

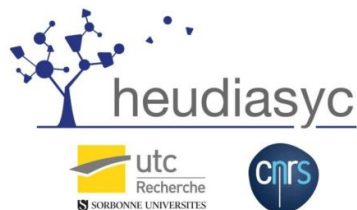
Presentation of graduation internship



Implementing Machine Learning Methods In Credit Risk Management

Inter: Zinan ZHOU
Supervisor: Dritan NACE

Introduction and state of art



- Build a model that borrowers can use to help make the best financial decisions.
- Classify clients as good clients and bad clients
- Machine Learning models

Techniques used for credit risk analysis could be categorized:

- Statistical methods:
Logistic regression, classification tree, etc
- Artificial Intelligence techniques:
Support Vector Machine, artificial neural networks, etc
- Hybrid approaches and Ensemble methods

Methods	German Dataset			Australian Dataset		
	OA (%)	Se(%)	Sp(%)	OA(%)	Se(%)	Sp(%)
LDA	72.00	72.43	71.00	85.80	92.50	80.42
QDA	68.00	67.71	68.67	80.59	66.48	91.91
LogR	76.40	88.14	49.00	86.53	88.28	85.12
DT	71.80	79.57	53.67	82.18	80.41	83.56
k-NN	69.90	89.85	23.33	69.13	54.39	80.94
DSISsvm	77.10	88.86	49.67	86.96	89.25	85.12

Table 1. Performance comparisons of different classifiers*

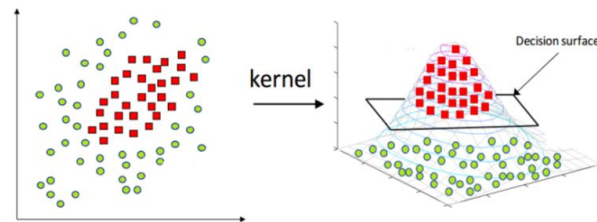
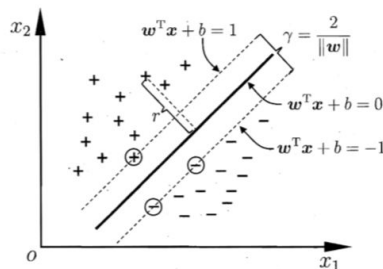
* Bio-inspired credit risk analysis: Computational intelligence with support vector machines. Lean Yu, Kin Keung Lai

Principle of SVM

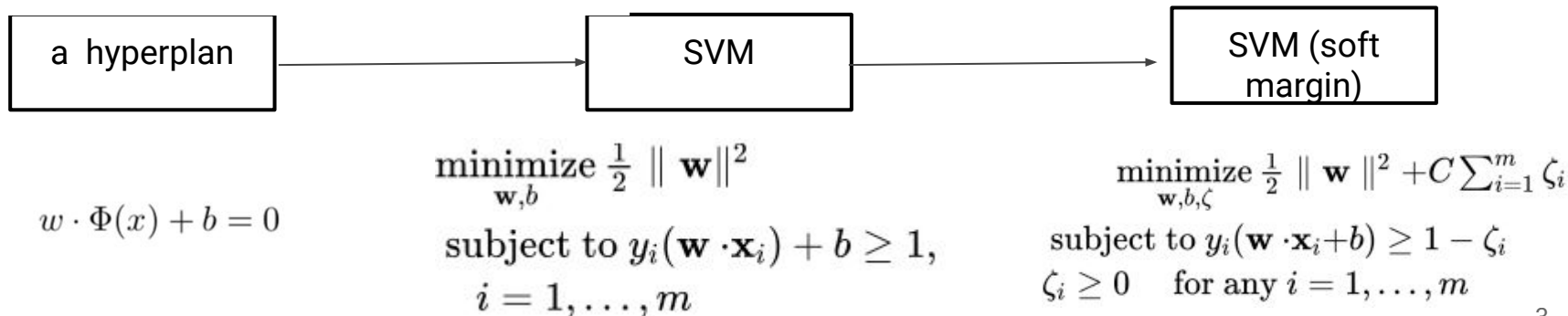
Support Vector Machine

SVM decision function:

$$f(x) = \text{sign}(w \cdot \Phi(x) + b) = \text{sign}\left(\sum_{i=1}^l \alpha_i y_i K(x_i, x) + b\right)$$



Kernel Trick converts non-linear classification to linear classification by shifting the lower dimension space to higher dimensional space.



Model Selection

- Data processing
 - Principal Components Analysis
- Model selection
 - SVM and SVM based models
- Model evaluation
 - Bad precision
 - correctly predicted as bad / classified as bad
 - Good precision
 - Total accuracy
 - AUC

Data imbalance problem :

- Over sampling and under sampling methods

Our results with advanced SVM methods :

	Bad precision	Good precision	Total accuracy	AUC
SVM	0.294	0.887	0.855	0.59
Fuzzy SVM	0.625	0.716	0.711	0.67
Bilateral Fuzzy SVM	0.549	0.781	0.769	0.665
Least Square Fuzzy SVM	0.37	0.891	0.866	0.63
Weighted LSSVM	0.412	0.852	0.828	0.632
LS Bilateral FuzzySVM	0.444	0.845	0.826	0.645



CVXOPT: Python Software for Convex Optimization
sklearn.svm.SVC

- ❖ The data are retrieved from the loan inventory of a bank out of France.

