

# Classifying Socially Sensitive Data Without Discrimination: An Analysis of a Crime Suspect Dataset

Faisal Kamiran

King Abdullah University of Science and Technology, KSA  
faisal.kamiran@kaust.edu.sa

Asim Karim

Lahore University of Management Sciences, Pakistan  
akarim@lums.edu.pk

Sicco Verwer

Ministry of Security and Justice, the Netherlands  
Radboud University Nijmegen, the Netherlands  
s.verwer@cs.ru.nl

Heike Goudriaan

Statistics Netherlands, the Netherlands  
h.goudriaan@cbs.nl

**Abstract**—Social discrimination against certain sensitive groups within society (e.g., females, blacks, minorities) is prohibited by law in many countries. To prevent discrimination arising from the use of discriminatory data, recent data mining research has focused on methods for making classifiers learned over discriminatory data discrimination-aware. Most of these methods have been tested on standard classification datasets that have been tweaked for discrimination analysis rather than over actual discriminatory data. In this paper, we study discrimination-aware classification when applied to a real-world dataset of *Statistics Netherlands*, which is a census body in the Netherlands. Specifically, we consider the use of classifiers for predicting whether an individual is a crime suspect, or not, to support law enforcement and security agencies' decision making. Our results show that discrimination does exist in real world datasets and blind use of classifiers learned over such datasets can exacerbate the discrimination problem. We demonstrate that discrimination-aware classification methods can mitigate the discriminatory effects and that they lead to rational and legally acceptable decisions.

**Keywords**-discrimination; classification

## I. INTRODUCTION

Social discrimination refers to biased decision making in favor of, or against, a person or a thing on the basis of affiliation of that person or thing to a certain group, class, or category rather than on merit. Discriminatory practices exist in employment, income, education, finance, and other benefits/services when decisions are made on the basis of sensitive attributes like age, gender, skin color, religion, race, language, culture, marital status, economic condition, and other non-merit factors. Discrimination is increasingly often considered unacceptable from social, ethical, and legal perspectives. Many anti-discrimination laws [1], [2], [3], [4] have been enacted and several anti-discrimination organizations (e.g., ENAR [5]) are working for the eradication of discrimination. The consequences of discriminatory practices can range from legal conviction to a variety of social problems like high unemployment rate, frustration, low productivity, and disputes.

The discrimination-aware classification problem studies the construction and use of classifiers learned from discriminatory or biased data. This problem is relatively new and was recently introduced by Pedreschi et al. [6], however, it was warmly welcomed by the data mining research community with the development of many novel discrimination detection and prevention methods. Moreover the recent debate at the European Parliament on the reform of the privacy directive that includes a new article on profiling and non-discrimination, makes it a hot topic for the data mining and legal research community [7], [8].

Recently, we collaborated with the Dutch Research and Documentation Center (WODC) associated with the Ministry of Security and Justice. One of the responsibilities of this center is analyzing and modeling demographic and crime data from *Statistics Netherlands*, the national census body, to support decision and policy making. The center showed interest in implementing discrimination-aware techniques to ensure that it gives non-discriminatory recommendations (w.r.t. different sensitive attributes like ethnicity and gender) to the decision makers, e.g., the Minister of Security and Justice. This interest emerged from the realization that removing the sensitive attribute ethnicity from the classification model does not remove the correlation between attributes ethnicity and crime suspect because of indirect discrimination, and usage of standard classifiers learned on such data can lead to discriminatory recommendations. In fact, the discrimination problem can be exacerbated by the blind use of standard discrimination-ignorant classifiers. This is troublesome as often it is assumed that learned classifiers provide accurate and unbiased decisions. The WODC was interested in investigating the effect that a non-discriminatory view of their data would have on policy making. Thus, we proposed the discrimination-aware data mining paradigm in order to provide appropriate solutions for avoiding discrimination.

In this paper, we study discrimination-aware classification

using a real world dataset of *Statistics Netherlands* with the aim of controlling the propagation of discrimination effects to automated decision making process. This dataset contains demographic, economic, and crime information on the whole Dutch population. We investigate the use of this dataset for the automatic profiling of persons as crime suspects or not. Clearly, it is undesirable for these profiles to be discriminatory. In our experiments, we show that if we learn a standard classifier without taking the discrimination effect into account, the learned classifier carries this discrimination effects to future decision making. In fact, since the classifier uses discriminatory correlations in its decisions, the discrimination effect is even increased. We present and discuss the results of discrimination-aware classification techniques on this dataset, demonstrating that the use of discrimination-aware classification techniques for a large part neutralizes the discriminatory effects while only incurring a small cost in predictive accuracy. The resulting models are in our opinion much more useful for policy making than blindly learned standard classifiers.

