# Data challenge & SHS: Logistic regression and linear model

Julie Josse, Gaël Varoquaux, and Bénédicte Colnet

February 2021

**Abstract**

In this tutorial, you will perform a logistic regression with `R`. This is the first exercice and we will do it together in class. At the end you can find an exercice with a simple linear regression you should be able to do alone at home (solutions will be given later).

## Contents

*Credits for this lab*: **An Introduction to Statistical Learning: With Applications in R** book from Garet James, Daniela Witten, Trevor Hastie, Robert Tibshirani (in particular for the exercice on Logistic Regression and stock market) The exercice on linear model comes from Imke Mayer's labs. Thanks to them.

## Logistic regression: stock market data

In this part we use the `Smarket` data, which is part of the `ISLR` library. This data set consists of percentage returns for the S&P 500 stock index over 1250 days, from the beginning of 2001 until the end of 2005.

The S&P 500,or simply the S&P, is a stock market index that measures the stock performance of 500 large companies listed on stock exchanges in the United States. It is one of the most commonly followed equity indices. (I guess we can compare it with the French CAC 40)

Therefore you have 1250 observations on the following 9 variables.

`Year` The year that the observation was recorded

`Lag1` Percentage return for previous day

`Lag2` Percentage return for 2 days previous

`Lag3` Percentage return for 3 days previous

`Lag4` Percentage return for 4 days previous

**Lag5** Percentage return for 5 days previous

**Volume** The number of shares traded

**Today** The percentage return on the date in question

**Direction** A factor with levels Down and Up indicating whether the market had a positive or negative return on a given day

## Question 1: Data exploration

Load the library `ISLR` and inspect the data set. Do you see a link between returns? For example you can also look at correlation. What can you say on the volume of shares traded over year?

**Solution**

```
library(ISLR)
Smarket <- Smarket
names(Smarket)
```

```
## [1] "Year"      "Lag1"      "Lag2"      "Lag3"      "Lag4"      "Lag5"
## [7] "Volume"    "Today"     "Direction"
```

```
summary(Smarket)
```

```
##       Year           Lag1                Lag2                Lag3
##  Min.   :2001   Min.   :-4.922000   Min.   :-4.922000   Min.   :-4.922000
##  1st Qu.:2002   1st Qu.:-0.639500   1st Qu.:-0.639500   1st Qu.:-0.640000
##  Median :2003   Median : 0.039000   Median : 0.039000   Median : 0.038500
##  Mean   :2003   Mean   : 0.003834   Mean   : 0.003919   Mean   : 0.001716
##  3rd Qu.:2004   3rd Qu.: 0.596750   3rd Qu.: 0.596750   3rd Qu.: 0.596750
##  Max.   :2005   Max.   : 5.733000   Max.   : 5.733000   Max.   : 5.733000
##       Lag4                Lag5               Volume          Today
##  Min.   :-4.922000   Min.   :-4.92200   Min.   :0.3561   Min.   :-4.922000
##  1st Qu.:-0.640000   1st Qu.:-0.64000   1st Qu.:1.2574   1st Qu.:-0.639500
##  Median : 0.038500   Median : 0.03850   Median :1.4229   Median : 0.038500
##  Mean   : 0.001636   Mean   : 0.00561   Mean   :1.4783   Mean   : 0.003138
##  3rd Qu.: 0.596750   3rd Qu.: 0.59700   3rd Qu.:1.6417   3rd Qu.: 0.596750
##  Max.   : 5.733000   Max.   : 5.73300   Max.   :3.1525   Max.   : 5.733000
##  Direction
##  Down:602
##  Up  :648
##
##
##
##
```

```
cor(Smarket[,-9])
```

```
##               Year         Lag1         Lag2         Lag3         Lag4
## Year    1.00000000  0.029699649  0.030596422  0.033194581  0.035688718
## Lag1    0.02969965  1.000000000 -0.026294328 -0.010803402 -0.002985911
## Lag2    0.03059642 -0.026294328  1.000000000 -0.025896670 -0.010853533
## Lag3    0.03319458 -0.010803402 -0.025896670  1.000000000 -0.024051036
## Lag4    0.03568872 -0.002985911 -0.010853533 -0.024051036  1.000000000
## Lag5    0.02978799 -0.005674606 -0.003557949 -0.018808338 -0.027083641
## Volume  0.53900647  0.040909908 -0.043383215 -0.041823686 -0.048414246
## Today   0.03009523 -0.026155045 -0.010250033 -0.002447647 -0.006899527
```
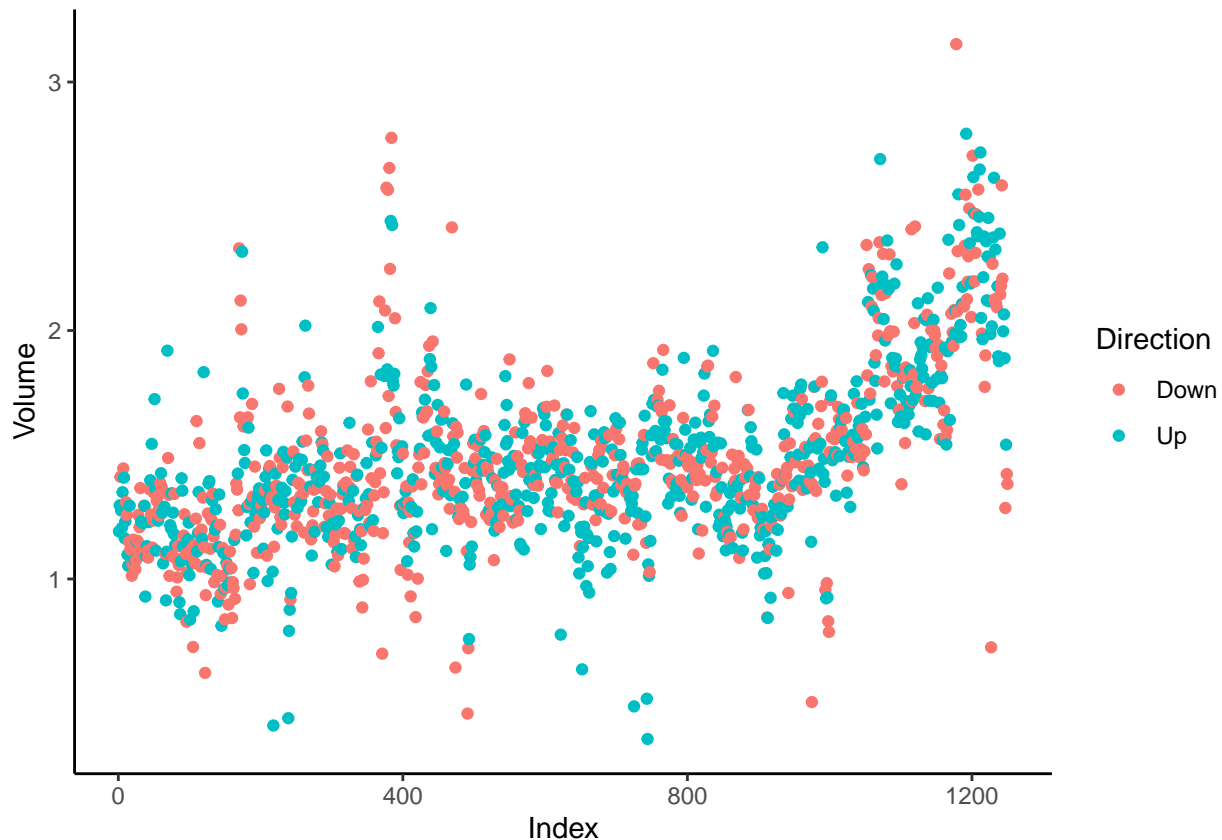
```
##                   Lag5       Volume        Today
## Year     0.029787995  0.53900647  0.030095229
## Lag1    -0.005674606  0.04090991 -0.026155045
## Lag2    -0.003557949 -0.04338321 -0.010250033
## Lag3    -0.018808338 -0.04182369 -0.002447647
## Lag4    -0.027083641 -0.04841425 -0.006899527
## Lag5     1.000000000 -0.02200231 -0.034860083
## Volume -0.022002315  1.00000000  0.014591823
## Today  -0.034860083  0.01459182  1.000000000
```

There appears to be little correlation between today's returns and previous days' returns. We can observe a correlation on the `Year` and `Volume`.

```
library(ggplot2)
```

```
## Warning: package 'ggplot2' was built under R version 3.6.2
```

```
ggplot(Smarket, aes(x = as.numeric(row.names(Smarket)), y = Volume, color = Direction)) +
  geom_point() +
  theme_classic() +
  xlab("Index")
```



**End of solution**

## Question 2: Logistic regression

Fit a logistic regression model in order to predict `Direction` using all the other available variables.

For this you can use `glm()`, a class of models that includes logistic regression.

Interpret the result. What is the coefficient that is the most linked to the outcome according to this model?

**Solution**

```
glm.fit = glm(Direction~Lag1+Lag2+Lag3+Lag4+Lag5+Volume, data=Smarket, family = binomial)
summary(glm.fit)
```

```
##
## Call:
## glm(formula = Direction ~ Lag1 + Lag2 + Lag3 + Lag4 + Lag5 +
##     Volume, family = binomial, data = Smarket)
##
## Deviance Residuals:
##    Min      1Q  Median      3Q     Max
## -1.446  -1.203   1.065   1.145   1.326
##
## Coefficients:
##             Estimate Std. Error z value Pr(>|z|)
## (Intercept) -0.126000   0.240736  -0.523    0.601
## Lag1        -0.073074   0.050167  -1.457    0.145
## Lag2        -0.042301   0.050086  -0.845    0.398
## Lag3         0.011085   0.049939   0.222    0.824
## Lag4         0.009359   0.049974   0.187    0.851
## Lag5         0.010313   0.049511   0.208    0.835
## Volume       0.135441   0.158360   0.855    0.392
##
## (Dispersion parameter for binomial family taken to be 1)
##
##     Null deviance: 1731.2  on 1249  degrees of freedom
## Residual deviance: 1727.6  on 1243  degrees of freedom
## AIC: 1741.6
##
## Number of Fisher Scoring iterations: 3
```

The smallest p-value here is associated with `Lag1`. The negative coefficient for this predictor suggests that if the market had a positive return yesterday, then it is less likely to go up today. However, at a value of 0.15, the p-value is still relatively large, and so there is no clear evidence of a real association between `Lag1` and `Direction`.

Be careful to look at which variable is the 1 or the 0. `R` automatically creates so-called dummy variables you can inspect with `contrasts()`.

```
contrasts(Smarket$Direction)
```

```
##      Up
## Down  0
## Up    1
```

**End of solution**

## Question 3: Prediction

You can use the `predict()` function to perform prediction that the market will go up given other values. Remember that it corresponds to the quantity:

$$\mathbb{P}(Direction = Up | Lag1, \ldots, Volume)$$

If no data set is supplied to the `predict()` function, then the probabilities are computed for the training data used to fit the logistic regression model.

