

Exploring Data using Patterns: A Survey

Vargha Dadvar^a, Lukasz Golab^a, Divesh Srivastava^b

^a*University of Waterloo, 200 University Avenue, Waterloo, Ontario, N2L3G1, Canada*

^b*AT&T Chief Data Office, 1 AT&T Way, Bedminster, NJ, 07921, USA*

Abstract

We present a survey of data exploration methods that extract multidimensional patterns from datasets consisting of dimension and measure attributes. These patterns are designed to summarize common properties of tuples sharing the same values of the measure attributes. We review motivating applications, we provide a categorization of the characteristics of patterns produced by various solutions to this problem, we categorize and experimentally evaluate commonly used performance optimizations, and we suggest directions for future research.

Keywords: Data exploration, Data summarization, Data explanation, Data cube, Pattern mining

1. Introduction

2 Data volumes have been growing rapidly in recent years. As a result,
3 data-intensive methods are now common in many contexts, including busi-
4 ness, science, and public governance. This motivates the need for tools that
5 allow users who are not necessarily data management experts to explore large
6 datasets. Such tools range from data visualization and aggregation [14, 19]
7 to flexible search interfaces such as keyword search in structured databases
8 [26].

9 In this paper, we focus on the exploration of datasets containing dimen-
10 sion attributes and binary or numeric measure attributes. In traditional
11 business datasets, dimension attributes often describe products or employees,
12 and measure attributes indicate sales totals or salaries. In Internet-of-Things
13 (IoT) and infrastructure monitoring, dimension attributes may describe de-
14 vice properties and measure attributes correspond to performance statistics.

15 In Web datasets, dimension attributes may describe products, with aggregate
16 user ratings as measure attributes. Additionally, in any of these applica-
17 tions, derived measure attributes may exist, e.g., a binary attribute denoting
18 whether a given record was determined to be an outlier or to contain an
19 error.

20 The *data cube* [12] has traditionally been used to explore these kinds of
21 datasets, by allowing users to aggregate, roll-up and drill-down using vari-
22 ous subsets of group-by attributes. However, in large-scale databases, the
23 data cube may be very large and may not immediately reveal interesting
24 trends. As argued in a recent vision paper by Vassiliadis and Marcel [22],
25 next-generation Business Intelligence (BI) tools require new concepts and op-
26 erators to help users discover information, among them those for automatic
27 mining of models and patterns. This motivates the need for richer data ex-
28 ploration tools that can operate over data cubes and other multi-dimensional
29 data models.

30 We observe that recent work on exploring multi-dimensional and OLAP
31 datasets proposed a variety of methods to identify interesting fragments of
32 the data, described using combinations of values of the dimension attributes
33 [1, 2, 5, 7, 9, 11, 15, 18, 20, 21, 24, 25]. These value combinations are
34 referred to as patterns. Below, we give examples to show that this class of
35 methods provides *interpretable* summaries and explanations of trends in the
36 data, aligning well with the anticipated needs of next-generation BI and data
37 exploration tools.

38 Table 1 shows a flight dataset that will serve as a running example. For
39 each flight, the dataset includes a record id, followed by three dimension
40 attributes, Day of the week, flight Origin and flight Destination, as well as
41 two measure attributes, a numeric attribute denoting how late the flight
42 was and a binary attribute denoting whether the flight was full. First, note
43 that the following two patterns summarize most of the tuples corresponding
44 to full flights: (Day=*, Origin=*, Dest=London) and (Day=*, Origin=*,
45 Dest=Frankfurt), where a star matches all the values of the corresponding
46 attribute. In other words, full flights are mainly those that arrive in London
47 or Frankfurt. Next, consider the pattern (Day=Mon, Origin=*, Dest=*),
48 corresponding to flights scheduled on Mondays. This pattern may be inter-
49 esting because none of these flights are full, which differs significantly from
50 the fraction of full flights in the entire table. Finally, suppose a data analyst
51 is surprised by the high average delay of flights in Table 1. Here, the pattern
52 (Day=*, Origin=*, Dest=London) can serve as a potential explanation since

Table 1: A flight dataset

id	Day	Origin	Dest.	Late	Full
1	Fri	SF	London	20	1
2	Fri	London	LA	16	1
3	Sun	Tokyo	Frankfurt	10	1
4	Sun	Chicago	London	15	1
5	Sat	Beijing	Frankfurt	13	1
6	Sat	Frankfurt	London	19	1
7	Tue	Chicago	LA	5	0
8	Wed	London	Chicago	6	0
9	Thu	SF	Frankfurt	15	1
10	Mon	Beijing	SF	4	0
11	Mon	SF	London	7	0
12	Mon	SF	Frankfurt	5	0
13	Mon	Tokyo	Beijing	6	0
14	Mon	Frankfurt	Tokyo	4	0

53 flights arriving in London have some of the longest delays in the table.

54 We survey methods that automatically identify such patterns. Users may
55 then inspect the patterns and explore tuples covered by the patterns. They
56 may then extract patterns corresponding to smaller subsets of the data found
57 to be interesting in the earlier exploration step, and so on. Of course, there
58 are other approaches to data exploration and analysis, such as visualization
59 and fitting models to predict the measure attributes based on the dimension
60 attributes (we will comment on the role of prediction models in data cube
61 exploration in Section 3.2). We restrict the scope of this survey to pattern-
62 based exploration and refer the reader to [19] for a survey of data visualization
63 methods. Furthermore, visualization may be thought of as complementary
64 to pattern-based exploration. For example, interesting patterns may be iden-
65 tified, followed by visualizing the statistical properties of the corresponding
66 tuples.

67 We make the following contributions towards an understanding of pattern-
68 based data exploration methods.

- 69 1. We survey recent work on data exploration using multi-dimensional
70 patterns and propose a categorization based on the properties of pat-
71 terns suggested for exploration: *coverage*, *contrast*, and *information*.

- 72 2. We categorize and experimentally evaluate performance optimizations
73 frequently used in pattern-based data exploration: *top-down pruning*,
74 *row pruning*, *column pruning* and *parallel processing*. Notably, we ex-
75 perimentally demonstrate that some performance optimizations origi-
76 nally proposed for one method are also effective when applied to other
77 methods. This observation should be of interest to researchers working
78 in this area since it points out existing performance optimizations that
79 may be beneficial for newly developed data exploration techniques.
- 80 3. We suggest open problems for future research.

81 In the remainder of this paper, we present the required background in
82 Section 2, we categorize existing solutions in Section 3, we classify and ex-
83 perimentally evaluate performance optimizations in Sections 4 and 5, respec-
84 tively, and we conclude in Section 6 with directions for future work.

85 A short version of this survey was presented at DOLAP 2021 [10]. This
86 extended version includes two new sections with a classification (Section 4)
87 and an experimental study (Section 5) of performance optimizations.

88 2. Background

89 We are given a dataset S with a set D of dimension attributes and a set
90 M of measure attributes (also referred to as outcomes in some prior work [6]).
91 Let D_1, D_2, \dots, D_d be the d dimension attributes and let M_1, M_2, \dots, M_m be
92 the m measure attributes. For now, we assume, as in the majority of previous
93 work, that the dimension attributes are categorical, and we will comment on
94 ordinal and numeric dimension attributes in Section 6. Measure attributes
95 may be binary or numeric.

96 Let $dom(D_i)$ be the active domain of the i th dimension attribute. A pat-
97 tern p is a tuple from $dom(D_1) \cup \{*\} \times \dots \times dom(D_d) \cup \{*\}$, i.e., from the
98 data cube over the dimension attributes, with ‘*’ denoting all possible values
99 of that attribute. A tuple $t \in S$ matches p , denoted by $t \simeq p$, if $p[D_j] =$
100 ‘*’ or $t[D_j] = p[D_j]$ for each dimension attribute D_j . For example, tuple
101 4 from Table 1 matches the patterns $(*, *, *)$ and $(*, *, London)$, but does
102 not match the pattern $(Fri, *, *)$; to simplify the notation, we drop attribute
103 names from patterns and only include attribute values. Some approaches
104 (e.g., [11]) support richer patterns with disjunctions and dimension hierar-
105 chies. Other methods described in this paper can also support disjunctions
106 and dimension hierarchies, at the cost of a larger search space of candidate

107 patterns. However, for simplicity of presentation, in the remainder of this pa-
 108 per, we will illustrate the methods using simple conjunctive patterns without
 109 hierarchies.

110 Let $sup(p)$ be the support of p in S , i.e., the number of tuples matching p ,
 111 and let $sup_r(p)$ be the number of tuples matching p and satisfying a predicate
 112 r over the measure attributes. For example, $sup(*, *, London) = 4$ and
 113 $sup_{Full=0}(*, *, London) = 1$. Furthermore, let $\theta_r(p) = \frac{sup_r(p)}{sup(p)}$, which is the
 114 fraction of tuples matching p that also satisfy the predicate r . For example,
 115 $\theta_{Full=1}(*, *, London) = \frac{3}{4}$.

