kaggle          🔍 Search                    Competitions   Datasets   Notebooks   Discussion   Courses   •••   🔔  🖼

Code                          This Notebook has been released under the Apache 2.0 open source license.          Download Code

```
 1   #### Introduction ####
 2   # In the past, the best submissions for this data mining challenge
 3   # (the challenge not the data set) have been simple logistic
 4   # regression models that use variables that make intuitive sense. It is absolutely paramount for a data
 5   # scientist to understand the domain the data is from. There are examples where 10 Phd guys build
 6   # some ridiculous NN and then there is this guy who regresses three variables out of the 100 or so
 7   # and comes up with a better model.
 8   # Having said that, I am lazy and intrigued at the same time. I don't want to dig into the data
 9   # too much, but just want to build a classifier for the hell of it. Classic mistake.
10   # There are two reasons why in this particular case this is a bad idea. First, the data is not
11   # very complex. Second, it is a very intuitive case. If these conditions are met, it is likely better
12   # to build a simple classifier, instead of going ham on the methods that are available, like
13   # building a random forest. This approach generally means you have no idea what you're dealing
14   # with and/or are just enjoying the academic challenge.
15   # This should be the lesson of this competition, KISS - Keep It Simple Stupid. Too bad it won't stick.
16   # So let's go down the rabbit hole.
17   require(Matrix)
18   require(data.table)
19
20   # read data
21   d <- read.csv('../input/carInsurance_train.csv')
22   # classify <- read.csv('~/Desktop/ML/DSS_DMC_2017/Working Material/test_set.csv')
23
24   ################################################################################
25   #### data exploration ####
26   ## should probably do outliers and data exploration
27   boxplot(d) # wow, balance. Could think about normalizing it and kicking the guy up top out. Let's check
28   # if balance correlates with the Car Insurance to be classified
29
30   cor.test(d$Balance, d$CarInsurance)
31   # pretty siginificant. What if we remove some "outliers"?
32   cor.test(d[d$Balance <= 20000,"Balance"], d$CarInsurance[d$Balance <= 20000])
33   # stops being significant. I am definitely overlooking some mathematical truth that would be apparent if I took the
34
35   ## let's look at the others without balance throwing everything in disarray
36   boxplot(d[,-c(1,7)]) # obviously we don't need ID either
37   # days passed still messed up
38
39   cor.test(d$DaysPassed, d$CarInsurance)
40   # pretty relevant
41
42   #### We'll leave it at that. It should give you an idea. Caret is a good resource as well http://topepo.github.io/c
43
44   ################################################################################
45   #### Cleaning ####
46
47   ### Cleaning has two parts. Getting rid of the usual, like handling NAs and dealing with nonsensical stuff (e.g. ID
48   ## My assumptions are that the day and time of day for the calls influence conversion rates.
49   # ******* should probably look at the columns for LastCallDay and LastCallMonth. They don't make sense in their cur
50   # Actually, finding out which weekday calls were made to what time would be very helpful. Because cold calls tend t
51   # We start with the 2nd part first, because NA imputation also somewhat depends on the steps we will take here. You
52
53   ## Get the weekdays in a workable format
54   # What are we doing here? We believe there is a chance that calls made at a certain time on certain days, e.g. Thur
55   # However, we don't know which year this data is from. Which does not matter. We just need to have a uniform distri
