



Customer Churn Reduction

Contents

Introduction.....	2
1.1 Problem Statement.....	2
1.2 Dataset.....	2
Methodology	0
2.1 Exploratory Data Analysis (EDA).....	0
2.1.1 The Target Variable - Churn	0
2.1.2 Missing Value Analysis	1
2.1.3 Multicollinearity	1
2.1.4 Analysis of churn ratio with different predictors	2
2.2 Modeling.....	10
2.2.1 Logistic regression.....	10
R implementation	10
Python implementation	13
2.2.2 Random Forest.....	13
R Implementation	13
Python Implementation	14
Result and performance measure.....	15
3.1 Python Implementation Results.....	15
Logistic Regression.....	15
Random Forest.....	16
3.2 R Implementations results	18
Logistic Regression.....	18
Random Forest.....	18
Logistic regression diagnostics.....	20
R code	22
Python code	35

Chapter 1

Introduction

Customer churn is defined as the rate when a customer leaves or stop paying for a product of service. Reduction customer churn is important because cost of acquiring a new customer is higher than retaining an existing one. The full cost of customer churn includes both lost revenue and the marketing costs involved with replacing those customers with new ones. Reducing customer churn is a key business goal of every business. Predicting and preventing customer churn represents a huge additional potential revenue source for every business.

1.1 Problem Statement

In this problem statement, we were provided with a train and test dataset of a telecom company. The data set consists of 20 variables describing various services and charges associated with them, duration and any service calls made by customer. The dataset also contains geographic location of customers in form of state. 'Churn' is the target variable, which tells us weather the customer has churned out or not.

1.2 Dataset

Two different datasets were provided as train data and test data. Data contains 20 predictor variables and 1 target variables. Train data had 3333 observations while test data consist of 1667. Total 21 variables were present in data. All the variables are described in table 1 below.

Variables	Description
State *	State to which customer belongs
Account length	Service usage period
Area code	Telephone area code
Phone number *	Customer's phone number
International plan *	'yes' if customer opted for international plan else 'no'
Voice mail plan *	'yes' if customer opted for voice mail plan else 'no'
Number vmail messages	Number of voice messages stored or received by customer
Total day minutes	Total minutes in day time usage
Total day calls	Total calls made in day time
Total day charges	Charges for services used during day time
Total eve minutes	Total minutes in evening time usage
Total eve calls	Total calls made in evening time
Total eve charges	Charges for services used during evening time
Total night minutes	Total minutes in night time usage
Total night calls	Total calls made in night time
Total night charges	Charges for services used during night time
Total intl minutes	Total international minutes used
Total inil calls	Total international calls made
Total intl charges	Charges for international calls
Number customer services call	Services call made by customer
Churn	Target - 'True.' If customer churned else 'False'

Table1. Variables description

Out of 20 variables, 04 were categorical and 16 were continuous. Categorical variables are marked as * in table. Sample data is shown below.

state	account length	area code	phone number	international plan	voice mail plan	number vmail messages	total day minutes	total day calls	total day charge	total eve minutes
KS	128	415	382-4657	no	yes	25	265.1	110	45.07	197.4
OH	107	415	371-7191	no	yes	26	161.6	123	27.47	195.5
NJ	137	415	358-1921	no	no	0	243.4	114	41.38	121.2
OH	84	408	375-9999	yes	no	0	299.4	71	50.9	61.9
OK	75	415	330-6626	yes	no	0	166.7	113	28.34	148.3
AL	118	510	391-8027	yes	no	0	223.4	98	37.98	220.6
MA	121	510	355-9993	no	yes	24	218.2	88	37.09	348.5
MO	147	415	329-9001	yes	no	0	157	79	26.69	103.1
LA	117	408	335-4719	no	no	0	184.5	97	31.37	351.6
WV	141	415	330-8173	yes	yes	37	258.6	84	43.96	222
IN	65	415	329-6603	no	no	0	129.1	137	21.95	228.5
RI	74	415	344-9403	no	no	0	187.7	127	31.91	163.4

total eve calls	total eve charge	total night minutes	total night calls	total night charge	total intl minutes	total intl calls	total intl charge	number customer service calls	Churn
99	16.78	244.7	91	11.01	10	3	2.7	1	False.
103	16.62	254.4	103	11.45	13.7	3	3.7	1	False.
110	10.3	162.6	104	7.32	12.2	5	3.29	0	False.
88	5.26	196.9	89	8.86	6.6	7	1.78	2	False.
122	12.61	186.9	121	8.41	10.1	3	2.73	3	False.
101	18.75	203.9	118	9.18	6.3	6	1.7	0	False.
108	29.62	212.6	118	9.57	7.5	7	2.03	3	False.
94	8.76	211.8	96	9.53	7.1	6	1.92	0	False.
80	29.89	215.8	90	9.71	8.7	4	2.35	1	False.
111	18.87	326.4	97	14.69	11.2	5	3.02	0	False.
83	19.42	208.8	111	9.4	12.7	6	3.43	4	True.
148	13.89	196	94	8.82	9.1	5	2.46	0	False.

Figure 1. Sample train data

Chapter 2

Methodology

Customer churn reduction is a business scenario in which a company is trying to retain a customer which is more likely to leave the services. For reducing churn rate, we need to identify which customers are most likely to churn and which are not. So Churn reduction is a classification problem.

The solution is divided into 3 parts.

1. Exploratory data analysis(EDA) was performed to explore the structure of data. Some of the basic assumptions were made about the data ie. Which variables are most likely causing churn. During exploration dataset was checked for missing values, multi collinearity and other model/algorithm specific assumptions.
2. After EDA, for learning two models were used, logistic regression and random forest. Some data pre-processing was done to prepare training data for learning model.
3. In the last part, performance tuning was done to increase the accuracy of models.

Both the algorithms and EDA were implemented in R and python. Both implementations were similar with little difference due to difference in learning algorithm implementation.

2.1 Exploratory Data Analysis (EDA)

Exploratory data analysis a.k.a. EDA was performed on training data using R and python. We looked at the structure of training data and found 20 predictors, 1 target variable and 3333 observations.

2.1.1 The Target Variable - Churn

The target variable was 'churn'. Initially if the customers churned out it flagged as 'True.', otherwise 'False.'. Later during it was changes as True = 1 and False = 0. Out of 3333 customers, 483 customers churned out and 2850 didn't churned out.

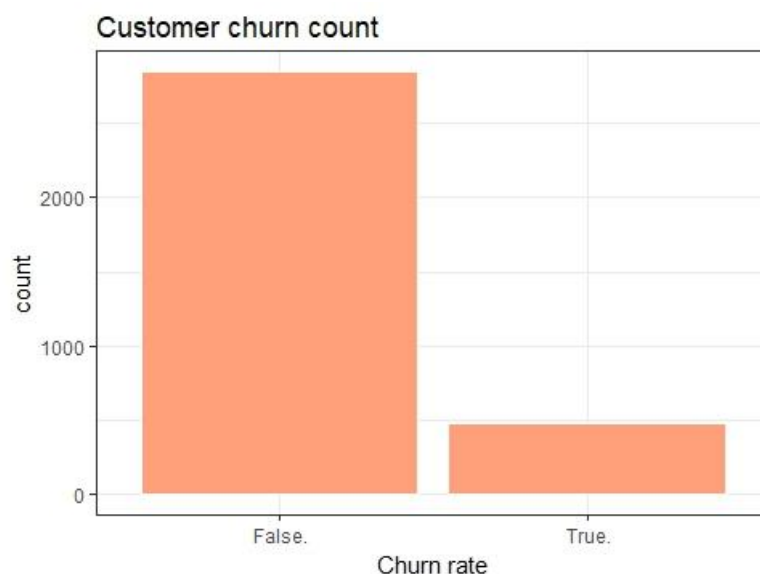


Fig2. Customer churn frequency

Since the churn event rate is approx. 14 % so it was concluded that it was not a highly imbalanced class. It won't affect our analysis significantly.

Also since our target variable is binary class, it is following the first assumption of logistic regression which is required dependent variable to be binary.

2.1.2 Missing Value Analysis

No missing values were present in the training and test dataset.

account length	0	international plan	0
voice_mail_plan	0	number_vmail_messages	0
total_day_minutes	0	total_day_calls	0
total_day_charge	0	total_eve_minutes	0
total_eve_calls	0	total_eve_charge	0
total_night_minutes	0	total_night_calls	0
total_night_charge	0	total_intl_minutes	0
total_intl_calls	0	total_intl_charge	0
number_customer_service_calls	0	churn	0

Table2. Missing value count

2.1.3 Multicollinearity

Multicollinearity exists whenever two or more of the predictors in a regression model are moderately or highly correlated. Multicollinearity is the condition when one predictor can be used to predict other. The basic problem is multicollinearity results in unstable estimation of coefficients which makes it difficult to access the effect of independent variable on dependent variable. Correlation plot was used in R and python to detect highly collinear variables.

From the correlation plot we can see that –

1. 'Total day minutes' and 'total day charges' are highly collinear
2. 'Total eve minutes' and 'total eve charges' are highly collinear
3. 'Total night minutes' and 'total night charges' are highly collinear
4. 'Total intl minutes' and 'total intl charges' are highly collinear

Random forest can handle multi collinear variables because of bagging approach. So we will not remove these variables in random forest implementation.

One of the assumptions of logistic regression is that logistic regression requires there to be little or no multicollinearity among the independent variables. This means that the independent variables should not be too highly correlated with each other. Due to this assumption, one the predictors from each set was removed when logistic learner was trained.

Multicollinearity matrix is visualized in figure 3.

Variables are highly collinear are highlighted with red colour with their corresponding score. Apart from correlation matrix, multicollinearity was again checked during logistic regression diagnostics using VIF (Variation inflation factor).

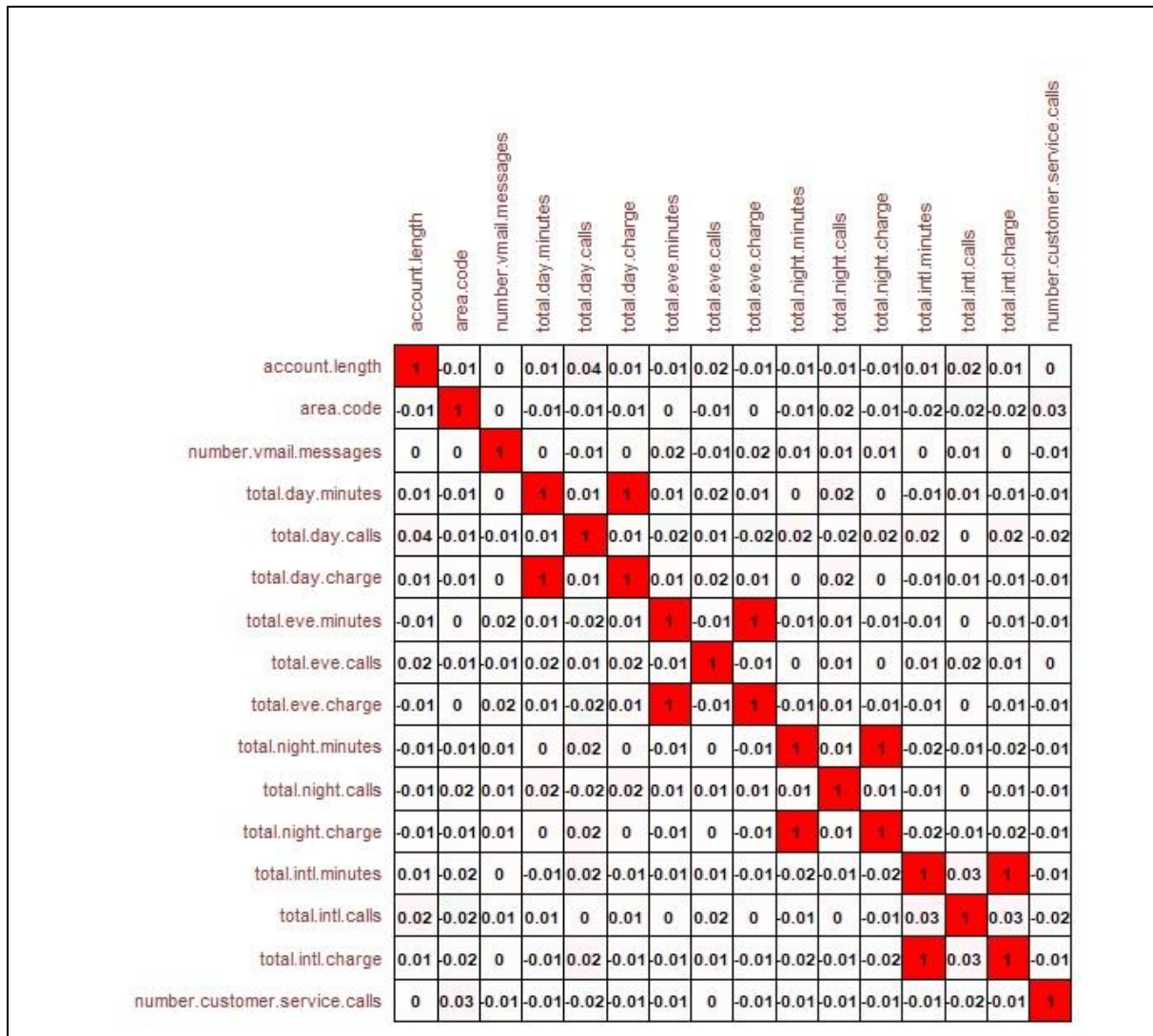


Fig 3. Multicollinearity matrix

2.1.4 Analysis of churn ratio with different predictors

The churn ratio with each predictors was analysed using boxplot and bar chart.

1. State

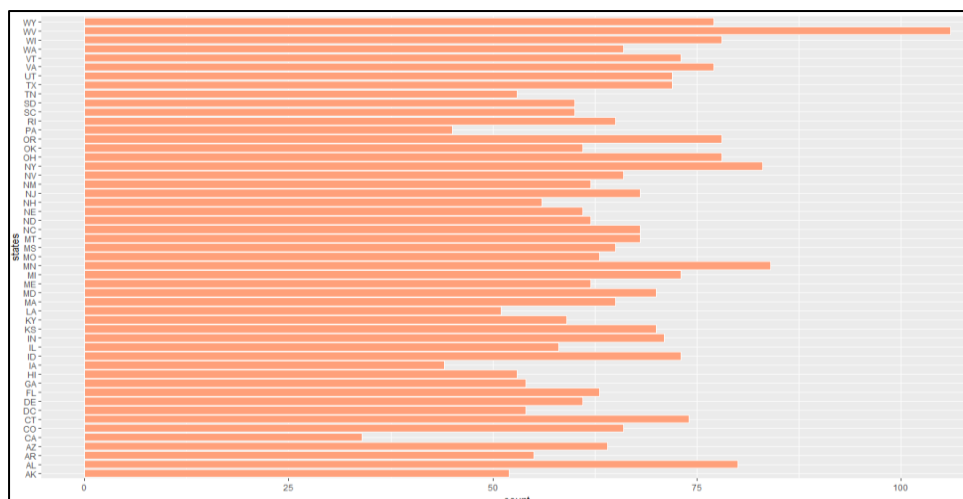


Fig. 4 State wise customer count

From the plot we can see that maximum customers are from west virginia and lowest are from California.

