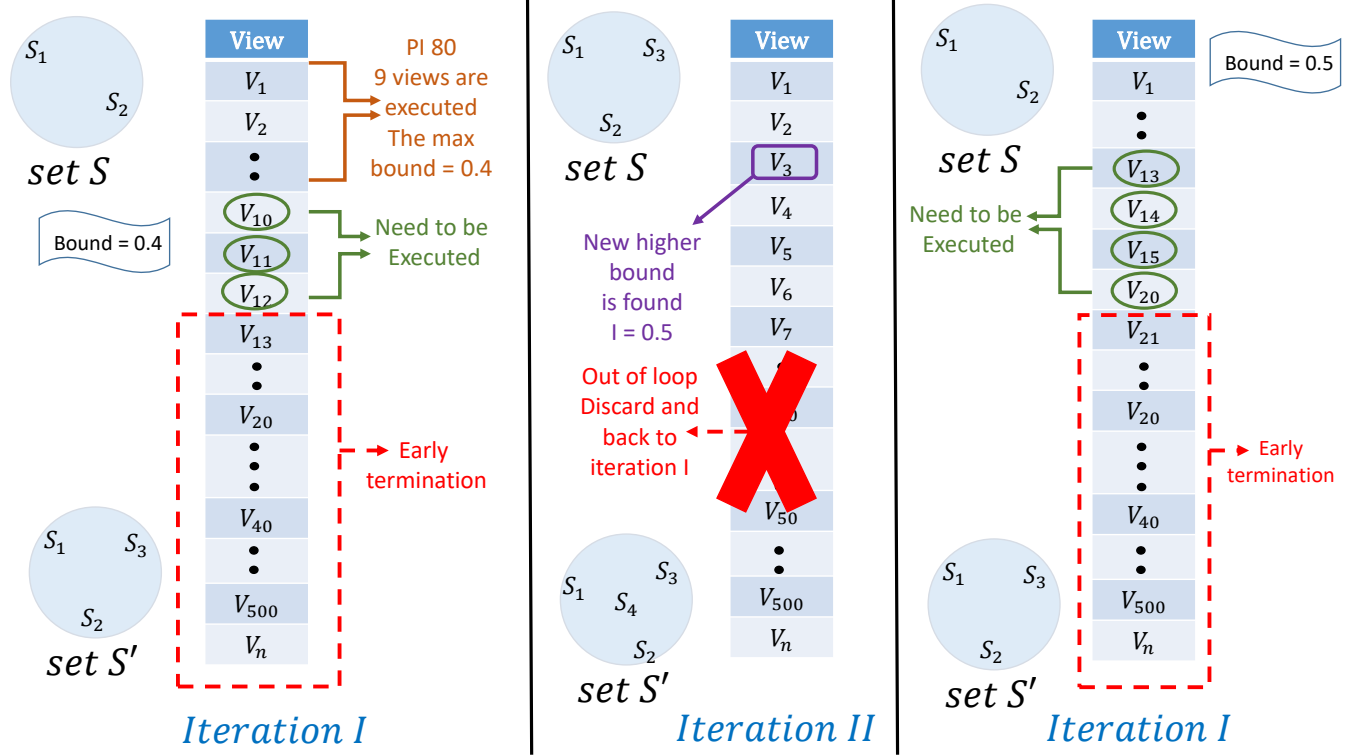


# **Rectifying Bound**

*22 August 2018*

Figure 1: Rectifying maximum bound in DiVE-Greedy-Adaptive while  $k = 5$ 

## 1 Correcting wrong maximum bound in adaptive pruning schemes

In order to reduce costs, our DiVE schemes utilize the importance score bound to do pruning. There are two pruning techniques proposed: 1) static bound approach and 2) adaptive bound approach. In the static bound, the theoretical maximum bound ( $\sqrt{2}$ ) is used and this bound will not be changed until the end of running. Meanwhile, adaptive used the estimation of maximum bound as the first, then this bound is updated where there is a higher maximum bound found.

