

Subgroup Discovery – Advanced Review

Martin Atzmueller

Knowledge and Data Engineering Group, University of Kassel, Germany
atzmueller@cs.uni-kassel.de

Keywords

subgroup discovery, local exceptionality detection, interestingness measures, algorithms, exploratory data mining

Abstract

Subgroup discovery is a broadly applicable descriptive data mining technique for identifying *interesting* subgroups according to some property of interest. This article summarizes fundamentals of subgroup discovery, before it reviews algorithms and further advanced methodological issues. In addition, we briefly discuss tools and applications of subgroup discovery approaches. In that context, we also discuss experiences and lessons learned and outline future directions in order to show the advantages and benefits of subgroup discovery.

Introduction

Subgroup discovery (66; 120; 80; 22; 101) has been established as a general and broadly applicable technique for descriptive and exploratory data mining: It aims at identifying descriptions of subsets of a dataset that show an interesting behavior with respect to certain interestingness criteria, formalized by a quality function, e. g., (120). This article summarizes fundamental concepts of subgroup before it provides an advanced review on algorithms, methodological issues, and applications. Overall, subgroup discovery and analytics are important tools for descriptive data mining: They can be applied, for example, for obtaining an overview on the relations in the data, for automatic hypotheses generation, and for data exploration. Prominent application examples include knowledge discovery in medical and technical domains, e.g., (49; 80; 22; 63). Typically, the discovered patterns are especially easy to interpret by the users and domain experts, cf. (49; 60). Standard subgroup discovery approaches commonly focus on a *single* target concept as the property of interest (66; 80; 60), while the quality function framework also enables *multi-target concepts*, e. g., (67; 14). Furthermore, more complex target properties (45; 81) can be formalized as *exceptional models*, cf. (81).

The remainder of this article is organized as follows. First, we present fundamentals of subgroup discovery. After that, we discuss state of the art algorithms. In the following sections, we outline methods for subgroup set selection, discuss applications and experiences, and provide an outlook on future directions and challenges. Finally, we conclude with a summary.

Preprint of: Martin Atzmueller. Subgroup Discovery – Advanced Review. WIREs Data Mining Knowl Discov 2015, 5:35–49. doi: 10.1002/widm.1144

Fundamentals of Subgroup Discovery

Subgroup discovery (66; 120; 67; 80; 22; 101) has been established as a versatile and effective method in descriptive and exploratory data mining. Similar to other methods for mining supervised local patterns, e.g., discriminative patterns (39), contrast sets (29), and emerging patterns (43), subgroup discovery aims at identifying *interesting* groups of individuals, where “interestingness is defined as distributional unusualness with respect to a certain property of interest” (120). Subgroup discovery has been well investigated concerning binary and nominal target concepts, i. e., properties of interest with a finite number of possible values (66; 120; 17). Furthermore, numeric target concepts have received increasing attention in subgroup discovery recently, and several approaches for using numeric attributes have been proposed, e.g., (97; 63; 56; 9). In the scope of this paper, we will adopt a broad definition of subgroup discovery, including single binary, nominal, and numeric target variables, but also extending to multi-target concepts, and to *exceptional model mining* (45; 81), as a variant of subgroup discovery that especially focuses on complex target properties.

In the remainder of this section, we first summarize the idea of local exceptionality detection employed by subgroup discovery. After that, we provide some necessary definitions and notation, before we formally tackle quality functions, optimistic estimates, and *top-k* pruning strategies.

Local Exceptionality Detection

Subgroup discovery is based on the idea of *local exceptionality detection*, that is, how locally exceptional, relevant, and thus interesting patterns can be detected, so-called nuggets in the data (cf. (66)). Local pattern mining, e. g., (98; 73) aims to detect such locally interesting patterns, in contrast to global models. Related approaches include, for example, frequent pattern mining (58), mining association rules (2; 75) and closed representations (31; 32). In contrast to those, however, subgroup discovery allows for a flexible definition of the applied quality function or interestingness measure, respectively. Therefore, many of the mentioned techniques can be captured in a subgroup discovery setting, e. g., (13) for a description-oriented community discovery method.

As sketched above, the exceptionality of a pattern is measured by a certain quality function. According to the type of the property of the subgroup, that we are interested in, we can distinguish between simple concepts such as a minimal frequency/size of the subgroup (also known as support for association rules), a deviating target share (confidence) of a binary target property of the subgroup, a significantly different subgroup mean of a numeric target concept, or more complex models, e. g., based on several target attributes for which their distribution significantly differs comparing the subgroup and the whole dataset. Using a quality function, a set of subgroups is then identified using a given subgroup discovery algorithm, e. g., using a heuristic or exhaustive search (66; 120; 80; 17) strategy, or a direct sampling approach (33). Typically, the *top-k* subgroups, or those above a minimal quality threshold are determined.

Basic Definitions

Formally, a *database* $D = (I, A)$ is given by a set of individuals I and a set of attributes A . For nominal attributes, a *selector* or *basic pattern* ($a_i = v_j$) is a Boolean function $I \rightarrow \{0, 1\}$ that is true if the value of attribute $a_i \in A$ is equal to v_j for the respective individual. For a numeric attribute a_{num} selectors ($a_{num} \in [min_j; max_j]$) can be defined analogously for each interval $[min_j; max_j]$ in the domain of a_{num} . The Boolean function is then set to true if the value of attribute a_{num} is within the respective range. The set of all basic patterns is denoted by Σ .

A subgroup is described using a description language, cf. (120), typically consisting of attribute–value pairs, e. g., in conjunctive or disjunctive normal form. Below, we present an exemplary conjunctive pattern description language; internal disjunctions can also be generated by appropriate attribute–value construction methods, if necessary, cf. (17). A *subgroup description* or (complex) *pattern* P is then given by a set of basic patterns $P = \{sel_1, \dots, sel_l\}$, $sel_i \in \Sigma, i = 1, \dots, l$, which is interpreted as a conjunction, i. e., $P(I) = sel_1 \wedge \dots \wedge sel_l$, with $length(P) = l$. A pattern can thus also be interpreted as the *body* of a *rule*. The *rule head* then depends on the property of interest. A *subgroup* $S_P := ext(P) := \{i \in I | P(i) = true\}$, i. e., a *pattern cover* is the set of all individuals that are covered by the subgroup description P .

The set of all possible subgroup description, and thus the possible search space is then given by 2^Σ , that is, all combinations of the basic patterns contained in Σ . In this context, the pattern $P = \emptyset$ covers all instances contained in the database.

Quality Functions

In general, quality and interestingness measures can be grouped into two categories: *Objective* and *subjective* measures (114; 46). Typically, a quality measure is determined according to the requirements and objectives of the analysis. Then also combinations of objective and subjective measures into hybrid quality measures are usually considered, cf. (6) for rules.

Common subjective interestingness measures are understandability, unexpectedness (new knowledge or knowledge contradicting existing knowledge), interestingness templates (describing classes of interesting patterns), and actionability (patterns which can be applied by the user to his or her advantage (103)). Objective measures are data driven and are derived using structure and properties of the data, e. g., based on statistical tests. In the following, we focus on such objective interestingness measures formalized by quality functions.

A *quality function*

$$q: 2^\Sigma \rightarrow \mathbb{R}$$

maps every pattern in the search space to a real number that reflects the interestingness of a pattern (or the pattern cover, respectively). The result of a subgroup discovery task is the set of k subgroup descriptions res_1, \dots, res_k with the highest interestingness according to the selected quality function.

In the binary, nominal and numeric setting a large number of quality functions has been proposed in literature, cf. (52; 66). In general, quality functions utilize the statistical distribution of the target concept(s) to score a subgroup pattern P . More complex quality functions compare a set of distributions, e. g., by utilizing the concept of exceptional models (81) discussed below. We can consider, for example, a pair of variables, or a whole set of variables arranged in a Bayesian network.

In addition to testing the (statistical) validity of the patterns, the (syntactical) complexity and simplicity of a pattern can also be considered. Commonly, simpler patterns are easier to understand and to interpret (108). Then, (6) describes a combination of quality measures for rules concerning the validity, i.e., the accuracy and the simplicity of the contained patterns. Furthermore, cost-based quality functions, e. g., (74), and cost-sensitive approaches (99) allow the modeling of costs for the quality assessment.

