



**HAL**  
open science

## **DEvIANT : Discovering significant exceptional (dis)agreement within groups**

Adnene Belfodil, Wouter Duivesteijn, Marc Plantevit, Sylvie Cazalens,  
Philippe Lamarre

► **To cite this version:**

Adnene Belfodil, Wouter Duivesteijn, Marc Plantevit, Sylvie Cazalens, Philippe Lamarre. DEvIANT : Discovering significant exceptional (dis)agreement within groups. Joint European Conference on Machine Learning and Knowledge Discovery in Databases (ECML PKDD 2019), Sep 2019, Würzburg, Germany. pp.3-20, 10.1007/978-3-030-46150-8\_1 . hal-02961093

**HAL Id: hal-02961093**

**<https://hal.science/hal-02961093>**

Submitted on 8 Oct 2020

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# DEvIANT : Discovering significant exceptional (dis)agreement within groups

Adnene Belfodil<sup>1</sup>, Wouter Duivesteijn<sup>2</sup>, Marc Plantevit<sup>3</sup>, Sylvie Cazalens<sup>1</sup>, and Philippe Lamarre<sup>1</sup>

<sup>1</sup> Univ Lyon, INSA Lyon, CNRS, LIRIS UMR 5205, F-69621, Lyon, France

<sup>2</sup> Technische Universiteit Eindhoven, Eindhoven, The Netherlands

<sup>3</sup> Univ Lyon, CNRS, LIRIS UMR 5205, F-69622, Lyon, France

**Abstract.** We consider any type of data featuring individuals (e.g., parliamentarians, customers) performing observable actions (e.g., votes, ratings) on entities (e.g., legislative procedures, movies). In such data, we aim to find contexts (i.e. subgroup of entities) for which an exceptional (dis)agreement is observed among a group of individuals. To this end, we introduce the novel problem of discovering statistically significant exceptional contextual intra-group agreement patterns. To handle the data sparsity, we use the Krippendorff’s Alpha measure to assess the agreement among individuals. We devise a branch-and-bound algorithm, named DEvIANT, to discover such patterns. DEvIANT exploits both closure operators and tight optimistic estimates. We derive analytic approximations for the confidence intervals (CIs) associated to patterns for a computationally efficient significance assessment. We prove that these approximate CIs are nested along specialization of patterns. This makes it possible to incorporate pruning properties in the algorithm to early discard non-significant patterns. Empirical study on several datasets demonstrates the efficiency and the usefulness of DEvIANT.

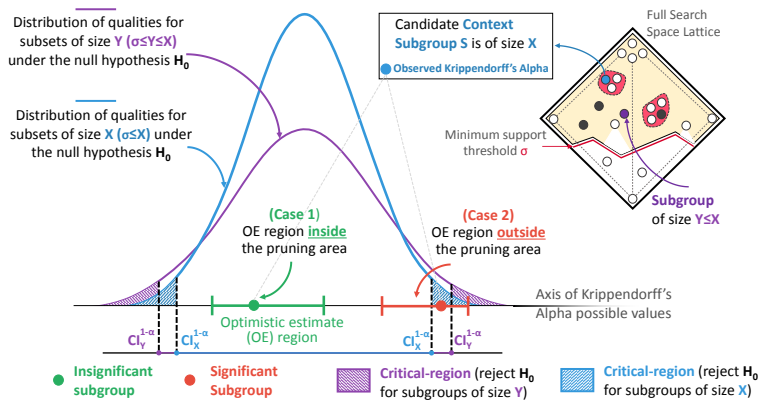
## 1 Introduction

Consider data describing the organization and votes of the European Parliament (EP). Such dataset records the votes of each member (MEP) in voting sessions held in the parliament, as well as the information on the parliamentarians (e.g., gender, national party, European party alliance) and the sessions (e.g., topic, date). This data offers interesting opportunities to study the agreement or disagreement of coherent subgroups, especially to highlight some unexpected ones. It is to be expected that on the majority of voting sessions, MEPs will vote along the lines of their European party alliance. However, when matters are of interest to a specific nation within Europe, alignments may change and agreements can be formed or dissolved. For instance, when a ballot on fishing rights is put before the MEPs, the island nation of the UK can be expected to agree on a specific course of action regardless of their party alliance, fostering an exceptional agreement where strong polarization exists otherwise.

We aim to discover such exceptional (dis)agreements. This is not limited to just EP or voting data: members of the US congress also votes on bills while

Amazon-like customers post ratings or reviews of products. A challenge with this approach when considering voting or rating data, however, is to effectively handle the absence of outcomes (sparsity) which is high in such datasets. For instance, in the European parliament data, MEPs votes on average on only a three-quarter of all sessions. These outcomes are not missing at random: special workgroups are often formed of MEPs who are tasked with properly studying a specific topic, and members of these workgroups are more likely to vote on their topic of study. Hence, present values are likely associated with more pressing votes, which means that missing values need to be treated carefully. This problem becomes much worse when looking at Amazon or Yelp rating data: the vast majority of customers will not have rated the vast majority of products/places.

In this paper, we introduce the problem of discovering significantly exceptional contextual intra-group agreement patterns, rooted in the Subgroup Discovery (SD) [42]/ Exceptional Model Mining (EMM) [8] framework. To tackle the data sparsity issue, we measure the agreement among groups with *Krippendorff's alpha* [27], a measure developed in the context of content analysis [28] which is well-known to handle missing outcomes elegantly. We develop a branch-and-bound algorithm to automatically find subgroups featuring statistically significant exceptional (dis)agreement among groups. This algorithm enables to early discard non significant subgroups by pruning unpromising branches of the search space. Fig. 1 illustrates this. Suppose, we are interested by subgroups of entities (e.g. voting sessions) whose sizes are greater than a support threshold  $\sigma$ . For a given subgroup of size  $X \geq \sigma$ , we gauge its exceptionality by its *p-value*, i.e. the probability of observing a quality (i.e. Krippendorff's alpha) for a random subset of entities is at least as extreme as the one observed for the subgroup. Hence avoiding to report subgroups which observe a high quality that appears due to chance only (i.e. accepting the null hypothesis). For this task to be achieved, we can estimate the empirical distribution of the quality of random



**Fig. 1:** Main DEVIANT properties for safe sub-search space pruning. A subgroup is reported as significant if its related Krippendorff's Alpha falls in the critical region of the corresponding empirical distribution of random subsets (DFD). When traversing the search space downward (decreasing support size) the approximate confidence intervals are nested. If the optimistic estimates region falls into the confidence interval computed on the related DFD, the sub-search space can be safely pruned.

subsets (DFD: Distribution of False Discoveries, cf. [9,34]) and establish, for a certain critical value  $\alpha$ , a confidence interval  $CI_X^{1-\alpha}$  over the corresponding distribution under the null hypothesis. If the subgroup quality is outside  $CI_X^{1-\alpha}$ , this means that the subgroup is statistically significant ( $p\text{-value} \leq \alpha$ ), otherwise the subgroup is a spurious finding. We prove by approximation that the confidence intervals are nested:  $\sigma \leq Y \leq X \Rightarrow CI_X^{1-\alpha} \subseteq CI_Y^{1-\alpha}$ . Moreover, we compute a tight optimistic estimate (OE) [18] to have a lower-bound and an upper bound of the quality of any specialization of a subgroup having its size greater than  $\sigma$ . Combining this two properties, if the OE region falls into the corresponding CI, we can safely prune large parts of the sub-search space that do not contain significant subgroups. In summary, the main contributions are:

- We introduce the novel problem of discovering statistically significant exceptional contextual intra-group agreement patterns (Section 3).
- We derive an analytical approximation of the confidence intervals associated to subgroups. This allows a computationally efficient assessment of the statistical significance of the findings. Furthermore, we show that approximate confidence intervals are nested (Section 4). Particular attention is also paid to the variability of outcomes among raters (Section 5).
- We devise a branch-and-bound algorithm to discover exceptional contextual intra-group agreement patterns (Section 6). It exploits tight optimistic estimates on Krippendorff’s alpha and nested approximate CIs property to early discard non significant patterns.
- We report an empirical evaluation (Section 7) which studies the performance and the potential of the proposed approach.

## 2 Background and Related Work

**Measuring Agreement.** Several measures of agreement focus on two targets (Pearson’s correlation coefficient, Spearman’s  $\rho$ , Kendall’s  $\tau$ , Association) and cannot handle missing values well. As pointed out by Krippendorff [28, p.244], using association and correlation measures to assess agreement leads to particularly misleading conclusions. When all data falls along a line  $Y = aX + b$ , correlation is perfect, but agreement requires that  $Y = X$ . Cohen’s  $\kappa$  [4] is a seminal measure of agreement between two raters who classify items into a fixed number of mutually exclusive categories. The Fleiss  $\kappa$  [14] extends this notion to multiple raters and requires that each item receives the exact same number of ratings. Krippendorff’s alpha generalizes these measures while handling multiple raters, missing outcomes and several metrics [28, p.232].

**Discovering Significant Patterns.** The statistical assessment of patterns have received a considerable attention for a decade [41,20,21], especially for association rules [19,40,36]. Some works focused on the statistical significance of results in Subgroup Discovery/Exceptional Model Mining during enumeration [9,34] or a posteriori [10] for statistical validation of the found subgroups.

**Voting and Rating data Analysis.** In a previous work [2], we proposed a method to discover exceptional *inter-group agreement* in voting or rating data.

This method does not allow to discover intra-group agreement. Several works in the literature addressed the problem of uncovering groups in rating datasets whose members exhibit an agreement or discord [6,38] or a specific rating distribution [1] (e.g. polarized, homogeneous) given upfront by the end-user. This is done by aggregating the ratings by using an arithmetic mean or a rating distribution. However, none of these methods allow to discover automatically exceptional (dis)agreement within groups. Moreover, these methods may output misleading hypotheses over the intra-group agreement. This is due to two main factors: aggregating ratings in a distribution is (i) highly affected by the data sparsity (e.g. two reviewees may have significantly different number of expressed ratings) and (ii) may conceal the true nature of the underlying intra-group agreement. For instance, a rating distribution computed for a collection of movies may highlight a polarized distribution of ratings (interpreted as a disagreement) while ratings over each movie may describe a consensus between raters (movies are either highly or lowly rated or by the majority of the group). These two issues are addressed by Krippendorff’s alpha.

### 3 Problem Definition

The data we are interested in consists of a set of individuals (e.g. *social network users, parliamentarians*) who give outcomes (e.g. *ratings, votes*) on entities (e.g. *movies, ballots*). This type of data is called behavioral dataset (e.g. Tab. 1).

**Definition 1 (Behavioral Dataset).** A behavioral dataset  $\mathcal{B} = \langle G_I, G_E, O, o \rangle$  is defined by (i) a finite collection of Individuals  $G_I$ , (ii) a finite collection of Entities  $G_E$ , (iii) a domain of possible Outcomes  $O$ , and (iv) a function  $o : G_I \times G_E \rightarrow O$  that gives the outcome of an individual  $i$  over an entity  $e$ .

