

GMP, gross metropolitan product, is the value of final goods and services produced within a metropolitan statistical area. Those metropolitan statistical areas, or MSAs, are determined by U.S. Statistical agencies and the U.S. Bureau of Economic Analysis estimates these MSAs' to the country's gross metropolitan products. In the paper by Bettencourt et al, the GMP and population size of the area has a special relationship

$$(GMP) \approx c * (population\ size)^b$$

for some rational number $c > 0$ and $b > 1$, which is also called “supra-linear power law scaling”.^[1] Unlike the linear model, $y \approx a_0 + a_1x$, or the quadratic model, $y \approx a_0 + a_1x + a_2x^2$, supra-linear scaling model uses a positive rational number $b > 1$ as an exponent. In this paper, we will first verify the previous theorem that ‘GMP and population size have supra-linear power law scaling relationship,’ and investigate the alternative linear model that uses population size and other variables as variables from the U.S. Bureau of Economic Analysis describing MSAs in 2006 consists of GMP, population size, finances, professional and technical services, information, communication and technology, and management of firms and enterprises for each area. Comparing to the previous supra-linear model, the alternative linear models can reflect more variables for estimation and will be easier to analyze intuitively because the models are linear.

Supra-linear power law scaling model can be easily converted into a linear model using variable transformation. Since we have

$$(GMP) \approx c * (population\ size)^b$$

if we take logarithms on both sides,

$$\log(GMP) \approx \log(c) + b * \log(population\ size)$$

which represents the linear relationship between log-scale GMP and log-scale population size.

We used the squared-error loss on a log scale as a loss function,

$$L(z, \theta) = [\log(Y)^2 - \mu_{\theta}(N)]^2$$

to calculate the in-sample loss and evaluate the supra-linear model. Not only used the loss function, but we also used the adjusted R^2 to evaluate how the model fits well to the data.

For alternative models, we used multi-variable linear models to predict the response. The alternative hypothesis basically assumed the linear relationship between per-capita GMP and population size, as well as using finances, professional and technical services, information, communication and technology as additional variables that helps to determine the per-capita GMP. We assumed that the higher per-capita GMP implies the probability of the higher economic level, so we used economic variables that represent the economic level.

We used log scale GMP versus log scale population size model instead of using normal GMP versus population size to evaluate the supra-linear law scaling model, from the mathematical backgrounds described above. On the figure 1, the first model, supra-linear law scaling model, showed clear linear relationship when we transformed the variables and responses into log scale. Typically, most of the data were clustered in the lower left corner of the figure 1. The estimate intercept for the linear model was 8.796 and the slope of the log-scale population was 1.123. From the estimates, we can calculate the c and b in the previous equation

$$\log(GMP) \approx \log(c) + b * \log(population\ size)$$

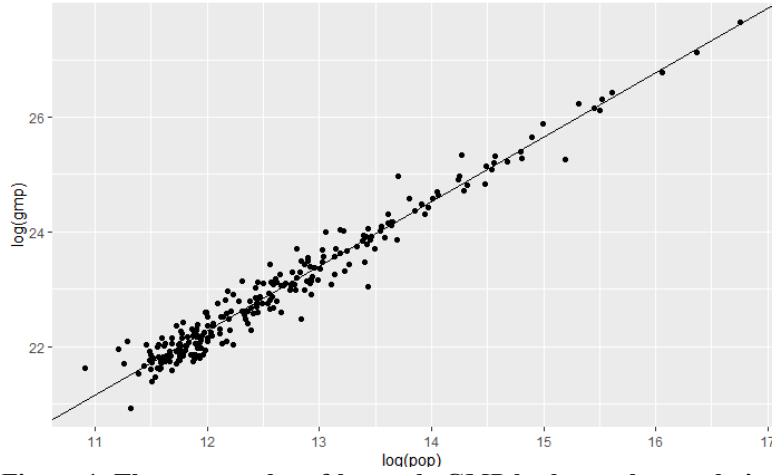


Figure 1. The scatter plot of log-scale GMP by log-scale population size, with linear model through data.

since $c = e^{8.796} = 6607.76$, $b = 1.123$. If we wanted to see the relationship between per-capita GMP and population size

$$\begin{aligned} \log(\text{per-capita GMP}) &= \log\left(\frac{(GMP)}{(\text{population size})}\right) \\ &\approx \log(c) + (b - 1) * \log(\text{population size}) \end{aligned}$$

so we can easily check that log scale per-capita GMP also has a linear relationship with log-scale population size, with coefficients $c = e^{8.796} = 6607.76$, $b' = b - 1 = 0.123$. t-value for each estimation was 47.94 and 77.54, the variance of residuals was $5.64 * 10^{-2}$, and the adjusted R^2 value for the model was 0.961. Since our log-scale population size had high t-value, we cannot reject the null hypothesis: ‘log-scale population size is not a meaningless variable to estimate the log-scale GMP’. From high adjusted R^2 value, we can trust the model. Double-checking with the loss function, the in-sample loss for the model was $5.62 * 10^{-2}$, which is quite low. Therefore, we can definitely say that the supra-linear law scaling model is plausible.

Figure 2, 3 and 4 represents the relationship between log scale population and log scale GMP, coloring the points with finances, finances, professional and technical services, and information, communication and technology. Typically, a considerable number of data was

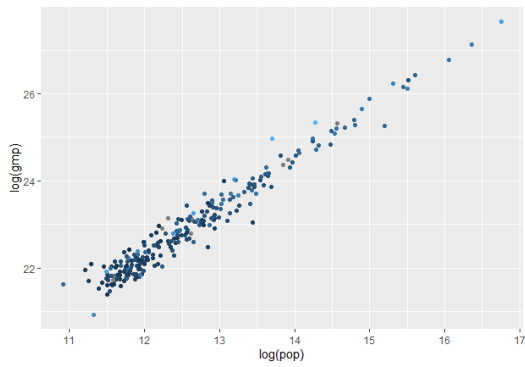


Figure 2. The scatter plot of log-scale GMP by log-scale population size, with points colored according to the level of finances.

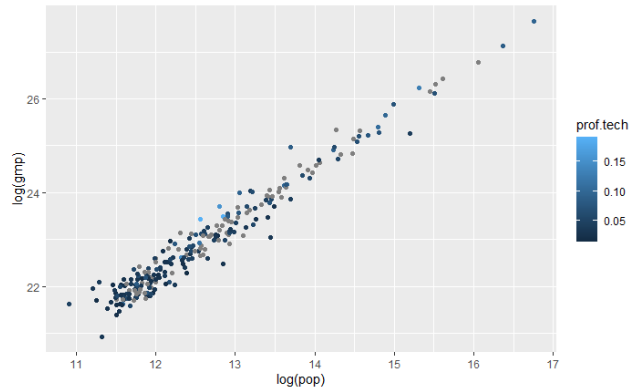


Figure 3. The scatter plot of log-scale GMP by log-scale population size, with points colored according to the level of professional and technical services.