```

```r
56
57   clean_d <- d
58
59   max(clean_d[clean_d$LastContactMonth == 'feb', 'LastContactDay']) # check if the year the calls were made isn't a l
60   clean_d$DateCall <- as.Date(paste(clean_d$LastContactDay, clean_d$LastContactMonth, "2015", sep = '/'), "%d/%b/%Y")
61   clean_d$Weekday <- factor(weekdays(clean_d$DateCall))
62
63   # Next, we want to know what time people were called during the day.
64   # Let's see when the calls were made and what the working hours are.
65   plot(table(clean_d$CallStart)) # not very informative let's take the minutes and seconds off
66   plot(table(call_hr <- gsub("(:\\d{2})", "", clean_d$CallStart))) # ok... they are pretty diligent in calling people
67
68   # We could take the times as they are given. However, that would be too much noise, in my opinion. Therefore, I opt
69   clean_d$CallDayTime <- as.numeric(gsub("(:\\d{2})", "", clean_d$CallStart))
70   require(car)
71   clean_d$CallDayTime <- factor(recode(clean_d$CallDayTime, "c('9', '10', '11')='morning'; c('12', '13', '14')='midda
72
73
74   ##
75   # Convert last phone call time to minutes
76   require(chron)
77   clean_d$call_dur_min <- 60 * 24 * as.numeric(times(clean_d$CallEnd)-times(clean_d$CallStart))
78
79   ## find NAs
80   # I didn't come up with that. Google + Stackoverflow = Bliss
81   na_count <-sapply(clean_d, function(y) sum(length(which(is.na(y)))))
82   na_count <- data.frame(na_count)
83
84   ## The issue is that we have a lot of NAs. Question is, what to do with them? Kick the missing values? Include them
85   # But most classifiers can't handle missing values. So we should replace them. For a more holistic
86   # overview of the how's and why's read the article above. We are going to go with the cool solution.
87   # We replace NAs with factors falling out of a k-Nearest-Neighbor
88
89   ## Before we can do that, we need to take away all the columns that we won't consider in our analysis. Reason being
90   summary(clean_d)
91   # ID is the first to go, redundant a.f.. Outcome probably too much noise, has to go. CallStart and CallEnd, LastCon
92   # Kick everything we don't need. This is a little complicated in R if you want to do it by column name. We do it an
93
94   sub_clean_d <- subset(clean_d, select = -c(Id, LastContactDay, LastContactMonth, Outcome, CallStart, CallEnd, DateC
95
96   # We can do the NA replacement now
97   ## Impute NAs https://www.r-bloggers.com/missing-value-treatment/
98   # this might leave us in Hell's Kitchen. Well... actually it's Freedman's Kitchen https://www.r-bloggers.com/freedm
99
100  require(DMwR)
101  set.seed(42)
102
103  clean_d_imputeknn <- knnImputation(sub_clean_d)  # perform knn imputation.
104  # check if it missed any NAs
105  anyNA(clean_d_imputeknn)
106
107
108  #################################################################################
109  #### Building the actual model ####
110  ## We start with a random forest.
111  ## Using the caret package we cross validate before we train the model
112  # Cross validation is pretty much the most imortant thing, because it reduces overfitting. It means we partition th
113  # Basically with the caret package, we define a cross validation object which we pass to the actual function to cre
114
115  ## Making a nice data frame for the modelling operations.
116  model_d <- clean_d_imputeknn
117
118
119  # Split data into test and train: so later we can validate our model
120  require(caret)
121  set.seed(42)
122
123  # We will use this index to define the rows, which will go into train and test data sets respectively
124  train_index <- createDataPartition(model_d$CarInsurance, p = 0.75, list = FALSE, times = 1)
125  training <- model_d[train_index, ]
```

