

# Data Driven Contact Strategy

*Machine Learning applied to Marketing*

 LUISS x **Deloitte.**

Fabiana Caccavale  
Marco Amadori  
Lorenzo Meloncelli





# INDEX



1. Goals
2. Project overview
3. Data comprehension
4. Data preparation
5. Propensity model
  - Features selection
  - Model training/testing
6. Eligibility model
7. Contact strategy
8. Conclusions and future developments

# GOALS



To estimate the likelihood of positive responses to a specific campaign (cross-selling/solution).

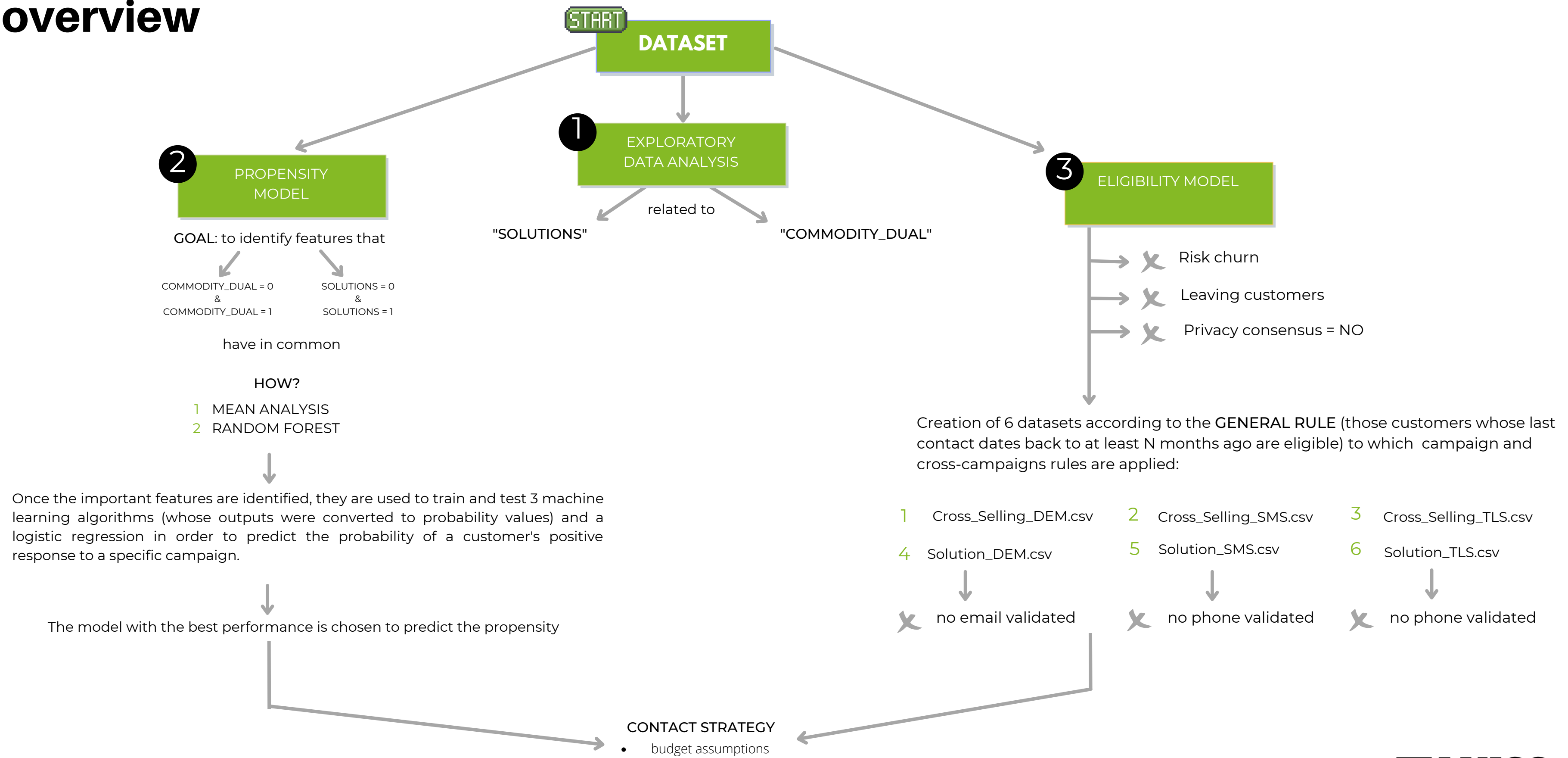


To analyze customers' contactability.

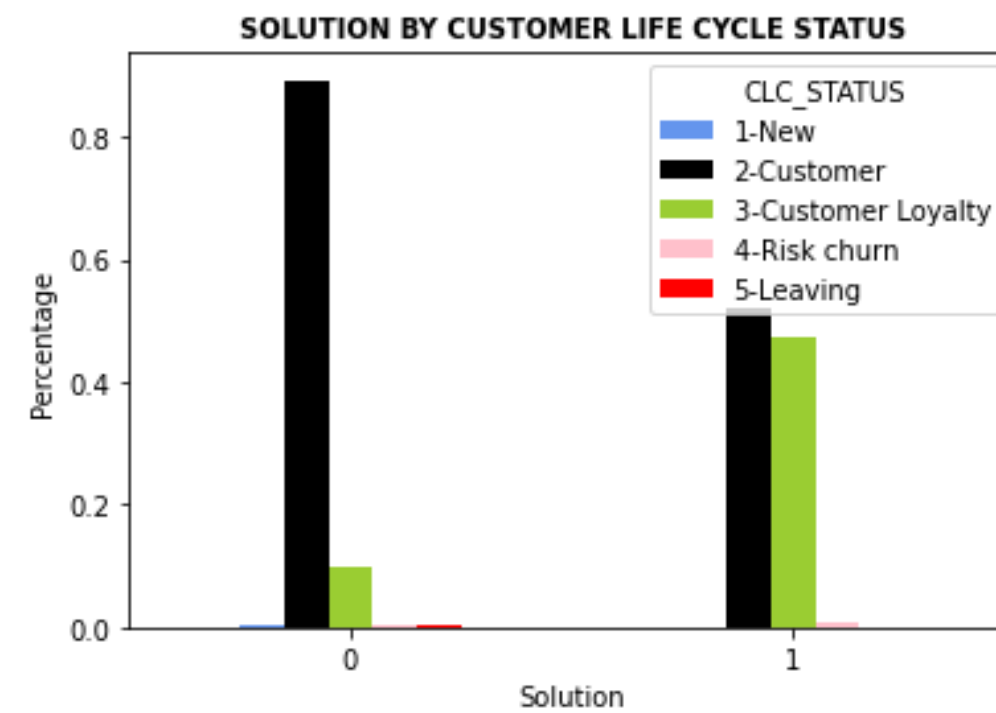
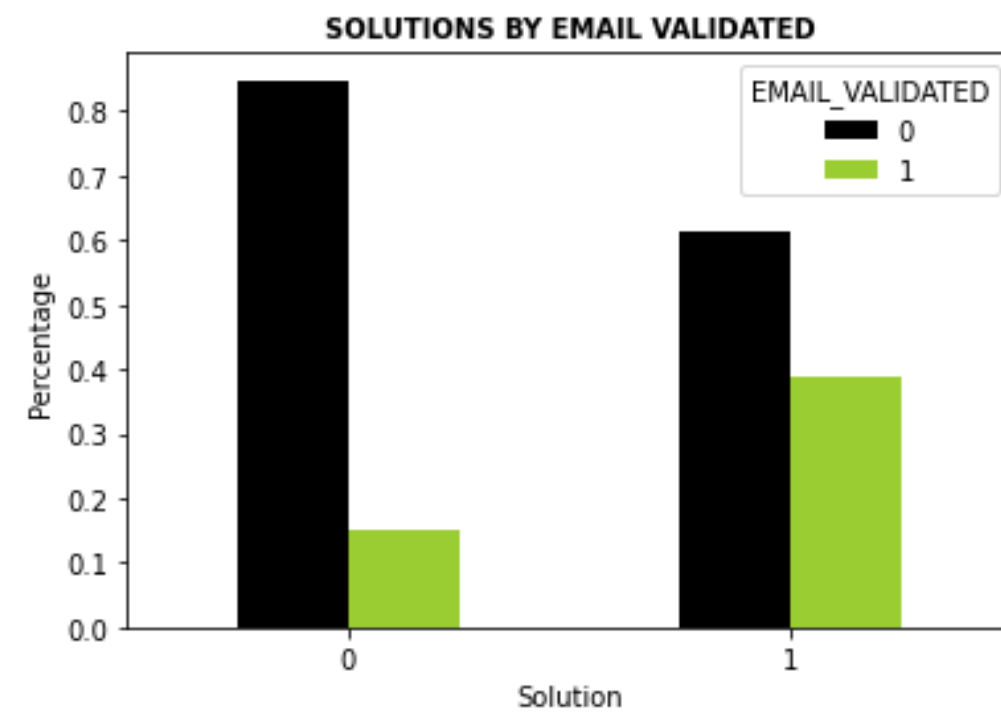
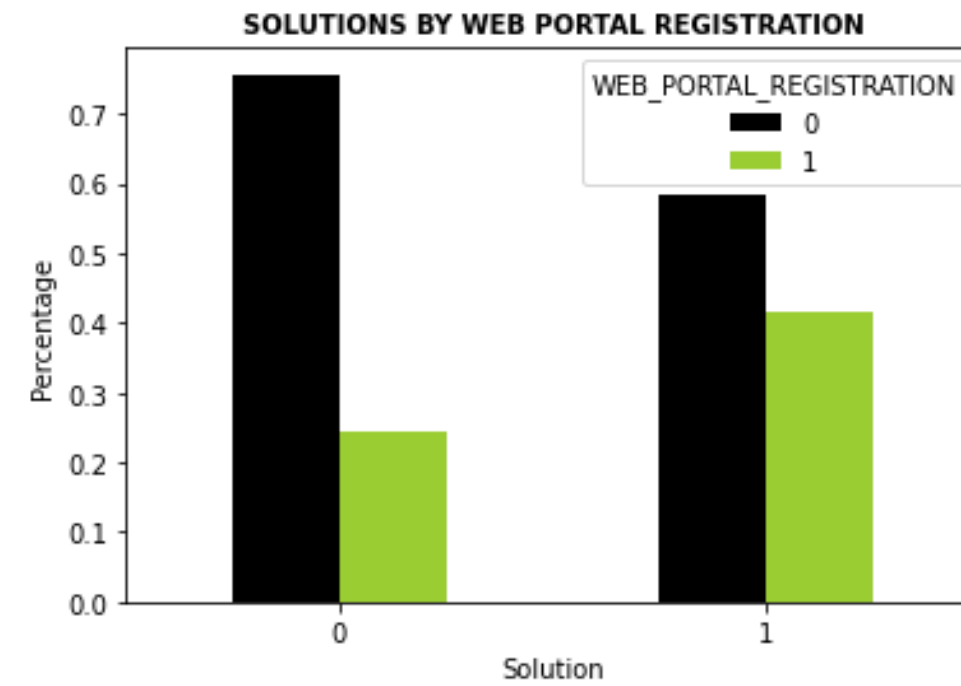
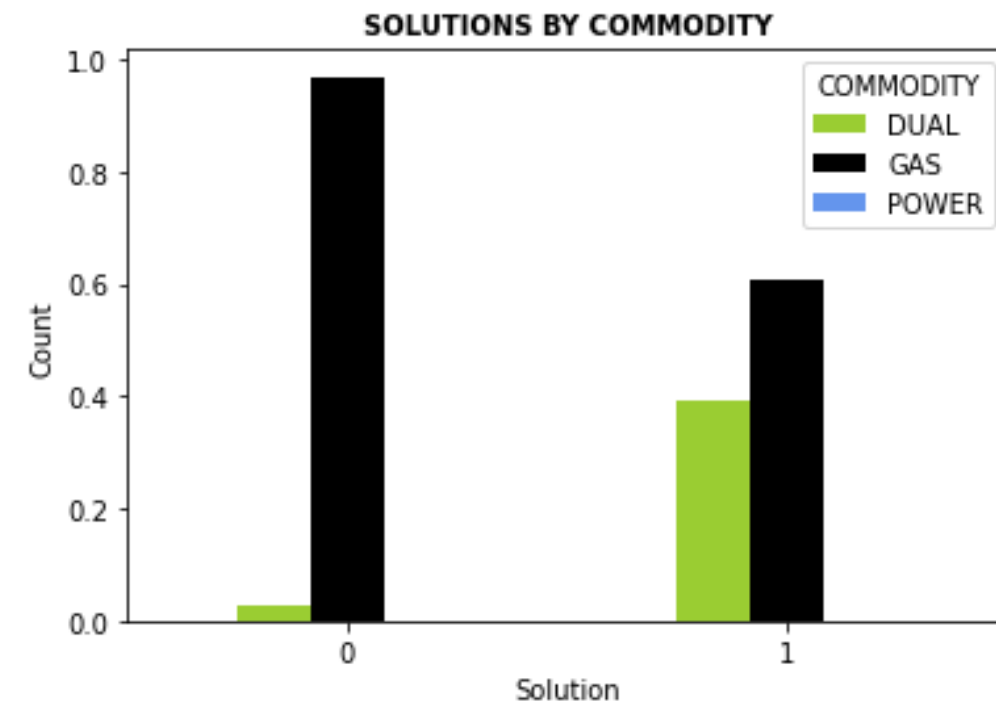
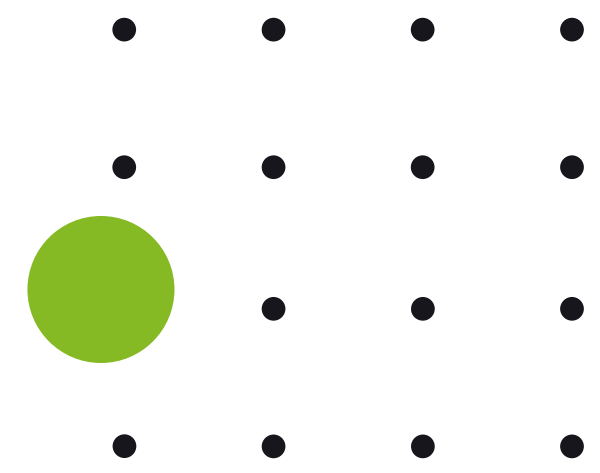