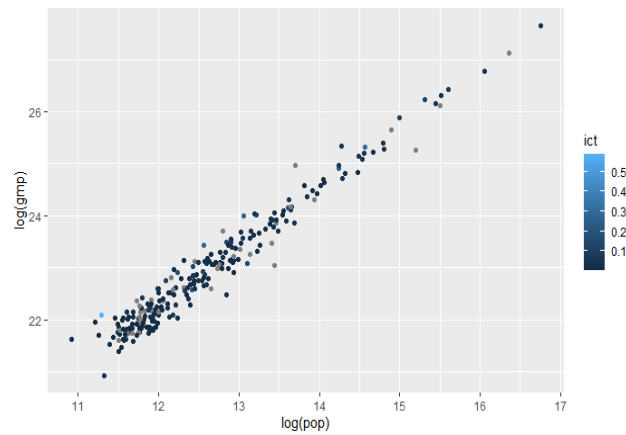


Figure 4. The scatter plot of log-scale GMP by log-scale population size, with points colored according to the level of information, communication and technology.

Missing Value; 3.69% of the finances data, 32.79% of the professional and technical services data, and 16.80% of the information, communication and technology data were missing values. We were not able to find the characteristic properties of information, communication and technology data on the graph in figure 4. However, the higher finances data tends to have higher log-scale population value and higher log-scale GMP in figure 2 and the higher – professional and technical services data tends to have higher log-scale population value and higher log-scale GMP in figure 3. However, it was not enough to say that those were noticeable levels.

For the first alternative linear model (per-capita GMP) \sim (population size) + (finances),

the estimated intercept was $2.365 * 10^4$, estimate coefficient for population size was $1.249 * 10^{-3}$, and estimate coefficient for finances was $5.189 * 10^4$. t-values for each coefficient were 18.109, 4.056, and 6.159 which are quite high so that we cannot reject any variables in the model. However, the residual standard error was 7821, which is quite large even considering the graph is not a log-scale, and the adjusted R-square value is 0.2615. Moreover, the in-sample loss for the model was $6.038 * 10^8$; Comparing quickly in terms of scale, its log value is 17.916, and our supra-linear model's in-sample loss value was 0.056. These results support that we cannot trust our first alternative model.

For the second alternative linear model (per-capita GMP) \sim (population size) + (professional and technical services), the estimated intercept was $2.433 * 10^4$, estimate coefficients for population size was $1.090 * 10^{-3}$, and estimated coefficient for professional and technical services was $1.419 * 10^5$. t-value for each coefficient was 20.193, 3.186 and 6.482 so we also cannot reject any variables in the model. The residual standard error was 7948 and the adjusted R-square value was 0.292. In sample loss for the model was $6.202 * 10^8$. So, we cannot trust our second alternative model also.

The third alternative model (per-capita GMP) \sim (population size) + (information, communication and technology) had estimate intercept $2.915 * 10^4$, estimated coefficient for population size was $1.990 * 10^{-3}$, and estimated coefficient for information, communication and technology was $5.236 * 10^4$. t-value for each coefficient was 45.703, 6.350 and 6.126 so we cannot reject any variables in the model. The residual standard error was 7378 and the adjusted R-square value was 0.296. In sample loss for the model was $5.364 * 10^8$. So, we cannot trust all three alternative models.

In conclusion, the supra-linear law scaling model performed much better than any alternative models we considered. It showed a higher adjusted R-square value and lower in-sample loss. The other variables, finances, professional and technical service, information,

communication and technology, showed a weak relationship with log GMP and log population size, but it was not significant.

One concern we had was whether to include graphs that express the relationship between per-capita GMP and other variables. Those relationships are approximately described in figure 2,3,4 by coloring points with other variables, but there is a lack of an intuitive understanding of the relationship between response and those variables. But we decided not to contain those graphs since we have roughly confirmed that the relationships between those do not exist strongly with figure 2, 3, 4. Also, these relationships only has an indirect connection with the models we consider; In the multi-variable linear model $Y \sim X_1 + X_2$, even if we find the strong relationship between Y and X_2 in the model $Y \sim X_2$, this does not guarantee the relationship between Y and X_2 in the model $Y \sim X_1 + X_2$.

In addition, the in-sample loss value through loss function was only comparable to the absolute number in the supra-linear model and the alternative models, and it was really hard to determine how much difference this really was, because one was log scale, and others were normal scale. We would like to make up for it if possible, in the future.

References

- [1] Luís M. A. Bettencourt, José Lobo, Dirk Helbing, Christian Kühnert, Geoffrey B. West, *Growth, innovation, scaling, and the pace of life in cities*, Proceedings of the National Academy of Sciences Apr 2007, 104 (17) 7301-7306; DOI: 10.1073/pnas.0610172104

Unit 2 Paper Technical Appendices

youngwoo Kwon

2021 3 4

#Summary In Appendix 1, the theoretical background for the data modification was proven. Also, the hypothesis selection was done.

In Appendix 2, basic data analysis was done. The code calculated the proportion of missing values and displayed some scatter plots explaining the relationship between GMP and population size. The code also plotted other variables to find the connection with the previous relationship.

In Appendix 3, the code chose the linear model and plotted that model. It calculated the loss function outcome and residual variances.

In Appendix 4, the alternative models were written. Also, the loss function outcomes for those models were calculated.

#Appendix 1: Detail of Statistical models

1. If $Y \approx cN^b$ for some $c > 0, b > 1$, then $\log(\frac{Y}{N}) \approx \beta_0 + \beta_1 \log(N)$ for some $\beta_0 \in (-\infty, \infty), \beta_1 > 0$, and also $\log(Y) \approx \beta_0 + (1 + \beta_1) \log(N)$.

Let $Y \approx cN^b$. Then, $\log(Y) \approx \log(c) + b \log(N)$. Therefore, for $\beta_0 = \log(c), \beta_1 = b - 1$, $\log(Y) \approx \beta_0 + (1 + \beta_1) \log(N)$. Since $c > 0, b > 1$, we can say that $\beta_0 \in (-\infty, \infty)$ and $\beta_1 > 0$.