The elements from  $G_I$  (resp.  $G_E$ ) are augmented with descriptive attributes  $\mathcal{A}_I$  (resp.  $\mathcal{A}_E$ ). Attributes  $a \in \mathcal{A}_I$  (resp.  $\mathcal{A}_E$ ) may be Boolean, numerical or categorical attributes potentially organized among a taxonomy. Subgroups (subsets) of  $G_I$  (resp.  $G_E$ ) can be defined using descriptions from  $\mathcal{D}_I$  (resp.  $\mathcal{D}_E$ ). These descriptions are formalized by conjunctions of conditions on the values of the attributes. Descriptions of  $\mathcal{D}_I$  are called *groups*, denoted  $g$ . Descriptions of  $\mathcal{D}_E$  are called *contexts*, denoted  $c$ . From now on,  $G$  (resp.  $\mathcal{D}$ ) denotes both collections  $G_I$  (resp.  $\mathcal{D}_I$ ) and  $G_E$  (resp.  $\mathcal{D}_E$ ) if no confusion can arise. We denote by  $G^d$  the subset of records characterized by the description  $d \in \mathcal{D}$ . Descriptions from

ide themes			date	idi country group age				idi ide o(i,e)   idi ide o(i,e)					
$e_1$	1.20	Citizen’s rights	20/04/16	$i_1$	France	S&D	26	$i_1$	$e_2$	Against	$i_3$	$e_1$	For
$e_2$	5.05	Economic growth	16/05/16	$i_2$	France	PPE	30	$i_1$	$e_5$	For	$i_3$	$e_2$	Against
$e_3$	1.20	Citizen’s rights; 7.30 Judicial Coop	04/06/16					$i_1$	$e_6$	Against	$i_3$	$e_3$	For
$e_4$	7	Security and Justice	11/06/16	$i_3$	Germany	S&D	40	$i_2$	$e_1$	For	$i_3$	$e_5$	Against
$e_5$	7.30	Judicial Coop	03/07/16					$i_2$	$e_3$	Against	$i_4$	$e_1$	For
$e_6$	7.30	Judicial Coop	29/07/16	$i_4$	Germany	ALDE	45	$i_2$	$e_4$	For	$i_4$	$e_4$	For
								$i_2$	$e_5$	For	$i_4$	$e_6$	Against

(a) Entities

(b) Individuals

(c) Outcomes

**Table 1:** Example of behavioral dataset - European Parliament Voting dataset

$\mathcal{D}$  are partially ordered by a specialization operator denoted  $\sqsubseteq$ . A description  $d_1$  is a specialization of  $d_2$ , denoted  $d_1 \sqsubseteq d_2$ , iff  $d_2 \Rightarrow d_1$  from a logical point of view. It follows that  $G^{d_2} \subseteq G^{d_1}$ .

### 3.1 Krippendorff's Alpha

Krippendorff's Alpha ( $A$ ) measures the agreement among raters. This measure has several properties that make it attractive in our setting, namely: (i) it is applicable to any number of observers; (ii) it handles various domains of outcomes (ordinal, numerical, categorical, time series); (iii) it handles behavioral data with missing values; (iv) it takes into consideration the agreement expected by chance [28]. In its most general form,  $A$  is defined by:

$$A = 1 - \frac{D_o}{D_e} \quad (1)$$

where  $D_o$  (resp.  $D_e$ ) is a measure of the disagreement observed (resp. by chance). Hence, when  $A = 1$ , the agreement is as large as it can possibly be (given the class prior), and when  $A = 0$ , the agreement is indistinguishable to agreement by chance. We can also have  $A < 0$ , where disagreement is larger than expected by chance and which corresponds to systematic disagreement.

Given a behavioral dataset  $\mathcal{B}$ , we want to measure Krippendorff's alpha for a given context  $c \in \mathcal{D}_E$  characterizing a subset of entities  $G_E^c \subseteq G_E$ , which indicates to what extent the individuals who comprise some selected group are in agreement  $g \in \mathcal{D}_I$ . From Eq. (1), we have  $A(S) = 1 - \frac{D_o(S)}{D_e}$  for any  $S \subseteq G_E$ . Note that the measure takes into consideration only entities where at least two individuals expressed an outcome. We assume that the entities that do not fulfil this requirement are removed in preprocessing.

$$D_o(S) = \frac{1}{\sum_{e \in S} n_e} \sum_{c, k \in O^2} \delta_{ck} \cdot \sum_{e \in S} \frac{n_{ec} \cdot n_{ek}}{n_e - 1} \quad (2)$$

$n_e$  is the number of expressed outcomes for the entity  $e$  and  $n_{ec}$  (resp.  $n_{ek}$ ) represents the number of outcomes equal to  $c$  (resp.  $k$ ) expressed for the entity  $e$ .  $\delta_{ck}$  is a distance measure between outcomes, which can be defined according to the domain of the outcomes (e.g.  $\delta_{ck}$  can correspond to the Kronecker delta for categorical outcomes or distance between ordinal values for ratings. Choices for the distance measure are discussed in [30]). We define below  $D_e$  that represents the disagreement expected by chance in Krippendorff's alpha:

$$D_e = \frac{1}{\sum_{e \in G_E} n_e \cdot (\sum_{e \in G_E} n_e - 1)} \sum_{c, k \in O^2} \delta_{ck} \cdot n_c \cdot n_k \quad (3)$$

With  $n_c$  (resp.  $n_k$ ) the number of expressed outcomes equal to  $c$  (resp.  $k$ ) observed in the entire behavioral data. This corresponds to the disagreement by chance observed on the overall marginal distribution of outcomes.

*Example:* Tab. 2 summarizes the behavioral data given in Tab. 1. The disagreement expected by chance is equal to (given:  $nb(F) = 8$ ,  $nb(A) = 6$ ):  $D_e = \frac{1}{13 \times 14} \times ((8 \times 6) + (6 \times 8)) = 48/91$ . If we want to evaluate the intra-agreement between the four individuals in the global context (considering all entities), we need to compute, first, the observed disagreement  $D_o(G_E)$ . Simply put,  $D_o(G_E)$  is the weighted average of the two last lines by considering the quantities  $n_e$  as the weights:  $D_o(G_E) = \frac{4}{14}$ . Hence, for the global context,  $A(G_E) = 1 - \frac{4}{14} / \frac{48}{91} = 0.46$ . Consider the context  $c = \langle \text{themes} \supseteq \{7.30 \text{ Judicial Coop.}\} \rangle$ , having as support:  $G_E^c = \{e_3, e_5, e_6\}$ . The observed disagreement is obtained by computing the weighted average, only considering the entities belonging to the context:  $D_o(G_E^c) = \frac{4}{7}$ . Hence, the contextual intra-agreement is:  $A(G_E^c) = 1 - \frac{4}{7} / \frac{48}{91} = -0.08$ .

Comparing  $A(G_E^c)$  and  $A(G_E^*)$  leads to the following statement: “*while parliamentarians are slightly in agreement in overall terms, judicial cooperation related questions create systematic disagreement among the parliamentarians*”.

	[F]or		[A]gainst			
	$e_1$	$e_2$	$e_3$	$e_4$	$e_5$	$e_6$
$i_1$		A			F	A
$i_2$	F		A	F	F	
$i_3$	F	A	F		A	
$i_4$	F			F		A
$n_e$	3	2	2	2	3	2
$D_o(e)$	0	0	1	0	$\frac{2}{3}$	0

**Table 2:** Summarized Behavioral Data  $D_o(e) = \frac{n_{ec} \cdot n_{ek}}{\sum_{c,k \in O^2} \delta_{ck} n_e \cdot (n_e - 1)}$

### 3.2 Mining Significant Patterns with Krippendorff’s Alpha

Considering a selected group of individuals  $g \in \mathcal{D}_I$ , we are interested in finding contexts  $c \in \mathcal{D}_E$  where the observed intra-agreement, denoted  $A^g(c)$ , *significantly* differs from the expected intra-agreement (i.e. that appear due to chance alone). In the same spirit as [9,34,41], we evaluate the quality of patterns by statistical significance of the contextual intra-agreement. This choice is motivated by: (i) the desire to not specify to the algorithm an arbitrary threshold on the distance from the overall intra-agreement observed (fixing the critical value  $\alpha$  is more intuitive), (ii) the recommendations of Krippendorff [22] to provide a confidence interval on the alpha metric rather than a point-value.

In this work, we are interested in finding patterns of the form  $(g, c) \in \mathcal{D}_I \times \mathcal{D}_E$  highlighting an exceptional intra-agreement between members of a group of individuals  $g$  over a context  $c$ . To perform such task, we formalize the problem using the well-established framework of SD/EMM [8], while giving particular attention to the statistical significance and soundness of the discovered patterns [21].

**Problem Statement.** (*Discovering Exceptional Contextual Intra-group Agreement patterns*). Given a behavioral dataset  $\mathcal{B} = \langle G_I, G_E, O, o \rangle$ , a minimum group support threshold  $\sigma_I$ , a minimum context support threshold  $\sigma_E$ , a significance critical value  $\alpha \in ]0, 1]$  and the null hypothesis  $H_0$  corresponding to that the observed Krippendorff’s alpha is generated by a *distribution of false discoveries* (cf. [9]). The goal is to find the pattern set  $P \subseteq \mathcal{D}_I \times \mathcal{D}_E$  such that:

$$\forall (g, c) \in P : |G_I^g| \geq \sigma_I \wedge |G_E^c| \geq \sigma_E \text{ we have: } p\text{-value}^g(c) \leq \alpha$$

## 4 Exceptional Contexts: Evaluation and Pruning

From now and on, to avoid overloading notation and for the sake of simplicity, we omit the exponent  $g$  if no confusion can arise, while keeping in mind a selected group of individual  $g \in \mathcal{D}_I$  related to a subset  $G_I^g \subseteq G_I$ .

### 4.1 Gauging Exceptionality of a Subgroup

To evaluate to what extent our findings are exceptional, we follow the significant pattern mining paradigm. That is, we consider each context  $c$  as a hypothesis test which returns a  $p$ -value. The  $p$ -value is the probability of obtaining an intra-agreement at least as extreme as the one observed over the current context  $A(G_E^c)$ , assuming the truth of the null hypothesis  $H_0$ . The pattern is accepted if  $H_0$  is rejected. This happens if the  $p$ -value is under a critical significance value  $\alpha$  which amounts to test if the observed intra-agreement  $A(G_E^c)$  is outside the confidence interval  $CI^{1-\alpha}$  established using the distribution assumed under  $H_0$ .

$H_0$  corresponds to a baseline finding. i.e. the observed contextual intra-agreement is generated by the distribution of random subsets equally likely to occur, a.k.a: *Distribution of False Discoveries* (DFD, c.f. [9]). We resolve to evaluate the  $p$ -value of the observed  $A$  against the distribution of random subsets of a cardinality equal to the size of the observed subgroup  $G_E^c$ . The subsets are issued by a uniform sampling without replacement from the entire entities collection. The rationale behind using sampling without replacement is that the observed subgroup does not contain multiple instances of the same entity. Moreover, drawing samples only from the collection of subsets of size equal to  $|G_E^c|$  allows to drive more judicious conclusions: the variability of the statistic  $A$  is impacted by the size of the considered subgroups, since smaller subgroups are more likely to observe low/high values of  $A$ . The same reasoning was followed in [34]

We define  $\theta_k : F_k \rightarrow \mathbb{R}$  as the random variable corresponding to the observed intra-agreement  $A$  of  $k$ -sized subsets  $S \in G_E$ . i.e. for any  $k \in [1, n]$  with  $n = |G_E|$  we have  $\theta_k(S) = A(S)$  and  $F_k = \{S \in G_E \text{ s.t. } |S| = k\}$ .  $F_k$  is then the set of possible subsets which are equally likely to occur under the null hypothesis  $H_0$ . That is  $\mathbb{P}(S \in F_k) = \binom{n}{k}^{-1}$ . We denote by  $CI_k^{1-\alpha}$  the  $(1 - \alpha)$  confidence interval related to the probability distribution of  $\theta_k$  under the null hypothesis  $H_0$ . To easily manipulate  $\theta_k$ , we reformulate  $A$  using equations (1), (2) and (3) as such:

$$A(S) = \frac{\sum_{e \in S} v_e}{\sum_{e \in S} w_e} \quad \text{with } w_e = n_e \text{ and } v_e = n_e - \frac{1}{D_e} \sum_{c, k \in O^2} \delta_{ck} \cdot \frac{n_{ec} \cdot n_{ek}}{(n_e - 1)} \quad (4)$$

Considering the null hypothesis  $H_0$  and under the assumption that the underlying distribution of intra-agreements is a normal distribution<sup>4</sup>  $\mathcal{N}(\mu_k, \sigma_k^2)$ , one

<sup>4</sup> In the same line of reasoning of [7,35], one can assume that the underlying distribution can be derived from what prior beliefs the end-user may have on such distribution. If only the observed expectation  $\mu$  and variance  $\sigma^2$  are given as constraints which must hold for the underlying distribution, the maximum entropy distribution (*the one that takes into account no other prior information than the given constraints*) is known to be the normal distribution  $\mathcal{N}(\mu, \sigma^2)$  [5, p.413]



can define  $CI_k^{1-\alpha}$  by computing  $\mu_k = E[\theta_k]$  and  $\sigma_k^2 = \text{Var}[\theta_k]$ . Doing so, requires either calculating estimator of such moments empirically by drawing a large number  $r$  of uniformly generated samples from  $F_k$  or deriving analytically the formula of  $E[\theta_k]$  and  $\text{Var}[\theta_k]$ . In the former case, the confidence interval  $CI_k^{1-\alpha}$  endpoints are given by [17, p.9]:  $\mu_k \pm t_{1-\frac{\alpha}{2}, r-1} \sigma_k \sqrt{1 + (1/r)}$  with  $\mu_k$  and  $\sigma_k$  being estimated empirically on the  $r$  samples and  $t_{1-\frac{\alpha}{2}, r-1}$  the  $(1 - \frac{\alpha}{2})$  percentile of Student's t-distribution with  $r - 1$  degrees of freedom. In the latter case ( $\mu_k$  and  $\sigma_k$  are known/derived analytically), the  $(1 - \alpha)$  confidence interval can be computed in its most basic form, that is  $CI_k^{1-\alpha} = [\mu_k - z_{(1-\frac{\alpha}{2})} \sigma_k, \mu_k + z_{(1-\frac{\alpha}{2})} \sigma_k]$  with  $z_{(1-\frac{\alpha}{2})}$  the  $(1 - \frac{\alpha}{2})$  percentile of  $\mathcal{N}(0, 1)$ .

