

Institutional sector classification

A Machine Learning Application



Paolo **Massaro**

Divisione Informazioni Anagrafiche
Dipartimento **ECS**

Oliver **Giudice**

Divisione Ricerca sulle Tecnologie Avanzate
Dipartimento **IT**



BANCA D'ITALIA

EUROSISTEMA

Problem statement

Given

a set of **features**
of a company

Numeric and non-numeric: name, number of employees, balance sheet data, whether publicly held or not, etc.

Determine

the appropriate **SAE** code
to assign to it

SAE="SETTORE DI ATTIVITA' ECONOMICA" is a code defined by Circ. 140/97 meant to cluster companies into one of 116 "institutional sectors" (e.g., public institution, productive company, financial holding, etc.)

Problem statement

Given

a set of **features**
of a company

Numeric and non-numeric: name, number of employees, balance sheet data, whether publicly held or not, etc.

Determine

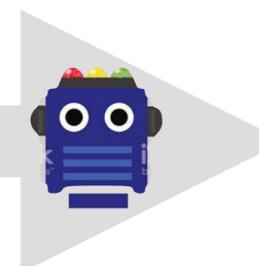
the appropriate **SAE** code
to assign to it

SAE="SETTORE DI ATTIVITA' ECONOMICA" is a code defined by Circ. 140/97 meant to cluster companies into one of 116 "institutional sectors" (e.g., public institution, productive company, financial holding, etc.)



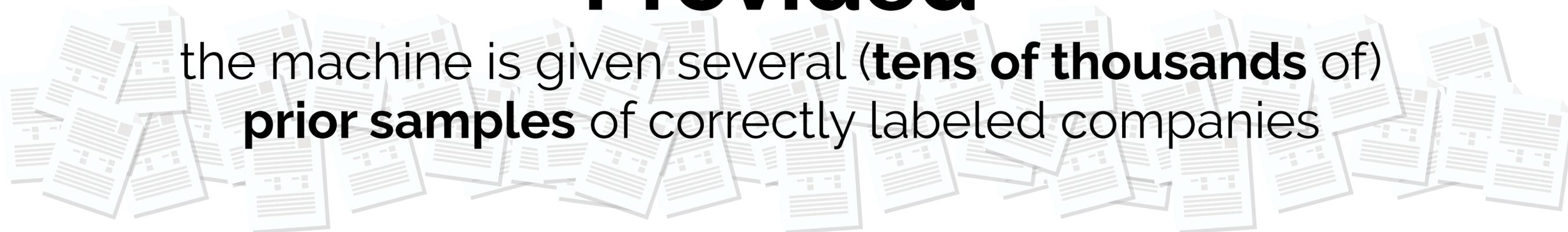
Machine Learning approach

We start from existing **data**; a "**machine learning model**" is trained from companies already labeled (by hand); on the basis of this "past experience" it learns to predict what SAE any new company belongs to



Provided

the machine is given several (**tens of thousands** of)
prior samples of correctly labeled companies



Why and when should ML help here?



Machine Learning approach



We start from existing **data**; a "**machine learning model**" is trained from companies already labeled (by hand); on the basis of this "past experience" it learns to predict what SAE any new company belongs to



AS-IS: Who classifies companies into SAEs?

Type of company	Classified by Institution	Data trustworthiness
Public Administrations	ISTAT	authoritative
Supervised Entities	Bank of Italy	authoritative
Other (Productive Companies, non-supervised Companies, etc...)	Financial Intermediaries	may be: <ul style="list-style-type: none">incorrect ← fixinconsistent ← spotstale ← updatemissing (~30%) ← autofill

Why and when should ML help here?



Machine Learning approach



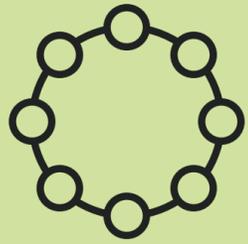
We start from existing **data**; a "**machine learning model**" is trained from companies already labeled (by hand); on the basis of this "past experience" it learns to predict what SAE any new company belongs to



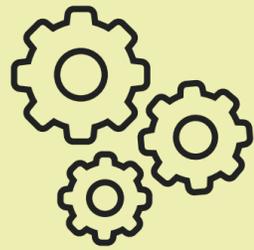
AS-IS: Who classifies companies into SAEs?

Type of company	Classified by Institution	Data trustworthiness
Public Administrations	ISTAT	authoritative
Supervised Entities	Bank of Italy	authoritative
Other (Productive Companies, non-supervised Companies, etc...)	Financial Intermediaries	may be: ▶ incorrect ▶ inconsistent ▶ stale ▶ missing (~30%)

fix
spot
update
autofill



Data



Preprocessing



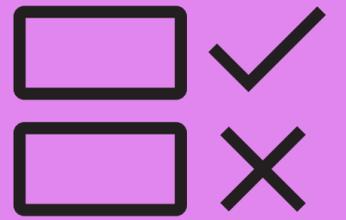
Feature extraction



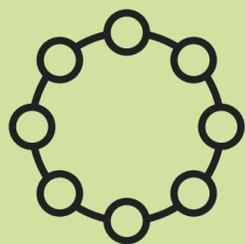
Imbalanced learning



Classification



Results



Data

Preprocessing

Feature extraction

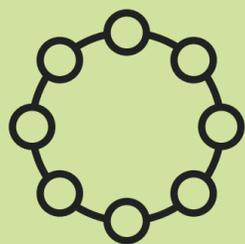
Imbalanced learning

Classification

Results

Original datasets

Dataset	#	Origin
Anagrafe Soggetti	42M	Bank of Italy
Listed Companies	1K	Bank of Italy
ATECO	3.6M	Ag. Entrate
Balance Sheet et al.	2.2M	CERVED
Info Imprese	2.2M	INFOCAMERE



Data

Preprocessing

Feature extraction

Imbalanced learning

Classification

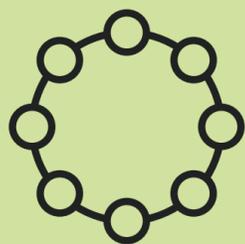
Results

Original datasets



Platform

Dataset	#	Origin
Anagrafe Soggetti	42M	Bank of Italy
Listed Companies	1K	Bank of Italy
ATECO	3.6M	Ag. Entrate
Balance Sheet et al.	2.2M	CERVED
Info Imprese	4.8M	INFOCAMERE



Data

Preprocessing

Feature extraction

Imbalanced learning

Classification

Results

Data ingestion (ETL)



Extract



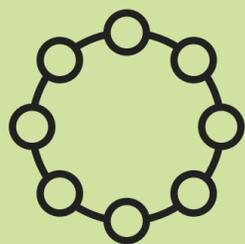
Transform

Load

Platform

Op

Dataset	#	Origin
Anagrafe Soggetti	42M	Bank of Italy
Listed Companies	1K	Bank of Italy
ATECO	3.6M	Ag. Entrate
Balance Sheet et al.	2.2M	CERVED
Info Imprese	4.8M	INFOCAMERE



Data

Preprocessing

Feature extraction

Imbalanced learning

Classification

Results

Data ingestion (ETL)



Extract

Transform

Load



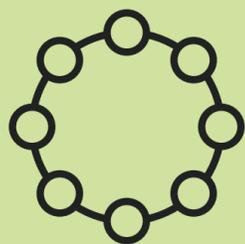
Platform

Op

Tools

Dataset # Origin

Anagrafe Soggetti	42M	Bank of Italy
Listed Companies	1K	Bank of Italy
ATECO	3.6M	Ag. Entrate
Balance Sheet et al.	2.2M	CERVED
Info Imprese	4.8M	INFOCAMERE



Data

Preprocessing

Feature extraction

Imbalanced learning

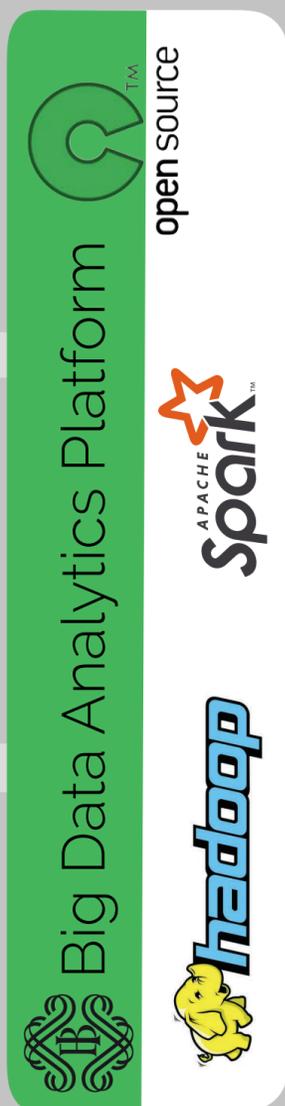
Classification

Results

Data ingestion (ETL)



Extract



Transform



Load



Platform

Op

Tools

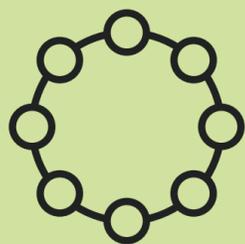
Input to ML machinery

Single text file

1.4M records, **400M**Bytes

Each record contains info about:

- a. Company structure
- b. Balance sheet
- c. Other info
- d. **SAE**



Data

Preprocessing

Feature extraction

Imbalanced learning

Classification

Results

Data ingestion (ETL)