After doing this prediction, you will have the probability of having $Y = 1$. Now, create a confusion matrix with the function `table()` to determine how many observations were correctly or incorrectly classified.
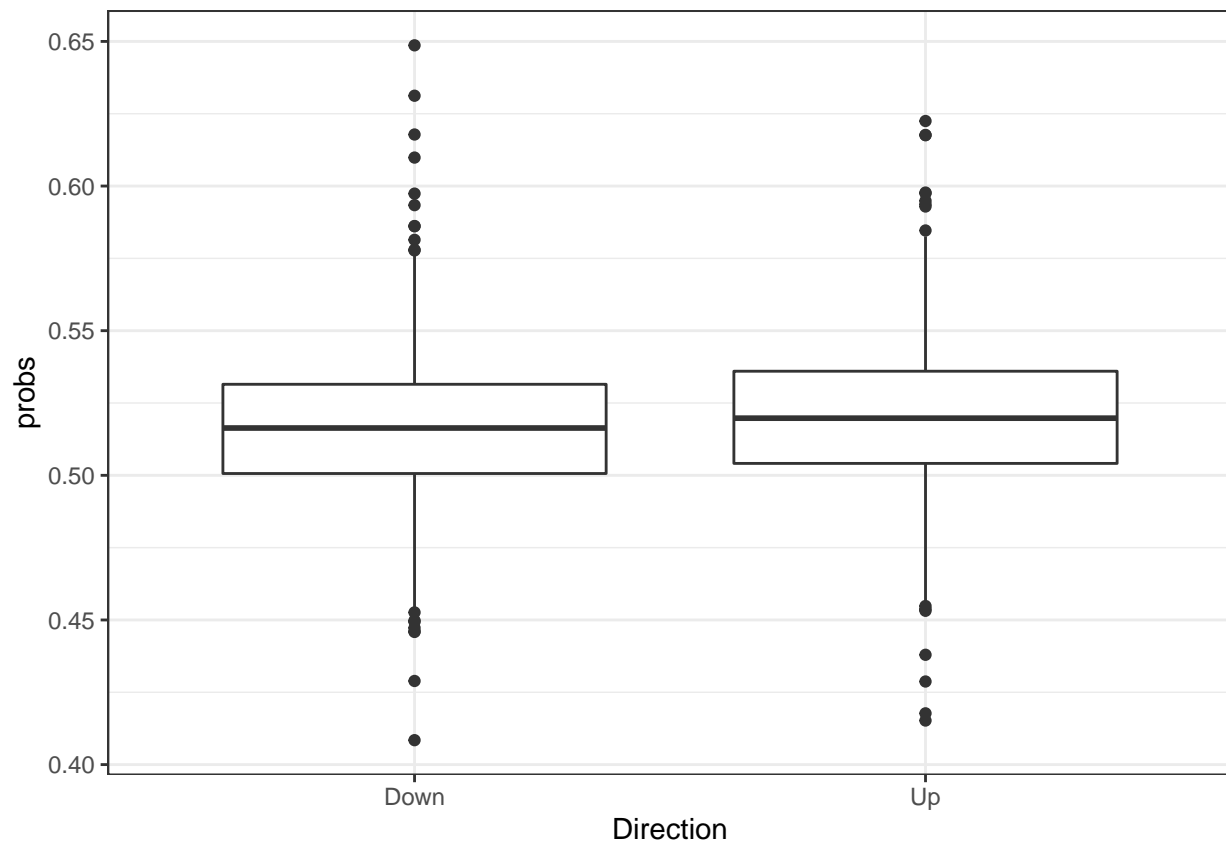
Conclude on this model efficacy. What would you do to better assess this model efficacy?

**Solution**

```
glm.probs <- predict(glm.fit, type = "response")

Smarket$probs <- glm.probs


ggplot(Smarket, aes(y = probs, x = Direction, group = Direction)) +
  geom_boxplot() +
  theme_bw()
```



The following two commands create a vector of class predictions based on whether the predicted probability of a market increase is greater than or less than 0.5.

```
glm.pred = rep(0, 1250)
```

```
glm.pred[glm.probs > .5] = 1
table(glm.pred, Smarket$Direction)
```

```
##
## glm.pred Down  Up
##        0  145 141
```

```
##           1  457 507
```

The diagonal elements of the confusion matrix indicate correct predictions, while the off-diagonals represent incorrect predictions. Hence our model correctly predicted that the market would go up on 507 days and that it would go down on 145 days, for a total of $507 + 145 = 652$ correct predictions. Logistic regression correctly predicted the movement of the market 52.2% of the time.

This value is a little better than the random classifier. Remember that without knowing anything on your data, when you want to classify something you can still use a random classifier that will say 0 or 1 at each new guess without any **a priori** on the data. You can notice that the current performance is very bad because 52.2% is our training error, which is clearly optimistic (you will see this with Gaël class and the MOOC on Scikit Learn).

The next part of the solution is a bonus part. You can look at it if you want. To implement this strategy, we will first create a vector corresponding to the observations from 2001 through 2004. We will then use this vector to create a held out data set of observations from 2005.

```
test = Smarket[Smarket$Year == 2005,]
glm.fit = glm(Direction~Lag1+Lag2+Lag3+Lag4+Lag5+Volume, family=binomial, data = Smarket[Smarket$Year <
glm.probs = predict(glm.fit, newdata = test, type = "response")
glm.pred = rep(0, nrow(test))
glm.pred[glm.probs >.5]=1
table(glm.pred, test$Direction)
```

```
##
## glm.pred Down Up
##        0   77 97
##        1   34 44
```

The results are rather disappointing: the test error rate is 48%, which is worse than random guessing! Note that if it was possible to accurately predicts day return with previous days, it would be easier to be a trader ;) !

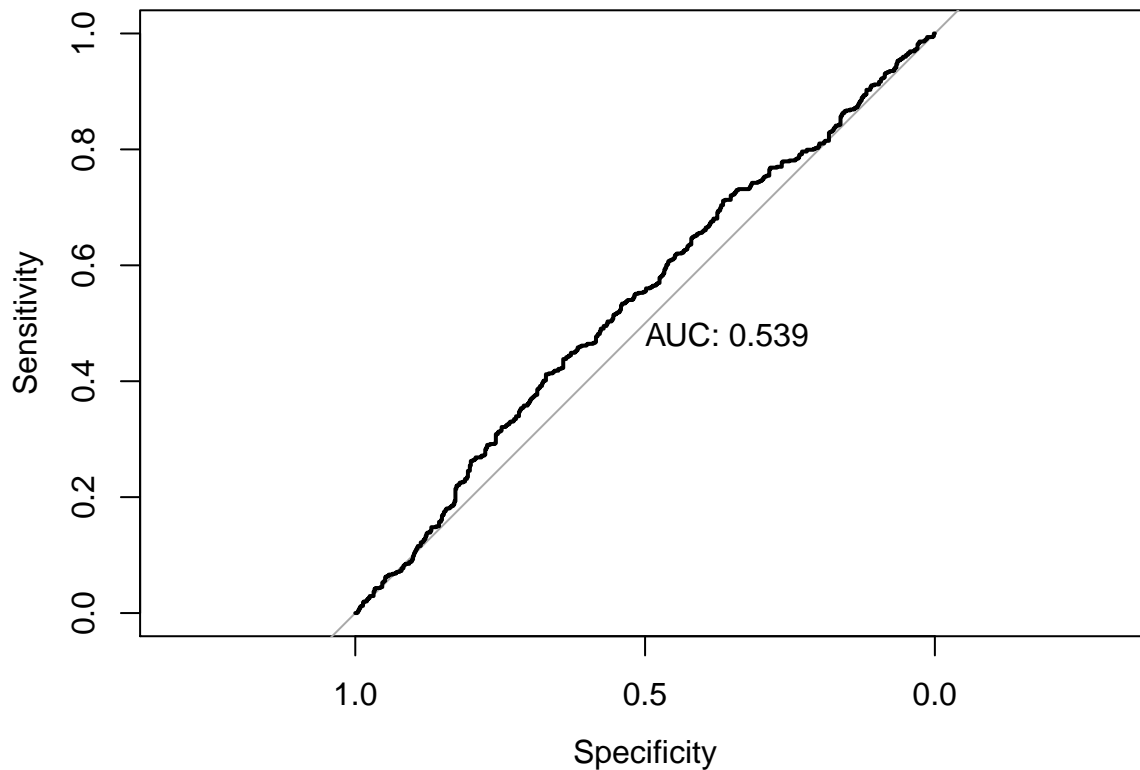**End of solution**

## Question 4: ROC curves

The prediction performed before is by default made with a cutoff at 0.5. But maybe another threshold would help to have a better performance. Using the library `pROC`, screen for the best cutoff. The function you will use is the function `roc()`.

**Solution**

```
library(pROC)
```

```
## Type 'citation("pROC")' for a citation.
```

```
##
## Attaching package: 'pROC'
```

```
## The following objects are masked from 'package:stats':
##
##     cov, smooth, var
```

```
glm.fit = glm(Direction~Lag1+Lag2+Lag3+Lag4+Lag5+Volume, data=Smarket, family =binomial)
glm.probs <- predict(glm.fit, type = "response")
test_roc = roc(Smarket$Direction ~ glm.probs, plot = TRUE, print.auc = TRUE)
```

```
## Setting levels: control = Down, case = Up
```

```
## Setting direction: controls < cases
```

```
as.numeric(test_roc$auc)
```

```
## [1] 0.5387341
```

A good model will have a high AUC, that is as often as possible a high sensitivity and specificity.

**End of solution**

## Exploratory data analysis and simple regression

### The database

The data are stored in the file 'bea-2006.csv'. It contains information about the economies of the 366 metropolitan statistical areas" (cities) of the US in 2006. In particular, it lists, for each city:

- the population,
- the total value of all goods and services produced for sale in the city that year per person (per capita gross metropolitan product", pcgmp),
- and the share of economic output coming from *four* selected industries.

### Question 1: load data

Load the data and perform a summary analysis.

*Solution 1*

```
data <- read.csv('bea-2006.csv', row.names=1)
```

```
summary(data)
```

```
##      pcgmp             pop              finance          prof.tech
## Min.   :14920   Min.   :   54980   Min.   :0.03845   Min.   :0.01474
```

```
##   1st Qu.:26532    1st Qu.:   135625    1st Qu.:0.10403    1st Qu.:0.02932
##   Median :31615    Median :   231500    Median :0.14140    Median :0.04212
##   Mean   :32923    Mean   :   680898    Mean   :0.15082    Mean   :0.04905
##   3rd Qu.:38212    3rd Qu.:   530875    3rd Qu.:0.18122    3rd Qu.:0.05932
##   Max.   :77860    Max.   :18850000     Max.   :0.38480    Max.   :0.19080
##                                         NA's   :12         NA's   :112
##       ict             management
##   Min.   :0.00349   Min.   :0.00042
##   1st Qu.:0.01215   1st Qu.:0.00294
##   Median :0.02218   Median :0.00651
##   Mean   :0.03910   Mean   :0.00908
##   3rd Qu.:0.04072   3rd Qu.:0.01191
##   Max.   :0.58600   Max.   :0.05431
##   NA's   :76        NA's   :157
```
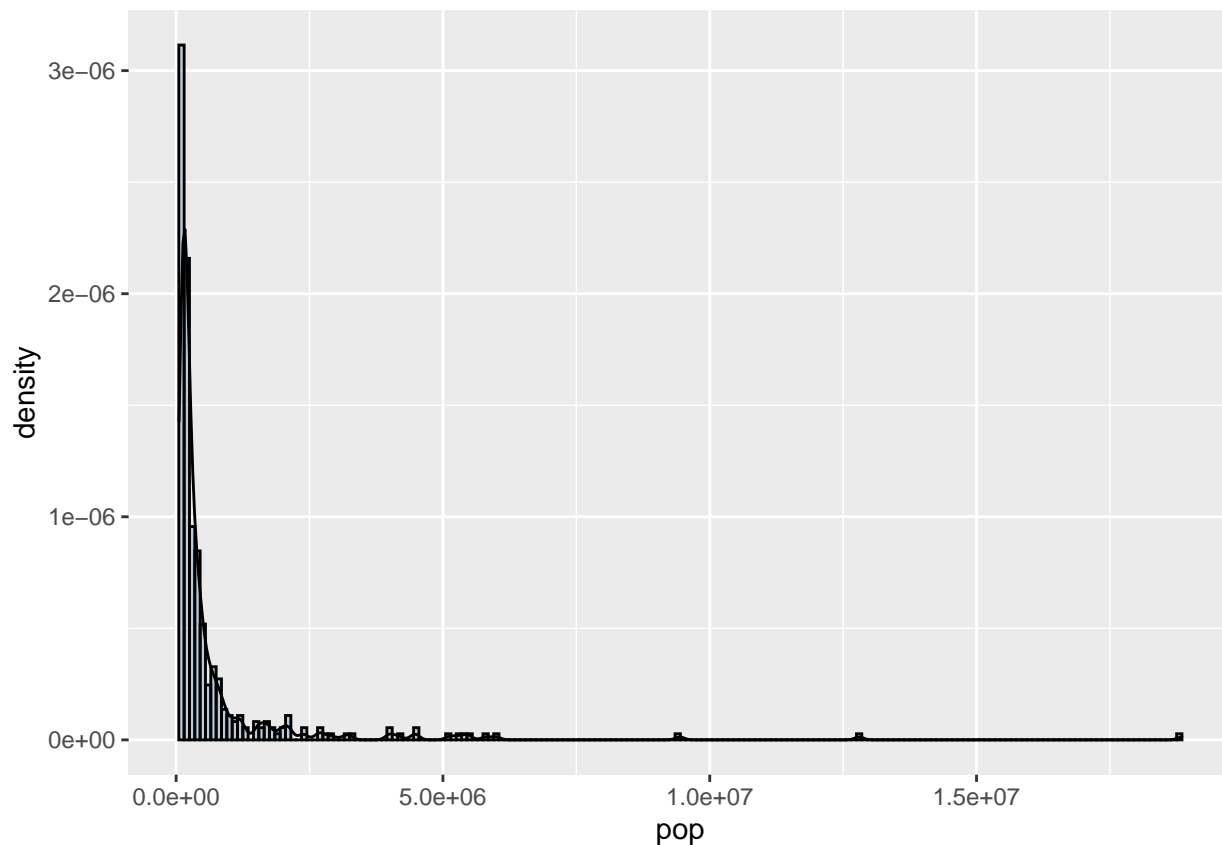
*End of solution 1*