Binary and Nominal Target Quality Functions

Most of the quality functions for binary target concepts are based on the parameters contained in a four-fold table, e. g., (66) covering the positive/negative instances for the pattern P , its complement, and the general population, respectively. Many of the quality measures proposed in (66) trade-off the size $n = |ext(P)|$ of a subgroup and the deviation $t_P - t_0$, where t_P is the average value of a given target concept in the subgroup identified by the pattern P and t_0 the average value of the target concept in the general population. For binary (and nominal value) target concepts this relates to the *share* of the target concept in the subgroup and the general population. Thus, typical quality functions are of the form

$$q_S^a(P) = n^a \cdot (t_P - t_0), a \in [0; 1]. \quad (1)$$

For binary target concepts, this includes for example the *weighted relative accuracy* (q_S^1) for the size parameter $a = 1$, a simplified binomial function ($q_S^{0.5}$), for $a = 0.5$, or the *added value* function (q_S^0) for $a = 0$, which is order-equivalent to the lift (9) and the relative gain quality (22) function. Further examples for quality functions are given by the binominal test quality function q_B and the Chi-Square quality function q_C :

$$q_B(P) = \frac{(t_P - t_0) \cdot \sqrt{n}}{\sqrt{t_0 \cdot (1 - t_0)}} \cdot \sqrt{\frac{N}{N - n}}, \quad q_C(P) = \frac{n}{N - n} \cdot (t_P - t_0)^2,$$

where $N = |D|$ denotes the size of the database (general population).

Nominal valued target concepts (given by basic patterns) can be analyzed as in the binary case (one vs. all). For nominal attributes for which the set of the different nominal values needs to be analyzed, the binary case can be generalized analogously to multi-class settings, such that the whole distribution of the different values is assessed, cf. (66; 1). As an example, the quality function q_S^a can be generalized as follows for the general nominal setting:

$$q_M^a(P) = n^a \cdot \sum_{v_i} (t_P^{v_i} - t_0^{v_i})^2, a \in [0; 1], \quad (2)$$

where $t_P^{v_i}$ and $t_0^{v_i}$ denote the target shares in the subgroup and the general population, respectively, for each of the respective nominal values v_i contained in the value domain of the nominal target concept.

Other alternatives to the quality functions presented above include, for example, functions adapted from the area of *association rules*, e. g., (2) concerning the support and confidence parameters, as well as adaptations of measures from *information retrieval*, e. g., precision and recall and their combination in the F-measure, cf. (12). For more details, we refer to, e. g., (60) which provides a broad overview on quality functions used for subgroup discovery. In principle, many measures for *subgroup analysis* in epidemiology can also be utilized for subgroup discovery, especially in the medical domain. For example, the Odds Ratio function, sensitivity, specificity, significance, false alarm rate etc., see e. g., (53; 80; 77) for a survey and discussion. Furthermore, (88) provide an in-depth discussion for using the odds ratio and define statistically non-redundant subgroups utilizing the error bounds of the odds ratio measure.

Numeric Target Quality Functions

Quality functions for numeric target concepts, i. e., numeric attributes can be formalized by slightly adapting the quality functions q_a for binary targets presented above, cf. (66). The target shares t_P, t_0 of the subgroup and the general population, are replaced by the mean values of the target variable m_P, m_0 , respectively.

For the analog to the quality function q_S^a this results in:

$$q_M^a(P) = n^a(m_P - m_0), a \in [0; 1], \quad (3)$$

It is easy to see, that this function includes the binary formalization as a special case when we set $m = 1$ if the boolean target concept is *true* and $m = 0$, if it is *false*.

Using the parameter a , Equation 3 can be utilized for formalizing (order-)equivalent functions for several typically applied quality functions:

- The *mean gain* function q_M^0 ranks subgroups by the respective means m_P (order-equivalently) of the target concept, without considering the size of the subgroup. Therefore, a suitable minimal subgroup size threshold is usually required.
- Another simple example is given by the *mean test* (66; 54) with $q_M^{0.5}$. Furthermore, $q_M^{0.5}$ is also order-equivalent to the *z-score* quality function (104), given by $q_M^{0.5} \cdot \sigma_0$, where σ_0 is the standard deviation in the total population.
- Analogously to the weighted relative accuracy, the *impact* quality function (119) is given by q_M^1 .

Further quality functions consider, for example, the median (104) or the variance (26) of the target concepts in the subgroup. For more details on quality functions based on statistical tests (e. g., Student t-test or Mann-Whitney U-test) we refer to (66; 104).

Multi-Target Quality Functions

For multi-target quality functions, we consider functions that take into account a set of target concepts, e. g., (67). It is possible to extend single-target quality functions accordingly, for example, by extending an univariate statistical test to the multivariate case, e. g., (14): We then need to compare the multivariate distributions of a subgroup and the general population in order to identify interesting (and exceptional) patterns. For comparing multivariate means, for example, for a set of m numeric attributes T_M , with $m = |T_M|$ we can make use of Hotelling's T-squared test (61), for the quality measure q_H :

$$q_H(P) = \frac{n(n-m)}{m(n-1)} (\mu_P^{T_M} - \mu_0^{T_M})^\top CV_P^{T_M^{-1}} (\mu_P^{T_M} - \mu_0^{T_M}),$$

where $\mu_P^{T_M}$ is the vector of the means of the model attribute in the subgroup S_P , $CV_P^{T_M}$ is the respective covariance matrix, and $\mu_0^{T_M}$ is the vector of the means of the (numeric) target concepts in D .

As another option for a disjunction of target concepts, (122) propose convex quality functions for discovering cluster groups: Their approach does not use a single target concept, but allows for a disjunction of several target concepts/variables.

A more general framework for multi-target quality functions is given by *exceptional model mining* (81): It tries to identify interesting patterns with respect to a local model derived from a set of attributes. The interestingness can be defined, e.g., by a significant deviation from a model that is derived from the total population or the respective complement set of instances within the population. In general, a model consists of a specific *model class* and *model parameters* which depend on the values of the model attributes in the instances of the respective pattern cover. The quality measure q then determines the interestingness of a pattern according to its model parameters. Following (82), we outline some examples below, and refer to (81) for a detailed description.

- A simple example for an exceptionality measure for a set of attributes considers the task of identifying subgroups in which the correlation between two numeric attributes is especially strong, e. g., as measured by the Pearson correlation coefficient. This *correlation model class* has exactly one parameter, i.e., the correlation coefficient.
- Furthermore, using a *simple linear regression model*, we can compare the slopes of the regression lines of the subgroup to the general population or the subgroups' complement. This *simple linear regression model* shows the dependency between two numeric variables x and y : It is built by fitting a straight line in the two dimensional space by minimizing the squared residuals e_j of the model:

$$y_i = a + b \cdot x_i + e_j$$

As proposed in (81), the slope $b = \frac{cov(x,y)}{var(x)}$ computed given the covariance $cov(x, y)$ of x and y , and the variance $var(x)$ of x can then be used for identifying interesting patterns.

- The *logistic regression model* is used for the classification of a binary target attribute $y \in T$ from a set of independent binary attributes $x_j \in T \setminus y, j = 1, \dots, |T| - 1$.

The model is given by:

$$y = \frac{1}{1 + e^{-z}}, z = b_0 + \sum_j b_j x_j .$$

Interesting patterns are then those, for example, for which the model parameters b_j differ significantly from those derived from the total population.

- Another example is given by a *Bayesian network* as a rather complex target model. Then, a quality function for assessing the differences between two Bayesian networks can be defined. As proposed in (45), for example, models can then be compared based on the edit distance (112). Then networks induced by a subgroup and the general population (or the subgroups' complement), respectively, can be analyzed for identifying interesting patterns.

Top-K Pruning

As the result of subgroup discovery, the applied subgroup discovery algorithm can return a result set containing those subgroups above a certain minimal quality threshold, or only the top- k subgroups, that can then also be postprocessed further. While both options have their relevance depending on the analysis goals, the top- k approach provides more flexibility for applying different pruning options in the subgroup discovery process.

Basically, in a top- k setting, the set of the top- k subgroups is determined according to a given quality function. Then different pruning strategies can be applied for restricting the search space of a subgroup discovery algorithm. A simple option is given by *minimal support pruning* based on the antimonotone constraint of the subgroup size analogously to the Apriori algorithm for mining association rules, cf. (2). Furthermore, properties of certain quality functions enable more powerful approaches.

For several quality functions, for example, *optimistic estimates* (57; 9) can be applied for determining upper quality bounds: Consider the search for the k best subgroups: If it can be proven, that no subset of the currently investigated hypothesis is interesting enough to be included in the result set of k subgroups, then we can skip the evaluation of any subsets of this hypothesis, but can still guarantee the optimality of the result.

Another pruning mechanism is given by *generalization-aware pruning* (83), such that the quality of a subgroup is estimated against the qualities of its generalizations. Below, we discuss these two options in more detail.

