

Dataset balancing and noise reduction

Rafael Baeta

Universidade Federal de Minas Gerais
Departamento de Ciência da Computação

2 de março de 2017

DCC
DEPARTAMENTO DE
CIÊNCIA DA COMPUTAÇÃO

UF *m* G



Introduction



- Data in the real world is dirty!
 - incomplete: lacking attribute values, lacking certain attributes of interest, or containing only aggregate data
 - noisy: containing errors or outliers
 - inconsistent: containing discrepancies in codes or names

Introduction



- Many problems in data occurs due to:
 - Data is not always available
 - Equipment malfunction
 - Data not entered due to misunderstanding
 - Certain data may not be considered important at the time of entry

Quality decisions must be based on quality data!

Introduction



- And these problems can lead to:
 - Unbalanced datasets
 - Some classes have a much smaller number of examples than other classes
 - Noise data
 - Some data can be similar and belong to different classes.
 - Data attributes can be corrupted or missing.

Introduction

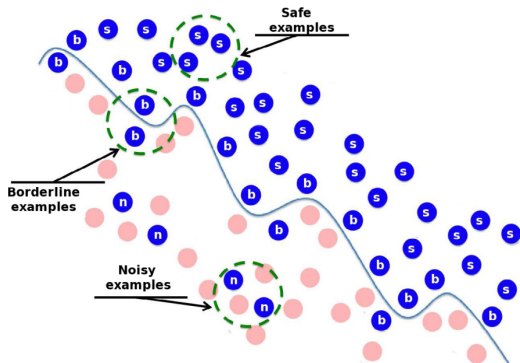


Figura : Unbalanced and noisy data

Motivation



- Datasets with high unbalanced data can lead to privilege some classes
- Some data can be wrong and can penalize the process of learning
- Some redundant data can be removed raising the training speed.
- The classifier becomes more general.

Dataset balancing



What to do if get more data for minority classes is not possible?

Dataset balancing



What to do if get more data for minority classes is not possible?

- Undersampling
 - Reduction of data that belongs to majority classes. How to select data to be removed?

Dataset balancing



What to do if get more data for minority classes is not possible?

- Undersampling
 - Reduction of data that belongs to majority classes. How to select data to be removed?
- Oversampling
 - Generate more data for the minority classes. How to generate this data?

Dataset balancing



What to do if get more data for minority classes is not possible?

- Undersampling
 - Reduction of data that belongs to majority classes. How to select data to be removed?
- Oversampling
 - Generate more data for the minority classes. How to generate this data?
- Weighting
 - Assigning larger misclassifying cost for minority class samples

Dataset balancing



What to do if get more data for minority classes is not possible?

- Undersampling
 - Reduction of data that belongs to majority classes. How to select data to be removed?
- Oversampling
 - Generate more data for the minority classes. How to generate this data?
- Weighting
 - Assigning larger misclassifying cost for minority class samples
- Ensemble methods
 - A combination of a variety of base classifiers

Undersampling



- An undersampling technique removes samples from the majority class
- It can remove relevant samples.
- A common technique is random undersampling.
- Other techniques: Neighborhood Cleaning Rule (NCR), Condensed Nearest Neighbour (CNN), One-Sided-Selection (OSS),¹ etc.

¹Miroslav Kubat, Stan Matwin et al. "Addressing the curse of imbalanced training sets: one-sided selection". Em: *ICML*. vol. 97. Nashville, USA. 1997, pp. 179–186.

Undersampling



- This approach² define k clusters.
- Calculate the number of majority examples extracted from each cluster based on a proportion between the majority and minority classes.
- Randomly select majority class samples within each cluster.
- The majority samples can be selected using heuristics of distances.

²Show-Jane Yen e Yue-Shi Lee. "Cluster-based under-sampling approaches for imbalanced data distributions". *Em: Expert Systems with Applications* 36.3 (2009), pp. 5718–5727.

