

Tutorial 2: Properties of Random Variables

Anh Le

September 4, 2015

Agenda (and learning goals)

1. Implement formulas for Expected Values, Variance, etc. in R
 - learn vectorized operation
2. Download data automatically from the web
 - learn `help()` in R
 - learn reproducible analysis even at the downloading data step
3. Draw the plots you saw from lectures in R (histograms, density plots, boxplot, normal quantile plot, scatterplot)
 - learn how to generate random sample
 - learn how to inspect the distribution of real data
4. Tips and tricks

1. Implement expected value and variance formula

Calculate Expected Value:

Use `sum()` (to get the sum) and `length()` (to get the number of elements in a vector). Calculate:

$$E(X) = \frac{1}{n} \sum_{i=1}^n X_i$$

```
X <- rnorm(1000)
sum(X) / length(X)
```

```
## [1] 0.04186559
```

```
mean(X)
```

```
## [1] 0.04186559
```

Calculate Variance:

$$Var(X) = \frac{1}{n-1} \sum_{i=1}^n (X_i - E(X))^2$$

Let's break down this formula. Mathematically, the formula mean that for each element X_i in the vector X : - subtract $E(X)$ from X_i , square the result - then we add up all the results and divide by $n - 1$

So we can naively translate that into code as follows:

```
myVec <- rnorm(1000, mean = 2, sd = 5)

myVar1 <- function(X) {
  n <- length(X)

  sum = 0
  # For each element X_i
  for (i in 1:n) {
    # Subtract E(X), square the result, then add the results together
    sum = sum + (X[i] - mean(X)) ** 2
  }

  return(sum / (n - 1))
}

myVar1(myVec)
```

```
## [1] 25.54645
```

```
var(myVec)
```

```
## [1] 25.54645
```

But loops in R are notoriously slow! We should use vectorized operation instead. For example,

```
X <- 1:5

# To subtract E(X) from each element
X - mean(X)
```

```
## [1] -2 -1  0  1  2
```

```
# To square all elements
X ** 2
```

```
## [1]  1  4  9 16 25
```

```
# To calculate the sum of squares
sum(X ** 2)
```

```
## [1] 55
```

Let's use this to rewrite `myVar1` so that it's faster:

```
myVar2 <- function(X) {
  return(sum((X - mean(X)) ** 2) / (length(X) - 1))
}
```

```
myVar2(myVec)
```

```
## [1] 25.54645
```

```
myVar1(myVec)
```

```
## [1] 25.54645
```

```
var(myVec)
```

```
## [1] 25.54645
```

Let's compare the speed:

```
library(rbenchmark) # install.packages if you don't have the package
benchmark(myVar1(myVec), myVar2(myVec))
```

```
##           test replications elapsed relative user.self sys.self
## 1 myVar1(myVec)           100   0.618         618   0.618       0
## 2 myVar2(myVec)           100   0.001          1   0.001       0
##   user.child sys.child
## 1           0         0
## 2           0         0
```

In-class exercise: Implement covariance formula

You'll learn about the properties of covariance next week. For now, you can implement the following formula of covariance in R.

$$\text{cov}(X, Y) = \frac{1}{N-1} \sum_{i=1}^N (X_i - \bar{X})(Y_i - \bar{Y})$$

```
X <- rnorm(100)
Y <- X + rnorm(10)
myCov(X, Y)
```

```
## [1] 1.028888
```

```
cov(X, Y)
```

```
## [1] 1.028888
```

2. Download data automatically from the web

```
# install.packages("WDI")
library(WDI)
```

```
## Loading required package: RJSONIO
```

```
help(WDI)
```