## Question 2: data exploration

Produce histogram of population (density and the histogram with "bar") and the box plot of the pgmp column.

Tips: Don't hesitate to do an histogram without the outliers.

*Solution 2*

```
library(ggplot2)
ggplot(data, aes(x = pop)) +
 geom_histogram(aes(y=..density..), alpha=0.3,
                position="identity", binwidth = 100000, fill = "steelblue", color = "black")+
 geom_density(alpha=0.9)
```

```
theme_bw()
```

```
## List of 93
##  $ line                    :List of 6
##   ..$ colour       : chr "black"
##   ..$ size         : num 0.5
##   ..$ linetype     : num 1
##   ..$ lineend      : chr "butt"
##   ..$ arrow        : logi FALSE
##   ..$ inherit.blank: logi TRUE
##   ..- attr(*, "class")= chr [1:2] "element_line" "element"
##  $ rect                    :List of 5
##   ..$ fill         : chr "white"
##   ..$ colour       : chr "black"
##   ..$ size         : num 0.5
##   ..$ linetype     : num 1
##   ..$ inherit.blank: logi TRUE
##   ..- attr(*, "class")= chr [1:2] "element_rect" "element"
##  $ text                    :List of 11
##   ..$ family       : chr ""
##   ..$ face         : chr "plain"
##   ..$ colour       : chr "black"
##   ..$ size         : num 11
##   ..$ hjust        : num 0.5
##   ..$ vjust        : num 0.5
##   ..$ angle        : num 0
##   ..$ lineheight   : num 0.9
```

```
##   ..$ margin       : 'margin' num [1:4] 0pt 0pt 0pt 0pt
##   .. ..- attr(*, "valid.unit")= int 8
##   .. ..- attr(*, "unit")= chr "pt"
##   ..$ debug        : logi FALSE
##   ..$ inherit.blank: logi TRUE
##   ..- attr(*, "class")= chr [1:2] "element_text" "element"
## $ title                    : NULL
## $ aspect.ratio             : NULL
## $ axis.title               : NULL
## $ axis.title.x             :List of 11
##   ..$ family       : NULL
##   ..$ face         : NULL
##   ..$ colour       : NULL
##   ..$ size         : NULL
##   ..$ hjust        : NULL
##   ..$ vjust        : num 1
##   ..$ angle        : NULL
##   ..$ lineheight   : NULL
##   ..$ margin       : 'margin' num [1:4] 2.75pt 0pt 0pt 0pt
##   .. ..- attr(*, "valid.unit")= int 8
##   .. ..- attr(*, "unit")= chr "pt"
##   ..$ debug        : NULL
##   ..$ inherit.blank: logi TRUE
##   ..- attr(*, "class")= chr [1:2] "element_text" "element"
## $ axis.title.x.top         :List of 11
##   ..$ family       : NULL
##   ..$ face         : NULL
##   ..$ colour       : NULL
##   ..$ size         : NULL
##   ..$ hjust        : NULL
##   ..$ vjust        : num 0
##   ..$ angle        : NULL
##   ..$ lineheight   : NULL
##   ..$ margin       : 'margin' num [1:4] 0pt 0pt 2.75pt 0pt
##   .. ..- attr(*, "valid.unit")= int 8
##   .. ..- attr(*, "unit")= chr "pt"
##   ..$ debug        : NULL
##   ..$ inherit.blank: logi TRUE
##   ..- attr(*, "class")= chr [1:2] "element_text" "element"
## $ axis.title.x.bottom      : NULL
## $ axis.title.y             :List of 11
##   ..$ family       : NULL
##   ..$ face         : NULL
##   ..$ colour       : NULL
##   ..$ size         : NULL
##   ..$ hjust        : NULL
##   ..$ vjust        : num 1
##   ..$ angle        : num 90
##   ..$ lineheight   : NULL
##   ..$ margin       : 'margin' num [1:4] 0pt 2.75pt 0pt 0pt
##   .. ..- attr(*, "valid.unit")= int 8
##   .. ..- attr(*, "unit")= chr "pt"
##   ..$ debug        : NULL
##   ..$ inherit.blank: logi TRUE
```

```
##   ..- attr(*, "class")= chr [1:2] "element_text" "element"
## $ axis.title.y.left       : NULL
## $ axis.title.y.right      :List of 11
##   ..$ family     : NULL
##   ..$ face       : NULL
##   ..$ colour     : NULL
##   ..$ size       : NULL
##   ..$ hjust      : NULL
##   ..$ vjust      : num 0
##   ..$ angle      : num -90
##   ..$ lineheight : NULL
##   ..$ margin     : 'margin' num [1:4] 0pt 0pt 0pt 2.75pt
##   .. ..- attr(*, "valid.unit")= int 8
##   .. ..- attr(*, "unit")= chr "pt"
##   ..$ debug      : NULL
##   ..$ inherit.blank: logi TRUE
##   ..- attr(*, "class")= chr [1:2] "element_text" "element"
## $ axis.text               :List of 11
##   ..$ family     : NULL
##   ..$ face       : NULL
##   ..$ colour     : chr "grey30"
##   ..$ size       : 'rel' num 0.8
##   ..$ hjust      : NULL
##   ..$ vjust      : NULL
##   ..$ angle      : NULL
##   ..$ lineheight : NULL
##   ..$ margin     : NULL
##   ..$ debug      : NULL
##   ..$ inherit.blank: logi TRUE
##   ..- attr(*, "class")= chr [1:2] "element_text" "element"
## $ axis.text.x             :List of 11
##   ..$ family     : NULL
##   ..$ face       : NULL
##   ..$ colour     : NULL
##   ..$ size       : NULL
##   ..$ hjust      : NULL
##   ..$ vjust      : num 1
##   ..$ angle      : NULL
##   ..$ lineheight : NULL
##   ..$ margin     : 'margin' num [1:4] 2.2pt 0pt 0pt 0pt
##   .. ..- attr(*, "valid.unit")= int 8
##   .. ..- attr(*, "unit")= chr "pt"
##   ..$ debug      : NULL
##   ..$ inherit.blank: logi TRUE
##   ..- attr(*, "class")= chr [1:2] "element_text" "element"
## $ axis.text.x.top         :List of 11
##   ..$ family     : NULL
##   ..$ face       : NULL
##   ..$ colour     : NULL
##   ..$ size       : NULL
##   ..$ hjust      : NULL
##   ..$ vjust      : num 0
##   ..$ angle      : NULL
##   ..$ lineheight : NULL
```

```
##    ..$ margin       : 'margin' num [1:4] 0pt 0pt 2.2pt 0pt
##    .. ..- attr(*, "valid.unit")= int 8
##    .. ..- attr(*, "unit")= chr "pt"
##    ..$ debug        : NULL
##    ..$ inherit.blank: logi TRUE
##    ..- attr(*, "class")= chr [1:2] "element_text" "element"
##  $ axis.text.x.bottom       : NULL
##  $ axis.text.y              :List of 11
##    ..$ family       : NULL
##    ..$ face         : NULL
##    ..$ colour       : NULL
##    ..$ size         : NULL
##    ..$ hjust        : num 1
##    ..$ vjust        : NULL
##    ..$ angle        : NULL
##    ..$ lineheight   : NULL
##    ..$ margin       : 'margin' num [1:4] 0pt 2.2pt 0pt 0pt
##    .. ..- attr(*, "valid.unit")= int 8
##    .. ..- attr(*, "unit")= chr "pt"
##    ..$ debug        : NULL
##    ..$ inherit.blank: logi TRUE
##    ..- attr(*, "class")= chr [1:2] "element_text" "element"
##  $ axis.text.y.left         : NULL
##  $ axis.text.y.right        :List of 11
##    ..$ family       : NULL
##    ..$ face         : NULL
##    ..$ colour       : NULL
##    ..$ size         : NULL
##    ..$ hjust        : num 0
##    ..$ vjust        : NULL
##    ..$ angle        : NULL
##    ..$ lineheight   : NULL
##    ..$ margin       : 'margin' num [1:4] 0pt 0pt 0pt 2.2pt
##    .. ..- attr(*, "valid.unit")= int 8
##    .. ..- attr(*, "unit")= chr "pt"
##    ..$ debug        : NULL
##    ..$ inherit.blank: logi TRUE
##    ..- attr(*, "class")= chr [1:2] "element_text" "element"
##  $ axis.ticks               :List of 6
##    ..$ colour       : chr "grey20"
##    ..$ size         : NULL
##    ..$ linetype     : NULL
##    ..$ lineend      : NULL
##    ..$ arrow        : logi FALSE
##    ..$ inherit.blank: logi TRUE
##    ..- attr(*, "class")= chr [1:2] "element_line" "element"
##  $ axis.ticks.x             : NULL
##  $ axis.ticks.x.top         : NULL
##  $ axis.ticks.x.bottom      : NULL
##  $ axis.ticks.y             : NULL
##  $ axis.ticks.y.left        : NULL
##  $ axis.ticks.y.right       : NULL
##  $ axis.ticks.length        : 'unit' num 2.75pt
##    ..- attr(*, "valid.unit")= int 8
```

```
##    ..- attr(*, "unit")= chr "pt"
## $ axis.ticks.length.x       : NULL
## $ axis.ticks.length.x.top   : NULL
## $ axis.ticks.length.x.bottom: NULL
## $ axis.ticks.length.y       : NULL
## $ axis.ticks.length.y.left  : NULL
## $ axis.ticks.length.y.right : NULL
## $ axis.line                 : list()
##    ..- attr(*, "class")= chr [1:2] "element_blank" "element"
## $ axis.line.x               : NULL
## $ axis.line.x.top           : NULL
## $ axis.line.x.bottom        : NULL
## $ axis.line.y               : NULL
## $ axis.line.y.left          : NULL
## $ axis.line.y.right         : NULL
## $ legend.background         :List of 5
##    ..$ fill         : NULL
##    ..$ colour       : logi NA
##    ..$ size         : NULL
##    ..$ linetype     : NULL
##    ..$ inherit.blank: logi TRUE
##    ..- attr(*, "class")= chr [1:2] "element_rect" "element"
## $ legend.margin             : 'margin' num [1:4] 5.5pt 5.5pt 5.5pt 5.5pt
##    ..- attr(*, "valid.unit")= int 8
##    ..- attr(*, "unit")= chr "pt"
## $ legend.spacing            : 'unit' num 11pt
##    ..- attr(*, "valid.unit")= int 8
##    ..- attr(*, "unit")= chr "pt"
## $ legend.spacing.x          : NULL
## $ legend.spacing.y          : NULL
## $ legend.key                :List of 5
##    ..$ fill         : chr "white"
##    ..$ colour       : logi NA
##    ..$ size         : NULL
##    ..$ linetype     : NULL
##    ..$ inherit.blank: logi TRUE
##    ..- attr(*, "class")= chr [1:2] "element_rect" "element"
## $ legend.key.size           : 'unit' num 1.2lines
##    ..- attr(*, "valid.unit")= int 3
##    ..- attr(*, "unit")= chr "lines"
## $ legend.key.height         : NULL
## $ legend.key.width          : NULL
## $ legend.text               :List of 11
##    ..$ family       : NULL
##    ..$ face         : NULL
##    ..$ colour       : NULL
##    ..$ size         : 'rel' num 0.8
##    ..$ hjust        : NULL
##    ..$ vjust        : NULL
##    ..$ angle        : NULL
##    ..$ lineheight   : NULL
##    ..$ margin       : NULL
##    ..$ debug        : NULL
##    ..$ inherit.blank: logi TRUE
```