To know when the bound should be updated, sampling based on prediction interval is used. Before running the program, user needs to defined what PI that she wants to use. For instance, while users set PI to 80, it means after 9 views are executed, the current bound will be updated to the maximum importance score which have seen so far. In the algorithm 1 and 2,  $getMaxPI(L)$  is the function to get the estimated maximum bound from some number of executed views based on PI. Generally, PI can be defined as following:

- PI80: need to execute 9 views
- PI85: need to execute 12 views
- PI90: need to executes 20 views
- PI95: need to executes 40 views
- PI97: need to executes 60 views

Our experiment results show that Adaptive scheme has the best pruning performance while PI80 is used. However, it reduces the effectiveness of recommended views due to only small number of executed views are needed for PI80. It causes the increasing of probability to have wrong bound. The safest way is to use higher PI such as PI95 or PI97. However, if there is a way to keep using PI80 without reducing effectiveness, it will definitely be very good. In fact, *the goal of pruning scheme is to minimize query view execution (i.e., use low PI) without reducing the quality of recommended views*. In order to overcome this issue, rectifying bound of adaptive pruning is proposed. The algorithms of our rectifying bound can be seen in algorithm 1 for DiVE-Greedy-Adaptive and 2 for DiVE-dSwap-Adaptive.

To get the importance score of view, the function ***get\_I\_score(Vi)*** will refer to this Data Frame

View	I (Vi)
$V_1$	0.32
$V_2$	0.56
$V_3$	0.41
$V_4$	unknown
$V_5$	unknown
$V_6$	unknown
⋮	⋮
$V_{55}$	unknown
⋮	⋮
$V_{500}$	$V_{500}$

**Data Frame (DF)**



**Database**

While ***get\_I\_score(Vi)*** finds the importance score of view is unknown, then query view will be executed. the importance score will be stored to Data Frame and returned to the function.

In the end of running, the number of executed views and unexecuted views in this Data Frame are counted

As a result, the **number of pruned queries can be known**

Figure 2: *get\_I\_score* function

The idea behind the rectifying bound algorithm is to keep track: 1) the position of maximum *setDist* score in  $L$  while it gets early termination, 2) the set  $S$  (i.e., the initial set  $S$  while iteration start ) and 3)  $S'$  (i.e., the set  $S$  in the end of iteration) in each Greedy and Swap iteration. For instance, Figure 1 shows the rectifying algorithm of DiVE-Greedy-Adaptive while  $k = 5$ . Firstly, there are two most distant views in the set  $S$  as the initialization. The first iteration of Greedy is to add one most optimal view to the set  $S$ , while in the end of the first iteration, size of  $S = 3$ . As shown in the Figure, some views are executed in advanced (i.e.,  $V_1 - V_9$ ) due to PI80 is currently used (Algorithm 1, line 11 ). The one of view from 9 executed views which has the highest improvement to set  $S$  will be added to set  $S$  (i.e.,  $S \cup V_1$ ), this  $F(S \cup V_1)$  will be the current  $F(S)$ . Moreover, the  $max_b$  is obtained from the highest importance score of these executed views, e.g.,  $max_b = 0.4$ .

Notice that up to this step, the current  $F(S)$  and  $max_b$  are known. The next step is to calculate  $maxF(S \cup V_i)$  using actual diversity score ( $setDist(V_i, S)$ ) and upper bound of importance score ( $max_b$ ). If  $maxF(S \cup X_i) > F(S)$ , then the  $X_i$  can “potentially” improve  $F(S)$ . However, if  $maxF(S \cup X_i) < F(S)$ , then  $X_i$  can be “pruned”. As shown in the *Iteration I*, three views ( $V_{10}, V_{11}, V_{12}$ ) have  $maxF > F$ . Hence, these three views need to be executed. Otherwise, views below  $V_{13}$  can be ignored (i.e., early termination start on  $V_{13}$ ). In the end of the *iteration I*, the result set  $S'$  will be generated which consists of three views (set  $S \cup X_i$ ). Moreover, in the end of each iteration,  $S, S', L, max_d$ , and *iteration count* are stored,  $max_d$  is the maximum view while it gets early termination (i.e.,  $V_{13}$ ).

While in the next iteration (i.e., *Iteration II*), there is a higher bound compared to the current  $max_b$ , then the current  $max_b$  will be updated to the new one and the loop will be stoped. Afterwards, the *Iteration I* will be visited again and before that  $S, S', L$  are fetched from the previous result of *iteration I*. As shown in the Figure, after revisit the *iteration I* with the new  $max_b$ , there are some views that need to be executed (i.e.,  $V_{13}, V_{14}, V_{15}, V_{20}$ ) and the early termination start on  $V_{21}$ . If in this revisit, there is a view that can improve the  $S'$  from the previous result then  $F(S)$  will be updated, otherwise,  $F(S)$  is still same and go to the next iteration with the current  $max_b$ .