Undersampling

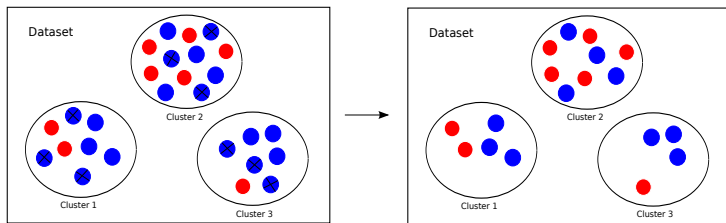


Figura : Cluster undersampling

Oversampling



- SMOTE³ is a popular oversampling technique for binary problems
- Minority class is over-sampled by creating "synthetic examples".
- Extra training data is created by operating in the "feature space".

³Nitesh V. Chawla et al. "SMOTE: synthetic minority over-sampling technique". *Em: Journal of artificial intelligence research* 16 (2002), pp. 321–357.

Oversampling



- Synthetic samples are generating in the following way:
 - Take the difference between the feature vector (sample) under consideration and its neighbor.
 - Multiply this difference by a random number between 0 and 1 and add it to the feature vector under consideration
 - This causes the selection of random points along the line segment between specific features.

Oversampling

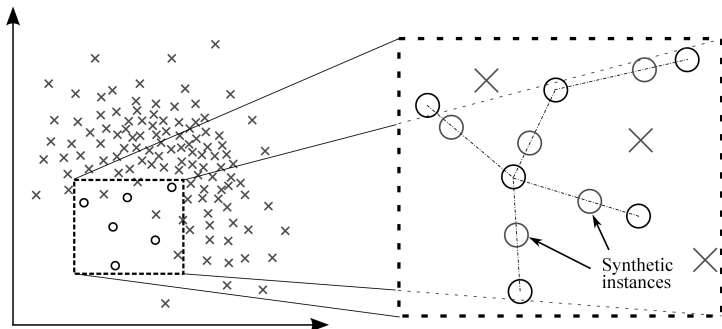


Figura : Smote example

Oversampling



- In⁴ is created a approach to deal with multiclass inbalance
- This approach uses the SMOTE algorithm to generate syntethic samples
- The synthetic samples are created based in some configurations that can be "outlier", "rare", "safe", "bordeline"

⁴José A Sáez, Bartosz Krawczyk e Michał Woźniak. “Analyzing the oversampling of different classes and types of examples in multi-class imbalanced datasets”. *Em: Pattern Recognition* 57 (2016), pp. 164–178.

Oversampling

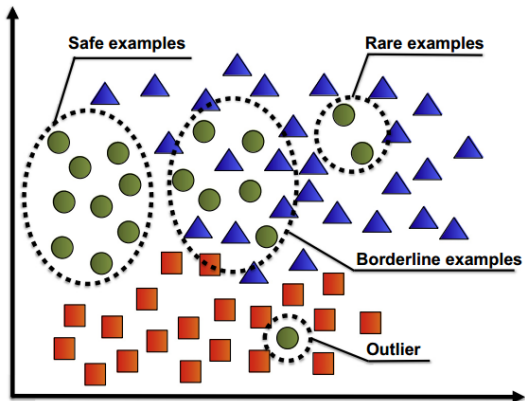


Figura : Smote example

Oversampling



- The samples can be classified as "rare", "outlier", "safe", "borderline" depending on the number of examples ($cn(e)$) that share the class label.
 - Safe examples: $cn(e) \geq 4$
 - Borderline examples: $2 \leq cn(e) \leq 3$
 - Rare examples: $cn(e) = 1$
 - Outliers examples: $cn(e) = 0$

Oversampling



- The synthetic examples are created as follow:
 - A random example x of a class and type are iteratively selected
 - Then, a random example y within the $k=5$ nearest neighbors is selected
 - With these two examples, x and y , the new synthetic example is created by interpolation as SMOTE does.

Weighting

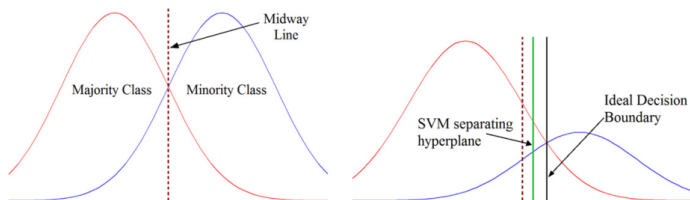


- Focus on modifying existing classification algorithms to strengthen the learning of minority class.⁵
- Most of algorithms in this family are based on SVM and Neural Network⁶

⁵Shounak Datta e Swagatam Das. “Near-Bayesian Support Vector Machines for imbalanced data classification with equal or unequal misclassification costs”. *Em: Neural Networks 70* (2015), pp. 39–52.

