

Transcription

Case Study: Analysis of Credit Defaulters

CASE STUDY INTRODUCTION

upGrad



Golden Corp.

- Golden Corp. is a financial corporation which offers a wide range of credit cards to its customers
- Ms. Gocha, Head, Decision Analytics Division of Raipur is worried about the increasing default rate in her region and wants to get insights on default patterns before modifying existing policies for issuing credit cards

In this video, we will discuss application of interval estimation and hypothesis testing through a practical example. We will look at the case of Golden Corp, which is a financial institution offering wide range of credit cards.

Within Golden Corp, we have Ms. Gocha who heads a decision analytics team in Raipur. She is worried about the default rate in her region and thinks that the existing policies of the company need to be modified. Before doing that, she wants to analyse the historical data and get some insight on default and spending patterns of credit card holders.

CASE STUDY : OBJECTIVE

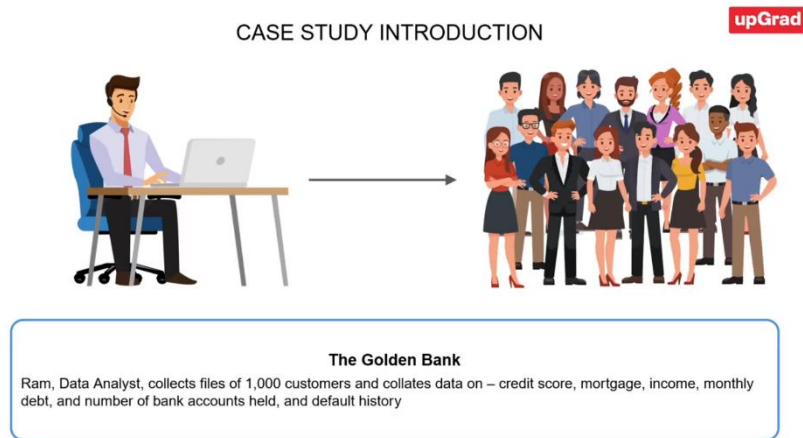
upGrad

- 1 What is the average credit score of defaulters?
- 2 What should be the cut-off score for second level of scrutiny?
- 3 Does the population default rate of those who have more than 10 accounts exceed 20%?
- 4 Which group is a better target for Golden Corp: customers with home mortgages or those who rent?

She goes on to discuss the project with her colleague Ram, and they decide on some of the more important objectives. The objectives are:

- One, what is the average credit score of defaulters?
- Two, what should be the cut-off score for a second scrutiny?
- Three, if the population default rate of those who have more than 10 bank accounts exceed 20%.

- And four, which group is a better target for Golden Corp, customers with home mortgages or customers who rent their houses?



We look at each of these objectives in greater detail. After discussing the objectives, Ram collects random files of 1000 credit card holders from the data team. He collates information on credit score, mortgages, income, monthly spending, number of bank accounts held and customers default history.

CustomerID	HomeOwnership	CreditScore	Income	MonthlyDebt	Number of Open Accounts	Defaults
1	Rent	748	83,311	8,257	10	-
2	Rent	695	86,335	21,118	5	-
3	Home Mortgage	51,810	12,952	12	1	-
4	Home Mortgage	70,039	6,101	12	-	-
5	Home Mortgage	719	5,26,568	39,491	7	-
6	Home Mortgage	688	1,02,629	13,726	6	-
7	Rent	731	68,527	8,950	4	-
8	Home Mortgage	689	1,11,953	23,371	7	-
9	Home Mortgage	738	1,26,413	23,360	7	-
10	Home Mortgage	738	1,26,413	23,360	7	-
11	Home Mortgage	733	1,95,548	15,214	25	1
12	Home Mortgage	731	2,01,481	31,695	24	-
13	Rent	724	1,13,438	28,574	16	1
14	Home Mortgage	733	1,95,548	15,214	25	1
15	Home Mortgage	731	2,01,481	31,695	24	-
16	Rent	724	1,13,438	28,574	16	1
17	Home Mortgage	694	66,454	16,946	10	-
18	Home Mortgage	700	1,01,501	29,841	3	-
19	Rent	746	1,12,330	9,952	5	-
20	Rent	746	1,12,330	9,952	5	-
21	Own Home	644	59,499	19,516	7	-
22	Rent	746	1,12,330	9,952	5	-
23	Home Mortgage	643	2,74,409	72,444	38	-
24	Rent	746	1,12,330	9,952	5	-
25	Rent	746	1,12,330	9,952	5	-
26	Rent	717	55,744	5,155	6	-
27	Rent	716	1,09,077	20,288	9	-
28	Rent	697	46,952	1,214	4	-
29	Home Mortgage	679	1,39,384	29,970	9	1
30	Rent	741	1,10,401	22,963	9	-
31	Rent	746	1,12,330	9,952	5	-
32	Rent	746	1,12,330	9,952	5	-
33	Rent	746	1,12,330	9,952	5	-
34	Home Mortgage	729	1,49,830	15,860	10	-
35	Home Mortgage	717	1,87,600	45,496	16	-
36	Rent	739	70,609	13,345	5	-

