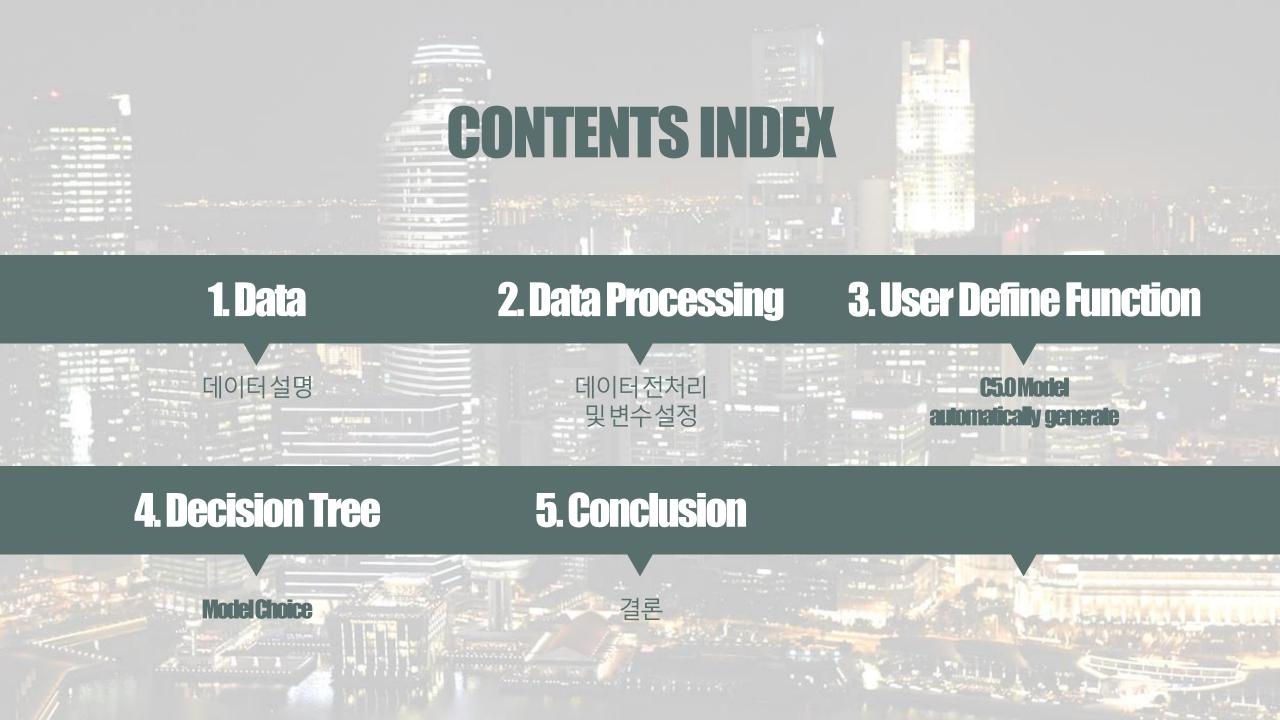


박데이터 경영 MBA 5기 U2016054 이병준 2016년 10월 29일 23시



1. Data

○ 사용데이터

- 2조 Customer Signatures 중 수업시간에만들었던 기본 15개 파생변수와 조별 과제에서 생성한 11개 데이터 사용.
 - 기본파생변수:환불금액,환불횟수,구매상품다양성,내점일수와구매주기,내점당구매건수,주중주말구매패턴,12/6/3개월 구매금액및횟수
 - 2조파생변수:가격선호도,시즌선호도,주구매상품,휴면/이탈고객,주구매시간대,선호할부,선호구매파트,주이용지점,선호할인율,
- -6조CS Data "instMonth" 평균할부기간데이터사용.

```
custid sex rf_amt rf_cnt buy_brd visits API
                                                NPPV wk_amt we_amt
                                                                      wk_pat
                                                                                amt12 nop12
                                                                                               amt6 nop6
                                                                                                           amt3
                                      4 91 1.500000 254300 127000 weekdays
                                                                               381300
                                                                                          6 381300
                                                                                                          89700
                                                            26100 weekdays
                                                                               155700
                                                                                         4 155700
                 0
                                     3 121 1.333333 129600
                                                                                                      4 155700
                                     15 24 1.933333 460900 1556500 weekend 2017400
           -221730
                                                                                         29 786060
                                                                                                     14 18060
                                         52 1.571429 767000 1903000 weekend 2670000
                                                                                        11 2670000
                                                                                                     11 2670000
                                    114 3 2.131579 6917729 3465820 weekdays 10383549
        2 -1934341
                                                                                        243 6040055 114 2425630
        2 -820800
                                      9 40 3.111111 457580 554100
                                                                        none 1011680
                                                                                        28 1011680
                                                                                                     28 1011680
nop3 pr_pref sea_pref stay_out fav_time fav_paymthd fav_part_mean_amt fav_part_cnt main_store avg_disc fav_part_index
  2 3-grade
              Spring
                                 19A
                                                          90533.33
                                                                                   Sinchon
                                                                                                4%
                         stay
  4 2-grade
              Spring
                         out
                                 19A
                                                          43200.00
                                                                                   Sinchon
                                                                                               10%
  3 3-grade
                Fall
                                 18人
                                                          45424.17
                                                                                   Chunho
                                                                                                3%
                         stay
 11 4-grade
              Spring
                                 19周
                                                         183860.00
                                                                                  Sinchon
                                                                                                7%
                         out
 47 1-grade
              Summer
                         out
                                 18人
                                                          64108.94
                                                                                  Sinchon
                                                                                                6%
 28 1-grade
              Spring
                         out
                                 18A
                                                          20718.57
                                                                                    Muyuk
                                                                                                4%
