# Association Analysis of Semi-structured Data for Discrimination Discovery in Business

**Luong Thanh Binh[1], Franco Turini[2]**
[1]IMT Institute for Advanced Studies Lucca, Italy
[2]Department of Informatics, University of Pisa, Italy

**Abstract -** *Data mining techniques have taken a critical role in life in numerous domains such as consumer analytics, finance, banking, medicine, biology, and astronomy… Recently, data mining techniques have found their application also in discovering illegal discriminatory treatment on the bases of sensitive attributes such as race, color, religion, nationality, gender, age… In this paper, we propose a framework to solve the discrimination matter in the context of semi-structured business data, and in particular in the calculation of taxes for imported goods. This framework is able to discover possibly discriminatory relations among data by finding discriminatory association rules with the support of a common sense knowledge base and text mining techniques. The framework has applied to the problem of HTS (US Harmonized Tariff Schedule) showing some satisfactory results.*

**Keywords:** association analysis, data mining, discrimination, HTS, text mining

## 1   Introduction

In human beings societies, discriminatory behaviour happens in situations when members of a minority are treated unequally or worse than the ones of majority group(s) without regard to individual merits. This problem has been surveyed for a long time by economists, sociologists and legislators. We can find studies of unfair treatment for instance in racial profiling and redlining [16], [17], personnel selection [7], [16], and mortgage granting [10]…

Typical cause of discrimination is sensitive attributes such as gender, age, racial, colors… For instance, female employees were often paid less or had narrowed opportunities to promotion or bonus, e.g. unjust treatments in the Merrill Lynch case in 2004 [15], Wal-mart case in 2007 [11]. A recent case attracted attention of the community especially of apparel importers in the United States market: Totes-Isotoner sued US government for violating their right to equal protection under the law. They complained that there was discrimination in tariff rates for certain men's and women's gloves on the basis of gender. For example, the Harmonized Tariff Schedule (HTS) imposed a tax of 14 percent on men's gloves whereas 12.6 percent on gloves for others. Yet, this complaint was dismissed by the U.S. Court of International Trade since the Court found that the argument and evidence Totes alleged were inadequate and not persuasive. Above incidents and others motivate an interest for implementing knowledge discovery models or more particularly association analysis to extract possibly illegal relations between the discriminatory treatments and sensitive attributes e.g. gender, age, color... However, data are often not well-structured or clean due to noises, incompleteness or inconsistency which are caused by a variety of information sources, ambiguity by synonyms or negative meanings, etc or mistakes of users… For example, a database of merchandise may provide varied descriptions for the same item that differ in producers or in materials. Direct mining on such data is impossible since the provided information is understandable for human beings but it is not in the form accepted by mining algorithms.

In this research, we propose a framework to discover discrimination in the context of semi-structured business data, and, in particular in the calculation of taxes for imported goods. The scope of "semi-structured business data" used herein covers all business data which are not completely well-structured as pairs of attribute/value such as tables in relational database that also contain textual parts which can be split into smaller components. The framework can find possibly discriminatory relations in the form of discriminatory association rules based on a common sense knowledge base and text mining techniques like syntax analysis… It is then applied to the HTS data to find some discriminatory association rules between tariffs and some attributes including gender.

The problem of discriminatory treatments in sensitively social applications data mining has been studied much recently was first mentioned in [3] about classifiers that might execute racial discrimination. We foresee three non mutually-exclusive strategies towards discrimination prevention. The first one is to adapt the preprocessing approaches of data sanitization [6, 21] and hierarchy-based generalization [18], [22]. Along this line, [8] adopts a controlled distortion of the training set. The second one is to post-process the produced classification model. Along this line, [12], [13, [14] propose a confidence-

altering approach for classification rules inferred by the CPAR algorithm of [23]. The third one is to modify the classification learning algorithm by integrating some specialized steps of discrimination analysis and possibly discrimination measures calculation within it. This approach inspires our research.

The remainder of this paper is structured as follows. Section 2 contains details of the HTS case. Section 3 defines the problem and represents the discriminatory association analysis framework on semi-structured business data for discrimination discovery. Experimental results are shown in section 4. Finally, discussion on the approach and future development will be presented in section 5.

