

# Holistic and Compact Selectivity Estimation for Hybrid Queries over RDF Graphs<sup>\*</sup>

Andreas Wagner<sup>1</sup>, Veli Bicer<sup>2</sup>, Thanh Tran<sup>3</sup>, and Rudi Studer<sup>1</sup>

<sup>1</sup> Karlsruhe Institute of Technology

<sup>2</sup> IBM Research Centre Dublin

<sup>3</sup> San Jose State University

{a.wagner,rudi.studer}@kit.edu, velibice@ie.ibm.com,  
ducthanh.tran@sjsu.edu

**Abstract.** Many RDF descriptions today are *text-rich*: besides *structured* data they also feature much *unstructured* text. Text-rich RDF data is frequently queried via predicates matching structured data, combined with string predicates for textual constraints (*hybrid queries*). Evaluating hybrid queries efficiently requires means for *selectivity estimation*. Previous works on selectivity estimation, however, suffer from inherent drawbacks, which are reflected in efficiency and effectiveness issues. We propose a novel estimation approach, *TopGuess*, which exploits topic models as data synopsis. This way, we capture correlations between structured and unstructured data in a *holistic and compact* manner. We study TopGuess in a theoretical analysis and show it to guarantee a linear space complexity w.r.t. text data size. Further, we show selectivity estimation time complexity to be independent from the synopsis size. In experiments on real-world data, TopGuess allowed for great improvements in estimation accuracy, without sacrificing efficiency.

## 1 Introduction

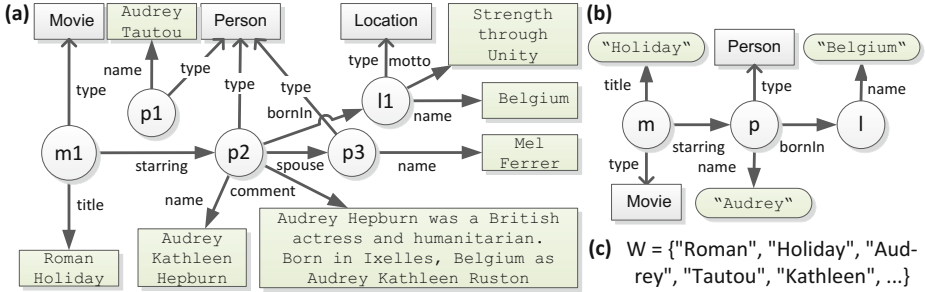
RDF data contains descriptions of entities, with each description being a set of *triples*. Many RDF descriptions feature *textual data*: On the one hand, structured RDF data often comprises text via predicates such as `comment` or `description`. On the other hand, unstructured Web documents are frequently annotated with structured data (e.g., via RDFa or Microformats).

**Hybrid Queries.** Such text-rich RDF descriptions are often queried with queries, which comprise predicates that match *structured data* as well as *words* in text data (*hybrid queries*). Consider the following example (cf. Fig. 1-a/b):

```
SELECT * WHERE {
  ?m ex:title ?title .      ?p ex:name ?name .      ?l ex:name ?name2 .
  ?m ex:starring ?p .      ?p ex:bornIn ?l .
  ?m a Movie .             ?p a Person .
  FILTER (contains(?title,"Holiday") && contains(?name,"Audrey") &&
          contains(?name2,"Belgium")) }
```

---

<sup>\*</sup> This work was supported by the European Union through project XLike (FP7-ICT-2011-288342).



**Fig. 1.** (a) RDF graph about “Audrey Hepburn” and her movie “Roman Holiday”. (b) Hybrid query graph asking for movies with title “Holiday” and starring a person with name “Audrey”, who was born in “Belgium”. (c) Vocabulary  $W$  comprising all words from attributes values in the data graph.

Hybrid queries are highly relevant for RDF stores with SPARQL fulltext extension, e.g., LARQ<sup>1</sup> or Virtuoso<sup>2</sup>. In fact, every store that supports `FILTER` clauses on texts faces hybrid queries.

**Selectivity Estimation.** For finding an optimal query plan, RDF stores rely on *selectivity estimates* to approximate the result size of a query (fragment) [19]. Various selectivity estimation techniques have been proposed for relational data [1,7,9,17,18,21] as well as for RDF data [10,16,19,20,23]. These techniques summarize the data via a *data synopsis* (e.g., histograms, join synopses, tuple-graph synopses, or probabilistic relational models (PRM)) to capture data correlations/statistics. Based on these synopses, different estimation approaches have been proposed to efficiently approximate the query’s selectivity.

However, when applying state-of-the-art selectivity estimation techniques for hybrid queries *effectiveness and efficiency issues* (I.1 and I.2) arise:

(I.1) Effectiveness Issues. Queries over RDF data typically comprise a large number of joins. Thus, a data synopsis must capture statistics for multiple (joined) query patterns in order to allow effective estimates. Recent work on selectivity estimation for RDF data captured multiple patterns either via computing statistics for frequent (star/chain) join patterns [10,16], via heuristics [19,20], or by conditional probabilities [23]. These strategies work well for structured query patterns (i.e., class, relation, and attributes), because the number of structured elements (i.e., class, relation, and attributes) is usually *small and independent of the instance data size*. However, in the presence of textual data and hybrid queries, capturing multiple patterns is much harder, since the number of words is oftentimes very large. For instance, the DBLP dataset (used in our evaluation) has 25 million words vs. 56 structure elements.

*Example.* In Fig. 1-a, there can be many entities of type `Person` (i.e., bindings for `?p a Person`), while only few entities have a `name` “Audrey”. So, in order

<sup>1</sup> <http://jena.sourceforge.net/ARQ/lucene-arq.html>

<sup>2</sup> <http://virtuoso.openlinksw.com>

to estimate the # bindings for  $?p$  in Fig. 1-b, a synopsis has to capture statistics for any word associated (via *name*) with *Person* entities.

Moreover, the number of words is strongly dependent on the data and text size, respectively. So, with growing data sets, statistics become increasingly complex and space consuming. To address this, all previous works summarized the words via a small synopsis, which is constructed either by hashing [19], heuristics [20], categorization [16], discretization [10], or via string synopses [23]. Unfortunately, such coarse-grained synopses result in an “information loss”. That is, oftentimes heuristics must be employed to estimate selectivities of query keywords (e.g., “Audrey”), which leads to severe miss-estimations.

(I.2) Efficiency Issues. All previous works [10,16,19,20,23], aim at constructing a *query-independent* data synopsis at offline time. In fact, previous approaches directly use this offline constructed data synopsis to estimate the query selectivity. Thus, a large synopsis would influence the estimation efficiency. In order to guarantee an efficient selectivity estimation at runtime, existing approaches either construct only *small* synopses [10,16,19,23] or purely rely on heuristics for selectivity estimation [20].

**Contributions.** (1) In this paper, we propose a novel approach (*TopGuess*), which utilizes relational topic models as data synopsis. Such synopses are well-suited to summarize text data, because they provide statistics for the *complete vocabulary* of words by means of topics. So, no “information loss” can occur due to coarse-grained synopses. Furthermore, correlations between structured query patterns (e.g., `?m ex:starring ?p` and `?m rdf:type Movie`, see Fig. 1-b) can also be captured. Thus, we have an effective and holistic synopsis for hybrid queries (effectiveness issues, I.1). The *TopGuess* approach also constructs a query-independent data synopsis at offline time. However, in contrast to previous approaches, we do not directly use this large synopsis at runtime. Instead, we only employ a small and compact synopsis (Bayesian network), which is constructed specifically for the current query (efficiency issues, I.2).

(2) We provide a theoretical analysis: *TopGuess* achieves a linear space complexity w.r.t. text data size (cf. Thm. 1). Further, *TopGuess* has an estimation time complexity that is independent of the synopsis size (given the number of topics), cf. Thm. 4.

(3) We conducted experiments on real-world data: *TopGuess* could improve the effectiveness by up to 88% – without sacrificing runtime performance.

**Outline.** First, in Sect. 2 we outline preliminaries. We introduce the *TopGuess* approach in Sect. 3. In Sect. 4, we present evaluation results, before we outline related work in Sect. 5. We conclude with Sect. 6.

