

## Question1

Q1)

$$\begin{aligned} 1.1) \quad \bar{x} &= 10.467 \\ \sigma^2 &= 7.559 \\ \sigma &= 2.749 \end{aligned} \quad \left. \vphantom{\begin{aligned} \bar{x} &= 10.467 \\ \sigma^2 &= 7.559 \\ \sigma &= 2.749 \end{aligned}} \right\} \text{ by using R}$$

For CI with unknown mean and unknown variance, we use:

$$\mu \in \left( \bar{x} - t_{\alpha/2, n-1} \frac{s}{\sqrt{n}}, \bar{x} + t_{\alpha/2, n-1} \frac{s}{\sqrt{n}} \right)$$

$$95\% \Rightarrow \frac{\alpha}{2} = 0.025$$

$$df = n - 1 = 19$$

$$t_{\alpha/2, n-1} = t_{0.025, 19} = 2.09302 \quad (\text{by using } t \text{ table})$$

Then,

$$\begin{aligned} \mu \in \left( 10.467 - (2.09302) \frac{2.749}{\sqrt{20}}, 10.467 + (2.09302) \frac{2.749}{\sqrt{20}} \right) \\ = (9.18, 11.75) \end{aligned}$$

∴ The estimated mean fuel efficiency of vehicles that are all-wheel driver (sample size  $n = 20$ ) is 10.467. We are 95% confident the population mean fuel efficiency for this group is between 9.18 and 11.75.

1.1

R code:

```
fueldb <- read.csv("fuel.efficiency.csv", header = TRUE)
print(head(fueldb))

a_fuel <- fueldb[fueldb$Type == "A", "FA"]
a_mean <- mean(a_fuel)
a_var <- var(a_fuel) #default is the unbiased variance
a_sd <- sqrt(a_var)
a_df <- length(a_fuel)-1

#alpha/2 = 0.025, df = 20-1 = 19
# t = 2.09302

low_a_ci <- a_mean - (2.09302 * (a_sd/sqrt(a_df+1)))
high_a_ci <- a_mean + (2.09302 * (a_sd/sqrt(a_df+1)))
```

$$\begin{array}{lcl}
 1.2) \quad \hat{\mu}_A = 10.467 & \left. \begin{array}{l} \hat{\mu}_P = 8.772 \\ \hat{\sigma}_A^2 = 7.559 \\ \hat{\sigma}_P^2 = 9.387 \end{array} \right\} & \begin{array}{l} n_A = 20 \\ n_P = 25 \\ 95\% \Rightarrow 1.96 \text{ in value} \end{array} \\
 & & \text{by using R}
 \end{array}$$

difference in fuel efficiency between A and P:

$$\begin{aligned}
 \hat{\mu}_A - \hat{\mu}_P &= 10.467 - 8.772 \\
 &= 1.695
 \end{aligned}$$

interval:

$$\left( \hat{\mu}_A - \hat{\mu}_P - z_{\alpha/2} \sqrt{\frac{\hat{\sigma}_A^2}{n_A} + \frac{\hat{\sigma}_P^2}{n_P}}, \hat{\mu}_A - \hat{\mu}_P + z_{\alpha/2} \sqrt{\frac{\hat{\sigma}_A^2}{n_A} + \frac{\hat{\sigma}_P^2}{n_P}} \right)$$

hence,

$$\begin{aligned}
 &\left( 1.695 - 1.96 \sqrt{\frac{7.559}{20} + \frac{9.387}{25}}, 1.695 + 1.96 \sqrt{\frac{7.559}{20} + \frac{9.387}{25}} \right) \\
 &= (-0.006287, 3.3963) \\
 &\quad -0.00576
 \end{aligned}$$

$\therefore$  The estimated difference in mean fuel efficiency between all-wheel drive vehicles (sample size  $n = 20$ ) and part-time four-wheel drive vehicles (sample size  $n = 25$ ) is 1.695. We are 95% confident the population mean difference in fuel efficiency is between -0.0062.87 up to 3.3963. As the interval includes zero, we cannot rule out the possibility of there being no difference at a population level between all-wheel drive vehicles and part-time four-wheel drive vehicles.