               PAPV P_group p_trend instCnt instMonth instRatio
fav_good_index
               90000
                                                           0.17
                                new
                                                  3.0
            2 50000
                                                  0.0
                                                           0.00
                                new
            3 130000
                                                  3.0
                                                           0.31
                                keep
            4 380000
                                         10
                                                  5.1
                                                           0.91
                                new
              90000
                            loyalty
                                                  2.9
                                                           0.26
            6 110000
                                                  3.4
                                                           0.75
                                 new
```

2. Data Processing [데이터 전처리]

C Library

```
library(plyr) # Join 및 데이터 전처리를 위해 호출
library(dplyr) # 데이터 전처리를 위해 호출
library(C50) # C50, caret, ROCR Decision Tree 생성을 위해 호출
library(caret)
library(ROCR)
library(combinat) # 조합을 사용하기 위해 호출
library(gtools) # 순열을 사용하기 위해 호출
library(Epi) # 시각화
```

Processing

```
user <- read.csv("2_group.csv", stringsAsFactors = T)
user_tree <- user[,c(1,2,37:length(user))]
# 0 : 무효 1 : 남성 2: 대
user_tree <- user_tree[user_tree$sex!=0,] #무효는 제외
user_tree$sex <- factor(user_tree$sex)
#명목형 변수 한글사용 시 Model이 생성되지 않는 점과 Factor화
user_tree$wk_pat <- as.character(user_tree$wk_pat)
user_tree[user_tree$wk_pat=="주울형",]$wk_pat <- "weekend"
user_tree[user_tree$wk_pat=="주형형",]$wk_pat <- "weekdays"
user_tree[user_tree$wk_pat=="주형헌음",]$wk_pat <- "none"
user_tree[user_tree$wk_pat=="유형헌음",]$wk_pat <- "none"
user_tree[user_tree$main_store=="신촌점",]$main_store <- "Sinchon"
user_tree[user_tree$main_store=="신촌점",]$main_store <- "Chunho"
user_tree[user_tree$main_store=="건호점",]$main_store <- "Muyuk"
user_tree[user_tree$main_store=="본점",]$main_store <- "Bon"
user_tree$main_store <- factor(user_tree$main_store)
```

2. Data Processing [데이터 전처리]

Processing

```
user_tree$group_member <- as.character(user_tree$group_member)
user_tree[user_tree$group_member=="그룹사",]$group_member <- "group"
user_tree[user_tree$group_member=="일반회원",]$group_member <- "normal"
user_treesgroup_member <- factor(user_treesgroup_member)
user_tree[user_tree$stay_out=="",]$stay_out <- "뮤지" # 해당사항이 안될 경우 뮤지로 판단.
user_tree$stay_out <- as.character(user_tree$stay_out)
user_tree[user_tree$stay_out=="이탈/휴면",]$stay_out <- "out"
user_tree[user_tree$stay_out=="유지",]$stay_out <- "stay"
user_tree$stay_out <- factor(user_tree$stay_out)
#선호 할부 개월에서 한글인 개월 삭제.
user_tree$fav_paymthd <- as.character(user_tree$fav_paymthd)
user_tree$fav_paymthd <- gsub(pattern = "개월",replacement = "",user_tree$fav_paymthd)
user_tree$fav_paymthd <- factor(user_tree$fav_paymthd)
                                                                                                            fav_part_index
                                                                                                                              fav_part
                                                                                                                        1 스포츠캐주얼
                                                                                                                              패션잡화
 한글 벡터를 포함하고 있는 데이터 항목들 전체 인덱스화 (우유 : 1, 치즈 : 2와 같은형태)
                                                                                                                              열라이브
df <- data.frame(fav_part_index = 1:length(unique(user_tree$fav_part)), fav_part =unique(user_tree$fav_part))
user_tree <- join(user_tree,df,by="fav_part")
user_tree$fav_part_index <- factor(user_tree$fav_part_index)
df2 <- data.frame(fav_good_index = 1:length(unique(user_tree$fav_goodcd)), fav_goodcd =unique(user_tree$fav_goodcd))
user_tree <- join(user_tree,df2,by="fav_goodcd")
                                                                                                            fav_good_index
                                                                                                                               fav_goodcd
user_tree$fav_good_index <- factor(user_tree$fav_good_index)
                                                                                                                                골프웨머
                                                                                                                          수입종합화장품
user_tree <- user_tree[,-c(12,13,22,26,31)]
                                                                                                                              미확인코너
cust_6 <- read.csv("6custsig.csv", stringsAsFactors = F)</pre>
tmp_cust <- cust_6[,c("custid","PAPV","P_group","p_trend","instCnt","instMonth","instRatio")]</pre>
user_tree <- join(user_tree,tmp_cust,by="custid")
```

