

Question 1. First, use the R command

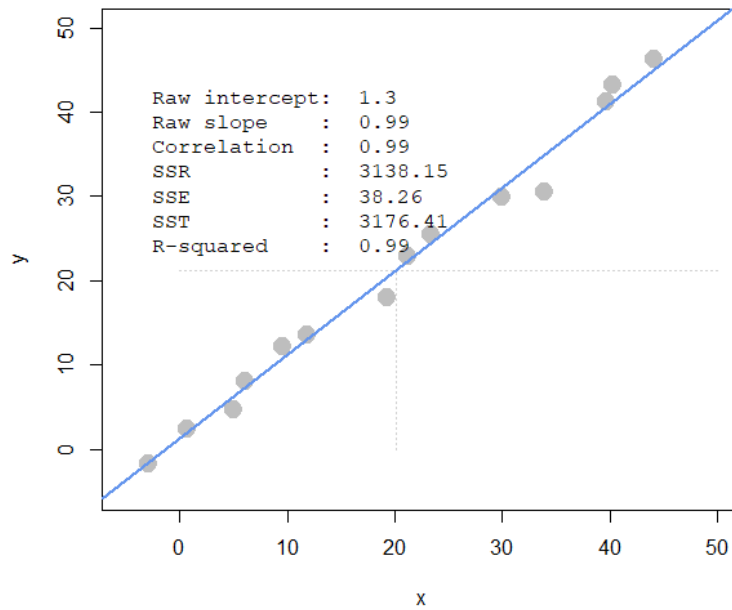
```
> source("demo_simple_regression_rsqr.R")
```

then, call the function

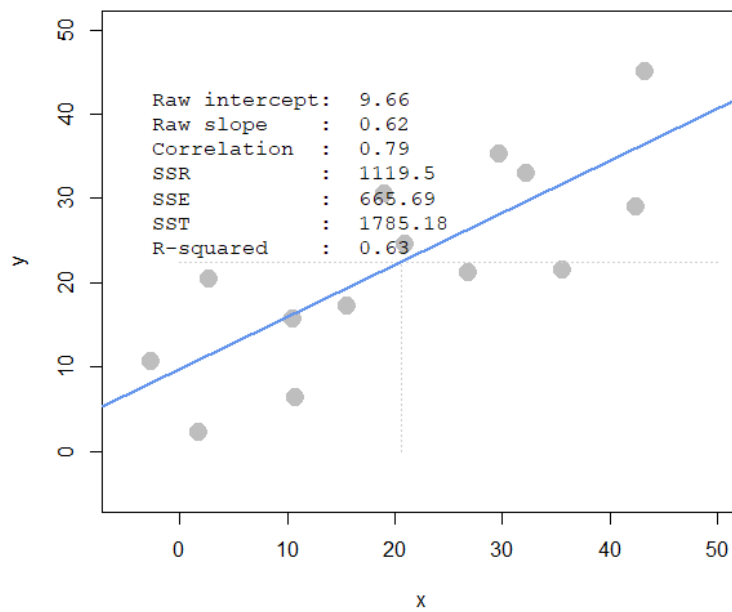
```
> interactive_regression_rsqr()
```

four times to generate four scenarios.

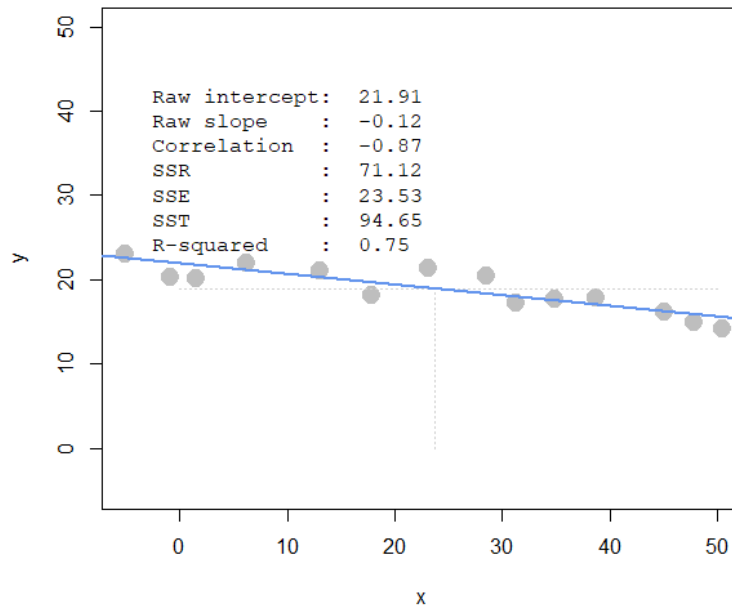
Scenario 1:



