

# Thurstonian model for object ranking after determinization of object using query aware approach

<sup>#1</sup>Ms. Sagare Priyanka B., <sup>#2</sup>Mrs. Bhokare Uma M.

<sup>#1</sup>Assistant Professor,  
Yashoda Technical Campus, Wadhe Satara.

<sup>#2</sup>Assistant Professor,  
Yashoda Technical Campus, Wadhe Satara.



## ABSTRACT

Use of cloud computing as well as web based applications is growing and hence this allows end users to frequently save their data in different present web applications. Frequently, user data is automatically created through different methods such as data enrichment methods, signal processing methods etc. and stored it into the web applications. These automatically generated contents are frequently resulted as ambiguous as well as objects with probabilistic attributes. In this project we are presenting the efficient, accurate and ranking based method for solving the determinization problem. The main aim of this project is to present technique that generates the deterministic probabilistic data representation which optimizes end application quality over deterministic data. The determinization problem will be explored in selection queries and triggers data processing tasks. In existing methods, the objects retrieved without ranking, therefore to overcome this problem we introduced ranking method based on Thurstonian Model. In this method, ranking of objects will be done according the scores of objects. For performance evaluation, performance measure the in terms of F-score, precision and recall metrics.

**Keywords:** Determinization, uncertain data, data quality, branch and bound algorithm, query workload.

## ARTICLE INFO

### Article History

Received: 25<sup>th</sup> March 2017

Received in revised form :  
25<sup>th</sup> March 2017

Accepted: 25<sup>th</sup> March 2017

### Published online :

4<sup>th</sup> May 2017

## I. INTRODUCTION

Today's Era of web, computer network individuals uses web based applications and Cloud computing also, users conjointly have to be compelled to store their data on their varied net applications. From the spread of signal processing, web applications analysis or improvement methods of various kinds of information are generated mechanically before stored on within web applications. For example, modern cameras come with features like vision analysis with tags like landscape or portrait, indoors or outdoors, scenery etc. Trendy pic cameras have microphones feature with a group of tags for users to talk sentences that which are processed by a speech recognizer[2]. To web applications the pic are often streamed in time victimization wireless property like Flickr. [1].

Within the past the determinization downside has not been examined extensively. The foremost connected analysis efforts unit, that explore the thanks to give settled answers to an issue (e.g. conjunctive selection question ) over probabilistic data , our main target is to confirm

determinization of the change in data it store on in traditional databases which optimizes the expected cost of queries among all the determinized illustration, To the determinization problem many approaches are designed. Two basic ways that one is Top-1 and other each one techniques, whereby we have a tendency to elect the foremost probable worth / attribute with non zero value generally.

Different context information into web applications has several issues relating to mechanically generated content and it can't be determined from its context and it's going to lead to objects with probabilistic values of attributes. for e.g., vision analysis could lead to tags with possibilities and like ,recognizer for automatic generated speech (ASR) could turn out Associate in Nursing best-N list or a mixed network of utterances[2][5]. Probabilistic information should be "determinized" before start hold on in gift Web applications. We have a tendency to ask drawback the matter of mapping probabilistic information into their respected settled illustration cause problem with determinization. for e.g., Associate in application for Nursing finish which supports

alerts and triggers on mechanically derived described content. Such samples Associate in Nursing end-application which includes publish/subscribe systems like GoogleAlert, whereby users can notify their respective subscriptions with keywords (e.g., "Ind vs Aus Cricket") and predicates among over information (here., information is of video category). Currently concerning/contemplate/take into account a video about Ind vs. Aus cricket that's may be revealed on YouTube. The video can have a collection of tags that were extracted victimisation either machine-controlled vision processing, techniques for data extraction applied for transcribed speech. Such machine-controlled tools could turn out tags with possibilities (e.g., "Ind": 0.5, "vs": 0.8, "Aus": 0.2, "Cricket": zero.1), comparing the actual tags of the video may be of "Ind", "Aus" and "Cricket": . The determination process ought to video association with applicable tags for such subscribers WHO square measure extremely inquisitive about the video (i.e., whose subscribed list contains the words "Ind vs. Aus Cricket") don't seem to verified whereas alternative square measure not swamped by unsuitable information[9][10]. Thus, within the example higher than, the determinization process should minimize metrics like false +ves and false -ves that result from a determinized illustration of information[1].