2. International plan

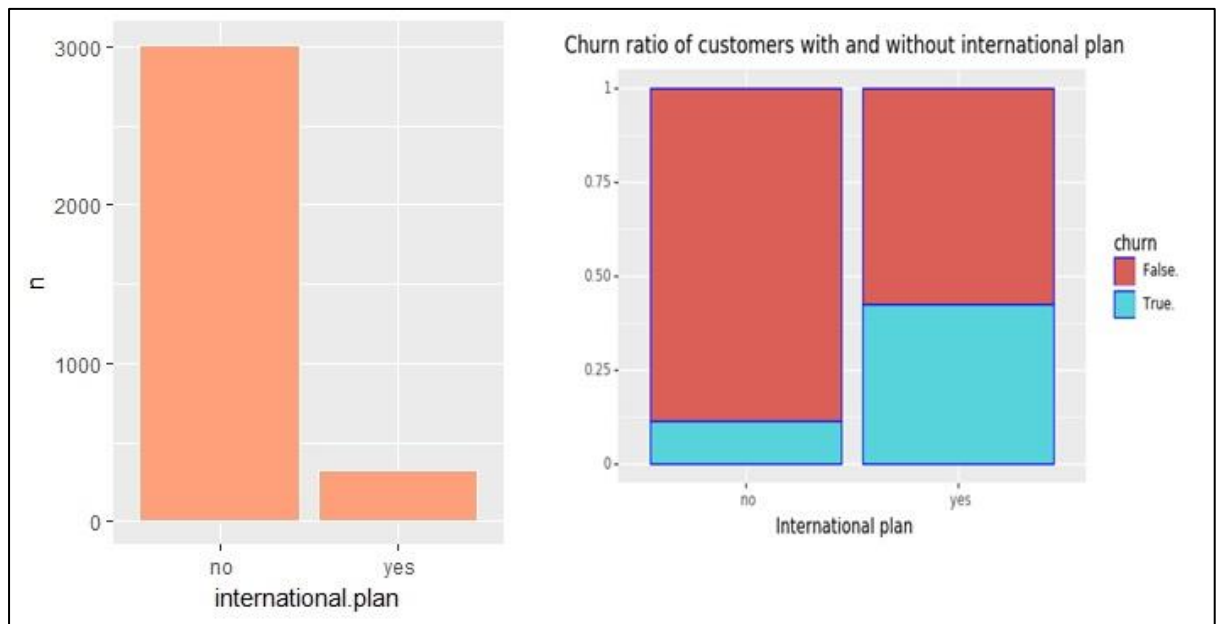


Fig 4. Churn ratio with and without international plan subscription

Most of the customers did not opt for the international plan subscription. Customer churn is more with customers with international plan. Further details can be examined by looking at percentage.

Churn	N	Percentage
False	2664	88.5
True	346	11.5

Table 3. Churn rate % for customers without international plan subscription

Churn	N	Percentage
False	186	57.6
True	137	42.4

Table 4. Churn rate % for customers with international plan subscription

As it is evident that only 11.5 % customer churned in without international plan category. 42.4 % customer churn out with international plan subscription. There may be some issue with international plan as customer churn rate is higher with international plan.

3. Voice mail plan

922 customers subscribed for voice mail plan, 2411 did not.

Churn	N	Percentage
False	842	91.3
True	80	8.68

Table 5. Churn rate % for customers with voice mail plan subscription

Churn	N	Percentage
False	2008	83.3
True	403	16.7

Table 6. Churn rate % for customers without voice mail plan subscription

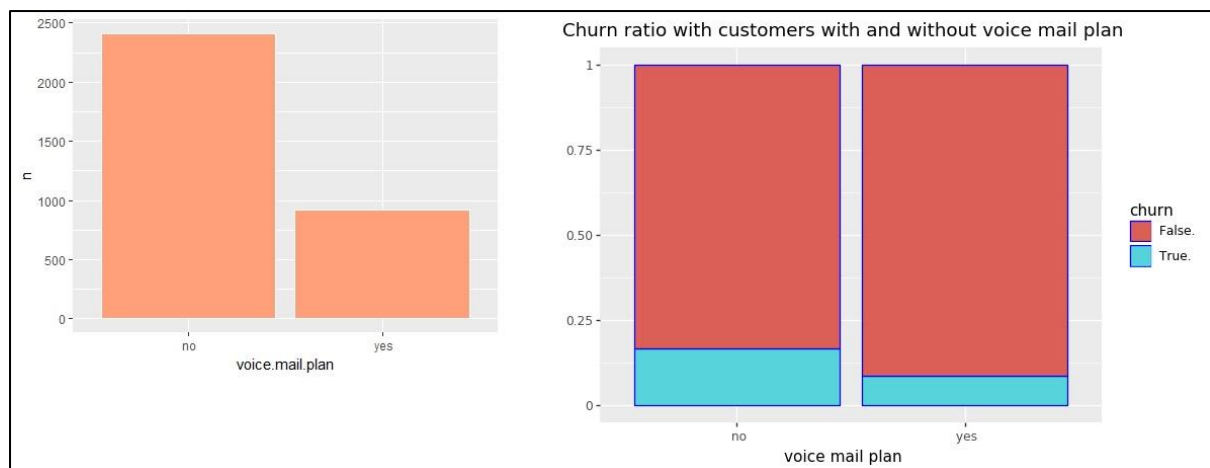


Fig 5. Churn ratio with and without voice mail plan subscription

922 customers have voice mail plan and 80 (8.68 %) customers out of 922 churned out. 2411 customers don't have voice mail plan and 403 (16.7 %) out of 2411 churn out. So customers without voice mail plan have higher churn rate.

4. number customer service calls

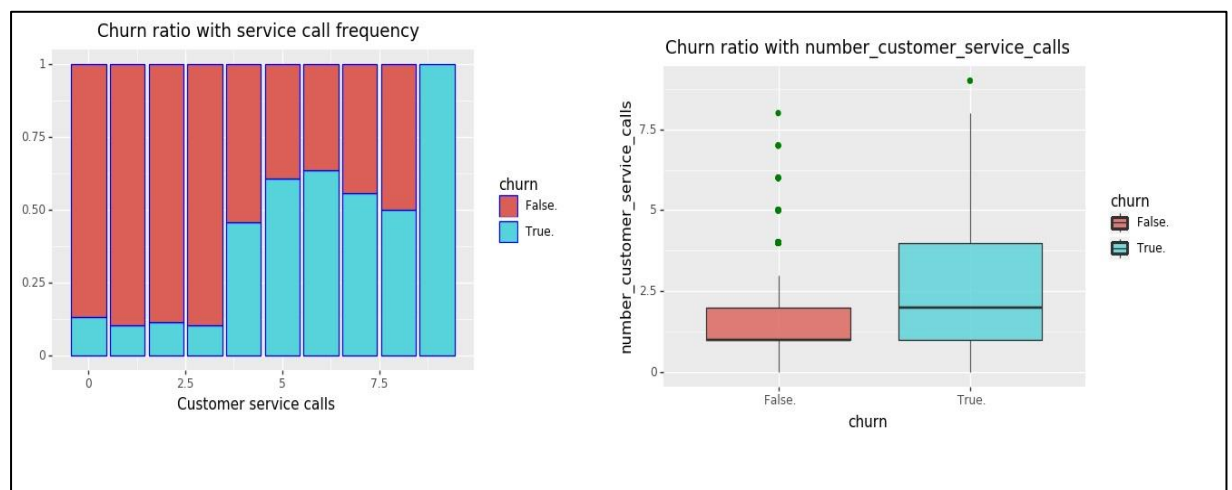


Fig. 6 Churn ratio with service call frequency

The churn rate is increasing with an increase in frequency of service calls.

5. Total day minutes, total day calls and total day charges

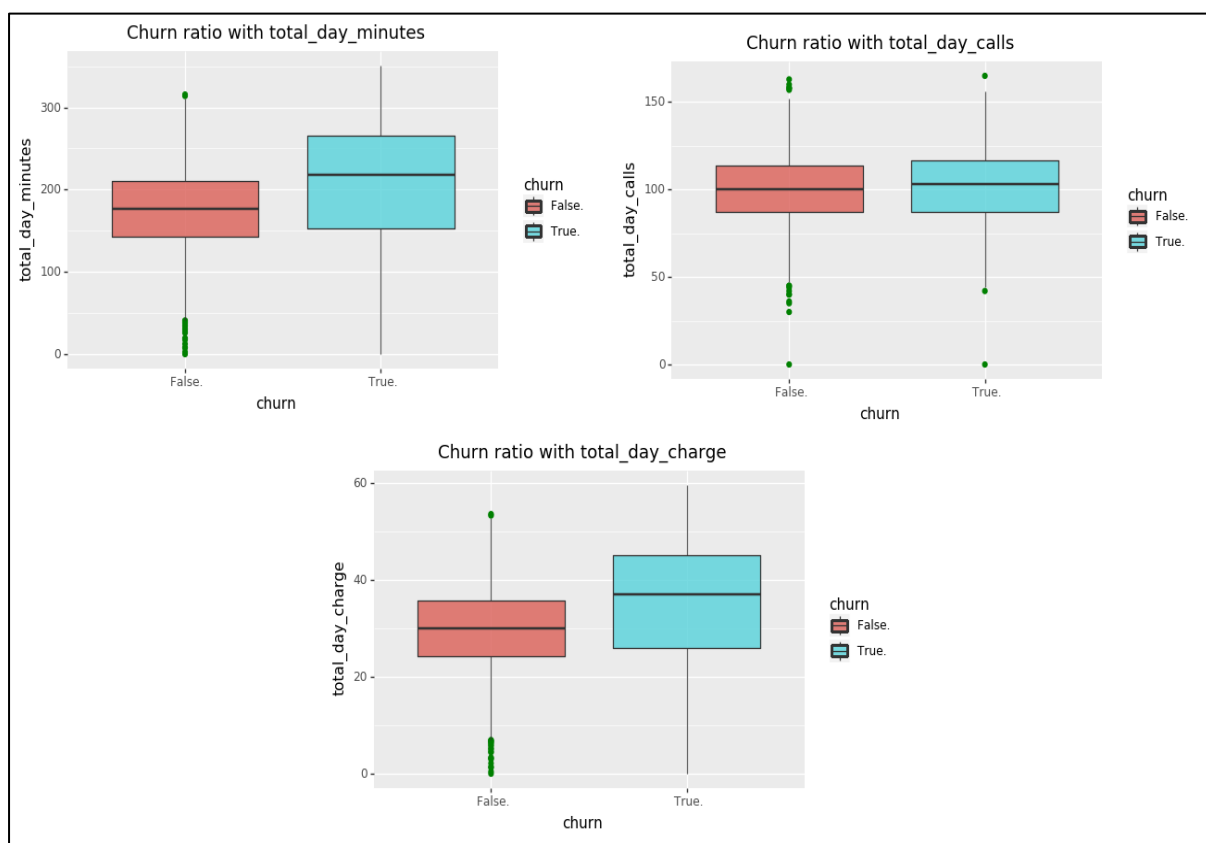


Fig 7. Total day minutes, total day calls and total day charge

	Total day minutes	Total day calls	Total day charge
Count	3333	3333	3333
Mean	179.775	100.435	30.562
Standard deviation	54.467	20.069	9.259
Min	0.0	0.0	0.0
25 %	143.70	87	24.43
50 %	179.40	101	30.50
75 %	216.40	114	36.79
Max	350.80	165	59.64

Table 7. summary statistics for Total day minutes, total day calls and total day charge

From the summary statistics, no anomalies are visible. Every value seems to be in standard range. From the plots we can say that total number of day calls are approximately same for across churn and non-churn customers. But day call duration and day call charges are slightly higher in case of customers who are churned out. Here 'total day minutes' and 'total day charges' are collinear.

6. Total eve minutes, total eve calls and total eve charge

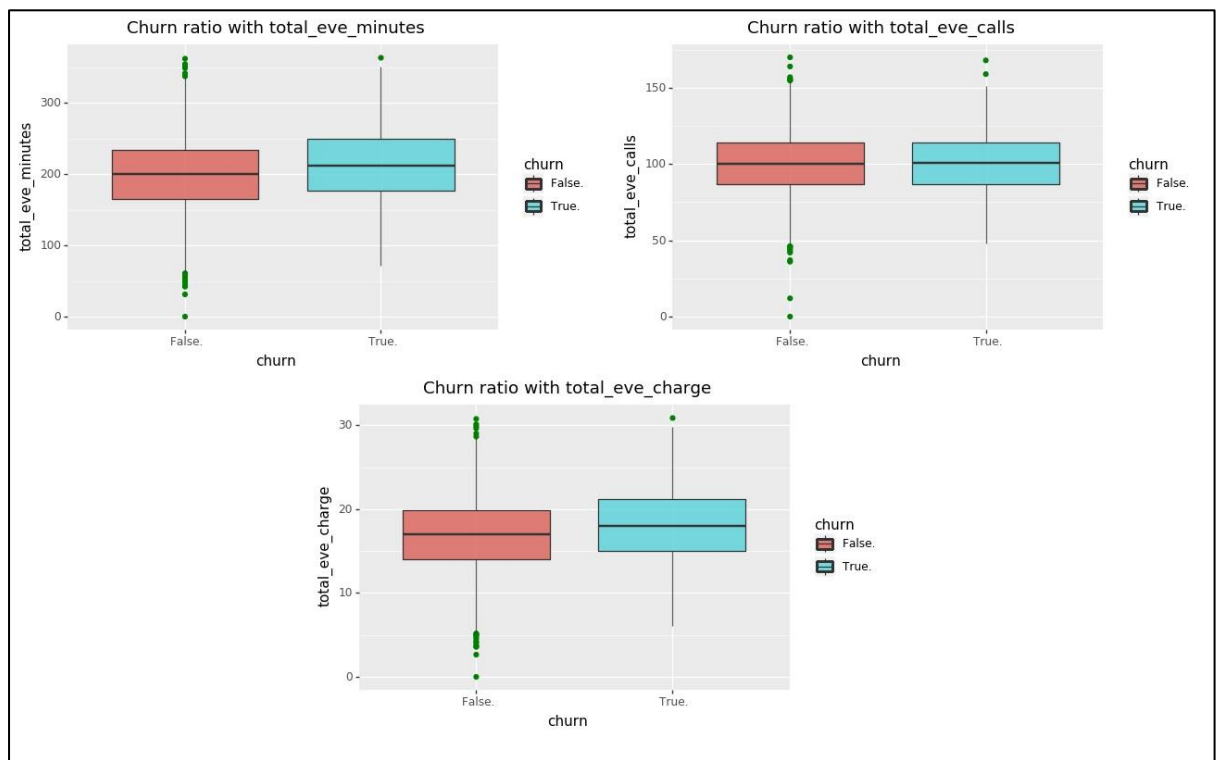


Fig 8. Total eve minutes, total eve calls and total eve charge

	Total eve minutes	Total eve calls	Total eve charge
Count	3333	3333	3333
Mean	200.98	100.114	17.083
Standard deviation	50.71	19.922	4.31
Min	0.0	0.0	0.0
25 %	166.60	87	14.16
50 %	201.40	100	17.12
75 %	235.30	114	20.00
Max	363.70	170	30.91

Table 8. summary statistics for Total eve minutes, total eve calls and total eve charge

No anomalies are visible in statistic summary. Every value seems to be in standard range. From the plots we can say that total number of eve calls are approximately same for across churn and non-churn customers. But eve call duration and eve call charge are slightly higher in case of customers who are churned out. Here 'total eve minutes' and 'total eve charge' are collinear.

