

Big Homework OSDA

Mukhanko Artem

19.12.2023

Table of contents

Introduction	2
1) Selected Datasets	2
1.1) Income dataset.....	2
1.2) Stroke prediction Dataset.....	3
1.3) Employee Future Prediction.....	3
2) Classification with standard ML tools	4
3) Lazy-FCA with pattern structures.....	5
Conclusion.....	6

Introduction

In this project will be introduced implementation of a Lazy FCA classification algorithm based on pattern structures. Proposed algorithm was compared with popular ML models: Logistic Regression, Random Forest, KNeighbors, Decision Tree. You can find all code in my GitHub repository: <https://github.com/artemm26/Lazy-FCA>

1) Selected Datasets

I used the following datasets in the project:

Income dataset: [Income Dataset \(kaggle.com\)](#)

Stroke prediction dataset: [Stroke Prediction Dataset \(kaggle.com\)](#)

Employee future prediction: [Employee Future Prediction \(kaggle.com\)](#)

1.1) Income dataset

The dataset provided predictive feature like education, employment status, marital status to predict if the salary is greater than \$50K.

Attribute Information:

- 1) age: (17 - 90 y.o.)
- 2) workclass: type of job
- 3) fnlwgt: final weight
- 4) education: (HS-grad, some-college...)
- 5) educational-num: education as Integer
- 6) marital-status: marital status
- 7) occupation
- 8) relationship: relationship status (husband, not-in-family...)
- 9) race (white, black)
- 10) gender (male, female)
- 11) capital-gain
- 12) capital-loss
- 13) hours-per-week
- 14) native-country

1.2) Stroke prediction Dataset

According to the World Health Organization (WHO) stroke is the 2nd leading cause of death globally, responsible for approximately 11% of total deaths. This dataset is used to predict whether a patient is likely to get stroke based on the input parameters like gender, age, various diseases, and smoking status. Each row in the data provides relevant information about the patient.

Attribute Information:

- 1) id: unique identifier
- 2) gender: "Male", "Female" or "Other"
- 3) age: age of the patient
- 4) hypertension: 0 if the patient doesn't have hypertension, 1 if the patient has hypertension
- 5) heart_disease: 0 if the patient doesn't have any heart diseases, 1 if the patient has a heart disease
- 6) ever_married: "No" or "Yes"
- 7) work_type: "children", "Govt_jov", "Never_worked", "Private" or "Self-employed"
- 8) Residence_type: "Rural" or "Urban"
- 9) avg_glucose_level: average glucose level in blood
- 10) bmi: body mass index
- 11) smoking_status: "formerly smoked", "never smoked", "smokes" or "Unknown"*
- 12) stroke: 1 if the patient had a stroke or 0 if not

1.3) Employee Future Prediction

A company's HR department wants to predict whether some customers would leave the company in next 2 years.

Attribute Information:

- 1) Education: education level
- 2) JoiningYear: year of joining company
- 3) City: city office, where posted
- 4) PaymentTier: (1- highest, 2-mid level, 3- lowest)
- 5) Age: current age
- 6) Gender (male, female)
- 7) EverBenched: ever kept out of projects for 1 month or more

- 8) ExperienceInCurrentDomain: experience in current field
- 9) LeaveOrNot: whether employee leaves the company in next 2 years

2) Classification with standard ML tools

After selecting 3 datasets I proceed Machine learning models. As the preprocessing I fixed obvious errors in data, like deliting nan values and fixed datatypes in the datasets: formatted numerical to int and float datatypes; converted categorical features. For encoding I used LabelEncoder. This preprocessing step aimed to enhance the overall quality of the datasets, ensuring they were well-structured and ready for consumption by a variety of classification algorithms.

The models were used: Logistic Regression, Random Forest, KNeighbors, Decision Tree. They have been selected for their widespread usage and proven performance in classification tasks. In this project, I employed accuracy score and F1-score as key metrics for evaluating the performance of classification models across three datasets. Accuracy score represents the ratio of correctly predicted instances to the total instances, providing a general overview of a model's correctness. On the other hand, F1-score balances precision and recall, making it particularly valuable when dealing with imbalanced datasets. It considers both false positives and false negatives, offering a more comprehensive assessment of a model's effectiveness in scenarios where different misclassification types bear varying consequences.

