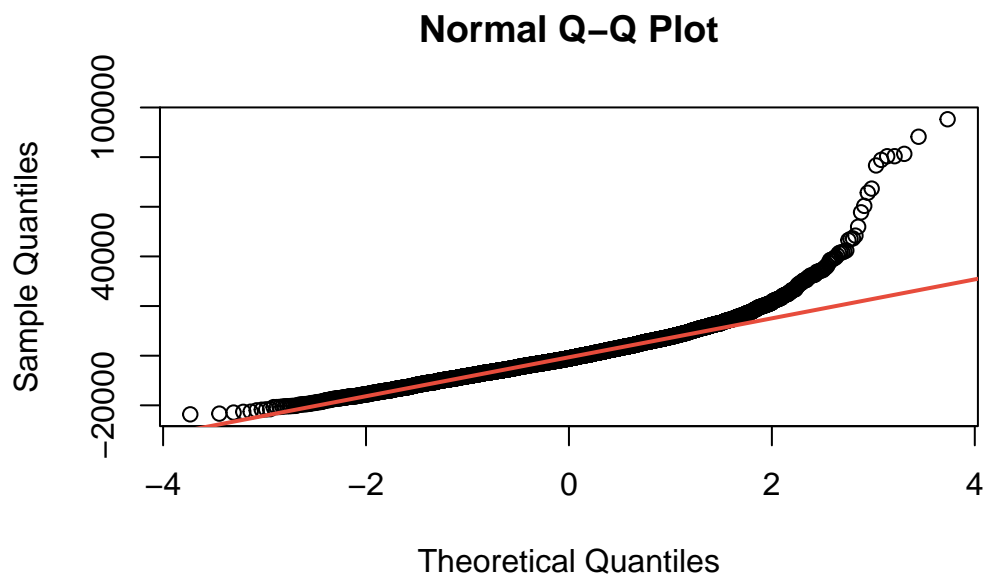
Residual standard error: 9517 on 5252 degrees of freedom

Multiple R-squared: 0.1766, Adjusted R-squared: 0.1744

F-statistic: 80.44 on 14 and 5252 DF, p-value: < 2.2e-16

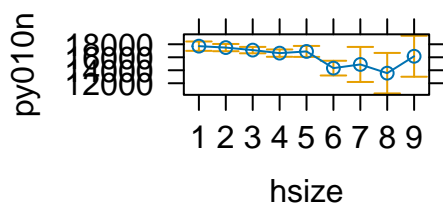
```
qqnorm(residuals(model_int))
```

```
qqline(residuals(model_int), col = "#E74C3C", lwd = 2)
```

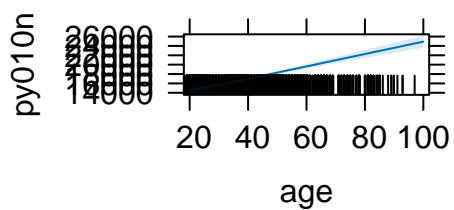


```
plot(allEffects(model_int))
```

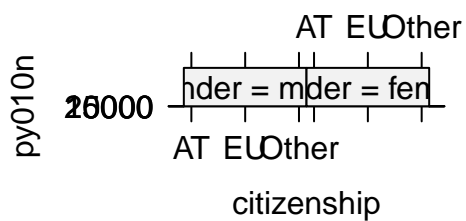
hsize effect plot



age effect plot



gender*citizenship effect plot



6 6. Descriptive statistics

6.1 6.1 Numeric summaries

```
summary(dat$py010n)
```

| Min. | 1st Qu. | Median | Mean | 3rd Qu. | Max. |
|------|----------|----------|----------|----------|-----------|
| 1.93 | 10066.01 | 16225.84 | 16952.35 | 21939.78 | 118362.27 |

```
summary(dat$age)
```

| Min. | 1st Qu. | Median | Mean | 3rd Qu. | Max. |
|-------|---------|--------|-------|---------|-------|
| 16.00 | 29.00 | 40.00 | 39.73 | 49.00 | 97.00 |

```
summary(dat$hsize)
```

| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|-----|------|------|------|-----|-----|----|----|----|
| 679 | 1140 | 1454 | 1105 | 496 | 274 | 48 | 36 | 35 |

