

model_CA_LM.R

atchirc

Mon May 22 16:49:40 2017

```
library(MASS)
library(car)
library(DataCombine)    # Pair wise correlation
library(stargazer)
library(dplyr)          # Data aggregation
library(glmnet)
source('../atchircUtils.R')

data    <- read.csv('../intrim/eleckart.csv')

# KPI selection
# units, product_mrp, list_mrp, COD, Prepaid are factors
# Insig : Affiliates corr OnlineMarketing
# Insig : Radio corr Other
# Insig : Digital, ContentMarketing corr SEM
# delivery(b/c)days are corr, lets choose deliverydays
# will use marketing levers rather TotalInvestment

# Filter significant KPIs
model_data <- subset(data, product_analytic_sub_category=='CameraAccessory',
  select = -c(product_analytic_sub_category,product_mrp,
    units,COD,Prepaid,deliverybdays,
    TotalInvestment,Affiliates,Radio,Digital,
    ContentMarketing,sla,procurement_sla))

model_data_org <- model_data
model_data[,c(8:12)] <- model_data[,c(8:12)]*10000000

# # *****
# #           FEATURE ENGINEERING -PASS2 ----
# # *****
#
# # . . . . List Price Inflation ----
model_data$chnghlist <- c(0,diff(model_data$list_mrp))
#
# # . . . . Discount Inflation ----
model_data$chnghdisc <- c(0,diff(model_data$discount))
#
```

*

****PROCs:****

Linear, Ridge and Lasso Model are wrapped with abstract functions. This would facilitate readable code for model building and Model optimization. Set Class definitions

```
setOldClass('elnet')
setClass(Class = 'atcglmnet',
  representation (
    R2 = 'numeric',
    mdl = 'elnet',
    pred = 'matrix'
  )
)
```

```
setOldClass('lm')
setClass(Class = 'atclm',
  representation (
    R2 = 'numeric',
    mdl = 'lm',
    pred = 'matrix'
  )
)
```

Finding min lambda from 1000 iterations Function to find Min Lambda using bootstrap method. minlambda identified over 1000 cross validation trails. observed minlambda used for Ridge and Lasso regression.

```
findMinLambda <- function(x,y,alpha,folds) {
  lambda_list <- list()
  for (i in 1:1000) {
    cv.out <- cv.glmnet(as.matrix(x), as.vector(y), alpha=alpha,
                       nfolds=folds)
    lambda_list <- append(lambda_list, cv.out$lambda.min)
  }
  return(min(unlist(lambda_list)))
}
```

Linear Model with Regularization Wrapper function for Ridge and Lasso regression. functions performs Ridge/Lasso regression and returns R2, Model and Predicted values as `atcglmnet` object

```
atcLmReg <- function(x,y,l1l2,folds) {
  # l1l2 = 0 for L1, 1 for L2

  if (l1l2) { # Lasso/L2
    min_lambda <- findMinLambda(x,y,1,folds)
  } else { # Ridge/L1
    min_lambda <- findMinLambda(x,y,0,folds)
  }
  mdl <- glmnet(x,y,alpha=l1l2,lambda = min_lambda)
```

```

pred      <- predict(mdl,s= min_lambda,newx=x)

# MSE
mean((pred-y)^2)
R2 <- 1 - (sum((y-pred )^2)/sum((y-mean(pred))^2))
return(new('atcglmnet', R2 = R2, mdl=mdl, pred=pred))
}

```

*

MODELING

```
# Prune KPI as part of model optimization
model_data <- na.omit(model_data)
model_data <- subset(model_data,select=-c(TV,SEM,discount))
```