If we subtract $\log(N)$ in both sides, $\log(Y) - \log(N) \approx \log(c) + (b - 1) \log(N)$. So $\log(\frac{Y}{N}) \approx \beta_0 + \beta_1 \log(N)$ for $\beta_0 = \log(c), \beta_1 = b - 1$.

2. Three hypothesis about how these other variables might influence per-capita GMP (pcgmp).
 - 1) There is a linear relationship between Per-Capita GMP and population + finance. (pcgmp ~ pop + finance)
 - 2) There is a linear relationship between Per-Capita GMP and population + information, communication and technology. (pcgmp ~ pop + ict)
 - 3) There is a linear relationship between Per-Capita GMP and population + professional and technical services. (pcgmp ~ pop + prof.tech)

#Appendix 2: Exploratory analyses

1. Read and modify the data

```
library(ggplot2)
library(dplyr)
```

```
##
## 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
```

```
mydata = read.csv("http://dept.stat.lsa.umich.edu/~bbh/s485/data/gmp-2006.csv")
head(mydata)
```

```
##           MSA pcgmp    pop finance prof.tech    ict
## 1           Akron, OH 32890 699300 0.12940  0.05440    NA
## 2           Albany, GA 24270 163000 0.08217    NA 0.00708
## 3  Albany-Schenectady-Troy, NY 36840 850300 0.15780  0.09399 0.04511
## 4           Albuquerque, NM 37660 816000 0.15990  0.09978 0.20500
## 5           Alexandria, LA 25490 152200 0.09152  0.03790 0.01134
## 6 Allentown-Bethlehem-Easton, PA-NJ 30160 794400 0.13670    NA 0.03384
## management
## 1    0.054310
## 2         NA
## 3         NA
## 4    0.006509
## 5    0.015210
## 6         NA
```

```
newdata <- mydata
newdata$pcgmp <- as.double(newdata$pcgmp)
newdata$pop <- as.double(newdata$pop)
newdata$gmp <- newdata$pop * newdata$pcgmp
head(newdata)
```

```
##           MSA pcgmp    pop finance prof.tech    ict
## 1           Akron, OH 32890 699300 0.12940  0.05440    NA
## 2           Albany, GA 24270 163000 0.08217    NA 0.00708
## 3  Albany-Schenectady-Troy, NY 36840 850300 0.15780  0.09399 0.04511
## 4           Albuquerque, NM 37660 816000 0.15990  0.09978 0.20500
## 5           Alexandria, LA 25490 152200 0.09152  0.03790 0.01134
## 6 Allentown-Bethlehem-Easton, PA-NJ 30160 794400 0.13670    NA 0.03384
## management      gmp
## 1    0.054310 22999977000
```

```
## 2      NA 3956010000
## 3      NA 31325052000
## 4 0.006509 30730560000
## 5 0.015210 3879578000
## 6      NA 23959104000
```

2. Missing Values

```
nrow(newdata)
```

```
## [1] 244
```

```
Finance_prop = nrow(newdata[4] %>% na.omit())/nrow(newdata)
Prof.tech_prop = nrow(newdata[5] %>% na.omit())/nrow(newdata)
ict_prop = nrow(newdata[6] %>% na.omit())/nrow(newdata)
management_prop = nrow(newdata[7] %>% na.omit())/nrow(newdata)
Finance_prof.tech_prop = nrow(newdata[4:5] %>% na.omit())/nrow(newdata)
```

```
Finance_prop
```

```
## [1] 0.9631148
```

```
Prof.tech_prop
```

```
## [1] 0.6721311
```

```
ict_prop
```

```
## [1] 0.8319672
```

```
management_prop
```

```
## [1] 0.5532787
```

```
Finance_prof.tech_prop
```

```
## [1] 0.6557377
```

```
nrow(newdata %>% na.omit())/nrow(newdata)
```

```
## [1] 0.3729508
```

96.31148% of data have no missing value in finance section.

67.21311% of data have no missing value in professional and technical services section.

83.19672% of data have no missing value in information, communication and technology section.

55.32787% of data have no missing value in and enterprises section.

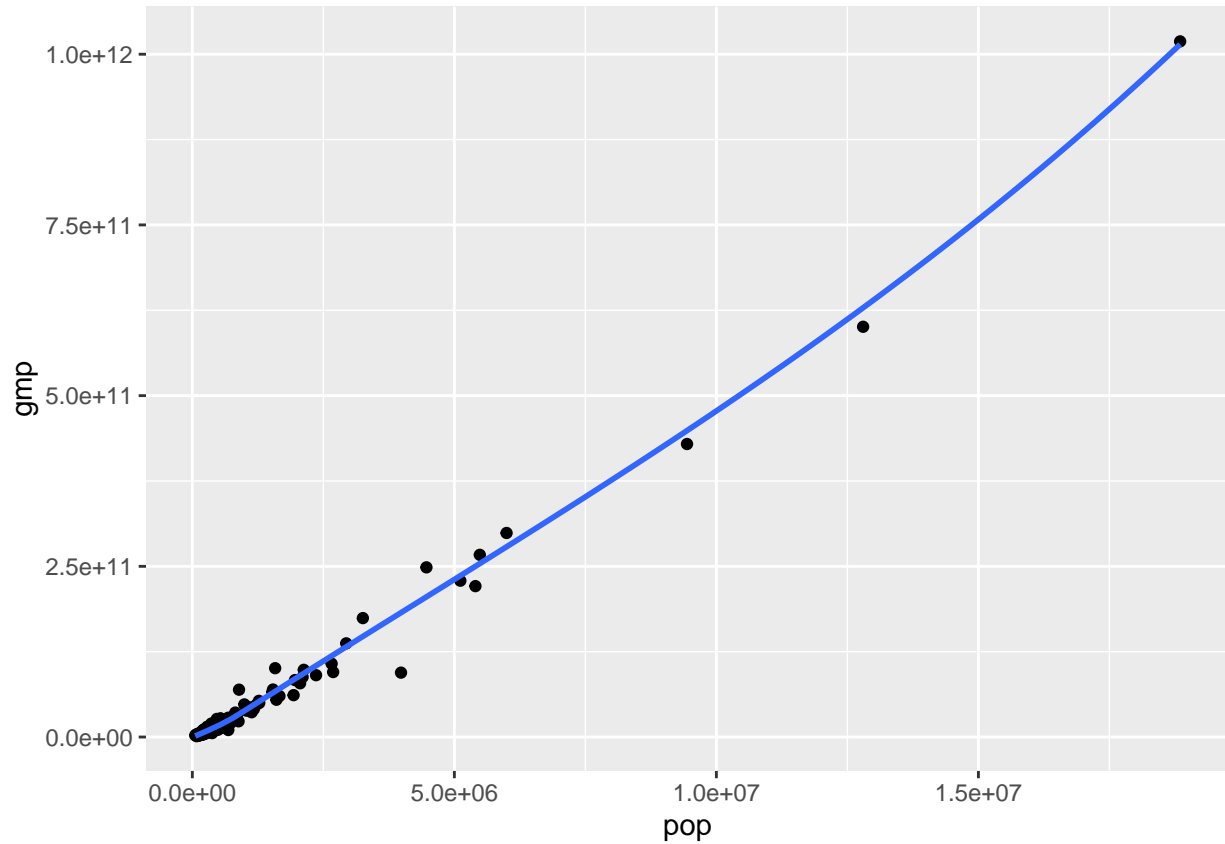
65.57377% of data have no missing value in finance and professional and technical services section.