```
##    ..- attr(*, "class")= chr [1:2] "element_text" "element"
## $ legend.text.align        : NULL
## $ legend.title             :List of 11
##   ..$ family       : NULL
##   ..$ face         : NULL
##   ..$ colour       : NULL
##   ..$ size         : NULL
##   ..$ hjust        : num 0
##   ..$ vjust        : NULL
##   ..$ angle        : NULL
##   ..$ lineheight   : NULL
##   ..$ margin       : NULL
##   ..$ debug        : NULL
##   ..$ inherit.blank: logi TRUE
##   ..- attr(*, "class")= chr [1:2] "element_text" "element"
## $ legend.title.align       : NULL
## $ legend.position          : chr "right"
## $ legend.direction         : NULL
## $ legend.justification     : chr "center"
## $ legend.box               : NULL
## $ legend.box.just          : NULL
## $ legend.box.margin        : 'margin' num [1:4] 0cm 0cm 0cm 0cm
##   ..- attr(*, "valid.unit")= int 1
##   ..- attr(*, "unit")= chr "cm"
## $ legend.box.background     : list()
##   ..- attr(*, "class")= chr [1:2] "element_blank" "element"
## $ legend.box.spacing        : 'unit' num 11pt
##   ..- attr(*, "valid.unit")= int 8
##   ..- attr(*, "unit")= chr "pt"
## $ panel.background          :List of 5
##   ..$ fill         : chr "white"
##   ..$ colour       : logi NA
##   ..$ size         : NULL
##   ..$ linetype     : NULL
##   ..$ inherit.blank: logi TRUE
##   ..- attr(*, "class")= chr [1:2] "element_rect" "element"
## $ panel.border              :List of 5
##   ..$ fill         : logi NA
##   ..$ colour       : chr "grey20"
##   ..$ size         : NULL
##   ..$ linetype     : NULL
##   ..$ inherit.blank: logi TRUE
##   ..- attr(*, "class")= chr [1:2] "element_rect" "element"
## $ panel.spacing             : 'unit' num 5.5pt
##   ..- attr(*, "valid.unit")= int 8
##   ..- attr(*, "unit")= chr "pt"
## $ panel.spacing.x           : NULL
## $ panel.spacing.y           : NULL
## $ panel.grid                :List of 6
##   ..$ colour       : chr "grey92"
##   ..$ size         : NULL
##   ..$ linetype     : NULL
##   ..$ lineend      : NULL
##   ..$ arrow        : logi FALSE
```

```
##   ..$ inherit.blank: logi TRUE
##   ..- attr(*, "class")= chr [1:2] "element_line" "element"
## $ panel.grid.major       : NULL
## $ panel.grid.minor       :List of 6
##   ..$ colour    : NULL
##   ..$ size      : 'rel' num 0.5
##   ..$ linetype  : NULL
##   ..$ lineend   : NULL
##   ..$ arrow     : logi FALSE
##   ..$ inherit.blank: logi TRUE
##   ..- attr(*, "class")= chr [1:2] "element_line" "element"
## $ panel.grid.major.x     : NULL
## $ panel.grid.major.y     : NULL
## $ panel.grid.minor.x     : NULL
## $ panel.grid.minor.y     : NULL
## $ panel.ontop            : logi FALSE
## $ plot.background        :List of 5
##   ..$ fill      : NULL
##   ..$ colour    : chr "white"
##   ..$ size      : NULL
##   ..$ linetype  : NULL
##   ..$ inherit.blank: logi TRUE
##   ..- attr(*, "class")= chr [1:2] "element_rect" "element"
## $ plot.title             :List of 11
##   ..$ family    : NULL
##   ..$ face      : NULL
##   ..$ colour    : NULL
##   ..$ size      : 'rel' num 1.2
##   ..$ hjust     : num 0
##   ..$ vjust     : num 1
##   ..$ angle     : NULL
##   ..$ lineheight: NULL
##   ..$ margin    : 'margin' num [1:4] 0pt 0pt 5.5pt 0pt
##   .. ..- attr(*, "valid.unit")= int 8
##   .. ..- attr(*, "unit")= chr "pt"
##   ..$ debug     : NULL
##   ..$ inherit.blank: logi TRUE
##   ..- attr(*, "class")= chr [1:2] "element_text" "element"
## $ plot.title.position    : chr "panel"
## $ plot.subtitle          :List of 11
##   ..$ family    : NULL
##   ..$ face      : NULL
##   ..$ colour    : NULL
##   ..$ size      : NULL
##   ..$ hjust     : num 0
##   ..$ vjust     : num 1
##   ..$ angle     : NULL
##   ..$ lineheight: NULL
##   ..$ margin    : 'margin' num [1:4] 0pt 0pt 5.5pt 0pt
##   .. ..- attr(*, "valid.unit")= int 8
##   .. ..- attr(*, "unit")= chr "pt"
##   ..$ debug     : NULL
##   ..$ inherit.blank: logi TRUE
##   ..- attr(*, "class")= chr [1:2] "element_text" "element"
```

```
##  $ plot.caption             :List of 11
##   ..$ family      : NULL
##   ..$ face        : NULL
##   ..$ colour      : NULL
##   ..$ size        : 'rel' num 0.8
##   ..$ hjust       : num 1
##   ..$ vjust       : num 1
##   ..$ angle       : NULL
##   ..$ lineheight  : NULL
##   ..$ margin      : 'margin' num [1:4] 5.5pt 0pt 0pt 0pt
##   .. ..- attr(*, "valid.unit")= int 8
##   .. ..- attr(*, "unit")= chr "pt"
##   ..$ debug       : NULL
##   ..$ inherit.blank: logi TRUE
##   ..- attr(*, "class")= chr [1:2] "element_text" "element"
##  $ plot.caption.position    : chr "panel"
##  $ plot.tag                 :List of 11
##   ..$ family      : NULL
##   ..$ face        : NULL
##   ..$ colour      : NULL
##   ..$ size        : 'rel' num 1.2
##   ..$ hjust       : num 0.5
##   ..$ vjust       : num 0.5
##   ..$ angle       : NULL
##   ..$ lineheight  : NULL
##   ..$ margin      : NULL
##   ..$ debug       : NULL
##   ..$ inherit.blank: logi TRUE
##   ..- attr(*, "class")= chr [1:2] "element_text" "element"
##  $ plot.tag.position        : chr "topleft"
##  $ plot.margin              : 'margin' num [1:4] 5.5pt 5.5pt 5.5pt 5.5pt
##   ..- attr(*, "valid.unit")= int 8
##   ..- attr(*, "unit")= chr "pt"
##  $ strip.background         :List of 5
##   ..$ fill        : chr "grey85"
##   ..$ colour      : chr "grey20"
##   ..$ size        : NULL
##   ..$ linetype    : NULL
##   ..$ inherit.blank: logi TRUE
##   ..- attr(*, "class")= chr [1:2] "element_rect" "element"
##  $ strip.background.x       : NULL
##  $ strip.background.y       : NULL
##  $ strip.placement          : chr "inside"
##  $ strip.text               :List of 11
##   ..$ family      : NULL
##   ..$ face        : NULL
##   ..$ colour      : chr "grey10"
##   ..$ size        : 'rel' num 0.8
##   ..$ hjust       : NULL
##   ..$ vjust       : NULL
##   ..$ angle       : NULL
##   ..$ lineheight  : NULL
##   ..$ margin      : 'margin' num [1:4] 4.4pt 4.4pt 4.4pt 4.4pt
##   .. ..- attr(*, "valid.unit")= int 8
```
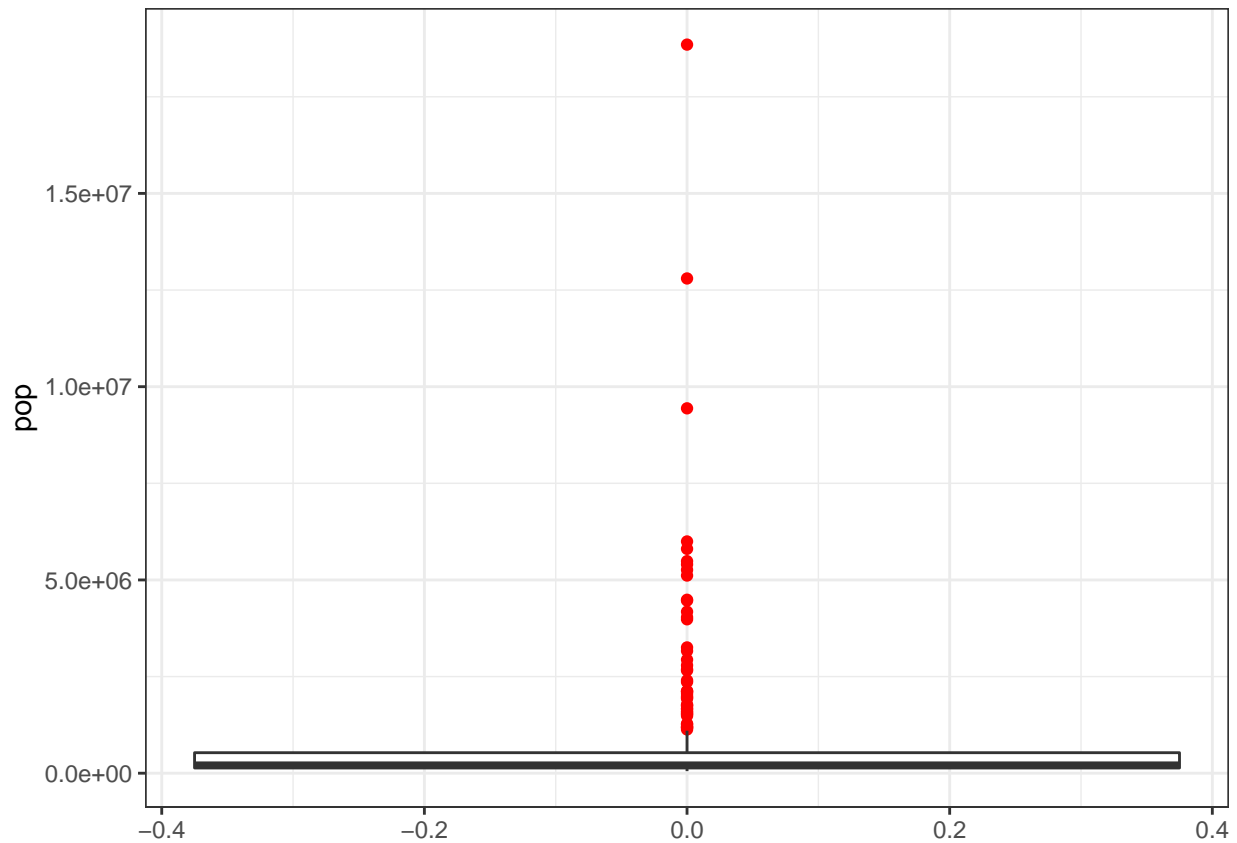
```
##     .. ..- attr(*, "unit")= chr "pt"
##     ..$ debug        : NULL
##     ..$ inherit.blank: logi TRUE
##     ..- attr(*, "class")= chr [1:2] "element_text" "element"
##   $ strip.text.x               : NULL
##   $ strip.text.y               :List of 11
##     ..$ family       : NULL
##     ..$ face         : NULL
##     ..$ colour       : NULL
##     ..$ size         : NULL
##     ..$ hjust        : NULL
##     ..$ vjust        : NULL
##     ..$ angle        : num -90
##     ..$ lineheight   : NULL
##     ..$ margin       : NULL
##     ..$ debug        : NULL
##     ..$ inherit.blank: logi TRUE
##     ..- attr(*, "class")= chr [1:2] "element_text" "element"
##   $ strip.switch.pad.grid     : 'unit' num 2.75pt
##     ..- attr(*, "valid.unit")= int 8
##     ..- attr(*, "unit")= chr "pt"
##   $ strip.switch.pad.wrap     : 'unit' num 2.75pt
##     ..- attr(*, "valid.unit")= int 8
##     ..- attr(*, "unit")= chr "pt"
##   $ strip.text.y.left         :List of 11
##     ..$ family       : NULL
##     ..$ face         : NULL
##     ..$ colour       : NULL
##     ..$ size         : NULL
##     ..$ hjust        : NULL
##     ..$ vjust        : NULL
##     ..$ angle        : num 90
##     ..$ lineheight   : NULL
##     ..$ margin       : NULL
##     ..$ debug        : NULL
##     ..$ inherit.blank: logi TRUE
##     ..- attr(*, "class")= chr [1:2] "element_text" "element"
##   - attr(*, "class")= chr [1:2] "theme" "gg"
##   - attr(*, "complete")= logi TRUE
##   - attr(*, "validate")= logi TRUE
```
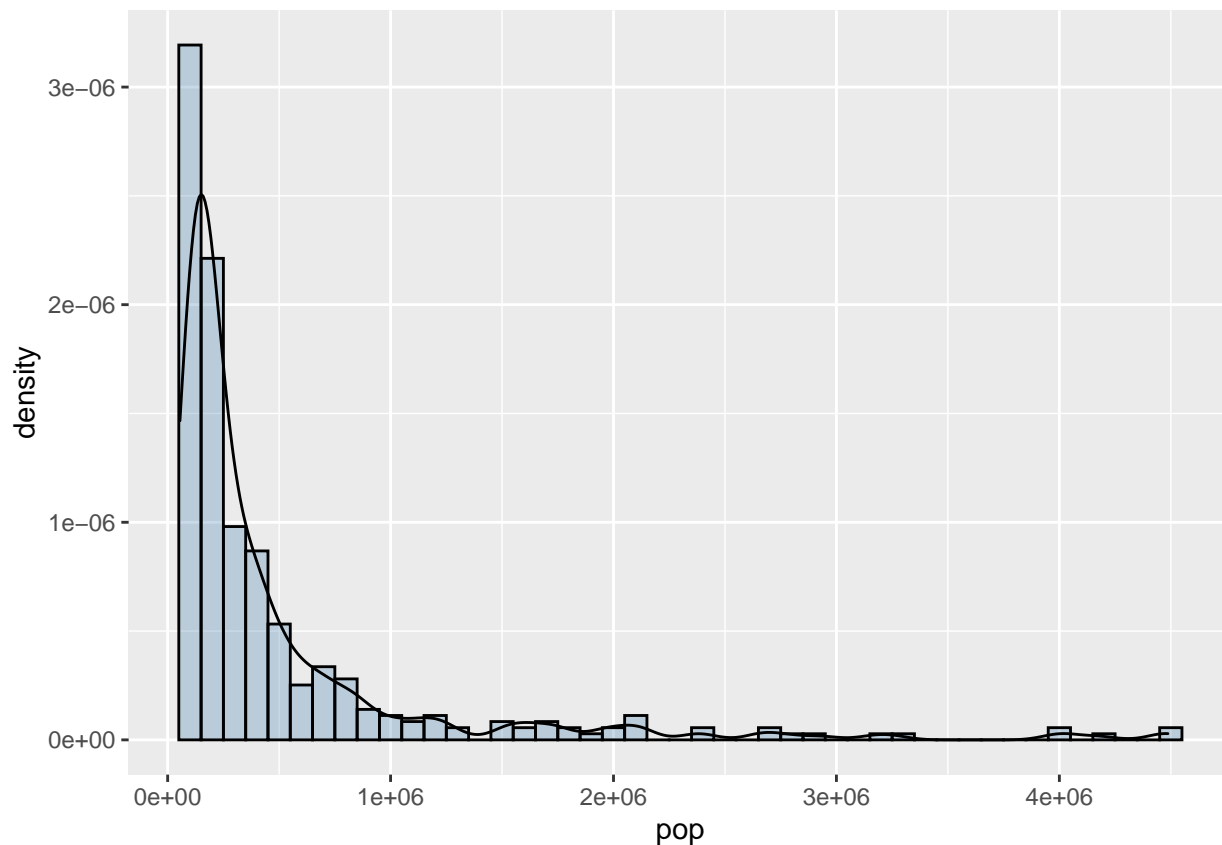
Note that with this plot we can see that outliers seem to be present. You can always use a boxplot to have evidence of it.

```
ggplot(data, aes(y = pop)) +
  geom_boxplot(outlier.colour = "red") +
  theme_bw()
```

