

(22)

Categorical: State, Int.l.Plan, Vmail.Plan, Churn.(target variable)

Numerical: Account.Length, Area.Code (more likely to be categorical), Vmail.Message, CustServ.Calls

Day.Mins	Eve.Mins	Night.Mins	Intl.Mins
Day.Calls	Eve.Calls	Night.Calls	Intl.Calls
Day.Charge	Eve.Charge	Night.Charge	Intl.Charge

(دقت شود که Phone در اینجا بصورت ID رفتار میکند و بهتر است متغیر در نظر گرفته نشود و اینکه Area.code نیز اگر چه مقادیر عددی گرفته ولی هر عدد نمایانگر یک area است (رسته ای))

```

1 data_set = read.csv(file.choose(), header = T)[, -1]
2 str(data_set)

```

2:14 (Top Level)

Console Terminal Background Jobs

R 4.2.2 ~ /

* DONE (ggplot2)

The downloaded source packages are in
'C:\Users\amirmahdi\AppData\Local\Temp\Rtmpm0kei2\downloaded_packages'

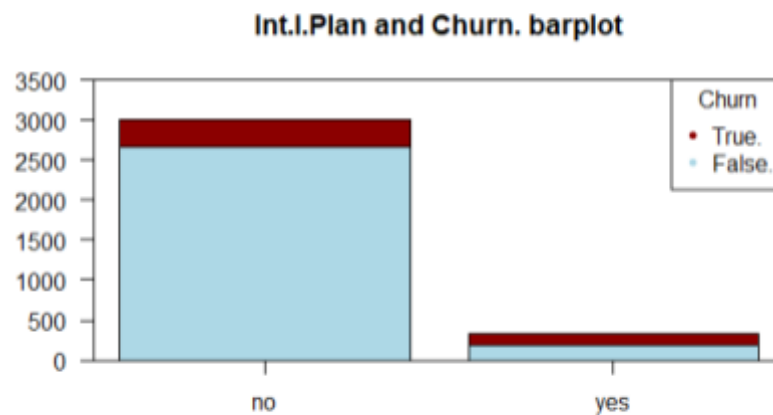
```

> data_set = read.csv(file.choose(), header = T)[, -1]
> str(data_set)
'data.frame': 3333 obs. of 21 variables:
 $ State      : chr  "KS" "OH" "NJ" "OH" ...
 $ Account.Length: int  128 107 137 84 75 118 121 147 117 141 ...
 $ Area.Code   : int  415 415 415 408 415 510 510 415 408 415 ...
 $ Phone       : chr  "382-4657" "371-7191" "358-1921" "375-9999" ...
 $ Int.l.Plan  : chr  "no" "no" "no" "yes" ...
 $ VMail.Plan  : chr  "yes" "yes" "no" "no" ...
 $ VMail.Message: int  25 26 0 0 0 0 24 0 0 37 ...
 $ Day.Mins    : num  265 162 243 299 167 ...
 $ Day.Calls   : int  110 123 114 71 113 98 88 79 97 84 ...
 $ Day.Charge  : num  45.1 27.5 41.4 50.9 28.3 ...
 $ Eve.Mins    : num  197.4 195.5 121.2 61.9 148.3 ...
 $ Eve.Calls   : int  99 103 110 88 122 101 108 94 80 111 ...
 $ Eve.Charge  : num  16.78 16.62 10.3 5.26 12.61 ...
 $ Night.Mins  : num  245 254 163 197 187 ...
 $ Night.Calls : int  91 103 104 89 121 118 118 96 90 97 ...
 $ Night.Charge: num  11.01 11.45 7.32 8.86 8.41 ...
 $ Intl.Mins   : num  10 13.7 12.2 6.6 10.1 6.3 7.5 7.1 8.7 11.2 ...
 $ Intl.Calls  : int  3 3 5 7 3 6 7 6 4 5 ...
 $ Intl.Charge : num  2.7 3.7 3.29 1.78 2.73 1.7 2.03 1.92 2.35 3.02 ...
 $ CustServ.Calls: int  1 1 0 2 3 0 3 0 1 0 ...
 $ Churn       : chr  "False." "False." "False." "False." ...

```

(25)

A) Churn and Int.l.Plan



Non-normalize barchart

data_set\$Churn.	data_set\$Int.l.Plan		Row Total
	no	yes	
False.	2664 0.885	186 0.576	2850
True.	346 0.115	137 0.424	483
Column Total	3010 0.903	323 0.097	3333

Statistics for All Table Factors

Pearson's Chi-squared test

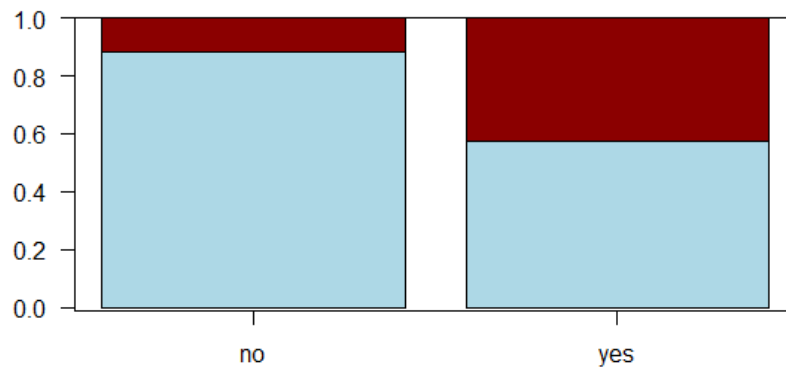
 $\chi^2 = 225.0541$ d.f. = 1 p = 7.14514e-51

Pearson's Chi-squared test with Yates' continuity correction

 $\chi^2 = 222.5658$ d.f. = 1 p = 2.493108e-50

Contingency table

Int.I.Plan and Churn. barplot

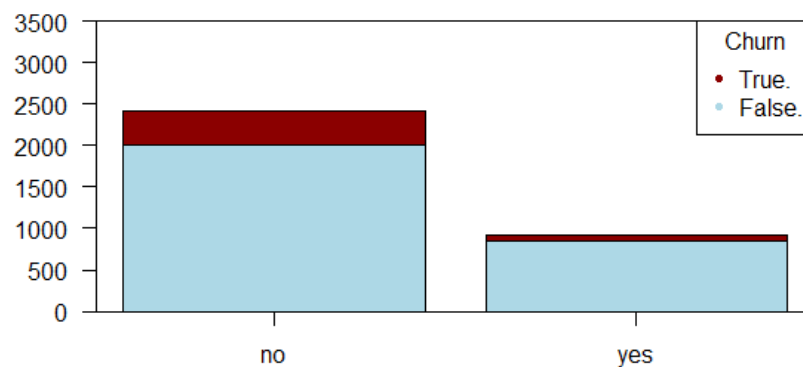


Normalized barchart

با بررسی barplot و contingency table میتوان نتیجه گرفت که درصد زیادی از مشتریانی که Int.I.Plan دارند، churn کردند (حدود 42%، چهار برابر کسانی که این سرویس را ندارند). پس گویا این سرویس برای مشتریان قابل قبول نبوده. البته باید توجه کرد که فراوانی مشتریانی که Int.I.Plan دارند نسبت به آنان که ندارند نیز بسیار کم است (323 / 3333) ضمناً p-value مقدار کمی دارد و فرض صفر (استقلال) رد میشود

B) Churn and Vmail.Plan

Vmail.Plan and Churn. barplot



Non-normalized barchart

data_set\$Churn.	data_set\$Vmail.Plan		Row Total
	no	yes	
False.	2008 0.833	842 0.913	2850
True.	403 0.167	80 0.087	483
Column Total	2411 0.723	922 0.277	3333

Statistics for All Table Factors

Pearson's Chi-squared test

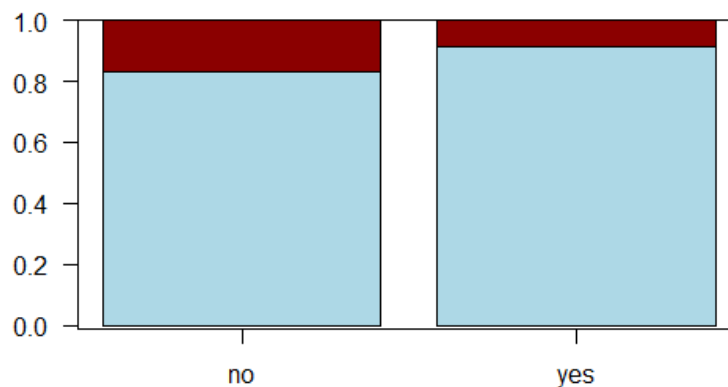
Chi^2 = 34.77733 d.f. = 1 p = 3.696527e-09