3. User Defined Function

사용자 함수 생성 (자동 독립변수 할당 및 옵션 적용 함수)

- 함수생성이유
 - C5모델생성의경우다양한방법과조건, Formula내변수설정등다각도로접근이가능하여사용자또는분석가가일일이적용할때많은 시간과노력이필요하다고판단되어사용자함수를생성.
 - 파라미터
 - ✓ Train Data:트레이닝용데이터
 - ✓ Test Data:테스트용데이터
 - ✓ Formula 내 독립 변수: 문자벡터
 - ✓ Control Parameters: 컨트롤및 C5.0 모델링시사용되는 옵션 값들 (Winnowing, Global Pruning, Boosting)
 - ✓ CF: Pruning Severity 가지치기의강도
 - 사용법

```
str <- c("fav_time","buy_brd","avg_disc","amt12")
# Sex ~ 독립변수 값 들에 할당 결과 : sex ~ fav_time + buy_brd + avg_disc + amt12 의 formula로 변환
sub <- c(T,F) # T,F의 값들 중 3개를 선택하여 표시 TTT,TFT 등등.
count <- c(seq(0,40,10)) # trials = 0 ~ 40 step 10 형태로 입력.
train_data <- user.train[,-1] # Training Data
test_data <- user.test # Testing Data
result <- make(user.train[,-1],user.test,str,sub,count) # 함수적용.
```

결과값.

3. User Defined Function

사용자 함수 (자동 독립변수 할당 및 옵션 적용 함수)

```
make <- function(train_data,test_data,str,detail,cnt,CF=0.25){</pre>
 df <- data.frame(Accur = integer(), Param = character(), Sub = character(), Cnt = integer())</pre>
 permn_str <- permn(str) # 입력받은 독립변수들을 순열로 생성
 sub_comb <- permutations(2,3,detail,repeats=TRUE) # T.F로 들어온 값을 2개 중 3번 선택, 중복 허용
 length_permn <- length(permn_str) # 순열의 크기
 length_sub <- NROW(sub_comb) # 생성된 Options의 크기
 length_count <- length(cnt) # trials의 크기.
 rows <- 0
 for(i in 1:length_permn){
   # 순열의 크기만큼 반복 (반복시키는 미유 - winnowing=T에서 입력되는 독립변수 순서에 따라 사용률이 달라지고 모델도 달라짐)
   param <- paste(permn_str[[i]],collapse = "+") # 독립변수들끼리 +로 문자열 합치기 a + b + c
   for(j in 1:length_sub){ # Options만큼 반복.
     sub_1 <- sub_comb[i,1] # Winnowing</pre>
     sub_2 <- sub_comb[j,2] # noGlobalPruning</pre>
     sub_3 <- sub_comb[j,3] # Trials</pre>
     c5_options <- C5.0Control(winnow = sub_1, noGlobalPruning = sub_2, CF = CF) # 옵션 생성 Default CF = 0.25
     for(k in 1:length_count) { # trials 횟수 만큼 반복
       rows <- rows + 1
       params <- as.formula(gsub("\\\"","",paste("sex",param, sep=" ~ "))) # Formula 생성.
       if(cnt[k]>0){ # trials > 0 미상일때 5
        c5_model <- C5.0(params, data=train_data, control = c5_options, rules=sub_3, trials=cnt[k])
         c5_model <- C5.0(params, data=train_data, control = c5_options, rules=sub_3)
       test_data$c5_pred <- predict(c5_model,test_data,type="class")</pre>
       test_data$c5_pred_prob <- round(predict(c5_model,test_data,type="prob"),2)
       Accur <- confusionMatrix(test_data$c5_pred, test_data$sex)
       AccurResult <- Accur$overall[[1]] # 해당 모델의 정확도 확인
       param_in <- deparse(params) # Formula 를 문자열 형태로 변환
       sub <- toString(sub_comb[j,]) # List에 있는 T,F 값들을 문자열 형태로 변환
       df2 <- data.frame(Accur = AccurResult, Param = param_in, Sub = sub, Cnt = cnt[k])
       # DataFrame에 동적으로 삽입하기 위해 DF생성
       df <- rbind(df,df2) #df에 합치기.
   print(rows)
 return(df) #모든 반복이후 결과값 출력
```

