```
Jiancong Chen -Jerry.ai
    • Question 1
          • Step 0: Construct synthetic tables
          • Step 1: What information would you derive from it, and how will you derive it?
                • 1a Whether a referral lead to success acquisition of new customer, which is defined as a customer who makes purchase
                  through Jerry.

    1b The percentage of success acquisition

                1c Risk perspective.

    1d Distribution of total between users who were referred and who were not referred

    1e Dsitribution of final bill between users who were refered and who were not refered

    1f Are there any differences for the final bill among groups

                • 1g Month specific information? Do we observe any unique signal based on months?
          • Step 2: How would you make a recommendation
                Answer:

    2b whether the rewards program help to boost the total amount

                • 2c Should we be considerate about certain users who may risky and only take the promotion?

    Question 2

    Code block

This is Jiancong (Nigel) Chen's home assignment submission to Jerry. Al's interview. Question #1 was coded in R due to the richness of statistical
packages as well as data analysis scripting packages. Question #2 was coded in Python and typed as text in the Rconsole.
Question 1
As no original/synthetic data was given, I constructed a synthetic dataset, that will enable me to show some data wrangling results, data
visulization and statistical tools to make recommendations on whether the referral program should continue or stop rather than speaking on top of
air about "what I want to do" or "how I can do".
Step 0: Construct synthetic tables
Note that this step is not needed when real data becomes available. This is only to make my further analysis more interpretative and
representable. I generated a synthetic pool of ids (expressed as emails), name, and referring user id (also expressed as emails). The constructed
table include situations where one user refer several new users.
For the Table Purchase, it includes ids that were referred by other users and ids who registered without referral. As the monthly premium for auto
insurance may range between 0 (no need situation) and 120 (based on my own experience), I set the maximum total monthly cost for a user at
120. As I was not given any other coupon situations, so I assumed the discount is the referral (10/referral) + a random small number.
Note: there are 592 entries of User table.
 rm(list = ls())
 library(randNames)
 set.seed(111)
```

## Generate the User table. The pool dataframe has N\_user unique user names Pool <- rand\_names(n = N\_user)[, c('email', 'name.first')]</pre> ## Warning: `as\_data\_frame()` is deprecated as of tibble 2.0.0. ## Please use `as\_tibble()` instead. ## The signature and semantics have changed, see `?as\_tibble`.

## This warning is displayed once every 8 hours. ## Call `lifecycle::last\_warnings()` to see where this warning was generated.

names(Pool) <- c('id', 'name')</pre> ## Assuming a rand% of the users were actually referred by other users. ## Note rand% is generated randomly. rand <- runif(1, min = 0, max = 1) # 0.59 for this case idx usr <- sample(N user, size = N user \* rand, replace = TRUE) idx\_refer <- sample(N\_user, size = N\_user \* rand, replace = FALSE)</pre> id <- Pool\$id[idx\_usr]</pre> refer <- Pool\$id[idx\_refer]</pre> User <- data.frame(id = id, name = Pool\$name[idx\_usr], referring\_user\_id = refer)</pre>

The following code generates the data for the Table: purchase. It allows one user to have multiple purchase records. The total amount of purchase was randomly generated between 0 - 100. The discount was only simulated for the ones that successfully refer other users as found in the User table. For the date variable, currently I am only simulating 2020 January - 2020 March for simplicity N\_purchase= 1000 total <- runif(N\_purchase, 0, 10) \* 10</pre> daterange <- c(as.numeric(as.Date('2020-01-01')):as.numeric(as.Date('2020-03-31')))</pre> date <- as.Date(sample(daterange, N\_purchase, replace = TRUE), origin = '1970-01-01') # Sample from the user pool. Also enable any users to have multple purchase records Purchase\_id <- Pool[sample(1:N\_user, N\_purchase, replace = TRUE),]</pre> Purchase <- data.frame(Purchase\_id, total, date)</pre> # For discount, only consider the ones with referral and give them \$10 Purchase\$discount =  $0 + runif(N_purchase, 0, 5)$ Purchase discount[which(Purchase id %in% User id)] = 10 + Purchase discount[which(Purchase id %in% User id)]