The rest of the paper is organized as follows. Section II gives an overview of existing discrimination-aware data mining works and a brief review of the discrimination-aware classification techniques we use in this study. Section III introduces *Statistics Netherlands* and gives general information on the available data. Section IV provides detailed information on the data used and presents experimental evaluations of discrimination-aware classification techniques when applied on this data. We conclude our study in Section V.

## II. DISCRIMINATION-AWARE DATA MINING

The concept of discrimination has been studied in social sciences for a long time. These studies led to the rise of many anti-discrimination organizations (e.g., ENAR [5]) and to the enactmentment of many anti-discrimination laws [1], [2], [3], [4]. However, the concept of discrimination-aware data mining is quite new and got attention from the computer science research community only in the recent years. Discrimination-aware data mining works can be divided into two main categories, i.e. the detection of discriminatory patterns from a given dataset (discrimination detection) and learning of discrimination-aware classifiers to avoid biased decision making in future (discrimination prevention).

Direct discrimination arises when sensitive or discriminatory attributes are utilized in learning and prediction. Nonetheless, it has been shown that discrimination is not removed by simply removing these attributes from the dataset [9]. That is, discriminatory decisions can still be made due to correlation of sensitive attributes with other attributes (indirect discrimination, also known as the *redlining*<sup>1</sup>). This issue has been studied in greater detail in [10], [11].

<sup>1</sup><http://en.wikipedia.org/wiki/Redlining>, May. 12th, 2012

Pedreschi et al. [6], [12], [13], [14] mainly focus on discrimination detection by discovering discriminatory classification rules from biased datasets following a frequent itemset mining approach coupled with a measure of discrimination. A central notion in the works on identifying discriminatory rules is that of the *context* of the discrimination. This context shows those regions where the discriminatory practices are more obvious. [14] addresses both the discrimination detection and the discrimination prevention problem. It proposes a variant of k-NN classification for the discovery of discriminated objects. The data object of a deprived community (e.g., female) with a treatment that is significantly different from its neighbors (objects from the favored community, e.g., males) are considered discriminated objects. The discrimination is prevented by changing the class labels of these discriminated objects.

Proposed methods for discrimination prevention are either based on data preprocessing or algorithm/model tweaking. Data preprocessing methods modify the biased data by removing discriminatory patterns from it before learning a prediction model from it. In works on discriminatory rule protection [15], [16], [17], data transformations are proposed for making discovered discriminatory classification rules discrimination-free according to a discrimination measure. The key limitation of these methods is their applicability to rule based classifiers only that may not be the best classifier for a given problem. In [18], [19], [20], [9], data sampling and massaging techniques are presented for removing discrimination w.r.t. a single sensitive attribute. We discuss these methods in more detail later in this section.

Proposed methods for discrimination prevention requiring learning model adaptation include those for decision trees [21], naive Bayes classifiers [22], and logistic regression [23]. All these methods require that the learning model or algorithm is tweaked, and the first two methods are specific to their respective classifiers. For example in [21], the authors propose a strategy for relabeling the leaf nodes of a decision tree to make it discrimination-free.

In this paper we present the results obtained by massaging and reweighing techniques of [9] and modifying the decision thresholds as explained in [22] on a real world dataset. Now we give an overview of the massaging, reweighing, and threshold modification methods in more detail.

1) *Massaging*: In *massaging* we change the class labels in the training data to make it discrimination-free; some objects of the deprived community (e.g., females) change from class  $-$  to  $+$ , and the same number of objects of the favored community (e.g., males) change from  $+$  to  $-$ . In this way the discrimination decreases, yet the overall class distribution is maintained; the same number of people have the positive class as before. This strategy reduces the discrimination to the desirable level with the least number of changes to the dataset while keeping the overall class distribution fixed. It is important to notice that we do not randomly pick the

objects to relabel. Instead, we use a ranker to rank the objects of favored and deprived community separately with respect to the class probability. Any probabilistic classifier can be used as a ranker. Based on this ranking we can see, for the deprived and favored communities separately, which instances are closest to the *decision boundary*. The objects close to the decision boundary are those with a probability close to 0.5. We discovered in our experiments that the data objects close to the decision boundary are more vulnerable to the effect of discrimination. Therefore, we select these objects first to relabel. Figure 1 shows the framework of massaging technique.