## Simple Random Forest on Insurance Call Forecast

R script using data from Car Insurance Cold Calls · 985 views · 2y ago

Copy and Edit   6

```r
130
131
132   #### Random Forest ####
133   # Training the actual model. We have to pass CarInsurance, i.e. the variable to be classified, as a factor, not as
134   set.seed(42)
135   model_rf = train(factor(CarInsurance)~., data=training, trControl=train_control, method="rf")
136
137   # We make a frame and fill it with the predicted values. This allows us to the test the quality of the model in the
138   prediction_rf = predict(model_rf, subset(testing, select=-c(CarInsurance)))
139
140   #Compute the accuracy of predictions with a confusion matrix
141   confusionMatrix(prediction_rf, testing$CarInsurance)
142
143   #### Logistic Regression ####
144   ## simple first, then bagging and boosting
145   set.seed(42)
146   model_logreg <- glm(factor(CarInsurance) ~., family=binomial(link='logit'), data=training)
147
148   prediction_logreg = predict(model_logreg, subset(testing, select=-c(CarInsurance)), type='response') # by choosing
149
150   table(testing$CarInsurance, prediction_logreg > 0.5) # LogReg gives the results as probabilities, so we can't use t
151
152   ## so Random Forest results are actually better. However, the RF takes about a minute or so to calculate the model.
153
154   #### LogitBoost ####
155   set.seed(42)
156   model_logitboost <- train(factor(CarInsurance)~., data=training, trainControl=train_control, method="LogitBoost", n
157
158   prediction_logitboost = predict(model_logitboost, subset(testing, select=-c(CarInsurance)))
159   confusionMatrix(prediction_logitboost, testing$CarInsurance)
160
161   #### XGB Trees ####
162   set.seed(42)
163   model_xgbtrees <- train(factor(CarInsurance)~., data=training, method='xgbTree', trainControl=train_control, metric
164
165   prediction_xgbtrees = predict(model_xgbtrees, subset(testing, select= -c(CarInsurance)))
166   confusionMatrix(prediction_xgbtrees, testing$CarInsurance) # worst result yet :(
```

**Did you find this Notebook useful?**
Show your appreciation with an upvote

0

## Run Info

| | | | |
|---|---|---|---|
| Succeeded | False | Run Time | 1200.5 seconds |
| Exit Code | 137 | Queue Time | 0 seconds |
| Docker Image Name | kaggle/rstats (Dockerfile) | Output Size | 0 |
| Timeout Exceeded | True | Used All Space | False |
| Failure Message | The kernel was killed for running longer than 1200 seconds. | | |