Optimistic Estimate Pruning

The basic principle of *optimistic estimates* is to safely prune parts of the search space, e. g., as proposed in (120) for binary target variables. This idea relies on the intuition that if the k best hypotheses so far have already been obtained, and the optimistic estimate of the current subgroup is below the quality of the worst subgroup contained in the k best, then the current branch of the search tree can be safely pruned. More formally, an optimistic estimate oe of a quality function qf is a function such that $P' \supset P \Rightarrow oe(P) > qf(P')$, i.e., that no refinement P' of the pattern P can exceed the quality $oe(P)$.

An optimistic estimate is considered *tight* if there is a subset $S' \subseteq D$, such that $oe(P) = q(P'_S)$. While this definition requires the existence of a subset S' , there is not necessarily a pattern P'_S that describes S' , cf. (57).

For several quality functions, including many for the binary and numeric case described above, there exist (tight) optimistic estimates that can be applied for pruning. For a detailed overview, we refer to (57; 9).

Generalization-Aware Pruning

In general, a pattern can be compared to its generalizations, e. g., in order to fulfill a minimal improvement constraint (106), such that subgroups with a lower target share than any of its generalizations are removed. This can analogously be applied for mean-based numeric target measures (26), such that the mean observed in a subgroup needs to deviate significantly from the mean values induced by its generalizations. Accordingly, Lemmerich et al. (83) propose a pruning technique utilizing generalization-aware techniques, and present a set of optimistic estimates in this setting. These estimates take into account subgroup covers of the generalizations of a pattern, and allow for a rather efficient pruning approach for generalization-aware techniques.

Algorithms for Subgroup Discovery

Many algorithms for subgroup discovery focus on binary or nominal target attributes, for which heuristic (cf. (66; 49; 76)) as well as exhaustive methods (cf. (66; 120; 65; 17; 9)) are applied. Algorithms for multi-relational data include those by the *MIDOS* (120) and *SubgroupMiner* (68) system, for which the latter also includes analysis options for spatio-temporal data. Essentially, heuristic approaches like beam search trade agility for completeness and are often applied in certain domains for which exhaustive methods take rather long to explore the complete search space, e. g., in dense numerical data for which numeric features are discretized on the fly such that numeric subspaces can be explored heuristically, cf. (92). However, due to efficient pruning techniques exhaustive methods can both achieve sufficiently good runtimes and guarantee completeness even in complex domains like ubiquitous social data, e. g., (11).

Furthermore, for the multi-target concept setting including exceptional model mining, there exist heuristic (71; 116) approaches as well efficient exhaustive discovery algorithms (82). As another option, sampling (e.g., (120; 113; 47)) can be used for estimating the quality of a subgroup on a (potentially significantly) smaller subset of the case base. The sequential sampling algorithm *GSS* (110), for example, discovers the k best subgroups according to a given confidence level, for quality functions that can be estimated with bounded error. Also, local pattern sampling methods, e. g., (33) can be seen as an alternative to exhaustive mining approaches utilizing direct sampling procedures with confidence bounds.

Further subgroup discovery approaches apply evolutionary techniques, i. e., genetic algorithms, e. g., (41; 37; 90; 38) for their discovery and refinement strategy. For removing uninteresting subgroups, expectation-driven measures, e. g., (85) can be considered such that (expected) interactions between variables are captured.

In general, exhaustive algorithms typically make use of the proposed pruning options for an efficient processing of the search space. Combining intelligent sampling approaches with fast exhaustive methods, e.g., with *SD-Map* (17) or *SD-Map** (9) can then be seen as a promising option for efficiently mining potentially arbitrarily large databases. Below, we summarize the main characteristics of heuristic and exhaustive approaches, and discuss exemplary algorithms.

Heuristic Algorithms

For heuristic approaches, commonly a beam search (89) strategy is used because of its efficiency. The search starts with a list of subgroup hypotheses of size w (corresponding to the *beam width*), which may be initially empty. The w subgroup hypotheses contained in the beam are then expanded iteratively, and only the best w expanded subgroups are kept implementing a hill-climbing greedy search. Lavrac et al. (80), for example, describe the application of the beam-search based *CN2-SD* algorithm adapted for subgroup discovery. To improve upon simple greedy approaches, other alternatives such as the *PRIM* algorithm (48) have been proposed, which employ a *patient* search strategy.

Beam search traverses the search space non-exhaustively and thus does not guarantee to discover the complete set of the top- k subgroups, or all subgroups above a minimal quality threshold. It can also be regarded as a variant of an anytime algorithm, since the search process can be stopped at any point such that the *currently best* subgroups are available. It is also possible to apply beam search to larger description spaces, e. g., including richer descriptions for numeric attributes, cf. (92). Furthermore, subgroup set selection can also be integrated into such heuristic approaches (116) as described below.

Alternatives include genetic algorithms, e. g., (41; 37; 90; 38), that cast the subgroup discovery task into an evolutionary optimization problem. This can also be applied for subgroup discovery in continuous domains, e. g., (107).

Efficient Exhaustive Algorithms

In contrast to heuristic methods, exhaustive approaches guarantee to discover the best solutions. However, the runtime costs of a (naive) exhaustive algorithm usually prohibit its application for larger search spaces. Examples of exhaustive algorithms include Apriori-based methods (2), for example, the Apriori-SD (65) algorithm; more alternatives are mentioned below.

Depending on the applied algorithm, there are different pruning options that can be applied. Many state-of-the-art algorithms apply extensions of frequent pattern trees (FP-trees) (59) in a pattern-growth fashion. Then, typically optimistic estimate pruning is applied, while generalization-aware pruning is better supported by layer-wise algorithms based on the Apriori (2) principle. Furthermore, (122; 123) have proposed branch-and-bound algorithms that require special (convex) quality functions for pruning the search space.

As efficient exhaustive algorithms, the *BSD* and the *SD-Map** algorithms, for example, allow the efficient handling of binary, nominal and numeric target properties. Both algorithms apply optimistic estimate pruning, but utilize different core data structures, bitset representations vs. extended FP-trees, cf. (59). FP-trees are also used in other subgroup discovery algorithms, e. g., by DpSubgroup (57; 56). As an extension of *SD-Map**, the GP-Growth algorithm (82) allows subgroup discovery for single target and multi-target concepts, e. g., for exceptional model mining; several model classes and quality functions can be implemented using the algorithm. In the following, we briefly review those algorithms in more detail.

SD-Map, SD-Map* and GP-Growth

*SD-Map** (9) is based on the efficient *SD-Map* (17) algorithm utilizing an FP-tree data structure, cf. (59) i. e., an extended prefix-tree-structure that stores information for pattern refinement and evaluation. *SD-Map** applies a divide and conquer method, first mining patterns containing one selector and then recursively mining patterns of size 1 conditioned on the occurrence of a (prefix) 1-selector. For the binary case, an FP-tree node stores the subgroup size and the true positive count of the respective subgroup description. In the continuous case, it considers the sum of values of the target variable, enabling us to compute the respective quality functions value accordingly. Therefore, all the necessary information is locally available in the FP-tree structure.

For extending the FP-tree structure towards multi-target concepts, we utilize the concept of *evaluation bases* introduced by (82). Then, all information required for the *evaluation* of the respective quality functions is stored in the nodes of the FP-tree, as the basis of the GP-Growth algorithm extending *SD-Map/SD-Map**. With this technique, a large number of single and multi-target concept quality functions can be implemented, cf. (82).

BSD

The BSD algorithm (86) utilizes a vertical bitset (bitvector) based data structure. Vertical data representations have been proposed for other data mining tasks, e.g., by (121). (36) used a bitset representation for maximal frequent itemset mining. As a general search strategy, BSD uses a depth-first-search approach with one level look-ahead (similar to the DpSubgroup (57; 56) algorithm). BSD uses a vertical data layout utilizing bitsets (vectors of bits), for the selectors, the instances reflecting the current subgroup hypothesis, and an additional array for the (numeric) values of the target variable. Then, the search, i. e., the refinement of the patterns can be efficiently implemented using logical *AND* operations on the respective bitsets, such that the target values can be directly retrieved.

Subgroup Set Selection

Due to multi-correlations between the selectors, some of the subgroups contained in the result set can overlap significantly. Therefore, usually a *high-quality* set of *diverse* subgroups should be retrieved. Subgroup set selection is one of the critical issues for removing redundancy and improving the interestingness of the overall subgroup discovery result, e. g., as first described in (66). Constraints denoting redundancy filters, for example, can be used to prune large regions of the search space. This is especially important for certain search strategies, which do not constrain the search space themselves, e.g., exhaustive search compared to beam search. Klösgen (66) distinguishes two types of redundancy filters: Logical and heuristic filters. The filters include either logical or heuristic implications for the truth value of a constraint condition with respect to a predecessor/successor pair of subgroups. Logical filters can be described as *strong filters*; they can be used to definitely exclude a region of the search space. Heuristic filters can be used as *weak filters*; these are applied as a first step in a brute force search, where the excluded regions of the search space can be refined later.

