

Hypothesis Testing Problem using R

Eriko Funasato

October 31, 2017

Introduction

In elementary of statistic, we have learnt about hypothesis testing. A statistical hypothesis test is a method of statistical inference. Commonly, two statistical data sets are compared, or a data set obtained by sampling is compared against a synthetic data set from an idealized model([Wikipedia](#)). We compare one alternative statement against a null hypothesis that is supposed to be true. In hypothesis testing, we determined if alternative claim is true by calculating whether the null hypothesis should be rejected or not.

In this process, we could introduce terms such as types of errors, significant level, p-value, and some graphs to make us easier to analyze the dataset. In this post, I am going to apply R into hypothesis testing problem by some example using the dataset "earn-of-collegemajors-all-ages"([Github](#)). -import the dataset and packages to be used

```
dat <- read.csv("C:/Users/eriko/stat133/stat133-hws-fall17/post01/data/earn-of-collegemajors-all-ages.csv", stringsAsFactors=FALSE)
library(dplyr)
```

```
## Warning: package 'dplyr' was built under R version 3.4.2
```

```
##
## 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
```

```
library(ggplot2)
```

```
## Warning: package 'ggplot2' was built under R version 3.4.2
```

```
library(BSDA)
```

```
## Warning: package 'BSDA' was built under R version 3.4.2
```

```
## Loading required package: lattice
```

```
##
## Attaching package: 'BSDA'
```

```
## The following object is masked from 'package:datasets':
##
##   Orange
```

```
setwd("/Users/eriko/stat133/stat133-hws-fall17/post01")
```

#Graphing and Analyze * Before we start doing a hypothesis testing problem, let's organize the data best fit to be worked in R studio, and do some analysis and graphing process to know more about the dataset.

* check the structure of the dataset.

```
str(dat)
```

* Rename the columns The name of columns in the dataset is too long and complicated, it is better to replace them by shorter abbreviation.

* Add column of employment rate

* Picking data of “Computers & Mathematics” and “engineering”

* summary of data of “Computers & Mathematics” and “engineering”

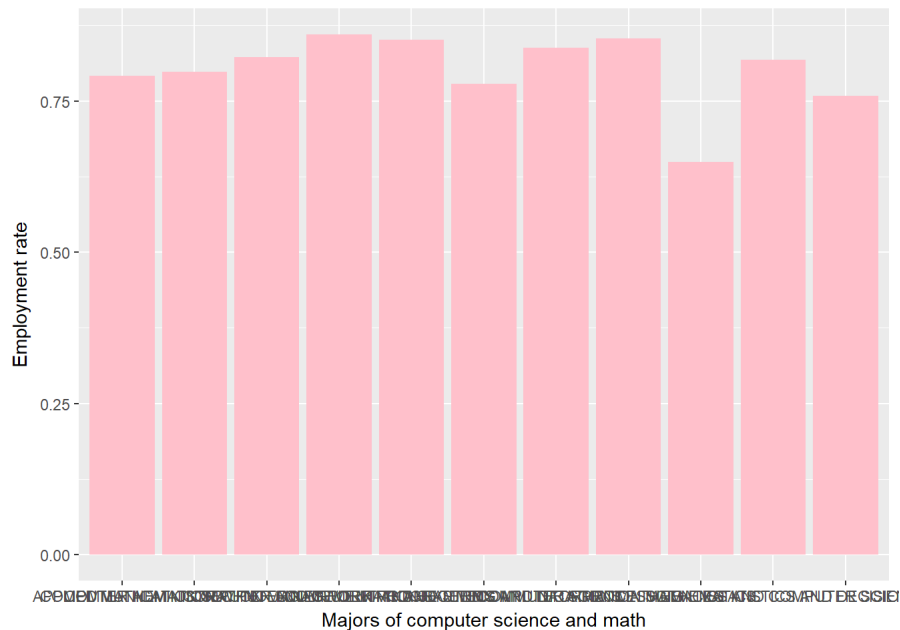
```
sink()
sink(file = "/Users/eriko/stat133/stat133-hws-fall17/post01/output/summary-engin.txt")
summary(engin)
```

```
##      major_code      major      major_cat      total
## Min.   :1401      Length:29      Length:29      Min.    : 6264
## 1st Qu.:2406      Class :character      Class :character      1st Qu.: 18347
## Median :2413      Mode  :character      Mode  :character      Median : 37382
## Mean   :2490                                     Mean   :123311
## 3rd Qu.:2499                                     3rd Qu.:138366
## Max.   :5008                                     Max.   :671647
##      employed      employed_fulltime      unemployed      unemployed_rate
## Min.   : 4120      Min.   : 3350      Min.   : 0      Min.   :0.00000
## 1st Qu.:12876      1st Qu.: 9226      1st Qu.: 617      1st Qu.:0.04384
## Median :27275      Median :22104      Median :1521      Median :0.04985
## Mean   :90413      Mean   :76414      Mean   :5048      Mean   :0.05063
## 3rd Qu.:101273      3rd Qu.:85014      3rd Qu.:5498      3rd Qu.:0.05882
## Max.   :489965      Max.   :422317      Max.   :26064      Max.   :0.08599
##      median      q1      q3      employed_rate
## Min.   :60000      Min.   :40000      Min.   :82000      Min.   :0.5413
## 1st Qu.:67000      1st Qu.:46900      1st Qu.:96000      1st Qu.:0.7001
## Median :75000      Median :50000      Median :102000      Median :0.7295
## Mean   :77759      Mean   :52459      Mean   :108534      Mean   :0.7278
## 3rd Qu.:85000      3rd Qu.:60000      3rd Qu.:116000      3rd Qu.:0.7702
## Max.   :125000      Max.   :75000      Max.   :210000      Max.   :0.8351
```