To monthly assign each customer to a campaign and a contact channel.

# Project overview

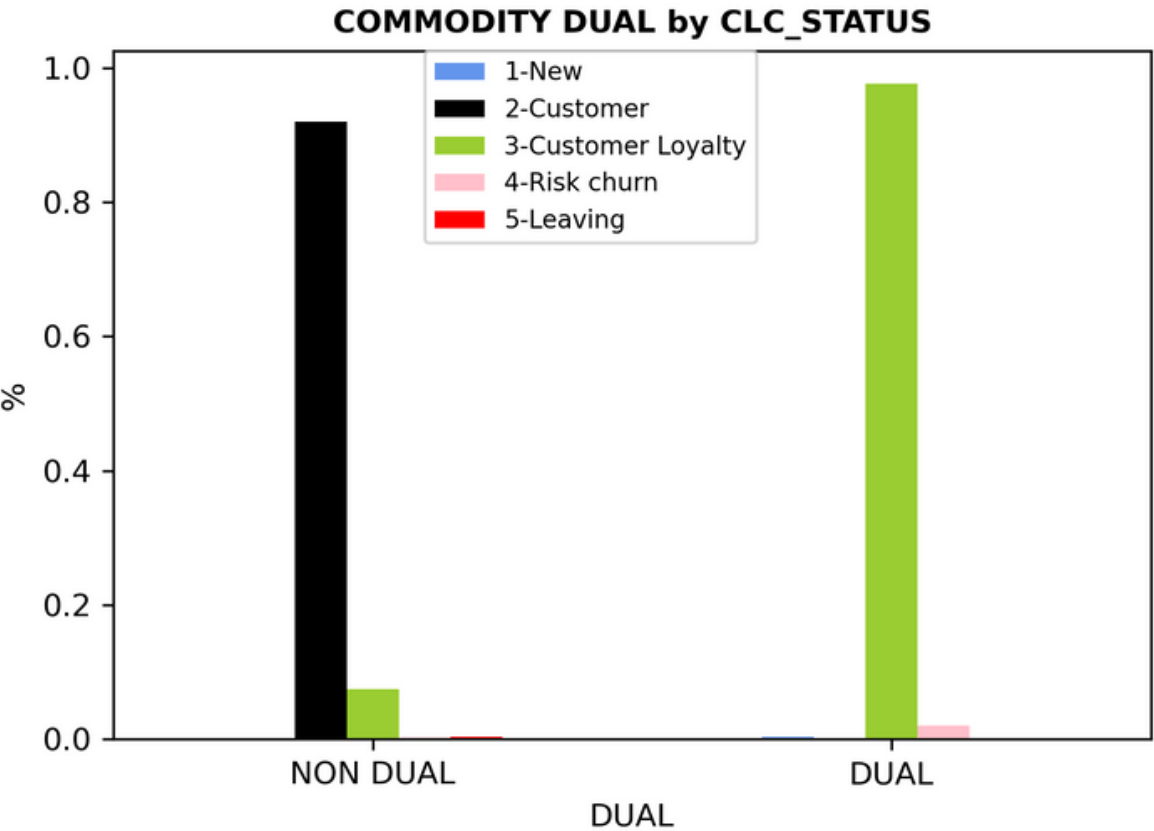
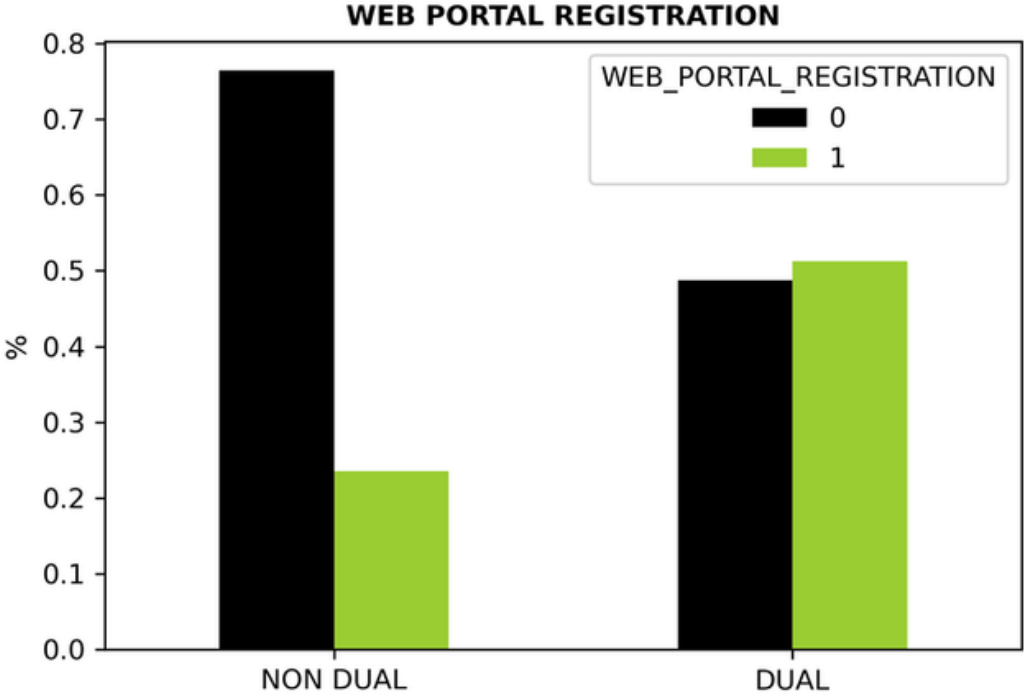
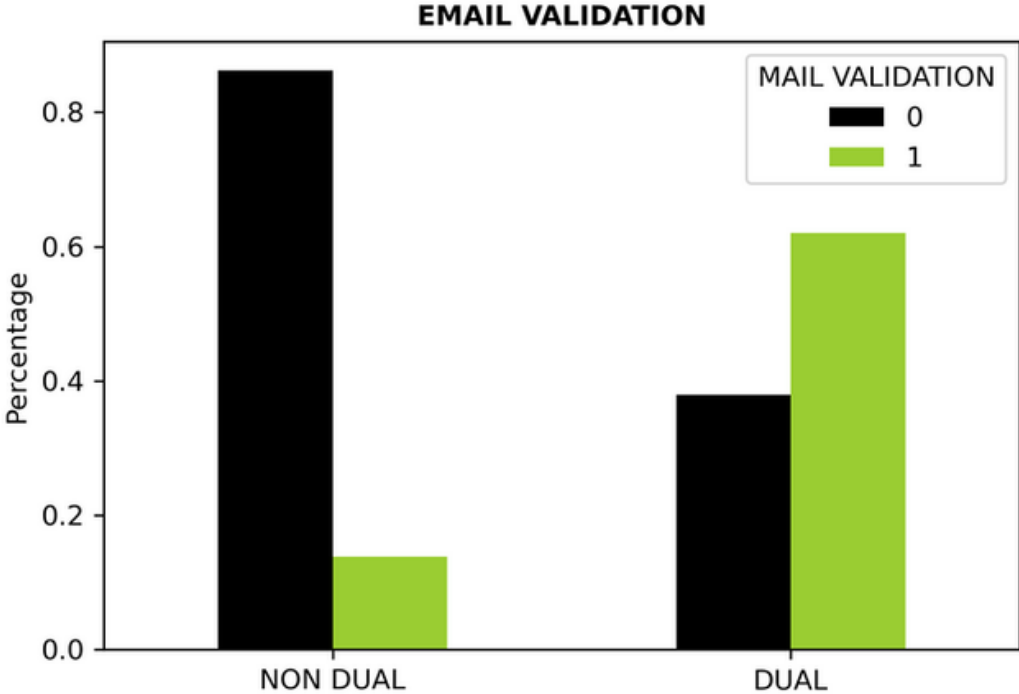
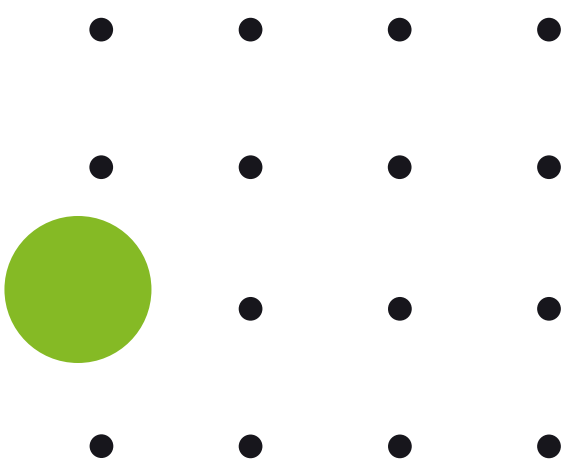