## 2 Preliminaries

**Data and Query Model.** We use RDF as data model:

**Definition 1 (RDF Graph).** Let  $\ell_a$  ( $\ell_r$ ) be a set of attribute (relation) labels. A RDF graph is a directed labeled graph  $\mathcal{G} = (V, E, \ell_a, \ell_r)$ , with nodes  $V =$

$V_E \uplus V_A \uplus V_C$  where  $V_E$  are entity nodes,  $V_A$  are attribute value nodes, and  $V_C$  are class nodes. Edges (so-called triples)  $E = E_R \uplus E_A$  and type are a disjoint union of relation edges  $E_R$  and attribute edges  $E_A$ . Relation edges connect entity nodes:  $\langle s, r, o \rangle \in E_R$ , with  $s, o \in V_E$  and  $r \in \ell_r$ . Attribute edges connect an entity with an attribute value:  $\langle s, a, o \rangle \in E_A$ , with  $s \in V_E, o \in V_A$  and  $a \in \ell_a$ . Triple  $\langle s, \text{type}, o \rangle \in E$  models that entity  $s \in V_E$  belongs to class  $o \in V_C$ .

We conceive an attribute value in  $V_A$  as a *bag-of-words*. Further, let a vocabulary  $W$  comprise all such words. That is,  $W$  is derived from words in attribute values: for each triple  $\langle s, p, o \rangle \in E_A$  we add all words in  $o$  to  $W$ . See also Fig. 1.

Conjunctive queries resemble the basic graph pattern (BGP) feature of SPARQL. In this work, we use hybrid queries:

**Definition 2 (Hybrid Query).** A query  $Q$  is a directed graph  $G_Q = (V_Q, E_Q)$ , with  $V_Q = V_{Q_V} \uplus V_{Q_C} \uplus V_{Q_K}$ ,  $V_{Q_V}$  as variables,  $V_{Q_C}$  as constants, and  $V_{Q_K}$  as keywords. Edges  $E_Q$  are called predicates: (1) Class predicates  $\langle s, \text{type}, o \rangle$ , with  $s \in V_{Q_V}, o \in V_{Q_C}$ . (2) Relation predicates  $\langle s, r, o \rangle$ , with  $s \in V_{Q_V}, o \in V_{Q_C} \uplus V_{Q_V}$ , and  $r \in \ell_r$ . (3) String predicates  $\langle s, a, o \rangle$ , with  $s \in V_{Q_V}, o \in V_{Q_K}$ , and  $a \in \ell_a$ .

Fig. 1-b shows an example query. Query semantics follow those for BGPs. That is, results are subgraphs of the data graph, which match all query predicates. The only difference is due to keyword nodes: a value node  $o \in V_A$  matches a keyword  $w \in V_{Q_K}$ , if the bag-of-words from  $o$  contains word  $w$ .

We rely on two data synopses: topic models and Bayesian networks (BNs):

**Topic Models.** Topic models assume that texts are mixtures of “hidden” topics, where a topic is a probability distribution over words. These topics are abstract clusters of words – formed according to word co-occurrences. More formally, a text collection can be represented by  $k$  topics  $\mathcal{T} = \{t_1, \dots, t_k\}$ , where  $W$  is the vocabulary (see above definition) and each topic  $t \in \mathcal{T}$  is a multinomial distribution of words in  $W$ :  $\mathcal{P}(w | t) = \beta_{tw}$  and  $\sum_{w \in W} \beta_{tw} = 1$ .

*Example.* Three topics are depicted in Fig. 2-c:  $\mathcal{T} = \{t_1, t_2, t_3\}$ . Every topic  $t$  assigns a probability (represented by vector  $\beta_t$ ) to each word in the vocabulary. Probabilities in  $\beta_t$  indicate the importance of words within topic  $t$ . For instance, “Belgium” is most important for topic  $t_3$  ( $\beta_{t_3w} = 0.014$ ), cf. Fig. 2-c.

**Bayesian Networks.** A Bayesian network (BN) is a directed graphical model, which compactly represents a joint probability distribution via its structure and parameters [12]. The structure is a directed acyclic graph, where nodes stand for random variables and edges represent dependencies. Given a node  $X_i$  and its parents  $\mathbf{Pa}(X_i) = \{X_j, \dots, X_k\}$ ,  $X_i$  depends on  $\mathbf{Pa}(X_i)$ , but is *conditionally independent* of all non-ancestor random variables (given  $\mathbf{Pa}(X_i)$ ).

BN parameters are given by conditional probability distributions (CPDs). That is, each random variable  $X_i$  is associated with a CPD capturing the conditional probability  $\mathcal{P}(X_i | \mathbf{Pa}(X_i))$ . The joint distribution  $\mathcal{P}(X_1, \dots, X_n)$  can be estimated via the chain rule [12]:  $\mathcal{P}(X_1, \dots, X_n) \approx \prod_i \mathcal{P}(X_i | \mathbf{Pa}(X_i))$ .

*Example.* A BN is shown in Fig. 3. Nodes such as  $X_m$  and  $X_{\text{holiday}}$  stand for random variables. Edges stand for dependencies between those nodes. For instance, the edge  $X_m \rightarrow X_{\text{holiday}}$  denotes a dependency between the parent,

$X_m$ , and the child,  $X_{holiday}$ . In fact, given its parent,  $X_{holiday}$  is conditionally independent of all non-ancestor variables, e.g.,  $X_p$ . Every node has a CDP. For example,  $X_{holiday}$  has a CDP for  $\mathcal{P}(X_{holiday} | X_m)$ , cf. Fig. 3-b.

**Problem.** Given a hybrid query  $Q$ , we aim at a result size estimation function  $\mathcal{F}(Q)$  as [9]:  $\mathcal{F}(Q) \approx \mathcal{R}(Q) \cdot \mathcal{P}(Q)$ .

Let  $\mathcal{R}$  be a function  $\mathcal{R} : Q \rightarrow \mathbb{N}$  that provides an upper bound cardinality for a result set for query  $Q$ . Further, let  $\mathcal{P}$  be a *probabilistic component*, which assigns a probability to query  $Q$  that models whether  $Q$ 's result set is non-empty.

$\mathcal{R}(Q)$  can be easily computed as product over ‘‘class cardinalities’’ of  $Q$  [9]. That is, for each variable  $v \in V_{Q_V}$  we bound the number of its bindings,  $R(v)$ , as number of entities belonging  $v$ 's class  $c$ :  $|\{s | \langle s, type, c \rangle \in E\}|$ . If  $v$  has no class, we use the number of all entities (i.e.,  $|V_E|$ ) as a bound. Finally,  $\mathcal{R}(Q) = \prod_v R(v)$ .

In the remainder of the paper, we provide an effective (I.1) and efficient (I.2) instantiation of the probabilistic component  $\mathcal{P}$ .