37.29508% of data have no missing value.

3. Scatter plot

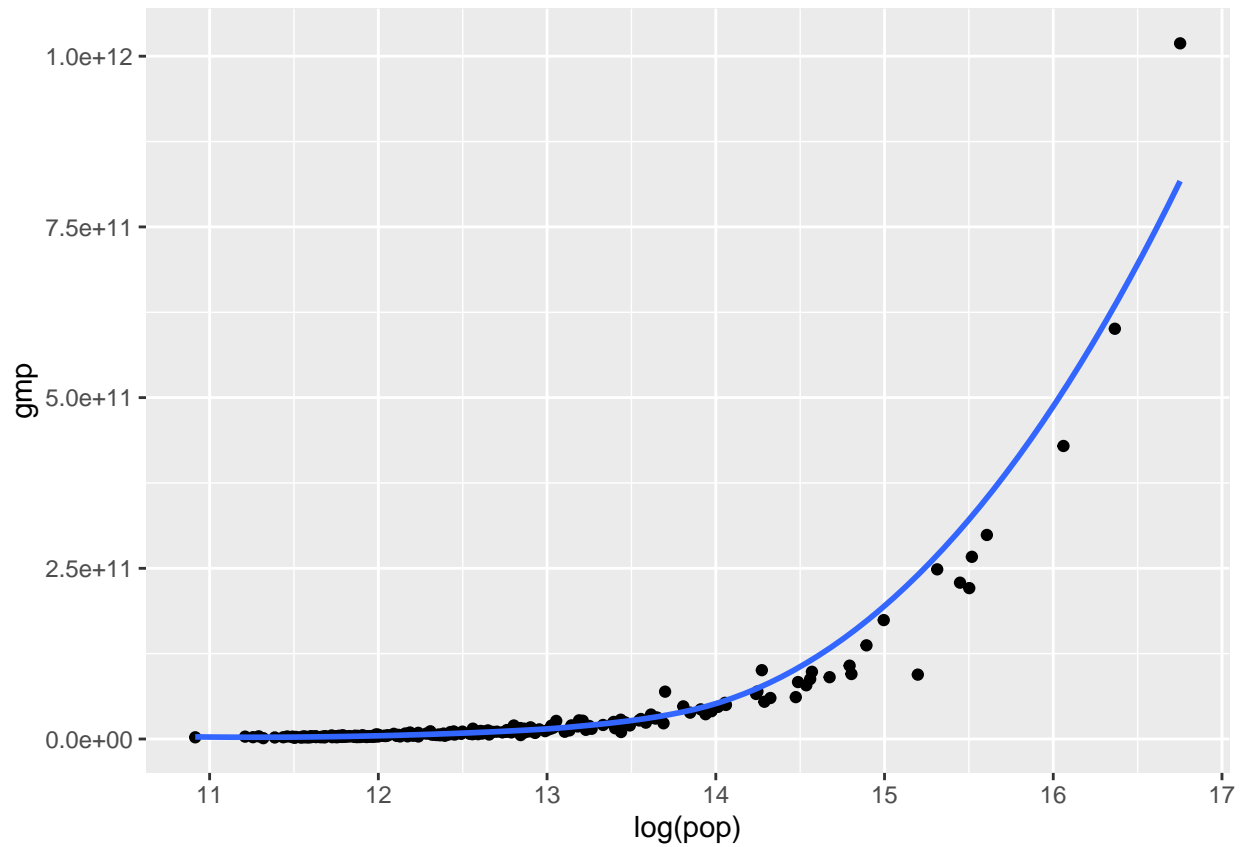
```
gmp_pop = ggplot(newdata, aes(y=gmp, x=pop)) +
  geom_point() +
  geom_smooth(se=FALSE)
gmp_pop
```

```
## 'geom_smooth()' using method = 'loess' and formula 'y ~ x'
```



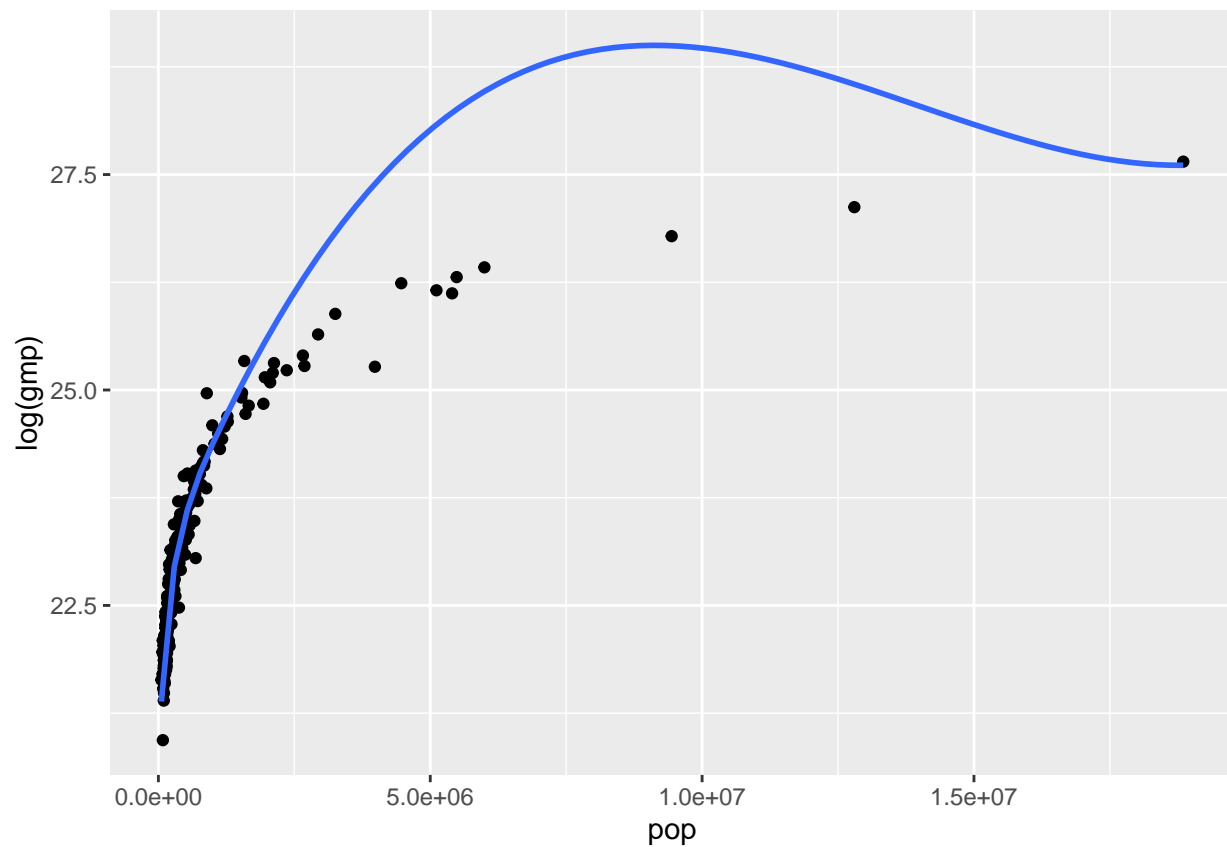
```
loggmp_pop = ggplot(newdata, aes(y=gmp, x=log(pop))) +
  geom_point() +
  geom_smooth(se=FALSE)
loggmp_pop
```

```
## 'geom_smooth()' using method = 'loess' and formula 'y ~ x'
```



