# Predicting customer behavior using Machine Learning models

## 1 Break-Even Response Rate & ROI

```r
cost_per_offer <- 1.5
amazon_prime_fee <- 8.99
avg_revenue <-40
COGS <- 0.7
shipping <- 6

profit_per_customer <- avg_revenue*(1-COGS) + amazon_prime_fee – shipping #We multiply the average
#revenue with the percentage that is not COGS, then add revenue generated by the prime membership
#and subtract the shipping cost supported by Amazon.


breakeven_response_rate <- cost_per_offer/profit_per_customer
```

```r
# pls do not modify the codes below
# these are for TAs to check results
print(paste("cost_per_offer is ", cost_per_offer))
```

```
[1] "cost_per_offer is  1.5"
```

```r
print(paste("profit_per_customer is", profit_per_customer))
```

[1] "profit_per_customer is 14.99"

```r
print(paste("breakeven_response_rate is", breakeven_response_rate))
```

[1] "breakeven_response_rate is 0.10006671114076"

```r
sub_sum <- sum(data_full$subscribe == "yes", na.rm = TRUE) #Filtered and summed only
#the customers who are subscribers (ignored the null values).
total_costs_of_mailing_blanket <- cost_per_offer * 10000 ##The total cost of mailing
#all 10000 customers. Whether they will subscribe or not, we still incur this cost.


total_profit_blanket <- profit_per_customer * sub_sum #Determined the expected profit.



ROI_blanket <- (total_profit_blanket-total_costs_of_mailing_blanket)/
  total_costs_of_mailing_blanket
```

```r
print(paste("total_costs_of_mailing_blanket is ", total_costs_of_mailing_blanket))
```

[1] "total_costs_of_mailing_blanket is  15000"

```r
print(paste("total_profit_blanket is ", total_profit_blanket))
```

[1] "total_profit_blanket is  12561.62"

```r
print(paste("ROI_blanket is ", ROI_blanket))
```

[1] "ROI_blanket is  -0.162558666666667"

# 2 Unsupervised Learning for Segmentation and Targeting

```r
# RFM variables

rfm <- data_full%>%

  mutate(recency = last,  # days since last purchase

          frequency = home + sports + clothes + health + books + digital + toys,

          #total purchases across categories

          monetary_value = electronics+nonelectronics  #total spending


  )


summary(rfm[c("recency", "frequency", "monetary_value")])
```

```
     recency          frequency       monetary_value

 Min.   : 1.00    Min.   : 1.00    Min.   : 15.0

 1st Qu.: 7.00    1st Qu.: 1.00    1st Qu.:128.0

 Median :11.00    Median : 2.00    Median :209.0

 Mean   :12.26    Mean   : 3.85    Mean   :208.2

 3rd Qu.:15.00    3rd Qu.: 6.00    3rd Qu.:284.0

 Max.   :35.00    Max.   :12.00    Max.   :478.0
```

```r
#Scaling data

data_kmeans <- rfm %>%

    mutate(

        recency = scale(recency),

        frequency = scale(frequency),

        monetary_value = scale(monetary_value)

    )

#K-means clustering

data_kmeans <- data_kmeans %>%

    select(recency, frequency, monetary_value) %>%

    mutate(
```
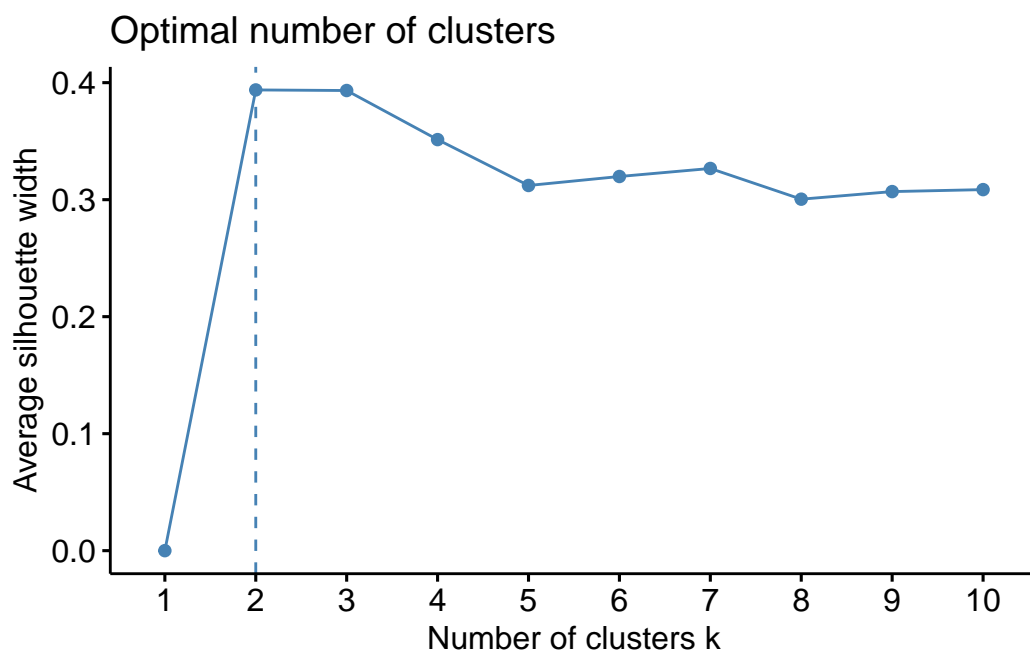
```
12          recency = scale(recency),

13          frequency = scale(frequency),

14          monetary_value = scale(monetary_value)

15      )
```

```
1   # Determine the optimal number of clusters using the Silhouette method below

2   set.seed(888)

3   pacman::p_load(factoextra, cluster)

4   fviz_nbclust(data_kmeans, kmeans, method = "silhouette")
```