6.2 6.2 Frequency tables

```
table(dat$gender)
```

| male | female |
|------|--------|
| 3004 | 2263 |

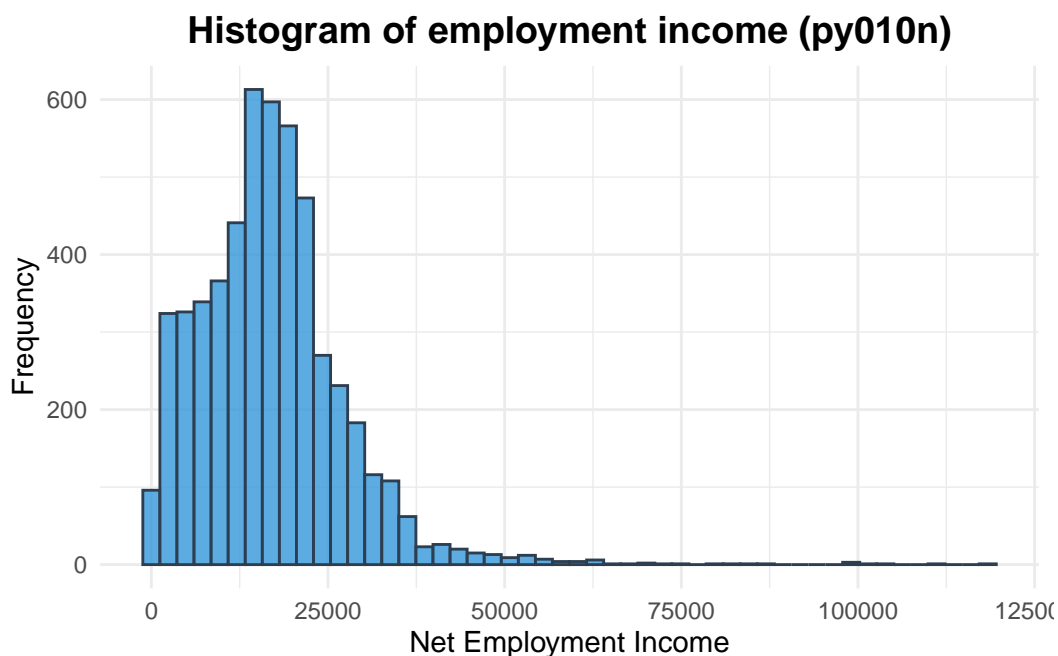
```
table(dat$citizenship)
```

| AT | EU | Other |
|------|----|-------|
| 5021 | 79 | 167 |

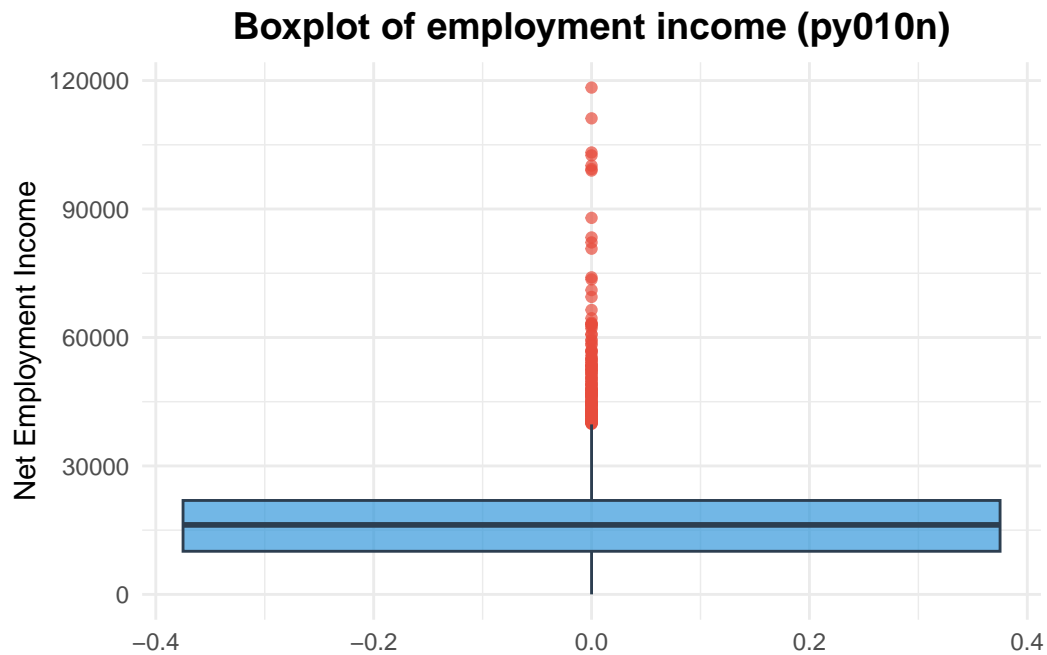
7 7. Univariate visualizations

7.1 7.1 Employment income (py010n)

```
ggplot(dat, aes(x = py010n)) +  
  geom_histogram(bins = 50, fill = "#3498DB", color = "#2C3E50", alpha = 0.8) +  
  labs(title = "Histogram of employment income (py010n)",  
       x = "Net Employment Income",  
       y = "Frequency") +  
  theme_minimal() +  
  theme(plot.title = element_text(hjust = 0.5, face = "bold", size = 14))
```

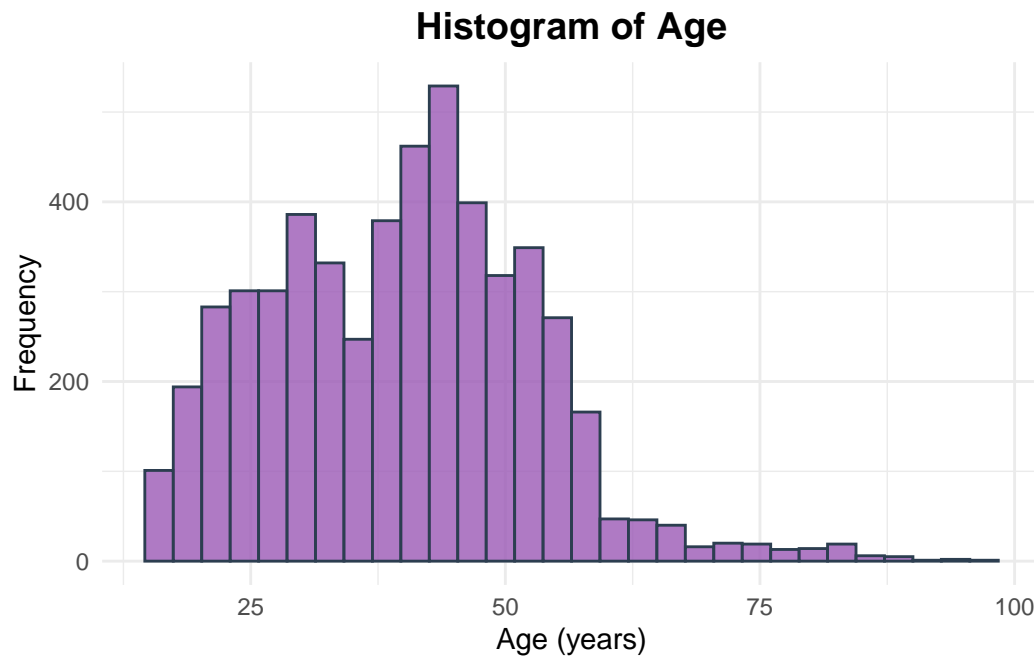


```
ggplot(dat, aes(y = py010n)) +  
  geom_boxplot(fill = "#3498DB", color = "#2C3E50", alpha = 0.7, outlier.color = "#E74C3C") +  
  labs(title = "Boxplot of employment income (py010n)",  
       y = "Net Employment Income") +  
  theme_minimal() +  
  theme(plot.title = element_text(hjust = 0.5, face = "bold", size = 14))
```

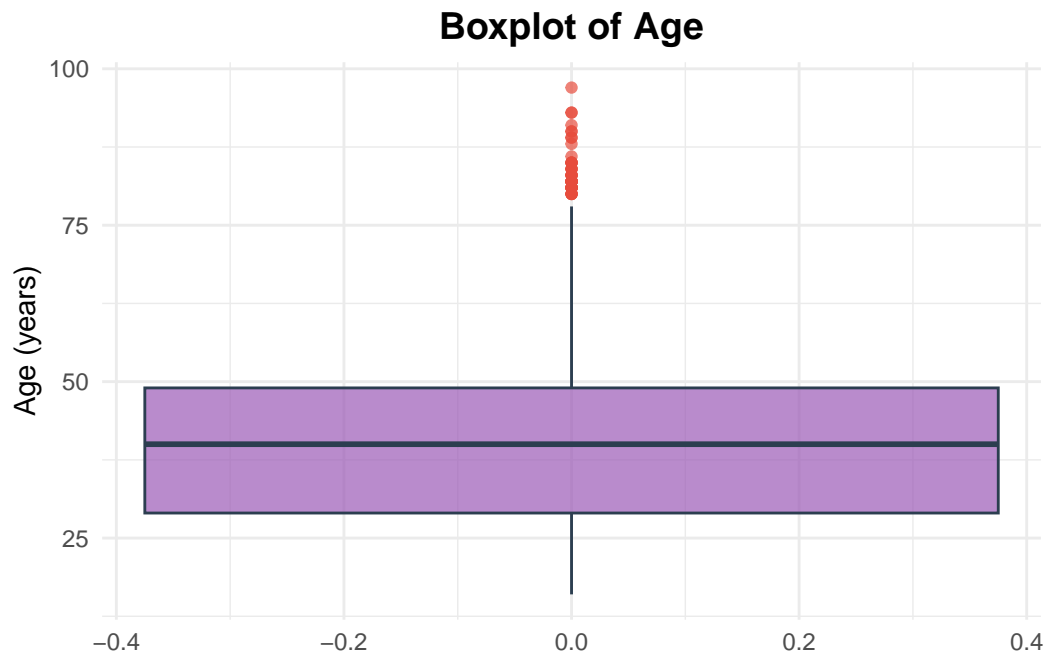


7.2 7.2 Age

```
ggplot(dat, aes(x = age)) +
  geom_histogram(bins = 30, fill = "#9B59B6", color = "#2C3E50", alpha = 0.8) +
  labs(title = "Histogram of Age",
       x = "Age (years)",
       y = "Frequency") +
  theme_minimal() +
  theme(plot.title = element_text(hjust = 0.5, face = "bold", size = 14))
```