(a)	<table style="border-collapse: collapse;"> <tr> <td></td> <td style="text-align: center;"><math>t_1</math></td> <td style="text-align: center;"><math>t_2</math></td> <td style="text-align: center;"><math>t_3</math></td> </tr> <tr> <td style="border-right: 1px solid black; padding-right: 5px;"><math>\lambda_{movie}</math></td> <td style="border: 1px solid black; text-align: center;">3</td> <td style="border: 1px solid black; text-align: center;">0</td> <td style="border: 1px solid black; text-align: center;">1</td> </tr> <tr> <td style="border-right: 1px solid black; padding-right: 5px;"><math>\lambda_{person}</math></td> <td style="border: 1px solid black; text-align: center;">0</td> <td style="border: 1px solid black; text-align: center;">5</td> <td style="border: 1px solid black; text-align: center;">2</td> </tr> </table>		$t_1$	$t_2$	$t_3$	$\lambda_{movie}$	3	0	1	$\lambda_{person}$	0	5	2	(b)	<table style="border-collapse: collapse;"> <tr> <td></td> <td style="text-align: center;"><math>t_1</math></td> <td style="text-align: center;"><math>t_2</math></td> <td style="text-align: center;"><math>t_3</math></td> </tr> <tr> <td style="border-right: 1px solid black; padding-right: 5px;"><math>\omega_{starring}^{t_1}</math></td> <td style="border: 1px solid black; text-align: center;">0</td> <td style="border: 1px solid black; text-align: center;">7</td> <td style="border: 1px solid black; text-align: center;">2</td> </tr> <tr> <td style="border-right: 1px solid black; padding-right: 5px;"><math>\omega_{starring}^{t_2}</math></td> <td style="border: 1px solid black; text-align: center;">0</td> <td style="border: 1px solid black; text-align: center;">1</td> <td style="border: 1px solid black; text-align: center;">0</td> </tr> <tr> <td style="border-right: 1px solid black; padding-right: 5px;"><math>\omega_{starring}^{t_3}</math></td> <td style="border: 1px solid black; text-align: center;">1</td> <td style="border: 1px solid black; text-align: center;">3</td> <td style="border: 1px solid black; text-align: center;">2</td> </tr> </table>		$t_1$	$t_2$	$t_3$	$\omega_{starring}^{t_1}$	0	7	2	$\omega_{starring}^{t_2}$	0	1	0	$\omega_{starring}^{t_3}$	1	3	2																																												
	$t_1$	$t_2$	$t_3$																																																																								
$\lambda_{movie}$	3	0	1																																																																								
$\lambda_{person}$	0	5	2																																																																								
	$t_1$	$t_2$	$t_3$																																																																								
$\omega_{starring}^{t_1}$	0	7	2																																																																								
$\omega_{starring}^{t_2}$	0	1	0																																																																								
$\omega_{starring}^{t_3}$	1	3	2																																																																								
(c)	<table style="border-collapse: collapse; width: 100%;"> <tr> <td style="border-bottom: 1px solid black; padding: 2px;"><math>W</math></td> <td style="border-bottom: 1px solid black; padding: 2px;"><math>\beta_1</math></td> <td style="border-bottom: 1px solid black; padding: 2px;"><math>W</math></td> <td style="border-bottom: 1px solid black; padding: 2px;"><math>\beta_2</math></td> <td style="border-bottom: 1px solid black; padding: 2px;"><math>W</math></td> <td style="border-bottom: 1px solid black; padding: 2px;"><math>\beta_3</math></td> </tr> <tr> <td style="padding: 2px;">film</td> <td style="padding: 2px;">0.024</td> <td style="padding: 2px;">born</td> <td style="padding: 2px;">0.027</td> <td style="padding: 2px;">city</td> <td style="padding: 2px;">0.025</td> </tr> <tr> <td style="padding: 2px;">play</td> <td style="padding: 2px;">0.023</td> <td style="padding: 2px;">woman</td> <td style="padding: 2px;">0.026</td> <td style="padding: 2px;">location</td> <td style="padding: 2px;">0.024</td> </tr> <tr> <td style="padding: 2px;">...</td> <td style="padding: 2px;">...</td> <td style="padding: 2px;">...</td> <td style="padding: 2px;">...</td> <td style="padding: 2px;">...</td> <td style="padding: 2px;">...</td> </tr> <tr> <td style="padding: 2px;">...</td> <td style="padding: 2px;">...</td> <td style="padding: 2px;">audrey</td> <td style="padding: 2px;">0.013</td> <td style="padding: 2px;">belgium</td> <td style="padding: 2px;">0.014</td> </tr> <tr> <td style="padding: 2px;">holiday</td> <td style="padding: 2px;">0.011</td> <td style="padding: 2px;">hepburn</td> <td style="padding: 2px;">0.012</td> <td style="padding: 2px;">...</td> <td style="padding: 2px;">...</td> </tr> <tr> <td style="padding: 2px;">roman</td> <td style="padding: 2px;">0.010</td> <td style="padding: 2px;">...</td> <td style="padding: 2px;">...</td> <td style="padding: 2px;">...</td> <td style="padding: 2px;">...</td> </tr> <tr> <td style="padding: 2px;">...</td> <td style="padding: 2px;">...</td> <td style="padding: 2px;">belgium</td> <td style="padding: 2px;">0.009</td> <td style="padding: 2px;">holiday</td> <td style="padding: 2px;">0.004</td> </tr> <tr> <td style="padding: 2px;">hepburn</td> <td style="padding: 2px;">0.004</td> <td style="padding: 2px;">...</td> <td style="padding: 2px;">...</td> <td style="padding: 2px;">...</td> <td style="padding: 2px;">...</td> </tr> <tr> <td style="padding: 2px;">...</td> <td style="padding: 2px;">...</td> <td style="padding: 2px;">holiday</td> <td style="padding: 2px;">0.002</td> <td style="padding: 2px;">hepburn</td> <td style="padding: 2px;">0.002</td> </tr> <tr> <td style="padding: 2px;">belgium</td> <td style="padding: 2px;">0.001</td> <td style="padding: 2px;">...</td> <td style="padding: 2px;">...</td> <td style="padding: 2px;">...</td> <td style="padding: 2px;">...</td> </tr> <tr> <td style="text-align: center; padding: 2px;"><math>t_1</math></td> <td></td> <td style="text-align: center; padding: 2px;"><math>t_2</math></td> <td></td> <td style="text-align: center; padding: 2px;"><math>t_3</math></td> <td></td> </tr> </table>			$W$	$\beta_1$	$W$	$\beta_2$	$W$	$\beta_3$	film	0.024	born	0.027	city	0.025	play	0.023	woman	0.026	location	0.024	...	...	...	...	...	...	...	...	audrey	0.013	belgium	0.014	holiday	0.011	hepburn	0.012	...	...	roman	0.010	...	...	...	...	...	...	belgium	0.009	holiday	0.004	hepburn	0.004	...	...	...	...	...	...	holiday	0.002	hepburn	0.002	belgium	0.001	...	...	...	...	$t_1$		$t_2$		$t_3$	
$W$	$\beta_1$	$W$	$\beta_2$	$W$	$\beta_3$																																																																						
film	0.024	born	0.027	city	0.025																																																																						
play	0.023	woman	0.026	location	0.024																																																																						
...	...	...	...	...	...																																																																						
...	...	audrey	0.013	belgium	0.014																																																																						
holiday	0.011	hepburn	0.012	...	...																																																																						
roman	0.010	...	...	...	...																																																																						
...	...	belgium	0.009	holiday	0.004																																																																						
hepburn	0.004	...	...	...	...																																																																						
...	...	holiday	0.002	hepburn	0.002																																																																						
belgium	0.001	...	...	...	...																																																																						
$t_1$		$t_2$		$t_3$																																																																							

**Fig. 2.** Synopsis parameters (stored on disk): (a)  $\lambda_{movie}$  and  $\lambda_{person}$  parameter for three topics. (b)  $\omega$  matrix for **starring** relation and three topics – rows (columns) represent source (target) topics of the relation. (c) Words from the vocabulary  $W$  and their corresponding probabilities for each topic,  $\beta_t$ . Note, data is taken from Fig. 1-a/c.

### 3 TopGuess

Targeting the effectiveness (I.1) and efficiency (I.2) issues of existing works w.r.t. hybrid queries (cf. Sect. 1), we now introduce our novel TopGuess approach.

More precisely, we present an uniform data synopsis based on relational topic models in Sect. 3.1 (I.1), and show in Thm. 1 that this synopsis has a linear space complexity w.r.t. vocabulary  $W$  (I.2). Further, we introduce a probabilistic component  $\mathcal{P}$  in Sect. 3.2 and 3.3, and show in Thm. 4 selectivity computation complexity to be independent of the synopsis size (I.2).

Note, the topic model (data synopsis) is learned at offline time and may be stored on disk. At runtime, we construct a small BN for each given query – reflecting our data synopsis as well as query characteristics via topic mixtures.

#### 3.1 Relational Topic Models as Data Synopsis

**Synopsis Parameters.** For an effective synopsis over text-rich RDF data, TopGuess exploits *relational topic models* [2,4,14,25]. These topic models

summarize the data by means of one *uniform synopsis* – considering structured and text data. More precisely, our synopsis comprises of two parts:

(1) First, the TopGuess synopsis captures text data in a low-dimensional representation via a set of  $k$  topics  $\mathcal{T} = \{t_1, \dots, t_k\}$ .

*Example.* Fig. 2-c groups words from the vocabulary  $W = \{\text{“Roman”}, \text{“Holiday”}, \dots\}$ , cf. Fig. 1-c, via three topics,  $\mathcal{T} = \{t_1, t_2, t_3\}$ . This way, text data in Fig. 1-a, e.g., associated via attribute *comment*, is compactly represented.