## 2 HTS problem

The Harmonized Tariff Schedule-HTS provides a tariff classification system for merchandise imported in US, including nomenclatures (names), descriptions for goods, and formulae for calculating tariff rates consisting of ad valorem, specific and estimated ad valorem equivalent (AVE) tariffs. Even though this system is carefully built and updated conforming to the law, Michael Barbaro has uncovered an oddity in the tariff system [1], [2]: duties on men's and women's garments are different for no apparent reason as the example shown in Figure 1. He calculated that the government imposes a 14 percent tariff on women's, but only 9 percent on men's on overall. And according to evaluation of Matt Gersper and Tom Gould [5], US importers have overpaid more than 1.3 billion dollars for discriminatory duties. Therefore, some clothing importers backed a lawsuit that would force the tariffs on similar items to be equalized, on the ground of gender like Totes Isotones, Steve Madden, Asics and Columbia Sportswear [1], [2], [5]. From the point of view of knowledge discovery, there is a possibility for applying data mining techniques to find possible relations between attribute(s) of apparel (gender, kind of apparel, materials…) and discrimination in tariff.

## 3 Finding potential discrimination by association analysis

### 3.1 Problem formulation

We use a convention that an *a*-item is an expression $a = v$, where *a* is an attribute and $v \in dom(a)$. An item is any *a*-item. Let *I* be a database of itemsets referring to a set of attributes: $A = \{a_1, ..., a_n\}$. Among these attributes there is often a special attribute $a_{tg}$ used for classifying itemsets which is one of the purposes of the application such as the cost of service, level of health care service… Its value depends on other attributes but this dependence cannot be represented as a simple function. Its domain is: $dom(a_{tg}) = \{c_1, ... , c_m\}$. Besides, there is an attribute, called description *des*, which contains a number of elementary attributes and is expressed in a quite free form. For instance, an item of a merchandise



Fig. 1. Discriminatory tariff in HTS

database: "*Men's or boys' suits, of synthetic fibers, not knitted or crocheted, containing 36 percent or more by weight of wool or fine animal hair*" which can be mapped into simpler sub-items:

- *Gender* = male
- *Name of goods* = {suits}
- *Form of production* = {not knitted, not crocheted}
- *Quantity* = 36%
- *Materials* = {synthetic fibers, wool, fine animal hair}

This kind of data is still useful since it contains precious information. However, directly digging knowledge from this item is impossible since the information is not well-organized. The matter turns simple if the description can be extracted as pairs of attribute/value which are items.

Summing up, we have to map these semi-structured items onto a set of well-defined transactions. This process may need help from some form of text mining. Along this line, *des* is split to sub-items:

$$\{a_{il},..., a_{ip}\}, dom(des) = \{d_1,...,d_m\}$$

Those sub-items can be divided into sensitive group $A_1$ such as the *Gender* item and non-sensitive group $A_2$ such as *Name of goods*, *Materials* items. We define *S*, a subset of $A \cup A_1$ as the *Sensitive Attribute Set*, which is a set of potentially discriminatory-PD attributes (due to categorization of [21,22,23]) such as gender, race… Our task is to verify the effect of these PD attributes *S* on the target attribute $a_{tg}$. The other set $\overline{S}$ is composed of potentially non-discriminatory-PND attributes, building the background information. From now on, the transactions of the database will have the form of a three-tuple:

$$(\{u_i\}, \{v_i\}, c_i) \qquad (1)$$

where $\{u_i\}$ is the set of non-sensitive items, $\{v_i\}$ is the set of sensitive items, and $c_i$ is the target item. While inspecting the database, if the variances in the target attributes $a_{tg}$ considerably depend on sensitive attributes *S*, it can be said that this correlation is discriminatory, or represented in association rules:

$$\{u_k\}, \{v\} \to c, \qquad (2)$$
$$\{u_k\}, \{v'\} \to c'$$

where:

- $\{u_k\}$ are items referring to $\overline{S}$ that defines the context.
- $\{v\}$, $\{v'\}$ refer to the same sensitive attributes with different values.
- $c, c'$ are target items with different values.

