



A model to Determine Business Opportunities

Ricardo Moncayo

B.Sc. Electronic Engineer , M.Sc. Biomedical Engineering

June 3, 2022

THE PROBLEM:

The challenge: To build a predictive model to determine which business opportunities are more likely to win, using data integration from different sources

PROPOSAL

Three Approaches were tested:

- ▶ A model trained from the raw data (Raw model)
- ▶ A model enhanced by selecting relevant features
- ▶ A model with relevant features without the country index (naive decision)
- ▶ A model-based on a supervised clustering of the country data relationships

PROPOSAL

Three Approaches were tested:

- ▶ A model trained from the raw data (Raw model)
- ▶ A model enhanced by selecting relevant features
- ▶ A model with relevant features without the country index (naive decision)
- ▶ A model-based on a supervised clustering of the country data relationships

An AdaBoost classifier is trained using all data without the opportunity number

PROPOSAL

Three Approaches were tested:

- ▶ A model trained from the raw data (Raw model)
- ▶ **A model enhanced by selecting relevant features**
- ▶ A model with relevant features without the country index (naive decision)
- ▶ A model-based on a supervised clustering of the country data relationships

PROPOSAL

Three Approaches were tested:

- ▶ A model trained from the raw data (Raw model)
- ▶ A model enhanced by selecting relevant features
- ▶ A model with relevant features without the country index (naive decision)
- ▶ A model-based on a supervised clustering of the country data relationships

Relevant features are selected using mRMR (minimum Redundancy Maximum Relevance), and an AdaBoost validated in the complete dataset

PROPOSAL

Three Approaches were tested:

- ▶ A model trained from the raw data (Raw model)
- ▶ A model enhanced by selecting relevant features
- ▶ A model with relevant features without the country index (naive decision)
- ▶ A model-based on a supervised clustering of the country data relationships

PROPOSAL

Three Approaches were tested:

- ▶ A model trained from the raw data (Raw model)
- ▶ A model enhanced by selecting relevant features
- ▶ A model with relevant features without the country index (naive decision)
- ▶ A model-based on a supervised clustering of the country data relationships

Without country index data, relevant features are selected, and an AdaBoost validated

PROPOSAL

Three Approaches were tested:

- ▶ A model trained from the raw data (Raw model)
- ▶ A model enhanced by selecting relevant features
- ▶ A model with relevant features without the country index (naive decision)
- ▶ A model-based on a supervised clustering of the country data relationships

PROPOSAL

Three Approaches were tested:

- ▶ A model trained from the raw data (Raw model)
- ▶ A model enhanced by selecting relevant features
- ▶ A model with relevant features without the country index (naive decision)
- ▶ A model-based on a supervised clustering of the country data relationships

A model is trained (with the most relevant features) using the data from a country and evaluating in the remaining countries independently. The higher results correspond to countries with similar decision processes

A MODEL TRAINED FROM THE RAW DATA (RAW MODEL)

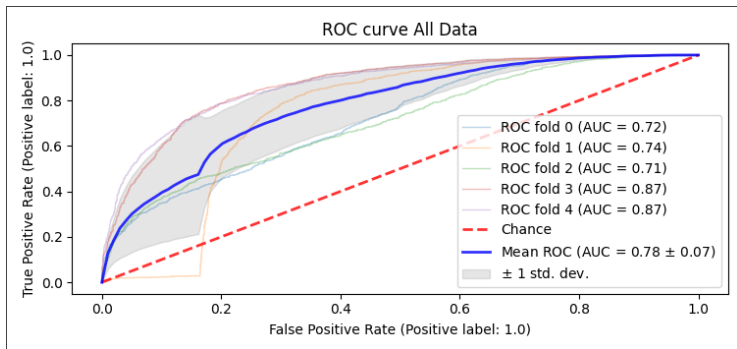
- ▶ The opportunity number is removed from the data
- ▶ An AdaBoost classifier is trained and testing using a 5-fold cross-validation

A MODEL TRAINED FROM THE RAW DATA (RAW MODEL)

- ▶ The opportunity number is removed from the data
- ▶ An AdaBoost classifier is trained and testing using a 5-fold cross-validation

