

I'm Toshihiro Kamishima.

Today, we would like to talk about fairness-aware classification.



Fairness-aware data mining is a data analysis taking into account potential issues of fairness, discrimination, neutrality, or independence.

There are several kinds of tasks: fairness-aware classification, detection of unfair events, and fairness-aware data publication.

These are examples of fairness-aware data mining applications.

Outline	
Applications fairness-aware data mining applications 	
 Difficulty in Fairness-aware Data Mining Calders-Verwer's discrimination score, red-lining effect 	
 Fairness-aware Classification fairness-aware classification, three types of prejudices 	
Methods prejudice remover regularizer, Calders-Verwer's 2-naïve-Bayes 	
 Experiments experimental results on Calders&Verwer's data and synthetic data 	
Related Work	
 privacy-preserving data mining, detection of unfair decisions, explainability, fairness-aware data publication 	
Conclusion	
	3

This is an outline of our talk.

We first give examples of fairness-aware data mining applications.

After showing difficulty in fairness-aware data mining, we define a fairness-aware classification task.

We then propose our method for this task, and empirically compare it with an existing method. Finally, we review this new research area, and conclude our talk.



Let me start by giving examples of applications.



The first application is the elimination of discrimination.

Due to the spread of data mining technologies, data mining is being increasingly applied for serious decisions. For example, credit scoring, insurance rating, employment application, and so on.

Additionally, accumulation of massive data enables to reveal personal information.

To cope with these situation, fairness-aware data mining techniques are used for excluding the influence of socially sensitive information from serious decisions: information considered from the viewpoint of social fairness, and information restricted by law or contracts.



The second application is related to the filter bubble problem, which is a concern that personalization technologies narrow and bias the topics of information provided to people. Pariser shows an example of a friend recommendation list in Facebook. To fit for his preference, conservative people are eliminated form his recommendation list, while this fact is not notified to him.



To cope with this filter bubble problem, an information neutral recommender system enhances the neutrality from a viewpoint specified by a user and other viewpoints are not considered. For this purpose, fairness-aware data mining techniques are applied. In the case of Pariser's Facebook example, a system enhances the neutrality in terms of whether conservative or progressive, but it is allowed to make biased recommendations in terms of other viewpoints, for example, the birthplace or age of friends.

Non-Redundant Clustering [Gondek+ 04] **non-redundant clustering** : find clusters that are as independent from a given uninteresting partition as possible a conditional information bottleneck method, which is a variant of an information bottleneck method clustering facial images • Simple clustering methods find two clusters: one contains only faces, and the other contains faces with shoulders • Data analysts consider this clustering is useless and uninteresting • A non-redundant clustering method derives more useful male and female clusters. which are independent of the above clusters 8

The third application is excluding uninteresting information.

Outline	
Applications applications of fairness-aware data mining 	
Difficulty in Fairness-aware Data Mining • Calders-Verwer's discrimination score, red-lining effect	
 Fairness-aware Classification fairness-aware classification, three types of prejudices 	
Methods prejudice remover regularizer, Calders-Verwer's 2-naïve-Bayes 	
Experiments experimental results on Calders&Verwer's data and synthetic data 	
Related Work privacy-preserving data mining, detection of unfair decisions, explainability, situation testing, fairness-aware data publication 	
Conclusion	
	9

We then show the difficulty in fairness-aware data mining.

Diffic	ulty in Fa	irness-awa	re Data Mini	
US Census Data : predict whether their income is high or low				
		Male	Female	
Hi	gh-Income	3,256 <u>fe</u> r	wer 590	
Lo	ow-income	7,604	4,831	
of Hig /hile 3 <mark>Occ</mark>	Females are r h-Male data is 5 0% of Male data	ninority in the high 5.5 times # of High-F 1 are High income, c ning techniques pre	Female data Female data only 11% of Females a	are
	Minor pa and thus mir	atterns are frequentl norities tend to be tr	y ignored eated unfairly	

This is Caldars & Verwer's example to show why fairness-aware data mining is needed.

This is a task to predict whether their income is high or low.

In this US census data, females are minority in the high-income class.

Because mining techniques prefer simple hypothesis, minor patterns are frequently ignored; and thus minorities tend to be treated unfairly.



Caldars & Verwer proposed a score to quantify the degree of unfairness.

This is defined as the difference between conditional probabilities of High-income decisions for males and females. The larger score indicates the unfairer decision.

The baseline CV score for the US Census Data Set is 0.19.