Let us take a quick glimpse of the generated data

3 agnetha.donkersloot@example.com

4 toni.dupuis@example.com

5 harry.zhang@example.com

## [1] 0.6165541

library(dplyr)

library(ggplot2)

mutate(count = n()) %>%

'%!in%' <- **function**(x,y)!('%in%'(x,y))

mutate(refered = case\_when(id %in% User\$id ~ 'Yes',

# Let us first take a look at the merged data frame.

Purchase2 <- Purchase %>%

3 agnetha.donkersloot@example.com

# Now let us visualize the distribution

No

# Now let us visualize for the discount part

ggplot(data = Purchase2) + geom\_boxplot(aes(x = refered, y = total))

4 toni.dupuis@example.com

5 harry.zhang@example.com

6 rows

100 -

0 -

group.

15 **-**

75 **-**

25 -

total - discount

6 cohen.hughes@example.com

6 cohen.hughes@example.com

head(User) id name referring\_user\_id <fct> <fct> <fct> 1 mhrd.glshn@example.com nznyn.sdr@example.com مهراد 2 bojoura.kentie@example.com guillaume.berger@example.com Bojoura 3 ella.lehto@example.com Ella debra.dunne@example.com 4 baer.kauffman@example.com julia.murto@example.com Baer 5 martin.denis@example.com Martin tasmira.rocha@example.com 6 eric.tosse@example.com Eric mhrd.sdr@example.com

6 rows head(Purchase) name date discount id total <chr> <dpl> <date> <chr> <dbl> 1 lima.desouza@example.com Lima 97.54454 2020-02-05 13.194779 2 afsar.tokatlioglu@example.com 37.36030 2020-02-02 12.459178 Afşar

Agnetha

Toni

Harry

Cohen

2020-02-09

2020-01-12

2020-02-23

2020-01-09

1.729834

14.571302

13.501691

1.347864

25.11911

41.22814

48.47339

45.72990

Step 1: What information would you derive from it, and how will you derive it? 1a Whether a referral lead to success acquisition of new customer, which is defined as a customer who makes purchase through Jerry. # Filter the list of customers who were referred Success\_refer\_usr <- unique(Purchase\$id[which(Purchase\$id %in% User\$referring\_user\_id)]) # Glimpse the list head(Success\_refer\_usr) ## [1] "lima.desouza@example.com" "afsar.tokatlioglu@example.com" ## [3] "harry.zhang@example.com" "cohen.hughes@example.com" ## [5] "kadir.agaoglu@example.com" "eren.karaer@example.com" 1b The percentage of success acquisition

Ideally, we would want this metrix to be as high as possible. Based on the current synthetic dataset, the success rate is around 65%. I

print(length(Success\_refer\_usr)/length(User\$referring\_user\_id))

## 1c Risk perspective. The percenrage of existing users who gambled the promotion, defined as a user who refered to multiple referrals but the new referrals are 'robot' who never make purchases. Based on the synthetic data, most users only have one referal, only 4 users have multiple referrals. If these referals are not success users, these 4 users should be flagged and less promotions should be aviable for them for risk purposes and customer relations. # Add a column to count the number of total refers for any specific user

## 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

1d Distribution of total between users who were refered and who were not refered

the group of referred and non-referred, it indicates the referral program is contributing. Otherwise, we should dig dipper into the cause.

First, let us wrangle the data to create a new column for the purchase table that reflect whether this user is referred or not.

id %!**in**% User\$id ~ 'No'))

This can help us identify whether the referral system is actually helping us gain revenus. If there's significant differenes in the total cost between

print(paste('Potential risky user:', User\_Agg\$id[which(User\_Agg\$count > 1)])) ## [1] "Potential risky user: ryan.levesque@example.com" ## [2] "Potential risky user: baer.kauffman@example.com"

mutate(count = replace(count, which(referring\_user\_id == 'No refer'), 0))

User\_Agg <- User %>% group\_by(referring\_user\_id) %>%

head(Purchase2) id name total date discount refered <chr> <chr> <dbl> <date> <dbl> <chr> 1 lima.desouza@example.com Lima 97.54454 2020-02-05 13.194779 Yes 2 afsar.tokatlioglu@example.com Afşar 37.36030 2020-02-02 12.459178 Yes