## II. LITERATURE SURVEY

In this section we discussed the review of different references for Uncertain Objects Determinization and Ranking using Query Aware.

In [2], D. V. Kalashnikov, J. Xu S. Mehrotra, N. Venkatasubramanian represents the, over each real speech recognizer's output as well as artificial data sets empirical analysis has been conducted. To boost quality of speech recognition it explores within the type of co-occurrence between pictures tags however linguistics information is exploited. For mutually exclusive image information this can be not appropriate.

In [3], J. Li and J. Wang, represents, in real-time the annotation of a picture will be provided. It provides content based mostly image retrieval and storage of classified pictures.

In [4], C. Wangand, L. Zhang, F. Jing, and H. Zhang represents the, effectiveness of the proposed methodology as per the Experimental results on each non-Web images of Corel dataset and web images of photograph forum sites demonstrate.

In [5], J. Li and A. Deshpande represents for computing the consensus answers for various distance metrics, we tend to obtain polynomial optimal time or approximation algorithms. Here XOR model is used for computation.

In [7], R. Cheng, J. Chen, and X. Xie represents the, for the purpose of cleaning data purposes a metric is used for a probabilistic database and explores a higher quality metric.

In [11]. B. Sigurbjörnsson and R. V. Zwol, represents the, with completely different levels of comprehensiveness of firstly given tagging, for effectively which gives suggestion for relative tags for a set of photos.

## III. PROPOSED APPROACH FRAMEWORK AND DESIGN

### Problem Definition

Under the important time environments, uncertain objects might occur in several kinds of data sources with goal that such unresolved objects might have quite numerous descriptions over the net applications and cloud computing domains. thus it is required to check determinization of object before uploading. Recently completely several methods for solving determinization problem are developing. Top-1 and all techniques are ways of solving such problem ,for choosing the all attainable values or most probable values severally. however these types of ways are agnostic to the end-application oftentimes resulted into suboptimal results. to boost the results of such ways, an additional recent technique introduced that planned the efficient determinization algorithms those are quicker as compared to enumeration primarily based best solution however achieves nearly identical quality because the best solution. The limitation of this technique is that, it doesn't supporting the objects ranking extraction.

### System Architecture

The system architecture of project is as shown in figure where

Image is taken as input whose features are extracted in feature extraction module where we get idea about all set of tags along with visual features. Those tags are used as input for branch and bound algorithm which gives result as set of tags with optimal result. And then after for rank order results Thurstonian model can be used which is based on mean and deviation calculus also uses precision and recall to check the rank of object.

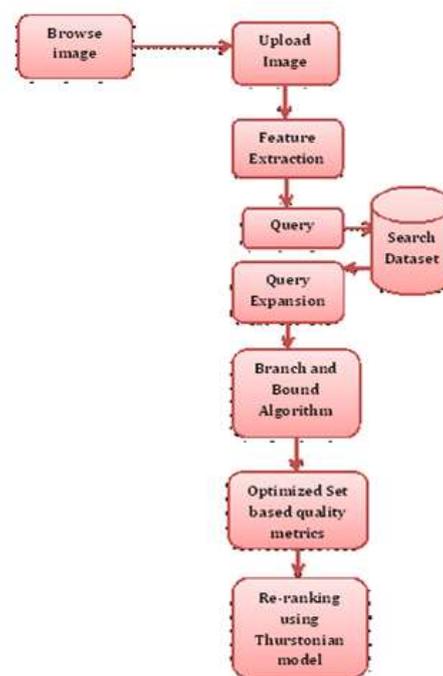


Figure 1: System architecture

## Mathematical model

Let,  $D = \{s, e, N, Y, F, Fme, DD, NDD, \emptyset\}$

//  $D$  is determinization technique.