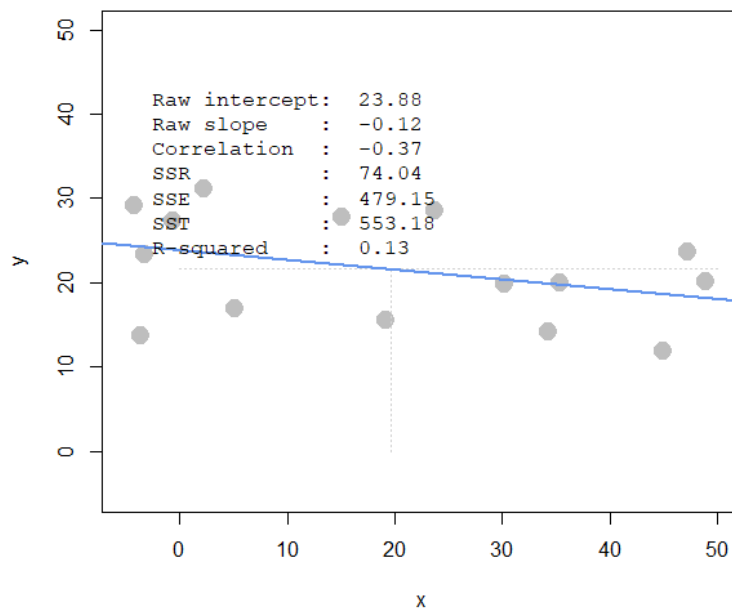
Scenario 2:



Scenario 3:



Scenario 4:



- (a) Scenario 1 has a stronger R^2 .
- (b) Scenario 3 has a stronger R^2 .
- (c) SSE: Scenario 2 > Scenario 1. SSR: Scenario 1 > Scenario 2. SST: Scenario 1 > Scenario 2.
- (d) SSE: Scenario 4 > Scenario 4. SSR: Scenario 4 > Scenario 4. SST: Scenario 4 > Scenario 4.

Question 2. (a) First, read the dataset and directly do the linear regression via `lm()` function..

```

1 # Question 2 (a)
2 salary <- read.csv("programmer_salaries.txt", sep="\t")
3 salary_regression <- lm(salary$Salary ~
4     salary$Experience +
5     salary$Score +
6     salary$Degree) # do linear regression
7 summary(salary_regression, data=salary)

```

Here's partial results. One can obtain R^2 and the first 5 values of \hat{y} and ϵ here.

```

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept)   7.9448    7.3808   1.076   0.2977
salary$Experience  1.1476    0.2976   3.856   0.0014 **
salary$Score    0.1969    0.0899   2.191   0.0436 *
salary$Degree   2.2804    1.9866   1.148   0.2679
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 2.396 on 16 degrees of freedom
Multiple R-squared:  0.8468,    Adjusted R-squared:  0.8181
F-statistic: 29.48 on 3 and 16 DF,  p-value: 9.417e-07

> head(salary_regression$fitted.values)
      1      2      3      4      5      6
27.89626 37.95204 26.02901 32.11201 36.34251 38.24380
> head(salary_regression$residuals)
      1      2      3      4      5      6
-3.8962605  5.0479568 -2.3290112  2.1879860 -0.5425072 -0.2437966

```

(b) The hand-craft linear regression code:

```

1 # Question 2 (b)
2 ones <- replicate(length(salary$Salary), 1) # create ones column vector
3 # Combine column vectors to a matrix
4 X <- cbind(ones, salary$Experience, salary$Score, salary$Degree)
5 y <- salary$Salary
6 beta_hat <- solve(t(X) %*% X) %*% t(X) %*% y # some linear algebra
7 y_hat <- X %*% beta_hat # predicted values
8 res <- y - y_hat # residuals
9 SSR <- sum((y_hat - mean(y))^2)
10 SSE <- sum((y - y_hat)^2)
11 SST <- SSR + SSE

```

(iii) $\hat{\beta} = (7.944849, 1.147582, 0.196937, 2.280424)^T$.