Agnetha

Toni

Harry

Cohen

25.11911

41.22814

48.47339

45.72990

2020-02-09

2020-01-12

2020-02-23

2020-01-09

1.729834 No

14.571302 Yes

13.501691 Yes

1.347864 No

75 total 20 -25 **-**

refered

 $ggplot(data = Purchase2) + geom_boxplot(aes(x = refered, y = discount))$ 

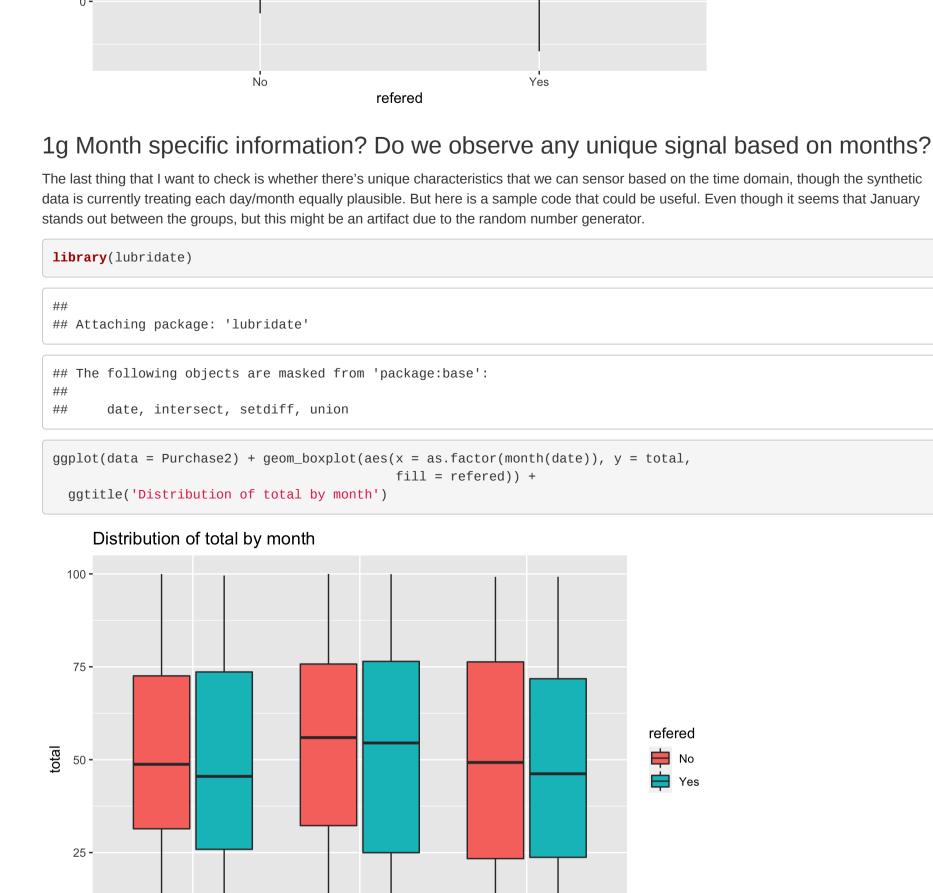
Note: as the total is generated randomly following uniform distribution, there should not be any significant differences among the groups.

1e Dsitribution of final bill between users who were refered and who were not refered

Note: given the promotion and constructed dataset, it is expected that the promotion for people who use refer will be higher than the non-refer

Yes





as.factor(month(date))

as.factor(month(date))

 $ggplot(data = Purchase2) + geom_boxplot(aes(x = as.factor(month(date)), y = total - discount,$ 

fill = refered)) +

ggtitle('Distribution of discount by month')

ggtitle('Distribution of final bill by month')