	Credit Score	Income	Monthly Debt
Mean	735	1,06,614	19,130
Min	594	1,009	533
Max	751	6,34,657	97,150

Let us take a look at the data. We have customer ID. We have home ownership, whether the customer owns a house, rents it or has a home loan. Then, we have credit score of customers, his or her monthly income, monthly debt or average credit card bill of the customer.

Then, there is number of open accounts, which is the total number of bank accounts held by the customer. Finally, we have default, which is one if the customer has defaulted on three or more consecutive bill payments and zero otherwise.



What is the average credit score of defaulters?

Coming back to the objectives, here is the first one, what is the average credit score of defaulters? How do we go about this? We have a sample data set, and based on the sample, we can estimate the 95% confidence interval for population mean of defaulters' credit score.

```
1 data<- read.csv("CreditCards.csv")
2
3 defaulters <- data[data$Default == 1, ]
4 cs <- defaulters$CreditScore
5 cs<-na.omit(cs)
```

Environment: Global Environment

Data:

- data: 1000 obs. of 7 variables
- defaulters: 142 obs. of 7 variables

Values:

cs: int [1:142] 714 713 678 729 714 716 725 716 787 65...

Console:

```
> defaulters <- data[data$Default == 1, ]
> head(defaulters)
  CustomerID HomeOwnership CreditScore Income MonthlyDebt Number.of.Open.Accounts
3      7619 Home Mortgage      NA      51809.06      12052.49              12
12     3318      Rent      714 123697.92      28574.10              16
13     0564 Home Mortgage      713 195548.00      15213.68              25
21     2825      Rent      NA      5891.52      2546.76              7
28     9829 Home Mortgage      678 139393.50      29969.65              9
37     5199 Home Mortgage      NA      61812.80      15253.20              16
Default:
3      1
12     1
13     1
21     1
28     1
37     1
> cs <- defaulters$CreditScore
> cs<-na.omit(cs)
>
```

Let's look at that in R. So, we start with reading the data, we again use read.csv function here. The data is called credit card dot csv. Once we run it, we can do head data. We can also see the data here in this window. It has 1000 observation, which is on sample size.

Now, we want to estimate the confidence interval for credit score of defaulters. So, first we would need to subset the data. We only need the data of defaulters and defaulters would be all the data points or all the customers for whom default is equal to 1.

So, net subset R data, we have formed defaulters, which is data such that data dollar default equals to 1, and we can again do head defaulters. So, we have all the customers for whom default is 1.

Next, we want to estimate confidence interval for credit score, and you can see there are NAs here. So, before estimating the confidence interval, we would need to get rid of these NAs.

Let us default cs as the credit score of defaulters or the series of these credit scores, and remember how we get rid of the NAs. We will use na.omit function. So, we will write na.omit cs. So, once we run it, we have a sample of defaulters of credit score.

```

1 data<- read.csv("CreditCards.csv")
2
3 defaulters <- data[dataDefault == 1, ]
4 cs <- defaulters$CreditScore
5 cs<-na.omit(cs)
6
7 n<- length(cs)
8 m<-mean(cs)
9 s<-sd(cs)
10 se<-s/sqrt(n)
11
12 t<-qt(1- 0.05/2, n-1)
13

```

Environment: Global Environment
 Data: data (1000 obs. of 7 variables), defaulters (142 obs. of 7 variables)
 Values: cs (int [1:116] 714 713 678 729 714 716 725 716 707 65...), m (699.172413793183), n (116), s (27.8803852413862), se (2.58119354161224), t (1.9888754118391)

Now, let us estimate the 95% confidence interval for mean credit score amongst defaulters. So, for sample size would be length of cs and the mean can be calculated using mean cs, which is defined as n here. We have s or standard deviation equal to sdcs, and then we compute the standard error, which is s divided by square root of n.