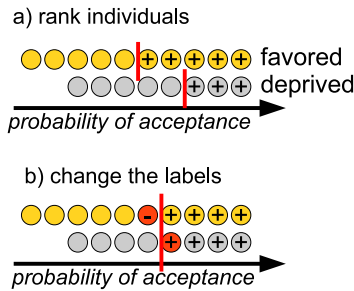


Figure 1. Massaging technique

2) *Reweighting*: The *reweighing* approach is less intrusive than *massaging* as it does not change the class labels of the objects. Instead of relabeling the objects, different weights are attached to them. For example, the deprived community objects with + class get higher weights than the deprived community objects with - class and the favored community objects with + get lower weights than the favored community objects with - class. We refer the reader to [9] to have detailed idea of weight calculation of objects to make the training dataset discrimination-free. The procedure for reweighing is presented in Figure 2.

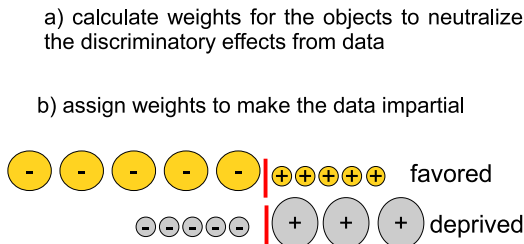


Figure 2. Sampling technique

3) *Modifying thresholds*: *Modifying decision thresholds* is a simple technique that works not on the dataset but on a predictive classification model such as the Naive Bayes classifier. The decision thresholds determine what

probability an object needs in order to be labeled as positive or negative (suspect or not). By changing these thresholds for the different communities (see Figure 3), we can influence the positive class probabilities of the favored and deprived communities. We increase the decision threshold of favored communities until their positive class probability equals to overall positive class probability. Similarly, we decrease the decision threshold for the deprived communities. The ratio of assigned positive labels in the data set is then equal for every community. However, since the sensitive value denoting the community of objects is unknown during classification, a second model is learned in order to predict the community value. The threshold corresponding to this community value is then used to classify an object.

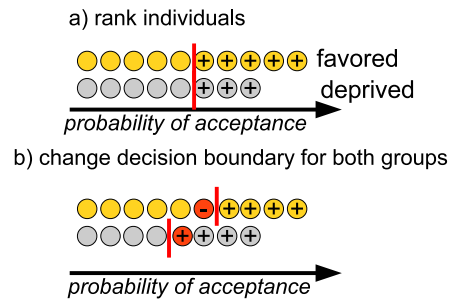


Figure 3. Modifying decision thresholds

### III. STATISTICS NETHERLANDS

Statistics Netherlands is a national body responsible for collecting and processing data, making it available to support policy making and scientific research. In addition to its responsibility for (official) national statistics, *Statistics Netherlands* is responsible for providing data on the Netherlands for the production of European (community) statistics. The legal basis for *Statistics Netherlands* and its work is the Act of 20 November 2003 last amended by the Act of 15 December 2004.

The information *Statistics Netherlands* publishes encompass a multitude of societal aspects, from macro-economic indicators such as economic growth and consumer prices, to the incomes of people and households, to for instance butterfly populations and the number and type of registered crimes. *Statistics Netherlands*' statistical programs (the long-term statistical program and the annual work program) are set by the Central Commission for Statistics. This is an independent commission that watches over the independence, impartiality, relevance, quality, and continuity of the statistical program. *Statistics Netherlands* decides autonomously which methods to use to make these statistics, and whether or not to publish results.

Statistics Netherlands is allowed to conduct supplementary surveys among companies and private persons. Companies are usually obliged by law to supply information to

*Statistics Netherlands* and can be forced to cooperate under certain circumstances. For its part *Statistics Netherlands* is obliged to keep all individual data confidential.