7. Total night minutes, total night calls and total night charge

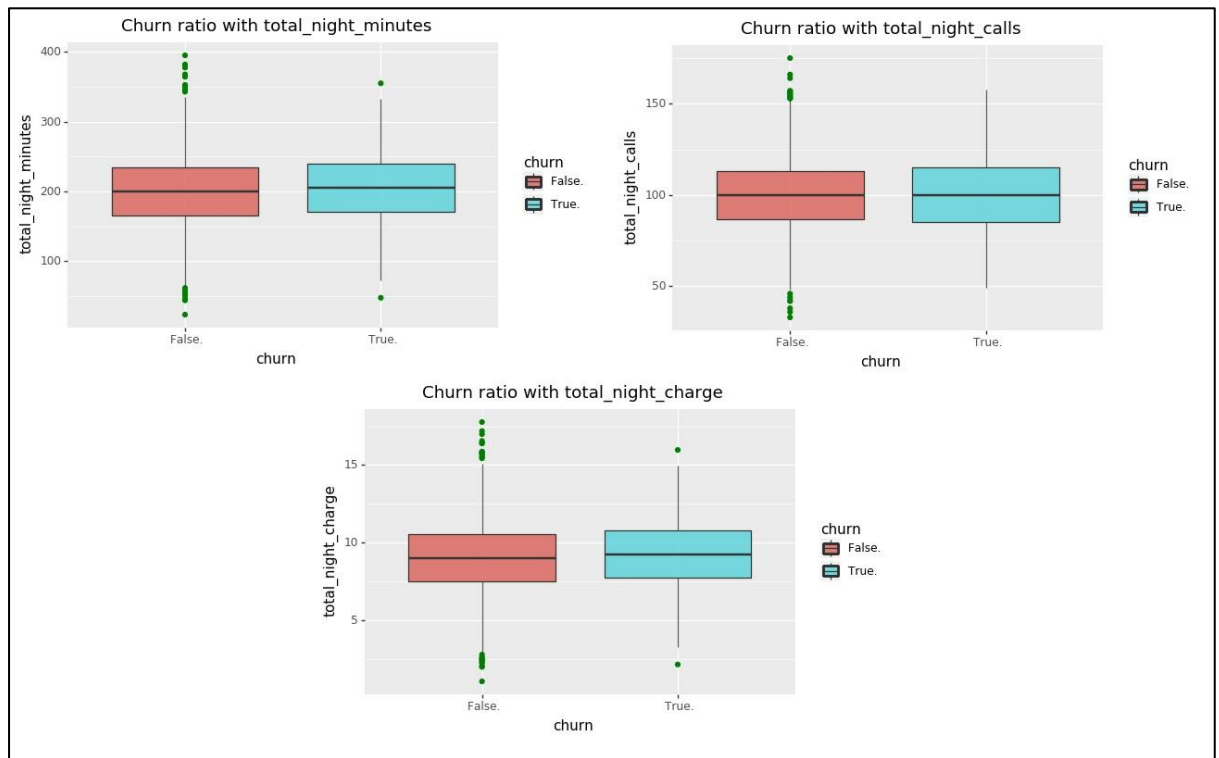


Fig 9. Total night minutes, total night calls and total night charge

	Total night minutes	Total night calls	Total night charge
Count	3333	3333	3333
Mean	200.87	100.107	9.039
Standard deviation	50.57	19.56	2.275
Min	23.20	33	1.04
25 %	167	87	7.52
50 %	201	100	9.05
75 %	235	113	10.59
Max	395	175	10.77

Table 8. summary statistics for Total night minutes, total night calls and total night charge

No anomalies are visible in statistic summary. Every value seems to be in standard range. From the plots we can say that total number of night calls, night minutes and night charges are approximately same for across churn and non-churn customers. 'total night minutes' and 'total night charge' are collinear.

8. Total intl minutes, total intl calls and total intl charges

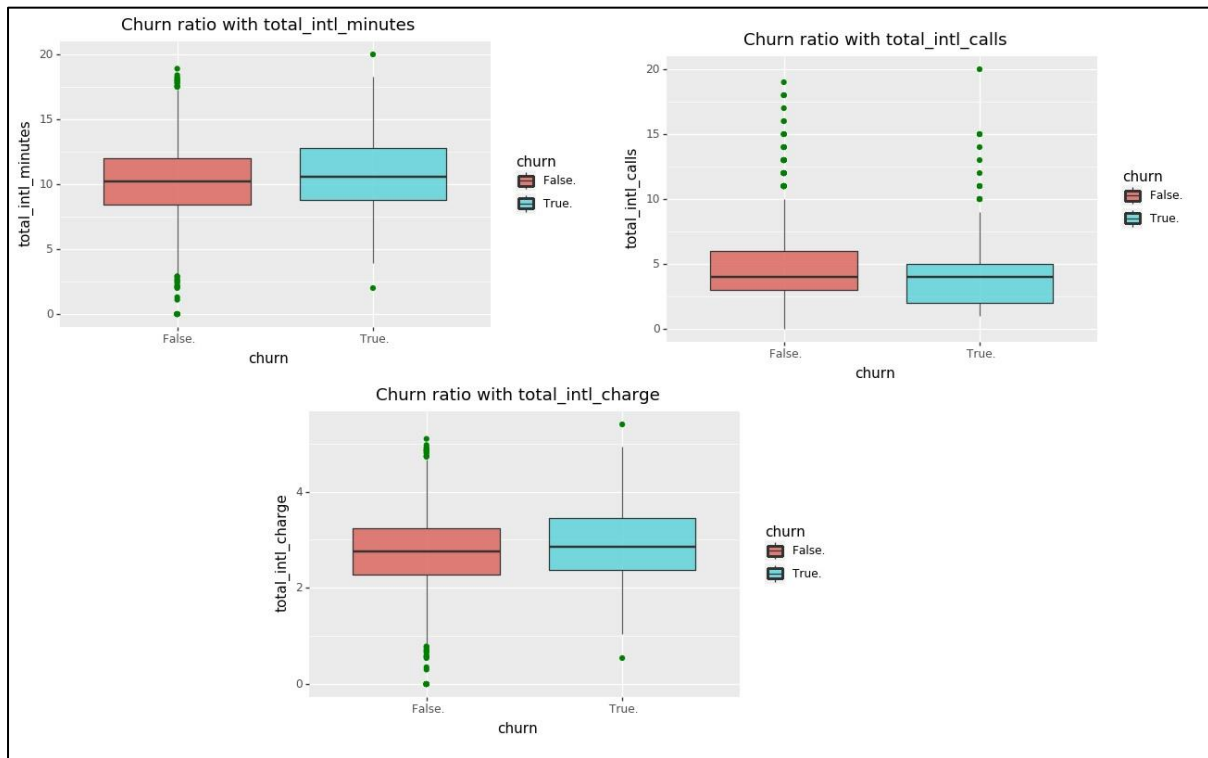


Fig 10. Total intl minutes, total intl calls and total intl charge

	Total intl minutes	Total intl calls	Total intl charge
Count	3333	3333	3333
Mean	10.23	4.47	2.764
Standard deviation	2.79	2.46	.753
Min	0	0	0
25 %	8.50	3	2.30
50 %	10.30	4	2.78
75 %	12.10	6	3.27
Max	20	20	5.40

Table 9. summary statistics for Total intl minutes, total intl calls and total intl charge

No anomalies are visible in statistic summary. Every value seems to be in standard range. From the plots we can say that total number of international calls, international minutes and international charges are approximately same for across churn and non-churn customers. 'total intl minutes' and 'total intl charge' are collinear.

9. Variable importance

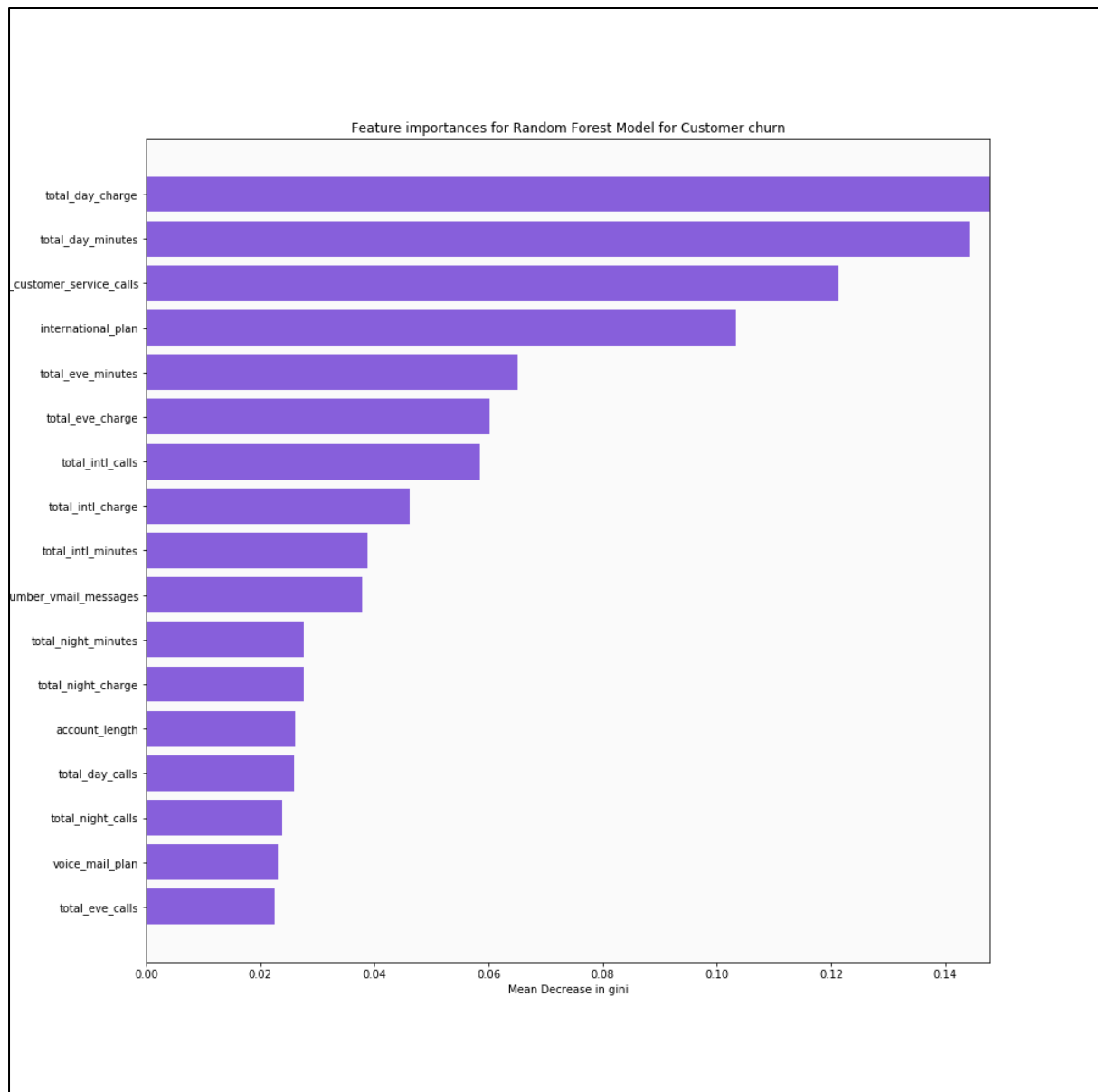


Figure 11. Variable importance

2.2 Modeling

Customer churn reduction is a binary classification problem. Here we have to build a model which can classify if a customer will churn out of not. In our dataset, dependent variable 'churn' was binary. So two machine learning algorithms were selected for learning.

1. Logistic regression, a member of generalized linear model family
2. Random forest, an ensemble tree based technique

Both training models logistic regression and random forest were implemented in R and python. After building an initial model, performance tuning was done using hyperparameter tuning for optimised parameters.

In R , 0/False was the positive class while in python implementation True/1 was positive class.

2.2.1 Logistic regression

First logistic algorithm was trained.

Logistic algorithm works on numeric data. So categorical data ie. International plan and voice mail plan and target predictor 'Churn' were encoded in to binary encoding.

In case if international plan and voice mail plan, yes was coded as 1 and no as 0. In case of Churn True was encoded as 1 and False as 0. True/1 means customer churned and False/0 means it stayed.

- Train data was divided into train dataset and validation set.
- Logistic regression models were trained on train dataset.
- Validation set and AIC score was used to select the best models out of all trained models.
- Final test and prediction was performed on test data which was provided separately.

R implementation

In R 4 different models were trained.

1. fullModel was trained with all the predictors.
2. Model2 was trained by removing collinear variables.
3. model_backward was trained using backward elimination approach.
4. A model with forward elimination was trained but it was equivalent to backward elimination approach.

fullModel

```
1. summary(fullModel)
2. Call:
3. glm(formula = Churn ~ ., family = binomial(link = "logit"), data = train_data)
4.
5. Deviance Residuals:
6.      Min       1Q   Median       3Q      Max
7. -2.1327  -0.5188  -0.3414  -0.1933   3.3265
8.
9. Coefficients:
10.              Estimate Std. Error z value Pr(>|z|)
11. (Intercept)    -8.836e+00  8.637e-01 -10.231  < 2e-16 ***
12. account.length    1.062e-04  1.658e-03   0.064  0.948912
13. international.plan1  1.950e+00  1.727e-01  11.295  < 2e-16 ***
14. voice.mail.plan1  -2.427e+00  7.056e-01  -3.439  0.000583 ***
15. number.vmail.messages  4.704e-02  2.213e-02   2.125  0.033562 *
```

```

16. total.day.minutes      2.582e+00  3.939e+00  0.655 0.512215
17. total.day.calls        4.543e-03  3.256e-03  1.395 0.162923
18. total.day.charge      -1.511e+01  2.317e+01 -0.652 0.514408
19. total.eve.minutes      6.128e-01  1.953e+00  0.314 0.753647
20. total.eve.calls        2.533e-03  3.310e-03  0.765 0.444154
21. total.eve.charge      -7.126e+00  2.297e+01 -0.310 0.756396
22. total.night.minutes   -1.731e-01  1.059e+00 -0.163 0.870202
23. total.night.calls      -8.406e-04  3.454e-03 -0.243 0.807717
24. total.night.charge     3.957e+00  2.354e+01  0.168 0.866526
25. total.intl.minutes    -4.260e+00  6.301e+00 -0.676 0.498984
26. total.intl.calls      -1.024e-01  3.027e-02 -3.382 0.000720 ***
27. total.intl.charge     1.607e+01  2.333e+01  0.689 0.490923
28. number.customer.service.calls  5.005e-01  4.820e-02 10.383 < 2e-16 ***
29. ---
30. Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
31.
32. (Dispersion parameter for binomial family taken to be 1)
33.
34. Null deviance: 1934.3 on 2333 degrees of freedom
35. Residual deviance: 1514.0 on 2316 degrees of freedom
36. AIC: 1550
37.
38. Number of Fisher Scoring iterations: 6

```

Model2

```

1. > summary(model2)
2.
3. Call:
4. glm(formula = Churn ~ . - total.day.minutes - total.eve.minutes -
5. total.night.minutes - total.intl.minutes, family = binomial(link = "logit"),
6. data = train_data)
7.
8. Deviance Residuals:
9. Min 1Q Median 3Q Max
10. -2.1280 -0.5209 -0.3397 -0.1929 3.3593
11.
12. Coefficients:
13. Estimate Std. Error z value Pr(>|z|)
14. (Intercept) -8.8607344 0.8636786 -10.259 < 2e-16 ***
15. account.length 0.0001257 0.0016568 0.076 0.939502
16. international.plan1 1.9437989 0.1723920 11.275 < 2e-16 ***
17. voice.mail.plan1 -2.4224099 0.7048692 -3.437 0.000589 ***
18. number.vmail.messages 0.0467668 0.0221136 2.115 0.034443 *
19. total.day.calls 0.0046369 0.0032528 1.426 0.154008
20. total.day.charge 0.0790371 0.0075897 10.414 < 2e-16 ***
21. total.eve.calls 0.0025136 0.0033060 0.760 0.447064
22. total.eve.charge 0.0833592 0.0159778 5.217 1.82e-07 ***
23. total.night.calls -0.0008455 0.0034513 -0.245 0.806482
24. total.night.charge 0.1098887 0.0294295 3.734 0.000188 ***
25. total.intl.calls -0.1020647 0.0301623 -3.384 0.000715 ***
26. total.intl.charge 0.2972813 0.0895991 3.318 0.000907 ***
27. number.customer.service.calls 0.4997406 0.0481675 10.375 < 2e-16 ***
28. ---
29. Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
30.
31. (Dispersion parameter for binomial family taken to be 1)
32.
33. Null deviance: 1934.3 on 2333 degrees of freedom
34. Residual deviance: 1515.0 on 2320 degrees of freedom
35. AIC: 1543
36.
37. Number of Fisher Scoring iterations: 6

```


If we compare fullModel and model2 summary, we can see that in full model due to variables with multi-collinear only **international plan, voice mail plan, number customer service call and total intl calls** are considered statistically significant with **AIC score 1550**.

After removing highly collinear variables in model2, variables **international plan, voice mail plan, number customer service call along with total day charge, total eve charge, total night charge, total intl charge, total intl call and total number vmail meassage** are considered statistically significant with **AIC score 1543**.

