

Community Search 实验报告

丁霄汉 2017312365 吴超月 2017213865 罗 瑶 2017213866

所选论文： Akbas Esra, Zhao Peixiang.Truss-based Community Search: a Truss-equivalence Based Indexing Approach. Proceedings of the VLDB Endowment, 10(11):1298-1309,2017.

目录

- Community Search 实验报告..... 1
- 1. 相关工作介绍..... 2
 - 1.1 Community Search..... 2
 - 1.2 Truss..... 2
 - 1.3 Graph Summarization..... 2
 - 1.4 TCP-Index..... 2
- 2. 问题定义..... 3
 - 2.1 图的类型..... 3
 - 2.2 社区模型..... 3
- 3. 算法详细分析..... 3
- 4. 实验结果..... 5
 - 4.1 索引构建..... 6
 - 4.2 社区搜索..... 6
 - 4.2.1 变化点的度数..... 6
 - 4.2.2 变化 truss k..... 7
 - 4.3 质量评估..... 8
- 5. 算法的优势和局限..... 9
 - 5.1 优势..... 9
 - 5.2 局限..... 9
- 参考文献..... 10

1. 相关工作介绍

1.1 Community Search

作为众所周知的社区检测问题的一个查询依赖变体,社区搜索就是在图中找到涉及给定查询顶点(或一组查询顶点)的内聚和密集连接的子图。**Community Search** 这个问题一开始被探索时被建模成具有距离和大小约束的 **k-core** 问题,使得这个问题成了 **NP 难**[29]。后续有进一步提出用 γ -quasi-k-cliques 方式在线搜索查询顶点的重叠社区[7],但是得到的社区可能不会内聚连接。同时也有工作做影响力社区搜索,致力于查询 **top-r** 最有影响力的 **k-core** 社区[21]。**Local Community Search**[34,1,7]是识别包含查询节点的 **k-core** 社区,同时也是最大化或者最小化某些指标。但是,这难以避免搭便车现象,并且也只能得到一个社区。同时也有工作做属性图[9,12]和空间图[8]的社区搜索。

1.2 Truss

密集的,有凝聚力的子图是揭示真实世界的大量图的潜在社区结构的关键组件。密集和凝聚图的建模和量化方面有丰富的文献,包括 clique 或 quasi-clique[6,31], **k-core**[25,1,15]和 nucleus[26,27]。所有上述模型除 **k-core** 外都受到计算困难问题的困扰[28],而 **k-core** 可能导致不一致的子图[37,5]。**Truss** 是在高阶图形三角形上定义的,在社区建模和计算中享有众多优势[5]。**Truss** 具体定义可以参见论文。

1.3 Graph Summarization

随着图的规模和复杂度的增加,图概括技术已经被用于将大量图简化为足够的和保存质量的总结,从而显着减少图的存储和计算成本[17]。**GraSS** [19]通过贪婪地将顶点分组成一个概率邻接矩阵来对图进行总和,在这个矩阵上邻域查询可以被有效地近似。在[24]中,图形总结为超级节点和超级边缘,保证重构误差。压缩的摘要用于近似查询,包括邻接度,度数,特征向量中心性和子图计数。在[22]中,作者设计了贪心和随机算法来压缩有界最小描述长度(MDL)错误的图。**VOD** [17]是一种基于词汇的图表汇总方法,旨在最大限度地减少图的信息论编码成本。**SNAP** [30]根据属性对顶点进行分组,迭代地分组,直到达到最大的属性和关系兼容的分组。图汇总是问题驱动的,并且通常针对与应用程序相关的目标进行优化。但是,本文所讨论的社区搜索问题并没有先验的图的概括方法。目前还没有研究探索 **k-Truss** 社区边缘的内在关系的工作。

1.4 TCP-Index

基于 **Truss** 的社区搜索的最先进的解决方案是 **TCP-Index** (三角连接预服务索引) [11]。本文主要也是与 **TCP-Index** 比较。

