

Multi-class Imbalanced Learning

DingHao

December 5, 2016

Contents

Introduction

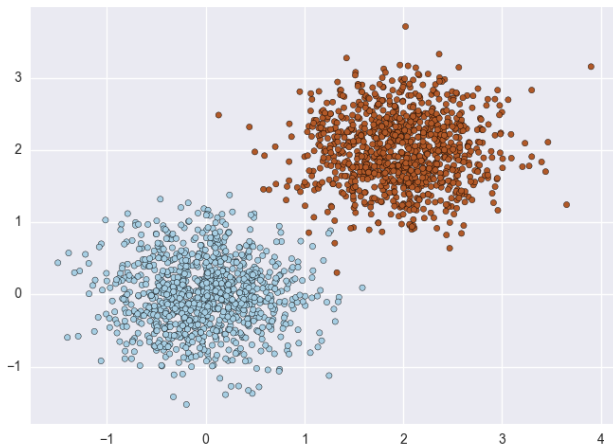
Evaluation Criteria

Approaches

Future Work

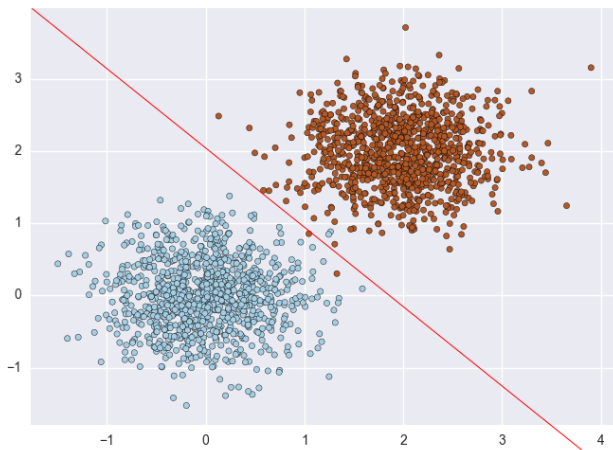
Introduction

Classification Problem



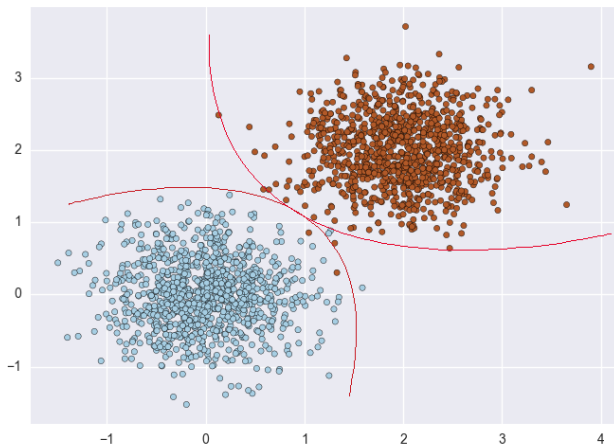
Introduction

classifiers



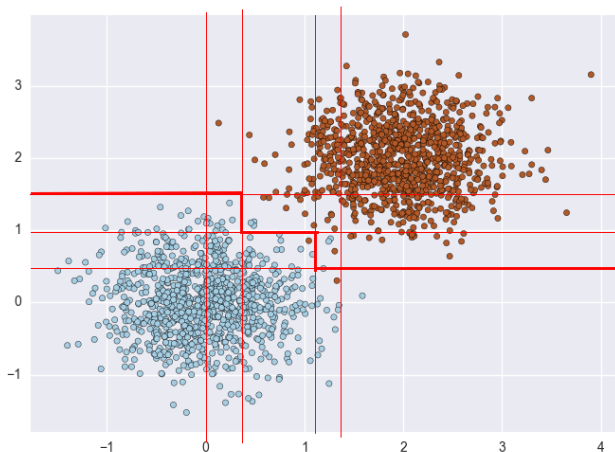
Introduction

classifiers



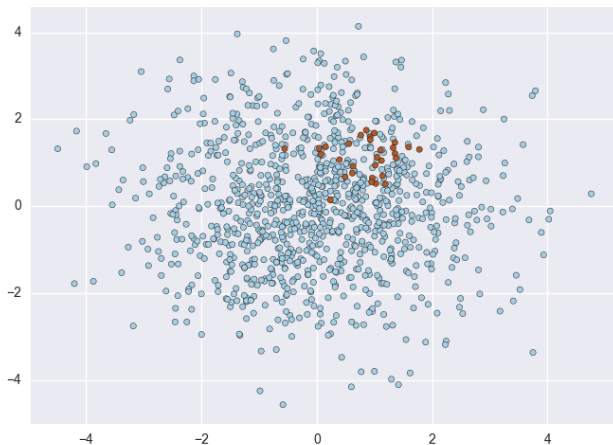
Introduction

classifiers



Introduction

Real World



Introduction

Doggy



Introduction

Rare Word

叅爰碇叁収支叵妄叔呶噀恹愉囹囿囑囁囁團
壘壘壘壘壘壘壘壘壘壘壘壘壘壘壘壘壘壘壘壘
香宜弃寓封对渺允佻膾牖龙馗嶸嶸嶸嶸嶸嶸
笑配渫叁悵欄櫛腩牖牖腩并廌廌廌廌廌廌廌
迴巡井弃弍彊弼彊弼彊彊彊彊彊彊彊彊彊彊
𠂔械掣舉掀擣攢檠敦𦉳𦉳𦉳𦉳𦉳𦉳𦉳𦉳𦉳𦉳𦉳
吁叵𩚑𩚑𩚑𩚑𩚑𩚑𩚑𩚑𩚑𩚑𩚑𩚑𩚑𩚑𩚑𩚑𩚑𩚑𩚑𩚑
洵涓瀾𦉳𦉳𦉳𦉳𦉳𦉳𦉳𦉳𦉳𦉳𦉳𦉳𦉳𦉳𦉳𦉳𦉳𦉳𦉳𦉳
洵涓瀾𦉳𦉳𦉳𦉳𦉳𦉳𦉳𦉳𦉳𦉳𦉳𦉳𦉳𦉳𦉳𦉳𦉳𦉳𦉳𦉳