Let's download GDP data:

```
d_gdp <- WDI(country = "all", indicator = "NY.GDP.MKTP.KD",
             extra = TRUE, start = 2010, end = 2011)
head(d_gdp)
```

```
##      iso2c      country NY.GDP.MKTP.KD year
## 1      1A      Arab World  2.103825e+12 2010
## 2      1A      Arab World  2.173896e+12 2011
## 3      1W      World      6.770244e+13 2011
## 4      1W      World      6.564782e+13 2010
## 5      4E East Asia & Pacific (excluding high income) 8.480167e+12 2011
## 6      4E East Asia & Pacific (excluding high income) 7.820048e+12 2010
##      iso3c      region capital longitude latitude      income      lending
## 1      ARB Aggregates Aggregates Aggregates
## 2      ARB Aggregates Aggregates Aggregates
## 3      WLD Aggregates Aggregates Aggregates
## 4      WLD Aggregates Aggregates Aggregates
## 5      <NA>      <NA>      <NA>      <NA>      <NA>      <NA>
## 6      <NA>      <NA>      <NA>      <NA>      <NA>      <NA>
```

Note how the dataset includes regions' aggregate data as well. We can exclude those rows as follows:

```
# Note that the region variable is available because we specified WDI(extra=TRUE)
d_gdp <- d_gdp[d_gdp$region != "Aggregates", ]
head(d_gdp)
```

```
##      iso2c      country NY.GDP.MKTP.KD year iso3c
## NA      <NA>      <NA>      NA      NA      <NA>
## NA.1    <NA>      <NA>      NA      NA      <NA>
## NA.2    <NA>      <NA>      NA      NA      <NA>
## NA.3    <NA>      <NA>      NA      NA      <NA>
## 11      AD Andorra    3346317329 2010      AND
## 12      AD Andorra    3185604582 2011      AND
##
##      region      capital longitude
## NA      <NA>      <NA>      <NA>
## NA.1    <NA>      <NA>      <NA>
## NA.2    <NA>      <NA>      <NA>
## NA.3    <NA>      <NA>      <NA>
## 11      Europe & Central Asia (all income levels) Andorra la Vella 1.5218
## 12      Europe & Central Asia (all income levels) Andorra la Vella 1.5218
##      latitude      income      lending
## NA      <NA>      <NA>      <NA>
```

```
## NA.1      <NA>                <NA>                <NA>
## NA.2      <NA>                <NA>                <NA>
## NA.3      <NA>                <NA>                <NA>
## 11      42.5075 High income: nonOECD Not classified
## 12      42.5075 High income: nonOECD Not classified
```

3. Draw the plots you saw from lectures in R (histograms, density plots)

We can generate random samples from various distributions in R, using `rbinom`, `rnorm`, `rpois`, etc.

Binomial distribution:

```
binomdraws <- rbinom(n=1000, size=100, prob=0.33)
length(binomdraws)
```

```
## [1] 1000
```

```
mean(binomdraws)
```

```
## [1] 33.073
```

Normal (Gaussian) distribution:

Draw normal samples

```
normdraws <- rnorm(n = 1000, mean = 10, sd = 5)
length(normdraws)
```

```
## [1] 1000
```

```
mean(normdraws)
```

```
## [1] 9.973859
```

```
var(normdraws)
```

```
## [1] 24.34949
```

Inspecting distribution with Histogram, Density plots, and Box plot

```
par(mfrow = c(1, 3))

normdraws <- rnorm(n = 1000, mean = 10, sd = 5)

# Histogram
```

```
hist(normdraws, main="Histogram")
# Density plot
normdensity <- density(normdraws)
plot(normdensity, main="Density plot")
# Box plot
boxplot(normdraws, main="Boxplot")
```