Model Selection

- Good clients : Bad clients = 20 : 1
- Data extremely imbalancing
- Random Forest models
 - Balanced Random Forest (BRF)
 - Improved Balanced Random Forest

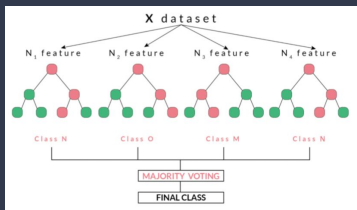
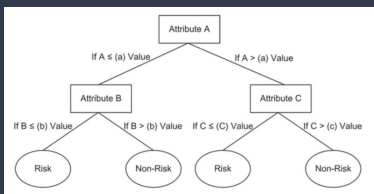


Fig 2. An example of decision tree and an example of random forest

* An experimental comparison of classification algorithms for imbalanced credit scoring data sets. Iain Brown and Christophe Mues. 2012

	Bad precision	Good precision	Total accuracy	AUC
Fuzzy SVM	0.625	0.716	0.711	0.67
Balanced Random Forest (BRF)	0.838	0.764	0.768	0.801
Improved BRF	0.908	0.701	0.708	0.799

❖ imblearn: imbalanced learn

Algorithm of Improved Balanced Random Forest model:

1. Set the threshold :
 $0 < \text{threshold_low}, \text{threshold_high} < 1$
2. Training the model and consider the predicted probability $[\text{threshold_low}, \text{threshold_high}]$ as uncertain.
3. Put the uncertain part into the next BRF model for training, get new predicted probability.
4. Repeat step 2,3.
5. Output: the average of the predicted probability Prob_final .
 $\text{Prob_final} > 0.5$: Good clients; $\text{Prob_final} < 0.5$: Bad clients

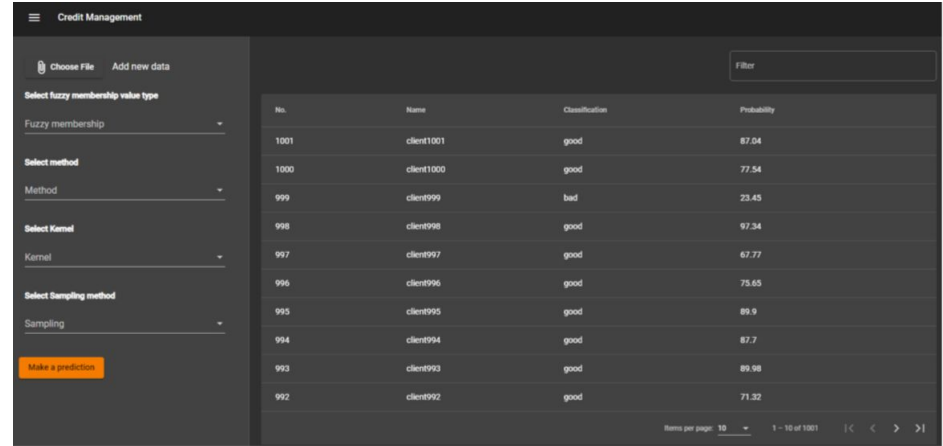
Conclution

1. Programming 5 algorithms based on SVM in Python :

- Fuzzy SVM, Bilateral Fuzzy SVM,
- Least Square Fuzzy SVM,
- Weighted Least Square SVM,
- Least Square Bilateral Fuzzy SVM.

2. Data imbalancing degrades SVM model performance. Verifying that Random Forest models perform better when the data is extremely imbalanced.

3. We improved the Balanced Random Forest model, which greatly improved the accuracy of minority class.



The screenshot shows a web application titled "Credit Management". On the left is a sidebar with configuration options: "Choose File" and "Add new data" at the top; "Select fuzzy membership value type" with a dropdown for "Fuzzy membership"; "Select method" with a dropdown for "Method"; "Select Kernel" with a dropdown for "Kernel"; and "Select sampling method" with a dropdown for "Sampling". At the bottom of the sidebar is an orange button labeled "Make a prediction". The main area on the right displays a table of client data with columns: "No.", "Name", "Classification", and "Probability". The table contains 10 rows of data. At the bottom right of the table area, there is a pagination control showing "Items per page: 10" and "1 - 10 of 1001".

No.	Name	Classification	Probability
1001	client1001	good	87.04
1000	client1000	good	77.54
999	client999	bad	23.45
998	client998	good	97.34
997	client997	good	67.77
996	client996	good	75.65
995	client995	good	89.9
994	client994	good	87.7
993	client993	good	89.98
992	client992	good	71.32

Fig. Frontend of the application

4. Save the model and cooperate with colleagues to complete the user interface.
5. Solutions to classification problems can also be applied to medical and other fields.

Thank you for your
attention !