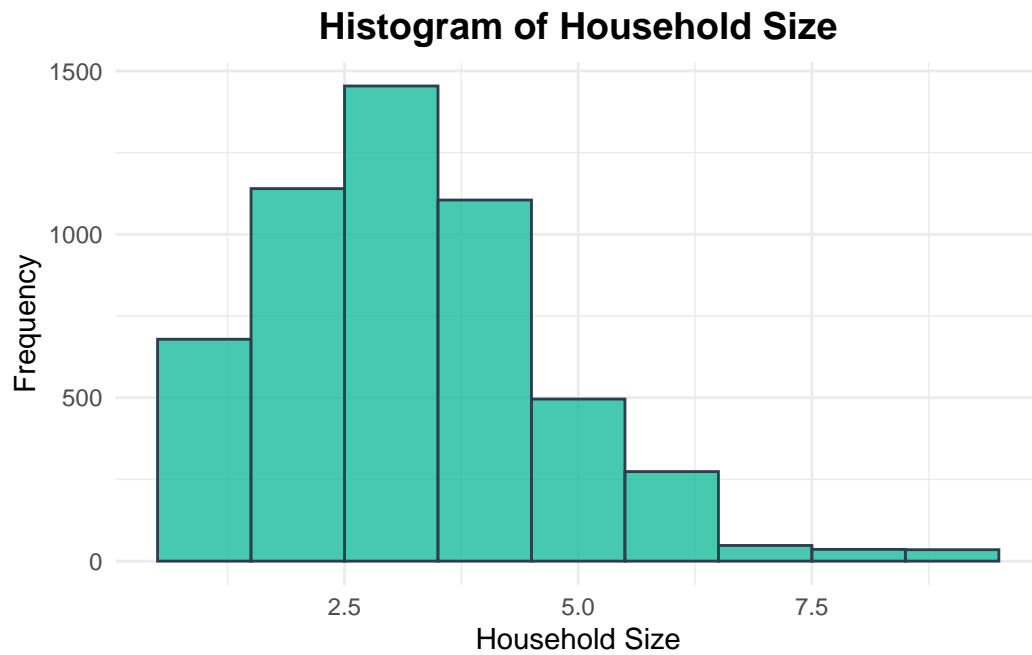


```
ggplot(dat, aes(y = age)) +  
  geom_boxplot(fill = "#9B59B6", color = "#2C3E50", alpha = 0.7, outlier.color = "#E74C3C") +  
  labs(title = "Boxplot of Age",  
        y = "Age (years)") +  
  theme_minimal() +  
  theme(plot.title = element_text(hjust = 0.5, face = "bold", size = 14))
```

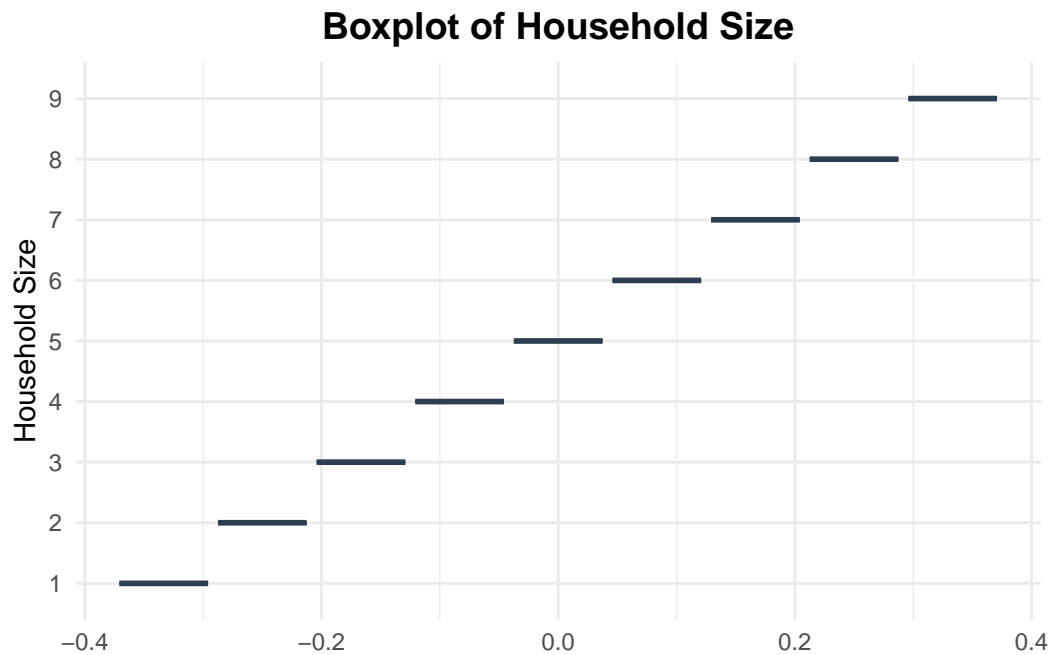


7.3 7.3 Household size (hsize)

```
ggplot(dat, aes(x = as.numeric(hsize))) +
  geom_histogram(binwidth = 1, fill = "#1ABC9C", color = "#2C3E50", alpha = 0.8) +
  labs(title = "Histogram of Household Size",
       x = "Household Size",
       y = "Frequency") +
  theme_minimal() +
  theme(plot.title = element_text(hjust = 0.5, face = "bold", size = 14))
```