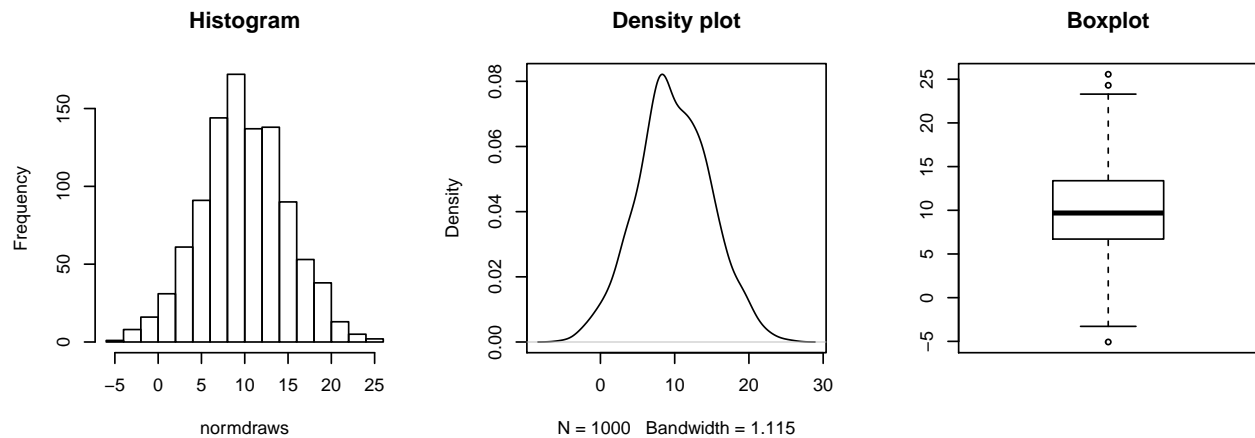
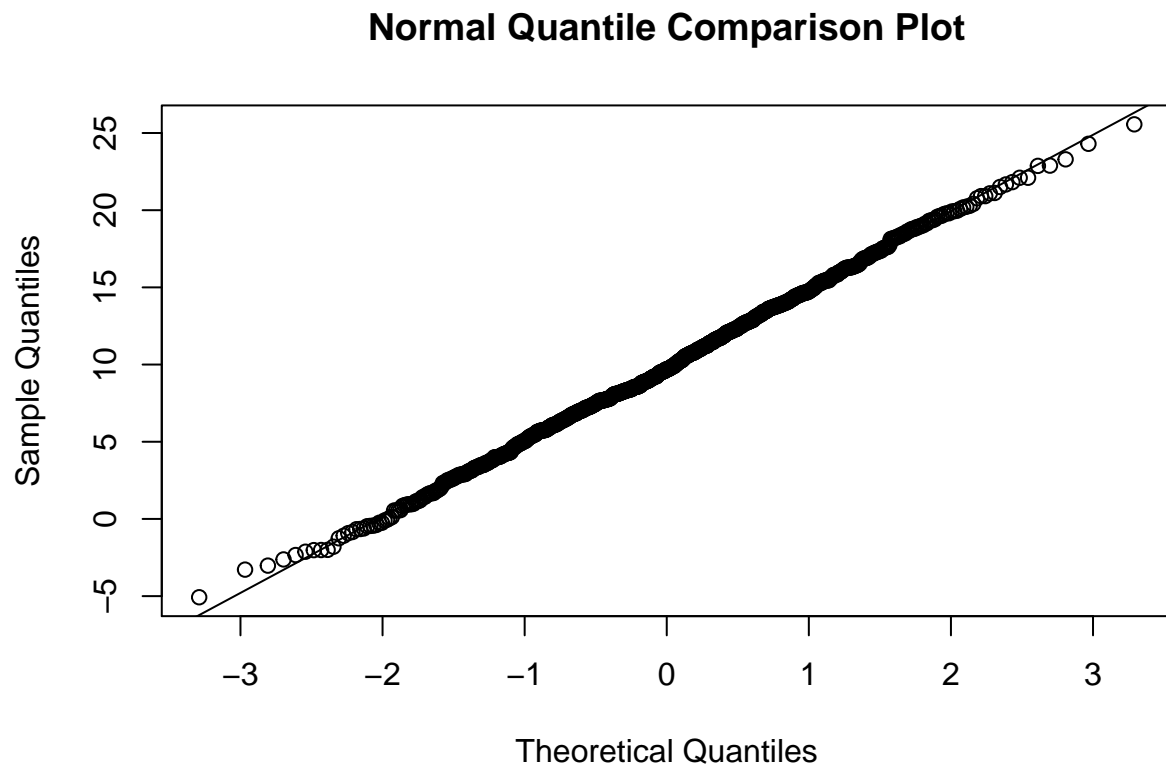