The general task of diverse subgroup set discovery is described in (116), for which different types of redundancy and according selection heuristics are proposed. There are several heuristic approaches for pattern set selection in general, e. g., (35), as well as for subgroup discovery (72; 116). In particular, several quality functions for selecting *pattern teams* are proposed in (72), which can be applied in a post-processing step. Furthermore, a method for semi-automatic retrieval of a set of subgroups is described in (18), for which mixed-initiative case-based techniques are applied. In the following, we outline several options for *diversity-aware* and *redundancy-aware* subgroup set discovery and selection in more detail, focusing on condensed representations, relevance criteria, covering approaches, as well as causal subgroup analysis for identifying causally related sets of subgroups.

Condensed Representations

In the field of association rules, *condensed* representations of frequent item sets have been developed for reducing the size of the set of association rules that are generated, e. g., (102; 27). These representations are used for the (implicit) redundancy management, since then the condensed patterns also describe the specifically interesting patterns, and can significantly reduce the size of the result sets. The efficiency of the association rule discovery method is also increased significantly. Such techniques can also be generalized for frequent patterns (c.f., (95; 34)). For subgroup discovery, target-closed representations can be formalized, cf. (55) for details. In that case, also an implicit redundancy management based on the subgroup descriptions is performed.

Relevance of Subgroups

As a quite simple method for redundancy management of subgroups for binary targets we can consider the (*ir-*)*relevance* (e.g., (49)) of a subgroup with respect to a set of subgroups: A (specialized) subgroup hypothesis S_N is *irrelevant* if there exists a (generalized) subgroup hypothesis S_P such that the true positives of S_N are a subset of the true positive of S_P and the false positives of S_N are a superset of the false positives of S_P . This concept is also closely related to the concept of closed patterns (51) and according relevance criteria (78; 51).

Embedding this redundancy management technique into the search process is straightforward: When considering a subgroup hypothesis for inclusion into the set of the k best subgroups, the test for (strict) irrelevancy can be applied. In addition, such a method can also be implemented in an optional post-processing step. Furthermore, (55) introduces delta-relevance which provides a more relaxed definition of coverage (essentially trading off precision vs. simplicity) with the overall goal of summarizing relevant patterns even more.

Covering Approaches

Similar to covering approaches for rule learning, a subgroup set can also be selected according to its overall coverage of the dataset. The *weighted covering algorithm* (49; 80) is such an approach that works by example reweighting. It iteratively focuses the subgroup selection method on the space of target records not covered so far, by reducing the weights of the already covered data records. As discussed in (64), example reweighting can also be used as a search heuristic – in combination with a suitable quality function. In this way, weighted covering is integrated in the subgroup discovery algorithm, i.e., the search step directly (e.g., (49)): In each search iteration only the best subgroup is considered, then the instances are reweighted, focusing the subgroup discovery method on the not yet covered target class cases.

Causal Subgroup Analysis

For identifying causal subgroups efficiently, constraint-based causal discovery algorithms, e. g., (40) can be applied. These try to limit the possible causal models by analyzing observational data. In causal subgroup analysis (e. g., (69; 20)) subgroups are identified which are causal for the target concept; for a causal subgroup, the manipulation of an instance to belong to the subgroup would also affect the probability of the instance to belong to the target group (40). In accordance with the general principle that correlation does not imply causation, constraint-based algorithms apply statistical tests, e. g., the χ^2 -test for independence in order to test the (conditional) dependence and independence of variables to *exclude* certain causal relations. After causal subgroups have been detected, the user can retain these (important) subgroups, which have a direct dependency relation to the target concept, in contrast to the remaining non-causal subgroups, which are often redundant given the causal subgroups.

Tools and Applications

Subgroup discovery is a powerful and broadly applicable data mining approach, in particular, for descriptive data mining tasks. It is typically applied, for example, in order to obtain an overview on the relations in the data and for automatic hypotheses generation. Furthermore, also predictive tasks can be tackled, e. g., by stacking approaches (70), or by applying the *LeGo* framework for combining local patterns into global models.

From a tool perspective, there exist several software packages for subgroup discovery, e. g., (15; 94; 10; 42). As open source options, there are, for example, subgroup discovery modules for the data mining systems *Orange* (42) and *RapidMiner* (96), the *Cortana* (94) system for discovering local patterns in data, as well as the specialized subgroup discovery and analytics system VIKAMINE (15; 10). Using the latter a number of successful real-world subgroup discovery applications have been implemented. These cover, for example, knowledge discovery and quality control setting in the medical domain (22; 23; 105), industrial applications (9), as well as pattern analytics in the social media domain (11). The system is targeted at a broad range of users, from industrial practitioners to ML/KDD researchers, students, and users interested in knowledge discovery and data analysis in general. Especially the visual mining methods enable the direct integration of the user to overcome major problems of automatic data mining methods, cf. (15; 10).

Applications include, for example, knowledge discovery in the medical domain, technical (fault) analysis, e. g., (9; 62), or mining social data, e. g., (13; 4; 14). We discuss these exemplarily below. Furthermore, for heterogenous data exceptional model mining (81; 45; 44; 82) opens up a wide range of options. There are also applications in related fields, e. g., in software engineering (30) for requirements engineering and design.

Knowledge Discovery in the Medical Domain

Subgroup discovery is a prominent approach for mining medical data, e. g., (49; 79; 50; 22; 23; 15; 105). Using the VIKAMINE system, for example, subgroup discovery has been applied for large-scale knowledge discovery and quality control in the clinical application SONOCONSULT, cf., (105). For this, several data sources including structured clinical documentation and unstructured documents, e.g., (7), were integrated. The main data source was given by the SONOCONSULT system, which has been in routine use since 2002 as the only documentation system for ultrasound examinations in the DRK-hospital of Berlin-Köpenick; since 2005, it is in routine use at the university hospital of Würzburg. The physicians considered statistical analysis as one of the most desirable features. In the analysis and results, e. g., (22; 23; 24), subgroup discovery was applied on a large set of clinical features together with laboratory values and diagnostic information from several systems.

According to the physicians, subgroup discovery and analysis was quite suitable for examining common medical questions, e.g. whether a certain pathological state is significantly more frequent if combinations of other pathological states exist, or if there are diagnoses, which one physician documents significantly more or less frequently than the average. Furthermore, VIKAMINE also provided an intuitive overview on the data, in addition to the knowledge discovery and quality control functions. Then, subgroup discovery can be performed in a semi-automatic approach, first generating hypothesis using automatic methods that are then inspected and refined using visual analytics techniques, cf. (15; 19).

Technical Fault Analysis

Technical applications of subgroup discovery include, for example, mining service processes (100), analysis of smart electrical meter data (62), or fault analysis of production processes (20). The latter, for example, has been implemented using VIKAMINE (9): It aimed at large-scale fault detection and analysis using subgroup discovery. Specifically, the task required the identification of subgroups (as combination of certain factors) that cause a significant increase/decrease in, e.g., the fault/repair rates of certain products. Similar problems in industry concern, for example, the number of service requests for a certain technical component, or the number of calls of customers to service support.

Such applications of subgroup discovery often require the utilization of continuous parameters. Then, the target concepts can often not be analyzed sufficiently using the standard discretization techniques, since the discretization of the variables causes a loss of information. As a consequence, the interpretation of the results is often difficult using standard data mining tools. In this context, VIKAMINE provided state-of-the-art algorithmic implementations, cf. (9), and enabled a semi-automatic involvement of the domain experts for effectively contributing in a discovery session.

Subgroup Discovery in Social Media

In mining social media (3), subgroup discovery methods are a versatile tool for focusing on different facets of the application domain. Subgroup discovery was applied, for example, for obtaining descriptive profiles of spammers in social media, specifically social bookmarking systems. Here, subgroup discovery was applied for the characterization of spammers, i. e., to describe them by their most prominent features (12). The mined patterns capturing certain spammer subgroups provide explanations and justifications for marking or resolving spammer candidates.

In addition, in social bookmarking systems, it is usually useful to identify high-quality tags, i.e., tags with a certain maturity, cf. (8). Subgroup discovery was applied for obtaining maturity profiles of tags based on a set of graph centrality features on the tag – tag cooccurrence graph, which are simple to compute and to assess. Then, they can be applied for tag recommendations, faceted browsing, or for improving search.