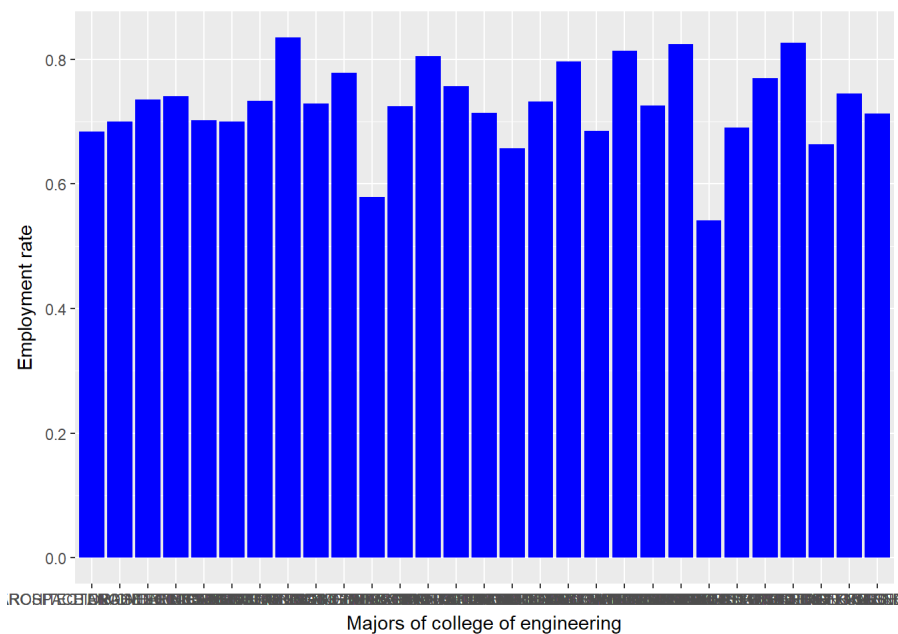
```
sink()
```

* Plot and store employment rate of each dataframe

```
ggplot(cs_math)+geom_col(aes(x=major,y=employed_rate),color=NA,fill="pink")+labs(x="Majors of computer science and math",y="Employment rate")
```



```
ggplot(engin)+geom_col(aes(x=major,y=employed_rate),color=NA,fill="blue")+labs(x="Majors of college of engineering",y="Employment rate")
```



```
#save
pdf(file = "C:/Users/eriko/stat133/stat133-hws-fall17/post01/image/csmath-ggplot.pdf")
ggplot(cs_math)+geom_col(aes(x=major,y=employed_rate),color=NA,fill="pink")+labs(x="Majors of computer science and math",y="Employment rate")
dev.off()
```

```
## png
## 2
```

```
#save
pdf(file = "C:/Users/eriko/stat133/stat133-hws-fall17/post01/image/engin-ggplot.pdf")
ggplot(engin)+geom_col(aes(x=major,y=employed_rate),color=NA,fill="blue")+labs(x="Majors of college of engineering",y="Employment rate")
dev.off()
```

```
## png
## 2
```

Hypothesis * To build up my own hypothesis testing problem, I would like to compare which one of "Computer science & math" and "Engineering" has higher income salary, since I was wondering for a long time if which of the field is paid more.

