

# Supervised Descriptive Rule Discovery: A Unifying Survey of Contrast Set, Emerging Pattern and Subgroup Mining



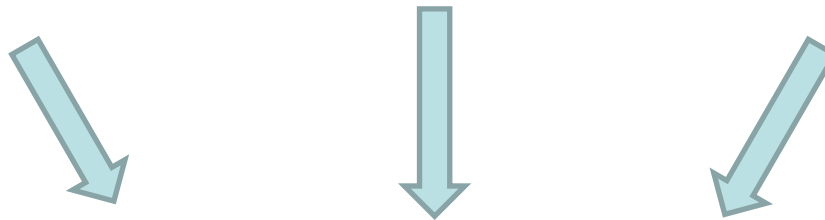
TECHNISCHE  
UNIVERSITÄT  
DARMSTADT

Petra Kralj Novak, Nada Lavrac, Geoffrey I. Webb

**C**ontrast  
**S**et  
**M**ining

**E**merging  
**P**attern  
**M**ining

**S**ubgroup  
**D**iscovery



## Supervised Descriptive Rule Discovery

1. Einführung
2. Descriptive Rule Discovery Verfahren
  - 2.1 Contrast Set Mining
  - 2.2 Emerging Pattern Mining
  - 2.3 Subgroup Discovery
3. Einheitliches Framework für Supervised Descriptive Rule Induction
4. Fazit

## Symbolic Data Analysis Techniques



### ■ Descriptive Induction

- Unlabeled Data
- Erkennen von Mustern

→ Contrast Set Mining  
Emerging Pattern Mining



### ■ Predictive Induction

- Labeled Data
- Klassifizierung von Beispielen

→ Subgroup Discovery

→ Hauptziel: menschlich interpretierbare Unterschiede zwischen Gruppen finden

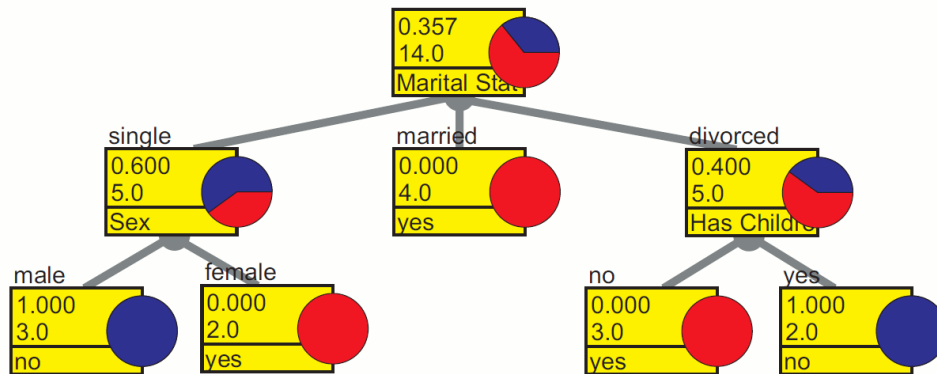
# Beispiel Datenbank



TECHNISCHE  
UNIVERSITÄT  
DARMSTADT

| Education  | Marital Status | Sex    | Has Children | Approved |
|------------|----------------|--------|--------------|----------|
| primary    | single         | male   | no           | no       |
| primary    | single         | male   | yes          | no       |
| primary    | married        | male   | no           | yes      |
| university | divorced       | female | no           | yes      |
| university | married        | female | yes          | yes      |
| secondary  | single         | male   | no           | no       |
| university | single         | female | no           | yes      |
| secondary  | divorced       | female | no           | yes      |
| secondary  | single         | female | yes          | yes      |
| secondary  | married        | male   | yes          | yes      |
| primary    | married        | female | no           | yes      |
| secondary  | divorced       | male   | yes          | no       |
| university | divorced       | female | yes          | no       |
| secondary  | divorced       | male   | no           | yes      |