Furthermore, subgroup discovery has been utilized for community pattern analytics in social media, e. g., (5), as well as semi-automatic approaches for pattern detection in geo-referenced tagging data, exemplified by an application using Flickr data, cf., (11). In this domain of “Big Data”, subgroup discovery can also provide suitable solutions using the efficient automatic algorithms; the combination of automatic and interactive visualization methods complemented each other for a successful discovery process. Especially subgroup introspection, and pattern explanation capabilities, e. g., (19; 18; 25) proved essential during pattern analytics and assessment.

Future Directions and Challenges

Overall, there is already a large body of works on algorithmic as well as methodological issues on subgroup discovery. Major challenging points include the algorithmic performance, the redundancy of the result set of subgroups, adequate comprehensive visualization, and the processing and integration of heterogenous data. Larger search spaces, like those encountered for numerical data, (complex) multi-relational datasets, e. g., encountered in social networks, or spatio-temporal data require efficient algorithms that can handle those different types of data, e. g., (91; 14). Also combinations of such different data characteristics, for example, temporal pattern mining for event detection (28), or temporal subgroup analytics (111) provide further challenges, especially considering sophisticated exceptional model classes in that area.

Typically, heuristic approaches are first established before advanced processing methods like sophisticated pruning and suppression heuristics enable exhaustive subgroup discovery techniques. Furthermore, processing large volumes of data (i. e., big data) is another challenge. In that area, extensions of techniques for parallelizing the computation, e. g., (86) or techniques from the field of association rules, e. g., (87) can provide interesting options for further improvements in that area. Also, sampling approaches, e. g., (33; 110; 109) can be applied for addressing these issues.

In addition, the integration of (rich) background knowledge in a knowledge-intensive approach, e. g., (21; 22; 16; 117) is a prerequisite for the analysis of large datasets for which relations and prior information needs to be utilized. This also tackles the area of automatic subgroup discovery, recent search strategies, e. g., (93) and the applied significance filtering in many methods and tools (84). For a full-scale approach, these issues need to be addressed such that suitable methods can be integrated comprehensively, from automatic to interactive approaches, e. g., (15; 10; 115), which can also be applied for generating appropriate explanations (118). Then, the combination, integration, and further development of such techniques and methods will lead to novel techniques embedded in a comprehensive approach for subgroup discovery and analytics, towards robust tool implementations, and finally to further successful applications of subgroup discovery.

Conclusions

In this paper, we have reviewed key areas of subgroup discovery considering algorithmic issues for discovering subgroups as well as for refining the discovery results. Thus, we covered the fundamentals of subgroup discovery, provided an overview from an algorithmic point of view, briefly discussed efficient algorithms, and summarized approaches for selecting a final subgroup set. Furthermore, we presented different tools and applications of subgroup discovery and outlined several interesting and challenging directions for future work. Overall, subgroup discovery is a versatile and powerful method that can be tuned to many application characteristics. It provides a comprehensive approach integrating different techniques for providing solutions in many domains.

References

- [1] T. Abudawood and P. Flach. Evaluation Measures for Multi-Class Subgroup Discovery. In *Proc. ECML/PKDD*, volume 5782 of *LNAI*, pages 35–50, Heidelberg, Germany, 2009. Springer Verlag.
- [2] R. Agrawal and R. Srikant. Fast Algorithms for Mining Association Rules. In J. B. Bocca, M. Jarke, and C. Zaniolo, editors, *Proc. 20th Int. Conf. Very Large Data Bases, (VLDB)*, pages 487–499. Morgan Kaufmann, 1994.
- [3] M. Atzmueller. Mining Social Media: Key Players, Sentiments, and Communities. *WIREs: Data Mining and Knowledge Discovery*, 1069, 2012.
- [4] M. Atzmueller. Data Mining on Social Interaction Networks. *Journal of Data Mining and Digital Humanities*, 1, June 2014.
- [5] M. Atzmueller. Social Behavior in Mobile Social Networks: Characterizing Links, Roles and Communities. In A. Chin and D. Zhang, editors, *Mobile Social Networking: An Innovative Approach*, Computational Social Sciences, pages 65–78. Springer Verlag, Heidelberg, Germany, 2014.

- [6] M. Atzmueller, J. Baumeister, and F. Puppe. Quality Measures and Semi-Automatic Mining of Diagnostic Rule Bases (extended version). In *Proc. 15th International Conference on Applications of Declarative Programming and Knowledge Management (INAP 2004)*, number 3392 in LNAI, pages 65–78, Heidelberg, Germany, 2005. Springer Verlag.
- [7] M. Atzmueller, S. Beer, and F. Puppe. Data Mining, Validation and Collaborative Knowledge Capture. In S. Brüggemann and C. d’Amato, editors, *Collaboration and the Semantic Web: Social Networks, Knowledge Networks and Knowledge Resources*. IGI Global, Hershey, PA, USA, 2012.
- [8] M. Atzmueller, D. Benz, A. Hotho, and G. Stumme. Towards Mining Semantic Maturity in Social Bookmarking Systems. In *Proc. Workshop Social Data on the Web, 10th Intl. Semantic Web Conference*, 2011.
- [9] M. Atzmueller and F. Lemmerich. Fast Subgroup Discovery for Continuous Target Concepts. In *Proc. 18th International Symposium on Methodologies for Intelligent Systems (ISMIS 2009)*, volume 5722 of LNCS, pages 1–15, Heidelberg, Germany, 2009. Springer Verlag.
- [10] M. Atzmueller and F. Lemmerich. VIKAMINE - Open-Source Subgroup Discovery, Pattern Mining, and Analytics. In *Proc. ECML/PKDD*, volume 7524 of LNAI, Heidelberg, Germany, 2012. Springer Verlag.
- [11] M. Atzmueller and F. Lemmerich. Exploratory Pattern Mining on Social Media using Geo-References and Social Tagging Information. *IJWS*, 2(1/2), 2013.
- [12] M. Atzmueller, F. Lemmerich, B. Krause, and A. Hotho. Who are the Spammers? Understandable Local Patterns for Concept Description. In *Proc. 7th Conference on Computer Methods and Systems*, Krakow, Poland, 2009. Oprogramowanie Nauko-Techniczne.
- [13] M. Atzmueller and F. Mitzlaff. Efficient Descriptive Community Mining. In *Proc. 24th International FLAIRS Conference*, pages 459 – 464, Palo Alto, CA, USA, 2011. AAAI Press.
- [14] M. Atzmueller, J. Mueller, and M. Becker. *Mining, Modeling and Recommending ‘Things’ in Social Media*, volume 8940 of LNAI, chapter Exploratory Subgroup Analytics on Ubiquitous Data. Springer Verlag, Heidelberg, Germany, 2015.
- [15] M. Atzmueller and F. Puppe. Semi-Automatic Visual Subgroup Mining using VIKAMINE. *Journal of Universal Computer Science*, 11(11):1752–1765, 2005.
- [16] M. Atzmueller and F. Puppe. A Methodological View on Knowledge-Intensive Subgroup Discovery. In *Proc. 15th International Conference on Knowledge Engineering and Knowledge Management (EKAW 2006)*, number 4248 in LNAI, pages 318–325, Heidelberg, Germany, 2006. Springer Verlag.

- [17] M. Atzmueller and F. Puppe. SD-Map - A Fast Algorithm for Exhaustive Subgroup Discovery. In *Proc. PKDD*, volume 4213 of *LNAI*, pages 6–17, Heidelberg, Germany, 2006. Springer Verlag.
- [18] M. Atzmueller and F. Puppe. A Case-Based Approach for Characterization and Analysis of Subgroup Patterns. *Journal of Applied Intelligence*, 28(3):210–221, 2008.
- [19] M. Atzmueller and F. Puppe. Semi-Automatic Refinement and Assessment of Subgroup Patterns. In *Proc. 21th International Florida Artificial Intelligence Research Society Conference (FLAIRS-2008)*, pages 518–523, Palo Alto, CA, USA, 2008. AAAI Press.
- [20] M. Atzmueller and F. Puppe. *Knowledge Discovery Enhanced with Semantic and Social Information*, chapter A Knowledge-Intensive Approach for Semi-Automatic Causal Subgroup Discovery. Springer Verlag, Heidelberg, Germany, 2009.
- [21] M. Atzmueller, F. Puppe, and H.-P. Buscher. Towards Knowledge-Intensive Subgroup Discovery. In *Proc. LWA 2004, Germany*, pages 117–123, 2004.
- [22] M. Atzmueller, F. Puppe, and H.-P. Buscher. Exploiting Background Knowledge for Knowledge-Intensive Subgroup Discovery. In *Proc. 19th International Joint Conference on Artificial Intelligence (IJCAI-05)*, pages 647–652, Edinburgh, Scotland, 2005.
- [23] M. Atzmueller, F. Puppe, and H.-P. Buscher. Profiling Examiners using Intelligent Subgroup Mining. In *Proc. 10th International Workshop on Intelligent Data Analysis in Medicine and Pharmacology (IDAMAP-2005)*, pages 46–51, Aberdeen, Scotland, 2005.
- [24] M. Atzmueller, F. Puppe, and H.-P. Buscher. A Semi-Automatic Approach for Confounding-Aware Subgroup Discovery. *International Journal on Artificial Intelligence Tools (IJAIT)*, 18(1):1 – 18, 2009.
- [25] M. Atzmueller and T. Roth-Berghofer. The Mining and Analysis Continuum of Explaining Uncovered. In *Proc. 30th SGAI International Conference on Artificial Intelligence (AI-2010)*, 2010.
- [26] Y. Aumann and Y. Lindell. A Statistical Theory for Quantitative Association Rules. *J. Intell. Inf. Syst.*, 20(3):255–283, 2003.
- [27] Y. Bastide, N. Pasquier, R. Taouil, G. Stumme, and L. Lakhal. Mining Minimal Non-Redundant Association Rules Using Frequent Closed Itemsets. In J. Lloyd, V. Dahl, U. Furbach, M. Kerber, K.-K. Lau, C. Palamidessi, L., M. Pereira, Y. Sagiv, P., and J. Stuckey, editors, *Computational Logic - CL 2000. Proc. CL '00*, pages 972–986, Heidelberg, Germany, 2000. Springer Verlag.