#### **Social Statistics Database**

The data used in this study is derived from the Social Statistics Database (SSD) from *Statistics Netherlands*. The SSD contains information from many different registered sources on the whole Dutch population where the Dutch population is defined as all persons who have been (temporarily) registered as a resident of the Netherlands. Since data from the different registered sources is coupled, thus information on different domains (e.g., demographic, socioeconomic, education, health) becomes available for every person in the Dutch population.

Demographic information (e.g., gender, date of birth, country of birth, marital status, town of residence) is mainly derived from the municipal personal records database (Gemeentelijke Basisadministratie persoonsgegevens; GBA). All changes in the GBA, like moving houses or divorces, are also recorded. Much of the socioeconomic information has been made available through the Tax and Customs Administration (Belastingdienst) and the Institute for the Execution of Employee Insurances (Uitvoeringsinstituut Werknemersverzekeringen; UWV) – the organization that is responsible for the implementation of employee insurances and provision of labor market and data services.

The SSD also contains information on whether or not a member of the population has been a crime suspect. This information is available from 1999 onwards and is derived from a nationwide database used by the Dutch Police for the registration of criminal suspects, Herkenningsdienstsysteem (HKS). A crime suspect is someone against whom a police report has been filed and who is at least 12 years old. For every crime suspect, additional information is available on the crime incident(s) of which one is accused (e.g., the type of crime).

## IV. EXPERIMENTAL EVALUATION

### A. Dataset

In the present study, a dataset has been created from the available information in the SSD, consisting of a random sample of 5% of the people between 12 and 79 years who have been arrested as a crime suspect at least once in 2006 (no. of suspects=10,239) and a random sample of 0.15% of the people between 12 and 79 who have not been arrested as a crime suspect in 2006 (no. of non\_suspect=20,478). This resampling is done to overcome the severe imbalance between number of crime suspects and non-suspects in the original data, which can cause classifiers to label all persons as non-suspects. In the end, this sampling procedure resulted in a dataset with 30,717 unique persons, of which exactly one third had been a crime suspect in 2006.

Each person is described by 39 attributes. These attributes can be divided into four main categories: demographic information, family information, socioeconomic information, and neighborhood information. The detail of attributes in each category is as follows:

**Demographic Information:** gender, age, number of relocations within the municipality in the previous 5 years, number of relocations outside the municipality in the previous 5 years, degree of urbanization of the municipality, type of household, position within the household, regional code, ethnic group, first or second generation immigrant, has been in the Netherlands less than 5 years, and has been in the Netherlands less than 10 years.

**Family Information:** (legal) parents live in the same address, one or both (legal) parents receive social benefits, one or both (legal) parents have a job, and dummies indicating whether one or more of the previous three variables is missing.

**Socioeconomic information:** monthly income of the household before taxes, monthly income of parents before taxes, monthly personal income from employment before taxes, socioeconomic status in categories, living in a rented or owner-occupied dwelling, value of the dwelling, and percentile score of household income.

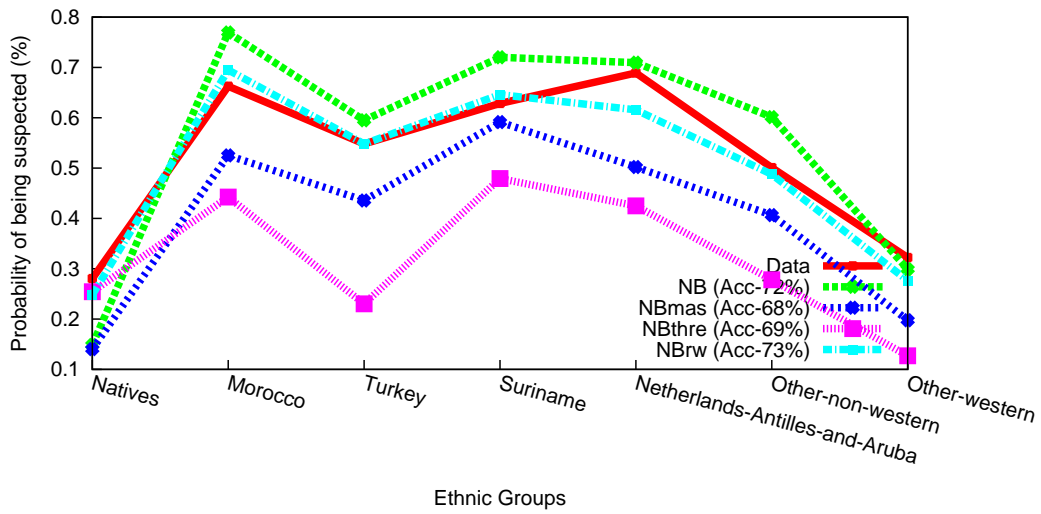
**Information on Neighborhood:** total population, number of natives, number of non-western immigrants, number of people aged 15-64 with employment, number of natives aged 15-64 with employment, number of non-western immigrants aged 15-64 with employment, number of persons receiving social benefits, number of persons receiving social benefits since over a year, number of rented residencies, and number of owner-occupied residencies.

