

# Energy effectiveness of a mortgage portfolio

Jakub Vondráček, Daniel Till

2.1.2023

# Overview

## 1 Data analysis

- Data exploration
- Augmenting the dataset
- Datasets

## 2 Models

- Used models
- Model evaluation

## 3 Results

- Imbalanced data
- Predefined distribution
- Feature Importance
- Balanced data

## 4 Appendix

# Data exploration

Value	Count	Frequency
D	8650444	43.8%
C	5789245	29.3%
E	3602636	18.3%
F	919699	4.7%
B	476051	2.4%
G	281359	1.4%
A	11011	0.1%

Table 1: CURRENT\_ENERGY\_RATING – frequencies in the cleaned dataset.

$$A < B < C < D < E < F < G$$

Client has the following data available:

- BUILDING\_AGE\_CLASS,
- FLOOR\_AREA,
- NUMBER\_HABITABLE\_ROOMS,
- PROPERTY\_TYPE,
- POSTCODE (ADDRESS).

# Augmenting the dataset

- Left out the variables: POSTCODE.
- Included the variables:
  - A\_PROP\_POSTCODE, ..., G\_PROP\_POSTCODE,
  - LOCAL\_AUTHORITY\_LABEL,
  - BUILT\_FORM,
  - POSTCODE\_PROPORTIONS\_ARE\_RELIABLE\_IND,
  - IS\_EPC\_LABEL\_BEFORE\_2008\_INCL,
  - POSTCODE\_COUNT.

# Datasets

- Small dataset - client given data + augmentation.
- Big dataset - all available features.

# Splitting the dataset

- Training set - 25 % of the observations.
- Validation set - 25 % of the observations.
- Test set - 50 % of the observations.

# Adjusting class weights

- Imbalanced dataset.
- Predefined distribution (predefined by the client).
- Balanced dataset.



# Used models

- Ordinal regression
- Support vector machine
- Gradient boosted decision trees

# Model evaluation

## Used metrics

- Accuracy
  - Proportion of correctly classified examples
- Ranked Probability Score
  - Metric that takes into account ordinality
- Confusion matrix
  - Precision, Recall, F1 Score

# Imbalanced data

	Light Gradient Boosting Machine Big	Baseline Model Small	Light Gradient Boosting Machine Small	Ordinal Regression Small	Support Vector Machine Small
A	0.311	0.413	<b>0.358</b>	0.369	0.405
B	0.093	0.247	<b>0.127</b>	0.132	0.154
C	0.040	0.088	0.056	<b>0.054</b>	0.077
D	0.029	<b>0.027</b>	0.032	0.035	0.041
E	0.060	0.113	0.080	0.097	<b>0.061</b>
F	0.081	0.259	0.170	0.203	<b>0.162</b>
G	0.071	0.421	<b>0.293</b>	0.343	0.295
Weighted Avg	0.042	0.082	<b>0.060</b>	0.066	0.067

Table 2: Ranked Probability Score for models with imbalanced dataset.

Light Gradient Boosting Machine Big	Baseline Model Small	Light Gradient Boosting Machine Small	Ordinal Regression Small	Support Vector Machine Small
0.657	0.438	<b>0.565</b>	0.533	0.519

Table 3: Accuracy for models with imbalanced dataset.

# Imbalanced data

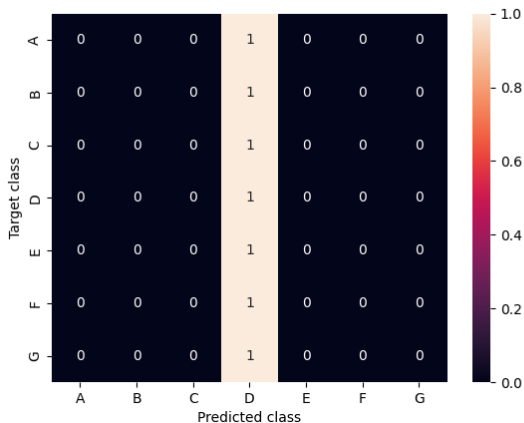


Figure 1: Confusion Matrix for Baseline Small Model with imbalanced dataset.

