

# Ejercicio 3.3

Silvia Pineda

## Carga de Datos y Librerías

```
library(naniar)
library(tidyverse)

-- Attaching core tidyverse packages ----- tidyverse 2.0.0 --
v dplyr     1.1.4     v readr     2.1.5
v forcats   1.0.1     v stringr   1.5.2
v ggplot2   4.0.0     v tibble    3.3.0
v lubridate 1.9.4     v tidyr    1.3.1
v purrr    1.1.0

-- Conflicts -----
x dplyr::filter() masks stats::filter()
x dplyr::lag()    masks stats::lag()
i Use the conflicted package (<http://conflicted.r-lib.org/>) to force all conflicts to beco
```

```
library(rio)

data <- read.csv("students_FP.csv",
  na.strings = c("", "NA", "NaN", "NULL"),
  stringsAsFactors = TRUE
)
```

## Imputación múltiple

```
library(mice)
```

```
Attaching package: 'mice'
```

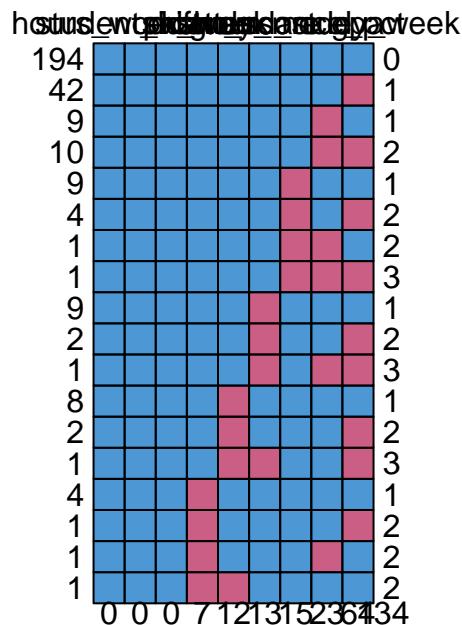
```
The following object is masked from 'package:stats':
```

```
filter
```

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

```
cbind, rbind
```

```
md.pattern(data)
```



|     | student_id | hours_work_week | shift | program | exam_score | study_mode |
|-----|------------|-----------------|-------|---------|------------|------------|
| 194 | 1          |                 | 1     | 1       | 1          | 1          |
| 42  | 1          |                 | 1     | 1       | 1          | 1          |
| 9   | 1          |                 | 1     | 1       | 1          | 1          |
| 10  | 1          |                 | 1     | 1       | 1          | 1          |
| 9   | 1          |                 | 1     | 1       | 1          | 1          |
| 4   | 1          |                 | 1     | 1       | 1          | 1          |
| 1   | 1          |                 | 1     | 1       | 1          | 1          |

|     |                |                  |     |     |    |   |
|-----|----------------|------------------|-----|-----|----|---|
| 1   | 1              | 1                | 1   | 1   | 1  | 1 |
| 9   | 1              | 1                | 1   | 1   | 1  | 0 |
| 2   | 1              | 1                | 1   | 1   | 1  | 0 |
| 1   | 1              | 1                | 1   | 1   | 1  | 0 |
| 8   | 1              | 1                | 1   | 1   | 0  | 1 |
| 2   | 1              | 1                | 1   | 1   | 0  | 1 |
| 1   | 1              | 1                | 1   | 1   | 0  | 0 |
| 4   | 1              | 1                | 1   | 0   | 1  | 1 |
| 1   | 1              | 1                | 1   | 0   | 1  | 1 |
| 1   | 1              | 1                | 1   | 0   | 1  | 1 |
| 1   | 1              | 1                | 1   | 0   | 0  | 1 |
| 0   | 0              | 0                | 7   | 12  | 13 |   |
|     | attendance_pct | hours_study_week | gpa |     |    |   |
| 194 | 1              | 1                | 1   | 0   |    |   |
| 42  | 1              | 1                | 0   | 1   |    |   |
| 9   | 1              | 0                | 1   | 1   |    |   |
| 10  | 1              | 0                | 0   | 2   |    |   |
| 9   | 0              | 1                | 1   | 1   |    |   |
| 4   | 0              | 1                | 0   | 2   |    |   |
| 1   | 0              | 0                | 1   | 2   |    |   |
| 1   | 0              | 0                | 0   | 3   |    |   |
| 9   | 1              | 1                | 1   | 1   |    |   |
| 2   | 1              | 1                | 0   | 2   |    |   |
| 1   | 1              | 0                | 0   | 3   |    |   |
| 8   | 1              | 1                | 1   | 1   |    |   |
| 2   | 1              | 1                | 0   | 2   |    |   |
| 1   | 1              | 1                | 0   | 3   |    |   |
| 4   | 1              | 1                | 1   | 1   |    |   |
| 1   | 1              | 1                | 0   | 2   |    |   |
| 1   | 1              | 0                | 1   | 2   |    |   |
| 1   | 1              | 1                | 1   | 2   |    |   |
| 15  |                | 23               | 64  | 134 |    |   |