# Predictive Induction



- C4.5 (1993)
- Hoher Wert auf Vollständigkeit

Sex = female  $\rightarrow$  Approved = yes

MaritalStatus = single AND Sex = male  $\rightarrow$  Approved = no

MaritalStatus = married  $\rightarrow$  Approved = yes

MaritalStatus = divorced AND HasChildren = yes  $\rightarrow$  Approved = no

MaritalStatus = divorced AND HasChildren = no  $\rightarrow$  Approved = yes

MaritalStatus = single AND Sex = male  $\rightarrow$  Approved = no

Sex = male  $\rightarrow$  Approved = no

Sex = female  $\rightarrow$  Approved = yes

MaritalStatus = married  $\rightarrow$  Approved = yes

MaritalStatus = divorced AND HasChildren = yes  $\rightarrow$  Approved = no

MaritalStatus = single  $\rightarrow$  Approved = no

- Class Attribut entspricht normalem Attribut
  - Regeln redundant
  - Regeln fehlen
- }
- Aufgabenstellung
  - Algorithmus
  - Constraints des Benutzers



1. Einführung
2. Descriptive Rule Discovery Verfahren
  - 2.1 Contrast Set Mining
  - 2.2 Emerging Pattern Mining
  - 2.3 Subgroup Discovery
3. Einheitliches Framework für Supervised Descriptive Rule Induction
4. Fazit

## 2.1 Contrast Set Mining

- **Ziel:** Finde ContrastSets mit denen eine Gruppe von anderen unterschieden werden kann

$$\textit{ContrastSet}_i \rightarrow G_i \quad \textit{ContrastSet}_j \rightarrow G_j \quad \dots \dots$$

- **Qualitätsbewertung:**

$$\textit{SuppDiff}(X, G_i, G_j) = |\textit{support}(X, G_i) - \textit{support}(X, G_j)| \geq \delta$$

## 2.2 Emerging Pattern Mining

- **Ziel:** ItemSets deren Support signifikant von einem Datenset zu einem anderen steigt

$$ItemSet_1 \rightarrow D_1 \quad ItemSet_2 \rightarrow D_2$$

- aufkommende Trends in zeitlich gemessenen Datenbanken oder differenzierende Merkmale zwischen Klassen von Daten finden
- **Qualitätsbewertung:**

$$GrowthRate(ItemSet, D_1, D_2) = \frac{support(ItemSet, D_1)}{support(ItemSet, D_2)}$$

## 2.3 Subgroup Discovery

- **Ziel:** finde UG, die statistisch gesehen am interessantesten sind bezüglich einer Property of Interest

*SubgroupDescription*  $\rightarrow$  *Class*

- **Qualitätsbewertung:**

$$WRAcc(X, C) = P(X) \cdot (P(Y|X) - P(Y))$$

$$q_g(X, C) = \frac{p}{n+g}$$

- Zielt auf ein gutes Maß zwischen Rule Coverage und Precision ab



1. Einführung
2. Descriptive Rule Discovery Verfahren
  - 2.1 Contrast Set Mining
  - 2.2 Emerging Pattern Mining
  - 2.3 Subgroup Discovery
3. Einheitliches Framework für Supervised Descriptive Rule Induction
4. Fazit

### 3. Ein einheitliches Framework - Inhalt

- Terminologie
- Task Definitions
- Rule Learning Heuristics
- Rule Selection Mechanism

## 3.1 Vereinheitlichung der Terminologie

### *Definition 1: Kompatibilität von Termen*

*Terme sind kompatibel, wenn sie ...*

- *... in äquivalente logische Ausdrücke umgeformt werden können*
- *... die selbe Bedeutung haben*

## 3.1 Vereinheitlichung der Terminologie

*Lemma 1: Die Terme in CSM, EPM und SD sind kompatibel*

| Contrast Set Mining                            | Emerging Pattern Mining                      | Subgroup Discovery               | Rule Learning                  |
|--|--|----------------------------------|--------------------------------|
| contrast set                                   | itemset                                      | subgroup description             | rule condition                 |
| groups $G_1, \dots, G_n$                       | data sets $D_1$ and $D_2$                    | class/property $C$               | class/concept $C_i$            |
| attribute-value pair                           | item   | logical (binary) feature         | condition                      |
| examples in groups<br>$G_1, \dots, G_n$        | transactions in data sets<br>$D_1$ and $D_2$ | examples of<br>$C$ and $\bar{C}$ | examples of<br>$C_1 \dots C_n$ |
| examples for which<br>the contrast set is true | transactions containing<br>the itemset       | subgroup of instances            | covered examples               |
| support of contrast set on $G_i$               | support of EP in data set $D_1$              | true positive rate               | true positive rate             |
| support of contrast set on $G_j$               | support of EP in data set $D_2$              | false positive rate              | false positive rate            |