1.2

R code:

```
p_fuel <- fueldb[fueldb$Type == "P", "FA"]
```

```
p_mean <- mean(p_fuel)
```

```
p_var <- var(p_fuel) #default is the unbiased variance
```

```
# alpha/2 in z score = 1.96 (95%)
```

```
low_mean_diff_ci <- (a_mean - p_mean) -
1.96*(sqrt((a_var/length(a_fuel))+(p_var/length(p_fuel))))
```

```
high_mean_diff_ci <- (a_mean - p_mean) +
1.96*(sqrt((a_var/length(a_fuel))+(p_var/length(p_fuel))))
```

1.3)

$$H_0 : \mu_A < \mu_P$$

$$H_A : \mu_A \geq \mu_P$$

$$\hat{\mu}_A = 10.467 \quad n_A = 20$$

$$\hat{\mu}_P = 8.772 \quad n_P = 25$$

$$\hat{\sigma}_A^2 = 7.559$$

$$\hat{\sigma}_P^2 = 9.387$$

we know that  $\sigma_A^2 \neq \sigma_P^2$

we get test statistic :

$$Z(\hat{\mu}_A - \hat{\mu}_P) = \frac{\hat{\mu}_A - \hat{\mu}_P}{\sqrt{\frac{\hat{\sigma}_A^2}{n_A} + \frac{\hat{\sigma}_P^2}{n_P}}}$$

$$\begin{aligned} \Rightarrow Z(\hat{\mu}_A - \hat{\mu}_P) &= \frac{10.467 - 8.772}{\sqrt{\frac{7.559}{20} + \frac{9.387}{25}}} \\ &= \frac{1.695}{0.868} \\ &= 1.9528 \end{aligned}$$

$$p = 1 - P(Z < Z(\hat{\mu}_A - \hat{\mu}_P))$$

$$= 1 - 0.9746$$

$$= 0.0254$$

∴ The p-value is smaller than the significance level 0.05 (default). There is enough evidence to reject the null hypothesis that all-wheel-drives are less efficient than part-time four-wheel-drive vehicles.

1.3

R code:

```
my_z <- (a_mean - p_mean) / (sqrt((a_var/length(a_fuel)) +  
(p_var/length(p_fuel))))
```

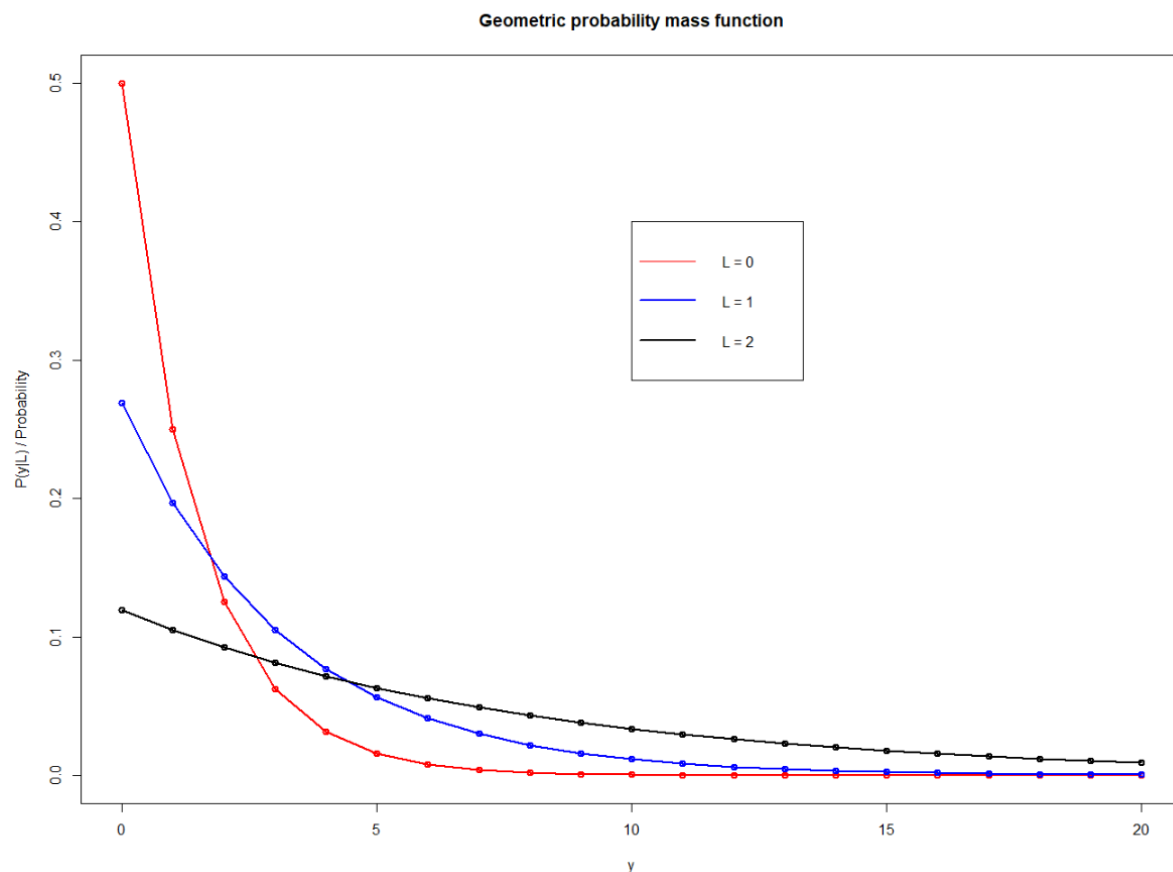
```
z_p = 1-pnorm(my_z)
```

```
pnorm(1.9528)
```