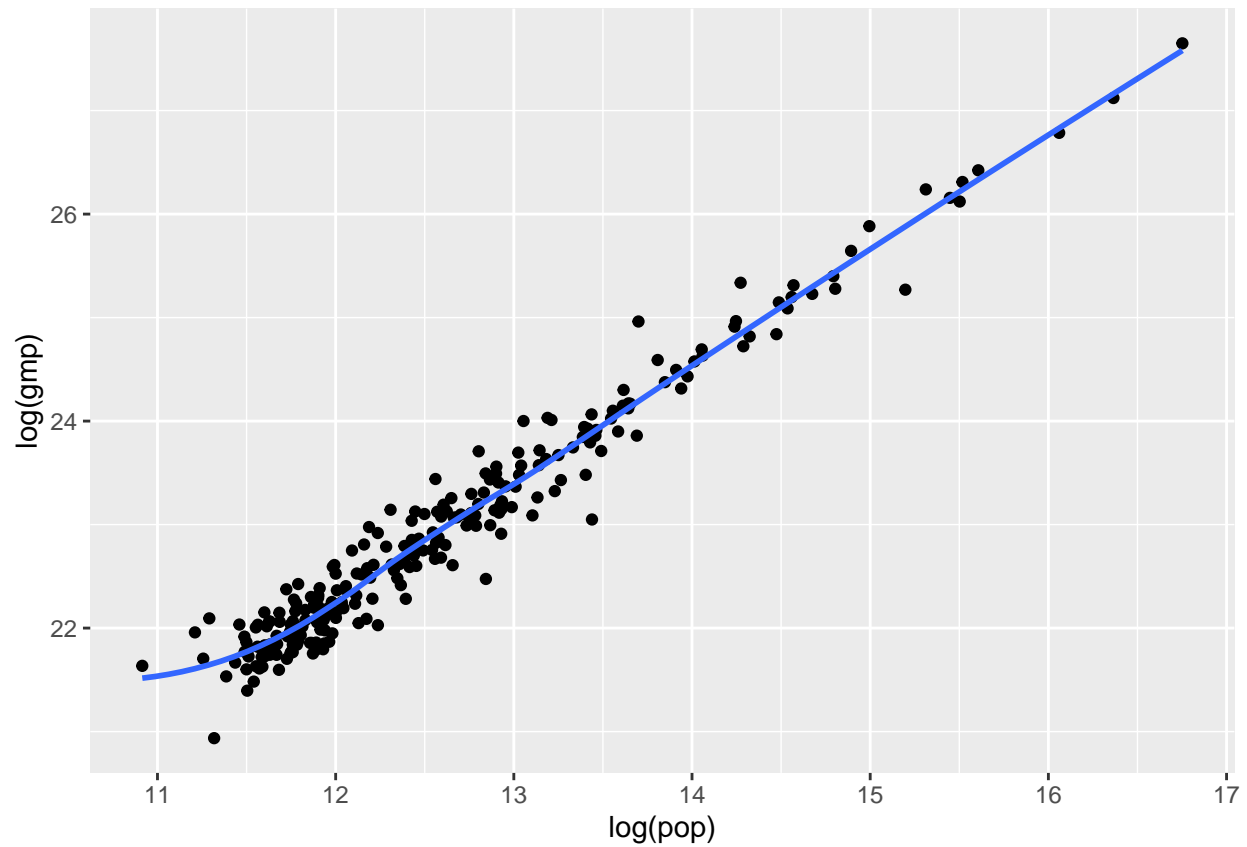
```
gmp_logpop = ggplot(newdata, aes(y=log(gmp), x=pop)) +  
  geom_point() +  
  geom_smooth(se=FALSE)  
gmp_logpop
```

```
## 'geom_smooth()' using method = 'loess' and formula 'y ~ x'
```



```
loggmp_logpop = ggplot(newdata, aes(y=log(gmp), x=log(pop))) +  
  geom_point() +  
  geom_smooth(se=FALSE)  
loggmp_logpop
```