Distribution of final bill by month

Distribution of discount by month

15 **-**

10 -

100

75 **-**

0.620

0.615

0.610

0.605

0.600

210

205

200

195

190

185

30

count

10 -

0 -

test1

##

##

m\_tukey

## \$refered

##

##

##

1.0

Refered user

1.0

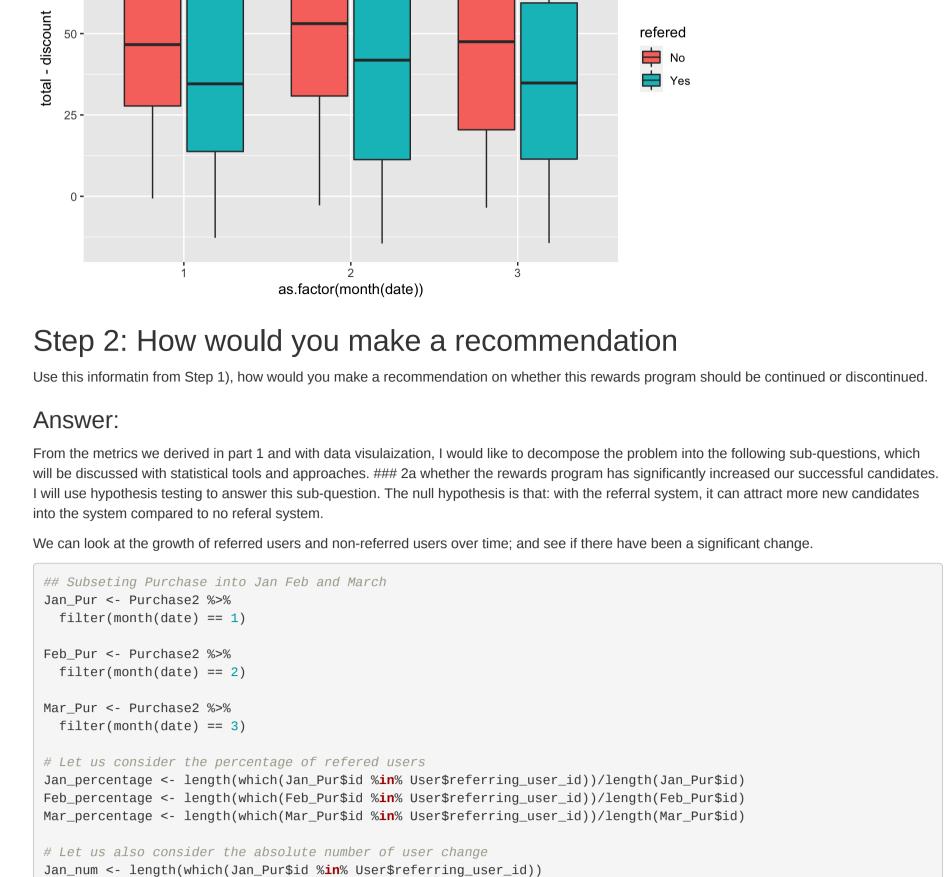
Percentage of refered user

 $ggplot(data = Purchase2) + geom_boxplot(aes(x = as.factor(month(date))), y = discount,$ 

fill = refered)) +

discount ■ No Yes 5 -

refered



Feb\_num <- length(which(Feb\_Pur\$id %in% User\$referring\_user\_id))</pre> Mar\_num <- length(which(Mar\_Pur\$id %in% User\$referring\_user\_id))</pre>

1.5

1.5

# First let us make histogram for the two groups

25

# Now let us set up the hypothesis testing (t-test) test1 <- t.test(total ~ refered, data = Purchase2)</pre>

## t = 1.0452, df = 935.35, p-value = 0.2962

mean in group No mean in group Yes

model <- lm(total ~ refered, data = Purchase2)</pre>

Tukey multiple comparisons of means

95% family-wise confidence level

## Yes-No -1.882567 -5.401025 1.635892 0.2939904

lwr

p adj

Welch Two Sample t-test

## 95 percent confidence interval:

51.04327

# Let us also try tukey's test

 $m_{tukey} < - TukeyHSD(x = m_anova)$ 

## Fit: aov(formula = model)

diff

def add(self, fr, to): newItvs = []start = frend = to

> for interval in self.intervals: if interval[0] < fr:</pre>

> > else:

else:

else:

def remove(self, fr, to):

else:

else:

self.intervals = newItvs

else:

newItvs = []

if interval[1] < fr:</pre>

newItvs.append(interval)

elif interval[0] >= fr and interval[0] <= to:</pre>

start = interval[0]

end = interval[1]

end = interval[1]

newItvs.append(interval)