2. 问题定义

2.1 图的类型

本文关注的网络图种类为无向无权无属性图。

2.2 社区模型

本篇工作是基于 **truss** 做的，由于 **k-truss** 得到的图可能不是连通图，所以在 **k-truss** 上加了 **triangle-connectivity** 约束来定义 **k-truss Community**。本文的社区搜索问题就是找到搜索节点的所有 **k-truss Community**。此外，如果网络图存在边的增加和删除操作时，文章给出的方法可以进行有效的更新。

3. 算法详细分析

算法的总体思路是在原始图 **G** 上建立一个索引 **EquiTruss**，然后 **k-truss Community** 的搜索可以只在索引上搜索，而不用访问原图 **G**。这涉及到两个部分：建立索引的过程和在索引上搜索的过程。

1. 建立索引。

首先介绍几个概念：

- **Edge Trussness**: 包含这条边的子图具有的最大的 **truss**;
- **k-triangle**: 三条边的 **Edge Trussness** 都大于等于 **k** 的三角形;
- **k-triangle connectivity**: 两个三角形由多个 **k-triangle** 连接到一起;
- **k-truss equivalence**: 两条边 **Edge Trussness** 均为 **k**, 且两条边所在的三角形 **k-triangle connectivity**。

我们根据 **k-truss equivalence** 等价关系来建立索引，**G** 的每条边可以唯一的映射到 **index** 的一个等价类中。然后我们可以把索引定义为图 $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ ，其中节点 $v \in \mathcal{V}$ ， v 代表一个等价类，而边 $(\mu, v) \in \mathcal{E}$ ，代表两个等价类中各存在一条边，这两条边具有 **triangle connectivity** 的关系。

文中的算法一首先计算图 **G** 中每条边的 **Edge Trussness**，**Edge Trussness** 是利用 **edge support**（每条边被多少个三角形包含）来计算的，**k** 初始化为 2，根据 **edge support** 从小到大遍历每条边，将 **k** 赋给当前 **edge support** 最小的边，然后将与此条边构成三角形的所有边的 **edge support** 减一，然后将边重新排序。计算 **edge support** 的时间复杂度为： $O(|E|^{1.5})$ ，而总的时间复杂度为： $O(|E|^{1.5})$ ，空间复杂度为： $O(|V| + |E|)$ 。算法具体如下：

Algorithm 1: Truss Decomposition

Input: A graph $G(V_G, E_G)$ **Output:** The edge trussness $\tau(e)$ for each $e \in E_G$

```
1 Computer  $\text{sup}(e)$  for each edge  $e \in E_G$ ;
2 Sort all edges in the non-decreasing order of their support;
3  $k \leftarrow 2$ ;
4 while  $\exists e \in E_G, \text{sup}(e) \leq (k - 2)$  do
5    $e^*(u, v) \leftarrow \arg \min_{e \in E_G} \text{sup}(e)$ ;
6   assume w.l.o.g.  $d(u) \leq d(v)$ ;
7   foreach  $w \in N(u)$  and  $(v, w) \in E_G$  do
8      $\text{sup}(u, w) \leftarrow \text{sup}(u, w) - 1$ ;
9      $\text{sup}(v, w) \leftarrow \text{sup}(v, w) - 1$ ;
10    Reorder  $(u, w)$  and  $(v, w)$  w.r.t. new edge support;
11    $\tau(e^*) \leftarrow k$ , remove  $e^*$  from  $E_G$ ;
12 if  $\exists e \in E_G$  then
13    $k \leftarrow k + 1$ ;
14   goto Step 4;
15 return  $\{\tau(e) | e \in E_G\}$ ;
```

在计算了所有边的 Edge Trussness 之后，我们可以将边按照 Edge Trussness 的值分类，Edge Trussness 大小相同的边在同一个类里面。进一步的，我们可以计算 index，也就是得到图 $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ 。我们按照 k 从小到大遍历每个类里的每条边。我们利用一条边初始化节点 $v \in V$ 并将其压入队列 Q ，每次从 Q 中取一条边 e 加入 v ，我们访问与边 e 构成三角形的所有边 e' ，如果这条边的 Edge Trussness 为 k ，我们将其加入队列，如果小于 k ，不用处理，如果大于 k ，将 v 的 id 加入 e' 的 list，用于建立 \mathcal{E} 。算法二的时间复杂度为： $O(|E|^{1.5})$ ，空间复杂度为： $O(|E|)$ 。具体算法如下：