- [28] I. Batal, D. Fradkin, J. Harrison, F. Moerchen, and M. Hauskrecht. Mining Recent Temporal Patterns for Event Detection in Multivariate Time Series Data. In *Proc. ACM SIGKDD*, KDD '12, pages 280–288, New York, NY, USA, 2012. ACM.
- [29] S. D. Bay and M. J. Pazzani. Detecting group differences: Mining contrast sets. *Data Mining and Knowledge Discovery*, 5(3):213–246, July 2001.
- [30] K. Behrenbruch, M. Atzmueller, C. Evers, L. Schmidt, G. Stumme, and K. Geihs. A Personality Based Design Approach Using Subgroup Discovery. In *Human-Centred Software Engineering*, volume 7623 of *LNCS*, pages 259–266. Springer Verlag, Heidelberg, Germany, 2012.
- [31] M. Boley and H. Grosskreutz. Non-redundant Subgroup Discovery Using a Closure System. *Machine Learning and Knowledge Discovery in Databases*, pages 179–194, 2009.
- [32] M. Boley, T. Horvath, A. Poigné, and S. Wrobel. Listing Closed Sets of Strongly Accessible Set Systems with Applications to Data Mining. *Theoretical Computer Science*, 411(3):691–700, 2010.
- [33] M. Boley, C. Lucchese, D. Paurat, and T. Gärtner. Direct Local Pattern Sampling by Efficient Two-step Random Procedures. In *Proc. ACM SIGKDD*, KDD '11, pages 582–590, New York, NY, USA, 2011. ACM.
- [34] J.-F. Boulicaut. *Encyclopedia of Data Warehousing and Mining*, chapter Condensed Representations for Data Mining, pages 37–79. Idea Group, 2006.
- [35] B. Bringmann and A. Zimmermann. The Chosen Few: On Identifying Valuable Patterns. In *Proc. IEEE Intl. Conf. on Data Mining*, pages 63–72, Washington, DC, USA, 2007. IEEE Computer Society.
- [36] D. Burdick, M. Calimlim, and J. Gehrke. MAFIA: A Maximal Frequent Itemset Algorithm for Transactional Databases. In *Proc. 17th International Conference on Data Engineering (ICDE'01)*, pages 443–452, 2001.
- [37] C. J. Carmona, P. González, M. J. del Jesús, and F. Herrera. NMEEF-SD: Non-dominated Multiobjective Evolutionary Algorithm for Extracting Fuzzy Rules in Subgroup Discovery. *IEEE T. Fuzzy Systems*, 18(5):958–970, 2010.
- [38] C. J. Carmona, P. González, M. J. del Jesus, and F. Herrera. Overview on Evolutionary Subgroup Discovery: Analysis of the Suitability and Potential of the Search Performed by Evolutionary Algorithms. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, 4(2):87–103, 2014.
- [39] H. Cheng, X. Yan, J. Han, and P. S. Yu. Direct Discriminative Pattern Mining for Effective Classification. In *Proc. 24th Intl. IEEE Conference on Data Engineering*, pages 169–178, Washington, DC, USA, 2008. IEEE Comp. Soc.

- [40] G. F. Cooper. A Simple Constraint-Based Algorithm for Efficiently Mining Observational Databases for Causal Relationships. *Data Min. Knowl. Discov.*, 1(2):203–224, 1997.
- [41] M. J. del Jesus, P. González, F. Herrera, and M. Mesonero. Evolutionary Fuzzy Rule Induction Process for Subgroup Discovery: A Case Study in Marketing. *IEEE T. Fuzzy Systems*, 15(4):578–592, 2007.
- [42] J. Demšar, T. Curk, A. Erjavec, Črt Gorup, T. Hočevar, M. Milutinovič, M. Možina, M. P. andMarko Toplak, A. Starič, M. Štajdohar, L. Umek, L. Žagar, J. Žbontar, M. Žitnik, and B. Zupan. Orange: Data Mining Toolbox in Python. *Journal of Machine Learning Research*, 14:2349–2353, 2013.
- [43] G. Dong and J. Li. Efficient Mining of Emerging Patterns: Discovering Trends and Differences. In *Proc. ACM SIGKDD*, pages 43–52, New York, NY, USA, 1999. ACM.
- [44] W. Duivesteijn, A. Feelders, and A. J. Knobbe. Different slopes for different folks: mining for exceptional regression models with cook’s distance. In *Proc. ACM SIGKDD*, pages 868–876, New York, NY, USA, 2012. ACM.
- [45] W. Duivesteijn, A. Knobbe, A. Feelders, and M. van Leeuwen. Subgroup Discovery Meets Bayesian Networks—An Exceptional Model Mining Approach. In *Proc. IEEE Intl. Conference on Data Mining (ICDM)*, pages 158–167, Washington, DC, USA, 2010. IEEE.
- [46] A. A. Freitas. On Rule Interestingness Measures. *Knowledge-Based Systems*, 12(5-6):309–325, 1999.
- [47] Y. Freund. Self Bounding Learning Algorithms. In *COLT: Proc. 11th Annual Conference on Computational Learning Theory*, New York, NY, USA, 1998. ACM.
- [48] J. H. Friedman and N. I. Fisher. Bump Hunting in High-Dimensional Data. *Statistics and Computing*, 9(2), 1999.
- [49] D. Gamberger and N. Lavrac. Expert-Guided Subgroup Discovery: Methodology and Application. *Journal of Artificial Intelligence Research*, 17:501–527, 2002.
- [50] D. Gamberger, N. Lavrac, and G. Krstacic. Active Subgroup Mining: A Case Study in Coronary Heart Disease Risk Group Detection. *Artificial Intelligence in Medicine*, 28:27–57, 2003.
- [51] G. C. Garriga, P. Kralj, and N. Lavrac. Closed Sets for Labeled Data. *Journal of Machine Learning Research*, 9:559–580, 2008.
- [52] L. Geng and H. J. Hamilton. Interestingness Measures for Data Mining: A Survey. *ACM Computing Surveys*, 38(3), 2006.