## 3.2 Vereinheitlichung der Task Definitions

### *Definition 2: Kompatibilität von Task Definitionen*

*Definitionen sind kompatibel, wenn ...*

- *... ein Learning Task zu einem anderen umgewandelt werden kann, ohne das Lernziel zu verändern*

## 3.2 Task Definitions

### ■ CSM

- Finde ContrastSets mit denen man am besten Gruppen voneinander unterscheiden kann

### ■ EPM

- Finde ItemSets deren Support signifikant von einem DatenSet zu einem anderen steigt

### ■ SD

- Finde UG die so groß wie möglich sind und möglichst ausgefallene statistische Charakteristika besitzen bezüglich der Property of Interest

## 3.2 Vereinheitlichung der Task Definitions

*Lemma 2:*

*Die Definitionen von CSM, EPM und SD sind kompatibel*

$$\begin{aligned} EPM(D_1, D_2) &\Leftrightarrow CSM(G_i, G_j) \\ CSM(G_i, G_j) &\Leftrightarrow SD(G_i) \text{ und } SD(G_j) \end{aligned}$$

| Contrast Set Mining   | Emerging Pattern Mining   | Subgroup Discovery  | Rule Learning   |
|---|---|---|---|
| <b>Given</b><br>examples in $G_1$ vs. $G_j$<br>from $G_1, \dots, G_i$                     | <b>Given</b><br>transactions in $D_1$ and $D_2$<br>from $D_1$ and $D_2$           | <b>Given</b><br>in examples $C$<br>from $C$ and $\bar{C}$ | <b>Given</b><br>examples in $C_i$<br>from $C_1 \dots C_n$ |
| <b>Find</b><br>$ContrastSet_{i_k} \rightarrow G_i$<br>$ContrastSet_{j_l} \rightarrow G_j$ | <b>Find</b><br>$ItemSet_{1_k} \rightarrow D_1$<br>$ItemSet_{2_l} \rightarrow D_2$ | <b>Find</b><br>$SubgrDescr_k \rightarrow C$               | <b>Find</b><br>$\{RuleCond_{i_k} \rightarrow C_i\}$       |

## 3.3 Vereinheitlichung der Rule Learning Heuristics



TECHNISCHE  
UNIVERSITÄT  
DARMSTADT

### *Definition 3: Kompatibilität von Heuristiken*

*Heuristik  $h_1$  ist kompatibel mit Heuristik  $h_2$ , wenn ...*

- ...  $h_2$  von  $h_1$  abgeleitet werden kann*
- ... für zwei Regeln  $R$  und  $R'$  gilt:  $h_1(R) > h_1(R') \Leftrightarrow h_2(R) > h_2(R')$*

## 3.3 Rule Quality Measures

- CSM

$$\text{SuppDiff}(X, G_i, G_j) = |\text{support}(X, G_i) - \text{support}(X, G_j)| \geq \delta$$

- EPM

$$\text{GrowthRate}(X, D_1, D_2) = \frac{\text{support}(X, D_1)}{\text{support}(X, D_2)}$$