Model with lower AIC score is considered better. Next, backward and forward elimination method was used. Models for forward and backward method were same so, selecting backward elimination method.

```

1. > summary(model_backward)
2.
3. Call:
4. glm(formula = Churn ~ international.plan + voice.mail.plan +
5.     number.vmail.messages + total.day.minutes + total.eve.minutes +
6.     total.night.charge + total.intl.calls + total.intl.charge +
7.     number.customer.service.calls, family = binomial(link = "logit"),
8.     data = train_data)
9.
10. Deviance Residuals:
11.     Min       1Q   Median       3Q      Max
12. -2.1014  -0.5209  -0.3373  -0.1952   3.3215
13.
14. Coefficients:
15.              Estimate Std. Error z value Pr(>|z|)
16. (Intercept)    -8.211093   0.620068  -13.242  < 2e-16 ***
17. international.plan1    1.943798   0.171974   11.303  < 2e-16 ***
18. voice.mail.plan1    -2.439077   0.703021   -3.469  0.000522 ***
19. number.vmail.messages    0.047180   0.022048    2.140  0.032366 *
20. total.day.minutes    0.013504   0.001288   10.480  < 2e-16 ***
21. total.eve.minutes    0.007039   0.001357    5.186  2.14e-07 ***
22. total.night.charge    0.110219   0.029372    3.753  0.000175 ***
23. total.intl.calls    -0.102416   0.030138   -3.398  0.000678 ***
24. total.intl.charge    0.298313   0.089341    3.339  0.000841 ***
25. number.customer.service.calls    0.495978   0.047997   10.333  < 2e-16 ***
26. ---
27. Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
28.
29. (Dispersion parameter for binomial family taken to be 1)
30.
31.     Null deviance: 1934.3  on 2333  degrees of freedom
32. Residual deviance: 1517.7  on 2324  degrees of freedom
33. AIC: 1537.7
34.
35. Number of Fisher Scoring iterations: 6

```

We are rejecting fullModel due to multi-collinearity assumption of logistic regression. Selecting model_backward because its AIC score is lower than model2, fullModel and all the predictors in model_backward are statically significant.

Assumptions of logistic regressions were check to see if any violation had happened. Results are in Appendix.

Python implementation

In python, two models were trained.

1. classifier_logit_default was trained using default parameters.
2. classifier_logit_2 was trained by tuning 'C' with grid search. Scikit-learn does not support step wise regression.

classifier_logit_default

```
1. LogisticRegression(C=1.0, class_weight=None, dual=False, fit_intercept=True,
2.     intercept_scaling=1, max_iter=100, multi_class='ovr', n_jobs=1,
3.     penalty='l2', random_state=456, solver='liblinear', tol=0.0001,
4.     verbose=0, warm_start=False)
```

classifier_logit_2

```
1. LogisticRegression(C=3, class_weight='balanced', dual=False,
2.     fit_intercept=True, intercept_scaling=1, max_iter=100,
3.     multi_class='ovr', n_jobs=1, penalty='l2', random_state=None,
4.     solver='liblinear', tol=0.0001, verbose=0, warm_start=False)
```

Results of hyper parameter tuning –

1. Best parameter: {'C': 3} Best score :0.7711101585940848

Model classifier_logit_2 was selected by comparing its performance on validation set.

2.2.2 Random Forest

After linear regression, random forest was trained. It was implemented in both R and python.

In both implementations random forest was first trained with default setting and the hyper parameters tuning was used to find the best parameters.

R Implementation

- In R mlr library was used for training and performance tuning.
- k-fold was selected as validation strategy with 10 iterations.

Two random forest model were trained

1. modelRFBaseline - with default parameters
2. modelRF1 - with tuned parameters

summary of modelRF1 –

```
1. > modelRF1$learner.model
2.
3. Call:
```

```

4. randomForest(formula = f, data = data, classwt = classwt, cutoff = cutoff,      nt
   ree = 429L, importance = TRUE, mtry = 6L, nodesize = 18L)
5.           Type of random forest: classification
6.           Number of trees: 429
7. No. of variables tried at each split: 6
8.
9.           OOB estimate of  error rate: 5.66%
10. Confusion matrix:
11.      0   1  class.error
12. 0 1919  76  0.03809524
13. 1   56 283  0.16519174

```

ModelRF1 was select.

Different parameters were tuned to increase the performance. These parameters were –

	Type	len	Def	Constr	Req	Tunable	Trafo
1. mtry	integer	-	-	2 to 10	-	TRUE	-
3. nodesize	integer	-	-	10 to 50	-	TRUE	-
4. ntree	integer	-	-	100 to 600	-	TRUE	-

Python Implementation

In python random forest was trained and hyperparameters optimisation was done using following parameters.

Summary of fit_randomForest

```

1. RandomForestClassifier(bootstrap=True, class_weight=None, criterion='gini',
2.                        max_depth=10, max_features='auto', max_leaf_nodes=None,
3.                        min_impurity_decrease=0.0, min_impurity_split=None,
4.                        min_samples_leaf=1, min_samples_split=2,
5.                        min_weight_fraction_leaf=0.0, n_estimators=100, n_jobs=1,
6.                        oob_score=False, random_state=12345, verbose=0,
7.                        warm_start=False)

```

```

1. param_dist = {'max_depth': [2,4,6,8,10],
2.              'bootstrap': [True, False],
3.              'max_features': ['auto', 'sqrt', 'log2', None],
4.              'criterion': ["gini", "entropy"],
5.              "n_estimators" : [100 ,200 ,300 ,400 ,500]
6.              }

```

After hyper tuning best parameters were –

```

1. Best Parameters using random search:
2. {'n_estimators': 100, 'max_features': 'auto', 'max_depth': 10, 'criterion': 'gini'
   , 'bootstrap': False}
3. Time taken in random search: 247.70

```

Variable importance was plotted. It is present in EDA section.

Chapter 3

Result and performance measure

Test data set was pre-processed in the same way as the train data was processed.

3.1 Python Implementation Results

True/1 was positive class.

Logistic Regression

classifier_logit_2 was the selected model and predictions were made on test dataset.

Confusion matrix for test data using logistic regression.

Observed/ Predicted	False	True
False	1096	347
True	39	185

Table.10 Confusion matrix for test data using logistic regression.

Classification performance metric of logistic regression in python

	Precision	Recall	F1-score
False / 0	0.76	0.97	0.85
True / 1	0.83	0.35	0.49
Avg / total	0.78	0.77	0.74

Table.11 Performance matrix for test data using logistic regression.

Different performance measures for logistic regression –

```
1. population: 1667
2. P: 224
3. N: 1443
4. PositiveTest: 532
5. NegativeTest: 1135
6. TP: 185
7. TN: 1096
8. FP: 347
9. FN: 39
10. TPR: 0.825892857143
11. TNR: 0.759528759529
12. PPV: 0.347744360902
13. NPV: 0.96563876652
14. FPR: 0.240471240471
15. FDR: 0.652255639098
16. FNR: 0.174107142857
17. ACC: 0.768446310738
18. F1_score: 0.489417989418
19. MCC: 0.428323776007
20. informedness: 0.585421616672
21. markedness: 0.313383127422
22. prevalence: 0.134373125375
23. LRP: 3.43447663648
24. LRN: 0.22923048097
25. DOR: 14.9826350403
26. FOR: 0.0343612334802
```

Mean accuracy using logistic regression was

1. Mean accuracy on test set 0.768446310738

Test error rate using logistic regression

1. Test error rate on test set 0.143

Random Forest

fit_randomForest was used for prediction on test dataset.

Confusion matrix for test data using random Forest in python.

Observed/ Predicted	False	True
False	1433	10
True	67	157

Table.12 Confusion matrix for test data using random forest

Classification performance metric of random forest in python

	Precision	Recall	F1-score
False / 0	0.99	0.96	0.97
True / 1	0.70	0.94	0.80
Avg / total	0.96	0.95	0.96

Table.13 Performance matrix for test data using random forest

Different performance measures for random forest –

1. population: 1667
2. P: 224
3. N: 1443
4. PositiveTest: 167
5. NegativeTest: 1500
6. TP: 157
7. TN: 1433
8. FP: 10
9. FN: 67
10. TPR: 0.700892857143
11. TNR: 0.99306999307
12. PPV: 0.940119760479
13. NPV: 0.955333333333
14. FPR: 0.00693000693001
15. FDR: 0.059880239521
16. FNR: 0.299107142857
17. ACC: 0.953809238152
18. F1_score: 0.803069053708
19. MCC: 0.788296379044
20. informedness: 0.693962850213
21. markedness: 0.895453093812
22. prevalence: 0.134373125375
23. LRP: 101.138839286

```
24. LRN: 0.301194422291
25. DOR: 335.792537313
26. FOR: 0.0446666666667
```

Mean accuracy using random forest –

1. Mean accuracy on test set 0.953809238152

Test error rate using logistic regression

1. Test error rate on test set 0.059

Result –

If we compare the results of logistic regression and random forest, random forest outperforms logistic regression in overall accuracy. Here, precision and recall are very important performance measures.

Recall - When a customer churns, how often does our classifier predict that correctly. For churned customers in logistic regression it is .35 while in random forest it is .94

Precision - When our classifier predicts a customer will churn, how often does that customer actually churn. For churned customer's Logistic regression's precision is .83 and random forest forest's is .70

But F1 score of random forest is better than logistic algorithm.

As we are trying to find which customers are going to churn out, **true positive rate** is also important metric. For logistic regression, it is .82 while for random forest it is .70. Here logistic regression is better than random forest.

Another important factor for consideration is that, these classifications are based on the probabilities of .5 for each target case. Probability threshold or cut-off is a business context decision. The threshold is usually set 0.5 by default. This means that anyone with a probability of more than .5 is predicted to churn. If we reduce the probability threshold, more people will be predicted to churn. If we are interested in customers which have most likely to churn, we can increase the probability threshold.

So the threshold or probability limit is decided based on business decision/requirement, marketing budget etc.

After comparing all the required performance measure, random forest was performing overall better than logistic regression.

3.2 R Implementations results

For comparing performance of logistic regression and random forest in R, probability threshold .7 and 0.3 were used. These cut off values were estimated while implementing random forest trying with different combinations.

Here positive class is False/0. Since False/0 is the positive class, **specificity** is a very important metric here. Specificity is the proportion of actual negatives that are correctly identified and our negative class is the customers who actually churned.

Logistic Regression

Result and summary of logistic regression

```
1. Confusion Matrix and Statistics
2.
3.      predicted
4. observed    0    1
5.      0 1284  159
6.      1  109  115
7.
8.      Accuracy : 0.8392
9.      95% CI : (0.8207, 0.8566)
10.     No Information Rate : 0.8356
11.     P-Value [Acc > NIR] : 0.360535
12.
13.      Kappa : 0.3685
14.  McNemar's Test P-Value : 0.002761
15.
16.      Sensitivity : 0.9218
17.      Specificity : 0.4197
18.      Pos Pred Value : 0.8898
19.      Neg Pred Value : 0.5134
20.      Prevalence : 0.8356
21.      Detection Rate : 0.7702
22.      Detection Prevalence : 0.8656
23.      Balanced Accuracy : 0.6707
24.
25.      'Positive' Class : 0
```

Confusion matrix

Observed/ Predicted	False	True
False	1284	159
True	109	115

Table.14 Confusion matrix of logistic regression

Random Forest

Result and summary of random forest

```
1. Confusion Matrix and Statistics
2.
3.      predicted
4. observed    0    1
5.      0 1404  39
6.      1  43  181
7.
8.      Accuracy : 0.9508
9.      95% CI : (0.9393, 0.9607)
```

```

10.    No Information Rate : 0.868
11.    P-Value [Acc > NIR] : <2e-16
12.
13.                Kappa : 0.7869
14.    McNemar's Test P-Value : 0.7404
15.
16.                Sensitivity : 0.9703
17.                Specificity : 0.8227
18.                Pos Pred Value : 0.9730
19.                Neg Pred Value : 0.8080
20.                Prevalence : 0.8680
21.                Detection Rate : 0.8422
22.    Detection Prevalence : 0.8656
23.    Balanced Accuracy : 0.8965
24.
25.    'Positive' Class : 0

```

Confusion matrix of prediction done using random forest

Observed/ Predicted	False	True
False	1404	39
True	43	181

Table.15 Confusion matrix for random forest

Results-

Comparing results of logistic regression –

Metric / Algorithm	Logistic Regression	Random Forest
Accuracy	.83	.95
Sensitivity	.92	.97
Specificity	.41	.82
Pos Pred Value	.88	.97
Neg Pred Value	.51	.80
Precision	.88	.97
Recall	.92	.97
F1	.90	.97
Prevalence	.83	.86
Detection Rate	.77	.84
Detection Prevalence	.86	.86
Balanced Accuracy	.67	.89

Table.16 Comparing Logistic regression and random forest performance

As we can see that random forest is outperforming logistic regression in every performance metric. We were interested in classifying which customers were going to churned. Here Specificity tells us how well our classifier is predicting, which customers are churning out. Logistic regression's specificity is .41 while random forest's is .82.

Random forest's accuracy is .95 which is way better than logistic regression. Also from the confusion matrix, we can see that false positive rate is very less for random forest.

So random forest is a better classifier for this problem statement.

Appendix A

Logistic regression diagnostics

```
1. #Logistic regression diagnostics - now we are going to check for logistic regression assumptions
2.
3. #1. Linearity assumption
4. # Checking for a linear relationship between continuous variables and logit of the outcome. This
5. # can be done by visually inspecting the scatter plot between each predictor and logit values.
6.
7. # selecting continuous predictors
8. continuous_data <- test_dataset %>%
9.   select_if(is.numeric)
10. predictors <- colnames(continuous_data)
11.
12. # calculating logit and adding it in 'continuous_data'
13. continuous_data <- continuous_data %>%
14.   mutate(logit = log(pred_backmodel/(1-pred_backmodel))) %>%
15.   gather(key = "predictors", value = "predictor.value", -logit)
16. ggplot(data = continuous_data, mapping = aes(x = logit, predictor.value)) +
17.   geom_point(size = 0.5, alpha = 0.5) +
18.   geom_smooth(method = "loess") +
19.   theme_bw() +
20.   facet_wrap(~predictors, scales = "free_y")
21.
22. # The smoothed scatter plots show all the continuous predictors are very near
23. # linearly associated with the outcome in logit scale in figure 11.
```

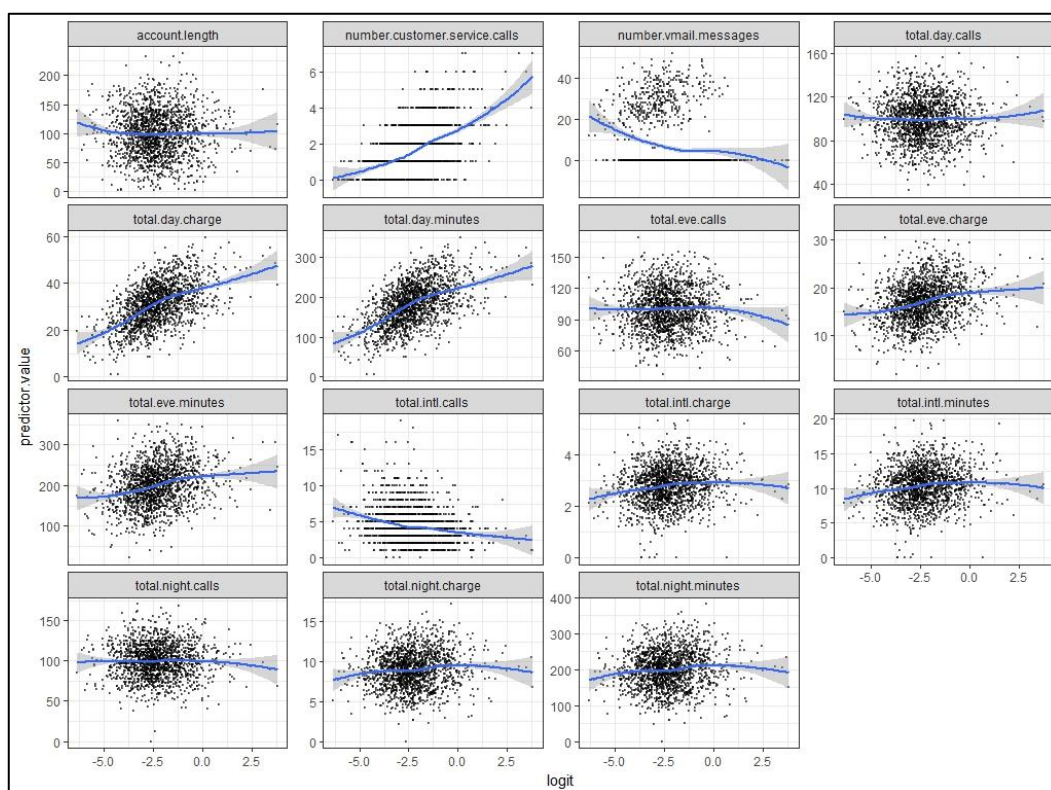


Figure 11 – Checking for linearity assumption

```

1. # 2. Multicollinearity
2. # we have removed multi-collinear variables while models.
3. #performing double check using variation inflation factor (VIF)
4. car::vif(model_backward)
5. #As a rule of thumb, a VIF value that exceeds 5 or 10 indicates a problematic amount of
6. #collinearity. In our model voice.mail.plan and number.vmail.messages is showing vif score of
7. # < 16. SO, we will remove one of the variable and retrain and check the final model after
8. #dignostics.
9. > car::vif(model_backward)
10.      international.plan      voice.mail.plan      number.vmail.messages
11.      1.050493      16.255768      16.219928
12.      total.day.minutes      total.eve.minutes      total.night.charge
13.      1.045576      1.031209      1.008662
14.      total.intl.calls      total.intl.charge      number.customer.service.calls
15.      1.012671      1.012673      1.0076539
16. #3. Influential values
17. # top 3 largest values
18. plot(model_backward, which = 4, id.n = 3)
19. #278, 1862, 3292