Extract



Transform



Load



Inside the ML machinery, for each company

Company structure

15 numeric features

- ▶ Num. of employees
- ▶ PA-owned shares
- ▶

Balance sheet

14 numeric features

- ▶ Share capital
- ▶ Personnel costs
- ▶

Other info

3 structured features

- ▶ Listed (y/n)
- ▶ ATECO
- ▶ Comune

Name & notes

2 textual features

- ▶ company name
- ▶ **balance notes**

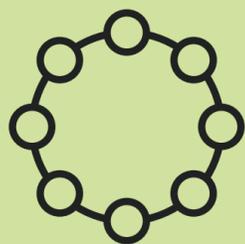


SAE

Platform

Op

Tools



Data

Preprocessing

Feature extraction

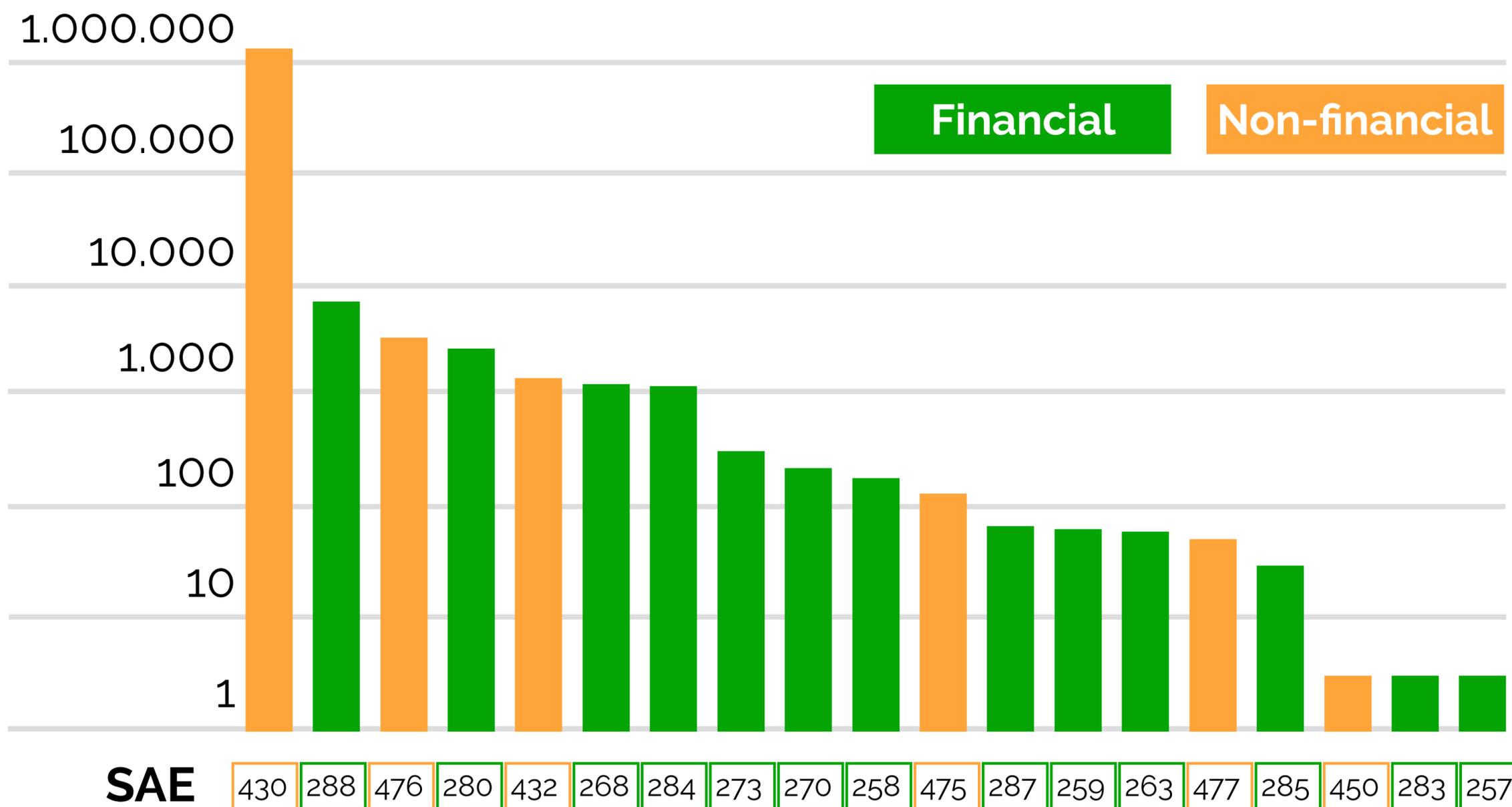
Imbalanced learning

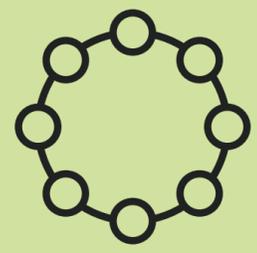
Classification

Results

SAE: (Un)balanced data

Number of companies per SAE





Data

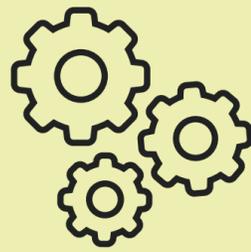
Preprocessing

Feature extraction

Imbalanced learning

Classification

Results



Preprocessing

Specific w.r.t data type

Dealing with Missing

Structured
Data

Normalization

Try to fix or infer:

- Use "zero" or average value
- Regression on other variables
- etc.

(un)structured
Data

Lemmatization
Stemming
Dictionary-checking

Ignore

It is textual and can be divided:

- ▶ company denomination (always present)
- ▶ **balance notes** (missing in almost 50% of the dataset)

Data

Preprocessing

Imbalanced learning

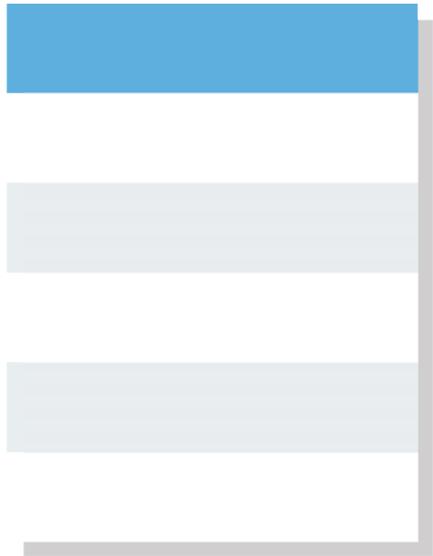
Classification

Results



Feature extraction

Types of features



Data

Preprocessing

Imbalanced learning

Classification

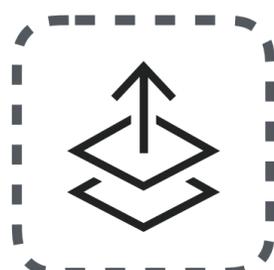
Results



Feature extraction

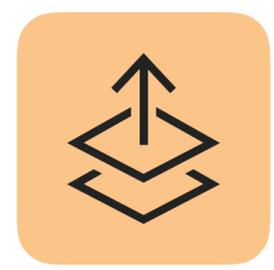
Types of features

numeric quantity



[direct]

categorical property



[one-hot-encoding]

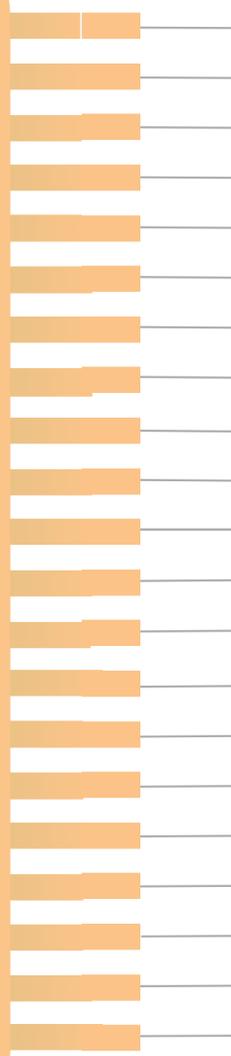
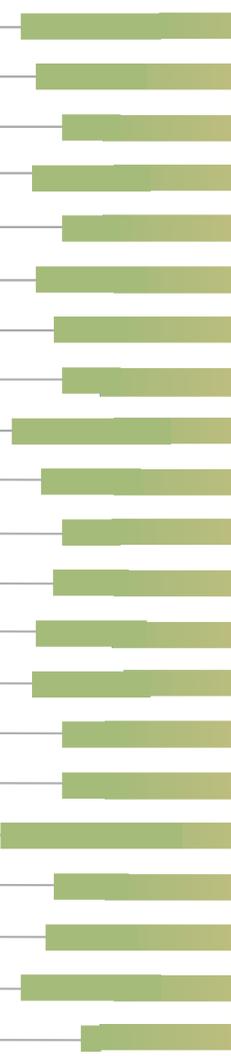
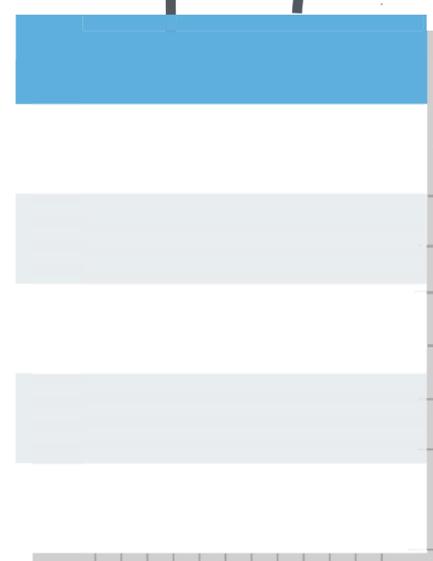
textual property

[tf-idf]