However, on one hand, due to the problem setting, establishing the confidence interval empirically is computationally expensive since it need to be calculated for each enumerated context which can become quickly unfeasible even for relatively small behavioral datasets. In the other hand, deriving analytically a computationally efficient form of  $E(\theta_k)$  is notoriously difficult, given that  $E[\theta_k] = \frac{1}{\binom{n}{k}} \sum_{S \in F_k} \frac{\sum_{e \in S} v_e}{\sum_{e \in S} w_e}$  and  $\text{Var}[\theta_k] = \frac{1}{\binom{n}{k}} \sum_{S \in F_k} \left( \frac{\sum_{e \in S} v_e}{\sum_{e \in S} w_e} - E[\theta_k] \right)^2$ .

Since  $\theta_k$  can be seen as a weighted arithmetic mean, one can model the random variable  $\theta_k$  as the ratio  $\frac{V_k}{W_k}$ . With  $V_k, W_k$  two random variables  $V_k : F_k \rightarrow \mathbb{R}$  and  $W_k : F_k \rightarrow \mathbb{R}$  with  $V_k(S) = \frac{1}{k} \sum_{e \in S} v_e$  and  $W_k(S) = \frac{1}{k} \sum_{e \in S} w_e$ . An elegant way to deal with a ratio of two random variables is to approximate its moments using the *Taylor series* following the line of reasoning of [11] and [25, p.351], since no easy analytic expression of  $E[\theta_k]$  and  $\text{Var}[\theta_k]$  can be derived. For the sake of brevity, the detailed computation of the formulas presented next are omitted. For more details, please refer to Appendix A.

**Proposition 1 (An approximate Confidence Interval  $\widehat{CI}_k^{1-\alpha}$  for  $\theta_k$ ).**  
Given  $k \in [1, n]$  and a significance critical value  $\alpha \in ]0, 1]$ ,  $\widehat{CI}_k^{1-\alpha}$  is given by:

$$\widehat{CI}_k^{1-\alpha} = \left[ \widehat{E}[\theta_k] - z_{1-\frac{\alpha}{2}} \sqrt{\widehat{\text{Var}}[\theta_k]}, \widehat{E}[\theta_k] + z_{1-\frac{\alpha}{2}} \sqrt{\widehat{\text{Var}}[\theta_k]} \right] \quad (5)$$

with:  $\widehat{E}[\theta_k]$  a Taylor approximation for the expectation  $E[\theta_k]$  expanded around  $(\mu_{V_k}, \mu_{W_k})$  and  $\widehat{\text{Var}}[\theta_k]$  a Taylor approximation for  $\text{Var}[\theta_k]$  given by:

$$\widehat{E}[\theta_k] = \binom{n}{k} - 1 \frac{\mu_v}{\mu_w} \beta_w + \frac{\mu_v}{\mu_w} \quad \widehat{\text{Var}}[\theta_k] = \binom{n}{k} - 1 \frac{\mu_v^2}{\mu_w^2} (\beta_v + \beta_w) \quad (6)$$

$$\text{with:} \quad \begin{aligned} \mu_v &= \frac{1}{n} \sum_{e \in G_E} v_e & \mu_w &= \frac{1}{n} \sum_{e \in G_E} w_e & n &= |G_E| \\ \mu_{v^2} &= \frac{1}{n} \sum_{e \in G_E} v_e^2 & \mu_{w^2} &= \frac{1}{n} \sum_{e \in G_E} w_e^2 & \mu_{vw} &= \frac{1}{n} \sum_{e \in G_E} v_e w_e \end{aligned}$$

$$\text{and:} \quad \beta_v = \frac{1}{n-1} \left( \frac{\mu_{v^2}}{\mu_v^2} - \frac{\mu_{vw}}{\mu_v \mu_w} \right) \quad \beta_w = \frac{1}{n-1} \left( \frac{\mu_{w^2}}{\mu_w^2} - \frac{\mu_{vw}}{\mu_v \mu_w} \right)$$

It is worth mentioning that, the complexity of the computation of the approximate confidence interval  $\widehat{CI}_k^{1-\alpha}$  is  $O(n)$  with  $n$  the size of entities collection  $G_E$ .

## 4.2 Pruning the Search Space

**Optimistic Estimate on Krippendorff’s Alpha:** to quickly prune unpromising areas of the search space, we define a tight optimistic estimate [18] on Krippendorff’s alpha. We leverage results of Eppstein and Hirschberg [13] who propose an elegant *linear algorithm* Random-SMWA<sup>5</sup> to find subsets with maximum weighted average. Remind that  $A$  can be seen as a weighted average from Eq. (4).

In a nutshell, Random-SMWA seeks to remove  $k$  values to find a subset of  $S$  having  $|S| - k$  values with maximum weighted average. The authors model the problem as such: given  $|S|$  values decreasing linearly with time, find the time at which the  $|S| - k$  maximum values add to zero. In the scope of this work and given  $\sigma_E$  a user defined support threshold on the minimum allowed size of context extents,  $k$  is fixed to  $|S| - \sigma_E$ . The obtained subset correspond to the smallest allowed subset having its support  $\geq \sigma_E$  maximizing the weighted average quantity  $A$ . The Random-SMWA algorithm can be tweaked<sup>6</sup> to retrieve the smallest subset of size  $\geq \sigma_E$  having analogously the minimum possible weighted average quantity  $A$ . We refer to the algorithm returning the maximum (resp. minimum) possible weighted average by RandomSMWA<sup>max</sup> (resp. RandomSMWA<sup>min</sup>).

**Proposition 2 (Upper and Lower bounds for  $A$ ).** *Given  $S \subseteq G_E$ , minimum context support threshold  $\sigma_E$ , and the following functions:*

$$UB(S) = A(\text{RandomSMWA}^{\max}(S, \sigma_E)) \quad LB(S) = A(\text{RandomSMWA}^{\min}(S, \sigma_E))$$

*we know that  $LB$  (resp.  $UB$ ) is a lower (resp. upper) bound for  $A$ , i.e.:*

$$\forall c, d \in \mathcal{D}_E : c \sqsubseteq d \wedge |G_E^c| \geq |G_E^d| \geq \sigma_E \Rightarrow LB(G_E^c) \leq A(G_E^d) \leq UB(G_E^c)$$

Using these results, we define the optimistic estimate for  $A$  as an interval bounded by the minimum and the maximum  $A$  measure that one can observe from the subsets of a given subset  $S \subseteq G_E$ , that is:  $OE(S, \sigma_E) = [LB(S), UB(S)]$ .

**Nested Confidence intervals for  $A$ :** the desired property between two confidence intervals of the same significance level  $\alpha$  related to respectively  $k_1, k_2$  with  $k_1 \leq k_2$  is that  $CI_{k_1}^{1-\alpha}$  encompasses the  $CI_{k_2}^{1-\alpha}$ . Colloquially speaking, larger samples lead to "narrower" confidence intervals. This property is intuitively plausible since the dispersion of the observed intra-agreement for smaller samples is likely to be higher than the dispersion for larger samples. Having such property allows to prune sub-search space related to a context  $c$  when traversing the search space downward if the optimistic estimate  $OE(G_E^c, \sigma_E) \subseteq CI_{|G_E^c|}^{1-\alpha}$ .

Proving  $CI_{k_2}^{1-\alpha} \subseteq CI_{k_1}^{1-\alpha}$  for  $k_1 \leq k_2$  for the exact confidence interval is an uneasy task, since it requires to derive analytically  $E[\theta_k]$  and  $\text{Var}[\theta_k]$  for any  $1 \leq k \leq n$ . It is worth mentioning that the expected value  $E[\theta_k]$  varies when  $k$  varies. We study such property for the approximate confidence interval  $\widehat{CI}_k^{1-\alpha}$ .

<sup>5</sup>Random-SMWA: Randomized algorithm - Subset with Maximum Weighted Average.

<sup>6</sup>Finding the subset having the minimum weighted average is a dual problem to finding the subset having the maximum weighted average. To solve the former problem using Random-SMWA, we modify the values of  $v_i$  to  $-v_i$  and keep the same weights  $w_i$ .

**Proposition 3 (Minimum cardinality constraint for nested approximate confidence intervals).** *Given a context support threshold  $\sigma_E$  and  $\alpha$ :*

$$\text{if } \sigma_E \geq C^\alpha = \frac{4n\beta_w^2}{z_{1-\frac{\alpha}{2}}^2(\beta_v + \beta_w) + 4\beta_w^2} \text{ then:}$$

$$\forall k_1, k_2 \in \mathbb{N} : \sigma_E \leq k_1 \leq k_2 \Rightarrow \widehat{CI}_{k_2}^{1-\alpha} \subseteq \widehat{CI}_{k_1}^{1-\alpha}$$

Combining Properties (1), (2) and (3), we formalize the pruning region property which answers to: *when to prune the sub-search space under a context  $c$ ?*

**Corollary 1 (Pruning regions).** *Given a behavioral dataset  $\mathcal{B}$ , a context support threshold  $\sigma_E \geq C^\alpha$ , a significance critical value  $\alpha \in ]0, 1]$ , For any  $c, d \in \mathcal{D}_E$  such that  $c \sqsubseteq d$  with  $|G_E^c| \geq |G_E^d| \geq \sigma_E$ , we have:*

$$OE(G_E^c, \sigma_E) \subseteq \widehat{CI}_{|G_E^c|}^{1-\alpha} \Rightarrow A(G_E^d) \in \widehat{CI}_{|G_E^d|}^{1-\alpha} \Rightarrow p\text{-value}(d) > \alpha$$

**Proofs:** all proofs of propositions and properties can be found in Appendix A.

## 5 On handling variability of outcomes among raters

In Section 4, we defined the confidence interval  $CI^{1-\alpha}$  established over the DFD. By taking into consideration the variability induced by the selection of a subset of entities, Such confidence interval enables to avoid reporting subgroups indicating an intra-agreement likely (w.r.t. the critical value  $\alpha$ ) to be observed by a random subset of entities. For a more statistically sound results, one should not only take into account the variability induced by the selection of subsets of entities, but also the variability induced by the outcomes of the selected group of individuals. This is well summarized by Hayes and Krippendorff [22] “The obtained value of  $A$  is subject to random sampling variability—specifically variability attributable to the selection of units (i.e. entities) in the reliability data (i.e. behavioral data) and the variability of their judgments”. To address these two questions, they recommend to employ a standard Efron & Tibshirani *bootstrapping approach* [12] to empirically generate the sampling distribution of  $A$  and produce an empirical confidence interval  $CI_{\text{bootstrap}}^{1-\alpha}$ .