No hay ninguna observación con todas las variables missing.

```
impData <- mice(select(data,-student_id),m=5,maxit=50,seed=500)
```

| iter | imp | variable         |                |     |            |         |            |  |
|------|-----|------------------|----------------|-----|------------|---------|------------|--|
| 1    | 1   | hours_study_week | attendance_pct | gpa | exam_score | program | study_mode |  |
| 1    | 2   | hours_study_week | attendance_pct | gpa | exam_score | program | study_mode |  |
| 1    | 3   | hours_study_week | attendance_pct | gpa | exam_score | program | study_mode |  |











```

44  4 hours_study_week attendance_pct gpa exam_score program study_mode
44  5 hours_study_week attendance_pct gpa exam_score program study_mode
45  1 hours_study_week attendance_pct gpa exam_score program study_mode
45  2 hours_study_week attendance_pct gpa exam_score program study_mode
45  3 hours_study_week attendance_pct gpa exam_score program study_mode
45  4 hours_study_week attendance_pct gpa exam_score program study_mode
45  5 hours_study_week attendance_pct gpa exam_score program study_mode
46  1 hours_study_week attendance_pct gpa exam_score program study_mode
46  2 hours_study_week attendance_pct gpa exam_score program study_mode
46  3 hours_study_week attendance_pct gpa exam_score program study_mode
46  4 hours_study_week attendance_pct gpa exam_score program study_mode
46  5 hours_study_week attendance_pct gpa exam_score program study_mode
47  1 hours_study_week attendance_pct gpa exam_score program study_mode
47  2 hours_study_week attendance_pct gpa exam_score program study_mode
47  3 hours_study_week attendance_pct gpa exam_score program study_mode
47  4 hours_study_week attendance_pct gpa exam_score program study_mode
47  5 hours_study_week attendance_pct gpa exam_score program study_mode
48  1 hours_study_week attendance_pct gpa exam_score program study_mode
48  2 hours_study_week attendance_pct gpa exam_score program study_mode
48  3 hours_study_week attendance_pct gpa exam_score program study_mode
48  4 hours_study_week attendance_pct gpa exam_score program study_mode
48  5 hours_study_week attendance_pct gpa exam_score program study_mode
49  1 hours_study_week attendance_pct gpa exam_score program study_mode
49  2 hours_study_week attendance_pct gpa exam_score program study_mode
49  3 hours_study_week attendance_pct gpa exam_score program study_mode
49  4 hours_study_week attendance_pct gpa exam_score program study_mode
49  5 hours_study_week attendance_pct gpa exam_score program study_mode
50  1 hours_study_week attendance_pct gpa exam_score program study_mode
50  2 hours_study_week attendance_pct gpa exam_score program study_mode
50  3 hours_study_week attendance_pct gpa exam_score program study_mode
50  4 hours_study_week attendance_pct gpa exam_score program study_mode
50  5 hours_study_week attendance_pct gpa exam_score program study_mode

