

DiVE: Diversifying View Recommendation for Visual Data Exploration

ABSTRACT

To support effective data exploration, there has been a growing interest in developing solutions that can automatically recommend data visualizations that reveal interesting and useful data-driven insights. In such solutions, a large number of possible data visualization views are generated and ranked according to some metric of importance (e.g., a deviation-based metric), then the top-k most important views are recommended. However, one drawback of that approach is that it often recommends similar views, leaving the data analyst with a limited amount of gained insights. To address that limitation, in this work we posit that employing diversification techniques in the process of view recommendation allows eliminating that redundancy and provides a good and concise coverage of the possible insights to be discovered. To that end, we propose a hybrid objective utility function, which captures both the importance, as well as the diversity of the insights revealed by the recommended views. While in principle, traditional diversification methods (e.g., Greedy Construction) provide plausible solutions under our proposed utility function, they suffer from a significantly high query processing cost. In particular, directly applying such methods leads to a “process-first-diversify-next” approach, in which all possible data visualization are generated first via executing a large number of aggregate queries. To address that challenge and minimize the incurred query processing cost, we propose an integrated scheme called *DiVE*, which efficiently selects the top-k recommended view based on our hybrid utility function. Specifically, *DiVE* leverages the properties of both the importance and diversity metrics to prune a large number of query executions without compromising the quality of recommendations. Our experimental evaluation on real datasets shows that DiVE can reduce the query processing cost by up to 40% compared to existing methods.

ACM Reference Format:

. 1997. DiVE: Diversifying View Recommendation for Visual Data Exploration. In *Proceedings of ACM conference (CIKM 2018)*. ACM, New York, NY, USA, Article 4, 13 pages. https://doi.org/10.475/123_4

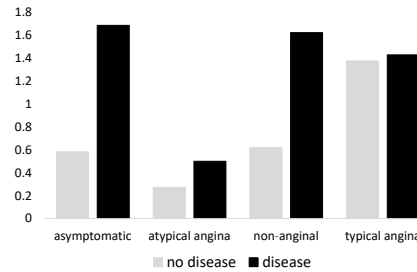
Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

CIKM 2018, October 2018,

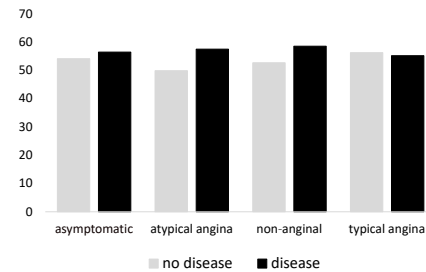
© 2018 Copyright held by the owner/author(s).

ACM ISBN 123-4567-24-567/08/06...\$15.00

https://doi.org/10.475/123_4



(a) Visualization of the avg. oldpeak vs. chest pain types



(b) Visualization of the average age vs. chest pain types

Figure 1: Important vs. less important view.

1 INTRODUCTION

In the recent years, visualization recommendation systems have become an integral part of data exploration systems. The users who are interested in finding some meaningful insights in data have neither time nor patience to manually generate all possible data visualizations. In addition to time, the domain knowledge is another key factor when generating visualizations manually. However, with an exponential growth of available data in various domains, there has been an increase in the number of *Data Enthusiasts*, people with little domain knowledge and technical expertise, looking for interesting trends in data.

For instance, consider a Cleveland heart disease dataset¹, which describes patients with and without a heart disease. A data enthusiast might be interested in conducting some comparison between people with heart disease (disease) and people without heart disease (no disease). Without any prior insights about data, she must manually specify different combinations of attributes, measures and aggregate functions

¹<http://archive.ics.uci.edu/ml/datasets/heart+Disease>

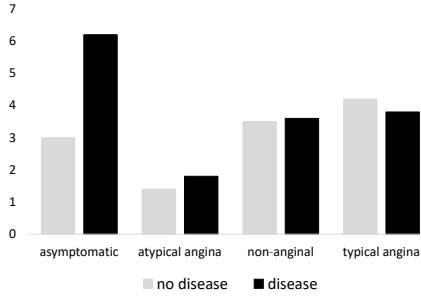


Figure 2: The visualization of maximum oldpeak vs. chest pain types

before finally generating a visualization that reveals some interesting information about the dataset. The user effort and time spent in that process increases exponentially with increase in the number of attributes and measures. Hence, several data-driven visualization recommendation tools have been proposed to reduce the user effort and time during data exploration [5, 15, 16]. The main goal of those recommendation systems is to provide the user with the most important visualizations (top-k views), which are selected from all possible visualizations. The top-k views are selected based on the most important views in the dataset. The importance of a view is defined on the basis of particular criteria. One of the widely used criteria for importance is based on the deviation between the queried subset of data (target view) with the reference subset of data (reference view). The reference subset can be another subset of the dataset, the rest of the dataset, or the whole dataset. The intuition behind deviation based approach is that views that reveal substantially different trends from the reference views are likely to be of higher interest to the user [15, 16].

Consider again the example of the heart disease dataset. Let the target subset be the data of people with heart disease and the reference subset be the data of people without heart disease. As shown in Figure 1a the average oldpeak (pressure of the ST segment) vs. chest pain types is more important view rather than the Figure 1b the average of age vs. chest pain types, due to the large deviation between target view (disease) and reference view (no disease) data. Figure 1a shows people with heart disease, especially who the chest pain types is asymptomatic tend to have much higher oldpeak rather than people without disease. To the contrary, the Figure 1b is potentially less important visualization compared to Figure 1a, even there is a deviation between disease and no disease but the deviation is very small and if it is compared to Figure 1a, it has lower deviation than Figure 1a. Figure 1b shows that there is no significant different in term of the average age that of the people with four types of chest pain.

Although the deviation based visualization recommendation systems automatically provide users with the most important visualizations, it is likely that the views in the top-k

set might be providing redundant information. For instance, Figure 2 provides the information which close to Figure 1a that the people with heart disease tend to have higher oldpeak values. Figure 2 has same attribute and attribute measure to Figure 1a, the only difference is the aggregate function, where Figure 1a uses AVG and Figure 2 uses MAX. Since both views have a high deviation from the reference subset, both will appear in the top-k set. This leads to an important observation that using only importance as the selection criteria may deliver redundant recommended views, which leads to presents not optimal insights.

However, novelty and diversity are one of the fundamental characteristics of any effective recommendation systems [3, 10, 20, 21]. Specifically, it is highly desirable that a visualization recommendation systems provides users with views that are both importance and also provide novel information that has not been revealed by the other views.

Towards designing an effective visualization recommendation systems that promotes both importance and novelty in recommended views, in this work, we propose an integrated approach called *DiVE*. In particular, *DiVE* aims to generate top-k visualizations that balance the tradeoff between importance and diversity. The main contributions of this paper are summarized as follows:

- We formulate the problem of evaluating recommended views that are both importance and diverse.
- We define a similarity measure to capture the distance between two visualizations.
- We present a hybrid objective function to balance the tradeoff between importance and diversity when ranking the visualizations.
- We propose the novel *DiVE* scheme, that employs various algorithms to evaluate the recommended visualizations based on the hybrid ranking/objective function.
- We present optimization techniques that leverage the hybrid objective function to substantially reduce the computational costs.
- We conduct an extensive experimental evaluation on real datasets, which compare the performance of various algorithms and illustrate the benefits achieved by *DiVE* both in terms of effectiveness and efficiency.