If objective variables are predicted by a naïve Bayes classifier, the CV score increases to 0.34, indicating unfair treatments.

Even if a feature, gender, is excluded, the CV score is improved to 0.28, but still being unfairer than its baseline.

Consequently, ignoring sensitive features is ineffective against the exclusion of their indirect influence. This red-lining effect makes fairness-aware data mining difficult.

	Outline
Ap	 applications applications of fairness-aware data mining
Di	 Friculty in Fairness-aware Data Mining Calders-Verwer's discrimination score, red-lining effect
Fa	 irness-aware Classification fairness-aware classification, three types of prejudices
M	ethods ● prejudice remover regularizer, Calders-Verwer's 2-naïve-Bayes
Ex	 experiments experimental results on Calders&Verwer's data and synthetic data
Re	 lated Work privacy-preserving data mining, detection of unfair decisions, explainability, situation testing, fairness-aware data publication
Сс	onclusion

We next define a task of fairness-aware classification



We introduce three types of variables.

An objective variable Y represents a result of serious decision. A sensitive feature S represents socially sensitive information. The other variables are non-sensitive features X.



A prejudice is one of causes of unfairness, which is defined as the statistical dependences of an objective variable or non-sensitive features on a sensitive feature.

Prejudices is classified into three types:

Direct prejudice is a clearly unfair state that a prediction model directly depends on a sensitive feature.

Indirect prejudice is the statistical dependence of an objective variable on a sensitive feature. Latent prejudice is the statistical dependence of non-sensitive features on a sensitive feature. We here focus on removing this indirect prejudice.



Fairness-aware classification is a variant of a standard classification.

In a case of a standard classification task, training data is sampled from a unknown true distribution.

A goal of this task is to learn a model approximating the true distribution.

In a case of a fairness-aware classification task, we assume a true fair distribution that satisfies two constraints.

This distribution should be similar to the true distribution, and should satisfy no indirect prejudice condition.

A goal of this task is to learn a model approximating this true fair distribution from training data sampled from the true distribution.

	Outline
	ations
Difficu • Ca	Ilty in Fairness-aware Data Mining
Fairne ● fai	ss-aware Classification rness-aware classification, three types of prejudices
Metho • pre	o <mark>ds</mark> ejudice remover regularizer, Calders-Verwer's 2-naïve-Bayes
Experi • ex	ments perimental results on Calders&Verwer's data and synthetic data
Relate • pri ex	d Work vacy-preserving data mining, detection of unfair decisions, plainability, situation testing, fairness-aware data publication
Concl	usion

We show our logistic regression with prejudice remover regularizer and Calders-Verwer's 2naïve-Bayes method.



We propose logistic regression with prejudice remover regularizer. We change an original logistic regression model in these two points. First, we add ability to adjust distribution of Y depending on the value of S. For this purpose, multiple logistic regression models are built separately, and each of these models corresponds to each value of a sensitive feature.



Second, we add a constraint of a no-indirect-prejudice condition, which is a main idea of this presentation.

For this purpose, we add this term to an original objective function.

This prejudice remover regularizer is designed so that the smaller value more strongly constraints the independence between Y and S.

Eta is a fairness parameter. The lager value more enforces the fairness.



As this prejudice remover regularizer, we adopt mutual information between S and Y so as to enforce the independence between Y and S.

This term can be computed by replacing with the summation over samples. But the computation of this term is rather complicated.



This distribution can be derived by marginalizing over X, but this is computationally heavy. We hence approximate by the sample mean over x for each pair of y and s. Unfortunately, this technique is applicable only if both Y and S are discrete.



We compared our method with Calders-Verwer's two-naïve-Bayes method. Unfair decisions are modeled by introducing of the dependence of X on S as well as on Y.



After parameters are estimated by the corresponding empirical distributions, joint distribution of Y and S is modified by this algorithm.

This algorithm updates the joint distribution so that its CV score decreases while keeping the updated distribution similar to the empirical distribution.

Outline	
Applications applications of fairness-aware data mining 	
 Difficulty in Fairness-aware Data Mining Calders-Verwer's discrimination score, red-lining effect 	
 Fairness-aware Classification fairness-aware classification, three types of prejudices 	
Methods prejudice remover regularizer, Calders-Verwer's 2-naïve-Bayes 	
Experiments • experimental results on Calders&Verwer's data and synthetic data	
 Related Work privacy-preserving data mining, detection of unfair decisions, explainability, situation testing, fairness-aware data publication 	
Conclusion	

We next show our experimental results.



We first test these four methods on Calders and Vewer's Test Data.