These association rules can then be used as evidence of discrimination. Yet, the results of (2) are hard to achieve and sometimes can not provide much in details about discrimination in the system due to the following possible reasons:

- The target item is not binary. For example, the "*rate type*" item in the HTS data has several values. It is difficult to clarify in which case(s) discrimination really happens for the target attribute and which attribute(s) has the negative effect causing the unequal target attribute.
- It is difficult to construct a context built of the PND items and some sub-items of the *des* to compare the effect of different values of sensitive attributes through the association analysis, especially when the database is very large and/or the target attribute is not binary.
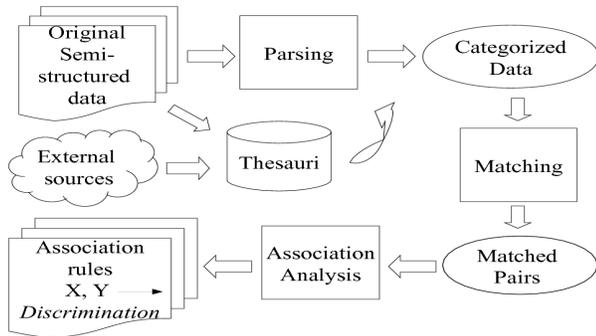


Fig. 2.  Framework of finding discriminatory base by association analysis

Therefore, the traditional rule mining process should be modified when applied to discrimination discovery process. This paper proposes a 3-step framework as presented in Figure 2 which is able to analyze semi-structured business data to find possible discriminatory associations. It is especially applicable to calculate discrimination in taxes for imported goods.

## 3.2    Proposed framework

### 3.2.1    Parsing

Figure 3a represents a typical case of not-well-structured raw data in which the "*description*" field can be still divided into more detailed subfields (sub-items) as shown in Figure 3c with categories of "*gender*", "*name of goods*"… In our framework, the original data will be firstly combined with external sources of information such as WordNet's thesaurus, Roget's thesaurus to build the common sense knowledge base. It is structured as a hierarchy of components (if any) and attributes of the sub-items extracted from *des*. From the

| ID | Description | Quantity2 | Code | Rate | Type |
|---|---|---|---|---|---|
| 62031210 | *Men's or boys' suits, of synthetic fibers, not knitted or crocheted, containing 36 percent or more by weight of wool or fine animal hair* | Kg | B | 17.5% | 7 |

(a)

| Gender | Name of goods | Materials | | | Forms | | Quantity | Limitation | Negative |
|---|---|---|---|---|---|---|---|---|---|
| male | suits | synthetic fibers | wool | fine animal hair | crocheted | knitted | 36% | more | not |

(b)

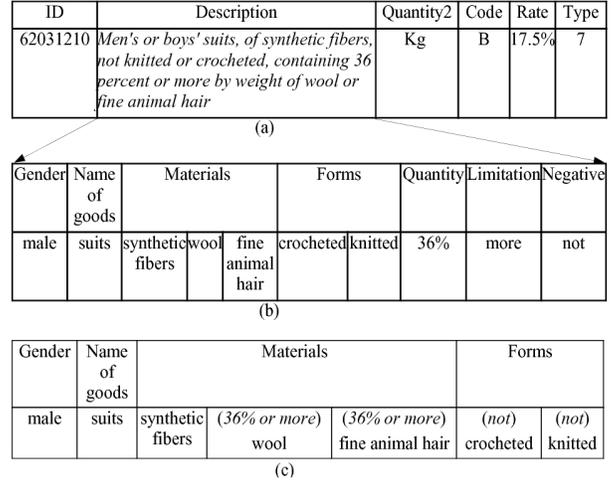| Gender | Name of goods | Materials | | | Forms | |
|---|---|---|---|---|---|---|
| male | suits | synthetic fibers | (*36% or more*) wool | (*36% or more*) fine animal hair | (*not*) crocheted | (*not*) knitted |