You can then reproduce the previous plot without the outliers:

```
ggplot(data[data$pop < 5000000 ,], aes(x = pop)) +
 geom_histogram(aes(y=..density..), alpha=0.3,
                position="identity", binwidth = 100000, fill = "steelblue", color = "black")+
 geom_density(alpha=0.9)
```

```
theme_bw()
```

```
## List of 93
##  $ line                    :List of 6
##   ..$ colour      : chr "black"
##   ..$ size        : num 0.5
##   ..$ linetype    : num 1
##   ..$ lineend     : chr "butt"
##   ..$ arrow       : logi FALSE
##   ..$ inherit.blank: logi TRUE
##   ..- attr(*, "class")= chr [1:2] "element_line" "element"
##  $ rect                    :List of 5
##   ..$ fill        : chr "white"
##   ..$ colour      : chr "black"
##   ..$ size        : num 0.5
##   ..$ linetype    : num 1
##   ..$ inherit.blank: logi TRUE
##   ..- attr(*, "class")= chr [1:2] "element_rect" "element"
##  $ text                    :List of 11
##   ..$ family      : chr ""
##   ..$ face        : chr "plain"
##   ..$ colour      : chr "black"
##   ..$ size        : num 11
##   ..$ hjust       : num 0.5
##   ..$ vjust       : num 0.5
##   ..$ angle       : num 0
##   ..$ lineheight  : num 0.9
```

```
##   ..$ margin        : 'margin' num [1:4] 0pt 0pt 0pt 0pt
##   .. ..- attr(*, "valid.unit")= int 8
##   .. ..- attr(*, "unit")= chr "pt"
##   ..$ debug         : logi FALSE
##   ..$ inherit.blank: logi TRUE
##   ..- attr(*, "class")= chr [1:2] "element_text" "element"
## $ title                    : NULL
## $ aspect.ratio             : NULL
## $ axis.title               : NULL
## $ axis.title.x             :List of 11
##   ..$ family      : NULL
##   ..$ face        : NULL
##   ..$ colour      : NULL
##   ..$ size        : NULL
##   ..$ hjust       : NULL
##   ..$ vjust       : num 1
##   ..$ angle       : NULL
##   ..$ lineheight  : NULL
##   ..$ margin      : 'margin' num [1:4] 2.75pt 0pt 0pt 0pt
##   .. ..- attr(*, "valid.unit")= int 8
##   .. ..- attr(*, "unit")= chr "pt"
##   ..$ debug       : NULL
##   ..$ inherit.blank: logi TRUE
##   ..- attr(*, "class")= chr [1:2] "element_text" "element"
## $ axis.title.x.top         :List of 11
##   ..$ family      : NULL
##   ..$ face        : NULL
##   ..$ colour      : NULL
##   ..$ size        : NULL
##   ..$ hjust       : NULL
##   ..$ vjust       : num 0
##   ..$ angle       : NULL
##   ..$ lineheight  : NULL
##   ..$ margin      : 'margin' num [1:4] 0pt 0pt 2.75pt 0pt
##   .. ..- attr(*, "valid.unit")= int 8
##   .. ..- attr(*, "unit")= chr "pt"
##   ..$ debug       : NULL
##   ..$ inherit.blank: logi TRUE
##   ..- attr(*, "class")= chr [1:2] "element_text" "element"
## $ axis.title.x.bottom      : NULL
## $ axis.title.y             :List of 11
##   ..$ family      : NULL
##   ..$ face        : NULL
##   ..$ colour      : NULL
##   ..$ size        : NULL
##   ..$ hjust       : NULL
##   ..$ vjust       : num 1
##   ..$ angle       : num 90
##   ..$ lineheight  : NULL
##   ..$ margin      : 'margin' num [1:4] 0pt 2.75pt 0pt 0pt
##   .. ..- attr(*, "valid.unit")= int 8
##   .. ..- attr(*, "unit")= chr "pt"
##   ..$ debug       : NULL
##   ..$ inherit.blank: logi TRUE
```

```
##   ..- attr(*, "class")= chr [1:2] "element_text" "element"
## $ axis.title.y.left        : NULL
## $ axis.title.y.right       :List of 11
##   ..$ family     : NULL
##   ..$ face       : NULL
##   ..$ colour     : NULL
##   ..$ size       : NULL
##   ..$ hjust      : NULL
##   ..$ vjust      : num 0
##   ..$ angle      : num -90
##   ..$ lineheight : NULL
##   ..$ margin     : 'margin' num [1:4] 0pt 0pt 0pt 2.75pt
##   .. ..- attr(*, "valid.unit")= int 8
##   .. ..- attr(*, "unit")= chr "pt"
##   ..$ debug      : NULL
##   ..$ inherit.blank: logi TRUE
##   ..- attr(*, "class")= chr [1:2] "element_text" "element"
## $ axis.text                :List of 11
##   ..$ family     : NULL
##   ..$ face       : NULL
##   ..$ colour     : chr "grey30"
##   ..$ size       : 'rel' num 0.8
##   ..$ hjust      : NULL
##   ..$ vjust      : NULL
##   ..$ angle      : NULL
##   ..$ lineheight : NULL
##   ..$ margin     : NULL
##   ..$ debug      : NULL
##   ..$ inherit.blank: logi TRUE
##   ..- attr(*, "class")= chr [1:2] "element_text" "element"
## $ axis.text.x              :List of 11
##   ..$ family     : NULL
##   ..$ face       : NULL
##   ..$ colour     : NULL
##   ..$ size       : NULL
##   ..$ hjust      : NULL
##   ..$ vjust      : num 1
##   ..$ angle      : NULL
##   ..$ lineheight : NULL
##   ..$ margin     : 'margin' num [1:4] 2.2pt 0pt 0pt 0pt
##   .. ..- attr(*, "valid.unit")= int 8
##   .. ..- attr(*, "unit")= chr "pt"
##   ..$ debug      : NULL
##   ..$ inherit.blank: logi TRUE
##   ..- attr(*, "class")= chr [1:2] "element_text" "element"
## $ axis.text.x.top          :List of 11
##   ..$ family     : NULL
##   ..$ face       : NULL
##   ..$ colour     : NULL
##   ..$ size       : NULL
##   ..$ hjust      : NULL
##   ..$ vjust      : num 0
##   ..$ angle      : NULL
##   ..$ lineheight : NULL
```

```