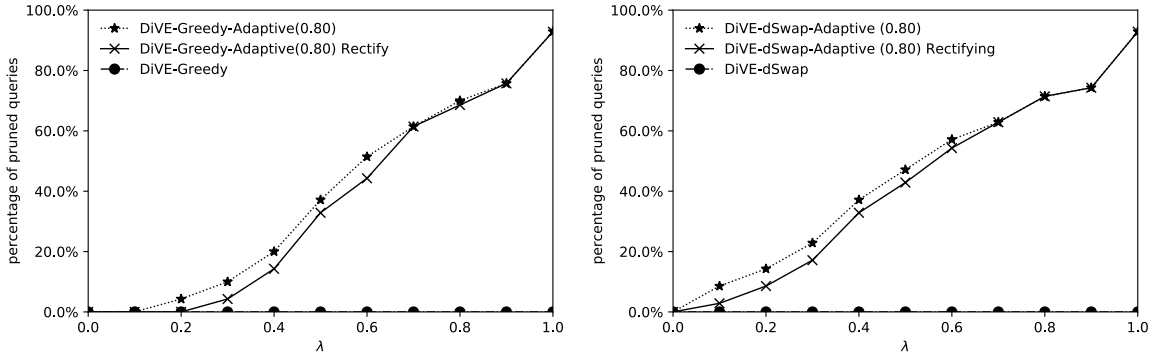


Figure 3: Rectifying vs. No Rectifying of adaptive scheme (PI80)

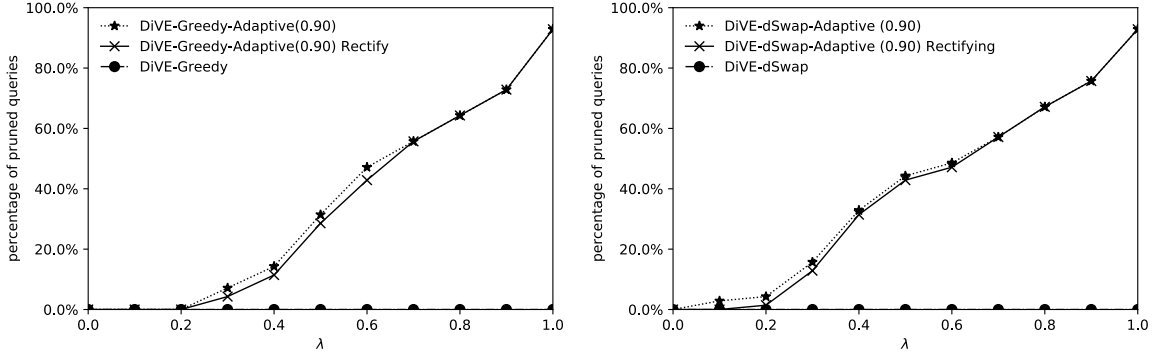


Figure 4: Rectifying vs. No Rectifying of adaptive scheme (PI90)

In the end of running, the number of executed and unexecuted queries are calculated. This algorithm uses book keeping strategy as shown in Figure 2.

The results of this rectifying bound strategy can be seen in Figure 3 to Figure 5. The results are average from five queries. These Figures show the performance of adaptive pruning scheme with rectifying bound strategy compared to without rectifying bound strategy. The pruning performance after applying rectifying bound strategy quite close to without rectifying bound strategy. Meanwhile, as shown in Figure 6 there is no effectiveness loss after rectifying bound is implemented especially for PI80. Moreover, the costs comparison of our rectifying algorithm is presented as well in Figure 7. This Figure shows the total cost of our pruning scheme with and without rectifying bound strategy which running on Flights dataset.

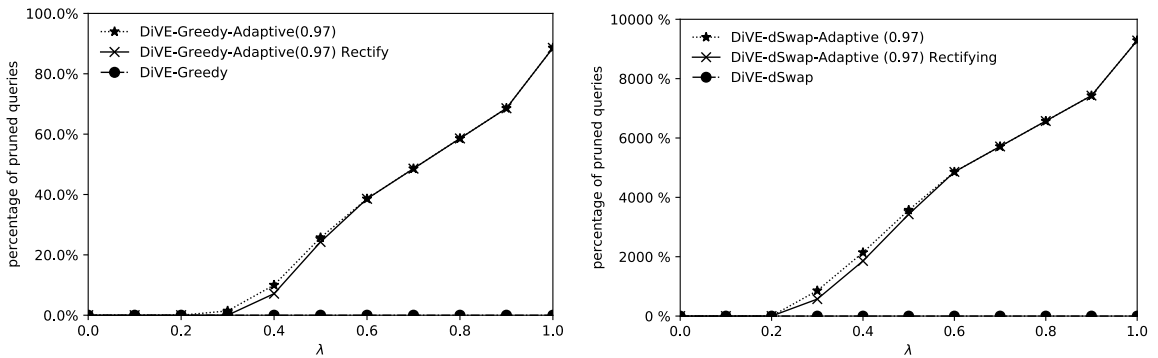


Figure 5: Rectifying vs. No Rectifying of adaptive scheme (PI97)