(c)

Fig. 3.  Example of transformation from semi-structured data to well-structured data.

structure of the sub-itemset, thesauri are built to help classify words into categories and support the extraction of sub-items from the *des*. There are two groups:

- *Form thesauri*: serve the purpose of normalizing all different words, e.g. thesaurus of synonyms, abbreviations. For instance, (*not*, *excluding*) is an entry of the negative thesaurus. If any word extracted from the description matches some relevant word of an entry of the thesaurus, e.g. "*excluding*", it will be replaced by its corresponding term, and in this case, that term is "*not*".
- *Category thesauri*: clarify categories (attributes) of terms. For example, thesaurus of materials of goods: *silk, fur*…

Therefore, redundant information is removed, such as stop words, for example, "to", "the", "of", "and"… or ambiguous words which can not be referred to any category (attribute). And all meaningful terms will be categorized into a particular thesaurus of the Category group. In other words, the original description item is transformed into pairs of attribute/value as well-defined sub-items (called clothing sub-itemset) as in Figure 3b.

Finally, possible semantic relationships among sub-items are searched with the help of formalized rules. In a database of clothing merchandise, for instance, it is observed that: *An amount is often related to some material of the clothing items, e.g. 15% of silk…* These rules may be formalized from the sub-items through syntactic analysis techniques borrowed from natural language processing such as: *Gender-name-materials 1-parts-materials 2 - forms of production.* The extracted categorized terms of each description are sequentially checked through rules to find their relationships. The above rule, for example, says that the *materials 1* belonging to the object level (*name*), have level 1; whereas the *materials 2* belonging to the component of object (*parts* item), have level 2 but not others if the checked sub-itemset satisfies the order of the rule. In the end, the returned result is a well-defined structure with semantics constraints as shown in Figure 3c.

### 3.2.2 Matching

Because the main purpose is checking discriminatory behaviors possibly caused by sensitive attributes, another task is required. That is finding itemsets having the same background information but some different (or opposite) sensitive attributes. Only when matching itemsets are found, mining is deployed on these items to check if their target items are different. In particular, for each itemset:

- Find itemsets that are different in the sensitive attributes $S$ (PD attributes)
- For each itemset found in the previous step, comparing their PND attribute. If there is any itemset which is completely the same in all non-sensitive attributes, it is said that they match each other.

For example, the data in Figure 3c will match the following original data: "*Women's or girls' suits, not knitted or crocheted, of synthetic fibers, containing 36 percent or more of wool or fine animal hair*" which has nearly the same sub-items as the itemset presented in Figure 3a except the gender sub-attribute (female vs male) as illustrated in Figure 4. We propose a simple algorithm as following but a convenient means for finding matching itemsets as specified in the two above steps.

```
   Algorithm MatchFinding(itemsetA)
{
  CandidateList = FindAllOppositeItemsets();
    //find all itemsets which are unequal in
    //values of PD attributes
    ForEach itemset in CandidateList
    {
        isMatch = true;
        ForEach attribute in AttributeSet
        {
            rate=
            CompareValuesOfTheCurrentAttribute(itemset,
            itemsetA);
            //calculate the level of matching
            //syntax rules and some heuristics can be
            //used
            If rate < 1
            //complete matching
            Then
              begin
                    isMatch = false;
                    break;
              end
        }
        If (isMatch)
        Then
            Add itemset to MatchedList;
    }
If IsEmpty(MatchedList)
Then
        return null;
Else
       return MatchedList;
}
```

DEFINITION 1. *A discriminatory indicator $\theta_i$ for the target item caused by the PD attribute $s_i$ on a given context $\{u_k\}$ is defined as following:*
*Given two tuples:*
$$(\{u_k\},\{v\},c), (\{u_k\},\{v'\},c')$$
*The value of $\theta_i$ is defined as:*

$$\theta_i = \begin{cases} 1 & if \quad c \neq c' \\ 0 & if \quad c = c' \end{cases}$$

where:

- $\{u_k\} \subseteq \overline{S}$ defines a context
- $\{v\},\{v'\}$: sensitive attributes with different values.
- $c, c'$: target attribute $a_{tg}$ with different values.