(iv) Use the command

`> head(y_hat)` and `> head(y_hat)`

Here's the result:

```

> head(y_hat) > head(res)
      [,1]      [,1]
[1,] 27.89626 [1,] -3.8962605
[2,] 37.95204 [2,]  5.0479568
[3,] 26.02901 [3,] -2.3290112
[4,] 32.11201 [4,]  2.1879860
[5,] 36.34251 [5,] -0.5425072
[6,] 38.24380 [6,] -0.2437966

```

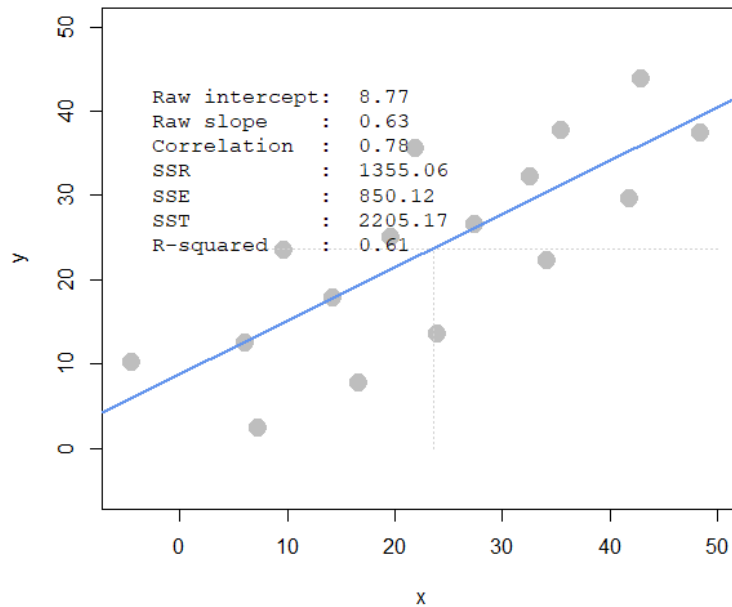
(v) $SSR = 507.896$, $SSE = 91.88949$, and $SST = 599.7855$.

(c) I generate another plot of scenario 2.

```

1 source("demo_simple_regression_rsqr.R")
2 points <- interactive_regression_rsqr()

```



Reuse the code from (b):

```

1 ones_c <- replicate(length(points[,1]), 1) # create ones
2 X_c <- cbind(ones_c, points[,1])
3 y_c <- points[,2]
4 beta_hat_c <- solve(t(X_c) %*% X_c) %*% t(X_c) %*% y_c
5 y_hat_c <- X_c %*% beta_hat_c
6 res_c <- y_c - y_hat_c
7 SSR_c <- sum((y_hat_c - mean(y_c))^2)
8 SSE_c <- sum((y_c - y_hat_c)^2)
9 SST_c <- SSR_c + SSE_c
10 R2_c_i <- SSR_c/SST_c
11 R2_c_ii <- cor(y_c, y_hat_c)^2

```

Both the two methods, say `R2_c_i` and `R2_c_ii` returns $R^2 = 0.6144897$. This meets the conclusion from the plot. ■

Question 3. Load the data.

```

1 # Question 3
2 auto <- read.table("auto-data.txt", header=FALSE, na.strings = "?")
3 names(auto) <- c("mpg", "cylinders", "displacement", "horsepower", "weight",
4                 "acceleration", "model_year", "origin", "car_name")