The rest of the paper is organized as follows: in Section II, we formulate the top-k diverse visualization problems and our related work; we present our proposed scheme *DiVE* in Section III; the experimental evaluation is reported in in Section IV and we conclude in Section V.

2 PRELIMINARIES AND RELATED WORK

2.1 Recommendation of Aggregate Views

Several recent research efforts have been directed to the challenging task of recommending aggregate views that reveal interesting data-driven insights (e.g., []). As in previous work, we assume a similar model, in which a visual data exploration session starts with an analyst submitting a query Q on a multi-dimensional database D_B . Essentially, Q selects a subset D_Q

from D_B by specifying a query predicate T . Hence, Q is simply defined as:

$Q: \text{SELECT } * \text{ FROM } D_B \text{ WHERE } T;$

Ideally, the analysts would like to generate some aggregate views (e.g., bar charts or scatter plots) that unearth some valuable insights from the selected data subset D_Q . However, achieving that goal is only possible if the analyst knows exactly what to look for! That is, if they know the parameters, which specify aggregate views that lead to those valuable insights (e.g., aggregate functions, grouping attributes, etc.). Meanwhile, such parameters only become clear in “hindsight” after spending long time exploring the underlying database. Hence, the goal of existing work, such as [5, 7, 15–17], is to automatically recommend such aggregate views.

To specify and recommend such views, as in previous work, we consider a multi-dimensional database D_B , which consists of a set of dimensional attributes \mathbb{A} and a set of measure attributes \mathbb{M} . Also, let \mathbb{F} be a set of possible aggregate functions over measure attributes, such as **COUNT**, **AVG**, **SUM**, **MIN** and **MAX**. Hence, specifying different combinations of dimension and measure attributes along with various aggregate functions, generates a set of possible views \mathbb{V} over the selected dataset D_Q . For instance, a possible aggregate view V_i is constituted by a tuple $\langle A_i, M_i, F_i \rangle$, where $A_i \in \mathbb{A}$, $M_i \in \mathbb{M}$, and $F_i \in \mathbb{F}$, and it can be formally defined as:

$V_i: \text{SELECT } A_i, F_i (M_i) \text{ FROM } D_B \text{ WHERE } T \text{ GROUP BY } A_i;$

Clearly, an analyst would be interested in those views that reveal interesting insights. However, manually looking for insights in each view $V_i \in \mathbb{V}$ is a labor-intensive and time-consuming process. For instance, consider again our example in the previous section. In that example, let D_B be the Cleveland Heart Disease data table (i.e., **tb_heart_disease**) and the analyst is selecting the subset of patients with heart disease (i.e., $D_Q = \text{disease}$ subset). Hence, the number of views to explore is equal to: $|\mathbb{V}| = |\mathbb{A}| \times |\mathbb{M}| \times |\mathbb{F}|$, where $|\mathbb{F}|$ is the number of SQL aggregate functions, and $|\mathbb{A}|$ and $|\mathbb{M}|$ are the number of attribute and measures in **tb_heart_disease**, respectively. For that medium-dimensionality dataset, that value of $|\mathbb{V}|$ goes up to 180 views, which is clearly unfeasible for manual exploration. Such challenge motivated multiple research efforts that focused on *automatic* recommendation of views based on some metrics that captures the importance of a recommended view (e.g., [5–7, 12, 15–17]). [Some of those works focus on recommending visualizations to facilitate a particular user intent or task \[6, 12, 12\].](#) For example: [explanations for a certain behavior, finding data anomalies or outliers and correlations among data attributes.](#) Hence, the criteria for ranking the visualizations is driven by the user intent. However, in visual data exploration, often the intent of the user is not clear. Towards that end, data driven metrics are employed to capture the interestingness or importance of a recommended visualization.

Among the data driven metrics, recent case studies have shown that [a deviation-based metric is effective in providing](#)

Table 1: Table of Symbols

Symbol	Description
k	no. of top recommended views
S	set of top-k recommended views
\mathbb{V}	set of all possible views
X	set of all candidate views
A	a dimensional attribute
M	a measure attribute
F	aggregate function
Q	a user query
D_B	a multi-dimensional database
D_Q	a target subset of D_B
D_R	a reference subset of D_B
V_i	a view query
$I(V_i)$	importance score of V_i
$I(S)$	importance score of views in S
$f(S, D)$	diversity score of views in S
$F(S)$	hybrid objective utility function value of S
$U(V_i)$	the utility score of each candidate view

[analysts with interesting visualizations that highlight some of the particular trends of the analyzed datasets \[11? \].](#)

In particular, the deviation-based metric measures the distance between $V_i D_Q$ and $V_i D_B$. That is, it measures the deviation between the aggregate view V_i generated from the subset data D_Q vs. that generated from the entire database D_B , where $V_i D_Q$ is denoted as *target* view, whereas $V_i D_B$ is denoted as *reference* view. The premise underlying the deviation-based metric is that a view V_i that results in a higher deviation is expected to reveal some important insights that are very particular to the subset D_Q and distinguish it from the general patterns in D_B .

While recommending views based on their importance has been shown to reveal some interesting insight, it also suffers from the drawback of recommending similar and redundant views, which leaves the data analyst with a limited scope of the possible insights. To address that limitation, in this work we posit that employing diversification techniques in the process of view recommendation allows eliminating that redundancy and provides a good and concise coverage of the possible insights to be discovered. In the next section, we discuss in details the formulation of both importance and diversity, and their impact on the view recommendation process.

[check table-text consistency and make shorter](#)

3 NEW TITLE

[short preamble](#)

3.1 Content-Driven Deviation

[As briefly described in the previous section, in this work we adopt a deviation-based metric to quantify the importance of an aggregate view \[15, 16\]. Essentially, the deviation-based](#)

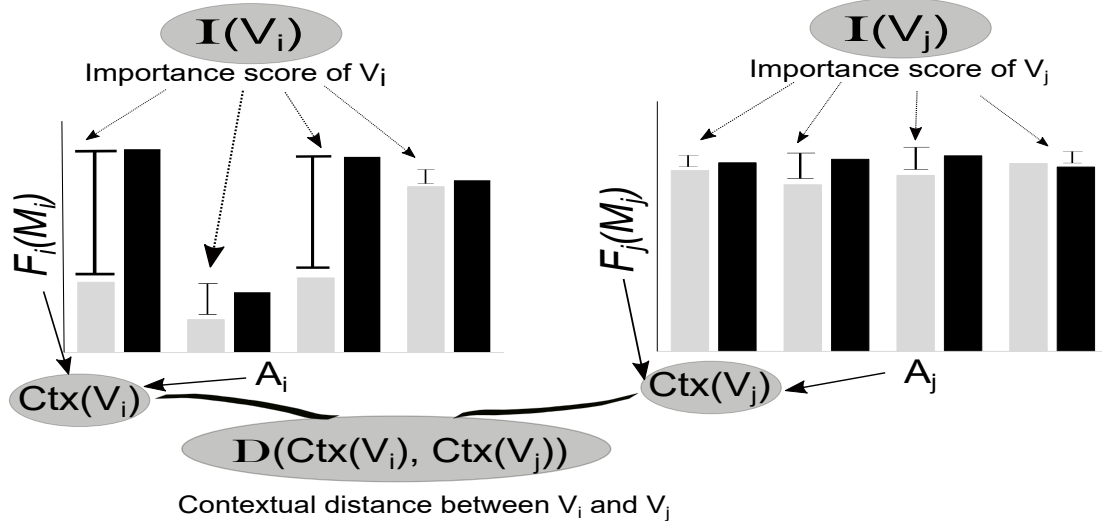


Figure 3: Content vs. Context of views.

metric compares an aggregate view generated from the selected subset dataset D_Q (i.e., target view $V_i D_Q$) to the same view if generated from a reference dataset D_R (i.e., reference view $V_i D_R$). That reference dataset could be the whole database (i.e., $D_R = D_B$) or a selected subset of the database.