Now, as population standard deviation is not known. We will use T distribution and for 95% confidence interval, alpha would be 0.05, and alpha by 2 would be 0.025. So, the T value at 95% confidence interval can be calculated using QT function. As we want to estimate two-side confidence interval, we will use alpha by 2 and the degree of freedom for this T distribution would be n minus 1, which is the sample size minus one.

```

14 library("BSDA")
15 tsum.test(mean.x = m, s.x = s, g.x = n)
16

```

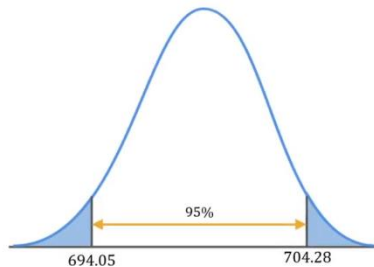
Environment: Global Environment
 Data: data (1000 obs. of 7 variables), defaulters (142 obs. of 7 variables)
 Values: ci (num [1:2] 694.704), cs (int [1:116] 714 713 678 729 714 716 725 716 707 65...), m (699.172413793183), n (116), s (27.8803852413862), se (2.58119354161224), t (1.9888754118391)

Now, we have T, we have standard error, and we have mean, so we can compute the confidence interval by adding and subtracting se into t from the mean. So, here is the confidence interval, which is 694.05 to 704.28.

Remember, we can also use t sum dot s function here. For that, we first need to call the library BSDA, and then we can call ts dot test function, mean dot x would be n, s dot x or standard deviation of sample would be s, and n dot x as the sample size.

If we run it, we have the confidence interval 694.05 to 704.28. So, based on this sample, we are 95% confident that the mean credit score of defaulters lie between 694 and 704.3.

INTERVAL ESTIMATION : CREDIT SCORE

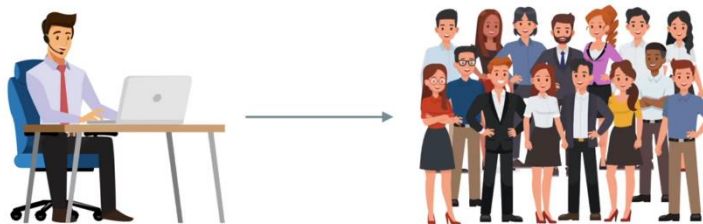


Credit Score

There is a 95% chance that mean credit score of defaulters lies between 694.05 to 704.28

We also have the plot here, there's a 95% chance that mean credit score of defaulters lie between 694 and 704.3.

CREDIT SCORE



Credit Score Scrutiny

Ram checks with the Due diligence team and learns that at present they conduct a 2nd round of scrutiny if an applicant's credit score is less than 700. The files with credit score greater than 700 are passed to the approval team.

Ram and Gocha also discuss about proposing a second scrutiny or inspection for those who have a very low credit score. Further examination of such applicants can help code and called identify potential defaulters, and it can then refuse such applications.

How do we define such a score or what should be the cut-off credit score for a second scrutiny? By cut-off score, we want to estimate the upper bound of credit score amongst defaulters, such that anyone below that score can be scrutinized again. In interval estimating term, we want to find one-sided confidence interval or the upper bound of credit score amongst defaulters.

```

10 se<-s/sqrt(n)
11
12 t<-qt(1-0.05/2, n-1)
13 ci<-m+c(se*t, se*t)
14
15 library("BSDA")
16 tsum.test(mean.x = m, s.x = s, n.x = n)
17
18 t<-qt(1-0.05, df = n-1)
19 m+t*se
20
21 tsum.test(mean.x = m, s.x = s, n.x = n, alternative = "less")
22
23:1 (Top Level) :
> m+t*se
[1] 703.4526
> tsum.test(mean.x = m, s.x = s, n.x = n, alternative = "less")
One-sample t-Test
data: Summarized x
t = 270.87, df = 115, p-value = 1
alternative hypothesis: true mean is less than 0
95 percent confidence interval:
NA 703.4526
sample estimates:
mean of x
699.1724
Warning message:
In tsum.test(mean.x = m, s.x = s, n.x = n, alternative = "less") :
argument 'var.equal' ignored for one-sample test.
>
  
```

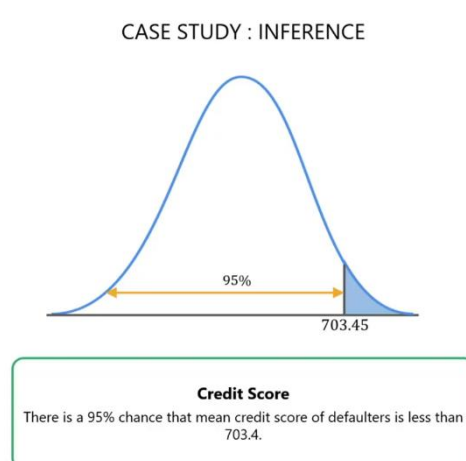
Environment pane details:

Global Environment	data	defaulters
data	1000 obs. of 7 variables	142 obs. of 7 variables
Values	ci: num [1:2] 694 704	ci: int [1:116] 714 713 678 729 714 716 725 716 707 65...
	m: 699.172413793103	n: 1161
	s: 27.8003052413862	se: 2.58119554161224
	t: 1.65822183083114	