```

```
summary(impData)
```

Class: mids

Number of multiple imputations: 5

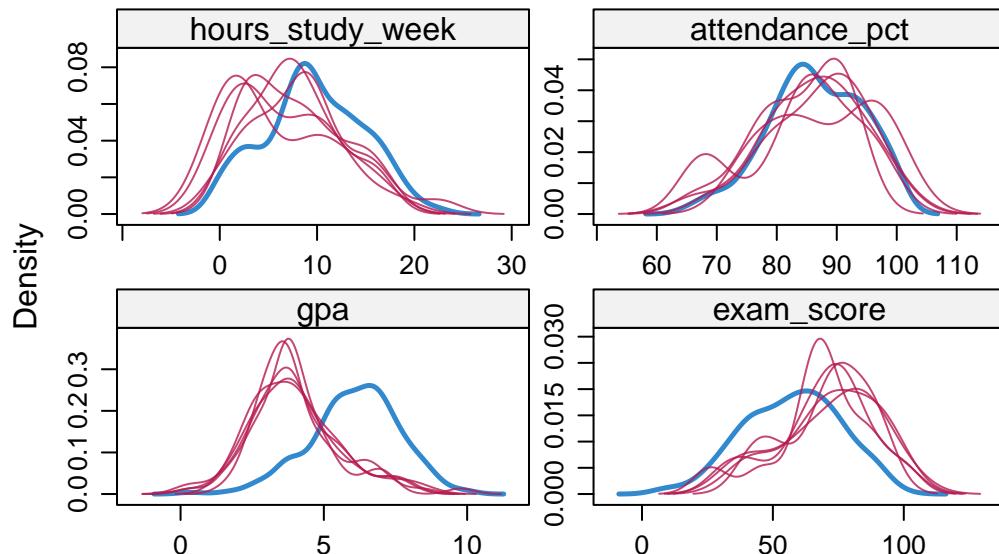
Imputation methods:

| hours_work_week | hours_study_week | attendance_pct | gpa   |
|-----------------|------------------|----------------|-------|
| ""              | "pmm"            | "pmm"          | "pmm" |
|                 |                  |                | shift |
| exam_score      | program          | study_mode     |       |
| "pmm"           | "polyreg"        | "polyreg"      | ""    |

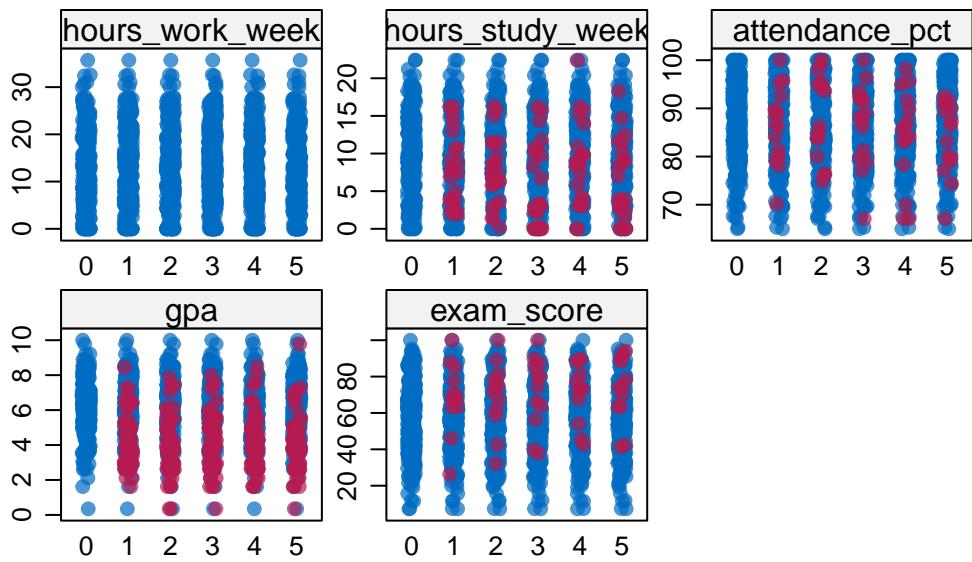
PredictorMatrix:

|                  | hours_work_week | hours_study_week | attendance_pct | gpa | exam_score |
|------------------|-----------------|------------------|----------------|-----|------------|
| hours_work_week  | 0               | 1                | 1              | 1   | 1          |
| hours_study_week | 1               | 0                | 1              | 1   | 1          |
| attendance_pct   | 1               | 1                | 0              | 1   | 1          |
| gpa              | 1               | 1                | 1              | 0   | 1          |
| exam_score       | 1               | 1                | 1              | 1   | 0          |
| program          | 1               | 1                | 1              | 1   | 1          |
|                  | program         | study_mode       | shift          |     |            |
| hours_work_week  | 1               | 1                | 1              |     |            |
| hours_study_week | 1               | 1                | 1              |     |            |
| attendance_pct   | 1               | 1                | 1              |     |            |
| gpa              | 1               | 1                | 1              |     |            |
| exam_score       | 1               | 1                | 1              |     |            |
| program          | 0               | 1                | 1              |     |            |