## Log

Download Log

| Time | Line # | Log Message |
|---|---|---|

| | | Code | Log | Data | Output | Comments |

```
3.1s          4                Pearson's product-moment correlation

                  data:  d$Balance and d$CarInsurance
3.1s          5   t = 2.6302, df = 3998, p-value = 0.008567
                  alternative hypothesis: true correlation is not equal to 0
                  95 percent confidence interval:
                   0.01058319 0.07245913
                  sample estimates:
3.1s          6         cor
                  0.04156101


                                Pearson's product-moment correlation

                  data:  d[d$Balance <= 20000, "Balance"] and d$CarInsurance[d$Balance <= 20000]
3.1s          7   t = 4.2517, df = 3974, p-value = 2.17e-05
                  alternative hypothesis: true correlation is not equal to 0
                  95 percent confidence interval:
                   0.03628313 0.09817159
                  sample estimates:
                         cor
                  0.06729209


3.2s          8                Pearson's product-moment correlation

                  data:  d$DaysPassed and d$CarInsurance
3.2s          9   t = 8.8714, df = 3998, p-value < 2.2e-16
                  alternative hypothesis: true correlation is not equal to 0
                  95 percent confidence interval:
                   0.1084183 0.1692057
                  sample estimates:
                         cor
                  0.1389429

                  [1] 27
3.3s         10   Loading required package: car
4.1s         11   Loading required package: chron
4.5s         12         Id              Age              Job            Marital
                  Min.   :   1    Min.   :18.00    management :893   divorced: 483
                  1st Qu.:1001    1st Qu.:32.00    blue-collar:759   married :2304
                  Median :2000    Median :39.00    technician :660   single  :1213
                  Mean   :2000    Mean   :41.21    admin.     :459
                  3rd Qu.:3000    3rd Qu.:49.00    services   :330
4.5s         13   Max.   :4000    Max.   :95.00    (Other)    :880
                                                   NA's       : 19
                     Education          Default           Balance          HHInsurance
                  primary  : 561    Min.   :0.0000    Min.   :-3058.0    Min.   :0.0000
                  secondary:1988    1st Qu.:0.0000    1st Qu.:  111.0    1st Qu.:0.0000
                  tertiary :1282    Median :0.0000    Median :  551.5    Median :0.0000
                  NA's     : 169    Mean   :0.0145    Mean   : 1532.9    Mean   :0.4928
                                    3rd Qu.:0.0000    3rd Qu.: 1619.0    3rd Qu.:1.0000
                                    Max.   :1.0000    Max.   :98417.0    Max.   :1.0000

                     CarLoan          Communication     LastContactDay    LastContactMonth
                  Min.   :0.000    cellular :2831    Min.   : 1.00     may    :1049
                  1st Qu.:0.000    telephone: 267    1st Qu.: 8.00     jul    : 573
                  Median :0.000    NA's     : 902    Median :16.00     aug    : 536
                  Mean   :0.133                      Mean   :15.72     jun    : 454
                  3rd Qu.:0.000                      3rd Qu.:22.00     nov    : 347
                  Max.   :1.000                      Max.   :31.00     apr    : 306
                                                                       (Other): 735
                     NoOfContacts       DaysPassed        PrevAttempts       Outcome
                  Min.   : 1.000    Min.   : -1.00    Min.   : 0.0000    failure: 437
                  1st Qu.: 1.000    1st Qu.: -1.00    1st Qu.: 0.0000    other  : 195
                  Median : 2.000    Median : -1.00    Median : 0.0000    success: 326
                  Mean   : 2.607    Mean   : 48.71    Mean   : 0.7175    NA's   :3042
                  3rd Qu.: 3.000    3rd Qu.: -1.00    3rd Qu.: 0.0000
                  Max.   :43.000    Max.   :854.00    Max.   :58.0000