List of company features



PCA



Data

Preprocessing



Feature extraction

Imbalanced learning

Classification

Results

Data

Preprocessing

Feature extraction

Classification

Results

A couple of unbalanced classes



Imbalanced learning



Data

Preprocessing

Feature extraction

Classification

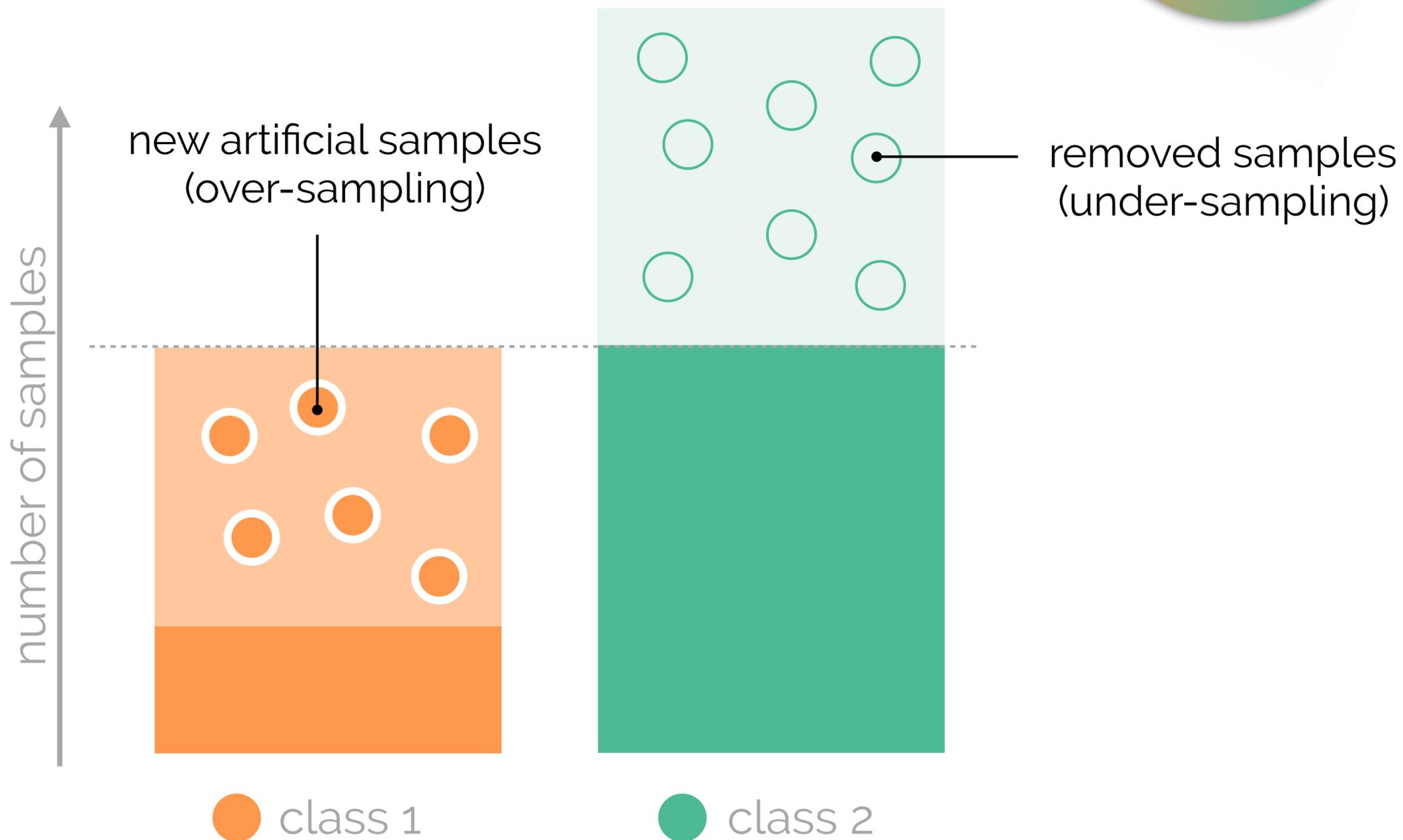
Results

Under-sampling & over-sampling



Imbalanced learning

SMOTE



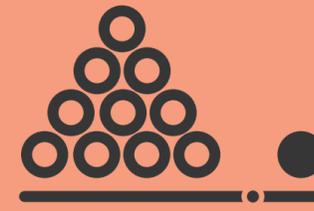