Clearly, the deviation between the target and reference views is a *data-driven* metric. That is, it measures the deviation between the aggregate result of $V_i D_Q$ and that of $V_i D_R$. Consequently, from a visualization point of view, it is a *content-based* metric as it captures the difference between the content of the visualization generated by $V_i D_Q$ vs. the visual content generated from $V_i D_R$. In the next, we formally describe the standard computation of that data-driven content-based metric, whereas the discussion of its counterpart context-driven metric is deferred to the next section.

To calculate that data-driven content-based deviation, each target view $V_i D_Q$ is normalized into a *probability distribution* $PV_i D_Q$ and similarly, each reference view into $PV_i D_R$. In particular, consider an aggregate view $V_i = A_i, M_i, F_i$. The result of that view can be represented as the set of tuples: $\langle a_j, g_j, a_j, g_j, \dots, a_t, g_t \rangle$, where t is the number of distinct values (i.e., groups) in attribute A_i , a_j is the j -th group in attribute A_i , and g_j is the aggregated value $F_i M_i$ for the group a_j . Hence, V is normalized by the sum of aggregate values $G = \sum_{j=1}^t g_j$, resulting in the probability distribution $PV = \langle \frac{g_1}{G}, \frac{g_2}{G}, \dots, \frac{g_t}{G} \rangle$.

Finally, the importance score of V_i is measured in terms of a score that captures the distance between $PV_i D_Q$ and $PV_i D_R$, and is simply defined as follows:

$$I(V_i) = \text{dist}(\mathcal{P}[V_i D_Q], \mathcal{P}[V_i D_R]) \quad (1)$$

where $I(V_i)$ is the importance score of V_i and dist is a distance function. In this work, we adopt the Euclidian distance, but other distance measures are also applicable (e.g., Earth Mover's distance, K-L divergence, etc.).

In current approaches for view recommendation, the importance value IV_i of each possible view V_i is computed, and the k views with the highest deviation are recommended (i.e., *top-k*) (e.g., [?]). However, in this work, our goal is to ensure that recommended views provide a good coverage of possible insights, which is achieved by considering the context of the recommended views, which is described next.

This metric measures how different a view is in relation to its reference. Hence, in order to compare the views, it needs to be a target subset and a reference subset of the D_B . Given target subset D_Q and D_R as the reference subset of D_B , we use a deviation-based metric as the content evaluation as employed in [15, 16]. The deviation value between target view $V_i D_Q$ and reference view $V_i D_R$ is called as importance score IV_i . The importance score of V_i is measured in terms of a score that captures the distance between probability distribution of target view $V_i D_Q$ and probability distribution of reference view $V_i D_R$, it can be seen in Figure 3 and it can be formally defined as follow:

$$I(V_i) = \text{dist}(\mathcal{P}[V_i D_Q], \mathcal{P}[V_i D_R]) \quad (2)$$

where $I(V_i)$ is the importance score of V_i and dist is the Euclidian distance or other distance functions. The higher distance between target view $V_i D_Q$ and reference view $V_i D_R$ implies to higher importance score I . When judging the importance score of a view, it is believed that a view with large deviation from a reference view tends to be more interested to the user [15, 16].

To confirm that both views have the same scale, the target view $V_i D_Q$ and reference view $V_i D_R$ are normalized. The results of target and reference views are summaries with two columns, which first column is A and the second column is the aggregate FM . Hence, the normalized probability distribution can be calculated by dividing g_n by G , where g_n is the second column of the result which is the aggregate value of F over M for the group-by A , and $G = \sum_{i=1}^n g_n$. G is the sum of aggregate of all groups n in attribute A . The result of normalized probability distribution can be shown as $\frac{g_1}{G}, \frac{g_2}{G}, \frac{g_3}{G}, \frac{g_4}{G}, \dots, \frac{g_n}{G}$.

stopped here. adjust the next section to consider the changes above

3.2 Context-Driven Deviation

Content-driven deviation only exposes the quality of individual view which is by comparing the view to its reference. However, since the views are recommended as a collection or set which consists of multiple views, the quality of the whole set is as important as the quality of individual views[14]. Hence, relying only on content-driven leads to generate the top-k views which are similar to each other and presenting views with limited amount of gained insights. To overcome this issue, we consider context-driven deviation that measures how different are the individual view within the set of views. In order to measures the contextual distance, we employing diversity function which allows to eliminate the redundancy. In fact, diversity has been well known and widely used in recommendation systems for maximizing information gain and minimizing redundancy[10, 18, 20, 21].

In this work, we employ a diversity measure based on the context of the views which is by calculating the distance between its context of two views.

For instance, let V_i and V_j be two views as shown in Figure 3, and Ctx is the function for obtaining context of the view. Also, let $CtxV_i = \langle A_i, M_i, F_i \rangle$ and $CtxV_j = \langle A_j, M_j, F_j \rangle$ where $A_i, A_j \in \mathbb{A}, M_i, M_j \in \mathbb{M}, F_i, F_j \in \mathbb{F}$.

In order to compare $CtxV_i$ and $CtxV_j$, it requires distance function. There have been a tons of similarity measurement which proposed in the literatures such as 1) binary similarity [2]; 2) cube similarity measure for multidimensional data [1]. There are also well known similarity measurements such as Eucliden, Jaccard, Manhattan, Hamming [2] that can be used for our work. However, since the context of view is the sets of categorical data, the Jaccard similarity can be the proper option. The Jaccard method measures similarity between finite sets as intersection of sets over union of sets.

Moreover, we assign the weights w to each components of the context of the view. As shown in the Figure 3, the context of views is a set that consist of three components which are the name of attribute A , measure M , and aggregate function F . The weights w is assigned to each component to define which component would be more prominent in measuring contextual distance between two views. For instance, analyst may consider that attribute A is more important than measure M and aggregate function F in terms of generating

top-k views. Thus, the weights ratio can be $3w_a : 2w_m : w_f$, where $w_a + w_m + w_f = 1$.

Using this weights, it affects the quality of the result of top-k views. The attribute A has three times weights of aggregate function F . As a consequence, attribute A has three times contribution compared to aggregate function F in measuring the contextual distance of two views. Hence, the Jaccard similarity between two views can be measured as,

where is the weights in equation? it should be weighted jaccard

$$J(CtxV_i, CtxV_j) = \frac{|CtxV_i \cap CtxV_j|}{|CtxV_i \cup CtxV_j|}$$

The contextual distance D between two views can be defined as:

$$D(CtxV_i, CtxV_j) = 1 - J(CtxV_i, CtxV_j) \quad (3)$$

As a summary, in the process to recommend users a set of important views, we work at two levels. At the first level, we evaluate that how “important” is the content of the view as compared to reference view. At the second level, we evaluate contextually how different a view is from other views in the recommended set.

3.3 Problem Definition

The proposed recommendation visualization systems should recommend set of views with the high importance score and maximum diversity as well as it should has high efficiency in terms of costs due to all computations are running on the fly. Hence, there will be two variables that can be used for evaluation, which are *effectiveness* and *efficiency*.

