

Measuring Discrimination in Socially-Sensitive Decision Records

a potentially discriminatory (PD) rule  $A, B \rightarrow C$ , where  $A$  is a non- empty PD itemset and  $B$  is a PND itemset.

- 1.Ratio Measures
- 2.Difference Measures
- 3.Discriminatory Classification Rules

1.define a family of formal measures of discrimination for classification rules

2.a notion of statistical significance

3.combine the discrimination measures with association rule mining

4.for rule-based classifiers, we propose a discrimination correction based on the measures

5.Experiment the theoretical definitions and results on the publicly available German credit granting dataset and on the CPAR rule-based classifier

a-discriminatory  
a-protective

1.Extraction and a-discrimination checking of PD classification rules  
2.potentially non-discriminatory (PND) and an inference model is proposed exploiting background knowledge with respect to the elift() measure

correction PD,PND

High values of the discrimination measures will occur when people in one or more of those categories is denied credit more often than people not in those categories.

unveiling all discriminatory decision patterns hidden in the historical data, combining discrimination analysis with association rule mining

unveiling discrimination in classifiers that learn over training data biased by discriminatory decisions

Our approach is validated on the German credit dataset and on the CPAR classifier.

a systematic framework for measuring discrimination