Introduction

叅爰碇叁収支段妄叔呷唢媵囡噉嘒嘑團
塹壅罌嚮鉅夢寅龔變鄴獫嫗燈熾娑享亞孳
香互弃寓討对𪛖允佻膾牖龙馐巘嶸巖幽𠂔𦏧
笑配潞沓幙欄櫜勝牖牖𣎵并薦廕庖廳傍膊龐
廻巡井奔式彊弼驅擯彙𪚩𪚩徬招扁籠慤忭憇
𠄎械掣舉掀揉攢縈𥽿𥽿𥽿𥽿𥽿𥽿𥽿𥽿𥽿𥽿
吁咍勔醵甌朞楨桃櫟歆欬欬步殛毆馳氺氺氺
洶涓瀾滄烈慟爬再爰寧𪚩𪚩𪚩𪚩𪚩𪚩𪚩𪚩𪚩𪚩

Introduction

Imbalanced Ratio

imbalanced ratio = majority class / minority class

ZooScan

$$427 / 28 = 15.25$$

Kaggle

$$1979 / 9 = 219.89$$

WHOI

$$2606720 / 4 = 651680$$

EVEN MORE THAN 10^8 !

WHY?

Evaluation Criteria

Name	Formula	Explanation
True Positive Rate (TP rate)	$TP / (TP + FP)$	The closer to 1, the better. TP rate = 1 when FP = 0. (No false positives)
True Negative Rate (TN rate)	$TN / (TN + FN)$	The closer to 1, the better. TN rate = 1 when FN = 0. (No false negatives)
False Positive Rate (FP rate)	$FP / (FP + TN)$	The closer to 0, the better. FP rate = 0 when FP = 0. (No false positives)
False Negative Rate (FN rate)	$FN / (FN + TP)$	The closer to 0, the better. FN rate = 0 when FN = 0. (No false negatives)

$$G - mean = \sqrt{TPr * TNr}$$

$$Precision = \frac{TP}{TP + FP}$$

$$Recall = \frac{TP}{TP + FN} = TPr$$

$$F - measure = \frac{2 * Precision * Recall}{Precision + Recall}$$

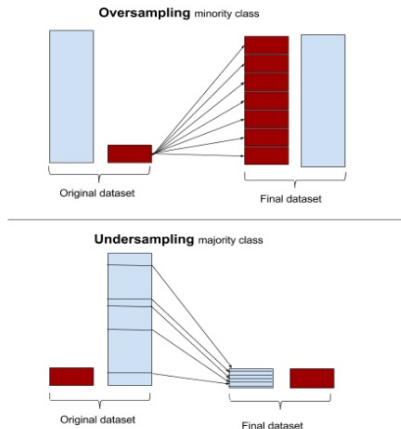
Approaches

Overview

- Sampling
 - Under-sampling*
 - Over-sampling*
- Cost-sensitive learning
- Ensembled classifier
 - EasyEnsemble*
 - BalanceCascade*

Approaches

Sampling



Best approach: SMOTE

Approaches

Cost-sensitive

$$L(x, i) = \sum_j P(j|x) c(i, j)$$

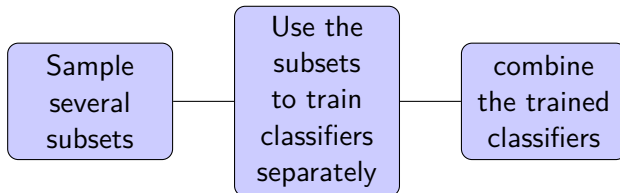
Minimize the overall cost.

- x : an example
- i : a class
- j : the j^{th} class
- P : Probability
- c : cost matrix

Best approaches: AdaCost, AsymBoost

Approaches

Ensembled Classifier



Best approaches: EasyEnsemble, BalanceCascade, SMOTEBoost

Approaches

EasyEnsemble.M \Rightarrow EasyEnsemble.D

- 1: Input: A set of minority class examples \mathcal{P} , $k-1$ sets of majority class examples \mathcal{N} , $|\mathcal{P}| < |\mathcal{N}_k|$, the number of subsets T to sample from \mathcal{N}_k , and s_i , the number of iterations to train an AdaBoost ensemble H_i
- 2: for $i \leftarrow 1:T$
- 3: $D_i = \mathcal{N}_1$
- 4: for $t \leftarrow 1:k$
- 5: Randomly sample a subset \mathcal{N}_{it} from \mathcal{N}_k , $N_{it}, |N_{it}| = |P| + \frac{C_1 * (|C_i| - |P|)}{|C_k|}$ in the t^{th}
- 6: $D_i = D_i \cup \mathcal{N}_{it}$
- 7: $H_t(x) = \text{sgn}(\sum_{d=1}^{s_i} \alpha_{t,d} h_{t,d}(x) - \theta_i)$
- 8: $H(x) = \text{sgn}(\sum_{t=1}^T \sum_{d=1}^{s_i} \alpha_{t,d} h_{t,d}(x) - \sum_{t=1}^T \theta_i)$

Future Work

- Optimize the algorithm to cost less runtime
- Use Kaggle and WHOI datasets
- Increase the amount of time in each dataset



Q&A