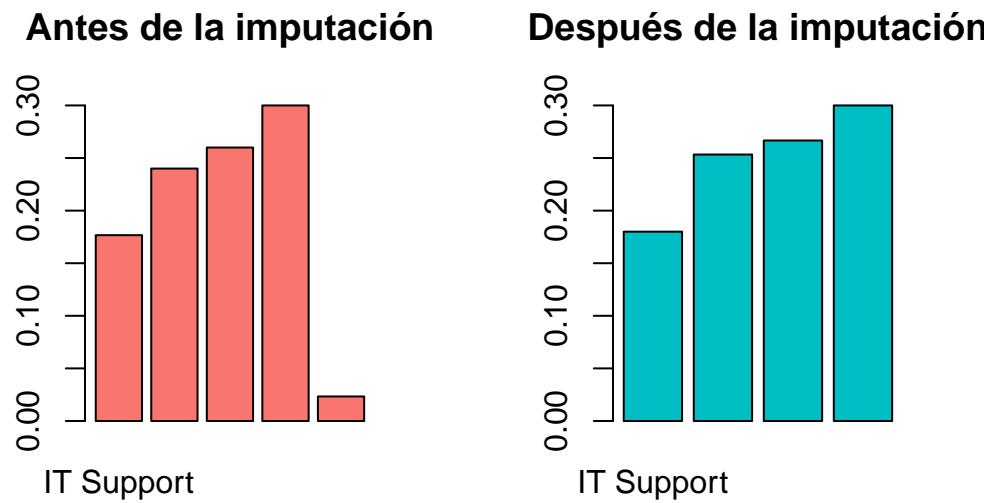
```
###Visualizar las variables cuantitativas  
densityplot(impData)
```



```
stripplot(impData, pch = 20, cex = 1.2)
```

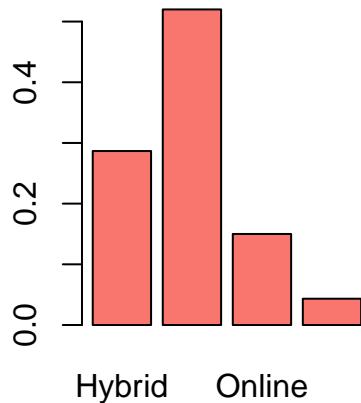


```
###Visualizar las variables cualitativas
completedData <- complete(impData,1)
par(mfrow = c(1, 2))
barplot(prop.table(table(data$program, useNA = "ifany")),
       main = "Antes de la imputación", col = "#F8766D" )
barplot(prop.table(table(completedData$program)),
       main = "Después de la imputación", col = "#00BFC4")
```

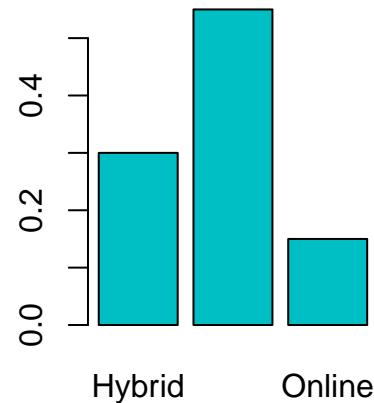


```
par(mfrow = c(1, 2))
barplot(prop.table(table(data$study_mode, useNA = "ifany")),
        main = "Antes de la imputación", col = "#F8766D" )
barplot(prop.table(table(completedData$study_mode)),
        main = "Después de la imputación", col = "#00BFC4")
```

**Antes de la imputación**



**Después de la imputación**



Excepto la variable attendance\_pct, ninguna variable se ha imputado demasiado bien, pero la que se ha imputado muy mal es la de gpa que todos los valores se han imputado por debajo del resto de la distribución, esto no necesariamente está mal, precisamente una hipótesis era que los datos faltantes de gpa era MNAR porque podrías estar relacionados con valores bajos de gpa, pero no explicado por el resto de variables. En este caso, al tener en cuenta de forma multivariante toda la base de datos, ha imputado los valores de la variable como bajos desplazando así toda la distribución.