3.3.1 Effectiveness. *DiVE* scheme using hybrid approach which designed to recommend a set of k views that considering importance and diversity. In this section, we formally define the problem of the quality of recommended views.

In order to generate the top-k views that considering both importance and diversity, each view is evaluated by its reference to measure the quality of individual view and all views in the set are evaluated in terms of diversity. For instance, given an example view V_i , $V_i \in S$, set of all possible views \mathbb{V} , our objective is to select a set of views $S \subseteq \mathbb{V}$ where size of $S = i = k$, such that the importance score of the views in S and the diversity of all views in the set S is maximized. To achieve this goal, there are two components that need to be considered, the importance score of a set of views S and the diversity score of the set S . Hence, we need to combine those two components as a hybrid objective utility function.

The importance score computation. The importance score of the set S is calculated as the average value of the importance measure of each view in S , it can be defined as:

$$I(S) = \frac{\frac{1}{k} \sum_{i=1}^k IV_i}{I_u}, V_i \in S$$

The average of the importance score of set S needs to be normalized due to the value is not as straight forward, it depends on the real data. The normalization works by dividing the average of importance score of set S by the absolute maximum of importance score or the absolute upper bound of the importance score I_u .

The diversity score computation. There are several different diversity functions have been employed in the literature [3, 9, 18], among which previous research has mostly focused on measuring diversity based on either the average or the minimum of the pairwise distances between elements of set [19]. We focus on the first of those variants (i.e., average), as it considers all the views in S . Given a distance matrix $D(V_i, V_j)$ as given in equation 2, the diversity of a set S can be measured by a diversity function $f(S, D)$ that captures the dissimilarity between the views in S , defined as:

$$f(S, D) = \frac{1}{k(k-1)} \sum_{i=1}^k \sum_{j>i}^k D(V_i, V_j), V_i, V_j \in S$$

There is no need to normalize the function for the diversity score of set S $f(S, D)$ due to the maximum diversity score of each pairs of views $div(V_i, V_j)$ never be more than 1 as in equation 3 and $f(S, D)$ uses $k(k-1)$ as the divisor to make sure that there will be no result that more than 1.

Hybrid objective utility function. In order to capture both importance and diversity in the set of recommended views, we define a hybrid objective function. Specifically, for a set of views $S \subseteq V$ an objective function is formulated as the linear weighted combination of the importance score, $I(S)$ and diversity function $f(S, D)$ which is defined as:

$$F(S) = (1 - \lambda) \cdot I(S) + \lambda \cdot f(S, D) \quad (4)$$

where $0 \leq \lambda \leq 1$ is employed to control the contribution between importance and diversity in the hybrid objective function. The higher values of λ result in a set of more diverse views whereas lower values of λ generate a set of the most important views that might be similar to each other. Given the hybrid objective function, our goal is to find an optimum set of views S^* that maximizes the objective function $F(S)$:

$$S^* = \underset{\substack{S \subseteq V \\ |S|=k}}{\operatorname{argmax}} F(S) \quad (5)$$

3.3.2 Efficiency. In the section before, we defined the problem in terms of quality of the results which hybrid utility function was proposed as the solution. This section defines the issue in terms of efficiency, as explained in the preliminaries section that in case of the modest data which has small number of dimensions, it can generate large number of views. Meanwhile, views that will be presented to users only in small number (top-k views) and all computation will be done on the fly. Hence, we need the scheme that robust in terms of the quality of results as well as in terms of efficiency (costs).

In fact, it has been mentioned a lot in the literatures [5, 15, 16] that the main issue of the visualization recommendation systems is the query execution. To deal with the query costs, there are several approaches that proposed, such as query results caching[9], shared computation among views, combine target and reference query, combination of multiple aggregates, combination of multiple GROUP BY, and parallel query and execution, [16, 19].

However, in this work, we propose *DiVE* scheme that leverages the properties of importance and diversity metrics to prune a large number of query executions without reducing the quality of the results. The detail is presented in the proposed methods section.

4 PROPOSED METHODS

As discussed in the introduction, the current view recommendations [5, 15, 16] generated solely on the basis of importance score suffer from the redundancy problem. The extreme solution to overcome the redundancy in the top-k views in the set S is to select views such that the diversity score of S is maximized.

This leads to two extreme baseline solutions, one based only on the importance score of the views which is called "**Linear-Importance**" and second based on only diversity score of the views which is called "**Greedy-Diversity**".

Instead of uses one of those two extrem, *DiVE* scheme employed hybrid function for evaluating top-k views which captures both the importance as well as diversity. Moreover, *DiVE* also equipped with pruning scheme which can prune a large number of queries without reducing the quality of the result. Our proposed *DiVE* scheme described below.

4.1 DiVE-Greedy Scheme

Algorithm 1: *DiVE* Greedy

Input: Set of views V and result set size k

Output: Result set $S \geq V$, $|S| = k$

```

1  $S \leftarrow [V_i, V_j]$  get two most distant views;
2  $X \leftarrow [V \setminus S]$ ;
3  $i \leftarrow \text{len}(S)$ ;
4 while  $i < k$  do
5   for  $j$  in set  $X$  do
6      $\max_v \leftarrow \operatorname{argmax}_F(X[j], S)$ ;
7   end
8    $S.add(\max_v)$ ;
9    $X.remove(\max_v)$ ;
10   $i \leftarrow i + 1$ ;
11 end
12 return  $S$ 
```

In order to capture both importance and diversity in the recommended top-k views, *DiVE* employs a greedy construction algorithm to iteratively select views that maximize the

hybrid objective function $F(S)$ as presented in equation 3, and it is called as **DiVE-Greedy** scheme.

The details of the *DiVE-Greedy* scheme are given in Algorithm 1. The key ingredient of any Greedy algorithm for solving an optimization problem is the objective function itself that needs to be maximized. In particular, *DiVE-Greedy* initializes the set S with two most distant views. The distance between all the views is calculated using the distance function as given in equation 2. In each iteration, a new candidate view is selected from remaining views X and added to S . Thus, *DiVE-Greedy* assigns a utility score to each candidate view which is based on the hybrid objective function $F(S)$ as defined in equation 4. The utility score of each candidate view V_i in X is computed as:

$$U(V_i) = (1 - \lambda) \cdot I(V_i) + \lambda \cdot \text{setDist}(V_i, S) \quad (6)$$

$$\text{Where } \text{setDist}(V_i, S) = \frac{1}{|S|} \sum_{V_j \in S} D(V_i, V_j)$$

Thus, the view with highest utility score in each iteration is selected and added to S .

DiVE-Greedy cost. The costs of Greedy Construction algorithm has two components which are the query execution cost C_Q that computing the importance score of view and the diversity cost C_D that computing set distance of each view from the views already in S . The complexity of query execution cost is $O(n)$ as the content of each view is generated only once. Meanwhile, the diversity cost C_D is $O(kn)$ where k is the size of subset of views S and n is the number of all possible views.

Traditional diversification method applied approach such "process-first-diversity-next" which leads to generating all possible views and computing the importance score in advance, then the diversity is executed next. Although, greedy algorithm is very efficient, however, as the number of attributes A , measures M and aggregate functions F increase, the number of views that need to be generated grows exponentially. Hence, without a good strategy, *DiVE-Greedy* suffers from the high costs of query execution.

Therefore, to overcome this issue, *DiVE-Greedy* employs static pruning strategy as elaborated next.