# Imbalanced data

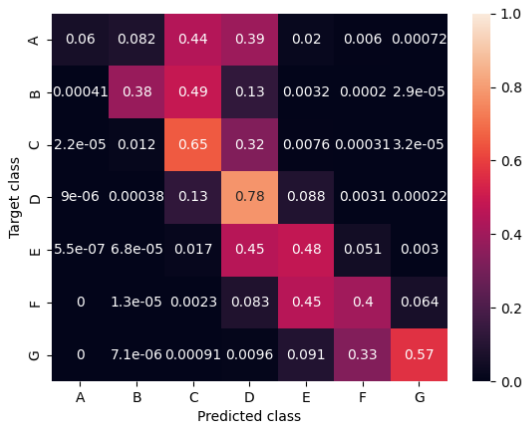


Figure 2: Confusion Matrix for Light Gradient Boosting Machine Model Big with imbalanced dataset.

# Imbalanced data

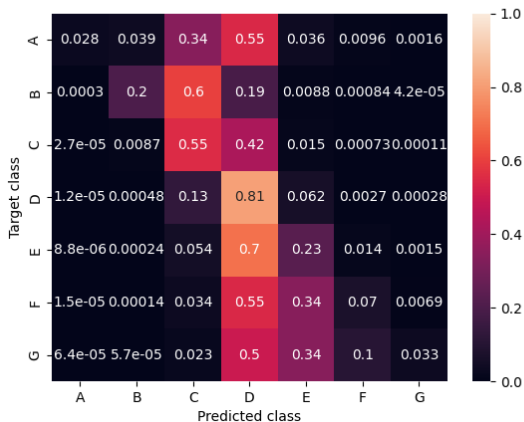


Figure 3: Confusion Matrix for Light Gradient Boosting Machine Small Model with imbalanced dataset.

# Imbalanced data

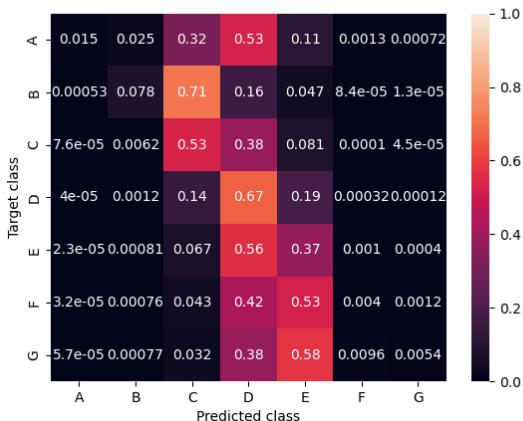


Figure 4: Confusion Matrix for Support Vector Machine Small Model with imbalanced dataset.

# Imbalanced data

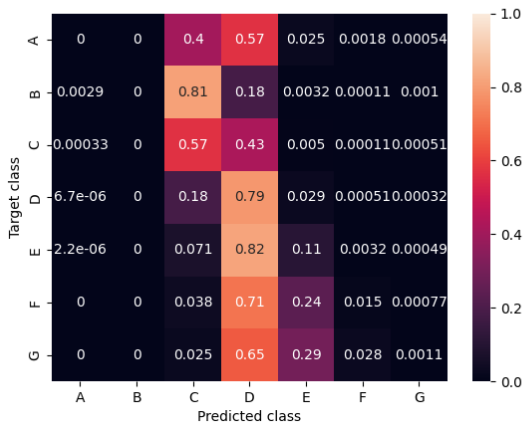


Figure 5: Confusion Matrix for Ordinal Regression Small Model with imbalanced dataset.



# Predefined distribution

	D	C	E	F	B	G	A
Predefined	38.9%	22.7%	17.2%	3.9%	16.3%	0.7%	0.3%
Our data	43.8%	29.3%	18.3%	4.7%	2.4%	1.4%	0.1%

Table 4: Distribution of target predefined by the client and in our dataset.

	Light Gradient Boosting Machine Big	Light Gradient Boosting Machine Small
Accuracy	0.632	0.543
Ranked Probability Score	0.045	0.063

Table 5: Accuracy and RPS for models with predefined distribution.

# Predefined distribution

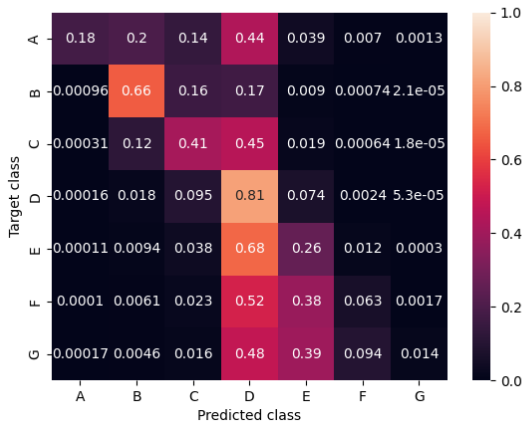


Figure 6: Confusion matrix for Light Gradient Boosting Machine Small trained with the predefined distribution.