Por otro lado las dos variables cualitativas parece que se imputa de forma correcta.

## Imputación con missForest

```
library(missForest)

set.seed(123)
data_missForest<-select(data,-student_id)
imp <- missForest(data_missForest)
imp$OOBerror
```

|           |           |
|-----------|-----------|
| NRMSE     | PFC       |
| 0.1441779 | 0.4063217 |

```

set.seed(123)
imp <- missForest(data_missForest, variablewise = TRUE)
imp$OOBerror

```

| MSE       | MSE        | MSE        | MSE       | MSE        | PFC       | PFC       |
|-----------|------------|------------|-----------|------------|-----------|-----------|
| 0.0000000 | 13.1874025 | 52.6990576 | 0.4157744 | 51.1884986 | 0.6928328 | 0.5261324 |
|           | PFC        |            |           |            |           |           |
| 0.0000000 |            |            |           |            |           |           |

```

# Calcular la Standard Deviation para normalizar solo en las cuantitativas
num_vars <- names(data_missForest)[sapply(data_missForest, is.numeric)]
id<-match(num_vars,names(data_missForest))

mse_num <- imp$OOBerror[id]
sd_num <- sapply(data_missForest[id], sd, na.rm = TRUE)

# Calcular el NMRSE
NRMSE <- sqrt(mse_num) / sd_num
names(NRMSE)<-colnames(data_missForest[id])
NRMSE

```

| hours_work_week | hours_study_week | attendance_pct | gpa       |
|-----------------|------------------|----------------|-----------|
| 0.0000000       | 0.7111160        | 0.9260997      | 0.4147127 |
|                 | exam_score       |                |           |
|                 | 0.3855255        |                |           |

De forma global los dos errores tanto para las cuantitativas como para las cualitativas, los errores son pequeños:

NRMSE = 0.14

PFC = 0.41

Si lo sacamos de forma individual por variable, para las cualitativas tenemos:

program (PFC = 0.69)

study\_mode (PFC = 0.53)

hours\_study\_week (NRMSE = 0.71)

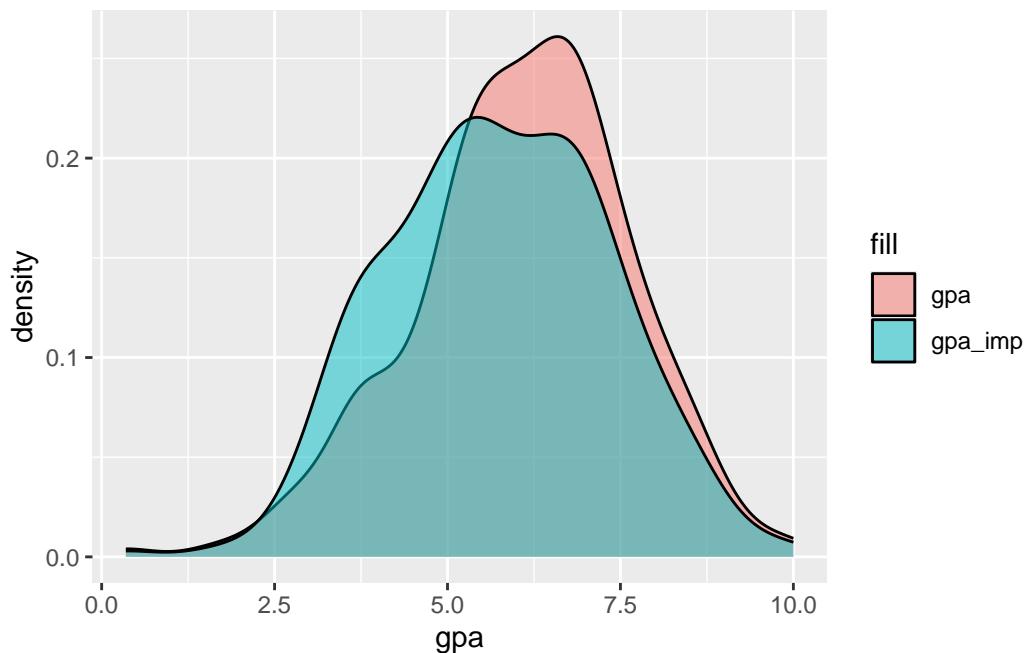
attendance\_pct (NRMSE = 0.93)

gpa (NRMSE = 0.41)

