

Maverick: Discovering Exceptional Facts from Knowledge Graphs

Gensheng Zhang*
gezhang@google.com
Google, Inc.

Damian Jimenez
damian.jimenez@mavs.uta.edu
The University of Texas at Arlington

Chengkai Li
cli@uta.edu
The University of Texas at Arlington

ABSTRACT

We present Maverick, a general, extensible framework that discovers exceptional facts about entities in knowledge graphs. To the best of our knowledge, there was no previous study of the problem. We model an exceptional fact about an entity of interest as a context-subspace pair, in which a subspace is a set of attributes and a context is defined by a graph query pattern of which the entity is a match. The entity is exceptional among the entities in the context, with regard to the subspace. The search spaces of both patterns and subspaces are exponentially large. Maverick conducts beam search on the patterns which uses a match-based pattern construction method to evade the evaluation of invalid patterns. It applies two heuristics to select promising patterns to form the beam in each iteration. Maverick traverses and prunes the subspaces organized as a set enumeration tree by exploiting the upper bound properties of exceptionality scoring functions. Results of experiments and user studies using real-world datasets demonstrated substantial performance improvement of the proposed framework over the baselines as well as its effectiveness in discovering exceptional facts.

ACM Reference Format:

Gensheng Zhang, Damian Jimenez, and Chengkai Li. 2018. Maverick: Discovering Exceptional Facts from Knowledge Graphs. In *SIGMOD'18: 2018 International Conference on Management of Data, June 10-15, 2018, Houston, TX, USA*. ACM, New York, NY, USA, 16 pages. <https://doi.org/10.1145/3183713.3183730>

1 INTRODUCTION

Knowledge graphs such as DBpedia [4], Freebase [6], Wikidata [37], and YAGO [31] record properties of and relationships between real-world entities. These data are used in numerous applications, including search, recommendation, and business intelligence. This paper introduces Maverick, a framework that, given an entity in a knowledge graph, discovers *exceptional facts* about the entity. Informally, such exceptional facts separate the entity from many other entities. Consider several factual statements in published news articles:

- (1) "Denzel Washington followed Sidney Poitier as only the second black to win the Best Actor award." (abcnews.go.com)

*The bulk of the work was done while the author was at UT-Arlington.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

SIGMOD'18, June 10-15, 2018, Houston, TX, USA

© 2018 Association for Computing Machinery.

ACM ISBN 978-1-4503-4703-7/18/06...\$15.00

<https://doi.org/10.1145/3183713.3183730>

- (2) "This was Brazil's first own goal in World Cup history ..." (yahoo.com)
- (3) "Hillary Clinton becomes first female presidential nominee." (chicagotribune.com)

An exceptional fact consists of three components: an *entity of interest*, a *context*, and a set of *qualifying attributes*. In each exceptional fact, among all entities in the context, the entity of interest is one of the few or even the only one that bears a particular value combination on the qualifying attributes. For example, in the above statement 1, the entity of interest is Denzel Washington, the context is the Academy Award Best Actor winners, and the qualifying attribute is ethnicity.

Discovery of exceptional facts is useful to important applications such as computational journalism [11, 12], recommendation systems, and data cleaning. a) In *fact-finding* [21, 22, 32, 41, 44], journalists are interested in monitoring data and discovering attention-seizing factual statements such as the aforementioned examples. These facts help make news stories substantiated and interesting, and they may even become leads to news stories. b) In *fact-checking* [20, 42], for vetting the statements made by humans, fact-checkers at news organizations such as The Washington Post, CNN, and PolitiFact can compare the statements with automatically-discovered facts. For example, an algorithm may find that Hillary Clinton is the second female presidential nominee, which contradicts with the statement 3 above.¹ c) Exceptional facts can help promote friends, news, products, and search results in various recommendation systems. d) When the discovered facts are inconsistent with known truth or apparent common sense, it reveals incomplete data or data errors. Such insights aid knowledge base cleaning and completion. For example, the above statement 3 may be generated using an incomplete source that misses the nomination of Victoria Woodhull.

Given an entity in a knowledge graph, an integer k , and an exceptionality scoring function, the objective of *exceptional fact discovery* is to find the top- k highest scored pairs of (context, attribute set). The entity is exceptional with regard to the attributes, while at the same time belonging to the context together with other entities. This description hinges upon two concepts—context and attribute—which we explain below.

- The *attributes* of an entity are the entity's incoming/outgoing edge labels, and the attribute values are the entity's direct neighbors. For example, Fig. 1 is an excerpt of a knowledge graph about FIFA World Cup, in which the edge labeled *awarded-to* from node G_1 to CRO captures the fact that the goal is awarded to the team Croatia. Entity G_1 has two attributes *scored-by* and *awarded-to*, with values S_1 and CRO , respectively.

¹The first female presidential nominee was Victoria Woodhull, according to <http://www.snopes.com/victoria-woodhull-hillary-clinton/>.

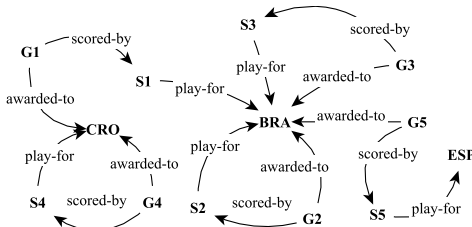


Figure 1: An excerpt of a knowledge graph.

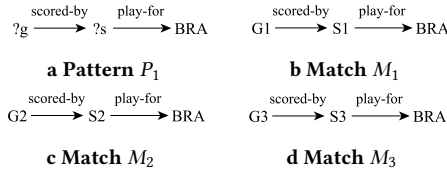


Figure 2: Pattern P_1 and variable $?g$ define a context consisting of all the goals scored by BRA players; M_1, M_2, M_3 are matches to P_1 in Fig. 1.

- A *context* is a set of entities sharing some common characteristics defined in a pattern query. In Fig. 2a, pattern P_1 and the variable $?g$ in it define a context C_1 of all the goals scored by players of team Brazil. Figs. 2b-2d show P_1 's matches in Fig. 1. For instance, match M_1 (Fig. 2b) is a subgraph of Fig. 1, in which $?g$ of P_1 is mapped to G_1 . Hence, G_1 belongs to context C_1 . Similarly, G_2 and G_3 in Fig. 1 also belong to C_1 based on M_2 and M_3 , while G_4 and G_5 are not part of C_1 .
- With respect to a *subspace* (i.e., a set of attributes), an entity is exceptional in a context if its attribute values deviate from the values of other entities in the same context. For example, the value of attribute *awarded-to* for G_1 is CRO, while the value is BRA for both G_2 and G_3 . The degree of exceptionality of an entity varies by different contexts and subspaces. For instance, one interpretation of statement 1 is that the context is the Academy Award Best Actor winners and the qualifying attribute is ethnicity; an alternative interpretation is that the context is all African Americans and the qualifying attribute is the award. Under some definitions of exceptionality, the second interpretation may render Denzel Washington more exceptional, since there are a lot more African Americans than winners of the award.

A holistic solution to exceptional fact discovery may be expected to synthesize whatever types of available data (structured databases, graphs, text, and so on), which is beyond the scope of this paper. Instead, our focus is on knowledge graphs which are becoming increasingly important to analytics and intelligence applications. To the best of our knowledge, there is no previous study on discovering exceptional facts about entities in knowledge graphs. The two most related areas are *outlier detection* in graphs [17, 28, 33, 38] and *outlying aspect mining* [2, 3, 15, 21, 32, 36, 40]. Duan et al. [15] and Vinh et al. [36] discussed the differences between these two areas. They achieve different goals. Outlier detection searches for all outlying objects among a set of objects. Outlying aspect mining, however, focuses on only one given object and returns the subspaces of attributes in which the object is relatively outlying, regardless of its true degree of outlyingness. In terms of objectives and problem modeling, the exceptional fact discovery problem formulated in this paper is closer to outlying aspect mining than outlier detection.

However, it focuses on graph data. In contrast, existing outlying aspect mining methods [3, 21, 32, 40] assume a single relational table. These methods take a tuple as input and returns two disjoint attribute sets. The first set of attributes define the context, i.e., the tuples having values identical to that of the input tuple on the attributes. On the second set of attributes, the input tuple has peculiar values compared to other tuples belonging to the context.

However, these methods for outlying aspect mining cannot be effectively applied to knowledge graphs, since they are specifically devised for single tables only. A seemingly plausible idea can be to represent a knowledge graph as a single table and then to apply the existing methods on the table. Consider the single-table model of RDF proposed in [8]. When adapting it for a knowledge graph, each tuple (row) is for an entity v and each attribute (column) corresponds to an edge label in the knowledge graph. The attribute is also associated with an edge direction—either incoming into or outgoing from v . The value at the junction of the row and the column is an entity or a set of entities adjacent to v via edges with the label and direction given by the column. Given this single-table representation of the knowledge graph, at least a few major problems render the existing outlying aspect mining methods inapplicable. *First*, in these methods a context, defined by a set of attributes, consists of the tuples having values identical to that of the input tuple. In other words, the context is the result of a conjunctive query over the attributes. For knowledge graphs, however, a context is defined by a graph pattern query, which cannot be captured by conjunctive queries on attributes in the aforementioned single-table representation. More specifically, an edge in the pattern may not be adjacent to the input entity and thus does not correspond to any of the entity's attributes. Hence, evaluating a pattern may involve self-joins of the single-table. Existing outlying aspect mining methods are not designed to accommodate joins. *Second*, the aforementioned set values in the single-table representation are not considered in the existing methods. An adaptation of the methods will thus require at least joins which, as mentioned above, are not supported by the methods. *Third*, due to the heterogeneity and scale of a large knowledge graph, such a single-table is extremely wide and sparse, which is well beyond the capacity of the existing methods because of the intrinsic exponential complexity of the problem's search space.

To discover the exceptional facts about an entity, we must explore two extremely large search spaces, one of patterns and the other of attribute subspaces. Section 5.1 shows that the number of patterns is at least exponential in the size of the graph. It is also clear that the number of subspaces is exponential in the number of attributes since a subspace is a combination of attributes. It is not computationally feasible to exhaustively enumerate all possible patterns and subspaces. Furthermore, it is challenging to prune patterns and subspaces, due to the non-existence of *downward closure property* (i.e., *anti-monotone property*) on typical exceptionality scoring functions.

To tackle these challenges, this paper introduces Maverick, a beam-search based framework. Given an input entity, Maverick discovers the top- k context-subspace pairs that give the entity the highest exceptionality scores. Maverick allows an application to plug in any exceptionality scoring function based on the application needs. Conceptually, Maverick organizes the search space of patterns as a partial order defined by the subsumption relation on

patterns and the search space of attribute subspaces as a set enumeration tree [30]. Intuitively, the search for top- k context-subspace pairs is performed in a nested-loop fashion in which the outer loop enumerates patterns and the inner loop enumerates subspaces. Maverick conducts breath-first beam search [43] on the space of patterns, starting from a pattern with a single variable node. On each visited pattern, Maverick applies a set of heuristics to prune its children so that Maverick visits at most w patterns at each level, where w is the beam width. Each visited pattern is evaluated over the knowledge graph to obtain the contexts it defines. For each context, Maverick calculates the input entity's exceptionality scores in different subspaces. It exploits an upper bound for exceptionality score to guide the traversal of the subspaces. The supersets of a subspace are pruned if their upper-bound scores are below the current top- k scores.

The paper reports the results of experiments on two real-world knowledge graphs, which verify Maverick's effectiveness in finding exceptional facts. The experiments compared the performance of a breath-first search method and the beam search method coupled with different candidate-selection heuristics. The experiment results establish that, even though the breath-first search method may evaluate more patterns in a fixed time frame than the beam search methods, it is not as effective as the beam search method using the proposed heuristics. We have also included some exceptional facts discovered by Maverick to demonstrate its practicality.

2 PROBLEM FORMULATION

In this section we formally define the data model of knowledge graphs, the concepts of context, attribute, and subspace, and the problem of exceptional fact discovery.

Knowledge Graphs

A knowledge graph $G(V_G, E_G)$ is a set of RDF [13] triples with node set $V_G \subseteq \mathcal{I}$ and edge set $E_G \subseteq V_G \times \mathcal{I} \times V_G$, where \mathcal{I} is the universe of IRIs.² In Fig. 1, there are three kinds of entities: goals (e.g., G_1), players (e.g., s_1), and teams (e.g., BRA). (Without loss of generality, we use an entity's name as its identifier (IRI) in the ensuing examples, assuming entity names are unique.) Three different types of edge labels represent different relationships: each player plays for a team (*play-for*), and each goal is scored by a player (*scored-by*) and is awarded to a team (*awarded-to*). For example, there is an own goal, as G_1 is scored by s_1 , a player of BRA , but awarded to CRO .

Patterns and Contexts

Definition 1 (Pattern P). A pattern is a weakly connected graph³ $P(V_P, E_P)$, where $V_P \subseteq \mathcal{I} \cup \mathcal{V}$, $E_P \subseteq V_P \times \mathcal{I} \times V_P$, and \mathcal{V} is the universe of variables. We also denote by $X_P \subseteq V_P$ the variables occurring in P . Δ

Definition 2 (Match M). A match $M(V_M, E_M)$ to a pattern $P(V_P, E_P)$ is a subgraph of G ($V_M \subseteq V_G$ and $E_M \subseteq E_G$) such that there exists a bijection $f : V_P \rightarrow V_M$ satisfying the following conditions:

- $|V_M| = |V_P|$, $|E_M| = |E_P|$;
- $\forall (v_i, l, v_j) \in E_P \Rightarrow (f(v_i), l, f(v_j)) \in E_M$;
- $\forall (u_i, l, u_j) \in E_M \Rightarrow (f^{-1}(u_i), l, f^{-1}(u_j)) \in E_P$;

²For the sake of simplicity, we do not consider blank nodes and literals.

³A weakly connected graph is a directed graph of which the corresponding undirected graph is connected.

- $\forall v \in \mathcal{I} \Rightarrow f(v) = v$.