The following tables present how the machine learning models performed on the three different datasets. It includes accuracy results, recall 0, recall 1 and F1 scores for each model.

Income prediction:

	Model	Accuracy	Recall_0	Recall_1	F_score
0	LogisticRegression	0.8167	0.9035	0.5417	0.5865
1	RandomForest	0.8333	0.8816	0.6806	0.6622
2	KNeighbors	0.8167	0.8553	0.6944	0.6452
3	DecisionTree	0.7700	0.8202	0.6111	0.5605

Stroke prediction:

	Model	Accuracy	Recall_0	Recall_1	F_score
0	LogisticRegression	0.8233	0.9383	0.3333	0.4176
1	RandomForest	0.8167	0.9383	0.2982	0.3820
2	KNeighbors	0.8100	0.9053	0.4035	0.4466
3	DecisionTree	0.7600	0.8313	0.4561	0.4194

Employee prediction:

	Model	Accuracy	Recall_0	Recall_1	F_score
0	LogisticRegression	0.7100	0.8614	0.3980	0.4727
1	RandomForest	0.8267	0.8911	0.6939	0.7234
2	KNeighbors	0.7900	0.9010	0.5612	0.6358
3	DecisionTree	0.8200	0.8861	0.6837	0.7128

3) Lazy-FCA with pattern structures

After binarization I performed 9 times Lazy Classification algorithm with pattern structured to tune the best hyperparameters: method (“standard”, “standard-support”, “ratio-support”) and alpha (0.1, 0.5, 0.9). As the result we can see values of accuracy score.

Results:

Income prediction:

```
Method: standard | Alpha: 0.1 | Accuracy: 0.71
Method: standard | Alpha: 0.5 | Accuracy: 0.41
Method: standard | Alpha: 0.9 | Accuracy: 0.0167
Method: standard-support | Alpha: 0.1 | Accuracy: 0.6767
Method: standard-support | Alpha: 0.5 | Accuracy: 0.2633
Method: standard-support | Alpha: 0.9 | Accuracy: 0.2467
Method: ratio-support | Alpha: 0.1 | Accuracy: 0.4833
Method: ratio-support | Alpha: 0.5 | Accuracy: 0.4767
Method: ratio-support | Alpha: 0.9 | Accuracy: 0.6533
Best score: 0.71
```

Best parameters:

Method = “standard”, alpha = 0.1

Stroke prediction:

```

Method: standard | Alpha: 0.1 | Accuracy: 0.5967
Method: standard | Alpha: 0.5 | Accuracy: 0.2267
Method: standard | Alpha: 0.9 | Accuracy: 0.0167
Method: standard-support | Alpha: 0.1 | Accuracy: 0.68
Method: standard-support | Alpha: 0.5 | Accuracy: 0.6267
Method: standard-support | Alpha: 0.9 | Accuracy: 0.61
Method: ratio-support | Alpha: 0.1 | Accuracy: 0.67
Method: ratio-support | Alpha: 0.5 | Accuracy: 0.6867
Method: ratio-support | Alpha: 0.9 | Accuracy: 0.7133
Best score: 0.7133

```

Best parameters:

Method = “ratio-support”, alpha = 0.9

Employee prediction:

```

Method: standard | Alpha: 0.1 | Accuracy: 0.7367
Method: standard | Alpha: 0.5 | Accuracy: 0.6533
Method: standard | Alpha: 0.9 | Accuracy: 0.0633
Method: standard-support | Alpha: 0.1 | Accuracy: 0.7367
Method: standard-support | Alpha: 0.5 | Accuracy: 0.6967
Method: standard-support | Alpha: 0.9 | Accuracy: 0.6767
Method: ratio-support | Alpha: 0.1 | Accuracy: 0.7533
Method: ratio-support | Alpha: 0.5 | Accuracy: 0.7467
Method: ratio-support | Alpha: 0.9 | Accuracy: 0.7767
Best score: 0.7767

```

Best parameters:

Method = “ratio-support”, alpha = 0.9

Conclusion

I applied the Lazy Formal Concept Analysis (FCA) algorithm to three datasets. In the overall evaluation, it falls in the middle of the scoreboard. However, it's worth noting that the pattern classifier tends to perform poorly when all attributes in the datasets are numerical. This observation suggests that the algorithm might have limitations or requires further adjustments when dealing with datasets composed entirely of numerical features.