A MODEL TRAINED FROM THE RAW DATA (RAW MODEL)

- ▶ The opportunity number is removed from the data
- ▶ An AdaBoost classifier is trained and testing using a 5-fold cross-validation



A MODEL ENHANCED BY SELECTING RELEVANT FEATURES

- ▶ mRmR is used to select 10 most relevant features
- ▶ Stable relevant features are finally selected
- ▶ An AdaBoost classifier is trained and testing using a 5-fold cross-validation

$$score_i(f) = \frac{relevance(f|target)}{redundancy(f|featuresselecteduntil i-1)}$$

A MODEL ENHANCED BY SELECTING RELEVANT FEATURES

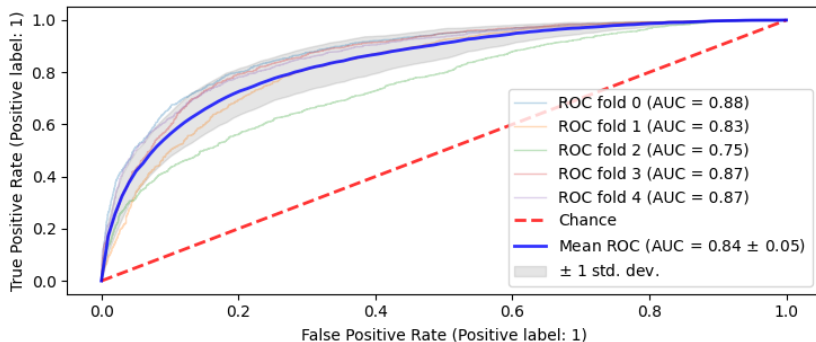
- ▶ mRmR is used to select 10 most relevant features
- ▶ Stable relevant features are finally selected
- ▶ An AdaBoost classifier is trained and testing using a 5-fold cross-validation

A MODEL ENHANCED BY SELECTING RELEVANT FEATURES

- ▶ mRmR is used to select 10 most relevant features
- ▶ Stable relevant features are finally selected
- ▶ An AdaBoost classifier is trained and testing using a 5-fold cross-validation

Best Features: CountryCode,Deal Size Category,Opportunity Amount USD,Revenue From Client Past Two Years,Sales Stage Change Count,Supplies Group,Supplies Subgroup,Total Days Identified Through Closing,Total Days Identified Through Qualified

RESULTS USING RELEVANT FEATURES



A MODEL WITHOUT COUNTRY INDEX AND THE MOST RELEVANT FEATURES

- ▶ Country indexes are removed from the data (mRmR)
- ▶ mRmR is used to select 10 most relevant features and stable features are finally selected
- ▶ An AdaBoost classifier is trained and testing using a 5-fold cross-validation in all data and in each country independently

A MODEL WITHOUT COUNTRY INDEX AND THE MOST RELEVANT FEATURES

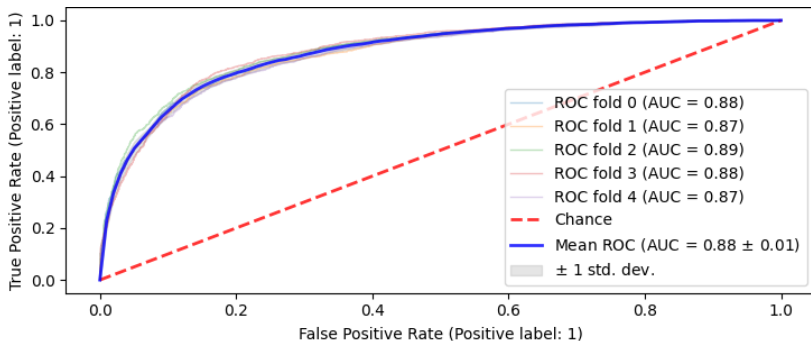
- ▶ Country indexes are removed from the data (mRmR)
- ▶ mRmR is used to select 10 most relevant features and stable features are finally selected
- ▶ An AdaBoost classifier is trained and testing using a 5-fold cross-validation in all data and in each country independently