1.4) For the estimated difference in fuel, we cannot rule out the possibility of there being no difference at population level between all-wheel drive vehicles and part-time four-wheel drive vehicles. This is because the interval includes zero. For estimated p-value, it's only an estimated p-value and the method we used is very close to true p-value with the size of population being moderate only. If the size of population grows larger the closure of p-value will be lower.

## Question2

2.1)



2.1

**R code:**

```
cal_prob <- function(L,y){  
  res <- ((exp(L) + 1)^(-y-1)) * exp(y*L)  
  return(res)  
}
```

```

my_y <- seq(0,20,1)
prob1 <- cal_prob(0,my_y)
plot(my_y, prob1, col = "red", lwd=2.5, type = "o", xlab = "y", ylab =
"P(y|L) / Probability", main = "Geometric probability mass
function" )
prob2 <- cal_prob(1,my_y)
lines(my_y, prob2, col = "blue", lwd=2.5, type = "o")
prob3 <- cal_prob(2, my_y)
lines(my_y, prob3, col = "black", lwd=2.5, type = "o")
legend(x=10,y=0.4,c("L = 0","L = 1","L = 2"), lty=c(1,1,1),
      pch=c("", "", "" ), col=c("red","blue","black"),
      lwd=c(1,2.5,2.5))

```

\* for entire  $\mathbb{Q}_2$ ,  $\log_e e = \log_e e$ , which is equal to 1

(2.2) Let  $L_s$  be the parameter:  $L$ , the log-odds of seeing a failure (tail) when the coin is tossed.

2.2) joint probability:

$$\begin{aligned}
 P(y|L_s) &= \prod_{i=1}^n P(y_i|L_s) \\
 &= (e^{L_s} + 1)^{-y_1 - 1} e^{y_1 L_s} \cdot (e^{L_s} + 1)^{-y_2 - 1} e^{y_2 L_s} \cdots (e^{L_s} + 1)^{-y_n - 1} e^{y_n L_s} \\
 &= (e^{L_s} + 1)^{-y_1 - y_2 - \cdots - y_n} (e^{L_s} + 1)^{-n} e^{n L_s} \cdot e^{y_1 + y_2 + \cdots + y_n} \\
 &= \frac{e^{n L_s}}{(e^{L_s} + 1)^n} \cdot \sum_{i=1}^n \frac{1}{(e^{L_s} + 1)^{y_i}} \cdot \sum_{i=1}^n e^{y_i}
 \end{aligned}$$

2.3)  $\therefore \frac{e^{n L}}{(e^{L+1})^n} \cdot \sum_{i=1}^n \frac{1}{(e^{L+1})^{y_i}} \cdot \sum_{i=1}^n e^{y_i}$ , where  $L$  is the parameter

$$\begin{aligned}
 L(y|L_s) &= -\log P(y|L_s) \\
 &= -\log \left[ \frac{e^{n L_s}}{(e^{L_s} + 1)^n} \cdot \sum_{i=1}^n \frac{1}{(e^{L_s} + 1)^{y_i}} \cdot \sum_{i=1}^n e^{y_i} \right] \\
 &= \log \frac{(e^{L_s} + 1)^n}{e^{n L_s}} + \sum_{i=1}^n \log (e^{L_s} + 1)^{y_i} + \sum_{i=1}^n -y_i \\
 &= \log (e^{L_s} + 1)^n - \log e^{n L_s} + \sum_{i=1}^n \log (e^{L_s} + 1)^{y_i} - \sum_{i=1}^n y_i \\
 &= n \log (e^{L_s} + 1) - n L_s + \sum_{i=1}^n y_i \log (e^{L_s} + 1) - \sum_{i=1}^n y_i
 \end{aligned}$$