- SD

$$\text{WRAcc}(X, C) = P(X) \cdot (P(Y|X) - P(Y)) \quad q_g(X, C) = \frac{p}{n+g}$$

$$\text{support}(X, Y) = \frac{\text{count}(X, Y)}{|Y|} = \text{TPr}(X, Y)$$

## 3.3 Vereinheitlichung der Rule Learning Heuristiken

*Lemma 3a: Support Difference Heuristik von CSM und Weighted Relative Accuracy von SD sind kompatibel*

$$\begin{aligned} WRAcc(X, Y) &= \\ &= P(X) \cdot [P(Y|X) - P(Y)] = P(Y \cdot X) - P(Y) \cdot P(X) \\ &= P(Y \cdot X) - P(Y) \cdot [P(Y \cdot X) + P(\bar{Y} \cdot X)] \\ &= (1 - P(Y)) \cdot P(Y \cdot X) - P(Y) \cdot P(\bar{Y} \cdot X) \\ &= P(\bar{Y}) \cdot P(Y) \cdot P(X|Y) - P(Y) \cdot P(\bar{Y}) \cdot P(X|\bar{Y}) \\ &= P(\bar{Y}) \cdot P(Y) \cdot [P(X|Y) - P(X|\bar{Y})] \\ &= P(Y) \cdot P(\bar{Y}) \cdot [TPr(X, Y) - FPr(X, Y)] \end{aligned}$$



konstant

➡  $WRAcc(X, C) = WRAcc(X, G_1) = P(G_1) \cdot P(G_2) \cdot [support(X, G_1) - support(X, G_2)].$

## 3.3 Vereinheitlichung der Rule Learning Heuristiken

*Lemma 3b: Die Growth Rate Heuristik vom EPM und die Generalization Quotient Heuristik sind kompatibel*

$$\begin{aligned} \text{GrowthRate}(X, D_1, D_2) &= \frac{\text{support}(X, D_1)}{\text{support}(X, D_2)} \\ &= \frac{\text{count}(X, D_1)}{\text{count}(X, D_2)} \cdot \frac{|D_2|}{|D_1|} = \frac{p}{n} \cdot \frac{N}{P} \end{aligned}$$

↑ konstant

➔  $\text{GrowthRate}(X, C, \bar{C}) = q_0(X, C) \cdot \frac{N}{P}$

## 3.3 Vereinheitlichung der Rule Learning Heuristiken

*Lemma 3: Die Definitionen der CSM, EPM und SD Heuristiken sind paarweise kompatibel*

| Contrast Set Mining     | Emerging Pattern Mining   | Subgroup Discovery | Rule Learning   |
|-------------------------|---------------------------|--------------------|---|
| $SuppDiff(X, G_i, G_j)$ |                           | $WRAcc(X, C)$      | Piatetski-Shapiro heuristic leverage                      |
|                         | $GrowthRate(X, D_1, D_2)$ | $q_g(X, C)$        | odds ratio for $g = 0$<br>accuracy/precision, for $g = p$ |



Zielen alle auf ein gutes Maß zwischen Rule Coverage und Precision ab

## 3.4 Vergleich der Rule Selection Mechanismen



### ■ Statistic Test

- Contrast Set Mining
- Statistical significance pruning
- Entfernt alle ContrastSets, die Spezialisierungen von generelleren CS sind, wenn sie einen ähnlichen Support haben

### ■ Weighted Covering Approach

- Subgroup Discovery
- Benutzt gewichtete Version von  $q_g$  und WRAcc
- Nach jeder Iteration wird die Gewichtung von pos. Bsp. Gesenkt

$$q'_g(X, Y) = \frac{p'}{n + g}$$



1. Einführung
2. Descriptive Rule Discovery Verfahren
  - 2.1 Contrast Set Mining
  - 2.2 Emerging Pattern Mining
  - 2.3 Subgroup Discovery
3. Einheitliches Framework für Supervised Descriptive Rule Induction
4. Fazit

- **Ziel:** Vereinheitlichung von Bereichen
- nur Terminologie, Task Definition und Heuristiken werden verglichen
- Rule Selection Mechanism unterscheiden sich!
  - ➔ andere Schwerpunkte

# Supervised Descriptive Rule Discovery: A Unifying Survey of Contrast Set, Emerging Pattern and Subgroup Mining



TECHNISCHE  
UNIVERSITÄT  
DARMSTADT

Petra Kralj Novak, Nada Lavrac, Geoffrey I. Webb

Vielen Dank für Eure Aufmerksamkeit!