Pearson's Chi-squared test with Yates' continuity correction

Chi^2 = 34.13166 d.f. = 1 p = 5.15064e-09

Contingency table

Vmail.Plan and Churn. barplot

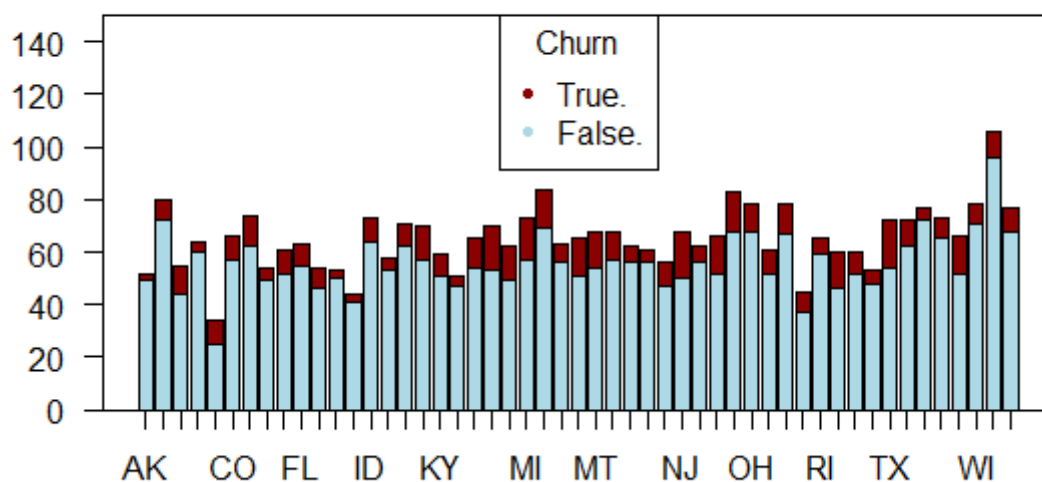


Normalized barchart

با بررسی barplot و contingency table میتوان نتیجه گرفت که سرویس Vmail.Plan برای مشتریان سرویس مناسبی بوده. زیرا درصد churn کنندگان در میان مشتریانی که Vmail.Plan را دارند حدود 8.7% و برای آنهایی که این سرویس را ندارند 17% می باشد. ضمناً p-value مقدار کمی دارد و فرض صفر (استقلال) رد میشود.

C) Churn and State

State and Churn. barplot

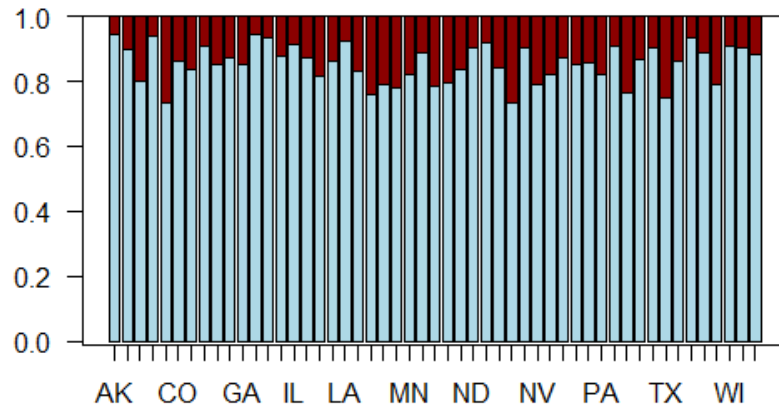


Non-normalized barchart

data_set\$State	data_set\$Churn.		Row Total
	False.	True.	
AK	49 0.017	3 0.006	52
AL	72 0.025	8 0.017	80
AR	44 0.015	11 0.023	55
AZ	60 0.021	4 0.008	64
CA	25 0.009	9 0.019	34
CO	57 0.020	9 0.019	66
CT	62 0.022	12 0.025	74
DC	49 0.017	5 0.010	54
DE	52	9	61

Contingency table

State and Churn. barplot



Normalized barchart

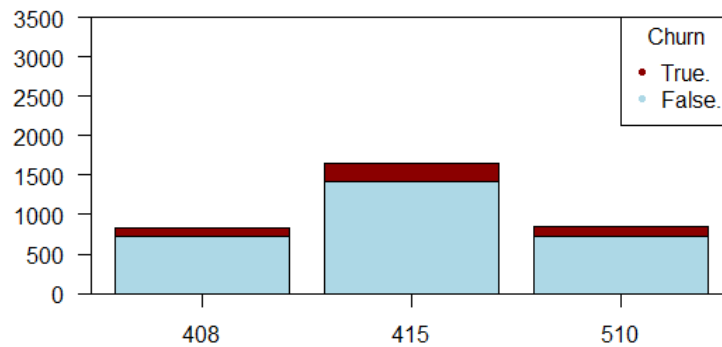
با بررسی barplot و contingency table نمیتوان با قطعیت در مورد استقلال یا عدم استقلال churn و State صحبت کرد. (p-value = 0.09).

تنها میتوان ایالت هایی که نسبت churn کنندگان در آنها به نسبت کل مشتری ها در همان ایالت بالا هست را یافت. (row percentage)

State	Percentage of Churn	Number of customers
CA	26%	34
MD	24%	70
NJ	26%	68
NV	21%	66
SC	23%	60
TX	25%	72
WA	21%	66

D) Churn and Area.Code

Area.Code and Churn. barplot



Non-Normalized barchart

data_set\$Churn.	data_set\$Area.Code			Row Total
	408	415	510	
False.	716 0.854	1419 0.857	715 0.851	2850
True.	122 0.146	236 0.143	125 0.149	483
Column Total	838 0.251	1655 0.497	840 0.252	3333

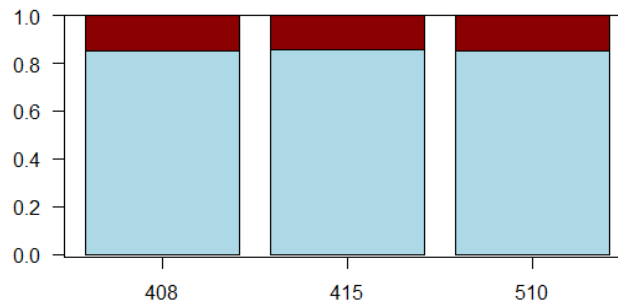
Statistics for All Table Factors

Pearson's Chi-squared test

Chi^2 = 0.1775407 d.f. = 2 p = 0.9150557

Contingency table

AreaCode and Churn. barplot



با بررسی barplot های بالا میبینیم که درصد Churn در هر سه AreaCode حدود 15% است و بنظر می رسد این دو متغیر از هم مستقل اند (البته تعداد مشتریان در AreaCode = 415 حدودا 2 برابر بقیه AreaCode هاست) ضمنا p-value مقدار زیادی دارد (90%) و فرض صفر (استقلال) قبول میشود

(b)

متغیر های Vmail.Plan و Int.l.Plan در هر الگوریتم data mining که در آینده قرار است بکار گرفته شود تاثیر بسزایی خواهند گذاشت

(26)

ارتباط متغیر Churn با دیگر متغیر های رسته ای

data_set\$Churn.	data_set\$Int.l.Plan		Row Total
	no	yes	
False.	2664 0.885	186 0.576	2850
True.	346 0.115	137 0.424	483
Column Total	3010 0.903	323 0.097	3333

Statistics for All Table Factors

Pearson's Chi-squared test

Chi² = 225.0541 d.f. = 1 p = 7.14514e-51

Pearson's Chi-squared test with Yates' continuity correction

Chi² = 222.5658 d.f. = 1 p = 2.493108e-50

data_set\$Churn.	data_set\$VMail.Plan		Row Total
	no	yes	
False.	2008 0.833	842 0.913	2850
True.	403 0.167	80 0.087	483
Column Total	2411 0.723	922 0.277	3333

Statistics for All Table Factors

Pearson's Chi-squared test

Chi² = 34.77733 d.f. = 1 p = 3.696527e-09

Pearson's Chi-squared test with Yates' continuity correction

Chi² = 34.13166 d.f. = 1 p = 5.15064e-09

data_set\$Churn.	data_set\$Area.Code			Row Total
	408	415	510	
False.	716 0.854	1419 0.857	715 0.851	2850
True.	122 0.146	236 0.143	125 0.149	483
Column Total	838 0.251	1655 0.497	840 0.252	3333

Statistics for All Table Factors

Pearson's Chi-squared test

Chi^2 = 0.1775407 d.f. = 2 p = 0.9150557

data_set\$State	data_set\$Churn.		Row Total
	False.	True.	
AK	49 0.017	3 0.006	52
AL	72 0.025	8 0.017	80
AR	44 0.015	11 0.023	55
AZ	60 0.021	4 0.008	64
CA	25 0.009	9 0.019	34
CO	57 0.020	9 0.019	66
CT	62 0.022	12 0.025	74
DC	49 0.017	5 0.010	54
DE	52	9	61

با بررسی درصد ها و همچنین فراوانی های crosstabulation های بالا میتوان دریافت که تفاوت درصد churn میان مشتریانی که int.l.plan دارند و ندارند بالاست (پس از churn مستقل نیست و pvalue نیز گویای این مطلب است). این تفاوت درصد در میان مشتریانی که vmail.plan دارند و ندارند نیز مشهود است. در مورد متغیر state بدلیل فراوانی حالت ها و زیاد نبودن p-value نمیتوان با اطمینان بالایی از مستقل بودن یا نبودن از churn صحبت کرد و تنها میتوان گفت برخی ایالت ها درصد و فراوانی churn بالایی دارند و باید به آنها توجه کرد. متغیر Area.Code باتوجه به درصد ها و pvalue از churn با اطمینان بالایی مستقل است.