The number of topics is dictated by the data characteristics. In particular, previous works allow to learn the optimal number of topics [8]. By means of this compact summary, TopGuess achieves a linear space complexity linear w.r.t. a vocabulary, see Thm. 1.

(2) Second, the TopGuess synopsis captures *correlations between those topics and structured data*. For our query model, we rely on two correlation parameters for selectivity estimation:  $\lambda$  and  $\omega$ . Note, for other kinds of queries, further types of correlation parameters may be considered.

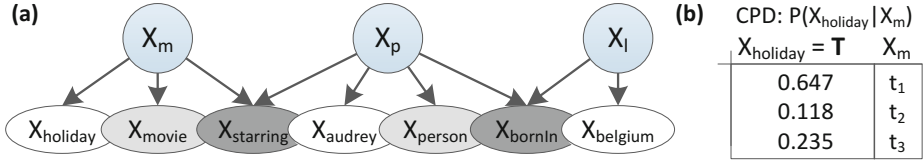
- Class-Topic Parameter  $\lambda$ . We capture correlations between a class  $c \in V_C$  and topics in  $\mathcal{T}$  via a vector  $\lambda_c$ , where each vector element,  $\lambda_c(t)$ , represents the weight between class  $c$  and topic  $t$ . A higher weight indicates that a class is more correlated with a topic.

*Example.* Fig. 2-a shows the two class-topic parameters for the classes *Movie* and *Person*:  $\lambda_{\text{movie}}$  and  $\lambda_{\text{person}}$ . Both indicate correlations w.r.t. topics  $\mathcal{T} = \{t_1, t_2, t_3\}$ . For instance,  $\lambda_{\text{movie}}$  states that class *Movie* is highly correlated with topic  $t_1$ , has some correlation with  $t_3$ , and has no correlation with  $t_2$ .

- Relation-Topic Parameter  $\omega$ . We measure correlations between a relation  $r$  and the topics in  $\mathcal{T}$  via a matrix  $\omega_r$ . Since relation  $r$  is observed “between” two entities, say  $(s, r, o)$ , the topic of its subject  $s$  and its object  $o$  is considered. Given  $k$  topics, matrix  $\omega_r$  has  $k \times k$  dimensions and entries such that: if entity  $s$  is associated with topic  $t_i$  and entity  $o$  has topic  $t_j$ , the weight of observing a relation  $r$  “between”  $s$  and  $o$  is given by the entry  $\langle i, j \rangle$ , denoted as  $\omega_r(t_i, t_j)$ . Note, TopGuess features a matrix  $\omega$  for each relation.

*Example.* Fig. 2-b depicts the relation-topic parameter for the *starring* relation:  $\omega_{\text{starring}}$ . According to  $\omega_{\text{starring}}$ , *starring* is most often observed (weight 7) if its subject (object) contains words from topic  $t_1$  ( $t_2$ ).

**Parameter Learning.** For training above parameters, we do not restrict TopGuess to a single topic model. Instead, different approaches can be used. For instance, classical topic models such as LDA [3] may be employed to learn the first part, i.e., word/topic probabilities, cf. Fig. 2-c. Then, correlations between those topics and classes/relations must be obtained. For this, topic models have been extended to consider structured data, so-called relational topic models [2,4,14,25]. Most notably, a recent approach, the *Topical Relational Model* (TRM) [2], trains topics as well as class/relation correlations simultaneously from RDF data. We used a TRM as data synopsis in our experiments.



**Fig. 3.** (a) Query-specific BN for the query in Fig.1-b. It contains three topical variables (color: blue, e.g.,  $X_m$ ), two class predicate variables (color: light gray, e.g.,  $X_{movie}$ ), two relation predicate variables (color: dark gray, e.g.,  $X_{starring}$ ), and three string predicate variables (color: white, e.g.,  $X_{holiday}$ ). Observed variables (color: white/light gray/dark gray) are independent from each other and only dependent on hidden topical random variables (color: blue) – as dictated by Def. 4. (b) CPD for random variable  $X_{holiday}$ , cf. parameters in Fig. 2-c.

**Discussion.** The TopGuess synopsis comes with key advantages: First, in contrast to existing work [23], we do not need separate synopses for structured and text data. This way, we may learn correlations in a uniform manner.

Moreover, TopGuess parameters are not required to be loaded in memory. This is a crucial advantage over state-of-the-art selectivity estimation systems [9,21,23], as memory is commonly a limited resource. So, TopGuess can utilize the *complete vocabulary*  $W$  for learning word/topic probabilities  $\beta$ .

Last, as empirically validated by our experiments, correlations between topics and structured data *suffice for an accurate selectivity estimation*. Since even a small number of topics can capture these correlations, our synopsis does not grow exponentially in its vocabulary size.

In fact, we can show that a topic-based data synopsis has linear space complexity w.r.t. its vocabulary:

**Theorem 1 (Synopsis Space Complexity).** *Given  $k$  topics, a vocabulary  $W$ , classes  $V_C$ , and relations  $\ell_r$ , TopGuess has a storage space complexity in  $O(|W| \cdot k + |V_C| \cdot k + |\ell_r| \cdot k^2)$ .*

*Proof.* See our report for a proof [22] ■

### 3.2 Probabilistic Component: Query-Specific BN

In this section, we exploit the synopsis parameters for an efficient probabilistic component (I.2, Sect. 1). For this, we first construct a small, query-specific BN and afterwards compute its joint probability in Sect. 3.3. Both steps are done at runtime. In contrast to [9,21,23], all synopsis parameters may be kept on disk, while only the query-specific BN is loaded into memory.

To construct a BN specifically for a query  $Q$ , we capture every query predicate in  $Q$  via a random variable: For a class predicate,  $\langle s, type, c \rangle$ , and relation predicate,  $\langle s, r, o \rangle$ , we create a binary random variable  $X_c$  and  $X_r$ . Similarly, for a string predicate,  $\langle s, a, w \rangle$ , we introduce a binary random variable  $X_w$ . Most importantly, we assume that each variable  $v$  in  $Q$  belongs to one or more topics in  $\mathcal{T}$ . So, we model variable  $v$  via a *topical random variable*,  $X_v$ . More formally,  $X_v$  has a multinomial distribution over the topics:

**Definition 3 (Topical Random Variable).** For a set of topics  $\mathcal{T}$ , a query  $Q$ , and its variable  $v \in V_{Q_V}$ , the random variable  $X_v$  is a multinomial topical random variable for  $v$ , with topics  $\mathcal{T}$  as sample space.

Based on topical random variables, we perceive query variables as topic mixtures. Thus,  $X_v$  captures query variable  $v$ 's "relatedness" to every topic. In the following, we denote the set of all string, class, relation, and topical random variables for query  $Q$  as  $\mathbf{X}_w$ ,  $\mathbf{X}_c$ ,  $\mathbf{X}_r$ , and  $\mathbf{X}_v$ .

We create a simple BN structure by means of a fixed structure assumption:

**Definition 4 (Topical Dependence Assumption).** Given a class/string predicate  $\langle v, *, * \rangle$ , the corresponding variable  $X$  depends only the topical random variable  $X_v$ . Given a relation predicate  $\langle v, *, y \rangle$ , the corresponding variable  $X$  depends only on two topical random variables:  $X_v$  and  $X_y$ .

The topical dependence assumption lies at the core of the TopGuess approach. It considers that query predicate probabilities depend on (and are governed by) the topics of their associated topical random variables. Further, the assumption allows us to model the query probability,  $\mathcal{P}(Q)$ , via a tree-shaped BN.

*Example.* Fig. 3-a depicts a BN for the query in Fig. 1-b. Adhering to Def. 4, each topical variable ( $X_m$ ,  $X_p$ , and  $X_l$ ) forms a small tree of dependent random variables. For instance, random variable  $X_{\text{holiday}}$  is only dependent on its topical variable,  $X_m$ . In fact, given  $X_m$ ,  $X_{\text{holiday}}$  is conditionally independent of all other variables, e.g.,  $X_{\text{audrey}}$ . This way, topic mixtures of  $X_m$ ,  $X_p$ , and  $X_l$  govern the overall query probability,  $\mathcal{P}(Q)$ .