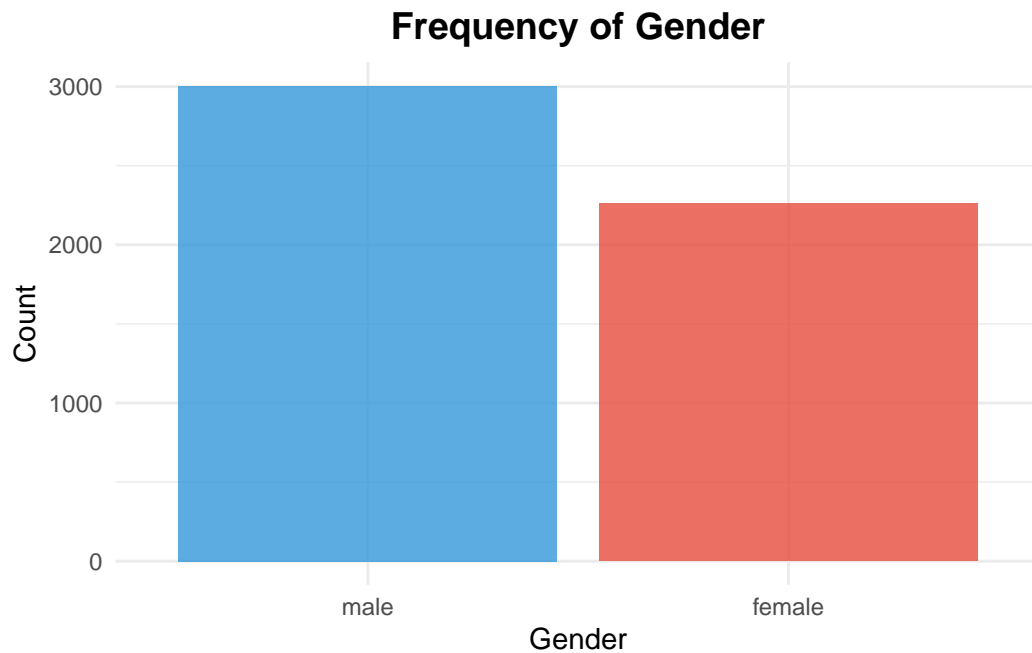


```
ggplot(dat, aes(y = hsize)) +  
  geom_boxplot(fill = "#1ABC9C", color = "#2C3E50", alpha = 0.7, outlier.color = "#E74C3C") +  
  labs(title = "Boxplot of Household Size",  
        y = "Household Size") +  
  theme_minimal() +  
  theme(plot.title = element_text(hjust = 0.5, face = "bold", size = 14))
```

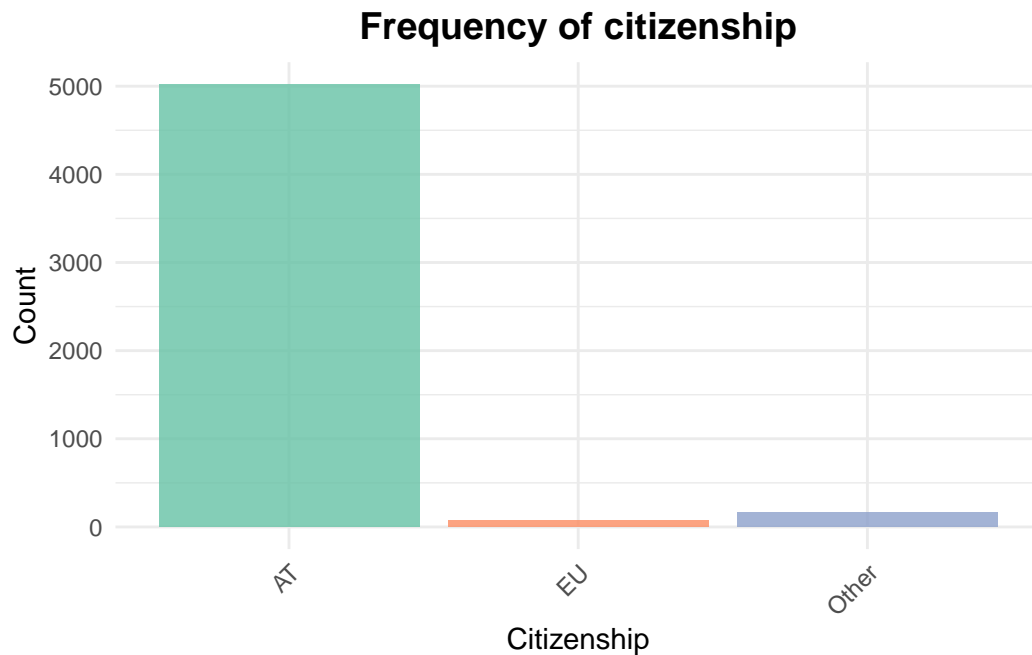
7.4 7.4 Gender

```
ggplot(dat, aes(x = gender, fill = gender)) +
  geom_bar(alpha = 0.8) +
  scale_fill_manual(values = c("male" = "#3498DB", "female" = "#E74C3C")) +
  labs(title = "Frequency of Gender",
       x = "Gender",
       y = "Count") +
  theme_minimal() +
  theme(plot.title = element_text(hjust = 0.5, face = "bold", size = 14),
        legend.position = "none")
```



7.5 7.5 Citizenship

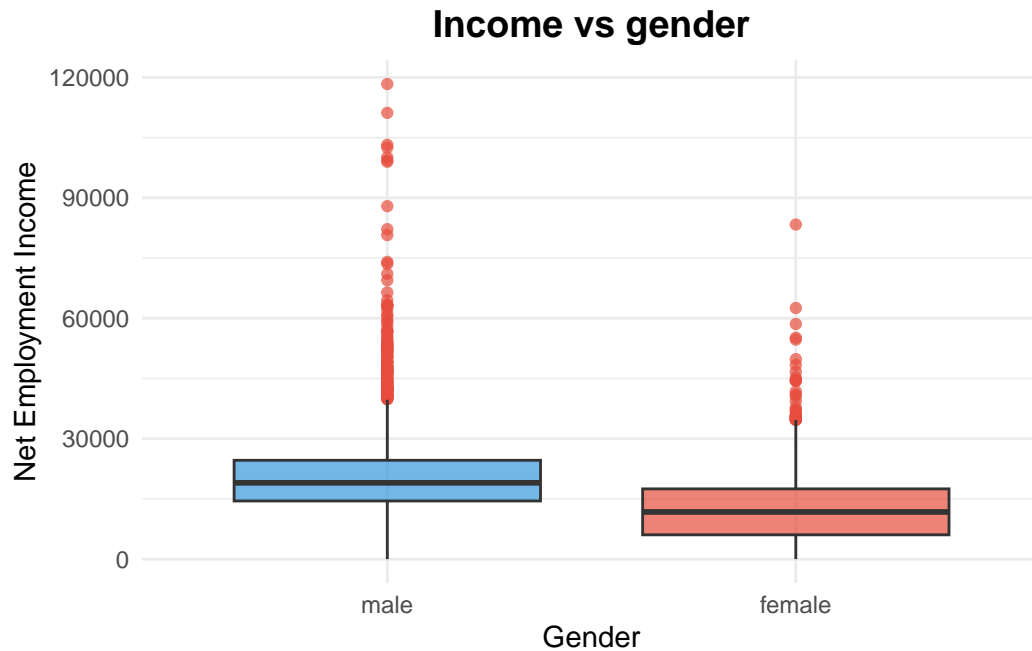
```
ggplot(dat, aes(x = citizenship, fill = citizenship)) +  
  geom_bar(alpha = 0.8) +  
  scale_fill_brewer(palette = "Set2") +  
  theme_minimal() +  
  theme(axis.text.x = element_text(angle = 45, hjust = 1),  
        plot.title = element_text(hjust = 0.5, face = "bold", size = 14),  
        legend.position = "none") +  
  labs(title = "Frequency of citizenship",  
       x = "Citizenship",  
       y = "Count")
```



8 8. Bivariate plots (predictors vs response)