Linear Model:

```
mdl <- lm(gmv~., data=model_data)
step_mdl <- stepAIC(mdl,direction = 'both',trace = FALSE)

stargazer(mdl,step_mdl, align = TRUE, type = 'text',
           title='Linear Regression Results', single.row=TRUE)
```

```
##
## Linear Regression Results
## =====
##                               Dependent variable:
##                               -----
##                               gmv
##                               (1)                (2)
## -----
## week                -6,523.196 (34,052.850)
## deliverycdays        266,593.100 (272,322.800)    187,512.800 (134,424.900)
## n_saledays           261,619.300 (163,878.600)    247,315.300 (156,402.500)
## Sponsorship          168,460.100** (66,427.240)    145,365.700*** (48,968.620)
## OnlineMarketing        0.034 (0.033)                0.035** (0.015)
## Other                 0.013 (0.017)
## NPS                   0.003 (0.018)
## list_mrp              0.0004** (0.0002)            0.0004*** (0.0001)
## chnglist              -0.00002 (0.0001)
## chngdisc              48,866.430 (29,654.780)    48,169.110* (28,188.840)
## Constant             -3,571,627.000 (11,186,860.000) -1,508,753.000 (1,176,816.000)
## -----
## Observations                52                52
## R2                          0.607                0.602
## Adjusted R2                  0.512                0.549
## Residual Std. Error    1,729,251.000 (df = 41)    1,662,414.000 (df = 45)
## F Statistic              6.342*** (df = 10; 41)    11.330*** (df = 6; 45)
## =====
## Note:                                *p<0.1; **p<0.05; ***p<0.01
```

```
knitr::kable(viewModelSummaryVIF(step_mdl))
```

var	Estimate	Std.Error	t-value	Pr(> t)	Significance	vif
chngdisc	4.817e+04	2.819e+04	1.709	0.094377	.	1.027431
deliverycdays	1.875e+05	1.344e+05	1.395	0.169886	NA	1.119227
list_mrp	3.709e-04	1.050e-04	3.532	0.000966	***	1.252244
n_saledays	2.473e+05	1.564e+05	1.581	0.120819	NA	1.133652
OnlineMarketing	3.529e-02	1.518e-02	2.325	0.024668	*	1.465913
Sponsorship	1.454e+05	4.897e+04	2.969	0.004782	**	1.444867

```
pred_lm <- predict(step_mdl, model_data)
```

Regularized Linear Model:

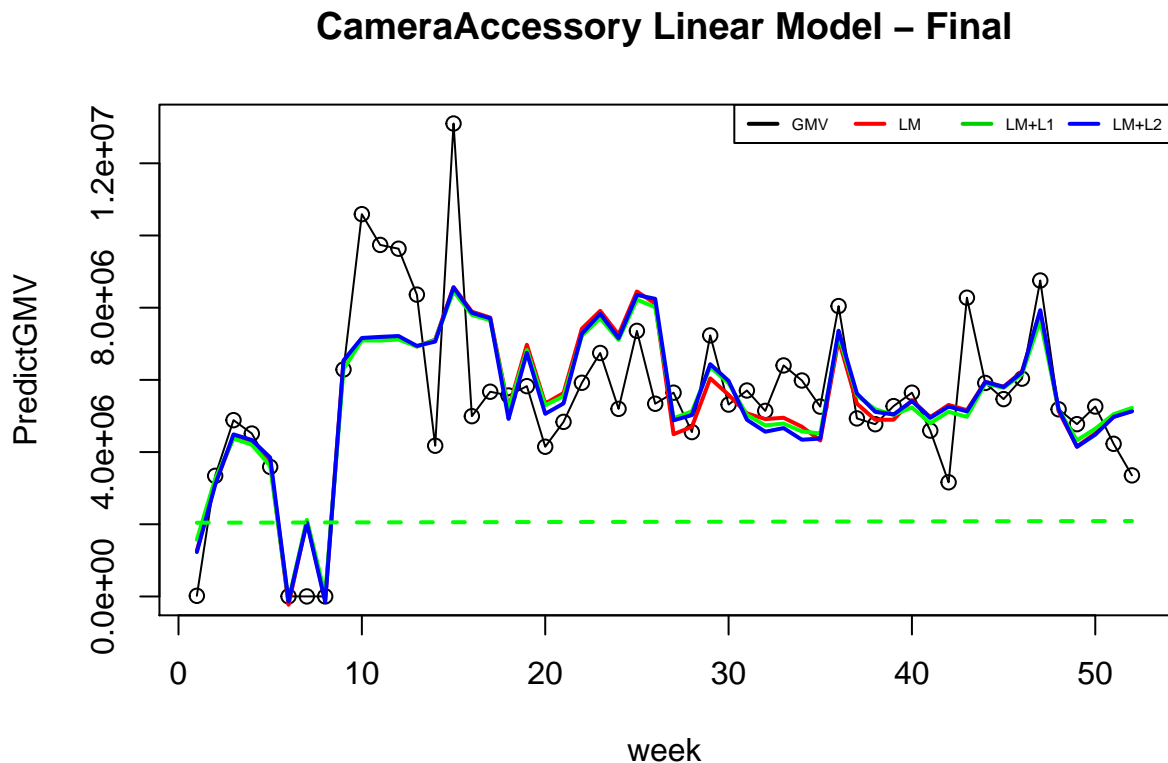
```
x = as.matrix(subset(model_data, select=-gmv))  
y = as.vector(model_data$gmv)  
  
ridge_out <- atcLmReg(x,y,0,3) # x, y, alpha, nfolds  
lasso_out <- atcLmReg(x,y,1,3) # x, y, alpha, nfolds
```