In short, a subgraph M of G is a match to pattern P if M is edge-isomorphic to P and, for each non-variable node v in P , $f(v)$ has the same identifier. Δ

Note that the semantics of patterns in our definition is similar to that of *basic graph patterns* in [19, 27]. However, there are two main differences. One is that patterns in this work are weakly connected. The other is that a match to a pattern is required to be edge-isomorphic to the pattern. Neither of them is enforced in [19, 27].

Definition 3 (Range of Variable R_x^P). Let \mathcal{M}_P be all the matches to pattern P in a knowledge graph G . (I.e., \mathcal{M}_P is $[[P]]_G$, the evaluation of P against G , using the terminology in [27].) For a variable $x \in X_P$, the range of x , denoted R_x^P , is a set of entities defined as

$$R_x^P = \{f(x) \mid M \in \mathcal{M}_P, f : V_P \rightarrow V_M\}. \quad \Delta$$

For example, P_1 in Fig. 2a has two variable nodes, $?g$ and $?s$. (To distinguish variables from entities, the names of variable nodes always start with the symbol $?$.) Figs. 2b-2d show P_1 's matches in Fig. 1. $R_{?g}^{P_1} = \{G_1, G_2, G_3\}$ and $R_{?s}^{P_1} = \{s_1, s_2, s_3\}$.

Definition 4 (Context $C_v^{P,x}$). Given an entity v , a pattern P , a variable $x \in X_P$ such that $v \in R_x^P$, the context of v defined by P and x is denoted $C_v^{P,x}$ and $C_v^{P,x} = R_x^P$. Δ

For example, the context of G_1 in the running example—goals scored by BRA players—is defined by pattern P_1 in Fig. 2a and variable $?g$ in the pattern: $C_{G_1}^{P_1, ?g} = R_{?g}^{P_1} = \{G_1, G_2, G_3\}$. On the other hand, since $G_1 \notin R_{?s}^{P_1}$, $?s$ in P_1 does not define a context of G_1 . Note that a pattern may define multiple contexts of v , since v may be mapped to different variables in the pattern. For example, consider pattern $P = \{(?g, \text{awarded-to}, ?a), (?g, \text{scored-by}, ?s), (?s, \text{play-for}, ?t)\}$. It defines two different contexts of BRA : $C_{BRA}^{P, ?a} = \{CRO, BRA\}$, $C_{BRA}^{P, ?t} = \{ESP, BRA\}$.

Entity Attributes and Subspaces

Given an entity of interest v , an attribute corresponds to the label of an edge incoming into or outgoing from v , and its value is the entity at the other end of the edge. Note that we need to distinguish between incoming attributes and outgoing attributes since an entity can be both sources and destinations of edges of the same label. For instance, a person can have a manager and meanwhile be the manager of someone else.

Definition 5 (Entity Attributes A_v). Given an entity v , its attributes A_v is the union of its incoming and outgoing attributes: $A_v = A_v^i \cup A_v^o$. The incoming attributes are a set of edge labels $A_v^i = \{(l, \leftarrow) \mid \exists(x, l, v) \in E_G\}$. Given an incoming attribute $a = (l, \leftarrow) \in A_v^i$, v 's value on attribute a is the set $v.a = \{x \mid (x, l, v) \in E_G\}$. Similarly, the outgoing attributes are $A_v^o = \{(l, \rightarrow) \mid \exists(v, l, x) \in E_G\}$. Given an outgoing attribute $a = (l, \rightarrow) \in A_v^o$, v 's value is $v.a = \{x \mid (v, l, x) \in E_G\}$. Δ

Definition 6 (Subspace A). A subspace A is a subset of v 's attributes, i.e., $A \subseteq A_v$. The projection of v 's attribute values onto subspace A is denoted $v.A$, and $v.\emptyset = \text{null}$. Δ

For example, in Fig. 1, $A_{CRO}^i = \{(\text{play-for}, \leftarrow), (\text{awarded-to}, \leftarrow)\}$; $CRO.(\text{awarded-to}, \leftarrow) = \{G_1, G_4\}$ and $A_{G_1}^o = \{(\text{scored-by}, \rightarrow), (\text{awarded-to}, \rightarrow)\}$; $G_1.(\text{awarded-to}, \rightarrow) = \{CRO\}$. Let subspace $A = \{(\text{scored-by}, \leftarrow), (\text{play-for}, \rightarrow)\}$. $A_{s_1} = A$ and $s_1.A = \{G_1, \{BRA\}\}$.

Exceptionality Score

Definition 7 (Exceptionality Scoring Function χ). An exceptionality scoring function $\chi(v, A, C) \in \mathbb{R}$ measures entity v 's degree of exceptionality with regard to subspace A in comparison with other entities in context C . Without loss of generality, we assume the range of χ is $[0, 1]$, with larger χ implying greater exceptionality. We also set $\chi(v, A, C) = 0$ if $A \not\subseteq A_v$ or $v \notin C$, to make χ a total function. \triangle

The Maverick framework is indifferent to the choice of the exceptionality scoring function. It can accommodate many different interestingness/outlyingness functions (see surveys such as [18, 24]). Hence, the focus of this paper is not on the design, evaluation and comparison of such exceptionality scoring functions. Rather, the goal is to develop a general framework for efficiently finding exceptional facts under various scoring functions. Nevertheless, to make the discussion concrete, we consider several representative functions, of which one is introduced below and two more are discussed in Section 4.2 and in Appendix A.2. To ensure consistency, the discussion uses our own notations and terminologies in presenting the adaptation of existing functions.

One-of-the-Few χ_f The one-of-the-few concept is adapted from [41]. The crux of the idea is that a factual claim about an entity is interesting when equally or more significant claims can be made about only few other entities. For example, in Fig. 1, it is interesting to claim “ G_1 is the only own goal among the goals scored by BRA players”, since such a unique claim cannot be made about any other goal scored by a BRA player. On the contrary, “ G_1 is the only goal scored by S_1 ” is not impressive, because the same kind of claim “ G_x is the only goal scored by S_y ” can be made for all 5 goals in Fig. 1.

The one-of-the-few measure [41] is based on multi-criteria dominance relationship which is irrelevant to this work. Our adaptation of [41] quantifies the rareness of attribute values based on frequency. Consider a context C , a subspace A , and any entity u in the context. We denote by p_S^A , or simply p_S when A is clear, the probability (or “frequency” as in [3]) of u taking values S in subspace A , i.e.,

$$p_S^A = p(u.A = S \mid u \in C) = |\{u \mid u \in C, u.A = S\}| / |C|. \quad (1)$$

Ranking facts directly by frequency is not robust, regarding which detailed analysis can be found in [41]. To intuitively understand the insight, consider an extreme example. Suppose in an organization everyone has a unique name. Given an particular individual x , a fact “ x is the only person with that name” has high exceptionality measured by frequency itself. However, it is not truly exceptional since the same kind of fact can be stated for everyone.

Based on the definition of p_S^A , the one-of-the-few χ_f quantifies the exceptionality of an entity of interest v by the pessimistic rank of the frequency of $v.A$. Specifically, the exceptionality of v is:

$$\chi_f(v, A, C) = |\{u \mid u \in C, p_{u.A} > p_{v.A}\}| / |C|. \quad (2)$$

For example, consider entity of interest $v_0 = G_1$ in Fig. 1 and context C defined by pattern P_1 and variable g in Fig. 2a, i.e., $C = C_{G_1}^{P_1, g} = \{G_1, G_2, G_3\}$. Table 1 shows the frequencies of attribute values in all subspaces. According to Table 1, $p_{G_2.A} = p_{G_3.A} = p_{\{BRA\}} = \frac{2}{3} > p_{G_1.A} = p_{\{CRO\}} = \frac{1}{3}$. Hence, $\chi_f(G_1, A, C) = \frac{|\{G_2, G_3\}|}{|C|} = \frac{2}{3}$. For $A = \{(awarded-to, \rightarrow), (scored-by, \rightarrow)\}$, $\chi_f(G_1, A, C) = 0$, since there exists no $u \in C$ such that $p_{u.A} > p_{G_1.A}$.

Table 1: The frequencies of attribute values in all subspaces for entity of interest G_1 with regard to context $C = \{G_1, G_2, G_3\}$.

| A | $v.A : p_{v.A}^C$ | $G_1.A$ |
|---|--|----------------------------------|
| $\{(awarded-to, \rightarrow)\}$ | $\langle\{CRO\}\rangle:1/3, \langle\{BRA\}\rangle:2/3$ | $\langle\{CRO\}\rangle$ |
| $\{(scored-by, \rightarrow)\}$ | $\langle\{S_1\}\rangle:1/3, \langle\{S_2\}\rangle:1/3, \langle\{S_3\}\rangle:1/3$ | $\langle\{S_1\}\rangle$ |
| $\{(awarded-to, \rightarrow), (scored-by, \rightarrow)\}$ | $\langle\{CRO\}, \{S_1\}\rangle:1/3, \langle\{BRA\}, \{S_2\}\rangle:1/3, \langle\{CRO\}, \{S_1\}\rangle$ | $\langle\{CRO\}, \{S_1\}\rangle$ |
| | $\langle\{BRA\}, \{S_3\}\rangle:1/3$ | |

Definition 8 (Top- k Exceptional Facts F_v). With regard to an entity v , the *rank* of a context-subspace pair (C, A) is the number of context-subspace pairs with greater exceptionality scores, i.e., $rank(C, A) = |\{(C', A') \in C_v \times A_v \mid \chi(v, A', C') > \chi(v, A, C)\}|$. C_v is the universe of v 's contexts: $C_v = \{C_v^{P, x} \mid P \in \mathcal{P}, x \in P, v \in R_x^P\}$, in which \mathcal{P} is the universe of patterns over G , i.e., $\mathcal{P} = \{P(V_P, E_P) \mid V_P \subseteq X \cup V_G, E_P \subseteq (X \cup V_G) \times L \times (X \cup V_G), P(V_P, E_P) \text{ is weakly connected}\}$ where X is the universe of variables. (C, A) is a *top- k exceptional fact* if its rank is lower than k . Hence, the set of top- k exceptional facts about v , F_v , is defined as $F_v = \{(C, A) \in C_v \times A_v \mid rank(C, A) < k\}$.⁴ \triangle

Problem Statement Given a knowledge graph G , an entity of interest v_0 , an integer k , and an exceptionality scoring function χ , the problem of *exceptional fact discovery* is to find F_{v_0} —the top- k exceptional facts about v_0 .

Continue the running example. With regard to G_1 , the context-subspace pair $(C_{G_1}^{P_1, g}, \{(awarded-to, \rightarrow)\})$ may be exceptional. The context $C_{G_1}^{P_1, g}$ is $\{G_1, G_2, G_3\}$, i.e., the goals scored by BRA players. An interpretation of G_1 's exceptionality with regard to the pair is: among all the goals scored by BRA players, G_1 is the only own goal.

Alternative Problem Modeling

There could be other ways of defining context and subspace. Definition 4 allows contexts based on arbitrary patterns. It is possible to adopt a more simplified and restricted definition that only allows such patterns to be in certain “shapes” such as paths, star graphs, and trees. Definition 6 dictates that a subspace must be a set of entity attributes. In other words, when comparing an entity with other entities in a context, the entity stands out with respect to a subspace if it satisfies the conjunctive condition formed on the attributes in the subspace (i.e., a star query) while most other entities do not. It is plausible to adopt a more complex and expressive definition that allows the framework to assess exceptionality of entities using more complex, general graph queries instead of only star queries.

The current choices of Definitions 4 and 6 are formed based on several considerations related to usability and practicality. Particularly, the exceptionality scoring functions in this section and Section 4.2, adapted from functions in the literature that define outlyingness of tuples in relational tables, are defined on the aforementioned star queries. It is thus unclear how to define a scoring function using more complex graph queries. While such is an interesting question to ponder, it falls outside this paper's scope. As mentioned earlier, the Maverick framework is indifferent to the choice of the exceptionality scoring function. The current simple definition of subspace also eases the task of ensuring the discovered facts can be intuitively expressed by the system and interpreted by

⁴The size F_v may be greater than k due to ties in exceptionality scores and thus ranks.

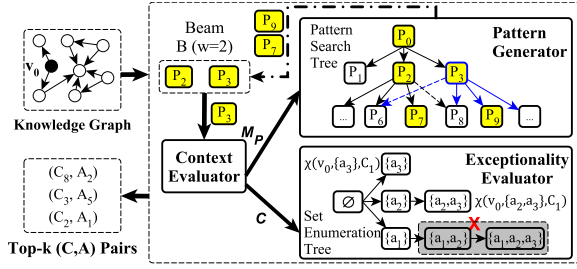


Figure 3: The framework of Maverick.

users. On a related note, while conducting the experiments (Section 6) we limited the sizes of the context-defining patterns and subspaces to be very small, only involving at most a handful of nodes and edges.

There could also be other ways of defining attributes. For example, one can define an entity's attributes as a vector of values either independent to or derived from the graph. For instance, [17, 28] consider, for each node, an associated mini-table containing information from external sources. A prevalent model of entities in knowledge graphs is embedding-based [7, 25, 39, 45], in which each entity is represented by a vector capturing its neighborhood information. Such vectors can also be used as entities' attributes. However, the vectors are indecipherable to human beings. Furthermore, in general knowledge graphs, an entity may have sub-properties, functional properties, and transitive properties [9, 14]. This work does not consider such models and thus is lack of reasoning capacity based on such properties. Some of such properties may be leveraged by pre-processing. For example, one may materialize the transitive properties. This can be an interesting future direction to explore.

3 OVERVIEW OF FRAMEWORK

We propose Maverick, an iterative framework for exceptional fact discovery. Intuitively, the process of discovering context-subspace pairs can be viewed as nested loops. The outer loop enumerates contexts, while the inner loop enumerates subspaces for each context. Given the entity of interest v_0 , while subspace enumeration in the inner loop enumerates the subsets of A_{v_0} , the outer loop enumerates contexts by patterns, since each context of an entity is defined by a pattern and one of the pattern's variables (c.f. Definition 4). Conceptually, Maverick organizes all the possible contexts as a partial order on patterns, i.e., a Hasse diagram, in which each node is a pattern and each edge represents the subsumption (subgraph-supergraph) relationship between the two patterns. The essence of the outer loop is thus a traversal of the search space of patterns.