                     CallStart         CallEnd          CarInsurance       DateCall
                  10:42:44:   3    10:22:30:   3    Min.   :0.000     Min.   :2015-01-08
                  11:48:25:   3    10:52:24:   3    1st Qu.:0.000     1st Qu.:2015-05-08
                  13:54:34:   3    11:27:46:   3    Median :0.000     Median :2015-06-05
                  15:27:56:   3    09:04:02:   2    Mean   :0.401     Mean   :2015-06-21
                  15:48:27:   3    09:06:42:   2    3rd Qu.:1.000     3rd Qu.:2015-08-11
                  17:02:39:   3    09:12:47:   2    Max.   :1.000     Max.   :2015-12-30
                  (Other) :3982    (Other) :3985
                     Weekday          CallDayTime        call_dur_min
                  Friday   :725    afternoon:1351    Min.   : 0.08333
                  Monday   :454    midday   :1365    1st Qu.: 2.10000
                  Saturday :380    morning  :1284    Median : 3.86667
                  Sunday   : 92                      Mean   : 5.84740
                  Thursday :849                      3rd Qu.: 7.66667
                  Tuesday  :724                      Max.   :54.21667
                  Wednesday:776
```

```
4.5s      14   Loading required package: DMwR
4.6s      15   Loading required package: methods
4.7s      16   Loading required package: lattice
4.7s      17   Loading required package: grid
45.8s     18   [1] FALSE
45.8s     19   Loading required package: caret
45.8s     20   Loading required package: ggplot2
46.2s     21   Loading required package: randomForest
46.2s     22   randomForest 4.6-12
               Type rfNews() to see new features/changes/bug fixes.

46.3s     23
               Attaching package: 'randomForest'

               The following object is masked from 'package:ggplot2':

                   margin

197.6s    24   Confusion Matrix and Statistics

                         Reference
               Prediction   0    1
                        0  483   79
                        1  106  332
197.6s    25
                         Accuracy : 0.815
                           95% CI : (0.7895, 0.8386)
               No Information Rate : 0.589
               P-Value [Acc > NIR] : < 2e-16

                            Kappa : 0.6216
               Mcnemar's Test P-Value : 0.05593

                      Sensitivity : 0.8200
                      Specificity : 0.8078
                   Pos Pred Value : 0.8594
                   Neg Pred Value : 0.7580
                       Prevalence : 0.5890
                   Detection Rate : 0.4830
             Detection Prevalence : 0.5620
                Balanced Accuracy : 0.8139

                   'Positive' Class : 0

197.7s    26
                      FALSE  TRUE
                  0    499    90
                  1    117   294
198.3s    27   Loading required package: caTools
211.5s    28   Confusion Matrix and Statistics

                         Reference
               Prediction   0    1
                        0  506  147
                        1   83  264
211.5s    29
                         Accuracy : 0.77
                           95% CI : (0.7426, 0.7958)
               No Information Rate : 0.589
               P-Value [Acc > NIR] : < 2.2e-16

                            Kappa : 0.5135
               Mcnemar's Test P-Value : 3.266e-05

                      Sensitivity : 0.8591
                      Specificity : 0.6423
                   Pos Pred Value : 0.7749
                   Neg Pred Value : 0.7608
                       Prevalence : 0.5890
                   Detection Rate : 0.5060
             Detection Prevalence : 0.6530
                Balanced Accuracy : 0.7507

                   'Positive' Class : 0

211.6s    30   Loading required package: xgboost
211.6s    31   Loading required package: plyr
211.6s    32
               Attaching package: 'plyr'

               The following object is masked from 'package:DMwR':

                   join

211.6s    33
211.6s    35   Failed. Exited with code 137.
```