ارتباط متغیر Vmail.Plan با دیگر متغیر های رسته ای

data_set\$VMail.Plan	data_set\$Int.l.Plan		Row Total
	no	yes	
no	2180 0.724	231 0.715	2411
yes	830 0.276	92 0.285	922
Column Total	3010 0.903	323 0.097	3333

Statistics for All Table Factors

Pearson's Chi-squared test

Chi^2 = 0.1202429 d.f. = 1 p = 0.7287712

Pearson's Chi-squared test with Yates' continuity correction

Chi^2 = 0.07913866 d.f. = 1 p = 0.7784681

data_set\$VMail.Plan	data_set\$Area.Code			Row Total
	408	415	510	
no	618 0.737	1184 0.715	609 0.725	2411
yes	220 0.263	471 0.285	231 0.275	922
Column Total	838 0.251	1655 0.497	840 0.252	3333

Statistics for All Table Factors

Pearson's Chi-squared test

Chi^2 = 1.368076 d.f. = 2 p = 0.5045755

data_set\$State	data_set\$Vmail.Plan		Row Total
	no	yes	
AK	36 0.015	16 0.017	52
AL	59 0.024	21 0.023	80
AR	41 0.017	14 0.015	55
AZ	45 0.019	19 0.021	64
CA	23 0.010	11 0.012	34
CO	47 0.019	19 0.021	66
CT	53 0.022	21 0.023	74
DC	36 0.015	18 0.020	54
DE	46 0.019	15 0.016	61
FL	43	20	63

از cross-tabulation های بالا میتوان دریافت که با اطمینان بالا میتوان گفت متغیر Vmail-plan با متغیر های Int.l.Plan و Area.Code مستقل است (با p-value های 0.7 و 0.5) در مورد رابطه Vmail.plan و State هم علی رغم زیاد بودن حالت های state با نگاه کردن به جدول و محاسبه p-value از آزمون chi-square میتوان دریافت که این دو متغیر مستقل اند (0.98).
تنها میتوان گفت تعداد مشتریانی که vmail-plan و int.l.plan را با هم ندارند بسیار زیاد (65%) است.

ارتباط متغیر Int.l.Plan و دیگر متغیرهای رسته ای

data_set\$Int.l.Plan	data_set\$Area.Code			Row Total
	408	415	510	
no	767 0.915	1505 0.909	738 0.879	3010
yes	71 0.085	150 0.091	102 0.121	323
Column Total	838 0.251	1655 0.497	840 0.252	3333

Statistics for All Table Factors			
Pearson's Chi-squared test			
Chi^2 = 7.936216	d.f. = 2	p = 0.01890917	

data_set\$State	data_set\$Int.l.Plan		Row Total
	no	yes	
AK	48 0.016	4 0.012	52
AL	72 0.024	8 0.025	80
AR	47 0.016	8 0.025	55
AZ	61 0.020	3 0.009	64
CA	30 0.010	4 0.012	34
CO	62 0.021	4 0.012	66
CT	66 0.022	8 0.025	74
DC	49 0.016	5 0.015	54
DE	51 0.017	10 0.031	61

حدود 45% از کل مشتریان، دارای Area.Code 405 و فاقد Int.l.Plan اند

ارتباط متغیر Area.Code با دیگر متغیر های رسته ای

data_set\$State	data_set\$Area.Code			Row Total
	408	415	510	
AK	14 0.017	24 0.015	14 0.017	52
AL	25 0.030	40 0.024	15 0.018	80
AR	13 0.016	27 0.016	15 0.018	55
AZ	15 0.018	36 0.022	13 0.015	64
CA	7 0.008	17 0.010	10 0.012	34
CO	25 0.030	29 0.018	12 0.014	66
CT	22 0.026	39 0.024	13 0.015	74
DC	14 0.017	27 0.016	13 0.015	54
DE	13 0.016	31 0.019	17 0.020	61
FL	12 0.014	31 0.019	20 0.024	63
GA	15 0.018	21 0.013	18 0.021	54
HI	15 0.018	30 0.018	8 0.010	53

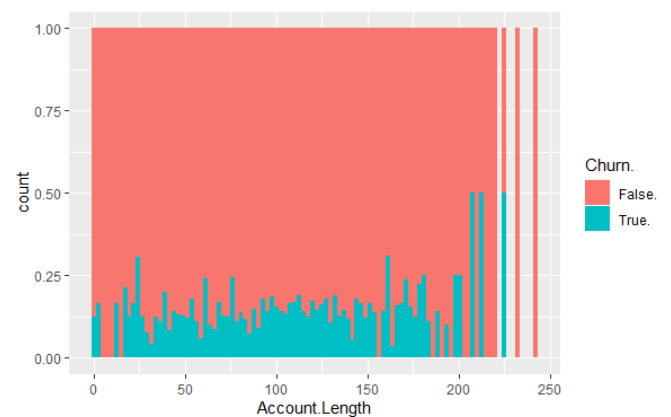
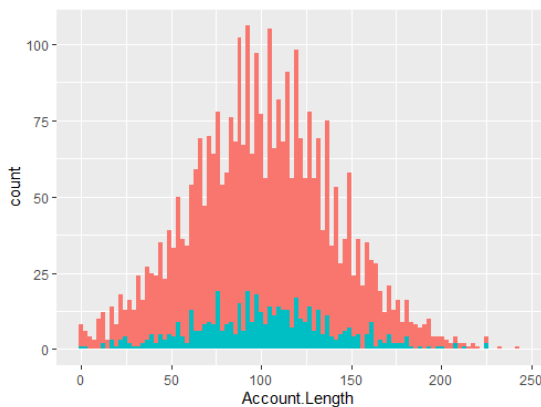
با بررسی درصد ها در هر سطر و $p\text{-value} = 0.6$ میتوان گفت این دو متغیر مستقلند.

(28)

مشتریانی که Churn کرده اند در AreaCode های مختلف بطور تقریباً یکسانی توزیع شده اند. ($p\text{-value} = 0.91$) همچنین در هر State مشتریان با AreaCode های مختلف نیز بطور تقریباً یکسان توزیع شده اند. ($p\text{-value} = 0.6$) همچنین متغیر Vmail_plan با متغیر های State و Int.l.Plan و Area.Code با مقادیر $p\text{-value} = 0.98, 0.77, 0.5$ (مستقل از Vmail.Plan اند) روابط بالا anomalous اند و باید قبل از مرحله ورودی دادن به الگوریتم با یک فرد expert در مورد این متغیرها صحبت کرد یا صحت data set را سوال کرد.

(30)

بررسی متغیر Account.length



```

19 churn.false = subset(data_set, data_set$Churn. == "False.")
20 churn.true = subset(data_set, data_set$Churn. == "True.")
21
22 t.test(churn.false$Account.Length, churn.true$Account.Length)
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100

```

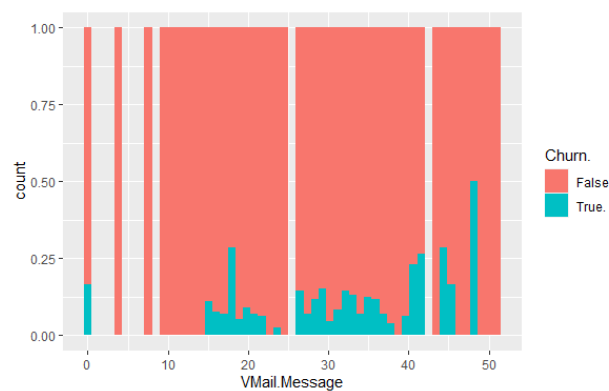
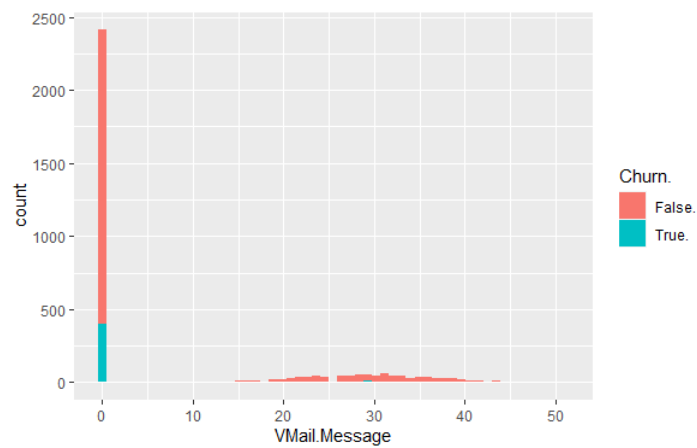
Welch Two Sample t-test

data: churn.false\$Account.Length and churn.true\$Account.Length
 $t = -0.96189$, $df = 659.91$, $p\text{-value} = 0.3365$
 alternative hypothesis: true difference in means is not equal to 0
 95 percent confidence interval:
 -5.690123 1.948299
 sample estimates:
 mean of x mean of y
 100.7937 102.6646

T-test for Churn and Account.Length relation

مقدار $p\text{-value}$ نشان میدهد میتوان با اطمینان خوبی گفت متغیر Account.Length از Churn مستقل است

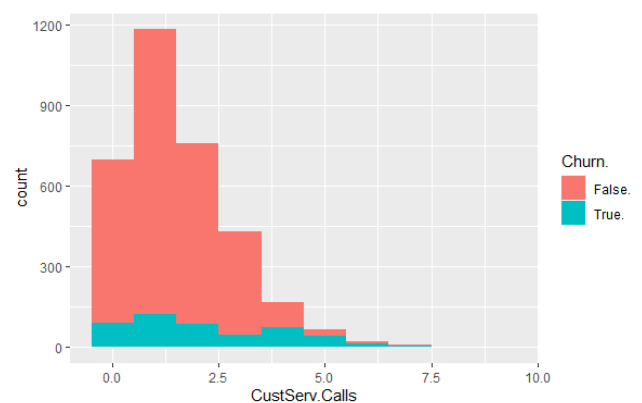
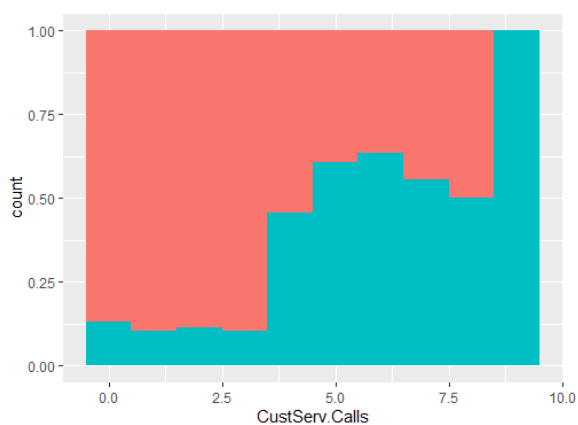
بررسی متغیر Vmail.Message



pvalue = 8.765e-09

تعداد بسیار زیادی از مشتریان این سرویس را ندارند که درصد churn نسبتاً پایینی هم دارند.