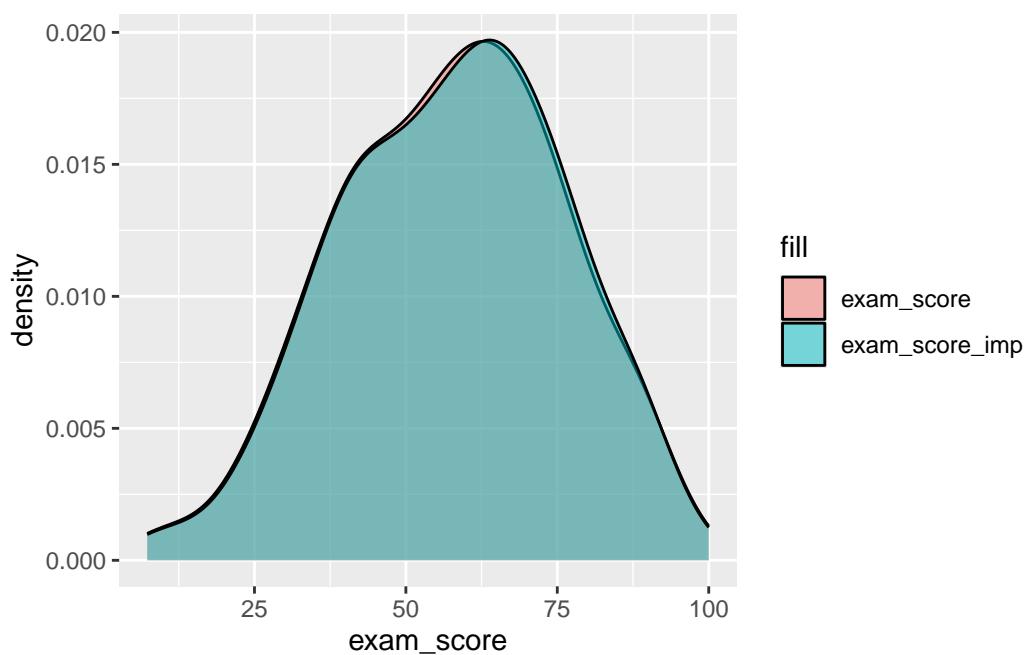
exam\_score (NRMSE = 0.38)

Ninguno parece demasiado alto, pero es curioso que gpa sea uno de los más bajos a diferencia del resto de imputaciones.