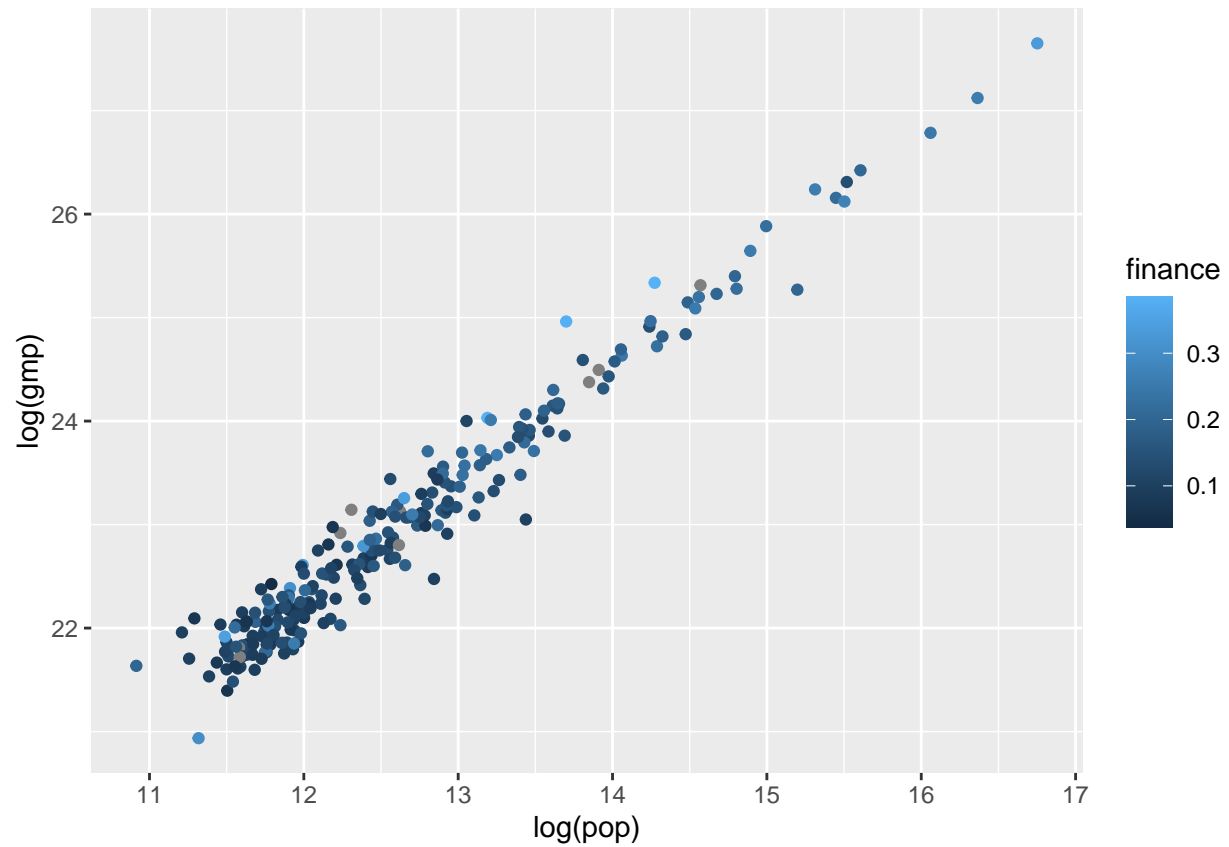
```
## 'geom_smooth()' using method = 'loess' and formula 'y ~ x'
```



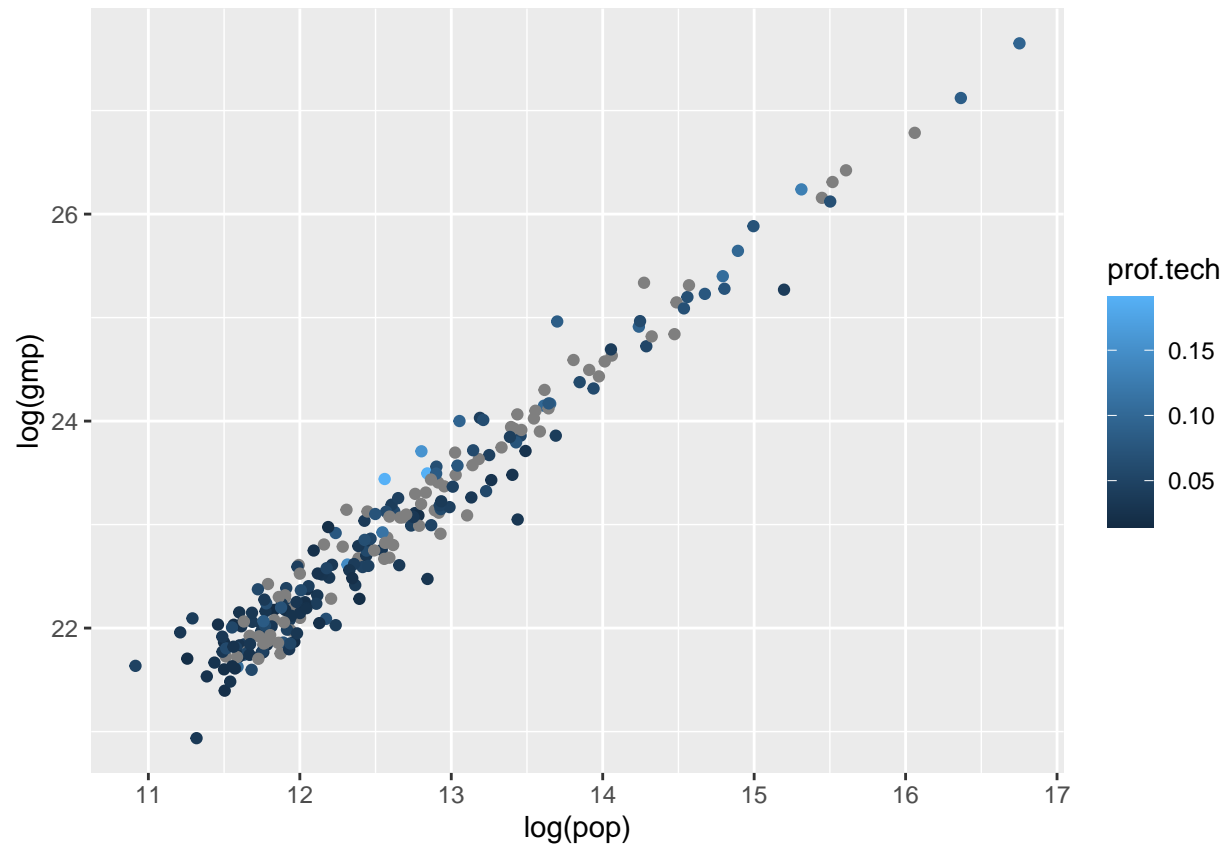
The results say that the $\log(y) \sim \log(x)$ is better scale for capturing patterns.

