

a hierarchical language model in enterprise corpora. Balog et al. [16] employed probabilistic models to study the problem of expert finding, which try to identify a list of experts for a query. However, these methods retrieve experts for each query independently, and cannot be directly used to deal with the expertise matching problem. Thus, several key questions arise for expertise matching, i.e., how to design a framework for expertise matching to guarantee an optimal solution under various constraints? How to develop an online algorithm so that it can incorporate user feedbacks in real time?

Fig. 1 shows an example of paper–reviewer matching problem (assigning reviewers to each paper). In the problem, the topics corresponding to a reviewer (i.e., the expertise of the reviewer) can be “machine learning”, “data mining”, “computational theory”, etc. Also, each paper has a distribution on different topics. There are some requirements to be a good assignment. First, for each paper, the assigned reviewers’ expertise should cover the topics of the paper, and all the reviewers should have a load balance (each reviewer can only review a certain number of papers). In addition, some reviewers might be senior and some might be average. We always hope that the review process of each paper can be “supervised” by at least one senior reviewer. Another example is the patient–doctor matching case, the topics corresponding to the doctor’s expertise include “pediatrics”, “rheumatology”, “neurology”, etc. Each doctor has different expertise degrees on different topics, while the disease of a patient also has a relevance distribution on the topics. Ideally, when arranging a consultation for a patient, the topics of the assigned doctors should contain the potential causes of the patient’s disease (e.g. a consultation for an SLE patient requires rheumatologist, nephrologist, cardiologist and neurologist), and all the doctors should have a load balance so that no doctor overstrains.

Contributions. In this paper, we formally define the problem of expertise matching and propose a constraint-based optimization framework to solve the problem. Specifically, the expertise matching problem is cast as a convex cost flow problem and the objective is then to find a feasible flow with minimum cost under certain constraints. We theoretically prove that the proposed framework can achieve an optimal solution under various constraints and develop an efficient algorithm to solve it. This paper is an extension and refinement of our previous conference paper [17], and differs the previous work in two aspects. (1) We reformalize our framework considering “multi-topic coverage” matching, which is very important to paper–reviewer assignment problem. We show that topic coverage measures (e.g. Coverage and Confidence) can also be incorporated into our framework.

(2) An additional dataset is introduced to evaluate the performance of our approach on multi-topic coverage. Experimental results substantiate the effectiveness and efficiency of the

proposed approach. We have applied the proposed method to help assign reviewers to papers for a top conference. Feedbacks from the conference organizers confirm the usefulness of the proposed approach.

The rest of the paper is organized as follows: Section 2 reviews the relevant literatures. Section 3 formally formulates the problem. Section 4 explains the proposed optimization framework. Section 5 gives experimental results that validate the effectiveness and the computational efficiency of our methodology. Finally, Section 6 concludes.

2. Related work

In general, existing methods for expertise matching mainly fall into two categories: probabilistic model and optimization model. The probabilistic model tries to improve the matching accuracy between experts and queries based on different probabilistic models such as keyword matching [1], latent semantic indexing [2], probabilistic topic modeling [3,4]. However, most of these methods do not consider the various constraints or simply consider the constraints by heuristics. The optimization model tries to incorporate the constraints as a component in an optimization framework such as integer linear programming [5] and minimum cost flow [6].

Most previous works cast expert matching or expert finding as an information-retrieval problem, in which every expert is represented as a “expertise” document and given a query, the goal is to retrieve most relevant experts. As a result, these methods mainly focus on two points: how to define the matching score between a query and a document; and how to represent each expert [18,19]. For example, Dumais and Nielsen [2] use latent semantic indexing (LSI) as the retrieval method and the abstracts provided by reviewer as expertise documents. Yu et al. [20] represent experts by analyzing text content and extracting related information. Basu et al. [1,21,22] integrate different sources of information for recommendation (e.g. publications, research interests, etc.). Yarowsky and Florian [23] assign a paper by computing its cosine similarity with a reviewer and choosing the one with the highest rank. Other expert finding work include [24,25].

In addition, different language models [14–16,26,27] and topic models [28] are used for expert matching/finding problem. In all of the language models, the matching score is the probability of a query given an expertise document, i.e., $p(q|d)$, but its definition varies. Mimno and McCallum [4] improve the matching accuracy by proposing a novel topic model Author-Persona-Topic (APT), in which experts are represented as independent distributions over topics. Karimzadehgan et al. also consider matching experts on multiple aspects of expertise [3]. Unlike the previous probabilistic