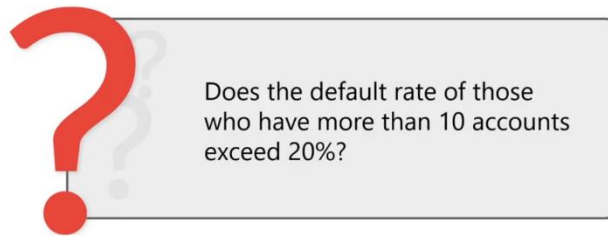
Earlier, we build a two-sided confidence interval. Now, we just want to get the critical value at the right side of the distribution. So, going back to R, first, we need the right T value at 95% confidence interval. So, we use QT function, and we pass 1 minus 0.05. Now, the 5% region is only at the right side of the distribution, and therefore, we use alpha 0.05 and not alpha by 2.

Next, we can calculate the upper bound as $m + t \times se$. So, this is the upward bound, 703.45. With 95% confidence, population mean of credit defaulters is less than 703.4. We can also estimate confidence interval using `t sum dot test` function.

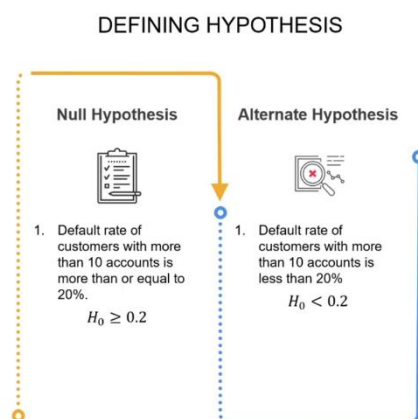
As we want to estimate the upper bound in the `t sum dot test` function, we will pass alternative as less. Once we run it, we have one-sided confidence interval, which is the upward bound of 703.4.



We can also see it graphically. Based on the given sample, we found that on average, there is a 95% chance that the credit score of defaulters is less than 703.4. Therefore, anyone with a score less than 703.4 should be inspected further.



Moving on, we also want to check if the population default rate of those who have more than 10 bank accounts exceed 20%. We can perform a hypothesis test for this.



So, a null would be that default rate of customers with more than 10 accounts is more than 20% or π is greater than an equal to 0.2. And our alternative hypothesis would be default rate of customers with more than 10 accounts is less than 20% or π is less than 0.2. This is a population proportion test. Let's do it in R.

The screenshot displays the RStudio interface with the following components:

- Source Editor:** Contains R code for a t-test:


```
21 tsm.test(mean.x = n, s.x = s, n.x = n, alternative = "less")
22
23 df <- data[nData$Number.of.Open.Accounts > 10,]
24 #H0: pi=0.2
25 #H1: pi= 0.2
26
27 pi <- 0.2
28 p <- nrow(df[df$Default == 1,])/nrow(df)
29 nc = nrow(df)
30 sigma <- sqrt(pi*(1-pi)/n)
31 Z.test <- (p-pi)/n
32 Z <- qnorm(0.95)
33
```
- Console:** Shows the execution of the code:


```
33:1 Top Level 2 R Script >
> Default
3 1
4 0
12 1
13 1
14 0
15 0
> p <- nrow(df[df$Default == 1,])/nrow(df)
> pi <- 0.2
> p
[1] 0.2412869
> sigma <- sqrt(pi*(1-pi)/n)
> Z.test <- (p-pi)/n
> Z <- qnorm(0.95)
[1] 1.644854
> Z.test
[1] 0.0003559212
>
```
- Environment:** Displays a data frame with the following values:

n	699.172413793103
m	1166
p	0.241286863270777
pi	0.2
s	27.8003852413862
se	2.58119354161224
sigma	0.0371390670355404
t	1.6582183003114
Z	1.64485362095147
Z.test	0.000355921235092909

For population proportion hypothesis test, we would want to identify the sample proportion of default amongst those who have more than 10 accounts. For that, we would subset the data and we'll define DF as data such that number of open accounts is greater than 10.

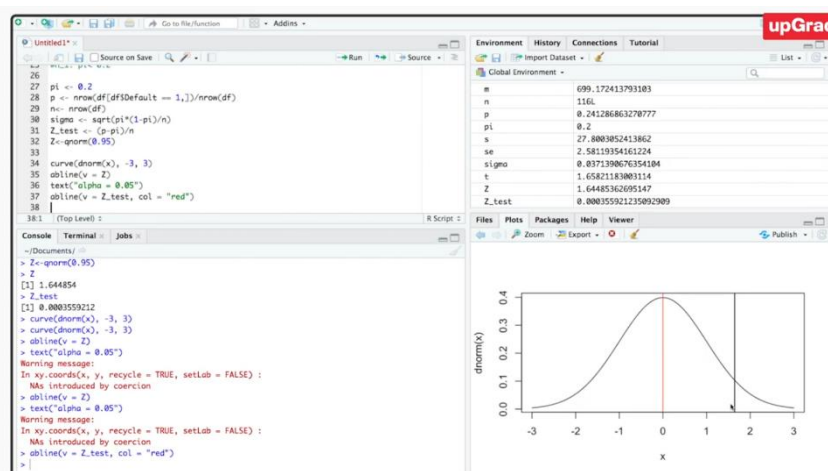
So, in DF, we have all the customers who hold more than 10 bank accounts in our sample. Based on the sample, we want to test the null or default rate of customers who have more than 10 accounts is 20% or 0.2 and the null hypothesis would be that π is less than 0.2. π is 0.2, so the π or population proportion in the hypothesis is 0.2.