Recall that we consider here a behavioral dataset  $\mathcal{B}$  reduced to the outcomes of a selected group of individual  $g$ . Following the bootstrapping scheme proposed by Krippendorff [29,22,43], the empirical confidence interval is computed by repeatedly performing the following steps: (1) resample  $n$  entities from  $G_E$  with replacement; (2) for each sampled entity, draw uniformly  $n_e \cdot (n_e - 1)$  pairs of outcomes according to the distribution of the observed pairs of outcomes; (3) compute the disagreement observed and calculate the Krippendorff alpha quantity on the resulting resample. This process, repeated  $b$  times leads to a vector of bootstrap estimates (sorted in ascending order)  $\hat{B} = [\hat{A}_1, \dots, \hat{A}_b]$ . Given the empirical distribution  $\hat{B}$ , the empirical confidence interval  $CI_{\text{bootstrap}}^{1-\alpha}$  is defined by the percentiles of  $\hat{B}$ . i.e.  $CI_{\text{bootstrap}}^{1-\alpha} = [\hat{B}_{\frac{\alpha}{2} \cdot b}, \hat{B}_{(1-\frac{\alpha}{2}) \cdot b}]$ . We denote by  $MCI^{1-\alpha}$  (Merged CI) the confidence interval that takes into consideration both  $CI^{1-\alpha} = [le_1, re_1]$  and  $CI_{\text{bootstrap}}^{1-\alpha} = [le_2, re_2]$ . We have  $MCI^{1-\alpha} = [\min(le_1, le_2), \max(re_1, re_2)]$ .

## 6 A Branch-and-bound solution: Algorithm DEvIANT

**Subgroup Enumeration.** In order to detect exceptional contextual intra-group agreement patterns, we need to enumerate candidate  $p = (g, c) \in (\mathcal{D}_I, \mathcal{D}_E)$ . For this task, different enumeration algorithms exist in the literature ranging from heuristic (e.g. beam-search [32]) to exhaustive techniques (e.g. GP-growth [33]). In this paper, we choose to exhaustively enumerate all candidate subgroups while leveraging closure operators [15] (since  $A$  computation only depends on the extent of a pattern). This makes it possible to avoid redundancy and to substantially reduce the number of visited patterns. With this aim in mind, and since the data we deal with are of the same format as those handled in our previous work [2], we apply EnumCC (Enumerate Closed Candidates) [2] to enumerate subgroups  $g$  (resp.  $c$ ) in  $\mathcal{D}_I$  ( resp.  $\mathcal{D}_E$ ). EnumCC goes in the same line of CloseByOne algorithm [31]. Given  $G$  a collection of records (which can be either  $G_E$  or  $G_I$ ), EnumCC traverses the search space in a DFS fashion and enumerates once and only once all the closed descriptions that fulfill the minimum support constraint  $\sigma$ . For more details, see Appendix B.

**DEvIANT (Algorithm 1)** implements an efficient branch-and-bound algorithm to Discover statistically significant Exceptional Intra-group Agreement patterns while leveraging closure, tight optimistic estimates and pruning properties. DEvIANT starts by selecting a group  $g$  of individual. Next, the corresponding behavioral dataset  $\mathcal{B}^g$  is established by reducing the original behavioral dataset  $\mathcal{B}$  to elements concerning solely the individuals comprising  $G_I^g$ . Subsequently, the bootstrap confidence interval  $\text{CI}_{\text{bootstrap}}^{1-\alpha}$  is calculated.

---

### Algorithm 1: DEvIANT( $\mathcal{B}, \sigma_E, \sigma_I, \alpha$ )

---

**Inputs :**  $\mathcal{B} = \langle G_I, G_E, O, o \rangle$  is a behavioral dataset,  
 $\sigma_E$  (resp.  $\sigma_I$ ) minimum support threshold of a context (resp. group),  
 $\alpha$  is a critical significance value (fixed to 0.05 in default setting).

**Output:**  $P$  is the set of exceptional intra-group agreement patterns

- 1  $P \leftarrow \{\}$
- 2 **foreach**  $(g, G_I^g, cont_g) \in \text{EnumCC}(G_I, *, \sigma_I, 0, True)$  **do**
- 3      $G_E(g) = \{e \in E \text{ s.t. } n_e^g \geq e\}$
- 4      $\mathcal{B}^g = \langle G_E(g), G_I^g, O, o \rangle$
- 5      $\text{CI}_{\text{bootstrap}}^{1-\alpha} = [\hat{B}_{\frac{\alpha}{2}} \cdot b, \hat{B}_{(1-\frac{\alpha}{2})} \cdot b]$      ▷ With  $\hat{B} = [\hat{A}_1^g, \dots, \hat{A}_b^g]$  computed on
- 6      $\sigma_E^g = \max(C^\alpha(g), \sigma_E)$      respectively  $b$  resamples of  $\mathcal{B}^g$
- 7     **foreach**  $(c, G_E^c, cont_c) \in \text{EnumCC}(G_E(g), *, \sigma_E^g, 0, True)$  **do**
- 8          $\text{MCI}_{|G_E^c|}^{1-\alpha} = \text{merge}(\widehat{CI}_{|G_E^c|}^{1-\alpha}, \text{CI}_{\text{bootstrap}}^{1-\alpha})$
- 9         **if**  $\text{OE}(G_E^c, \sigma_E^g) \subseteq \text{MCI}_{|G_E^c|}^{1-\alpha}$  **then**
- 10              $cont_c \leftarrow \text{False}$      ▷ Prune the unpromising search subspace under  $c$
- 11         **else if**  $A^g(G_E^c) \notin \text{MCI}_{|G_E^c|}^{1-\alpha}$  **then**
- 12              $p_{\text{new}} \leftarrow (g, c)$
- 13             **if**  $\nexists p_{\text{old}} \in P \text{ s.t. } \text{ext}(p_{\text{new}}) \subseteq \text{ext}(p_{\text{old}})$  **then**
- 14                  $P \leftarrow (P \cup p_{\text{new}}) \setminus \{p_{\text{old}} \in P \mid \text{ext}(p_{\text{old}}) \subseteq \text{ext}(p_{\text{new}})\}$
- 15              $cont_c \leftarrow \text{False}$      ▷ Prune the sub search space, generality concept
- 16 **return**  $P$

---

Before searching for exceptional contexts, the minimum context support threshold  $\sigma_E$  is adjusted to  $C^\alpha(g)$  (see Prop.3) if it is lower than  $C^\alpha(g)$ . Note that  $C^\alpha(g)$  is, in practice, much smaller than  $\sigma_E$ . Still, we keep this correction for algorithm soundness. Next, contexts are enumerated by EnumCC and for each candidate context  $c$ , the algorithm evaluates the optimistic estimate interval  $OE(G_E^c)$  (see Prop. 2). According to Corollary 1, if  $OE(G_E^c, \sigma_E^g)$  is inside  $MCI_{|G_E^c|}^{1-\alpha}$  the sub-search space under  $c$  can be pruned. Otherwise,  $A^g(G_E^c)$  is computed and evaluated against  $MCI_{|G_E^c|}^{1-\alpha}$ . If it is outside  $MCI_{|G_E^c|}^{1-\alpha}$ , this means that  $(g, c)$  is significant and should be kept in the final result set  $P$ . Eventually, to reduce the number of reported patterns, we choose to keep only the most general patterns while ensuring that each significant pattern in  $\mathcal{P}$  is represented by a pattern in  $P$ . This formally translates to:  $\forall p' = (g', c') \in \mathcal{P} \setminus P : p\text{-value}^{g'}(c') \leq \alpha \Rightarrow \exists p = (g, c) \in P$  s.t.  $\text{ext}(q) \subseteq \text{ext}(p)$ , with  $\text{ext}(q = (g', c')) \subseteq \text{ext}(p = (g, c))$  defined by  $G_I^{g'} \subseteq G_I^g$  and  $G_E^{c'} \subseteq G_E^c$ . This is based on the following postulate: the end-user is more interested by exceptional (dis)agreement within larger groups and/or for larger contexts rather than local exceptional (dis)agreement. Moreover, the end-user can always refine her analysis to obtain more fine-grained results by re-launching the algorithm starting from a specific context or group.

## 7 Empirical Evaluation

We report on both quantitative and qualitative experiments over the implemented algorithms. For reproducibility purposes, the source code and the data are made available in our companion page<sup>7</sup>. The following experiments aim to answer the following questions: **(Q<sub>1</sub>)** How well the Taylor approximate CI approaches the empirical CI? **(Q<sub>2</sub>)** How efficient is the Taylor approximate CI and the pruning properties? **(Q<sub>3</sub>)** Does DEVIANT provide interpretable patterns?

**Datasets:** experiments were carried on four real-world behavioral datasets (see Tab. 3). Two voting datasets (EPD8 and CHUS) and two rating datasets (Movielens and Yelp). Each dataset features entities and individuals that are described by categorical (C), numerical (N) attributes, or categorical attributes augmented with a taxonomy (H). We report also the equivalent number of items (in an item-set language) corresponding to the descriptive attributes (ordinal scaling [16]).

	$ G_E $	$\mathcal{A}_E$ (Items-Scaling)	$ G_I $	$\mathcal{A}_I$ (Items-Scaling)	Outcomes	Sparsity	$C^{0.05}$
EPD8 <sup>8</sup>	4704	1H + 1N + 1C (437)	848	3C (82)	3.1M (C)	78.6%	$\simeq 10^{-7}$
CHUS <sup>9</sup>	17350	1H + 2N (307)	1373	2C (261)	3M (C)	31.2%	$\simeq 10^{-6}$
Movielens <sup>10</sup>	1681	1H + 1N (161)	943	3C (27)	100K (O)	06.3%	$\simeq 0.06$
Yelp <sup>11</sup>	127K	1H + 1C (851)	1M	3C (6)	4.15M (O)	0.003%	$\simeq 1.73$

**Table 3:** Main characteristics of the behavioral datasets.  $C^{0.05}$  represent the minimum context support threshold over which we have nested approximate CI property.

<sup>7</sup><https://github.com/Adnene93/Deviant>

<sup>8</sup>Eight European Parliament Voting Dataset.

<sup>9</sup>102<sup>nd</sup>-115<sup>th</sup> congresses of the US House of representatives (Period: 1991-2015).

<sup>10</sup>Movie review dataset - <https://grouplens.org/datasets/movielens/100k/>

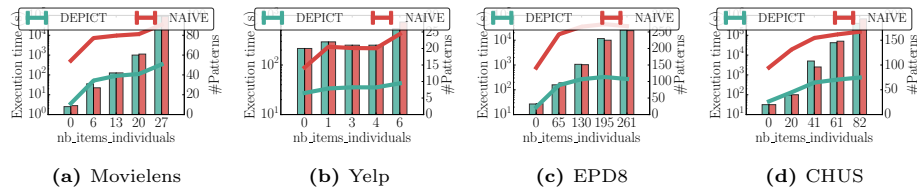
<sup>11</sup>Social network dataset - <https://www.yelp.com/dataset/challenge>

**Q<sub>1</sub>**- First, we evaluate to what extent the confidence interval computed empirically approximates the confidence interval computed by Taylor approximations. For this task, we run 1000 experiments for different subsets size  $k$  uniformly distributed in  $[1, n = |G_E|]$ . For each  $k$ , we compute the corresponding Taylor approximate  $\widehat{CI}_k^{1-\alpha} = [a^T, b^T]$  and empirical confidence interval  $ECI_k^{1-\alpha} = [a^E, b^E]$ . The latter CI is calculated by running  $10^4$  samples of size  $k$  from  $G_E$ , followed by the computation of the observed  $A$  on each sample which are used to estimates the moments of the empirical distribution required for the establishment  $ECI_k^{1-\alpha}$ . Once both CIs are computed, we measure the distance between them by using the Jaccard index, i.e.  $\text{dist}(ECI_k^{1-\alpha}, \widehat{CI}_k^{1-\alpha}) = 1 - \frac{(\min(b^E, b^T) - \max(a^E, a^T))}{(\max(b^E, b^T) - \min(a^E, a^T))}$ . We report in Tab. 4, the average  $\mu_{\text{err}}$  and the standard deviation  $\sigma_{\text{err}}$  of the observed distances (coverage error) over the 1000 experiments. We notice that the difference between the analytic Taylor approximation and the empirical approximation is negligible ( $\mu_{\text{err}}$  is less than  $10^{-2}$ ). Therefore, the CIs approximated by the two methods are so close, that it does not matter which method is used. Hence, the choice is guided by the computational efficiency.