Given that the search space of patterns can be extremely large (Section 5), it is impractical to adopt breath-first, depth-first, or heuristic search approaches due to memory and time constraints [29]. To address this challenge, we propose to traverse the search space by *beam search* [5]. Since beam search maintains a "beam" of heuristically w best nodes and prunes all other nodes, it is not guaranteed to be complete or optimal. However, good solutions can be found quickly if the heuristic is sound enough.

Fig. 3 and Alg. 1 illustrate the framework of Maverick, which has three main components: Context Evaluator (CE), Exceptionality Evaluator (EE), and Pattern Generator (PG). The beam search at

Algorithm 1: Discovering exceptional context-subspace pairs.

```

1 FACT-DISCOVER ( $G, v_0, \chi, k, w$ )
   Input:  $G$ : the knowledge graph;  $v_0 \in V_G$ : the entity of interest;
            $\chi$ : the exceptionality scoring function;  $k$ : the size of
           output;  $w$ : the beam width
   Output:  $H$ :  $k$  most exceptional context-subspace pairs
2  $P_0 \leftarrow (V_{P_0} = \{x_0\}, E_{P_0} = \emptyset)$ ; // Initial state.  $x_0$  is a variable.
3  $B \leftarrow \{P_0\}$ ; // Beam.
4  $i \leftarrow 1$ ; // Iteration number.
5 while  $B \neq \emptyset$  and  $i \leq \text{MAX\_ITERATION}$  do
6    $i \leftarrow i + 1$ ;  $B_{tmp} \leftarrow \emptyset$ ;
7   foreach  $P \in B$  do
8     // Obtain contexts of  $v_0$  and matches to  $P$ . (Section 3.1)
9      $C_{v_0}^P, M_P \leftarrow \text{CONTEXT-EVALUATOR}(P, v_0, G)$ ;
10    foreach  $C \in C_{v_0}^P$  do
11      // Exceptionality Evaluation. (Section 4)
12       $\mathcal{A} \leftarrow \text{EXCEPTIONALITY-EVALUATOR}(v_0, C, k, \chi)$ ;
13      foreach  $A \in \mathcal{A}$  do  $H \leftarrow H \cup \{(C, A)\}$ ;
14      // Find  $\mathcal{Y}$  — the children of  $P$ . (Section 5)
15       $\mathcal{Y} \leftarrow \text{PATTERN-GENERATOR}(v_0, P, M_P, w, G)$ ;
16       $B_{tmp} \leftarrow B_{tmp} \cup \mathcal{Y}$ ;
17    $B \leftarrow \text{top-}w \text{ of } B_{tmp} \text{ based on heuristics } h$ ; // Section 5.4
18 return top- $k$  pairs in  $H$  based on exceptionality scores;

```

the outer loop starts with a pattern P_0 with a single variable node x_0 (Lines 2–3 in Alg. 1). The search results in a *pattern search tree*, of which the root is P_0 . At each iteration, Maverick maintains a beam B of a fixed size w (Lines 6, 13, 14). The beam consists of heuristically the best w patterns (e.g., P_2, P_3 in Fig. 3 where $w = 2$) at the visited level of the pattern search tree. For each pattern P in B , component CE obtains the matches M_P to the pattern and the corresponding contexts $C_{v_0}^P$ of v_0 (Line 8). For each context C in $C_{v_0}^P$ (e.g. C_1 in Fig. 3), component EE finds the top- k scored subspaces according to a given exceptionality scoring function χ (Line 10, and Section 4). Component PG finds the children of the visited pattern based on its matches (Line 12, and Section 5). Since there are usually much more children than what the beam size w allows, PG applies a set of heuristics (Section 5.4) to prune the child patterns. Each child pattern is given a score that measures how promising it is according to the heuristics. The best w patterns among all the children of patterns in B will become the new beam B (Line 14), which is the input to the next iteration, e.g., $\{P_7, P_9\}$ in Fig. 3. The process ends when the limit on the number of iterations has reached. The limit is set to avoid overly-complex patterns which correspond to facts that are only convolutedly interesting. It also practically bounds the resource spent for running the algorithm. When the algorithm terminates, Maverick returns the k context-subspace pairs with the highest exceptionality scores (Line 15). Below, we discuss component CE in Section 3.1, EE in Section 4, and PG in Section 5.

3.1 Context Evaluator

The context evaluator (CE, Line 8 in Alg. 1) is responsible for obtaining the matches to a given pattern as well as the corresponding contexts. Its working is depicted in Alg. 2. We expect a graph query system to take a pattern as the input and return all the matches to the pattern (Line 3). The Maverick framework is agnostic to the

choice of the specific query processing system. According to Definition 4, for each variable in the pattern ($x \in X_P$), CE returns its range R_x^P as a context if the entity of interest v_0 is in the range (Line 5).

For example, consider graph G in Fig. 1, the entity of interest $v_0 = g_1$, and the pattern P_1 in Fig. 2a. $M_{P_1} = \{M_1, M_2, M_3\}$, where M_1, M_2 , and M_3 are in Figs. 2b–2d. P_1 has two variables, $?g$ and $?s$. Since $g_1 \in R_{?g}^{P_1} = \{g_1, g_2, g_3\}$ and $g_1 \notin R_{?s}^{P_1}$, P_1 defines one and only one context of g_1 , which is $C_{g_1}^{P_1, ?g} = R_{?g}^{P_1}$. Therefore, $C_{g_1}^{P_1} = \{C_{g_1}^{P_1, ?g}\}$.

Algorithm 2: Context evaluator.

```

1 CONTEXT-EVALUATOR ( $P, v_0, G$ )
2    $C_{v_0}^P \leftarrow \emptyset$ ; // The set of contexts defined by  $P$ .
3    $M_P \leftarrow \text{match}(G, P)$ ; // Matches to  $P$ .
4   foreach  $x \in X_P$  do
5     // Refer to Definitions 3 and 4 for  $R_x^P$  and  $C_{v_0}^{P,x}$ .
6     if  $v_0 \in R_x^P$  then  $C_{v_0}^P \leftarrow C_{v_0}^P \cup \{C_{v_0}^{P,x}\}$ ;
7   return ( $C_{v_0}^P, M_P$ );

```

4 EXCEPTIONALITY EVALUATOR

The Exceptionality Evaluator (EE) operates in the inner loop of the Maverick framework (function `EXCEPTIONALITY-EVALUATOR` (v_0, C, k, χ) at Line 10 of Alg. 1). For each context C of the entity of interest v , it finds the k subspaces A with the highest $\chi(v, A, C)$ scores. Note that it is sufficient to find these k subspaces, since the eventual output of Maverick is the top- k context-subspace pairs across all contexts of v . A naive solution of EE can exhaustively enumerate all possible subspaces of A_v and calculate the exceptionality score of v in each subspace. The apparent $O(2^{|A_v|})$ complexity of this approach renders it prohibitively expensive since many entities may have a lot of attributes. For instance, Denzel Washington has more than 40 attributes in the August 9, 2015 Freebase graph. It is thus crucial for Maverick to have an efficient subspace enumeration method in order to discover more exceptional context-subspace pairs. Section 4.1 discusses how Maverick uses a set enumeration tree to avoid exhaustively enumerating subspaces. Specifically, Maverick exploits the upper bound properties of exceptionality scoring functions to guide the traversal of the set enumeration tree. Section 4.2 introduces three representative exceptionality scoring functions along with their upper bound functions.

4.1 Finding Top- k Subspaces

EE applies a *set enumeration tree* (SE-tree) [30] to avoid exhaustively enumerating subspaces. Each node in the tree is a subspace—a subset of v 's attributes A_v . The children of a node correspond to various supersets of the node's associated attributes. Formally, let r be an (arbitrary) total order on A_v . The root of an SE-tree for A_v is the empty set. The children of a node $A \subset A_v$ in the tree form the set $\{A \cup \{a\} \mid a \in A_v \setminus A, \forall a' \in A, a' <_r a\}$. An SE-tree for $A_v = \{a_1, a_2, a_3\}$ is shown in Fig. 3. The gist is to explore the set enumeration tree using heuristic search methods such as best-first search and to prune branches that are guaranteed to not contain highly-scored subspaces.

What is particularly challenging is that an exceptionality scoring function χ usually does not have the *downward closure property* with respect to subspace inclusion, i.e., $\chi(v, A, C)$ can be greater than,

Algorithm 3: Exceptionality evaluator.

```

1 EXCEPTIONALITY-EVALUATOR ( $v, C, k, \chi$ )
2   // CS: current subspace; UA: attributes to visit; Tk: top- $k$  subspaces.
3    $CS \leftarrow \emptyset$ ;  $UA \leftarrow A_v$ ;  $Tk \leftarrow \emptyset$ ;
4   return EXPLORE-SUBSPACE( $v, C, k, \chi, CS, UA, Tk$ );
5 EXPLORE-SUBSPACE ( $v, C, k, \chi, CS, UA, Tk$ )
6   while  $UA \neq \emptyset$  do
7     // Calculate upper bounds.
8      $a_{\max} \leftarrow \arg \max_{a \in UA} \text{upper}(v, CS \cup \{a\}, C)$ ;
9      $A_{\max} \leftarrow CS \cup \{a_{\max}\}$ ;  $\text{upper}_{\max} \leftarrow \text{upper}(v, A_{\max}, C)$ ;
10     $UA \leftarrow UA \setminus \{a_{\max}\}$ ;
11    if  $|Tk| < k$  then
12      // -1 indicates the top- $k$  list  $Tk$  is not full.
13       $\text{score}_{\min} \leftarrow -1$ ;  $A_{\min} \leftarrow \emptyset$ ;
14    else ( $A_{\min}, \text{score}_{\min}$ )  $\leftarrow \arg \min_{(A, \text{score}) \in Tk} \text{score}$ ;
15    if  $\text{upper}_{\max} > \text{score}_{\min}$  then
16       $\text{score} \leftarrow \chi(v, A_{\max}, C)$ ;
17      if  $\text{score} > \text{score}_{\min}$  then
18        if  $\text{score}_{\min} \geq 0$  then
19           $Tk \leftarrow Tk \setminus \{(A_{\min}, \text{score}_{\min})\}$ ;
20           $Tk \leftarrow Tk \cup \{(A_{\max}, \text{score})\}$ ;
21        // Explore children subspaces.
22         $Tk \leftarrow \text{EXPLORE-SUBSPACE}(v, C, k, \chi, A_{\max}, UA, Tk)$ ;
23    return  $Tk$ ;

```

less than, or equal to $\chi(v, A', C)$ for any $A' \supseteq A$. As a matter of fact, none of the three representative functions that will be introduced in Section 4.2 satisfies the property (proof omitted). The lack of downward closure property makes it infeasible to prune the set enumeration tree based on exact exceptionality scores.

EE uses upper bounds on the exceptionality scoring function χ to allow for pruning of the set enumeration tree. Alg. 3 presents its pseudo code. The set enumeration tree nodes (i.e., subspaces) are visited in the descending order of their upper bounds (Line 6). If the upper bound score of a node is not greater than the score of the current k -th ranked subspace, the node and all its children are pruned (Line 12). Otherwise, the exact exceptionality score of the node is calculated (Line 13). The subspace is used to purge the current k -th subspace if its exact score is still greater (Lines 14–17). Regardless of whether the node makes into the top- k list, its children are enumerated recursively (Line 18).

The general upper bound function *upper* in Alg. 3 is defined as follows. By the definition, it is sound to prune a node and all its children if the condition in Line 12 is not satisfied.

Definition 9 (Upper bound of an exceptionality scoring function *upper*). Given an exceptionality scoring function χ , an upper bound of χ is a function that, for any entity v , context C , and subspace $A \subseteq A_v$, bounds the exceptionality score of v with respect to C and any superset of A , i.e.,

$$\text{upper}(v, A, C) \geq \max_{A \subseteq A' \subseteq A_v} \chi(v, A', C). \quad \Delta$$

The general upper bound function *upper* must be instantiated for specific exceptionality scoring functions χ . The Maverick framework expects an application developer to supply *upper* while specifying χ . Various outlying aspect mining methods [2, 3, 15] also devise upper bound functions for pruning set enumeration tree. They operate on the single-table data model and are thus inapplicable for graphs, as explained in Section 1. EE must use different

scoring functions and upper bound functions designed for knowledge graphs. The ensuing discussion in this section entails that.

4.2 Exceptionality Scoring Functions

As mentioned in Section 2, the general Maverick framework accommodates different exceptionality scoring functions beyond the one-of-the-few function χ_f . We discuss two more representative functions in this section and in Appendix A.2.

Outlyingness χ_o This measure, adopted from [3], is based on the distribution of attribute values. An entity receives a high score if it has rare attribute values while a lot of other entities share common attribute values. It quantifies the rareness of attribute values by $p_S^A = p(u.A = S \mid u \in C)$ (the same as for χ_f). Let \mathcal{S}_A be all possible attribute values on subspace A and in context C , i.e., $\mathcal{S}_A = \{u.A \mid u \in C\}$. The outlyingness score of an entity v is given by:

$$\chi_o(v, A, C) = \sum_{S \in \mathcal{S}_A} p_S \times (p_S - p_{v.A}) \times \mathbb{1}(p_S > p_{v.A}) \quad (3)$$

where $\mathbb{1}(\cdot)$ is the indicator function that returns 1 for a true condition and 0 otherwise. Essentially, the outlyingness score is the area above the accumulated frequency histogram of the context C with respect to the subspace A , starting from the frequency of $v.A$. The score is designed to quantify the “degree of unbalance” between the frequencies of entities in the context [3].