4.2 DiVE-Greedy-Static Pruning

In order to recommend a small subset of views, the importance score I of all possible views need to be computed. However, only few views are eventually recommended to the analyst and rest of the views are discarded. Clearly, this approach would hinder the performance of the view recommendation systems for high dimensional datasets. Thus, motivated by the need to reduce the number of views that need to be generated, *DiVE-Greedy* employs a static pruning technique.

The proposed pruning technique is based on the observation that the utility score of each view UV_i is a weighted sum of two different measures; 1) the importance score of view IV_i and 2) diversity score of a view from S $\text{setDist}(V_i, S)$. The

Views in $V \setminus S$	$\min U'$	$\max U'$
V_1	0.23	0.65
V_2	0.19	0.21 ✕
V_3	0.25	0.85
V_4	0.15	0.70
V_5	0.20	0.23 ✕
V_6	0.21	0.24 ✕
V_7	0.22	0.91

Figure 4: Max-Min Pruning: All views which has $\max U'$ less than the maximum of $\min U'$ will be pruned

diversity score of a view requires only CPU computations and is faster operation. Whereas, computing the importance score of a view by comparing the target view to the reference view incurs high cost which is dominated by I/O cost of the query executions.

DiVE-Greedy applied Max-Min pruning method which has been presented in [8], by leveraging the diversity score of a view to decide whether a view query should be executed or not, and it is called **DiVE-Greedy-Static**.

Max-Min utilize the maximum bound of importance score I_u and minimum bound of importance score $I_0 = 0$. In each iteration, the distance $\text{setDist}(V_i, S)$ between all remaining views in X to the current set is calculated. Hence, the utility score of each views is computed using the real value of distance $\text{setDist}(V_i, S)$ and the importance score I_u and $I_0 = 0$. Thus, it produces two version of utility score which are $\min U'$ and $\max U'$ of each candidate views, which can be formally defined as:

$$\begin{aligned} \min U'(V_i) &= (1 - \lambda) \cdot I_0(V_i) + \lambda \cdot \text{setDist}(V_i, S) \\ \max U'(V_i) &= (1 - \lambda) \cdot I_u(V_i) + \lambda \cdot \text{setDist}(V_i, S) \end{aligned}$$

If $\max U'$ of the candidate view less than the maximum of $\min U'$, then this view definitely will be pruned, the detail example can be seen in Figure 4. This approach generates same set of recommended views as generated by the *DiVE-Greedy* without pruning. This is due to the fact that in each iteration *DiVE-Greedy* selects the view with the highest utility score to be added to S .

This Max-Min pruning approach has been mentioned has a good performance in pruning [8]. However, in that literature, the dataset have a large of points which the value of each point is very diverse to others. To the contrary, this work is quite different. The context of view only has three dimensions, while using three dimensions and computing the diversity score of the small set of view S , the utility score may not diverse as in the literature. Perhaps the performance of pruning is not same as well as in the literature.

Moreover, although Greedy algorithms have been shown to be efficient and provide good approximations to the optimal

solutions [20], [18], [13]. However, since Greedy is constructive type algorithm, it constructs the set S by adding a new candidate view, there is no guarantee that the new view selected in each iteration is the best view for the objective function FS . It is because the view which has the highest utility score not necessary be the best one that improve the objective function FS (e.g: local optimum).

To overcome this Greedy limitation, in this work, we also proposed another scheme which based on swap technique. The discussion of the swap technique is elaborated next.

Algorithm 2: DiVE Swap

Input: Set of views V and result set size k
Output: Result set $S \geq V$, $|S| = k$

```

1  $S \leftarrow$  Result set of only importance or only diversity;
2  $X \leftarrow [V \setminus S]$ ;
3  $F_{current} \leftarrow 0$ ;
4  $improve \leftarrow True$ ;
5 while  $improve = True$  do
6   for  $i$  in set  $X$  do
7      $S' \leftarrow S$ ;
8     for  $j$  in set  $S$  do
9       if  $F(S') < F(S \setminus S_j \cup X_i)$  then
10         $S' \leftarrow S \setminus S_j \cup X_i$ ;
11      end
12    end
13    if  $F(S') > F(S)$  then
14       $S \leftarrow S'$ 
15    end
16  end
17  if  $F(S) > F_{current}$  then
18     $F_{current} \leftarrow F(S)$ ;
19     $improve \leftarrow True$ ;
20  else
21     $improve \leftarrow False$ ;
22  end
23 end
24 return  $S$ 
```

4.3 DiVE-Swap Scheme

Swap is local search type algorithm and it has been known and used to maximize diversity in the literature [4, 18]. This algorithm starts with a complete initial set S , and try to achieve better result by interchanging the remaining views in X to the current set S . If the views in X is able to give better objective function value $F(S)$, then this view will be joined to the current set and one view in the current set that has the lowest contribution to the $F(S)$ will be removed. The details of *DiVE-Swap* algorithm can be seen in Algorithm 2.

We proposed two types of Swaps which are: 1) **DiVE-iSwap**, the underlying behind this scheme is, it has the initial set from the result of Linear-Importance which is importance score maximized. Thus, there will be exchanging view from

X to the current set and view that can improve the $F(S)$ of current set S will be considered as the member of S . 2) **DiVE-dSwap** which is quite similar to DiVE-iSwap, however, this scheme is initialized by results of Greedy-Diversity, which is diversity maximized.

DiVE-Swap cost. The costs of Swap algorithm is also depend on the query execution time C_Q of all possible views and the diversity computation C_D . The query cost C_Q is executed only once but the cost is high due to it needs I/O cost. However, the complexity of diversity computation C_D is $O(k^2)$ and the number of distance computation depends on the number of iterations of the swap and the number of views in X . In the worst case, swap algorithm can perform $O(k^n)$ iterations.

DiVE-iSwap utilize the results of Linear-Importance as the first initialization. Hence, this algorithm cannot escape from executing all queries due to Linear-Importance needs to execute all possible views to get the results. However, the second proposed swap algorithm, *DiVE-dSwap* is initialized by the result of Greedy-Diversity. This algorithm does not execute any query to generate the results. Therefore, leveraging the properties of diversity for pruning in *DiVE-dSwap* scheme is possible. The proposed static pruning of *DiVE-dSwap* is discussed below.

4.4 DiVE-dSwap Static Pruning

As in the Greedy technique, *DiVE-dSwap* leverages the diversity score of a candidate view to decide whether a view query should be executed or not. The details of static pruning technique of *DiVE-dSwap* as following:

- *DiVE-dSwap* using initial set which is the result of Greedy-Diversity. All query views in the initial set need to be executed to get the objective function FS of the current set S .
- Before start an exchanging view from X to the current set S . All views in X will be sorted based on $setDist(V_i, S)$. This is to make sure that exchanging process starts from the view that has highest score of diversity to the current set. We called this technique as "top-1" method.
- In order to exchanging the candidate view in X , the query of the candidate view need to be executed to get the importance score. Instead of getting the actual importance score, the maximum bound of importance score is used to compute the utility score of each view as in Greedy technique. Hence, the result is not the actual utility score but $maxU'$, which defined as: $maxU'(V_i) = (1 - \lambda) \cdot I_u(V_i) + \lambda \cdot setDist(V_i, S)$.
- The exchanging process is started by comparing FS of the current set to FS of new set as can be seen in Algorithm 2 lines 9 - 10.