116 Let $sum^{M_i}(p)$ be the sum of the values of a measure attribute M_i over all
 117 the tuples matching p . Let $sum_r^{M_i}(p)$ be the sum of the values of a measure
 118 attribute M_i over all the tuples matching p and satisfying a predicate r over
 119 the measure attributes. For example, $sum^{Late}(*, *, London) = 20 + 15 + 19 +$
 120 $7 = 61$ and $sum_{Full=0}^{Late}(*, *, London) = 7$.

121 We survey solutions to the following data exploration problem: given a
 122 dataset S , produce a set or a list of patterns P over the dimension attributes
 123 of S , as defined above, that summarize common properties of tuples sharing
 124 the same values of the measure attribute(s). The number of patterns in
 125 P should be limited to direct the user’s attention to the most important
 126 or interesting regions of the data. This limit may be set explicitly by the
 127 user (in terms of the maximum number of patterns in P) or implicitly by
 128 retrieving the fewest possible patterns that jointly satisfy some property such
 129 as covering some fraction of the data.

130 This data exploration problem has the following applications.

- 131 • *Explaining the results of aggregate queries.* Suppose a data analyst
 132 issues the following query over Table 1: SELECT SUM(Late) FROM
 133 S. Suppose the analyst wishes to understand why the result, of 145, is so
 134 high. Here, interesting patterns are those which cover tuples that make
 135 a significant contribution to the result, i.e., those with a high $sum^{Late}()$
 136 such as $(*, *, London)$. The analyst may then zoom into flights landing
 137 in London and investigate potential reasons for the lengthy delays.
- 138 • *Analyzing outliers and data quality issues.* Suppose we have a binary
 139 measure attribute denoting whether a given tuple contains an error or is
 140 an outlier. This attribute could be created manually by domain experts
 141 or automatically by identifying tuples that violate data quality rules
 142 or deviate from the expected distribution. We may wish to produce

143 patterns that summarize the properties of erroneous tuples to help
144 determine the root cause of data quality problems.

- 145 • *Feature selection and explainable AI.* Before building prediction mod-
146 els, a data scientist may explore interesting patterns to understand
147 which dimension attributes are related to the measure attribute that is
148 to be predicted. Furthermore, suppose a data analyst wants to under-
149 stand how a black-box model makes classification decisions. Here, the
150 dimension attributes are the features given to the model as input, and,
151 as the measure attribute, the analyst records the predictions made by
152 the model. The analyst may then want to find interesting patterns that
153 explain the prediction decisions. For example, in Table 1, the pattern
154 (*Mon, *, **) is associated with tuples having $Full = 0$, suggesting that
155 flights scheduled on Mondays are usually not full¹.

156 3. Solutions

157 In this section, we provide a categorization of previous work on data
158 exploration using multi-dimensional patterns based on the pattern proper-
159 ties and ranking strategies used for pattern selection. We categorize these
160 properties into three types: those focusing on *coverage*, *contrast* and *infor-*
161 *mation*. Table 2 categorizes the surveyed methods and lists their motivating
162 applications, as mentioned in the corresponding papers.

163 Additionally, Table 3 lists the inputs and outputs of the surveyed meth-
164 ods. As explained in Section 2, these methods operate over datasets with
165 multiple dimension attributes and a measure attribute that needs to be cov-
166 ered, contrasted or explained. MRI and Smart Drilldown also require a
167 function that assigns pattern weights. CAPE is unique in its inputs in that
168 it requires a specific pattern that can be thought of as a starting point for
169 further exploration. In terms of outputs, most methods produce k best pat-
170 terns according to some properties. Methods based on contrast identify k
171 patterns with the highest contrast scores (details in Section 3.2). Methods
172 based on coverage and information typically use greedy heuristics to solve an

¹Model explanations may be global (to summarize how classification decisions are made) or local (to explain why a specific example was classified in a particular way); see [13] for a survey. The methods discussed in this paper are examples of global explanations.

Table 2: Methods surveyed

Method	Approach	Applications
CAPE [18]	Contrast	Explaining queries
Data Auditor [9]	Coverage	Data quality analysis
Data X-ray [24]	Contrast	Data quality analysis
DIFF [2]	Contrast	Outlier analysis
Explanation tables [7]	Information	Feature selection
Macrobase [1]	Contrast	Outlier analysis
MRI [5]	Coverage	Explaining queries over product ratings
RSExplain [20]	Contrast	Explaining queries
Scorpion [25]	Contrast	Outlier analysis
Shrink [11]	Information	Explaining queries
Smart Drilldown [15]	Coverage	Explaining queries, data cube exploration
SURPRISE [21]	Information	Explaining queries

173 underlying NP-hard problem related to maximizing the coverage or informa-
174 tion content of the selected patterns. Thus, their outputs are *approximately*
175 optimal with respect to the associated coverage, weighted coverage, or infor-
176 mation metric.

177 3.1. Methods Based on Coverage

178 The goal of these methods is to identify patterns that *cover* tuples of
179 interest; this may refer to covering tuples in the entire dataset, covering tuples
180 with a given value of a measure attribute, or covering tuples that contribute
181 to the result of a query. We discuss three coverage-based methods: Data
182 Auditor [9], MRI [5], and Smart Drilldown [15].

183 3.1.1. Method Details

184 Suppose we want to cover tuples having $Full = 1$ in Table 1 to sum-
185 marize the characteristics of full flights. A simple coverage-oriented ap-
186 proach is to sort the patterns according to $sup_{Full=1}()$ and output the top-
187 ranking patterns. Ignoring $(*, *, *)$, which always covers everything but is
188 not useful in data exploration, the top candidates are $(*, *, London)$ and
189 $(*, *, Frankfurt)$, which cover three full flights each, followed by the follow-
190 ing patterns that cover two such tuples each: $(Fri, *, *)$, $(Sun, *, *)$, $(Sat, *, *)$
191 and $(*, SF, *)$.

Table 3: Inputs and outputs of methods surveyed

Method	Input	Output
CAPE	Dataset, a pattern p	k patterns with the highest counterbalance score w.r.t. p
Data Auditor	Dataset, coverage thresholds $\theta_r(p)$	Fewest patterns that satisfy the coverage thresholds
Data X-ray	Dataset	k patterns with the highest diagnosis cost
DIFF	Dataset, contrast metric, support threshold	k patterns (that satisfy the support threshold) with the highest contrast
Explanation tables	Dataset	k most informative patterns
Macrobase	Dataset	k patterns with the highest risk ratio
MRI	Dataset, coverage threshold, pattern weighting function	k patterns that cover the required fraction of tuples with a minimal sum of weights
RSExplain	Dataset	k patterns with the highest intervention score
Scorpion	Dataset	k pattern with the highest influence score
Shrink	Dataset	k patterns that best summarize the distribution of the measure attribute
Smart Drilldown	Dataset, pattern weighting function	k patterns that maximize the product of coverage and sum of weights
SURPRISE	Dataset	k most informative patterns

192 There are two problems with this simple approach: there may be many
193 patterns with a nonzero $sup_{Full=1}(p)$, and some of these patterns may also
194 cover tuples with other values of the measure attribute (here, $Full = 0$). To
195 reduce the size of the output and to ensure that each pattern co-occurs with
196 the specified value of the measure attribute, Data Auditor solves the following
197 set cover problem. Continuing with our example, Data Auditor requires

198 a minimum threshold for $\theta_{Full=1}(p)$, i.e., the fraction of tuples covered by
 199 p that correspond to full flights. Suppose we require $\theta_r(p) \geq 0.75$. This
 200 threshold defines the candidate sets for the set cover problem, i.e., all the
 201 patterns p with $\theta_r(p) \geq 0.75$. The set cover objective is to select the fewest
 202 such patterns that together cover a specified fraction of tuples in S having
 203 $Full = 1$. Suppose this coverage fraction, which would again be set by the
 204 user, is 0.5. The goal is then to find the fewest patterns, as defined above,
 205 to cover half the full flights. Alternatively, Data Auditor may produce k
 206 patterns that (approximately) maximize coverage.