```

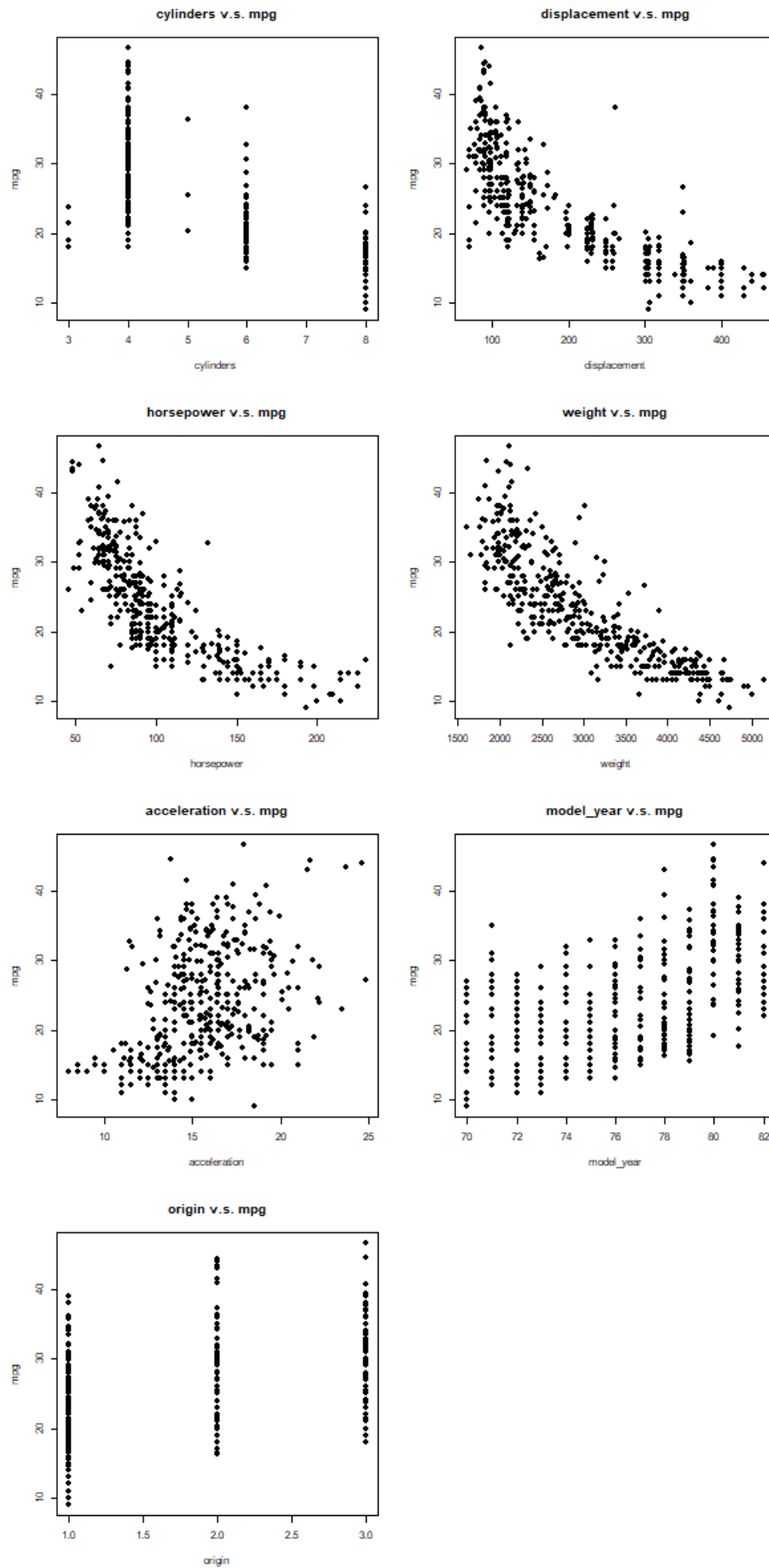
(a)(i) Plot the regression plot to all the other variables expect `car_name`.

```

1 # Question 3 (a-i)
2 png(filename = "3a.png", width = 600, height = 1200) # Subplots
3 par(mfrow=c(4,2))
4 cy <- plot(auto$cylinders, auto$mpg, # scatter
5           main="cylinders v.s. mpg",
6           xlab="cylinders ",
7           ylab="mpg",
8           pch=19)
9 di <- plot(auto$displacement, auto$mpg,
10          main="displacement v.s. mpg",