Next, we need to estimate the sample proportion or the default rate in DF. So, we can calculate sample proportion or P using this command. In the denominator, we have all the customers who are more than 10 accounts, which is the number of rows in DF. And in the numerator, we have those customers in DF who have defaulted, that is where default is equal to 1.

So, P or sample proportion of default rate for customers who hold more than 10 bank accounts is 0.24. Now to test this hypothesis, we would want to calculate sigma. We can define sample size as nrow DF, and this is how we would calculate sigma in proportion test. So, it is square root of $\pi \times 1 - \pi$ divided by N.

The test statistic that can be used is Z test. So, we'll define the test statistic as P minus π divided by N, and we can calculate the critical value using qnorm function. So, we will define C as to qnorm 0.95.

This is a one-tailed test, as our alternative hypothesis is π less than 0.2, and therefore, we use 5% as rejection region. The critical value is 1.64, while the test statistic is 0.0003, and since the Z test statistic is less than critical value, we do not reject the null that π is greater than equal to 0.2.



We can also plot the graph here, just remember how we will plot the graph using curve function. We will first plot the normal distribution and we will define the critical region or the rejection region.

So, this is the rejection region, and our test statistic is almost 0, it lies here. The test statistic does not line the rejection region. Hence, we don't reject the null.

Based on the **hypothesis test**, we **fail to reject** the null that for customers with more than 10 accounts, the **default rate exceeds 20%**.

So, we can conclude that based on the given sample, we fail to reject the null that customers who have more than 10 bank accounts, the default rate exceeds 20%. Golden Corp should thus be more cautious in approving credit cards to such applicants.

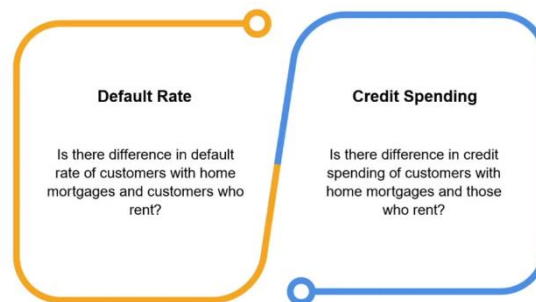
Which group is a better target for Golden Corp: Customers with home mortgages or those who rent?

Finally, the last objective is to identify which group should Golden Corp target in its new marketing strategy. They want to identify which group can generate more revenue for Golden Corp. Is it the customers with home mortgage or is it the customers who are living on rent?

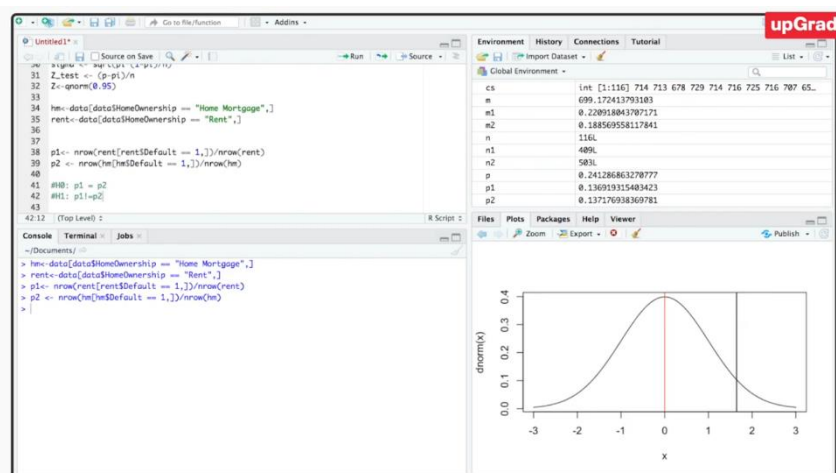
We do not include the customers who own houses here, as they are secured borrowers, and that any way given appropriate offers by Golden Corp, they want to make customized policy for the other group. The other two groups are also fundamentally comparable based on home ownership as they both pay a considerable proportion of income on houses, either in terms of EMI or in terms of rents.