$\therefore n \log (e^{L+1}) - n L + \sum_{i=1}^n y_i \log (e^{L+1}) - \sum_{i=1}^n y_i$ , where  $L$  is the parameter

2.4) maximum likelihood estimator  $\hat{L}$  for  $L$ :

let  $L_{lg}$  be the negative log-likelihood

$$\begin{aligned}\frac{dL_{lg}(y|L)}{dL} &= \sum_{i=1}^n y_i \frac{1}{L^{L+1}} + \frac{n}{L^{L+1}} - n \\&= \frac{\sum_{i=1}^n y_i}{L^{L+1}} + \frac{n}{L^{L+1}} - n \\&= \frac{n + \sum_{i=1}^n y_i}{L^{L+1}} - n \\&= \frac{n - n(L^{L+1}) + \sum_{i=1}^n y_i}{L^{L+1}} = \frac{n(-L^L) + \sum_{i=1}^n y_i}{L^{L+1}}\end{aligned}$$

Set derivative to 0:

$$\frac{n(-L^L) + \sum_{i=1}^n y_i}{L^{L+1}} = 0$$

$$L^{L+1} = 0 \quad \text{or} \quad n(-L^L) + \sum_{i=1}^n y_i = 0$$

$$L^L = -1$$

$$-nL^L = -\sum_{i=1}^n y_i$$

$$\log L^L = \log(-1)$$

hence, rejected

$$L^L = \frac{1}{n} \sum_{i=1}^n y_i$$

$$L(\log L) = \log \frac{1}{n} \sum_{i=1}^n y_i$$

$$L = \log \frac{1}{n} \sum_{i=1}^n y_i$$



2.5)

Assumptions : the population  $L$

i) have a mean of  $E[y_i] = \mu$

ii) have a variance of  $V[y_i] = \sigma^2$

iii) and are independent

bias :

$$b_L(\hat{L}) = E[\hat{L}(y)] - L$$

$$\text{Var}_L(\hat{L}) = E[(\hat{L}(y) - E[\hat{L}(y)])^2] = V[\hat{L}(y)]$$

$$= \frac{\sigma^2}{n}$$

$$= \frac{\sigma^2(e^L + 1)}{n}$$

$$\text{Let } X = \log \sum_{i=1}^n y_i - \log \sum_{i=1}^n (1), \quad X \sim e^L, \quad E(X) = e^L, \quad V(X) = e^L(e^L + 1)$$

$$f(X) = \log X$$

$$\frac{df(X)}{dX} = \frac{1}{X}$$

$$\frac{d^2f(X)}{dX^2} = -\frac{1}{X^2}$$

$$V[f(X)] = \left[ \frac{df(X)}{dX} \Big|_{X=\mu_X} \right]^2 \sigma_X^2$$

$$= \left( \frac{1}{e^L} \right)^2 e^L(e^L + 1)$$

$$= \frac{e^L + 1}{e^L}$$

Variance :

$$E[f(X)] = f(\mu_X) + \left[ -\frac{1}{X^2} \Big|_{X=\mu_X} \right] \frac{\sigma_X^2}{2}$$

$$= f(\mu_X) + \left[ -\frac{1}{\mu_X^2} \right] \frac{\sigma_X^2}{2}$$

$$= \log e^L + \left[ -\frac{1}{e^{2L}} \right] \frac{e^L(e^L + 1)}{2}$$

$$= L + \left[ -\frac{1}{e^{2L}} \right] \frac{e^{2L} + e^L}{2}$$

$$= L + \left[ -\frac{1 \cdot (e^{2L} + e^L)}{e^{2L} \cdot 2} \right]$$

$$= L - \frac{e^L + 1}{2e^L}$$

$$V_L(\hat{L}) = V\left[\log\left(\frac{1}{n} \sum_{i=1}^n y_i\right)\right] - L$$

$$= V[f(X)]$$

$$= \frac{e^L + 1}{e^L}$$

Bias :

$$b_L(\hat{L}) = E[\hat{L}(y)] - L$$

$$= E\left[\log\left(\frac{1}{n} \sum_{i=1}^n y_i\right)\right] - L$$

$$= E[\log(X)] - L$$

$$= L - \frac{e^L + 1}{2e^L} - L$$

$$= -\frac{e^L + 1}{2e^L}$$

## Question3

Q3

3.1)  $X \sim \text{Be}(\theta)$ ,  $\theta$  is the probability of a full moon day experiencing an above average number of dog bite admissions

$$n = 26$$

$$\hat{\theta} = \frac{11}{26}$$

$$V[X] = \left(\frac{11}{26}\right)\left(1 - \frac{11}{26}\right)$$
$$= 0.244 / \frac{165}{676}$$