# Predefined distribution

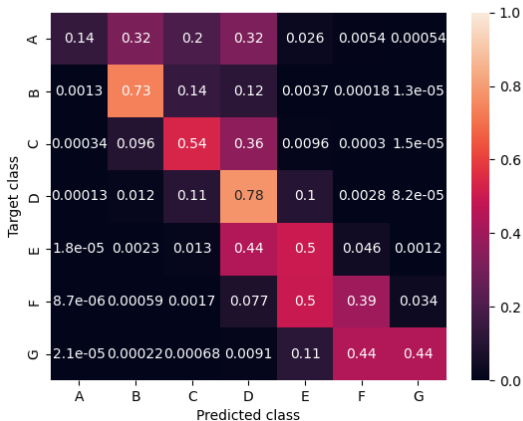


Figure 7: Confusion matrix for Light Gradient Boosting Machine Big trained with the predefined distribution.

# Feature Importance

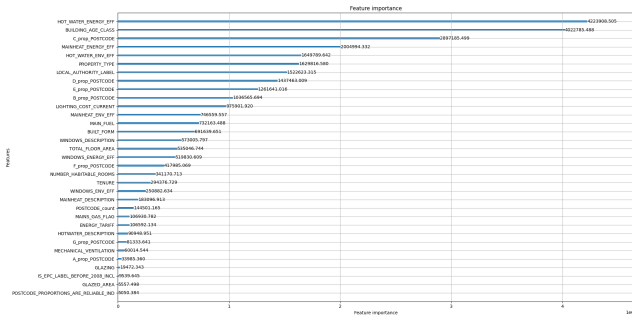
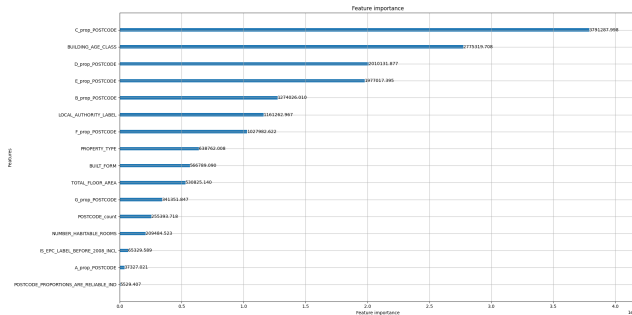


Figure 9: Feature Importance for Light Gradient Boosting Machine Big Model with imbalanced dataset.

# Feature Importance



**Figure 10:** Feature Importance for Light Gradient Boosting Machine Small Model with imbalanced dataset.

# Feature Importance

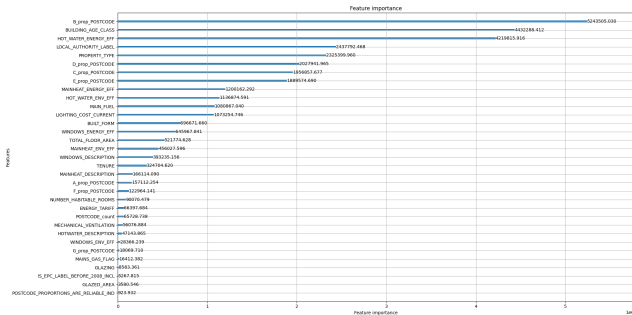
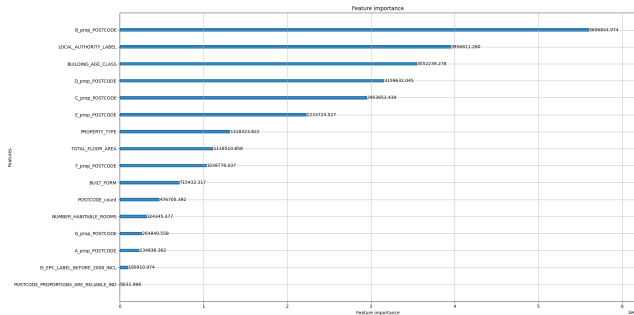


Figure 11: Feature Importance for Light Gradient Boosting Machine Big Model with predefined distribution.