if interval[1] <= to:</pre>

continue

newItvs.append([start, end]) self.intervals = newItvs return self.intervals

for interval in self.intervals: if interval[0] < fr:</pre>

if interval[1] < fr:</pre>

if interval[1] <= to:</pre> continue

newItvs.append(interval)

newItvs.append(interval)

elif interval[1] >= fr and interval[1] <= to:</pre> newItvs.append([interval[0], fr])

newItvs.append([interval[0], fr]) newItvs.append([to, interval[1]]) elif interval[0] >= fr and interval[0] <= to:</pre>

newItvs.append([to, interval[1]])

elif interval[1] >= fr and interval[1] <= to:</pre>

## data: total by refered

## -1.652090 5.417223 ## sample estimates:

m\_anova <- aov(model)</pre>

plot(1:3, c(Jan\_percentage, Feb\_percentage, Mar\_percentage), xlab = '2020-month', type = 'b', ylab = 'Percentage of refered user', main = 'Time change of refered user percentage')

Time change of refered user percentage

2.0

2020-month

Time change of refered user

2.0

2020-month

attract as many new users as it should be. This result is based on the synthetic data and is only used for illustration.

2b whether the rewards program help to boost the total amount

plot(1:3, c(Jan\_num, Feb\_num, Mar\_num), xlab = '2020-month', type = 'b', ylab = 'Refered user', main = 'Time change of refered user ')

2.5

2.5

data, the percentage decrease the most in Feburary. The absolute new refered users show a decling trend. These results indicate that we need to identify what are the main reasons that lead to the decrease in Feburary; also it seems the referal program is not working promptly, as it does not

The second factor to consider here is to determine whether Jerry can earn more revenus even if the referal program is giving out 10 dollar promotion. In this section, we will compare the revenue based on the two groups. This comparison can be done with hypothesis testing

3.0

3.0

refered

Yes

Based on the synthetic generated

## ggplot(data = Purchase2) + geom\_histogram(aes(total, group = refered, color = refered)) ## `stat\_bin()` using `bins = 30`. Pick better value with `binwidth`. 40 -

50

total

## alternative hypothesis: true difference in means is not equal to 0

49.16070

<del>7</del>5

100

Based on the synthetic data, we did not observe a significant difference of total cost between referred groups and non-referred group. The t-test returned a p-value of 0.4383, indicating the null hypothesis should not be rejected. There's no significant difference between the referred group versus the non-referred group. The ANOVA and Tukey's HSD test also indicate a p-value at 0.439. There is no significant differences between the two group. From this information, it seems by prompting the 10 dollar bonus/coupon will not really boost our revenue. There might be people who only take them as a way of lowering their expense. 2c Should we be considerate about certain users who may risky and only take the promotion? This category can be very interesting as there might be a group of people who post their referal link online and get the promo code. But in reality, they never really refer to people. In this severe case, the promotion is mostly wasted and will not result in any boost in revenue or increase in new users. We can take a look at the amount of referred people who may not be real user. length(which(User\$referring\_user\_id %in% Purchase\$id))/length(User\$referring\_user\_id) ## [1] 0.6182432 This simple algebra indicates only ~ 64.9% of the referred user are actually qualified users who will make a purchase through Jerry at a later time. If we believe 64.9% is not sufficient, then I will recommend canceling the program as it does not yield a satisfactory user attraction rate. Question 2 Code block The following code implements the add and remove function in Python that manages disjointed intervals of integers. I have tested all the exam sequence with the satisfactory result as shown. class AddnRemove(): def \_\_init\_\_(self): self.intervals = []

return self.intervals Now, we can execuate the codes. The results indicate the add/remove functions are working promptly. obj = AddnRemove() print(obj.add(1, 5))

## [[1, 5]] print(obj.remove(2, 3)) ## [[1, 2], [3, 5]] print(obj.add(6, 8)) ## [[1, 2], [3, 5], [6, 8]] print(obj.remove(4, 7)) ## [[1, 2], [3, 4], [7, 8]] print(obj.add(2, 7)) ## [[1, 8]] print(obj.add(10, 11)) ## [[1, 8], [10, 11]] print(obj.remove(9, 10)) ## [[1, 8], [10, 11]]