Best Features: Client Size By Employee Count, Deal Size Category, Opportunity Amount USD, Revenue From Client Past Two Years, Sales Stage Change Count, Supplies Group, Supplies Subgroup, Total Days Identified Through Closing, Total Days Identified Through Qualified

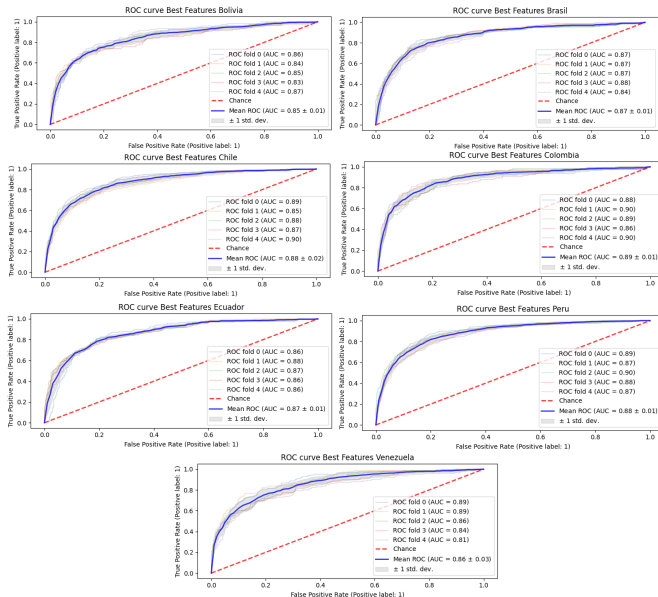
A MODEL WITHOUT COUNTRY INDEX AND THE MOST RELEVANT FEATURES

- ▶ Country indexes are removed from the data (mRmR)
- ▶ mRmR is used to select 10 most relevant features and stable features are finally selected
- ▶ An AdaBoost classifier is trained and testing using a 5-fold cross-validation in all data and in each country independently

RESULTS-RELEVANT FEATURES WITHOUT COUNTRY INDEX - ALL DATA



RESULTS RELEVANT FEATURES MODEL PER COUNTRY



A MODEL-BASED IN A SUPERVISED CLUSTERING OF THE COUNTRY DATA RELATIONSHIPS

- ▶ A model is trained using the data from a country and evaluating in the remaining countries (i.e. Training from Brasil evaluating in Colombia, Training from Brasil evaluating in Peru, Training from Colombia evaluating in Brasil...)
- ▶ The AUC is computed in each evaluation
- ▶ A matrix of all AUC is obtained, and a threshold is applied using as reference the AUC values obtained previously in the feature selection process

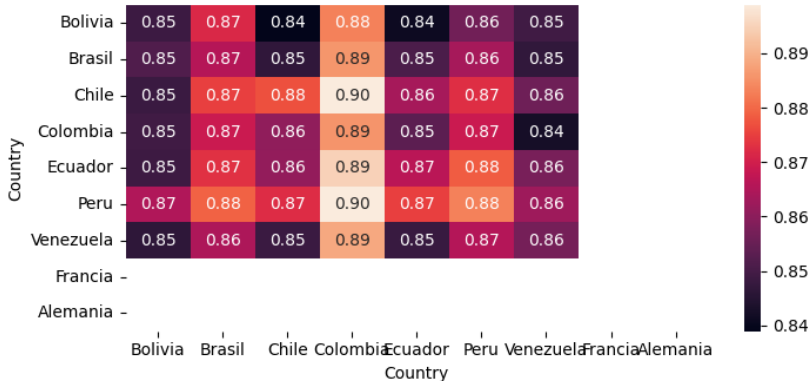
A MODEL-BASED IN A SUPERVISED CLUSTERING OF THE COUNTRY DATA RELATIONSHIPS

- ▶ A model is trained using the data from a country and evaluating in the remaining countries (i.e. Training from Brasil evaluating in Colombia, Training from Brasil evaluating in Peru, Training from Colombia evaluating in Brasil...)
- ▶ The AUC is computed in each evaluation
- ▶ A matrix of all AUC is obtained, and a threshold is applied using as reference the AUC values obtained previously in the feature selection process