**Algorithm 1:** *DiVE* Greedy Pruning Rectifying

---

**Input:** Set of views  $V$  and result set Size  $k$   
**Output:** Result set  $S \leq V$ , size  $S = k$

```

1  $S \leftarrow$  two most distant views
2  $X \leftarrow [V \setminus S]$ 
3 function  $\text{getL}(f, S, X, L)$ :
4   for  $X_i$  in set  $X$  do
5     for  $S_j$  in set  $S$  do
6        $d \leftarrow \text{setDist}(X_i, S)$ 
7        $L.append([X_i, d])$ 
8    $L \leftarrow \text{sorted\_by\_d}(L)$ 
9   return  $L$ 
10  $max_b \leftarrow \text{getMaxPI}(L)$ 
11  $rectify \leftarrow \text{False}$ 
12
13 while  $i < k$  do
14   if  $rectify = \text{False}$  then
15      $L \leftarrow \text{getL}(S, X)$ 
16      $S' \leftarrow S \cup L[X_1]$ 
17   for  $L_i$  in  $L$  do
18     if  $rectify = \text{True}$  then
19        $\text{start loop at } L[\text{min}_d]$ 
20     if  $F(S') < F(S \cup X_i, max_b)$  then
21        $I \leftarrow \text{get\_I\_score}(X_i)$ 
22       if  $F(S') < F(S \cup X_i, I)$  then
23          $S' \leftarrow S \cup X_i$ 
24       if  $I > max_b$  then
25          $max_b \leftarrow I$ 
26          $rectify = \text{True}$ 
27          $\text{break}(\text{Out of Loop})$ 
28       else
29          $rectify = \text{False}$ 
30   if  $rectify == \text{True}$  then
31      $G \leftarrow \text{fetchTempResult}(i - 2)$ 
32      $S, S' \leftarrow G[S], G[S']$ 
33      $L \leftarrow G[L]$ 
34      $i = i - 2$ 
35   else
36      $\text{storeTempResult}(i, S, S', L, \text{min}_d)$ 
37      $S \leftarrow S'$ 
38      $i = i + 1$ 
39 return  $S$ 