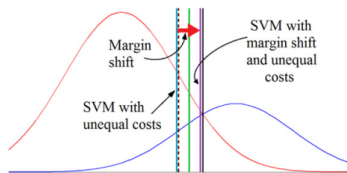
⁶Li Yijing et al. “Adapted ensemble classification algorithm based on multiple classifier system and feature selection for classifying multi-class imbalanced data”. *Em: Knowledge-Based Systems 94* (2016), pp. 88–104.

Weighting



(a) Class-conditional probability distributions for the majority class and minority class.

(b) Posterior probability distributions for the two classes. The separating hyperplane obtained by SVM is closer to the minority class, but does not coincide with the ideal decision boundary.



(c) Adding unequal costs for the two classes moves the hyperplane close to the midway line, this hyperplane can be shifted towards the ideal decision boundary using proper biasing.

Figura : SVM margin shift with unequal costs

Weighting



- Traditional SVM - $\min_{w,b} \frac{1}{2} \|w\|^2$
- Soft SVM - $\min_{w,\xi,b} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^N \xi$
- Soft SVM with unequal regularization - $\min_{w,\xi,b} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^N f(X_i) \xi$
 - $f(X_i) = \begin{cases} \frac{1}{\rho_+}, & \text{if } x_i \in I_+ \\ \frac{1}{\rho_-}, & \text{if } x_i \in I_- \end{cases}$

Ensemble



- This methodology⁷ is composed of the following steps:
 - The imbalanced training set is submitted to a feature selection

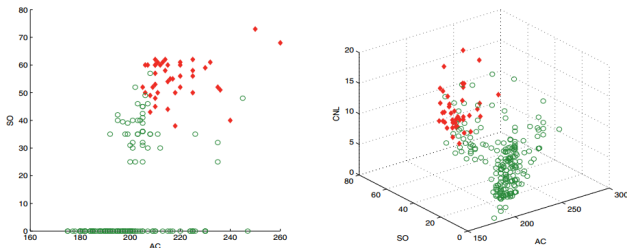


Figura : Feature selection

⁷Yijing et al., “Adapted ensemble classification algorithm based on multiple classifier system and feature selection for classifying multi-class imbalanced data”.

Ensemble



- Then the training set is divided in k subtraining sets
- Each training set is give to train a classifier.
- The answer of the classifiers are submited to a ensemble rule.

Ensemble



- This approach test 5 ensemble rules: max, min, product, majority vote and sum.
- Max - $\operatorname{argmax}_{1 \leq j \leq n} \{R_{C_1}, R_{C_2}, \dots, R_{C_m}\}, R_{C_j} = \max_{1 \leq i \leq k} AUC_{area_i} \cdot p_{ij}$
- Min - $\operatorname{argmax}_{1 \leq j \leq n} \{R_{C_1}, R_{C_2}, \dots, R_{C_m}\}, R_{C_j} = \min_{1 \leq i \leq k} AUC_{area_i} \cdot p_{ij}$
- Product - $\operatorname{argmax}_{1 \leq j \leq n} \{R_{C_1}, R_{C_2}, \dots, R_{C_m}\}, R_{C_j} = \prod_{i=1}^k AUC_{area_i} \cdot p_{ij}$
- Majority vote -
 $\operatorname{argmax}_{1 \leq j \leq n} \{R_{C_1}, R_{C_2}, \dots, R_{C_m}\}, R_{C_j} = \operatorname{count}_{C_1, C_2, \dots, C_m}(\operatorname{argmax}_{1 \leq i \leq k} AUC_{area_i} \cdot p_{ij})$
- Sum - $\operatorname{argmax}_{1 \leq j \leq n} \{R_{C_1}, R_{C_2}, \dots, R_{C_m}\}, R_{C_j} = \sum_{1 \leq i \leq k} AUC_{area_i} \cdot p_{ij}$

