

Tutorial 2 Solutions

STAT 3013/40278027

1. SI: 2.1, 2.2, 2.4. **Ans.** See the handwritten pages.

2. Consider the following:

- a. Visually display the data and discuss. Try taking the natural log of the data (when statisticians say “log” they mean natural log).
- b. Compute a six number summary of the data.
- c. Based on the “box plot rule”, determine if there are any outliers. Which countries are outliers? To use the rule examine the following: Are any values in the data below the 1st Quartile - 1.5 IQR? Are any values in the data above the 3rd Quartile + 1.5 IQR? IQR is the inter-quartile range.
- d. Let $Y = \log(\text{GDP})$. Suppose $Y \sim \text{normal}(\mu, \sigma^2)$. What is your best guess for μ and σ^2 as functions of Y (call these T_1 and T_2)? What are the means (expected values) of T_1 and T_2 ?

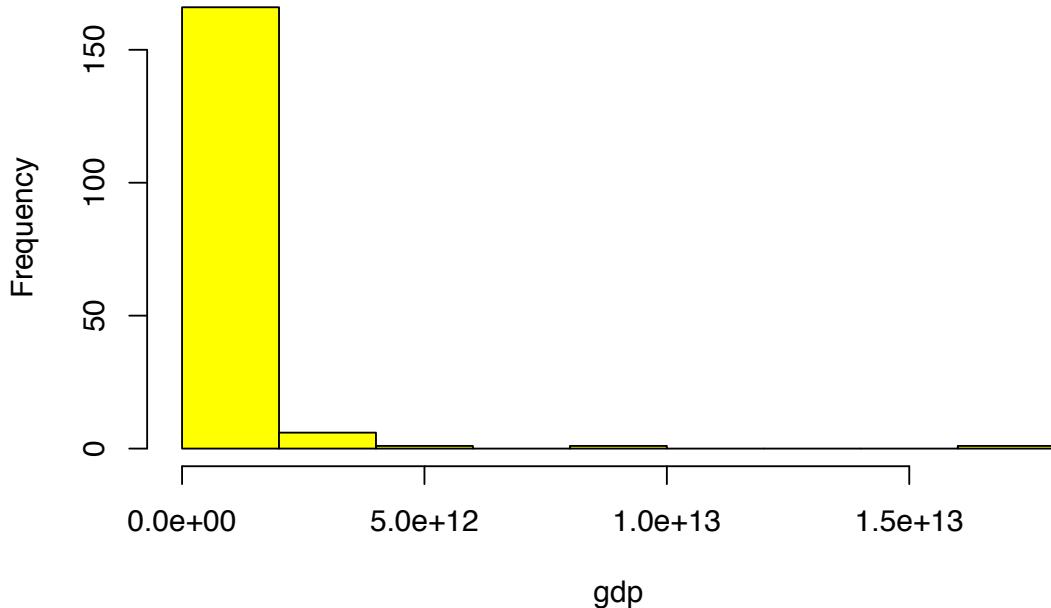
Ans. First let's load in the data. Note that I removed the missing values (NA's). Missing data is an extensive and important topic in statistics; as such be careful about removing missing data. At the very least, discuss your full sample of data and then which cases were removed due to missing data.

```
## GDP
gdp2013 <- read.table("gdp2013.txt", header=TRUE)
D <- gdp2013
D <- na.omit(D)
gdp <- D$Y2013
```

The histogram is unimodal, right skewed and appears to have some outliers.

```
## visual display
hist(gdp, col="yellow")
```