// This is used for object determinization.

$s$  = start of the system

$e$  = end of the system

$N$  = database containing number of images.

$F$  = It contains the feature vector.

$Fme$  = Algorithm of the system here branch and bound algorithm is used.

$DD$  = Deterministic data. i.e. number of instruction or space complexity.

$NDD$  = Nondeterministic data.

$Y$  = Processing on the input set with text and visual features we get output.

$\emptyset$  = failure or success of system. Success is desired output of the system i.e. determinization of object. Otherwise failure.

we choose one deterministic representation  $AO$  for  $O$ ,

Algorithm for branch and bound can be given as below:

// Here we consider object as image with specific tags as given in image.

Algorithm Branch\_and\_Bound(  $O, A, H, Q, N_{leaf}$  )

```
{
// O is set of object here it is permutation of given tags.
// Q is set of queries.
// Set of leaf nodes. Which affect the complexity.
1. //Set best answer set to empty.
   A ← ∅
2. c* ← ∞
   // Set lower cost to some value..here it is infinity.
3. Create and Assign empty priority queue H.
4. Add root to empty priority queue.
5. Repeat steps 6 to 12 until priority queue H is not empty
6. V ← Remove the best node with minimum cost from priority queue H.
7. // if c* is less than lower bound node can be pruned
   If  $lv \geq c^*$ 
   Continue;
8. if  $lv < c^*$ 
    $c^* \leftarrow lv$ 
    $A^* \leftarrow A(shv)$ 
9. if  $N_{leaf} \leq 0$ 
   break // as algorithm completed
10. if  $lv \neq hv$  then //need branching
11.  $w \leftarrow get\_next\_tag\_with\_min\_cost(v)$ 
12.  $N_{leaf} \leftarrow Branch(v, w, H, N_{leaf})$ 
   // v- node
   // w-tags
   // H-priority queue
   {
13. generate a copy vyes of a node.
14.  $N_{leaf} \leftarrow Update\ node(vyes, w, yes, Qv, N_{leaf})$ 
15. Add vyes to priority queue H.
16. Generate a copy vno of a node v.
17.  $N_{leaf} \leftarrow Update\ node(vno, w, no, Qv, N_{leaf})$ 
18. Add vno to priority queue H.
19. Return  $N_{leaf}$ .
20. Return  $A^*$ 
```

21. }

For a query  $Q$  if we compute expected cost WRT  $O$  is as:

$$E(cost(O, Q)) = \sum_{G \in W} cost(A_O, G, Q) \cdot \mathbb{P}(G = G_O).$$

And  $E(cost(O, Q))$  for a chosen answer set  $A$  can be derived as:

$$E(cost(O, Q)) = E(cost(A, G, Q))$$

Where

$$E(cost(A, G, Q)) = \begin{cases} c_Q^+ \cdot (1 - P_{O \in G_Q}) & \text{if } O \in A_Q; \\ c_Q^- \cdot P_{O \in G_Q} & \text{if } O \notin A_Q. \end{cases}$$

For query workload  $Q$  we have:

$$E(cost(O, Q)) = E(cost(A, G, Q)),$$

Where

$$E(cost(A, G, Q)) = \sum_{Q \in Q} f_Q \cdot E(cost(A, G, Q))$$

Branch and Bound (BB/B&B) is an algorithm design paradigm which is for discrete and combinational optimization based problem also generate real valued problems. A Branch-and-bound algorithm comprises of a systematic inventory of contestant solution by means of state space search: the set of candidate solution is forming a rooted tree.