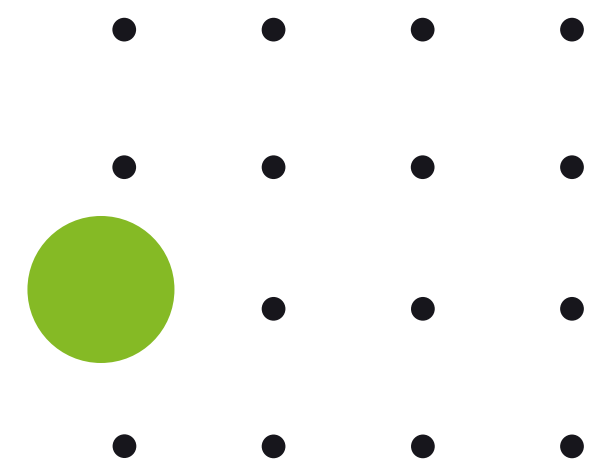
# Data comprehension (SOLUTIONS)



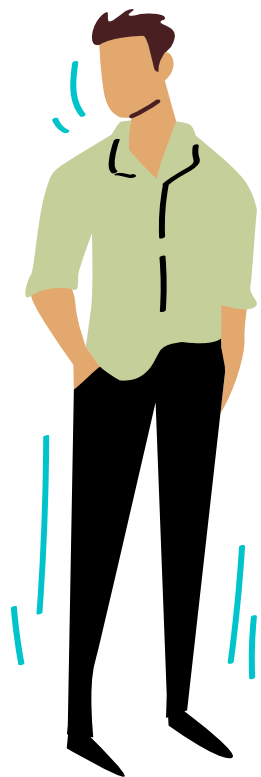
# Data comprehension

(COMMODITY\_DUAL)





# The typical customer:



## "COMMODITY\_DUAL" = 1

- 60% are registered to the web portal
- Loyal customers
- 60% validate their mail
- Area: NORTH

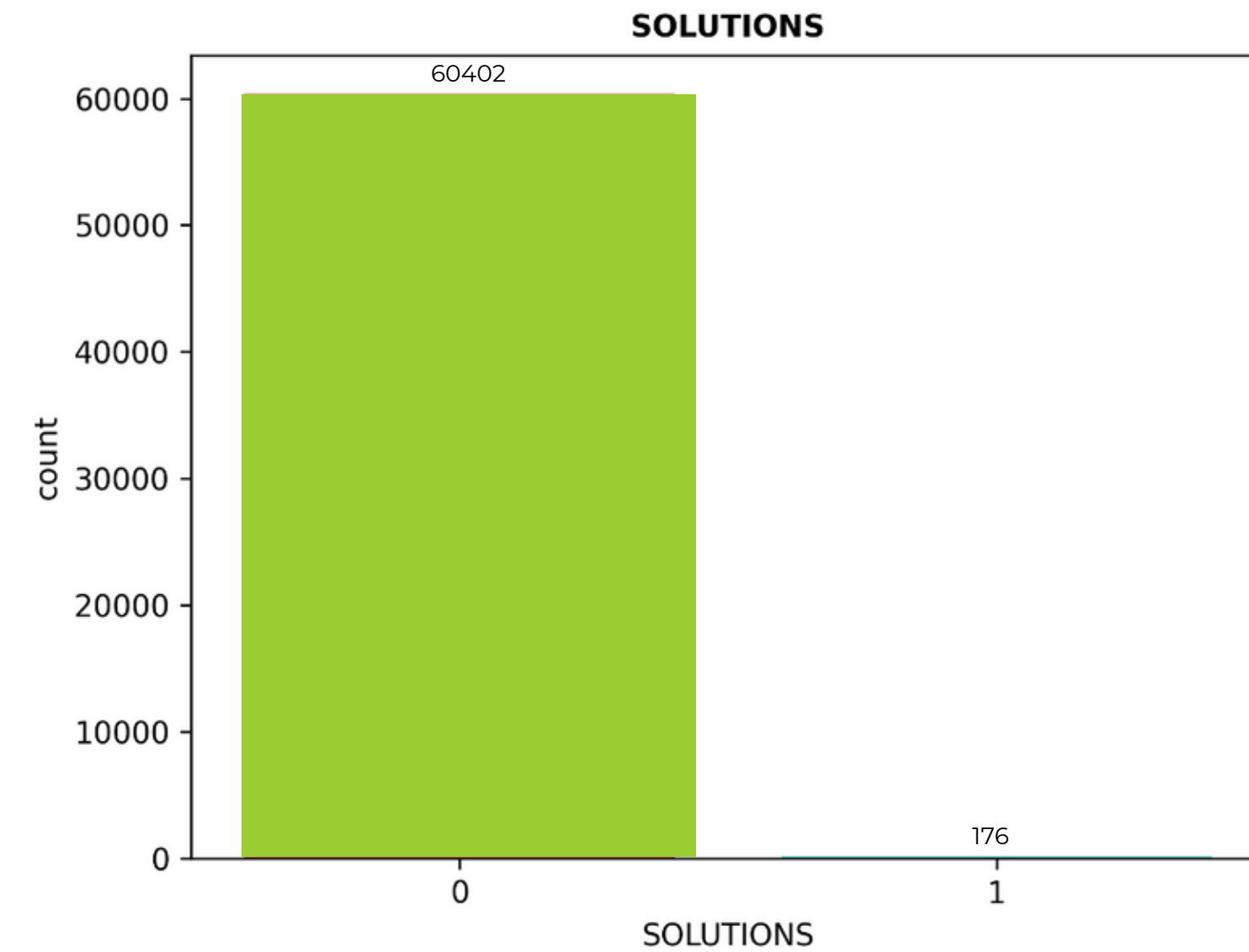
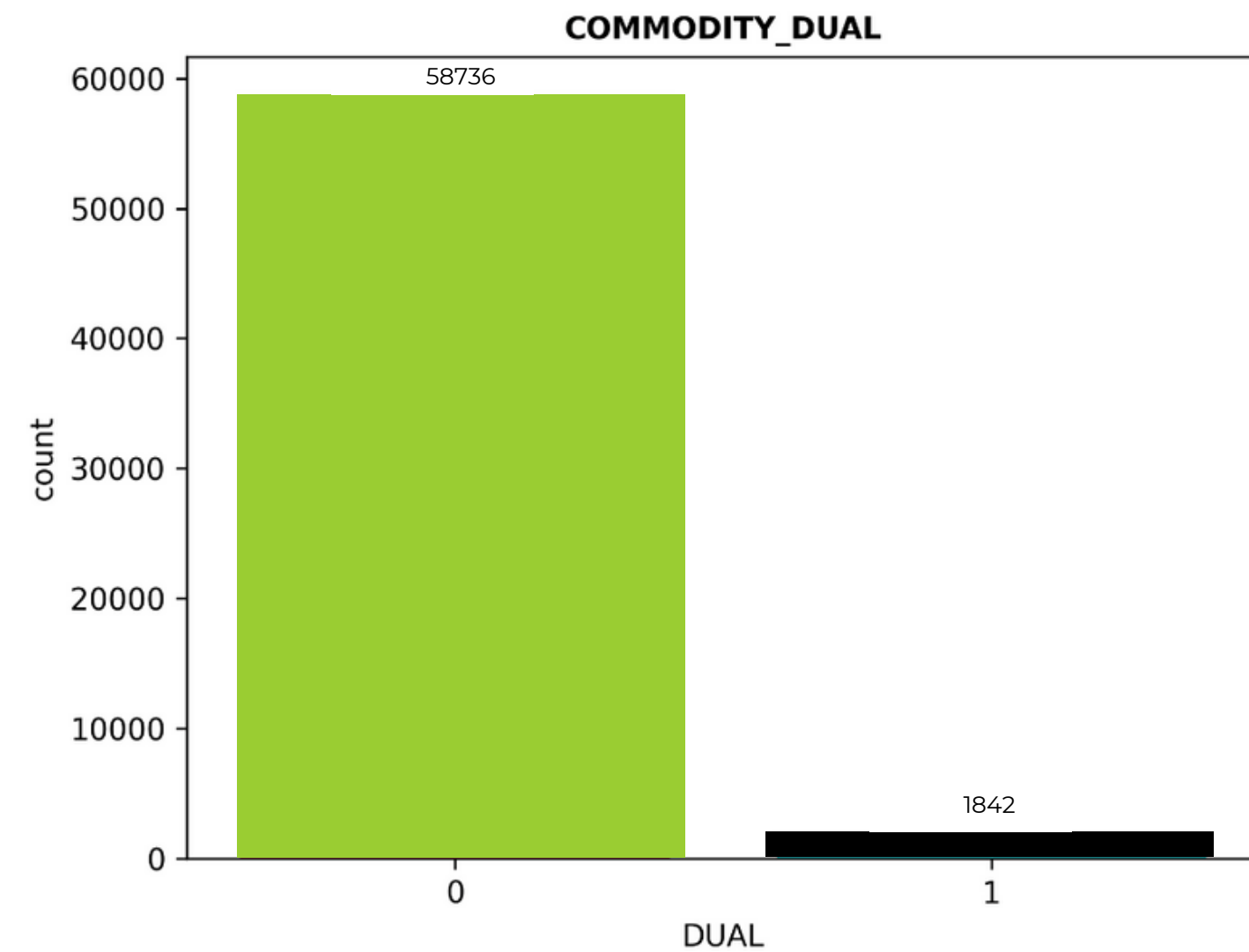
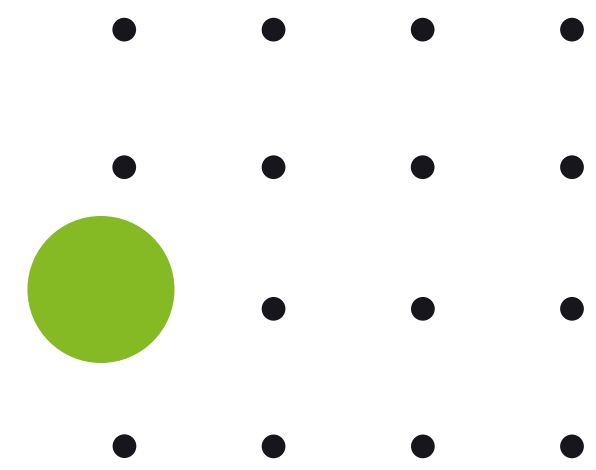


## "SOLUTIONS" = 1

- 40% are registered to the web portal
- Loyal customers
- 40% validate their mail
- Dual or gas contract
- More frequent inbound contacts
- Registered to the web portal

# Data comprehension:

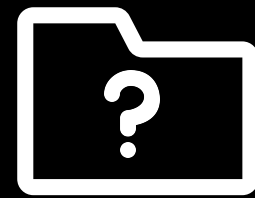
## Dataset imbalance



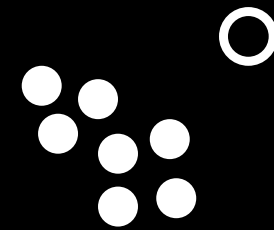


# Data preparation

Main issues:



Missing  
values



Outliers



Dataset  
imbalance

- Missing values were removed.
- Outliers were replaced with the appropriate median.
- To overcome the problem of dataset imbalance, we randomly selected observations from the majority class and deleted them from the training dataset.