Data

Preprocessing

Feature extraction

Classification

Results



Imbalanced learning

Data

Preprocessing

Feature extraction

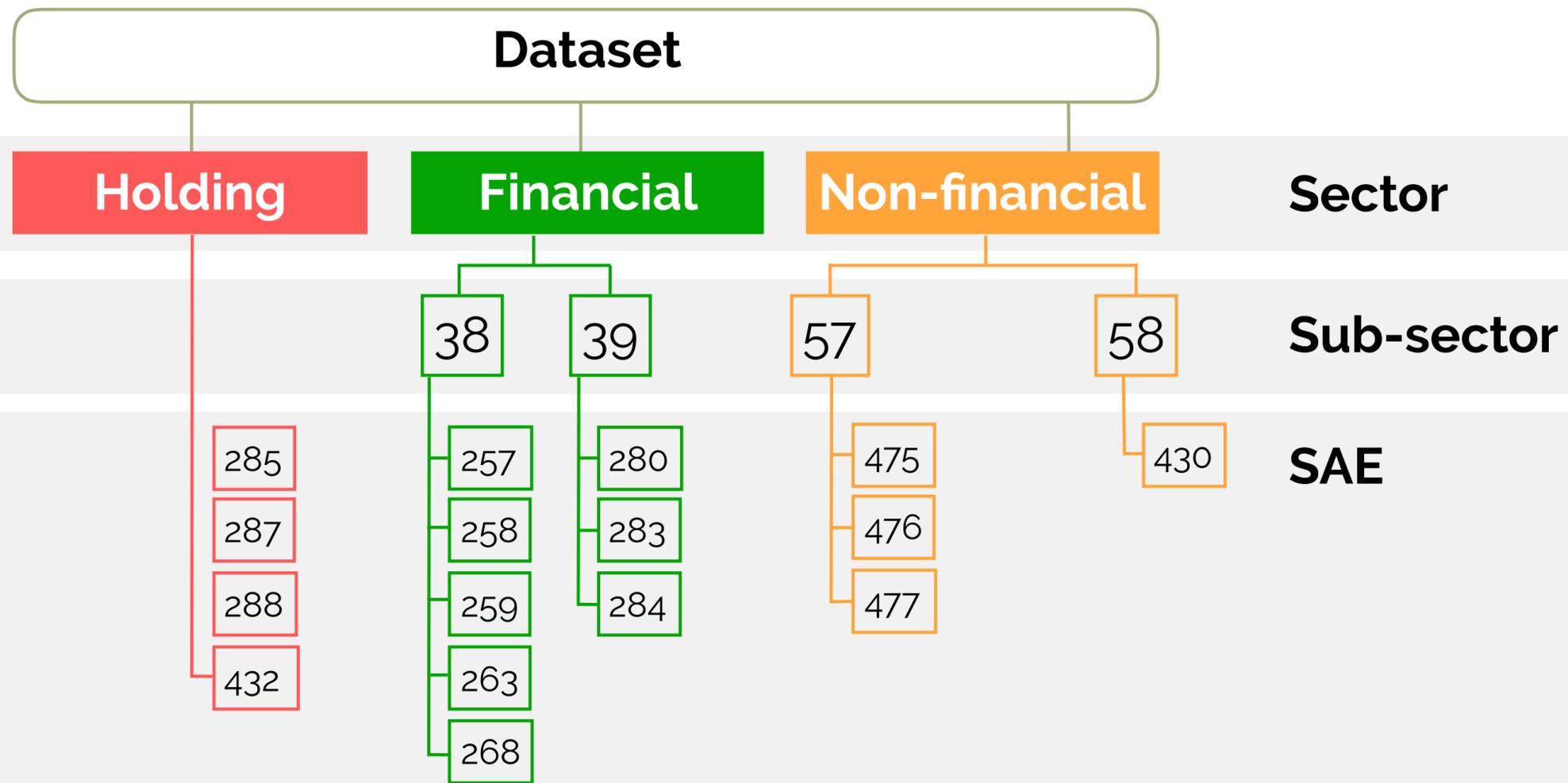
Imbalanced learning

Results



Classification

SAE hierarchy



Data

Preprocessing

Feature extraction

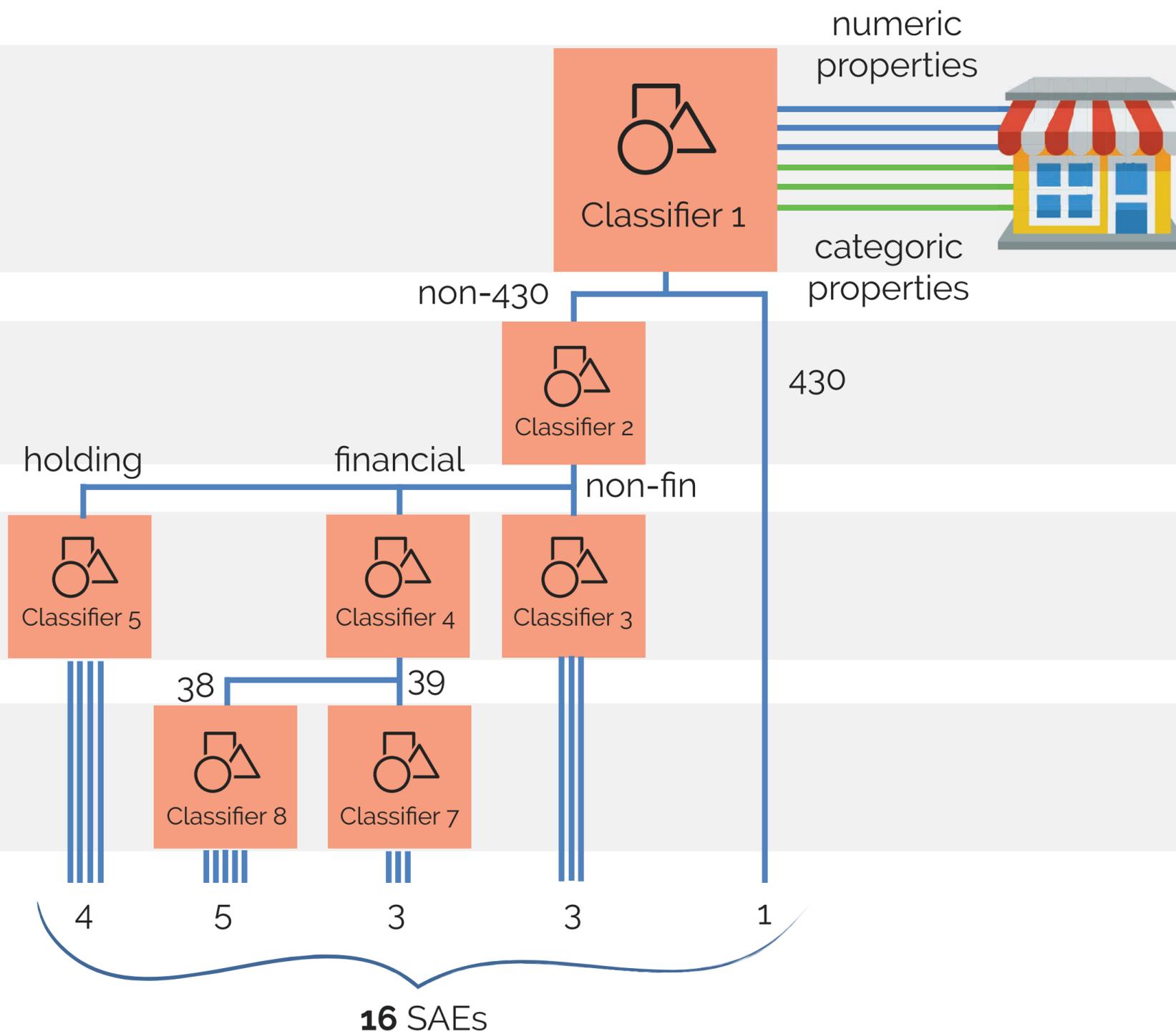
Imbalanced learning

Results

Classifier hierarchy



Classification



Data

Preprocessing

Feature extraction

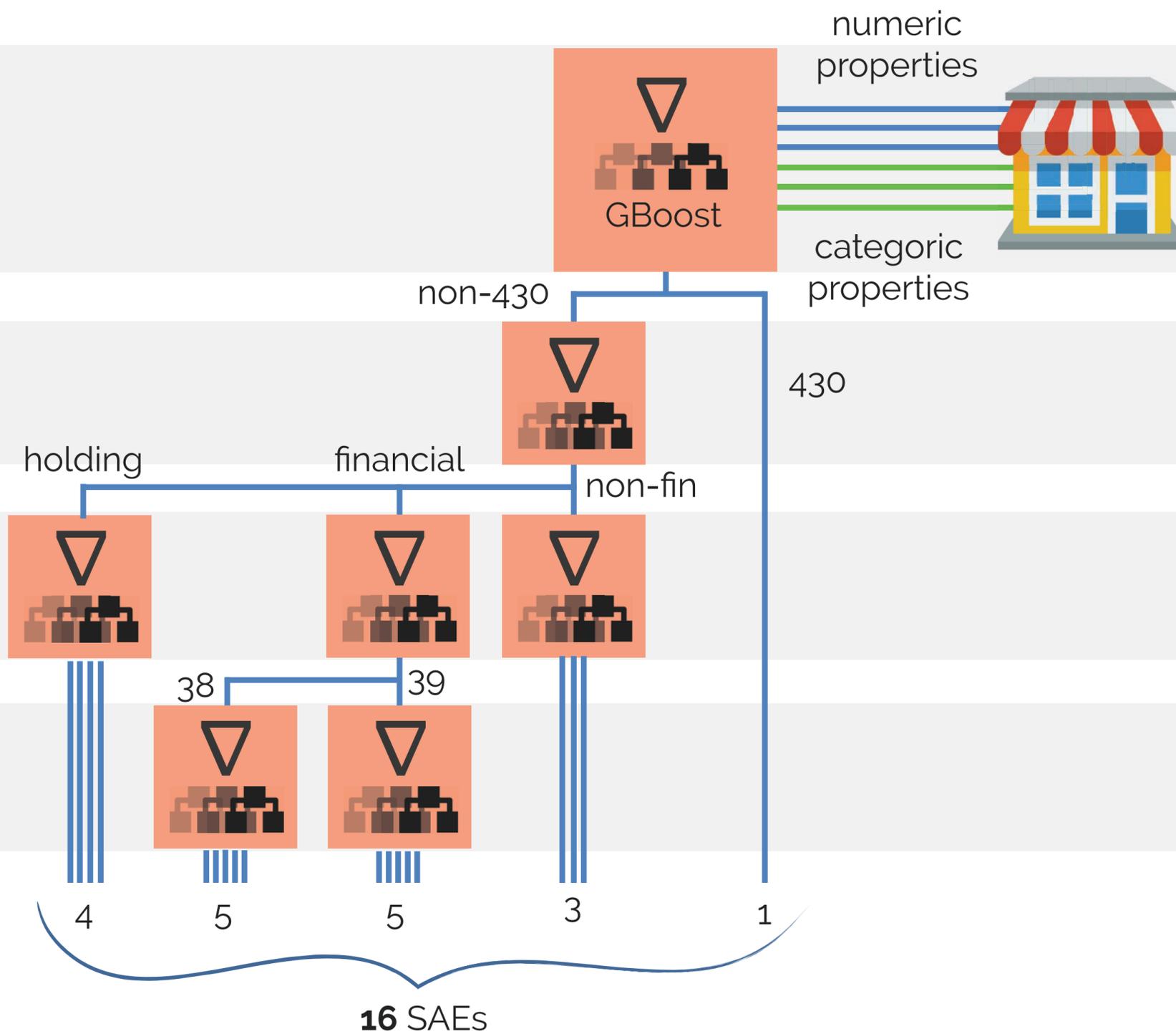