## Optimal number of clusters



```
1   # implement k-means clustering below

2

3   # do not modify seeds

4   set.seed(888)

5

6   result_kmeans <- kmeans(data_kmeans,

7       centers = 2,

8       nstart = 10
```

```
9    )
```

```
1    # use broom::tidy() to check the clusters.
2    pacman::p_load(broom)
3    tidy(result_kmeans)
```

```
# A tibble: 2 x 6
    recency frequency monetary_value  size withinss cluster
      <dbl>     <dbl>          <dbl> <int>    <dbl> <fct>
1 -0.00843    -0.547         -0.344  7263   13207. 1
2  0.0224      1.45           0.912  2737    5712. 2
```

```
1    data_full <- data_full %>%
2        mutate(subscribe = ifelse(subscribe == "yes", 1, 0))
3    data_full <- data_full %>%
4        mutate(segment = result_kmeans$cluster)
5    data_full %>%
6        group_by(segment) %>%
7        summarise(avg_subscribe_rate = mean(subscribe, na.rm = T)) %>%
8        ungroup()
```

```
# A tibble: 2 x 2
  segment avg_subscribe_rate
    <int>              <dbl>
1       1             0.0680
2       2             0.126
```

```
1    # ROI for cluster 1 and 2
2    cluster2 <- sum(result_kmeans$cluster == 2, na.rm = TRUE)
3    total_costs_of_mailing_kmeans <- cost_per_offer*cluster2
4    total_profit_kmeans <- profit_per_customer*cluster2* 0.12568506
5    ROI_kmeans <- (total_profit_kmeans-total_costs_of_mailing_kmeans)/
6        total_costs_of_mailing_kmeans
```

```
7
8    cluster1 <- sum(result_kmeans$cluster == 1, na.rm = TRUE)
9    total_costs_of_mailing_kmeans <- cost_per_offer*cluster1
10   total_profit_kmeans <- profit_per_customer*cluster1* 0.06801597
11   ROI_kmeans1 <- (total_profit_kmeans-total_costs_of_mailing_kmeans)/
12     total_costs_of_mailing_kmeans
13
14   ROI_kmeans
```

```
[1] 0.2560127
```

```
1    ROI_kmeans1
```

```
[1] -0.3202937
```

```
1    print(paste("ROI_kmeans is ", ROI_kmeans))
```

```
[1] "ROI_kmeans is  0.2560126996"
```