Figure 1: Density of normal distribution

Another way to check whether a variable is normally distributed is the “normal quantile comparison plot”. The more tightly our data points hug the diagonal line, the more normally distributed it is.

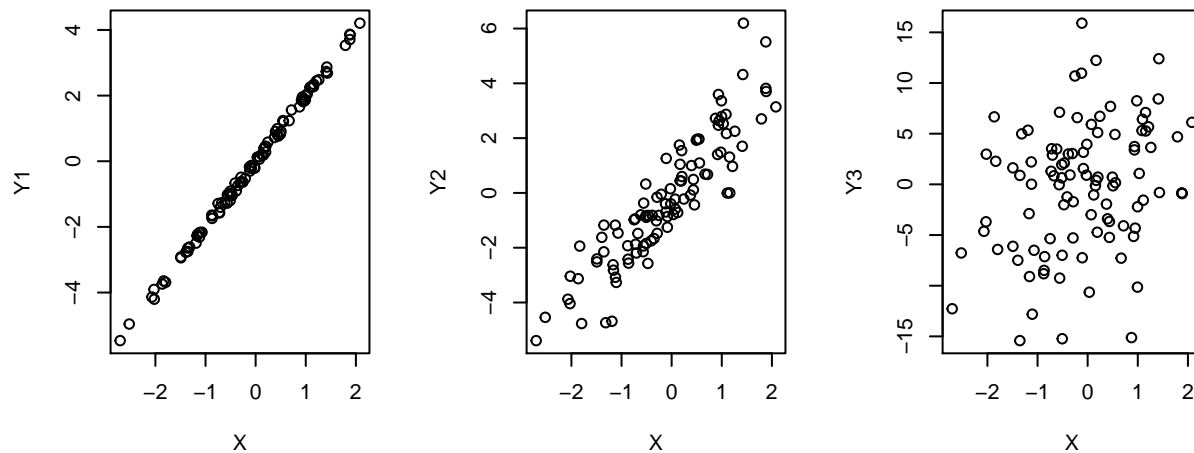
```
qqnorm(normdraws, main="Normal Quantile Comparison Plot")
qqline(normdraws)
```



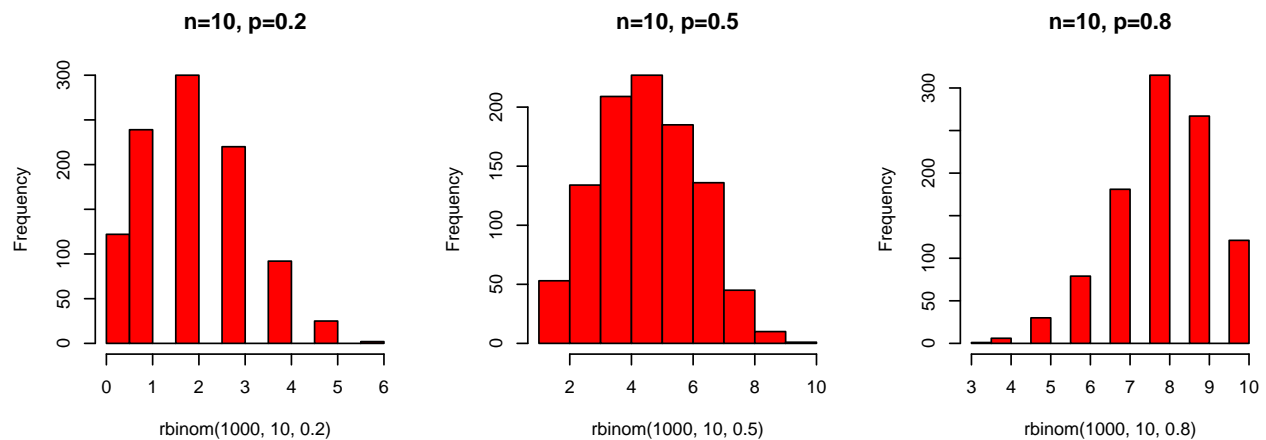
Inspecting relationship with scatterplot

```
X <- rnorm(n = 100)
Y1 <- 2 * X + rnorm(length(X), sd=0.1)
Y2 <- 2 * X + rnorm(length(X), sd=1)
Y3 <- 2 * X + rnorm(length(X), sd=5)

par(mfrow=c(1, 3))
plot(X, Y1)
plot(X, Y2)
plot(X, Y3)
```