##   ..$ margin        : 'margin' num [1:4] 0pt 0pt 2.2pt 0pt
##   .. ..- attr(*, "valid.unit")= int 8
##   .. ..- attr(*, "unit")= chr "pt"
##   ..$ debug         : NULL
##   ..$ inherit.blank: logi TRUE
##   ..- attr(*, "class")= chr [1:2] "element_text" "element"
## $ axis.text.x.bottom      : NULL
## $ axis.text.y             :List of 11
##   ..$ family       : NULL
##   ..$ face         : NULL
##   ..$ colour       : NULL
##   ..$ size         : NULL
##   ..$ hjust        : num 1
##   ..$ vjust        : NULL
##   ..$ angle        : NULL
##   ..$ lineheight   : NULL
##   ..$ margin       : 'margin' num [1:4] 0pt 2.2pt 0pt 0pt
##   .. ..- attr(*, "valid.unit")= int 8
##   .. ..- attr(*, "unit")= chr "pt"
##   ..$ debug        : NULL
##   ..$ inherit.blank: logi TRUE
##   ..- attr(*, "class")= chr [1:2] "element_text" "element"
## $ axis.text.y.left        : NULL
## $ axis.text.y.right       :List of 11
##   ..$ family       : NULL
##   ..$ face         : NULL
##   ..$ colour       : NULL
##   ..$ size         : NULL
##   ..$ hjust        : num 0
##   ..$ vjust        : NULL
##   ..$ angle        : NULL
##   ..$ lineheight   : NULL
##   ..$ margin       : 'margin' num [1:4] 0pt 0pt 0pt 2.2pt
##   .. ..- attr(*, "valid.unit")= int 8
##   .. ..- attr(*, "unit")= chr "pt"
##   ..$ debug        : NULL
##   ..$ inherit.blank: logi TRUE
##   ..- attr(*, "class")= chr [1:2] "element_text" "element"
## $ axis.ticks              :List of 6
##   ..$ colour       : chr "grey20"
##   ..$ size         : NULL
##   ..$ linetype     : NULL
##   ..$ lineend      : NULL
##   ..$ arrow        : logi FALSE
##   ..$ inherit.blank: logi TRUE
##   ..- attr(*, "class")= chr [1:2] "element_line" "element"
## $ axis.ticks.x            : NULL
## $ axis.ticks.x.top        : NULL
## $ axis.ticks.x.bottom     : NULL
## $ axis.ticks.y            : NULL
## $ axis.ticks.y.left       : NULL
## $ axis.ticks.y.right      : NULL
## $ axis.ticks.length       : 'unit' num 2.75pt
##   ..- attr(*, "valid.unit")= int 8
```

```
##    ..- attr(*, "unit")= chr "pt"
## $ axis.ticks.length.x        : NULL
## $ axis.ticks.length.x.top    : NULL
## $ axis.ticks.length.x.bottom: NULL
## $ axis.ticks.length.y        : NULL
## $ axis.ticks.length.y.left   : NULL
## $ axis.ticks.length.y.right  : NULL
## $ axis.line                  : list()
##    ..- attr(*, "class")= chr [1:2] "element_blank" "element"
## $ axis.line.x                : NULL
## $ axis.line.x.top            : NULL
## $ axis.line.x.bottom         : NULL
## $ axis.line.y                : NULL
## $ axis.line.y.left           : NULL
## $ axis.line.y.right          : NULL
## $ legend.background          :List of 5
##    ..$ fill         : NULL
##    ..$ colour       : logi NA
##    ..$ size         : NULL
##    ..$ linetype     : NULL
##    ..$ inherit.blank: logi TRUE
##    ..- attr(*, "class")= chr [1:2] "element_rect" "element"
## $ legend.margin              : 'margin' num [1:4] 5.5pt 5.5pt 5.5pt 5.5pt
##    ..- attr(*, "valid.unit")= int 8
##    ..- attr(*, "unit")= chr "pt"
## $ legend.spacing             : 'unit' num 11pt
##    ..- attr(*, "valid.unit")= int 8
##    ..- attr(*, "unit")= chr "pt"
## $ legend.spacing.x           : NULL
## $ legend.spacing.y           : NULL
## $ legend.key                 :List of 5
##    ..$ fill         : chr "white"
##    ..$ colour       : logi NA
##    ..$ size         : NULL
##    ..$ linetype     : NULL
##    ..$ inherit.blank: logi TRUE
##    ..- attr(*, "class")= chr [1:2] "element_rect" "element"
## $ legend.key.size            : 'unit' num 1.2lines
##    ..- attr(*, "valid.unit")= int 3
##    ..- attr(*, "unit")= chr "lines"
## $ legend.key.height          : NULL
## $ legend.key.width           : NULL
## $ legend.text                :List of 11
##    ..$ family       : NULL
##    ..$ face         : NULL
##    ..$ colour       : NULL
##    ..$ size         : 'rel' num 0.8
##    ..$ hjust        : NULL
##    ..$ vjust        : NULL
##    ..$ angle        : NULL
##    ..$ lineheight   : NULL
##    ..$ margin       : NULL
##    ..$ debug        : NULL
##    ..$ inherit.blank: logi TRUE
```

```
##    ..- attr(*, "class")= chr [1:2] "element_text" "element"
## $ legend.text.align        : NULL
## $ legend.title             :List of 11
##   ..$ family      : NULL
##   ..$ face        : NULL
##   ..$ colour      : NULL
##   ..$ size        : NULL
##   ..$ hjust       : num 0
##   ..$ vjust       : NULL
##   ..$ angle       : NULL
##   ..$ lineheight  : NULL
##   ..$ margin      : NULL
##   ..$ debug       : NULL
##   ..$ inherit.blank: logi TRUE
##   ..- attr(*, "class")= chr [1:2] "element_text" "element"
## $ legend.title.align       : NULL
## $ legend.position          : chr "right"
## $ legend.direction         : NULL
## $ legend.justification     : chr "center"
## $ legend.box               : NULL
## $ legend.box.just          : NULL
## $ legend.box.margin        : 'margin' num [1:4] 0cm 0cm 0cm 0cm
##   ..- attr(*, "valid.unit")= int 1
##   ..- attr(*, "unit")= chr "cm"
## $ legend.box.background     : list()
##   ..- attr(*, "class")= chr [1:2] "element_blank" "element"
## $ legend.box.spacing        : 'unit' num 11pt
##   ..- attr(*, "valid.unit")= int 8
##   ..- attr(*, "unit")= chr "pt"
## $ panel.background          :List of 5
##   ..$ fill        : chr "white"
##   ..$ colour      : logi NA
##   ..$ size        : NULL
##   ..$ linetype    : NULL
##   ..$ inherit.blank: logi TRUE
##   ..- attr(*, "class")= chr [1:2] "element_rect" "element"
## $ panel.border              :List of 5
##   ..$ fill        : logi NA
##   ..$ colour      : chr "grey20"
##   ..$ size        : NULL
##   ..$ linetype    : NULL
##   ..$ inherit.blank: logi TRUE
##   ..- attr(*, "class")= chr [1:2] "element_rect" "element"
## $ panel.spacing             : 'unit' num 5.5pt
##   ..- attr(*, "valid.unit")= int 8
##   ..- attr(*, "unit")= chr "pt"
## $ panel.spacing.x           : NULL
## $ panel.spacing.y           : NULL
## $ panel.grid                :List of 6
##   ..$ colour      : chr "grey92"
##   ..$ size        : NULL
##   ..$ linetype    : NULL
##   ..$ lineend     : NULL
##   ..$ arrow       : logi FALSE
```

```