```

```
11         xlab="displacement",
12         ylab="mpg",
13         pch=19)
14 ho <- plot(auto$horsepower, auto$mpg,
15           main="horsepower v.s. mpg",
16           xlab="horsepower",
17           ylab="mpg",
18           pch=19)
19 we <- plot(auto$weight, auto$mpg,
20           main="weight v.s. mpg",
21           xlab="weight",
22           ylab="mpg",
23           pch=19)
24 ac <- plot(auto$acceleration, auto$mpg,
25           main="acceleration v.s. mpg",
26           xlab="acceleration",
27           ylab="mpg",
28           pch=19)
29 my <- plot(auto$model_year, auto$mpg,
30           main="model_year v.s. mpg",
31           xlab="model_year",
32           ylab="mpg",
33           pch=19)
34 or <- plot(auto$origin, auto$mpg,
35           main="origin v.s. mpg",
36           xlab="origin",
37           ylab="mpg",
38           pch=19)
39 dev.off()
```



(ii) Write the correlation matrix to a .csv file:

```

1 # Question 3 (a-ii)
2 cor_matrix <- cor(auto[,colnames(auto)!="car_name"], # drop column car_name
3               use="pairwise.complete.obs") # omit NA's

```

	mpg	cylinders	displacement	horsepower	weight	acceleration	model_year	origin
mpg	1	-0.78	-0.8	-0.78	-0.83	0.42	0.58	0.56
cylinders	-0.78	1	0.95	0.84	0.9	-0.51	-0.35	-0.56
displacement	-0.8	0.95	1	0.9	0.93	-0.54	-0.37	-0.61
horsepower	-0.78	0.84	0.9	1	0.86	-0.69	-0.42	-0.46
weight	-0.83	0.9	0.93	0.86	1	-0.42	-0.31	-0.58
acceleration	0.42	-0.51	-0.54	-0.69	-0.42	1	0.29	0.21
model_year	0.58	-0.35	-0.37	-0.42	-0.31	0.29	1	0.18
origin	0.56	-0.56	-0.61	-0.46	-0.58	0.21	0.18	1

```

4 cor_matrix <- round(cor_matrix, digits=2)
5 write.table(cor_matrix, file="3a.csv")

```

Then by some magic,

I'll answer (iii)-(v) at the same time. First, by (i) and (ii), I found that displacement, horsepower and weight seem to (highly) related to mpg. However these relations seems not linear. Also the variable cylinder has a high correlation ($r = -0.78$) to mpg. However based on the scatter plot, I think it should be viewed as a discrete type data. It is not suit for a linear regression model.

(b) Though some of the variable are not suit for linear regression, I still create a model.

```

1 # Question 3 (b)
2 auto_lr_model <- lm(mpg ~ cylinders+displacement+horsepower+
3                     weight+acceleration+model_year+factor(origin), auto)
4 summary(auto_lr_model)

```

Here's the partial results:

```

              Estimate Std. Error t value Pr(>|t|)
(Intercept) -1.795e+01  4.677e+00  -3.839 0.000145 ***
cylinders    -4.897e-01  3.212e-01  -1.524 0.128215
displacement  2.398e-02  7.653e-03   3.133 0.001863 **
horsepower   -1.818e-02  1.371e-02  -1.326 0.185488
weight       -6.710e-03  6.551e-04  -10.243 < 2e-16 ***
acceleration  7.910e-02  9.822e-02   0.805 0.421101
model_year   7.770e-01  5.178e-02  15.005 < 2e-16 ***
factor(origin)2 2.630e+00  5.664e-01   4.643 4.72e-06 ***
factor(origin)3 2.853e+00  5.527e-01   5.162 3.93e-07 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 3.307 on 383 degrees of freedom
(因為不存在，6 個觀察量被刪除了)
Multiple R-squared:  0.8242,    Adjusted R-squared:  0.8205
F-statistic: 224.5 on 8 and 383 DF,  p-value: < 2.2e-16

```

(i) The variables displacement, weight, model_year, factor(origin)2 and factor(origin)3.

(ii) By (i), based on the plot and the p -values, I believe weight are the most effective at increasing mpg.