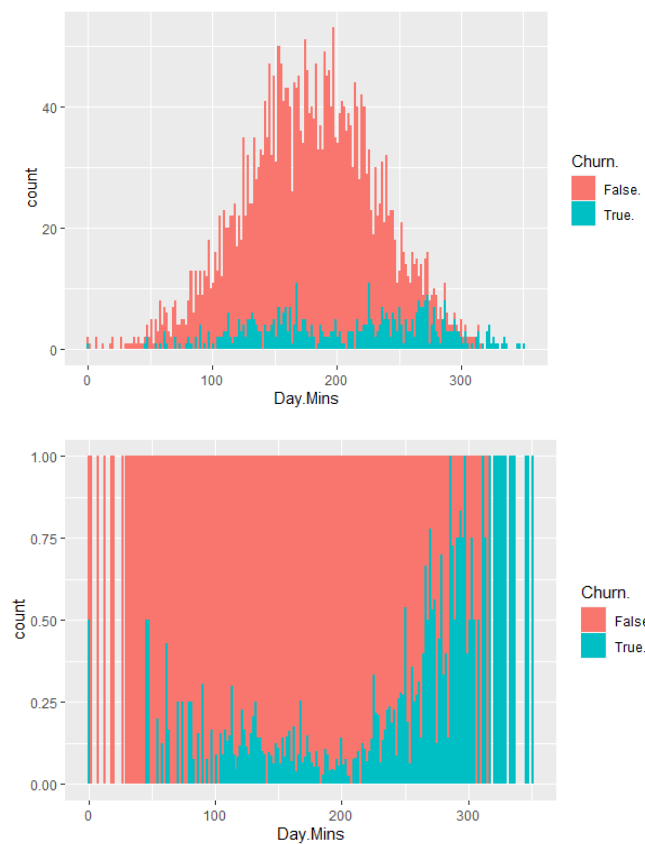
بررسی متغیر CustServ.Calls



p-value < 2.2e-16

با افزایش custServ.Calls، درصد churn افزایش می یابد (البته فراوانی مشتریان با custServ.Calls بالا نیز کم است)

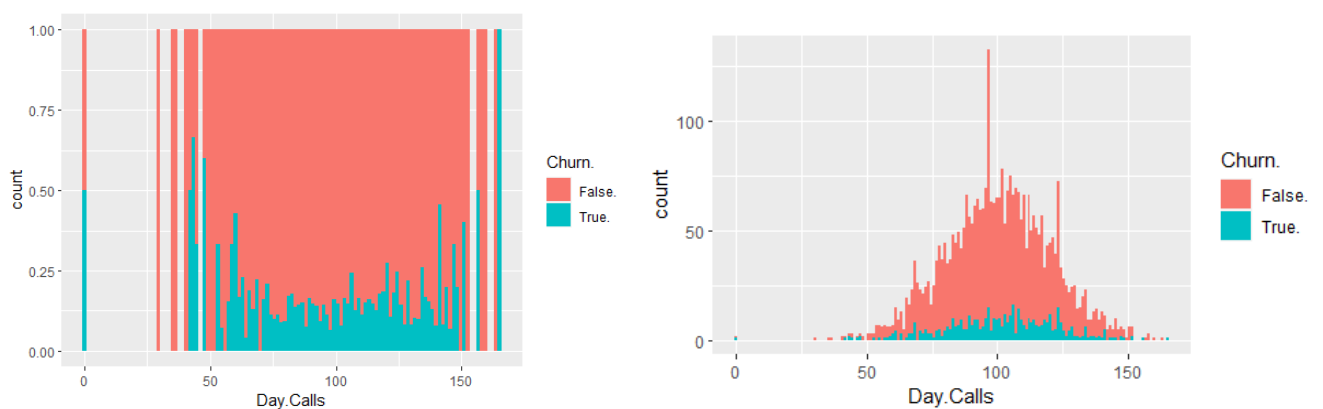
بررسی متغیر Day.Mins



p-value < 2.2e-16

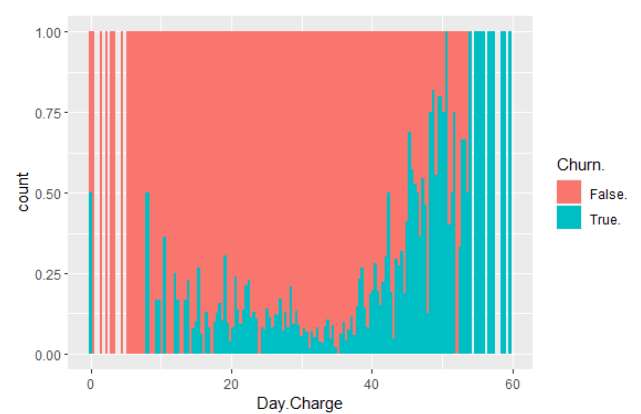
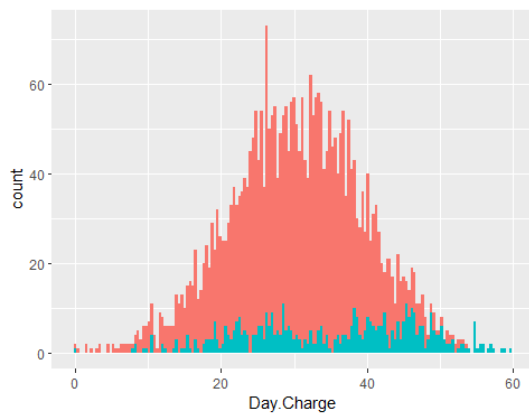
مشتریان با DayMins بیشتر از 250 احتمال churn بالایی دارند (البته اکثریت مشتریان در این بازه نیستند)

بررسی متغیر Day.Calls



p-value: 0.3165

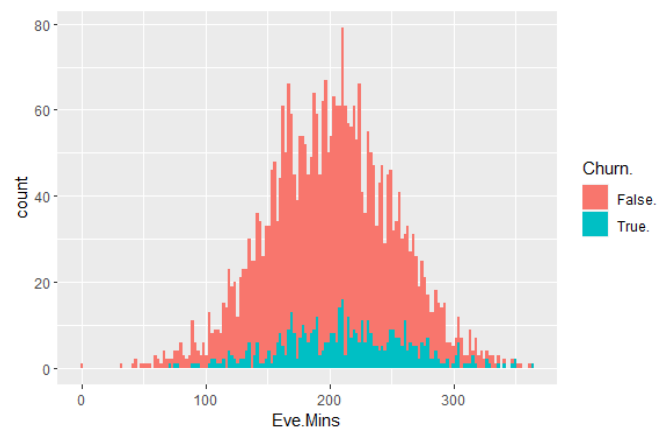
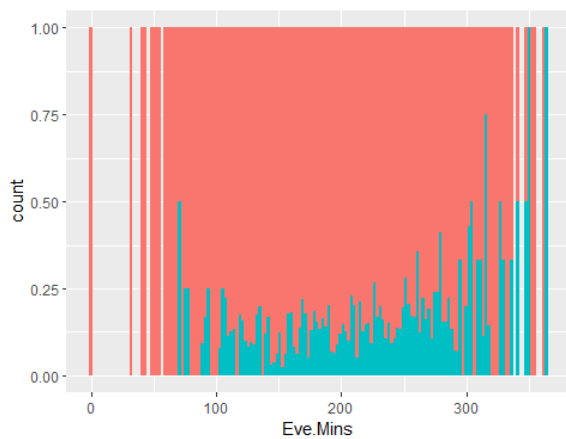
بررسی متغیر Day.Charge



P-value < 2.2e-16

این متغیر با متغیر Day.Mins ارتباط مستقیم دارد ($\text{cor}(\text{Day.Charge}, \text{Day.Mins}) = 1$) حاکی این مطلب است). مشتریان بالای 45 دلار Day.Charge، فراوانی کم و درصد churn بالایی دارند.