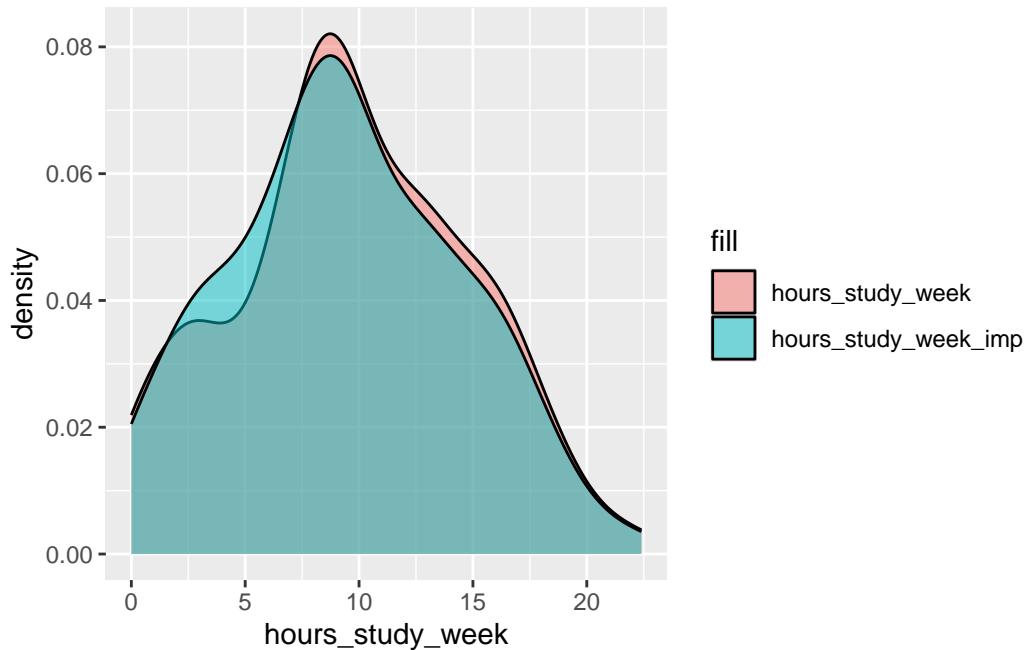
```
g1<-ggplot(data, aes(x =gpa, fill = "gpa")) +  
  geom_density(alpha = 0.5, na.rm = TRUE) +  
  geom_density(aes(x = imp$ximp$gpa,  
                  fill = "gpa_imp"), alpha = 0.5)  
  
g2<-ggplot(data, aes(x = exam_score, fill = "exam_score")) +  
  geom_density(alpha = 0.5, na.rm = TRUE) +  
  geom_density(aes(x = imp$ximp$exam_score,  
                  fill = "exam_score_imp"), alpha = 0.5)  
  
g3<-ggplot(data, aes(x = hours_study_week,  
                  fill = "hours_study_week")) +  
  geom_density(alpha = 0.5, na.rm = TRUE) +  
  geom_density(aes(x = imp$ximp$hours_study_week,  
                  fill = "hours_study_week_imp"), alpha = 0.5)  
  
g4<-ggplot(data, aes(x = attendance_pct,  
                  fill = "attendance_pct")) +  
  geom_density(alpha = 0.5, na.rm = TRUE) +  
  geom_density(aes(x = imp$ximp$attendance_pct,  
                  fill = "attendance_pct_imp"), alpha = 0.5)  
  
g1
```



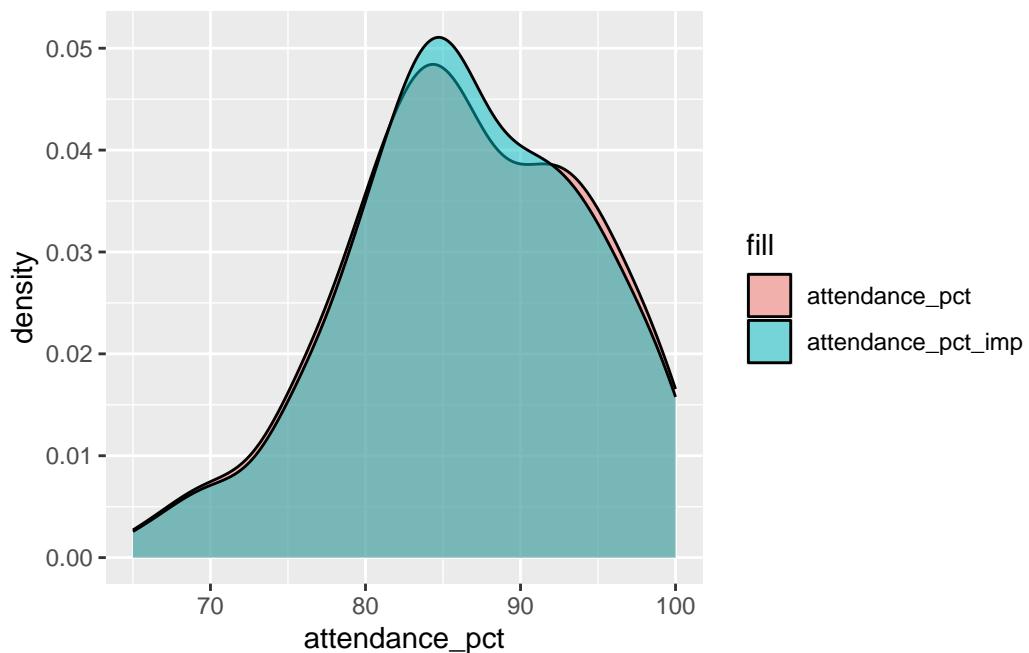
g2



g3



g4



El gráfico de gpa es el que peor se ajusta a su distribución, pero mejora mucho a lo que veíamos con el resto de imputaciones.