# Propensity model

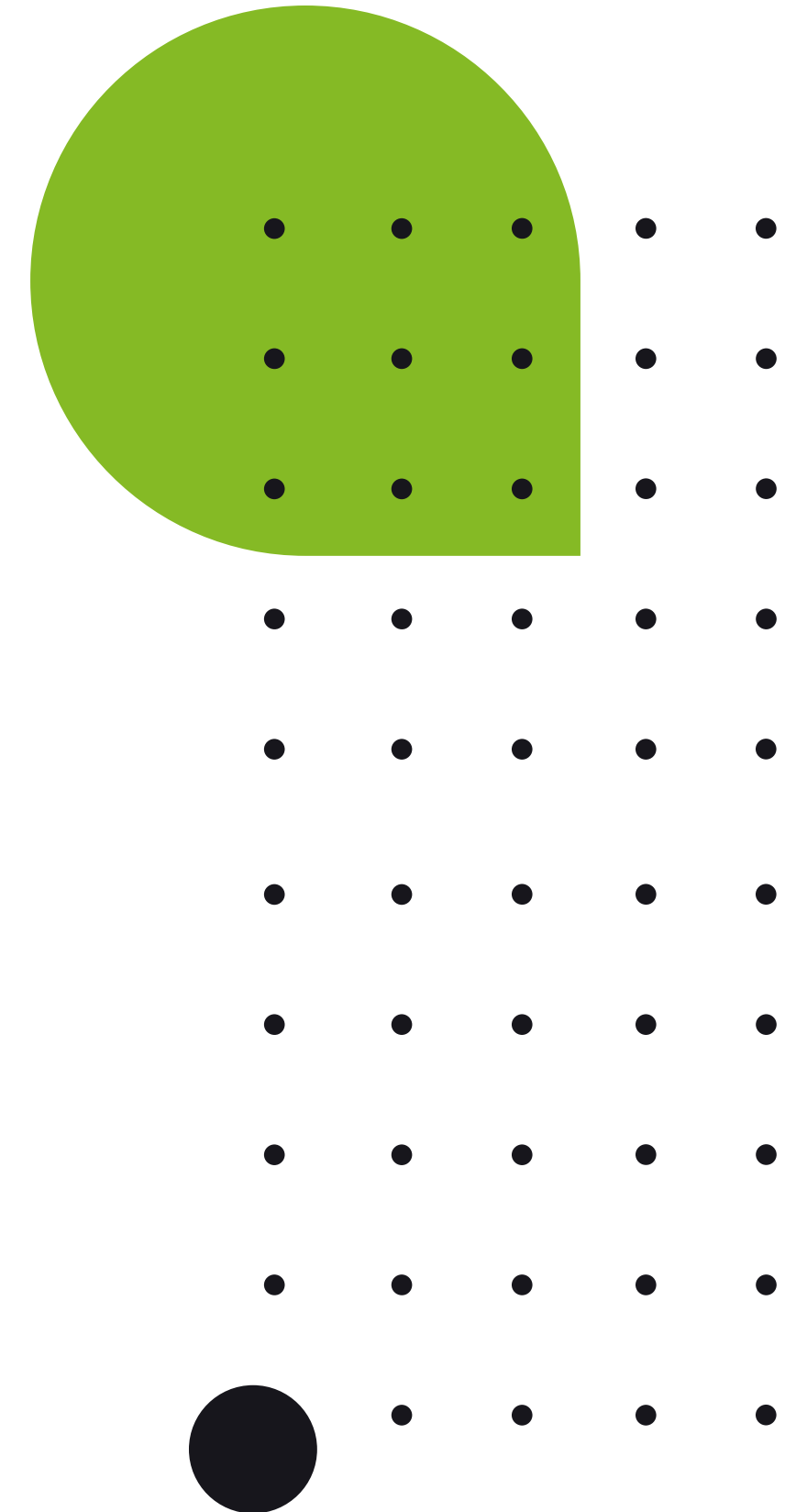
GOAL: to develop a model that would have helped to select which customers are most likely to purchase a particular contract.



To do that, features that could have most influenced consumers' choice were identified by comparing customers who had already purchased the contract to those who had not.

This was done through:

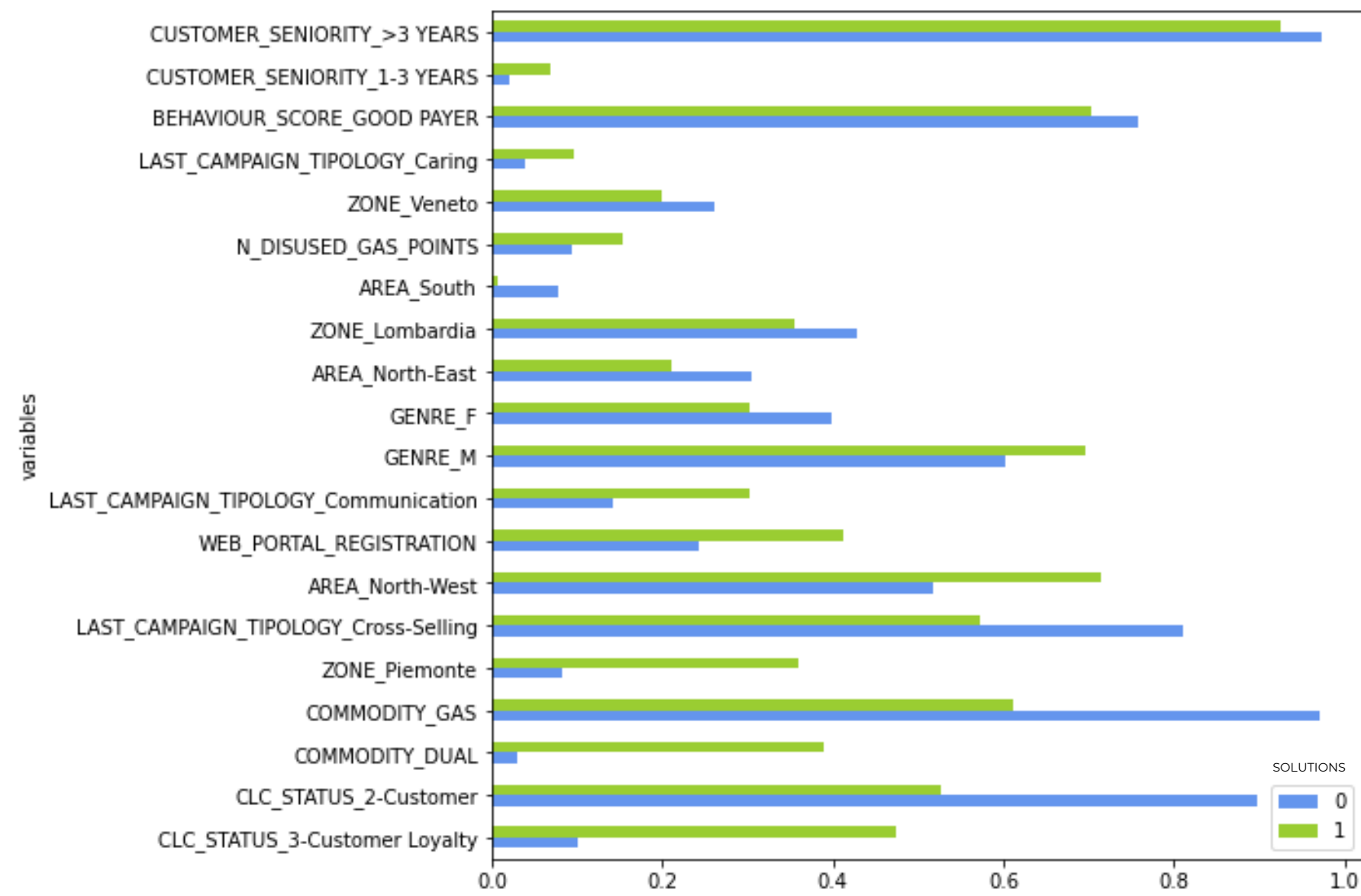
- Mean analysis
- Random Forest feature importance