ID = 62031210

| Gender | Name of goods | Materials | | | | Forms | |
|---|---|---|---|---|---|---|---|
| male | suits | synthetic fibers | *(36% or more)* wool | *(36% or more)* fine animal hair | | *(not)* crocheted | *(not)* knitted |

ID = 62041310

| Gender | Name of goods | Materials | | | | Forms | |
|---|---|---|---|---|---|---|---|
| female | suits | synthetic fibers | *(36% or more)* wool | *(36% or more)* fine animal hair | | *(not)* crocheted | *(not)* knitted |

Fig. 4. Matching itemset

The above definition can be explained as follows: For any itemset, the discriminatory indicator is activated when another itemset which is different in values of PD attribute(s) and target attribute on the same context is found. If only one discriminatory indicator is set, it is hard to precisely identify the effect of each PD attribute (if there are more than one PD attributes) on the target attribute. Thus, each PD attribute has its own discriminatory indicator. The individual effect of each PD attribute on the target attribute is then calculated by the support of that discriminatory indicator on the whole database. In particular, the following situations may happen:

- there are several itemsets satisfying the triple:
$$(\{u_k\}, v_i, \theta_i = 1)$$
In this case, the following association is generated:
$$\{u_k\}, v_i \rightarrow Discrimination \qquad (3)$$
where $\{u_k\}$ is a set of PND items forming the context, whereas $v_i$ is a single PD item. When $\{u_k\}$ is empty, it means that the discrimination caused by the PD attribute does not depend on a specific case, which is the simplest case.

- there are several itemsets satisfying the triple:
$$\{u_k\}, \{v_i\}, \{\theta_i = 1\}$$
In this case the following associations are generated, one for each PD item in the set $\{v_i\}$:
$$\{u_k\}, \{v_i\} \rightarrow Discrimination \qquad (4)$$
where $\{v_i\}$ is a set of more than 1 PD items.

We call (3) and (4) discriminatory association rules. In the two cases, the multiple-valued target item is replaced by a binary item, the discriminatory indicator ($\theta$ is 0 or 1). This result is extremely important since it helps to directly clarify whether the association between the PD attribute(s) and the target attribute is discriminatory commonly by using association analysis techniques.

### 3.2.3 Association analysis

Finally, the discriminatory association mining is implemented to compute the confidence of each of the possible discriminatory associations mentioned in the previous step. Any association rules mining algorithm can be applied, e.g., Apriori, FPGrowth... Given a user-specified $\alpha$-threshold (for the minimum support or minimum confidence), if any discriminatory association rule is retrieved, it will: i) prove that there is unjust discrimination in the given system on the bases of PD attributes and ii) reveal which attribute(s) causes that discriminatory treatment. Moreover, among PND attributes causing discrimination, it might be found that some of them are not reasonable. If no convincing argument is given for the negative effect of such attributes in the discrimination matter, they should be considered wrong/illegal and removed in the decision-making processes for business. Also, the achieved results can then be used in reasoners for further analysis as mentioned in [14]. Examples of possibly archived association rules will be presented in the Section 4.

## 3.3 Discrimination measures

A general principle mentioned in [9] is to consider group under-representation as a quantitative measure of the qualitative requirement that people in a group are treated "less favourably" (European Union Legislation 2009 [4]); UK Legislation (2009) [19] than others, or such that "a higher proportion of people without the attribute comply or are able to comply" (Australian Legislation 2009) to a qualifying criterion. For the purpose of quantifying the level of discrimination, we also propose some measures applied on the discriminatory association rules. Whereas [12], [13] use discrimination measure as a way to achieve the expected level of discrimination in the system, our relatively similar measures are used to obtain a precise vision of the actual discrimination in the decision system in which different PD attributes have different discriminatory effects and even their mutual action on each other.