* But to compare the income salary by hypothesis test, I need the mean and standard deviation of the datasets.

* To get the mean and sd, here is the formula ([Wan et al.\(2014\)](#)),

$$\bar{x} = \frac{q1+m+q3}{3}$$

$$s = \frac{q3 - q1}{1.35}$$

where q1 is first quartile, m is the median, q3 is the third quartile + Add the mean and standard deviation to each datasets

```
cs_math <- mutate(cs_math,mean=(q3+median+q1)/3,std=(q3-q1)/1.35)
engin <- mutate(engin,mean=(q3+median+q1)/3,std=(q3-q1)/1.35)
```

- Calculate the mean and the std of whole datasets

```
Mean_cs_math <- sum(cs_math$mean * cs_math$total) / sum(cs_math$total)
Mean_engin <- sum(engin$mean * engin$total) / sum(engin$total)
std_cs_math <- mean(cs_math$std)
std_engin <- mean(engin$std)
```

Hypothesis: Computer science & math earn more money than the engineering major.

1. Claim

H_0 (null): mean(cs_math income)-mean(engineering income) = 0

H_1 (alternative claim): mean(cs_math income)-mean(engineering income) > 0

2. since the number of both independent sample are huge enough, suppose they are approximately normally distributed. Standard deviation are known by calculation, we could use 2-sample Z-test.
3. Set the significant value as $\alpha=0.05$, which means confidence level is 0.95
4. Program calculation

- Calculate using `z.test`([RDocumentation](#))

```
z.test(cs_math$mean,engin$mean, alternative = "greater", mu = 0, sigma.x = std_cs_math,sigma.y = std_engin, conf.level = 0.95)
```

```
##
## Two-sample z-Test
##
## data: cs_math$mean and engin$mean
## z = -0.79045, p-value = 0.7854
## alternative hypothesis: true difference in means is greater than 0
## 95 percent confidence interval:
## -34129.93 NA
## sample estimates:
## mean of x mean of y
## 68506.06 79583.91
```

- Calculate using `t.test`([Phil Spector](#))

```
t.test(cs_math$mean,engin$mean,alternative = "greater",conf.level = 0.95)
```

```
##
## Welch Two Sample t-test
##
## data: cs_math$mean and engin$mean
## t = -2.5243, df = 23.49, p-value = 0.9906
## alternative hypothesis: true difference in means is greater than 0
## 95 percent confidence interval:
## -18592.52 Inf
## sample estimates:
## mean of x mean of y
## 68506.06 79583.91
```

5. Conclusion: fail to reject null

- For the result of calculation using z-test, p-value=0.7854 is greater than the $\alpha=0.05$, therefore fail to reject the null hypothesis, we don't have enough confidence to say that people who major in "Computer science & math" earn more than "Engineering" major.
- For the result of calculation using t-test, p-value=0.9906 is greater than the $\alpha=0.05$, therefore fail to reject the null hypothesis, we don't have enough confidence to say that people who major in "Computer science & math" earn more than "Engineering" major. Therefore, we cannot conclude which of the major("CS & Math" and "Engineering") can earn more money after graduate.
- The difference between the Z-test and t-test is, when the mean of the sample can be known, the standard deviation is unknown, then we use the **t-test** to estimate the population. If the standard deviation is known or calculated, then use **z-test** to do the evaluation.

More about hypothesis test.

- Not only the problems like the example I give, there are still many kind of hypothesis testing, such as:
 - one-sample z-test(left-sided,right-sided,both-sided)
 - one-sample t-test(left-sided,right-sided,both-sided)
 - two-proportion z test(inference of two proportion)
 - two-proportion t-test
 - two-sample t/z-test(inference of two mean)
 - chi-square test(testing a claim about a standard deviation of population)
 - Goodness-of-fit(test if sample data with k categories is "good fit" to an assumed distribution)
 - One-way Analysis of Variance(ANOVA):test for equality of more than three sample data.
 - And to do these in R, check ([Rtutorial](#))
- For the concepts of hypothesis test, see ([Pennstate](#)).

Take Home Message!!

Please try to make some function of t-test or z-test, for one-sample and two-sample.

And try to compare the employment rate and unemployment rate of the majors or the major categories you are interested in.