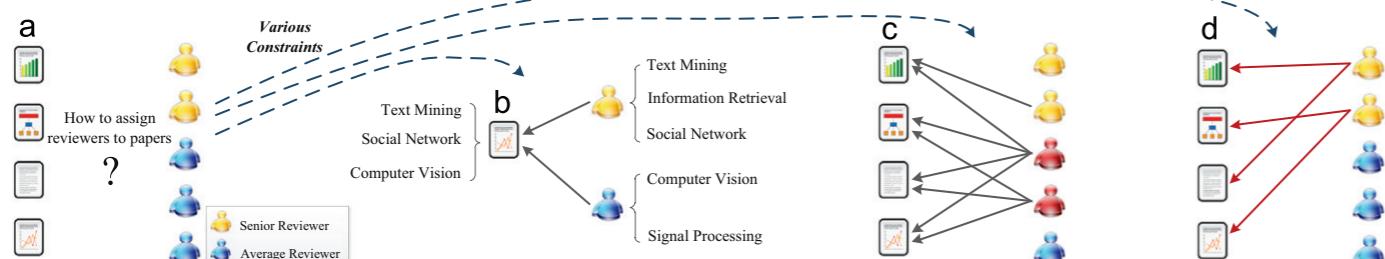


Fig. 1. (a) An illustration of the paper–reviewer assignment problem. To make an ideal assignment, some constraints should be considered. For example: (b) Each paper is related to one or several topics, each reviewer also has expertise on different topics. For a paper, the combination of the expertise of all assigned reviewers should cover all topics of the paper. (c) The assignment should be “load balanced”. A bad example is shown, where the two reviewers marked with red are assigned with too many works, meanwhile the blue one does not get any assignment. (d) The authoritative degree of reviewers may vary largely, it is also desirable that each paper is assigned with at least one senior reviewer, so that she/he can guide the review process.

企业语料库中的分层语言模型。Balog等人。[16]采用概率模型来研究专家发现的问题，该问题试图确定查询专家名单。但是，这些方法独立检索每个查询的专家，不能直接用于处理专业匹配问题。因此，有关专业知识匹配的几个关键问题，即如何设计专业知识匹配的框架，以保证在各种限制下的最佳解决方案？如何开发在线算法，以便它可以实时合并用户反馈？

提出的方法。我们已应用提出的方法，以帮助将审阅者分配给首脑会议的论文。会议组织者的反馈确认了拟议方法的有用性。其余部分组织如下：第2节审查相关文献。第3节正式制定了漂白剂。第4节解释了所提出的优化框架。第5节给出了验证了我们方法的有效性和计算效率的实验结果。最后，第6节结束了。

2. Related work

一般来说，现有的专业知识匹配方法主要分为两类：概率模型和优化模型。概率模型试图根据关键字匹配[1]，潜在语义索引[2]，概率主题建模[3,4]等不同概率模型来提高专家和查询之间的匹配准确性。然而，这些方法中的大多数不考虑各种限制或简单地通过启发式来控制约束。优化模型试图将约束包含在优化框架中的组件，例如整数线性编程[5]和最小成本流程[6]。最先前的工程投射专家匹配或专家查找作为信息检索问题，其中每个专家都被称为“专业知识”文件，并给出了查询，目标是检索大多数相关专家。结果，这些方法主要关注两点：如何在查询和文档之间定义匹配分数；以及如何代表每个专家[18,19]。例如，Dumais和Nielsen[2]使用潜在语义索引（LSI）作为审阅者提供的检索方法和摘要作为专业知识文件。yu等人。

[20]通过分析文本内容和提取相关信息来代表专家。Basu等人。[1,21,22]整合不同的信息来源以供建议（例如出版物，研究兴趣等）。Yarowsky和Florian

[23]通过将其与审阅者的余弦相似度计算并选择具有最高等级的余弦相似度分配纸张。其他专家寻找工作包括[24,25]。此外，不同的语言模型[14–16,26,27]和主题模型[28]用于专家匹配/查找问题。在所有语言模型中，匹配分数是给定专业知识文件的查询的概率，即 $P(Q|D)$ ，但其定义变化。MIMNO和MCALLUM[4]通过提出一个新颖的主题模型作者—Persons-pompelea-主题（APT）来提高匹配的准确性，其中专家表示为与主题的独立发行版。Karimzadehgan等。还要考虑匹配专门知识的多个方面的专家[3]。与以前的概率不同

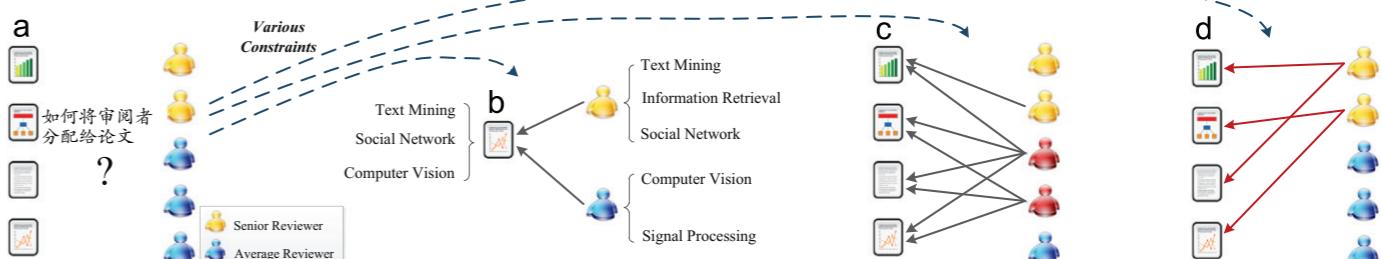


图1。(a) 纸张审阅者分配问题的示例。为了做出理想的作业，应考虑一些约束。例如：(b) 每篇论文与一个或多个主题有关，每个审稿人也具有不同主题的专业知识。对于一篇论文，所有指定审核人员的专业知识的组合应涵盖本文的所有主题。(c) 作业应该是“负载平衡”。显示了一个糟糕的例子，其中标有红色的两个评论员被分配了太多的作品，同时蓝色没有得到任何分配。(d) 审阅者的权威程度可能很大程度上，也有意义上，每篇论文至少有一个高级评论员分配，以便她/他可以指导审查流程。

models in which a query is matched as a whole unit, it tries to find “comprehensive” matchings to cover all subtopics of a query. New measures *Coverage* and *Confidence* are defined to evaluate multi-aspect expertise matching results. Several matching methods are proposed and show better multi-aspect performance than traditional ones.