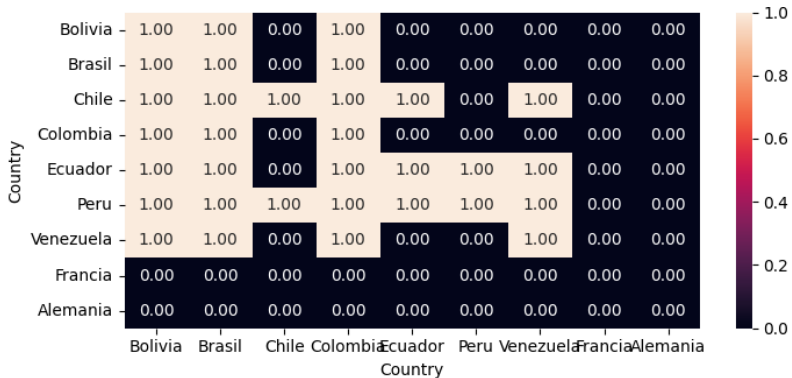
A MODEL-BASED IN A SUPERVISED CLUSTERING OF THE COUNTRY DATA RELATIONSHIPS

- ▶ A model is trained using the data from a country and evaluating in the remaining countries (i.e. Training from Brasil evaluating in Colombia, Training from Brasil evaluating in Peru, Training from Colombia evaluating in Brasil...)
- ▶ The AUC is computed in each evaluation
- ▶ A matrix of all AUC is obtained, and a threshold is applied using as reference the AUC values obtained previously in the feature selection process

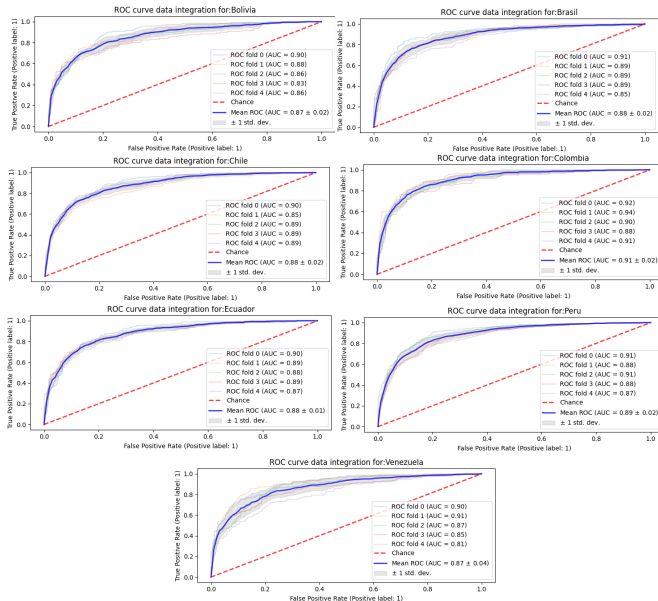
RESULTS OF THE AUC MATRIX- DATA RELATIONS



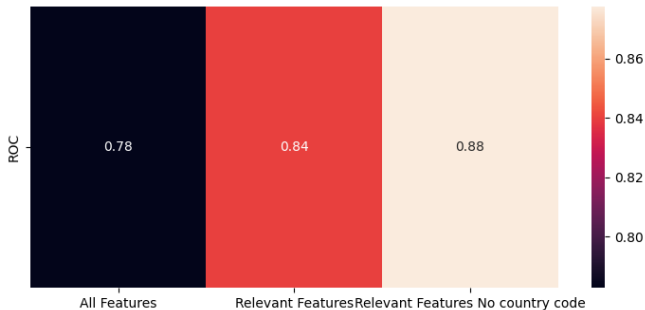
RELATIONSHIPS ARE REVEALED AFTER A THRESHOLD OPERATION



RESULTS PER COUNTRY



SUMMARY - RESULTS RELEVANT FEATURES SELECTION



SUMMARY - RESULTS DATA INTEGRATION PER COUNTRY



RESULTS OF THE BLIND DATA

	Relevant Features	Data Integration
Win	5247	5390
Lose	26576	26433

CONCLUSIONS

- ▶ Relevant Features enhance the data representation
- ▶ Data integration show an improvement in the representation
- ▶ Setting a manual threshold can be improved