4. Decision Tree

Model Choice

- 모델선택
 - ✓ 돌렸던 변수및 조건들에대해기장정확도가높은모델선정.
 - ✓ 옵션 winnowing:F, noGlobalPruning=F, rules=T, trials=20
 - ✓ 독립변수: pr_pref+sea_pref+API+amt12+instMonth+avg_disc+fav_part_index+fav_time+nop6+fav_good_index+buy_brd+ main_store

```
boost
              7167(23.9%) <<
                      <-classified as
          (a) (b)
         2858 6240 (a): class 1
          927 19959
                      (b): class 2
       Attribute usage:
       100.00% pr_pref
       100.00% sea_pref
       100.00% amt12
       100.00% instMonth
       100.00% avg_disc
       100.00% fav_part_index
       100.00% fav_time
       100.00% main_store
        96.93% fav_good_index
        95.28% API
        91.84% buy_brd
        80.66% nop6
```

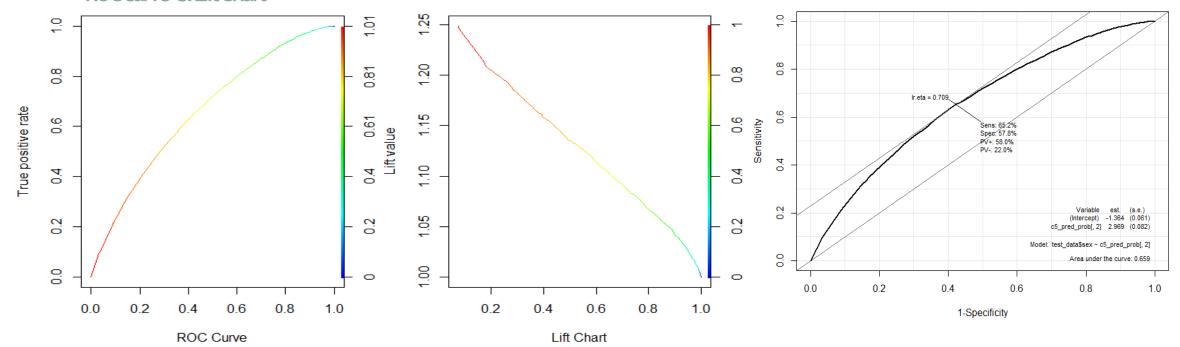
```
Confusion Matrix and Statistics
         Reference
Prediction 1
        1 1296 1019
        2 4768 12904
              Accuracy: 0.7105
                95% CI : (0.7041, 0.7167)
   No Information Rate: 0.6966
   P-Value [Acc > NIR] : 9.687e-06
                 Карра : 0.1702
 Mcnemar's Test P-Value : < 2.2e-16
           Sensitivity: 0.21372
           Specificity: 0.92681
        Pos Pred Value: 0.55983
        Neg Pred Value: 0.73019
            Prevalence: 0.30340
        Detection Rate: 0.06484
   Detection Prevalence: 0.11583
     Balanced Accuracy: 0.57027
       'Positive' Class : 1
```

- No Information Rage: 0.6966
- Accuracy: 0.7105

4. Decision Tree

Model Choice

ROCcurve & Lift Chart



• AUC:0.66

> performance(c5_pred,"auc")@y.values[[1]]
[1] 0.6592522

✓ 일반적으로덜정확한모델의값이출력되었다.0.5〈AUC〈=0.7

5. Conclusion

결론

- 느낀점
 - 구매 내역으로고객의성별을구분하는것에 있어서 많은 어려움이 있는 것으로 판단됨.
 - 남성 또는 여성의 구매 패턴이 비슷하며 뚜렷한 구분이 되는 데이터가 존재 하지 않는 것으로 판단됨.
 - CS데이터내그룹사및나이정보를가지고 모델을생성했을경우83%이상의 성별을매칭할수있었습니다.