However, most of the aforementioned works treat each query independently and ignore the certain constraints (e.g. load balance), thus they cannot be directly adapted to an expertise matching system. In the real world, expertise matching is a highly constrained problem, and some existing works study this constrained optimization problem using various methods. For example, Guervós et al. [29] combine a greedy and an evolutionary algorithm [30,31] to assign papers to reviewers. Karimzadehgan et al. [5] and Taylor [32] cast it as an integer linear programming (ILP) problem so approximate solutions can be found by any ILP solver. Sun et al. [7] solves the reviewer assignment by hybrid approach of domain knowledge. Recently, a few systems [8–10,33,34,11–13] have also been developed to help proposal-reviewer and paper-reviewer assignments. However, the expertise matching problem is still treated as an information-retrieval problem, which obviously cannot result in an optimal solution.

In this paper, we aim to formalize the problem of expertise matching in a constraint-based optimization framework and propose an efficient algorithm to solve the framework. The differences of our work from existing work are: (a) we offer an optimization framework that incorporates the expertise matching and various constraints together; (b) the framework can be easily extended since new constraint can be combined into the optimization framework by simply defining a new (hard or soft) constraint; and (c) the framework can guarantee an optimal solution.

3. Problem formulation

In this section, we first give several necessary definitions and then present a formal definition of the problem.

Given a set of experts $V = \{v_i\}$, each expert has different expertise over all topics. Formally, we assume that there are T aspects of expertise (called topics) and each expert v_i has different expertise degrees on different topics. Further, given a set of queries $Q = \{q_j\}$, each query is also related to multiple topics. Given this, we first define the concept of topic model.

Definition 1 (*Topic model*). A topic model θ of an expert (or a query) is a multinomial distribution of words $\{p(w|\theta)\}$. Each expert (query) is considered as a mixture of multiple topic models. The assumption of this model is that words associated with the expert (query) are sampled according to the word distributions corresponding to each topic, i.e., $p(w|\theta)$. Therefore, words with the highest probability in the distribution would suggest the semantics represented by the topic.

Assuming we have T topics, the expertise degree of expert v_i on topic $z \in \{1 \dots T\}$ is represented as a probability $\theta_{v,z}$ with $\sum_z \theta_{v,z} = 1$. Similarly, for each query, we also have a T -dimensional topic distribution with $\sum_z \theta_{q,z} = 1$. Notations are summarized in Table 1.

It is easy to understand that each query q_j can be represented as a sequence of words, i.e., $d_{q,j}$. To represent every expert v_i , without loss of generality, we also consider it as a sequence of words, i.e., $d_{v,i}$. Based on this representation, we can calculate the similarity (or relevance score) between each query and every expert using measures such as cosine similarity or language model. Given this, we can define our problem of expertise matching with various constraints.

Table 1
Notations.

Symbol	Description
M	Number of experts
N	Number of queries
T	Number of topics
V	The set of candidate experts
Q	The set of queries
v_i	One expert
q_j	One query
$\theta_{v,z}$	The probability of topic z given expert v_i
$\theta_{q,z}$	The probability of topic z given query q_j
$T(v_i)$	The set of major related topics of expert v_i
$T(q_j)$	The set of major related topics of query q_j

Problem 1 (*Expertise matching with constraints*). Given a set of experts V and a set of queries Q , the objective is to assign m experts to each query by satisfying certain constraints, such as (1) the number of assigned queries with each expert should be in a range $[n_1, n_2]$, where $n_1 \leq n_2$; (2) the experts' major topics should cover the query's related topics; (3) the assignment should avoid some conflict-of-interest (COI).

Actually in some applications, satisfying the constraints is more important than matching expertise with the queries. For example, in the conference paper-reviewer assignment, the authors of a paper should not be assigned to review their own papers. This must be a hard constraint. While in some other scenario, the constraint is relatively soft, for example the load balance among experts. The number of assigned queries to each expert can be in a range between n_1 and n_2 . In existing works, Dumais et al. [2] and Mimno et al. [4] mainly focus on improving the accuracy of expertise matching, but ignore how to obtain an optimal matching satisfying the various constraints. Karimzadehgan et al. [5] use integer linear programming to find the solution for expertise matching with constraints. However, the proposed model cannot guarantee an optimal solution. In this work, we propose a generalizable optimization framework to solve this problem. Various constraints can also be incorporated in the framework.

4. The constrained optimization framework

In this section, we propose a constraint-based optimization framework for expertise matching. We develop an efficient algorithm to solve the optimization framework based on the theory of convex cost flow, and also present an online matching algorithm to incorporate user feedbacks in real time.

Basic idea. The basic idea of our approach is to formulate this problem in a constrained optimization framework. Different constraints can be formalized as penalty in the objective function or be directly taken as the constraints in the optimization solving process. For solving the optimization framework, we transform the problem to a convex cost network flow problem, and present an efficient algorithm which guarantees an optimal solution.