207 The set cover problem is NP-hard, and Data Auditor uses the standard
 208 greedy heuristic that achieves a logarithmic approximation ratio in the size of
 209 the solution: it iteratively chooses the pattern that covers the most uncovered
 210 tuples (having the desired value of the measure attribute), until the desired
 211 fraction of such tuples has been covered (or until k patterns have been cho-
 212 sen). In our example, the first pattern added to P is either $(*, *, London)$ or
 213 $(*, *, Frankfurt)$ – they each cover three full flights and their $\theta_{Full=1}(p)$ val-
 214 ues are 0.75 each. Suppose the set cover algorithm selects $(*, *, Frankfurt)$.
 215 Since there are seven full flights in the dataset and our coverage threshold is
 216 0.5, we need to cover one more full flight. In the next iteration, the pattern
 217 that (has $\theta_{Full=1}(p) \geq 0.75$ and) covers the most remaining full flights is
 218 $(*, *, London)$ and the algorithm terminates, with P consisting of these two
 219 patterns.

220 Data Auditor solves a set cover problem in which patterns are prioritized
 221 by their coverage (of tuples that have a given value of the measure attribute
 222 and have not yet been covered). The other two coverage-based methods take
 223 into account pattern weights in addition to coverage, as described below.

224 The next method, Smart Drilldown, produces k patterns. In each of the
 225 k iterations of the algorithm, the chosen pattern maximizes the following
 226 objective: the number of tuples not yet covered multiplied by the *weight* of
 227 the pattern corresponding to some measure of interestingness. One simple
 228 weighting function proposed in [15] is the number of non-star values in the
 229 pattern; i.e., more specific patterns are considered to be more desirable. In
 230 Table 1 for example, this weighting function prefers $(Tue, Chicago, LA)$ over
 231 $(Tue, *, *)$ – both of these patterns have a support of one, but the former has
 232 more non-star values.

233 The third method based on coverage, MRI, finds k patterns that cover
 234 a user-specified fraction of the data and satisfy additional properties related
 235 to the variance of the measure attribute within each pattern. Here, the

236 weight of a pattern corresponds to this notion of variance, and patterns
237 with smaller weights are preferred. The motivating example for MRI was to
238 explain queries over product reviews, with the dimension attributes corre-
239 sponding to information about the reviewers (such as their gender and age)
240 and the numeric measure attribute corresponding to the average rating. Min-
241 imizing variance amounts to returning patterns (having high coverage and)
242 describing reviewers with similar opinions. For example, when applying MRI
243 to Table 1 with *Late* as the measure attribute, the pattern $(*, *, Frankfurt)$
244 is preferred over $(*, *, London)$. Both patterns cover four tuples, but the
245 former has a lower variance of the *Late* attribute within the covered tuples.

246 Both Smart Drilldown and MRI prioritize patterns with high coverage (of
247 tuples that have not yet been covered) and take pattern weights into account.
248 However, while Smart Drilldown uses a modified version of the greedy set
249 cover heuristic also used by Data Auditor, MRI uses a *hill climbing* heuristic.
250 First, MRI selects k patterns at random. Next, it makes small changes to the
251 patterns, if necessary, to ensure that the user-specified coverage threshold is
252 satisfied. A small change to increase coverage may correspond to replacing
253 one attribute value in a pattern with a ‘*’. Finally, MRI makes a second
254 round of small changes to the patterns in an attempt to reduce the sum of
255 their weights.

256 3.1.2. Pros and Cons

257 One advantage of coverage-based methods is conciseness: by design, they
258 identify (approximately) the fewest patterns that cover the desired fraction of
259 tuples of interest, or they produce k patterns that (approximately) maximize
260 coverage. On the other hand, in Data Auditor, some trial-and-error may be
261 required on the user’s part to select good values for the two required thresh-
262 olds. For example, a high value of θ will ensure that each selected pattern
263 co-occurs mainly with the specified value of the measure attribute, but will
264 disqualify more patterns from consideration, possibly leading to lower cover-
265 age. Similarly, MRI and Smart Drilldown require users to set parameters for
266 coverage and other pattern properties.

267 3.2. Methods based on Contrast

268 This group of methods includes CAPE [18], Data X-ray [24], DIFF [2],
269 Macrobase [1], RSExplain [20] and Scorpion [25]. Contrast-based methods
270 often assume a binary measure attribute, and select patterns co-occurring
271 with one value of the measure attribute but not the other. These patterns

272 reflect the contrast between tuples having different values of the measure
273 attribute.

274 One could argue that interpretable classifiers such as decision trees and
275 rule-based methods (see, e.g., [16]) can also be used for contrast-based data
276 exploration. These methods identify patterns of values of the feature at-
277 tributes that have high discriminative power in terms of the class variable
278 (in our case, the binary measure attribute). These patterns are therefore
279 likely to provide contrast as well. However, classification algorithms usu-
280 ally focus on out-of-sample predictive power and include optimizations such
281 as rule pruning to avoid overfitting. On the other hand, the methods cov-
282 ered in this survey focus explicitly on identifying a concise set of interesting
283 fragments of the data for user exploration.

284 Contrast-based methods can also explain the results of aggregate queries.
285 Consider the following query over Table 1: `SELECT SUM(Late) FROM S`
286 `WHERE Full=1`. Here, one measure attribute corresponds to the quantity
287 being aggregated. We then set the other (binary) measure attribute to one for
288 all tuples that participate in the query (i.e., tuples that match the `WHERE`
289 predicate), and we select patterns of tuples that contribute to the result of
290 the query (but would not contribute had the query been issued against the
291 other tuples in the dataset).

292 3.2.1. Method Details

293 DIFF is a recent solution that generalizes earlier contrast-based methods,
294 including Data X-ray, Macrobase, RSExplain and Scorpion. The authors of
295 DIFF observe that these methods all use a similar algorithmic framework but
296 different contrast metrics. DIFF supports these different contrast metrics.
297 To summarize the pattern mining framework used in DIFF to generalize
298 prior work, a contrast metric is calculated for each candidate pattern and the
299 highest-ranking patterns are returned. DIFF additionally implements several
300 performance optimizations that will be discussed in Section 4. Below, we give
301 several examples of contrast metrics originally used in earlier methods and
302 now supported by DIFF. We again use the running example in Table 1, with
303 *Full* as the binary measure attribute and *Late* as the additional numeric
304 measure attribute when needed.

305 *Risk ratio* was originally used by Macrobase; a related metric called Di-
306 agnosis Cost is used by Data X-ray. It is the ratio of the following two prob-
307 abilities: 1) the probability that a tuple with a particular value is covered
308 by the given pattern, and 2) the probability that a tuple with this particular

309 value occurs outside this pattern. In our example,

$$risk_{Full=1}(p) = \frac{\theta_{Full=1}(p)}{\frac{sup_{Full=1}(*,*,*) - sup_{Full=1}(p)}{(sup_{Full=1}(*,*,*) - (sup_{Full=1}(p)) + (sup_{Full=0}(*,*,*) - sup_{Full=0}(p))}}.$$

310 For instance, $risk_{Full=1}(*, *, London) = \frac{0.75}{0.4} = 1.875$, and $risk_{Full=1}(*, SF, *) =$
 311 $\frac{0.5}{0.5} = 1$. This indicates that $(*, *, London)$ represents full flights better than
 312 $(*, SF, *)$.

Mean shift computes the ratio of the mean of the measure attribute values co-occurring with the two values of the binary measure attribute. In our example,

$$mean_{Full=1}^{Late}(p) = \frac{sum_{Full=1}^{Late}(p) / sup_{Full=1}(p)}{sum_{Full=0}^{Late}(p) / sup_{Full=0}(p)}.$$

313 For instance, $mean_{Full=1}^{Late}(*, *, London) = \frac{54/3}{7/1} = 2.57$, indicating that full
 314 flights to London have delays that are 2.57 times longer than non-full flights
 315 to London.

Intervention was originally used by RSExplain; a related metric called Influence is used by Scorpion. It measures the ratio of contribution towards the numeric measure attribute for tuples occurring with the different values of the binary measure attribute. In our example,

$$intervention_{Full=1}^{Late}(p) = \frac{sum_{Full=1}^{Late}(*, *, *) - sum_{Full=1}^{Late}(p)}{sum_{Full=0}^{Late}(*, *, *) - sum_{Full=0}^{Late}(p)}.$$

316 For instance, $intervention_{Full=1}^{Late}(*, *, London) = \frac{108-54}{37-7} = 1.8$. In other
 317 words, if flights to London were removed from the dataset then full flights
 318 would have delays on average 1.8 times longer than non-full flights. On the
 319 other hand, $intervention_{Full=1}^{Late}(*, SF, *) = \frac{108-35}{37-12} = 2.92$, meaning that re-
 320 moving flights departing from SF from the dataset would create a greater
 321 *contrast* between the delays of full and non-full flights.

322 Finally, we discuss CAPE. Given a specific pattern p as input, CAPE
 323 finds patterns whose tuples have measure attribute values that *counterbalance*
 324 those of p . Thus, instead of ranking patterns according to some contrast
 325 metric, CAPE ranks patterns according to a counterbalance metric with
 326 respect to p , and outputs the highest-ranking such patterns.