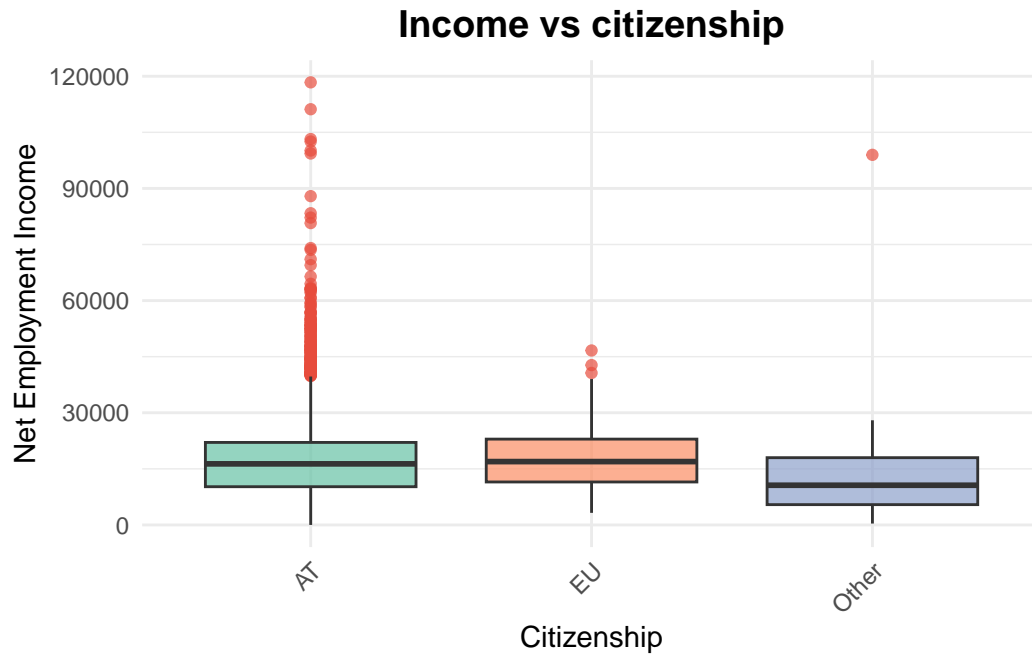
8.1 8.1 Gender vs income

```
ggplot(dat, aes(x = gender, y = py010n, fill = gender)) +  
  geom_boxplot(alpha = 0.7, outlier.color = "#E74C3C") +  
  scale_fill_manual(values = c("male" = "#3498DB", "female" = "#E74C3C")) +  
  labs(title = "Income vs gender",  
       x = "Gender",  
       y = "Net Employment Income") +  
  theme_minimal() +  
  theme(plot.title = element_text(hjust = 0.5, face = "bold", size = 14),  
        legend.position = "none")
```



8.2 8.2 Citizenship vs income

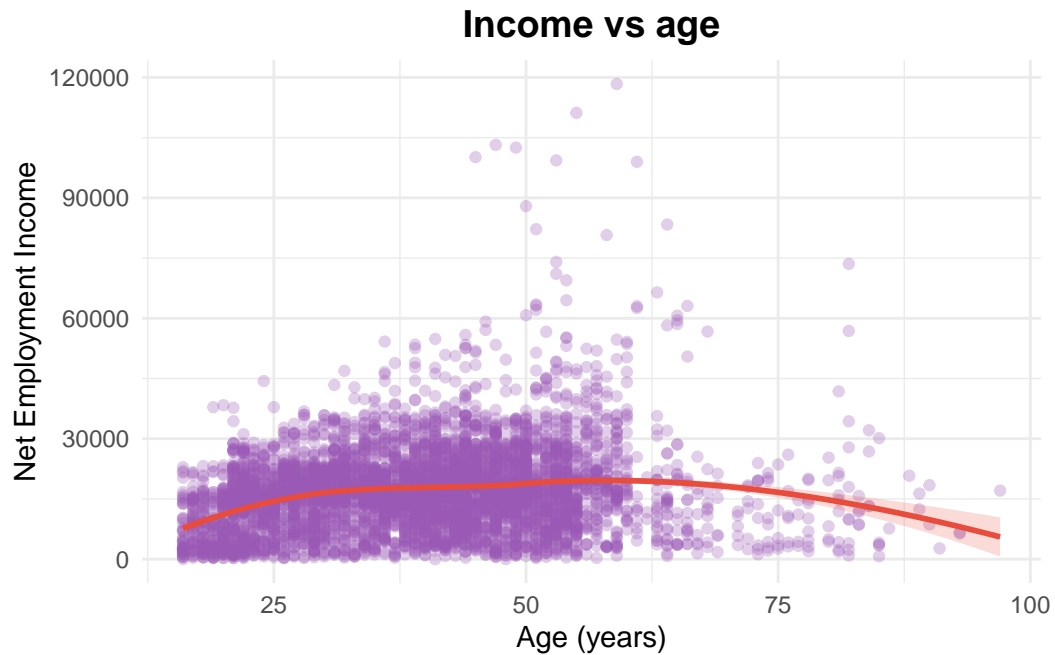
```
ggplot(dat, aes(x = citizenship, y = py010n, fill = citizenship)) +
  geom_boxplot(alpha = 0.7, outlier.color = "#E74C3C") +
  scale_fill_brewer(palette = "Set2") +
  theme_minimal() +
  theme(axis.text.x = element_text(angle = 45, hjust = 1),
        plot.title = element_text(hjust = 0.5, face = "bold", size = 14),
        legend.position = "none") +
  labs(title = "Income vs citizenship",
       x = "Citizenship",
       y = "Net Employment Income")
```



8.3 8.3 Age vs income

```
ggplot(dat, aes(x = age, y = py010n)) +
  geom_point(alpha = 0.3, color = "#9B59B6") +
  geom_smooth(method = "loess", color = "#E74C3C", fill = "#E74C3C", alpha = 0.2) +
  labs(title = "Income vs age",
       x = "Age (years)",
       y = "Net Employment Income") +
  theme_minimal() +
  theme(plot.title = element_text(hjust = 0.5, face = "bold", size = 14))
```