4.1. The framework

Now, we explain the proposed approach in detail. In general, our objective can be viewed from two perspectives: maximizing the matching score between experts and queries and satisfying

查询作为整个单位匹配的模型，它试图找到“全面”匹配以覆盖查询的所有子主题。新措施覆盖和信心被定义为评估多方面专业知识匹配结果。提出了几种匹配方法，并显示出比传统的更好的多方面性能。然而，大多数上述作品独立地处理每个查询并忽略某些约束（例如，负载平衡），因此它们不能直接适应专业匹配系统。在现实世界中，专业匹配是一个高度约束的问题，一些现有的作品使用各种方法研究了这种相应的优化问题。对于考试，Guervós等人。

[29]结合了贪婪和进化算法[30,31]将文件分配给审阅者。

Karimzadehgan等。[5]和Taylor

[32]将其投射为整数线性编程（ILP）问题，因此任何ILP求解器都可以找到近似解决方案。太阳等。

[7]通过域知识的混合方法解决审阅者分配。最近，一些系统

Table 1
Notations.

Symbol	Description
m	次数
n	次疑问题
v	题目
v_i	候选专家
q	查询
q_j	一个专家
y	一个查询
y_i	主题
z	给定查询
t	概率
$t(v_i)$	给定查询
$T(v_i)$	关于专家VI
$T(Q)$	$T(Q)$ 的主要相关主题集的一组主要相关主题QJ

问题1（与约束匹配的专业知识）。给定一组专家 v 和一组查询 Q ，目标是通过满足某些约束来将 M 专家分配给每个查询，例如（1）每个专家的分配查询的数量应该在一个范围内 $[n_1, n_2]$ ，其中 n_1 到 n_2 ；（2）专家的主要主题应涵盖查询的相关主题；（3）任务应避免某种利益冲突（COI）。

实际上在某些应用中，满足约束比与查询的匹配专业知识更重要。例如，在会议论文审阅者分配中，不应分配文件的作者来审查自己的论文。这必须是一个艰难的约束。虽然在一些其他场景中，约束相对较软，例如专家之间的负载平衡。对每个专家的分配查询的数量可以在 N_1 和 N_2 之间的范围内。在现有的作品中，Dumais等人。[2]和mimmo等。

[4]主要专注于提高专业知识匹配的准确性，但忽略了如何获得满足各种约束的最佳匹配。Karimzadehgan等。

[5]使用整数线性编程来找到与约束匹配的专业知识的解决方案。但是，所提出的模型无法保证最佳解决方案。在这项工作中，我们提出了一个普遍的优化框架来解决这个问题。各种约束也可以包含在框架中。

3. Problem formulation

在本节中，我们首先给出几个必要的定义，然后提出问题的正式定义。鉴于一组专家 V ，每个专家对所有主题都有不同的专业知识。从形式上讲，我们假设有关方面的专业知识（称为主题），每个专家 v_i 在不同主题上具有不同的专业知识水平。此外，鉴于第1/4季度/ q_j 的一组查询，每个查询还与多个主题相关。鉴于此，我们首先定义了主题模型的概念。

定义1（主题模型）。专家（或查询）的主题模型 Y 是单词 w_1, w_2, \dots, w_g 的多项分布。每个专家（查询）被视为多主题模型的混合。该模型的假设是根据与专家（查询）相关的单词根据关于每个主题的词分布，即 $P(w_i|v_i)$ 。因此，分布中具有最高概率的单词将建义主题所代表的神学。

4.受限的优化框架

在本节中，我们提出了一种基于约束的优化框架，用于专业匹配。我们开发了一种高效的算法，基于凸成本流理论来解决优化框架，并呈现一个在线匹配算法，以实时合并用户反馈。基本想法。我们方法的基本思想是在约束优化框架中制定这个问题。不同的共同之躯可以在目标函数中正式化，或者直接被视为优化解决过程中的约束。为了解决优化框架，我们将问题转换为凸起成本网络流量问题，并提出了一种保证最佳解决方案的有效算法。

假设我们有 T 主题，主题 Z 和 A 。专业 v_i 的专业知识学位 za_f 表示为具有 p_{zyvz} 的概率 $YVIZ$ 。类似地，对于每个查询，我们也具有与 p_{zyqjz} 的 t 维主题分布。1.总结了符号

表1很容易理解，每个查询 Q_j 都可以表示为单词序列，即 $d_{q,j}$ 。要代表每个专家 v_i ，没有普遍性，我们也将其视为一系列单词，即 DVi 。基于此表示，我们可以计算每个查询和每个专家之间的相似性（或相关性分数）使用余弦相似性或语言模型等措施。鉴于此，我们可以定义我们与各种限制的专业知识问题。

4.1. The framework

现在，我们详细解释了提出的方法。一般来说，我们的目标可以从两个角度看，最大化专家和疑问之间的匹配分数和令人满意

the given constraints. Formally, we denote the set of experts to answer query q_j as $V(q_j)$, and the set of queries assigned to expert v_i as $Q(v_i)$. Further, we denote the matching score between expert v_i and query q_j as R_{ij} . Therefore, a basic objective function can be defined as follows:

$$\text{Max} \sum_{v_i \in V} \sum_{q_j \in Q(v_i)} R_{ij} \quad (1)$$

The objective function can be equivalently defined as $\sum_{q_j \in Q} \sum_{v_i \in V(q_j)} R_{ij}$. In different applications, the constraints will be different. Here we use several general constraints to explain how the proposed framework can incorporate different constraints.

The first constraint is that each query should be assigned with exactly m experts. For example, in the paper-reviewer assignment task, each paper should be assigned with 3 or 5 reviewers. This constraint can be directly added into the optimization problem. Formally, we have

$$\text{ST1: } \forall q_j \in Q, |V(q_j)| = m \quad (2)$$

The second constraint is called as *expert load balance*, indicating that each expert can only answer a limited number of queries. There are two ways to achieve this purpose: define a *strict constraint* or add a *soft penalty* to the objective function.