- [53] A. S. Glas, J. G. Lijmer, M. H. Prins, G. J. Bonsel, and P. M. Bossuyt. The Diagnostic Odds Ratio: A Single Indicator of Test Performance. *Journal of Clinical Epidemiology*, 56(11):1129 – 1135, 2003.
- [54] H. Grosskreutz. Cascaded Subgroup Discovery with an Application to Regression. In *Proc. ECML/PKDD*, volume 5211 of *LNAI*, Heidelberg, Germany, 2008. Springer Verlag.
- [55] H. Großkreutz, D. Paurat, and S. Rüping. An Enhanced Relevance Criterion for More Concise Supervised Pattern Discovery. In *Proc. ACM SIGKDD*, KDD '12, pages 1442–1450, New York, NY, USA, 2012. ACM.
- [56] H. Grosskreutz and S. Rüping. On Subgroup Discovery in Numerical Domains. *Data Mining and Knowledge Discovery*, 19(2):210–226, Oct 2009.
- [57] H. Grosskreutz, S. Rüping, and S. Wrobel. Tight Optimistic Estimates for Fast Subgroup Discovery. In *Proc. ECML/PKDD*, volume 5211 of *LNAI*, pages 440–456, Heidelberg, Germany, 2008. Springer Verlag.
- [58] J. Han, H. Cheng, D. Xin, and X. Yan. Frequent Pattern Mining: Current Status and Future Directions. *Data Mining and Knowledge Discovery*, 15:55–86, 2007.
- [59] J. Han, J. Pei, and Y. Yin. Mining Frequent Patterns Without Candidate Generation. In W. Chen, J. Naughton, and P. A. Bernstein, editors, *2000 ACM SIGMOD Intl. Conference on Management of Data*, pages 1–12. ACM Press, 05 2000.
- [60] F. Herrera, C. Carmona, P. Gonzalez, and M. del Jesus. An Overview on Subgroup Discovery: Foundations and Applications. *Knowledge and Information Systems*, 29(3):495–525, 2011.
- [61] H. Hotelling. The Generalization of Student's Ratio. *Ann. Math. Statist.*, 2(3):360–378, 1931.
- [62] N. Jin, P. Flach, T. Wilcox, R. Sellman, J. Thumim, and A. Knobbe. Subgroup Discovery in Smart Electricity Meter Data. *Industrial Informatics, IEEE Transactions on*, 10(2):1327–1336, May 2014.
- [63] A. M. Jorge, F. Pereira, and P. J. Azevedo. Visual Interactive Subgroup Discovery with Numerical Properties of Interest. In *Proceedings of the 9th International Conference on Discovery Science (DS 2006)*, volume 4265 of *LNAI*, pages 301–305, Barcelona, Spain, October 2006. Springer.
- [64] B. Kavsek and N. Lavrac. Analysis of Example Weighting in Subgroup Discovery by Comparison of Three Algorithms on a Real-Life Data Set. In *Proc. Workshop on Advances in Inductive Rule Learning, at ECML/PKDD*, 2004.

- [65] B. Kavsek, N. Lavrac, and V. Jovanoski. APRIORI-SD: Adapting Association Rule Learning to Subgroup Discovery. In *Proc. 5th Intl. Symposium on Intelligent Data Analysis*, pages 230–241, Heidelberg, Germany, 2003. Springer Verlag.
- [66] W. Klösgen. Explora: A Multipattern and Multistrategy Discovery Assistant. In U. M. Fayyad, G. Piatetsky-Shapiro, P. Smyth, and R. Uthurusamy, editors, *Advances in Knowledge Discovery and Data Mining*, pages 249–271. AAAI Press, 1996.
- [67] W. Klösgen. *Handbook of Data Mining and Knowledge Discovery*, chapter 16.3: Subgroup Discovery. Oxford University Press, New York, 2002.
- [68] W. Klösgen and M. May. Census Data Mining - An Application. In D. Malerba and P. Brito, editors, *Proc. Workshop Mining Official Data, 6th European Conference, PKDD 2002*, Helsinki, 2002. Helsinki Univ. Printing House.
- [69] W. Klösgen and M. May. Spatial Subgroup Mining Integrated in an Object-Relational Spatial Database. In T. Elomaa, H. Mannila, and H. Toivonen, editors, *Proc. PKDD*, volume 2431 of *LNCS*, pages 275–286, Heidelberg, Germany, 2002. Springer Verlag.
- [70] P. Klügl, M. Toepfer, F. Lemmerich, A. Hotho, and F. Puppe. Collective Information Extraction with Context-Specific Consistencies. In *Proc. ECML/PKDD*, volume 7523 of *LNAI*, pages 728–743, 2012.
- [71] A. Knobbe, A. Feelders, and D. Leman. *Data Mining: Foundations and Intelligent Paradigms. Vol. 2*, chapter Exceptional Model Mining, pages 183–198. Springer Verlag, Heidelberg, Germany, 2011.
- [72] A. Knobbe and E. Ho. Pattern teams. In *Knowledge Discovery in Databases: PKDD 2006*, pages 577–584, Heidelberg, Germany, 2006. Springer Verlag.
- [73] A. J. Knobbe, B. Cremilleux, J. Fürnkranz, and M. Scholz. From Local Patterns to Global Models: The LeGo Approach to Data Mining. In *From Local Patterns to Global Models: Proceedings of the ECML/PKDD-08 Workshop (LeGo-08)*, pages 1 – 16, 2008.
- [74] R. M. Konijn, W. Duivesteijn, M. Meeng, and A. J. Knobbe. Cost-Based Quality Measures in Subgroup Discovery. In J. Li, L. Cao, C. Wang, K. C. Tan, B. Liu, J. Pei, and V. S. Tseng, editors, *Trends and Applications in Knowledge Discovery and Data Mining - PAKDD 2013 International Workshops: DMAApps, DANTh, QIMIE, BDM, CDA, CloudSD, Gold Coast, QLD, Australia, April 14-17, 2013, Revised Selected Papers*, volume 7867 of *Lecture Notes in Computer Science*, pages 404–415. Springer, 2013.
- [75] L. Lakhal and G. Stumme. Efficient Mining of Association Rules Based on Formal Concept Analysis. In *Formal Concept Analysis*, pages 180–195, Heidelberg, Germany, 2005. Springer.

- [76] N. Lavrac, B. Cestnik, D. Gamberger, and P. Flach. Decision Support Through Subgroup Discovery: Three Case Studies and the Lessons Learned. *Machine Learning*, 57(1-2):115–143, October 2004.
- [77] N. Lavrac, P. A. Flach, B. Kasek, and L. Todorovski. Rule Induction for Subgroup Discovery with CN2-SD. In N. L. D. M. M. Bohanec, B. Kasek, editor, *ECML/PKDD'02 workshop on Integration and Collaboration Aspects of Data Mining, Decision Support and Meta-Learning*, pages 77—87, 2002.
- [78] N. Lavrac and D. Gamberger. Relevancy in Constraint-based Subgroup Discovery. In H. M. Jean-Francois Boulicaut, Luc de Raedt, editor, *Constraint-based Mining and Inductive Databases*, volume 3848 of *LNCS*. Springer Verlag, Heidelberg, Germany, 2004.
- [79] N. Lavrac, D. Gamberger, and P. Flach. Subgroup Discovery for Actionable Knowledge Generation: Shortcomings of Classification Rule Learning and the Lessons Learned. In N. Lavrac, H. Motoda, and T. Fawcett, editors, *Proc. ICML 2002 workshop on Data Mining: Lessons Learned*, July 2002.
- [80] N. Lavrac, B. Kavsek, P. Flach, and L. Todorovski. Subgroup Discovery with CN2-SD. *Journal of Machine Learning Research*, 5:153–188, 2004.
- [81] D. Leman, A. Feelders, and A. Knobbe. Exceptional Model Mining. In *Proc. ECML/PKDD*, volume 5212 of *Lecture Notes in Computer Science*, pages 1–16. Springer, 2008.
- [82] F. Lemmerich, M. Becker, and M. Atzmueller. Generic Pattern Trees for Exhaustive Exceptional Model Mining. In *Proc. ECML/PKDD*, volume 7524 of *LNAI*, pages 277–292, Heidelberg, Germany, 2012. Springer Verlag.
- [83] F. Lemmerich, M. Becker, and F. Puppe. Difference-Based Estimates for Generalization-Aware Subgroup Discovery. In *Proc. ECML/PKDD*, volume 8190 of *Lecture Notes in Computer Science*, pages 288–303, Heidelberg, Germany, 2013. Springer Verlag.
- [84] F. Lemmerich and F. Puppe. A Critical View on Automatic Significance-Filtering in Pattern Mining. In *In Proc. Workshop Statistically Sound Data Mining, ECML/PKDD 2014*, Nancy, France.
- [85] F. Lemmerich and F. Puppe. Local Models for Expectation-Driven Subgroup Discovery. In *Proc. IEEE International Conference on Data Mining (ICDM)*, pages 360–369, Washington, DC, USA, 2011. IEEE.
- [86] F. Lemmerich, M. Rohlfs, and M. Atzmueller. Fast Discovery of Relevant Subgroup Patterns. In *Proc. 23rd International FLAIRS Conference*, pages 428–433, Palo Alto, CA, USA, 2010. AAAI Press.
- [87] H. Li, Y. Wang, D. Zhang, M. Zhang, and E. Y. Chang. PFP: Parallel Fp-Growth for Query Recommendation. In *Proc. RecSys*, pages 107–114, New York, NY, USA, 2008. ACM.