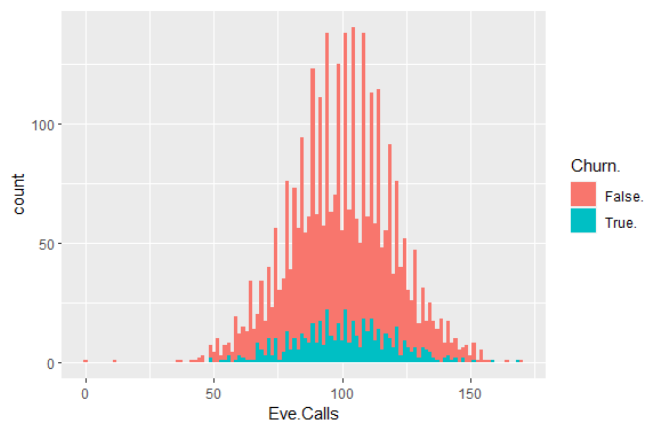
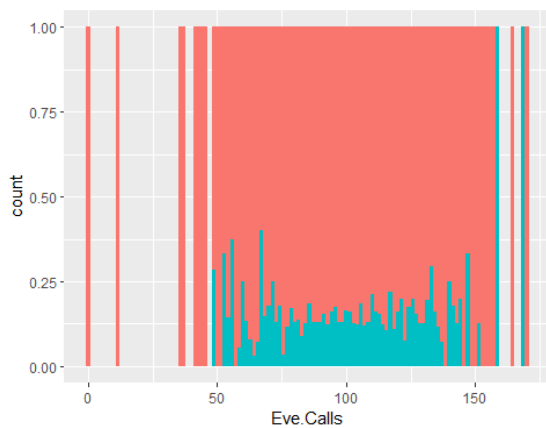
بررسی متغیر EVE.Mins



p-value: 1.839e-07

مشتریان بالای 300 EVE.Mins دارای churn rate بالایی هستند

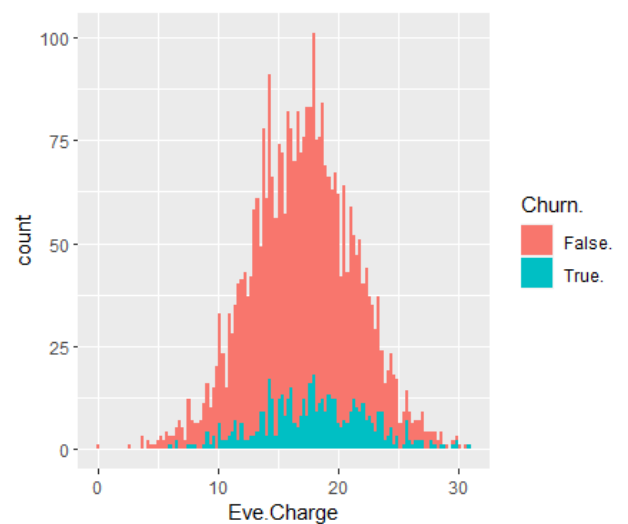
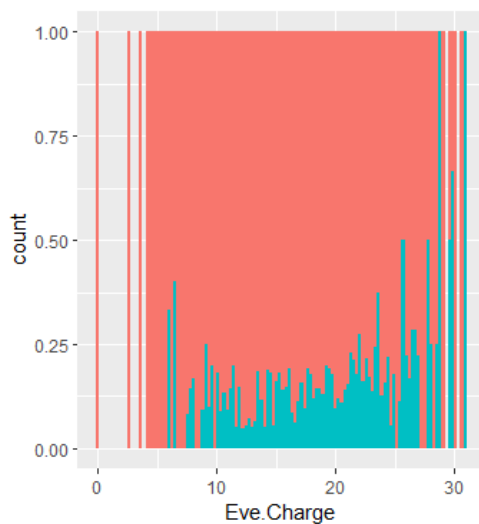
بررسی متغیر EVE.Calls:



p-value = 0.5912

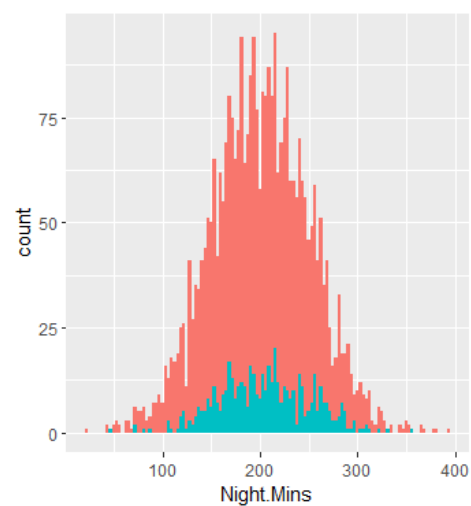
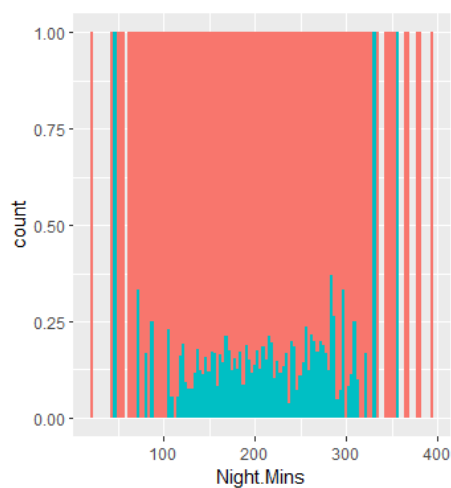
با اطمینان بالای میتوان گفت EVE.Calls و churn مستقل اند

بررسی متغیر EVE.Charge



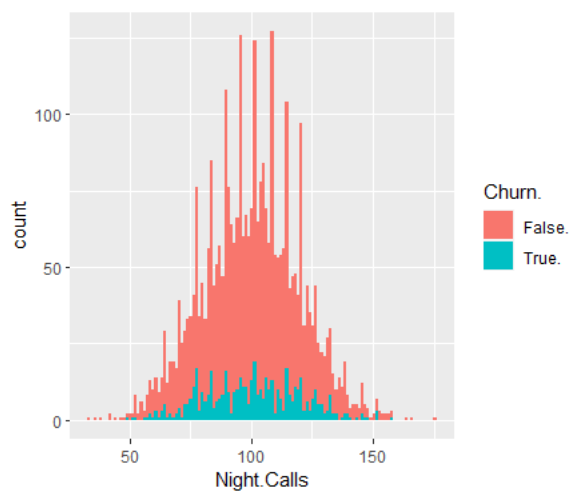
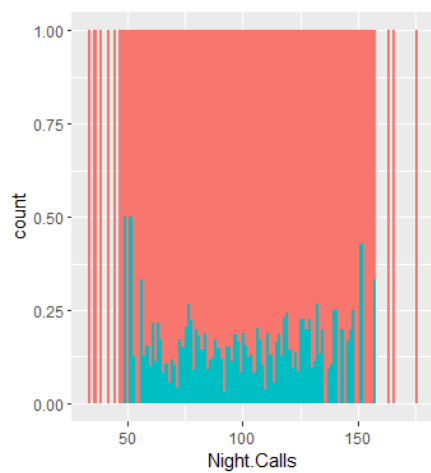
p-value: 1.843e-07

بررسی متغیر Nigth.Mins



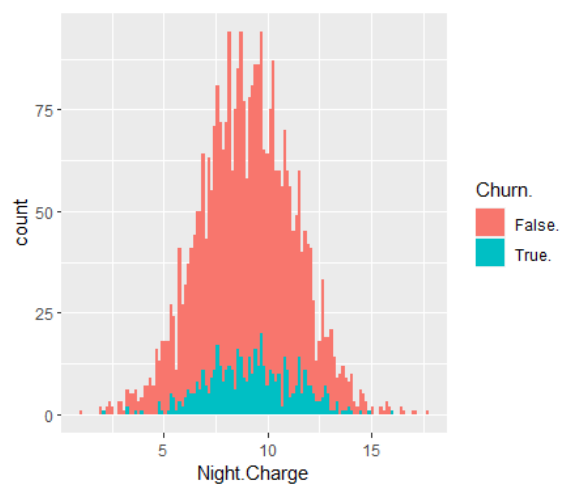
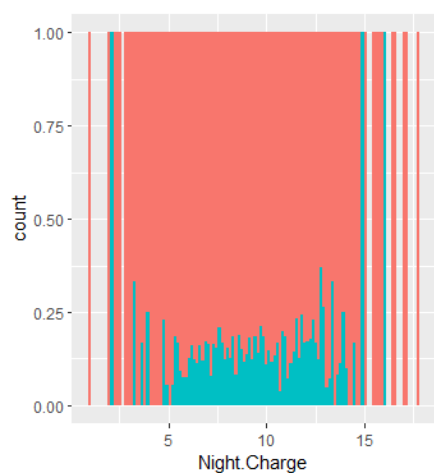
p-value: 0.03028

بررسی متغیر Night.Calls



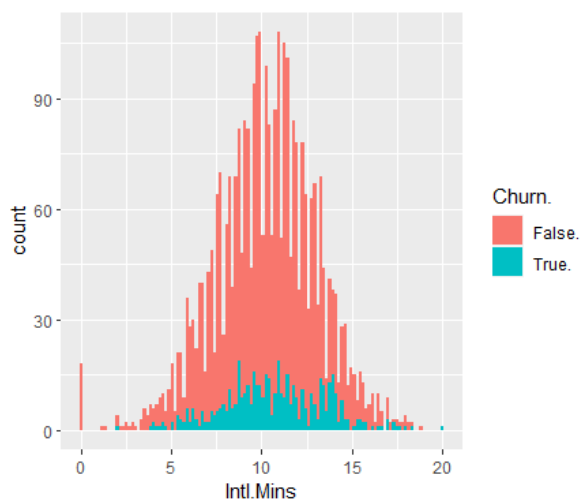
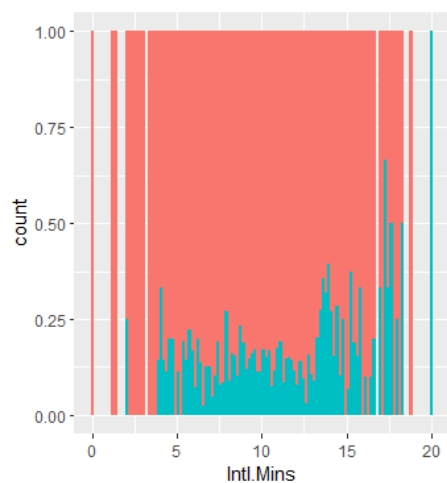
p-value = 0.7273