If candidate views in X , while using importance score I_u cannot perform better in terms of improving the objective function FS to the current set S , those views will be pruned. This method is valid due to if the maximum score of importance is used and that view cannot improve FS of the

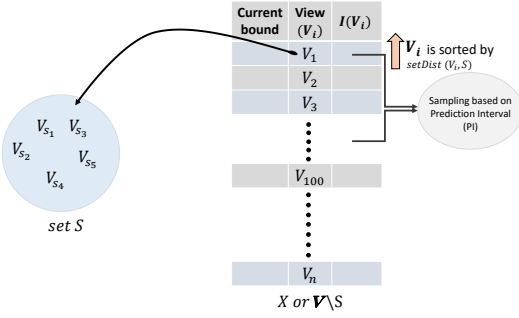


Figure 5: dSwap-Pruning: Each candidate view has importance score and *currentbound*. The importance score of view will be generated only for view which its utility score is able to improve set S while using *currentbound*

current set, then there is no reason to execute the view query to get the importance score. This pruning scheme is called **DiVE-dSwap-Static**.

All proposed pruning techniques including DiVE-Greedy-Static and DiVE-dSwap-Static are using static value I_u as the bound. Hence, the pruning performance of this technique will not optimal while the value of I_u is far away from the real maximum value of importance in database.

To overcome this issue, we proposed adaptive pruning scheme as described in the next section.

4.5 Adaptive Pruning Scheme

The idea behind pruning scheme is to minimize the query execution, which is by early prune low quality views. There are two factors that improve pruning performance: 1) weight of λ and 2) the max bound value I_u .

For instance, assume that an analyst wants to get set of views from view recommendation and she uses $\lambda = 0.7$. The λ value equal to 0.7 means that the diversity contribution to the utility score is 70% and the contribution of importance score will be 30 %, as it can be seen in equation 4. Thus, by setting the λ to higher value, the contribution of importance score will be lower and there will be more prune queries. This λ value is determined by analyst, however, in this experiment 0.5 is used as the default.

Moreover, *DiVE schemes* utilize the max bound I_u to check whether the query view need to be executed or not. In case that a view cannot improve the current set while using I_u then this view will not be executed. Hence, the optimal pruning can be achieved while the max bound value is close to the real maximum value of importance score in the dataset. In order to overcome this issue, instead of using static bound I_u , we proposed adaptive pruning scheme that automatically adapts the bound to the real maximum importance score in the dataset.

Adaptive pruning in Greedy is called *DiVE-Greedy-Adaptive* whereas in Swap is called *DiVE-dSwap-Adaptive*.

Table 2: Parameters testbed in the experiments

Parameter	Range (default)
datasets	Heart disease, Flights
sample queries	10
diversity weight ratio	3(A) : 2(M) : 1(F)
tradeoff weight λ	0.0, 0.2, 0.4, 0.5 , 0.6, 0.8, 1.0
result set (size of k)	5 , 15, 25, 35
prediction interval %	80 , 85, 90, 95, 97 , 98

5 EXPERIMENTAL EVALUATION

Experiment Setup. This experiment running on Windows 10 64 bit, Intel Core i7-7700 CPU @ 3.60 GHz, RAM 16 GB. The experiment is developed using Python and run on Python version 3.6.3 with PostgreSQL as the database engine. Two real datasets are used, as follows:

- Heart Disease Dataset ², it has 14 attributes in total which are 9 attributes A and 5 measure attributes M , A = sex, cp, fbs, restecg, exang, slope, ca, thal, num; M = age, trestbps, chol, thalach, oldpeak. This dataset has small number of rows (299 rows).
- Airline (Flights) Dataset ³, it has 7 attributes, 4 measure attributes. The detail as follows: A = year, month, week, day, carrier, origin, destination; M = arrivaldelay, departurdelay, weatherdelay, distance. This dataset has 855,632 number of rows.

5.1 Effectiveness Evaluation

In order to test the performance of our proposed approach, we run experiments using two real datasets: flights dataset and heart disease dataset. For each dataset, we execute ten random queries. Afterwards, the average of all results is calculated as the final result. We examine the performance of our proposed schemes in term of the quality of the result (effectiveness) and the efficiency. All parameters used in this experiments are summarized in Table 2 and the detail results of the experiments is elaborated next.

The impact of parameter weight λ to the overall $F(S)$. One of the advantage *DiVE* scheme is users can decide what kind of the results that they want by changing the value of weight λ . The main function of λ is this weight can be used to tradeoff between importance and diversity, it is explained in equation 4. Figure 6 shows the impact of λ to the overall objective values $F(S)$ by fixed k equal to 5, where $0 \leq \lambda \leq 1$. The figure shows that in the beginning when the λ is close to 0, means the role is dominated by the importance score, Linear-Importance has high overall objective function value $F(S)$ but it will decline gradually while increasing of λ . To the contrary, the Greedy-Diversity is the opposite, it has the high overall objective function value when the λ equal to 1, which means the diversity gets the full role. Hence, there is a crossover between Linear-Importance and Greedy-Diversity,

²<http://archive.ics.uci.edu/ml/datasets/heart+Disease>

³<http://stat-computing.org/dataexpo/2009/the-data.html>

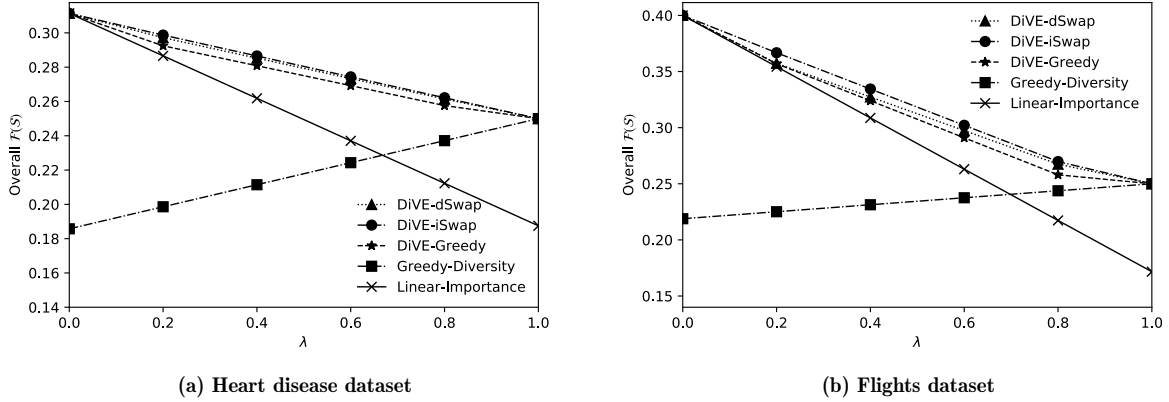


Figure 6: Impact of λ to overall objective function value $F(S)$ while $k = 5$ and running on three real datasets

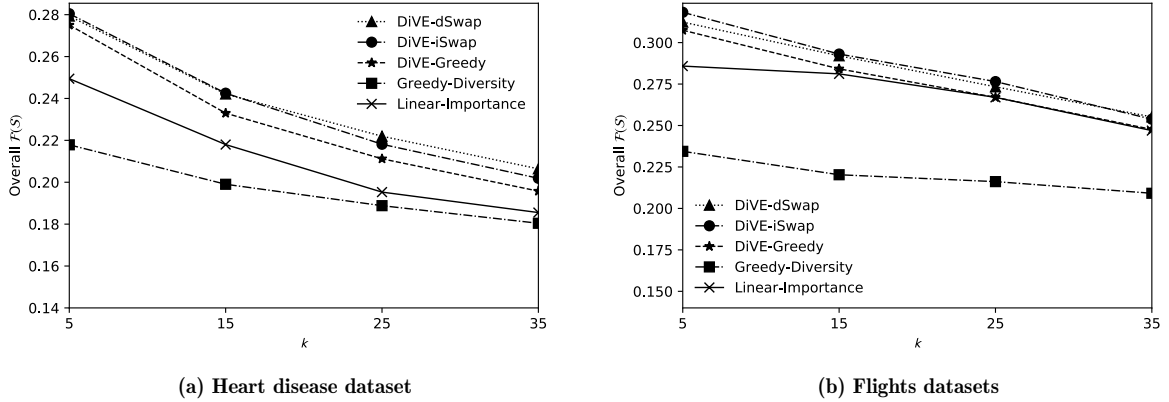


Figure 7: Overall objective function value $F(S)$ using different value of k and running on three real datasets

this crossover shows the role balancing between importance score and diversity score. Furthermore, our proposed schemes have stable performance in all various λ value and its $F(S)$ value always better than Linear-Importance and Greedy-Diversity.