LRns and NBs are respectively pure logistic regression and naïve Bayes without sensitive features.

PR is our logistic regression with prejudice remover regularizer.

CV2NB is Calders and Verwer's two-naïve-Bayes method.



We use two types of evaluation measures.

Accuracy measures how correct predicted classes are.

Normalized mutual information measures how fair predicted classes are.



These are experimental results.

X-axes correspond to fairness parameters, the lager value more enhances the fairness.

This chart (left) shows the change of prediction accuracy.

This chart (right) shows the change of the degree of fairness.

As the increase of a fairness parameter $\eta,$ prediction accuracy is worsened and the fairness is enforced.

Our metod could make fairer decisions than pure logistic regression and naïve Bayes.

Additionally, our method could make more accurate prediction than naïve Bayes and CV2NB. Unfortunately, CV2NB achieved near-zero NMI, but our method could not achieve it.



To investigate why our method failed to make a fairer prediction than that made by CV2NB, we tested our method on synthetic data.

Each sample composed of one sensitive feature, two non-sensitive features, and one objective variable.

Non-sensitive feature X is independent from S, but W depends on S.

And, both X and W equally influence an objective variable.

	M[Y X, S=0]		M[Y X, S=0]			
	X	W	bias	X	W	bias
η=0	11.3	11.3	-0.0257	11.3	11.4	0.0595
η=150	55.3	-53.0	-53.6	56.1	54.1	53.6
When η= for <i>X</i> and	=0, PR reg d <i>W</i> are alr	ularizer do nost equal	esn't affec	t weights,	and these	weights
When η= for X and When η= PR igno	=0, PR reg d <i>W</i> are alr =150, abso pres feature	ularizer do nost equal plutes of w es dependi	esn't affec eights for 2 ing on <i>S</i> if t	t weights, K are large he influenc	and these r than thos es to <i>Y</i> are	weights se for <i>W</i> equal

A prejudice remover regularizer doesn't affect weights if $\eta=0$, but it heavily does if $\eta=150$. When $\eta=0$, weights for x and w are almost equal.

When η =150, absolutes of weights for X are larger than those for W.

This indicates that our prejudice remover regularizer ignores features depending on s if the influences to y are equal.

On the other hand, CV2NB method treats all features equally.

Consequently, our prejudice remover can consider the differences among individual features, but Calders-Verwer's 2-naïve-Bayes can improve fairness more drastically.

Applications • applications of fair	rness-aware data mining
Difficulty in Fairness ● Calders-Verwer's (-aware Data Mining discrimination score, red-lining effect
airness-aware Clas ● fairness-aware cla	sification ssification, three types of prejudices
Vlethods ● prejudice remover	regularizer, Calders-Verwer's 2-naïve-Bayes
Experiments • experimental result	Its on Calders&Verwer's data and synthetic data
Related Work	g data mining, detection of unfair decisions,

We finally review this new research area.



Here, we'd like to point out the relation between an indirect prejudice and PPDM.

From information theoretic perspective, an indirect prejudice implies that mutual information between Y and S is non-zero.

From the viewpoint of privacy-preservation, this can be interpreted as the leakage of sensitive information when an objective variable is known.

On the other hand, there are some differences from PPDM like this.



Pedreschi et al. firstly addressed the problem of discrimination in data mining. Their algorithm enumerates unfair association rules.



Zliobaite et al. proposed a notion of explainability.



Luong et al. proposed another notion, situation testing.



Dwork et al. proposed a framework for fairness-aware data publishing.

Conclusion

Contributions

- the unfairness in data mining is formalized based on independence
- a prejudice remover regularizer, which enforces a classifier's independence from sensitive information
- experimental results of logistic regressions with our prejudice remover

Future Work

- improve the computation of our prejudice remover regularizer
- fairness-aware classification method for generative models
- another type of fairness-aware mining task

Socially Responsible Mining

 Methods of data exploitation that do not damage people's lives, such as fairness-aware data mining, PPDM, or adversarial learning, together comprise the notion of socially responsible mining, which it should become an important concept in the near future.

Our contributions are as follows.

In future work, we have to improve computation of prejudice remover regularizer Methods of data exploitation that do not damage people's lives, such as fairness-aware mining, PPDM, or adversarial learning, together comprise the notion of socially responsible mining, which it should become an important concept in the near future.



Program codes and data sets are available at these URLs.

We finally wish to thank Dr. Sicco Verwer for providing detail information about their work and program codes.

That's all I have to say. Thank you for your attention.