# Propensity model

## Features selection: mean analysis

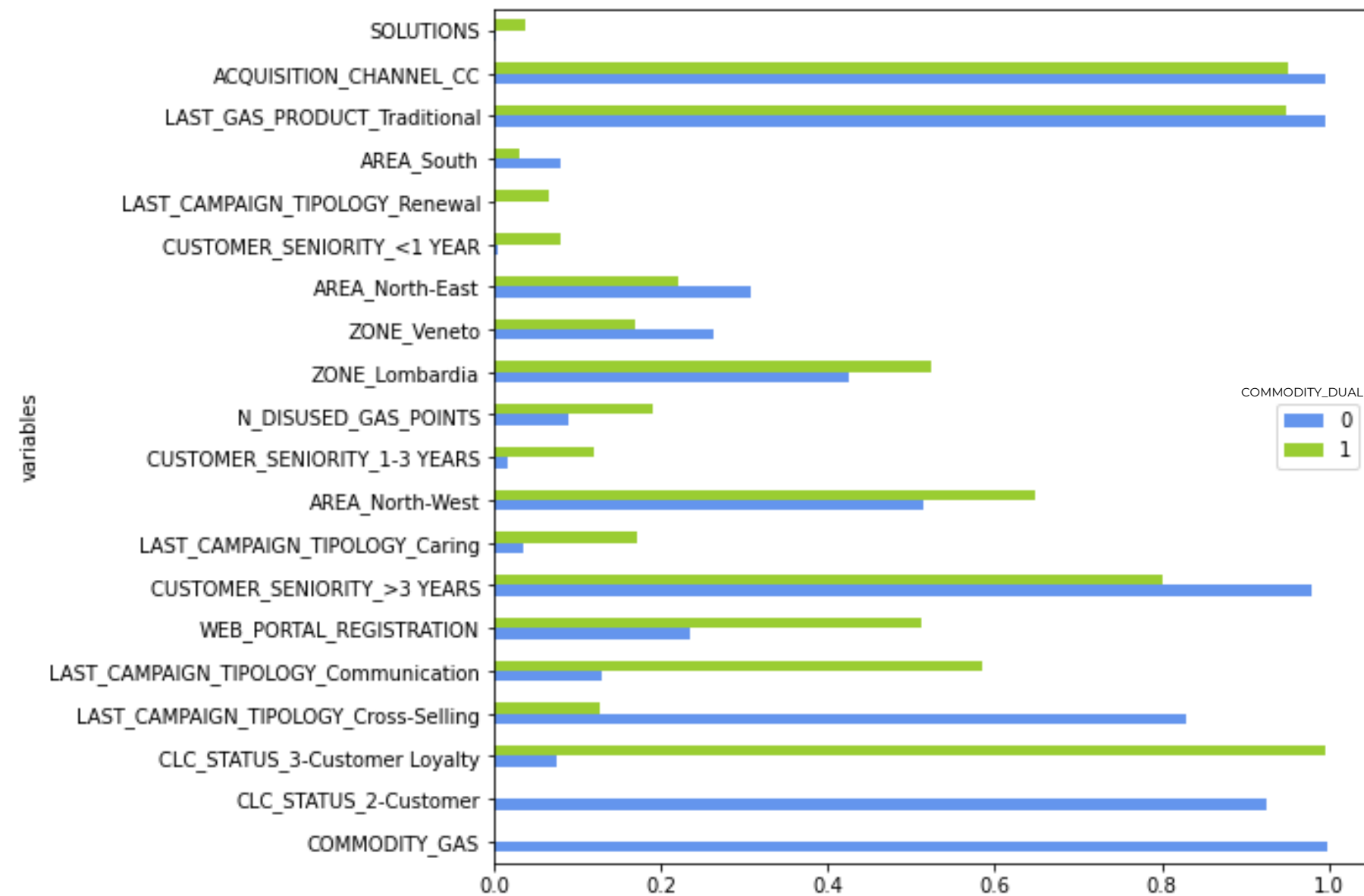
"SOLUTIONS"



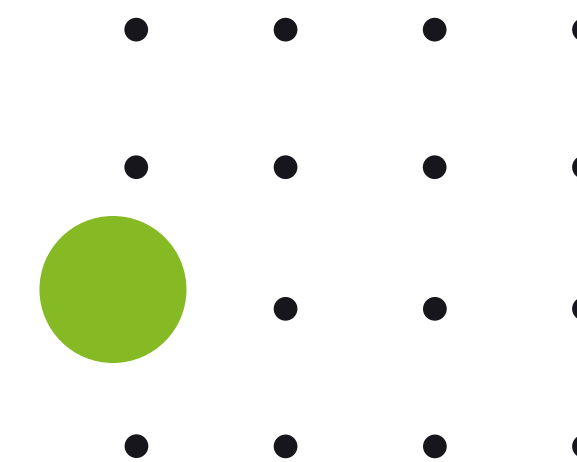
# Propensity model

## Features selection: mean analysis

"COMMODITY\_DUAL"