DECISION CRITERIA

upGrad

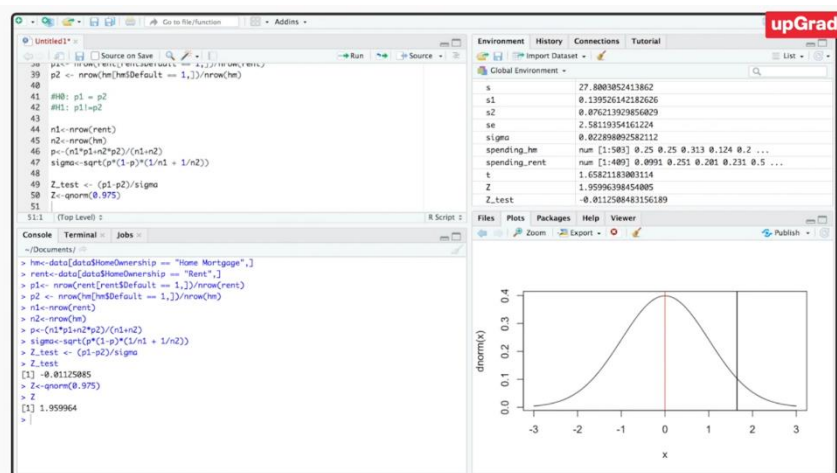


So, in order to identify the better target or the higher revenue generating group amongst the two, Ram decides to check for two factors, default rate and credit spending. So, he wants to check if there is a difference between the default rate of two groups. For this, we use hypothesis testing for difference in proportion.



So again, we will first need to subset the data to compute the proportion or default rate of the two groups. So, we'll define HM as the customers who have home mortgage, and we define rent as customers who live on rent. So, we subset on the home ownership column, and we create two subsets HM and rent here.

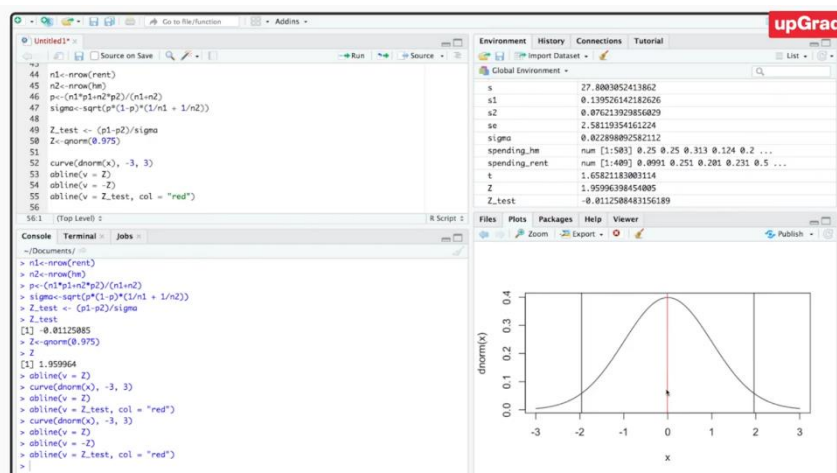
Next, we compute the default rate of each group. So, the default rate would simply be the number of defaulters in that group divided by the total number of customers and then 2. So, for example, P1 or default rate of those who rent houses would be equal to nrow, number of defaulters in the rent group divided by and nrow of rent group. And similarly, we will compute the default rate for those who have home mortgage.



Now, we want to test the null hypothesis that P1 is equal to P2, and the alternative would be P1 not equal to P2. We can define the sample size here, using again nrow command, and further for hypothesis testing of difference in proportion, we would need to compute P which is the pool proportion.

And if you remember, this is the formula for computing pool proportion, $n1$ into $P1$ plus $n2$ into $P2$ divided by $n1$ plus $n2$, and the sigma would be square root of $P \times 1 - P$ multiplied by 1 upon $n1$ plus 1 upon $n2$.

We have seen this formula in lectures. We have the sigma. So, we can calculate the test statistic. It would be $P1$ minus $P2$ divided by sigma. So, the test statistic is minus 0.01, and we can compute the critical value using qnorm function. Because it is a 2-tailed test, we'll use alpha by two, which is 2.5 and 1 minus 2.5 is 97.5. So, we will use qnorm of 0.975.

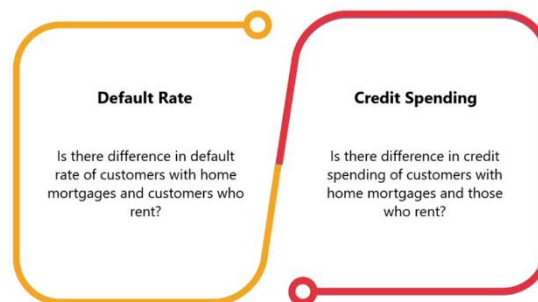


Critical value is 1.9. The test value is less than the critical value. We can also see it in the graph here. We are again using the curve function. This is a 2-tailed test. So, the rejection region would be on both side of the curves, and our test statistics lies here. It does not lie in the rejection region.