# Feature Importance



**Figure 12:** Feature Importance for Light Gradient Boosting Machine Small Model with predefined distribution.

# Balanced data

	Light Gradient Boosting Machine Big	Baseline Model Small	Light Gradient Boosting Machine Small	Support Vector Machine Small
A	0.312	0.413	<b>0.260</b>	0.303
B	0.096	0.247	<b>0.088</b>	0.111
C	0.041	0.088	<b>0.085</b>	0.101
D	0.029	<b>0.027</b>	0.089	0.083
E	0.061	0.113	0.079	<b>0.065</b>
F	0.083	0.259	0.095	<b>0.094</b>
G	0.073	0.421	<b>0.170</b>	0.177
Weighted Avg	0.043	<b>0.082</b>	0.087	0.087

**Table 6:** Ranked Probability Score for models with balanced dataset.

Light Gradient Boosting Machine Big	Baseline Model Small	Light Gradient Boosting Machine Small	Support Vector Machine Small
0.574	<b>0.438</b>	0.422	0.422

**Table 7:** Accuracy for models with balanced dataset.



# Balanced data

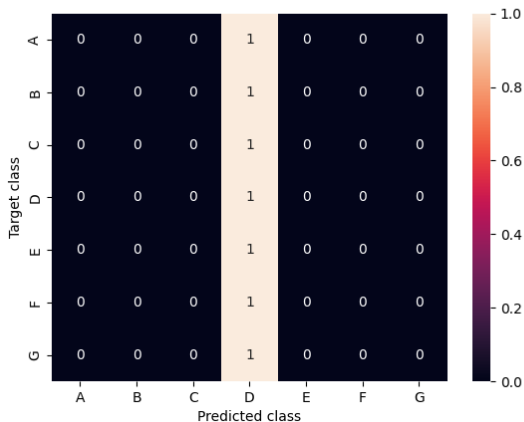


Figure 13: Confusion Matrix for Baseline Small Model (same as with unbalanced data)

# Balanced data

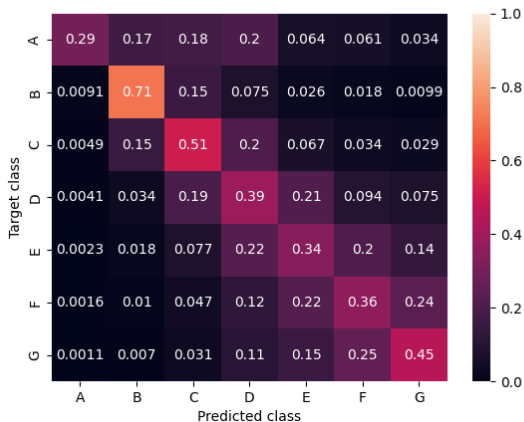


Figure 14: Confusion Matrix for Light Gradient Boosting Machine Small Model with balanced dataset.

# Balanced data

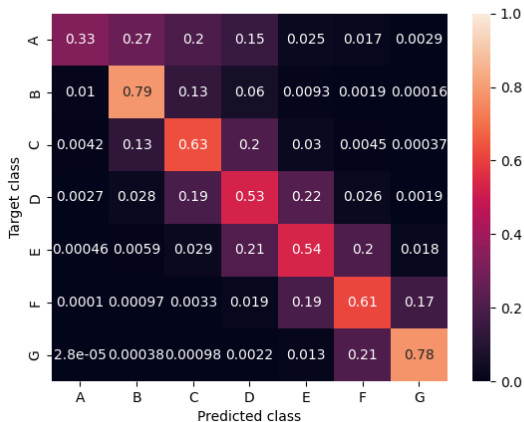


Figure 15: Confusion Matrix for Light Gradient Boosting Machine Model Big with balanced dataset.

# Balanced data

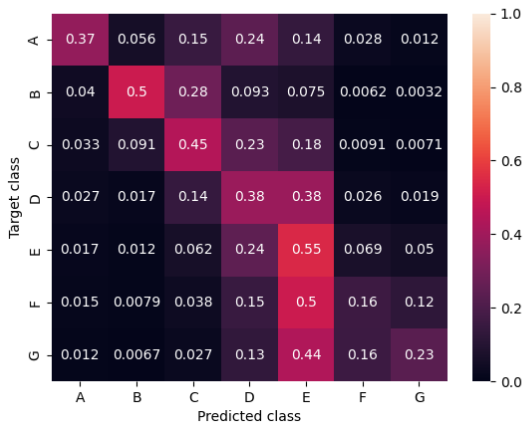


Figure 16: Confusion Matrix for Support Vector Machine Small Model with balanced dataset.

