

# Challenge 2

Code ▾

Hide

```
# Damit setzen wir das Working Directory auf den Ordner dieser Datei
if (!is.null(parent.frame(2)$ofile)) {
  this.dir <- dirname(parent.frame(2)$ofile)
  setwd(this.dir)}
```

## Get Data

Hide

```
library(data.table)

machines = fread(file = "machine_data.csv",
                 sep = ',',
                 header = TRUE,
                 stringsAsFactors = FALSE,
                 dec = '.')

products = fread(file = "product_data.csv",
                 sep = ',',
                 header = TRUE,
                 stringsAsFactors = FALSE,
                 dec = '.')

transaction = fread(file = "transactional_data.csv",
                   sep = ',',
                   header = TRUE,
                   stringsAsFactors = FALSE,
                   dec = '.')

mfailures = fread(file = "machine_failures.csv",
                  sep = ',',
                  header = TRUE,
                  stringsAsFactors = FALSE,
                  dec = '.')

mfailures = data.table(mfailures)
transaction = data.table(transaction)
products = data.table(products)
machines = data.table(machines)
```

### 1. Merge the transactional dataset with the machine failures data set setting failure variable to 0 when no failure is recorded

Hide

```
# merge the datasets
transactional = merge(transaction, mfailures, by = c('machine', 'timestamp', 'column'), all.x=T)
setDT(transactional)
# set the failure value of normal records to 0
transactional[is.na(failure), failure:=0]
```

**2. In the transactional data table, create a variable called “last\_vend” containing the timestamp of the previous sale of each machine**

[Hide](#)

```
#order by machine and date
transactional = transactional[order(machine,date)]

#create "last stand colum" with shift (one down):
transactional[, last_vend := shift(timestamp)]

# Empty first row for each machine since first selling point doesnt have a previous sale:
transaction[, last_vend := c(NA, timestamp[-.N]), by=machine]
summary(transactional)
```

**3. Create a new variable in the transactional data table called “deltahours” containing, for every sale, the hours that passed since the last sale**

[Hide](#)

```
help("difftime")
transactional[, deltahours := as.numeric(difftime(timestamp,last_vend,units = 'hours'
))]
# Only rows without NAs
transactional = transactional[complete.cases(transactional)]
head(transactional)
```

**4. Create an auxiliary data table called “machine\_daily\_average” 2ith the average daily sales per machine. Use this auxiliary table to attach to every row of the transactional data table to the average daily sales per machine.**

[Hide](#)

```
machine_daily_average = transactional[, .(daily_sales_machine=length(product_name)/uniqueN(date)),by=machine]
transactional = merge(transactional,machine_daily_average, by = 'machine', all.x=T)
head(transactional)
```

**5. Create a new variable called “delta” in the transactional data table containing a normalized version of deltahours consisting on the deltahours associated with each sale divided by the average deltahours of each machine i.e.  $\text{delta} = \text{deltahours} / (24/\text{daily\_sales\_machine})$ . The interpretation of delta is the amount of “missed sales” if the machine was selling at a constant rate**

[Hide](#)

```
transactional$delta <- transactional[, .(delta=as.numeric(deltahours*daily_sales_machine/24, units='hours'))]
head(transactional)

# Checking if high deltahours have a strong correlation with failure:
high_deltas = mean(transactional$failure[transactional$deltahours >
                                         unname(quantile(transactional$deltahours,
0.99,na.rm = TRUE))],na.rm = TRUE)

low_deltas = mean(transactional$failure[transactional$deltahours <
                                         unname(quantile(transactional$deltahours,0.
99, na.rm = TRUE))], na.rm = TRUE)

influence_deltas = data.table(high_deltas,low_deltas)
influence_deltas
```

Through this correlation analysis we intended to identify the relationship between the number of failures of a machine and lower sales. The higher correlation coefficient of machines with higher number of failures confirmed that they tend to have lower sales (which makes intuitively completely sense)

**6. Select 30% of the machines in the transactional data for testing and 70% of the machines for training and train a linear logistic regression model called “m” to predict whether a machine has a failure as a function of variable delta. What is the value of the intercept and the coefficient accompanying variable delta?**

[Hide](#)

```
set.seed(42)
idx <- sample(x=unique(transactional$machine), size=round(0.7*length(unique(transactional$machine))), 0))

train <- transactional[machine %in% idx, ]
test <- transactional[!machine %in% idx, ]

m <- glm(failure ~ delta, train, family="binomial")
summary(m)
```

**Answer:** Intercept : -6.911617; coef,delta: 0.562049

**a) What's the AUC, a measure of quality, of the model you have built on the train set? and on test set?**

[Hide](#)

```
library(pROC)
pred_train = predict(m, newdata = train, type = "response")
pred_test = predict(m, newdata = test, type = "response")

auc(train$failure, pred_train)
auc(test$failure, pred_test)
```