327 For example, suppose a user runs the following query on Table 1: SE-
 328 LECT AVG(Late) FROM T WHERE Day = 'Mon'. This query outputs

329 the value 5.2 and touches tuples identified by the pattern (*Mon*, *, *). This
330 average delay is lower than the average delay in the entire table, which is
331 10.4. The user may then input this pattern to CAPE and request patterns
332 that counterbalance the lower delays seen in this pattern. CAPE may then
333 return patterns such as (*Fri*, *, *) and (*Sat*, *, *), whose average delays are
334 higher than average (18 and 16, respectively). Thus, counterbalancing refers
335 to finding related patterns whose measure attribute values are “outliers” in
336 the other direction (in our example, higher than average) compared to the
337 input pattern (in our example, lower than average).

338 In CAPE, the pattern mining algorithm has two main modules.

339 The first module finds regression relationships in the data; then counter-
340 balancing can be used to find outliers with respect to these relationships. In
341 our example above, counterbalancing was based on a trivial regression rela-
342 tionship involving the *Late* attribute, namely that $AVG(Late)=10.5$. Since
343 the input pattern had lower than average delays, CAPE searched for patterns
344 having higher than average delays. To see an intuitive example of a more
345 complex regression pattern, paraphrased from [18], consider a database with
346 authors and their publications such as DBLP. In this example, a regression
347 pattern may hold for most authors such that they publish more papers over
348 time. Here, an unusual (with respect to the identified regression relationship)
349 pattern corresponding, say, to a lower-than-expected number of publications
350 in a given year for given author, can be counterbalanced by finding pat-
351 terns corresponding to years in which this author published more than the
352 expected number of papers.

353 After identifying regression relationships, the second module finds pat-
354 terns to counterbalance a given input pattern. Such patterns must satisfy
355 two objectives. First, they must be “outliers” in the opposite direction to
356 the input pattern with respect to some aggregate function over the measure
357 attribute. Second, they must be “close” to the input pattern in terms of the
358 values of the dimension attributes. In the DBLP example above, if the input
359 pattern refers to, say, year 2015, then patterns with similar years, say 2014
360 or 2016, would be preferred for counterbalancing.

361 3.2.2. *Pros and Cons*

362 By design, contrast-based methods are useful when exploring differences
363 between data subsets – something that coverage-based methods do not di-
364 rectly optimize for. On the other hand, contrast-based methods may not
365 guarantee concise results. One exception is Data X-ray, which performs a

Table 4: An explanation table of size four for the binary measure attribute Full

Day	Origin	Dest.	Full
*	*	*	0.5
Mon	*	*	0
	*	London	0.75
	*	Frankfurt	0.75

366 set-cover-like operation on the extracted patterns as a post-processing step
 367 to eliminate redundant patterns.

368 *3.3. Methods based on Information*

369 Finally, we discuss three methods that select patterns based on the infor-
 370 mation they provide about the distribution of the measure attribute: Explan-
 371 ation Tables [7], SURPRISE [21] and Shrink [11]. We show two examples of
 372 explanation tables, one with a binary measure attribute and one with a nu-
 373 meric measure attribute, followed by a greedy heuristic to construct (almost)
 374 optimal explanation tables, and an overview of SURPRISE and Shrink.

375 *3.3.1. Method Details*

376 Table 4 shows an explanation table of size four (i.e., containing four pat-
 377 terns) for the binary measure attribute Full based on Table 1. In addition
 378 to values of the dimension attributes, each explanation table pattern also in-
 379 cludes the fraction of matching tuples that have $Full = 1$. The first pattern
 380 in an explanation table is always the all-stars pattern, and, in this example,
 381 it indicates that half the flights in the entire dataset are full. The next pat-
 382 tern suggests that no flights on Mondays are full, and the last two patterns
 383 indicate that three-quarters of flights to London and Frankfurt are full.

384 In Table 5, we show an explanation table of size four for the measure
 385 attribute Late based on Table 1. Here, each pattern includes the average
 386 value of Late across its matching tuples. Again, we start with the all-stars
 387 pattern, which states that flights are 10.4 minutes late on average. The next
 388 pattern indicates that flights to London are 15.3 minutes late on average,
 389 and so on.

390 The greedy heuristic for constructing explanation tables used in [6, 7, 8]
 391 iteratively selects patterns that contain the most information about the dis-
 392 tribution of the measure attribute. To do so, the algorithm maintains a

Table 5: An explanation table of size four for the numeric measure attribute Late

Day	Origin	Dest.	Late
*	*	*	10.4
	*	London	15.3
Fri	*	*	18
Sat	*	*	16

393 *maximum-entropy* estimate of the distribution based only on the patterns
394 that have been added to the explanation table so far, and without assum-
395 ing any other information; recall that the entropy of a random variable X
396 with outcomes x_1 through x_n is defined as $E(X) = \sum_1^n -P(x_i) \log P(x_i)$.
397 The algorithm also keeps track of the distance between the estimated distri-
398 bution and the true distribution by computing their *Kullback-Leibler* (KL)
399 divergence. Given two random variables, X and Y , with the same space
400 of outcomes, x_1 through x_n and y_1 through y_n , respectively, the KL diver-
401 gence of their distributions is defined as $KL(X, Y) = \sum_1^n P(x_i) \log \frac{P(x_i)}{P(y_i)}$. To
402 quantify the information contained in a candidate pattern p , the algorithm
403 computes the reduction in KL-divergence if p were to be added to the expla-
404 nation table.

405 Returning to Table 4, the greedy algorithm starts by inserting the pattern
406 $(*, *, * || 0.5)$. At this point, knowing only this one piece of information, the
407 maximum-entropy estimate of the distribution of $Full$ is to assign $Full = 0.5$
408 to every tuple in Table 1². That is, knowing only that the average value of
409 Full in the entire table is 0.5, we obtain maximum entropy by assigning 0.5 to
410 each tuple. Next, it turns out that $(Mon, *, * || 0)$ gives the greatest reduction
411 in KL-divergence. Based on this new pattern, the maximum-entropy estimate
412 for tuples 10 through 14 in Table 1 changes to $Full = 0$. This revision causes
413 the estimates of the first nine tuples to change (from 0.5 to $\frac{7}{9}$) in order
414 to maintain consistency with the first pattern, which asserts that $Full =$
415 0.5 on average over the entire table. Given the updated maximum-entropy
416 estimate, the next pattern with the greatest reduction in KL-divergence is

² $Full$ is a binary attribute that can only be zero or one. However, for the purpose of measuring the divergence between the true and the estimated distributions, the maximum-entropy estimates are allowed to be real numbers between zero and one.

417 $(*, *, London || 0.75)$, and so on.

418 Similar reasoning can explain how Table 5 was created. The first pattern
419 asserts that flights are late by 10.4 minutes on average. Given this estimate,
420 $(*, *, London || 15.3)$ provides the most information about the distribution of
421 *Late*. The maximum-entropy estimate of *Late* is now updated accordingly,
422 That is, tuples 1, 4, 6 and 11, corresponding to flights to London, receive
423 an estimate of 15.3, and the remaining tuples receive an estimate of 8.4
424 to maintain consistency with the first pattern. The next most-informative
425 pattern is then selected, and so on.

426 SURPRISE is a similar method, whose goal is to identify surprising frag-
427 ments of a dataset where the distribution of the measure attribute is different
428 than what the user has seen so far. Suppose the user queries Table 1 and
429 finds out that flights are 10.4 minutes late on average. SURPRISE finds the
430 most informative *non-overlapping* and *contained* patterns, i.e., those which
431 lead to the greatest reduction in KL-divergence between the true distribu-
432 tion of *Late* and the maximum-entropy estimated distribution. Restricting
433 the output to such patterns makes it easier to update the estimated distri-
434 bution. In our example, the most informative pattern is $(*, *, London)$ and
435 its most informative subset is $(*, SF, London)$.

436 Finally, we discuss Shrink. The motivation behind Shrink was to create a
437 new OLAP operator to reduce the size of a data cube after a user has drilled
438 down to a fine-grained level. This is done by merging similar slices while
439 balancing the tradeoff between the number of merged slices returned and
440 accuracy of the corresponding aggregate function over the measure attribute.
441 In the context of pattern-based exploration, we can position Shrink as a
442 method that starts with a full data cube (i.e., all possible patterns) and
443 produces k non-overlapping patterns that summarize the distribution of the
444 measure attribute with (approximately) a minimal sum of squared errors.