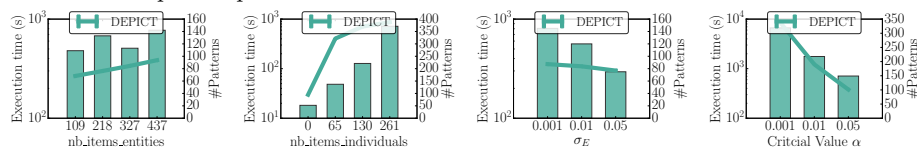
**Q<sub>2</sub>**- In order to evaluate how efficient are the proposed properties ((i) Taylor approximate CI, (ii) optimistic estimates and (iii) nested approximates CIs), we choose to compare DEVIANT against a **Naive** approach where the three aforementioned properties are disabled. For a fair comparison, **Naive** pushes monotonic constraints (minimum support threshold) and employs closure operators while estimating empirically the confidence interval by successive random

$\mathcal{B}$	$\mu_{\text{err}}$	$\sigma_{\text{err}}$	$\ \mathcal{B}$	$\mu_{\text{err}}$	$\sigma_{\text{err}}$	$\ \mathcal{B}$	$\mu_{\text{err}}$	$\sigma_{\text{err}}$	$\ \mathcal{B}$	$\mu_{\text{err}}$	$\sigma_{\text{err}}$
CHUS	0.007	0.004	EPD8	0.007	0.004	Movielens	0.0075	0.0045	Yelp	0.008	0.007

**Table 4:** Coverage error between empirical CIs and Taylor CIs.



**Fig. 2:** Comparison between DEVIANT and **Naive** when varying the size of the description space  $\mathcal{D}_I$ . Lines correspond to the execution time and bars correspond to the number of outputted patterns. Parameters:  $\sigma_E = \sigma_I = 1\%$  and  $\alpha = 0.05$ .



**Fig. 3:** Effectiveness of DEVIANT on EPD8 when varying sizes of both search spaces  $\mathcal{D}_E$  and  $\mathcal{D}_I$ , minimum context support threshold  $\sigma_E$  and the critical value  $\alpha$ . Default parameters: full search spaces  $\mathcal{D}_E$  and  $\mathcal{D}_I$ ,  $\sigma_E = 0.1\%$ ,  $\sigma_I = 1\%$  and  $\alpha = 0.05$ .

trials from  $F_k$ . For this study, we choose to run both algorithms while disabling the bootstrap  $CI_{\text{bootstrap}}^{1-\alpha}$  computation since the overhead induced by the computation of  $CI_{\text{bootstrap}}^{1-\alpha}$  is the same for both algorithms. We vary the size of the descriptions space related to groups of individuals  $\mathcal{D}_I$  while considering the whole description space of entities. Results of this experiment are reported in Fig.2. We observe that DEvIANT outperforms Naive in terms of runtime by nearly two orders of magnitude while outputting the same number of the desired patterns.

Furthermore, we report in Fig.3 the performance of DEvIANT in terms of runtime and the number of outputted patterns. We observe that when varying the descriptions space sizes, DEvIANT requires more time to finish. It is worth mentioning that the size of individuals search space  $\mathcal{D}_I$  substantially affects the runtime of DEvIANT. This is mainly due to the fact that larger  $\mathcal{D}_I$  leads to more candidate group of individuals  $g$  which require: (i) DEvIANT to generate the bootstrapping confidence interval  $CI_{\text{bootstrap}}^{1-\alpha}$  and (ii) to mine for exceptional contexts  $c$  concerning the candidate group  $g$ . Finally, we observe that when  $\alpha$  decreases, the execution time required for DEvIANT to finish increases while returning more patterns. This may at first, seem counter-intuitive since less patterns are naturally considered significant when alpha decreases, but this is supported by the fact that DEvIANT considers in the resulting pattern set only the most general patterns. Hence, when  $\alpha$  decreases, DEvIANT goes deeper in the contexts search space, implying thus much more candidate patterns to be tested and thus a larger results set. The same conclusions can be drawn as well from experiments performed on Yelp, MovieLens and CHUS (see Appendix. C).

**Q<sub>3</sub>**- We illustrate a number of examples of the outputted exceptional contextual intra-group agreement patterns in the benchmark datasets. Table 5 reports the exceptional contexts observed among the republicans party during the 115<sup>th</sup> congress. For instance, Pattern  $p_1$  illustrated in Fig. 4, highlights a collection of voting sessions addressing Government and Administrative issues where a clear polarization is observed between two clusters of House republicans. Notable roll call vote of this context in which a significant disagreement was observed between republicans is “**House Vote 417**”<sup>12</sup> which was closely watched by the media<sup>13</sup>.

Table 6 depicts some patterns returned by DEvIANT when carried on MovieLens Datasets. For instance, pattern  $p_2$  reports that “Middle-aged Men” group observe a significantly higher intra-group agreement compared to the overall intra-group agreement for movies labeled with both adventure and musical genres (e.g. The Wizard of Oz (1939)).

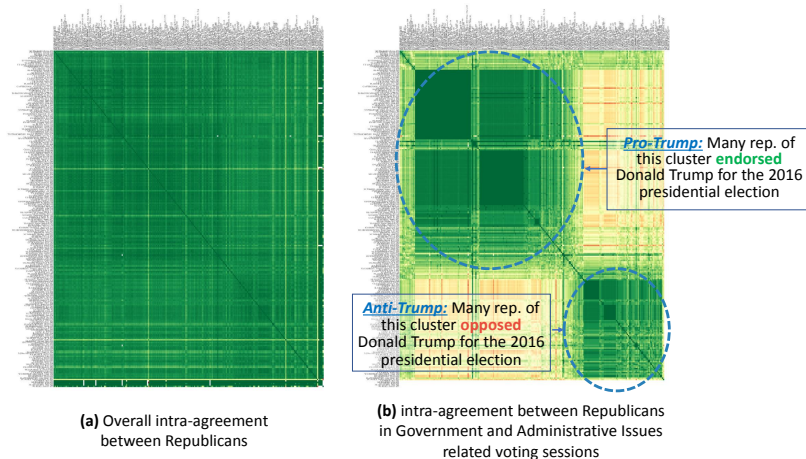
id	group ( $g$ )	context ( $c$ )	$A^g(*)$	$A^g(c)$	$p$ -value	IA
$p_1$	Republicans	20.11 Government Branch Relations, Admin. Issues, and Constitutional Reforms	0.83	0.32	< .001	Conflict
$p_2$	Republicans	5 Labor	0.83	0.63	< .01	Conflict
$p_3$	Republicans	20.05 Nominations and Appointments	0.83	0.92	< .001	Consensus

**Table 5:** Exceptional consensual/conflictual subjects among Republicans Party representatives in the 115<sup>th</sup> congress of the US House of Representatives.  $\alpha = 0.01$

<sup>12</sup><https://projects.propublica.org/represent/votes/115/house/1/417>

<sup>13</sup>Washington Post:<https://wapo.st/2W32I9c>; Reuters:<https://reut.rs/2TF0dgV>





**Fig. 4:** Illustrating Pattern 1 from Tab. 5 with a similarity matrix between Republicans. Each cell represents the ratio of voting sessions in which both Rep. agreed. A green cell reports a strong agreement whereas a red cell highlights a strong disagreement.

id	group ( $g$ )	context ( $c$ )	$A^g(*)$	$A^g(c)$	$p$ -value	IA
$p_1$	Old	1.Action & 2.Adventure & 6.Crime Movies	-0.06	-0.29	< 0.01	Conflict
$p_2$	Middle-aged Men	2.Adventure & 12.Musical Movies	0.05	0.21	< 0.01	Consensus
$p_3$	Old	4.Children & 12.Musical Movies	-0.06	-0.21	< 0.01	Conflict

**Table 6:** Top3-Exceptionally consensual/conflictual movies genres between MovieLens’ raters,  $\alpha=0.01$ . Patterns are ranked by the absolute difference between  $A^g(c)$  and  $A^g(*)$ .

## 8 Conclusion and Future Directions

In this paper we introduced the novel problem of discovering statistically significant exceptional contextual intra-group agreement patterns. We devised a branch-and-bound algorithm, named DEVIANT, which efficiently search for the desired patterns while leveraging closure operators, approximate confidence intervals (CIs), tight optimistic estimates on Krippendorff’s Alpha measure and the property of nested CIs. The empirical experiments demonstrated both the efficiency and the usefulness of DEVIANT over multiple behavioral datasets relevant to various domains ranging from political analysis to rating data analysis. In future research, we plan (i) to incorporate FDR (False Discovery Rate) control to tackle the multiple comparison problem [21] (ii) to investigate exceptional intra-group agreement compared to the one observed between all individuals over the same context and (iii) to integrate the option of choosing which kind of exceptional consensus the end-user want, i.e. is the exceptional consensus is observed because the group members “liked/voted for” (or “disliked/voted against”) the context related entities? All this being done within the perspective to provide a comprehensive framework and tool<sup>14</sup>) for behavioral data analysis alongside exceptional inter-group agreement pattern discovery implemented in [2].

<sup>14</sup>A prototype is available online in <http://contentcheck.liris.cnrs.fr>



## References

1. S. Amer-Yahia, S. Kleisarchaki, N. K. Kolloju, L. V. Lakshmanan, and R. H. Zamar. Exploring rated datasets with rating maps. In *WWW*, pages 1411–1419. International World Wide Web Conferences Steering Committee, 2017.
2. A. Belfodil, S. Cazalens, P. Lamarre, and M. Plantevit. Flash points: Discovering exceptional pairwise behaviors in vote or rating data. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 442–458. Springer, 2017.
3. M. Boley, T. Horváth, 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.
4. J. Cohen. A coefficient of agreement for nominal scales. *Education and Psychological Measurement*, 20:37–46, 1960.
5. T. M. Cover and J. A. Thomas. *Elements of information theory*. John Wiley & Sons, 2012.
6. M. Das, S. Amer-Yahia, G. Das, and C. Yu. Mri: Meaningful interpretations of collaborative ratings. *PVLDB*, 4(11):1063–1074, 2011.
7. T. De Bie. An information theoretic framework for data mining. In *KDD*, pages 564–572. ACM, 2011.
8. W. Duivesteijn, A. J. Feelders, and A. Knobbe. Exceptional model mining. *Data Mining and Knowledge Discovery*, 30(1):47–98, 2016.
9. W. Duivesteijn and A. Knobbe. Exploiting false discoveries—statistical validation of patterns and quality measures in subgroup discovery. In *Data Mining (ICDM), 2011 IEEE 11th International Conference on*, pages 151–160. IEEE, 2011.
10. W. Duivesteijn, A. J. Knobbe, A. Feelders, and M. van Leeuwen. Subgroup discovery meets bayesian networks - an exceptional model mining approach. *ICDM*, 2010.
11. F. Duris, J. Gazdarica, I. Gazdaricova, L. Strieskova, J. Budis, J. Turna, and T. Szemes. Mean and variance of ratios of proportions from categories of a multinomial distribution. *Journal of Statistical Distributions and Applications*, 5(1):2, 2018.
12. B. Efron and R. J. Tibshirani. *An introduction to the bootstrap*. CRC press, 1994.
13. D. Eppstein and D. S. Hirschberg. Choosing subsets with maximum weighted average. *J. Algorithms*, 24(1):177–193, 1997.
14. J. L. Fleiss. Measuring nominal scale agreement among many raters. *Psychological Bulletin*, 76(5):378–382, 1971.
15. B. Ganter and S. Kuznetsov. Pattern structures and their projections. *ICCS*, 2001.
16. B. Ganter and R. Wille. *Formal concept analysis - mathematical foundations*. Springer, 1999.
17. S. Geisser. *Predictive Inference*, volume 55. CRC Press, 1993.
18. H. Grosskreutz, S. Rüping, and S. Wrobel. Tight optimistic estimates for fast subgroup discovery. In *ECML/PKDD (1)*, volume 5211 of *Lecture Notes in Computer Science*, pages 440–456. Springer, 2008.
19. W. Hämaläinen. Statapriori: an efficient algorithm for searching statistically significant association rules. *Knowl. Inf. Syst.*, 23(3):373–399, 2010.
20. W. Hämaläinen and G. I. Webb. Statistically sound pattern discovery. In *KDD*, page 1976. ACM, 2014.
21. W. Hämaläinen and G. I. Webb. A tutorial on statistically sound pattern discovery. *Data Min. Knowl. Discov.*, 33(2):325–377, 2019.