**b) Plot the function of probability of failure with respect to delta to gain intuition:**

[Hide](#)

```

intercept = m$coefficients["(Intercept)"]
coefficient = m$coefficients["delta"]
f = function(x){(1/(1+exp(-(intercept + (coefficient*x))))))}
curve(f, 0, 20,
      xlab = "Delta", ylab="Prob of Failure", col="red",
      main = "Probability of Failure with respect to Delta")

```

**c) Let us create alarms with two levels of priority: med-risk and high-risk. Med-risk alarms will fire when the probability of failure is  $\geq 60\%$  and High-risk when that probability is  $\geq 80\%$ .**

**i. What are the threshold deltas for each type of alarm to fire? #Hint (check the case)**

[Hide](#)

```

# Can be solved using the function and log it:
log_f = function(x){
  (log((1/x)-1) + intercept) / -coefficient}

threshold_med_risk = log_f(0.6)
threshold_high_risk = log_f(0.8)

threshold_med_risk
threshold_high_risk

```

**ii. How many of these alarms would be fired per day on average according to your model?**

[Hide](#)

```

daily_alarm_med_risk = length(transactional$delta[transactional$delta >= threshold_med_risk]) / length(unique(transactional$date))

daily_alarm_high_risk = length(transactional$delta[transactional$delta >= threshold_high_risk]) / length(unique(transactional$date))

daily_alarm_med_risk

daily_alarm_high_risk

```

**iii. What % of these will be “false alarms” i.e. failure variable is equal to 0, for each level of priority?**

[Hide](#)

```

transactional$alarm_med_risk = 0
transactional$alarm_med_risk[transactional$delta >= threshold_med_risk] = 1

transactional$alarm_high_risk = 0
transactional$alarm_high_risk[transactional$delta >= threshold_high_risk] = 1

false_alarm_rate_med_risk = sum(transactional$alarm_med_risk[transactional$failure == 0]) / sum(transactional$alarm_med_risk)
false_alarm_rate_high_risk = sum(transactional$alarm_high_risk[transactional$failure == 0]) / sum(transactional$alarm_high_risk)

sprintf("%1.2f%%", 100*false_alarm_rate_med_risk)
sprintf("%1.2f%%", 100*false_alarm_rate_high_risk)

```

**d) In this exercise we will estimate the profit impact of our EWS system vs the current system:**

**i. If we set the EWS only with the med-risk alarms, what is the annual profit we will generate vs the current system as a % of the total profit?**

**ii. And if we set the EWS only with the high-risk alarms?**

[Hide](#)

```
# Assumptions:  
margin_per_item = 1.7  
cost_of_checking = 10  
time_to_fix = 1.5  
current_false_alarms = 2.2
```

Function for calculating increase in profit based on risk threshold:

[Hide](#)

```

profit_increase = function(threshold) {
  if (threshold == threshold_med_risk){
    dt_alarms = transactional[transactional$alarm_med_risk == 1,]
  }
  else {
    dt_alarms = transactional[transactional$alarm_high_risk == 1,]
  }

  ## Additional costs from false alarms
  # New costs
  new_costs = nrow(dt_alarms[dt_alarms$failure == 0]) * cost_of_checking

  # Current costs
  current_costs = transactional[, .(max = max(timestamp), min = min(timestamp)), by =
machine]
  current_costs$cost = as.numeric(difftime(current_costs$max,
                                           current_costs$min,
                                           units = "days"))/ 365 * current_false_alar
ms * cost_of_checking
  # Total current costs
  current_costs = sum(current_costs$cost)

  # Net costs
  net_costs = new_costs - current_costs
#-----

# Additional revenue
dt_alarms[, threshold_hours := threshold * (24 / daily_sales_machine)]
dt_alarms[, threshold_hours_fixed := threshold_hours + time_to_fix]
dt_alarms[, delta_fixed := threshold_hours_fixed * daily_sales_machine / 24 ]

# Additional sales:
dt_alarms[, won_sales := failure * (delta - delta_fixed)]
add_margin = sum(dt_alarms$won_sales) * margin_per_item

# Profit increase
delta_profit = add_margin - net_costs
current_profit = nrow(transactional) * margin_per_item
profit_increase = delta_profit/current_profit

  return(profit_increase)
}

```

[Hide](#)

```

# i) increase Profit of medium risk:
profit_increase_medium_risk = profit_increase(threshold_med_risk)
sprintf("%1.2f%%", 100*profit_increase_medium_risk)

# ii) Increase Profit of high risk
profit_increase_high_risk = profit_increase(threshold_high_risk)
sprintf("%1.2f%%", 100*profit_increase_high_risk)

```