

Datasættets første parameter "customer\_id" kan udelades, da det er en identifikator, der ikke tjener noget formål i en model, der kan forudsige, om en bruger er konverteret eller ej.

Parameteren "credit\_account\_id" er en hash og vil derfor være en kategorisk variabel. Den har 148 unikke værdier, hvor de fleste opstår 1-4 gange, men en enkelt af værdierne kan findes i datasættet 687 gange. Det er umiddelbart svært at se, hvordan den skulle kunne bidrage til modellen, og parameteren ekskluderes derfor fra datasættet.

Det er fundet at parameteren "age" mangler 177 værdier. Ud af 891 rækker, er det en ret stor andel, og derfor er de manglende værdier i parameteren imputed. Da mean imputation vil gøre at 177 af rækkerne får samme alder, hvilket formegentlig ikke vil være retvisende i et datasæt hvor parameteren spænder mellem 0.42 og 80, er missForest valgt som imputation-metode, da den imputer værdier ved at kigge på hele datasættet.

Da target-parameteren er binær, laves der en logistisk regressionsmodel på datasættet. Et summary af modellen viser at parameterene "customer\_segment", "gender", "age" og "related\_customers" har lave p-værdier, hvilket antyder, at der er en stærk relation mellem disse parametre og hvorvidt en bruger er konverteret eller ej.

```
> model <- glm(dataset$converted ~ ., data = dataset, family = binomial(link = "logit"))
>
> summary(model)
```

Call:  
glm(formula = dataset\$converted ~ ., family = binomial(link = "logit"),  
data = dataset)

Deviance Residuals:

Min	1Q	Median	3Q	Max
-2.7572	-0.5950	-0.4023	0.6242	2.5094

Coefficients:

	Estimate	Std. Error	z value	Pr(> z )
(Intercept)	3.060e+01	6.042e+02	0.051	0.959602
customer_segment	-1.239e+00	1.516e-01	-8.171	3.06e-16 ***
gendermale	-2.668e+00	2.022e-01	-13.190	< 2e-16 ***
age	-4.857e-02	8.169e-03	-5.945	2.77e-09 ***
related_customers	-3.853e-01	1.111e-01	-3.468	0.000524 ***
family_size	-8.818e-02	1.211e-01	-0.728	0.466488
initial_fee_level	7.343e-04	1.181e-03	0.622	0.534218
branchHelsinki	-1.276e+01	6.042e+02	-0.021	0.983152
branchTampere	-1.239e+01	6.042e+02	-0.021	0.983642
branchTurku	-1.240e+01	6.042e+02	-0.021	0.983625

---  
Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Figure 1: Summary af den logistiske regressionsmodel, som viser, at modellen har en lav p-værdi for parametrene "customer\_segment", "gender", "age" og "related\_customers".

En ANOVA test, altså en analyse af variansen, viser, at forskellen mellem *null deviance* og *residual deviance*<sup>1</sup> er størst, når parameterne "customer\_segment", "gender", "age" og "related\_customers" tilføjes modellen.

```
> anova(model, test = "Chisq")
Analysis of Deviance Table

Model: binomial, link: logit

Response: dataset$converted

Terms added sequentially (first to last)
```

	Df	Deviance	Resid. Df	Resid. Dev	Pr(>Chi)
NULL			890	1186.66	
customer_segment	1	102.254	889	1084.40	< 2.2e-16 ***
gender	1	257.206	888	827.20	< 2.2e-16 ***
age	1	28.946	887	798.25	7.442e-08 ***
related_customers	1	19.861	886	778.39	8.328e-06 ***
family_size	1	0.485	885	777.90	0.4862
initial_fee_level	1	1.023	884	776.88	0.3117
branch	3	3.707	881	773.17	0.2948

---  
 Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Figure 2: ANOVA test af den logistiske regressionsmodel, som viser at modellens *residual deviance* falder mest, når parametrene "customer\_segment", "gender", "age" og "related\_customers" introduceres.

En anden metode er at kigge på forskellen i modellens forklaringsgrad. Hvis forklaringsgradens forøgelse er stor, når en parameter tilføjes, sammenlignet med en model, hvor alle parametre undtagen den undersøgte er inkluderet, så fortæller det, at parameteren er vigtig for modellen. Da det ikke er muligt at lave en standard  $R^2$ -mål på en logistisk regressionsmodel, bruges i stedet McFadden  $R^2$  indeks, som er den mest brugte erstatning herfor.

customer_segment	gender	age	related_customers	family_size	initial_fee_level	branch
0.05819323	0.1820019	0.03276349	0.012201	0.0004543226	0.0003408775	0

Figure 3: Tabel som viser, hvor meget McFadden  $R^2$  indekset stiger, når den enkelte parameter er tilføjet modellen, sammenlignet med en model, hvor den undersøgte parameter er holdt ude.

<sup>1</sup>Deviance er et "goodness-of-fit"-mål, som man bruger til statistiske modeller, som man ofte bruger i generaliserede lineære modeller, som en logistisk regression er.

Resultaterne bekræfter, hvad der blev fundet i ANOVA testen; at de vigtigste parametre for forudsigelse om en bruger er konverteret eller ej er "gender", "customer\_segment", "age" og "related\_customers" i den nævnte rækkefølge.

## Kode

```
1 setwd(dirname(rstudioapi::getActiveDocumentContext()$path))
2
3
4 library(dplyr)
5 library(pscl)
6 library(missMethods)
7 library(randomForest)
8 library(missForest)
9
10 #Load data and exclude the parameter "customer_id"
11 dataset <- read.csv("case_data.csv", header = T)
12 dataset <- as.data.frame(dataset[,2:10])
13
14 #make sure, that the target variable is a boolean
15 dataset[,1] <- as.factor(dataset[,1])
16
17 #Investigate the dataset
18 nrow(dataset)
19 sapply(dataset, n_distinct)
20
21 data.frame(table(dataset[,8]))
22 #Exclude the parameter "credit_account_id" according to the reasons
   mentioned in the report
23 dataset <- dataset[,-8]
24
25 #Handle missing data
26 colnames(dataset)[colSums(is.na(dataset)) > 0]
27 nrow(dataset[is.na(dataset$age),])
28 #dataset <- impute_mean(dataset, type = "columnwise")
29 set.seed(1234)
30 dataset <- missForest(dataset)$ximp
31
32 #Build the logistic regression model
33 model <- glm(dataset$convertd~., data = dataset, family = binomial(link
   = "logit"))
34
35 summary(model)
36 anova(model, test = "Chisq")
37
38 #Find the McFadden R^2 index for when all parameters are used
39 McFadden_full <- pR2(model)[4]
40
41 #Initiate a table for the increasements in the McFadden R^2 indexes
```

```
42 r_sqrd_increasement <- matrix(0, nrow=1, ncol = ncol(dataset)-1)
43 colnames(r_sqrd_increasement) <- colnames(dataset[,2:ncol(dataset)])
44
45 #Loop that build a logistic regression model where each of the parameters
   are excluded and then find the increasement i the McFadden R^2 index
46 for(i in 2:(ncol(dataset)-1)){
47   reduced_model <- glm(converted~., data = dataset[, -i], family =
     binomial(link = "logit"))
48   #print( McFadden_full - pR2(reduced_model)[4])
49   r_sqrd_increasement[1,i-1] <- McFadden_full - pR2(reduced_model)[4]
50 }
51
52 r_sqrd_increasement
```

case\_code.R