For instance, consider the same example used in explaining χ_f : $v_0 = G_1$, $C = C_{G_1}^{P_1, ?g} = \{G_1, G_2, G_3\}$, and $A = \{(awarded-to, \rightarrow)\}$. According to Table 1, $\chi_o(G_1, \{(awarded-to, \rightarrow)\}, C) = p_{\{CRO\}} \times (p_{\{CRO\}} - p_{\{CRO\}}) \times 0 + p_{\{BRA\}} \times (p_{\{BRA\}} - p_{\{CRO\}}) \times 1 = \frac{2}{3}(\frac{2}{3} - \frac{1}{3}) = \frac{2}{9}$. Another example is, for $A = \{(awarded-to, \rightarrow), (scored-by, \rightarrow)\}$, $\chi_o(G_1, A, C) = 0$ since there exists no $u \in C$ such that $p_{u.A} > p_{G_1.A}$.

4.3 Upper Bound Functions

In this section we devise upper bound functions for the three representative exceptionality functions introduced in Section 2 (χ_f) and Section 4.2 (χ_o and χ_i). We prove that these designs satisfy Definition 9 and thus ensure the soundness of Alg. 3, with regard to any given entity v , context C , and subspaces $A \subseteq A' \subseteq A_v$. Recall that we denote by p_S^A , or simply p_S , the frequency of entity’s attribute value S in subspace A (Eq. (1)).

Theorem 1 (Upper bound of χ_f). $upper_f(v, A, C) \geq \chi_f(v, A', C)$, given the following definition in which $\overline{C_v} = C \setminus \{v\}$:

$$upper_f(v, A, C) = |\{u \mid u \in \overline{C_v}, p_{u.A} > 1/|C|\}| / |C| \quad (4)$$

The theorem holds because $\frac{1}{|C|} \leq p_{u.A'} \leq p_{u.A}$ for any $A' \supseteq A$. We omit the detailed proof here.

Theorem 2 (Upper bound of χ_o). $upper_o(v, A, C) \geq \chi_o(v, A', C)$, given the following definition where $\mathcal{S}_A = \{u.A \mid u \in C\}$:

$$upper_o(v, A, C) = \sum_{S \in \mathcal{S}_A} (p_S)^2 - \frac{(2p_{v.A} + 1) \times |C| - 2}{|C|^2}. \quad (5)$$

Our final note is that an upper bound function may have limited pruning power when it gives loose bounds on exceptionality scores, resulting in exponential complexity in subspace enumeration. Our empirical results in Section 6.2, though, verified that the several upper bound functions proposed in this paper (Eqs. (5)–(7)) substantially reduced the overhead of subspace enumeration.

5 PATTERN GENERATOR

The Pattern Generator (PG) is used in Line 12 of Alg. 1 in the Maverick framework. Its pseudo code is in Alg. 4. At each iteration of the beam search on patterns, it finds the children of each visited pattern P (Line 3, see Alg. 5) in the current beam. A child pattern, if not pruned (see Section 5.3), is given a score that measures how promising it is according to a few heuristics (Line 5, see Section 5.4). Among all the children of the patterns in the current beam, the w children with the highest scores are returned to form the new beam (Line 14 in Alg. 1), where w is the predefined beam width. The new beam becomes the input to the next iteration. This section first describes the search space of patterns (Section 5.1) and then discusses how to efficiently explore the space by applying pruning rules (Section 5.3) and selection heuristics (Section 5.4).

Algorithm 4: Pattern generator.

```

1 PATTERN-GENERATOR( $v_0, P, M_P, w, G$ )
2    $\mathcal{Y} \leftarrow \emptyset$ ; // Promising children of  $P$ .
3   // Find  $P$ ’s children, see Alg. 5.
4   children  $\leftarrow$  FIND-CHILDREN( $v_0, P, M_P, G$ );
5   foreach child  $\in$  children do
6      $\mathcal{Y} \leftarrow \mathcal{Y} \cup \{(child, h(v_0, child))\}$ ; // See Section 5.4 for  $h$ .
7   return top- $w$  of  $\mathcal{Y}$  based on score;
```

5.1 Search Space of Patterns

The search space of patterns is a Hasse diagram of valid patterns, where a pattern is *valid* if it contains at least one variable node and it has a match (Definition 2) in the knowledge graph G . We exclude invalid patterns since they cannot lead to relevant facts. For example, pattern $\{(?g, scored-by, ?s_1), (?g, scored-by, ?s_2)\}$ does not have a match and is thus invalid because no goal is scored by more than one player. Formally, the search space of patterns is a Hasse diagram $\mathbb{P}(V_P, E_P)$, where V_P is the set of valid patterns and $E_P \subseteq V_P \times V_P$ is the set of edges. There exists an edge from *parent* pattern P_i to *child* pattern P_j if P_i is an immediate subgraph of P_j , i.e., P_i has exactly one edge less than P_j . A pattern can have multiple children and multiple parents. Fig. 4 shows an excerpt of the search space of patterns over the data graph in Fig. 1. In the figure, P_6 and P_7 are the children of P_2 , and both P_2 and P_3 are the parents of P_7 .

One may realize already that \mathbb{P} can be extremely large. We prove in Theorem 3 that the *order* of \mathbb{P} (i.e., the cardinality of V_P) is exponential to the orders of G ’s weakly connected components (WCCs). A WCC is a maximal subgraph of a directed graph, in which every pair of vertices are connected, ignoring edge direction. Given that knowledge graphs are all well connected, it is impossible to exhaustively enumerate the patterns. For example, according to Theorem 3, the data graph in Fig. 1 has at least $2^{13+1} - 2 - 13 + 13 =$

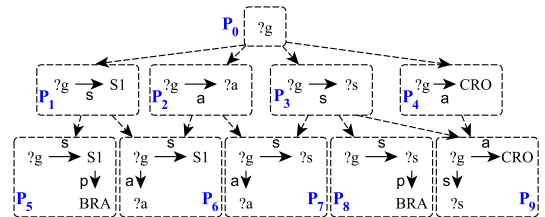


Figure 4: An excerpt of the search space of patterns over Fig. 1. Edge labels: *a*: awarded-to, *p*: play-for, *s*: scored-by.

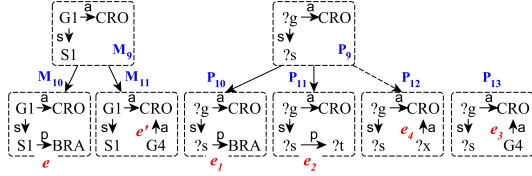


Figure 5: Illustration of how the child patterns of a pattern are constructed. P_{10} and P_{11} are obtained based on M_{10} and edge e . P_{12} and P_{13} can be obtained based on M_{11} and edge e' , but they are pruned based on rules in Section 5.3.

16,382 patterns. (The graph itself is the only WCC, with 13 nodes.) Note that Theorem 3 only provides a loose bound. In practice, the number is even much larger, exacerbating the challenge. Section 6.2 shows that the tiny graph has more than 69,000 patterns with merely no more than 5 edges.

Theorem 3. Let \mathcal{W} be the set of WCCs in a knowledge graph G , a lower bound on \mathbb{P} 's order is:

$$|\mathbb{P}| \geq \sum_{W \in \mathcal{W}} (2^{|V_W|+1} - 2) - |V_G| + \max_{W \in \mathcal{W}} |V_W|.$$

5.2 Match-based Construction of Patterns

Given the current beam of patterns, Maverick finds top context-subspace pairs using its context evaluator (Section 3.1) and exceptionality evaluator (Section 4). Among the child patterns of the evaluated patterns, the promising ones are chosen to form the new beam for the next iteration. While Section 5.4 discusses how to select the promising patterns, this section proposes an efficient way of generating the child patterns. Note that the aforementioned Hasse diagram of patterns is not pre-materialized. Rather, the patterns need to be constructed before we can evaluate them.

To construct the child patterns of an evaluated pattern P , a simple approach is to enumerate all possible ways of expanding P by adding one more edge. A major drawback of this approach is it may construct many invalid patterns that do not have any match. Some invalid patterns can be easily recognized by referring to the schema graph of the data. However, chances are most of the schema-abiding patterns are still invalid because they do not have matching instances in the data graph, given the sheer diversity of a knowledge graph. The system will evaluate such patterns in vain to get empty results in order to realize they are invalid.

To avoid evaluating invalid patterns, we propose a *match-based pattern construction method*. Instead of constructing the child patterns by directly expanding P , this method expands the matches of P and constructs the child patterns from the expanded matches. It guarantees to construct only valid patterns and evade the evaluation of invalid patterns. The method is based on the following theorem.

Theorem 4. Suppose P' is a child of $P \in \mathbb{P}$, i.e., $(P, P') \in E_{\mathbb{P}}$ and thus P' is a valid pattern with matches. Given any match M' to P' , there exists a match M to P that is a subgraph of M' , i.e., $\forall M' \in \mathcal{M}_{P'}, \exists M \in \mathcal{M}_P$ s.t. $V_M \subseteq V_{M'}$ and $E_M \subseteq E_{M'}$.

Based on Theorem 4, the method that constructs the child patterns of P is illustrated in Alg. 5. For a match M of P , it finds each of its weakly connected supergraphs by adding an edge that exists in the data graph G and is adjacent to a node in M (Line 6). Given each such resulting supergraph M' , let $(u, l, w) = E_{M'} \setminus E_M$ and, without loss of generality, assume $u \in V_M$. If $w \in V_M$, then the only child of P obtained from M' is $P + (f^{-1}(u), l, f^{-1}(w))$ (Line 13), where we

denote by $P + e$ the supergraph of P by adding edge e for brevity. If $w \notin V_M$, then two child patterns are obtained: $P + (f^{-1}(u), l, w)$ and $P + (f^{-1}(u), l, z)$, where z is a variable and $z \notin X_P$ (Line 16; Line 19 for the symmetric case). Fig. 5 shows an example of obtaining a pattern's children. For instance, P_{10} can be obtained by adding e_1 , which is obtained by replacing s_1 of edge e with variable $?s$.

Algorithm 5: Find all the children of a given pattern.

```

1 FIND-CHILDREN ( $v_0, P, \mathcal{M}_P, G$ )
2    $D \leftarrow \emptyset$ ; // The set of  $P$ 's children.
3    $\mathbb{M} \leftarrow \{M \in \mathcal{M}_P \mid f: V_P \rightarrow V_M \text{ and } \exists x \in X_P \text{ s.t. } f(x) = v_0\}$ ; // Rule 1
4   foreach  $M \in \mathbb{M}$  do
5     Let  $f$  be the bijection  $f: V_P \rightarrow V_M$ ;
6      $\bar{E}_M = \{(u, l, w) \in E_G \setminus E_M \mid u \in V_M \text{ or } w \in V_M\}$ ;
7     foreach  $(u, l, w) \in \bar{E}_M$  do
8        $z \leftarrow$  a new variable and  $z \notin X_P$ ;
9        $x \leftarrow f^{-1}(u), y \leftarrow f^{-1}(w)$ ;
10      if  $\nexists x \in X_P$  s.t.  $f(x) = u$  or  $f(x) = w$  then
11        continue; // Rule 2
12      else if  $u \in V_M$  and  $w \in V_M$  then
13         $P_1 \leftarrow P + (x, l, y)$ ;
14         $D \leftarrow D \cup \{P_1\}$ ;
15      else if  $w \notin V_M$  then //  $\exists x \in X_P$  s.t.  $f(x) = u$ 
16         $P_1 \leftarrow P + (x, l, w); P_2 \leftarrow P + (x, l, z)$ ;
17         $D \leftarrow D \cup \{P_1, P_2\}$ ;
18      else //  $\exists y \in X_P$  s.t.  $f(y) = w$ 
19         $P_1 \leftarrow P + (u, l, y); P_2 \leftarrow P + (z, l, y)$ ;
20         $D \leftarrow D \cup \{P_1, P_2\}$ ;
21   return  $D$ ;
```

5.3 Pattern Pruning Strategies

The search space \mathbb{P} of patterns as defined in Section 5.1 and constructed using the match-based pattern construction method in Section 5.2 has an enormous size. To ensure efficiency, the Pattern Generator (PG) employs two pruning rules to exclude irrelevant patterns from \mathbb{P} and to avoid repeated constructions of patterns from certain type of parent patterns.

Rule 1 (RelevantOnly). Exclude a pattern if it does not define any context for the entity of interest v_0 .

The rationale behind Rule 1 is, for discovering exceptional facts about v_0 , a pattern is relevant only if it defines a context for v_0 . By this rule, the match-based pattern construction method only expands a match in which v_0 is an image of a variable in P . It is guaranteed that the patterns obtained define v_0 's contexts.

Rule 2 (VarOnly). Expand a pattern only if the new edge has at least one variable.

Let P' be a child pattern of P . The extra edge in P' , i.e., $e = E_{P'} \setminus E_P$, belongs to one of the 7 types in Fig. 6. Rule 2 avoids constructing P' from P if e belongs to types 6-7. This rule is based on Theorem 5. Simply put, enforcing Rule 2 will not miss any contexts of v_0 .

Theorem 5. Let P' be a child of $P \in \mathbb{P}$, $e = E_{P'} \setminus E_P$, $C_{v_0}^P$ be all the contexts of v_0 defined by P : $C_{v_0}^P = \{R_{x'}^P \mid x' \in X_P, v_0 \in R_{x'}^P\}$, then $C_{v_0}^{P'} = C_{v_0}^P$, if e belongs to types 6-7.

5.4 Pattern Selection Heuristics (h)

Even with the rules proposed in Section 5.3, there are still too many patterns. In this section, we propose two scoring heuristics for selecting promising patterns to visit, to substantiate the function

| | | | |
|-----------------|-----------------|-----------------|-----------------|
| (1) (x, l, y) | (2) (x, l, z) | (3) (x, l, w) | (4) (u, l, y) |
| (5) (u, l, z) | (6) (u, l, v) | (7) (u, l, w) | |

Figure 6: Consider a pattern P and its child pattern P' . The 7 types of the extra edge $e = E_{P'} \setminus E_P$. x, y, z are variables, $x, y \in X_P, z \notin X_P$. u, v, w are non-variables, $u, v \in V_P \cap I$, and $w \in V_G \setminus V_P$.

h in Line 5 of Alg. 4. A heuristic gives each pattern a score, based on which the w patterns with the highest scores form the beam for the next iteration of beam search.