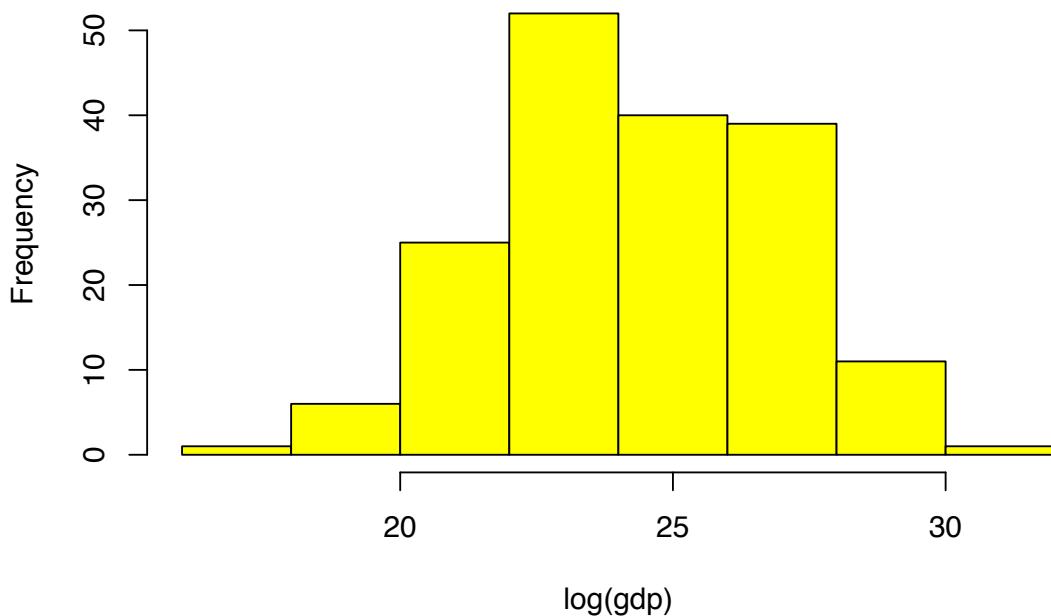
Histogram of gdp



If we take the log of the data it is far more symmetric.

```
## visual display  
hist(log(gdp), col="yellow")
```

Histogram of log(gdp)



Let's get a six number summary (using the logged data)

```
log.gdp <- log(gdp)
summary(log.gdp)

##      Min. 1st Qu. Median    Mean 3rd Qu.    Max.
## 17.46   22.73   24.15   24.23   26.13   30.45
```

Let's get the quartiles and the IQR.

```
##
Q <- quantile(log.gdp, prob=c(0.25, 0.5, 0.75), type=6)
Q

##      25%      50%      75%
## 22.71874 24.14522 26.13676
IQR <- unname(Q[3]-Q[1])
IQR

## [1] 3.418021
```

Let's identify outliers based on the box plot method.

```
##
high.outliers <- log.gdp[log.gdp > Q[3] + 1.5*IQR]
low.outliers <- log.gdp[log.gdp < Q[1] - 1.5*IQR]

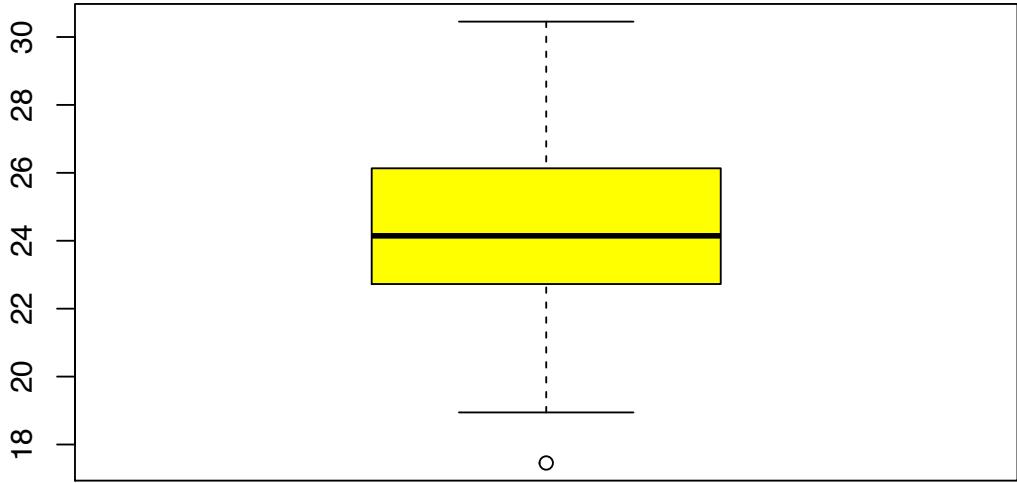
high.outliers

## numeric(0)
low.outliers

## [1] 17.45664
# let's find the countries that have high outliers
D[log.gdp >= min(high.outliers),1]

## Warning in min(high.outliers): no non-missing arguments to min; returning
## Inf
## factor(0)
## 214 Levels: Afghanistan Albania Algeria American Samoa Andorra ... Zimbabwe
# let's find the countries that have high outliers
D[log.gdp <= max(low.outliers),1]

## [1] Tuvalu
## 214 Levels: Afghanistan Albania Algeria American Samoa Andorra ... Zimbabwe
# let's make a boxplot.
boxplot(log.gdp, col="yellow")
```



- If we assume the data $Y_1, \dots, Y_n \sim iid n(\mu, \sigma^2)$ then reasonable guesses for the population mean and variance are the sample mean and variances. So we have:

$$\bar{y} = 24.2263843 = \hat{\mu}$$

$$S^2 = 6.0934187 = \hat{\sigma}^2$$

A nice property of these estimators is that $E[\hat{\mu}] = E[\bar{Y}] = \mu$ and $E[\hat{\sigma}^2] = E[S^2] = \sigma^2$. We proved the latter in lecture and for fun let's do the former:

$$\begin{aligned}
 E[\hat{\mu}] &= E[\bar{Y}] &= E[(1/n)(Y_1 + \dots + Y_n)] \\
 &= (1/n)E[(Y_1 + \dots + Y_n)] \\
 &= (1/n)(E[Y_1] + \dots + E[Y_n]) \\
 &= (1/n)nE[Y_1] = (1/n)n\mu = \mu
 \end{aligned}$$