22. A. F. Hayes and K. Krippendorff. Answering the call for a standard reliability measure for coding data. *Communication methods and measures*, 1(1):77–89, 2007.
23. M. Kaytoue, S. O. Kuznetsov, A. Napoli, and S. Duplessis. Mining gene expression data with pattern structures in formal concept analysis. *Information Sciences*, 181(10):1989–2001, 2011.
24. G. M. P. van Kempen and L. J. van Vliet. Mean and variance of ratio estimators used in fluorescence ratio imaging. *Cytometry: The Journal of the International Society for Analytical Cytology*, 39(4):300–305, 2000.
25. M. Kendall, A. Stuart, and J. Ord. Kendall’s advanced theory of statistics. v. 1: Distribution theory. 1994.
26. A. Kleiner, A. Talwalkar, P. Sarker, and M. I. Jordan. The big data bootstrap. In *ICML*. icml.cc / Omnipress, 2012.
27. K. Krippendorff. Estimating the reliability, systemic error and random error of interval data. *Educational and Psychological Measurement*, 30(1):61–70, 1970.
28. K. Krippendorff. Content analysis, an introduction to its methodology. 2004.
29. K. Krippendorff. Bootstrapping distributions for krippendorff’s alpha. *Available at web. asc. upenn. edu/usr/krippendorff/boot. c-Alpha. pdf*, 2006.
30. K. Krippendorff. Computing krippendorff’s alpha-reliability. 2011.
31. S. O. Kuznetsov and S. A. Obiedkov. Comparing performance of algorithms for generating concept lattices. *Journal of Experimental & Theoretical Artificial Intelligence*, 14(2-3):189–216, 2002.
32. M. van Leeuwen and A. J. Knobbe. Diverse subgroup set discovery. *Data Min. Knowl. Discov.*, 25(2):208–242, 2012.
33. F. Lemmerich, M. Becker, and M. Atzmueller. Generic pattern trees for exhaustive exceptional model mining. In *ECML/PKDD*, pages 277–292, 2012.
34. F. Lemmerich, M. Becker, P. Singer, D. Helic, A. Hotho, and M. Strohmaier. Mining subgroups with exceptional transition behavior. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 965–974. ACM, 2016.
35. J. Lijffijt, B. Kang, W. Duivesteyn, K. Puolamäki, E. Oikarinen, and T. De Bie. Subjectively interesting subgroup discovery on real-valued targets. In *ICDE*, pages 1352–1355. IEEE Computer Society, 2018.
36. S. Minato, T. Uno, K. Tsuda, A. Terada, and J. Sese. A fast method of statistical assessment for combinatorial hypotheses based on frequent itemset enumeration. In *ECML/PKDD (2)*, volume 8725 of *Lecture Notes in Computer Science*, pages 422–436. Springer, 2014.
37. N. Mukhopadhyay. *Probability and statistical inference*. CRC Press, 2000.
38. B. Omidvar-Tehrani, S. Amer-Yahia, P.-F. Dutot, and D. Trystram. Multi-objective group discovery on the social web. In *ECMLPKDD*, 2016.
39. J. M. Steele. *The Cauchy-Schwarz master class: an introduction to the art of mathematical inequalities*. Cambridge University Press, 2004.
40. A. Terada, M. Okada-Hatakeyama, K. Tsuda, and J. Sese. Statistical significance of combinatorial regulations. *Proceedings of the National Academy of Sciences*, 110(32):12996–13001, 2013.
41. G. I. Webb. Discovering significant patterns. *Machine learning*, 68(1):1–33, 2007.
42. S. Wrobel. An algorithm for multi-relational discovery of subgroups. PKDD, 1997.
43. A. Zapf, S. Castell, L. Morawietz, and A. Karch. Measuring inter-rater reliability for nominal data— which coefficients and confidence intervals are appropriate? *BMC Medical Research Methodology*, 16(1):93, 2016.

## A Appendix: Proofs

Recall that  $\theta_k : F_k \rightarrow \mathbb{R}$  is the random variable corresponding to the observed intra-agreement  $A$  (Krippendorff's alpha) of subsets  $S \in G_E$  of size  $k$ . i.e. for any  $k \in [1, n]$  with  $n = |G_E|$  we have  $\theta_k(S \in F_k) = A(S)$  and  $F_k = \{S \in G_E \text{ s.t. } |S| = k\}$ .  $F_k$  is then the set of possible outcomes which are equally likely to occur under the null hypothesis  $H_0$ .  $G_E$  contains  $n$  records (i.e.  $|G_E| = n$ ). Each record  $e \in G_E$  is associated to a value  $v_e$  and  $w_e$ .  $\theta_k$  can be expressed as a ratio of two random variable  $\frac{V_k}{W_k}$ . With  $V_k, W_k$  two random variables  $V_k : F_k \rightarrow \mathbb{R}$  and  $W_k : F_k \rightarrow \mathbb{R}$  with  $V_k(S) = \frac{1}{k} \sum_{e \in S} v_e$  and  $W_k(S) = \frac{1}{k} \sum_{e \in S} w_e$ .

*Proof (Proposition 1).* For any  $f(x, y)$ , the bivariate second order Taylor expansion about any  $\lambda = (\lambda_x; \lambda_y)$  is (a concise lecture note<sup>15</sup> follows the same reasoning and explains the derivations) :

$$\begin{aligned} f(x, y) &= f(\lambda) + f'_x(\lambda)(x - \lambda_x) + f'_y(\lambda)(y - \lambda_y) \\ &+ \frac{1}{2} (f''_{xx}(\lambda)(x - \lambda_x)^2 + f''_{xy}(\lambda)(x - \lambda_x)(y - \lambda_y) + f''_{yy}(\lambda)(y - \lambda_y)^2) + \epsilon \end{aligned} \quad (7)$$

with  $\epsilon$  is a remainder of smaller order than the term of the equation.

An approximation of the expectation  $E[f(x, y)]$  expanded around  $\lambda = (\lambda_x; \lambda_y)$  is:

$$E[f(x, y)] \approx f(\lambda) + \frac{1}{2} [f''_{xx}(\lambda)\text{Var}[X] + f''_{xy}(\lambda)\text{Cov}[X, Y] + f''_{yy}(\lambda)\text{Var}[Y]] \quad (8)$$

Given that  $f(x, y) = \frac{x}{y}$  and using the fact that  $E[X - \mu_x] = 0$  (This is valid for both  $V$  and  $W$ ). We have:  $\text{Var}[X] = E[(X - \mu_x)^2]$  and  $\text{Cov}[X, Y] = E[(X - \mu_x)(Y - \mu_y)]$ . We can derive an approximation of  $E[\theta_k] = E[\frac{V_k}{W_k}]$  around  $(\mu_{V_k}, \mu_{W_k})$ .

$$E[\theta_k] = E\left[\frac{V_k}{W_k}\right] = E[f(V_k, W_k)] \approx \frac{\mu_{V_k}}{\mu_{W_k}} - \frac{\text{Cov}[V_k, W_k]}{\mu_{W_k}^2} + \frac{\text{Var}[W_k]\mu_{V_k}}{\mu_{V_k}^3} \quad (9)$$

The formulas of  $E[V_k]$  (resp.  $E[W_k]$ ) and  $\text{Var}[V_k]$  (resp.  $\text{Var}[W_k]$ ) can be derived analytically. We denote by  $\mu_v$  (resp.  $\mu_w$ ) the arithmetic mean of the values (resp. weights) corresponding to each entity  $e \in G_E$ . i.e:  $\mu_v = \frac{1}{n} \sum_{e \in G_E} v_e$  and  $\mu_w = \frac{1}{n} \sum_{e \in G_E} w_e$  with  $n = |G_E|$ .

$$E[V_k] = \frac{1}{\binom{n}{k}} \sum_{S \in F_k} \frac{1}{k} \sum_{e \in S} v_e = \frac{1}{n} \sum_{e \in G_E} v_e = \mu_v \quad (10)$$

$$\begin{aligned} \text{Var}[V_k] &= \frac{1}{\binom{n}{k}} \sum_{S \in F_k} \left( \frac{1}{k} \sum_{e \in S} v_e - E[V_k] \right)^2 = \frac{1}{\binom{n}{k}} \sum_{S \in F_k} \left( \frac{1}{k} \sum_{e \in S} v_e - \mu_v \right)^2 \\ &= \frac{1}{k} \left( \frac{n}{n-1} (\mu_{v^2} - \mu_v^2) \right) - \frac{1}{n-1} (\mu_v^2 - \mu_{v^2}) \text{ with } \mu_{v^2} = \frac{1}{n} \sum_{e \in G_E} v_e^2 \end{aligned} \quad (11)$$

<sup>15</sup>see <http://www.stat.cmu.edu/~hseltman/files/ratio.pdf>

Same reasoning applies to compute the expected value and the variance related to  $W_k$ :

$$E[W_k] = \frac{1}{n} \sum_{e \in G_E} w_e = \mu_w \quad (12)$$

$$\begin{aligned} \text{Var}[W_k] &= \frac{1}{\binom{n}{k}} \sum_{S \in F_k} \left( \frac{1}{k} \sum_{e \in S} w_e - E[W_k] \right)^2 \\ &= \frac{1}{k} \left( \frac{n}{n-1} (\mu_{w^2} - \mu_w^2) \right) - \frac{1}{n-1} (\mu_w^2 - \mu_w^2) \quad \text{with } \mu_{w^2} = \frac{1}{n} \sum_{e \in G_E} w_e^2 \end{aligned} \quad (13)$$

We derive now the formula for  $\text{Cov}(V_k, W_k)$ . The same line of reasoning for the computation of the variance of  $V_k$  and  $W_k$  applies. We obtain:

$$\begin{aligned} \text{Cov}[V_k, W_k] &= \frac{1}{\binom{n}{k}} \sum_{S \in F_k} \left( \frac{1}{k} \sum_{e \in S} v_e - E[V_k] \right) \left( \frac{1}{k} \sum_{e \in S} w_e - E[W_k] \right) \\ &= \frac{1}{k} \left( \frac{n}{n-1} (\mu_{vw} - \mu_v \mu_w) \right) - \frac{1}{n-1} (\mu_v \mu_w - \mu_{vw}) \\ &\quad \text{with } \mu_{vw} = \frac{1}{n} \sum_{e \in G_E} w_e v_e \end{aligned} \quad (14)$$

Using equations (10), (11), (12), (13), (14), we derive the approximation of  $E[\theta_k]$  after simplifications of (9).

$$E[\theta_k] \approx \widehat{E}[\theta_k] = \left( \frac{n}{k} - 1 \right) \frac{\mu_v}{\mu_w} \beta_w + \frac{\mu_v}{\mu_w} \quad \text{with } \beta_w = \frac{1}{n-1} \left( \frac{\mu_{w^2}}{\mu_w^2} - \frac{\mu_{vw}}{\mu_v \mu_w} \right) \quad (15)$$