```

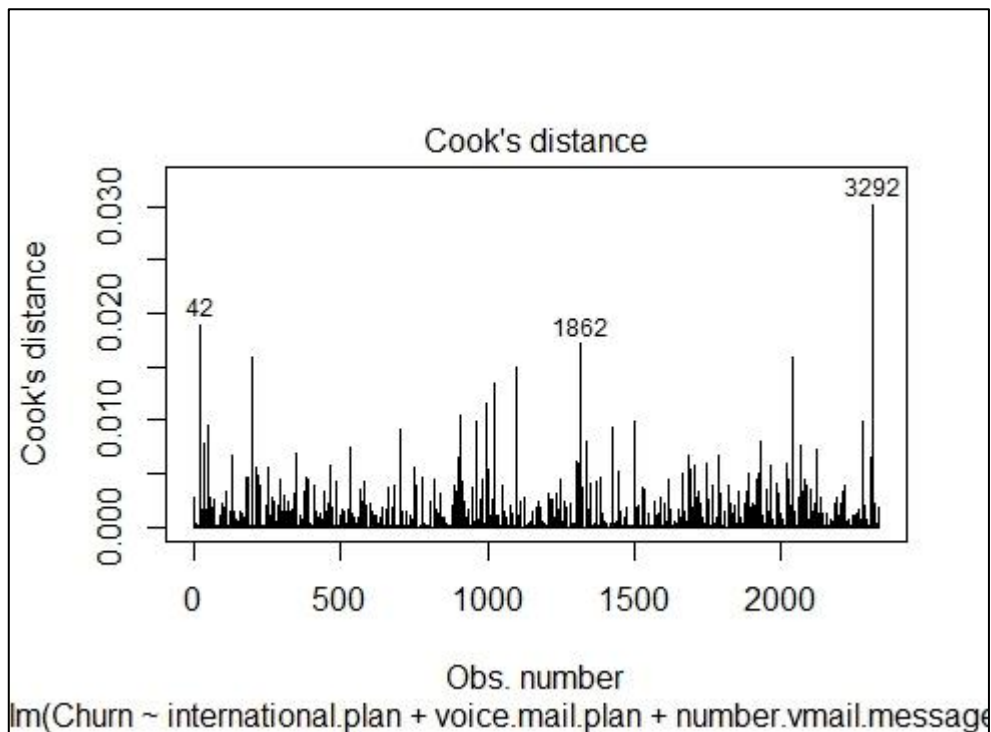


Fig12. Identifying influential & extreme points

```

1. # Data points with an absolute standardized residuals above 3 represent possible outliers
2. #and may deserve closer attention
3. #augment - add columns to the original dataset such as predictions, residuals and cluster assignments

```

```

4.
5. model.data <- augment(model_backward) %>%
6.   mutate(index = 1:n())
7. model.data %>% top_n(3, .cooks_d)
8. # plotting standardized residuals
9. ggplot(model.data, aes(index, .std.resid)) +
10.   geom_point(aes(color = Churn), alpha = .5) +
11.   theme_bw()
12. #
13. model.data %>%
14.   filter(abs(.std.resid) > 3)
15. # one possible influential observations was found on our training set.
16. # row 1890

```

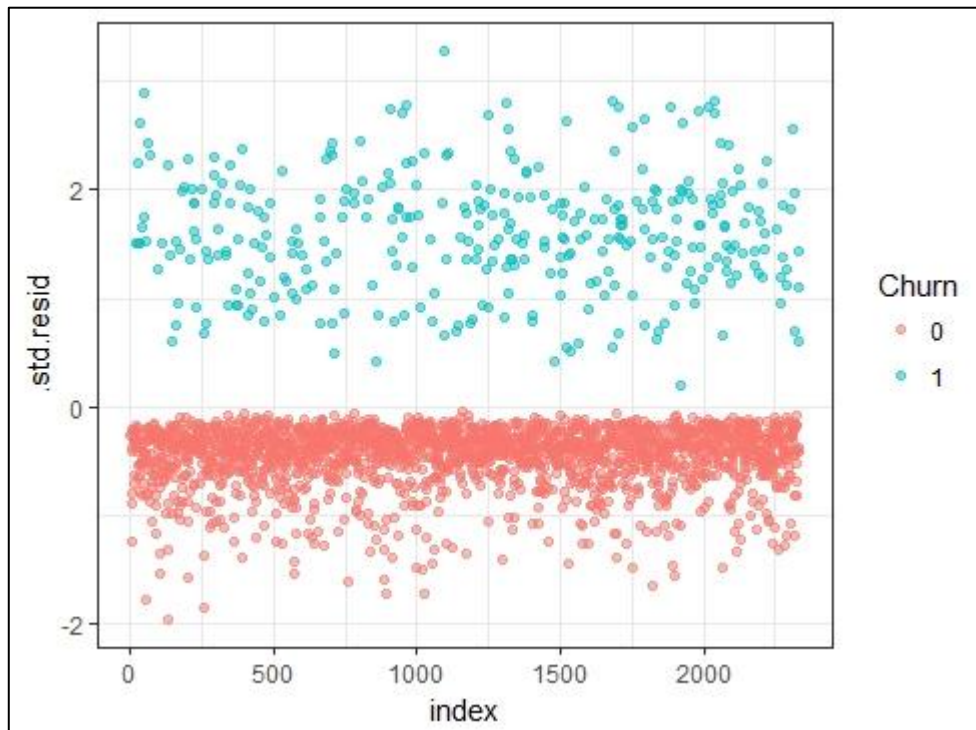


Fig13. standardized residuals

R code

ChurnReductionEDA.R

```

1. # Customer churn reduction EDA
2.
3. rm(list=ls())
4. #loading requiried libraries
5.
6. library(dplyr)
7. library(ggplot2)
8. library(stringr)
9. library(corrplot)
10.
11.
12. fillColor = "#FFA07A"
13. fillColorRed = "#56B4E9"
14.
15. train_data <-
16.   read.csv("Train_data.csv",
17.           sep = ',',
18.           na.strings = c(" ", "NA"))

```

```

19.
20. #looking at dimensions
21. dim(train_data)
22.
23. #Train data set consist of 3333 observations and 21 variabiles
24.
25. #checking structure of dataset
26. str(train_data)
27.
28. # Visualizing target class frequencies
29. train_data %>%
30.   count(Churn) %>%
31.   ggplot(aes(x = Churn,
32.             y = n)) +
33.     geom_bar(stat = 'identity',
34.             colour = "white",
35.             fill = fillColor) +
36.     labs(x = 'Churn rate', y = 'count ', title = 'Customer churn count') +
37.     theme_bw()
38.
39. table(train_data$Churn)
40. #Looking at the frequencies of churn , it is not looking like highly imbalance prob
    lem.
41.
42. # Now looking for any missing values
43. sapply(train_data, function(x) {
44.   sum(is.na(x))
45. }) # There are no missing values in dataset
46.
47.
48.
49. #selecting numeric variables
50. numCols <- unlist(lapply(train_data,is.numeric))
51. numVarDataset <- train_data[,numCols]
52.
53. # Visualizing correlation
54. par(mfrow = c(1, 1))
55. corr <- cor(numVarDataset)
56. corrplot(
57.   corr,
58.   method = "color",
59.   outline = TRUE,
60.   cl.pos = 'n',
61.   rect.col = "black",
62.   tl.col = "indianred4",
63.   addCoef.col = "black",
64.   number.digits = 2,
65.   number.cex = 0.60,
66.   tl.cex = 0.70,
67.   cl.cex = 1,
68.   col = colorRampPalette(c("green4", "white", "red"))(100)
69. )
70.
71. # From corrplot we can see that dataset consist of multicollinearity
72. # total.day.minutes and total.day.charge are highly collinear
73. # total.eve.minutes and total.eve.charge are highly collinear
74. # total.night.minutes and total.night.charge are highly collinear
75. # total.intl.minutes and total.intl are highly collinear
76. # we can exclude one of these predictors later during modeling
77.
78. ##### Generic EDA function for continous variables
79. plot_continuous <- function(dataset, variable,targetVariable) {
80.   var_name = eval(substitute(variable), eval(dataset))
81.   target_var = eval(substitute(targetVariable), eval(dataset))
82.   par(mfrow = c(1, 2))
83.   print(summary(var_name))

```

```

84. print(summary(target_var))
85. possible_outliers <- (boxplot.stats(var_name)$out)
86. print(possible_outliers)
87. print(paste("Total possible outliers", length(possible_outliers)))
88. table(possible_outliers)
89. ggplot(train_data, aes(target_var, var_name, fill = target_var)) +
90.   geom_boxplot(alpha = 0.8) + theme(legend.position = "null")
91. }
92.
93.
94. ##### looking at 'state' variable. It is a factor variable
95. train_data %>%
96.   count(state) %>%
97.   ggplot(mapping = aes(x = state, y = n)) +
98.   geom_bar(stat = 'identity',
99.     colour = 'white',
100.     fill = fillColor) +
101.     labs(x = "states", y = "count", "Customers per state") +
102.     coord_flip()
103.   # From the plot we can see that maximum customers are from west virginia and lowe
   st are from California
104.
105.
106.
107.
108.   # looking at each variable
109.   plot_continuous(train_data, account.length, Churn)
110.   # As we can see, that there are some possible outliers but they are not very
   extreme. Ignoring them
111.
112.
113.   ##### analysing international.plan #####
   #####
114.   str(train_data$international.plan) # it is a categorical variable
115.   table(train_data$international.plan)
116.   train_data %>%
117.     count(international.plan) %>%
118.     ggplot(mapping = aes(x = international.plan, y = n)) +
119.     geom_bar(stat = 'identity',
120.       colour = 'white',
121.       fill = fillColor)
122.   # From the plot we can see that most customers don't have international plan.
123.
124.   # next examining for the churn rate percentage of customers with national and
   international plan
125.
126.   national_cust_churnRate <- train_data %>%
127.     select(international.plan, Churn) %>%
128.     filter(str_detect(international.plan, "no")) %>%
129.     group_by(Churn) %>%
130.     summarise(n = n()) %>%
131.     mutate(percentage = (n / sum(n)) * 100)
132.   #Only 11.49 % customer with national plan churn out.
133.
134.
135.   international_cust_churnRate <- train_data %>%
136.     select(international.plan, Churn) %>%
137.     filter(str_detect(international.plan, "yes")) %>%
138.     group_by(Churn) %>%
139.     summarise(n = n()) %>%
140.     mutate(percentage = (n / sum(n)) * 100)
141.   # 42.42 % customers with international plan had churn out. It means that the
   telecom company
142.   # is mainly losing customers with international plans.

```

```

143.
144. #####Analysing voice.mail.plan #####
#####
145. table(train_data$voice.mail.plan)
146.
147. # customers with voice plan and their churn rate
148. voice_plan_churnRate <- train_data %>%
149.   select(voice.mail.plan, Churn) %>%
150.   filter(str_detect(voice.mail.plan, "yes")) %>%
151.   group_by(Churn) %>%
152.   summarise(n = n()) %>%
153.   mutate(churnRatePercentage = (n / sum(n)) * 100)
154.
155. ggplot(data = voice_plan_churnRate,
156.   mapping = aes(x = Churn, y = churnRatePercentage)) +
157.   geom_bar(stat = 'identity',
158.     colour = 'white',
159.     fill = fillColorRed) +
160.   labs(title = 'Voice main plan customers churn rate')
161.
162. # 922 customers have voice mail plan and 80 (8.68 %) customers out of 922 c
hurn out.
163.
164. #customers without voice plan and their churn rate
165. non_voice_plan_churnRate <- train_data %>%
166.   select(voice.mail.plan, Churn) %>%
167.   filter(str_detect(voice.mail.plan, "no")) %>%
168.   group_by(Churn) %>%
169.   summarise(n = n()) %>%
170.   mutate(churnRatePercentage = (n / sum(n)) * 100)
171.
172. ggplot(data = non_voice_plan_churnRate,
173.   mapping = aes(x = Churn, y = churnRatePercentage)) +
174.   geom_bar(stat = 'identity',
175.     colour = 'white',
176.     fill = fillColor) +
177.   labs(title = 'Non voice plan Customer churn rate')
178.
179. # 2411 customers dont have voice mail plan and 403 (16.7 %) out of 2411 chu
rn out
180.
181. #So customers without voice plan have higher churn rate
182.
183.
184.
185. # removing parameters that dosn't seem to be logical parameter for customer
churn.
186. #So removing state, area code and phone number
187. train_data$state <- NULL
188. train_data$area.code <- NULL
189. train_data$phone.number <- NULL
190.
191.
192. ##### Analysing number.vmail.messages #####
#####
193. str(train_data$number.vmail.messages)
194. plot_continous(train_data, number.vmail.messages,Churn)
195.
196. # no extreme outliers detected.
197.
198. ##### Analysing total.day.minutes #####
#####
199. str(train_data$total.day.minutes)
200. plot_continous(train_data, total.day.minutes,Churn)
201.
202. # no extreme outliers detected

```

```

203.
204. ##### Analysing total.day.calls #####
#####
205. str(train_data$total.day.calls)
206. plot_continuous(train_data, total.day.calls, Churn)
207.
208. # no extreme outliers detected
209.
210. ##### Analysing total.day.charge #####
#####
211. str(train_data$total.day.charge)
212. plot_continuous(train_data, total.day.charge, Churn)
213.
214. # no extreme outliers detected
215.
216. ##### Analysing total.eve.minutes #####
#####
217. str(train_data$total.eve.minutes)
218. plot_continuous(train_data, total.eve.minutes, Churn)
219.
220. # no extreme outliers detected
221.
222. ##### Analysing total.eve.calls #####
#####
223. str(train_data$total.eve.calls)
224. plot_continuous(train_data, total.eve.calls, Churn)
225.
226. # no extreme outliers detected
227.
228. ##### Analysing total.eve.charge #####
#####
229. str(train_data$total.eve.charge)
230. plot_continuous(train_data, total.eve.charge, Churn)
231.
232. # no extreme outliers detected
233.
234. ##### Analysing total.night.minutes #####
#####
235. str(train_data$total.night.minutes)
236. plot_continuous(train_data, total.night.minutes, Churn)
237.
238. # no extreme outliers detected
239.
240. ##### Analysing total.night.calls #####
#####
241. str(train_data$total.night.calls)
242. plot_continuous(train_data, total.night.calls, Churn)
243.
244. # no extreme outliers detected
245.
246. ##### Analysing total.night.charge #####
#####
247. str(train_data$total.night.charge)
248. plot_continuous(train_data, total.night.charge, Churn)
249.
250. # no extreme outliers detected
251.
252. ##### Analysing total.intl.minutes #####
#####
253. str(train_data$total.intl.minutes)
254. plot_continuous(train_data, total.intl.minutes, Churn)
255.
256. # no extreme outliers detected
257.
258. ##### Analysing total.intl.calls #####
#####