### B. Experimental Setup

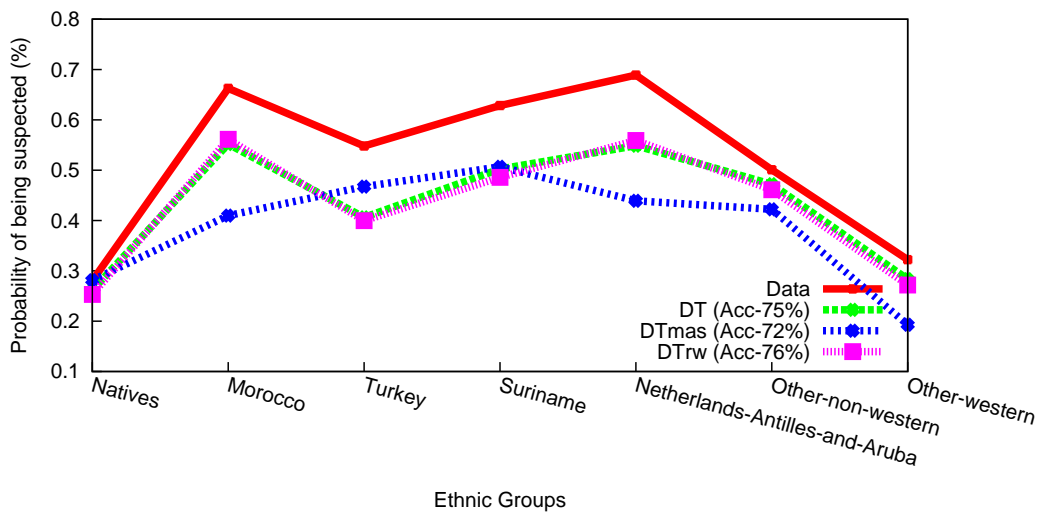
We consider the discrimination-aware classification problem and evaluate discrimination arising from classifiers trained on the discriminatory dataset. We assume ethnic group (Native, Morocco, Turkey, Suriname, Netherlands-Antilles and Aruba, Other Non-Western, and Other Western) to be the sensitive attribute and crime suspect (yes or no) to be the class attribute. In our experiments, we compare the following classification methods (standard implementations available in the `e1071` and `rpart` packages in R<sup>2</sup>):

- Traditional Naive Bayes Classifier (label NB).
- Naive Bayes Classifier learned over massaged training data (label NBmas).
- Naive Bayes Classifier learned over reweighed training data (label NB<sub>rw</sub>).
- Naive Bayes Classifier with modified decision thresholds (label NB<sub>thre</sub>).
- Traditional Decision Tree Classifier (label DT).

<sup>2</sup><http://www.r-project.org/>



(a) Naive Bayes



(b) Decision Tree

Figure 4. Crime suspect probability of different ethnic groups produced by classifiers and in the original data.

- Decision Tree learned over massaged training data (label DTmas).
- Decision Tree learned over reweighed training data (label DT<sub>rw</sub>).

Since the learned decision trees did not provide a smooth positive class distribution, we could not use the threshold modification method on decision trees. Furthermore, note that this is a multi-valued sensitive attribute problem. Discrimination is therefore gauged by comparing the probability of being a crime suspect produced by the methods across the different ethnic groups. The results reported in the paper are obtained on unaltered test set using *10-fold cross-validation*, i.e., no preprocessing is applied to test set. We also use standard classifier parameters in our experiments.

### C. Results and Discussion

The results of our experiments are shown in Figures 4 and 5. In these figures the X-axis shows different ethnic groups, while the Y-axis of Figure 4 shows the probability of being a crime suspect and the Y-axis of Figure 5 shows the false positive rate. The accuracy scores for each method are given in the legend.

