

NATIONAL RESEARCH UNIVERSITY

Query-based classification with interval pattern structures: application to credit scoring

S.Kuznetsov, Y.Kashnitskiy, A.Masyutin School of Computer Sciences

Higher School of Economics , 2017 www.hse.ru



Risk Management in Retail Banking

Banks and credit organizations take risks in order to make profit.

- The higher the risk the higher is return.
- The balance between the risks and return is necessary for staying afloat.





Scoring is a Key Tool to Manage Credit Risk in Retail Banking





FCA in Classification Tasks

Formal Concept Analysis

provides a framework to represent the relations between objects and their attributes. There are several algorithms to use this framework such as JSM classification algorithm:

$$(\bigvee_{i=1}^{p} [\delta(g_{test}) \sqsubseteq h_{i}^{+}]) \land \neg(\bigvee_{j=1}^{n} [\delta(g_{test}) \sqsubseteq h_{j}^{-}])$$

where h is hypothesis and g_{test} is unlabeled object, $\delta(g_{test})$ is its description. Basically, the classification can be based on several *concepts* that are relevant for the unlabeled object. The concept-based learning model for standard objectattribute representation (i.e., formal contexts) is naturally extended to pattern structures.



Query-based Classification versus Training Based Methods

Training Based Methods

Learning Step: The training data is used to construct a model which relates the feature variables. <u>Test Step</u>: The trained model is used to predict the class variable for test instances. In case of pattern structures the learning step implies the full lattice construction, which can be NP-hard in general case.



Accept!

Training

data

Query based Classification

<u>No Learning Step</u>: the test object is processed online in a kNN manner. No prior model is designed. <u>Test Step</u>: the test object is classified based on its description.



Definitions

Let us give several definitions:

Let G be a set (of objects), let (D, \sqcap) be a meet-semi-lattice (of all possible object descriptions) and let $\delta: G \to D$ be a mapping. Then $(G, \underline{D}, \delta)$, where $\underline{D} = (D, \sqcap)$, is called a *pattern structure*

A pattern structure gives rise to the following derivation operator:

$$A^{\diamond} = \prod_{g \in A} \delta(g) \quad \text{for } A \in G,$$
$$d^{\diamond} = \{g \in G \mid d \sqsubseteq \delta(g)\} \quad \text{for } d \in (D, \ \Box).$$

Instead of strict notion of hypothesis we introduce alpha-weak premises:

A pattern $c \in D$ is an α - weak positive premise (classifier) iff:

$$\frac{||c^{\diamond} \cap G_{-}||}{||G_{-}||} \leq \alpha \text{ and } \exists A \subseteq G_{+} : c \sqsubseteq A^{\diamond}$$



The classification of a test object is based on a voting scheme among premises. In most general case voting scheme F is a mapping:

$$F(g_{test}, h_1^+, ..., h_p^+, h_1^-, ..., h_n^-) \to [-1, 1, \emptyset]$$

where g_{test} is the test object with unknown class, h_i^+ is a positive premise $\forall i = \overline{1, p}$ and h_j^- is a negative premise $\forall j = \overline{1, n}$, -1 is a label for negative class, and 1 is a label for positive class (i.e. defaulters).

The voting scheme is built upon weighting function ω (), aggregation operator A() and comparing operator \bigotimes :

$$F(\omega(\cdot), A(\cdot), \otimes) =$$

= $(A_{i=1}^{p}[\omega(h_{i}^{+})]) \otimes (A_{j=1}^{n}[\omega(h_{j}^{-})])$



Query Based Classification Algorithm

The algorithm is designed to be run iteratively and includes randomization.

<u>Inputs:</u> positive and negative contexts (two parts of training data with defaulters and non-defaulters).