# Feature Importance

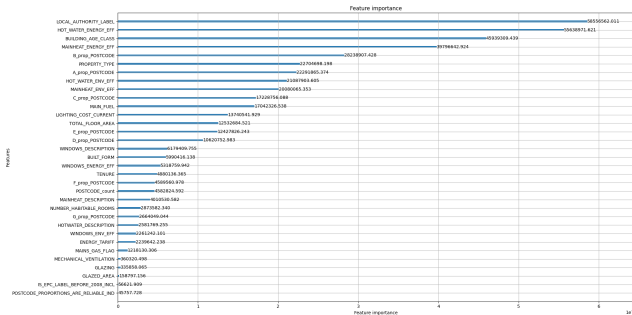
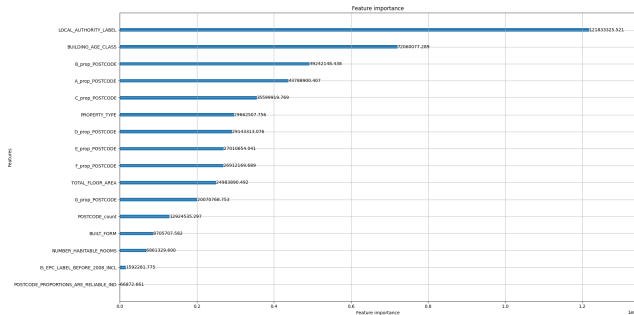


Figure 17: Feature Importance for Light Gradient Boosting Machine Model Big with balanced dataset.

# Feature Importance



**Figure 18:** Feature Importance for Light Gradient Boosting Machine Small Model with balanced dataset.

Thank you for your attention

# Appendix

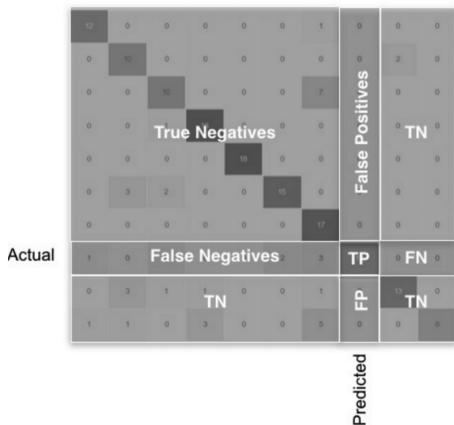


Figure 19: Schematic representation of a confusion matrix and the terms True Positive, True Negative, False Positive and False Negative.



# Appendix

$$\text{RPS}((x^T, y)^T) = \frac{1}{R+1} \sum_{k=0}^R \left[ \left( \sum_{j=0}^k p_j(x) \right) - \left( \sum_{j=0}^k y^{(j)} \right) \right]^2,$$

$$\text{RPS}_i = \frac{1}{|i|} \sum_{((x_l^T, y_l)^T) \in i} \text{RPS}((x_l^T, y_l)^T),$$

$$\text{Accuracy} = \frac{\sum_{i=0}^R c_{ii}}{\sum_{i,j=0}^R c_{ij}},$$

# Appendix

- True Positives ( $TP_i$ ):  $c_{ii}$ ,
- True Negatives ( $TN_i$ ):  $\sum_{j,k \in \{0, \dots, R\} \setminus \{i\}} c_{jk}$ ,
- False Positives ( $FP_i$ ):  $\sum_{j \in \{0, \dots, R\} \setminus \{i\}} c_{ji}$ ,
- False Negatives ( $FN_i$ ):  $\sum_{j \in \{0, \dots, R\} \setminus \{i\}} c_{ij}$ .

# Appendix

$$\text{Precision}_i = \frac{TP_i}{TP_i + FP_i},$$

$$\text{Recall}_i = \frac{TP_i}{TP_i + FN_i},$$

$$\text{F1 Score}_i = 2 \frac{\text{Precision}_i \cdot \text{Recall}_i}{\text{Precision}_i + \text{Recall}_i},$$

where  $i = 0, \dots, R$ .