445 Notably, Shrink explicitly allows patterns with disjunctions of values and
446 dimension hierarchies. Table 6 shows an example of a summary with two
447 patterns that may be considered by Shrink based on Table 1. We assume
448 that the aggregate function is $AVG(Late)$. Given the average values of the
449 *Late* measure attribute reported by the summary, the sum of squared errors
450 is 77.1.

451 Merging slices is done greedily. In each iteration, pairs of slices are chosen
452 for merging if merging the corresponding aggregate values of the measure
453 attributes leads to the smallest increase in sum of squared errors.

Table 6: A summary of size two considered by Shrink [11]

Day	Origin	Dest.	Late
Fri,Sun,Sat,Thu	*	*	15.4
Tue,Wed,Mon	*	*	5.3

454 *3.3.2. Pros and Cons*

455 By design, information-based methods produce *informative* patterns that
 456 highlight fragments of the data with surprising distributions of the measure
 457 attribute. However, these methods tend to be expensive, especially as the
 458 number of dimension attributes grows.

459 **4. Performance Optimizations**

460 A critical challenge in pattern-based data exploration arises from the
 461 size of the search space: the number of possible patterns is exponential in
 462 the number of dimension attributes. Coverage-based methods therefore deal
 463 with a large number of candidate patterns when constructing a set cover,
 464 contrast-based methods must compute contrast scores for many candidate
 465 patterns, and information-based methods must keep track of the information
 466 content of many patterns. In this section, we categorize frequently used
 467 performance optimizations to address these challenges and enable interactive
 468 data exploration.

469 We identify the following categories of optimizations: *top-down pruning*,
 470 *row pruning*, *column pruning* and *parallel processing*. Table 7 lists the meth-
 471 ods that originally used these optimizations, indicated by a ‘*’; for brevity,
 472 we do not explicitly list methods whose contrast metrics are supported by
 473 DIFF since the optimizations implemented in DIFF also apply to those meth-
 474 ods. As we will explain throughout this section, some optimizations may
 475 potentially be applicable to other methods, indicated by ‘a’ in Table 7. In
 476 particular, optimizations that will be experimentally evaluated in Section 5
 477 are indicated by an ‘E’.

478 *4.1. Top-Down Pruning*

479 Top-down pruning refers to traversing the space of candidate patterns
 480 from general (all-stars) to specific and pruning candidates along the way.

Table 7: Performance optimizations, methods that use them (indicated by a ‘*’), methods that may potentially use them (indicated by an ‘a’), and new method-optimization combinations that will be tested experimentally in Section 5 (indicated by an ‘E’). Methods are ordered by category, divided by horizontal lines: coverage-based, followed by contrast-based, followed by information-based.

Method	Top-down pruning	Row pruning	Column pruning	Parallel processing
Data Auditor [9]	*	a	E	a
Smart Drilldown [15]	*	*	a	a
CAPE [18]	a	a	*	a
DIFF [2]	*	E	*	*
Explanation tables [7]	E	*	*	*
Shrink [11]		a	a	a
SURPRISE [21]	a	a	a	a

481 The idea is similar to that of the Apriori algorithm for frequent itemset
 482 mining [3].

483 In the context of data exploration using patterns, we define the *an-*
 484 *cestors* of a pattern p as all patterns p' such that $p \subseteq p'$, and the *de-*
 485 *scendants* of p as all patterns p' such that $p' \subseteq p$. We can generate the
 486 ancestors of p by replacing non-star values with stars. Similarly, we can
 487 generate the descendants of p by replacing one or more stars with val-
 488 ues from the corresponding column. For example, in Table 1, the ances-
 489 tors of $(*, SF, London)$ are $(*, *, London)$, $(*, SF, *)$ and $(*, *, *)$. Simi-
 490 larly, the descendants of $(*, *, London)$ are $(*, SF, London)$, $(Fri, *, London)$,
 491 $(Mon, SF, London)$, and so on.

492 Note that if p is a descendant of p' (equivalently, if p' is an ancestor of p)
 493 then $sup(p) \leq sup(p')$. DIFF exploits this observation by requiring the user
 494 to set a minimum support threshold for a candidate pattern. Then, similar
 495 to Apriori frequent itemset mining, DIFF traverses the space of candidate
 496 patterns from the top down. When a pattern p is encountered with a sup-
 497 port below the threshold, p and all of its descendants can be ignored (i.e.,
 498 their contrast measure will not be calculated). Since DIFF supports vari-
 499 ous contrast metrics used in other contrast-based methods, this performance
 500 optimization also applies to other contrast-based methods.

501 Coverage-based methods such as Data Auditor and Smart Drilldown also

502 use top-down pruning. In each iteration of the greedy set cover heuristic,
503 these methods select the next best pattern based on the number of uncovered
504 tuples it covers (multiplied by the weight in case of Smart Drilldown). Thus,
505 a performance bottleneck results from having to keep track of the number
506 of uncovered tuples that can be covered by the candidate patterns – this
507 changes in every iteration, whenever a new pattern is added to the solution.

508 Here, an important observation is that a descendant of a pattern p cannot
509 cover more uncovered elements than p . Data Auditor exploits this observa-
510 tion and does not generate all the patterns that serve as input to greedy set
511 cover beforehand. Instead, patterns are generated on-demand, only after all
512 of their ancestors have already been considered. For example, a pattern such
513 as $(Fri, *, London)$ would only be generated after all of its ancestors - includ-
514 ing $(Fri, *, *)$ and $(*, *, London)$ - have already been considered. Until then,
515 $(Fri, *, London)$ can be safely ignored: it cannot cover more elements than
516 its ancestors and therefore is guaranteed to not be selected by the greedy
517 set cover heuristic at this time. This optimization can greatly reduce the
518 number of generated patterns whose coverage (of uncovered elements) must
519 be re-computed while constructing the solution.

520 Top-down pruning has not been applied to information-based methods.
521 One problem is that a descendant pattern may be more informative, even if
522 its support is lower. In Section 5, we will investigate whether the approach
523 used in DIFF, which is to require the user to set a minimum support threshold
524 for a candidate pattern, is beneficial for information-based methods.

525 Finally, as we indicated in Table 7, there does not appear to be an obvious
526 way to apply top-down pruning to Shrink. During pattern mining, Shrink
527 starts with all possible patterns and iteratively merges patterns until the
528 final result contains only k merged patterns, for some user-supplied value
529 of k . This is a bottom-up approach. Furthermore, ignoring patterns with
530 low support does not appear to be helpful since these patterns would not
531 be merged and instead would appear individually in the final output. This
532 would increase the size of the output, which is the opposite of the intended
533 goal of reducing the number of patterns.

534 4.2. Row Pruning

535 A simple example of row pruning is sampling: we draw a random sample
536 from the input dataset and run a data exploration algorithm on the sample
537 instead of the entire dataset. In principle, sampling applies to all the methods
538 discussed in this survey. In particular, sampling is used by information-based

539 methods such as Explanation Tables to reduce the number of patterns whose
540 information gain must be computed, and by coverage-based methods such as
541 Smart Drilldown to reduce the number of patterns whose coverage must be
542 computed. Furthermore, sampling may speed up the computation of infor-
543 mation gain or coverage of candidate patterns since there is less data to scan.
544 For example, the explanation table construction algorithm in [6, 7, 8] draws
545 a random sample in every iteration. Next, the set of candidate patterns cor-
546 responds to only those patterns that have a non-zero support in the sample.
547 The intuition is that patterns with frequently occurring combinations of di-
548 mension attribute values are likely to be sampled and also likely to contain
549 information about the distribution of the measure attribute.

550 However, a disadvantage of generating patterns from a sample is that
551 any pattern statistics calculated from a sample, such as contrast scores or
552 information gain, may not be accurate with respect to the full dataset. This
553 is in contrast to top-down pruning, which does not affect the quality of the
554 output: candidates are ignored only if they are guaranteed to not be included
555 in the output because their support is below the minimum support threshold
556 set by the user. Flashlight, which is one of the explanation table algorithms
557 proposed in [6], suggests the following compromise: sampling is used only to
558 restrict the candidate pattern space (to those with non-zero support in the
559 sample), but the information gain of candidate patterns is computed over
560 the full dataset, not the sample. This way, the resulting explanation table
561 is more informative since the patterns reflect the distribution of the measure
562 attribute in the entire dataset rather than the distribution within the sample.

563 *4.3. Column Pruning*

564 This category of optimizations removes some columns from consideration
565 in order to reduce the space of candidate patterns. We identified two exam-
566 ples of column pruning: removing correlated columns and limiting pattern
567 size. Again, both of these optimizations apply to all the methods surveyed
568 in this paper.

569 In terms of removing correlated columns, both CAPE and DIFF include a
570 pre-processing step that removes attributes that are functionally determined
571 by other attributes. The modified dataset is then used for exploration.