Last, note that the topical dependence assumption leads to a valid BN:

**Theorem 2.** A query-specific BN constructed according to Def. 4 is acyclic.

*Proof.* See [22] for a proof ■

### 3.3 Probabilistic Component: Query Probability Computation

Having formed the BN structure for a given query  $Q$ , we may compute the query probability,  $\mathcal{P}(Q)$ , via the chain rule (CR) [12]:

$$\mathcal{P}(Q) = \mathcal{P} \left( \bigwedge \mathbf{X}_w = \mathbf{T} \quad \bigwedge \mathbf{X}_c = \mathbf{T} \quad \bigwedge \mathbf{X}_r = \mathbf{T} \right) \quad (1a)$$

$$\begin{aligned} &\stackrel{\text{CR}}{\approx} \prod_{\langle v, a, w \rangle \in Q} \mathcal{P}(X_w = \mathbf{T} \mid X_v) \cdot \prod_{\langle v, \text{type}, c \rangle \in Q} \mathcal{P}(X_c = \mathbf{T} \mid X_v) \\ &\cdot \prod_{\langle v, r, y \rangle \in Q} \mathcal{P}(X_r = \mathbf{T} \mid X_v, X_y) \end{aligned} \quad (1b)$$

In order to solve Eq. 1, we require a CPD for each random variable, cf. Fig. 3-b. We rely on TopGuess parameters as well as distributions of topical random variables to approximate these CPDs. As topical variables  $\mathbf{X}_v$  are hidden, we learn their distributions from observed random variables ( $\mathbf{X}_w$ ,  $\mathbf{X}_c$ ,  $\mathbf{X}_r$ ).



In the following, we first discuss CPD estimation for observed random variables, given topical random variables (topic distributions). Subsequently, we present learning of topic distributions for hidden topical variables.

**Query Predicate Probabilities.** Probabilities for query predicates are influenced by their associated topical random variables and their TopGuess parameters. In other words, we may compute the CPDs for  $\mathbf{X}_w$ ,  $\mathbf{X}_c$ , and  $\mathbf{X}_r$  using topic distributions of topical variables and probabilities estimated by the corresponding  $\beta$ ,  $\lambda$ , or  $\omega$  parameter:

(1) Class Predicates. Adhering to the topical dependence assumption, the probability of observing a class,  $\mathcal{P}(X_c = \mathbf{T})$ , is only dependent on its topical variable  $X_v$ . We use the class-topic parameter  $\lambda$  to obtain the weight  $\lambda_c(t)$ , which indicates the correlation between topic  $t$  and class  $c$ :

$$\mathcal{P}(X_c = \mathbf{T} \mid X_v, \lambda) = \sum_{t \in \mathcal{T}} \mathcal{P}(X_v = t) \frac{\lambda_c(t)}{\sum_{t' \in \mathcal{T}} \lambda_c(t')}$$

*Example.* Fig. 3-a shows two random variables,  $X_{movie}$  and  $X_{person}$ , which dependent on their topical variables  $X_m$  and  $X_p$ . For computing  $\mathcal{P}(X_{movie} = \mathbf{T})$  and  $\mathcal{P}(X_{person} = \mathbf{T})$ , the parameters  $\lambda_{movie}$  and  $\lambda_{person}$  are used, cf. Fig. 2-a. Assuming  $\mathcal{P}(X_m = t_1) = 0.6$ ,  $\mathcal{P}(X_m = t_2) = 0.1$ , and  $\mathcal{P}(X_m = t_3) = 0.3$ , we get:  $\mathcal{P}(X_{movie} = \mathbf{T}) = 0.6 \cdot 0.75 + 0.1 \cdot 0 + 0.3 \cdot 0.25 = 0.525$ .

(2) Relation Predicates. A relation predicate  $\langle v, r, y \rangle$  connects two query variables, which have the two topical variables  $X_v$  and  $X_y$ . Random variable  $X_r$  solely depends on the topics of  $X_v$  and  $X_y$ . The correlation between relation  $r$  and these topics is given by the relation-topic parameter  $\omega_r$ :

$$\mathcal{P}(X_r = \mathbf{T} \mid X_v, X_y, \omega_r) = \sum_{t, t' \in \mathcal{T}} \frac{\mathcal{P}(X_v = t) \omega_r(t, t') \mathcal{P}(X_y = t')}{\sum_{t'', t''' \in \mathcal{T}} \omega_r(t'', t''')}$$

*Example.* In Fig. 3-a, we have the variables  $X_{starring}$  and  $X_{bornIn}$  – both dependent on two topical variables. For instance,  $X_{starring}$  depends on  $X_m$  and  $X_p$ .  $\mathcal{P}(X_{starring} = \mathbf{T})$  is estimated via matrix  $\omega_{starring}$ , cf. Fig. 2-b.

(3) String Predicates. For each string predicate  $\langle v, a, w \rangle$ , there is a random variable  $X_w$ . The word-topic parameter  $\beta_{tw}$  represents the probability of observing word  $w$  given topic  $t$ . Thus,  $\mathcal{P}(X_w = \mathbf{T})$  is calculated as probability of observing  $w$ , given the topics of  $v$ 's topical variable,  $X_v$ :

$$\mathcal{P}(X_w = \mathbf{T} \mid X_v, \beta_{1:K}) = \sum_{t \in \mathcal{T}} \mathcal{P}(X_v = t) \frac{\beta_{tw}}{\sum_{t' \in \mathcal{T}} \beta_{t'w}}$$

*Example.* Fig. 3-a depicts three random variables for string predicates. Given  $\mathcal{P}(X_m)$  as in the above example, the probability for “holiday” is (cf. Fig. 3-b):

$$\mathcal{P}(X_{holiday} = \mathbf{T}) = 0.6 \cdot \frac{0.011}{0.017} + 0.1 \cdot \frac{0.002}{0.017} + 0.3 \cdot \frac{0.004}{0.017} = 0.47$$

**Learning Topic Distributions.** Finally, we wish to estimate topic distributions for the hidden topical variables based on  $\mathbf{X}_w$ ,  $\mathbf{X}_c$ , and  $\mathbf{X}_r$ . We aim at finding a topic distribution for every topical variable, so that the query probability in Eq. 1 is maximized. Thus, this optimal topic distribution directly gives

us  $\mathcal{P}(Q)$ . Let  $\theta_{vt}$  denote a set of topic parameters for topical random variable  $X_v$ . Further, let  $\theta = \{\theta_{vt} \mid v \in V_{Q_V}, t \in \mathcal{T}\}$  be the set of all parameters  $\theta_{vt}$ . Then, we search for parameter  $\theta$  that maximizes the log-likelihood of Eq. 1:

$$\arg \max_{\theta} L(\theta : \mathbf{X}_w, \mathbf{X}_c, \mathbf{X}_r) \tag{2}$$

where  $L(\theta : \mathbf{X}_w, \mathbf{X}_c, \mathbf{X}_r)$  is the log-likelihood defined as:

$$\begin{aligned} L(\theta : \mathbf{X}_w, \mathbf{X}_c, \mathbf{X}_r) &= \mathcal{P}(\mathbf{X}_w, \mathbf{X}_c, \mathbf{X}_r \mid \theta, \beta, \omega, \lambda) \\ &= \sum_v \sum_{X_w \in \mathbf{X}_w^v} \log \mathcal{P}(X_w \mid X_v, \beta) + \sum_v \sum_{X_c \in \mathbf{X}_c^v} \log \mathcal{P}(X_c \mid X_v, \lambda) \\ &\quad + \sum_{v,y} \sum_{X_r \in \mathbf{X}_r^{v,y}} \log \mathcal{P}(X_r \mid X_v, X_y, \omega) \end{aligned}$$

where  $\mathbf{X}_w^v$  and  $\mathbf{X}_c^v$  is the set of all string/class random variables having  $X_v$  as parent.  $\mathbf{X}_r^{v,y}$  is the set of all relation random variables with parents  $X_v$  and  $X_y$ .