For *strict*, we add a constraint indicating that the number of assigned queries to every expert v_i should be equal or larger than a minimum number n_1 , but be equal or smaller than a maximum number n_2 . The *strict* constraint can be written as

$$\text{ST2 (strict): } \forall v_i \in V, n_1 \leq |Q(v_i)| \leq n_2 \quad (3)$$

The other way is to add a soft penalty to the objective function (Eq. (1)). For example, we can define a square penalty as $|Q(v_i)|^2$. By minimizing the sum of the penalty $\sum_i |Q(v_i)|^2$, we can achieve a soft load balance among all experts, i.e.:

$$\text{softpenalty: Min} \sum_{v_i \in V} |Q(v_i)|^2 \quad (4)$$

These two methods can be also used together. Actually, in our experiments, soft penalty method gives better results than strict constraint. Combining them together can yield a further improvement.

The third constraint is called *authority balance*. In real application, experts have different expertise level (authoritative level). Take the paper-reviewer assignment problem as an example. Reviewers may be divided into two levels: senior reviewers and average reviewers. Intuitively, we do not hope that the assigned reviewers to a paper are all average reviewers. It is desirable that the senior reviewers can cover all papers to guide (or supervise) the review process. Without loss of generality, we divide all experts into K levels, i.e., $V^1 \cup V^2 \cup \dots \cup V^K = V$, with V^1 representing experts of the highest authoritative level. Similar to *expert load balance*, we can define a strict constraint like $|V^1 \cap V(q_j)| \geq 1$, and also add a penalty function to each query q_j over the k -level experts. Following, we give a simple method to instantiate the penalty function:

$$\text{ST3: Min} \sum_{k=1}^K \sum_{j=1}^N |V^k \cap V(q_j)|^2 \quad (5)$$

Besides the above constraints, we also wish to assign experts that can cover all topics in the query. In the paper-reviewer assignment problem, for example, a paper may be related to several research areas, thus ideally in a comprehensive assignment, the paper is reviewed by a group of experts which cover all of the topics.

We introduce *related topics* of queries/experts. The related topics of query q_j and expert v_i are denoted as $T(q_j)$ and $T(v_i)$, indicating the most relevant aspects to the query and expert, respectively. $T(q_j)$ and $T(v_i)$ can be determined in different ways. Again in the paper-reviewer assignment problem, authors may be required to select the related topics of their papers from pre-defined categories in the submission, and reviewers can also select their expertise topics. In addition, it is also possible to estimate related topics from the learned topic distributions: (a) select top- k topics in θ_{v_i} , θ_{q_j} as the corresponding related topics; or (b) use thresholds τ_v and τ_q to prune topics, i.e., related topics are determined by $\theta_{v_i,z} > \tau_v$ and $\theta_{q_j,z} > \tau_q$.

Follow the work [3], we can incorporate different evaluation measures for topic covering. One measure is called *Coverage*, as we hope that assigned experts can cover different topics of a given query, i.e.,

$$\text{Coverage}(q_j) = \frac{|T(q_j) \cap \bigcup_{v_i \in V(q_j)} T(v_i)|}{|T(q_j)|} \quad (6)$$

In this way, an optimal assignment should maximize the *Coverage*, e.g., $T(q_j) \subseteq \bigcup_{v_i \in V(q_j)} T(v_i)$. But this becomes the NP-hard set cover problem which is intractable. Consequently, we choose to find less optimal solutions by making further assumptions: for a specific query q_j , each assigned expert $v_i \in V(q_j)$ can select only one *responsible topic* $\hat{T}_{q_j}(v_i) \in T(v_i)$, and covers this topic for the query. The optimal solution under the assumption provides a lower bound of the original problem. Finding the matching maximizing the coverage with responsible topics actually opens another optimization problem, but fortunately this can be incorporated into our framework. We leave the discussion to Section 5.2.

Another measure proposed is called *Confidence*, as we prefer the related topics to be covered by as many experts as possible, i.e.,

$$\text{Confidence}(q_j) = \frac{1}{m} \sum_{v_i \in V(q_j)} \frac{|T(q_j) \cap T(v_i)|}{|T(q_j) \cap \bigcup_{v_i \in V(q_j)} T(v_i)|} \quad (7)$$

Generally there is a *Coverage-Confidence* tradeoff. To achieve both high coverage and high confidence, a measure *Average Confidence* is accordingly defined as:

$$\text{AverageConfidence}(q_j) = \frac{1}{m} \sum_{v_i \in V(q_j)} \frac{|T(q_j) \cap T(v_i)|}{|T(q_j)|} \quad (8)$$

We finally choose the average confidence as the fourth constraint in our optimization framework.

$$\text{ST4: Max} \sum_{q_j \in Q} \sum_{v_i \in V(q_j)} \frac{|T(q_j) \cap T(v_i)|}{|T(q_j)|} \quad (9)$$

The last constraint is called *COI avoidance*. In many cases, we need to consider the conflict-of-interest (COI) problem. For example, an author, of course, should not review his own or his coauthors' paper. This can be accomplished through employing a binary $M \times N$ matrix U . An element with value of 0, i.e., $U_{ij}=0$, represents expert v_i has the conflict-of-interest with query q_j . A simple way is to multiply the matrix U with the matching score R in (Eq. (1)).

