

# 08 Feature Engineering

Fabian Blasch

06/14/2022

## 1 Load Data

```
# source AUX
source("../Misc/Auxilliary.R")

# packages
get.packages(c("ggplot2", "patchwork", "modeldata", "recipes", "stringr"))

# load data
dat_bids <- readRDS("../Data/Bid Tab RDS/Bids_df.RDS")
dat_aucs <- readRDS("../Data/Bid Tab RDS/Aucs_df.RDS")
```

## 2 Bidder Interactions

Given that some bidders may interact when finding prices it makes sense to create a dummy that represents bidder groups. Accordingly, we may first identify the column via a prefix *Vend\_* to then find all the pairwise interactions utilizing the package *recipes*.

```
# change vendor ID names to subset easier with starts_with
indlv <- which(names(dat_aucs) == "969A")
names(dat_aucs)[11:indlv] <- paste0("Vend_", names(dat_aucs)[11:indlv])

# now use recipes to create vendor interactions
recip <- recipe(Winning_Bid ~., data = dat_aucs)

# interactions
recip |> step_interact(terms = ~ starts_with("Vend"):starts_with("Vend")) |>
  prep(dat_aucs) |>
  bake(dat_aucs) -> dat_aucs_vend_int

# apply over cols remove all columns that contain only 0, i.e., all interactions
# of firms that never bid in the same auction
sapply(dat_aucs_vend_int, \(x){

  # check if amount of occurrences is larger than
  (sum(abs(x)) > 0) |> tryCatch(error = \(e) TRUE) # catch all cases where summing
```

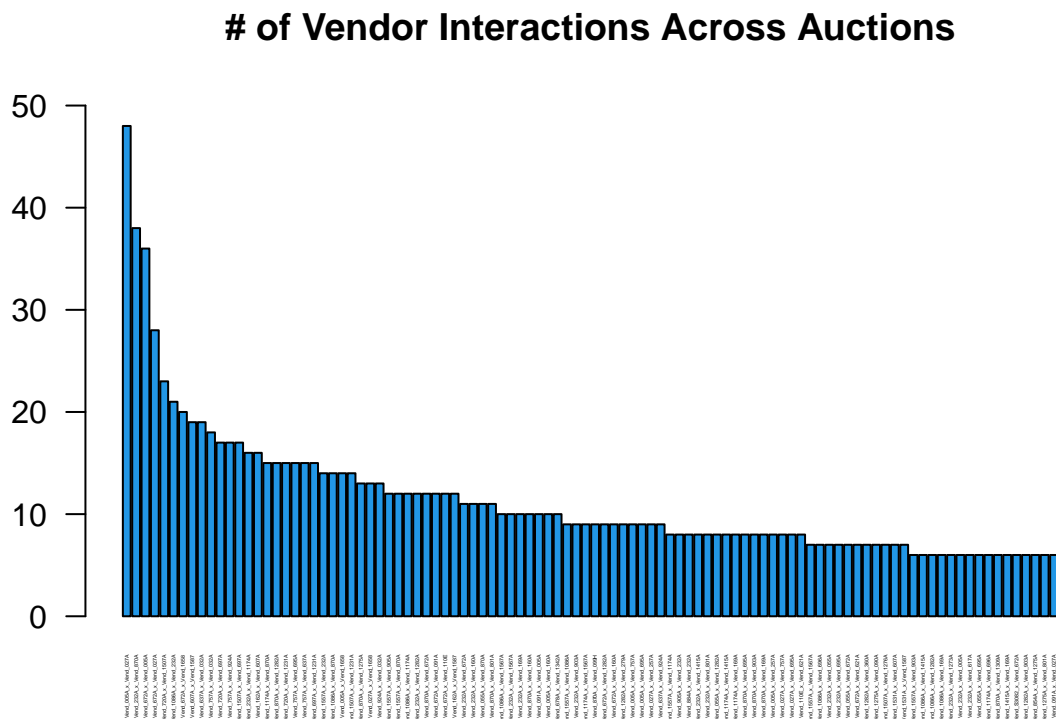
```

# makes no sense and return TRUE
}) -> tmp

# use to subset our data set
dat_aucsvend_int <- dat_aucsvend_int[, tmp]

# display which bidders interact most
ind <- str_detect(names(dat_aucsvend_int), "_x_")
sm <- sapply(dat_aucsvend_int[, ind], sum) |> sort(decreasing = TRUE)
barplot(sm[1:100], ylim = c(0, 50), las = 2, cex.names = 0.2, col = 4,
        main = "# of Vendor Interactions Across Auctions")

```



### 3 Sum of Auctions Won

Since the Data does not contain any firm specific information and estimating cost functions is not feasible on firm level with the small data set at hand, we may sum up the cumulative volume of contracts won by all firms that participate in an auction. Given that firms have no other advantage than their cost structure this should account for potential difference in cost structure among participating companies.

```
# find sum of auctions won grouped by year and Vendor ID
```

## 4 Train and Test

```
# splits
{set.seed(33)
ind <- sample(1:nrow(dat_aucsvend_int), replace = FALSE,
              size = floor(nrow(dat_aucsvend_int) * 0.2))
}

# train and test
dat_aucsvend_split <- list("Train" = dat_aucsvend_int[!(1:nrow(dat_aucsvend_int) %in% ind), ],
                           "Test" = dat_aucsvend_int[ind, ])

# write split
# saveRDS(dat_aucsvend_split, "../Data/Bid Tab RDS/Aucsvend_split.RDS")
# saveRDS(dat_aucsvend_int, "../Data/Bid Tab RDS/Aucsvend.RDS")
```