The impact of k to the overall objective function value $F(S)$. It has been explained the impact of λ to the overall objective function value $F(S)$. This section discuss the overall objective function value $F(S)$ by increasing k . Figure 7 shows the overall objective values $F(S)$ by using different of k . As shown in those figures, the value of $F(S)$ decrease while increasing the value of k . There are two reasons why this moment happens: 1) the importance score of the set $I(S)$ is generated by sum of all importance score of views and it is divided by the number of views= k , increasing the number of k means increasing the divisor. 2) the diversity score of a set of views $f(S, D)$ also depends on the number of k . While k is increasing, the diversity score of the set will be decreased,

due to the score of $f(S, D)$ is yielded from the sum of all distance of views divided by $k * k - 1$. Hence, the more views selected, the probability to get similar views will be high and the score of diversity will be declined.

The most important point in this results is our *DiVE* schemes always have better overall objective function value $F(S)$ compared to the two extreme baselines in various number of k .

5.2 Efficiency and Pruning Scheme Evaluation

The total time to run *DiVE* schemes. In order to start analyzing the efficiency, we need to know the main issue in term of costs. Figure 8 shows the example of exactly time that needed to run schemes on flights dataset. It shows Greedy-Diversity which only considering diversity and no query executions, it has very low costs. Meanwhile, Linear-Importance and *DiVE-Greedy* seems in the same line but that was not exactly same. The total of diversity computations C_D of *DiVE-Greedy* is very low, the total costs of *DiVE-Greedy* is dominated by

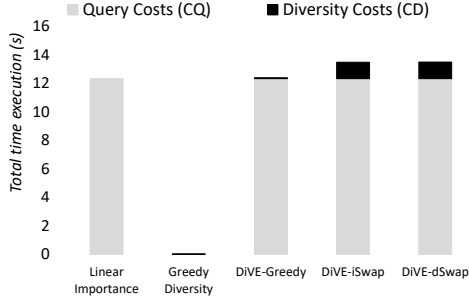


Figure 8: Total time (seconds) to execute schemes on flights dataset using $k=5$ and $\lambda = 0.5$. It shows that costs are dominated by query costs C_Q

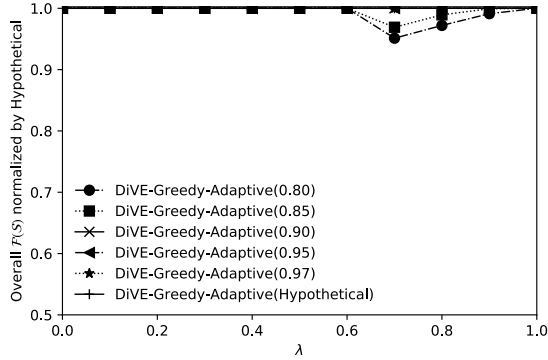


Figure 9: Impact of λ to the performance of *DiVE-Greedy-Adaptive-Pruning* with different of PI in terms of the effectiveness, running on Heart disease dataset using $k = 5$. Using PI 80 and PI 85 can reduce the quality of the result while the λ value is low

query costs C_Q . Due to of this reason, the total execution time of *DiVE-Greedy* closes to Linear-Importance. Those are the proof that the total costs of all schemes except Greedy-Diversity are dominated by query cost C_Q .

Due to the page limitation, all our experiments in three real dataset cannot be showed. Hence, for the next sections, we use heart disease dataset as our focus observation.

Impact of λ to the pruned queries of *DiVE-Greedy* scheme.

In this work, we proposed two kind of pruning schemes, that are using estimated static value of maximum value of the importance score $maxI$ and using adaptive maximum value of importance score.

The first our pruning scheme is using static value $maxI$. To check the performance of our pruning scheme, we run it to all real datasets using different value of λ . We analyze the result by comparing the result of schemes with pruning enable on it and the schemes without pruning enable on it.

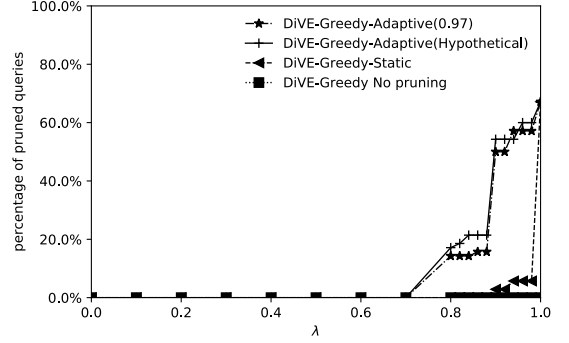


Figure 10: Impact of λ to the pruned queries of *DiVE-Greedy* scheme, running on Heart disease dataset using $k = 5$. Using PI-97 able to prune queries around 20 percent while λ higher than 0.8

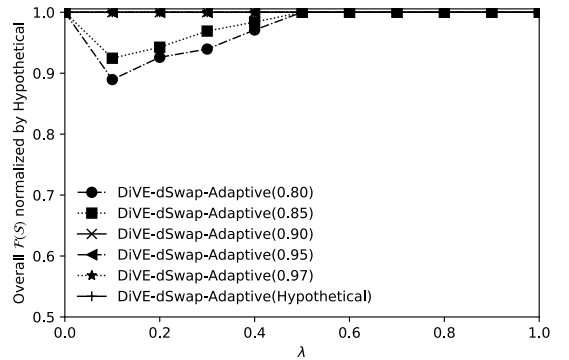


Figure 11: Impact of λ to the performance of *DiVE-dSwap-Adaptive-Pruning* with different of PI in terms of the effectiveness, running on Heart disease dataset using $k = 5$. Using PI 80 and PI 85 can reduce the quality of the result while the λ value is low

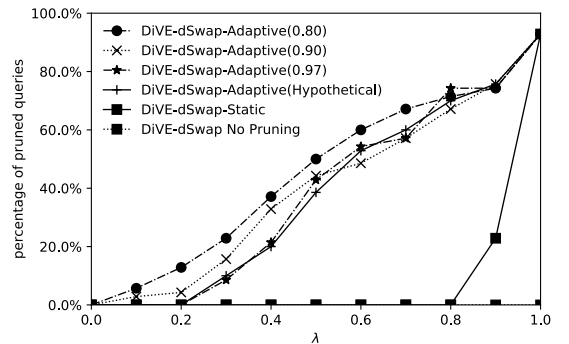


Figure 12: Impact of λ to the pruned queries of *DiVE-dSwap* scheme, running on Heart disease dataset using $k = 5$. Adaptive pruning scheme able to prune queries start from λ 0.1, it has more pruned queries while λ increased

The example result which running on heart disease dataset is shown in Figure 10. The static pruning scheme only able to prune queries while the value of λ is high, closes to 0.9. As we expected, it is because the value of $maxI$ that too far from the real value of importance score.

