

論文 / 著書情報
Article / Book Information

題目(和文)	
Title(English)	Learning under Class-Balance Change: Distribution Matching via Direct Divergence Estimation
著者(和文)	MARTHINUSC.DUPLES
Author(English)	Marthinus Christoffel du Plessis
出典(和文)	学位:博士(工学), 学位授与機関:東京工業大学, 報告番号:甲第9558号, 授与年月日:2014年3月26日, 学位の種別:課程博士, 審査員:杉山 将,佐藤 泰介,徳永 健伸,村田 剛志,瀬々 潤
Citation(English)	Degree:Doctor (Engineering), Conferring organization: Tokyo Institute of Technology, Report number:甲第9558号, Conferred date:2014/3/26, Degree Type:Course doctor, Examiner:,,,,,
学位種別(和文)	博士論文
Category(English)	Doctoral Thesis
種別(和文)	論文要旨
Type(English)	Summary

(博士課程)
Doctoral Program

論文要旨

THESIS SUMMARY

専攻:	Computer Science	専攻
Department of		
学生氏名:	Marthinus Christoffel du	
Student's Name	Plessis	

申請学位 (専攻分野):	博士 (Engineering)
Academic Degree Requested	Doctor of
指導教員 (主):	Masashi Sugiyama
Academic Advisor(main)	
指導教員 (副):	
Academic Advisor(sub)	

要旨 (和文 2000 字程度)

Thesis Summary (approx.2000 Japanese Characters)

In most machine learning algorithms, it is assumed that the training and test environment are the same and that the supervisor (teacher) assigns labels to all training samples. In many real-world datasets, these assumptions are however violated due to a changing environment or imperfect supervision.

Many of these situations in the classification scenario can be characterized as a change in class balance. In this thesis, three such problems are considered: classification under class-balance change, labeling of unsupervised datasets differing by class balance, and classification of partially labeled data.

Due to real-world effects such as biased selection or non-stationarity, the class balance between the training and test dataset may differ. Training a classifier on labeled training data and then applying it on test data with a different class prior may cause an excess misclassification rate. This can, however, be corrected for by reweighting training samples with the class prior of the test environment.

In practice, however, the test class priors are unknown. The focus of Chapter 4 is to estimate the class balance in a semi-supervised setup. In the semi-supervised setup, labeled training samples and unlabeled test samples are available. In this thesis, it is shown that the class balance may be estimated by matching a model of the test input density to the true test input density. We further show that the distributions can be matched by minimizing an f -divergence between the two distributions. This f -divergence in turn must be estimated from samples. It was shown that the existing method can be interpreted as indirectly estimating the Kullback-Leibler divergence via posterior estimation.

We show that the class balance can be estimated by directly estimating the f -divergence via Fenchel duality. Of the f -divergences, we are in particular interested in using the Pearson divergence since it is more robust and may be analytically estimated. Empirically we show that

the proposed method obtains an accurate estimate of the class balance. When a classifier is reweighted with the estimated class-balance, it leads to a lower misclassification rate. This method was also extended to match distributions by minimizing the L_2 distance between probability densities.

Secondly, the problem of labeling of a dataset without any supervisory information is considered. In this setting, we assume that two *unlabeled* datasets with different class balances are available. We show that it is possible to obtain a labeling of unlabeled samples using only the two unlabeled datasets. This labeling corresponds to classification up to class commutation. Furthermore, we show that this labeling may be obtained by estimating the sign of the density difference between the two distributions. Using the ramp-loss, this sign may be directly estimated.

Finally, the problem of classification using positive-labeled and unlabeled data is considered. This problem often occurs in outlier detection problems, where a dataset of nominal samples and a separate polluted dataset consisting of nominal samples and outliers are available. It has previously been shown that a classifier may be trained in this setting if the class balance is known. We introduce the framework of partially matching distributions in order to estimate the class prior. These distributions may be matched by minimizing f -divergences. Analysis of the existing method shows that it can be interpreted as indirectly matching the Pearson divergence. This Pearson divergence may be analytically estimated using Fenchel duality. More importantly, due to the analytic estimator of the Pearson divergence, we can obtain an analytic estimate of the class balance. Analysis also showed that both methods are positively biased, with the bias depending on the class overlap. On real-world datasets used in experiments however, the Pearson divergence gave a very accurate estimate of the class balance.

We conclude that many different learning problems characterized by a change in class balance can be formulated as distribution matching problems. Furthermore, by selecting an appropriate divergence and directly estimating it, practical algorithms can be obtained that yield excellent experimental results.

備考：論文要旨は、和文 2000 字と英文 300 語を 1 部ずつ提出するか、もしくは英文 800 語を 1 部提出してください。

Note: Thesis Summary should be submitted in either a copy of 2000 Japanese Characters and 300 Words (English) or 1copy of 800 Words (English).