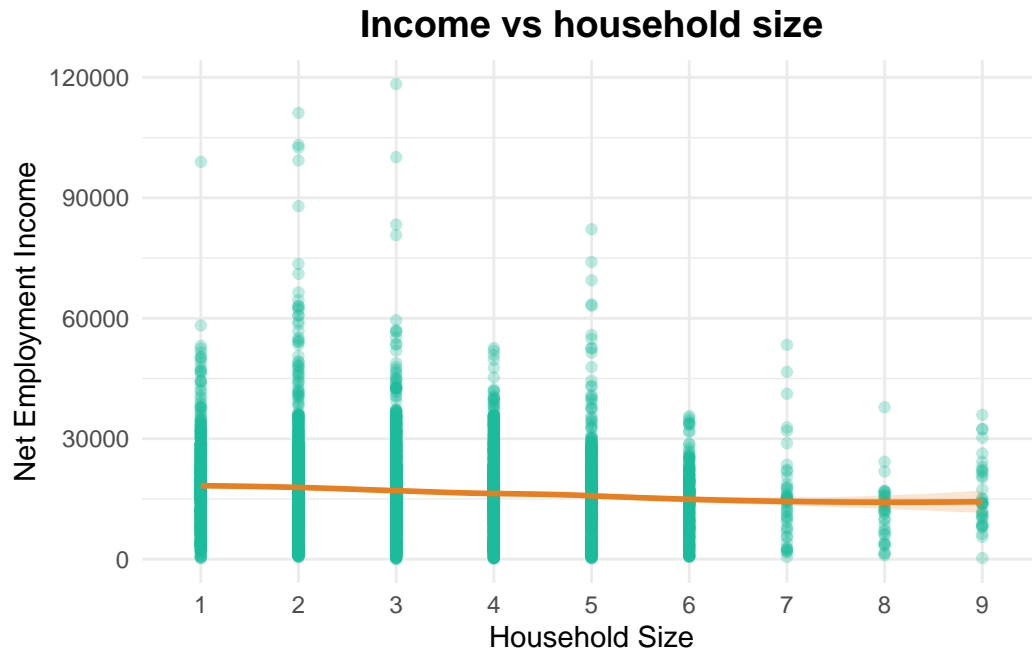
`geom_smooth()` using formula = 'y ~ x'



8.4 8.4 Household size vs income

```
ggplot(dat, aes(x = hsize, y = py010n, group = 1)) +
  geom_point(alpha = 0.3, color = "#1ABC9C") +
  geom_smooth(method = "loess", color = "#E67E22", fill = "#E67E22", alpha = 0.2) +
  labs(title = "Income vs household size",
       x = "Household Size",
       y = "Net Employment Income") +
  theme_minimal() +
  theme(plot.title = element_text(hjust = 0.5, face = "bold", size = 14))
```

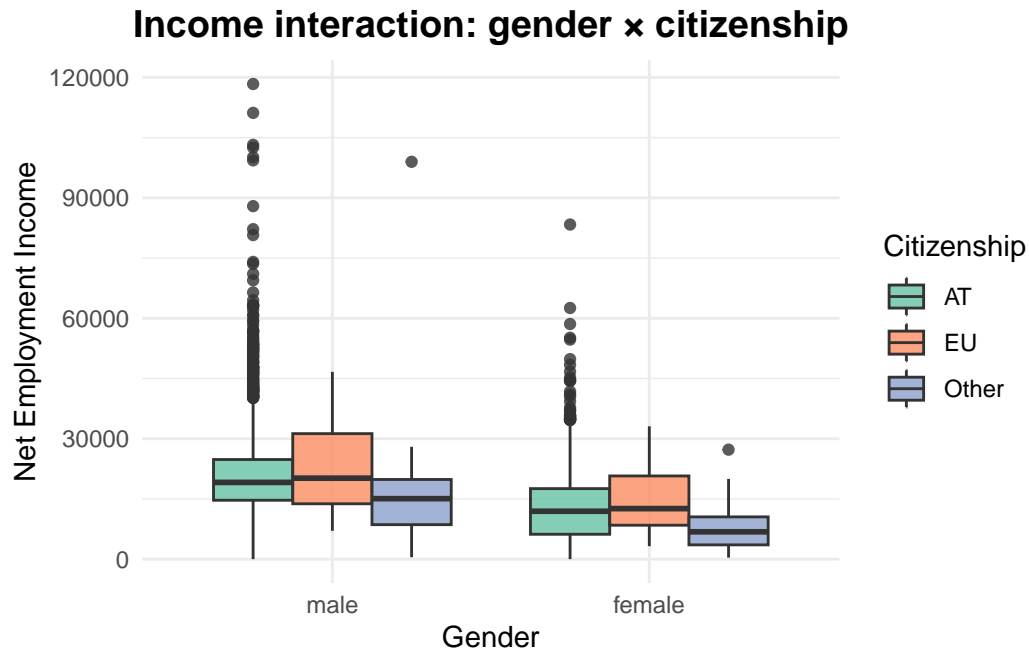
`geom_smooth()` using formula = 'y ~ x'



9 9. Interaction plots

9.1 9.1 Gender × Citizenship

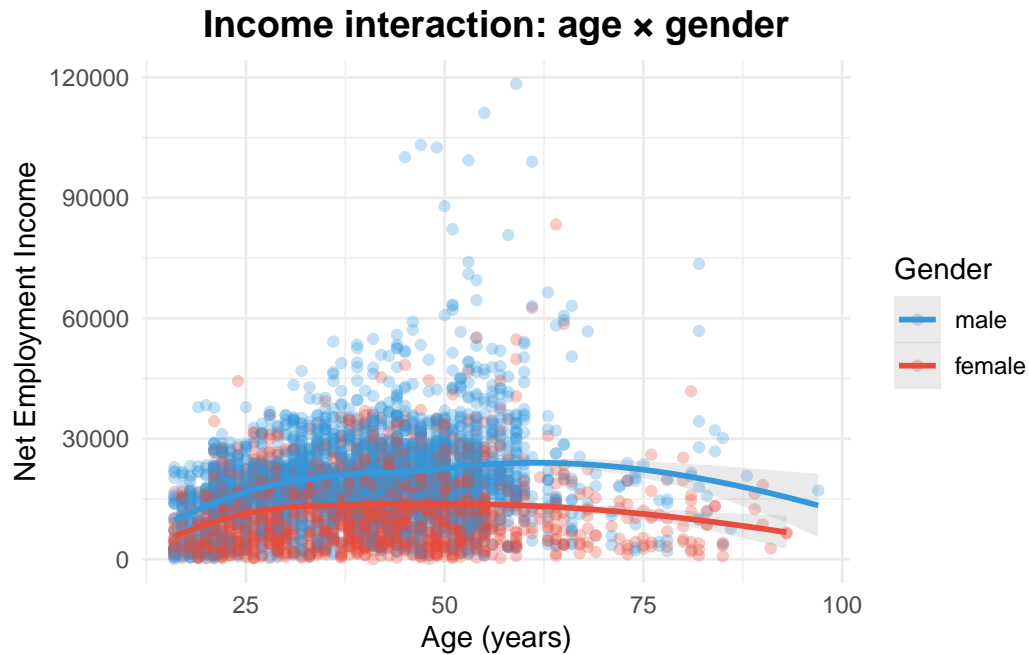
```
ggplot(dat, aes(x = gender, y = py010n, fill = citizenship)) +
  geom_boxplot(position = "dodge", alpha = 0.8) +
  scale_fill_brewer(palette = "Set2") +
  labs(title = "Income interaction: gender × citizenship",
       x = "Gender",
       y = "Net Employment Income",
       fill = "Citizenship") +
  theme_minimal() +
  theme(plot.title = element_text(hjust = 0.5, face = "bold", size = 14))
```