Heuristic 1 (Optimistic). Given a pattern P , the entity of interest v_0 , let $C_{v_0}^P$ be the set of contexts defined by P , i.e., $C_{v_0}^P = \{C_{v_0}^{P,x} | x \in X_P, v_0 \in R_x^P\}$, then

$$h_{opt}(v_0, P) = \max_{C \in C_{v_0}^P} upper(v_0, \emptyset, C)$$

where $upper(v_0, \emptyset, C)$ is an upper bound of χ with regard to C for any subspace (see Definition 9).

h_{opt} simply uses the exceptionality score upper bound of P . It optimistically assumes the ideal case for each pattern, where the entity of interest is most exceptional among the entities in a context defined by the pattern. In Section 4, we discussed the upper bound functions for various exceptionality functions. Note that we have $p_{v_0, \emptyset} = 1$ (Eq. (1)) since $v_0 \cdot \emptyset = \text{null}$ (Definition 6) and we consider all null values equal, $upper_o(v_0, \emptyset, C) = 1 - \frac{3 \times |C| - 2}{|C|^2}$, $upper_f(v_0, \emptyset, C) = 1 - \frac{1}{|C|}$, and $upper_i(v_0, \emptyset, C) = 1 - 2 - |C| \log_2 |C| / (|C| \log_2 |C| - |C - 1| - 1) \log_2 (|C - 1|)$ (cf. Appendix A.2). In sum, all the three upper bounds increase when the context size increases. In other words, h_{opt} selects the patterns that define large contexts. However, a large context may contain many entities of different characteristics, which may make the entity of interest less exceptional. Note that, since h_{opt} depends on context size $|C|$, all the child patterns of P need to be evaluated in order to get $|C|$. It is also required for heuristic h_{conv} below for the same reason.

Heuristic 2 (Convergent). Consider a pattern P and the entity of interest v_0 . Given P' , a parent of P in the pattern search tree, we define $r_x = |C_{v_0}^{P,x}| / |C_{v_0}^{P',x}|$. The score of P is

$$h_{conv}(v_0, P) = \max_{(P', P) \in E_{\mathbb{P}} \text{ and } P' \in B, C_{v_0}^{P',x} \in C_{v_0}^{P',x}} \left[r_x \times \max_{A \subseteq A_{v_0}} \chi(v_0, A, C_{v_0}^{P',x}) + (1 - r_x) \times upper(v_0, \emptyset, C_{v_0}^{P,x}) \right]$$

The h_{conv} score of P is a weighted sum of the upper bound of P (for any subspace) and the best score of the parent pattern P' . Note that Maverick performs a beam search and the patterns visited form a pattern search tree. P could be constructed from different parent patterns in the current beam B . The above equation thus uses the best score across all such parents. For this reason, the edge adjacent to P in the pattern search tree comes from the parent P' that gives it the best score. If P' possesses some highly-scored context-subspace pairs, h_{conv} gives favorable score to P if P and P' define similar contexts; otherwise, h_{conv} favors a P that defines smaller contexts. Compared with h_{opt} , h_{conv} is potentially both more efficient and more effective. It can be more efficient since it may favor child patterns that define smaller contexts. Such child patterns usually can be evaluated more efficiently since they have less matches. It can be more effective since it discards child patterns that define contexts where the entity of interest may not be exceptional, based on the highest score of the context-subspace pairs for the parent pattern. When h_{conv} is used for choosing patterns to form the

beams, the sizes of the contexts defined by the patterns in a path of the tree may gradually become smaller and eventually converge. We thus call h_{conv} Convergent.

6 EXPERIMENTS

6.1 Experiment Setup

The framework and algorithms of Maverick are implemented in Python. The experiments were conducted on a 16-core, 32GB-RAM node in Stampede—a cluster of the Extreme Science and Engineering Discovery Environment (XSEDE: <https://www.xsede.org>). All datasets used in experiments are available in Neo4j format at <https://doi.org/10.5281/zenodo.1185476>.

Datasets The experiments used the following two real-world graphs:

- **WCGoals**. It was constructed by crawling data from the FIFA World Cup website (<http://www.fifa.com/worldcup/index.html>). It consists of 49,078 nodes, 158,114 edges, 13 different edge labels, and 11 entity types: **WORLD_CUP**, **ROUND_CATEGORY**, **ROUND**, **STADIUM**, **TEAM**, **GAME**, **GROUP**, **PLAYER**, **BIBNUM**, **PARTICIPANT**, and **GOAL**.
- **OscarWinners**. This is a subgraph of Freebase. It has 42,148 nodes, 63,187 edges, 24 distinct edge labels, and 13 entity types including **PERSON**, **FILM_CREW**, **AWARD_WON**, **FILM_CHARACTER**, **AWARD_CATEGORY**, **PERFORMANCE**, **GENRE**, **AWARD**, **FILM**, **COUNTRY**, **FILM_CREW_ROLE**, **CEREMONY**, and **SPECIAL_PERFORMANCE_TYPE**. Each film in the graph has won at least one Academy Award (Oscar).

The two graphs were stored using Neo4j (<https://neo4j.com>) graph database. The patterns are expressed in Neo4j's query language Cypher. The experiment results using the two graphs and different exceptionality scoring functions are similar. Therefore, we only report our findings on WCGoals and exceptionality function χ_o , except that Section A.1 reports the discovered exceptional facts using both WCGoals and OscarWinners.

Methods Compared The experiments compared the performance of a breadth-first search method and several beam search methods (Section 3) coupled with different heuristics (Section 5.4):

- **Beam-Rdm**: Beam search that randomly selects child patterns.
- **Beam-Opt**: Beam search using h_{opt} in selecting child patterns.
- **Beam-Conv**: Beam search using h_{conv} in selecting child patterns.
- **Breadth-First**: The breadth-first search method that enumerates all possible patterns.

The family of beam search methods and Breadth-First differ in two ways. *Firstly*, beam search only visits a fixed number of patterns at each level of the pattern search tree, whereas Breadth-First visits all. *Secondly*, beam search visits the patterns by the decreasing order of their scores, whereas Breadth-First does not assume any order. The experiment results establish that, even though Breadth-First may evaluate more patterns than the beam search methods in a fixed time frame, it is not as effective as Beam-Conv which discovers more highly-scored context-subspace pairs using less time.

6.2 Efficiency

We measured how fast Maverick discovers highly-scored context-subspace pairs and how fast it explores the search space of patterns. We executed Maverick for multiple 2-minute runs and recorded a) the scores of discovered context-subspace pairs; b) the time when each context-subspace pair was discovered; and c) the number of visited patterns in the outer loop of the framework.

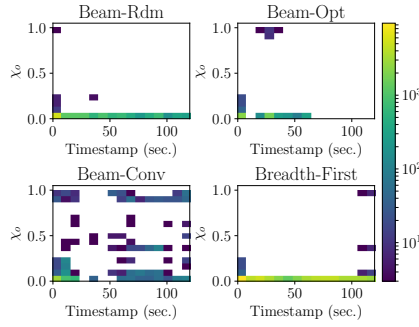
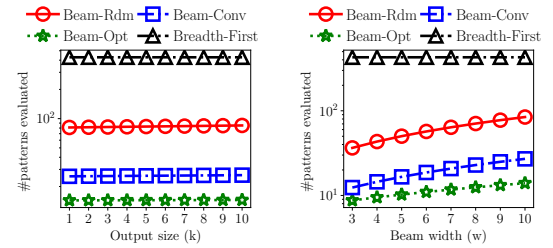


Figure 7: The heat map of exceptionality scores (χ_o) and timestamps of all the discovered context-subspace pairs during 2-minute runs for 10 entities of interest (v_0) in WCGoals ($k = 10$, $w = 10$).

Fig. 7 shows the heat map of the context-subspace pairs' exceptionality scores by their timestamps. It includes all the discovered context-subspace pairs during the 2-minute runs for 10 entities of interest in WCGoals. We run 10 times per entity for all the methods, since Beam-Rdm selects child patterns randomly. Both the output size k and the beam width w were set to 10. The 10 entities were randomly chosen from those that have highly-scored context-subspace pairs. Each bucket in the figure corresponds to a particular range of scores and a 8-second time frame in the 2-min run. The color of the bucket reflects how many context-subspace pairs (from all 100 runs for the 10 entities) discovered during the time frame fall into the corresponding score range. Intuitively, if the upper left portion of a heat map is more populated, the corresponding method performs better, since it means the method discovers highly-scored pairs faster. If the upper portion of a heat map is more populated, it means the method discovers more highly-scored pairs. The figure shows that Beam-Conv is both efficient and effective in discovering highly-scored context-subspace pairs. In contrast, Beam-Opt performed poorly. The results confirm the analysis in Section 5.4: preferring patterns that produce large contexts (h_{opt}) degrades not only the efficiency but also the effectiveness of Maverick. It is because such patterns are usually more expensive to evaluate and the produced contexts may include more varieties of entities, which makes the entity of interest less exceptional. With regard to Breadth-First, since it enumerates candidate patterns exhaustively, it may discover some highly-scored pairs that reside in the low levels of the pattern search tree. For example, some highly-scored pairs for entity Goal(46683) were found using the 2-edge pattern in Fig. 2a. Given that the number of patterns with no more than 2 edges is small (more details in the discussion of results regarding pruning strategies), Beam-Rdm is likely to hit such small patterns that define contexts in which the entity of interest is exceptional.

Fig. 8 shows the impact of output size k and beam width w on how many patterns Maverick can manage to evaluate, in other words, how many nodes in \mathbb{P} it can manage to visit. The y-axis is the average number of evaluated patterns across the aforementioned 10 runs. Since we observed similar results on the 10 entities from WCGoals, the figure only depicts the results on Goal(46683). Fig. 8a shows that varying k from 1 to 10 (fixing w at 10) barely had any impact on the number of evaluated patterns. Since k controls the number of context-subspace pairs that Maverick returns, it mainly affects EE, which is responsible for finding top- k subspaces with regard to each context. Thus k needs to be very large in order to have



a Varying k , fixing $w = 10$. b Varying w , fixing $k = 10$.

Figure 8: Effect of k and w on the number of evaluated patterns.

a significant impact on the number of evaluated patterns, since EE is the least time-consuming component, as explained as follows. Table 2 provides the breakdown of execution time of different search methods into the three components in the workflow—Context Evaluator (CE), Exceptionality Evaluator (EE), and Pattern Generator (PG). The results are the average of the runs which are the same as in Fig. 7. (The summation in each column is slightly less than 100%, since we do not include operations such as framework initialization in the breakdown.) Another observation from Table 2 is that the execution time of PG dominates more substantially in Beam-Opt and Beam-Conv than in Beam-Rdm and Breadth-First. The reason is PG in both Beam-Opt and Beam-Conv needs to compute h for each child pattern based on the pattern selection heuristics, which entails evaluating the child patterns to obtain the context sizes. In fact, on average, PG in Beam-Opt and Beam-Conv spent more than 99% and 96% of its time on applying the heuristics.

Table 2: Breakdown of execution time by components.

| | Beam-Rdm | Beam-Opt | Beam-Conv | Breadth-First |
|----|----------|----------|-----------|---------------|
| CE | 25.52% | 1.56% | 1.90% | 28.36% |
| EE | 0.41% | 0.65% | 0.32% | 2.79% |
| PG | 61.49% | 97.69% | 95.92% | 53.89% |

Fig. 8b depicts the results when w varied from 3 to 10 and k was fixed at 10. It shows the number of evaluated patterns increased by w in the three beam search methods. When w increases, the methods evaluate more patterns from lower levels in the pattern search tree, which have less edges and can be evaluated more efficiently than those from higher levels. In a fixed time frame, the methods can then evaluate more patterns in total, as shown in the figure. Fig. 10 shows the average time that Context Evaluator (Alg. 2) spends on pattern evaluation increases when the level of pattern (i.e., the number of edges) increases. Since Breadth-First does not need to calculate scores for patterns and does not have a limit on the number of patterns to visit at each level, it tends to evaluate more patterns but may only evaluate patterns at low levels. Fig. 9a compares the numbers of patterns evaluated at different levels by the four methods, when both k and w stayed at 10. Breadth-First evaluated patterns up to level 3 and spent most of its time on level 3. On the contrary, the beam search methods evaluated at most 10 patterns at each level and covered more levels.

Fig. 8b also suggests that Beam-Rdm evaluated more patterns than Beam-Conv and Beam-Opt. It is because Beam-Rdm (like Breadth-First) does not compute scores for child patterns, which is expensive. Since Beam-Opt favors patterns that define larger contexts, it evaluated the fewest patterns since it spent more time to calculate the sizes of the contexts. On the other hand, Beam-Conv prefers patterns defining smaller contexts, which allowed it

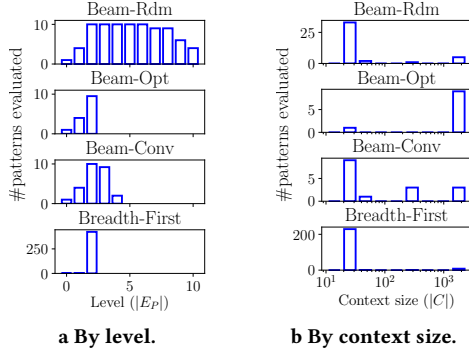


Figure 9: Number of evaluated patterns by level and context size.

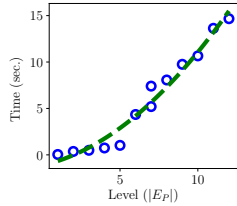


Figure 10: Time spent by Context Evaluator (Alg. 2) on evaluating patterns at different levels.