Algorithm 2: Index Construction for EquiTruss

Input: A graph $G(V_G, E_G)$ **Output:** EquiTruss: $\mathcal{G}(\mathcal{V}, \mathcal{E})$

/* Initialization */

```
1 Truss Decomposition ( $G$ );
2 foreach  $e \in E_G$  do
3    $e.\text{processed} \leftarrow \text{FALSE}$ ;
4    $e.\text{list} \leftarrow \emptyset$ ;
5   if  $\tau(e) = k$  then
6      $\Phi_k \leftarrow \Phi_k \cup \{e\}$ ;
7 snID  $\leftarrow 0$ ; /* Super-node ID initialized to 0 */
/* Index Construction */
8 for  $k \leftarrow 3$  to  $k_{max}$  do
9   while  $\exists e \in \Phi_k$  do
10     $e.\text{processed} = \text{TRUE}$ ;
11    Create a super-node  $\nu$  with  $\nu.\text{snID} \leftarrow ++\text{snID}$ ;
12     $\mathcal{V} \leftarrow \mathcal{V} \cup \{\nu\}$ ; /* A new super-node for  $\mathcal{C}_e$  */
13     $Q.\text{enqueue}(e)$ ;
14    while  $Q \neq \emptyset$  do
15       $e(u, v) \leftarrow Q.\text{dequeue}()$ ;
16       $\nu \leftarrow \nu \cup \{e\}$ ; /* Add  $e$  to super-node  $\nu$  */
17      foreach  $id \in e.\text{list}$  do
18        Create a super-edge  $(\nu, \mu)$  where  $\mu$  is an
        existing super-node with  $\mu.\text{snID} = id$ ;
19         $\mathcal{E} \leftarrow \mathcal{E} \cup \{(\nu, \mu)\}$ ; /* Add super-edge */
20      foreach  $w \in N(u) \cap N(v)$  do
21        if  $\tau(u, w) \geq k$  and  $\tau(v, w) \geq k$  then
22          ProcessEdge( $u, w$ );
23          ProcessEdge( $v, w$ );
24       $\Phi_k \leftarrow \Phi_k - \{e\}$ ;  $E \leftarrow E - \{e\}$ ;
25 return  $\mathcal{G}(\mathcal{V}, \mathcal{E})$ ;
```

```

26 Procedure ProcessEdge( $u, v$ )
27 if  $\tau(u, v) = k$  then /*  $k$ -triangle connectivity */
28   if  $(u, v).processed = \text{FALSE}$  then
29      $(u, v).processed = \text{TRUE}$ ;
30      $Q.enqueue(u, v)$ ;
31 else /*  $\tau(u, v) > k$  */
32   if  $snID \notin (u, v).list$  then
33      $(u, v).list \leftarrow (u, v).list \cup \{snID\}$ ;

```

2. 利用索引来搜索 k -truss community，算法 3 是社区搜索的过程，非常简单且直观，毕竟我们已经得到了等价类，作者直接选择用证明来说明算法的正确性。在算法 3 中，时间复杂度仅仅被搜索到的社区集合的大小决定，具体算法如下：

Algorithm 3: Community Search Based on EquiTruss

Input: $\mathcal{G}(\mathcal{V}, \mathcal{E})$, the truss value $k \geq 3$, the query vertex q
Output: \mathcal{A} : all k -truss communities containing q

```

/* Initialization */
1 foreach  $v \in \mathcal{V}$  do
2    $v.processed \leftarrow \text{FALSE}$ ;
3  $l \leftarrow 0$ ;
/* BFS traversal for community search */
4 foreach  $v \in \mathcal{H}(q)$  do
5   if  $\tau(v) \geq k$  and  $v.processed = \text{FALSE}$  then
6      $v.processed \leftarrow \text{TRUE}$ ;
7      $l \leftarrow l + 1$ ;  $\mathcal{A}_l \leftarrow \emptyset$ ;
8      $Q \leftarrow \emptyset$ ;  $Q.enqueue(v)$ ;
9     while  $Q \neq \emptyset$  do
10       $v \leftarrow Q.dequeue()$ ;
11       $\mathcal{A}_l \leftarrow \mathcal{A}_l \cup \{e | e \in v\}$ ;
12      foreach  $(v, \mu) \in \mathcal{E}$  do
13        if  $\tau(\mu) \geq k$  and  $\mu.processed = \text{FALSE}$  then
14           $\mu.processed \leftarrow \text{TRUE}$ ;
15           $Q.enqueue(\mu)$ ;
16 return  $\{\mathcal{A}_1, \dots, \mathcal{A}_l\}$ ;

```

4. 实验结果

我们用 Java 语言实现了上述算法，并在四个数据集上进行了实验。

实验环境：

处理器: Intel(R) Core(TM) i7-7700 CPU @ 3.60GHz 3.60 GHz
 已安装的内存(RAM): 16.0 GB (15.9 GB 可用)
 系统类型: 64 位操作系统，基于 x64 的处理器

```

C:\Users\Administrator>java -version
java version "9.0.1"
Java(TM) SE Runtime Environment (build 9.0.1+11)
Java HotSpot(TM) 64-Bit Server VM (build 9.0.1+11, mixed mode)

```

数据集：

四个来自于真实世界的数据集，数据集的基本情况见下表。其中 d_{\max} 表示图中点的最大度数。

Network	$ V_G $	$ E_G $	d_{\max}
Amazon	334863	925872	549

DBLP	317080	1049866	343
YouTube	1134890	2987624	28754
LiveJournal	399762	34681189	14815

4.1 索引构建

实验从对给定图建立索引开始，这个过程通常在进行社区搜索前完成。我们将建立好的索引保留在主存中，来作为后续进行大型图表社区搜索的有效工具。在实验中主要对索引构建两个方面的指标进行衡量：（1）索引构建的时间（2）构建的索引的大小。相关实验结果如下表所示：

Graph	Graph size (MB)	Index space (MB)	Construction time (Sec.)
Amazon	12	6.4	1.5
DBLP	13.2	6.92	2
YouTube	36.9	18.3	8.4
LiveJournal	478	312	366.3

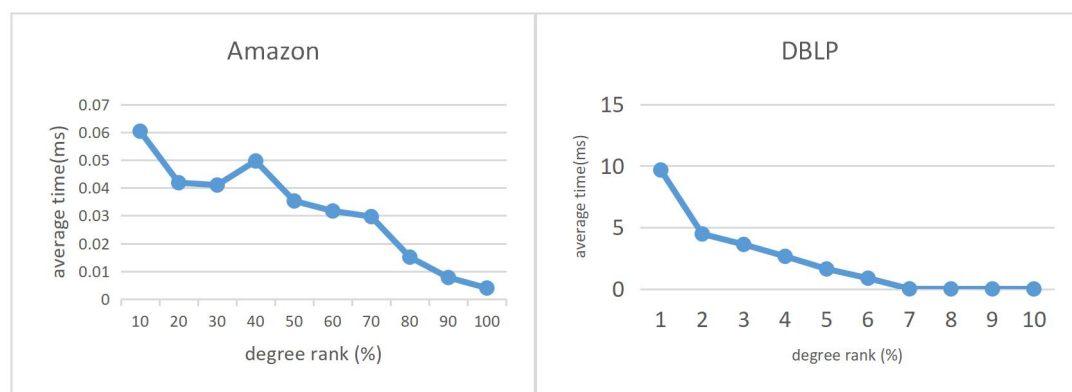
从表中可以看出，随着图大小的增加，构建的索引大小和所需时间也随之增加，但是索引的大小均比原图要小，这是因为在索引中没有存储冗余的信息。

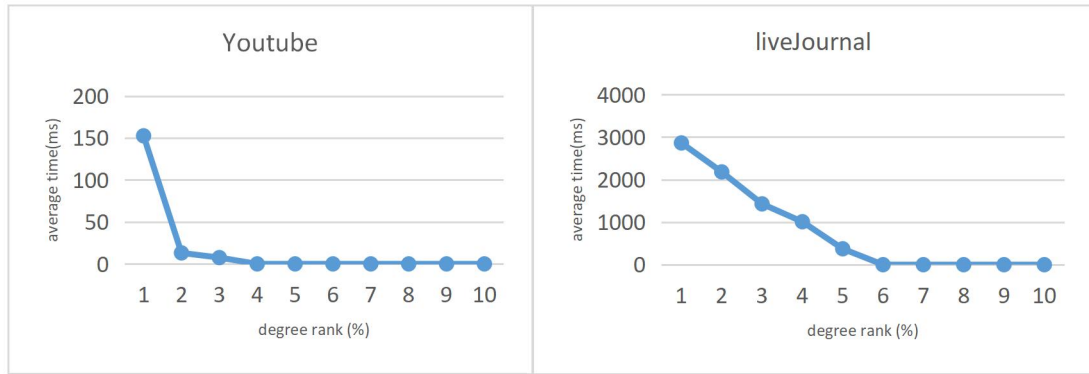
4.2 社区搜索

一旦建立好图对应的索引后，就可以利用这些索引进行社区搜索。在实验中考虑两种不同的设置：（1）查询具有不同的度数的点，即考察点的度数对社区搜索运行效率的影响情况（2）考察所设置的 $\text{truss } k$ 的不同对社区搜索运行效率的影响情况。

4.2.1 变化点的度数

对每个数据集的点按照度数降序排序，并按照度数将这些点分成十份，使得第一份数据中的点拥有最高的度数。然后从十份数据中各随机选择 100 个点进行社区搜索，并计算每份数据的平均搜索时间。实验在 Amazon、DBLP、YouTube、LiveJournal 四个数据集上进行，其中 Amazon 数据集中设置的 $\text{truss } k$ 为 4，DBLP 中的 $\text{truss } k$ 为 5，YouTube 的 $\text{truss } k$ 设置为 4，LiveJournal 的 $\text{truss } k$ 为 6。实验结果如下所示。

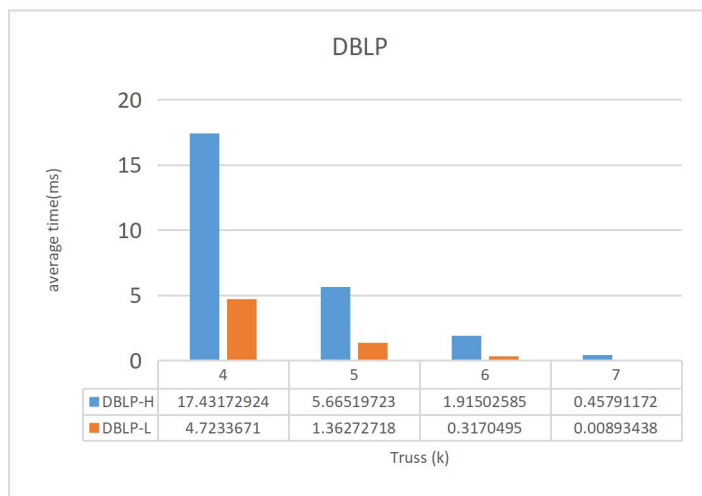
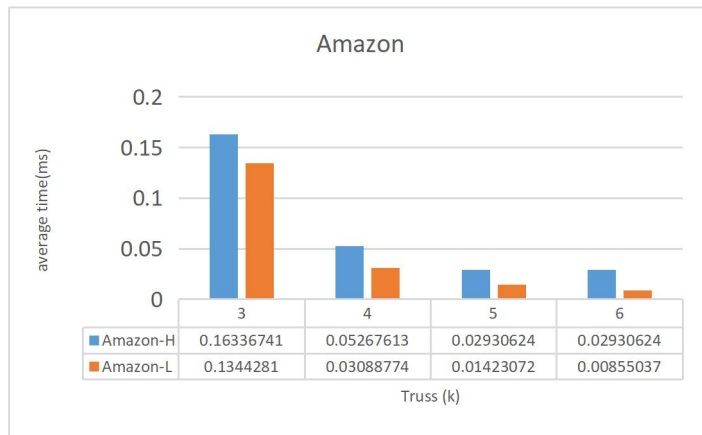


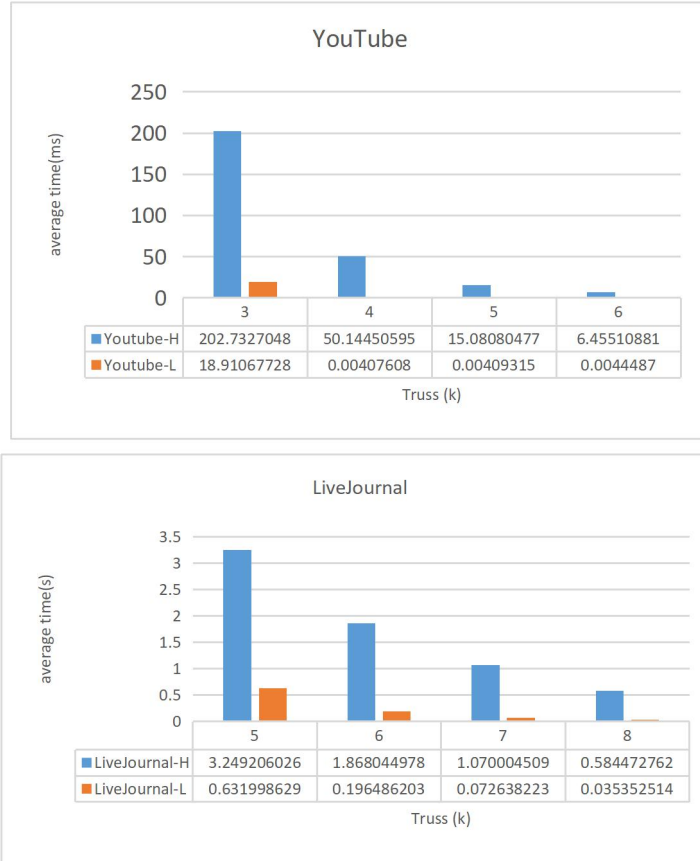


从折线图中可以看出，对拥有较高度数的点进社区搜索所需的时间要明显高于度数较低的点。特别的，在 **YouTube** 数据集中，随着点度数的降低，所需的时间迅速减少，这是因为对于那些度数少的点来说，基本上很少存在对应的 **k-truss** 社区。

4.2.2 变化 truss k

对每个数据集的点按照度数降序排序，并按照度数将这些点分成十份，使得第一份数据中的点拥有最高的度数。然后从中选取两个待测数据集：（1）从度数在前 30%的点中任意选取 100 个点（2）从后 70%的点中任意选取 100 个点，对这些点进行社区搜索并计算每份数据的平均搜索时间。实验在 **Amazon**、**DBLP**、**YouTube**、**LiveJournal** 四个数据集上进行，实验结果如下所示。





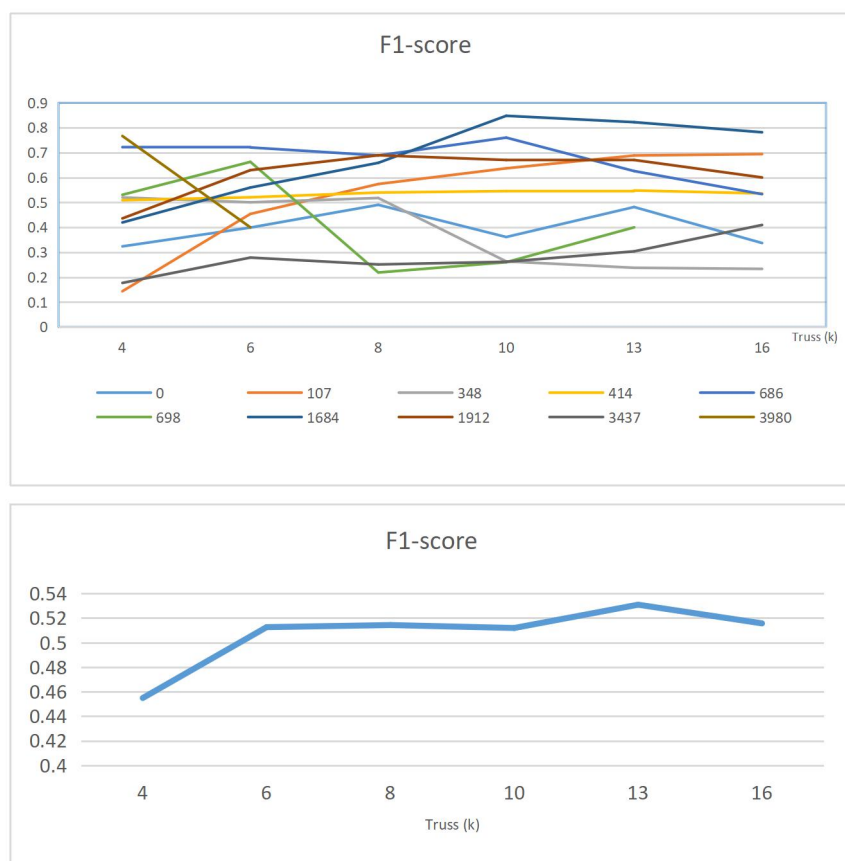
从上面的数据图表中可以看出，随着 **truss k** 的增加，所需要的搜索时间越来越少，这种情况在大型图中的搜索边的尤其显著。并且在每组实验中，拥有较高度数的点对应的社区搜索时间也越长，这与我们上一个实验的结论保持一致。

4.3 质量评估

在评价 **k-truss community model** 的有效性的时候，我们采用了 **Facebook** 数据集，因为这个数据集无向无环无属性，并且为单个点提供 **ground truth**，**Facebook** 数据集为 10 个点分别提供了它们所在的所有社区，我们可以直接利用我们得到的社区集合与点本身的社区集合计算 **F1-Score**，社区集合的 **F1-Score** 计算公式为：

$\max_{f: \mathcal{C} \rightarrow \hat{\mathcal{C}}} \frac{1}{|\mathcal{C}|} \sum_{C \in \text{dom}(f)} F_1(C, f(C))$ 。下面分别

展示在 **Facebook** 数据集上进行社区搜索的 10 个点（用点的 **index** 表示）不同 **k** 下各自和平均的 **F1-score**。可以看出总体的 **F1-score** 在 0.45 以上，随着 **k** 的不同，**F1-score** 基本保持在一个稳定的数值，浮动较小。



5. 算法的优势和局限

5.1 优势

首先本文介绍的算法证明完整，理论扎实，让人信服。其次这是一个确定性的算法，不像有些解决方案属于近似算法。其实在建立索引后，节点的社区搜索时间复杂度非常低，这适合与图更新较少且搜索较多的情况。并且 k 值是可以设置的，这在宽泛意义上使得使用者对得到社区的紧密程度可调。此外，如果网络图存在边的增加和删除操作时，文章给出的方法较先前的工作可以进行有效的更新。

5.2 局限

本篇工作查询节点的时候需要设置 k 值，虽然作者解释说 k 值的大小一定程度上反映着得到的社区的紧密程度，这的确在更宽泛的意义上是 **make sense** 的，但是当我们将这个算法用于真正的社区搜索的时候，我们发现 k 值会直接影响搜索的好坏程度，并且针对每个节点的好的 k 值都可以不一样，好的 k 值也没有肉眼可见的规律可调，这使得这个算法的真实的使用性降低了。

参考文献

- [1] N. Barbieri, F. Bonchi, E. Galimberti, and F. Gullo. Efficient and Effective Community Search. *Data Min. Knowl. Discov.*, 29(5):1406–1433, 2015.
- [5] J. Cohen. Trusses: Cohesive Subgraphs for Social Network Analysis. NSA:Technical report, 2008.
- [6] W. Cui, Y. Xiao, H. Wang, Y. Lu, and W. Wang. Online Search of Overlapping Communities. In *SIGMOD’13*, pages 277–288, 2013.
- [7] W. Cui, Y. Xiao, H. Wang, and W. Wang. Local Search of Communities in Large Graphs. In *SIGMOD’14*, pages 991–1002, 2014.
- [8] Y. Fang, R. Cheng, X. Li, S. Luo, and J. Hu. Effective Community Search over Large Spatial Graphs. *Proc. VLDB Endow.*, 10(6):709–720, 2017.
- [9] Y. Fang, R. Cheng, S. Luo, and J. Hu. Effective Community Search for Large Attributed Graphs. *Proc. VLDB Endow.*, 9(12):1233–1244, 2016.
- [11] X. Huang, H. Cheng, L. Qin, W. Tian, and J. X. Yu. Querying K-truss Community in Large and Dynamic Graphs. In *SIGMOD’14*, pages 1311–1322, 2014.
- [12] X. Huang and L. V. S. Lakshmanan. Attribute truss community search. *Proc. VLDB Endow.*, 10(9):949 – 960, 2017.
- [15] W. Khaouid, M. Barsky, V. Srinivasan, and A. Thomo. K-core Decomposition of Large Networks on a Single PC. *Proc. VLDB Endow.*, 9(1):13–23, 2015.
- [17] D. Koutra, U. Kang, J. Vreeken, and C. Faloutsos. Summarizing and Understanding Large Graphs. *Stat. Anal. Data Min.*, 8(3):183–202, 2015.
- [19] K. LeFevre and E. Terzi. GraSS: Graph Structure Summarization. In *SDM*, 2010.
- [21] R.-H. Li, L. Qin, J. X. Yu, and R. Mao. Influential Community Search in Large Networks. *Proc. VLDB Endow.*, 8(5):509–520, 2015.
- [22] S. Navlakha, R. Rastogi, and N. Shrivastava. Graph Summarization with Bounded Error. In *SIGMOD’08*, pages 419–432, 2008.
- [24] M. Riondato, D. García-Soriano, and F. Bonchi. Graph Summarization with Quality Guarantees. *Data Min. Knowl. Discov.*, 31(2):314–349, 2017.
- [25] A. E. Sariyuce, B. Gedik, G. Jacques-Silva, K.-L. Wu, and U. V. Catalyurek. Incremental K-core Decomposition: Algorithms and Evaluation. *The VLDB Journal*, 25(3):425–447, 2016.

- [26] A. E. Sariyüce and A. Pinar. Fast Hierarchy Construction for Dense Subgraphs. *Proc. VLDB Endow.*, 10(3):97–108, 2016.
- [27] A. E. Sariyüce, C. Seshadhri, A. Pinar, and U. V. Catalyürek. Finding the Hierarchy of Dense Subgraphs Using Nucleus Decompositions. In *WWW’15*, pages 927–937, 2015.
- [28] Y. Shao, L. Chen, and B. Cui. Efficient Cohesive Subgraphs Detection in Parallel. In *SIGMOD’14*, pages 613–624, 2014.
- [29] M. Sozio and A. Gionis. The Community-search Problem and How to Plan a Successful Cocktail Party. In *KDD’10*, pages 939–948, 2010.
- [30] Y. Tian, R. A. Hankins, and J. M. Patel. Efficient aggregation for graph summarization. In *SIGMOD’08*, pages 567–580, 2008.
- [31] C. Tsourakakis, F. Bonchi, A. Gionis, F. Gullo, and M. Tsiarli.

Denser Than the Densest Subgraph: Extracting Optimal Quasi-cliques with Quality Guarantees. In *KDD’13*, pages 104–112, 2013.
- [34] Y. Wu, R. Jin, J. Li, and X. Zhang. Robust Local Community Detection: On Free Rider Effect and Its Elimination. *Proc. VLDB Endow.*, 8(7):798–809, 2015.
- [37] F. Zhao and A. K. H. Tung. Large Scale Cohesive Subgraphs Discovery for Social Network Visual Analysis. *Proc. VLDB Endow.*, 6(2):85–96, 2012.