##   ..$ inherit.blank: logi TRUE
##   ..- attr(*, "class")= chr [1:2] "element_line" "element"
## $ panel.grid.major        : NULL
## $ panel.grid.minor        :List of 6
##   ..$ colour     : NULL
##   ..$ size       : 'rel' num 0.5
##   ..$ linetype   : NULL
##   ..$ lineend    : NULL
##   ..$ arrow      : logi FALSE
##   ..$ inherit.blank: logi TRUE
##   ..- attr(*, "class")= chr [1:2] "element_line" "element"
## $ panel.grid.major.x      : NULL
## $ panel.grid.major.y      : NULL
## $ panel.grid.minor.x      : NULL
## $ panel.grid.minor.y      : NULL
## $ panel.ontop             : logi FALSE
## $ plot.background         :List of 5
##   ..$ fill       : NULL
##   ..$ colour     : chr "white"
##   ..$ size       : NULL
##   ..$ linetype   : NULL
##   ..$ inherit.blank: logi TRUE
##   ..- attr(*, "class")= chr [1:2] "element_rect" "element"
## $ plot.title              :List of 11
##   ..$ family     : NULL
##   ..$ face       : NULL
##   ..$ colour     : NULL
##   ..$ size       : 'rel' num 1.2
##   ..$ hjust      : num 0
##   ..$ vjust      : num 1
##   ..$ angle      : NULL
##   ..$ lineheight : NULL
##   ..$ margin     : 'margin' num [1:4] 0pt 0pt 5.5pt 0pt
##   .. ..- attr(*, "valid.unit")= int 8
##   .. ..- attr(*, "unit")= chr "pt"
##   ..$ debug      : NULL
##   ..$ inherit.blank: logi TRUE
##   ..- attr(*, "class")= chr [1:2] "element_text" "element"
## $ plot.title.position     : chr "panel"
## $ plot.subtitle           :List of 11
##   ..$ family     : NULL
##   ..$ face       : NULL
##   ..$ colour     : NULL
##   ..$ size       : NULL
##   ..$ hjust      : num 0
##   ..$ vjust      : num 1
##   ..$ angle      : NULL
##   ..$ lineheight : NULL
##   ..$ margin     : 'margin' num [1:4] 0pt 0pt 5.5pt 0pt
##   .. ..- attr(*, "valid.unit")= int 8
##   .. ..- attr(*, "unit")= chr "pt"
##   ..$ debug      : NULL
##   ..$ inherit.blank: logi TRUE
##   ..- attr(*, "class")= chr [1:2] "element_text" "element"
```

```
##  $ plot.caption              :List of 11
##   ..$ family      : NULL
##   ..$ face        : NULL
##   ..$ colour      : NULL
##   ..$ size        : 'rel' num 0.8
##   ..$ hjust       : num 1
##   ..$ vjust       : num 1
##   ..$ angle       : NULL
##   ..$ lineheight  : NULL
##   ..$ margin      : 'margin' num [1:4] 5.5pt 0pt 0pt 0pt
##   .. ..- attr(*, "valid.unit")= int 8
##   .. ..- attr(*, "unit")= chr "pt"
##   ..$ debug       : NULL
##   ..$ inherit.blank: logi TRUE
##   ..- attr(*, "class")= chr [1:2] "element_text" "element"
##  $ plot.caption.position     : chr "panel"
##  $ plot.tag                  :List of 11
##   ..$ family      : NULL
##   ..$ face        : NULL
##   ..$ colour      : NULL
##   ..$ size        : 'rel' num 1.2
##   ..$ hjust       : num 0.5
##   ..$ vjust       : num 0.5
##   ..$ angle       : NULL
##   ..$ lineheight  : NULL
##   ..$ margin      : NULL
##   ..$ debug       : NULL
##   ..$ inherit.blank: logi TRUE
##   ..- attr(*, "class")= chr [1:2] "element_text" "element"
##  $ plot.tag.position         : chr "topleft"
##  $ plot.margin               : 'margin' num [1:4] 5.5pt 5.5pt 5.5pt 5.5pt
##   ..- attr(*, "valid.unit")= int 8
##   ..- attr(*, "unit")= chr "pt"
##  $ strip.background          :List of 5
##   ..$ fill        : chr "grey85"
##   ..$ colour      : chr "grey20"
##   ..$ size        : NULL
##   ..$ linetype    : NULL
##   ..$ inherit.blank: logi TRUE
##   ..- attr(*, "class")= chr [1:2] "element_rect" "element"
##  $ strip.background.x        : NULL
##  $ strip.background.y        : NULL
##  $ strip.placement           : chr "inside"
##  $ strip.text                :List of 11
##   ..$ family      : NULL
##   ..$ face        : NULL
##   ..$ colour      : chr "grey10"
##   ..$ size        : 'rel' num 0.8
##   ..$ hjust       : NULL
##   ..$ vjust       : NULL
##   ..$ angle       : NULL
##   ..$ lineheight  : NULL
##   ..$ margin      : 'margin' num [1:4] 4.4pt 4.4pt 4.4pt 4.4pt
##   .. ..- attr(*, "valid.unit")= int 8
```
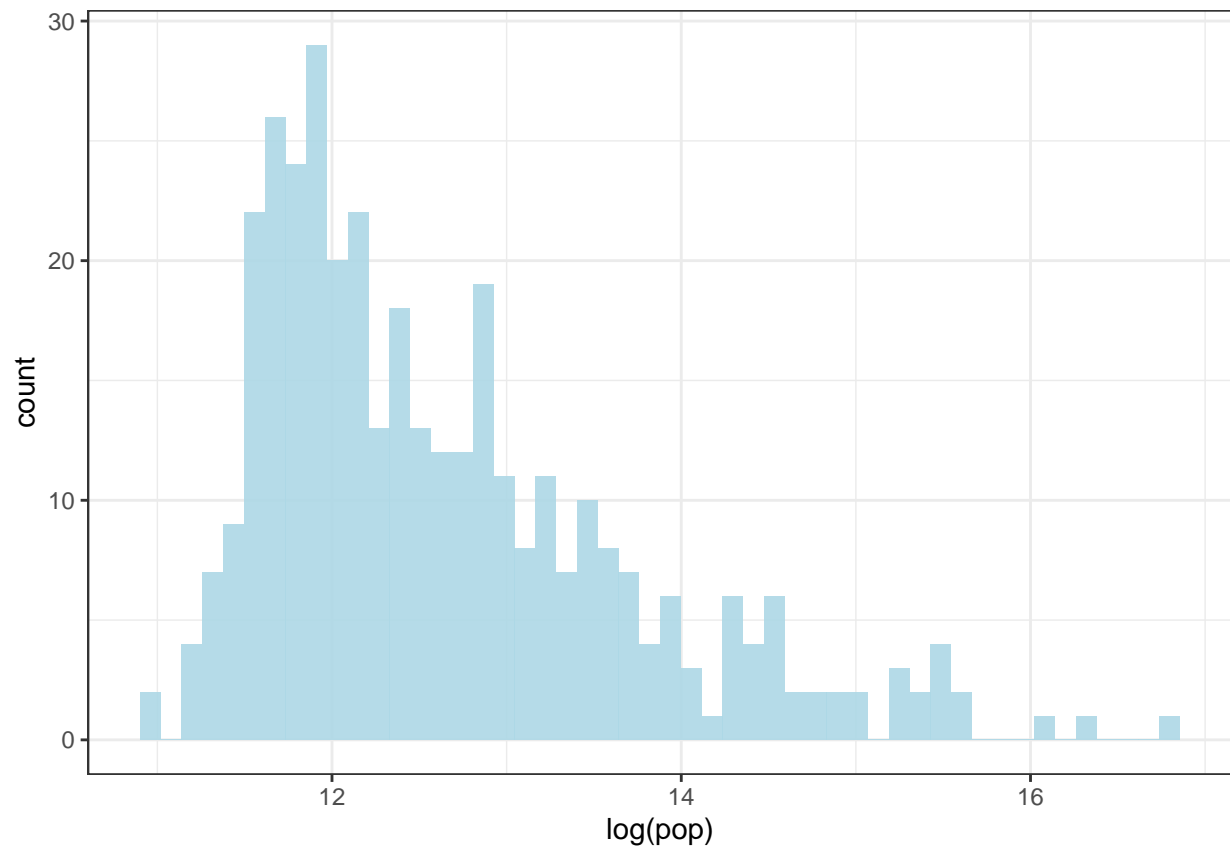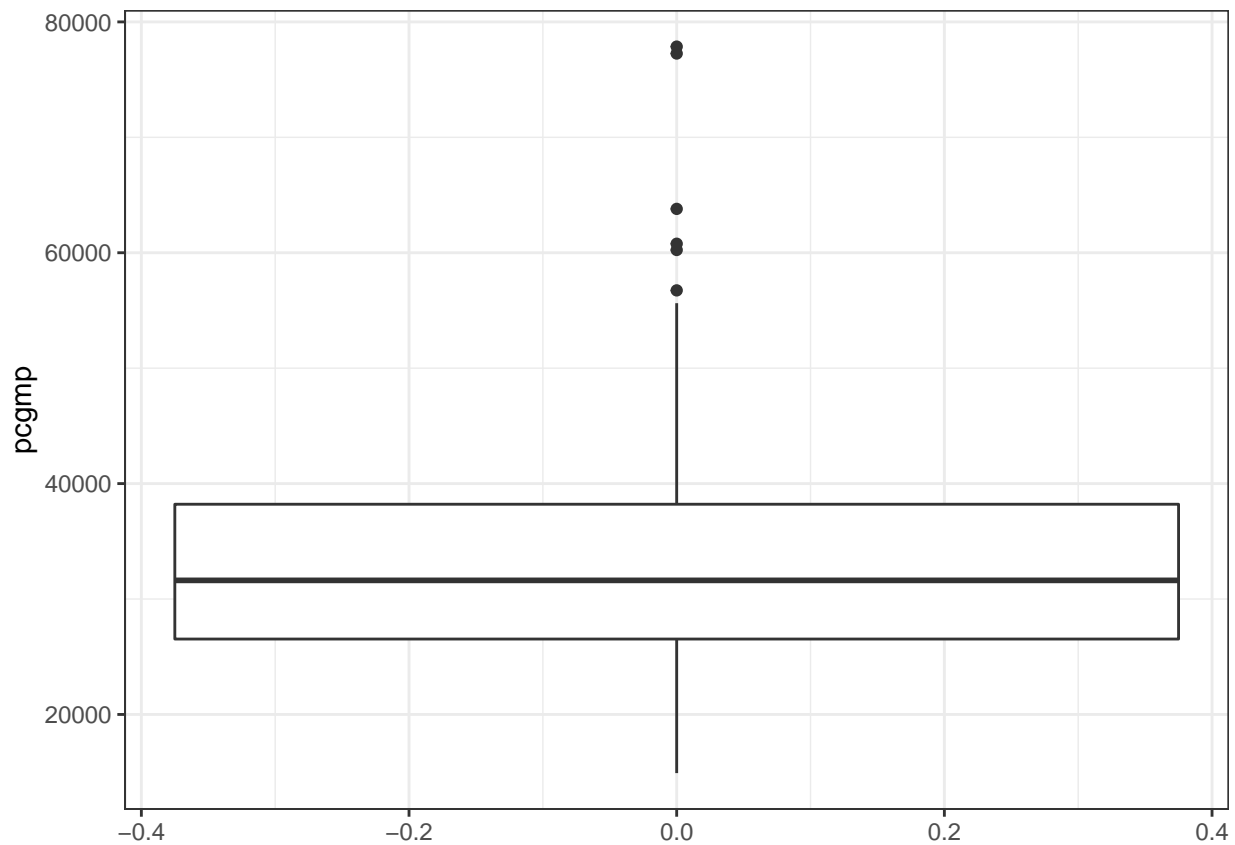
```
##   .. ..- attr(*, "unit")= chr "pt"
##   ..$ debug        : NULL
##   ..$ inherit.blank: logi TRUE
##   ..- attr(*, "class")= chr [1:2] "element_text" "element"
## $ strip.text.x            : NULL
## $ strip.text.y            :List of 11
##   ..$ family       : NULL
##   ..$ face         : NULL
##   ..$ colour       : NULL
##   ..$ size         : NULL
##   ..$ hjust        : NULL
##   ..$ vjust        : NULL
##   ..$ angle        : num -90
##   ..$ lineheight   : NULL
##   ..$ margin       : NULL
##   ..$ debug        : NULL
##   ..$ inherit.blank: logi TRUE
##   ..- attr(*, "class")= chr [1:2] "element_text" "element"
## $ strip.switch.pad.grid   : 'unit' num 2.75pt
##   ..- attr(*, "valid.unit")= int 8
##   ..- attr(*, "unit")= chr "pt"
## $ strip.switch.pad.wrap   : 'unit' num 2.75pt
##   ..- attr(*, "valid.unit")= int 8
##   ..- attr(*, "unit")= chr "pt"
## $ strip.text.y.left       :List of 11
##   ..$ family       : NULL
##   ..$ face         : NULL
##   ..$ colour       : NULL
##   ..$ size         : NULL
##   ..$ hjust        : NULL
##   ..$ vjust        : NULL
##   ..$ angle        : num 90
##   ..$ lineheight   : NULL
##   ..$ margin       : NULL
##   ..$ debug        : NULL
##   ..$ inherit.blank: logi TRUE
##   ..- attr(*, "class")= chr [1:2] "element_text" "element"
## - attr(*, "class")= chr [1:2] "theme" "gg"
## - attr(*, "complete")= logi TRUE
## - attr(*, "validate")= logi TRUE
```

```r
ggplot(data, aes(x = log(pop))) +
  geom_histogram(bins = 50, alpha = 0.9, fill = "lightblue") +
  theme_bw()
```

```
ggplot(data, aes(y = pcgmp)) +
  geom_boxplot() +
  theme_bw()
```
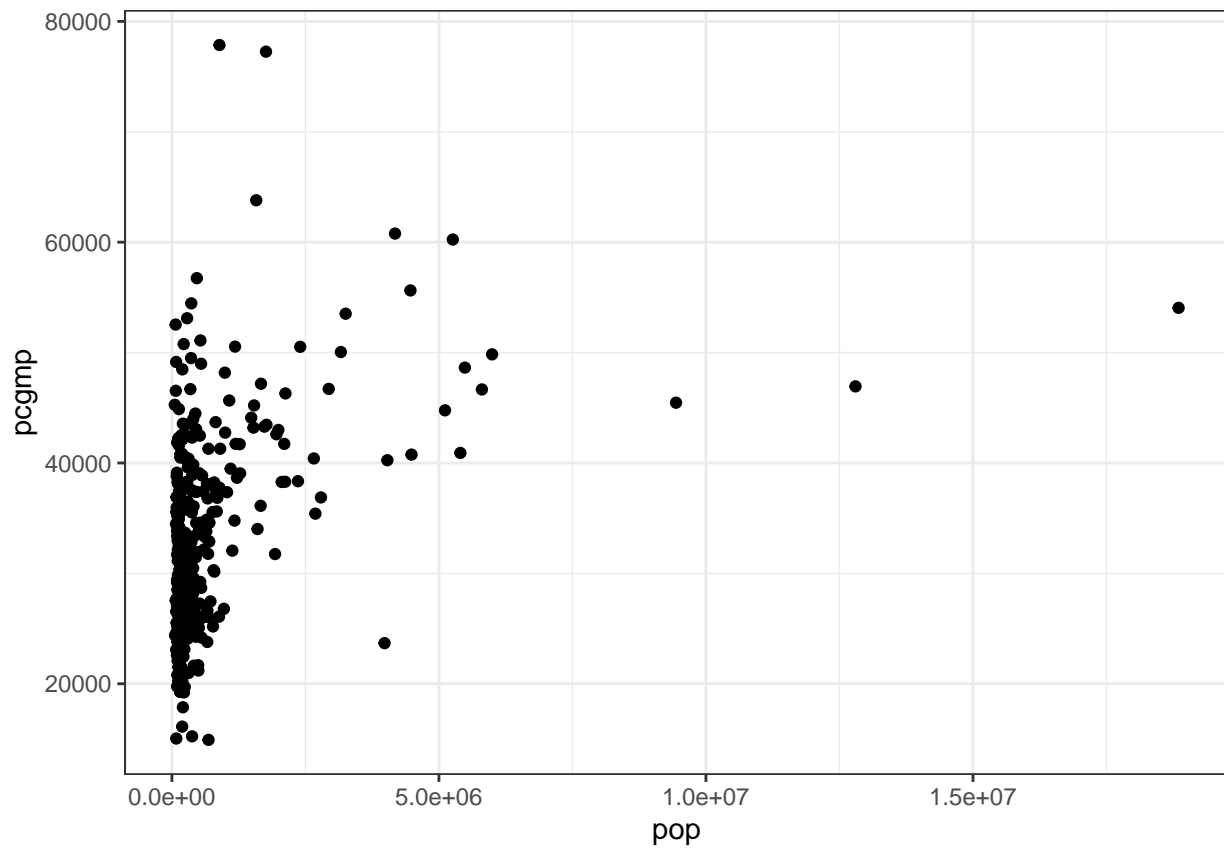
*End of solution 2*
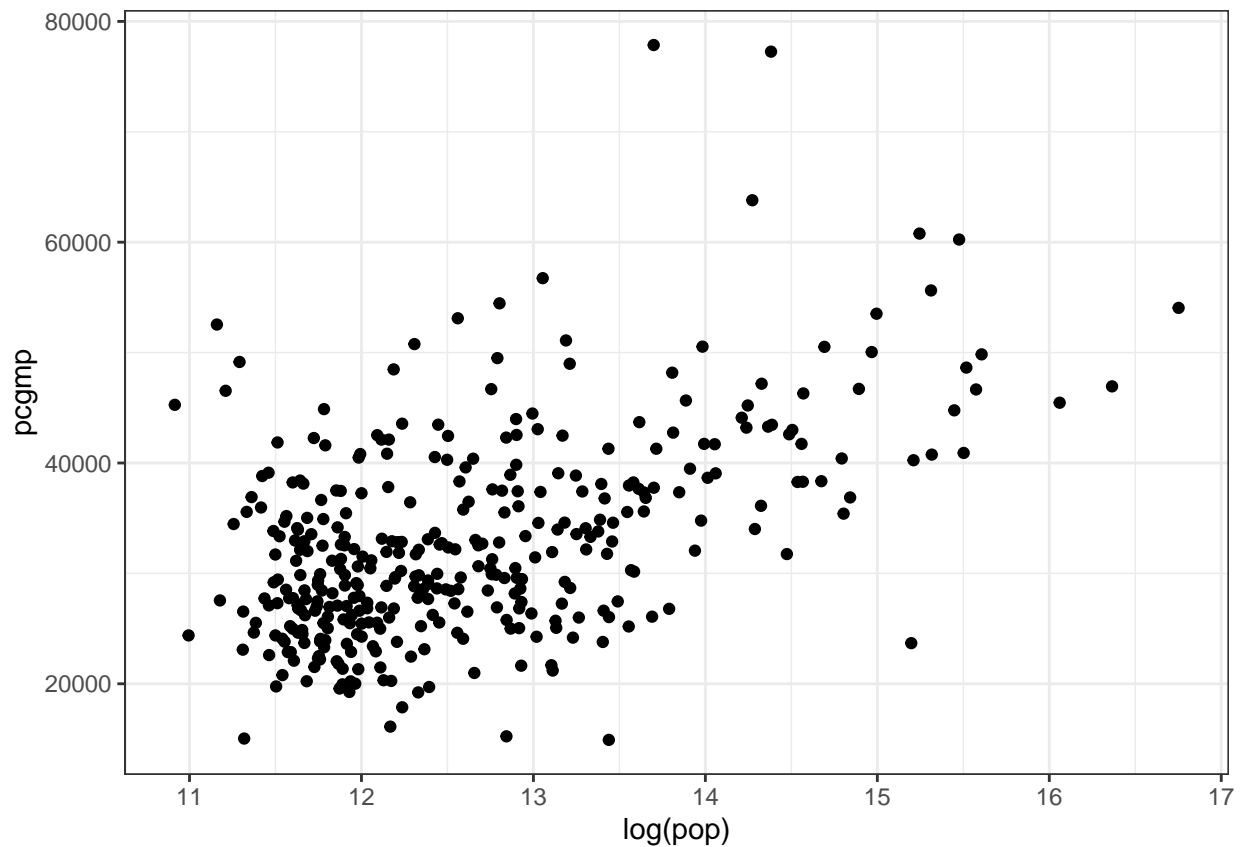
## Question 3: GMP and population