بررسی متغیر Night.Charge



p-value: 0.03027

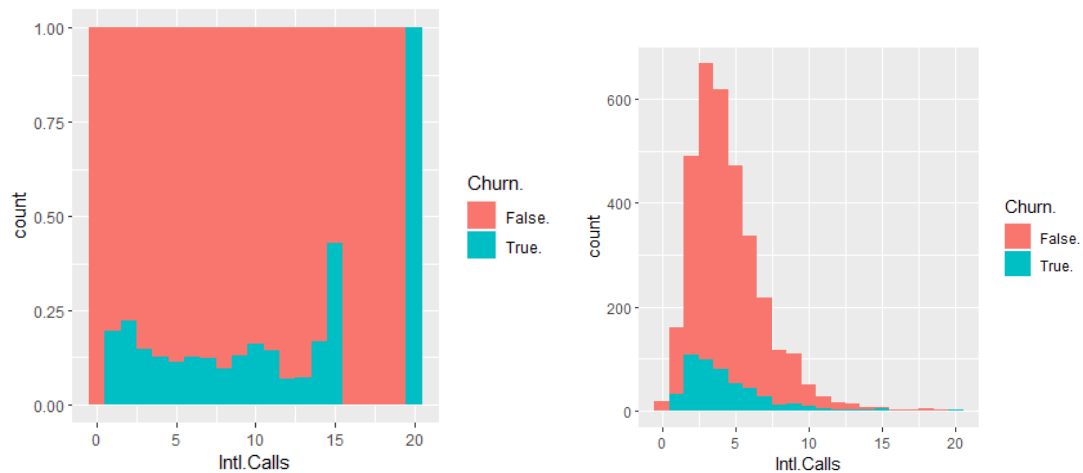
بررسی متغیر Intl.Mins :



p-value: 9.066e-05

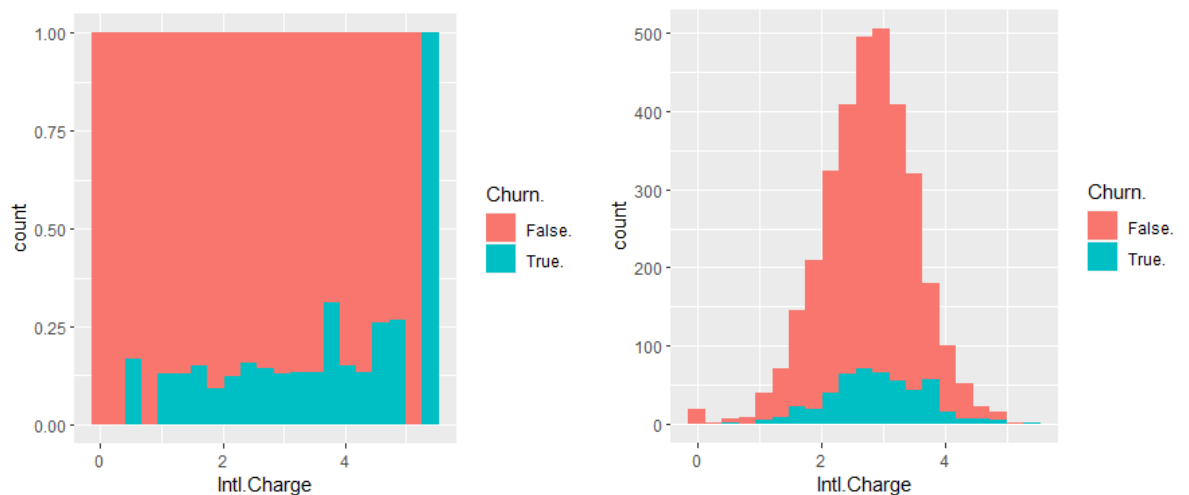
مشتریان دارای Intl.Mins بالای 15 Churn بالا و فراوانی کمی دارند

بررسی متغیر Intl.Calls:



p-value: 0.003186

بررسی متغیر Intl.Charge:



p-value: 9.026e-05

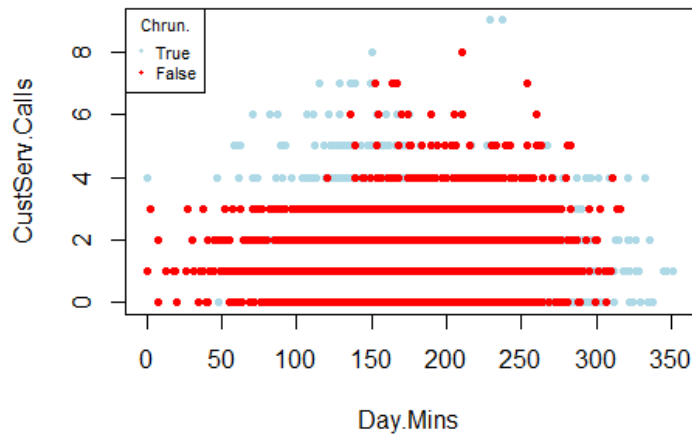
(b)

با مشاهده p-value ها متوجه میشویم که احتمالاً متغیرهای مربوط به Calls غیر international (یعنی Eve/Day/Night Calls) با churn مستقل اند و تاثیر بسزایی در الگوریتم ما نخواهند گذاشت (البته آنها را حذف نمیکنیم، بلکه برای مراحل بعد کاندید میکنیم).

با محاسبه cor بین متغیرهای عددی فوق متوجه میشویم که بین هر جفت Mins و Charge رابطه تقریباً خطی وجود دارد (در تمامی جفت متغیرهای Day.Charge و Day.Mins یا Night.Charge و Night.Mins و ...) پس بهتر است از هر جفت، یکی از متغیرها را در الگوریتم لحاظ کنیم.

بقیه متغیرها را باید در الگوریتم مان لحاظ کنیم

(32)



```

18
19 churn.false = subset(data_set, data_set$Churn. == "False." &
20                           data_set$Day.Mins < 150 & data_set$CustServ.Calls >= 4)
21 churn.true = subset(data_set, data_set$Churn. == "True." &
22                       data_set$Day.Mins < 150 & data_set$CustServ.Calls >= 4)
23
24 t.test(churn.true$Intl.Charge, churn.false$Intl.Charge)
25
26
16:55 (Top Level) R Scri
Console Terminal Background Jobs
R 4.2.2 ~ /
+ data_set$Day.Mins < 150 & data_set$CustServ.Calls >= 4)
> churn.true = subset(data_set, data_set$Churn. == "True." &
+                       data_set$Day.Mins < 150 & data_set$CustServ.Calls >= 4)
> t.test(churn.true$Intl.Charge, churn.false$Intl.Charge)

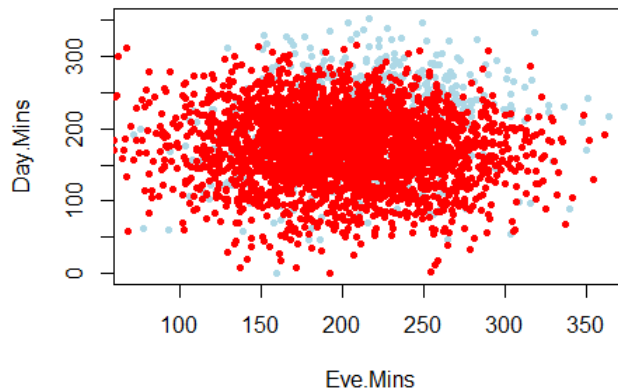
Welch Two Sample t-test

data: churn.true$Intl.Charge and churn.false$Intl.Charge
t = -0.31876, df = 7.652, p-value = 0.7584
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
-0.6491370 0.4925655
sample estimates:
mean of x mean of y
2.724571 2.802857

```

از سوال 30 میتوان دریافت که مشتریان با CustServ.Calls بالای 3 churn rate بالایی دارند. اما وقتی این متغیر با Day.Mins بصورت همزمان بررسی میشوند میتوان دریافت که آن دسته از مشتریان که CustServ.Calls بیشتر از 3 دارند اما Day.Mins بالای 150 هم دارند کمتر churn میکنند.

از طرفی آن دسته از مشتریانی که CustServ.Calls کمتر از 3 و Day.Mins شان از 300 بالاتر است churn rate بالایی دارند.



```

19 churn.false = subset(data_set, data_set$Churn. == "False." &
20   data_set$Day.Mins > 250 & data_set$Eve.Mins >= 200)
21 churn.true = subset(data_set, data_set$Churn. == "True." &
22   data_set$Day.Mins > 250 & data_set$Eve.Mins >= 200)
23
24 t.test(churn.true$Intl.Charge, churn.false$Intl.Charge)
25
26
25:1 (Top Level) :-

```

```

Console Terminal Background Jobs
R 4.2.2 ~ /
+ data_set$Day.Mins > 250 & data_set$Eve.Mins >= 200)
> churn.true = subset(data_set, data_set$Churn. == "True." &
+ data_set$Day.Mins > 250 & data_set$Eve.Mins >= 200)
> t.test(churn.true$Intl.Charge, churn.false$Intl.Charge)

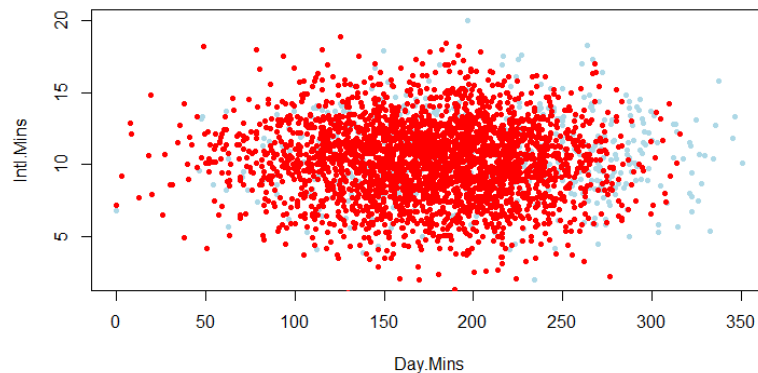
Welch Two Sample t-test

data: churn.true$Intl.Charge and churn.false$Intl.Charge
t = -0.38317, df = 133.36, p-value = 0.7022
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
 -0.2559487  0.1728745
sample estimates:
mean of x mean of y
 2.793540  2.835077

```

از سوال 30 میتوان دریافت که مشتریان با Day.Mins بالای 250 churn بالایی دارند، اما آن دسته از مشتریان با Day.Mins بالای 250 و Eve.Mins کمتر از 200 churn بالایی ندارند. در نتیجه مشتریان با Day.Mins و Eve.Mins بالا بطور قابل توجهی churn rate بالایی دارند.