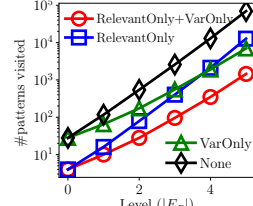


Figure 11: The pruning power of different pruning strategies.

to evaluate more patterns. This is verified in Fig. 9b, which shows the numbers of evaluated patterns with different context sizes, when k and w were both 10.

Effect of pruning strategies We examined the effectiveness of the two pruning rules from Section 5.3 by comparing the following pruning strategies. In order to comprehensively compare these strategies, we used Breadth-First as the search method since it exhaustively enumerates all possible candidate patterns at all levels.

- None: No child pattern pruning rule is applied;
- RelevantOnly (Rule 1);
- VarOnly (Rule 2);
- RelevantOnly+VarOnly: Apply both *RelevantOnly* and *VarOnly*.

Fig. 11 shows the number of patterns visited by Breadth-First on the data graph in Fig. 1. The figure reveals that both rules can significantly reduce the number of candidate patterns. For instance, there are 69,582 candidate patterns at level 5 when no pruning rule is applied (*None*). The number is reduced to 12,740 and 6,963 by following *RelevantOnly* and *VarOnly*, respectively. It is further reduced to 1,448 with both rules applied (*RelevantOnly+VarOnly*). The figure also shows that the number of patterns still grows exponentially to the level of the pattern search tree even with both pruning rules applied, which suggests an enormous search space of patterns. Since *VarOnly* is stricter than *RelevantOnly*, as *VarOnly* only allows expanding on variable nodes, the growth rate of *VarOnly* can be smaller than *RelevantOnly*. Fig. 11 also confirms that.

Effect of upper bound of exceptionality functions Fig. 12 depicts the effect of using upper bound functions in pruning subspaces (Section 4.3). It shows the time and the number of subspaces visited for Game(903)—one of the 10 entities used in Fig. 7—with/without applying upper bound functions under varying k . (Results for the other 9 entities are similar.) The measures are averages of 10 runs. The figure verifies that the upper bound functions significantly

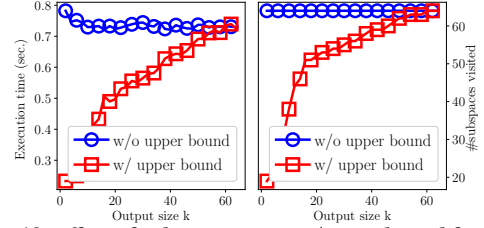


Figure 12: Effect of subspace pruning (upper bound functions).

improve the performance of exceptionality score calculations. Under relatively small k (e.g. 10), the execution time was reduced by more than half when the upper bound was applied. As k increased, the upper bound function's pruning power gradually diminished. Eventually, it was no longer able to prune any subspaces after $k = 60$.

6.3 Effectiveness

Section A.1 reports a few examples of discovered exceptional facts using both WCGoals and OscarWinners. We also conducted experiments to verify if Maverick can effectively discover highly-scored context-subspace pairs. Fig. 13 shows the score distributions of the top-10 context-subspace pairs for the same 10 entities in Section 6.2. There were 10 2-minute runs per entity. Both k and w were set to 10. The results in Fig. 13 are averaged over all entities and all runs. The last row are the results when Maverick uses patterns mined by a frequent pattern (FP) mining algorithm [16] as candidate patterns, instead of the ones discovered in \mathbb{P} . We set the minimum support to be 1,000, as lower value leads to excessive execution time.⁵

The results show that the output of Beam-Rdm mainly consists of pairs scored low. It is expected because the chance of hitting a promising pattern by a random method is very low due to the large search space of patterns. It is not surprising either to observe Beam-Opt performed badly as explained in Section 6.2. In contrast, Beam-Conv significantly outperformed other beam search methods, as it found much more highly-scored context-subspace pairs. It also found substantially more highly-scored pairs than Breadth-First in score range $[0.8-1.0]$. This observation confirms that a wide pattern search tree hinders Breadth-First's performance. Using FPs was not effective in discovery of exceptional facts. There are mainly two reasons. 1) Due to practical considerations such as efficiency and resources, FP mining techniques usually consider only node types (e.g., Team) but not node IDs (e.g., BRA). 2) An FP mining algorithm is not designed for individual entities. This leads to two consequences. One is that there may be no FP that defines some context for a given entity. The other is that the exceptionality of the entity may not be revealed in the contexts defined by the FPs. In fact, the experiments on WCGoals yielded only 12 FPs and all of them are about only two node types: PARTICIPANT and GOAL.

We also use a variation of *coverage error* [35] to measure the effectiveness of the four methods. For each method, we evaluated the result of its 2-minute run, using the result of its 10-hour run as the ground truth. The ground truth is the list of discovered context-subspace pairs during the 10-hour run, ranked by their exceptionality scores. Given the set of discovered context-subspace pairs in a 2-minute run, H , the coverage error is the

⁵The algorithm did not finish after more than 10 hours on graph WCGoals when the minimum support was set to 500.

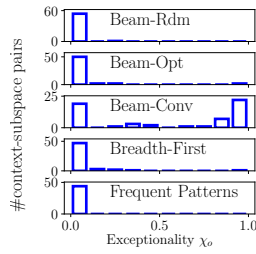


Figure 13: Score distributions of top-10 context-subspace pairs for 10 entities, 10 2-minute runs per entity.

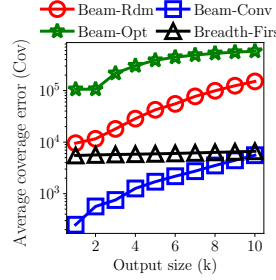


Figure 14: Average coverage error on 10 entities. Beam width 10.

average rank position of the pairs in the ground truth, defined by $Cov = \frac{1}{|H|} \sum_{(C,A) \in H} \text{rank}(C,A)$. Fig. 14 reports the average coverage error of each method under varying output size k . Table 3 shows the average and median coverage errors under varying beam width w . In Fig. 14, the coverage error of Beam-Conv is less than other methods by orders of magnitude, which suggests that Beam-Conv found highly-scored context-subspace pairs. Table 3 shows that coverage error decreases when beam width increases. The reason is that a wider beam leads to more patterns visited at every level and thus a better coverage of patterns. It is especially beneficial when highly-scored pairs reside in patterns at lower levels.

Table 3: The effect of beam width (w) on the coverage errors of top-10 context-subspace pairs of 10 entities. In each cell: the average and the median coverage errors. Both numbers are the smaller the better.

| w | Beam-Rdm | Beam-Opt | Beam-Conv | Breadth-First |
|-----|---------------|---------------|------------------|---------------|
| 3 | 3375.7/2636.9 | 2151.0/1071.5 | 49.7/12.5 | 383.0/390.5 |
| 4 | 3293.2/1607.3 | 2675.7/1622.1 | 52.1/12.5 | 383.0/390.5 |
| 5 | 2743.2/1871.2 | 2418.3/1550.8 | 30.2/26.0 | 383.0/390.5 |
| 6 | 2890.9/1809.2 | 2288.7/1259.7 | 20.6/1.0 | 383.0/390.5 |
| 7 | 2821.4/1398.9 | 1789.3/1259.7 | 21.8/1.0 | 383.0/390.5 |
| 8 | 2646.4/1818.6 | 1721.5/1168.8 | 78.3/3.0 | 383.0/390.5 |
| 9 | 2262.8/1653.4 | 1365.5/1107.3 | 36.6/4.2 | 383.0/390.5 |
| 10 | 2720.8/1619.9 | 1365.5/1107.3 | 58.4/22.1 | 383.0/390.5 |

7 RELATED WORK

In exceptional fact discovery, the output context-subspace pairs can be viewed as a way of explaining outliers. Most conventional outlier detection solutions, including those for graphs, focus on finding outliers but do not explain why they are outlying. For example, CODA [17] finds a list of community outliers, and FocusCO [28] clusters an attributed graph and then discovers outliers in the clusters. Besides the limitation that both approaches are only suitable for homogeneous graphs, it is up to users to figure out the explanations of the outliers. Although these two systems make such explanations easier by providing the communities or clusters in which the outliers reside, it still requires substantial expertise to summarize the communities/clusters' characteristics. A few works improve the interpretation of outliers' outlyingness [1, 23, 33]. For instance, systems such as [1] and [23] use visualization to help users identify outliers and potentially discover their outlying aspects.

Although most existing outlying aspects mining approaches focus on finding global outlying aspects and do not consider contexts [15, 36], there are a few attempts to find contextual outlying aspects [2, 3, 32, 40]. The general Maverick framework allows users

to adopt any exceptionality measure in the literature such as outlierness [3]. Although Maverick focuses on categorical attributes at this stage, it can be extended for numerical attributes so that measures such as skyline points [32], promotiveness [40], outlierness [2], outlyingness rank [15], and z-score [36] can be adopted in the framework.

Trummer et al. [34] developed the SURVEYOR system to mine the dominant opinion on the Web regarding whether a subjective property (e.g., "safe cities") applies to an entity. This is useful for populating a knowledge base with ground truth for answering subjective queries. While they focus on deriving entities' hidden properties which may or may not be exceptional, Maverick focuses on finding exceptional entities using existing data in knowledge graphs.

8 CONCLUSION AND FUTURE WORK

In this paper, we study the problem of discovering exceptional facts about entities in knowledge graphs. Each exceptional fact consists of a pair (*context*, *subspace*). To tackle the challenge of exploring the exponential large search spaces of both contexts and subspaces, we propose a beam search based framework, Maverick, which applies a set of heuristics during the discovery. The experiment results show that our proposed framework is both efficient and effective for discovering exceptional facts.

Interesting future work can be pursued along several directions on both *efficiency* and *usability*. With regard to *efficiency*, the current system focuses on finding exceptional facts given a specific entity. What is more appealing is a discovery mode in which the system automatically finds facts for all entities. Straightforwardly applying the system on a large knowledge graph will thus lead to exhaustive and repetitive computations for a huge number of entities. Devising algorithms for sharing computations across different entities can significantly increase the system's capability over large knowledge graphs. Furthermore, the system currently considers a static knowledge graph which in reality constantly evolves and grows. To produce up-to-date facts, one has to repeatedly apply the system, which is not practical given the sheer size and change frequency of real-world knowledge bases. Hence, another substantial improvement of the system will be adding incremental exceptional fact discovery algorithms. With regard to *usability*, how to present exceptional facts poses intriguing challenges related to user interface, data visualization, and exceptionality measures. For instance, it is appealing for the system to produce natural language descriptions of the facts. While exceptionality measures such as one-of-the-few may lend itself to simple template-based translations, it is much more challenging to precisely convey facts ranked by more complex measures such as outlyingness and isolation score. Moreover, while our user study helps gain insights into different exceptionality scoring functions, to thoroughly understand their strengths and limitations, a more comprehensive and larger scale user study is worth doing.

ACKNOWLEDGMENTS

The work is partially supported by NSF grants IIS-1408928 and IIS-1719054. Any opinions, findings, and conclusions or recommendations expressed in this publication are those of the authors and do not necessarily reflect the views of the funding agencies.