4. Other variances and gmp~pop

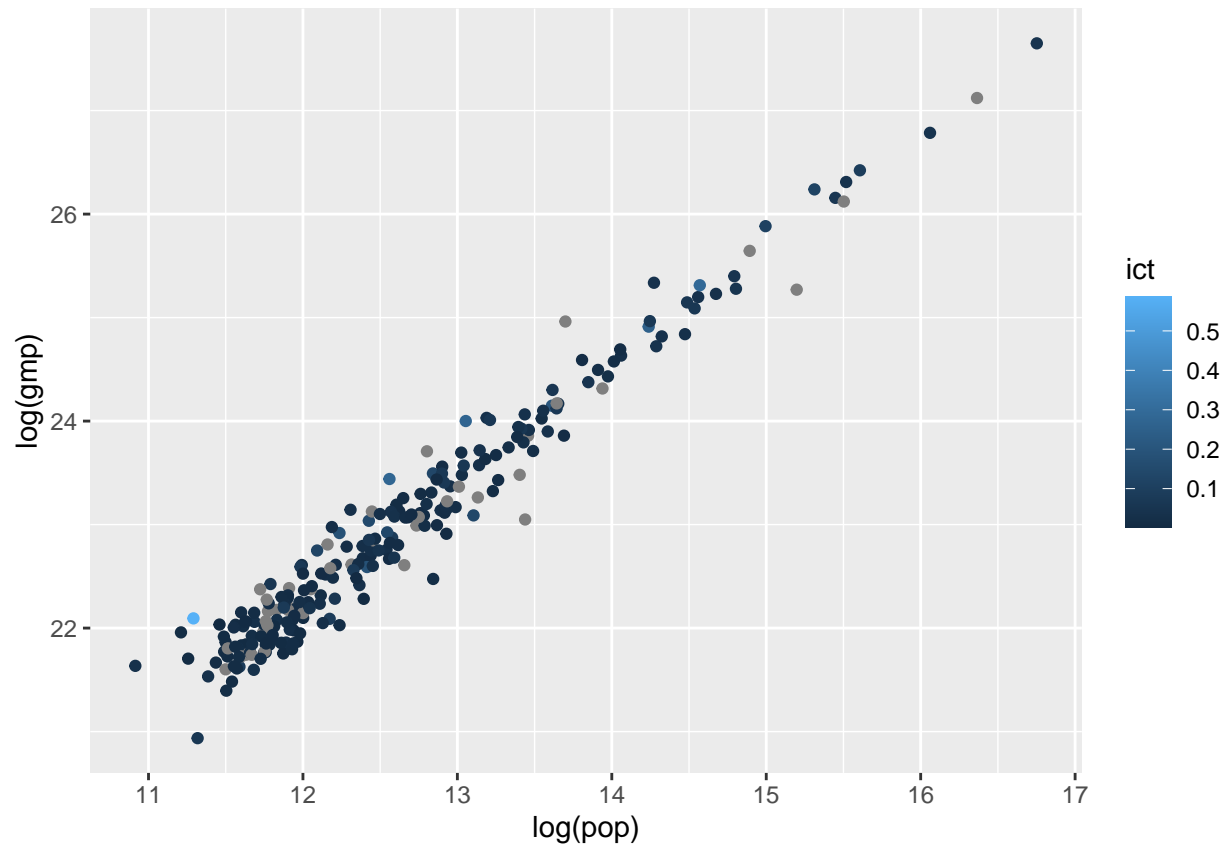
```
loggmp_logpop_finance = ggplot(newdata, aes(y=log(gmp), x=log(pop))) +  
  geom_point(aes(colour = finance))  
loggmp_logpop_finance
```



```
loggmp_logpop_prof.tech = ggplot(newdata, aes(y=log(gmp), x=log(pop))) +  
  geom_point(aes(colour = prof.tech))  
loggmp_logpop_prof.tech
```



```
loggmp_logpop_ict = ggplot(newdata, aes(y=log(gmp), x=log(pop))) +  
  geom_point(aes(colour = ict))  
loggmp_logpop_ict
```

I didn't remove the NA values because ggplot would automatically neglect and do not colour the data that have NA value

#Appendix 3: Fitting the power law model

1. Basic linear model

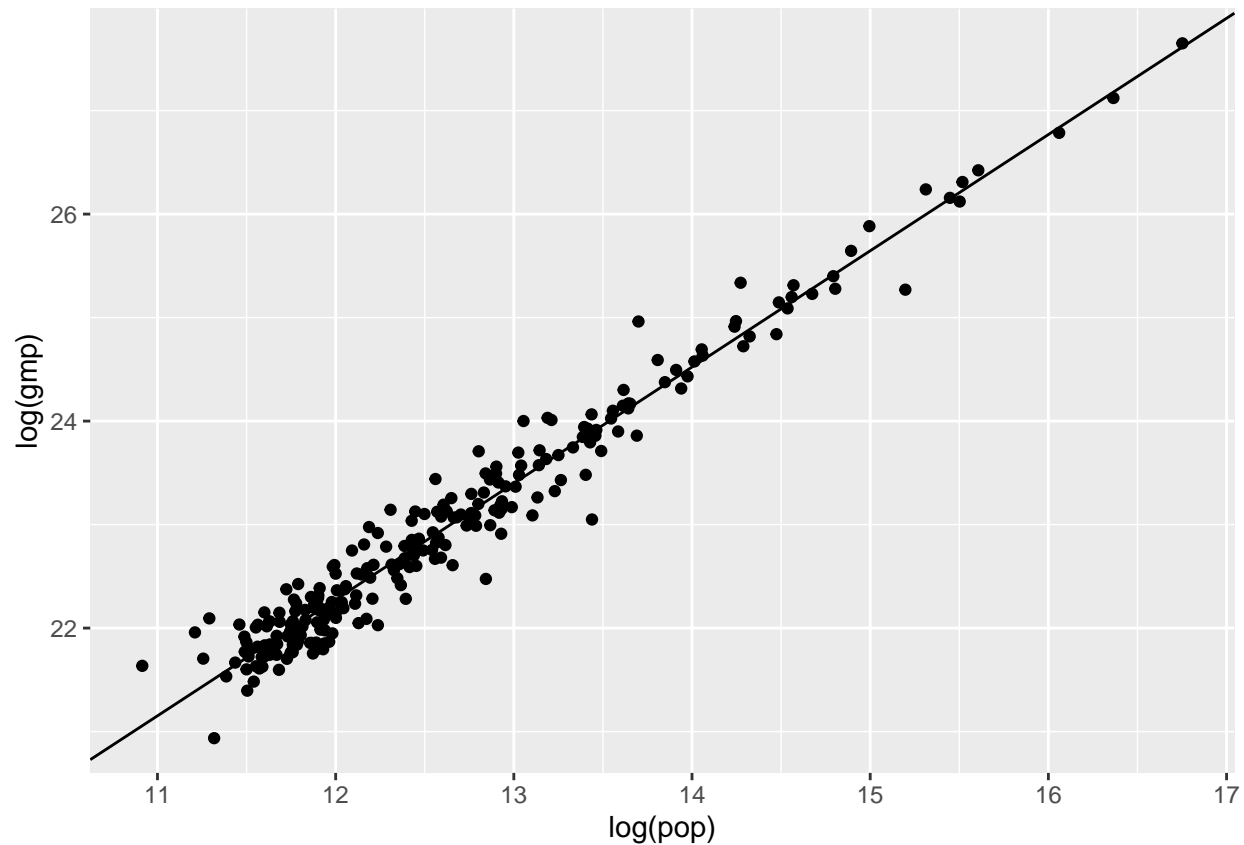
```
lm_loggmp_logpop = lm(log(gmp)~log(pop), data = newdata)
summary(lm_loggmp_logpop)

##
## Call:
## lm(formula = log(gmp) ~ log(pop), data = newdata)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.84226 -0.13993  0.00157  0.12942  0.77779
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  8.79623    0.18350   47.94  <2e-16 ***
## log(pop)      1.12326    0.01449   77.54  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.238 on 242 degrees of freedom
## Multiple R-squared:  0.9613, Adjusted R-squared:  0.9611
## F-statistic: 6012 on 1 and 242 DF, p-value: < 2.2e-16
```