```

---

**Algorithm 2:** *DiVE* dSwap Pruning Rectifying

---

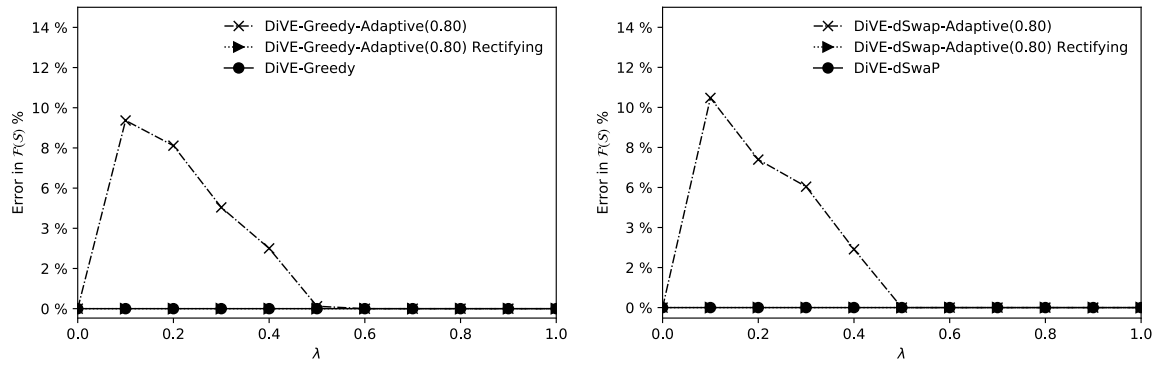
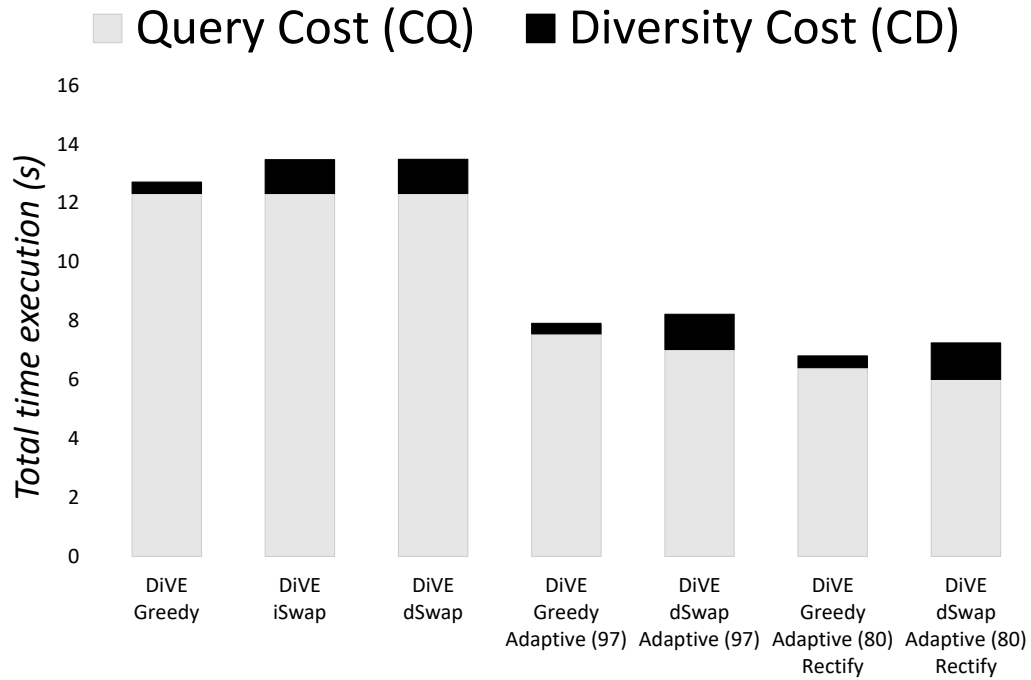
**Input:** Set of views  $V$  and result set Size  $k$   
**Output:** Result set  $S \leq V$ , size  $S = k$

```

1  $S \leftarrow$  Result set of only diversity
2  $X \leftarrow [V \setminus S]$ 
3 function  $\text{getL}(f, S, X, L)$ :
4   for  $X_i$  in set  $X$  do
5     for  $S_j$  in set  $S$  do
6        $d \leftarrow \text{setDist}(X_i, S)$ 
7        $L.\text{append}([S_j, X_i, d])$ 
8    $L \leftarrow \text{sorted\_by\_d}(L)$ 
9   return  $L$ 
10  $F_{\text{current}}, \text{counter} \leftarrow 0, 0$ 
11  $\text{improve}, \text{rectify} \leftarrow \text{True}, \text{False}$ 
12  $\text{max}_b \leftarrow \text{getMaxPI}(L)$ 
13
14 while  $\text{improve} = \text{True}$  do
15    $\text{counter} = \text{counter} + 1$ 
16   if  $\text{rectify} = \text{False}$  then
17      $L \leftarrow \text{getL}(S, X)$ 
18      $S' \leftarrow S$ 
19   for  $L_i$  in  $L$  do
20     if  $\text{rectify} = \text{True}$  then
21        $\text{start loop at } L[\text{min}_d]$ 
22     if  $F(S') < F(S \setminus S_j \cup X_i, \text{max}_b)$  then
23        $I \leftarrow \text{get\_I\_score}(X_i)$ 
24       if  $F(S') < F(S \setminus S_j \cup X_i, I)$  then
25          $S' \leftarrow S \setminus j \cup X_i$ 
26       if  $I > \text{max}_b$  then
27          $\text{max}_b \leftarrow I$ 
28          $\text{rectify} = \text{True}$ 
29          $\text{break}(\text{Out of Loop})$ 
30     else
31        $\text{rectify} = \text{False}$ 
32   if  $\text{rectify} == \text{True}$  then
33      $G \leftarrow \text{fetchTempResult}(\text{counter} = 1)$ 
34      $S, S' \leftarrow G[S], G[S']$ 
35      $L \leftarrow G[L]$ 
36      $\text{counter} = 0$ 
37      $\text{improve} \leftarrow \text{True}$ 
38   else
39      $\text{storeTempResult}(\text{counter}, S, S', L, \text{min}_d)$ 
40     if  $F(S') > F(S)$  then
41        $S \leftarrow S'$ 
42     if  $F(S) > F_{\text{current}}$  then
43        $F_{\text{current}} \leftarrow F(S)$ 
44        $\text{improve} \leftarrow \text{True}$ 
45     else
46        $\text{improve} \leftarrow \text{False}$ 
47 return  $S$ 

```

---

Figure 6: Error  $F(S)$  after rectifying boundFigure 7: Total costs of schemes running on Flights dataset,  $k = 5$ , and  $\lambda = 0.5$