Finally, by incorporating different constraints in Eqs. (4)–(9) and the COI matrix U into the basic objective function (Eq. (1)), we can result in a constraint-based optimization framework, e.g.:

$$\begin{aligned} \text{Max} \sum_{v_i \in V} \sum_{q_j \in Q(v_i)} U_{ij} R_{ij} - \sum_{k=1}^K \left(\mu_k \sum_{j=1}^N |V^k \cap V(q_j)|^2 \right) \\ - \beta \sum_{v_i \in V} |Q(v_i)|^2 + \lambda \sum_{q_j \in Q} \sum_{v_i \in V(q_j)} \frac{|T(q_j) \cap T(v_i)|}{|T(q_j)|} \\ \text{s.t. } \forall q_j \in Q, |V(q_j)| = m \\ \forall v_i \in V, n_1 \leq |Q(v_i)| \leq n_2 \end{aligned} \quad (10)$$

给定的约束。正式, 我们表示, 将查询 Q_j 作为 v_i (Q_j) 以及分配给 Expert VI 的查询集作为 $Q(v_i)$ 。此外, 我们表示专家 VI 和查询 Q_j 之间的匹配分数作为 R_{ij} 。因此, 基本目标函数可以定义如下:

$$\text{Max} \sum_{v_i \in V} \sum_{q_j \in Q(v_i)} R_{ij} \quad (1)$$

客观函数可以等同地定义为 $P(Q_j) \cap Q(v_i)$, 其中 v_i 是不同的应用中的约束

会是不同的。在这里, 我们使用几个一般约束来解释所提出的框架如何包含不同的约束。第一个约束是每个查询都应以完整的专家分配。例如, 在纸张审阅者分配任务中, 应使用3或5名审核员分配每份纸张。可以直接添加该约束在优化问题中。正式, 我们有

我们介绍查询/专家的相关主题。查询 Q_j 和专家 VI 的相关主题表示为 $T(Q_j)$ 和 $T(VI)$, 分别表示查询和专家的最相关方面。 $T(Q_j)$ 和 $T(VI)$ 可以以不同的方式确定。再次在纸张审阅者分配问题中, 可能需要提交人从提交中预定义的类别选择文件的相关主题, 审阅者也可以选择他们的专业知识。此外, 还可以从学习主题分发中估计相关主题: (a) 在 YVI , YQ_j 中选择 Top-K 主题作为相应的相关主题; 或 (b) 使用阈值电视和 tq to prune 主题, 即相关主题由 $YVIZ$ 4TV 和 $YQJZ$ 4TQ 确定。按照工作[3], 我们可以纳入不同的评估措施, 以了解主题覆盖。一项措施称为覆盖, 我们希望指定的专家可以涵盖给定查询的不同主题, 即,

$$\text{Coverage}(q_j) = \frac{|T(q_j) \cap \bigcup_{v_i \in V(q_j)} T(v_i)|}{|T(q_j)|} \quad (6)$$

通过这种方式, 最佳任务应最大化覆盖率, 例如, tq_j dsvia q_j tvi 。但这成为了难以相容的NP-Hard Set涵盖问题。因此, 我们选择通过进一步假设找到更少的最佳解决方案: 对于特定的查询 Q_j , 每个分配的专家 V_i 只能选择一个响应的主题 tq_j vi at vi , 并涵盖查询的本主题。假设下的最佳解决方案提供了原始问题的下限。找到匹配最大化的覆盖范围与负责任主题实际上打开另一个优化产品, 但幸运的是, 这可以纳入我们的框架。我们将讨论留给5.2节。提出的另一项措施被称为信心, 因为我们更喜欢与尽可能多的专家涵盖的相关主题, 即,

$$\text{ST1: } \forall q_j \in Q, |V(q_j)| = m \quad (2)$$

第二个约束被称为专家负载余额, 表明每个专家只能应答有限数量的查询。有两种方法可以实现此目的: 定义严格的约束或向目标函数添加软罚款。对于严格, 我们添加了一个约束, 指示每个专家 VI 的分配查询的数量应该等于或大于最小数字 $N1$, 但是等于或小于最大数量 $N2$ 。严格的约束可以写成

$$\text{ST2 (strict): } \forall v_i \in V, n_1 \leq |Q(v_i)| \leq n_2 \quad (3)$$

另一种方法是向目标函数 (EQ. (1)) 添加软罚款。例如, 我们可以定义 AS9Q VI 92 通过最小化罚款 $Pi9q vi 92$ 的总和, 我们可以在所有专家中实现软载余额, 即:

$$\text{Confidence}(q_j) = \frac{1}{m} \sum_{v_i \in V(q_j)} \frac{|T(q_j) \cap T(v_i)|}{|T(q_j) \cap \bigcup_{v_i \in V(q_j)} T(v_i)|} \quad (7)$$

通常存在覆盖置信措施。为实现高覆盖率和高信任, 因此衡量平均置信度被定义为:

$$\text{AverageConfidence}(q_j) = \frac{1}{m} \sum_{v_i \in V(q_j)} \frac{|T(q_j) \cap T(v_i)|}{|T(q_j)|} \quad (8)$$

我们终于选择了我们优化框架中的第四个限制的平均信心。

$$\text{ST4: Max} \sum_{q_j \in Q} \sum_{v_i \in V(q_j)} \frac{|T(q_j) \cap T(v_i)|}{|T(q_j)|} \quad (9)$$