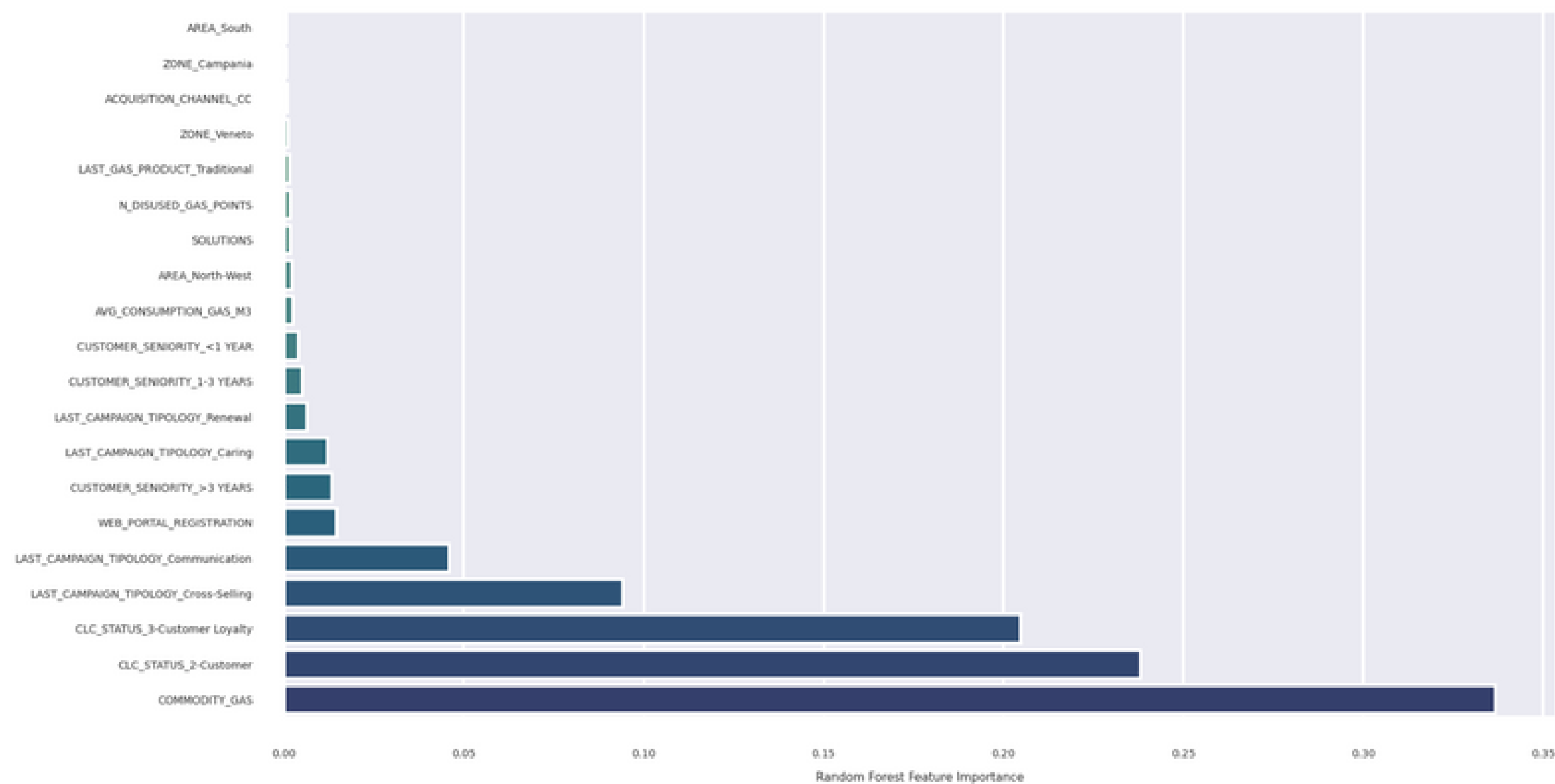


# Propensity model

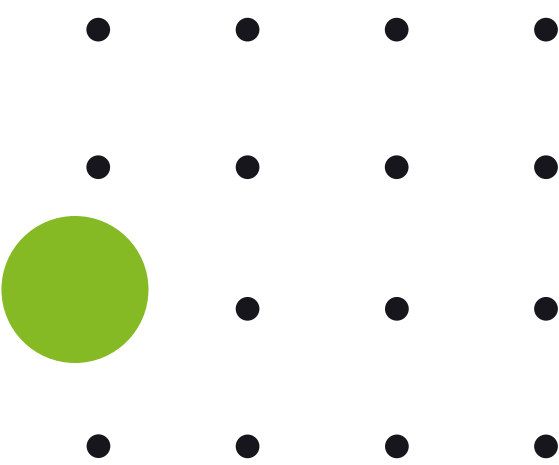


Features selection: Random Forest features importance

y: COMMODITY\_DUAL

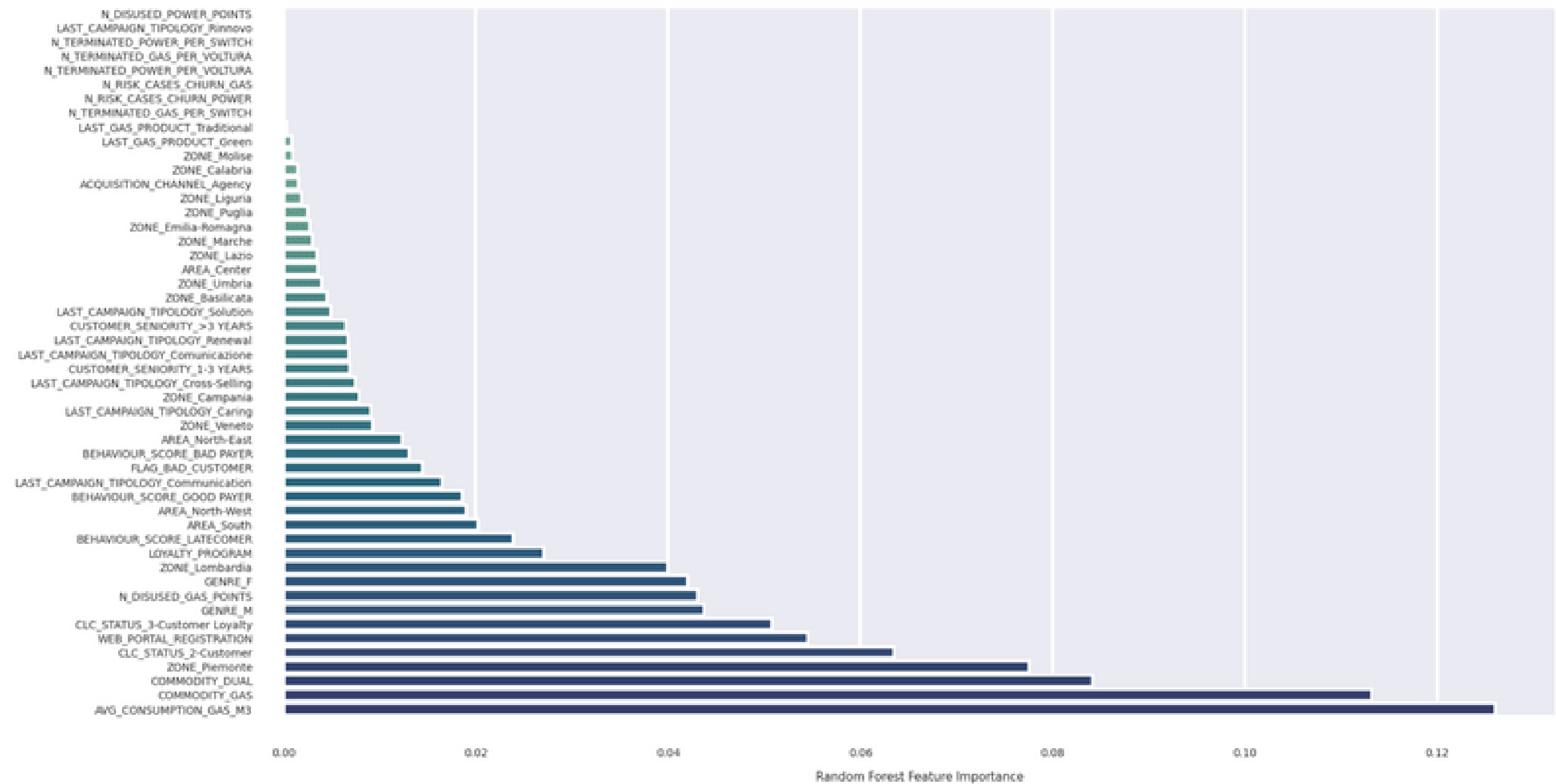


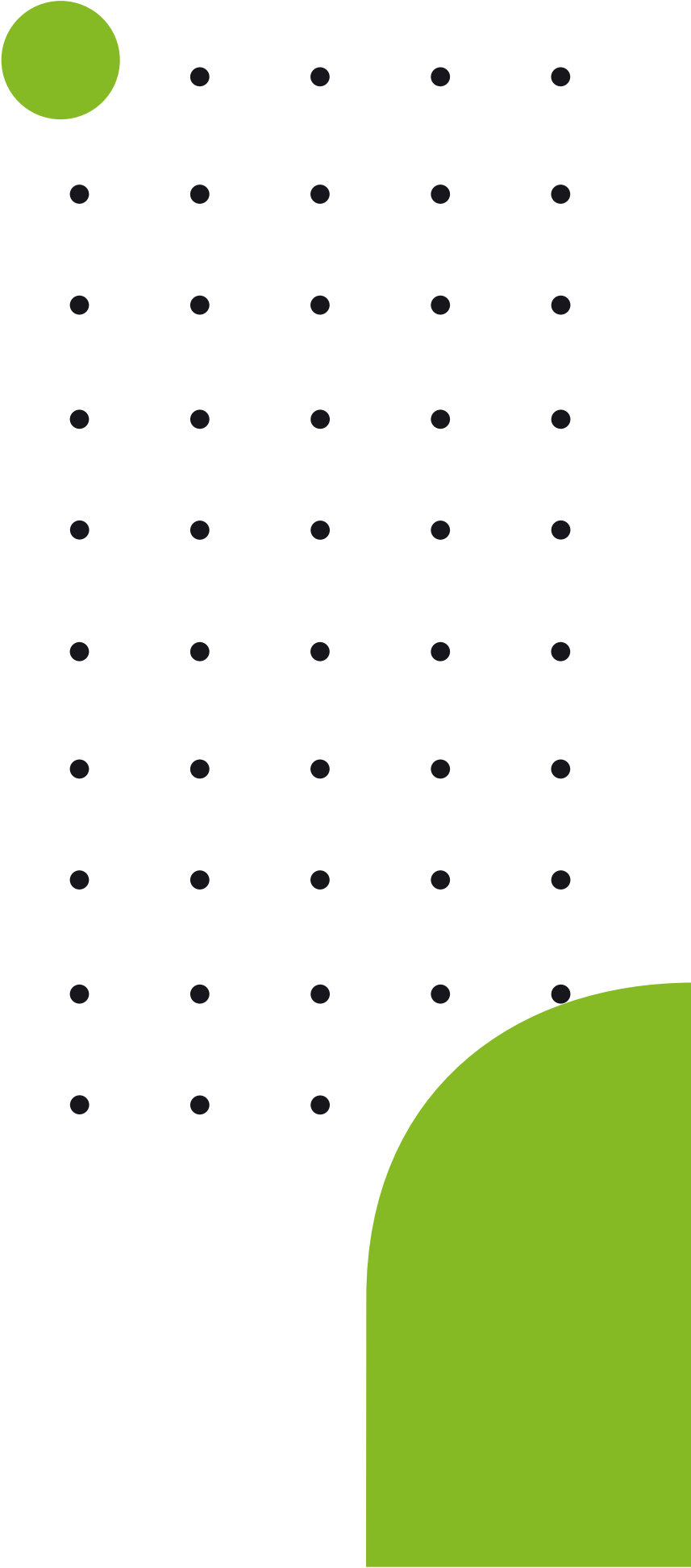
# Propensity model



## Features selection: Random Forest features importance

y: SOLUTIONS





After removing the variables that lacked an economic meaning together with the strongly correlated ones, the resulting selected features were the following:

### SOLUTIONS

['AVG\_CONSUMPTION\_GAS\_M3',  
'COMMODITY\_DUAL', 'ZONE\_Piemonte',  
'WEB\_PORTAL\_REGISTRATION', 'AREA\_North-West',  
'CLC\_STATUS\_3-Customer Loyalty',  
'BEHAVIOUR\_SCORE\_GOOD PAYER',  
'LOYALTY\_PROGRAM', 'AREA\_SOUTH',  
'ZONE\_VENETO',  
'LAST\_CAMPAGN\_TIPOLOGY\_Caring', 'AREA\_North-East',  
'LAST\_CAMPAGN\_TIPOLOGY\_Cross-Selling', 'CUSTOMER\_SENIORITY\_>3 YEARS',  
'CUSTOMER\_SENIORITY\_<1 YEAR',  
'ACQUISITION\_CHANNEL\_CC',  
'BEHAVIOUR\_SCORE\_BAD PAYER', 'AREA\_CENTER' ]

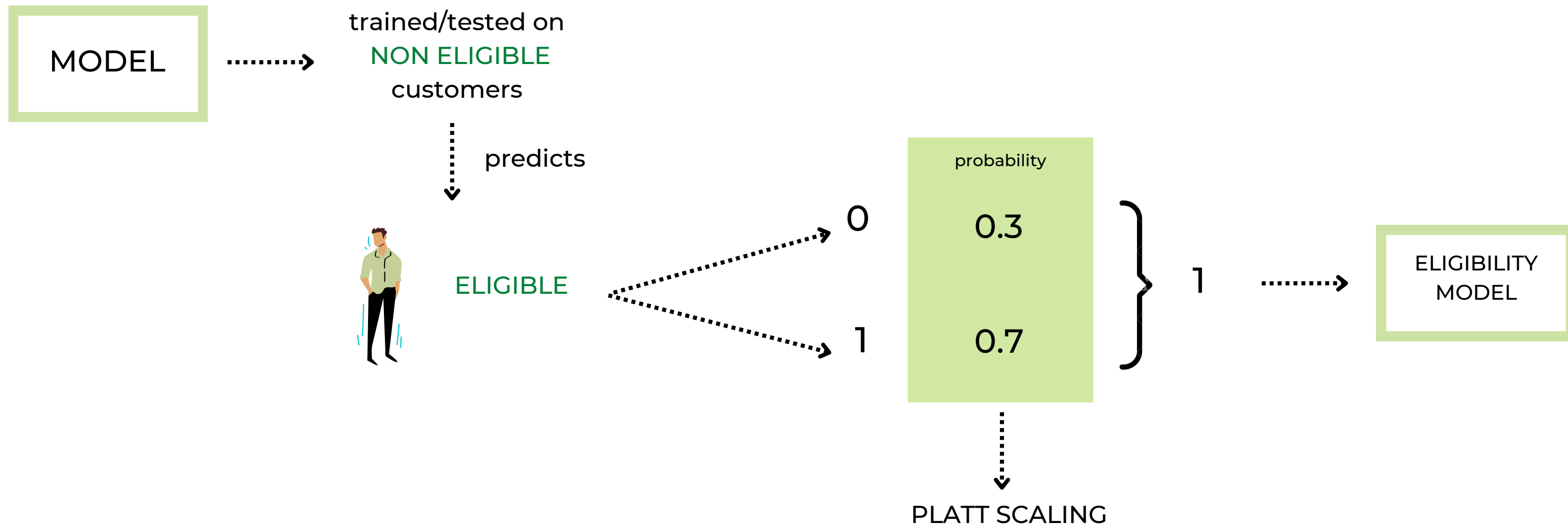
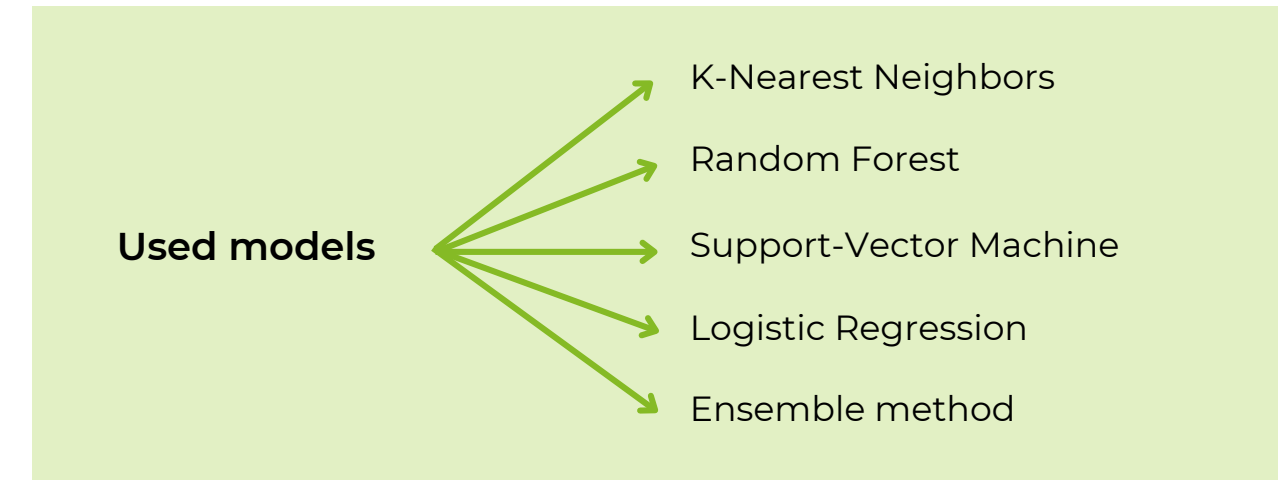
### COMMODITY\_DUAL