DEFINITION 2. *Let* $\bar{s}, s \rightarrow \theta_i = 1$ *(discrimination) and* $\bar{s} \rightarrow \theta_i = 1$ *be association rules with confidences correspondingly* $\rho$ *and* $\varphi$. *We have:*
- *Absolute difference:*     $abs\_lift = \varphi - \rho$     (5)
- *Relative difference:*     $rev\_lift = (\varphi - \rho)/\varphi$     (6)
- *Ratio gain of difference:*  $rg\_lift = \varphi/\rho$     (7)

where:
- $\bar{s} \subseteq \bar{S}$ are PND items
- $s \subset S$ is a set of PD items

Many legislative documents regulate a threshold which is the acceptable maximum difference. For example, in the UK, a difference of 5% in confidence between female and male treatment is assumed by courts to be significant of discrimination against women. The absolute difference reveals whether the discrimination caused by $s_i$ will be gained or decreased when some extra PD attribute is added to the context which means that the discrimination will be worse when these two PD attribute go in tandem. For the same meanings, the relative difference computes the relative reduction of equality with the presence of another PD attribute besides the context. The US Legislation [20], for instance, states that "a selection rate for any race, sex, or ethnic group which is less than four-fifths (80%) of the rate of the highest rate will be generally regarded as evidence of adverse impact." The ratio gain of difference tells how much the degree of difference will be amplified if more sensitive part is added to the context.

When $|S_i| = 1$ and both the following two discriminatory association rules are achieved:

$$\bar{s}, s_i = v_1 \rightarrow \theta_i = 1,$$
$$\bar{s}, s_i = v_2 \rightarrow \theta_i = 1$$

where $v_1$ and $v_2$ are two opposite values of the PD attribute $s_i$, and the corresponding confidences are $\rho$ and $1 - \rho$. And through the matching step, $\theta_i$ is specified as the possible benefit such as a favored action like accepting a loan. We can compare by then the disparity between different values of the sensitive items such as what is the difference of income between male and female:

- Absolute opposite difference:
  $$abs\_lift = 1 - 2\rho \qquad (8)$$
- Relative opposite difference:
  $$rev\_lift = (1 - 2\rho)/(1 - \rho) \qquad (9)$$
- Ratio gain opposite difference:
  $$rg\_lift = (1 - 2\rho)/(1 - \rho) \qquad (10)$$

## 4 Experimental results in HTS problem

We applied the framework to the mentioned HTS problem. Our experimental goals are:
- Finding possible relationships between attributes of apparel and discriminatory tariff in the form of discriminatory association rules.
- Then, verifying the charge that the discrimination in tariff for certain apparels is based on the gender attribute.

To realize this component, we have built a common sense knowledge base as illustrated in Figure 5 from the original data as presented in Figure 3a. Additionally, thesauri of synonyms, abbreviations, negative meanings are also built.
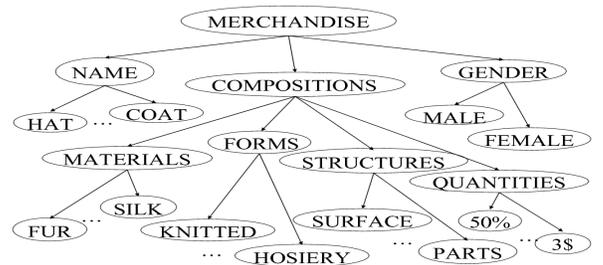


Fig. 5. Common sense knowledge base for the HTS problem

Data are cleaned and standardized as in Figure 3b and hierarchically structured as in Figure 3c: for each apparel, information is categorized into the name of the apparel, its gender (which group of sex it is produced for), material, form of production, quantities of each material, and quantities for its components (structure) as built in the common sense knowledge base. It is easily seen that the tariff attribute is the target attribute, the gender attribute is the PD attributes and all the others form the background (the context). After the matching, the association analysis is performed by the Apriori algorithm to dig out possibly discriminatory relations among attributes. When the minimum support is set equal to 5 percent, we found the following maximal length discriminatory rules:

```
form = "knitted", "crocheted", material = "fine
animal hair", "wool"➜discrimination (conf = 73.53%)
```

```
form = "not_knitted", "not_crocheted", material =
"fine animal hair", "wool"➜discrimination
(conf=50%)
```

which means that on the context specified in the antecedent of the rules, the tariff will be different if gender is different from merchandise to merchandise with the confidences correspondingly are 73.53% and 50%.