6JJ Q 2.1.) x_1, \dots, x_n iid $\text{Unif}(0, \theta)$

$$f(x) = \frac{1}{\theta} = \frac{1}{\theta}$$

$Z \sim \text{Unif}(a, b)$

$$f(z) = \frac{1}{b-a}$$

See the table of distributions

a.) Find the density of $x_{(n)} = \max(x_1, \dots, x_n)$

- Let $Y = \max(x_1, \dots, x_n)$
- $P(Y \leq y) = P(x_{(n)} \leq y)$ If the largest x_i is less than y , then all x_i s are less than y .

$$= P(x_1 \leq y, x_2 \leq y, \dots, x_n \leq y)$$

$$= P(x_1 \leq y) \times \dots \times P(x_n \leq y)$$
 B/c of independence

$$= \underbrace{[P(x_i \leq y)]^n}_{\text{CDF}}$$
 B/c identically distributed

$$F(y) = P(x_i \leq y) = \int_0^y \frac{1}{\theta} dx = \frac{y}{\theta}$$

$$\therefore P(x_{(n)} \leq y) = \left[\frac{y}{\theta} \right]^n \quad \therefore f(y) = \frac{n y^{n-1}}{\theta^n}$$

b.) If $MSE(T(\underline{x})) \rightarrow 0$ as $n \rightarrow \infty$, then $T(\underline{x})$ is a Consistent estimator.

$$\bullet T_1(\underline{x}) = z \bar{x}$$

$$\Rightarrow E(z \bar{x}) = z E\left(\frac{1}{n}(x_1 + \dots + x_n)\right) = \frac{z}{n} E(x_1 + \dots + x_n)$$

$$= \frac{z}{n} \sum_{i=1}^n E(x_i) = \frac{z}{n} n \left[\frac{0+\theta}{2} \right] = \theta \quad \therefore \text{unbiased}$$

$$\Rightarrow V(z \bar{x}) = 4 V\left(\frac{1}{n}(x_1 + \dots + x_n)\right) = \frac{4}{n^2} V(x_1 + \dots + x_n)$$

Independence means there are no Covariance terms

$$= \frac{4}{n^2} n V(x) = \frac{4}{n} \left[\frac{(\theta-0)^2}{12} \right] = \frac{\theta^2}{3n} \rightarrow 0 \quad \text{as } n \rightarrow \infty$$

$$\Rightarrow MSE(\bar{x}) = v(\bar{x}) + [Bias(\bar{x})]^2$$

$$= \frac{\theta^2}{3n} + \sigma^2 \rightarrow 0 \text{ as } n \rightarrow \infty$$

$\therefore \bar{x}$ is a consistent estimator of θ

$$\bullet T_2(\underline{x}) = \frac{(n+1)x_{(n)}}{n}$$

$$\Rightarrow E\left(\frac{(n+1)x_{(n)}}{n}\right) = \frac{(n+1)}{n} E(x_{(n)})$$

$$\Rightarrow E(x_{(n)}) = \int_0^\theta x \frac{n x^{n-1}}{\theta^n} dx = \frac{n}{\theta^n} \int_0^\theta x^n dx$$

$$= \frac{n}{\theta^n} \left[\frac{1}{n+1} x^{n+1} \right] \Big|_0^\theta = \frac{n}{n+1} \theta$$

$$\therefore E(T_2) = \frac{(n+1)}{n} \frac{n}{(n+1)} \theta = \theta \therefore \text{unbiased}$$

$$\Rightarrow V(T_2(\underline{x})) = \frac{(n+1)^2}{n^2} V(x_{(n)}) \quad | \quad v(x) = E(x^2) - [E(x)]^2$$

$$\Rightarrow E(x_{(n)}^2) = \int_0^\theta x^2 \frac{n x^{n-1}}{\theta^n} dx = \frac{n}{\theta^n} \int_0^\theta x^{n+1} dx$$

$$= \frac{n}{\theta^n} \frac{1}{n+2} x^{n+2} \Big|_0^\theta = \frac{n}{\theta^n} \frac{\theta^{n+2}}{n+2} = \frac{n}{n+2} \theta^2$$

$$\Rightarrow v(T_2) = E(T_2^2) - [E(T_2)]^2$$

$$= \frac{(n+1)^2}{n(n+2)} \theta^2 - \theta^2 = \frac{1}{n(n+2)} \theta^2$$

$$\therefore MSE(T_2) = v(T_2) + [Bias(T_2)]^2$$

$$= \frac{1}{n(n+2)} \theta^2 + \sigma^2 \rightarrow 0 \text{ as } n \rightarrow \infty$$

$\therefore T_2$ is a consistent estimator of θ

• Which estimator should we choose?