The same reasoning applies to approximate  $\text{Var}[\theta_k]$  using Taylor expansion. We will confine ourselves to a first order Taylor expansion around  $(\mu_v, \mu_w)$  to make the analytic derivation of the approximation of  $\text{Var}[\theta_k]$  feasible. Same observation have been made by [24,11] and [25, p.351] to approximate the variance for a ratio random variable. We obtain:

$$\text{Var}[\theta_k] = \text{Var}[f(V_k, W_k)] \approx \frac{\text{Var}[V_k]}{\mu_{W_k}^2} - 2 \frac{\mu_{V_k} \text{Cov}[V_k, W_k]}{\mu_{W_k}^3} + \frac{\mu_{V_k}^2 \text{Var}[W_k]}{\mu_{W_k}^4} \quad (16)$$

After simplifications and by using the same line of reasoning when deriving the expected value approximation reported in equation (15), we obtain:

$$\begin{aligned} \text{Var}[\theta_k] &\approx \widehat{\text{Var}}[\theta_k] = \left( \frac{n}{k} - 1 \right) \frac{\mu_v^2}{\mu_w^2} (\beta_v + \beta_w) \\ &\quad \text{with } \beta_w = \frac{1}{n-1} \left( \frac{\mu_{w^2}}{\mu_w^2} - \frac{\mu_{vw}}{\mu_v \mu_w} \right) \quad \text{and } \beta_v = \frac{1}{n-1} \left( \frac{\mu_{v^2}}{\mu_v^2} - \frac{\mu_{vw}}{\mu_v \mu_w} \right) \end{aligned} \quad (17)$$

We denote by  $\widehat{CI}_k^{1-\alpha}$  the approximate confidence interval calculated using the approximations of the expected value  $\widehat{E}[\theta_k]$  (15) and the variance  $\widehat{\text{Var}}[\theta_k]$  (17) that is:

$$\widehat{CI}_k^{1-\alpha} = \left[ \widehat{E}[\theta_k] - z_{1-\frac{\alpha}{2}} \sqrt{\widehat{\text{Var}}[\theta_k]}, \widehat{E}[\theta_k] + z_{1-\frac{\alpha}{2}} \sqrt{\widehat{\text{Var}}[\theta_k]} \right] \quad (18)$$

It is worth mentioning that, the complexity of the computation of the approximate confidence interval (18) is linear to the size  $n$ .  $\square$

*Proof (Proposition 2).* To alleviate the text, we will omit  $\sigma_E$  as a parameter in the proof and keep in mind that we consider the minimum support threshold  $\sigma_E$ . Given that  $c \sqsubseteq d$ , with  $c, d$  two descriptions from  $\mathcal{D}$ , we have  $G_E^d \subseteq G_E^c$ . The proposition stems from the fact that:

1.  $A(G_E^c) \leq UB(G_E^c)$  since **RandomSMWA**<sup>max</sup> computes the subset  $S_{\max}^c$  having the maximum weighted average  $A$  as proven by Epstein and Hirschberg [13].
2.  $UB$  is monotonic w.r.t. the partial order  $\subseteq$  between sets. That is:

$$\forall S, S' \subseteq G_E : S' \subseteq S \Rightarrow UB(S') \leq UB(S)$$

This can be proven by reduction to absurdity. We denote by  $S'_{\max} \subset S'$  (resp.  $S_{\max} \subset S$ ) the optimal subset of  $S'$  (resp.  $S$ ) having its size  $\geq \sigma_E$  and the maximum possible weighted average  $A$ . Suppose that  $\exists S, S' \subseteq G_E : S' \subseteq S \wedge UB(S') > UB(S)$  ( $A(S'_{\max}) > A(S_{\max})$ ). Since  $S' \subseteq S$ , this means that there is another subset, namely  $S'_{\max}$ , in  $S$  that observes a greater weighted average  $A$  than the actual optimal subset  $S_{\max}$ , which is absurd.

From (1) and (2) we have:  $A(G_E^d) \leq UB(G_E^d) \leq UB(G_E^c)$ . Same reasoning hold to prove that  $LB$  is a lower bound.  $\square$

*Proof (Proposition 3).* In order to prove the desired property for the approximate confidence intervals, we need to determine first if the variance decreases when  $k$  increases.

$$k_1, k_2 \in \mathbb{N} : \text{if } k_1 \leq k_2 \Rightarrow \widehat{\text{Var}}[\theta_{k_1}] \geq \widehat{\text{Var}}[\theta_{k_2}] \quad (19)$$

From (17),  $\widehat{\text{Var}}[\theta_k] = \left(\frac{n}{k} - 1\right) \frac{\mu_v^2}{\mu_w^2} (\beta_v + \beta_w)$ . Given that  $\frac{n}{k} - 1$  is a decreasing function w.r.t.  $k$ , proving (19) requires that  $\beta_v + \beta_w$  is a positive quantity. This stems from the fact that the original formula of the approximate variance given in (16) is positive. This can be proved by a direct application of the Covariance inequality [37, p, 149], which itself is an application of the Cauchy-Schwarz inequality [39]. Since  $\beta_v + \beta_w$  is of the same sign of (17), we have  $\beta_v + \beta_w \geq 0$ . For the sake of a self-contained proof of this assertion below:

We have, from (16) is of the same sign of:

$$\frac{\text{Var}[V_k]}{\mu_{V_k}^2} - 2 \frac{\text{Cov}[V_k, W_k]}{\mu_{V_k} \mu_{W_k}} + \frac{\text{Var}[W_k]}{\mu_{W_k}^2} \quad (20)$$

From the Covariance inequality, we have  $\text{Cov}[V_k, W_k] \leq \sigma[V_k]\sigma[W_k]$  with  $\sigma^2[V_k] = \text{Var}[V_k]$  and  $\sigma^2[W_k] = \text{Var}[W_k]$ , hence we have the quantity (20) is greater than:

$$\frac{\sigma^2[V_k]}{\mu_{V_k}^2} - 2\frac{\sigma[V_k]\sigma[W_k]}{\mu_{V_k}\mu_{W_k}} + \frac{\sigma^2[W_k]}{\mu_{W_k}^2} \quad (21)$$

This quantity can be rewritten as such:

$$\frac{\sigma[V_k]}{\mu_{V_k}} \left( \frac{\sigma[V_k]}{\mu_{V_k}} - \frac{\sigma[W_k]}{\mu_{W_k}} \right) - \frac{\sigma[W_k]}{\mu_{W_k}} \left( \frac{\sigma[V_k]}{\mu_{V_k}} - \frac{\sigma[W_k]}{\mu_{W_k}} \right) \quad (22)$$

This quantity is clearly positive since (22) is equal to  $\left( \frac{\sigma[V_k]}{\mu_{V_k}} - \frac{\sigma[W_k]}{\mu_{W_k}} \right)^2 \geq 0$ . Hence  $\beta_v + \beta_w \geq 0$  which confirms that the variance is decreasing w.r.t. increasing size  $k$  as stated in equation (19).

Recall that, we want to ensure by approximation that for  $\sigma_E \leq k_1 \leq k_2$  with  $\sigma_E$  a threshold on the context support, we have  $\widehat{CI}_{k_2}^{1-\alpha} \subseteq \widehat{CI}_{k_1}^{1-\alpha}$ . Hence, we need to find the minimum  $\sigma_E$  above which such property is valid. This amounts to find a lower bound for  $\sigma_E$  such that:

$$z_{1-\frac{\alpha}{2}} \sqrt{\widehat{\text{Var}}[\theta_{k_1}]} - z_{1-\frac{\alpha}{2}} \sqrt{\widehat{\text{Var}}[\theta_{k_2}]} \geq \left| \widehat{E}[\theta_{k_1}] - \widehat{E}[\theta_{k_2}] \right| \quad (23)$$

By using the definitions of  $\widehat{\text{Var}}[\theta_k]$  and  $\widehat{E}[\theta_k]$  from equations (15) and (17), the inequality (23) can be rewritten as such:

$$\left( \sqrt{\frac{n}{k_1} - 1} + \sqrt{\frac{n}{k_2} - 1} \right) \leq z_{1-\frac{\alpha}{2}} \sqrt{\frac{\beta_v + \beta_w}{\beta_w^2}} \quad (24)$$

And since  $\sigma_E \leq k_1 \leq k_2$ :

$$2\sqrt{\frac{n}{\sigma_E} - 1} \leq z_{1-\frac{\alpha}{2}} \sqrt{\frac{\beta_v + \beta_w}{\beta_w^2}} \quad (25)$$

After simplifications, we obtain:

$$\sigma_E \geq C^\alpha = \frac{4n\beta_w^2}{z_{1-\frac{\alpha}{2}}^2(\beta_v + \beta_w) + 4\beta_w^2} \quad (26)$$

□

*Proof (Corollary 1).* The proof is a straightforward. From proposition 2, we have that for any  $c, d \in \mathcal{D}_E$  s.t.  $c \sqsubseteq d$ , if  $G^c \geq G^d \geq \sigma_E$  then:

$$A(G_E^d) \in OE(G_E^c, \sigma_E) \quad (27)$$

From proposition 3, if  $\sigma_E \geq C^\alpha$  we have:

$$CI_{|G_E^c|}^{1-\alpha} \subseteq \widehat{CI}_{|G_E^d|}^{1-\alpha} \quad (28)$$

From (27) and (28) and given that  $OE(G_E^c, \sigma_E) \subseteq \widehat{CI}_{|G_E^c|}^{1-\alpha}$ , it follows that  $A(G_E^d) \in OE(G_E^c, \sigma_E) \subseteq \widehat{CI}_{|G_E^c|}^{1-\alpha} \subseteq \widehat{CI}_{|G_E^d|}^{1-\alpha}$  hence  $p\text{-value}(d) > \alpha$ . □

## B Appendix: Enumeration Algorithm

Given a collection of records  $G$  which descriptive attributes are  $\mathcal{A} = \{a_1, \dots, a_m\}$  which can be boolean, numerical or categorical potentially organized among a taxonomy. Attributes  $\mathcal{A}$  allow to structure the search space  $\mathcal{D}$  by considering descriptions  $d \in \mathcal{D}$  which are conjunction of conditions over the attributes domains of interpretation. A condition over a categorical attribute is an equality test while a condition over a numerical attributes is a membership test in an interval.  $G^d$  denotes the set of records of  $G$  covered by the description  $d$ .

EnumCC algorithm enumerates once and only once all closed descriptions whose associated subgroups fulfill the minimum support constraint  $\sigma_E$ . The algorithm follows the same reasoning of most common SD algorithms and goes in the same line of CloseByOne Algorithm (CbO) [31] and Divide-And-Conquer Algorithm [3]. It traverses the search lattice  $\mathcal{D}$  in a top-down (DFS) fashion starting from the most general description  $*$  whose support is the entire collection  $G$ . It proceeds by atomic refinements to progress, step by step, toward more specific descriptions. This is enabled by a refinement operator denoted  $\eta_j$  for the  $j^{th}$  attribute.  $\eta_j$  keeps all conditions related to attributes  $a_i$  for  $i \neq j$  intact, and refines only the  $j^{th}$  condition. If the condition is related to a numerical attribute a left or a right minimal change is performed [23]. If the condition is related to a categorical attribute, return an equality test for all possible values of the domain (if the condition was never refined before), otherwise no refinement is possible. If the attribute is an HMT (categorical attribute augmented with a taxonomy) only one tag is refined to its child or an additional tag is appended [2]. In a nutshell, for each parameter description  $d$ , EnumCC starts by assessing if the subgroup  $G^d$  is valid ( $|G^d| \geq \sigma$ ) and. In this case, the closed description `closure_d` is computed and returned only if the canonicity test is passed (cf. [16, p.66-68]). `closure_d` corresponds to the tightest description of  $G^d$  (maximum in the sense of partial order  $\sqsubseteq$  ordering descriptions in  $\mathcal{D}$ ) which is the

---

### Algorithm 2: EnumCC( $G, d, \sigma_G, f, cnt$ )

---

**Inputs :**  $G$  is the collection of records depicted each by  $m$  attributes,  
 $d$  a description from  $\mathcal{D}$ ,  $\sigma_G$  a support threshold,  
 $f \in [1, m]$  a refinement flag,  $cnt$  a boolean.