از آنجا که رابطه مستقیمی بین Eve.Mins و Eve.Charge وجود دارد، در نتیجه نتایج فوق برای نمودار Day.Mins و Eve.Charge نیز برقرار است.



```

19 churn.false = subset(data_set, data_set$Churn. == "False." &
20                             data_set$Day.Mins > 200 & data_set$Intl.Mins >= 13)
21 churn.true = subset(data_set, data_set$Churn. == "True." &
22                     data_set$Day.Mins > 200 & data_set$Intl.Mins >= 13)
23
24 t.test(churn.true$Intl.Charge, churn.false$Intl.Charge)
25
26
25:1 (Top Level) :

```

```

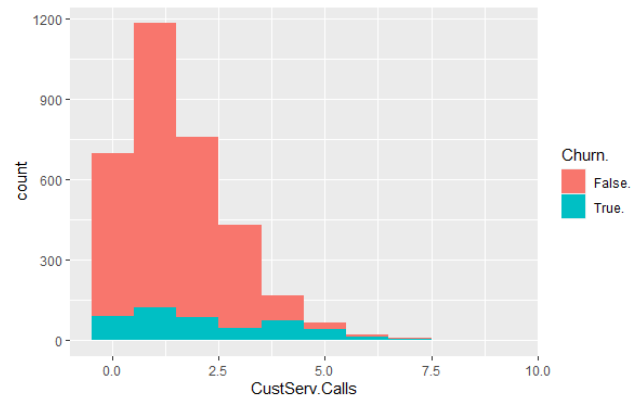
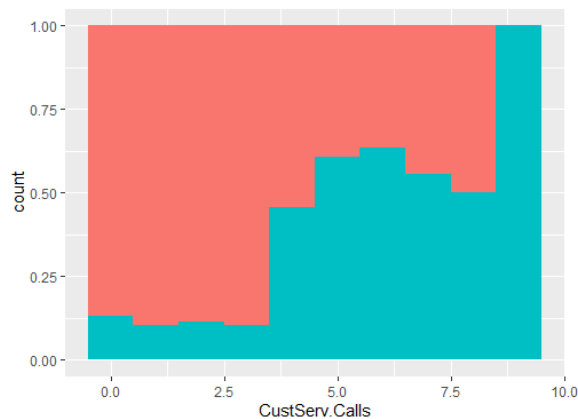
R 4.2.2 ~ /
+ data_set$Day.Mins > 200 & data_set$Intl.Mins >= 13)
> churn.true = subset(data_set, data_set$Churn. == "True." &
+ data_set$Day.Mins > 200 & data_set$Intl.Mins >= 13)
> t.test(churn.true$Intl.Charge, churn.false$Intl.Charge)

Welch Two Sample t-test

data: churn.true$Intl.Charge and churn.false$Intl.Charge
t = 1.3857, df = 87.634, p-value = 0.1694
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
 -0.03280006  0.18386589
sample estimates:
mean of x mean of y
 3.896897  3.821364

```

از سوال 30 بنظر می رسید آن دسته مشتریانی که Intl.Mins بالای 13 دارند با احتمال بیشتری Churn میکنند. اما با نگاه به ترکیب این متغیر با Day.Mins میتوان متوجه شد که مشتریان با Intl.Mins بالای 13 و Day.Mins بالای 200 هستند که درصد بیشتر churn کنندگان حالت univariate را تشکیل میدهند.



مثال کتاب: تقسیم CustServ.Calls به دو دسته low (کمتر از 4) و high (بیشتر از 4).

```

29 new_data_set = cbind(CustServ = data_set$CustServ.Calls,
30                       ChurnRate = data_set$Churn.)
31
32 # nrow(new_data_set)
33 # as.numeric(new_data_set[, 1])
34
35 for(i in 1 : nrow(new_data_set)){
36   if(as.numeric(new_data_set[i, 1]) >= 4){
37     new_data_set[i, 1] = "high";
38   }
39   else{
40     new_data_set[i, 1] = "low";
41   }
42 }
43
44 library(gmodels)
45 CrossTable(y = new_data_set[, 1], x = new_data_set[, 2],
46            prop.r = FALSE, prop.t = FALSE,
47            chisq = TRUE, prop.chisq = FALSE,
48            dnn = c("Churn", "Cust.Serv.Calls"))

```

Churn	high	low	Row Total
False.	129	2721	2850
	0.483	0.887	
True.	138	345	483
	0.517	0.113	
Column Total	267	3066	3333
	0.080	0.920	

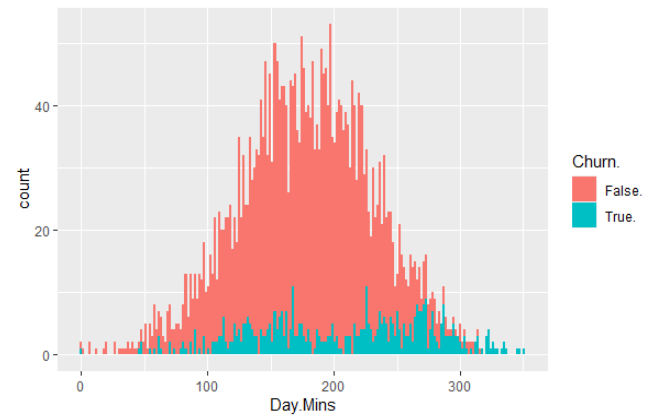
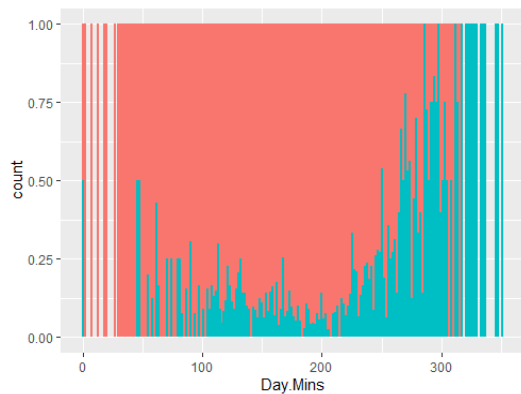
Statistics for All Table Factors

Pearson's Chi-squared test

Chi^2 = 324.0392 d.f. = 1 p = 1.910278e-72

درصد churn در این تقسیم بندی جدید حدود 4 برابر (0.517 / 0.113) می باشد و شرط استقلال قویا رد میشود که نشان میدهد تقسیم بندی خوبی است (Maximize the effect of classes)

یک مثال دیگر



```

29 new_data_set = cbind(CustServ = data_set$Day.Mins,
30                       ChurnRate = data_set$Churn.)
31
32 # nrow(new_data_set)
33 # as.numeric(new_data_set[, 1])
34
35 for(i in 1 : nrow(new_data_set)){
36   if(as.numeric(new_data_set[i, 1]) >= 250){
37     new_data_set[i, 1] = "high";
38   }
39   else{
40     new_data_set[i, 1] = "low";
41   }
42 }
43
44 library(gmodels)
45 CrossTable(y = new_data_set[, 1], x = new_data_set[, 2],
46            prop.r = FALSE, prop.t = FALSE,
47            chisq = TRUE, prop.chisq = FALSE,
48            dnn = c("Churn", "Day.Mins"))
49
44:17 [Untitled]

```

Console Terminal Background Jobs

R 4.2.2

Total Observations in Table: 3333

Churn	Day.Mins		Row Total
	high	low	
False.	174 0.530	2676 0.891	2850
True.	154 0.470	329 0.109	483
Column Total	328 0.098	3005 0.902	3333

Statistics for All Table Factors

Pearson's Chi-squared test

Chi^2 = 309.3387 d.f. = 1 p = 3.042938e-69

تقسیم بندی Day.Mins به دو دسته low (کمتر از 250) و high (بیشتر از 250) درصد churn کردن Day.Mins = high حدود 5 برابر Day.Mins = low می باشد. ضمناً p-value نیز استقلال را رد میکند (Maximize the effect of classes)