We observe that when we apply standard Naive Bayes classifier without taking discrimination into account, it predicts a high crime suspect probability for the minority ethnic groups and a low probability for the native citizens. It is important to recall the discrimination calculation method used by [9], where the difference between probabilities of favored and deprived communities for the favored class is considered discrimination. In our case, the difference between the probabilities of being a crime suspect between different ethnic groups will be considered discrimination. The ultimate goal is to make these differences as small as possible.

In the original data the crime rate for the minorities is higher than for native people (see Figure 4), but this difference in crime rate is exaggerated by the standard discrimination-ignorant classifiers. Using such classifiers would therefore be equivalent to committing discrimination against ethnic minority groups. Also, it would help some of the actual suspects among natives to escape, due to unbalanced attention to (non-suspect) minority groups. This is evident from Figure 5 that gives the false positive rate for each ethnic group. For the standard Naive Bayes classifier, the probability of being falsely accused is only 9% for a native person and 62% for a person from Morocco. Such a huge difference in false positives for different ethnic groups is a clear example of discrimination.

Figure 4(a) also shows the results of discrimination-aware Naive Bayes classification methods. We observe that threshold modification provides the best control over the discrimination problem by reducing the differences in the predicted probabilities (of being a crime suspect) between the minority groups and natives, while maintaining high

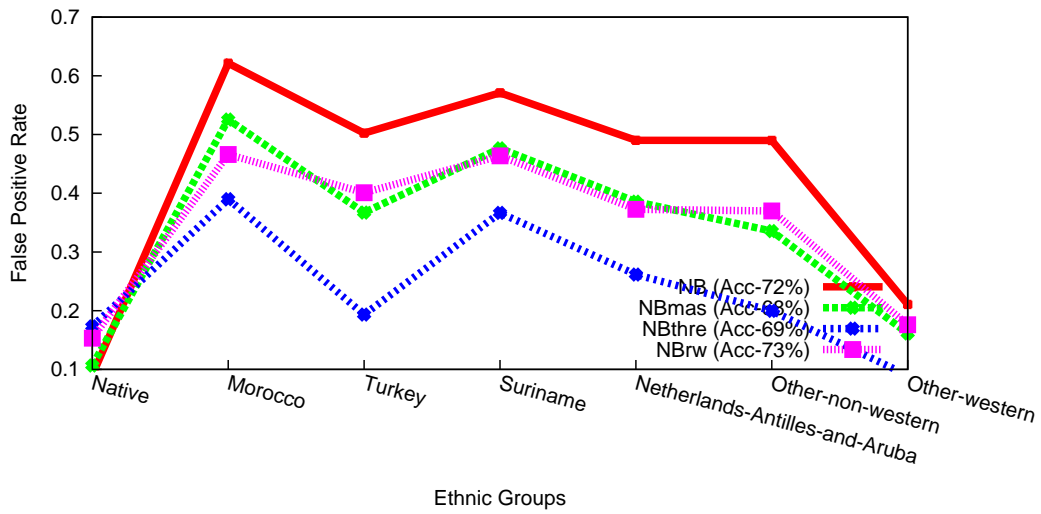
accuracy. The reweighing technique does not perform well, most likely due to its uniform weight calculation. Figure 4(b) shows similar results for discrimination-aware Decision Trees: reweighing has little effect, and massaging lowers the discrimination, at least for the two most discriminated communities (Morocco and Netherlands-Antilles and Aruba). An interesting observation is that the standard Decision Tree produces much lower discrimination as compared to the standard Naive Bayes classifier. Furthermore, it achieves a higher accuracy, highlighting that Decision Trees are clearly to be preferred over Naive Bayes classifiers on this dataset. We also note that the reweighing and massaging techniques work much better with the Naive Bayes classifier. This can be attributed to the use of a Naive Bayes classifier as a ranker in the massaging and reweighting procedures.

Figure 5 also shows the false positive rates for different discrimination-aware methods. Again, the threshold modification method provides the best (lowest and most similar) false positive rates. However, people from Morocco are still twice as likely to be falsely suspected than native Dutch people. The difference in false positive rates between the massaging and reweighing techniques are negligible. The false positive rates for the decision tree classifier are difficult to compare due to the well known unstable behavior of this classifier. The massaging technique results in higher rates for most of the ethnic groups, but the differences between the different groups are similar for all methods.