**Output:** yields all closed descriptions, i.e.  $clo[\mathcal{D}] = \{clo(d) \text{ s.t. } d \in \mathcal{D}\}$

```

1 if  $|G^d| \geq \sigma$  then
2   closure_d  $\leftarrow \delta(G^d)$            ▷ compute the most specific description of  $G^d$ 
3   if  $d \prec_f \text{closure\_d}$  then
4     cnt_c  $\leftarrow \text{copy}(cnt)$            ▷ cnt_c value can be modified by a caller algorithm
5     yield (closure_d,  $G^{\text{closure\_d}}$ , cnt_c)   ▷ yield results and wait for next call
6     if cnt_c then
7       foreach  $j \in [f, m]$  do
8         foreach  $d' \in \eta_j(\text{closure\_d})$  do
9           foreach  $(nc, G^{nc}, cnt\_nc) \in \text{EnumCC}(G, d', \sigma_G, j, cnt\_c)$  do
10            yield ( $nc, G^{nc}, cnt\_nc$ )

```

---

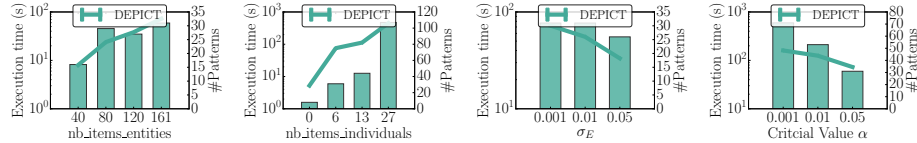
conjunction of all descriptions (conjunction of conditions) related to the records  $g \in G^d$ . Next, if the caller-algorithm allows the algorithm to continue (boolean `cnt_c` kept `True`), the description `closure_d` is refined by starting from the last refined attribute (pointed out by the flag  $f \in [1..m]$ ), since refining preceding attributes will certainly causes the next canonicity test to fail causing the algorithm to backtrack. Eventually, a recursive call is done to explore the sub-search space related to  $d$  (`closure_d`).



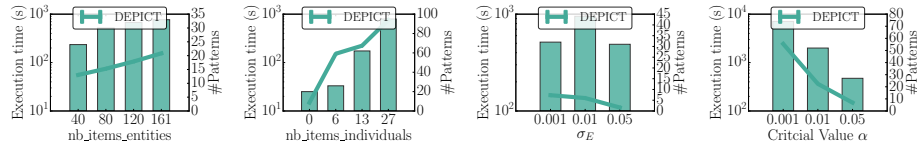
## C Appendix: Additional Experiments

### C.1 Performance evaluation

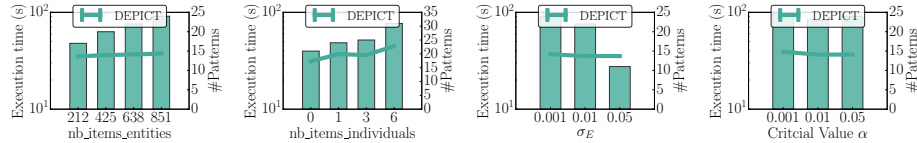
Additional experiments reporting the execution time and the number of reported significant pattern by DEVIANT on Movielens, Yelp, CHUS and EPD8. In this



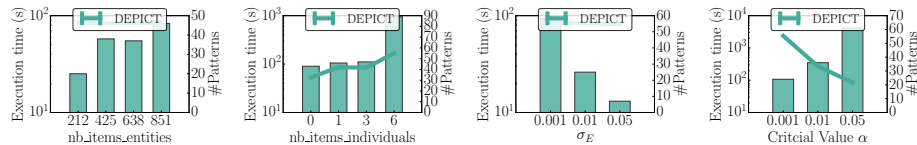
**Fig. 5:** Effectiveness of DEVIANT on Movielens when varying sizes of both search spaces  $\mathcal{D}_E$  and  $\mathcal{D}_I$ , minimum context support threshold  $\sigma_E$  and the critical value  $\alpha$ . Default parameters: full search spaces  $\mathcal{D}_E$  and  $\mathcal{D}_I$ ,  $\sigma_E = 0.1\%$ ,  $\sigma_I = 1\%$  and  $\alpha = 0.05$ . Bootstrapping Confidence intervals for handling variability of outcomes is disabled



**Fig. 6:** Effectiveness of DEVIANT on Movielens when varying sizes of both search spaces  $\mathcal{D}_E$  and  $\mathcal{D}_I$ , minimum context support threshold  $\sigma_E$  and the critical value  $\alpha$ . Default parameters: full search spaces  $\mathcal{D}_E$  and  $\mathcal{D}_I$ ,  $\sigma_E = 0.1\%$ ,  $\sigma_I = 1\%$  and  $\alpha = 0.05$ . Bootstrapping Confidence intervals for handling variability of outcomes is enabled

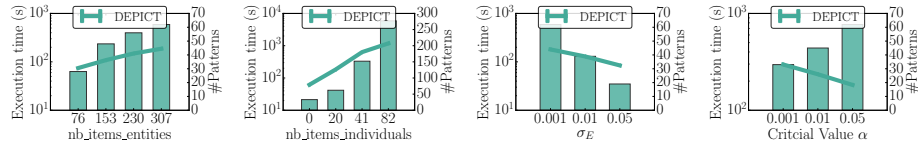


**Fig. 7:** Effectiveness of DEVIANT on Yelp when varying sizes of both search spaces  $\mathcal{D}_E$  and  $\mathcal{D}_I$ , minimum context support threshold  $\sigma_E$  and the critical value  $\alpha$ . Default parameters: full search spaces  $\mathcal{D}_E$  and  $\mathcal{D}_I$ ,  $\sigma_E = 0.1\%$ ,  $\sigma_I = 1\%$  and  $\alpha = 0.05$ . Bootstrapping Confidence intervals for handling variability of outcomes is disabled

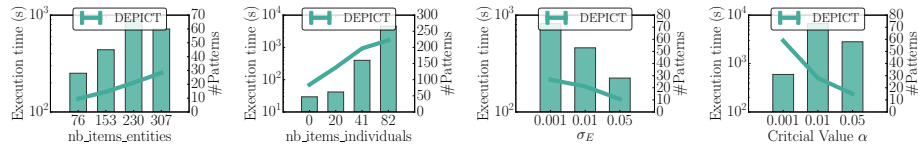


**Fig. 8:** Effectiveness of DEVIANT on Yelp when varying sizes of both search spaces  $\mathcal{D}_E$  and  $\mathcal{D}_I$ , minimum context support threshold  $\sigma_E$  and the critical value  $\alpha$ . Default parameters: full search spaces  $\mathcal{D}_E$  and  $\mathcal{D}_I$ ,  $\sigma_E = 0.1\%$ ,  $\sigma_I = 1\%$  and  $\alpha = 0.05$ . Bootstrapping Confidence intervals for handling variability of outcomes is enabled

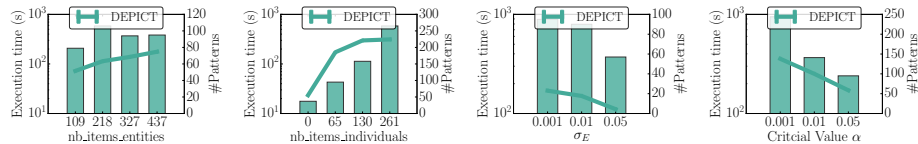
experiments we study also the overhead induced by the computation of the bootstrapping confidence interval required to handle the variability of outcomes and evaluated for each generated group of individuals.



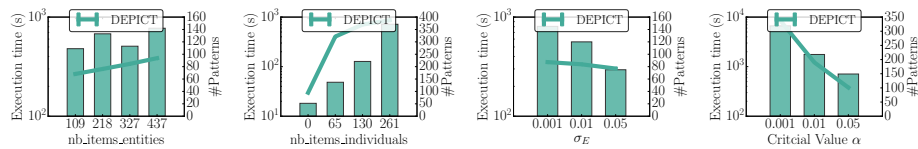
**Fig. 9:** Effectiveness of DEvIANT on CHUS when varying sizes of both search spaces  $\mathcal{D}_E$  and  $\mathcal{D}_I$ , minimum context support threshold  $\sigma_E$  and the critical value  $\alpha$ . Default parameters: full search spaces  $\mathcal{D}_E$  and  $\mathcal{D}_I$ ,  $\sigma_E = 0.1\%$ ,  $\sigma_I = 1\%$  and  $\alpha = 0.05$ . Bootstrapping Confidence intervals for handling variability of outcomes is disabled



**Fig. 10:** Effectiveness of DEvIANT on CHUS when varying sizes of both search spaces  $\mathcal{D}_E$  and  $\mathcal{D}_I$ , minimum context support threshold  $\sigma_E$  and the critical value  $\alpha$ . Default parameters: full search spaces  $\mathcal{D}_E$  and  $\mathcal{D}_I$ ,  $\sigma_E = 0.1\%$ ,  $\sigma_I = 1\%$  and  $\alpha = 0.05$ . Bootstrapping Confidence intervals for handling variability of outcomes is enabled



**Fig. 11:** Effectiveness of DEvIANT on EPD8 when varying sizes of both search spaces  $\mathcal{D}_E$  and  $\mathcal{D}_I$ , minimum context support threshold  $\sigma_E$  and the critical value  $\alpha$ . Default parameters: full search spaces  $\mathcal{D}_E$  and  $\mathcal{D}_I$ ,  $\sigma_E = 0.1\%$ ,  $\sigma_I = 1\%$  and  $\alpha = 0.05$ . Bootstrapping Confidence intervals for handling variability of outcomes is disabled.



**Fig. 12:** Effectiveness of DEvIANT on EPD8 when varying sizes of both search spaces  $\mathcal{D}_E$  and  $\mathcal{D}_I$ , minimum context support threshold  $\sigma_E$  and the critical value  $\alpha$ . Default parameters: full search spaces  $\mathcal{D}_E$  and  $\mathcal{D}_I$ ,  $\sigma_E = 0.1\%$ ,  $\sigma_I = 1\%$  and  $\alpha = 0.05$ . Bootstrapping Confidence intervals for handling variability of outcomes is enabled.

## C.2 Qualitative evaluation

In this appendix, we report additional illustrative examples depicting the significant patterns discovered by DEvIANT when carried on the Eighth European Parliament (EPD8) dataset, US House of representatives (CHUS) dataset and Yelp dataset.

id	group ( $g$ )	context ( $c$ )	$A^g(*)$	$A^g(c)$	$p$ -value	IA
$p_1$	S&D	8.10 Revision of the Treaties and intergovernmental conferences	0.81	0.44	< 0.001	Conflict
$p_2$	*	2 Internal market, single market 6 External relations of the Union	0.27	0.54	< 0.001	Consensus
$p_3$	S&D	8.30 Treaties in general	0.81	0.55	< 0.001	Conflict
$p_4$	*	2 Internal market, single market, 4.15 Employment policy, act. combat unemployment	0.27	0.53	< 0.001	Consensus
$p_5$	ALDE	1.20.09 Protection of privacy and data protection 8 State and evolution of the Union	0.73	0.48	< 0.001	Conflict

**Table 7:** Top-5 Exceptional consensual/conflictual subjects among European Political Groups in the 8<sup>th</sup> EU parliament.  $\alpha = 0.01$ . Patterns are ranked by the absolute difference between  $A^g(c)$  and  $A^g(*)$ .

id	group ( $g$ )	context ( $c$ )	$A^g(*)$	$A^g(c)$	$p$ -value	IA
$p_1$	*	10 Health & Medical	0.05	-0.33	< 0.001	Conflict
$p_2$	*	03 Automotive	0.05	-0.29	< 0.001	Conflict
$p_3$	*	14 Local Services	0.05	-0.25	< 0.001	Conflict
$p_4$	newcommer	21.177 American (Traditional) and 21.24 Breakfast & Brunch	-0.06	0.24	< 0.01	Consensus
$p_5$	*	11 Home Services	0.05	-0.25	< 0.01	Conflict

**Table 8:** Top-5 Exceptional consensual/conflictual subjects among Yelp users.  $\alpha = 0.01$ . Patterns are ranked by the absolute difference between  $A^g(c)$  and  $A^g(*)$ .