REFERENCES

- [1] Leman Akoglu, Mary McGlohon, and Christos Faloutsos. 2010. Oddball: Spotting anomalies in weighted graphs. In *PAKDD*. 410–421.
- [2] Fabrizio Angiulli, Fabio Fasseti, Giuseppe Manco, and Luigi Palopoli. 2016. Outlying property detection with numerical attributes. *DMKD* (2016), 1–30.
- [3] Fabrizio Angiulli, Fabio Fasseti, and Luigi Palopoli. 2009. Detecting Outlying Properties of Exceptional Objects. *TODS* 34, 1 (April 2009), 7:1–7:62.
- [4] Sören Auer, Christian Bizer, Georgi Kobilarov, Jens Lehmann, Richard Cyganiak, and Zachary Ives. 2007. DBpedia: A Nucleus for a Web of Open Data. In *ISWC*.
- [5] Roberto Bisiani. 1992. Beam Search. In *Encyclopedia of Artificial Intelligence* (2nd ed.), Stuart C Shapiro (Ed.). Wiley-Interscience, 1467–1468.
- [6] Kurt Bollacker, Colin Evans, Praveen Paritosh, Tim Sturge, and Jamie Taylor. 2008. Freebase: a collaboratively created graph database for structuring human knowledge. In *SIGMOD*. 1247–1250.
- [7] Antoine Bordes, Nicolas Usunier, Alberto Garcia-Duran, Jason Weston, and Oksana Yakhnenko. 2013. Translating embeddings for modeling multi-relational data. In *NIPS*. 2787–2795.
- [8] Mihaela A. Bornea, Julian Dolbey, Anastasios Kementsietsidis, Kavitha Srinivas, Patrick Dantressangle, Octavian Udrea, and Bishwaranjan Bhattacharjee. 2013. Building an Efficient RDF Store over a Relational Database. In *SIGMOD*.
- [9] Dan Bricklayer and R.V. Guha. 2014. RDF Schema 1.1. *W3C Recommendation* (2014).
- [10] J. Cohen. 1988. *Statistical Power Analysis for the Behavioral Sciences*. Lawrence Erlbaum Associates.
- [11] Sarah Cohen, James T. Hamilton, and Fred Turner. 2011. Computational Journalism. *Commun. ACM* 54, 10 (Oct. 2011), 66–71.
- [12] Sarah Cohen, Chengkai Li, Jun Yang, and Cong Yu. 2011. Computational Journalism: A Call to Arms to Database Researchers. In *CIDR*. 148–151.
- [13] Richard Cyganiak, David Wood, and Markus Lanthaler. 2014. RDF 1.1 concepts and abstract syntax. *W3C Recommendation* (2014).
- [14] Mike Dean, Guus Schreiber, Sean Bechhofer, Frank van Harmelen, Jim Hendler, Ian Horrocks, Deborah L McGuinness, Peter F Patel-Schneider, and Lynn Andrea Stein. 2004. OWL web ontology language reference. *W3C Recommendation* (2004).
- [15] Lei Duan, Guanting Tang, Jian Pei, James Bailey, Akiko Campbell, and Changjie Tang. 2015. Mining outlying aspects on numeric data. *DMKD* 29, 5 (2015), 1116–1151.
- [16] Mohammed Elseidy, Ehab Abdelhamid, Spiros Skiadopoulos, and Panos Kalnis. 2014. Grami: Frequent subgraph and pattern mining in a single large graph. *Proceedings of the VLDB Endowment* 7, 7 (2014), 517–528.
- [17] Jing Gao, Feng Liang, Wei Fan, Chi Wang, Yizhou Sun, and Jiawei Han. 2010. On community outliers and their efficient detection in information networks. In *KDD*. 813–822.
- [18] Liqiang Geng and Howard J Hamilton. 2006. Interestingness measures for data mining: A survey. *ACM Computing Surveys (CSUR)* 38, 3 (2006), 9.
- [19] Steve Harris and Andy Seaborne. 2013. SPARQL 1.1 Query Language. *W3C recommendation* 15 (2013).
- [20] Naeemul Hassan, Bill Adair, James T. Hamilton, Chengkai Li, Mark Tremayne, Jun Yang, and Cong Yu. 2015. The Quest to Automate Fact-Checking. In *Proceedings of the 2015 Computation+Journalism Symposium*.
- [21] Naeemul Hassan, Afroza Sultana, You Wu, Gensheng Zhang, Chengkai Li, Jun Yang, and Cong Yu. 2014. Data In, Fact Out: Automated Monitoring of Facts by FactWatcher. *PVLDB, demonstration description* 7, 13 (2014), 1557–1560.
- [22] Xiao Jiang, Chengkai Li, Ping Luo, Min Wang, and Yong Yu. 2011. Prominent streak discovery in sequence data. In *KDD*. 1280–1288.
- [23] U Kang, Jay-Yoon Lee, Danai Koutra, and Christos Faloutsos. 2014. Net-ray: Visualizing and mining billion-scale graphs. In *PAKDD*. 348–361.
- [24] Hans-Peter Kriegel, Peer Kröger, and Arthur Zimek. 2010. Outlier detection techniques. In *KDD*.
- [25] Yankai Lin, Zhiyuan Liu, Maosong Sun, Yang Liu, and Xuan Zhu. 2015. Learning Entity and Relation Embeddings for Knowledge Graph Completion.. In *AAAI*. 2181–2187.
- [26] Fei Tony Liu, Kai Ming Ting, and Zhi-Hua Zhou. 2008. Isolation forest. In *ICDM*. IEEE, 413–422.
- [27] Jorge Pérez, Marcelo Arenas, and Claudio Gutierrez. 2009. Semantics and complexity of SPARQL. *TODS* 34, 3 (2009), 16.
- [28] Bryan Perozzi, Leman Akoglu, Patricia Iglesias Sánchez, and Emmanuel Müller. 2014. Focused Clustering and Outlier Detection in Large Attributed Graphs. In *KDD*.
- [29] Stuart Russell and Peter Norvig. 2010. *Artificial Intelligence: A Modern Approach (3rd Edition)*. Prentice Hall.
- [30] Ron Rymon. 1992. Search through systematic set enumeration. *Technical Reports MS-CIS-92-66, Department of Computer and Information Science, University of Pennsylvania* (1992).
- [31] Fabian M Suchanek, Gjergji Kasneci, and Gerhard Weikum. 2007. Yago: a core of semantic knowledge. In *WWW*. 697–706.
- [32] Afroza Sultana, Naeemul Hassan, Chengkai Li, Jun Yang, and Cong Yu. 2014. Incremental Discovery of Prominent Situational Facts. In *ICDE*. 112–123.
- [33] Hanghang Tong and Ching-Yung Lin. 2012. Non-negative residual matrix factorization: problem definition, fast solutions, and applications. *Statistical Analysis and Data Mining* 5, 1 (2012), 3–15.
- [34] Immanuel Trummer, Alon Halevy, Hongrae Lee, Sunita Sarawagi, and Rahul Gupta. 2015. Mining Subjective Properties on the Web. In *SIGMOD*. 1745–1760.
- [35] Grigorios Tzoumakas, Ioannis Katakis, and Ioannis Vlahavas. 2009. Mining multi-label data. In *Data mining and knowledge discovery handbook*. Springer, 667–685.
- [36] Nguyen Xuan Vinh, Jeffrey Chan, Simone Romano, James Bailey, Christopher Leckie, Kotagiri Ramamohanarao, and Jian Pei. 2016. Discovering outlying aspects in large datasets. *DMKD* 30, 6 (2016), 1520–1555.
- [37] Denny Vrandečić and Markus Krötzsch. 2014. Wikidata: a free collaborative knowledgebase. *Commun. ACM* 57, 10 (2014), 78–85.
- [38] Xiang Wang and Ian Davidson. 2009. Discovering contexts and contextual outliers using random walks in graphs. In *ICDM*. 1034–1039.
- [39] Zhen Wang, Jianwen Zhang, Jianlin Feng, and Zheng Chen. 2014. Knowledge Graph Embedding by Translating on Hyperplanes.. In *AAAI*. 1112–1119.
- [40] Tianyi Wu, Dong Xin, Qiaozhu Mei, and Jiawei Han. 2009. Promotion analysis in multi-dimensional space. *PVLDB* 2, 1 (2009), 109–120.
- [41] You Wu, Pankaj K. Agarwal, Chengkai Li, Jun Yang, and Cong Yu. 2012. On “One of the Few” Objects. In *KDD*. 1487–1495.
- [42] You Wu, Pankaj K. Agarwal, Chengkai Li, Jun Yang, and Cong Yu. 2017. Computational Fact Checking through Query Perturbations. *TODS* 42, 1 (Jan. 2017), 4:1–4:41.
- [43] Yuehua Xu and Alan Fern. 2007. On learning linear ranking functions for beam search. In *ICML*. 1047–1054.
- [44] Gensheng Zhang, Xiao Jiang, Ping Luo, Min Wang, and Chengkai Li. 2014. Discovering General Prominent Streaks in Sequence Data. *TKDD* 8, 2 (June 2014), 9:1–9:37.
- [45] Miao Zhou, Chunhong Zhang, Xiao Han, Yang Ji, Zheng Hu, and Xiaofeng Qiu. 2016. Knowledge Graph Completion for Hyper-relational Data. In *International Conference on Big Data Computing and Communications*. 236–246.

A APPENDIX

A.1 Case Study

To illustrate the effectiveness of Maverick, we present below some examples of exceptional facts discovered by Maverick in both graph WCGoals and graph OscarWinners.

Goal(46683) is the only own goal in Brazil’s World Cup history.

Exceptionality $\chi_o = 0.986$
Subspace $\{(awared-to, \rightarrow)\}$
Context $C_{Goal(46683)}^{P, x_0}$, where $P = \{(x_0, scored-by, x_1), (x_1, play-for, BRA)\}$

Indeed, among all the 221 goals that were scored by Brazil players in the *FIFA World Cup Finals* tournaments, Goal(46683), which was awarded to Croatia, was the only goal not awarded to Brazil. This exceptional fact has a very high score.

Among all the crew members of Oscar winning films, Paul J. Franklin (FilmCrew(7674)) is the only crew member with role Computer Animation.

Exceptionality $\chi_f = 0.784$
Subspace $\{(film-crew-role, \rightarrow)\}$
Context $C_{FilmCrew(7674)}^{P_0, x_0}$

This example demonstrates the utility of Maverick in revealing data errors, as motivated in Section 1. While investigating why this entity is exceptional, an analyst will realize the exceptional fact is due to a data error. An edge mistakenly links from node Paul J. Franklin to a genre node Computer Animation which is incorrectly used in this case as a role node. The correct crew role node should have been Computer Animator.

Goal(23464) is the only goal awarded to Paraguay, among all the goals scored in matches hosted in Mexico City that had at least two goals.

Exceptionality $\chi_f = 0.983$
Subspace $\{(awarded-to, \rightarrow)\}$
Context $C_{Goal(24227)}^{P, x_1}$, where $P = \{(x_0, goal, x_1), (x_0, goal, x_2), (x_0, venue, Mexico\ City)\}$

There are in total 62 goals scored in matches hosted in Mexico City, among which 58 were scored in 18 multiple-goal matches. These 58 goals were awarded to 12 different teams. Paraguay is the only team that was awarded only one of the 58 goals.

Game(899) is one of the only two games in which the home team was Senegal, among all the games where there was a player wearing the number 21 shirt.

Exceptionality $\chi_f = 0.959$
Subspace $\{(home, \rightarrow)\}$
Context $C_{Game(899)}^{P, x_1}$, where $P = \{(x_0, bibnum, Bibnum(21)), (x_0, participate-in, x_1)\}$

In 761 games some player wore number 21. Game(899) is one of the only two such games in which the home team was Senegal.

Among the Oscar winning films produced in the United States, The Lord of the Rings: The Return of the King (Film(31768)) is one of the only 7 films that were also produced in New Zealand.

Exceptionality $\chi_o = 0.676$
Subspace $\{(country, \rightarrow)\}$
Context $C_{Film(31768)}^{P, x_0}$, where $P = \{(x_0, country, USA)\}$

There are in total 662 Oscar winning films produced in the United States, of which 545 were produced solely in the United States. Only 7 of the co-produced films were co-produced in New Zealand. However, the score of this fact is not as high as that of the last two facts, because China, Brazil, and a few other countries co-produced even less films.

A.2 Exceptionality Scoring Function: Isolation Score χ_i

Isolation Score χ_i The isolation score χ_i is inspired by iForest [26] and iPath [36]. Both iForest and iPath are applicable on real value attributes. By randomly choosing a pivot in the range of an attribute, both methods iteratively split a set of entities into two disjoint subsets, until the entities in each set have an identical value. The iForest score and iPath score are defined using the number of splits applied. iPath only iteratively splits the subsets containing the entity of interest. The entity's score is the number of splits until the entity has the same value as other entities in its subsuming set. iForest splits all subsets. Essentially, both iForest and iPath follow the *minimum description length* principle, and both scores are functions of the estimated description length of the entity's attribute value, which is the number of splits. Inspired by iForest score, we define isolation score χ_i as follows:

$$\chi_i(v, A, C) = 1 - 2^{-\frac{-\log_2 p_{v,A}}{-\sum_{S \in S_A} (p_S \times \log_2 p_S)}} \quad (6)$$

where the numerator in the exponent is the description length of v 's attribute value, while the denominator is the average description length of attribute values in subspace A . Intuitively, if $v.A$ is peculiar, then the description length of $v.A$ is longer than average and $\chi_i(v, A, C)$ is closer to 1.

For example, let the conditions be the same as the ones used in explaining χ_o and χ_f : $v_0 = G_1$, $C = C_{G_1}^{P_1, ?g} = \{G_1, G_2, G_3\}$, and $A = \{(awarded-to, \rightarrow)\}$. According to Table 1, $-\sum_{S \in \{\{BRA\}, \{CRO\}\}} (p_S \times \log_2 p_S) = -(\frac{1}{3} \times \log_2 \frac{1}{3} + \frac{2}{3} \times \log_2 \frac{2}{3}) = \log_2 3 - \frac{2}{3} = 0.91$, $-\log_2 p_{G_1.A} = 1.58$, then $\chi_f(G_1, A, C) = 0.7$. Similarly, $\chi_f(G_1, \{(awarded-to, \rightarrow), (scored-by, \rightarrow)\}, C) = 0$.

Theorem 6 (Upper bound of χ_i). $upper_i(v, A, C) \geq \chi_i(v, A', C)$, given the following definition:

$$upper_i(v, A, C) = 1 - 2^{-\frac{-\log_2 \frac{1}{|C|}}{-q_{v.A} - \sum_{S \in S_A} (p_{v.A} \setminus \{v.A\}) (p_S \times \log_2 p_S)}} \quad (7)$$

where $q_{v.A} = \frac{1}{|C|} \times \log_2 \frac{1}{|C|} + (p_{v.A} - \frac{1}{|C|}) \times \log_2 (p_{v.A} - \frac{1}{|C|})$.

PROOF. Please refer to Appendix A.5 for the proof. ■

A.3 Complexity Analysis of the Beam Search Method

This section presents a brief analysis of the complexity of our beam search method. Let w be the beam size, if Maverick stops at level l of the Hasse diagram \mathbb{P} , which is the search sapce of patterns, then it evaluates exceptionality in at least $(l-1)w+1$ contexts (patterns). For each context, the complexity of computing exceptionality is $O(2^{|A_{v_0}|})$, as we discussed in Section 4. Assume the average degree of an entity is d , the average number of variables in a pattern of size k is $\frac{k}{2}$, then a pattern of size k has $(d+1)^{k/2}$ children. For each child, Maverick needs to compute h scores for selecting promising candidates (Section 5.4). The main computational cost of h scores is the calculation of context sizes, which requires pattern evaluation. Given a pattern of size k , its evaluation can be done in $O(|E(G)|^k)$, since it may require k self-join operations on E_G . In sum, the estimated complexity is $O((w(l-1)+1)(2^{|A_{v_0}|} + \sum_{k=0}^l (d+1)^{k/2} \times |E(G)|^k)) = O(2^{|A_{v_0}|} + d^l |E(G)|^l)$.

A.4 User Study for Comparing Exceptionality Scoring Functions

We conducted a user study to assess the quality of four different exceptionality scoring functions χ_f (one-of-the-few, Section 2), χ_o (outlyingness, Section 4.2), χ_i (isolation score, Appendix A.2), and *frequency rank* (Eq. (1)). The measure *frequency rank* is simply to rank the exceptionality of entities in context C with respect to subspace A based on their attribute value frequency $p_{v.A}$, as defined in Eq. (1). The lower the frequency, the more exceptional it is. We compared the rankings of exceptional facts based on these functions as well as actual user preferences.