最后的约束被称为 COI 避免。在许多情况下, 我们需要考虑利益冲突 (COI) 问题。例如, 当然, 作者不应该审查自己或他的共同咨询纸。这可以通过采用二进制 $M \times N$ 矩阵 U 来完成。具有 0, 1, 1, 0 的元素, 表示专家 VI 对查询 Q_j 具有兴趣冲突。一种简单的方法是将矩阵 U 乘以匹配得分 r (eq. (1))。最后, 通过在 EQ S 中纳入不同的约束。(4) – (9) 和 COI 矩阵 U 进入基本目标函数 (EQ. (1)) , 我们可以导致基于约束的优化框架, 例如:

$$\text{ST3: Min} \sum_{k=1}^K \sum_{j=1}^N |V^k \cap V(q_j)|^2 \quad (5)$$

除了上述约束之外, 我们还希望分配可以涵盖查询中所有主题的专家。例如, 在纸张审阅者分配问题中, 一篇论文可能与若干研究领域有关, 因此理想情况下, 该论文由一组专家审查, 涵盖了所有主题。

$$\begin{aligned} \text{Max} \sum_{v_i \in V} \sum_{q_j \in Q(v_i)} U_{ij} R_{ij} - \sum_{k=1}^K \left(\mu_k \sum_{j=1}^N |V^k \cap V(q_j)|^2 \right) \\ - \beta \sum_{v_i \in V} |Q(v_i)|^2 + \lambda \sum_{q_j \in Q} \sum_{v_i \in V(q_j)} \frac{|T(q_j) \cap T(v_i)|}{|T(q_j)|} \\ \text{s.t. } \forall q_j \in Q, |V(q_j)| = m \\ \forall v_i \in V, n_1 \leq |Q(v_i)| \leq n_2 \end{aligned} \quad (10)$$

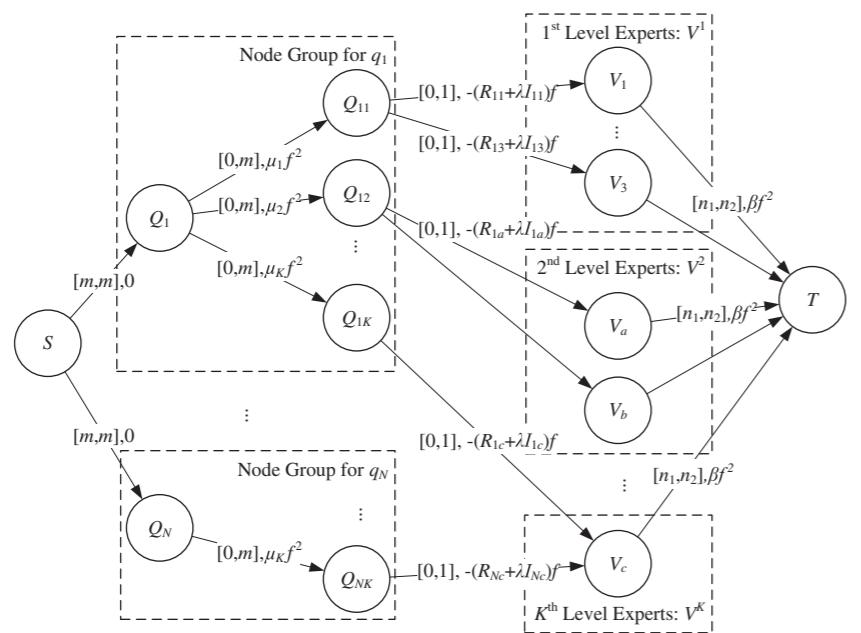


Fig. 2. The construction of convex cost network flow according to objective function (10).

where λ , β and μ_k are lagrangian multipliers, used to tradeoff the importance of different components in the objective function.

Now the problem is how to define the topic distribution θ , how to calculate the pairwise matching score R_{ij} , and how to optimize the framework.

4.2. Modeling multiple topics

The goal of topic modeling is to associate each expert v_i with a vector $\theta_{v_i} \in \mathbb{R}^T$ of T -dimensional topic distribution, and as well to associate each query q_j with a vector $\theta_{q_j} \in \mathbb{R}^T$. The topic distribution can be obtained in many different ways. For example, in the paper-reviewer assignment problem, each reviewer can select their expertise topics from a predefined categories. In addition, we can use statistical topic modeling [35,36] to automatically extract topics from the input data. In this paper, we use the topic modeling approach to initialize the topic distribution of each expert and each query.

To extract the topic distribution, we can consider that we have a set of M expert documents and N query documents (each representing an expert or a query). An expert's document can be obtained by accumulating the content information related to the expert. For example, we combine all publication papers as the expert document of a reviewer, thus expert v_i 's document can be represented as $d_{v_i} = \{w_{ij}\}$. Each query can also be viewed as a document. Then we can learn these T topic aspects from the collection of expert documents and query documents using a topic model such as LDA [36]. Specifically, let $D = \{d_{v_1}, \dots, d_{v_M}\}$ be the set of experts' documents. The log-likelihood of the whole collection according to LDA is

$$\log p(D|\theta, \phi) = \sum_{d \in D} \sum_{w \in d} c(w, d) \log \left(\sum_{z=1}^T p(w|z, \phi_z) p(z|d, \theta_d) \right) \quad (11)$$

where $c(w, d)$ is the count of word w in document d , $p(w|z, \phi_z)$ is the probability of topic z generating word w , and $p(z|d, \theta_d)$ is the probability of document d containing topic z .