We use gradient ascent optimization to learn the parameter  $\theta$ . First, we parametrize each  $\mathcal{P}(X_v = t)$  with  $\theta_{vt}$  such that

$$\mathcal{P}(X_v = t) = \frac{e^{\theta_{vt}}}{\sum_{t' \in \mathcal{T}} e^{\theta_{vt'}}$$

to obtain a valid probability distribution over the topics. Obtaining the gradient requires dealing with the log of the sum over the topics of each topical variable. Therefore, we make use of theorem [12]:

**Theorem 3.** *Given a BN and  $\mathcal{D} = \{\mathbf{o}[1], \dots, \mathbf{o}[M]\}$  as a partially observed dataset. Let  $X$  be a variable in that BN with  $\mathbf{Pa}(X)$  as its parents. Then:*

$$\frac{\partial L(\theta : \mathcal{D})}{\partial \mathcal{P}(x \mid \mathbf{pa})} = \frac{1}{\mathcal{P}(x \mid \mathbf{pa})} \sum_{m=1}^M \mathcal{P}(x, \mathbf{pa} \mid \mathbf{o}[m], \theta),$$

This provides the necessary form of the gradient. Now, the gradient of the log-likelihood w.r.t. parameter  $\theta_{vt}$  is:

$$\frac{\partial L(\theta : \mathbf{X}_w, \mathbf{X}_c, \mathbf{X}_r)}{\partial \theta_{vt}} = \underbrace{\frac{\partial L(\theta : \mathbf{X}_w, \mathbf{X}_c, \mathbf{X}_r)}{\partial \mathcal{P}(X_v = t)}}_{(i)} \cdot \underbrace{\frac{\partial \mathcal{P}(X_v = t)}{\partial \theta_{vt}}}_{(ii)} \tag{3}$$

The (i)-part of the gradient, Eq. 3, may be obtained via Theorem 3:

$$\begin{aligned} \frac{\partial L(\theta : \mathbf{X}_w, \mathbf{X}_c, \mathbf{X}_r)}{\partial \mathcal{P}(X_v = t)} &= \frac{1}{\mathcal{P}(X_v = t)} \left( \sum_{X_w \in \mathbf{X}_w^v} \mathcal{P}(X_w = t \mid X_w, \beta) \right. \\ &\quad \left. + \sum_{X_c \in \mathbf{X}_c^v} \mathcal{P}(X_c = t \mid X_c, \lambda) + \sum_y \sum_{X_r \in \mathbf{X}_r^{v,y}} \mathcal{P}(X_r = t \mid X_r, X_y, \omega) \right) \end{aligned}$$

Using Bayes rule we get:

$$\begin{aligned} \frac{\partial L(\theta : \mathbf{X}_w, \mathbf{X}_c, \mathbf{X}_r)}{\partial \mathcal{P}(X_v = t)} &= \sum_{X_w \in \mathbf{X}_w^v} \frac{\mathcal{P}(X_v = t) \mathcal{P}(X_w | \beta, t)}{\sum_{t'} \mathcal{P}(X_v = t') \mathcal{P}(X_w | \beta, t')} \\ &+ \sum_{X_c \in \mathbf{X}_c^v} \frac{\mathcal{P}(X_v = t) \mathcal{P}(X_c | \lambda, t)}{\sum_{t'} \mathcal{P}(X_v = t') \mathcal{P}(X_c | \lambda, t')} \\ &+ \sum_y \sum_{X_r \in \mathbf{X}_r^{v,y}} \frac{\mathcal{P}(X_v = t) \sum_{t'} \mathcal{P}(X_r | X_y, \omega, t')}{\sum_{t''} \mathcal{P}(X_v = t'') \sum_{t'''} \mathcal{P}(X_r | X_y, \omega, t''')} \end{aligned}$$

Finally, the (ii)-part of the gradient in Eq. 3 is given by:

$$\frac{\partial \mathcal{P}(X_v = t)}{\partial \theta_{tv}} = \frac{e^{\theta_{tv}} \sum_{t'-t} e^{\theta_{t'v}}}{(\sum_{t'} e^{\theta_{t'v}})^2}$$

**Time Complexity.** Query probability estimation has a complexity bound:

**Theorem 4 (Time Complexity for  $\mathcal{P}(Q)$  Estimation).** *Given  $k$  topics and a query  $Q$ , the time for computing  $\mathcal{P}(Q)$  in Eq. 1 is in  $O(\psi \cdot |Q| \cdot k)$ , with  $\psi$  as number of iterations needed for optimization and  $|Q|$  as # predicates in  $Q$ .*

*Proof.* A proof is given in [22] ■

Note,  $\psi$  is determined by the specific algorithm used for optimization. So, overall complexity for computing  $\mathcal{P}(Q)$  is independent of the synopsis size given the topics.

## 4 Evaluation

We conducted experiments to analyze the effectiveness (I.1) and efficiency (I.2) of TopGuess. Overall, our results are very promising: we achieved up to 88% more accurate selectivity estimates, while runtime was comparable to the baselines. Further, in contrast to the baselines, we noted TopGuess’s runtime behavior to be much less influenced by the synopsis size – thereby confirming Thm. 4.

### 4.1 Evaluation Setting

**Systems.** We employ two categories of baselines: (1) String predicates are combined with structured predicates via an *independence assumption*: IND baseline. That is, the selectivity of string predicates and structured predicates is estimated using two separate synopses: a string synopsis (explained below) and a synopsis for structured RDF data based on histograms [10]. Obtained probabilities are combined in a greedy fashion while assuming independence. (2) We reuse our previous work on BNs for text-rich data graphs [23]: BN baseline. Here, all query predicates are captured uniformly via a single BN. To handle string predicates, we employ *n-gram string synopses* [24]

A n-gram synopsis summarizes the vocabulary by “omitting” certain n-grams. Thus, a synopsis represents a subset of all possible n-grams occurring in the

**Table 1.** Data synopsis memory/disk space in MB

	IMDB		DBLP	
	BN/IND	TopGuess	BN/IND	TopGuess
Mem.	{2, 4, 20, 40}	$\leq 0.1$	{2, 4, 20, 40}	$\leq 0.1$
Disk	0	281.7	0	229.9

data. A simplistic strategy is to choose random n-gram *samples* from the data. Another approach is to construct a *top-k* n-gram synopsis. For this, n-grams are extracted from the data together with their counts. Then, the  $k$  most frequent n-grams are included in the synopsis. Further, a *stratified bloom filter* (SBF) synopsis has been proposed [24], which uses bloom filters as a heuristic map that projects n-grams to their counts. Note, we refer to omitted n-grams as “missing”. The probability for missing n-grams cannot be estimated with a probabilistic framework, as such strings are not included in a sample space. So, a string predicate featuring a missing n-gram is assumed to be independent from the remainder of the query. Its probability is computed via a heuristic. We employ the *leftbackoff* strategy, which finds the longest known n-gram that is the pre- or postfix of the missing n-gram. Then, the probability of the missing n-gram is approximated based on statistics for its pre-/postfix [24].

Combining string synopses with the two categories of baselines yields six systems:  $\text{IND}_{\text{sample}}$ ,  $\text{IND}_{\text{top-k}}$ , and  $\text{IND}_{\text{SBF}}$  rely on the independence assumption, while  $\text{BN}_{\text{sample}}$ ,  $\text{BN}_{\text{top-k}}$ , and  $\text{BN}_{\text{SBF}}$  represent BN approaches.

**Data.** We employ two real-world RDF datasets: DBLP [15] and IMDB [6]. From both datasets we have large vocabularies: 25,540,172 (DBLP) and 7,841,347 (IMDB) words. Note, while DBLP and IMDB feature text-rich attributes, they differ in their overall amount of text. On average an attribute in DBLP contains 2.6 words, with a variance of 2.1 words. In contrast, IMDB attributes contain 5.1 words, with a variance of 95.6 words. Moreover, we observed during learning of the BN baseline that there are more data correlations in IMDB than in DBLP. We expect correlations between query predicates to have a strong influence on the system effectiveness.

**Queries.** We used IMDB [6] and DBLP [15] keyword search benchmarks: We generated 54 DBLP queries from [15]. Further, we constructed 46 queries for IMDB based on queries in [6]. We omitted 4 queries from [6], as they could not be translated to conjunctive queries. Overall, our load features 100 queries with: 0-4 relation, 1-7 string, 1-4 class predicates, and 2-11 predicates in total. Further query statistics and a complete query listing can be found in [22].