حال تقسیم بندی را طوری انجام می دهیم که اثر کلاس ها minimum شوند (یعنی با این کلاس بندی متغیر Day.Mins و Churn مستقل میشوند)

```

29 new_data_set = cbind(CustServ = data_set$Day.Mins,
30                       ChurnRate = data_set$Churn.)
31 for(i in 1 : nrow(new_data_set)){
32   if(as.numeric(new_data_set[i, 1]) <= 100){
33     new_data_set[i, 1] = "low";
34   }
35   else{
36     new_data_set[i, 1] = "high";
37   }
38 }
39
40 CrossTable(y = new_data_set[, 1], x = new_data_set[, 2],
41            prop.r = FALSE, prop.t = FALSE,
42            chisq = TRUE, prop.chisq = FALSE,
43            dnn = c("Churn", "Day.Mins"))
44
45
46
47
48
49
50

```

Churn	Day.Mins		Row Total
	high	low	
False.	2639 0.853	211 0.883	2850
True.	455 0.147	28 0.117	483
Column Total	3094 0.928	239 0.072	3333

Statistics for All Table Factors

Pearson's Chi-squared test

Chi^2 = 1.601107 d.f. = 1 p = 0.2057464

همانطور که دیده میشود، اگر مرز میان low و high را در Day.Mins مساوی 100 قرار دهیم، بخوبی نتوانسته ایم تفاوت churn rate را در تقسیم بندی مان نشان دهیم (churn rate ها در هر ستون مقدار تقریباً برابری با سطر نظیرشان دارند و p-value بزرگتر از 0.05)

(34)

تقسیم بندی به دو دسته high و low بر اساس روش equal-frequency binning دقت شود که باید data_set را قبل از binning مرتب (sort) کنیم

همانطور که دیده میشود، تفاوت قابل توجهی میان درصد churn کنندگان و آنانی که churn نمیکنند در DayMins = high / low دیده نمیشود. پس تقسیم بندی مناسبی نیست

```

2 numOfBins = 2
3 freq1 = round(numOfRecords / numOfBins, 0);
4 freq2 = freq1 + 1
5
6 sortedDayMins = new_data_set[order(as.numeric(new_data_set[, 1]),
7                                   decreasing = FALSE), ]
8 level = c("low", "high")
9
10 for(i in 1:2){
11   for(j in 1:numOfRecords){
12     if((i - 1) * freq < j && j <= i * freq){
13       sortedDayMins[j] = level[i]
14     }
15   }
16 }
17 sortedDayMins[3333, 1] = "high"
18 library(gmodels)
19 CrossTable(sortedDayMins[, 2], sortedDayMins[, 1],
20           prop.r = FALSE, prop.t = FALSE)

```

2-14 (Top Level)

Console Terminal Background Jobs

R 4.2.2 ~/

Churn	Day.Mins high	low	Row Total
False.	1375 0.825	1475 0.885	2850
True.	292 0.175	191 0.115	483
Column Total	1667 0.500	1666 0.500	3333

Statistics for All Table Factors

Pearson's Chi-squared test

Chi^2 = 24.62856 d.f. = 1 p = 6.951394e-07

روش equal-width binning:

```

25 (range = (max(as.numeric(new_data_set[, 1])) -
26         min(as.numeric(new_data_set[, 1])) + 1) / numOfBins);
27
28 equal_length_data_set = new_data_set;
29
30 for(i in 1:numOfBins){
31   for(j in 1:numOfRecords){
32     if((i - 1) * range <= as.numeric(new_data_set[j, 1]) &&
33       as.numeric(new_data_set[j, 1]) <= i * range){
34       equal_length_data_set[j, 1] = level[i]
35     }
36   }
37 }
38
39 library(gmodels)
40 CrossTable(equal_length_data_set[, 2], equal_length_data_set[, 1],
41           prop.r = FALSE, prop.t = FALSE
42           , chisq = TRUE, prop.chisq = FALSE,|
43           dnn = c("Churn", "Day.Mins"))

```

Churn	Day.Mins		Row Total
	high	low	
False.	1454 0.830	1396 0.882	2850
True.	297 0.170	186 0.118	483
Column Total	1751 0.525	1582 0.475	3333

Statistics for All Table Factors

Pearson's Chi-squared test

chi^2 = 18.16722 d.f. = 1 p = 2.023317e-05

در این روش نیز همانطور که دیده میشود، تفاوت درصد میان churn = False و churn = True در مقایسه سطری به طرز قابل توجهی وجود ندارد، و درصد های churn در گروه های low و high یکسان است.

(6) نظر گرایی جدول تعاملی چیست؟ clustered bar chart

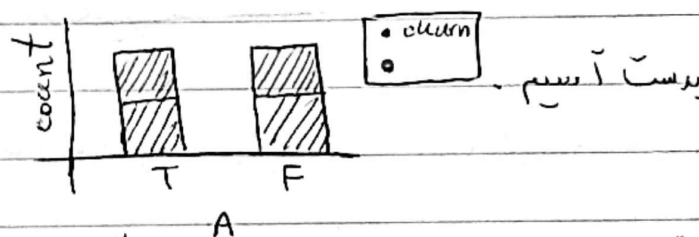
(7) با یک مثال، اثر متقابل (interaction) بین دو متغیر کیفی را بیان کنید.

نمایش دهنده دو متغیر کیفی (A و B) بصورت دایره‌های رابطه‌ها با متغیر هدف بررسی

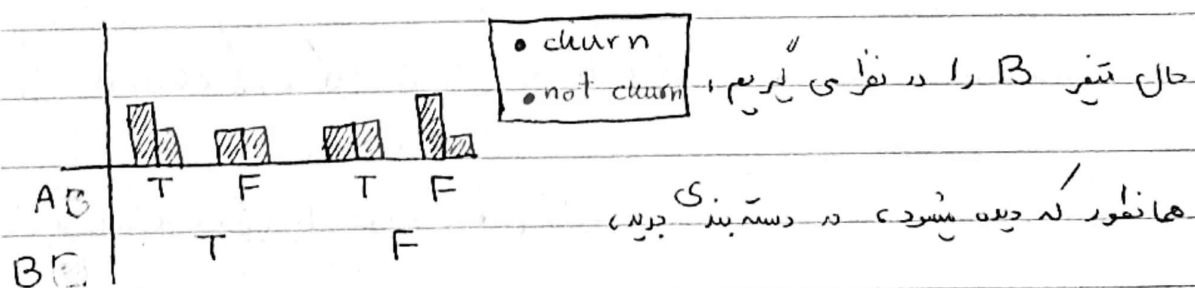
میشود، در صورتیکه رابطه جدید بدست آمده با رابطه هر کدام از متغیرها A و B

بصورت کلی با متغیر هدف تفاوت داشته باشد، می‌گوئیم A و B با هم interaction دارند.

بنوان مثال، فرض کنید رابطه متغیر A با دو مقدار T و F، با churn



با نگاه به bar chart بالا، گویا متغیر پیشرو A ارتباطی با target ندارد. اما



در هر کدام از چهار حالت (A و B)، بین churn کردن و آنبانی churn

تفاوت تفاوت معناداری وجود دارد، در نتیجه متغیر A و B با هم interaction

دارند.

می‌دانند
(8) overlay histogram
رابطه میان متغیر هدف و یک متغیر

numerical را نشان دهد. برای نشان دادن proportion متغیر هدف
به طور واضح از حالت normalized استفاده می‌کنیم

(9) مزیت normalized histogram : نسبت متغیر هدف و در هر ستون متغیر

بطلو بطور واضح نشان میدهد و ستون ارتباط میان متغیر هدف و ستون را

راحت‌تر و visible مشاهده کرد. (برای مثال با زیاد شدن مقدار متغیر ستون، درصد متغیر هدف churn

نیز بیشتر می‌شود)

عیب normalized histogram : محدود فرادانی متغیر ستون اطلاعاتی نمی‌دهد. برای

مثال میلست تنها دو نفر در یک ستون از histogram هستند که اگر آن دو

نفر churn نشد، آن ستون در حالت normalized و churn می‌شود و میلست

برداشت درستی از نمودار ابقام نشود

عزدار
با توجه به توضیحات داده شده، بهتر است برای overlay histogram هم

normalized و هم non normalized را رسم کرد.

(13) روش Binning بر اساس متغیر بشماره در این روش متغیر کیفی A

را بر اساس درصدهای متغیر هدف (که به تعداد گروه تقسیم می کنیم) (داده با)

تبدیل به متغیر کیفی می کنیم. برای مثال، متغیر A مقادیر 0 تا 10

می گیرد، اما در صد churn در record های متغیر A با مقدار 3 از 3،

چندین برابر record های با مقدار 3 است. بنابراین، می توان متغیر A را به

دو گروه slow (کمتر از 3) و high (بیشتر از 3) قرار داد.

(15) باید قبل از derive کردن متغیر جدید به scale متغیرهای numerical قبلی

ترتیب کنیم. برای مثال، اگر متغیر عددی A، میانگین و انحراف معیار بیشتر از متغیر B دارد،

میانگین گرفتن از record های A و B باعث bias نسبت متغیر A می شود. پس بهتر است

قبل از میانگین گرفتن، هر دو متغیر A و B را به Z-score استاندارد کنیم.

(19) میان متغیرهای بشماره عددی، Pearson correlation را بدست می آوریم.

چنانچه این مقدار 1 بود، یکی از دو متغیر بشماره را حذف می کنیم (این کار را

می توان با matrix plot و یا با جیل "correlation" انجام داد) برای متغیرهای بشماره

با correlation بالا و p-value کم) باید برای مایل بندی

Data Mining: کاربرد برای حذف شواهد

(20

Bar charts: (لینی) رسته ای

Histograms: هردو

Summary statistics: هردو (برای لینی ها) و ^{barplot} نمودار را (نشان میدهد)

Crosstabulations: (رسته ای) و ~~نمودار~~

Correlation analysis: عددی (رسته)

Scatter plot: هردو

Binning: عددی