*

PLOTTING MODEL RESULTS

Plot Model prediction and base sales:

```
plot(model_data$gmv, main = 'CameraAccessory Linear Model - Final',
     xlab='week', ylab='PredictGMV')
lines(model_data$gmv)
lines(pred_lm, col='red', lwd=2)
lines(ridge_out@pred, col='green', lwd=2)
lines(lasso_out@pred, col='blue', lwd=2)
lines(step_mdl$coefficients['(Intercept)'] + step_mdl$coefficients['week'] * model_data$week,
     lty=2, lwd=2, col='red')
lines(ridge_out@mdl$a0 + ridge_out@mdl$beta['week', 1] * model_data$week,
     lty=2, lwd=2, col='green')
lines(lasso_out@mdl$a0 + lasso_out@mdl$beta['week', 1] * model_data$week,
     lty=2, lwd=2, col='blue')
legend('topright', inset=0, legend=c('GMV', 'LM', 'LM+L1', 'LM+L2'), horiz = TRUE,
     lwd = 2, col=c(1:4), cex = 0.5)
```



*

*Model Coefficients:**

```
coeff_lm <- as.data.frame(as.matrix(coef(step_md1)))
coeff_l1 <- as.data.frame(as.matrix(coef(ridge_out@mdl)))
coeff_l2 <- as.data.frame(as.matrix(coef(lasso_out@mdl)))
```

```
lm_df=data.frame('x'=rownames(coeff_lm),'y'=coeff_lm)
colnames(lm_df) = c('coeff','lm')
l1_df=data.frame('x'=rownames(coeff_l1),'y'=coeff_l1)
colnames(l1_df)= c('coeff','l1')
l2_df=data.frame('x'=rownames(coeff_l2),'y'=coeff_l2)
colnames(l2_df) <- c('coeff','l2')
```

```
smry <- merge(lm_df,l1_df,all = TRUE)
smry <- merge(smry,l2_df,all=TRUE)
```

```
print(smry)
```