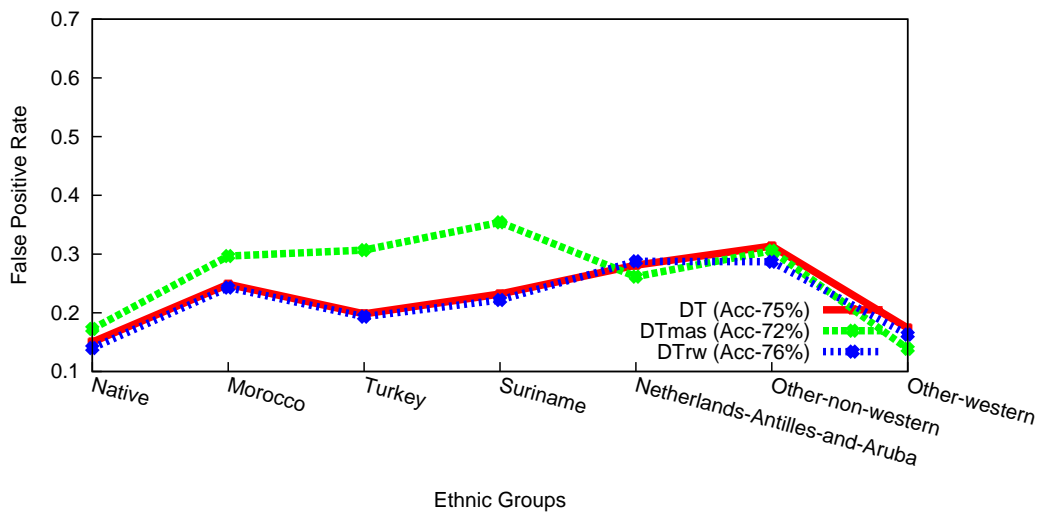
Overall, we find that discrimination-aware classification methods produce lower discrimination and more uniform false positive rates. Thus, they are strongly recommended over standard discrimination-ignorant methods. Among the different discrimination-aware methods evaluated here, the Naive Bayes classifier with threshold modification appears most appropriate for this dataset. Decision trees, although being better in their standard form, are less robust for discrimination-aware classification of this dataset.

## V. CONCLUSION

This paper presents an analysis of a real-world dataset w.r.t. social discrimination. In particular, we studied discrimination-aware classification by testing on an actual dataset from *Statistics Netherlands*, which maintains demographic, economic, and crime information of all Dutch citizens. Our results show that using a standard discrimination-ignorant classifier exacerbates the discrimination problem by increasing the probability difference of being a crime suspect between people from minority and from non-minority groups. Furthermore, people from the minority groups are more likely to be incorrectly classified as a crime suspect when using such methods. These results highlight the importance of discrimination-aware classifiers in practice. Among the three discrimination-aware techniques evaluated, we find that modifying the decision threshold of a Naive Bayes classifier produces good discrimination control, and that



(a) Naive Bayes



(b) Decision Tree

Figure 5. False positive rates for different ethnic groups produced by the classifier.

data preprocessing methods (massaging and/or reweighing) reduce discrimination for both the Naive Bayes classifier and Decision Trees. For Decision Trees, however, the reduction in discrimination is much smaller and thus are not advocated for discrimination-aware decision making. Possible explanations are that the preprocessing methods use a different (a Naive Bayes) classifier to rank the objects, or that Decision Trees already result in less discrimination. Investigating these explanations is an interesting direction for future work. For the Naive Bayes classifier, there is a large reduction in discrimination and we therefore recommend using this classifier on this data.

In future work, we would like to investigate whether it is possible to remove discrimination and to minimize false positive rates without using (or knowing) the ethnicity of a person in the prediction model and to determine the effect on the accuracy of the resulting classifier.

We conclude by saying that this study validates the usefulness of discrimination-aware data mining works in practical settings.