Imbalanced learning

Results

Classifier hierarchy



Classification



Data

Preprocessing

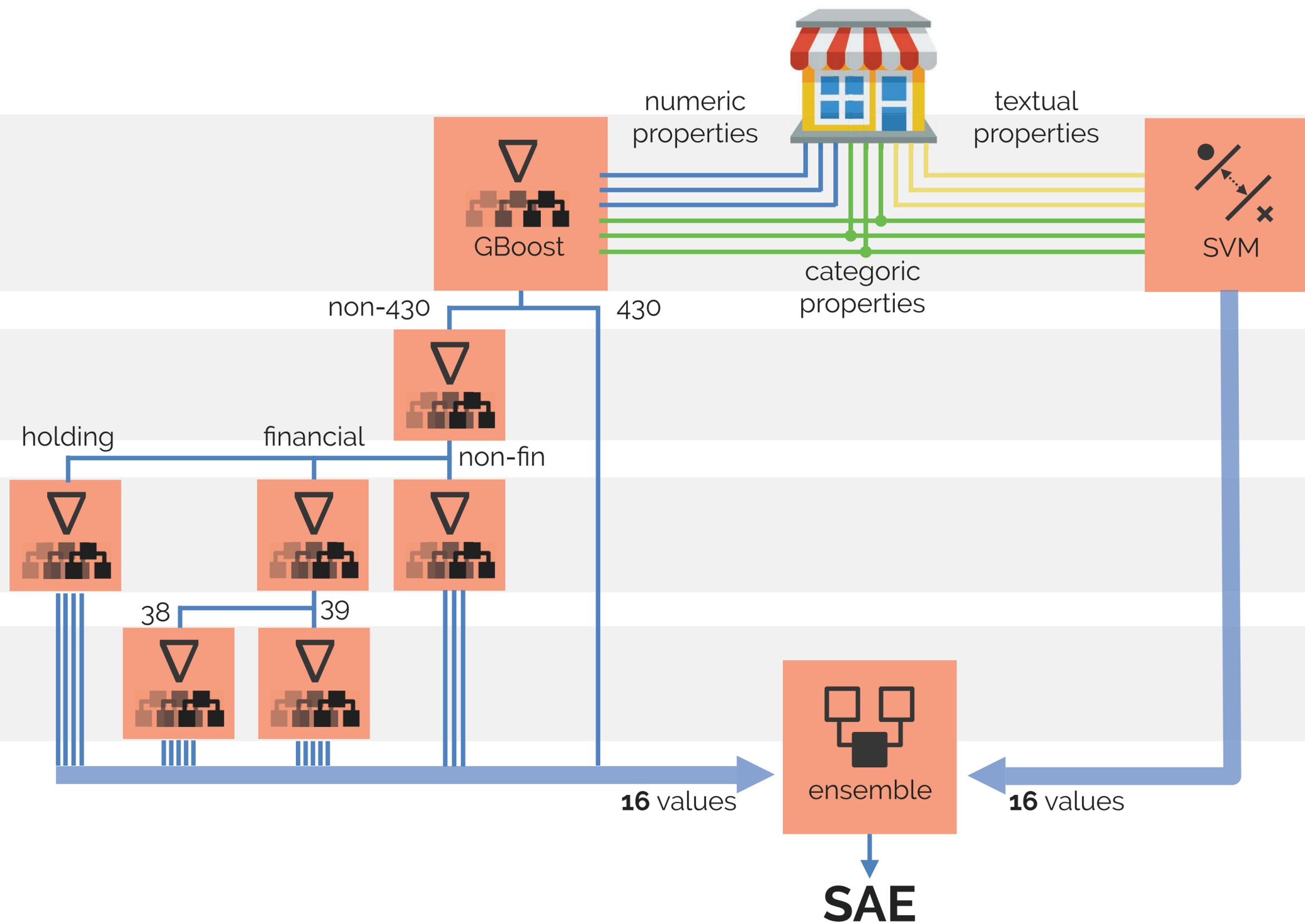
Feature extraction

Imbalanced learning

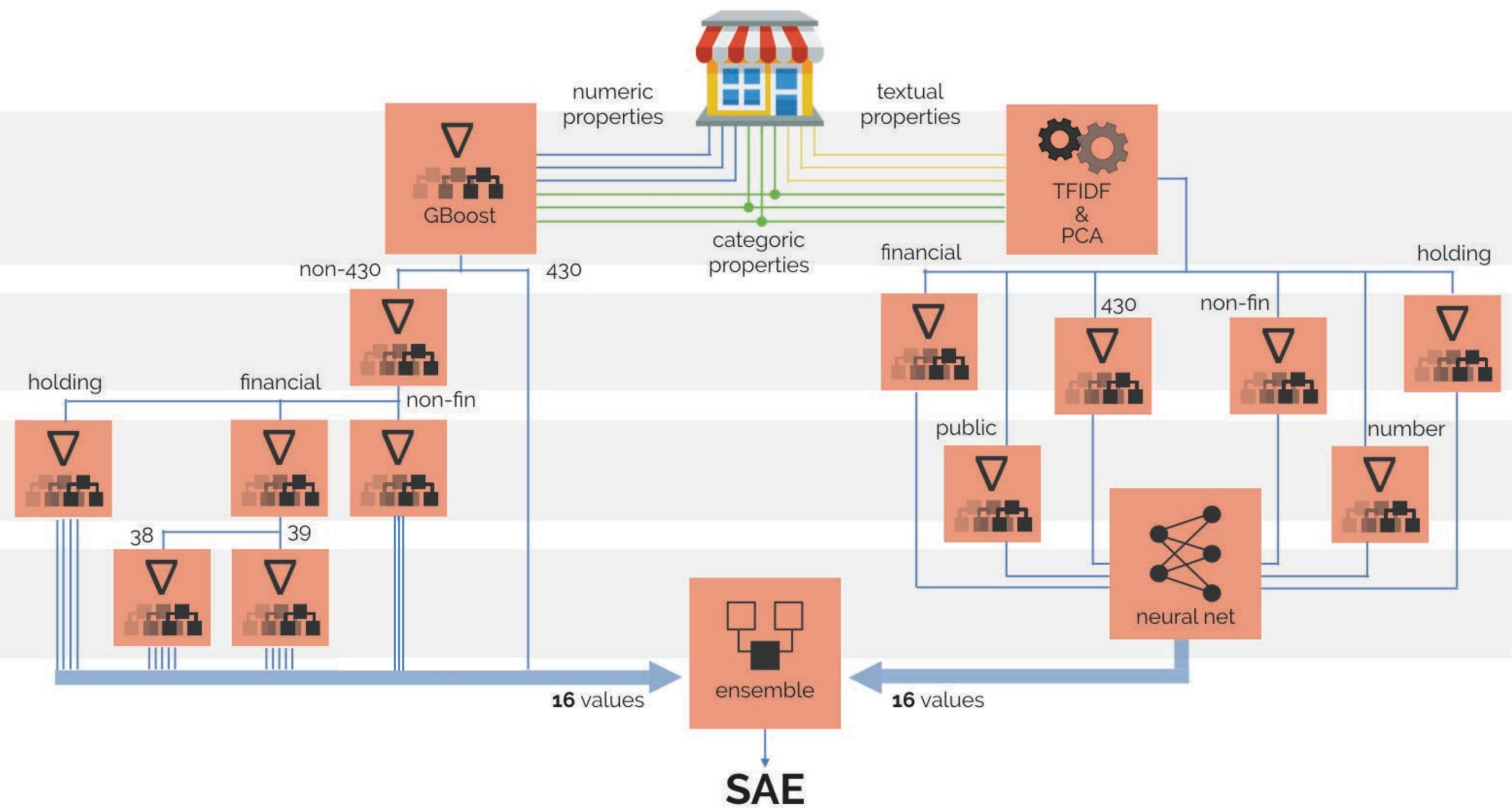
Results

Ensemble classifier

Classification



Ensemble Neural classifier



Data

Preprocessing

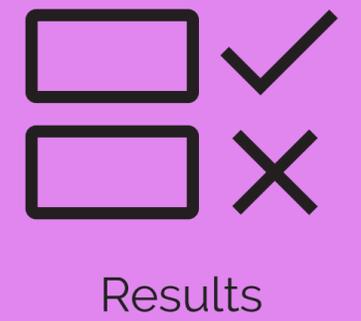
Feature extraction

Imbalanced learning

Results

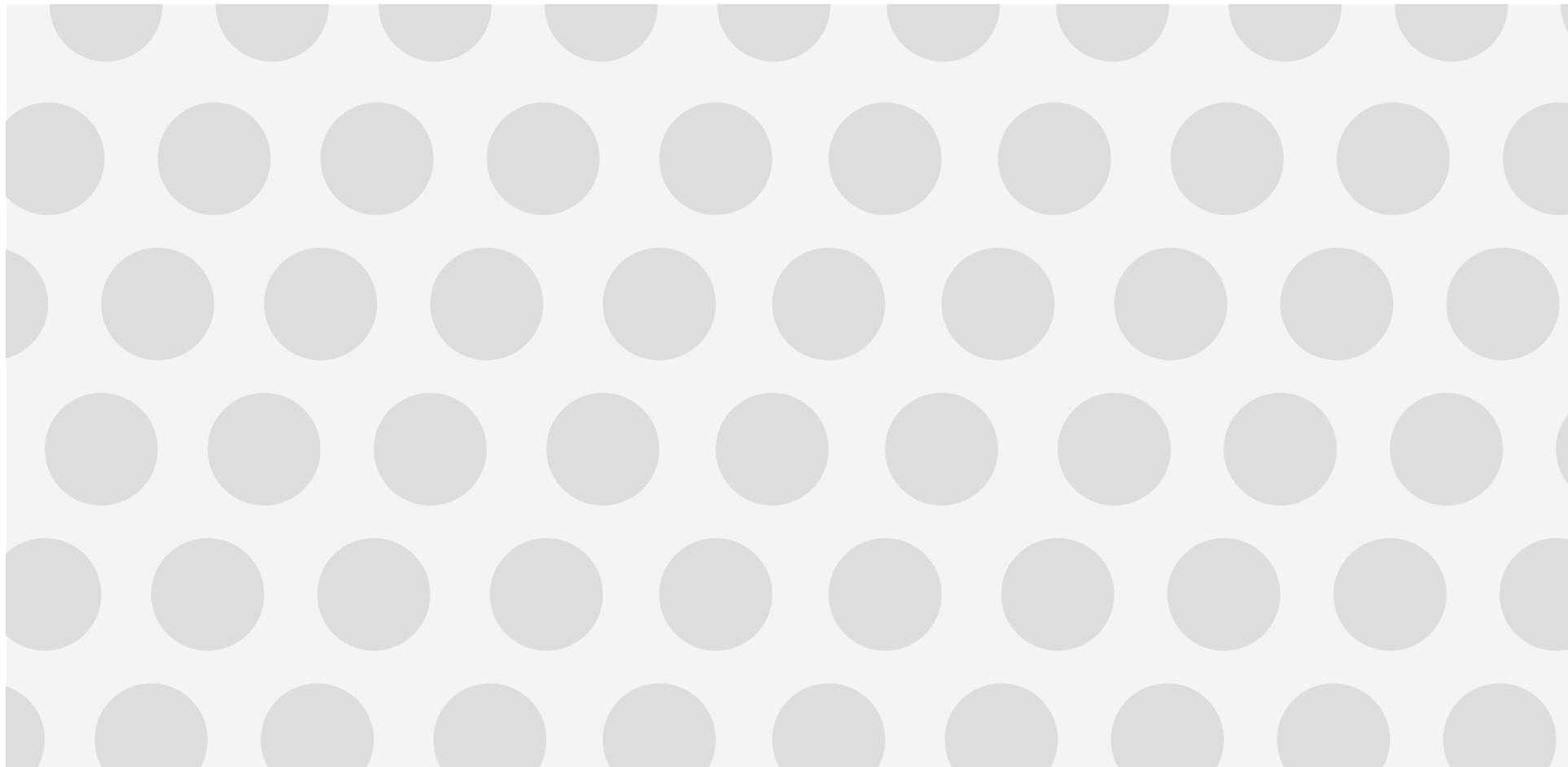