◦ Both are unbiased and consistent

∴ Let's compare based on their variances,
which in this case is the same as
comparing their MSEs.

$$V(T_1) = \frac{\theta^2}{3n} \quad \text{vs} \quad V(T_2) = \frac{\theta^2}{n(n+2)}$$

$$\begin{aligned} V(T_2) - V(T_1) &= \frac{\theta^2}{n(n+2)} - \frac{\theta^2}{3n} \\ &= \left[\frac{1}{n(n+2)} - \frac{1}{3n} \right] \theta^2 \\ &\stackrel{\text{brace}}{=} 0 \text{ for } n=1 \\ &< 0 \text{ for } n \geq 2. \end{aligned}$$

$V(T_2) < V(T_1)$ for $n \geq 2$ ∴ choose T_2 .

GJJ Q 2.2) $\hat{\theta}_1, \hat{\theta}_2$ are independent and unbiased estimators of θ .

$$V(\hat{\theta}_1) = \sigma_1^2, V(\hat{\theta}_2) = \sigma_2^2$$

- Consider: $\hat{\theta} = k_1 \hat{\theta}_1 + k_2 \hat{\theta}_2$

\rightarrow we want an estimator which is unbiased but has the smallest variance.

$$1.) E(\hat{\theta}) = k_1 E(\hat{\theta}_1) + k_2 E(\hat{\theta}_2) = \theta$$

$$= k_1 \theta + k_2 \theta = \theta$$

$$\therefore k_1 + k_2 = 1$$

$$2.) \text{Now } V(\hat{\theta}) = k_1^2 V(\hat{\theta}_1) + k_2^2 V(\hat{\theta}_2)$$

$$= k_1^2 \sigma_1^2 + k_2^2 \sigma_2^2$$

$$\Rightarrow \min_{k_1, k_2} k_1^2 \sigma_1^2 + k_2^2 \sigma_2^2 \text{ s.t. } k_1 + k_2 = 1$$

such that

$$\Rightarrow \min_{k_1} k_1^2 \sigma_1^2 + (1-k_1)^2 \sigma_2^2$$

$$\frac{d}{dk_1} = 2k_1 \sigma_1^2 - 2(1-k_1) \sigma_2^2 = 0$$

$$\Rightarrow k_1 \sigma_1^2 - \sigma_2^2 + k_1 \sigma_2^2 = 0$$

$$k_1 (\sigma_1^2 + \sigma_2^2) = \sigma_2^2$$

$$\therefore k_1 = \frac{\sigma_2^2}{\sigma_1^2 + \sigma_2^2} \quad \& \quad k_2 = \frac{\sigma_1^2}{\sigma_1^2 + \sigma_2^2}$$

Do we have a minimum?

$$\frac{d^2}{dk_1^2} = 2\sigma_1^2 + 2\sigma_2^2 > 0 \quad \therefore \text{we found the minimum.}$$

GJJ Q 2.4) Recall Def 2.2] $\hat{\theta}$ is weakly consistent if $P(|\hat{\theta} - \theta| > \varepsilon) \rightarrow 0$ as $n \rightarrow \infty$.

\rightarrow we can show: If $v(\hat{\theta}) \rightarrow 0$ & $\text{Bias}(\hat{\theta}) \rightarrow 0$

$\therefore \text{MSE}(\hat{\theta}) \rightarrow 0$ as $n \rightarrow \infty$

then $\hat{\theta}$ is a consistent estimator of θ .

$$\begin{aligned} \rightarrow P(\hat{\theta} = \theta) &= \frac{(n-1)}{n} \\ P(\hat{\theta} = \theta+n) &= \frac{1}{n} \end{aligned} \quad \left. \begin{array}{l} \text{we only have} \\ \text{two options for } \hat{\theta} \end{array} \right\}$$

- Let's consider the direct approach:

$$\begin{aligned} \lim_{n \rightarrow \infty} P(|\hat{\theta} - \theta| > \varepsilon) &= \lim_{n \rightarrow \infty} P(\hat{\theta} = \theta+n) \\ &= \lim_{n \rightarrow \infty} \frac{1}{n} = 0 \end{aligned}$$

$\therefore \hat{\theta}$ is consistent estimator

- Note: $E(\hat{\theta}) = \theta P(\hat{\theta} = \theta) + (\theta+n) P(\hat{\theta} = \theta+n)$

$$\begin{aligned} &\text{Possibility} \quad \text{Probability} \\ &= \theta \frac{(n-1)}{n} + (\theta+n) \frac{1}{n} \\ &= \theta + 1 \end{aligned}$$

$\rightarrow 0$ as $n \rightarrow \infty \therefore \text{Bias} \rightarrow 0$.