##		coeff	lm	l1	l2
## 1	(Intercept)	-1.508753e+06	2.042501e+06	-2.730584e+06	
## 2	chnngdisc	4.816911e+04	4.513849e+04	4.859238e+04	
## 3	chnnglist		NA	7.272228e-06	-1.849719e-05
## 4	deliverycdays	1.875128e+05	1.656200e+05	2.483629e+05	
## 5	list_mrp	3.709176e-04	3.277707e-04	3.937760e-04	
## 6	n_saledays	2.473153e+05	2.318695e+05	2.589813e+05	
## 7	NPS		NA	-5.349811e-03	1.793655e-03
## 8	OnlineMarketing	3.529442e-02	2.337185e-02	3.158909e-02	
## 9	Other		NA	7.735092e-03	1.239416e-02
## 10	Sponsorship	1.453657e+05	1.416823e+05	1.660920e+05	
## 11	week		NA	9.151802e+02	-4.611292e+03

```
print(paste0('Ridge regression R2 : ',ridge_out@R2))
```

```
## [1] "Ridge regression R2 : 0.602792651908966"
```

```
print(paste0('Lasso regression R2 : ',lasso_out@R2))
```

```
## [1] "Lasso regression R2 : 0.607243602040244"
```

```
print(paste0('Linear Mode R2 : ',getModelR2(step_md1)))
```

```
## [1] "Multiple R-squared: 0.6017,\tAdjusted R-squared: 0.5486 "
```

```
## [1] "Linear Mode R2 : Multiple R-squared: 0.6017,\tAdjusted R-squared: 0.5486 "
```

*

Significant KPI

Lasso(LM+L1) regression results a simple explainable model with significant KPIs as Discount Inflation, Deliverycday, sale days, Sponsorship Discount, week, NPS

```
# Model Optimization
```

```
# coeff      lm      l1      l2
# 1      (Intercept) -4.205266e+06  3.743013e+06 -2.335133e+06
# 2      chngdisc      NA  3.544890e+04  2.297922e+04
# 3      chnglist      NA  1.274977e-05 -2.125097e-06
# 4      deliverycdays      NA  1.399561e+05  9.078950e+04
# 5      discount  6.485938e+04  6.976909e+03  2.857188e+04
# 6      list_mrp  3.520229e-04  2.898529e-04  3.339852e-04
# 7      n_saledays  2.494251e+05  2.376959e+05  2.589315e+05
# 8      NPS      NA -8.022442e-03  0.000000e+00
# 9      OnlineMarketing  4.147731e-02  2.946905e-02  4.207859e-02
# 10     Other      NA  6.919302e-03  1.216733e-02
# 11     SEM -5.362909e-02 -3.241843e-02 -4.862319e-02
# 12     Sponsorship  2.619984e+05  2.082814e+05  2.920367e+05
# 13     TV      NA -1.952227e+05 -5.558398e+05
# 14     week      NA -6.411466e+03 -1.947268e+03
# [1] "Ridge regression R2 : 0.635910648911486"
# [1] "Lasso regression R2 : 0.648390286764186"
# [1] "Multiple R-squared:  0.6301,\tAdjusted R-squared:  0.5808 "
# [1] "Linear Mode      R2 :
#      Multiple R-squared:  0.6301,\tAdjusted R-squared:  0.5808 "
```

```
# coeff      lm      l1      l2
# 1      (Intercept) -1.508753e+06  2.042501e+06 -2.805416e+06
# 2      chngdisc  4.816911e+04  4.513849e+04  4.861658e+04
# 3      chnglist      NA  7.272228e-06 -1.891686e-05
# 4      deliverycdays  1.875128e+05  1.656200e+05  2.500288e+05
# 5      list_mrp  3.709176e-04  3.277707e-04  3.945720e-04
# 6      n_saledays  2.473153e+05  2.318695e+05  2.592140e+05
# 7      NPS      NA -5.349811e-03  1.912645e-03
# 8      OnlineMarketing  3.529442e-02  2.337185e-02  3.179100e-02
# 9      Other      NA  7.735092e-03  1.246340e-02
# 10     Sponsorship  1.453657e+05  1.416823e+05  1.662989e+05
# 11     week      NA  9.151802e+02 -4.788093e+03
# [1] "Ridge regression R2 : 0.602792651908966"
# [1] "Lasso regression R2 : 0.607259895391884"
# [1] "Multiple R-squared:  0.6017,\tAdjusted R-squared:  0.5486 "
# [1] "Linear Mode      R2 : Multiple R-squared:  0.6017,\tAdjusted R-squared:  0.5486 "
# >
```