```



```

259.     str(train_data$total.intl.calls)
260.     plot_continuous(train_data, total.intl.calls, Churn)
261.
262.     # no extreme outliers detected
263.
264.     ##### Analysing total.intl.charge #####
265.     str(train_data$total.intl.charge)
266.     plot_continuous(train_data, total.intl.charge, Churn)
267.
268.     # no extreme outliers detected
269.
270.     ##### Analysing number.customer.service.calls #####
271.     str(train_data$number.customer.service.calls)
272.
273.     plot_continuous(train_data , number.customer.service.calls, Churn)
274.     table(train_data$number.customer.service.calls)

```

inputPrep.R

```

1. # preparing train,validation and test dataset for random forest and logistic regres
   sion
2.
3. library(caret)
4. library(stringr)
5.
6. # reading train set
7. inputdata <- read.csv("Train_data.csv", sep = ',', header = TRUE, na.strings = c("
   ", "NA"))
8. inputTest <- read.csv("Test_data.csv", sep = ',', header = TRUE, na.strings = c(" "
   ", "NA"))
9.
10. #removing variables which not helpfull in churn reduction. Using these variables do
    sn't make sense.
11. #for training data
12. inputdata$state <- NULL
13. inputdata$area.code <- NULL
14. inputdata$phone.number <- NULL
15.
16. #for test data
17. inputTest$state <- NULL
18. inputTest$area.code <- NULL
19. inputTest$phone.number <- NULL
20.
21. # We are going to implement random forest and logistic regression and compare perfor
    mance of
22. # both models. For comparision we need same training and test data.
23. # random forest can use input both character and numeric data but logistic regressi
    on can only
24. # work on numeric data.
25. # Transforming and re-
    encoding the data so that both random forest and logistic regression can train and
    predict
26. # on same dataset.
27.
28. #for training data
29. inputdata[, 'international.plan'] <- ifelse(str_detect(inputdata[, 'international.pla
    n'], "yes"), 1, 0)
30. inputdata$international.plan <- as.factor(inputdata$international.plan)
31. inputdata[, 'voice.mail.plan'] <- ifelse(str_detect(inputdata[, 'voice.mail.plan'], "y
    es"), 1, 0)
32. inputdata$voice.mail.plan <- as.factor(inputdata$voice.mail.plan)
33. inputdata[, "Churn"] <- str_replace(inputdata[, "Churn"], "\\.", "")
34. inputdata[, "Churn"] <- ifelse (str_detect(inputdata[, "Churn"], "True"), 1, 0)

```



```

35. inputdata$Churn <- as.factor(inputdata$Churn)
36.
37. #for test data
38. inputTest[, 'international.plan'] <- ifelse(str_detect(inputTest[, 'international.pla
n'], "yes"), 1, 0)
39. inputTest$international.plan <- as.factor(inputTest$international.plan)
40. inputTest[, 'voice.mail.plan'] <- ifelse(str_detect(inputTest[, 'voice.mail.plan'], "y
es"), 1, 0)
41. inputTest$voice.mail.plan <- as.factor(inputTest$voice.mail.plan)
42. inputTest[, "Churn"] <- str_replace(inputTest[, "Churn"], "\\.", "")
43. inputTest[, "Churn"] <- ifelse (str_detect(inputTest[, "Churn"], "True"), 1, 0)
44. inputTest$Churn <- as.factor(inputTest$Churn)
45.
46. set.seed(987)
47. # splitting input train data into train and validation set
48. churnIndex <- createDataPartition(inputdata$Churn, p = 0.7, times = 1, list = FALSE
)
49. train_data <- inputdata[churnIndex,]
50. validation_data <- inputdata[-churnIndex,]
51.
52. test_dataset <- inputTest

```

logisticRegression.R

```

1. # Churn Reduction Logistic regression implementation
2.
3. library(stringr)
4. library(corrplot)
5. library(caret)
6. library(dplyr)
7. library(tidyr)
8. library(broom)
9.
10.
11.
12. # One of the assumptions of logistic regression is that there is no high intercorre
lations
13. #among predictors.
14. # From EDA we found some variables have multi-collinearity
15.
16.
17.
18. # Now we will build our model.
19. # For the first model we will select all the predictors
20.
21. fullModel <- glm(Churn ~., data = train_data, family = binomial(link = 'logit'))
22. summary(fullModel)
23.
24. # Predicting and checking model performance with full model
25. pred_fullmodel_val <- predict(fullModel, validation_data[, -
18], type = "response")
26. pred_fullmodel_valT <- ifelse(pred_fullmodel_val < 0.5, 1, 0)
27. xtab=table(observed=validation_data[,18], predicted=pred_fullmodel_valT)
28. fullmodel_confmat_val <- confusionMatrix(xtab)
29. print(fullmodel_confmat_val)
30.
31. # now in second model we will remove variables affected by multi-collinearity
32. model2 <- glm(Churn ~. -total.day.minutes-total.eve.minutes-
33.               total.night.minutes-total.intl.minutes,
34.               data = train_data, family = binomial(link = 'logit'))
35. summary(model2)
36.
37. # Predicting and checking model performance with model2
38. pred_model2_val <- predict(model2, validation_data[, -18], type='response')
39. pred_model2_val <- ifelse(pred_model2_val > 0.5, 1, 0)

```

```

40. xtab1 <- table(observed = validation_data[,18], predicted = pred_model2_val)
41. model2_confmat_val <- confusionMatrix(xtab1)
42.
43. # if we compare fullModel and model2 summary, we can see that
44. # In full model due to variables with multi-
    collineary only international.plan, voice.mail.plan,
45. # number.customer.service.call are considered statistically significant with AIC score 1540.
46.
47. # After removing highly collinear variables in model2, variables nternational.plan,
    voice.mail.plan,
48. # number.customer.service.call along with total.day.charge, total.eve.charge, total
    .night.charge,
49. # and total.intl.charge are considered statistically significant with AIC score 1543.
50.
51. # model with lower AIC score is considered better
52.
53. # NOW we will use Stepwise Procedures
54.
55. #1 backwards elimination
56. backward <- step(fullModel)
57. # from backwards elimination approach, we got a model with AIC 1529.02 which is better
58. #than model2
59. # training model using best result from 'backward'
60. model_backward <- glm(Churn ~ international.plan+voice.mail.plan+number.vmail.messages+
61.                        total.day.minutes+total.eve.minutes+total.night.charge+
62.                        total.intl.calls+total.intl.charge+number.customer.service.
    calls,
63.                        data = train_data, family = binomial(link = 'logit'))
64. summary(model_backward)
65.
66. #predicting on validation set
67. pred_backmodel_val <- predict(model_backward, validation_data[,18], type = 'response')
68. pred_backmodel_valT <- ifelse(pred_backmodel_val > 0.5,1,0)
69. xtab2 <- table(observed = validation_data[,18], predicted = pred_backmodel_valT)
70. backward_confmat_val <- confusionMatrix(xtab2)
71.
72.
73.
74. # model_backward is slightly better than model2
75.
76. # 2. forward elimination
77. nullModel <- glm(Churn ~ 1, data =train_data, family = binomial)
78. summary(nullModel)
79.
80. forward <- step(nullModel, scope = list(lower =formula(nullModel), upper = formula(
    fullModel)),
81.                        direction = 'forward')
82.
83. formula(backward)
84. formula(forward)
85. # both the steps are giving the same formula with same AIC score ie. 1537.7
86.
87. # We are rejecting fullModel due to multi-
    collinearity assumption of logistic regression.
88. # so selecting model_backward because its AIC score is lower than model2 and all the
    predictors
89. # in model_backward are stastically significant.
90.
91.
92.
93. #predicting using model_backward on test set

```

```

94. pred_backmodel <- predict(model_backward, test_dataset[, -18], type = 'response')
95. pred_backmodelT <- ifelse(pred_backmodel > 0.5, 1, 0)
96. xtab_backmodel <- table(observed = test_dataset[, 18], predicted = pred_backmodelT)

97. backmodel_confmat <- confusionMatrix(xtab_backmodel)
98. backmodel_performance <- data.frame(backmodel_confmat$byClass)
99. backmodel_performance <- rbind(accuracy = backmodel_confmat$overall, backmodel_perfor-
    mance)
100.    #comparing model_backward performance on validation and test set for overfit-
    ting
101.    z <- as.data.frame(rbind(backmodel_confmat$byClass, backward_confmat_val$byCla-
    ss))
102.    #model_backward is performing with acceptable variation on both dataset, so
    no overfitting.
103.
104.
105.    # plotting histogram of prediction
106.    pred_hist <- data.frame(pred_backmodel)
107.    pred_hist %>%
108.    ggplot(mapping = aes(x = pred_backmodel)) +
109.    geom_histogram(bins = 50, fill = 'grey40') +
110.    labs(title = " Prediction histograms")
111.
112.    # range of predictions
113.    round(range(pred_backmodel), 2)
114.    median(pred_backmodel)
115.    #The prediction range from 0 to 0.97 with median 0.0872
116.
117.
118.    # Selecting probability threshold value is business context decision and a
119.    # tradoff between true positive and false positive classifications
120.    # The threshold here is set to .5 . This means that anyone with a probabili-
    ty of
121.    # more than .5 is predicted to churn. If we reduce the probability threshold
    , more people will
122.    # be predicted to churn, this gives us a higher number of "at risk customers
    " to target.
123.    # However, this increases the likelihood that customers who are not at risk
    will pass the
124.    # threshold and be predicted to churn.
125.    # If we are concerned with marketing expenditure, then a higher threshold s
    hould be
126.    # targeted (above 0.8 or 0.9). Otherwise, lower thresholds can be targeted,
127.    # so the company can target larger amounts of customers who are at risk of c
    hurning.
128.
129.
130.    #Logistic regression diagnostics - now we are going to check for logistic re-
    gression assumptions
131.
132.    #1. Linearity assumption
133.    # Checking for a linear relationship between continuous variables and logit o
    f the outcome. This
134.    #can be done by visually inspecting the scatter plot between each predictor
    and logit values.
135.
136.    # selecting continuous predictors
137.    continuous_data <- test_dataset %>%
138.    select_if(is.numeric)
139.    predictors <- colnames(continuous_data)
140.
141.    # calculating logit and adding it in 'continuous_data'
142.    continuous_data <- continuous_data %>%
143.    mutate(logit = log(pred_backmodel/(1-pred_backmodel))) %>%

```

```

144.     gather(key = "predictors", value = "predictor.value", -logit)
145.     ggplot(data = continous_data, mapping = aes(x = logit, predictor.value)) +
146.     geom_point(size = 0.5, alpha = 0.5) +
147.     geom_smooth(method = "loess") +
148.     theme_bw() +
149.     facet_wrap(~predictors, scales = "free_y")
150.
151.     # The smoothed scatter plots show all the continous predictors are very near
152.     # linearly associated with the outcome in logit scale.
153.
154.     # 2. Multicollinearity
155.     # we have removed multi-collinear variables while models.
156.     #performing double check using variation inflation factor (VIF)
157.     car::vif(model_backward)
158.     #As a rule of thumb, a VIF value that exceeds 5 or 10 indicates a problemati
c amount of
159.     #collinearity. In our model voice.mail.plan and number.vmail.messages is sho
wing vif score of
160.     # < 16. SO, we will remove one of the variable and retrain and check the fin
al model after
161.     #dignostics.
162.
163.     #3. Influential values
164.     # top 3 largest values
165.     plot(model_backward, which = 4, id.n = 3)
166.     #278, 1862, 3292
167.     # Data points with an absolute standardized residuals above 3 represent poss
ible outliers
168.     #and may deserve closer attention
169.     #augment - add columns to the original dataset such as predictions, residual
s and cluster assignments
170.
171.     model.data <- augment(model_backward) %>%
172.     mutate(index = 1:n())
173.     model.data %>% top_n(3, .cooksd)
174.     # plotting standardized residuals
175.     ggplot(model.data, aes(index, .std.resid)) +
176.     geom_point(aes(color = Churn), alpha = .5) +
177.     theme_bw()
178.     #
179.     model.data %>%
180.     filter(abs(.std.resid) > 3)
181.     # one possible influential observations was found on our training set.
182.     # row 1890
183.
184.     #-----
-----
185.     # creating a final model by removing one of the variables detected by VIF an
d removing
186.     # one influential observations
187.
188.     train_data <- train_data[-1890,]
189.     formula(model_backward)
190.
191.     final_model <- glm(Churn ~ international.plan + voice.mail.plan +
192.     total.day.minutes + total.eve.minutes + total.night.cha
rge +
193.     total.intl.calls + total.intl.charge + number.customer.
service.calls,
194.     data = train_data, family = binomial(link = 'logit'))
195.
196.     #checking VIF
197.     car::vif(final_model)
198.     # no multicollinear variable found in model
199.

```

```

200.     pred_final <- predict(final_model, test_dataset, type = 'response')
201.     pred_finalT <- ifelse(pred_final>0.5, 1,0)
202.     xtab_final <- table(observed = test_dataset[,18], predicted = pred_finalT)
203.     final_confmat <- confusionMatrix(xtab_final)
204.     finalModel_performance <- data.frame(final_confmat$byClass)
205.     finalModel_performance <- rbind(accuracy = final_confmat$overall, finalModel
_performance)
206.
207.
208.
209.     # After comparing confusion matrix of model_backward and final_modal, we fou
nd that
210.     # both models are performing very similar.
211.
212.     #-----
-----
213.
214.     # for comparision with random forest performance with probabiltly threshlod .
75,.25
215.
216.     pred_compareRF <- predict(final_model, test_dataset[, -
18], type = 'response')
217.     pred_compareRFT <- ifelse(pred_compareRF>0.25, 1,0)
218.
219.     xtab_compare <- table(observed = test_dataset[,18], predicted = pred_compare
RFT)
220.     confMatLogit <- confusionMatrix(xtab_compare)
221.     compare_performance <- data.frame(confMatLogit$byClass)
222.     compare_performance <- rbind(accuracy = confMatLogit$overall, compare_perfor
mance)
223.     print(confMatLogit)
224.     print(compare_performance)
225.
226.
227.     # Model input and output for logistic regression
228.     write.csv(test_dataset, file = "InputLogisticRegressionR.csv")
229.     write.csv(pred_compareRFT, file="outputLogisticRegressionR.csv")

```

randomForestMLR.R

```

1. # Randomforest and parameter tuning unis MLR package
2.
3. library(mlr)
4. #library(caret)
5.
6.
7. trainTask <- makeClassifTask(data = train_data, target = "Churn")
8. validationTask <- makeClassifTask(data = validation_data, target = "Churn")
9.
10. #creating randomoforest classifier/learner
11. randomFOrest.learner.baseline <- makeLearner("classif.randomForest")
12. randomFOrest.learner <- makeLearner("classif.randomForest")
13.
14.
15. # setting validaion staertgy using cv with 10 folds
16. cvstr <- makeResampleDesc("CV", iters = 10L)
17.
18. # Our main aim is churn reduction. So we are interest in a model which can classify
customers
19. # who are going to churn out. There for model should reduce false positive rate. As
positive class
20. #for the model is 'false'
21.
22.
23. randomFOrest.learner$par.vals <-