The user study participants were asked to choose a fact from a pair of facts that they deemed to be more exceptional. For quality assurance we manually crafted a set of trivial facts that are clearly non-exceptional or dull. We then formulated test pairs, of which each is composed of a regular fact and a trivial fact. The participants are expected to choose the regular fact as more exceptional. A participant's quality is thus gauged by their accuracy on the test pairs,

which were mixed together with regular pairs without disclosure to the participants.

We used 10 regular facts and thus 45 pairs of these facts. We randomly selected 10 entities from graph OscarWinners, and run Maverick using Beam-Conv to discover exceptional facts of the entities. The exceptionality scoring function used in the discovery was outlieriness (χ_o). Among all the facts discovered by Maverick, 10 facts are picked so that the scores of facts are roughly evenly distributed in $[0, 1]$. For each selected fact, χ_f , χ_i , and *frequency rank* are also calculated. We crafted 8 trivial facts and formulated 8 test pairs by pairing up these trivial facts with regular ones. Hence a participant responded to at most 53 pairs. The facts were presented to the participants in their natural language descriptions which we manually generated in the form of one-of-the-few facts. The study was conducted on line, for which the participants were solicited from computer science graduate students in the authors' institution. 4, 212 responses from 84 participants were recorded in total.

For each exceptionality scoring function, we constructed a vector X of 45 values corresponding to the 45 pairs of regular facts. For each pair, the value in X is the difference between the two facts' ranks according to the scoring function. We also constructed another vector Y , in which a value is the difference between how many participants favored one fact versus another in the corresponding pair. The correlation between the scoring function and the participants is calculated using the Pearson Correlation Coefficient (PCC) which is defined as $(E(XY) - E(X)E(Y)) / (\sqrt{E(X^2) - (E(X))^2} \sqrt{E(Y^2) - (E(Y))^2})$. A PCC value in the ranges of $[0.5, 1.0]$, $[0.3, 0.5]$ and $[0.1, 0.3]$ indicates a strong, medium and small positive correlation, respectively [10].

Table 4: User study results at different participant quality levels.

| # of pairs correct (\geq) | # of participants | χ_i | χ_f | χ_o | <i>frequency</i> <i>rank</i> |
|----------------------------------|----------------------|----------|----------|----------|---------------------------------|
| 0 | 84 | 0.370 | 0.564 | 0.295 | 0.429 |
| 1 | 84 | 0.370 | 0.564 | 0.295 | 0.429 |
| 2 | 78 | 0.392 | 0.585 | 0.317 | 0.412 |
| 3 | 74 | 0.410 | 0.597 | 0.335 | 0.400 |
| 4 | 66 | 0.449 | 0.642 | 0.378 | 0.338 |
| 5 | 53 | 0.491 | 0.649 | 0.422 | 0.377 |
| 6 | 43 | 0.614 | 0.730 | 0.558 | 0.283 |
| 7 | 22 | 0.738 | 0.831 | 0.696 | 0.061 |
| 8 | 9 | 0.750 | 0.864 | 0.711 | -0.007 |

The results of the user study are presented in Table 4, in which each row shows the PCC values calculated using participants at a different quality level (measured by number of test pairs the participants got correct, i.e., the first column in the table). We can make a few observations on the results. 1) As the quality level increases the number of participants that are accounted for at that level decreases, showing that the test pairs were successful in filtering out low performing participants. 2) For scoring functions χ_i , χ_f and χ_o , the correlation with participants steadily increases when the participants' quality increases. 3) In general, the scoring function that performed the best was χ_f , followed by χ_i and χ_o . The results show a strong correlation between these three functions and high-quality human participants, which suggests these

functions are effective in ranking the facts. The observation that χ_f performed the best could be due to a bias: the natural language descriptions of the facts were in the form of one-of-the-few facts. (On a side note, this suggests a strength of χ_f related to usability, as there is no clear way of directly expressing facts in line with χ_i and χ_o .) 4) On the contrary, *frequency rank* displayed a decrease in correlation as participant quality increases and its correlation was never strong. We reason that this could be due to low performing participants directly using frequency to hastily assess whether a fact is exceptional or not, without carefully examining the nature of the fact. The fact that *frequency rank* attains stronger correlation with lower-quality participants verifies that it cannot be used as a robust exceptionality scoring function, as explained in Section 2.

A.5 Proofs of Theorems

Proof of Theorem 2

PROOF. Let $\{p_{v.A}, p_{S_1}, \dots, p_{S_N}\}$ be the probability distribution of attribute values in subspace A . According to [3], for any $A' \supseteq A$, $\chi_o(v, A', C)$ is maximized when the additional attributes in $A' \setminus A$ preserve the current attribute value distribution, except that the additional attributes make v different from all other entities, i.e., the optimal distribution of attribute values in subspace A' is $\{p_{v.A'}, p_{v.A} - p_{v.A'}, p_{S_1}, \dots, p_{S_N}\}$, where $p_{v.A'} = \frac{1}{|C|}$. (Note that $p_S \geq \frac{1}{|C|}$ for any S .) In other words, the entities having value $v.A$ on subspace A are partitioned into v itself (having value $v.A'$ on subspace A') and the rest (having identical value on A'). Based on Eq. (3), after a few polynomial manipulations, which we omit here, we have $\chi_o(v, A', C) \leq \sum_{S \in S_A} p_S^2 - \frac{(2p_{v.A} + 1) \times |C| - 2}{|C|^2}$. ■

Proof of Theorem 6

PROOF. By Eq. (6), $\chi_i(v, A', C)$ is maximized when the denominator in the exponent is minimized and the numerator is maximized. Let $\{p_{v.A}, p_{S_1}, \dots, p_{S_N}\}$ be the probability distribution of attribute values in subspace A . Similar to the proof of Theorem 2, we prove that $\chi_i(v, A', C)$ is maximized when the distribution in subspace A' is $\{p_{v.A'}, p_{v.A} - p_{v.A'}, p_{S_1}, \dots, p_{S_N}\}$, where $p_{v.A'} = \frac{1}{|C|}$. Partition the entities having value S in A into two disjoint subsets that have values S_1 and S_2 in A' , respectively, i.e., $p_S^A = p_{S_1}^{A'} + p_{S_2}^{A'}$. Without loss of generality, assume $p_{S_1}^{A'} \leq p_{S_2}^{A'}$. We have

$$\begin{aligned} (1) \quad & -p_{S_1}^{A'} \log_2 p_{S_1}^{A'} - p_{S_2}^{A'} \log_2 p_{S_2}^{A'} \geq -p_S^A \log_2 p_S^A, \\ (2) \quad & -p_{S_1}^{A'} \log_2 p_{S_1}^{A'} - p_{S_2}^{A'} \log_2 p_{S_2}^{A'} \geq \\ & -p_{v.A'}^{A'} \log_2 p_{v.A'}^{A'} - (p_S^A - p_{v.A'}^{A'}) \log_2 (p_S^A - p_{v.A'}^{A'}). \end{aligned}$$

They can be derived by

$$\begin{aligned} (1) \quad & -p_{S_1}^{A'} \log_2 p_{S_1}^{A'} - p_{S_2}^{A'} \log_2 p_{S_2}^{A'} \geq -p_{S_1}^{A'} \log_2 p_{S_1}^{A'} - p_{S_2}^{A'} \log_2 p_{S_2}^{A'} \\ & = -(p_{S_1}^{A'} + p_{S_2}^{A'}) \log_2 p_{S_2}^{A'} = -p_S^A \log_2 p_{S_2}^{A'} \geq -p_S^A \log_2 p_S^A, \\ (2) \quad & \text{Let } p_S^A = n p_{v.A'}^{A'} = \frac{n}{|C|}, p_{S_1}^{A'} = m p_{v.A'}^{A'} = \frac{m}{|C|}, \text{ then } 1 \leq m < n, \\ & -p_{S_1}^{A'} \log_2 p_{S_1}^{A'} - p_{S_2}^{A'} \log_2 p_{S_2}^{A'} - (-p_{v.A'}^{A'} \log_2 p_{v.A'}^{A'} - (p_S^A - p_{v.A'}^{A'}) \log_2 (p_S^A - p_{v.A'}^{A'})) \\ & = -m p_{v.A'}^{A'} \log_2 m p_{v.A'}^{A'} - (n p_{v.A'}^{A'} - m p_{v.A'}^{A'}) \log_2 (n p_{v.A'}^{A'} - m p_{v.A'}^{A'}) + p_{v.A'}^{A'} \log_2 p_{v.A'}^{A'} + (n p_{v.A'}^{A'} - p_{v.A'}^{A'}) \log_2 (n p_{v.A'}^{A'} - p_{v.A'}^{A'}) \\ & = p_{v.A'}^{A'} ((n-1) \log_2 (n-1) - m \log_2 m - (n-m) \log_2 (n-m)) \\ & = p_{v.A'}^{A'} \log_2 \frac{(n-1)^{(n-1)}}{m^m (n-m)^{(n-m)}} \\ & \geq p_{v.A'}^{A'} \log_2 \min\left(\frac{(n-1)^{(n-1)}}{1^{1(n-1)}(n-1)}, \frac{(n-1)^{(n-1)}}{(n-1)^{(n-1)}(n-(n-1))^{(n-(n-1))}}\right) = 0. \end{aligned}$$

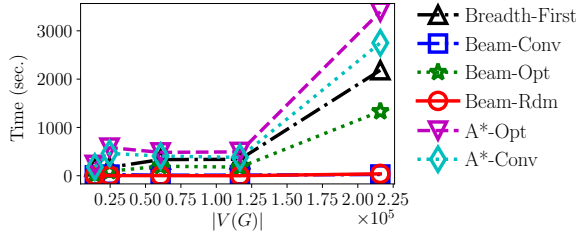


Figure 15: Execution time of enumerating all candidates up to level 2 by different orders of graphs.

As a result, since $-\log_2 p_{v.A'} \geq \log_2 p_i$ for any $p_i \geq p_{v.A'}$, by dividing $p_{v.A}$ to $p_{v.A'}$ and $p_{v.A} - p_{v.A'}$, $\chi_i(v, A', C)$ is maximized. In other words, the optimal distribution in subspace A' is $\{p_{v.A'}, p_{v.A} - p_{v.A'}, p_{S_1}, \dots, p_{S_N}\}$. ■

Proof of Theorem 3

PROOF. Given a WCC $W \in \mathcal{W}$, it has at least one subgraph of order i , for every $i \in [1, |V_W|]$. For each subgraph of size i , there are 2^i corresponding patterns that can be constructed by replacing some nodes with variables. Hence, for each W , there are at least $\sum_{i=1}^{|V_W|} 2^i = 2^{|V_W|+1} - 2$ patterns. Since every such pattern is isomorphic to a subgraph of W , it is guaranteed to be valid. Note that two patterns of the same order constructed from two subgraphs in two different WCCs can be equivalent if all their nodes are variables. Therefore, each $W \in \mathcal{W}$ has at most $|V_W|$ patterns that are equivalent to others. There are at least $\max_W |V_W|$ unique patterns in which all nodes are variables. Thus, after excluding double-counted patterns,

$$\begin{aligned} |V_{\mathbb{P}}| &\geq \sum_W (2^{|V_W|+1} - 2) - \sum_W |V_W| + \max_W |V_W| \\ &= \sum_W (2^{|V_W|+1} - 2) - |V_G| + \max_W |V_W|. \end{aligned}$$

Proof of Theorem 4

PROOF. Since P' is a child of P , P' has one edge more than P . Suppose $(u, l, w) = E_{P'} \setminus E_P$, and f' is the bijection $f' : V_{P'} \rightarrow V_{M'}$. We prove the theorem by constructing M . More specifically, let $E_M = E_{M'} \setminus \{(f'(u), l, f'(w))\}$, and $V_M = \cup_{(v_i, l', v_j) \in E_M} \{v_i, v_j\}$. We can construct a bijection $f : V_P \rightarrow V_M$ such that $f(u) = f'(u)$ for any $u \in V_P$. Since f' satisfies the edge isomorphism, f also satisfies it, i.e., $\forall (v_i, l', v_j) \in E_P, (f(v_i), l', f(v_j)) \in E_M$, and vice versa. By Definition 2, M is a match to P . ■

Proof of Theorem 5

PROOF. Since both ends of e are entities, we have $X_P = X_{P'}$. By Theorem 4, $\forall M' \in \mathcal{M}_{P'}$, there exists $M \in \mathcal{M}_P$ which is a subgraph of M' . Let $e' = E_{M'} \setminus E_M$, then $e' = e$ by Definition 2. Therefore, $\forall x \in X_{P'}, R_x^{P'} \subseteq R_x^P$. Similarly, $\forall M \in \mathcal{M}_P$, the graph $M + e$ is a match to P' . As a result, $\forall x \in X_P, R_x^P \subseteq R_x^{P'}$. In sum, $\forall x \in X_P, R_x^P = R_x^{P'}$, and $C_{v_0}^P = \{R_{x'}^P \mid x' \in X_P, v_0 \in R_{x'}^P\} = \{R_{x'}^{P'} \mid x' \in X_{P'}, v_0 \in R_{x'}^{P'}\} = C_{v_0}^{P'}$. ■

A.6 Scalability

To test the scalability of Maverick, we generated synthetic graphs using a benchmark data generator *BSBM* about products, vendors, consumers, and reviews (<http://wifo5-03.informatik.uni-mannheim.de/bizer/berlinsparqlbenchmark/spec/BenchmarkRules>). We varied the number of products from 100 to 2,000, which resulted in graphs of order $|V(G)|$ from 14,195 to 215,895. Fig. 15 shows how the execution time of all algorithms increased along with the order of the graphs. Given that Beam-Conv and Beam-Rdm need less computation in candidate selection, both scaled much more gracefully than others. The figure also shows the execution time of two A* algorithms guided by h_{opt} and h_{conv} , respectively. The results confirmed that computing h_{opt} is more expensive than h_{conv} . This is because computing h_{opt} requires examining all the variables in a pattern, while computing h_{conv} only requires that for context-defining variables in the parent. Due to the overhead of computing h scores, the A* algorithms are shown more expensive than Breadth-First.