Ensemble

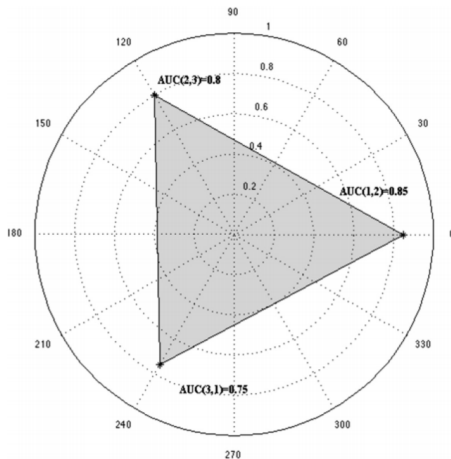


Figura : AUCarea metric

Ensemble

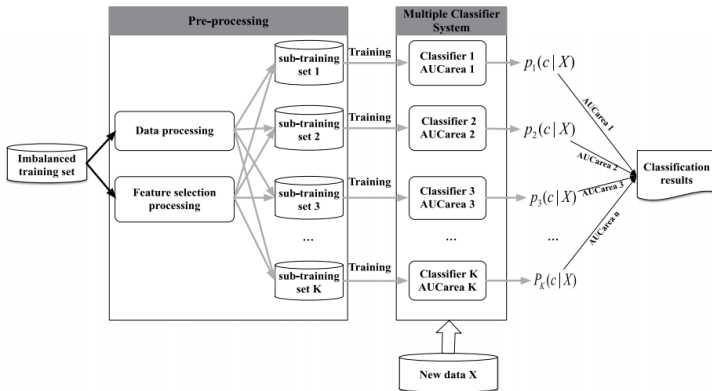


Figura : Architecture of ensemble

Ensemble

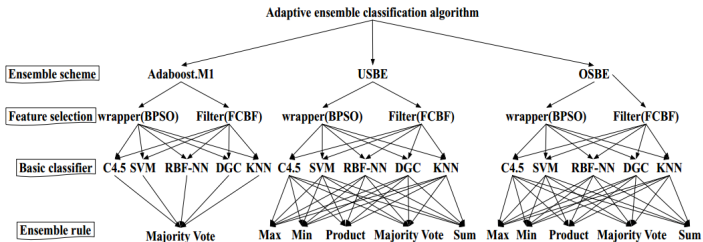


Figura : Adapting path of AMCS

Conclusion



- Dataset balancing is essential to avoid overfitting in classifiers.
- Remove useless samples
- Make classifiers more general

References I



- Batista, Gustavo EAPA, Ronaldo C Prati e Maria Carolina Monard. “A study of the behavior of several methods for balancing machine learning training data”. Em: *ACM Sigkdd Explorations Newsletter* 6.1 (2004), pp. 20–29.
- Cateni, Silvia, Valentina Colla e Marco Vannucci. “A method for resampling imbalanced datasets in binary classification tasks for real-world problems”. Em: *Neurocomputing* 135 (2014), pp. 32–41.
- Chawla, Nitesh V. et al. “SMOTE: synthetic minority over-sampling technique”. Em: *Journal of artificial intelligence research* 16 (2002), pp. 321–357.
- Datta, Shounak e Swagatam Das. “Near-Bayesian Support Vector Machines for imbalanced data classification with equal or unequal misclassification costs”. Em: *Neural Networks* 70 (2015), pp. 39–52.
- Kubat, Miroslav, Stan Matwin et al. “Addressing the curse of imbalanced training sets: one-sided selection”. Em: *ICML*. Vol. 97. Nashville, USA. 1997, pp. 179–186.

References II



- Rahman, M Mostafizur e DN Davis. “Addressing the class imbalance problem in medical datasets”. *Em: International Journal of Machine Learning and Computing* 3.2 (2013), p. 224.
- Sáez, José A, Bartosz Krawczyk e Michał Woźniak. “Analyzing the oversampling of different classes and types of examples in multi-class imbalanced datasets”. *Em: Pattern Recognition* 57 (2016), pp. 164–178.
- Yen, Show-Jane e Yue-Shi Lee. “Cluster-based under-sampling approaches for imbalanced data distributions”. *Em: Expert Systems with Applications* 36.3 (2009), pp. 5718–5727.
- Yijing, Li et al. “Adapted ensemble classification algorithm based on multiple classifier system and feature selection for classifying multi-class imbalanced data”. *Em: Knowledge-Based Systems* 94 (2016), pp. 88–104.