```

```

24. list(ntree = 500L,
25.       importance = TRUE,
26.       cutoff = c(0.70, 0.30))
27.
28. r <- resample(
29.   learner = randomF0rest.learner,
30.   task = trainTask,
31.   resampling = cvstr,
32.   measures = list(tpr, fpr, fnr, fpr, acc),
33.   show.info = TRUE
34. )
35.
36. # Probability threshold or cutoff is a bussiness context decision. The threshold is
   usually set
37. #to .5 by default. This means that anyone with a probability of more than .5 is pre
   dicted to churn.
38. #If we reduce the probability threshold, more people will be predicted to churn
39. # to reduce the false positive rate , we are tuning cut-off here.
40. # cut-off tuning is perfomed by cross-validation
41. # Tuning cutoff value starting with default .50,.50
42. #.70 and .30 is selected as cutoff value
43. # Main purpose here is tuning cutoff.
44. cutoff1 <- calculateConfusionMatrix(r$pred)
45.
46. # tuning mtry ,ntree, nodesize using hypertuning
47. params <- makeParamSet(
48.   makeIntegerParam("mtry", lower = 2, upper = 10),
49.   makeIntegerParam("nodesize", lower = 10, upper = 50),
50.   makeIntegerParam("ntree", lower = 100, upper = 600)
51. )
52.
53. #random search with 100 iterations
54. ctrl <- makeTuneControlRandom(maxit = 100L)
55.
56. tuneRF <- tuneParams(
57.   learner = randomF0rest.learner,
58.   task = trainTask,
59.   resampling = cvstr,
60.   measures = list(acc),
61.   par.set = params,
62.   control = ctrl,
63.   show.info = TRUE
64. )
65.
66. tunedRFmodel <- setHyperPars(learner = randomF0rest.learner,
67.                               par.vals = tuneRF$x)
68.
69. # tuning with default random forest
70. modelRFBaseline <- mlr::train(randomF0rest.learner.baseline, trainTask)
71.
72. # training random forest with tuned parameters
73. modelRF1 <- mlr::train(tunedRFmodel, trainTask)
74.
75. #prediction on validation set using baseline
76. predictBaseline <- predict(modelRFBaseline, validationTask)
77.
78. #performance measures for baseline model on validation set
79. accuracy_bl <- performance(pred = predictBaseline, measures = acc)
80. f1Measure_bl <- performance(pred = predictBaseline, measures = f1)
81. truePositiveRate_bl <-
82.   performance(pred = predictBaseline, measures = tpr)
83. trueNegativeRate_bl <-
84.   performance(pre = predictBaseline, measures = tnr)
85. falsePositiveRate_bl <-
86.   performance(pred = predictBaseline, measures = fpr)
87. falseNegativeRate_bl <-

```

```

88. performance(pred = predictBaseline, measures = fnr)
89.
90. baselinePerformanceMetric <-
91.   c(
92.     accuracy_bl,
93.     f1Measure_bl,
94.     truePositiveRate_bl,
95.     trueNegativeRate_bl,
96.     falsePositiveRate_bl,
97.     falseNegativeRate_bl
98.   )
99.
100.
101.
102.   #prediction on validation set using tuned random forest
103.   predictRF1 <- predict(modelRF1, validationTask)
104.
105.   #performance measures for tuned model on validation set
106.   accuracy_val <- performance(pred = predictRF1, measures = acc)
107.   f1Measure_val <- performance(pred = predictRF1, measures = f1)
108.   truePositiveRate_val <-
109.     performance(pred = predictRF1, measures = tpr)
110.   trueNegativeRate_val <-
111.     performance(pre = predictRF1, measures = tnr)
112.   falsePositiveRate_val <-
113.     performance(pred = predictRF1, measures = fpr)
114.   falseNegativeRate_val <-
115.     performance(pred = predictRF1, measures = fnr)
116.
117.
118.   validationPerformanceMetric <-
119.     c(
120.       accuracy_val,
121.       f1Measure_val,
122.       truePositiveRate_val,
123.       trueNegativeRate_val,
124.       falsePositiveRate_val,
125.       falseNegativeRate_val
126.     )
127.
128.
129.
130.
131.
132.   mainTestTask <- makeClassifTask(data = test_dataset, target = "Churn")
133.
134.
135.   #prediction on main testset using tuned RF
136.
137.   predictTest <- predict(modelRF1, mainTestTask)
138.   confusionMatrixTest <- calculateConfusionMatrix(predictTest)
139.
140.   #calculating performance measure for predictions
141.   accuracy <- performance(pred = predictTest, measures = acc)
142.   f1Measure <- performance(pred = predictTest, measures = f1)
143.   truePositiveRate <- performance(pred = predictTest, measures = tpr)
144.   trueNegativeRate <- performance(pre = predictTest, measures = tnr)
145.   falsePositiveRate <- performance(pred = predictTest, measures = fpr)
146.   falseNegativeRate <- performance(pred = predictTest, measures = fnr)
147.
148.   testPerformanceMetric = c(
149.     accuracy,
150.     f1Measure,
151.     truePositiveRate,
152.     trueNegativeRate,
153.     falsePositiveRate,

```

```

154.         falseNegativeRate
155.     )
156.
157.
158.
159.     xt <- table(observed = test_dataset[,18], predicted = predictTest$data$response
nse)
160.     final_confmatRF <- confusionMatrix(xt)
161.     compare_performanceRF <- data.frame(final_confmatRF$byClass)
162.     compare_performanceRF <- rbind(accuracy = final_confmatRF$overall, compare_p
formanceRF)
163.     print(final_confmatRF)
164.     print(compare_performanceRF)
165.
166.     # Model input and output for Random Forest
167.     write.csv(test_dataset, file = "InputRandomForestR.csv")
168.     write.csv(predictTest$data$response, file="outputRandomForestR.csv")

```

Python code

Customer Churn Reduction - logistic Regression.py

```

1. # coding: utf-8
2.
3. # # Customer Churn Reduction - Logistic Regression
4. # In the notebook, we will use logistic regression to predict whether a customer wi
ll churn or not.
5.
6. # ### importing librarires
7.
8. # In[1]:
9.
10.
11. import pandas as pd
12. import numpy as np
13. import seaborn as sns
14. import matplotlib.pyplot as plt
15. from plotnine import *
16. from sklearn.preprocessing import LabelEncoder
17. from sklearn.linear_model import LogisticRegression
18. from sklearn.model_selection import train_test_split, RandomizedSearchCV, GridSearc
hCV
19. from sklearn.metrics import classification_report
20. from pandas_ml import ConfusionMatrix
21.
22.
23. # In[2]:
24.
25.
26. # Read in train and test data
27. train_data = pd.read_csv("Train_data.csv")
28. test_data = pd.read_csv("Test_data.csv")
29.
30.
31. # In[3]:
32.
33.
34. # change column names for ease of use and display first 5 rows
35. train_data.columns = train_data.columns.str.lower().str.replace(' ', '_')
36. test_data.columns = test_data.columns.str.lower().str.replace(' ', '_')
37.
38.
39. # In[4]:
40.
41.

```



```

42. print(" number of rows and columns in train data ",train_data.shape)
43.
44.
45. # In[5]:
46.
47.
48. print(" number of rows and columns in test data ",test_data.shape)
49.
50.
51.
52. # In[6]:
53.
54.
55. train_data.describe()
56. train_data.head()
57.
58.
59. # In[7]:
60.
61.
62. #removing variables which not helpfull in churn reduction. Using these variables do
    sn't make sense.
63. #for training data
64. train_data = train_data.drop(['state','area_code','phone_number'], axis = 1)
65. test_data = test_data.drop(['state','area_code','phone_number'], axis = 1)
66.
67.
68. # In[8]:
69.
70.
71. # checking for missing values in train
72. train_data.isnull().sum()
73.
74.
75. # In[9]:
76.
77.
78. # checking for missing values in test dataset
79. test_data.isnull().sum()
80.
81.
82. # ##### There are no missing values in train and test data.
83.
84. # Now looking at target variable.
85. # churn: This is the target variable. Churn is defined as whether the customer leav
    es the services or not. churn = True means customer left ,churn = false means custo
    mer stays
86.
87. # In[10]:
88.
89.
90. plt.figure(figsize=(8,6))
91. sns.set_style('ticks')
92. sns.countplot(train_data.churn,palette='summer')
93. plt.xlabel('Customer churn')
94.
95.
96. # In[11]:
97.
98.
99. # churn ratio of customers with and without internation plan
100.     ggplot(train_data) +     aes('international_plan', fill = 'churn') +     geom
        _bar(position = "fill", color= 'blue') + labs(x = "International plan", y = "") +
        ggtitle("Churn ratio of customers with and without international plan") +     th
        eme(figsize=(6, 4))
101.

```

```

102.
103.     # #### Customers with international plan are churning out more as compare to
        domestic customers.
104.
105.     # In[12]:
106.
107.
108.     # churn ratio of customers with voice_mail_plan
109.     ggplot(train_data) + aes('voice_mail_plan', fill = 'churn') + geom_bar
        r(position = "fill", color = 'blue') + labs(x = "voice mail plan", y = "") +
        ggtitle("Churn ratio with customers with and without voice mail plan") +
        theme(figure_size=(6, 4))
110.
111.
112.     # #### Customers without voice mail plan are churning out more as compare to
        customers with voice mail plan.
113.
114.     # In[13]:
115.
116.
117.     # # churn ratio of customers with respect to service calls
118.     ggplot(train_data) + aes('number_customer_service_calls', fill = 'churn'
        ) + geom_bar(position = "fill", color = 'blue') + labs(x = "Customer service
        calls", y = "") + ggtitle("Churn ratio with service call frequency") +
        theme(figure_size=(6, 4))
119.
120.
121.     # #### customers with higher service calls ie > 3 are churning out more.
122.
123.     # ### 3. Correlation matrix for continuous predictors
124.
125.     # In[14]:
126.
127.
128.     churn_corr = train_data.corr()
129.     cmap = sns.diverging_palette(5, 250, as_cmap=True)
130.
131.     def magnify():
132.         return [dict(selector="th",
133.             props=[("font-size", "7pt")]),
134.             dict(selector="td",
135.                 props=[('padding', "0em 0em")]),
136.             dict(selector="th:hover",
137.                 props=[("font-size", "12pt")]),
138.             dict(selector="tr:hover td:hover",
139.                 props=[('max-width', '200px'),
140.                     ('font-size', '12pt')])
141.         ]
142.
143.     churn_corr.style.background_gradient(cmap, axis=1) .set_properties(**{'ma
        x-width': '90px', 'font-
        size': '10pt'}) .set_caption("Correlation matrix") .set_precision(2) .set_
        table_styles(magnify())
144.
145.
146.     # ### From corrplot we can see that dataset consist of multicollinearity
147.     # 1. total.day.minutes and total.day.charge are highly collinear
148.     # 2. total.eve.minutes and total.eve.charge are highly collinear
149.     # 3. total.night.minutes and total.night.charge are highly collinear
150.     # 4. total.intl.minutes and total.intl.charge are highly collinear
151.     #
152.     # Multi-
        collinearity violates the assumption of logistic regression. So we will be removing
        one of these predictors
153.     # from the model.
154.

```

```

155.     # ### 4. Exploring continous predictors
156.
157.     # In[15]:
158.
159.
160.     # function for exploring distributions by continuous predictors with there s
    ummary stats
161.     def countPred_eda(train_data, variableName, targetVariable):
162.         print(train_data[variableName].describe())
163.         return ggplot(train_data) + aes(targetVariable, variableName, fill =
    targetVariable) + geom_boxplot(alpha = .8, outlier_color = "green") + labs(x =
    targetVariable, y = variableName) + ggtitle("Churn ratio with " + variableName
    ) + theme(figure_size=(6, 4))
164.
165.
166.
167.
168.     # In[16]:
169.
170.
171.     # --- total_day_minutes --- #
172.     countPred_eda(train_data, 'total_day_minutes', 'churn')
173.
174.
175.     # ##### It is evident from above plot that churn rate is higher when count of
    total_day_minute is higher
176.
177.     # In[17]:
178.
179.
180.     # --- total_day_calls ---- #
181.     countPred_eda(train_data, 'total_day_calls', 'churn')
182.
183.
184.     # In[18]:
185.
186.
187.     # --- total_day_charge --- #
188.     countPred_eda(train_data, 'total_day_charge', 'churn')
189.
190.
191.     # In[19]:
192.
193.
194.     # --- total_eve_minutes --- #
195.     countPred_eda(train_data, 'total_eve_minutes', 'churn')
196.
197.
198.     # In[20]:
199.
200.
201.     # --- total_eve_calls --- #
202.     countPred_eda(train_data, 'total_eve_calls', 'churn')
203.
204.
205.     # In[21]:
206.
207.
208.     # --- total_eve_charge --- #
209.     countPred_eda(train_data, 'total_eve_charge', 'churn')
210.
211.
212.     # In[22]:
213.
214.
215.     # --- total_night_minutes --- #

```

```

216.     countPred_eda(train_data, 'total_night_minutes', 'churn')
217.
218.
219.     # In[23]:
220.
221.
222.     # --- total_night_calls --- #
223.     countPred_eda(train_data, 'total_night_calls', 'churn')
224.
225.
226.     # In[24]:
227.
228.
229.     # --- total_night_charge --- #
230.     countPred_eda(train_data, 'total_night_charge', 'churn')
231.
232.
233.     # In[25]:
234.
235.
236.     # --- total_intl_minutes --- #
237.     countPred_eda(train_data, 'total_intl_minutes', 'churn')
238.
239.
240.     # In[26]:
241.
242.
243.     # --- total_intl_calls --- #
244.     countPred_eda(train_data, 'total_intl_calls', 'churn')
245.
246.
247.     # In[27]:
248.
249.
250.     # --- total_intl_charge --- #
251.     countPred_eda(train_data, 'total_intl_charge', 'churn')
252.
253.
254.     # In[28]:
255.
256.
257.     # --- number_customer_service_calls --- #
258.     countPred_eda(train_data, 'number_customer_service_calls', 'churn')
259.
260.
261.     # In[29]:
262.
263.
264.     # Removing highly collinear variables from train and test
265.     train_data = train_data.drop(['total_day_charge', 'total_eve_charge', 'total_n
266.     ight_charge', 'total_intl_charge'], axis = 1)
267.     test_data = test_data.drop(['total_day_charge', 'total_eve_charge', 'total_nig
268.     ht_charge', 'total_intl_charge'], axis = 1)
269.
270.
271.
272.     # encoding target variables
273.     le = LabelEncoder()
274.     # for train data
275.     train_data.churn = le.fit_transform(train_data.churn)
276.     # for test data
277.     test_data.churn = le.fit_transform(test_data.churn)
278.
279.

```

```

280.     # In[31]:
281.
282.
283.     # Encoding categorical variables
284.     # for train data
285.     train_data.international_plan = le.fit_transform(train_data.international_pl
an)
286.     train_data.voice_mail_plan = le.fit_transform(train_data.voice_mail_plan)
287.
288.     # for test data
289.     test_data.international_plan = le.fit_transform(test_data.international_plan
)
290.     test_data.voice_mail_plan = le.fit_transform(test_data.voice_mail_plan)
291.
292.
293.     # In[32]:
294.
295.
296.     test_data.head()
297.
298.
299.     # In[33]:
300.
301.
302.     train_data.churn.value_counts()
303.
304.
305.     # In[34]:
306.
307.
308.     # selecting predictors
309.     train_feature_space = train_data.iloc[:,train_data.columns != 'churn']
310.     # selecting target class
311.     target_class = train_data.iloc[:,train_data.columns == 'churn']
312.
313.
314.     # In[35]:
315.
316.
317.     # creating training and validation set
318.     training_set, validation_set, train_taget, validation_target = train_test_sp
lit(train_feature_space,
319.                                     target_c
lass,
320.                                     test_siz
e = 0.30,
321.                                     random_s
tate = 456)
322.
323.     # Cleaning test sets to avoid future warning messages
324.     train_taget = train_taget.values.ravel()
325.     validation_target = validation_target.values.ravel()
326.
327.
328.     # In[36]:
329.
330.
331.     # logistic regression classifier
332.     classifier_logit_default = LogisticRegression(random_state=456)
333.
334.
335.     # In[37]:
336.
337.
338.     classifier_logit_default.fit(training_set, train_taget)
339.

```

```

340.
341.     # In[38]:
342.
343.
344.     # Predicting the validation set results
345.     validation_prediction = classifier_logit_default.predict(validation_set)
346.
347.
348.     # In[39]:
349.
350.
351.     # confusion matrix for validation set
352.     validation_logit_crosstab = pd.crosstab(index = validation_target,
353.                                             columns = validation_prediction)
354.     validation_logit_crosstab = validation_logit_crosstab.rename(columns= {0: 'False', 1: 'True'})
355.     validation_logit_crosstab.index = ['False', 'True']
356.     validation_logit_crosstab.columns.name = 'n = 1000'
357.
358.
359.     # In[40]:
360.
361.
362.     validation_logit_crosstab
363.
364.
365.     # In[41]:
366.
367.
368.     #classification report on validation set
369.     target_names =[0,1]
370.
371.     validation_report = classification_report(validation_prediction, validation_target, target_names )
372.     print(validation_report)
373.
374.
375.     # In[42]:
376.
377.
378.     mean_accuracy_validation = classifier_logit_default.score(validation_set, validation_target)
379.     print(' Mean accuracy on validation set',mean_accuracy_validation)
380.
381.
382.     # In[43]:
383.
384.
385.     # calculating test error rate on validation set
386.     test_error_rate = 1 - mean_accuracy_validation
387.     print(' Test error rate on validation set',test_error_rate)
388.
389.
390.     # From the confusion matrix we can see that high accuracy of model is due to disproportionate number of non-churn
391.     # customers predicted correctly. The This model is working great for identifying non churning customer but performing poorly for
392.     # churning customers. We will tune the model to increase accuracy on churning customer.
393.
394.     # In[44]:
395.
396.
397.     # model 2
398.     classifier_logit_2 = LogisticRegression(class_weight='balanced')
399.