95%  $\Rightarrow$  1.96 in value

Confidence Interval:

$$\left( \frac{11}{26} - 1.96 \sqrt{\frac{\frac{165}{676}}{26}}, \frac{11}{26} + 1.96 \sqrt{\frac{\frac{165}{676}}{26}} \right)$$

$$= (0.233, 0.613)$$

∴ The estimated probability of a full moon day experiencing an above average number of dog bite admissions (sample size  $n=26$ ) is 0.244.

We are 95% confident the population probability dog bite admissions for this group is between 0.233 and 0.613.

3.1

R code:

#3.1

```
my_var <- 11/26 * (1-11/26)
```

```
my_lowci <- 11/26 - 1.96*(sqrt(my_var/26))
```

```
my_highci <- 11/26 + 1.96*(sqrt(my_var/26))
```



3.2)

let  $A$  = number of dog bite admissions of full moon days

let  $B$  = number of dog bite admissions of non-full moon days

$$H_0 : \theta_A = \theta_B$$

$$H_A : \theta_A \neq \theta_B$$

$$\theta_A = \frac{11}{26}$$

$$\theta_B = 0.53$$

$$Z_{\hat{\theta}} = \frac{\frac{11}{26} - 0.53}{\sqrt{\frac{0.53(1-0.53)}{26}}}$$

$$= \frac{-139/1300}{0.09788}$$

$$= -1.0924$$

$$P = 2P(Z < -|z_{\hat{\theta}}|)$$

$$= 2P(Z < -1.0924)$$

$$= 2(0.13733) \quad \text{by using } \text{pnorm}(-1.0924) \text{ in R}$$

$$= 0.274$$

$\therefore$  The  $p$ -value is greater than the significant level 0.05 (default).  
There is insufficient evidence to reject the null hypothesis that  
there's no difference in the probability of experiencing an  
above average number of dog bite admissions between full moon  
and non-full moon days.

3.2

R code:

```
the_z <- (11/26 - 0.53)/sqrt(0.53*(1-0.53)/26)
```

```
new_p <- 2*pnorm(-1.0924)
```

3.3) To find exact p-value:

$$\text{binom.test}(x=11, n=26, 0.53)$$

$$P = 0.3275$$

$\therefore$  The exact p-value is slightly greater than the approximated p-value.  
Yet the p-value is still greater than the significant level 0.05 (default).  
There's insufficient evidence to reject the null hypothesis that there's no difference in the probability of experiencing an above average number of dog bite admissions between full moon and non-full moon days.

3.4) Let A = number of dogbite admissions of full moon days

Let B = number of dogbite admissions of new moon days

$$H_0: \theta_A = \theta_B$$

$$H_A: \theta_A \neq \theta_B$$

We use a pooled estimate of  $\theta$ :

$$\hat{\theta}_p = \frac{11 + 20}{26 + 26} = \frac{31}{52} \approx 0.5962$$

$$\begin{aligned} Z(\hat{\theta}_A - \hat{\theta}_B) &= \frac{\frac{11}{26} - \frac{20}{26}}{\sqrt{\frac{31}{52} \left(1 - \frac{31}{52}\right) \left(\frac{1}{26} + \frac{1}{26}\right)}} \\ &= \frac{-\frac{9}{26}}{\sqrt{(0.2408) \frac{2}{26}}} \\ &= -2.5437 \end{aligned}$$

$$P = 2P(Z < -|Z(\hat{\theta}_A - \hat{\theta}_B)|)$$

$$= 2P(Z < -2.5437)$$

$$= 2P(0.005485)$$

$$= 0.011$$

$\therefore$  The p-value is smaller than the significant level 0.05 (default)

there is enough evidence to reject

the null hypothesis that the probability of

experiencing an above average number of dogbite admissions does not differ between days falling on the new moon and full moon.

3.4

R code:

```
my_Z <- (11/26 - 20/26)/(sqrt(31/52*(1-31/52)*2/26))
```

```
my_P <- 2*pnorm(-2.543629)
```