## 3 Decision Tree Analysis

```
1    data_tree <- rfm %>%
2        select(-user_id) %>%
3        select(-gender) %>%
4        select(-city)
5    # set seed
6    set.seed(1314520)
7
8    data_tree <- rfm %>%
9        mutate(subscribe = ifelse(subscribe == "yes", 1, 0))
10   n_rows_data_tree <- nrow(data_tree)
```

```
11   training_set_index <- sample(
12     x = 1:n_rows_data_tree,
13     size = 0.75*n_rows_data_tree,
14     replace = FALSE
15   )
16
17
18   # create data_training and data_test
19   data_training <- data_tree %>%
20     slice(training_set_index)
21
22   data_test <- data_tree %>%
23     slice(-training_set_index)
```

```
1   # This is to print out first 5 customers
2   training_set_index[1:5]
```

```
[1] 3620    43 3574 4308 7387
```
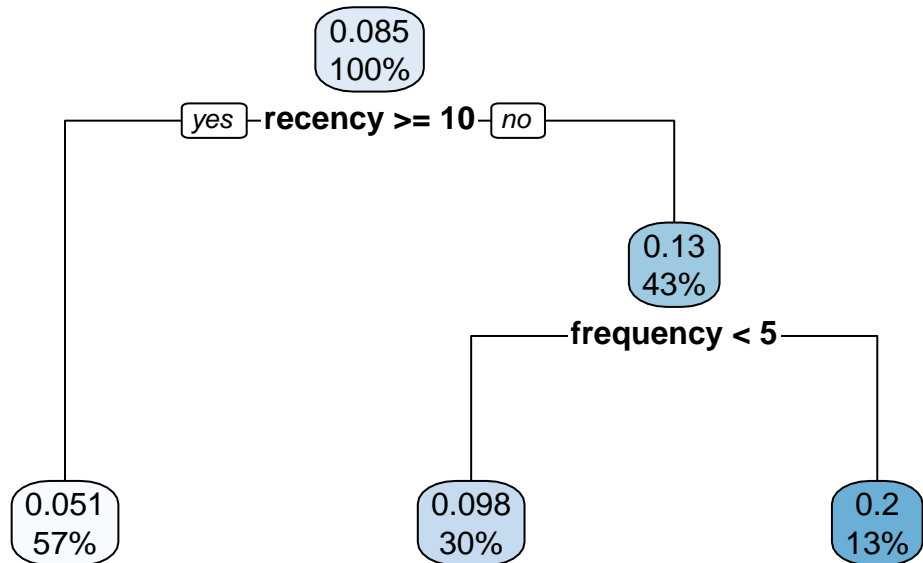
```
1   pacman::p_load(rpart, rpart.plot)
2
3   # train model tree1 below
4
5   tree1 <- rpart(
6     formula = subscribe ~ recency + frequency + monetary_value, #subscribe is the
7     #variable Tom tries to predict based on predictors
8     #(recency, frequency, monetary_value).
9     data = data_training, #training the model
10    method = "anova" #probabilities
11   )
12
13
14   # visualize tree1 below
```

```
15
16    rpart.plot(tree1)
```



```
1     # ROI for tree1
2     prediction_from_decision_tree <- predict(tree1, data_test)
3     data_test <- data_test %>%
4         mutate(predicted_prob_decisiontree = prediction_from_decision_tree)
5     data_test <- data_test %>%
6         mutate(is_target_decisiontree = ifelse(predicted_prob_decisiontree >
7                                        breakeven_response_rate, 1, 0))
8     total_costs_of_mailing_decisiontree <- cost_per_offer *
9       sum(data_test$is_target_decisiontree)
10    data_test_targeted_customers <- data_test %>%
11        filter(is_target_decisiontree == 1)
12    total_profit_decisiontree <- sum(data_test_targeted_customers$subscribe) *
13      profit_per_customer
14
15    # Compute ROI
```

```
16  ROI_decisiontree <-

17    (total_profit_decisiontree - total_costs_of_mailing_decisiontree)/

18    total_costs_of_mailing_decisiontree

19

20  ROI_decisiontree
```

```
[1] 0.7185987
```

## 4 Random Forest

```
1   pacman::p_load(ranger)

2   set.seed(888)

3   # train the random forest model below

4   randomforest <- ranger(

5     formula = subscribe ~ recency + frequency + monetary_value, #subscribe is the

6     #variable Tom tries to predict based on predictors

7     #(recency, frequency, monetary_value)

8     data = data_training, # training the model

9     probability = TRUE, #TRUE because we are interested in finding the probability

10    #of a customer subscribing

11    num.trees = 5000

12  )

13

14  prediction_from_randomforest <- predict(randomforest, data_test)

15

16  data_test <- data_test %>%

17      mutate(predicted_prob_randomforest = prediction_from_randomforest$predictions[, 2])

18

19  data_test <- data_test %>%

20      mutate(is_target_randomforest = ifelse(predicted_prob_randomforest >

21                                           breakeven_response_rate, 1, 0))
```

```
22
23  total_costs_of_mailing_randomforest <- cost_per_offer *
24    sum(data_test$is_target_randomforest)
25
26
27  data_responding_targeted_customers <- data_test %>%
28      filter(is_target_randomforest == 1) %>%
29      filter(subscribe == 1)
30
31  # total profits from responding customers
32  total_profit_randomforest <- nrow(data_responding_targeted_customers) * profit_per_customer
33
34  # Compute ROI
35  ROI_randomforest <-
36    (total_profit_randomforest - total_costs_of_mailing_randomforest) /
37    total_costs_of_mailing_randomforest
38
39  ROI_randomforest
```

    [1] 0.4295828

```
1  print(paste("ROI_randomforest is ", ROI_randomforest))
```

    [1] "ROI_randomforest is  0.429582760201743"