**Synopsis Size.** We employ baselines with varying synopsis size. For this, we varied # words captured by the string synopsis. The top-k and sample synopsis contained # words  $\in \{0.5K, 1K, 5K, 10K\}$ . The SBF string synopsis captured  $\{2.5K, 5K, 25K, 50K\}$  words for each attribute. Note, SBF systems featured most keywords occurring in our query load. Different string synopsis sizes translated to a memory consumption of baselines  $\in \{2, 4, 20, 40\}$  MB. IND and

BN baselines load their synopsis into main memory. In contrast, TopGuess keeps a large topic model at disk and constructs a small, query-specific BN in memory at runtime ( $\leq 100$  KBytes). Table 1 depicts further details.

**Implementation and Offline Learning.** For IND and BN baselines, we started by constructing their string synopses. Each synopsis was learned in  $\leq 1$ h.

Then, we constructed BN systems based on [23]. That is, we capture words and structured data elements using random variables and learn correlations between them, thereby forming a BN structure. For efficient selectivity estimation the network is reduced to a “lightweight” model capturing solely the most important correlations. Then, we calculate model parameters (CPDs) based on frequency counts. For IND systems, we do not need the model structure and merely keep the marginalized BN parameters. Structure and parameter learning combined took up to 3h. To compute query selectivities the BN systems need inferencing strategies. For this, we used a Junction tree algorithm [23].

TopGuess exploits an “off-the-shelf” TRM from [2]. The number of topics is an important factor – determining which correlations are discovered. We experimented with a varying number of topics  $\in [10, 100]$ . We found 50 topics are sufficient to capture all strong correlations in our datasets. The TopGuess learning took up to 5h and its parameters were stored on hard disk, cf. Table 1. At query time, we employed a greedy gradient ascent algorithm for learning the topic distributions. To avoid local maxima, we used up to 10 random restarts.

We implemented all systems in Java 6. Experiments were run on a Linux server with: 2 Intel Xeon CPUs at 2.33GHz, 16 GB memory assigned to the JVM, and a RAID10 with IBM SAS 10K rpm disks. Before each query execution we cleared OS caches. Presented values are averages over five runs.

## 4.2 Selectivity Estimation Effectiveness

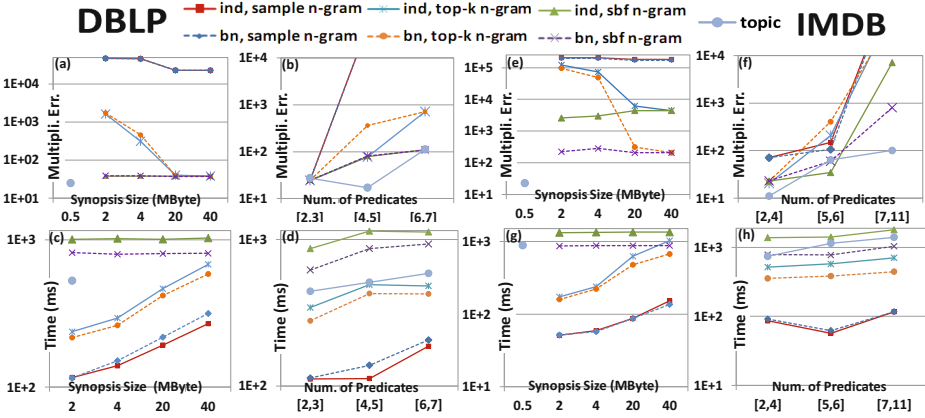
We employ the multiplicative error metric ( $me$ ) [7] for measuring effectiveness:

$$me(Q) = \frac{\max\{\mathcal{F}_e(Q), \mathcal{F}_a(Q)\}}{\min\{\mathcal{F}_e(Q), \mathcal{F}_a(Q)\}}$$

with  $\mathcal{F}_e(Q)$  and  $\mathcal{F}_a(Q)$  as exact and approximated selectivity. Intuitively,  $me(Q)$  is the factor to which  $\mathcal{F}_a(Q)$  under-/overestimates  $\mathcal{F}_e(Q)$ .

**Overall Results.** Figs. 4-a/e (b/f) show the multiplicative error vs. synopsis size (# predicates) for DBLP and IMDB. Baseline system effectiveness strongly varies with their synopsis size. In particular, for small synopses  $\leq 20$  MB IND and BN performed poorly. We explain this with missing words in their string synopses, which led to heuristics being used for probability computation. In simple terms, IND and BN systems traded synopsis space for estimation accuracy.

TopGuess, on the other hand, did not suffer from this issue. All its parameters (cf. Sect. 3.1) could be stored at disk and solely the query-specific BN was loaded at runtime. Thus, TopGuess could exploit very fine-grained probabilities and omitted any kind of heuristic. We observed that TopGuess reduced the error of the best BN system,  $BN_{\text{SBF}}$ , by 88% (33%) for IMDB (DBLP). Further, we outperformed the best IND system,  $IND_{\text{SBF}}$ , by 99% (35%) on IMDB (DBLP).



**Fig. 4.** Effectiveness: (a)+(b) for DBLP and (e)+(f) for IMDB. Efficiency: (c)+(d) for DBLP and (g)+(h) for IMDB. Y-axes are given in logarithmic scale.

**Synopsis Size.** Figs. 4-a/e show estimation errors w.r.t. in-memory synopsis size. An important observation is that the synopsis size is a key factor for effectiveness. Top-k and sample-based string synopsis systems were strongly affected by their (string) synopsis size. Given a small synopsis  $\leq 4$  MB, we observed that top-k/sample-based systems performed poorly. Here, many relevant query keywords were missed, leading to inaccurate heuristics being applied. With increasing synopsis size  $\in [4, 20]$  MB, the performance of top-k approaches converged to the most accurate baseline (SBF-based systems). For instance given 4 MB space, the  $\text{BN}_{\text{top-k}}$  approach preformed 95% worse than  $\text{BN}_{\text{SBF}}$  on IMDB, but only 33% worse given 20 MB. Further, we noted SBF-based approaches to perform fairly stable. We explain this with SBF systems using bloom filters as an effective summary. Such systems were able to capture most query keywords. Thus, few heuristic-based estimations were necessary. However, SBF-based systems also have a limited memory space and must eventually omit words.

In contrast, we observed TopGuess to use  $\leq 0.1$  MB memory for IMDB as well DBLP. We explain this extremely compact BN with: (1) TopGuess has a network size, which is bound by the query size. (2) The BN only contains random variables that either are binary or have a sample space, which is bounded by the number of topics. For example, over the DBLP query load TopGuess needed on average 40 KB. Yet, TopGuess resolves the issue of missing words completely: the TopGuess parameters (stored on disk) capture all words in the vocabulary. At runtime, TopGuess retrieves the necessary statistics for a particular query and constructs its query-specific BN. This way, TopGuess achieved up to by 88% (33%) better results on IMDB (DBLP) than the best baselines.

Overall, we can conclude that estimation effectiveness is driven by accurate string predicate probabilities. Thus, there is a strong need for a data synopsis allowing for extensive word/text data statistics.

**Correlations.** We found system performances to vary for IMDB and DBLP. For the IMDB dataset,  $\text{BN}_{\text{SBF}}$  could reduce errors of the  $\text{IND}_{\text{SBF}}$  approach by up

to 93%. On the other hand, for DBLP improvements were much smaller. These differences are due to the varying degree of correlations in our two datasets. While learning the BNs for BN, we observed less correlations in DBLP than in IMDB. For instance, for DBLP queries with string predicates `name` and `label`, we noted no significant correlations. Thus, the probabilities obtained from BN systems were almost identical to the ones from IND.

In contrast, even for the less correlated dataset DBLP, TopGuess outperforms the best baselines,  $\text{IND}_{\text{SBF}}$  and  $\text{BN}_{\text{SBF}}$ , by 35% and 33%. We explain this our a fine-grained, query-specific BN. More precisely, we observed that BN approaches exploited data correlations, which were observed in the data graph. However, TopGuess captured even minor correlations via its topic mixtures at query time – learned for each query individually.

**Query Size.** We depict the multiplicative error vs. # query predicates in Figs. 4-b/f. As expected, estimation errors increase for all systems in # predicates. For our baselines, we explain this behavior with: (1) Given an increasing # predicates, the chances of missing a query keyword increase. (2) When missing a single query keyword, the error is “propagated” throughout the computation.

However, while the TopGuess approach also led to more misestimates for larger queries, the degree of this increase was smaller. For instance, considering IMDB queries with 7-11 predicates, we could observe that TopGuess performs much more stable than BN or IND baselines, cf. Fig. 4-f.