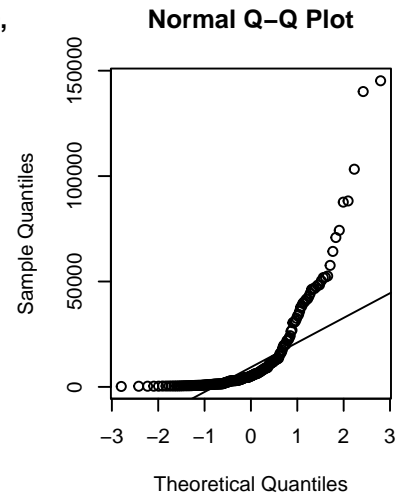
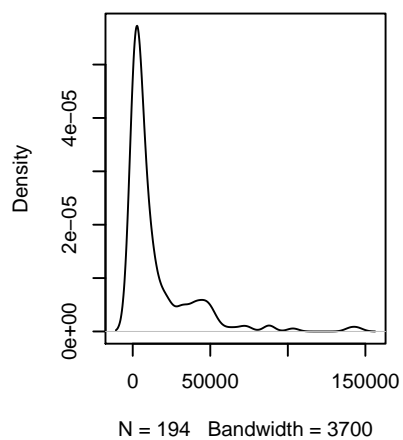
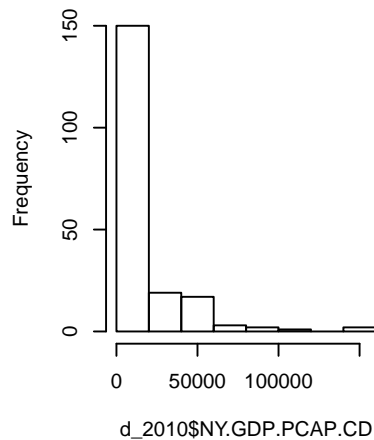
In-class exercise: Replicate binomial histogram in your lecture slides



In-class exercise: Plotting GDP per capita in 2010

Download GDP per capita data for all countries in 2010, using package WDI. Plot the histogram, density plot, and normal quantile comparison plot.

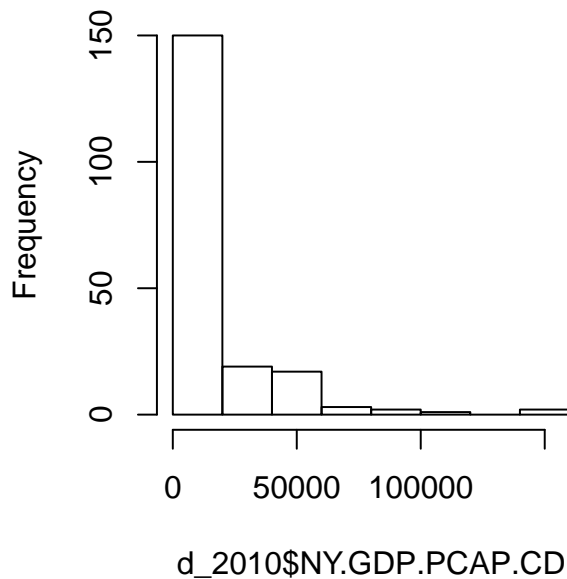
Histogram of d_2010\$NY.GDP.PCAP.Fault(x = d_2010\$NY.GDP.PCAP.CD,