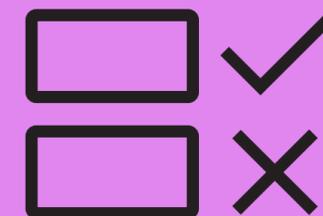
Classification



Datasets and performance

1.4 million records

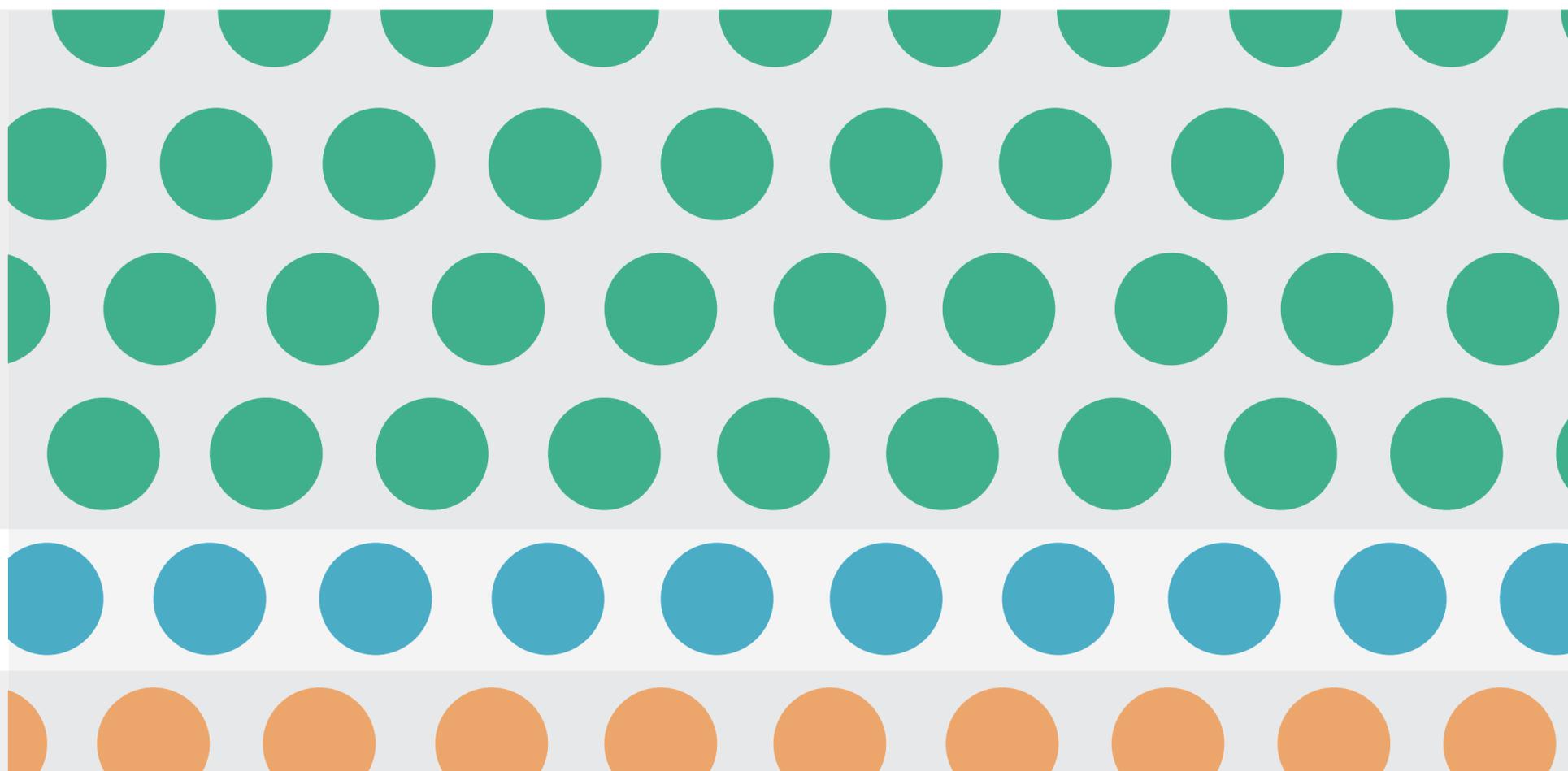




Results

Datasets and performance

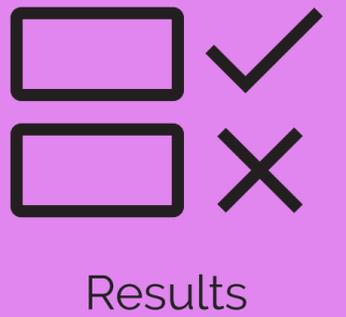
1.4 million records
of SAE-labeled data



430

288

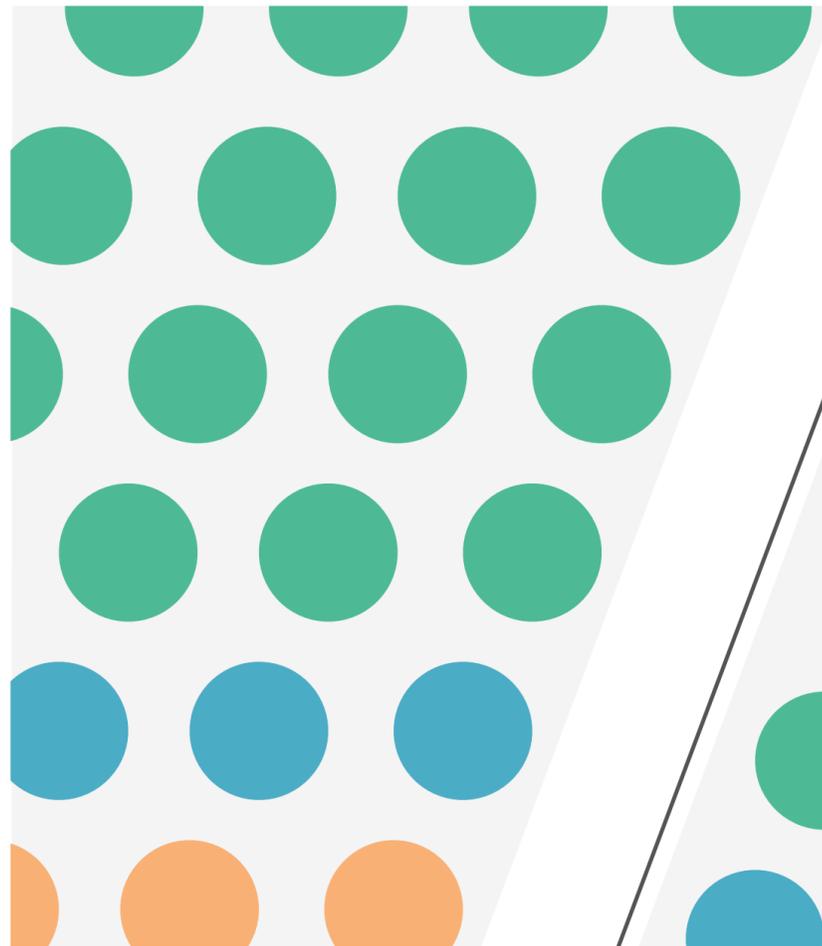
other SAEs



Datasets and performance

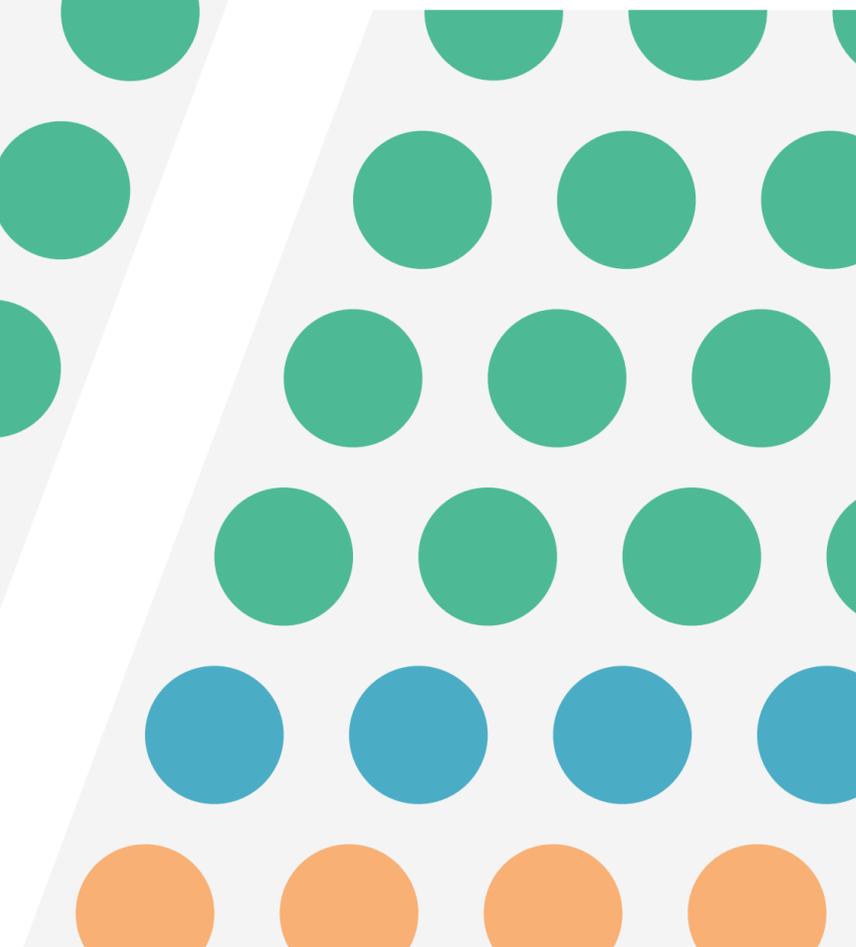
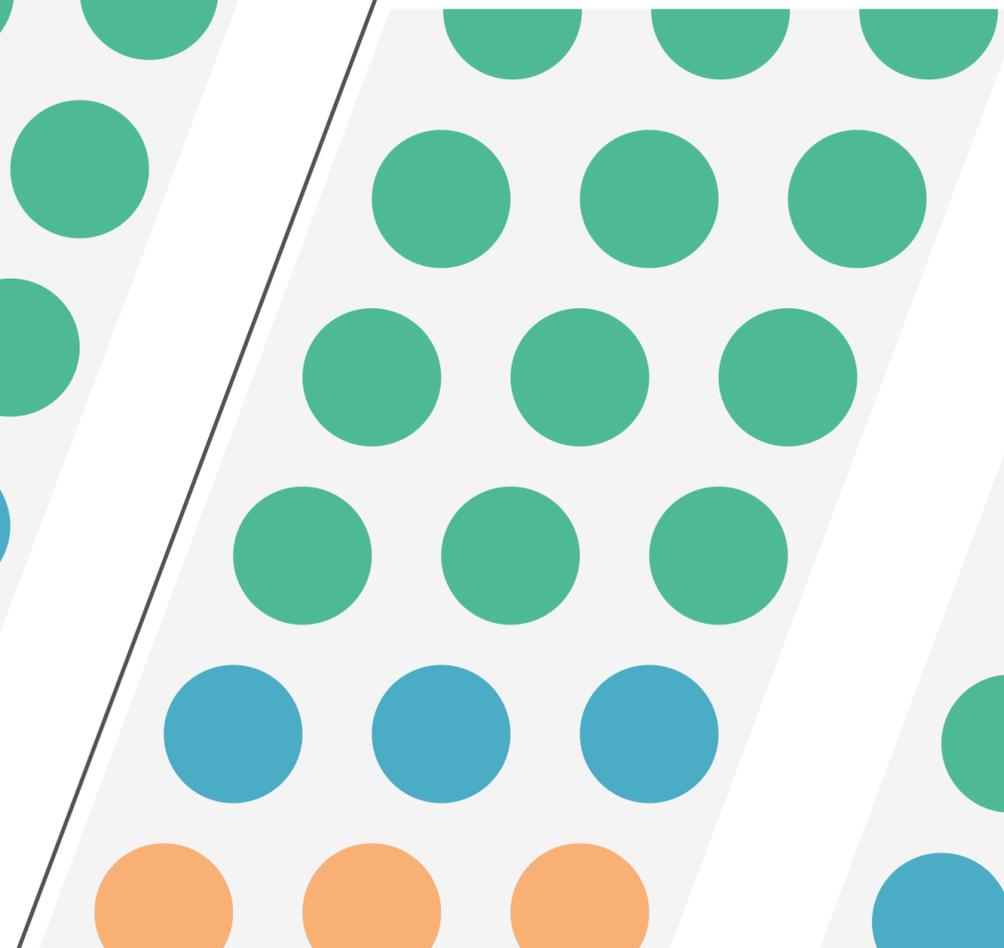
training set