To overcome the flaw in static pruning scheme, we also proposed adaptive pruning scheme which used dynamic value of the maximum importance score $maxI$. In order to get $maxI$ as close as possible to the real value of importance in the dataset, we applied sampling method. By using adaptive pruning, users also able to change the confidence and the margine error to tradeoff between time and precision.

Impact of λ to the pruned queries of DiVE-dSwap scheme.

If static pruning scheme only able to prune queries while λ closes to 0.9 and higher, in this section, we shows the adaptive pruning performance. Figure 12 shows the performance of adaptive pruning scheme by using different value of λ . It shows the impact of λ to the percentage of pruned queries. The adaptive pruning schemes especially *DiVE-dSwap-Adaptive-Pruning* is able to prune queries significantly, pruning start while λ closes to 0.2 and by increasing the λ , more queries can be pruned.

The impact of query load to the pruning performance. As shown in Figure 7, that the actual value of importance score affects to the overall value of objective function $F(S)$ and the shape of the graphs. The next question is do the actual value of importance score also affects the pruning performance while using adaptive pruning scheme? We did experiments using heart disease dataset to answer this question. We lists all subsets in the heart disease dataset and classify those subsets to three categories: low, middle, and high. Low category consist of subsets that have mostly low value of importance score, high category is the opposite, and the middle category is the middle of low and high. We compare the result in terms of pruning performance using input low, middle, and high categories. The result can be seen in Figure 13. It has different result while using low category and high category. Input using low category or low query load has more pruned queries compared to high importance query load (high category).

6 CONCLUSIONS

In this paper, we proposed *DiVE* scheme which the main purposes are to evaluates and optimizes the results of visualization recommendation systems with respect to importance and diversity. The advantage of *DiVE* is that analyst can set their preferences by changing the parameter to tradeoff between importance and diversity to get result set. We also performed an experimental study and present the results which focus on effectiveness and efficiency of our approach on real datasets. We proposed *DiVE* scheme which based on Greedy and Swap approach, *DiVE-iSwap* have the best performance in recommending result views but it has the highest costs due to this scheme executing all possible view from the dataset, this scheme can be used for the analyst who only cares about the results without worrying execution time. However, to the analyst who care about execution time, we

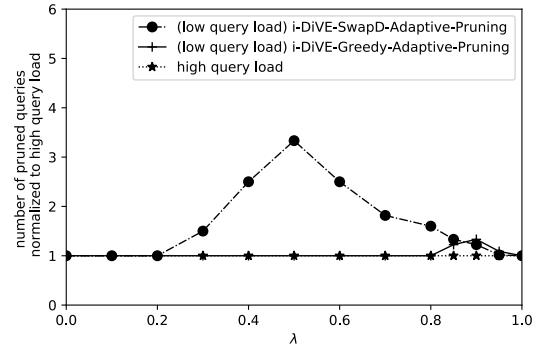


Figure 13: Number of pruned queries of high and low query load normalized by high importance query load using different value of λ , $k = 5$ and running on Heart disease dataset

proposed *DiVE-dSwap-Adaptive* and *DiVE-Greedy-Adaptive*, those schemes are able to decrease costs significantly without reducing the quality of results.

ACKNOWLEDGMENTS

This research is supported by the Indonesia Endowment Fund for Education (LPDP) as scholarship provider from the Ministry of Finance, Republic of Indonesia.

REFERENCES

- [1] E. Baikousi, G. Rogkakos, and P. Vassiliadis. [n. d.]. Similarity measures for multidimensional data. *ICDE 2011* ([n. d.]).
- [2] S. Choi, S. Cha, and C. C. Tappert. 2010. A survey of binary similarity and distance measures. *Journal of Systemics, Cybernetics and Informatics* (2010).
- [3] Charles L.A et al. Clarke. 2008. Novelty and diversity in information retrieval evaluation. *SIGIR* (2008).
- [4] M Drosou and E Pitoura. 2010. Search Result Diversification. *SIGMOD Record* (2010).
- [5] H Ehsan, M. Sharaf, and Panos K. Chrysanthos. 2016. MuVE: Efficient Multi-Objective View Recommendation for Visual Data Exploration. *ICDE* (2016).
- [6] Sean Kandel, Ravi Parikh, Andreas Paepcke, Joseph M Hellerstein, and Jeffrey Heer. 2012. Profiler: Integrated statistical analysis and visualization for data quality assessment. In *Proceedings of the International Working Conference on Advanced Visual Interfaces*. ACM, 547–554.
- [7] Alicia Key, Bill Howe, Daniel Perry, and Cecilia R. Aragon. 2012. VizDeck: self-organizing dashboards for visual analytics. *SIGMOD Conference* (2012).
- [8] Hina A. Khan and Mohamed A. Sharaf. 2015. Progressive diversification for column-based data exploration platforms. *ICDE* (2015).
- [9] Hina A Khan, Mohamed A Sharaf, and Abdullah Albarrak. 2014. DivIDE: Efficient Diversification for Interactive Data Exploration. *SSDBM* (2014).
- [10] Davood Rafiei, Krishna Bharat, and Anand Shukla. 2010. Diversifying web search results. *WWW* (2010).
- [11] Thibault Sellam and Martin L. Kersten. 2016. Ziggy: Characterizing Query Results for Data Explorers. *PVLDB* 9, 13 (2016), 1473–1476. <http://www.vldb.org/pvldb/vol9/p1473-sellam.pdf>
- [12] Jinwook Seo and Ben Shneiderman. 2006. Knowledge Discovery in High-Dimensional Data: Case Studies and a User Survey for the Rank-by-Feature Framework. *IEEE Trans. Vis. Comput. Graph.* 12, 3 (2006), 311–322. <https://doi.org/10.1109/TVCG.2006.50>
- [13] Barry Smyth and Paul McClave. 2001. Similarity vs . Diversity. (2001).

- [14] Manasi Vartak, Silu Huang, Tarique Siddiqui, Samuel Madden, and Aditya Parameswaran. 2017. Towards Visualization Recommendation Systems. *ACM SIGMOD Record* (2017).
- [15] Manasi Vartak and Samuel Madden. 2014. S EE DB : Automatically Generating Query Visualizations. *VLDB* (2014).
- [16] M Vartak, S Rahman, S Madden, A Parameswaran, and N Polyzotis. 2015. SEEDB : Efficient Data-Driven Visualization Recommendations to Support Visual Analytics. *VLDB* 8 (2015).
- [17] Fernanda B. Viegas, Martin Wattenberg, Frank Van Ham, Jesse Kriss, and Matt McKeon. 2007. Many Eyes: A site for visualization at internet scale. *IEEE Transactions on Visualization and Computer Graphics* (2007), 1121–1128.
- [18] Marcos R. et al. Vieira. 2011. On query result diversification. *ICDE* (2011).
- [19] Eugene Wu, Leilani Battle, and Samuel R. Madden. 2014. The case for data visualization management systems. *VLDB Endowment* (2014).
- [20] Cong Yu, Laks Lakshmanan, and Sihem Amer-Yahia. 2009. It takes variety to make a world: diversification in recommender systems. *EDBT* (2009).
- [21] Mi Zhang and Neil Hurley. 2008. Avoiding Monotony: Improving the Diversity of Recommendation Lists. *ACM Conference on Recommender Systems* (2008).