As we saw in the #Appendix 1, the $\log(c) = \log(8.79623)$ equals to the β_0 , and $b-1 = 1.12326 - 1 = 0.012326$ equals to the β_1 . Since the Adjusted R-squared value is over 0.96 and t value for each estimate is large, we can say that this model supports the supra-linear power-law scaling hypothesis.

2. Plot the data, errors and residuals

```
ggplot(newdata, aes(y=log(gmp), x=log(pop))) +
  geom_point() +
  geom_abline(intercept = lm_loggmp_logpop$coefficients[1], slope = lm_loggmp_logpop$coefficients[2])
```



```
var(lm_loggmp_logpop$residuals)
```

```
## [1] 0.05642693
```

```
0.238^2 #From Residual standard error at linear model summary
```

```
## [1] 0.056644
```

So the variance of the residuals are almost equal to the variance of the regression. Since we got high t-value and small p-value for each coefficient and high adjusted R-square value, we can trust the estimated coefficients.

3. Loss function, In-sample loss, estimated values of parameters

```
loss_log <-function(z, model){
  result = (log(z[1]) - predict(model, z[-1]))^2
  return(colMeans(result))
}
```

```
loss_log(newdata[c(8,3)], lm_loggmp_logpop)
```

```
##          gmp
## 0.05619567
```

(Used `log_e` instead of `log_10`. Essentially, $\log_e(x) = r \log_{10}(x)$ where $r = \log_e 10$, so nothing important changed.)

So the in-sample loss is 0.05619567. Since the in-sample loss is quite low, the expected values of the parameters make sense.

#Appendix 4: Fitting and assessment of alternate models

1, 2. Three alternate regression models & fit models

- 1) There is a linear relationship between Per-Capita GMP and population + finance. (pcgmp ~ pop + finance)
- 2) There is a linear relationship between Per-Capita GMP and population + information, communication and technology. (pcgmp ~ pop + ict)
- 3) There is a linear relationship between Per-Capita GMP and population + professional and technical services. (pcgmp ~ pop + prof.tech)

```
alt_model1 = lm(pcgmp~pop + finance, data = newdata)
alt_model2 = lm(pcgmp~pop + ict, data = newdata)
alt_model3 = lm(pcgmp~pop + prof.tech, data = newdata)
```

```
summary(alt_model1)
```

```
##
## Call:
## lm(formula = pcgmp ~ pop + finance, data = newdata)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -24223  -4509   -989    3878   33425
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  2.365e+04  1.306e+03  18.109 < 2e-16 ***
## pop          1.249e-03  3.080e-04   4.056 6.82e-05 ***
## finance      5.189e+04  8.425e+03   6.159 3.19e-09 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 7821 on 232 degrees of freedom
## (9 observations deleted due to missingness)
## Multiple R-squared:  0.2678, Adjusted R-squared:  0.2615
## F-statistic: 42.43 on 2 and 232 DF, p-value: < 2.2e-16
```

```
summary(alt_model2)
```

```
##
## Call:
## lm(formula = pcgmp ~ pop + ict, data = newdata)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -17450.6  -4995.6   -801.7   4334.7  29372.5
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  2.915e+04  6.379e+02  45.703 < 2e-16 ***
```

```
## pop          1.990e-03  3.135e-04  6.350 1.42e-09 ***
## ict          5.236e+04  8.547e+03  6.126 4.70e-09 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 7378 on 200 degrees of freedom
## (41 observations deleted due to missingness)
## Multiple R-squared:  0.3029, Adjusted R-squared:  0.296
## F-statistic: 43.46 on 2 and 200 DF, p-value: < 2.2e-16
```

```
summary(alt_model3)
```

```
##
## Call:
## lm(formula = pcgmp ~ pop + prof.tech, data = newdata)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -16233  -4599   -667    3905   40522
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  2.433e+04  1.205e+03  20.193 < 2e-16 ***
## pop          1.090e-03  3.420e-04   3.186  0.00173 **
## prof.tech    1.419e+05  2.188e+04   6.482  1.05e-09 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 7948 on 161 degrees of freedom
## (80 observations deleted due to missingness)
## Multiple R-squared:  0.3004, Adjusted R-squared:  0.2917
## F-statistic: 34.56 on 2 and 161 DF, p-value: 3.25e-13
```

All three models have very low adjusted r-squared value.

3. Evaluate the model based on the square-error loss function

```
loss <-function(z, model){
  result = (z[1] - predict(model, z[-1]))^2
  return(colMeans(result))
}

loss(newdata[c(2,3,4)] %>% na.omit(), alt_model1)
```

```
##      pcgmp
## 60380645
```

```
loss(newdata[c(2,3,6)] %>% na.omit(), alt_model2)
```

```
##      pcgmp
## 53635797
```

```
loss(newdata[c(2,3,5)] %>% na.omit(), alt_model3)
```

```
##      pcgmp  
## 62017277
```

```
log(loss(newdata[c(2,3,4)] %>% na.omit(), alt_model1))
```

```
##      pcgmp  
## 17.91618
```

```
log(loss(newdata[c(2,3,6)] %>% na.omit(), alt_model2))
```

```
##      pcgmp  
## 17.79773
```

```
log(loss(newdata[c(2,3,5)] %>% na.omit(), alt_model3))
```

```
##      pcgmp  
## 17.94292
```

All three models have very large loss function output.

Therefore, it is hard to trust our coefficients from the linear model.