used to automatically learn classifier **parameters**



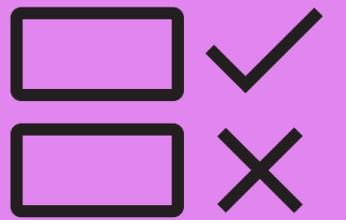
validation set

used to optimise classifier **hyperparameters**



test set

used to evaluate classifier **performance**



Results

Datasets and performance

training set

used to automatically learn classifier **parameters**

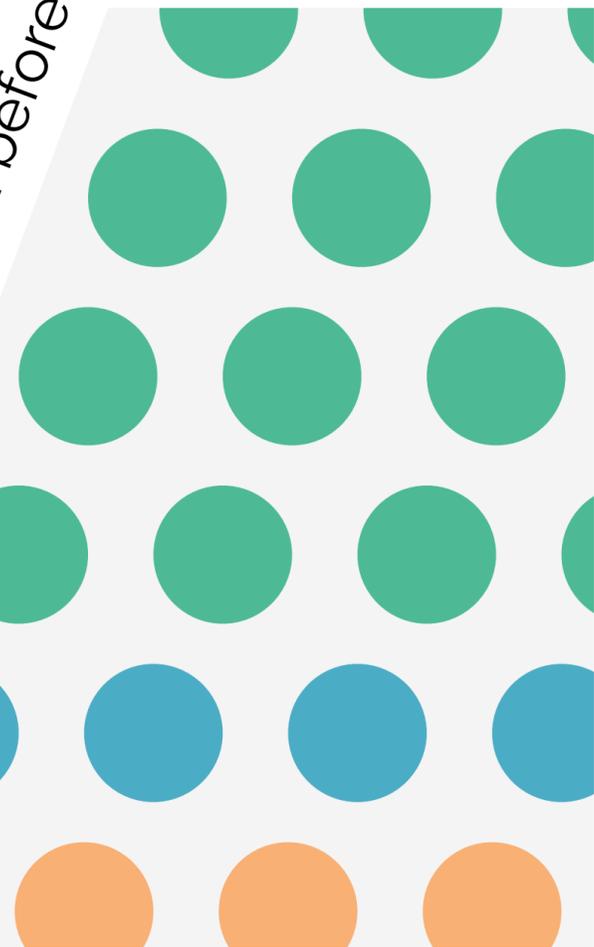


validation set

used to optimise classifier **hyperparameters**

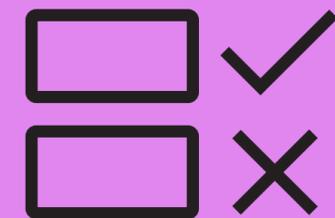


*420,000 records
the classifier has never seen before*



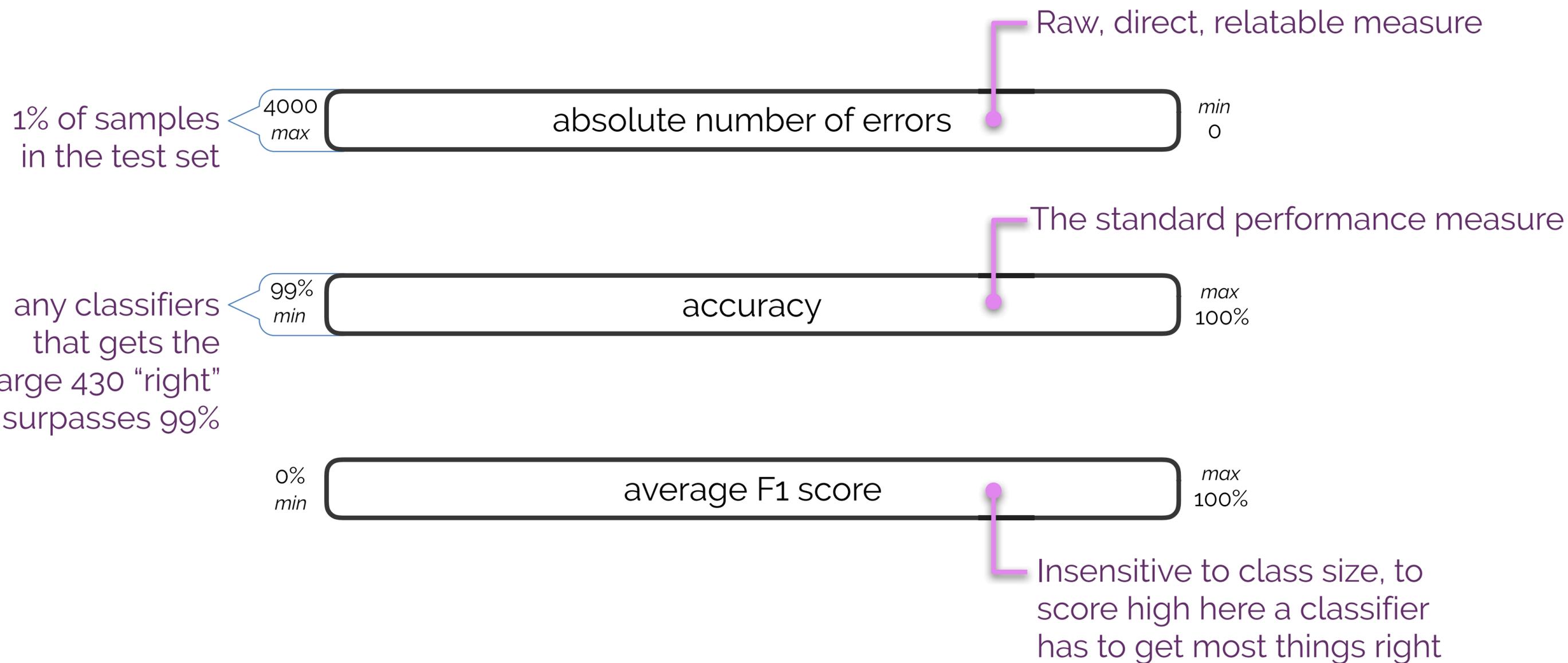
test set

used to evaluate classifier **performance**



Results

Performance metrics



Data

Preprocessing

Feature extraction

Imbalanced learning

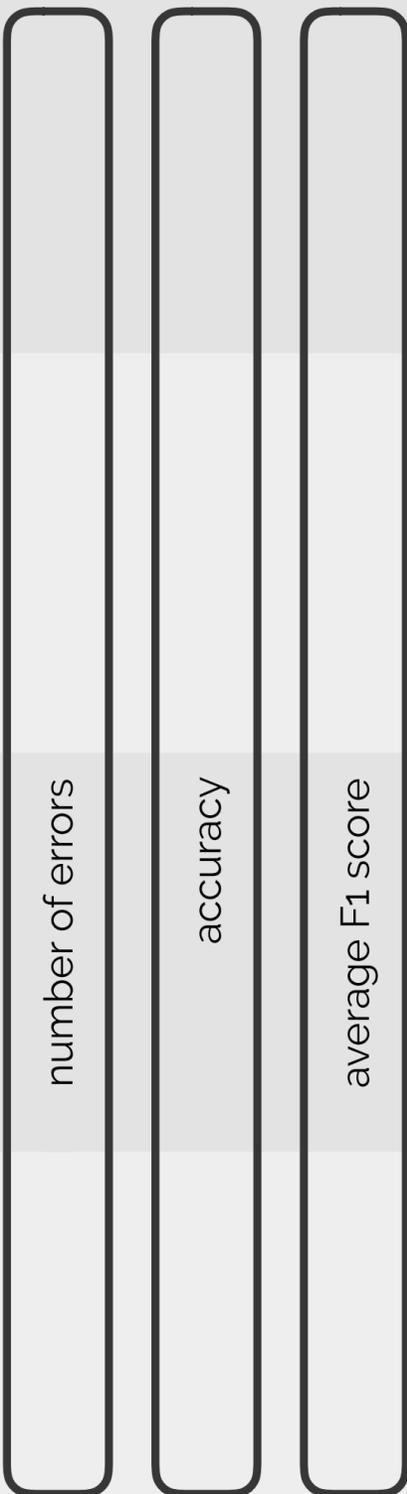
Classification



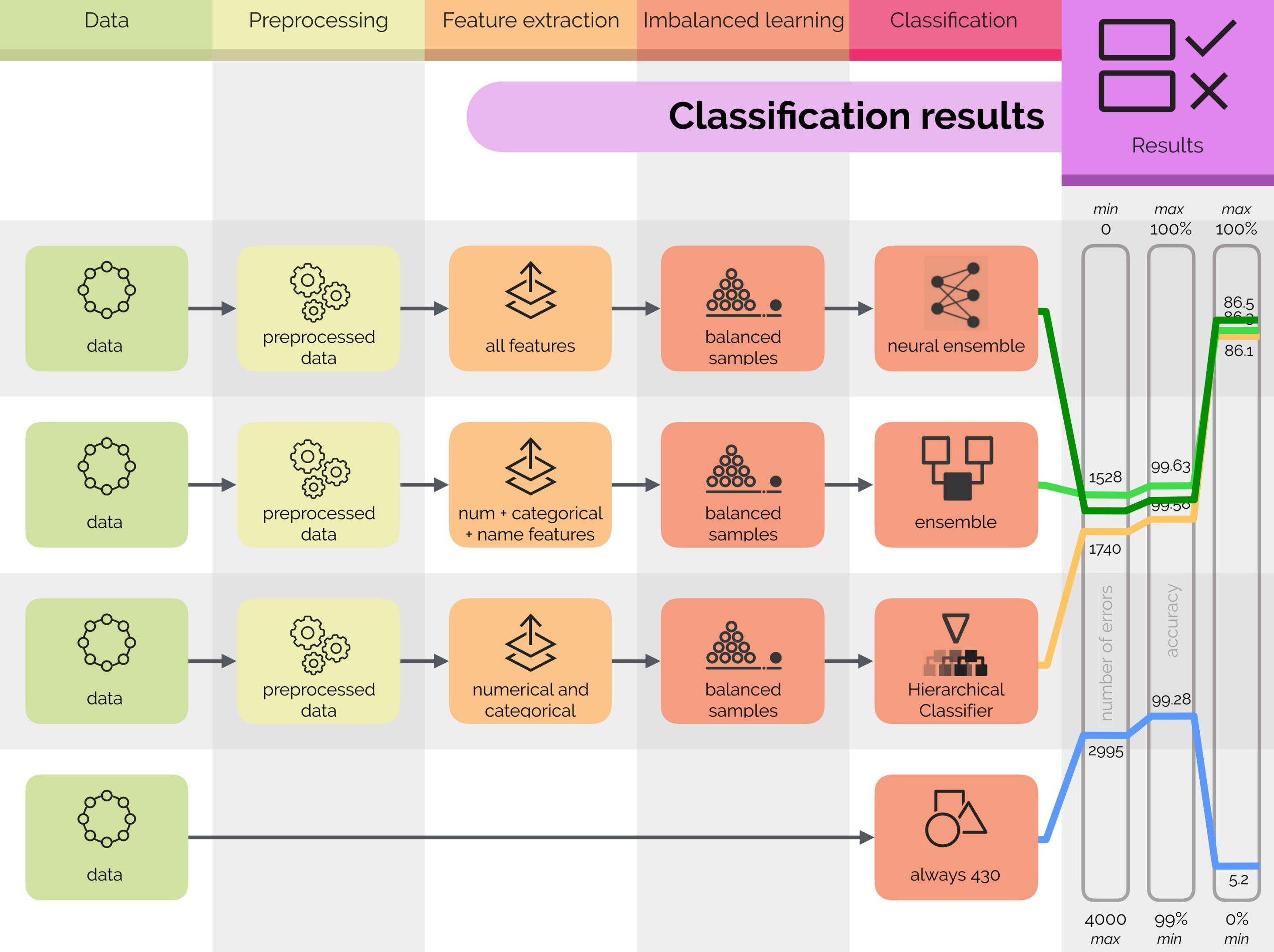
Results

Classification results

min
0 *max*
100% *max*
100%



4000 99% 0%
max *min* *min*



Data

Preprocessing

Feature extraction

Imbalanced learning

Classification

Classification results

✓
 ✗
 Results

data

preprocessed data

all features

balanced samples

neural ensemble

data

preprocessed data

num + categorical + name features

balanced samples

ensemble

data

preprocessed data