- [88] J. Li, J. Liu, H. Toivonen, K. Satou, Y. Sun, and B. Sun. Discovering Statistically non-redundant Subgroups. *Knowledge-Based Systems*, 67(0):315 – 327, 2014.
- [89] B. T. Lowerre. *The Harpy Speech Recognition System*. PhD thesis, Pittsburgh, PA, USA, 1976. AAI7619331.
- [90] J. M. Luna, J. R. Romero, C. Romero, and S. Ventura. Discovering Subgroups by Means of Genetic Programming. In K. Krawiec, A. Moraglio, T. Hu, A. S. Etaner-Uyar, and B. Hu, editors, *EuroGP*, volume 7831 of *Lecture Notes in Computer Science*, pages 121–132. Springer, 2013.
- [91] A. Magalhães and P. J. Azevedo. Contrast Set Mining in Temporal Databases. *Expert Systems*, pages n/a–n/a, 2014.
- [92] M. Mampaey, S. Nijssen, A. Feelders, and A. J. Knobbe. Efficient Algorithms for Finding Richer Subgroup Descriptions in Numeric and Nominal Data. In *Proc. IEEE International Conference on Data Mining (ICDM)*, pages 499–508, Washington, DC, USA, 2012. IEEE Computer Society.
- [93] M. Meeng, W. Duivesteijn, and A. J. Knobbe. ROCsearch - An ROC-guided Search Strategy for Subgroup Discovery. In M. J. Zaki, Z. Obradovic, P.-N. Tan, A. Banerjee, C. Kamath, and S. Parthasarathy, editors, *Proc. SIAM International Conference on Data Mining*, pages 704–712. SIAM, 2014.
- [94] M. Meeng and A. J. Knobbe. Flexible Enrichment with Cortana – Software Demo. In *Proc. Benelearn*, pages 117–119, 2011.
- [95] T. Mielikäinen. *Summarization Techniques for Pattern Collections in Data Mining*. PhD thesis, University of Helsinki, May 2005.
- [96] I. Mierswa, M. Wurst, R. Klinkenberg, M. Scholz, and T. Euler. YALE: Rapid Prototyping for Complex Data Mining Tasks. In *Proc. ACM SIGKDD, KDD '06*, pages 935–940, New York, NY, USA, 2006. ACM.
- [97] K. Moreland and K. Trueemper. Discretization of Target Attributes for Subgroup Discovery. In P. Perner, editor, *MLDM*, volume 5632 of *Lecture Notes in Computer Science*, pages 44–52. Springer, 2009.
- [98] K. Morik, J.-F. Boulicaut, and A. Siebes, editors. *Local Pattern Detection*. Springer Verlag, 2005.
- [99] M. Müller, R. Rosales, H. Steck, S. Krishnan, B. Rao, and S. Kramer. Subgroup Discovery for Test Selection: A Novel Approach and its Application to Breast Cancer Diagnosis. In A. S. J.-F. B. Niall M. Adams, Celine Robardet, editor, *Advances in Intelligent Data Analysis VIII, 8th International Symposium on Intelligent Data Analysis, IDA 2009*, volume 5772 of *Lecture Notes in Computer Science*, pages 119–130. Springer, 2009.

- [100] M. Natu and G. Palshikar. Interesting Subset Discovery and Its Application on Service Processes. In K. Yada, editor, *Data Mining for Service*, volume 3 of *Studies in Big Data*, pages 245–269. Springer Berlin Heidelberg, 2014.
- [101] P. K. Novak, N. Lavrač, and G. I. Webb. Supervised Descriptive Rule Discovery: A Unifying Survey of Contrast Set, Emerging Pattern and Subgroup Mining. *J. Mach. Learn. Res.*, 10:377–403, 2009.
- [102] N. Pasquier, Y. Bastide, R. Taouil, and L. Lakhal. Discovering Frequent Closed Itemsets for Association Rules. In C. Beeri and P. Buneman, editors, *Proc. 7th Intl. Conference on Database Theory (ICDT 99)*, volume 1540 of *Lecture Notes in Computer Science*, pages 398–416. Springer, 1999.
- [103] G. Piatetsky-Shapiro and C. J. Matheus. The Interestingness of Deviations. In *Proc. AAAI-94 Workshop on Knowledge Discovery in Databases (KDD-94)*, pages 25–36. ACM Press, New York, 1994.
- [104] B. Pieters, A. Knobbe, and S. Dzeroski. Subgroup Discovery in Ranked Data, with an Application to Gene Set Enrichment. In *Proc. Preference Learning Workshop (PL2010) at ECML/PKDD*, 2010.
- [105] F. Puppe, M. Atzmueller, G. Buscher, M. Huettig, H. Lührs, and H.-P. Buscher. Application and Evaluation of a Medical Knowledge-System in Sonography (SonoConsult). In *Proc. 18th European Conference on Artificial Intelligence (ECAI 2008)*, pages 683–687, 2008.
- [106] J. Roberto J. Bayardo. Efficiently Mining Long Patterns from Databases. In *SIGMOD '98: Proc. of the 1998 ACM SIGMOD Intl. Conference on Management of data*, pages 85–93, New York, NY, USA, 1998. ACM Press.
- [107] D. Rodríguez, R. Ruiz, J. C. Riquelme, and J. S. Aguilar-Ruiz. Searching for Rules to Detect Defective Modules: A Subgroup Discovery Approach. *Inf. Sci.*, 191:14–30, 2012.
- [108] W. Romao, A. A. Freitas, and I. M. de S. Gimenes. Discovering Interesting Knowledge from a Science & Technology Database with a Genetic Algorithm. *Applied Soft Computing*, 4(2):121–137, May 2004.
- [109] T. Scheffer and S. Wrobel. A Scalable Constant-Memory Sampling Algorithm for Pattern Discovery in Large Databases. In *Proc. PKDD*, pages 397–409, Heidelberg, Germany, 2002. Springer Verlag.
- [110] T. Scheffer and S. Wrobel. Finding the Most Interesting Patterns in a Database Quickly by Using Sequential Sampling. *Journal of Machine Learning Research*, 3:833–862, 2002.
- [111] C. Sáez, P. Rodrigues, J. Gama, M. Robles, and J. García-Gómez. Probabilistic Change Detection and Visualization Methods for the Assessment of Temporal Stability in Biomedical Data Quality. *Data Mining and Knowledge Discovery*, pages 1–26, 2014.

- [112] L. G. Shapiro and R. M. Haralick. A Metric for Comparing Relational Descriptions. *IEEE Trans. Pattern Anal. Mach. Intell.*, 7(1):90–94, 1985.
- [113] H. Toivonen. Sampling Large Databases for Association Rules. In T. M. Vijayaragaman, A. P. Buchmann, C. Mohan, and N. L. Sarda, editors, *Proc. 1996 Intl. Conference on Very Large Data Bases*, pages 134–145. Morgan Kaufman, 1996.
- [114] A. Tuzhilin. *Handbook of Data Mining and Knowledge Discovery*, chapter 19.2.2: Usefulness, Novelty, and Integration of Interestingness Measures. Oxford University Press, New York, 2002.
- [115] M. van Leeuwen. Interactive Data Exploration using Pattern Mining. In A. Holzinger and I. Jurisica, editors, *Lecture Notes in Computer Science*. Springer Berlin Heidelberg, 2014.
- [116] M. van Leeuwen and A. J. Knobbe. Diverse Subgroup Set Discovery. *Data Min. Knowl. Discov.*, 25(2):208–242, 2012.
- [117] A. Vavpetic and N. Lavrac. Semantic Subgroup Discovery Systems and Workflows in the SDM-Toolkit. *Comput. J.*, 56(3):304–320, 2013.
- [118] A. Vavpetic, V. Podpecan, and N. Lavrac. Semantic Subgroup Explanations. *J. Intell. Inf. Syst.*, 42(2):233–254, 2014.
- [119] G. I. Webb. Discovering Associations with Numeric Variables. In *Proc. ACM SIGKDD Intl. Conf. on Knowledge Discovery and Data Mining*, KDD '01, pages 383–388, New York, NY, USA, 2001. ACM.
- [120] S. Wrobel. An Algorithm for Multi-Relational Discovery of Subgroups. In *Proc. 1st Europ. Symp. Principles of Data Mining and Knowledge Discovery*, pages 78–87, Heidelberg, Germany, 1997. Springer Verlag.
- [121] M. J. Zaki. Efficient Enumeration of Frequent Sequences. In *CIKM '98: Proceedings of the seventh international conference on Information and knowledge management*, pages 68–75, New York, NY, USA, 1998. ACM.
- [122] A. Zimmermann and L. D. Raedt. Cluster-Grouping: From Subgroup Discovery to Clustering. In *Proc. ECML*, pages 575–577, 2004.
- [123] A. Zimmermann and L. D. Raedt. CorClass: Correlated Association Rule Mining for Classification. In E. Suzuki and S. Arikawa, editors, *Proc. 7th Intl. Conference on Discovery Science*, volume 3245 of *Lecture Notes in Computer Science*, pages 60–72, 2004.