When reducing the minimum support to 1 percent, the result is more surprising with more specific attributes:

```
form = "not_knitted", "not_crocheted", material =
"silk", "textile", "not man-made fiber", "not wool",
"not cotton", quantities = "70%", MFN[1]_rate_type = 7
➜discrimination (conf =  66.67%)
```

If the minimum support is set relatively high, i.e, greater than 5 percent, it will affect the number of selected candidates; therefore, decrease the number of rules as well as the number of items in the rules. For instance, when the minimum support is set to 10 percent, the achieved maximal length rule is relatively sparse:

```
MFN_rate_type  =  7,  form  =  "not_knitted",
"not_crocheted" ➜ discrimination (conf = 29.58%)
```

Besides, it shows that the ratio of apparels for male which are imposed a lower tariff than for female is often much higher than that one for female as illustrated in Figure 6 and Figure 7. It is immediate to see that, for a common kind of apparels, male often receives lower tariff than female. Moreover, the difference of male/female rate is not high if female is imposed lower tariff which is opposite to the situation when male retrieves lower tariff. In general, it can be said that apparels for male are favored in the tariff. This result comparatively fits researches of [1], [2], [5].

---

[1] MFN-Most Favored Nation implies the normal status in international trade. For a given merchandise, a MFN rate type is 7 means the duty for that merchandise is an *ad valorem* rate.
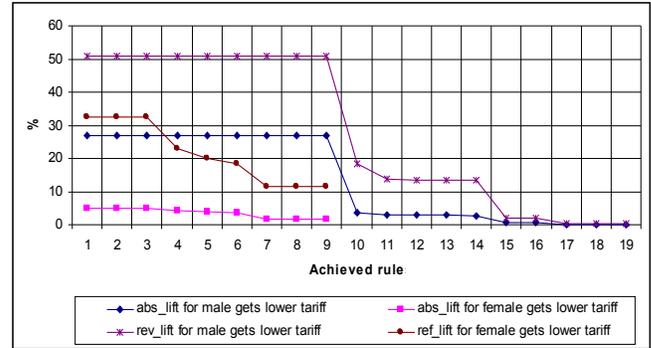


Fig. 6.  Distribution of discrimination in tariff between male and female in HTS w.r.t. *abs_lift* and *rev_lift*.



Fig. 7.  Distribution of discrimination in tariff between male and female in HTS w.r.t. *rg_lift*.

Looking back at the initial targets of our experiment, the following results are achieved:

- Finding relationships between attributes of apparel and discriminatory tariff in the form of discriminatory association rules. Through the achieved results, verify the evaluation that the discrimination in tariff for certain apparels is considerably affected on the basis of gender.

- Other likely hidden causes of this matter are the MFN rate type 7, "knitted", "crocheted", "not knitted", "not crocheted" for the form of production, and "silk", "textile", "wool" for materials with a relatively high confidence for different values of minimum support. This result strongly supports the findings of [12], [13] when PND attributes may have certain effect on discrimination. However, a new question arises: why the above three attributes have a strong connection to the disparity in the HTS? Is there any explanation for it?

- Verifying the potential of the framework. The proposed framework shows the ability to transform the ambiguous original data to the well categorized data which are necessary for mining discrimination in semi-structured data in business field.

And the framework is able to discover the discrimination in business data and present this relation in the form of discriminatory association rules from which we can draw conclusion of possible bases leading to discrimination.

# 5 Conclusions

Since data mining is significantly useful in analysis, we can expect that this field is capable of supporting the discrimination discovery process, especially reveal potentially unfair bases of disparate decisions. In this paper, we have formalized a general framework to discover discrimination in business data and, in particular in the calculation of discrimination in taxes for imported goods by means of association analysis. Possible discrimination is extracted in form of discriminatory association rules with support from a common knowledge base and some techniques of text mining such as cleaning. Experimental results on the HTS problem show the potential of the proposed framework. In the future, we will follow the approach of using association analysis methodologies to prevent discrimination caused by both multiple potentially discriminatory attributes and potentially non-discriminatory attributes. And if possible, de-discriminate those disparities by developing new measures and methods of measurement.