9.2 9.2 Age × Gender

```
ggplot(dat, aes(x = age, y = py010n, color = gender)) +
  geom_point(alpha = 0.3) +
  geom_smooth(method = "loess", se = TRUE, alpha = 0.2) +
  scale_color_manual(values = c("male" = "#3498DB", "female" = "#E74C3C")) +
  labs(title = "Income interaction: age × gender",
       x = "Age (years)",
       y = "Net Employment Income",
       color = "Gender") +
  theme_minimal() +
  theme(plot.title = element_text(hjust = 0.5, face = "bold", size = 14))
```

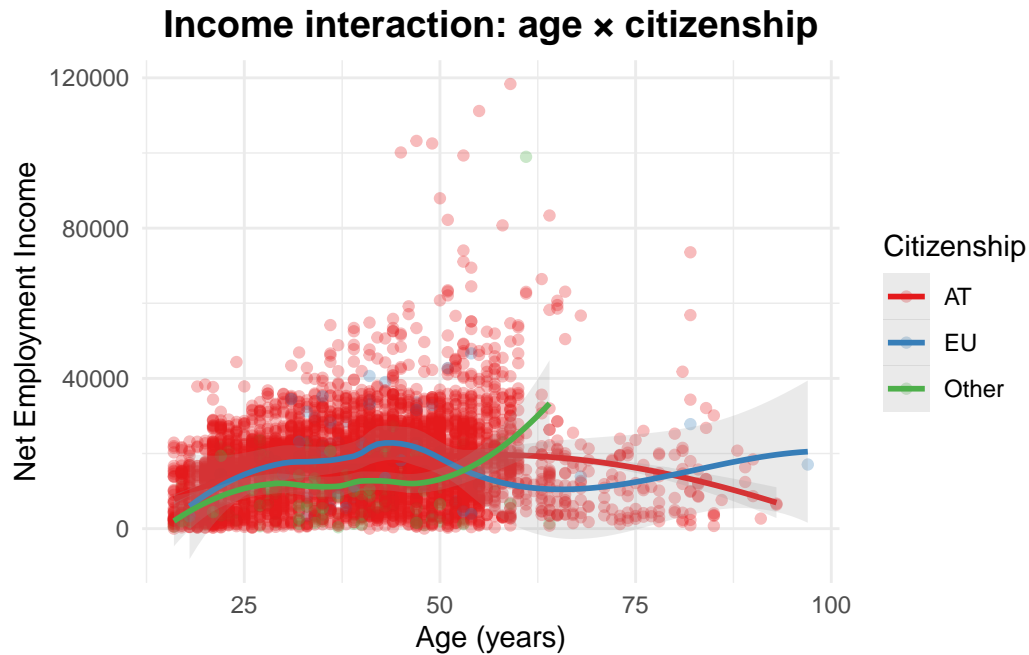
`geom_smooth()` using formula = 'y ~ x'



9.3 9.3 Age × Citizenship

```
ggplot(dat, aes(x = age, y = py010n, color = citizenship)) +
  geom_point(alpha = 0.3) +
  geom_smooth(method = "loess", se = TRUE, alpha = 0.2) +
  scale_color_brewer(palette = "Set1") +
  labs(title = "Income interaction: age × citizenship",
       x = "Age (years)",
       y = "Net Employment Income",
       color = "Citizenship") +
  theme_minimal() +
  theme(plot.title = element_text(hjust = 0.5, face = "bold", size = 14))
```

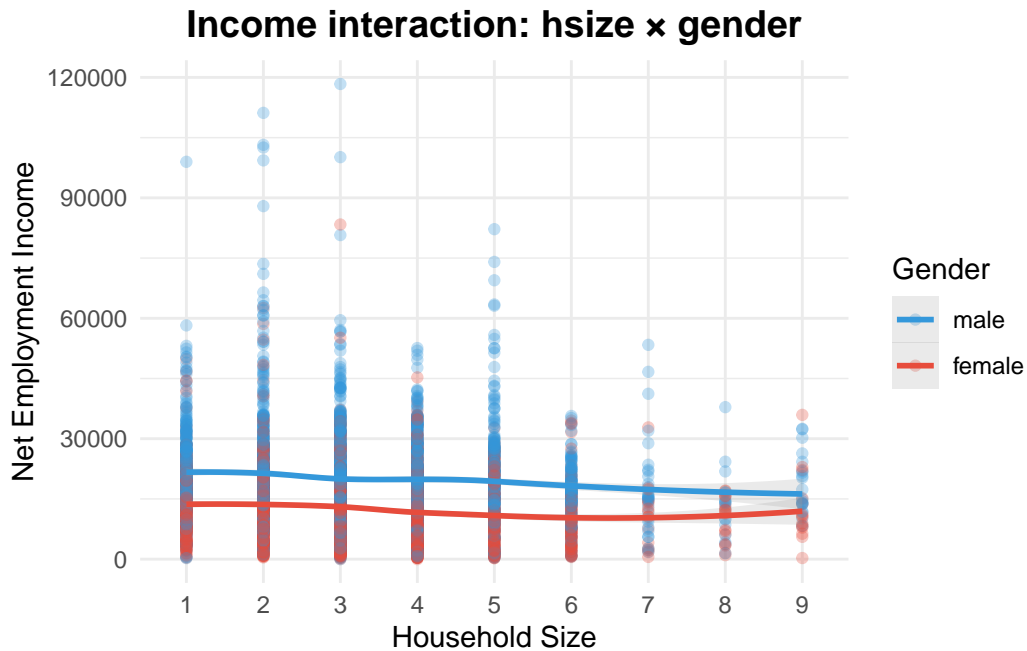
`geom_smooth()` using formula = 'y ~ x'



9.4 9.4 Household Size × Gender

```
ggplot(dat, aes(x = hsize, y = py010n, color = gender, group = gender)) +
  geom_point(alpha = 0.3) +
  geom_smooth(method = "loess", se = TRUE, alpha = 0.2) +
  scale_color_manual(values = c("male" = "#3498DB", "female" = "#E74C3C")) +
  labs(title = "Income interaction: hsize × gender",
       x = "Household Size",
       y = "Net Employment Income",
       color = "Gender") +
  theme_minimal() +
  theme(plot.title = element_text(hjust = 0.5, face = "bold", size = 14))
```

`geom_smooth()` using formula = 'y ~ x'



9.5 9.5 Household Size × Citizenship

```
ggplot(dat, aes(x = hsize, y = py010n, color = citizenship, group = citizenship)) +
  geom_point(alpha = 0.3) +
  geom_smooth(method = "loess", se = TRUE, alpha = 0.2) +
  scale_color_brewer(palette = "Set1") +
  labs(title = "Income interaction: hsize x citizenship",
       x = "Household Size",
       y = "Net Employment Income",
       color = "Citizenship") +
  theme_minimal() +
  theme(plot.title = element_text(hjust = 0.5, face = "bold", size = 14))
```

`geom_smooth()` using formula = 'y ~ x'

Warning in simpleLoess(y, x, w, span, degree = degree, parametric = parametric,
: pseudoinverse used at 3

Warning in simpleLoess(y, x, w, span, degree = degree, parametric = parametric,
: neighborhood radius 1

```
Warning in simpleLoess(y, x, w, span, degree = degree, parametric = parametric,  
: reciprocal condition number 0
```