Make a bivariate plot for per-capita GMP as a function of population. Describe the relationship in words. You can also try with $log(pop)$.

*Solution 3*

```
ggplot(data, aes(x = pop, y = pcgmp)) +
  geom_point() +
  theme_bw()
```

```
ggplot(data, aes(x = log(pop), y = pcgmp)) +
  geom_point() +
  theme_bw()
```

*End of solution 3*

## Question 4: simple linear model

Considering your previous plot, run the `lm()` function on the data and add the regression to the previous plot. Will you use `pop` or `log(pop)`? (would the last one still be a linear model?)
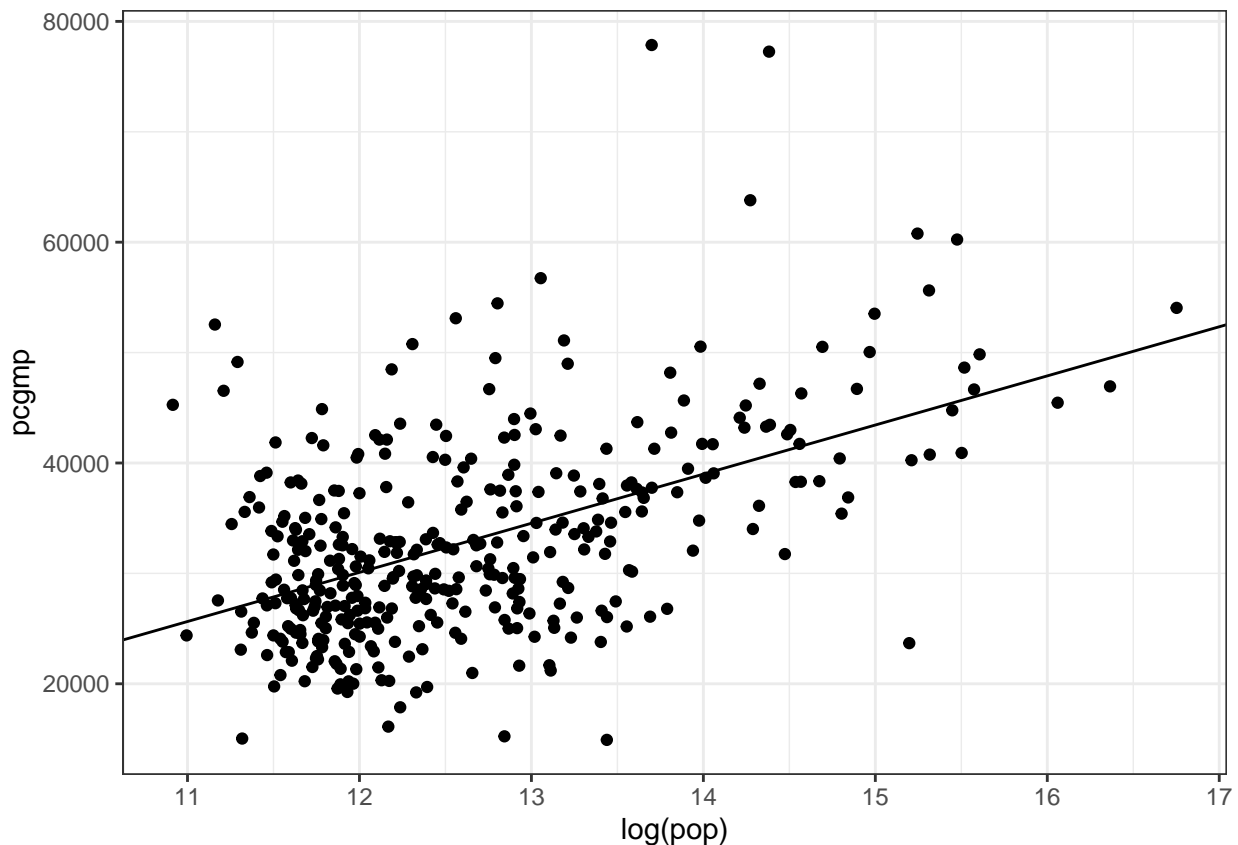
You can comment the result.

*Solution 4*

```
model = lm(pcgmp~log(pop), data = data)
print(model$coefficients)

## (Intercept)    log(pop)
##  -23306.199    4449.758

ggplot(data, aes(x = log(pop), y = pcgmp)) +
  geom_point() +
  theme_bw() +
  geom_abline(slope = model$coefficients[[2]], intercept = model$coefficients[[1]])
```

The fit does not seem to be a good one.

```
summary(model)
```

```
##
## Call:
## lm(formula = pcgmp ~ log(pop), data = data)
##
## Residuals:
##    Min     1Q Median     3Q    Max
## -21572  -4765  -1016   3686  40207
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -23306.2     4957.1  -4.702 3.67e-06 ***
## log(pop)      4449.8      390.9  11.383  < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 7929 on 364 degrees of freedom
## Multiple R-squared:  0.2625, Adjusted R-squared:  0.2605
## F-statistic: 129.6 on 1 and 364 DF,  p-value: < 2.2e-16
```

Note that it is not easy to interpret `log(pop)` in comparison with `pop`.

*End of solution 4*

# Question 5

Bonus question: could you do `log(pcgmp)` as a linear function of `pop` (or `log(pop)`)?

*Solution 5* This would be wrong! Because it implies a multiplication. Here to keep the model linear it should be $log(Y) = \beta X * \varepsilon$ *End of solution 5*