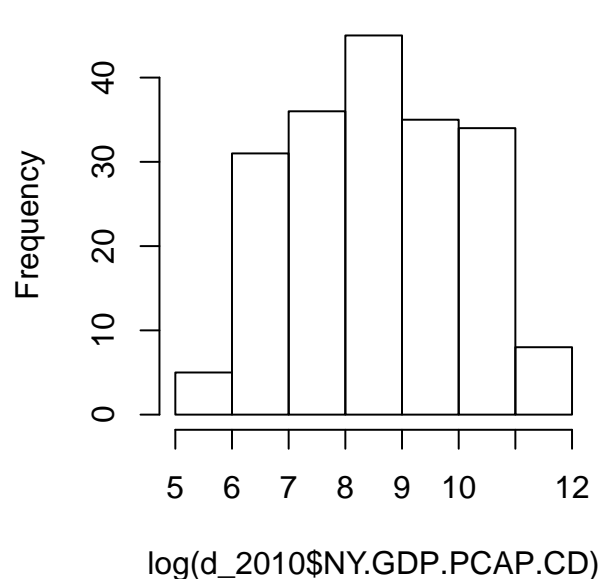
The distribution of GDP per capita has a long right tail. This is because a country's GDP per capita can go very high but cannot go lower than 0 (this phenomenon is called "left-censored"). Because of this, GDP per capita is NOT normally distributed, and can misbehave in models that assume normality. A common way to deal with this is to take the $\log(\text{GDP per capita})$ instead.

```
par(mfrow=c(1, 2))
hist(d_2010$NY.GDP.PCAP.CD, main="non log")
hist(log(d_2010$NY.GDP.PCAP.CD), main="log")
```

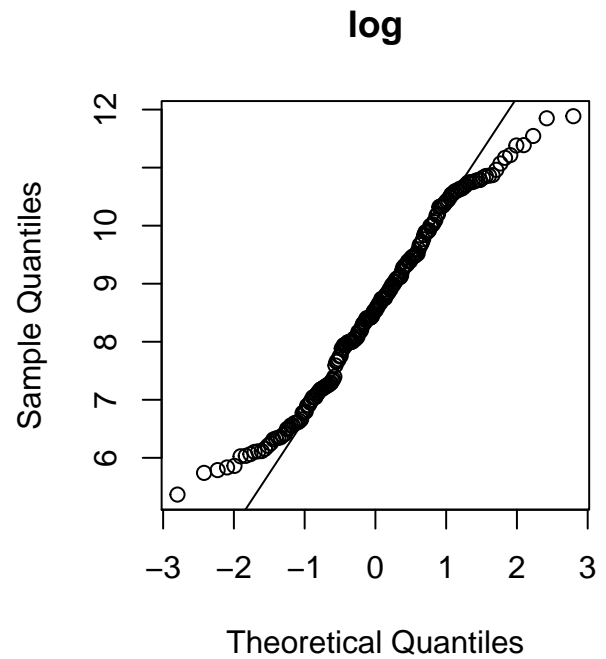
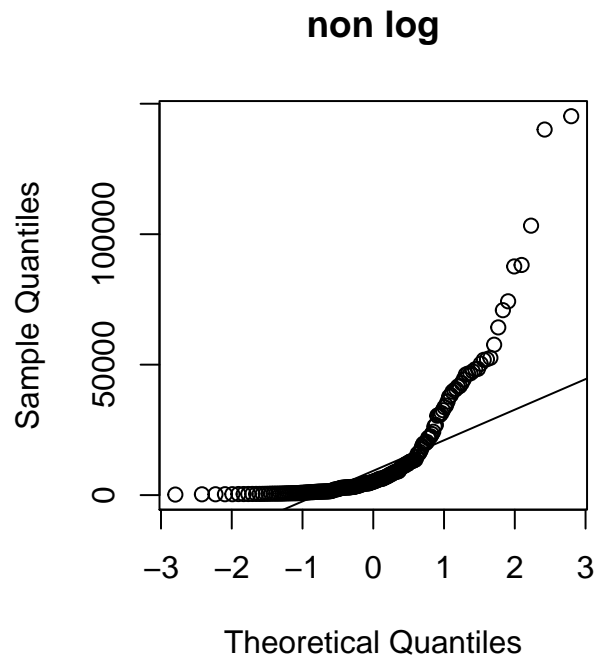
non log