(c)(i) Drop the column car_name of the dataset and standardize.

```

1 # Question 3 (c-i, ii)
2 auto_std <- data.frame(scale(auto[,colnames(auto)!="car_name"])) # Standardize
3 auto_lr_model_std <- lm(mpg ~ cylinders+displacement+horsepower+
4                         weight+acceleration+
5                         model_year, data=auto_std)
6 summary(auto_lr_model_std)

```

```

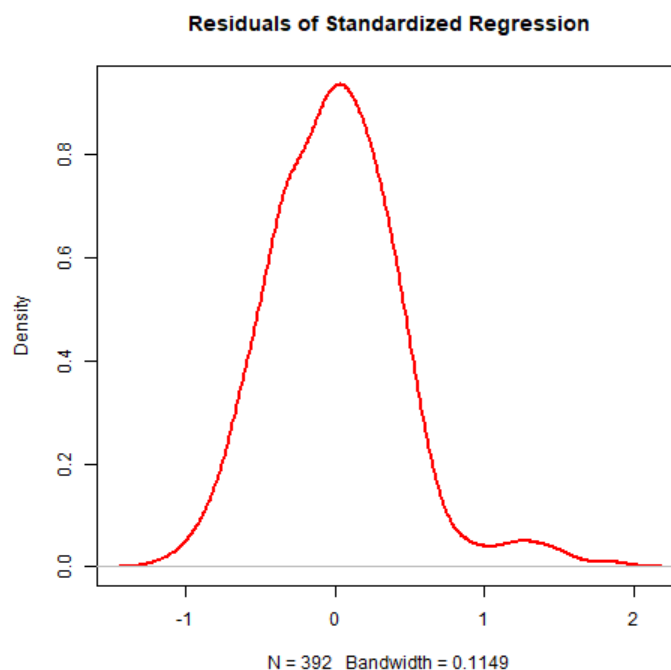
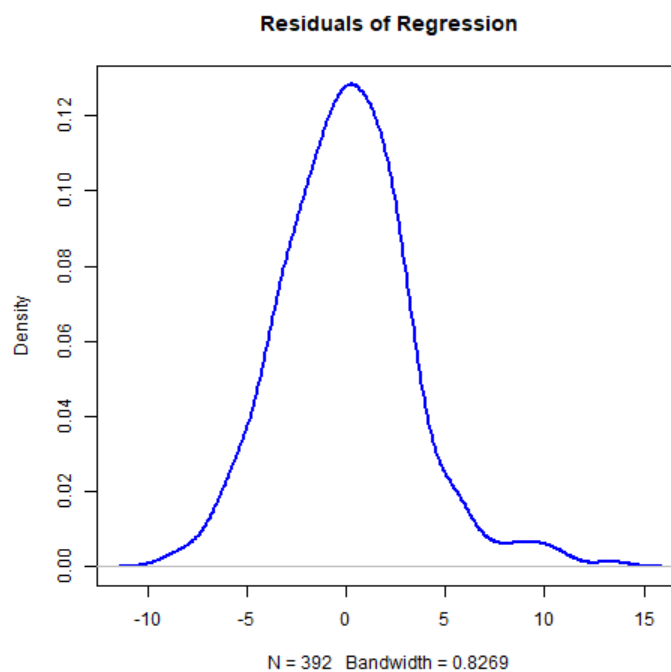
Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)  0.0004236  0.0222112   0.019  0.985
cylinders    -0.0717877  0.0722763  -0.993  0.321
displacement  0.1024348  0.0981565   1.044  0.297
horsepower   -0.0019273  0.0681403  -0.028  0.977
weight       -0.7361794  0.0725952 -10.141 <2e-16 ***
acceleration  0.0300867  0.0360009   0.836  0.404
model_year   0.3564069  0.0248929  14.318 <2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

(ii) Based on the result, it's still weight.

(iii) Plot the density of the residuals.

```
1 # Question 3 (c-iii)
2 png(filename = "3c-1.png")
3 plot(density(auto_lr_model$residuals),
4       main="Residuals of Regression",
5       col="blue", lwd=2)
6 dev.off()
7
8 png(filename = "3c-2.png")
9 plot(density(auto_lr_model_std$residuals),
10      main="Residuals of Standardized Regression",
11      col="red", lwd=2)
12 dev.off()
```



It looks like in both cases, the residuals are normally distributed and centered around zero, intuitively. ■