So based on the sample, we do not reject the null that default rate of customers with home mortgage is equal to default rate of customers who live on rents.

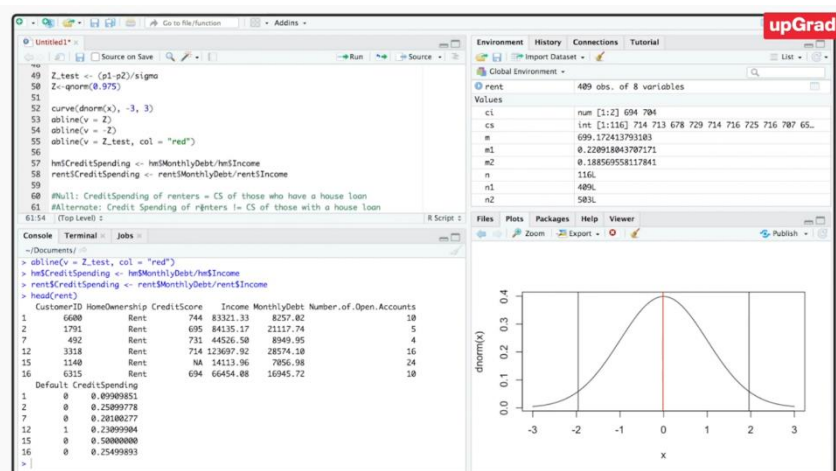
DECISION CRITERIA

upGrad



So, we have checked for the first factor default rate, we could not reject the null that the two groups have same default. Had they been different, and we could identify a group with a lower default rate, Golden Corp would have focused on that group to increase its revenue.

Next, we will check if there is difference in credit cards spending of customers with home loans and customers will live on rental properties. Generally speaking, customers or group who have a higher credit card bill are generally more valued by a bank or by a credit card company. It is a group which is generating more revenue for them.

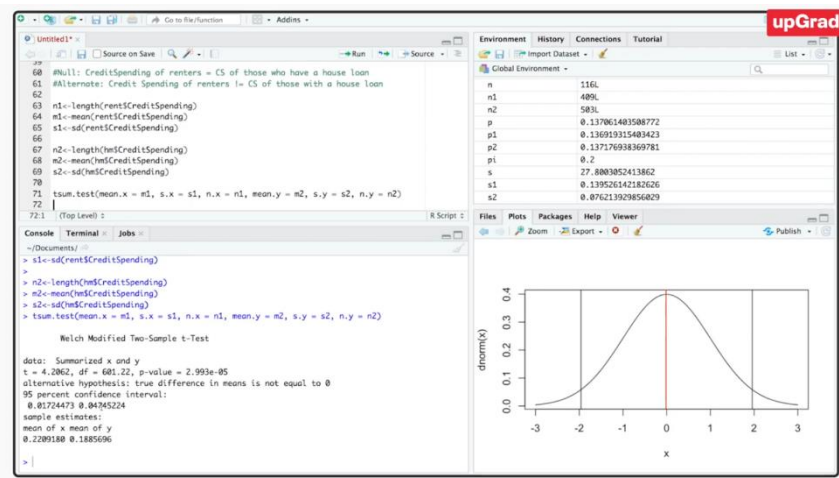


So, let's now check the hypothesis if we can see a difference in credit card spending of customers who have home loans versus customers who live on rent. First, we will approximate the credit spending using monthly debt or monthly credit card bill divided by the monthly income.

So, we will define a new column in the data set called credit spending, which is the monthly credit card bill as proportion of their income. So, we define this variable in both the groups. So, we have a new column called credit spending.

We want to conduct a difference in mean hypothesis. We want to test if credit spending of renters is equal to credit spending of those who have a house loan. So, we can define the null as credit spending of renters is equal to credit