['CLC\_STATUS\_3-Customer Loyalty',  
'AREA\_North-West',  
'WEB\_PORTAL\_REGISTRATION',  
'LAST\_CAMPAGN\_TIPOLOGY\_Cross-Selling',  
'N\_DISUSED\_GAS\_POINTS',  
'LAST\_CAMPAGN\_TIPOLOGY\_Caring',  
'SOLUTIONS', 'CUSTOMER\_SENIORITY\_>3 YEARS',  
'LAST\_CAMPAGN\_TIPOLOGY\_Renewal',  
'CUSTOMER\_SENIORITY\_1-3 YEARS',  
'AVG\_CONSUMPTION\_GAS\_M3',  
'LAST\_GAS\_PRODUCT\_Traditional',  
"COMMODITY\_DUAL"]

# How is the propensity measure captured?

Y → SOLUTIONS  
Y → COMMODITY\_DUAL

X → important  
X → features



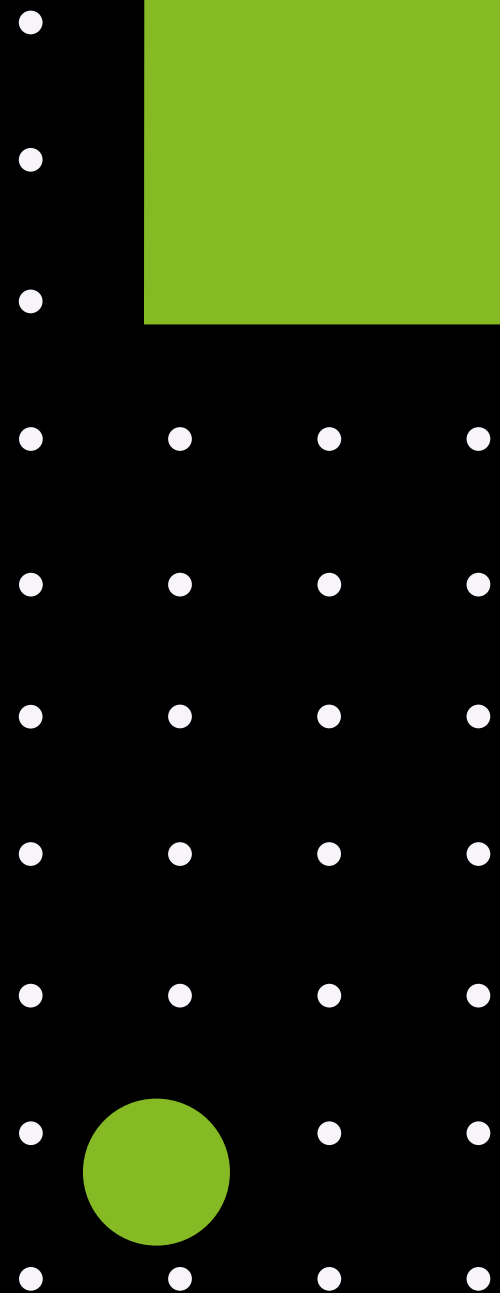


# Models' performance metrics

## COMMODITY\_DUAL

	Accuracy	Precision	Recall
Random Forest	0.89	0.83	0.97
KNN	0.86	0.81	0.94
SVM	0.89	0.83	0.97
Logistic Regression	0.89	0.83	0.97
<b>Ensemble method</b>	<b>0.89</b>	<b>0.83</b>	<b>0.98</b>

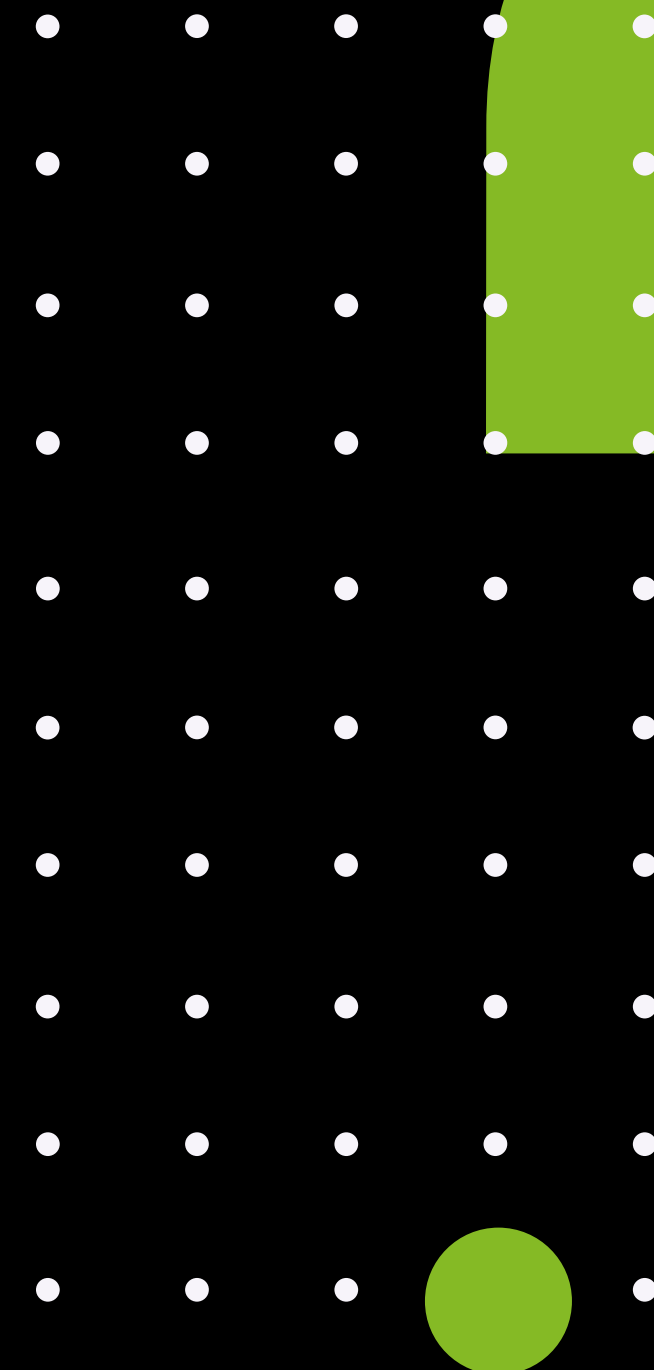
Most important evaluation metric when the cost of false negative is high (the cost of NOT contacting someone who is likely to give a positive response to the campaign is higher than the cost of contacting a customer who is not willing to sign the contract).



# Models' performance metrics

## SOLUTIONS

	Accuracy	Precision	Recall
Random Forest	0.68	0.72	0.57
KNN	0.67	0.68	0.62
SVM	0.68	0.67	0.69
Logistic Regression	0.65	0.65	0.63
<b>Ensemble method</b>	<b>0.7</b>	<b>0.68</b>	<b>0.71</b>



# Eligibility model

1

Customers who are at risk of being churned, who are leaving and who did not give the privacy consensus will not be contacted and, thus, they are deleted from the dataset.

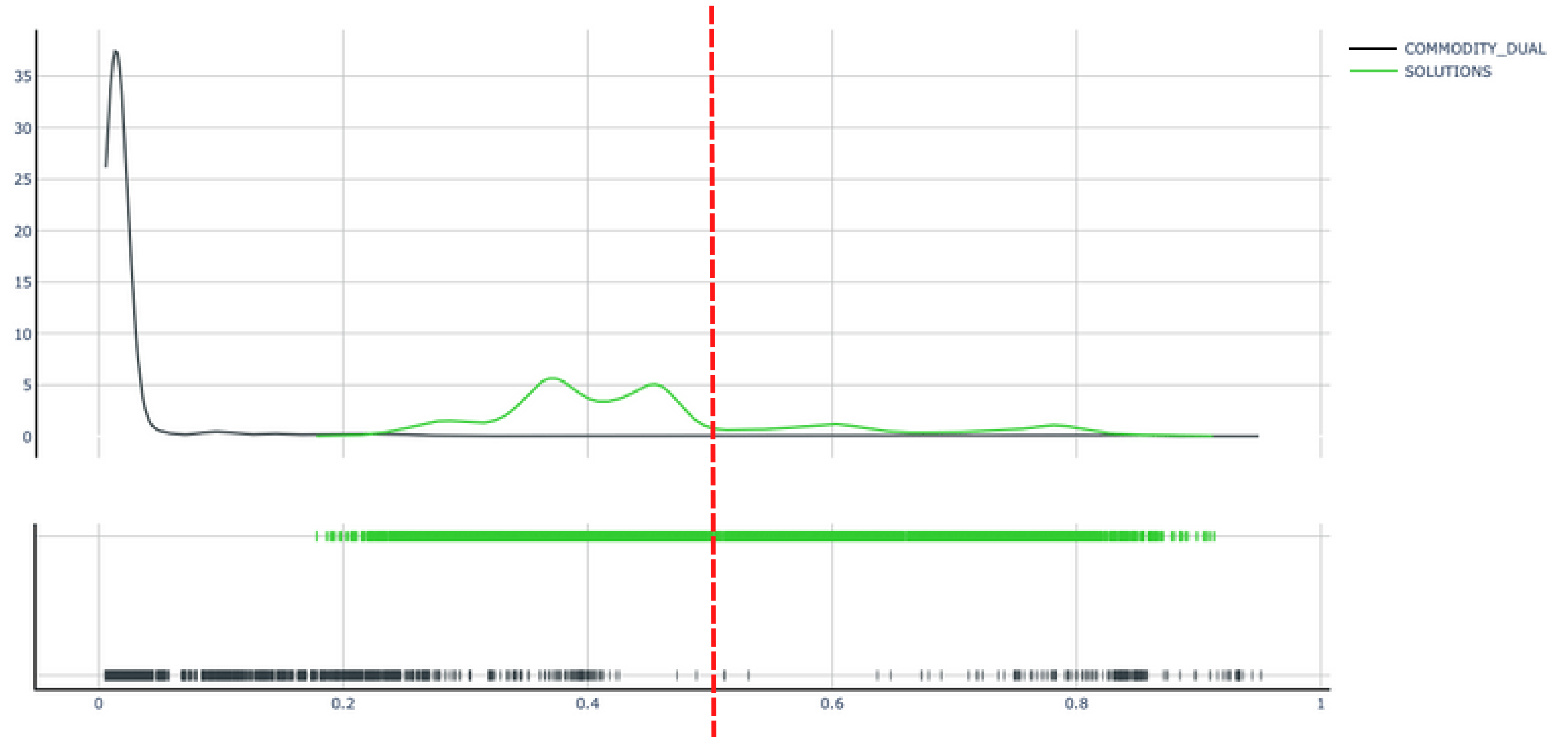
2

Six datasets were created according to the general rule (*those customers whose last contact dates back to at least N months ago are eligible*) and, according to the type of channel, those customers who did not validate their phone number or their e-mail are deleted.