The node selection algorithm maintains one priority queue  $H$  for selecting a node  $v$  that contains the most accurate sequence  $S_v$  to continue with to compare an upper bound on  $mv$ , it is sufficient to pick one answer sequence  $A$  from  $A_v$  and then set  $hv = E(cost(A, G, Q))$ . We can call such an answer sequence as the upper bound sequence  $S_{hy}$  for node  $v$ . The procedure for choosing  $S_{hy}$  determines the quality of the upper bound. Expected cost is calculated as.

$$E(cost(S, G, Q)) = c_Q^+ \cdot (1 - P_{O \in G_Q}) \cdot P_{O \in S_Q} + c_Q^- \cdot P_{O \in G_Q} \cdot (1 - P_{O \in S_Q}),$$

Where  $P_{oes_Q}$  is the probability of object  $O$  that satisfies query  $Q$  based on given sequence  $S$  which is associated.

To calculate a lower bound on  $mv$ , we can find the result set with minimum expected cost is costs as the lower bound. i.e., for a given node  $v$ .

$$\sum_{Q \in Q} f_Q \cdot \min_{A \in A_v} E(cost(A, G, Q)) \leq \min_{A \in A_v} E(cost(A, G, Q)).$$

Thus,  $Q \in Q$   $f_Q \cdot \min_{A \in A_v} E(cost(A, G, Q))$  this is same as lower bound on  $mv$ , and the task is to compute  $\min_{A \in A_v} E(cost(A, G, Q))$  for each  $Q \in Q$ .

Query-level optimizations can be applied to improve performance of B&B algorithm.

The approximation of F-measure for a defined  $AQ$  is can be calculated as:

$$\hat{E}(F_{\alpha}(A_Q, G_Q)) = \frac{(1 + \alpha) \sum_{O \in A_Q} P_{O \in G_Q}}{|A_Q| + \alpha \sum_{O \in G_Q} P_{O \in G_Q}}$$

Consider one example of image having 3 tags in it, and here we need to determinize those 3 tags for proceeding the algorithm. Algorithm proceeds with Object as an image then set of tags are permuted as there are 3 tags total 8 permutation can be formed which results as count of Nleaf will be 8, out of which best 1 need to be selected, Firstly one root get selected and then based on probability values of lower and upper bound of a node appropriate node will be selected and added to priority queue as an answer set to our query. The best with minimum cost get selected. And again if more than one images with same set is there then ranking algorithm is applied which help to get output with ranked order. For that thurstonian model will be used which is also work on mutually exclusive construct and work on principle of mean and variance calculus.

#### IV. EXPECTED RESULT

Figures given below shows the practical work to be done during implementation. Graphical representation for accuracy and data can be shown as in given figure. Time required for set of transactions is used to determine the performance of system.