572 Limiting pattern size refers to limiting the number of non-star values that
573 appear in a pattern, which again reduces the space of candidate patterns.
574 This optimization is used by CAPE, DIFF and SIRUM [8], which is one of
575 the explanation table construction algorithms.

576 *4.4. Parallel Processing*

577 Parallel processing is used by DIFF and explanation tables (SIRUM). In
578 a contrast-based method such as DIFF, there are many candidate patterns
579 whose contrast measure needs to be computed. These computations can
580 be done independently for each pattern, and therefore can be parallelized.
581 In explanation table construction, every iteration requires the computation
582 of information gain of candidate patterns. This can also be parallelized,
583 but only within an iteration, not across iterations, since information gain
584 of the remaining patterns may change when a new pattern is added to the
585 explanation table. Finally, in coverage-based methods, it may be possible to
586 parallelize the computation of coverage of candidate patterns.

587 **5. Experiments**

588 In Section 4, we suggested that some performance optimizations origi-
589 nally proposed in the context of one method may apply to other methods.
590 In this section, we experimentally verify this claim, by applying optimiza-
591 tions originally proposed for a method in one category to a representative
592 method from another category. In particular, we evaluate the following *new*
593 combinations of methods and optimizations.

- 594 1. DIFF (contrast-based) uses top-down pruning by ignoring patterns
595 whose support is below a user-specified threshold. To evaluate the im-
596 pact of this optimization on information-based methods, we add this
597 optimization to the Flashlight algorithm for explanation table construc-
598 tion. In other words, in each iteration, the modified Flashlight algo-
599 rithm does not compute the information gain of any patterns whose
600 support is below the threshold.
- 601 2. Row pruning (sampling) was originally used in Flashlight (information-
602 based) and Smart Drilldown (coverage-based). To evaluate the impact
603 of this optimization on contrast-based methods, we compare the per-
604 formance on DIFF on a full dataset and DIFF on a random sample of
605 a dataset.
- 606 3. DIFF (contrast-based) and SIRUM (information-based) limit the num-
607 ber of non-star values that a pattern may use, which is an example
608 of column pruning. We apply this optimization to the Data Auditor
609 coverage-based method. To do this, we modify the Data Auditor al-
610 gorithm to ignore patterns with more than a user-specified number of
611 non-star values.

612 Note that we do not investigate new applications of parallel processing
613 optimizations. These kinds of optimization usually require a redesign of the
614 underlying algorithm (e.g., the SIRUM algorithm for parallel construction of
615 explanation tables [8]), which is outside the scope of this survey and is an
616 interesting direction for future work.

617 5.1. Experimental Setup

618 **Overview:** Experiments were performed in Ubuntu Linux 16.04, on a
619 device with Intel Core i7-6700HQ 2.60 GHz processor and 12 GB of RAM.
620 Each experiment (i.e., each combination of method, dataset, and parameter
621 settings) was repeated three times. For each experiment, we report the aver-
622 age running time across the three runs, as well as a measure of the goodness
623 or quality of the output, such as coverage percentage or information gain.

624 **Data:** We use the following datasets:
625

- 626 • *Flight Upgrade:* This dataset includes 5794 United Airlines ticket records,
627 with 7 categorical flight information attributes (origin, destination, air-
628 plane model, etc.), and a binary measure attribute determining whether
629 or not the passenger had their ticket upgraded to business class. The
630 dataset was used to evaluate explanation tables in previous work [7]
631 and was obtained from the authors of the corresponding paper. We
632 use this dataset in Experiments 1 and 3 to evaluate optimizations of
633 the Flashlight algorithm for explanation tables and Data Auditor, re-
634 spectively.
- 635 • *Adult Income:* There are 32561 records in this US Census dataset,
636 where each record includes 8 categorical attributes about the adult
637 person demographics (age, sex, education, etc.), and a binary measure
638 attribute denoting whether the income level of the person is above or
639 below \$50,000. This dataset is accessible from the UCI repository³, and
640 we use this dataset for Experiment 2 to test optimizations of DIFF.
641 Note that we do not use the smaller Flight Upgrade dataset to test
642 DIFF because DIFF queries run much faster than the other methods,
643 and thus measuring DIFF runtimes on small datasets does not show
644 any meaningful differences.

³<https://archive.ics.uci.edu/ml/datasets/adult>

645 • *Bank Marketing*: This dataset, which can be accessed from UCI repos-
646 itory⁴, includes records of bank marketing campaign outcomes, includ-
647 ing 10 categorical attributes about the demographic information of each
648 customer and the campaign contacts with the customer. The binary
649 measure attribute determines whether the campaign was effective and
650 the customer subscribed to the service that was advertised. We use the
651 smaller version of the dataset with 4119 records for Experiments 1 and
652 3 (explanation tables and Data Auditor), and the larger full version
653 with 41188 records for Experiment 2 (DIFF).

654 **Source code:** We obtained the DIFF source code from the project
655 Github page⁵. We obtained the Flashlight source code, written in C++,
656 from the authors of [23], and the Data Auditor code, also written in C++,
657 from the authors of [9].

658 5.2. Experiment 1: Impact of top-down pruning on information-based meth- 659 ods

660 In this experiment, we compare the original Flashlight algorithm with
661 a modified algorithm that ignores patterns whose support is below a user-
662 specified threshold. This reduces the space of candidate patterns, at the
663 expense of solution quality since some of the pruned patterns may have been
664 informative. We thus investigate the tradeoff between running time and the
665 information gain of the resulting explanation tables for various minimum
666 support thresholds: 0.005, 0.01, 0.05. We fix the size of the explanation
667 table at ten patterns, and we set the sample size for pattern generation to
668 16 (which was shown to be effective in previous work [6]).

669 Figure 1 displays the results of this experiment on the Flight (left) and
670 Bank (right) datasets. The plots show how the information gain (y-axes)
671 and execution time (x-axes) of the Flashlight algorithm change for different
672 minimum support thresholds (the blue points), compared to the original
673 algorithm without any top-down pruning (the black point).

674 Our results show that top-down pruning does not reduce the running time
675 significantly, while the information gain of the resulting explanation tables
676 degrades. Further analysis showed that the main performance bottleneck
677 in the Flashlight algorithm is the generation of candidate patterns, and the

⁴<http://archive.ics.uci.edu/ml/datasets/Bank+Marketing>

⁵<https://github.com/stanford-futuredata/macrobases>

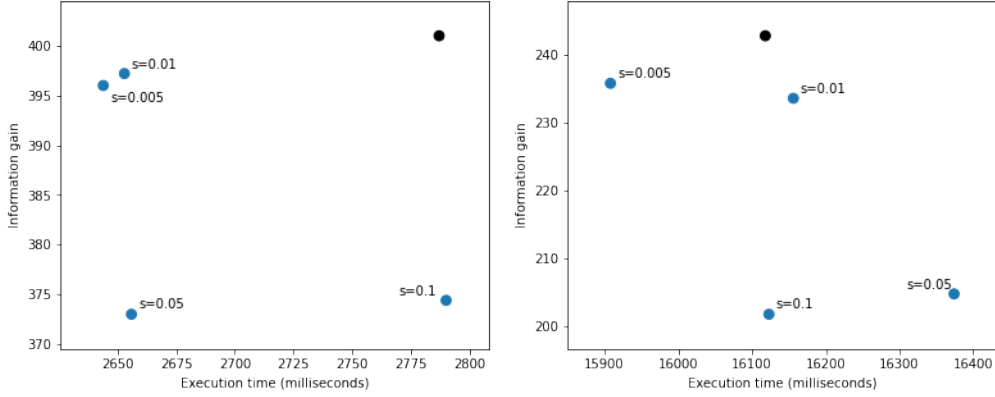


Figure 1: Experiment 1 results using Flashlight on Flight (left) and Bank (right) datasets

678 performance gain from not having to compute the information gain of some
 679 of the generated patterns (i.e., those below the given support threshold) is
 680 minimal. Furthermore, we noticed that patterns with large support increase
 681 the time taken by the algorithm to update the maximum-entropy estimate of
 682 the distribution of the measure attribute via iterative scaling. Thus, the time
 683 savings resulting from not having to compute the information gain of some
 684 patterns are often offset by the extra time taken to update the estimated
 685 distribution.

686 5.3. Experiment 2: Impact of row pruning on contrast-based methods

687 Next, we run DIFF on an entire dataset and on samples of the dataset.
 688 Reducing the sample size should make DIFF faster, at the expense of a larger
 689 error in the computed contrast scores. We thus evaluate the tradeoff between
 690 running time and the accuracy of contrast scores computed from a sample.