# 6 References

[1]  Barbaro, M., "In apparel, all tariffs aren't created equal". *New York Times*. Apr 28th 2007

[2]  Barbaro, M., "Clothing makers allege sex discrimination in U.S. tariffs". *International Herald Tribune*. Apr 29th 2007

[3]  Clifton, C., "Privacy preserving data mining: How do we mine data when we aren't allowed to see it?" *KDD 2003*. 2003

[4]  European Union Legislation 2009. (a) Racial Equality Directive, (b) Employment Equality Directive. http://www.ec.europa.eu/employment_social/fundamental_rights

[5]  Gersper, M., T. Gould, "Gender and Age discrimination costs U.S. importers billions".  US Custom House guide. Feb 4th 2008.

[6]  Hintoglu, A.A., A. Inan, Y. Saygin, "Suppressing data sets to prevent discovery of association rules". ICDM 2005, IEEE Computer Society, pp. 645-648. 2005

[7]  Holzer, H., S. Raphael, M. Stoll, "Black job applicants and the hiring officer's race". Industrial and Labor Relations Review, 57(2), pp.267–287. 2004

[8]  Katz, M.J., "The Economics of Discrimination: The three Fallacies of Croson". Yale Law Journal 100: 1033. 1991

[9]  Knopff, R., "On proving discrimination: Statistical methods and unfolding policy logics". Canadian Public Policy 12, pp.567-583. 1986

[10] LaCour-Little, M., "Discrimination in mortgage lending: A critical review of the literature". J. of Real Estate Literature, 7, pp.15–50. 1999

[11] Millman, J., "Largest discrimination case in history: Wal-mart case's appeal denied". Diversity Inc Magazine. Feb 7th 2007

[12] Pedreschi, D., S. Ruggieri, F. Turini, "Discrimination – aware Data Mining". KDD'08, pp.560-568. ACM, Aug 2008.

[13] Pedreschi, D., S. Ruggieri, F. Turini, "Integrating induction and deduction for Finding Evidence of Discrimination".  ICAIL  2009, pp.157-166. ACM, Jun 2009

[14] Pedreschi, D., S. Ruggieri, F. Turini, "DCUBE: Discrimination discovery in Databases".SIGMOD 2010, to be published.

[15] Sachdev, A., "Merrill Lynch to pay $2.2 million in damage in sexual discrimination case". *Chicago Tribune newspaper*. Apr 24th 2004.

[16] Squires, D., S. O'Connor, *Colour and Money*. Albany: SUNY Press. 2001

[17] Squires, D., "Racial profiling, insurance style: Insurance redlining and the uneven development of metropolitan areas". J. of Urban Affairs, 25(4), pp.391–410. 2003

[18] Sweeney L., "Achieving k-anonymity privacy protection using generalization and suppression". International Journal on Uncertainty and Fuzziness in Knowledge-Based Systems 10(5), pp.571-588. 2002

[19] UK Legislation, 2009. (a) Sex Discrimination Act, (b) Race Relation Act. http://www.statuelaw.gov.uk

[20] US Federal Legislation 2009. (a) Equal Credit Opportunity Act, (b) Fair Housing Discrimination Act, (c) Intentional Employment Discrimination, (d) Equal Pay Act, (e) Pregnancy Discrimination Act, (f) Civil Right Act. http://www.usdoj.gov

[21] Verykios V.S., A.K. Elmagarmid, E. Bertino, Y. Saygin, E. Dasseni, "Association rule hiding". IEEE Transactions on Knowledge and Data Engineering 16(4), pp.434-447. 2004

[22] Wang K., B.C.M. Fung, P.S. Yu, "Template-based privacy preservation in classification problems". ICDM 2005, IEEE Computer Society, pp.466-473. 2005