### 4.3 Selectivity Estimation Efficiency

We now analyze the estimation efficiency vs. synopsis size (# query predicates), cf. Figs. 4-c/g (d/h). For TopGuess, the reported times comprise parameter loading, BN construction, and topic learning. For BN and IND, the times represent only selectivity computation, i.e., no model learning or parameter loading.

**Overall Results.** Considering BN and IND systems, we saw that their string synopsis was a key performance factor. Intuitively, the more words were missed, the “simpler” and the more efficient these systems became. However, such gains came at the expense of effectiveness: the fastest baseline system,  $\text{IND}_{\text{sample}}$ , also computed the least accurate selectivity estimates.

Comparing the two systems with the best effectiveness, TopGuess and  $\text{BN}_{\text{SBF}}$ , TopGuess led to a better performance by up to 45%. Unfortunately, in comparison to top-k systems, TopGuess resulted in a performance decrease of 40%. We explain these drawbacks with the time-consuming disk I/O, which was needed for loading the statistics. However, while BN and IND clearly outperformed TopGuess w.r.t. small synopses  $\leq 4$  MB, TopGuess results are comparable for synopses  $\geq 20$  MB. We expect such effects to be more drastic for “large” BN/IND synopses  $\gg 100$  MB. So, TopGuess guarantees a much more “stable” behavior.

**Synopsis Size.** Figs. 4-c/g show time vs. synopsis size. For the baselines, we saw a correlation between synopsis size and runtime behavior: While BN and IND reach a high efficiency for synopses  $\leq 4$  MB, their performance decreases rapidly for synopses  $\geq 20$  MB. We explain this effect with the larger CPDs, which led to longer probability estimation times. We observed SBF-based approaches to

be less driven by their synopsis size. This is because their computational costs are mainly determined by bloom filters. In contrast, TopGuess did not suffer from this issue at all. That is, for a given query, TopGuess only loads/processes statistics necessary for that query. All others statistics are kept on disk.

**Query Size.** All systems had increasing estimation times in query size, cf. Figs. 4-d/h. This is because each additional query predicate translated to more probability computations. However, as TopGuess exploits a compact query-specific BN, we expected it’s performance to be less influenced by query size. To confirm this, we compared the standard deviation of the estimation time w.r.t. a varying # query predicates. For instance, the standard deviations was 82.48 ms (213.48 ms) for TopGuess (BN). The low deviation for TopGuess indicates that its probability estimation times varied less than those from BN systems.

## 5 Related Work

For selectivity estimation on *structured data*, existing works exploit various data synopses, e.g., join samples [1], graph synopses [18], or graphical models [9,21,23]. In particular, those synopses have been applied for selectivity estimation of structured SPARQL/BGP queries on RDF data, e.g., [7,10,16,19].

Unfortunately, such synopses can not effectively capture statistics for combinations of words and structured data elements. In order to guarantee a manageable synopsis size, various summarization techniques are commonly applied on words, e.g., hashing [19], heuristics [20], categorization [16], discretization [10], or string synopses [23] (I.1, Sect. 1). Moreover, since the offline constructed data synopsis is directly used for the selectivity estimation at runtime, existing works must either keep the data synopsis small [10,16,19,23] or solely employ heuristics [20] (I.2, Sect. 1). In contrast, our TopGuess approach overcomes both issues (I.1 and I.2) by relying on relational topic models as data synopsis and a query-specific BN for selectivity estimation.

With regard to selectivity estimation on *text data*, language models and other machine learning techniques have been employed [5,11,13,24]. More specifically, some works aim at substring or fuzzy string matching [5,13], while other approaches target “extraction” operators, e.g., dictionary-based operators [24]. However, such works do not consider correlations among multiple string predicates or correlations between string predicates and structured query predicates.

## 6 Conclusion

We proposed a holistic approach (TopGuess) for selectivity estimation of hybrid queries. We showed space and time complexity bounds for TopGuess. Further, we conducted empirical studies on real-world data and achieved strong effectiveness improvements, while not requiring additional runtime.

As a future work, we plan to extend TopGuess in order to become a more generic selectivity estimation approach for RDF data and BGP queries, respectively. For this, we replace the topic distributions in our data synopsis with different application-specific probability distributions, e.g., a continuous distribution for estimating the selectivity of range queries over numerical data.



## References

1. Acharya, S., Gibbons, P., Poosala, V., Ramaswamy, S.: Join synopses for approximate query answering. In: SIGMOD (1999)
2. Bicer, V., Tran, T., Ma, Y., Studer, R.: TRM - Learning Dependencies between Text and Structure with Topical Relational Models. In: Alani, H., et al. (eds.) ISWC 2013, Part I. LNCS, vol. 8218, pp. 1–16. Springer, Heidelberg (2013)
3. Blei, D., Ng, A., Jordan, M.: Latent dirichlet allocation. *The Journal of Machine Learning Research* 3, 993–1022 (2003)
4. Chang, J., Blei, D.: Relational Topic Models for Document Networks. In: AISTats (2009)
5. Chaudhuri, S., Ganti, V., Gravano, L.: Selectivity Estimation for String Predicates: Overcoming the Underestimation Problem. In: ICDE (2004)
6. Coffman, J., Weaver, A.C.: A framework for evaluating database keyword search strategies. In: CIKM (2010)
7. Deshpande, A., Garofalakis, M.N., Rastogi, R.: Independence is Good: Dependency-Based Histogram Synopses for High-Dimensional Data. In: SIGMOD (2001)
8. Doshi, F., Miller, K., Gael, J.V., Teh, Y.W.: Variational Inference for the Indian Buffet Process. *JMLR* 5, 137–144 (2009)
9. Getoor, L., Taskar, B., Koller, D.: Selectivity estimation using probabilistic models. In: SIGMOD (2001)
10. Huang, H., Liu, C.: Estimating Selectivity for Joined RDF Triple Patterns. In: CIKM (2011)
11. Jin, L., Li, C.: Selectivity Estimation for Fuzzy String Predicates in Large Data Sets. In: VLDB (2005)
12. Koller, D., Friedman, N.: Probabilistic graphical models. MIT Press (2009)
13. Lee, H., Ng, R.T., Shim, K.: Extending Q-Grams to Estimate Selectivity of String Matching with Low Edit Distance. In: VLDB (2007)
14. Liu, Y., Niculescu-Mizil, A., Gryc, W.: Topic-link LDA: Joint Models of Topic and Author Community. In: ICML (2009)
15. Luo, Y., Wang, W., Lin, X., Zhou, X., Wang, J., Li, K.: SPARK2: Top-k Keyword Query in Relational Databases. *TKDE* 23(12), 1763–1780 (2011)
16. Neumann, T., Moerkotte, G.: Characteristic sets: Accurate cardinality estimation for RDF queries with multiple joins. In: ICDE (2011)
17. Poosala, V., Haas, P., Ioannidis, Y., Shekita, E.: Improved histograms for selectivity estimation of range predicates. In: SIGMOD (1996)
18. Spiegel, J., Polyzotis, N.: Graph-based synopses for relational selectivity estimation. In: SIGMOD (2006)
19. Stocker, M., Seaborne, A., Bernstein, A., Kiefer, C., Reynolds, D.: SPARQL Basic Graph Pattern Optimization Using Selectivity Estimation. In: WWW (2008)
20. Tsialiamanis, P., Sidirouros, L., Fundulaki, I., Christophides, V., Boncz, P.: Heuristics-based Query Optimisation for SPARQL. In: EDBT (2012)
21. Tzoumas, K., Deshpande, A., Jensen, C.S.: Lightweight Graphical Models for Selectivity Estimation Without Independence Assumptions. In: PVLDB (2011)
22. Wagner, A., Bicer, V., Tran, D.T.: Topic-based Selectivity Estimation for Text-Rich Data Graphs, <http://www.aifb.kit.edu/web/Techreport3039>
23. Wagner, A., Bicer, V., Tran, T.D.: Selectivity estimation for hybrid queries over text-rich data graphs. In: EDBT (2013)
24. Wang, D.Z., Wei, L., Li, Y., Reiss, F., Vaithyanathan, S.: Selectivity estimation for extraction operators over text data. In: ICDE (2011)
25. Zhang, L., et al.: Multirelational Topic Models. In: ICDM (2009)