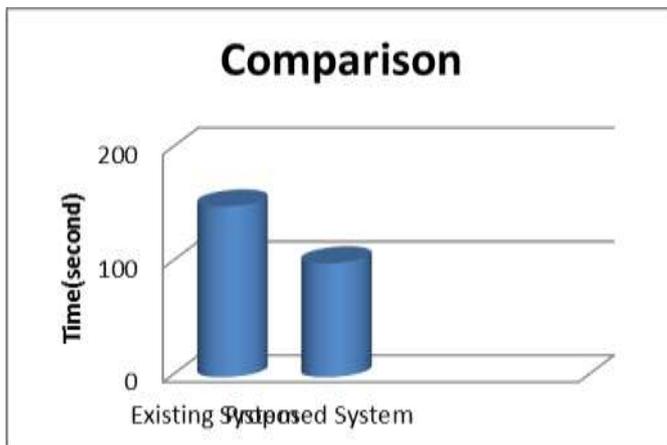


Figure 2: Time Comparison

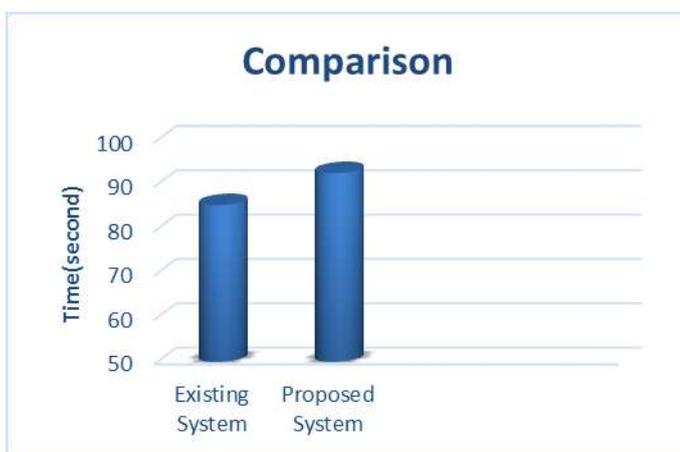


Figure 3: Accuracy Prediction

#### V. CONCLUSION

We projected unsure Objects Determinization and Ranking victimisation question Aware Approach and Thurstonian Model that generates the settled probabilistic knowledge illustration that optimizes finish application quality over settled knowledge. we've got projected economical determinization algorithms that area unit orders of magnitude quicker than the enumeration based mostly optimum answer however achieves identical quality because the optimum answer however achieves identical quality because the optimum answer and additionally retrieved objects during a ranked order.

#### REFERENCES

- [1] Jie Xu, Dmitri V. Kalashnikov, and Sharad Mehrotra, "Query aware determinization of uncertain object," *IEEE Trans. Knowl. Data Eng.*, vol. 27, no. 1, Jan 2015.
- [2] D. V. Kalashnikov, S. Melhotra, J. Xu, and N. Venkatasubramanian, "A semantics-based approach for speech annotation of images," *IEEE Trans. Knowl. Data Eng.*, vol. 23, no. 9, pp. 1373–1387, Sept. 2011.
- [3] J. Li and J. Wang, "Automatic linguistic indexing of pictures by a statistical modeling approach," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 25, no. 9, pp. 1075–1088, Sept. 2003.
- [4] C. Wangand, F. Jing, L. Zhang, and H. Zhang, "Image annotation refinement using random walk with restarts," in *Proc. 14th Annu. ACM Int. Conf. Multimedia*, New York, NY, USA, 2006.
- [5] R. Nuray-Turan, D. V. Kalashnikov, S. Mehrotra, and Y. Yu, "Attribute and object selection queries on objects with probabilistic attributes," *ACM Trans. Database Syst.*, vol. 37, no. 1, Article 3, Feb. 2012.
- [6] J. Li and A. Deshpande, "Consensus answers for queries over probabilistic databases," in *Proc. 28th ACM SIGMOD-SIGACTSIGARTSymp. PODS*, New York, NY, USA, 2009.
- [7] R. Cheng, J. Chen, and X. Xie, "Cleaning uncertain data with quality guarantees," in *Proc. VLDB*, Auckland, New Zealand, 2008.
- [8] S. Bhatia, D. Majumdar, and P. Mitra, "Query suggestions in the absence of query logs," in *Proc. 34th Int. ACM SIGIR*, Beijing, China, 2011.
- [9] D. V. Kalashnikov and S. Mehrotra, "Domain-independent datacleaning via analysis of entity-relationship graph," *ACM Trans. Database Syst.*, vol. 31, no. 2, pp. 716–767, Jun. 2006.
- [10] A. Rae, B. Sigurbjörnsson, and R. V. Zwol, "Improving tag recommendation using social networks," in *Proc. RIAO*, Paris, France, 2010.
- [11] B. Sigurbjörnsson and R. V. Zwol, "Flickr tag recommendation based on collective knowledge," in *Proc. 17th Int. Conf. WWW*, New York, NY, USA, 2008.