Data

## Data Sources

- ⌄ 🗁 Car Insurance Cold Calls
  - ⊞ carInsurance_test.csv                    19 columns
  - ⊞ carInsurance_train.csv                   19 columns
  - 🗎 DSS_DMC_Description.pdf

## Car Insurance Cold Calls

**We help the guys and girls at the front to get out of Cold Call Hell**

Last Updated: 2 years ago (Version 1)

**About this Dataset**

## Introduction

Here you find a very simple, beginner-friendly data set. No sparse matrices, no fancy tools needed to understand what's going on. Just a couple of rows and columns. Super simple stuff. As explained below, this data set is used for a competition. As it turns out, this competition tends to reveal a common truth in data science: KISS - Keep It Simple Stupid

What is so special about this data set is, given it's simplicity, it pays off to use "simple" classifiers as well. This year's competition was won by a C5.0 . Can you do better?
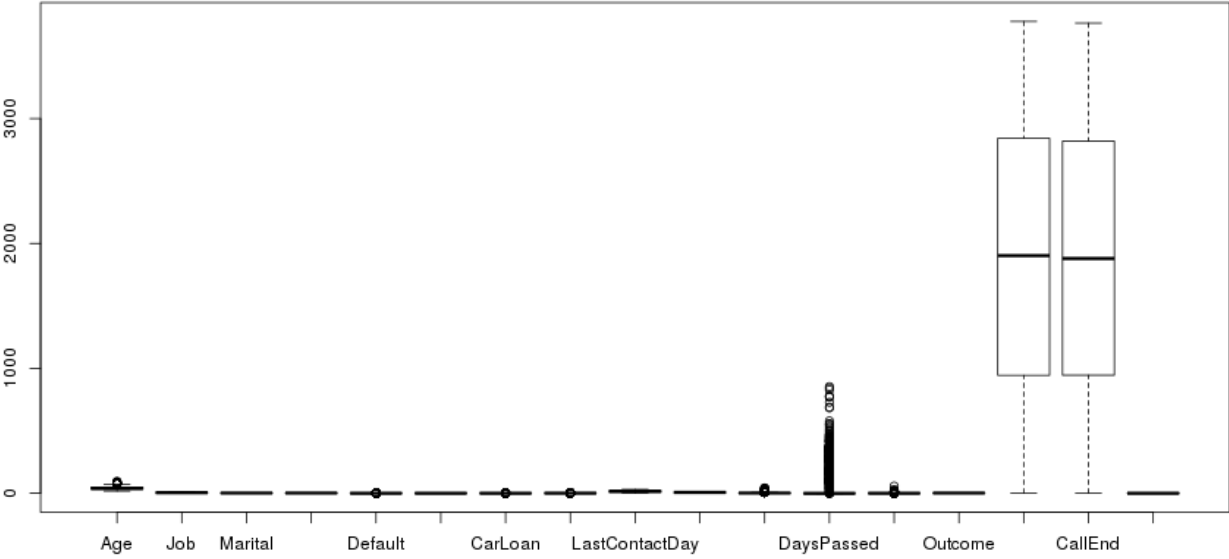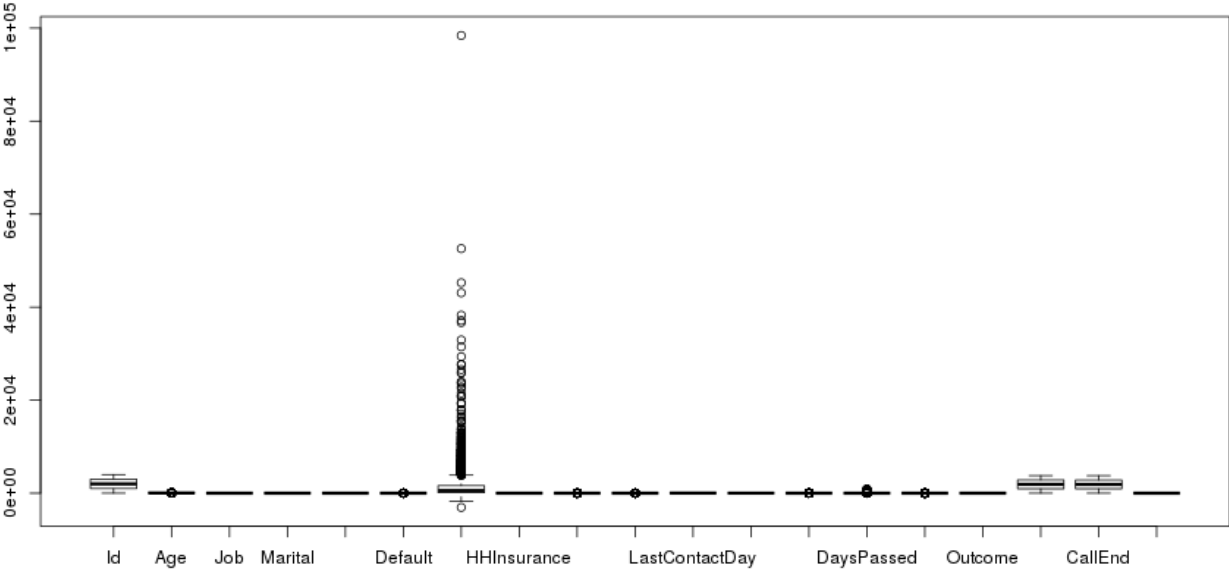
## Description

We are looking at cold call results. Turns out, same salespeople called existing insurance customers up and tried to sell car insurance. What you have are details about the called customers. Their age, job, marital status, whether the have home insurance, a car loan, etc. As I said, super simple.

What I would love to see is some of you applying some crazy XGBoost classifiers, which we can square off against some logistic regressions. It would be curious to see what comes out on top. Thank you for your time, I hope you enjoy using the data set.

## Acknowledgements

Thanks goes to the Decision Science and Systems Chair of Technical University of Munich (TUM) for getting the data set

Output Visualizations

Output visualization 3

## Comments (0)

Click here to comment...

Our Team   Terms   Privacy   Contact/Support