#### REFERENCES

- [1] C. Attorney-General's Dept., "Australian sex discrimination act 1984." 1984, via: <http://www.comlaw.gov.au/Details/C2010C00056>.
- [2] "European Union Legislation," 2012, via: [http://europa.eu/legislation\\_summaries/index\\_en.htm](http://europa.eu/legislation_summaries/index_en.htm).
- [3] "United Kingdom Legislation," 2012, via: <http://www.legislation.gov.uk/>.
- [4] "The US Federal Legislation," 2011, via: <http://www.justice.gov/crt>.
- [5] "European Network Against Racism," 1998, via: <http://www.enar-eu.org/>.
- [6] D. Pedreschi, S. Ruggieri, and F. Turini, "Discrimination-aware data mining," in *Proc. of ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, 2008.
- [7] P. De Hert and V. Papakonstantinou, "The proposed data protection regulation replacing directive 95/46/ec: A sound system for the protection of individuals," *Computer Law & Security Review*, vol. 28, no. 2, pp. 130–142, 2012.
- [8] L. Costa and Y. Pouillet, "Privacy and the regulation of 2012," *Computer Law & Security Review*, vol. 28, no. 3, pp. 254–262, 2012.
- [9] F. Kamiran and T. Calders, "Data preprocessing techniques for classification without discrimination," *Knowledge and Information Systems*, pp. 1–33, 2012.
- [10] I. Zliobaite, F. Kamiran, and T. Calders, "Handling conditional discrimination," in *Proc. of IEEE Int. Conf. on Data Mining (ICDM'11)*, 2011, pp. 992–1001.
- [11] F. Kamiran, I. Zliobaite, and T. Calders, "Quantifying explainable discrimination and removing illegal discrimination in automated decision making," *Knowledge and Information Systems*, pp. (Accepted) 1–33, 2013.
- [12] D. Pedreschi, S. Ruggieri, and F. Turini, "Measuring discrimination in socially-sensitive decision records," in *Proc. of the SIAM International Conference on Data Mining (SDM'09)*, 2009, pp. 581–592.
- [13] S. Ruggieri, D. Pedreschi, and F. Turini, "Data mining for discrimination discovery," *ACM Transactions on Knowledge Discovery from Data (TKDD)*, vol. 4, no. 2, p. 9, 2010.
- [14] B. Luong, S. Ruggieri, and F. Turini, "k-nn as an implementation of situation testing for discrimination discovery and prevention," in *Proc. of the 17th ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM, 2011, pp. 502–510.
- [15] S. Hajian, J. Domingo-Ferrer, and A. Martínez-Ballesté, "Rule protection for indirect discrimination prevention in data mining," *Modeling Decision for Artificial Intelligence*, pp. 211–222, 2011.
- [16] S. Hajian, J. Domingo-Ferrer, and A. Martinez-Balleste, "Discrimination prevention in data mining for intrusion and crime detection," in *IEEE Symposium on Computational Intelligence in Cyber Security (CICS)*. IEEE, 2011, pp. 47–54.
- [17] S. Hajian and J. Domingo-Ferrer, "A methodology for direct and indirect discrimination prevention in data mining," *IEEE Transactions on Knowledge and Data Engineering*, vol. accepted, 2012.
- [18] F. Kamiran and T. Calders, "Classifying without discriminating," in *Proc. of the 2nd Int. Conf. on Computer, Control and Communication (IC4)*, 2009, pp. 1–6.
- [19] T. Calders, F. Kamiran, and M. Pechenizkiy, "Building classifiers with independency constraints," in *IEEE ICDM Workshop on Domain Driven Data Mining (DDDM'09)*, 2009, pp. 13–18.
- [20] F. Kamiran and T. Calders, "Classification with no discrimination by preferential sampling," in *Proc. of the 19th Ann. Machine Learning Conf. of Belgium and the Netherlands (BENELEARN'10)*, 2010, pp. 1–6.
- [21] F. Kamiran, T. Calders, and M. Pechenizkiy, "Discrimination aware decision tree learning," in *Proc. of IEEE Int. Conf. on Data Mining (ICDM)*, 2010, pp. 869–874.
- [22] T. Calders and S. Verwer, "Three naive Bayes approaches for discrimination-free classification," *Data Mining and Knowledge Discovery*, vol. 21, no. 2, pp. 277–292, 2010.
- [23] T. Kamishima, S. Akaho, and J. Sakuma, "Fairness-aware learning through regularization approach," in *2011 IEEE 11th International Conference on Data Mining Workshops (ICDMW)*, 2011, pp. 643–650.