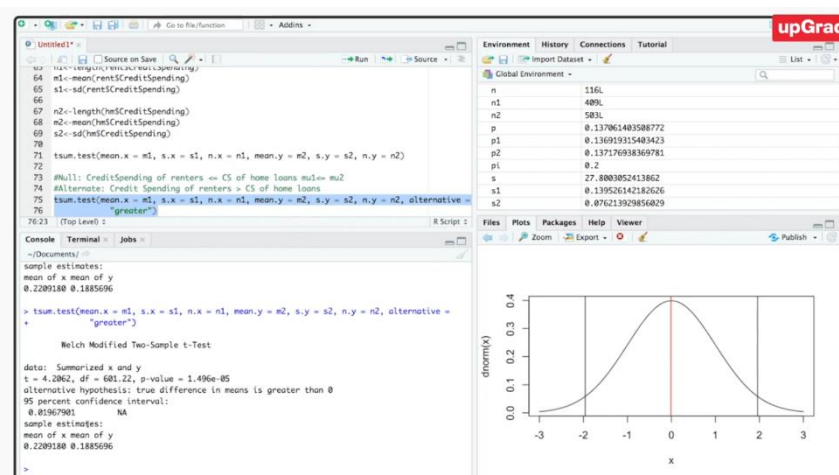
spending of those who have a house loan, and the alternative would be credit spending of renters is not equal to credit spending of those with a house loan.



Next, we can define the variables. So, we have n1 has lend of rent dollar credit spending, m1 as mean of rent dollar credit spending and standard deviation is s1, which is SD rent dollar credit spend.

Similarly, we have n2, m2, and s2 for those who have house loan. So, we can conduct the difference in mean test using t sum dot test function, we define mean dot x s dot x and dot x as n1, s1 and n1. And similarly for the second sample, we define m2, s2, and n2.

Once you run this, you can see that the 95% confidence interval for difference in means is 0.01 to 0.04. A null is the mean of first sample is equal to mean of second sample that is $\mu_1 - \mu_2$ is equal to zero. Since zero does not lie in this interval, we reject the null that credit spending of the two groups is same.

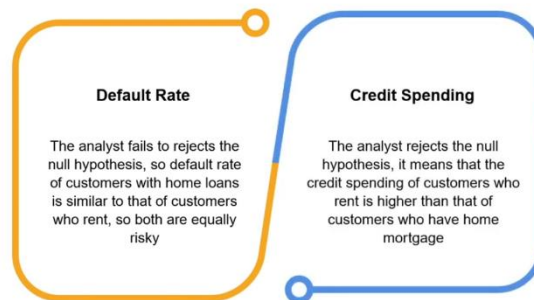


We can further test the hypothesis that spending of renters is less than spending of home loans. So, we will define the null as credit spending of renters less than equal to credit spending of home loans or μ_1 is less than equal to μ_2 , and the alternative will be credit spending of renters is greater than credit spending of home loans.

So, as our alternative is greater, in the t sum dot test function, we'll add another argument alternative equal to greater. So, the one-sided confidence interval or lower bound is 0.01, that is, based on this sample, we are 95% confident that the difference between credit spending of renters and credit spending of home loans is greater than or equal to 0.01. Since 0 does not lie in this interval, we reject the null that credit spending of those who rent is less than or equal to credit spending of those who have home loans.

DECISION CRITERIA

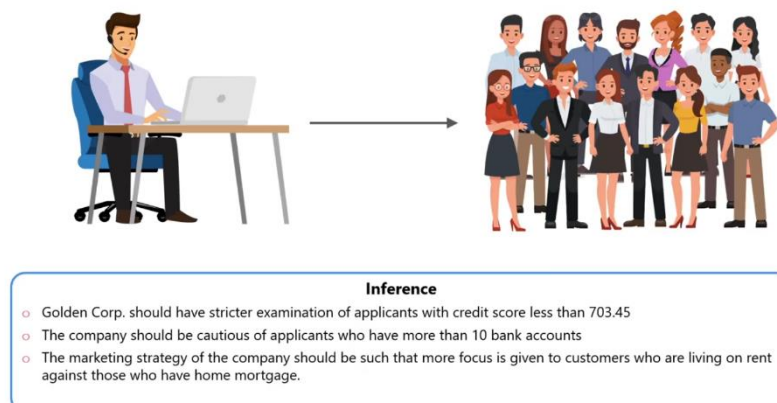
upGrad



So, coming back to choosing a better target, Ram concludes that based on the above sample, he cannot reject the hypothesis that default rate between the two groups is same. Further, he is able to reject the hypothesis that means spending of those who rent is less than means spending of those who have a house loan.

CASE STUDY

upGrad



Finally, based on the sample study, Ram concludes three points and reports them to Ms. Gocha.

1. Golden Corp should have stricter examination of applicants with credit score less than 703.4. We found that on average, there is a 95% chance that the credit score of defaulters is less than 703.4.
2. Secondly, Golden Corp should be cautious of applicants who have more than 10 bank accounts. Based on the sample, we fail to reject the null that default rate of customers with more than 10 accounts exceed 20%.

3. Finally, as we saw in the last test, it should target its marketing strategies to attract customers who are living on rent versus those who have a home mortgage.

No part of this publication may be reproduced, transmitted, or stored in a retrieval system, in any form or by any means, electronic, mechanical, photocopying, recording or otherwise, without the prior permission of the publisher.