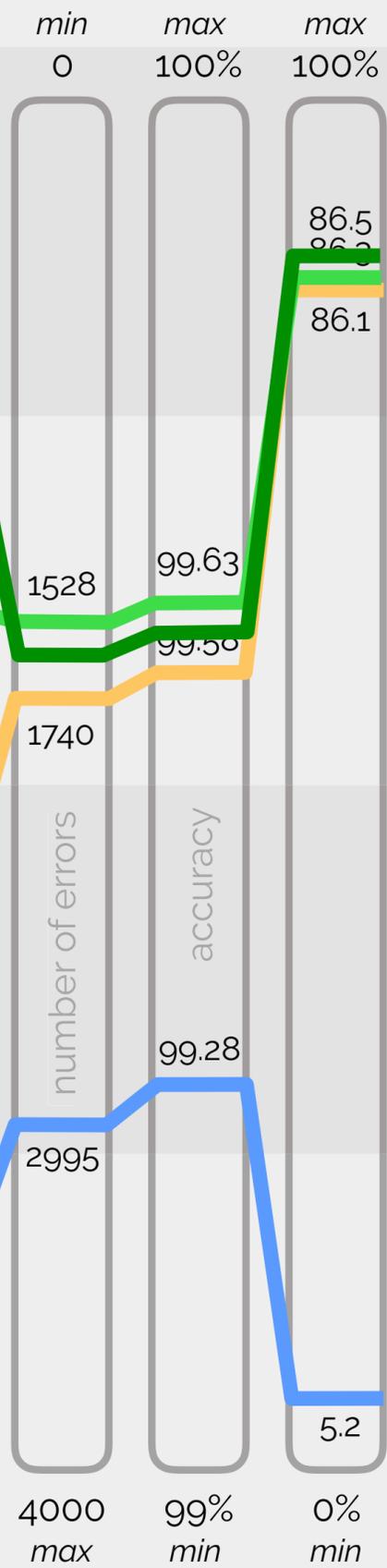
numerical and categorical

balanced samples

Hierarchical Classifier

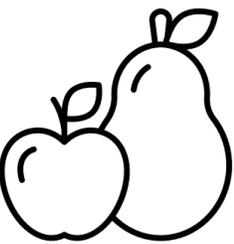
data

always 430



Conclusions

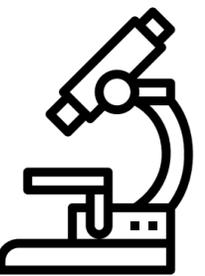
Dealing with **hybrid data** is complex and **different pipelines** (with ensemble techniques) are needed



Aa

Hierarchical structures give comfortable a-priori knowledge but are not well suited for “ambiguous” data

A **scientific paper** with details on all the techniques presented is currently under review and will be **published soon**.



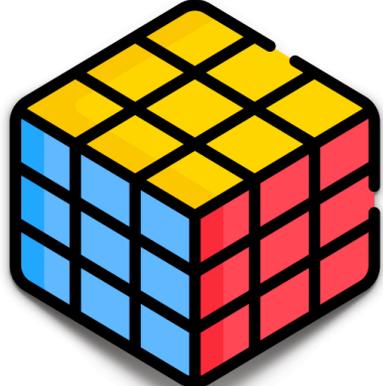
From problem

A **business necessity** to improve DQM activity efficiency

A **Machine Learning** solution could solve the problem

A **research activity** was carried out in order to find the best solution

A **final solution** is being developed as an integration in the enterprise SW

To  solution

Workshop on “Big Data & Machine Learning Applications for Central Banks”

October 22nd 2019

Centro Carlo Azeglio Ciampi



Thank you for your attention

Any questions?



BANCA D'ITALIA

EUROSISTEMA

The opinions expressed and conclusions drawn are those of the authors and do not necessarily reflect the views of the Bank of Italy.