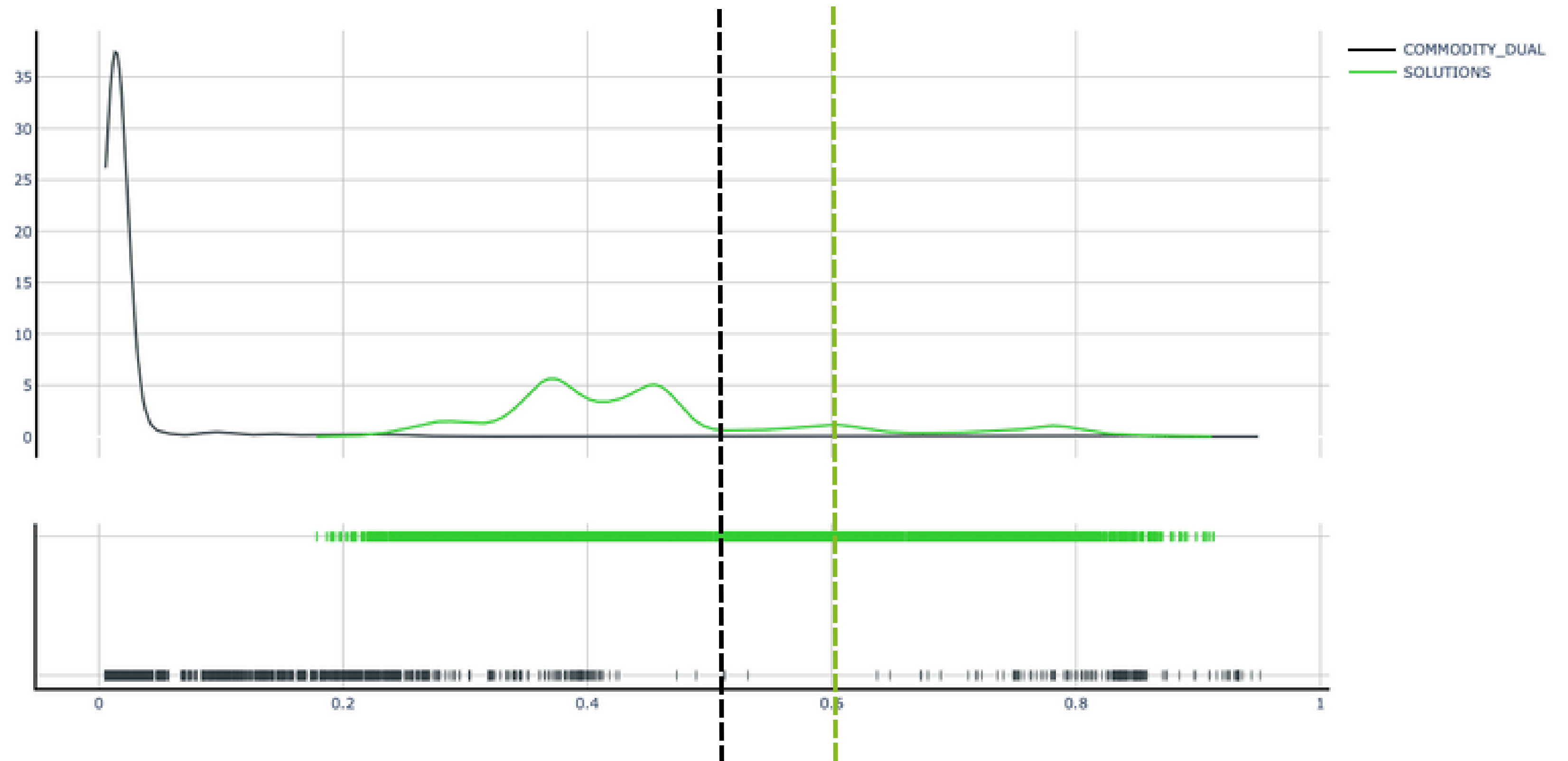
3

Campaign and cross-campaigns rules are applied to design the monthly contact strategy.

# MODEL FLEXIBILITY



# MODEL FLEXIBILITY





# Model Flexibility

## No budget

- Efficiency Oriented
- Only the most propense are selected



≈ 5000€

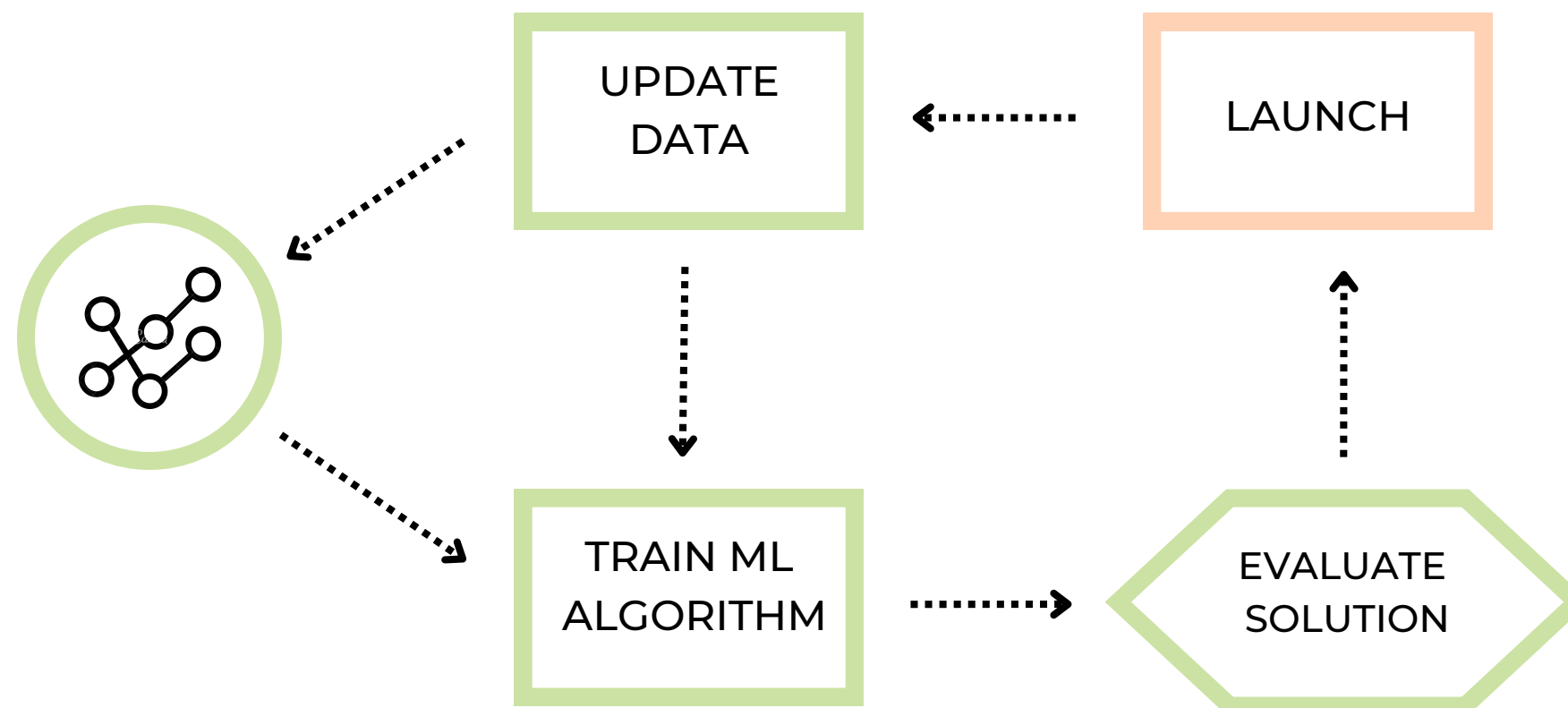
## Defined Budget

- Allows more variability
- Also customers with a lower probability may be contacted

# Conclusions and future developments

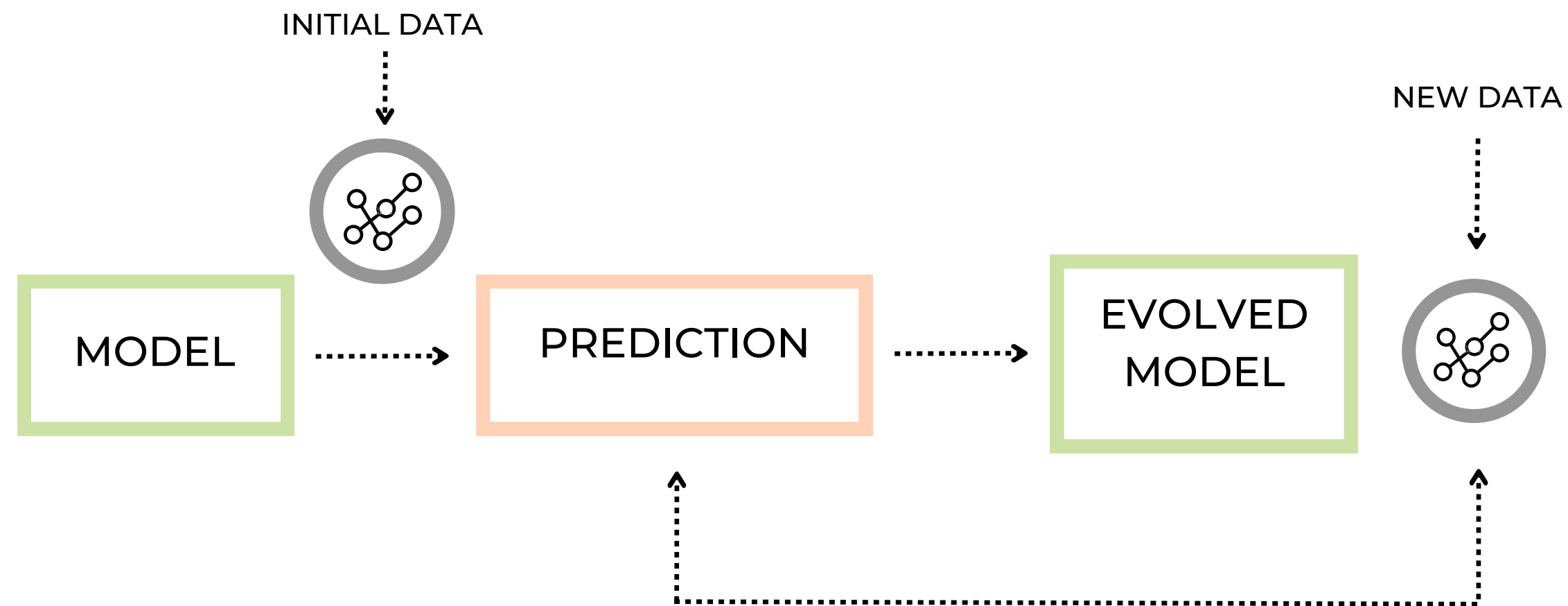
We have been able to optimize the marketing strategy of the company prioritizing the contact of a client based on their propensity to sign the contract

The resources and time saved may give the company a competitive advantage with respect to the competitors that still do not follow this data-driven approach.



# Conclusions and future developments

## ONLINE LEARNING





# Our Team



Fabiana Caccavale



Lorenzo Meloncelli



Marco Amadori

*Thank you for the attention!*