691 We test the following sample fractions of the full dataset: 0.01, 0.05,
 692 0.1 and 0.2. From the contrast metrics supported by DIFF, we test the
 693 risk ratio and the mean shift in two separate sets of experiments (recall
 694 Section 3.2.1). For the risk ratio scores, we compute the contrast between
 695 subsets of the data with different values of the binary measure (denoting
 696 high income in the Adult Income dataset or marketing campaign success in
 697 the Bank Marketing dataset). For the mean shift scores, we use numeric
 698 measure attributes (denoting the number of working hours per week in the
 699 Adult Income dataset and the number of times the customer was contacted
 700 in the Bank Marketing dataset). We use the default DIFF parameters for

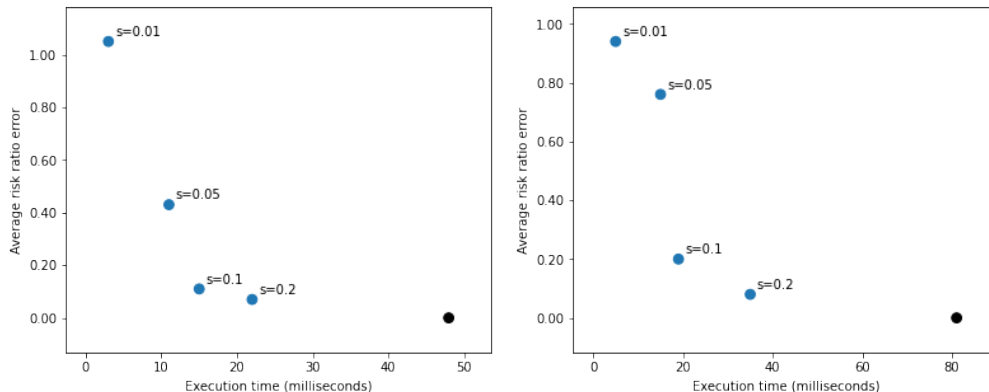


Figure 2: Experiment 2 results using DIFF risk ratio score on Income (left) and larger Bank (right) datasets

701 minimum support (10%) and minimum contrast score (1.5) thresholds.

702 To compare the DIFF contrast scores for full and sampled datasets, for
 703 each sample size, we compute the difference between the contrast scores of
 704 every pattern common between the full dataset result and the sampled result,
 705 and we take the average of these differences as the overall average error.

706 The results of the experiments with DIFF risk ratio score for the Income
 707 (left) and larger Bank (right) datasets are shown in Figure 2. The mean shift
 708 results are shown in Figure 3. The average contrast score errors for different
 709 sample size fractions are shown (blue points), compared with the full dataset
 710 (the black point).

711 Based on the four plots, we conclude that sampling is effective in reduc-
 712 ing the running time (x-axes), with the average contrast score errors becom-
 713 ing quite small at 0.1 and 0.2 sample fractions (y-axes). However, smaller
 714 samples appear to increase the error without a significant corresponding im-
 715 provement in running time.

716 *5.4. Experiment 3: Impact of column pruning based on pattern size on coverage-*
 717 *based methods*

718 Finally, we compare the original Data Auditor algorithm with a modified
 719 version that limits the number of non-star values a pattern may have. This
 720 optimization should reduce the space of candidate patterns and therefore
 721 the running time, possibly at the expense of solution quality (more patterns
 722 may be required to reach the desired coverage). However, patterns with

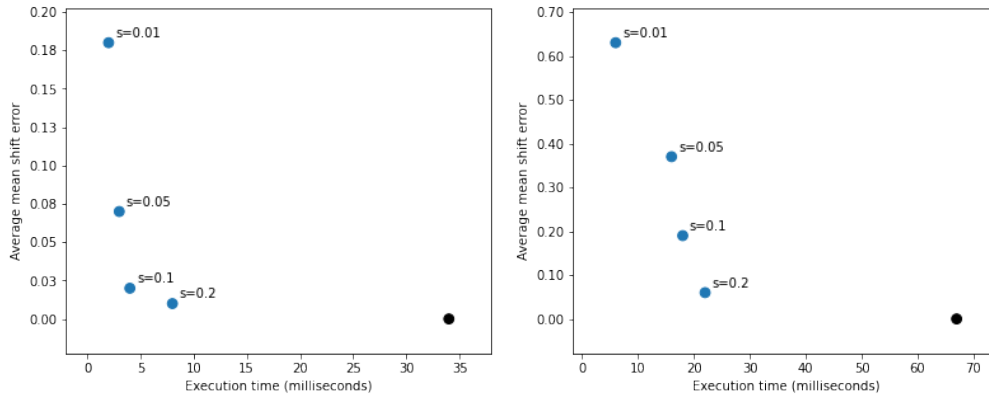


Figure 3: Experiment 2 results using DIFF mean shift score on Income (left) and larger Bank (right) datasets

723 many non-star values are expected to have low coverage, meaning that this
 724 optimization should be quite effective.

725 We test the following thresholds for pattern size, i.e., the maximum num-
 726 ber of non-star values in a pattern: 1, 2, 3 and 4. We set the $\theta_r(p)$ thresh-
 727 old to 0.9. By default, Data Auditor terminates when the generated patterns
 728 collectively have the user-specified coverage. However, in order to compare
 729 the results obtained using different maximum non-star thresholds based on
 730 their total coverage, we modified the code to output a fixed maximum num-
 731 ber of patterns, which in our experiments is set to 10, similar to the first
 732 experiment with explanation tables.

733 Figure 4 shows the results for the Flight (left) and Bank (right) datasets.
 734 The blue points represent the coverage and running time for different values
 735 of the pattern size threshold (denoted by m), while the black point shows the
 736 result of the original Data Auditor algorithm without any pattern pruning.

737 Limiting the number of non-star attributes significantly reduces the exe-
 738 cution time (by multiple orders of magnitude; x-axes) without degrading the
 739 total coverage significantly (y-axes). In the Flight dataset, setting the thresh-
 740 old to two results in the best trade-off between performance and accuracy,
 741 while in the Bank dataset, a higher threshold of four provides significant
 742 improvements in running time with a small decrease in coverage. In both
 743 plots, we see that having only one non-star attribute in the patterns reduces
 744 the total coverage significantly, likely because there are few such patterns
 745 that meet the $\theta_r(p)$ threshold, and therefore there are few candidate sets for

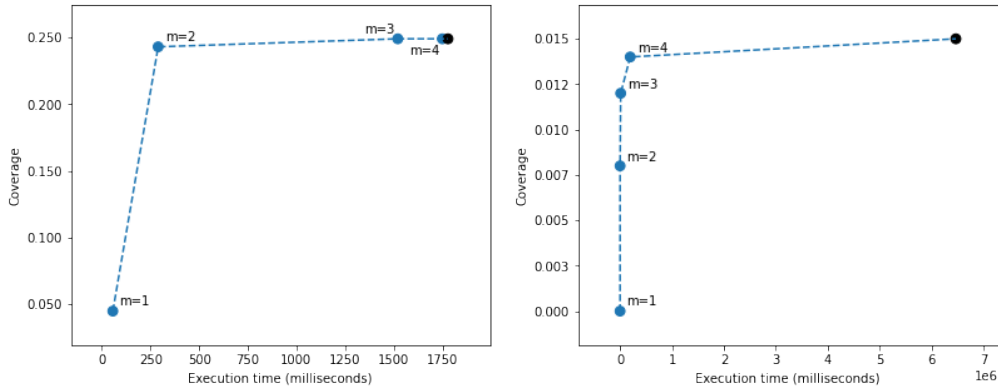


Figure 4: Experiment 3 results using Data Auditor on Flight (left) and Bank (right) datasets

746 the set cover solution produced by Data Auditor. In general, finding a good
 747 value of m for a given dataset is an interesting open problem.

748 6. Conclusions and Open Problems

749 We surveyed recent data exploration methods that extract interesting or
 750 informative fragments of the data, represented as patterns over the dimension
 751 attributes. We categorized these methods according to the properties of pat-
 752 terns they select, and we identified and experimentally evaluated frequently
 753 used performance optimizations. Our experimental study should be of in-
 754 terest to researchers working in the area of pattern-based data exploration
 755 as it suggests that many existing performance optimizations may apply to
 756 newly developed techniques. Specifically, our experiments showed that using
 757 a sample of a dataset is effective in improving the performance of contrast-
 758 based methods supported by DIFF, and we observed a performance gain
 759 after applying a limit on pattern size in the Data Auditor coverage-based
 760 method.

