# Task 3

## Ariella Fuzaylov and Candice Djorno

### 4/3/2021

We transform the data Employee_B_probs into a matrix.

```r
Employee.B.Branch=read.csv("Employee_B_by_Branch.csv", header=TRUE)
Employee.B.overall=read.csv("Employee_B_overall.csv", header=TRUE)
Employee.B.probs=as.matrix(read.csv("Employee_B_probs.csv", header=FALSE))
```

We add col and row names for human readability.

```r
knitr::kable(head(Employee.B.Branch))
```

| Year_Month | Branch | Total_of_Ratings | m_i |
|---|---|---:|---:|
| 2011-12 | Disneyland_California | 4.469767 | 215 |
| 2011-12 | Disneyland_HongKong | 4.134831 | 89 |
| 2011-12 | Disneyland_Paris | 4.118421 | 76 |
| 2011-7 | Disneyland_California | 4.232877 | 73 |
| 2011-7 | Disneyland_HongKong | 4.083333 | 24 |
| 2011-7 | Disneyland_Paris | 3.791667 | 72 |

```r
knitr::kable(head(Employee.B.overall))
```

| Year_Month | Total_of_Ratings | m_i |
|---|---:|---:|
| 2011-12 | 4.321053 | 380 |
| 2011-7 | 4.023669 | 169 |
| 2012-10 | 4.365155 | 419 |
| 2012-3 | 4.351351 | 370 |
| 2014-2 | 4.253968 | 252 |
| 2015-5 | 4.264012 | 678 |

```r
month<-Employee.B.overall$Year_Month

colnames(Employee.B.probs)<-month
rownames(Employee.B.probs)<-month

knitr::kable(Employee.B.probs)
```

| | 2011-12 | 2011-7 | 2012-10 | 2012-3 | 2014-2 | 2015-5 | 2015-6 | 2015-7 | 2016-11 | 2016-2 | 2016-8 | 2017-11 | 2017-4 | 2018-12 | 2018-8 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 2011-12 | 0.059431 | | | | | | | | | | | | | | |
| 2011-7 | 0.006198 | | | | | | | | | | | | | | |
| 2012-10 | 0.010696 | | | | | | | | | | | | | | |
| 2012-3 | 0.014932 | | | | | | | | | | | | | | |
| 2014-2 | 0.004887 | | | | | | | | | | | | | | |
| 2015-5 | 0.007769 | | | | | | | | | | | | | | |
| 2015-6 | 0.007395 | | | | | | | | | | | | | | |
| 2015-7 | 0.013308 | | | | | | | | | | | | | | |
| 2016-11 | 0.008163 | | | | | | | | | | | | | | |
| 2016-2 | 0.006139 | | | | | | | | | | | | | | |
| 2016-8 | 0.013588 | | | | | | | | | | | | | | |
| 2017-11 | 0.008045 | | | | | | | | | | | | | | |
| 2017-4 | 0.008440 | | | | | | | | | | | | | | |
| 2018-12 | 0.007199 | | | | | | | | | | | | | | |
| 2018-8 | 0.011491 | | | | | | | | | | | | | | |

## Subtask 1

**Estimate average rating**

```r
my_design<-svydesign(id=~Year_Month,prob=~diag(Employee.B.probs),
                     fpc=~rep(15/112,15),
                     data=Employee.B.overall,
                     pps=ppsmat(Employee.B.probs))
svymean(~Total_of_Ratings,my_design)
```

```
##                    mean     SE
## Total_of_Ratings 4.2182 0.0308
```

The estimated average satisfaction rating overall for the population of 40,041 reviews is 4.2182.

**Confidence interval**

```
conf= confint(svymean(x=~Total_of_Ratings,design = my_design))
conf
```

```
##                      2.5 %   97.5 %
## Total_of_Ratings 4.157762 4.278607
```

A 95% confidence interval is $[4.1577624, 4.278607]$.

## Subtask 2

**Calculate Mean by Branch**

```
knitr::kable(Employee.B.Branch%>%
  summarise(n= n(), Mean= mean(Total_of_Ratings),Var=sd(Total_of_Ratings)^2),caption = "Rating Summary
```

Table 4: Rating Summary Statistics

| n | Mean | Var |
|---|------|-----|
| 45 | 4.187407 | 0.0490165 |

```
knitr::kable(Employee.B.Branch%>% group_by(Branch)%>%
  summarise(n= n(), Mean= mean(Total_of_Ratings),StD=sd(Total_of_Ratings)), caption = "Rating Summarise
```

Table 5: Rating Summarised by Branch

| Branch | n | Mean | StD |
|--------|---|------|-----|
| Disneyland_California | 15 | 4.391302 | 0.1133633 |
| Disneyland_HongKong | 15 | 4.164182 | 0.1312240 |
| Disneyland_Paris | 15 | 4.006737 | 0.2094921 |

The estimated average rating for California is 4.391302, for HongKong is 4.164182, for Paris is 4.006737.

**Hypothesis Test**

We perform a hypothesis test to determine whether there is evidence that any of the ratings are statistically significantly different from each other in the population.

$H_0 : \mu_{california} = \mu_{hongkong} = \mu_{paris}$

$H_1 : \mu_{california} \neq \mu_{hongkong}$ or $\mu_{california} \neq \mu_{paris}$ or $\mu_{hongkong} \neq \mu_{paris}$ (i.e. the means are not all equal).

We perform an ANOVA.

```
rating_aov = aov(Total_of_Ratings~Branch,data=Employee.B.Branch)
summary(rating_aov)
```

```
##              Df Sum Sq Mean Sq F value   Pr(>F)
## Branch        2  1.121  0.5607   22.74 2.03e-07 ***
## Residuals    42  1.035  0.0247
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

We obtain $p-value < 2.03e-07$ so $p-value < \alpha$. Therefore, we reject the null hypothesis and we conclude that there is evidence that Employee B could achieve more precision for these estimates.

## Subtask 3

For Employee A: Overall estimated average rating: 4.2227 SE: 0.0125 95% confidence interval: [4.198217, 4.247116] California estimated average: 4.396533 HongKong estimated average: 4.213475 Paris estimated average: 3.976963 Result of ANOVA: the means are not all equal

For Employee B: Overall estimated average rating: 4.2182 SE: 0.0308 95% confidence interval: [4.157762, 4.278607] California estimated average: 4.391302 HongKong estimated average: 4.164182 Paris estimated average: 4.006737 Result of ANOVA: the means are not all equal

Let $\bar{y}_A$ be the estimated average for Employee A and $\bar{y}_B$ the estimated average for Employee B.

We observe that $SE(\bar{y}_A) < SE(\bar{y}_B)$ so $Var(\bar{y}_A) < Var(\bar{y}_B)$, therefore the estimate found by Employee A is more efficient than the estimate found by Employee B. Thus, the result found by Employee A provides the best answer. Because Employee B used months as clusters, this means people who went to the park in similar wheather would be in the same cluster. This would lead to homogeneity inside a single cluster, making cluster-sampling perform worse than SRSWOR.