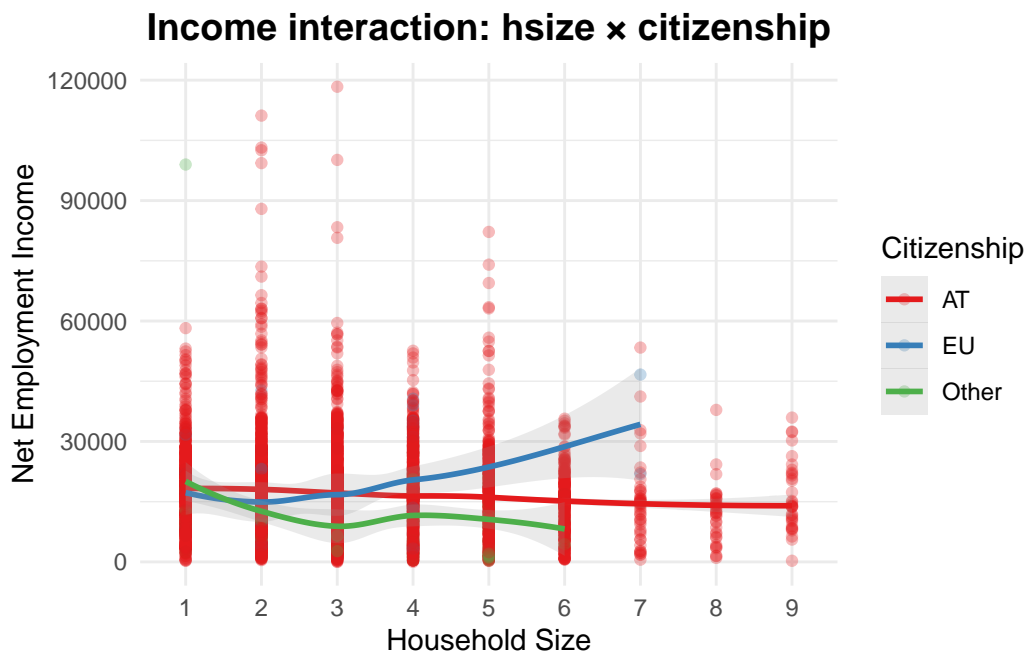
```
Warning in simpleLoess(y, x, w, span, degree = degree, parametric = parametric,  
: There are other near singularities as well. 4
```

```
Warning in predLoess(object$y, object$x, newx = if (is.null(newdata)) object$x  
else if (is.data.frame(newdata))  
as.matrix(model.frame(delete.response(terms(object))), : pseudoinverse used at 3
```

```
Warning in predLoess(object$y, object$x, newx = if (is.null(newdata)) object$x  
else if (is.data.frame(newdata))  
as.matrix(model.frame(delete.response(terms(object))), : neighborhood radius 1
```

```
Warning in predLoess(object$y, object$x, newx = if (is.null(newdata)) object$x  
else if (is.data.frame(newdata))  
as.matrix(model.frame(delete.response(terms(object))), : reciprocal condition  
number 0
```

```
Warning in predLoess(object$y, object$x, newx = if (is.null(newdata)) object$x  
else if (is.data.frame(newdata))  
as.matrix(model.frame(delete.response(terms(object))), : There are other near  
singularities as well. 4
```



10 10. Contingency tables (categorical × categorical)

```
table(dat$gender, dat$citizenship)
```

| | AT | EU | Other |
|--------|------|----|-------|
| male | 2867 | 40 | 97 |
| female | 2154 | 39 | 70 |

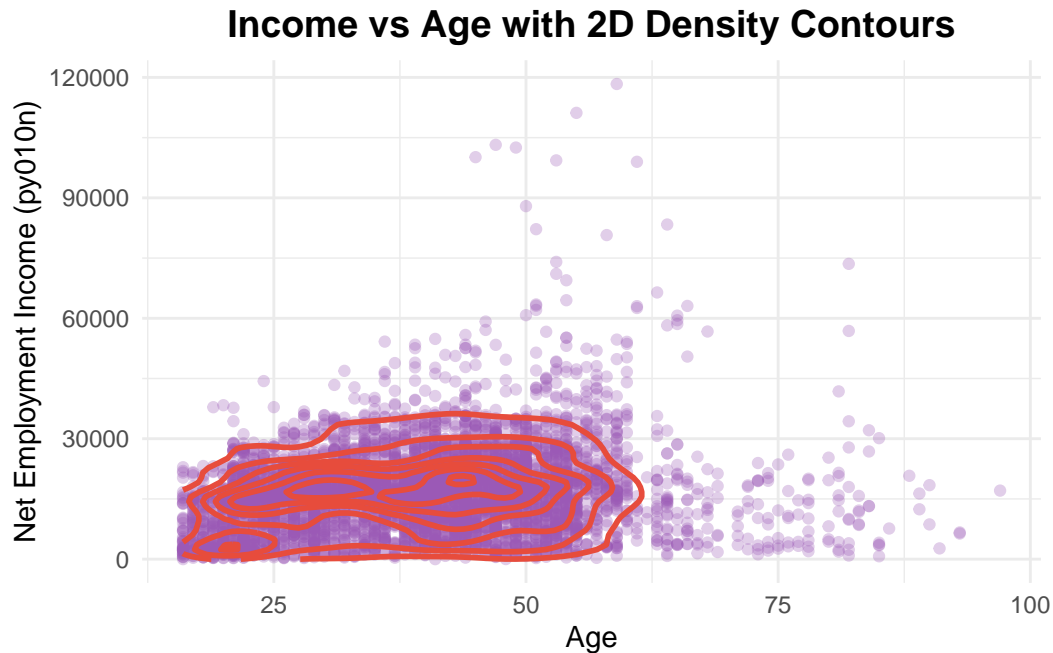
- If categories are too detailed, merge small ones:

```
dat$citizenship <- fct_lump(dat$citizenship, n = 3)
table(dat$gender, dat$citizenship)
```

| | AT | EU | Other |
|--------|------|----|-------|
| male | 2867 | 40 | 97 |
| female | 2154 | 39 | 70 |

10.1 Extra: 2D Density Contours

```
ggplot(dat, aes(x = age, y = py010n)) +
  geom_point(alpha = 0.3, color = "#9B59B6") +
  geom_density_2d(color = "#E74C3C", linewidth = 1) +
  labs(title = "Income vs Age with 2D Density Contours",
       x = "Age",
       y = "Net Employment Income (py010n)") +
  theme_minimal() +
  theme(plot.title = element_text(hjust = 0.5, face = "bold", size = 14))
```



11 11. Summary

- The income variable is highly right-skewed with outliers.
 - Men typically have higher median employment income than women.
 - Citizenship differences may indicate structural inequality in wages.
 - Age and income show a nonlinear increasing pattern.
 - Larger households do not clearly correlate with higher or lower income.
 - Some interaction effects appear visible (gender \times citizenship, etc.).
 - **Interaction insights:** Compares income distributions across citizenship groups, separately for men and women, highlighting potential interaction between gender and citizenship.
-

12 END OF DOCUMENT