761 Below, we offer suggestions for future work in this area.

762 *Benchmarks:* A performance comparison of contrast-based methods im-
 763 plemented within the DIFF frameworks appears in [2]. In terms of effective-
 764 ness, prior work reports that methods based on information provide more
 765 information about the distribution of the measure attribute than coverage-
 766 based methods [6, 7]; similarly, methods based on contrast provide more
 767 precise outlier explanations than methods based on coverage [24]. Some ap-

768 proaches were also evaluated through user studies against simple baselines
769 [18]. An interesting direction for future work is to develop a benchmark to
770 highlight the effectiveness of various types of methods in various applications.

771 *New applications:* Popular motivating applications that guided the de-
772 velopment of prior work were outlier and data error analysis, as well as query
773 result explanation. Recent interest in explainable AI motivates further stud-
774 ies on exploring the behaviour of black-box machine learning models such as
775 neural networks using multi-dimensional patterns, as was suggested in [7].
776 Since deep learning methods have been successful in the context of unstruc-
777 tured data such as text, images and graphs, future research should investigate
778 new ways of formulating interpretable patterns over these high-dimensional
779 unstructured datasets.

780 *Correlated measure attributes:* Another characteristic of prior work is
781 that it usually formulates exploration problems involving a single measure
782 attribute. Pattern-based exploration of multiple measure attributes is an
783 interesting area for future work.

784 *Feature reduction:* In terms of performance and scalability, the large num-
785 ber of possible patterns remains a challenge for many methods, especially
786 those based on information which cannot leverage Apriori-like pruning strate-
787 gies. This is an important challenge for interactive methods that allow users
788 to continuously issue new exploration tasks. As a result, some techniques
789 such as DIFF limit the number of dimension attributes for use in patterns
790 and discard redundant dimension attributes such as those functionally deter-
791 mined by other attributes. Distributed versions of some methods, including
792 DIFF [2] and Explanation Tables [8], have also been proposed to parallelize
793 the search for interesting patterns. In machine learning, there exists a vari-
794 ety of dimension reduction methods such as Principal Component Analysis
795 (PCA) and word embeddings. However, these methods are not known for
796 being interpretable and thus their suitability for pattern-based exploration
797 requires further study.

798 *Bringing order to dimension attributes:* Much of the previous work con-
799 siders categorical dimension attributes. However, there exist methods for
800 covering a multi-dimensional dataset using *hyper-rectangles* corresponding
801 to intervals over numeric dimension attributes [17], there exists a method to
802 cover data anomalies using intervals over numeric features [27], and expla-
803 nation tables have recently been extended to support ordinal and numeric
804 dimension attributes [23]. These extensions further increase the space of
805 candidate patterns and require additional performance optimizations. For

806 example, returning to Table 1, the Day attribute may lead to additional pat-
807 terns with ranges or intervals such as ($[Mon - Fri], *, *$) or ($[Sat - Sun], *, *$).
808 Techniques used to construct optimal histograms and optimal decision trees
809 may help to optimize the discovery of these types of patterns.

810 *Exploring data evolution:* Finally, recent work motivates the need for tools
811 to explore how data (and metadata) change over time [4]. Here, patterns may
812 summarize fragments of the data that have changed recently or are updated
813 often.

814 References

- 815 [1] F. Abuzaid, P. Bailis, J. Ding, E. Gan, S. Madden, D. Narayanan, K.
816 Rong, S. Suri: MacroBase: Prioritizing Attention in Fast Data. ACM
817 Trans. Database Syst. 43(4): 15:1-15:45 (2018)
- 818 [2] F. Abuzaid, P. Kraft, S. Suri, E. Gan, E. Xu, A. Shenoy, A. Anatha-
819 naraya, J. Sheu, E. Meijer, X. Wu, J. F. Naughton, P. Bailis, M. Zaharia:
820 DIFF: A Relational Interface for Large-Scale Data Explanation. PVLDB
821 12(4): 419-432 (2018)
- 822 [3] R. Agrawal, R. Srikant: Fast Algorithms for Mining Association Rules
823 in Large Databases. VLDB 1994: 487-499
- 824 [4] T. Bleifuss, L. Bornemann, T. Johnson, D. V. Kalashnikov, F. Nau-
825 mann, D. Srivastava: Exploring Change - A New Dimension of Data
826 Analytics. PVLDB 12(2): 85-98 (2018)
- 827 [5] M. Das, S. Amer-Yahia, G. Das, C. Yu: MRI: Meaningful Interpretations
828 of Collaborative Ratings. PVLDB 4(11): 1063-1074 (2011)
- 829 [6] K. El Gebaly, P. Agrawal, L. Golab, F. Korn, D. Srivastava: Inter-
830 pretable and Informative Explanations of Outcomes. PVLDB 8(1): 61-
831 72 (2014)
- 832 [7] K. El Gebaly, G. Feng, L. Golab, F. Korn, D. Srivastava: Explanation
833 Tables. IEEE Data Eng. Bull. 41(3): 43-51 (2018)
- 834 [8] G. Feng, L. Golab, D. Srivastava: Scalable Informative Rule Mining.
835 ICDE 2017: 437-448

- 836 [9] L. Golab, H. J. Karloff, F. Korn, D. Srivastava: Data Auditor: Explor-
837 ing Data Quality and Semantics using Pattern Tableaux. PVLDB 3(2):
838 1641-1644 (2010)
- 839 [10] L. Golab, D. Srivastava: Exploring Data using Patterns: A Survey and
840 Open Problems. DOLAP 2021: 116-120
- 841 [11] M. Golfarelli, S. Graziani, S. Rizzi: Shrink: An OLAP operation for
842 balancing precision and size of pivot tables. Data Knowl. Eng. 93: 19-41
843 (2014)
- 844 [12] J. Gray, S. Chaudhuri, A. Bosworth, A. Layman, D. Reichart, M. Venka-
845 traao, F. Pellow, H. Pirahesh: Data Cube: A Relational Aggregation
846 Operator Generalizing Group-by, Cross-Tab, and Sub Totals. Data Min.
847 Knowl. Discov. 1(1): 29-53 (1997)
- 848 [13] R. Guidotti, A. Monreale, S. Ruggieri, F. Turini, F. Giannotti, D. Pe-
849 dreschi: A Survey of Methods for Explaining Black Box Models. ACM
850 Comput. Surv. 51(5): 93:1-93:42 (2019)
- 851 [14] S. Idreos, O. Papaemmanouil, S. Chaudhuri: Overview of Data Explo-
852 ration Techniques. SIGMOD Conference 2015: 277-281
- 853 [15] M. Joglekar, H. Garcia-Molina, A. G. Parameswaran: Interactive data
854 exploration with smart drill-down. ICDE 2016: 906-917
- 855 [16] H. Lakkaraju, S. H. Bach, J. Leskovec: Interpretable Decision Sets: A
856 Joint Framework for Description and Prediction. KDD 2016: 1675-1684
- 857 [17] L. V. S. Lakshmanan, R. T. Ng, C. X. Wang, X. Zhou, T. Johnson: The
858 Generalized MDL Approach for Summarization. VLDB 2002: 766-777
- 859 [18] Z. Miao, Q. Zeng, C. Li, B. Glavic, O. Kennedy, S. Roy: CAPE: Explain-
860 ing Outliers by Counterbalancing. PVLDB 12(12): 1806-1809 (2019)
- 861 [19] X. Qin, Y. Luo, N. Tang, G. Li: Making data visualization more efficient
862 and effective: a survey. VLDB J. 29(1): 93-117 (2020)
- 863 [20] S. Roy, D. Suciú: A formal approach to finding explanations for database
864 queries. SIGMOD Conference 2014: 1579-1590

- 865 [21] S. Sarawagi: User-cognizant multidimensional analysis. VLDB J. 10(2-
866 3): 224-239 (2001)
- 867 [22] P. Vassiliadis, P. Marcel: The Road to Highlights is Paved with Good
868 Intentions: Envisioning a Paradigm Shift in OLAP Modeling. DOLAP
869 2018.
- 870 [23] M. Vollmer, L. Golab, K. Bohm, D. Srivastava: Informative Summa-
871 rization of Numeric Data. SSDBM 2019: 97-108
- 872 [24] X. Wang, X. L. Dong, A. Meliou: Data X-Ray: A Diagnostic Tool for
873 Data Errors. SIGMOD Conference 2015: 1231-1245
- 874 [25] E. Wu, S. Madden: Scorpion: Explaining Away Outliers in Aggregate
875 Queries. PVLDB 6(8): 553-564 (2013)
- 876 [26] J. Xu Yu, L. Qin, L. Chang: Keyword Search in Databases. Synthesis
877 Lectures on Data Management, Morgan & Claypool Publishers, 2009
- 878 [27] H. Zhang, Y. Diao, A. Meliou: EXstream: Explaining Anomalies in
879 Event Stream Monitoring. EDBT 2017: 156-167