```

```

400.
401.     param = {'C':[0.001,0.005,0.01,0.05,0.1,0.5,1,1.5,2,3]}
402.     #rs_cv = RandomizedSearchCV(estimator=classifier_logit_2, cv = 10,
403.                                #n_iter = 100,
404.                                #param_distributions=param, random_state=1234)
405.     rs_cv = GridSearchCV(estimator=classifier_logit_2, cv = 10,param_grid=param)

406.
407.     rs_cv.fit(training_set,train_taget)
408.
409.     print('Best parameter :{} Best score :{}'.format(rs_cv.best_params_,rs_cv.be
st_score_))
410.
411.
412.     # In[45]:
413.
414.
415.     classifier_logit_2.set_params(C = 3)
416.
417.
418.     # In[46]:
419.
420.
421.     classifier_logit_2.fit(training_set, train_taget)
422.
423.
424.     # In[47]:
425.
426.
427.     validation_prediction_tuned = classifier_logit_2.predict(validation_set)
428.
429.
430.     # In[48]:
431.
432.
433.     # confusion matrix for validation set
434.     validation_logit_crosstb1 = pd.crosstab(index = validation_target,
435.                                             columns = validation_prediction_tuned)
436.     validation_logit_crosstb1 = validation_logit_crosstb1.rename(columns= {0: 'F
alse', 1: 'True'})
437.     validation_logit_crosstb1.index = ['False', 'True']
438.     validation_logit_crosstb1.columns.name = 'n = 1000'
439.
440.
441.     # In[49]:
442.
443.
444.     validation_logit_crosstb1
445.
446.
447.     # In[50]:
448.
449.
450.     #classification report on validation set with hypertuning
451.     target_names =[0,1]
452.
453.     validation_report_tuning = classification_report(validation_prediction_tuned
, validation_target, target_names )
454.     print(validation_report_tuning)
455.
456.
457.     # ### Prediction and performance on test data
458.     #
459.     # Model classifier_logit_2 was selected.
460.
461.     # In[51]:

```

```

462.
463.
464.     # test set
465.     test_set = test_data.iloc[:,test_data.columns != 'churn']
466.     # selecting target class for test set
467.     test_set_target = test_data.iloc[:,test_data.columns == 'churn']
468.
469.
470.     # Predicting the test set results
471.     test_prediction = classifier_logit_2.predict(test_set)
472.
473.
474.     # In[52]:
475.
476.
477.     # confusion matrix for validation set
478.     test_logit_crosstb = pd.crosstab(index = test_data.churn,
479.                                     columns = test_prediction)
480.     test_logit_crosstb = test_logit_crosstb.rename(columns= {0: 'False', 1: 'True
481. e'})
482.     test_logit_crosstb.index = ['False', 'True']
483.     test_logit_crosstb.columns.name = 'n = onservation'
484.
485.     # In[53]:
486.
487.
488.     cm = ConfusionMatrix(test_data.churn, test_prediction)
489.
490.
491.     # In[61]:
492.
493.
494.     cm
495.
496.
497.     # In[60]:
498.
499.
500.     cm.print_stats()
501.
502.
503.     # In[55]:
504.
505.
506.     #classification report on test set with hypertuning
507.     target_names =[0,1]
508.
509.     test_report_tuning = classification_report(test_prediction, test_data.churn,
510. target_names )
511.     print(test_report_tuning)
512.
513.     # In[56]:
514.
515.
516.     mean_accuracy_test = classifier_logit_2.score(test_set, test_set_target)
517.     print(' Mean accuracy on test set',mean_accuracy_test)
518.
519.
520.     # In[57]:
521.
522.
523.     # calculating test error rate on test set
524.     test_error_rate_testset = 1 - mean_accuracy_test
525.     print(' Test error rate on test set',test_error_rate)

```



```

526.
527.
528.     # In[63]:
529.
530.
531.     test_set_target.churn.value_counts()
532.
533.
534.     # In[65]:
535.
536.
537.     # Model input and output
538.     test_set.to_csv('inputLogisticRegressionPython.csv', encoding = 'utf-
8', index = False)
539.     pd.DataFrame(train_target).to_csv('targetLogisticRegressionPython.csv', index
= False)
540.     pd.DataFrame(test_prediction, columns=['predictions']).to_csv('outputLogisti
cRegressionPython.csv')

```

customer Churn Reduction -Random forest.py

```

1. # coding: utf-8
2.
3. # # Customer Churn Reduction - Random forest
4. #
5. # Finding wether a customer will churn out or not. Random forest was use.
6.
7. # In[1]:
8.
9.
10. # importing requiried libraries
11. import pandas as pd
12. import numpy as np
13. from sklearn.preprocessing import LabelEncoder
14. from sklearn.model_selection import train_test_split, RandomizedSearchCV
15. from sklearn.ensemble import RandomForestClassifier
16. import time
17. import random
18. import matplotlib.pyplot as plt
19. from sklearn.model_selection import KFold, cross_val_score
20. from sklearn.metrics import confusion_matrix
21. from sklearn.metrics import classification_report
22. from pandas_ml import ConfusionMatrix
23.
24. get_ipython().run_line_magic('matplotlib', 'inline')
25.
26.
27. # In[2]:
28.
29.
30. # Reading train data
31. inputTrain = pd.read_csv("Train_data.csv")
32. # Reading test data
33. inputTest = pd.read_csv("Test_data.csv")
34.
35.
36. # In[3]:
37.
38.
39. # change column names for ease of use and display first 5 rows
40. inputTrain.columns = inputTrain.columns.str.lower().str.replace(' ', '_')
41. inputTest.columns = inputTest.columns.str.lower().str.replace(' ', '_')
42.
43.
44. # In[4]:

```

```

45.
46.
47. #dimensions
48. inputTrain.shape
49. inputTest.shape
50.
51.
52. # In[5]:
53.
54.
55. # calculating quick summary statistic for continous predictors
56. inputTrain.describe()
57.
58.
59. # In[6]:
60.
61.
62. #removing variables which not helpfull in churn reduction. Using these variables do
    sn't make sense.
63. #for training data
64. inputTrain = inputTrain.drop(['state', 'area_code', 'phone_number'], axis = 1)
65. inputTest = inputTest.drop(['state', 'area_code', 'phone_number'], axis = 1)
66.
67.
68. # In[7]:
69.
70.
71. # sanity check
72. print(inputTrain.shape)
73. print(inputTest.shape)
74.
75.
76. # In[8]:
77.
78.
79. print("Train dataset", inputTrain.churn.value_counts())
80. print("Test data", inputTest.churn.value_counts())
81.
82.
83. # In[9]:
84.
85.
86. # encoding categorical and target variable to binary
87. # converting international_plan, voice_mail_plan and churn
88.
89. le = LabelEncoder()
90.
91. inputTrain.international_plan = le.fit_transform(inputTrain.international_plan)
92. inputTrain.voice_mail_plan = le.fit_transform(inputTrain.voice_mail_plan)
93. inputTrain.churn = le.fit_transform(inputTrain.churn)
94.
95. inputTest.international_plan = le.fit_transform(inputTest.international_plan)
96. inputTest.voice_mail_plan = le.fit_transform(inputTest.voice_mail_plan)
97. inputTest.churn = le.fit_transform(inputTest.churn)
98.
99.
100.     # In[10]:
101.
102.
103.     # selecting predictors
104.     train_feature_space = inputTrain.iloc[:, inputTrain.columns != 'churn']
105.     # selecting target class
106.     target_class = inputTrain.iloc[:, inputTrain.columns == 'churn']
107.
108.
109.     # In[11]:

```

```

110.
111.
112.     # creating training and validation set
113.     training_set, validation_set, train_taget, validation_target = train_test_sp
lit(train_feature_space,
114.                                     target_c
lass,
115.                                     test_siz
e = 0.30,
116.                                     random_s
tate = 12345)
117.
118.     # Cleaning test sets to avoid future warning messages
119.     train_taget = train_taget.values.ravel()
120.     validation_target = validation_target.values.ravel()
121.
122.
123.     # ## Random forest Implementation
124.
125.     # In[12]:
126.
127.
128.     # using random forest classifier. setting a random state
129.     fit_randomForest = RandomForestClassifier(random_state=12345)
130.
131.
132.     # ### Hyperparameters Optimization
133.     #
134.     # Utilizing the RandomizedSearchCV functionality, we create a dictionary wit
h parameters we are looking to optimize to create the best model for our data.
135.
136.     # In[13]:
137.
138.
139.     np.random.seed(12)
140.     start = time.time()
141.
142.     # selecting best max_depth, maximum features, split criterion and number of
trees
143.     param_dist = {'max_depth': [2,4,6,8,10],
144.                   'bootstrap': [True, False],
145.                   'max_features': ['auto', 'sqrt', 'log2', None],
146.                   'criterion': ["gini", "entropy"],
147.                   "n_estimators" : [100 ,200 ,300 ,400 ,500]
148.                   }
149.     cv_randomForest = RandomizedSearchCV(fit_randomForest, cv = 10,
150.                                         param_distributions = param_dist,
151.                                         n_iter = 10)
152.
153.     cv_randomForest.fit(training_set, train_taget)
154.     print('Best Parameters using random search: \n',
155.           cv_randomForest.best_params_)
156.     end = time.time()
157.     print('Time taken in random search: {0: .2f}'.format(end - start))
158.
159.
160.     # In[39]:
161.
162.
163.     # Set best parameters given by random search
164.     fit_randomForest.set_params(criterion = 'gini',
165.                                 max_features = 'auto',
166.                                 max_depth = 10,
167.                                 n_estimators = 100
168.                                 )
169.

```

```

170.
171.     # In[40]:
172.
173.
174.     fit_randomForest.fit(training_set, train_taget)
175.
176.
177.     # In[16]:
178.
179.
180.     importances_rf = fit_randomForest.feature_importances_
181.     indices_rf = np.argsort(importances_rf)[::-1]
182.
183.
184.     # In[17]:
185.
186.
187.     def variable_importance_plot(importance, indices, training_set):
188.
189.         index = np.arange(len(training_set.columns))
190.
191.
192.         importance_desc = sorted(importance)
193.         feature_space = []
194.         for i in range(16, -1, -1):
195.             feature_space.append(training_set.columns[indices[i]])
196.
197.         fig, ax = plt.subplots(figsize=(14, 14))
198.
199.         ax.set_facecolor('#fafafa')
200.         plt.title('Feature importances for Random Forest Model for Customer chur
n')
201.         plt.barh(index,
202.                  importance_desc,
203.                  align="center",
204.                  color = '#875FDB')
205.         plt.yticks(index,
206.                    feature_space)
207.         plt.xlim(0, max(importance_desc))
208.         plt.xlabel('Mean Decrease in gini')
209.         plt.ylabel('Feature')
210.         plt.savefig('VarImp.png')
211.         #savefig('VarImp.pdf')
212.         plt.show()
213.
214.
215.         plt.close()
216.
217.
218.     # In[18]:
219.
220.
221.     variable_importance_plot(importances_rf, indices_rf, training_set)
222.
223.
224.     # ## perfoming Cross validation
225.
226.     # In[19]:
227.
228.
229.     # function to perform cross validation
230.     def cross_val_metrics(fit, training_set, train_taget, print_results = True):
231.
232.         n = KFold(n_splits=10)
233.         scores = cross_val_score(fit,

```

```

234.         training_set,
235.         train_taget,
236.         cv = n)
237.     if print_results:
238.         print("Accuracy: {0: 0.3f} (+/- {1: 0.3f})" .format(scores.mean(), scores.std() / 2))
239.     else:
240.         return scores.mean(), scores.std() / 2
241.
242.
243.     # In[20]:
244.
245.
246.     cross_val_metrics(fit_randomForest, training_set,
247.                       train_taget,
248.                       print_results = True)
249.
250.
251.     # ## prediction and performance measure on validation set
252.
253.     # In[41]:
254.
255.
256.     predictions_randomForest_validation = fit_randomForest.predict(validation_set)
257.
258.     validation_crosstb = pd.crosstab(index = validation_target,
259.                                     columns = predictions_randomForest_validation)
260.     validation_crosstb = validation_crosstb.rename(columns= {0: 'False', 1: 'True'})
261.     validation_crosstb.index = ['False', 'True']
262.     validation_crosstb.columns.name = 'n = 1000'
263.
264.
265.     # In[42]:
266.
267.
268.     # confusion matrix of validation set
269.     validation_crosstb
270.
271.
272.     # In[23]:
273.
274.
275.     # mean accuracy on validation set
276.     accuracy_randomForest_val = fit_randomForest.score(validation_set, validation_target)
277.     print(' Mean accuracy on validation set',accuracy_randomForest_val)
278.
279.
280.     # In[24]:
281.
282.
283.     # calculating test error rate on validation set
284.     test_error_rate = 1 - accuracy_randomForest_val
285.     print(' Test error rate on validation set',test_error_rate)
286.
287.
288.     # In[25]:
289.
290.
291.     #classification report on validation set
292.     target_names =[0,1]
293.
294.     validation_report = classification_report(predictions_randomForest_validation, validation_target, target_names )

```

```

295.     print(validation_report)
296.
297.
298.     # ## prediction and performance measure on test set
299.
300.     # In[43]:
301.
302.
303.     #selecting predictors
304.     test_set = inputTest.iloc[:,inputTest.columns != 'churn']
305.     # selecting target class
306.     test_target = inputTest.iloc[:,inputTest.columns == 'churn']
307.
308.     test_prediction = fit_randomForest.predict(test_set)
309.
310.
311.     # In[44]:
312.
313.
314.     #performing prediction on test set
315.     test_prediction = fit_randomForest.predict(test_set)
316.
317.
318.     # In[45]:
319.
320.
321.     # creating confusion matrix of test set
322.     confusion_matrix(test_target,test_prediction)
323.
324.
325.     # In[46]:
326.
327.
328.     test_rf_crosstb = pd.crosstab(index = test_target.churn,
329.                                   columns = test_prediction)
330.     test_rf_crosstb = test_rf_crosstb.rename(columns= {0: 'False', 1: 'True'})
331.     test_rf_crosstb.index = ['False', 'True']
332.     test_rf_crosstb.columns.name = 'n = 1667'
333.
334.
335.     # In[47]:
336.
337.
338.     test_rf_crosstb
339.
340.
341.     # In[31]:
342.
343.
344.     # mean accuracy on test set
345.     accuracy_randomForest_test = fit_randomForest.score(test_set, test_target.ch
urn)
346.     print(' Mean accuracy on test set',accuracy_randomForest_test)
347.
348.
349.     # In[32]:
350.
351.
352.     # calculating test error rate on test set
353.     test_error_rate_testset = 1 - accuracy_randomForest_test
354.     print(' Test error rate on test set',test_error_rate)
355.
356.
357.     # In[33]:
358.
359.

```

```

360.     #classification report on test set
361.     target_names =[0,1]
362.
363.     test_report = classification_report(test_prediction, test_target.churn, tar
get_names )
364.     print(test_report)
365.
366.
367.     # In[35]:
368.
369.
370.     cm = ConfusionMatrix(test_target.churn, test_prediction)
371.
372.
373.     # In[38]:
374.
375.
376.     cm
377.
378.
379.     # In[48]:
380.
381.
382.     cm.print_stats()
383.
384.
385.     # In[69]:
386.
387.
388.     # Model input and output
389.     test_set.to_csv('inputRandomForestPython.csv', encoding = 'utf-
8', index = False)
390.     pd.DataFrame(train_taget).to_csv('targetInputRandomForestPython.csv', index
= False)
391.     pd.DataFrame(test_prediction, columns=['predictions']).to_csv('outputRandomF
orestPython.csv')

```