Voting scheme:

$$F(g_{test}, h_1^+, ..., h_p^+, h_1^-, ..., h_n^-) = \\ = (\sum_{i=1}^p [\delta(g_{test}) \sqsubseteq h_i^+]) \otimes (\sum_{j=1}^n [\delta(g_{test}) \sqsubseteq h_j^-])$$

Training data

Bad obligors Good obligors

There are three parameters within the procedure:

- <u>Number of iterations</u>
- <u>Subsample size</u> percentage of the context randomly used for intersection with the test object (parameter)
- <u>Alpha-threshold</u> is the maximum allowable percentage of the opposite context for that the premise is not falsified



Query Based Classification Algorithm

<u>Step 1.</u> Extract *subsample size* of objects from the positive (negative) context.

<u>Step 2.</u> Intersect with the test object (g_{test}) and get a premise *h*.

<u>Step 3.</u> If *h* is an *alpha-threshold*-weak positive (negative) premise then add *h* to positive (negative) premises set.

Repeat Step 1 – Step 3 *number of iteration* times. These steps can be called a process of premises mining.

<u>Final step</u>. Having acquired the sets of positive and negative premises, use the voting scheme in order to produce the prediction (or margin).





- Open dataset devoted to the credit scoring. We considered the "Give Me Some Credit" contest held in 2012. (https://www.kaggle.com/c/GiveMeSomeCredit). The validation process requires calculation of performance metrics (ROC AUC and Gini coefficient) of the model based on the data sample that was retrieved from the same distribution but was not used to develop the model itself. All observations were randomly extracted from the contest dataset.
- The data represent the customers and their metrics assessed on the date of loan application. Each context consists of 1000 objects. The test dataset consists of 300 objects.
- Feature set represents various metrics such as loan amount, term, rate, payment-to-income ratio, age of the borrower, undocumented-to-documented income, credit history metrics etc.



"Give Me Some Credit" contest held in 2012

Variable Name Description Type SeriousDlqin2yrs Person experienced 90 days past due delinquency or worse Y/N Total balance on credit cards and personal lines of credit except real estate and no installment debt like car **RevolvingUtilizationOfUnsecuredLines** percentage loans divided by the sum of credit limits Age of borrower in years integer age Number of times borrower has been 30-59 days past due NumberOfTime30-59DaysPastDueNotWorse integer but no worse in the last 2 years. Monthly debt payments, alimony, living costs divided DebtRatio percentage by monthly gross income MonthlyIncome Monthly income real Number of Open loans (installment like car loan or NumberOfOpenCreditLinesAndLoans integer mortgage) and Lines of credit (e.g. credit cards) Number of times borrower has been 90 days NumberOfTimes90DaysLate integer or more past due. Number of mortgage and real estate loans including NumberRealEstateLoansOrLines integer home equity lines of credit Number of times borrower has been 60-89 days past NumberOfTime60-89DaysPastDueNotWorse integer due but no worse in the last 2 years. Number of dependents in family excluding themselves NumberOfDependents integer (spouse, children etc.)

Kaggle Data Description



Examples of one-factor distributions by target class



age by borrower class

MonthlyIncome by borrower class



NumberOfOpenCreditLinesAndLoans by borrower class



 One factor distributions show that the target variable is mixed within range of factors, although they are promising to make some discrimination between good and bad obligors.



Data and Experiments Results

• According to performed calculations the model with the highest Gini (70.8%) on the validation sample was obtained with XGBoost algorithm.

- Classical scorecard adopted in bank has shown moderate discriminatory power with 58.1% Gini.
- Query based algorithm reached accuracy of 66.3% in terms of Gini

| Metric | Query based classification | Scorecard | XGboost |
|------------|----------------------------|-----------|---------|
| Gini on | | | |
| validation | | | |
| dataset | 0.6624 | 0.5806 | 0.708 |









Conclusion

• When dealing with large numerical datasets, query based classification may be preferable to classification based on full set of hypotheses extracted from lattice, since it requires less time and memory resources

• We considered query based classification for Kaggle open data set devoted to problem of credit scoring

•The classification accuracy of the algorithm was compared both to conventionally adopted models used in the bank and to black-box models.



Thank you for your attention!



Appendix (feature importance)



Fig. 4. Feature importance estimated with Xgboost