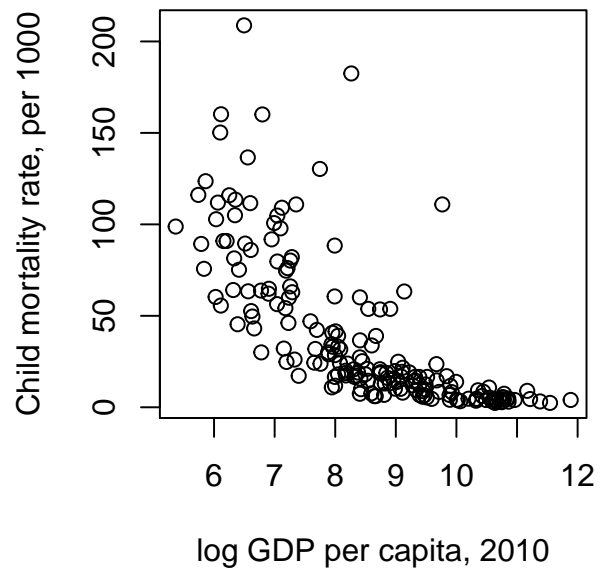
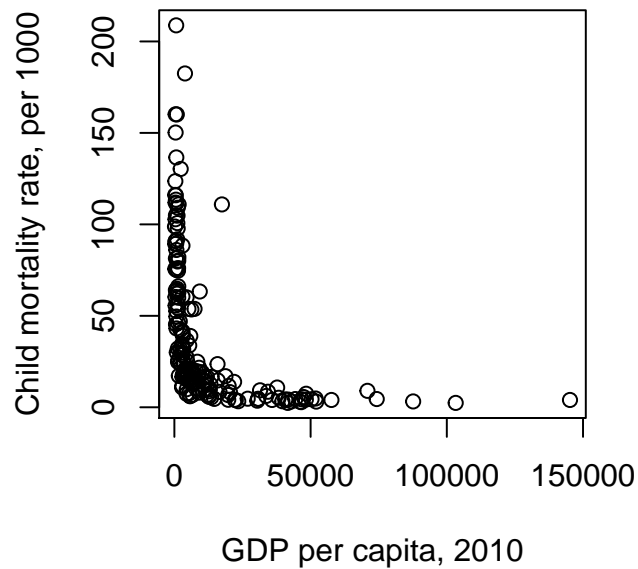
log



```
par(mfrow=c(1, 2))
qqnorm(d_2010$NY.GDP.PCAP.CD, main="non log")
qqline(d_2010$NY.GDP.PCAP.CD)
qqnorm(log(d_2010$NY.GDP.PCAP.CD), main="log")
qqline(log(d_2010$NY.GDP.PCAP.CD))
```

In-class exercise: Plot the relationship between GDP per capita and child mortality (“Mortality rate, under-5 (per 1000 live births)”)



4. Tips and tricks

1. You can name your knitr chunk
2. You can divide your R code into sections