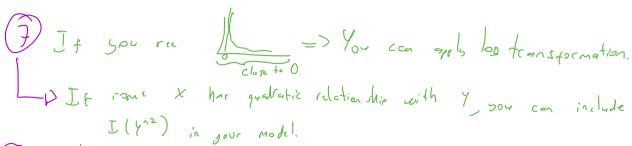
# BARML - Analytica Cup Plan

1) Check R2, it should be high Particularly Multiple R-Squarel. 2) Chech VIF (Variance Inflation Factor) In If VIFT are IT, which correlation coefficients LD If some columns are corelated drap one by one check VIF to If 2 variables are correlated at a rate greater than .6.7 to droppins the least theoretically important of the two.

3 Chede p-values [Po(>|t|)] of columns, they should be significant (\*\*\*) Apply Glescher Test to see if there is heteroscedasticity in a particular column lation.

Needo normally distributed residuals.

- 5) Apply White Test to see if there is heteroscedasticity
- (6) Use Durbin-Wation (DW) Method to see if there is antocorrelation in the column.



- 8) If data is "panel data", consider tenting enlageneits and use either Random Effects Model or Fixed Effects Model.
- 9) You can use Wald Statistics and remove irrelevant variables.
- (10) It son use Logistic Regression vice Wald-test for significance (similar to t-test for Linear Regression)
- 11) If there is a "rewre" ordinal column, the "rank" you can use as, factor (rank) in your model.
- (12) If you we Logiotic Regression, use McFad for understanding explanators power of covariates (similar to) R2)

  13) Pay y Hentian to data types! (Nominal, odda),

  Interval, Ratio). You can sometimes use ordinal as numeric (similar to ratio).
- (14) I of there are many attributes and attributes seen to be independent and equally important, you can use Naive Bayes Classifier.

- 15) If you were Decision Trees, first only consider aftributes with greater than information gain, then compare them on gain ratio.
  - 16 It son use a tree-based method, consider binning numeric affeibates so that they can be treated as nominal affributes.
  - 17 Especially in tree-barred methods, you can add an additional attribute as "attribute a > attribute b.

### **Data Preprocessing**

#### **COLUMNS REQUIRING "ONE HOT ENCODING" ONLY:**

- PRICE LIST -> ONE HOT ENCODING #DONE
- TECH -> ONE HOT ENCODING #DONE
- OFFER TYPE -> ONE HOT ENCODING #DONE
- BUSINESS\_TYPE -> ONE HOT ENCODING #DONE
- SALES\_OFFICE: ONE HOT ENCODING #DONE
- SALES\_BRANCH: ONE HOT ENCODING #DONE
- OWNERSHIP: ONE HOT ENCODING
- SALES\_LOCATION

#### **OTHER COLUMNS:**

- END\_CUSTOMER -> HAS END CUSTOMER
- ISIC -> HAS ISIC
- COSTS\_PRODUCT\_A to COSTS\_PRODUCT\_E:
  - COSTS\_PRODUCT\_\* -> INCLUDED\_COSTS\_PRODUCT\_\*
  - TOTAL\_COSTS\_PRODUCT: SUM OF ALL COSTS\_PRODUCT\_\*
- COUNTRY\_CODE: Convert to BINARY
- REV\_CURRENT\_YEAR.1 and REV\_CURRENT\_YEAR.2:
  - REV CURRENT YEAR.1 > EURO (with a fixed exchange rate)
  - REV\_CURRENT\_YEAR.2 -> EURO (with a fixed exchange rate)
  - DROP ONE OF THEM: REV CURRENT YEAR.1 or REV CURRENT YEAR

  - New column PREV\_YEAR\_PERCENTAGE\_INCREASE:
     ((REV\_CURRENT\_YEAR \*)
     REV\_CURRENT\_YEAR.2)/REV\_CURRENT\_YEAR)\*100)
- CREATION\_YEAR: Without null values, look at correlation with the target. If important, then think about filling missing values.
  - Extract only the year.
  - Calculate how long since CREATION\_YEAR
- OWNERSHIP:
  - NA <- No information:
    </p>
  - Without null values, look at correlation with the target. If important, then think about filling missing values or no information.
  - Without "no information" values, look at correlation with the target. If important, then think about filling missing values or no information.
  - Change it to one-hot encoding.

### Additional Advice from Stefan Heidekrüger

- -> Log transformations: log+1
- -> Create a BAC pipeline

- -> do not make data leakage
- -> different customer ids in train and validation set
- -> If Nominal columns are majority: tree-based solutions are better
- -> If numerical columns are majority: logistic solutions are better
- -> In date columns: you can extract Similar to unix timestamp, you can convert dates to
- -> LogReg,Random Forest: For most cases one-hot encoding is dealt implicitly
- -> If necessary: After one-hot encoding. Last 5% or last 10% to some "others"
- -> not more than 150 columns
- -> log transformations are kinda outlier prevention

### Filling Missing Values

- CREATION\_YEAR
- REV\_CURRENT\_YEAR.1
- REV\_CURRENT\_YEAR.2
- REV\_PERCENTAGE\_INCREASE
- REV\_CURRENT\_YEAR.1 and REV\_CURRENT\_YEAR.2: If both zero or NA, before replacing, be sure that it is converted to euro, then replace with mean
- OWNERSHIP
- OWNERSHIP\_NO\_INFO\_AS\_NA
- SALES\_LOCATION -> ONE HOT ENCODING
  - In Test set and has no SALES LOCATION
    - SALES\_LOCATION: NA TEST\_SET\_ID: 12359, Random Sales location in CH (?)
    - SALES\_LOCATION: NA TEST\_SET\_ID: 16396, Random Sales location in CH (?)

# Feature Importance/Check

- Correlation matrix between all columns
- R library for ranking which features are statistically significant
  - o library(Boruta) https://www.machinelearningplus.com/machine-learning/feature-selection/
  - VIP model about feature importanceShap values expensive to compute
- For numerical values, plot independent values with respect to the target, then consider applying log transformation to the column with many 0s.
- Isolate numerical values and try PCA
- Apply Logistic Regression
  - Check Variance Inflation Score
  - Apply McFad to mimic R^2.

- o For significance, use wald test
- You can try Random Forest, Gradient boosted trees etc.

\_\_\_\_\_

! Remember to set the seed before doing anything!:)

- + Joining, pre-processing
- + Extreme Values, Missing values
- + Data Augmentation: Feature extraction, combine features, separate features
- + Feature Importance, Explanatory

#### + Models:

- Ensemble: Majority Voting, bagging, boosting
- Start with: simple logistic regression
- Decision Trees
- Naive Bayes
- Random Forest
- Neural Networks
- AutoML: to narrow down optionsSplit validation: split by customer

### Task distribution:

# **Weekly Plan**

Week1: Data Preprocessing
Week2:
Week3: