

# Mining Statistically Significant Attribute Associations in Attributed Graphs

Jihwan Lee  
Department of Computer Science  
Purdue University  
West Lafayette, IN  
jihwan@purdue.edu

Keehwan Park  
Department of Computer Science  
Purdue University  
West Lafayette, IN  
park451@purdue.edu

Sunil Prabhakar  
Department of Computer Science  
Purdue University  
West Lafayette, IN  
sunil@purdue.edu

**Abstract**—Recently, graphs have been widely used to represent many different kinds of real world data or observations such as social networks, protein-protein networks, road networks, and so on. In many cases, each node in a graph is associated with a set of its attributes and it is critical to not only consider the link structure of a graph but also use the attribute information to achieve more meaningful results in various graph mining tasks. Most previous works with attributed graphs take into account attribute relationships only between individually connected nodes. However, it should be greatly valuable to find out which sets of attributes are associated with each other and whether they are statistically significant or not. Mining such significant associations, we can uncover novel relationships among the sets of attributes in the graph. We propose an algorithm that can find those attribute associations efficiently and effectively, and show experimental results that confirm the high applicability of the proposed algorithm.

## I. INTRODUCTION

Nowadays graphs have emerged as a powerful abstract data type to represent and analyze complex data in a broad range of commercial and scientific applications including social networks [1], [2], bioinformatics [3], world wide web [4], [5], and so on. Mining structured patterns in graphs have been actively studied in the literature and such patterns including cliques [6], subgraphs [7], [8], [9], paths [10] and trees [11] help us better understand the intrinsic characteristics of graph data. Also, when the graph data come with auxiliary information such as node attributes, such information can be applied to various application areas, e.g., community detection, link prediction, graph clustering, network modeling, and etc. Thus, attributed graphs are more important than ever before to complex mining tasks.

While node attributes can be successfully employed to augment various mining tasks, the node attributes themselves could give us interesting patterns for better understanding graphs. Given an attributed graph where each node is associated with its attribute values, one might be interested in a pattern of node attribute values which co-occur between connected nodes. Let's call such co-occurred attribute values between two connected nodes an attribute association. This information can tell us directly the attribute patterns shared by connected nodes over the entire graph. In large scale, one might be interested in which attribute associations are most frequently observed or which attribute vector is most

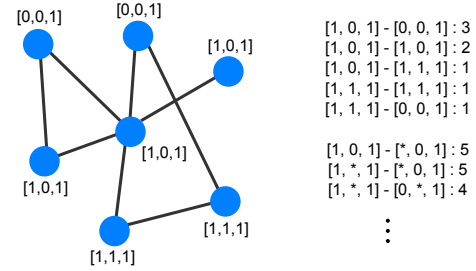


Fig. 1: Attribute associations in attributed graph

expected to be observed given another attribute vector in attribute associations. Looking at frequent attribute associations reveals the most dominant attribute associations in the graph by simply taking into account how many times they are held by connected nodes. Even though the frequent attribute associations give us which ones are dominant over the entire graph, they do not tell us which ones are really significant. That is because the frequency of an attribute association often does not depart from what we expect and therefore may not be meaningful actually if we already know the distributions of attribute values in the graph. Rather, identifying the statistically significant attribute associations where the pattern of the attribute association deviates from the expected, can potentially infer undiscovered possible relationships between nodes in the graph. The statistical significance of a pattern has been emphasized in various data mining problems [12], [13], [7], [14], [15] and the previous works already explored why a statistically significant pattern is more important rather than a frequent pattern. Thus, in this paper we define a statistically significant attribute association and address the problem of uncovering it in attributed graphs.

Fig. 1 shows an example that shows a list of possible attribute associations in an attribute graph. An attribute association is *frequent* if the number of pairs of nodes is above a given threshold which is determined by *freq\_support*. Unfortunately the frequency is not sufficient to measure the statistical significance of an attribute association since the frequency eventually depends on the actual distributions of the attribute values in the graph. We will closely see the set

difference between the two in Section V. Also when obtaining significant associations, each attribute value does not always have to take discrete attribute value, e.g., 0 or 1 in binary case, as long as the association has enough statistical significance. Accordingly, we introduce wildcard attribute notation (\*), which matches any value of the corresponding attribute.

The statistical significance of an attribute association with its frequency  $k$  is determined by the probability that it is observed at least  $k$  times or more, and the probability is called the  $p$ -value of the attribute association. By measuring  $p$ -value, we can identify the significant ones even though they are not frequent absolutely in the graph. Also, as shown in Fig. 1, we are interested in even associations of partial attribute values as long as they are statistically significant. The main challenge of the problem is how to estimate the probability that an attribute association occurs in a random graph. There are as many different attribute associations as the number of edges in a graph, and if we consider even the partial attribute associations then the number of possible attribute associations grows exponentially. We address the challenge by transforming a graph  $G$  into an alternative graph  $\mathcal{AG}$ , called *association graph*, where each vertex contains a subset of nodes in  $G$  that have the same or similar attribute values and each edge corresponds to a certain attribute association between two set of attribute values, each of which is represented by a cluster. During the process of transformation, we build  $\mathcal{AG}$  such that the edges (i.e., associations) are statistically significant.

To experimentally evaluate our work, we use two real world attributed graphs. One is the *DBLP co-authorship network* and the other is the *Yelp social network*. We present the statistically significant attribute associations extracted from the graphs and compare them against the frequent attribute associations qualitatively. In addition to that, we show quantitatively how the statistically significant attribute associations can be used for boosting the performance of the link prediction task.

We summarize the contributions of our work as follows:

- We formally define the novel problem of mining statistically significant attribute associations which aims to find patterns of co-occurred attribute values between nodes which deviate from the expected.
- We design and implement an algorithm that can find the statistically significant attribute associations efficiently and effectively.
- We conduct experiments using real world attributed graphs and show qualitative results as well as the actual application that can benefit from the results.

The paper is organized as follows. In Section II, we introduce previous works related to our problem and discuss how our problem differs from them. In Section III, we define the problem of mining statistically significant attribute association and provide basic background concepts. The novel algorithm to solve our problem is discussed in Section IV and we present our experimental findings in Section V. Finally, we conclude the paper in Sections VI.

## II. RELATED WORK

There are a number of previous works that have explored the statistical significance of patterns in various data mining and knowledge discovery tasks and have proposed efficient methods for mining the statistically significant patterns. [14], [7] study the statistical significance of subgraphs where the nodes of the graph are labeled. [14] addresses the problem of finding statistically significant connected subgraphs in a vertex-labeled graph where the labels are discrete and continuous. The statistical significance is quantified by using the *chi-square* statistic, which makes the naïve algorithm impractical because of the exponential number of subgraphs. They propose an efficient algorithm which converts the graph into a super-graph. In [7], the authors propose a technique for computing the statistical significance of frequent subgraphs in a graph database. In order to solve the difficulty of estimating the  $p$ -value of a subgraph directly in the graph space due to the flexible structures of graphs, they transform graphs into a feature space with predefined set of basis elements, and then approximate the significance of a feature vector in the feature space by using the binomial distribution. Although these two works explore the statistically significant patterns in graphs, they differ from our work in that they more focus on structured patterns, not attribute association patterns.

In addition to graphs, the statistical significance has been studied for other types of patterns as well. [12] extends the traditional association rule mining problem to searching statistically significant association rules such that some spurious rules are not included in the result set while considering statistical dependence. The significance of the observed frequency of an association rule is estimated by the binomial distribution. [15] solves the problem of mining statistically significant substrings in a string generated from a memoryless Bernoulli distribution and uses the chi-square statistic as a quantitative measure of statistical significance. The statistical significance is considered for the sequential pattern mining problem as well in [16]. The approach developed by the authors is able to efficiently mine unexpected patterns in sequence of itemsets without considering overlapping occurrences or conditioning the length of the sequence.

## III. PROBLEM STATEMENT

In this section, we give basic definitions of the attribute association, frequent association, statistically significant association, and define the problem of mining statistically significant attribute associations. Table I introduces the notations we use throughout the paper.

### A. Attribute Associations

Suppose we have an attributed graph  $G = (V, E, A)$  where  $V = \{u_1, u_2, \dots, u_{|V|}\}$  is a set of nodes,  $E = V \times V$  is a set of edges, and  $A = \{\vec{a}_{u_1}, \vec{a}_{u_2}, \dots, \vec{a}_{u_{|V|}}\}$  is a set of attribute vectors, each of which is associated with a node in  $V$ . The attribute vector  $\vec{a}_u$  of the node  $u$  that holds  $l$  different attributes is represented by a vector of  $l$  binary values in that each binary indicates whether the node  $u$  actually has a value for the

Notation	Meaning
$G = (V, E)$	attributed graph
$V = \{u_1, u_2, \dots, u_{ V }\}$	set of nodes in $G$
$E$	set of edges in $G$
$\mathcal{AG} = (\mathcal{V}, \mathcal{E})$	association graph
$\mathcal{V} = \{c_1, c_2, \dots, c_{ \mathcal{V} }\}$	set of clusters in $\mathcal{AG}$
$\mathcal{E}$	set of attribute associations in $\mathcal{AG}$
$\vec{a} = (a^1, a^2, \dots, a^l)$	attribute vector of size $l$
$\Delta$	attribute association
$\sigma$	<i>freq_support</i>
$\lambda$	<i>size_support</i>
$\delta_G$	density of graph $G$
$\Psi_c$	<i>p-value</i> of cluster $c$
$TS(u, v)$	tie-strength between node $u$ and $v$
$\Gamma(\cdot)$	set of neighbors
$\tilde{G}_c$	subgraph of nodes within cluster $c$

TABLE I: Basic notations

corresponding attribute (in case of an  $m$  multi-valued attribute, it can be transformed into  $m - 1$  dichotomous variables each with binary). Then we define an *attribute association* between a pair of attribute vectors  $\vec{a}_1$  and  $\vec{a}_2$  as follows:

*Definition 1:* Given two attribute vectors  $\vec{a}_1 = (a_1^1, a_1^2, \dots, a_1^l)$  and  $\vec{a}_2 = (a_2^1, a_2^2, \dots, a_2^l)$ , the attribute association between them, denoted by  $\Delta_{\vec{a}_1, \vec{a}_2}$ , is defined as a pair of two sets of attribute values,  $\{i | a_1^i = 1\}$  and  $\{i | a_2^i = 1\}$  where  $i \in \{1, 2, \dots, l\}$ .

Note that the attribute association is symmetric with respect to a given pair of attribute vectors  $\vec{a}_1$  and  $\vec{a}_2$ , that is,  $\Delta_{\vec{a}_1, \vec{a}_2} = \Delta_{\vec{a}_2, \vec{a}_1}$ . Every pair of nodes has its attribute association and therefore there are as many attribute associations as the number of edges in  $G$ . The attribute association information is widely used in many different applications. For example, the link prediction algorithms that aim to predict whether a link will be newly formed between two unconnected nodes in the future usually employ the link structure information around the two nodes but it could leverage from using the attributes of the nodes as well. Many previous researches have shown that nodes in a graph tend to establish homophily or heterophily relationships in terms of their attributes [17], [18], [19]. Another example of using attribute information is the community detection problem. Many early approaches to detect latent communities rely on only the link structure of a graph [20], [21], [22]. That is, they detect communities such that nodes within the same community interact with each other more frequently than with those outside the community. However more recent studies use the node attributes as well as the link structure and show that the attribute information is helpful for community detection [23], [24], [25].

If an attribute association  $\Delta$  is repeatedly observed and its frequency is over a given threshold  $\sigma$  that is referred as *freq\_support*, then we say  $\Delta$  is a frequent attribute association.

*Definition 2:* Given an attribute association  $\Delta$  and a support  $\sigma$ ,  $\Delta$  is called a frequent attribute association if  $fr(\Delta) \geq \sigma \times |E|$  where  $fr(\Delta)$  is the number of pairs of nodes with  $\Delta$ .

When a frequent attribute association is given, we can say that there are many pairs of nodes having the association but it does not necessarily mean that the attribute association is really interesting. For example, in a social network of *Purdue University Alumni*, it is not surprising to observe many connected nodes have the attribute association of  $\{\text{"Purdue"}, \text{"CS"}\} - \{\text{"Purdue"}, \text{"CS"}\}$ . So we are interested in statistically significant attribute associations rather than frequent ones, which will be discussed in the following section.

## B. Statistically Significance

The statistical significance of an object can be quantified by estimating the probability of the observed or rarer objects under the null hypothesis. Let  $\delta_G$  denote the density of  $G$  which is defined as the fraction of the number of edges in  $G$  over all pairs of nodes ( $\delta_G = \frac{|E|}{1/2 \cdot |V| \cdot (|V| - 1)}$ ). If we randomly select two groups of nodes no matter which attribute values they have, denoted by  $C_1$  and  $C_2$  respectively, then the expected number of edges between  $C_1$  and  $C_2$  is  $e(C_1, C_2) = |C_1| \cdot |C_2| \cdot \delta_G$  by assuming the probability of a pair of randomly selected nodes being connected to each other follows  $\delta_G$ . Also, assuming the edges are independent of each other, the actual number of edges  $M$  between  $C_1$  and  $C_2$  would follow the binomial distribution with parameters  $n = |C_1| \cdot |C_2|$  and  $p = \delta_G$ , and thus the probability of getting exactly  $k$  edges among  $n$  possible edges is given by the following probability mass function:

$$f(k; n, p) = P[M = k] = \binom{n}{k} p^k (1 - p)^{n-k} \quad (1)$$

If each of  $C_1$  and  $C_2$  is a group of nodes with the same attribute values in  $G$  which are specified by an attribute vector, then the attribute vectors  $\vec{a}_1$  and  $\vec{a}_2$  can be instantiated from  $C_1$  and  $C_2$  respectively and the attribute association between two attribute vectors is induced from the edges across the nodes of  $C_1$  and the nodes of  $C_2$ . So we can measure the statistical significance of a given attribute association  $\Delta_{\vec{a}_1, \vec{a}_2}$  based on the probability  $P[M \geq k]$  that the observed or higher number of edges occur between  $C_1$  and  $C_2$  in which the nodes have  $\vec{a}_1$  and  $\vec{a}_2$  respectively. The association is said to be statistically significant if the estimated probability  $P[M \geq k]$  is very small.

$$P[M \geq k] = 1 - \sum_{i=0}^{k-1} \binom{n}{i} p^i (1 - p)^{n-i} \quad (2)$$

*Definition 3:* An attribute association  $\Delta_{\vec{a}_1, \vec{a}_2}$  between  $C_1$  and  $C_2$  is statistically significant if the probability that the observed or more number of edges between  $C_1$  and  $C_2$  is less than  $\alpha$  which is called a significance level.

In order to show the assumption that the number of edges between two groups of nodes follows the binomial distribution is reasonable, we randomly sampled two groups of 50 nodes from the *DBLP co-authorship network* (the details of the network is described in Section V) 10,000 times and obtained

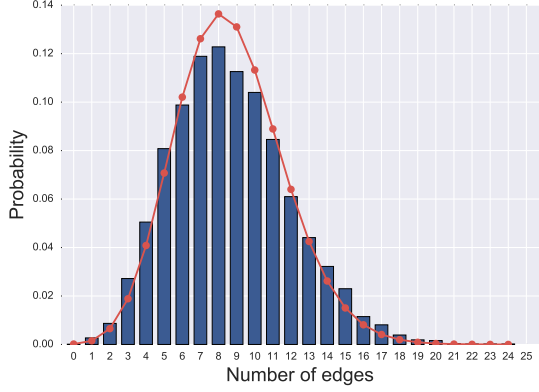


Fig. 2: Distribution of the number of edges between two groups of nodes

the empirical distribution of the number of edges residing between the two groups. As shown in Fig. 2, the empirical distribution (blue bar, mean: 8.68 / stddev: 3.29) is very closed to the actual binomial distribution (red line, mean: 8.64 / stddev: 2.93), which is verified by the chi-squared testing on the two distributions.

### C. Locality Preserving Significant Associations

An attribute association may reside in anywhere over the entire graph  $G$ . However, we expect that a certain attribute association could be observed more frequently among nodes which are closed to each other. For example, in the *DBLP co-authorship network*, some authors who have published papers in venues of data mining area are expected to have a certain attribute association with other authors in the same or similar area (e.g., the association of  $\{\text{ICDM}, \text{KDD}\} - \{\text{ICDM}, \text{NIPS}, \text{ICML}\}$ ). Any pair of authors in a relationship with the association could be seen in several locations of  $G$ , but some of them may be located very closely in terms of the hop distance in the graph and form a densely connected subgraph or community. Different communities that have the same venue pattern many times would be corresponding to different schools in different countries. That is, some attribute association patterns come with locality in the graph and such a pattern can be more statistically significant locally rather than globally. Besides, some attribute association patterns that are statistically significant locally may form another complex patterns (e.g., star or chain, not just pair) among them. One of the nice features of the algorithm we propose in Section IV is that it is able to effectively find all the statistically significant attribute associations while preserving the locality.

## IV. GRAPH TRANSFORMATION

In this section, we describe the algorithm that finds statistically significant attribute associations in a given attribute graph  $G$ . The basic approach for finding statistically significant attribute associations is to transform the original graph  $G$  into a new graph  $\mathcal{AG} = (\mathcal{V}, \mathcal{E}, \mathcal{A})$ , which is called *Association*

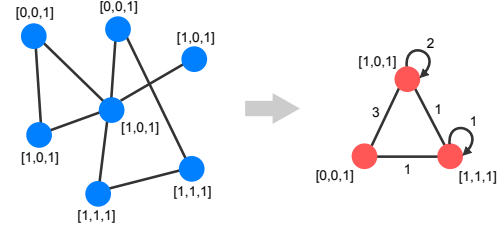


Fig. 3: Graph transformation

*Graph*, where each node in  $\mathcal{V}$  corresponds to a group of nodes in  $V$  which have the same or similar attribute values, each edge in  $\mathcal{E}$  is an attribute association  $\Delta$  between two attribute vectors, each attribute vector in  $\mathcal{A}$  represents one shared by a group of nodes in  $\mathcal{V}$ . To avoid confusion, from now on we call a node in  $\mathcal{V}$  a cluster and call an edge in  $\mathcal{E}$  an association. Each association  $\Delta$  is assigned a weight, referred as its strength  $w(\Delta)$ , that is given by the number of edges between nodes in the clusters forming the association. For a given association  $\Delta$  and its associated strength, defined as the number of edges between nodes in the clusters, we can determine whether  $\Delta$  is significant or not by looking at the strength and the size of the clusters to which  $\Delta$  is incident, which will be explained in detail in the following sections.

The graph transformation can be done through an iteration of two steps. We first start with a single cluster that contains all nodes of  $V$  in  $G$  and then the cluster is partitioned into several subclusters by applying two steps repeatedly and iteratively. For the first step, a cluster is split such that each subcluster contains a subset of  $V$  that have similar attribute values. This operation is able to be easily done using any clustering algorithms. In case of binary attributes, we just select one of the attributes and then do two-way split with respect to the attribute. In Section IV-A, we explain how to select the attribute. For the second step, we try to split a cluster such that each of the associations incident to the cluster has higher strength in order to obtain more significant associations between two sets of attributes. That is, the iteration of the two different splits alternate between performing the similarity-based split, which produces clusters with the same or similar attribute values, and the strength-based split, which makes associations more significant. It results in a new graph  $\mathcal{AG}$  where we can see groups of nodes with certain attribute values and significant associations between them, as shown in Fig. 3.

Algorithm 1 shows the whole structure of the graph transformation algorithm including the two steps of splits. and the following subsections describe how each split should be done in detail.

### A. Similarity-based split

As mentioned already, the goal of the first step is to maximize the similarity among attribute values in each cluster so that each cluster can represent a certain set of attribute values. Thus we select one of the clusters in  $\mathcal{AG}$  and then split it into two subclusters based on a certain attribute so that

---

**Algorithm 1** Algorithm for graph transformation

---

**Input:**  $G = (V, E, A)$ **Output:**  $\mathcal{AG} = (\mathcal{V}, \mathcal{E}, \mathcal{A})$ *Initialization :*

- 1:  $\mathcal{V} = \emptyset, \mathcal{E} = \emptyset$
- 2:  $c$  is initialized as a cluster containing all nodes in  $V$
- 3:  $\mathcal{V} = \mathcal{V} \cup c$

*Iterative Process*

- 4: **while** there exist at least one cluster to split **do**
  - 5:    $c = \text{findClusterForSimilaritySplit}(\mathcal{AG})$
  - 6:   **if** ( $c$  exists) **then**
  - 7:      $\text{similaritySplit}(\mathcal{AG}, c)$
  - 8:   **end if**
  - 9:    $c = \text{findClusterForStrengthSplit}(\mathcal{AG})$
  - 10:   **if** ( $c$  exists) **then**
  - 11:      $\text{strengthSplit}(\mathcal{AG}, c)$
  - 12:   **end if**
  - 13: **end while**
  - 14: **return**  $\mathcal{AG}$
- 

each subcluster contains a set of node that share the same value on the attribute. The way to select a cluster in  $\mathcal{AG}$  is based on the following idea. Basically we do not only want to maximize the similarity of attribute values in each subcluster after the split, but also want each subcluster to be statistically significant as much as possible in terms of the attribute values of its nodes.

To achieve the goal, we need to figure out which cluster should be split and which attribute should be used to split the cluster. Let  $p_i$  denote the probability that a value of 1 occurs at  $i$ -th attribute, which is the fraction of nodes with a value of 1 for the  $i$ -th attribute in  $G$ . So,  $p_i$  is considered an expectation of having the attribute value for a random node. First, an attribute on which a cluster should be split based is picked such that the probability of the attribute having the value of 1 in the cluster is least deviated from its corresponding  $p_i$ . It allows the subclusters to not only have higher similar attribute values among the nodes in them but also have the highest significance gain through the split. Once we decide which attribute should be used for the split of the clusters, we select one of the clusters to split. While assuming that the attributes are independent of each other and the number of times the value of 1 appears at the  $i$ -th attribute follows the binomial distribution with the probability  $p_i$ , the statistical significance  $\Psi_c$  of a cluster  $c$  is defined based on the product of p-values of the attribute values of the nodes in the cluster as follows,

$$\Psi_c = 1 - \prod_{i=1}^l \left( 1 - \sum_{j=0}^{k_i-1} \binom{|c|}{j} p_i^j (1-p_i)^{|c|-j} \right) \quad (3)$$

where  $k_i$  is the number of nodes having the value of 1 on the  $i$ -th attribute and  $|c|$  is the number of nodes in the cluster  $c$ . So for each cluster  $c$  we compute  $\Psi_{c'}$  of the subclusters  $c'$ . Remind that our goal is to split a cluster so

that its subclusters are most statistically significant. However, since the subclusters may have different significances (one can be highly significant but the others can be very low), we take subclusters with the lowest significance from each of the clusters in  $\mathcal{AG}$  and then select a cluster that will produce a subcluster with the highest significance among those subclusters, i.e.,

$$\arg \max_c \left( \min_{c' \in sb(c)} \Psi_{c'} \right) \quad (4)$$

where  $sb(c)$  is a set of subclusters that will be created after the split. In this way, we can avoid to split a cluster that will produce the least significant subclusters. By repeating this kind of split,  $\mathcal{AG}$  will have only clusters, in each of which the same attribute values are shared by its nodes, but we need to place one constraint while doing the split. Even though a cluster represents a certain set of attribute values shared in it, if it contains only a few nodes then its attribute values may not be meaningful at all when we look at an attribute association between clusters in  $\mathcal{AG}$ . Thus, we use *size\_support*, denoted by  $\lambda$ , to force a cluster not to split any more if all the subclusters that will be obtained after splitting the cluster have the sizes less than  $\lambda \cdot |V|$ . Thus, during the first step, we examine only clusters satisfying the  $\lambda$  threshold to determine which cluster should be split. Also, it is obvious that a cluster in which all its nodes have the same attribute values does not need to be split.

We do not only want nodes in the same cluster to have the same attribute values but also allow them to have similar attribute values. In other words, even though every node in a cluster does not agree on a certain attribute, if the distribution of the values of the attribute is statistically significantly deviated from the expectation, then those nodes are considered to have an identical value for the attribute.

Once a cluster is split at the first step, we move on to the second step to increase the significances of the attribute associations between clusters.

### B. Strength-based split

While the similarity-based split of the first step aims to increase the similarity of attribute values for a cluster, we try to maximize strengths of associations to which a cluster is incident through the strength-based split. Given an attribute association between two clusters, its strength is defined as the number of edges that connect the nodes of the clusters. The strength is not meaningful by itself because the significance depends on the sizes of the clusters as well as the strength. As we discussed the definition of a statistically significant attribute association in Section III-B, the stronger strength an attribute association has and the smaller the associated clusters are, the higher statistically significant the association is. Thus, in order to make an association more significant, a cluster that is one of the end points of the association needs to be split into subclusters such that nodes which have many common neighbor clusters belong to the same

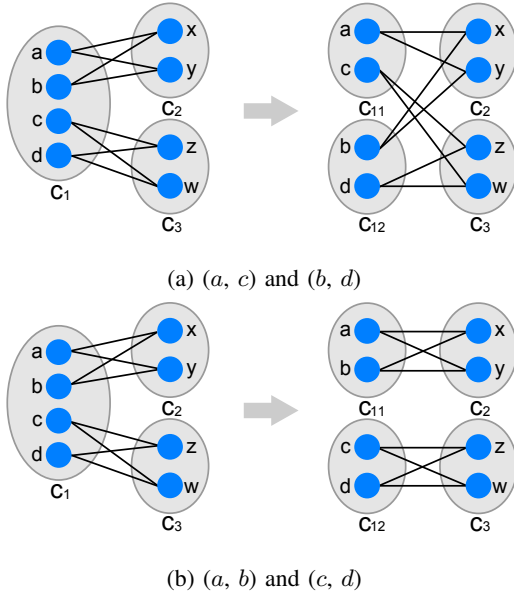


Fig. 4: Two different strength-based splits

subcluster. Fig. 4 illustrates the basic idea of the strength-based split. Suppose we want to maximize the significance of the associations held by the cluster  $c_1$  and we consider two different splits to do that as presented in Fig. 4a and Fig. 4b. The nodes  $a$  and  $b$  in  $c_1$  have edges, all of which are incident to other nodes in  $c_2$  while the nodes  $c$  and  $d$  are adjacent to only other nodes in  $c_3$ . Thus, in order for the subclusters obtained from splitting  $c_1$  to have associations of maximized significance, the split should produce two subclusters which contain the two nodes  $a$  and  $b$ , and the other two nodes  $c$  and  $d$ , respectively.

So we need to find the optimal split of a cluster so that its associations become more significant. For a given cluster  $c$  we try to split, we build a graph  $\tilde{G} = (\tilde{V}, \tilde{E})$  where  $\tilde{V} = \{u | u \in c\}$  and  $\tilde{E} = \{(u, v) | u, v \in c \wedge \exists c' \text{ s.t. } (u, w_1), (u, w_2) \in E \text{ and } w_1, w_2 \in c'\}$ , some of which are connected to each other if they have edges with some common neighbor clusters,  $\Gamma(c)$ . Those edges in  $\tilde{E}$  are weighted based on the fraction of edges to common neighbors among all of their edges. Then, we partition the graph  $\tilde{G}$  based on the weights of the edges in the graph and the subgraphs resulted from the partition become the subclusters we obtain through the strength-based split. For this task, we need to come up with a proper way to assign weights to the edges. We borrow the idea of tie-strength between individuals in social network. In the social science community, there are many different ways to define the tie-strength of an interpersonal relationship [26], and one widely used measure is the Jaccard index. That is, a tie-strength between two individuals  $u$  and  $v$  is determined by  $|\Gamma(u) \cap \Gamma(v)| / |\Gamma(u) \cup \Gamma(v)|$  where  $\Gamma(\cdot)$  is a set of neighbors of a node. In our setting, two nodes  $u$  and  $v$  in the cluster  $c$  may not have common neighbor nodes in  $G$  but some of their neighbor nodes may belong to the same neighbor cluster  $c'$  in

$\mathcal{AG}$ . Similarly, when  $u$  and  $v$  in  $c$  are connected to some of the nodes in a common neighbor cluster  $c'$  of  $c$ , there might not be common nodes in  $c'$  which are incident to both  $u$  and  $v$ . Thus, we modify the Jaccard index slightly so as to measure the tie-strength between  $u$  and  $v$  while capturing the common neighbor clusters.

$$TS(u, v) = \frac{\sum_{c' \in \Gamma(c)} \min\{\phi(u, c'), \phi(v, c')\}}{\sum_{c' \in \Gamma(c)} \max\{\phi(u, c'), \phi(v, c')\}} \quad (5)$$

where  $\phi(u, c') = |w | w \in c' \wedge (u, w) \in E|$ , that is the number of edges in  $E$  between  $u$  and any nodes in  $c'$ . Using this tie-strength measure, we can have nodes belong to the same subcluster after the split if they have many common neighbor clusters, regardless of whether they have common neighbor nodes in  $G$  or not (of course, it depends on the weight given by  $TS(\cdot, \cdot)$ ).

Once we have  $\tilde{G}$  for the cluster  $c$  then we perform graph partitioning on  $\tilde{G}$  to find optimal subclusters that can make the associations between  $c$  and  $c' \in \Gamma(c)$  more significant. Since all the edges  $\tilde{E}$  of  $\tilde{G}$  are assigned weights and  $\tilde{G}$  should be partitioned based on the weights, we take an approach to maximize the modularity of  $\tilde{G}$  [22]. The modularity  $Q(\tilde{G})$  is defined as

$$Q(\tilde{G}) = \frac{1}{2m} \sum_{u, v} \left[ A_{uv} - \frac{k_u k_v}{2m} \right] \delta(c_u, c_v) \quad (6)$$

where  $m = \tilde{E}$ ,  $k_u$  is the degree of  $u$ ,  $c_u$  is the group to which  $u$  belongs, and  $A_{uv}$  is 1 if there is an edge in  $\tilde{E}$  between  $u$  and  $v$  otherwise 0. That is, the modularity is the fraction of the edges that fall within the given groups minus the expected such fraction if edges were distributed at random. If we split the cluster  $c$  through the graph partitioning method as described, a set of nodes that share many common neighbor clusters is likely to fall within the same subcluster as much as possible, and different nodes that share only few common neighbors would be distributed to different subclusters. Thus, we can increase the statistical significances of the attribute associations.

During the second step, we enforce a couple of conditions to prune some clusters and associations in  $\mathcal{AG}$  and do not perform the strength-based split on them for both achieving computational efficiency and finding more meaningful results. As done in the first step, we use *size\_support*,  $\lambda$  because if the size of a cluster  $c$  is too small, we do not believe that  $c$  is representative of a certain set of attribute values. Thus, the strength-based split is run for a cluster  $c$  only when  $|c| \geq \lambda \cdot |V|$ . In addition to that, if a cluster has an attribute association with too weak strength, then we can safely discard it for the rest of the algorithm. Note that the strength of an attribute association between two clusters monotonically decreases as the two splits are performed iteratively while the statistically significance is not monotonic in either way. Since we consider only attribute associations between clusters satisfying the *size\_support* condition and the statistically

significance of an association depends on its strength and the sizes of the clusters at the end points, we can prune an attribute association from  $\mathcal{AG}$  as long as it meets the following condition.

*Lemma 1:* Given an attribute association  $\Delta_{c_1, c_2}$  and its two incident clusters  $c_1$  and  $c_2$ , if the strength of  $\Delta_{c_1, c_2}$  is less than  $\Phi^{-1}\left(1 - \alpha - \frac{C(p^2 + q^2)}{\sqrt{npq}}\right)\sqrt{npq} + np$ , then  $\Delta_{c_1, c_2}$  does not have a chance to be statistically significant any more, where  $n = |c_1| \cdot |c_2|$ ,  $p = \delta_G$ ,  $q = 1 - p$ ,  $\Phi(\cdot)$  is the error function, and  $C$  is a constant.

*Proof:* Given the size\_support  $\lambda$ , both the clusters  $c_1$  and  $c_2$  should have the size of at least  $|V| \cdot \lambda$  in order to make the attribute association  $\Delta_{c_1, c_2}$  considered as statistically significant. Also, let  $k$  denote the strength of  $\Delta_{c_1, c_2}$  and then according to the (2),  $P[X \geq k] \leq \alpha$ . If we approximate the binomial distribution using the normal distribution,

$$\begin{aligned} P[X \geq k] &= P\left[\frac{X - np}{\sqrt{npq}} \geq \frac{k - np}{\sqrt{npq}}\right] \\ &= P\left[Z \geq \frac{k - np}{\sqrt{npq}}\right] \leq \alpha \end{aligned} \quad (7)$$

Now we have the standard normal distribution and need to find the lower bound of  $k$  which satisfies the inequality (7). Using the error function  $\Phi(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-\frac{t^2}{2}} dt$  which is essentially identical to the standard normal cumulative distribution function [27],

$$\begin{aligned} \frac{k - np}{\sqrt{npq}} &\geq \Phi^{-1}(1 - \alpha) \\ k &\geq \Phi^{-1}(1 - \alpha)\sqrt{npq} + np \end{aligned} \quad (8)$$

The lower bound for  $k$  is originated from approximation based on the standard normal distribution and thus we need to get the error bound. According to the following *Berry-Essen theorem* [28],

$$\sup_{x \in \mathbb{R}} \left| P\left[\frac{B(p, n) - np}{\sqrt{npq}} - \Phi(x)\right] \right| \leq \frac{C(p^2 + q^2)}{\sqrt{npq}} \quad (9)$$

with  $C < 0.4748$ , we know that the error arising from the approximation is at most  $\frac{C(p^2 + q^2)}{\sqrt{npq}}$ . As a result, if we relax the lower bound for  $k$  in (8) to the extent of the error, then we obtain

$$k \geq \Phi^{-1}\left(1 - \alpha - \frac{C(p^2 + q^2)}{\sqrt{npq}}\right)\sqrt{npq} + np \quad (10)$$

Since  $p$  is very small and  $n$  is large for given  $\lambda$ , the error bound is small and we can still get a reasonably tight lower bound for  $k$ . Regarding the inverse of the error function, if we use  $\alpha = 0.01$  as the significance level, then  $\Phi^{-1}(1 - \alpha) = 1.8212$ . According to Lemma 1, we drop attribute associations if they are too weak to be able to be significant later on. In fact, such associations are noise and do not bring us any meanings. Rather, it prevents the strength-based split from running optimally. ■

	Original graph			Association graph	
	Nodes	Edges	Density	Nodes	Edges
DBLP	4,672	37,726	0.00346	195	6,302
Yelp	4,454	44,906	0.00453	202	8,388

TABLE II: Dataset statistics

Subarea	Conferences
DM/ML	ICDM, NIPS, ICML
OS	SOSP, OSDI
Theory	FOCS, STOC, SODA
Security	IEEE Symposium on Security and Privacy (S&P), ACM Conference on Computer and Communications Security (CCS)

TABLE III: DBLP subareas in computer science

## V. EXPERIMENTS

### A. Datasets

We ran the graph transformation algorithm on real co-authorship and social networks, and obtained the resulting association graphs. Using the association graphs, we analyzed qualitative differences between the statistically significant and frequent associations. Also we showed the application of the significant patterns to a link prediction problem, and synthetic graphs with attributes are considered to show the algorithm's scalability.

**DBLP.** We obtained a collection of bibliographic information from the DBLP website [29], an open bibliographic information provider of computer science journals and conferences. Each record of journal or conference paper has one or more authors, and the venue, on which it is published. We first filtered out any authors who appear in less than 3 papers. Then, we considered only papers published to the 10 conferences of 4 different subareas of computer science, i.e., data mining and machine learning (DM/ML), operating systems (OS), theory, and security. More details are shown in Table III. Then we built an attribute vector of length 10 for each node, i.e., if an author (or a node) published a paper to a conference in Table III, we set the corresponding vector value to 1. If not, we set the corresponding vector value to 0. Finally, an edge is formed if two authors (or nodes) have co-authored at least one paper in the dataset.

**Yelp.** Yelp is a provider of crowd-sourced reviews about local businesses, along with a social network. The Yelp challenge dataset [30] contains the social network, composed of the users (nodes) and their friend relations (edges). Also the sets of users' reviews are provided in the dataset. Each review is tied to a user and a business, and each business has a small set of business type categories. We first filtered out any users who have less than 10 reviews. Then, we considered only reviews for the restaurants, which has at least one of the 10 business categories, {Chinese, Japanese, Mediterranean, Thai, French, Greek, Vietnamese, Korean, Indian, British}. The node attributes are compiled similarly to the DBLP dataset. Note that we did not use some of the most popular restaurant categories, e.g., American, Mexican, and Italian. As the majority of users has left reviews on the restaurants of

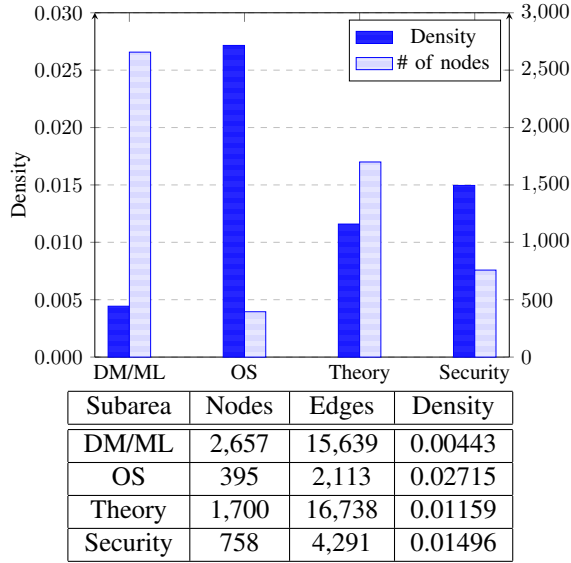


Fig. 5: DBLP graph statistics for different subareas

#	Association	
1	{SOSP, OSDI, S&P, CCS}	– {SOSP, OSDI}
2	{SOSP, OSDI, S&P, CCS}	– {S&P, CCS}
3	{ICML, ICDM, S&P(*)}	– {ICML, ICDM}
4	{SOSP, OSDI, S&P, CCS}	– {SOSP, S&P, CCS}
5	{FOCS(*), STOC(*), CCS}	– {S&P, CCS}
6	{ICML, ICDM, S&P(*)}	– {ICML(*), ICDM}
7	{SOSP, S&P, CCS}	– {SOSP, OSDI}
8	{SOSP, S&P, CCS}	– {S&P, CCS}
9	{ICML, ICDM, S&P(*)}	– {ICDM, OSDI(*)}
10	{ICML, ICDM}	– {ICML(*), ICDM}

TABLE IV: Significant associations minus Frequent associations for DBLP

such categories, they seem to appear in most of the attribute associations and carry little or no information.

### B. Effectiveness Analysis

To evaluate the effectiveness of our algorithm, we conducted the set difference between the statistically significant associations and the top-15 frequent associations. As the resulting significant associations contain wildcard attributes, it is not easy to make direct comparisons or set differences between the two. Thus, we took a conservative approach that as long as all attribute values of any top-15 frequent associations have exact or wildcard attribute match, we considered that there is a match. This approach is certainly in favor of the frequent associations, since it ignores that wildcard matches may lead to some other possible set of attribute values.

Table IV, we have the set difference between statistically significant and frequent associations for the DBLP dataset. First consider 4 subgraphs that only contains the nodes and their edges, whose attribute value for any conference of the corresponding subarea is 1. Fig. 5 describes the characteristics of each subgraph. The subgraph of DM/ML has a large number of nodes but its graph density is small, which means that the tie-strengths are weak. On the other hand, there are relatively

small numbers of nodes in the subgraph of OS and security but their densities are high, which means that the tie-strengths are strong among the nodes. We can easily identify the OS and security-related associations, which contain {SOSP, OSDI} and {S&P, CCS}, are appearing on top of the difference list. Also note that many frequent associations are related to DM/ML conferences since its subgraph contains the most number of edges while its density is low.

From Table IV, we can infer many interesting significant associations, which do not appear in the frequent association list. The association number 1, 2, 4, 7, and 8 clearly shows that the nodes who have authorship in the OS-related conferences tend to co-work with the authors in the security-related conferences. The association number 3, 5 and 6 shows that the nodes who have authorship in the security-related conferences frequently co-work with the authors in DM/ML and theory-related conferences. Interestingly enough, the association number 9 shows how the authors in DM/ML, security, and OS have frequent co-authorship relations in the graph. These results might look obvious to some of the readers who have a good understanding of co-authorship in computer science. However, when the relationships of attributes are little known, the discussed results may be intriguing.

Table V shows that the set difference between statistically significant and frequent associations for the Yelp dataset. Note that {Chinese, Japanese} appears very commonly in the association results due to their prevalence in node attributes. Thus, we will exclude them from the subsequent discussions. Also it turned out that the first 10 significant associations with the highest statistical significance are the same as the associations reported in Table V. That is, none of the first 10 significant associations are reported in the top-15 frequent association results, since the significant associations do not occur often in terms of frequency but do occur often in the dataset in a statistically significant manner.

Among the frequent visitors of {Mediterranean, Thai}, the association number 2 and 7 shows that the nodes with {Greek} attribute are strongly associated with the nodes with {Vietnamese, Korean} attributes, and the association number 4, 6 and 10 shows that the nodes with {Vietnamese, Korean, Indian} are strongly associated with the nodes with {Vietnamese, Korean} and {Greek} attributes. Also the association number 5 and 8 describes that the nodes with {Mediterranean} have statistically significant associations with the nodes with {Mediterranean, Thai, Greek}.

### C. Scalability Analysis

We evaluated the computation cost of our algorithm on synthetic attributed graphs of different sizes and densities. The experiments were carried on a machine with an Intel Xeon 3.1GHz CPU and 32GB memory, running 64bit Ubuntu 14.04. All algorithms are implemented in Python 2.7.

The graphs are generated based on the simplified version of Multiplicative Attribute Graph (MAG) model [17]. MAG is widely used in the literature to generate synthetic graphs with node attributes, and known to model real-world networks with

#	Association
1	{Chinese, Japanese, Mediterranean, Thai, Greek} – {Chinese, Mediterranean, Thai(*), Greek}
2	{Chinese, Japanese, Mediterranean, Thai, Vietnamese, Korean} – {Chinese, Japanese, Mediterranean, Thai, Greek}
3	{Chinese, Japanese, Thai, Vietnamese, Korean} – {Chinese, Japanese, Vietnamese, Korean}
4	{Chinese, Japanese, Mediterranean, Thai, Vietnamese, Korean, Indian} – {Chinese, Japanese, Mediterranean, Thai, Vietnamese, Korean}
5	{Chinese, Mediterranean, Thai(*), Greek} – {Chinese, Mediterranean}
6	{Chinese, Japanese, Mediterranean, Thai, Vietnamese, Korean, Indian} – {Chinese, Japanese, Thai}
7	{Chinese, Japanese, Mediterranean, Thai, Vietnamese, Korean} – {Chinese, Mediterranean, Thai(*), Greek}
8	{Chinese, Japanese, Mediterranean, Thai, Greek} – {Chinese, Mediterranean}
9	{Chinese, Japanese, Thai, Vietnamese, Korean} – {Chinese, Japanese, Thai}
10	{Chinese, Japanese, Mediterranean, Thai, Vietnamese, Korean, Indian} – {Chinese, Japanese, Mediterranean, Thai, Greek}

TABLE V: Significant associations minus Frequent associations for Yelp

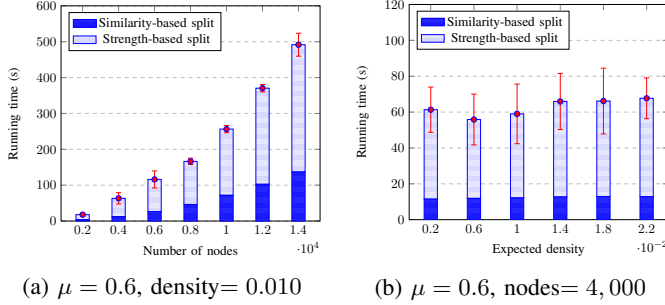


Fig. 6: Running time experiments on synthetic graph datasets

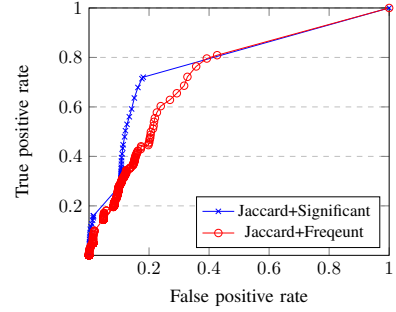


Fig. 7: Link prediction performance

flexibility. We conducted two sets of experiments with  $l = 5$  and  $\mu$ 's fixed, the probability of each attribute value being 1, i.e., each node has five binary attributes and the attributes are drawn from the same distribution, retaining the node attribute distribution throughout the experiments.

**Time complexity.** Our algorithm is divisive in nature and it splits at least one node of *Association Graph* in every iteration. First, the similarity-based split step will run  $\mathcal{O}(2^l)$  iterations. Usually the length of attribute vector is small,  $l \ll n$ , and the similarity-based split under reasonable settings takes much less time compared to that of the strength-based split. In the strength-based split step, it is not hard to see that the computation of tie-strengths between each pair of nodes,  $\mathcal{O}(n^2)$ , dominates the running time of the step. And we can notice that the algorithm will run  $\log n$  iterations of the strength-based steps on average. Accordingly the overall average time complexity of the algorithm is  $\mathcal{O}(n^2 \log n)$ .

**Results.** Fig. 6a shows the computation time over the number of nodes. We fixed the attribute link-affinity matrix [17], which determines the probability of edge formation between two sets of node attributes. Note that since we kept all parameters of the MAG model but the number of nodes, the graph density remained the same. We confirmed that the algorithm is of polynomial time in the number of nodes. This result is in line with the time complexity we discussed above.

In the second experiment, we fixed the number of nodes and the scale factor of the attribute matrix, which merely changes the expected number of edges. That is, we scaled the attribute matrix such that the resulting graphs have the graph densities as we desire, without changing any other properties of the graphs. In Fig. 6b, we can easily observe that the algorithm's

running time remains almost the same as we increase the expected graph density. The aforementioned time complexity should well explain the result.

Finally, both of the plots in Fig. 6 show that the running time of the strength-based split step dominates that of the similarity-based step. Also both plots describe that the running time of the similarity-based step remain the same as we add more edges with the number of nodes fixed, and the running time grows as we increase the number of nodes. This supports our intuition that the similarity-based split step is not relevant to the number of edges or graph density.

#### D. Application: Link prediction

As one of the application for which the statistically significant attribute associations are useful, the *link prediction* problem is considered. Many different approaches to the link prediction have been proposed for the past decade, but with the objective of showing the potential merit of the statistically significant attribute associations, we simply use the Jaccard coefficient proposed in [31] and compare the effects of using statistically significant attribute associations and frequent ones. Given a pair of nodes without an edge, we compute the prediction score by combining the Jaccard coefficient  $J(u, v)$  and the score  $S(u, v)$  resulted from either the significance or the normalized frequency of an attribute association between the nodes as follows

$$\text{pred}(u, v) = \tau \cdot J(u, v) + (1 - \tau) \cdot S(u, v) \quad (11)$$

and if it is over a given threshold then we predict that  $u$  and  $v$  will form a new link. We take two snapshots of the *DBLP co-authorship network* (Mar 2015 and Mar 2016) and all the

newly created links between the two snapshots are used for the positive samples. Similarly, a set of pairs of nodes that do not have an edge in both the snapshots are used for the negative samples. Since the number of negative samples far outweighs the number of positive samples, we do negative subsampling with the ratio of 1 : 5 (five negatives per one positive). In Fig. 7, we report the ROC curves for two different methods, Jaccard+Significant, and Jaccard+Frequent. As shown in Fig. 7, the link prediction can more benefit from employing the attribute information and the statistically significant attribute associations can achieve higher performance rather than the frequent ones.

## VI. CONCLUSION

We defined a problem of mining statistically significant attribute associations using *Association Graph*, which keeps the locality of attribute associations and carries the significant relationships between the sets of attribute values. And we proposed a novel, two-step iterative algorithm that efficiently and effectively generates an *Association Graph* from the original graph. The experiments are conducted on two real world datasets, and we ran some qualitative analysis on the results, confirming that our algorithm effectively finds the significant associations, which cannot be uncovered by conventional frequent association mining. Also we ran extensive scalability experiments on synthetic datasets, and confirmed that the algorithm is of polynomial running time in the number of nodes. Lastly, applying the results from one of the real world datasets to the link prediction task, and we showed how the statistically significant attribute associations can be used in practice.

For future work, we plan to investigate how we can exploit resulting *Association Graph* better, e.g., its locality preserving property, and if we can come up with a linear time algorithm or a distributed algorithm, which can be run on large-scale graphs.

## REFERENCES

- [1] S. Wasserman and K. Faust, *Social network analysis: Methods and applications*. Cambridge university press, 1994, vol. 8.
- [2] J. Scott, *Social network analysis*. Sage, 2012.
- [3] J. Hu, X. Shen, Y. Shao, C. Bystroff, and M. J. Zaki, "Mining protein contact maps," in *BIOKDD*, 2002, pp. 3–10.
- [4] A. Broder, R. Kumar, F. Maghoul, P. Raghavan, S. Rajagopalan, R. Stata, A. Tomkins, and J. Wiener, "Graph structure in the web," *Computer networks*, vol. 33, no. 1, pp. 309–320, 2000.
- [5] J. M. Kleinberg, R. Kumar, P. Raghavan, S. Rajagopalan, and A. S. Tomkins, "The web as a graph: measurements, models, and methods," in *Computing and combinatorics*. Springer, 1999, pp. 1–17.
- [6] J. Pei, D. Jiang, and A. Zhang, "Mining cross-graph quasi-cliques in gene expression and protein interaction data," in *21st International Conference on Data Engineering (ICDE'05)*. IEEE, 2005, pp. 353–356.
- [7] H. He and A. K. Singh, "Graphrank: Statistical modeling and mining of significant subgraphs in the feature space," in *ICDM'06. Sixth International Conference on Data Mining, 2006*. IEEE, 2006, pp. 885–890.
- [8] S. Ranu and A. K. Singh, "Graphsig: A scalable approach to mining significant subgraphs in large graph databases," in *2009 IEEE 25th International Conference on Data Engineering*. IEEE, 2009, pp. 844–855.
- [9] X. Yan, H. Cheng, J. Han, and P. S. Yu, "Mining significant graph patterns by leap search," in *Proceedings of the 2008 ACM SIGMOD international conference on Management of data*. ACM, 2008, pp. 433–444.
- [10] J. Scott, T. Ideker, R. M. Karp, and R. Sharan, "Efficient algorithms for detecting signaling pathways in protein interaction networks," *Journal of Computational Biology*, vol. 13, no. 2, pp. 133–144, 2006.
- [11] Y. Chi, Y. Yang, and R. R. Muntz, "Indexing and mining free trees," in *Third IEEE International Conference on Data Mining, 2003. ICDM 2003*. IEEE, 2003, pp. 509–512.
- [12] W. Hamalainen and M. Nykanen, "Efficient discovery of statistically significant association rules," in *ICDM'08. IEEE International Conference on Data Mining, 2008*. IEEE, 2008, pp. 203–212.
- [13] S. Gunnemann, P. Dao, M. Jamali, and M. Ester, "Assessing the significance of data mining results on graphs with feature vectors," in *2012 IEEE 12th International Conference on Data Mining (ICDM)*. IEEE, 2012, pp. 270–279.
- [14] A. Arora, M. Sachan, and A. Bhattacharya, "Mining statistically significant connected subgraphs in vertex labeled graphs," in *Proceedings of the 2014 ACM SIGMOD international conference on Management of data*. ACM, 2014, pp. 1003–1014.
- [15] M. Sachan and A. Bhattacharya, "Mining statistically significant substrings using the chi-square statistic," *Proceedings of the VLDB Endowment*, vol. 5, no. 10, pp. 1052–1063, 2012.
- [16] C. Low-Kam, C. Raïssi, M. Kaytoue, and J. Pei, "Mining statistically significant sequential patterns," in *2013 IEEE 13th International Conference on Data Mining (ICDM)*. IEEE, 2013, pp. 488–497.
- [17] M. Kim and J. Leskovec, "Multiplicative attribute graph model of real-world networks," *Internet Mathematics*, vol. 8, no. 1–2, pp. 113–160, 2012.
- [18] E. M. Rogers and D. K. Bhowmik, "Homophily-heterophily: Relational concepts for communication research," *Public opinion quarterly*, vol. 34, no. 4, pp. 523–538, 1970.
- [19] M. McPherson, L. Smith-Lovin, and J. M. Cook, "Birds of a feather: Homophily in social networks," *Annual review of sociology*, pp. 415–444, 2001.
- [20] M. Coscia, G. Rossetti, F. Giannotti, and D. Pedreschi, "Demon: a local-first discovery method for overlapping communities," in *Proceedings of the 18th ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM, 2012, pp. 615–623.
- [21] J. Yang and J. Leskovec, "Overlapping community detection at scale: a nonnegative matrix factorization approach," in *Proceedings of the sixth ACM international conference on Web search and data mining*. ACM, 2013, pp. 587–596.
- [22] V. D. Blondel, J.-L. Guillaume, R. Lambiotte, and E. Lefebvre, "Fast unfolding of communities in large networks," *Journal of statistical mechanics: theory and experiment*, vol. 2008, no. 10, p. P10008, 2008.
- [23] R. Balasubramanyan and W. W. Cohen, "Block-lda: Jointly modeling entity-annotated text and entity-entity links," in *SDM*, vol. 11. SIAM, 2011, pp. 450–461.
- [24] S. Günnemann, B. Boden, I. Färber, and T. Seidl, "Efficient mining of combined subspace and subgraph clusters in graphs with feature vectors," in *Advances in Knowledge Discovery and Data Mining*. Springer, 2013, pp. 261–275.
- [25] Y. Ruan, D. Fuhry, and S. Parthasarathy, "Efficient community detection in large networks using content and links," in *Proceedings of the 22nd international conference on World Wide Web*. International World Wide Web Conferences Steering Committee, 2013, pp. 1089–1098.
- [26] M. Gupte and T. Eliassi-Rad, "Measuring tie strength in implicit social networks," in *Proceedings of the 4th Annual ACM Web Science Conference*. ACM, 2012, pp. 109–118.
- [27] W. H. Greene, *Econometric analysis*. Pearson Education India, 2003.
- [28] S. Nagaev and V. Chebotarev, "On the bound of proximity of the binomial distribution to the normal one," *Theory of Probability & Its Applications*, vol. 56, no. 2, pp. 213–239, 2012.
- [29] "DBLP: computer science bibliography," <http://dblp.uni-trier.de/xml/>, April 2016.
- [30] "Yelp dataset challenge," <https://www.yelp.com/>, July 2014.
- [31] D. Liben-Nowell and J. Kleinberg, "The link-prediction problem for social networks," *Journal of the American society for information science and technology*, vol. 58, no. 7, pp. 1019–1031, 2007.