We use the Gibbs sampling algorithm [37,38] to learn the topic distribution θ_{v_i} for each expert. The topic distribution of query θ_{q_j} can be accordingly inferred from the produced θ_{v_i} .

4.3. Pairwise matching score

We employ a language model-based retrieval method to calculate the pairwise matching score. With language model, the matching score R_{ij} between expert v_i and query q_j is interpreted as a probability $R_{ij}^{LM} = p(q_j|d_i) = \prod_{w \in q_j} p(w|d_i)$, where

$$p(w|d_i) = \frac{N_{d_i}}{N_{d_i} + \lambda_D} \cdot \frac{tf(w, d_i)}{N_{d_i}} + \left(1 - \frac{N_{d_i}}{N_{d_i} + \lambda_D}\right) \cdot \frac{tf(w, D)}{N_D} \quad (12)$$

where N_{d_i} is the number of word tokens in document d_i , $tf(w, d_i)$ is the number of occurring times of word w in d_i , N_D is the number of word tokens in the entire collection, and $tf(w, D)$ is the number of occurring times of word w in the collection D . λ_D is the Dirichlet smoothing factor and is commonly set according to the average document length in the collection [26].

Our previous work extended LDA and proposed the ACT model [39] to generate a topic distribution. By considering the learned topic model, we can define another matching score as

$$R_{ij}^{ACT} = p(q_j|d_i) = \prod_{w \in q_j} \sum_{z=1}^T p(w|z, \phi_z) p(z|d, \theta_d) \quad (13)$$

Further, we can define a hybrid matching score by combining the two probabilities together

$$R_{ij}^H = R_{ij}^{LM} \times R_{ij}^{ACT} \quad (14)$$

4.4. Optimization solving

In order to maximize the objective function (Eq. (10)), we construct a convex cost network with lower and upper bounds imposed on the arc flows. Fig. 2 illustrates the constructing process as described in Algorithm 1.¹ Convex cost flow problem can be solved by transforming to an equivalent minimum cost flow problem [40]. The minimum cost flow of the network gives an optimal assignment with respect to (Eq. (10)).

¹ Every arc in the network is associated with lower and upper bound denoted as $[l, u]$ and a convex function of the arc flow f .

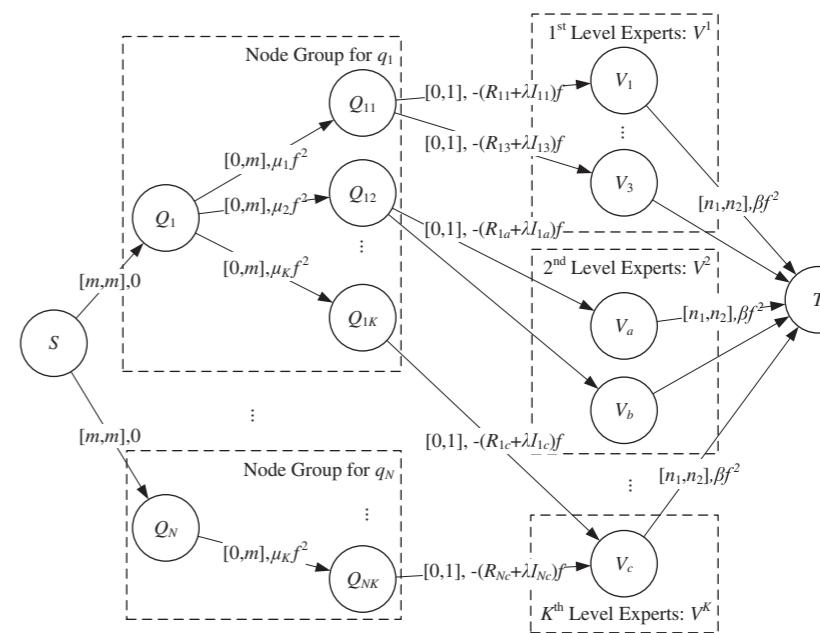


图2. 根据客观函数(10)的凸起成本网络流的构造。

其中L, B和MK是Lagrangian乘法器，用于在目标函数中对不同组件的重视进行权衡。现在问题是如何定义主题分发y，如何计算成对匹配分数RIJ，以及如何优化框架。

4.3. Pairwise matching score

4.3. Pairwise matching score

我们采用基于语言模型的检索方法来计算成对匹配分数。通过语言模型，Expert VI和查询QJ之间的匹配分数RIJ被解释为概率RLM IJ = P(qj|di) = qwaqjp / w9di

$$d_i \quad D \quad d_i \quad D \quad D$$

其中NDI是DID
DI, TF(W, DI)中的单词令牌的数量是DI中的单词W的发生时间，NDI是整个集合中的单词令牌的数量，而TF(W, D)是收集D。
LD中的单词W的发生次数是Dirichlet平滑因子，并且通常根据集合中的平均文档长度进行设置[26]。我们以前的工作扩展了LDA并提出了该法案

[39]要生成主题分发。通过考虑学习的主题模型，我们可以将另一个匹配分数定义为

$$R_{ij}^{ACT} = p(q_j|d_i) = \prod_{w \in q_j} \sum_{z=1}^T p(w|z, \phi_z) p(z|d, \theta_d) \quad (13)$$

此外，我们可以通过将这两个概率组合在一起定义混合匹配分数

$$R_{ij}^H = R_{ij}^{LM} \times R_{ij}^{ACT} \quad (14)$$

4.4. Optimization solving

为了最大化目标函数(EQ. (10))，我们构建一个凸起成本网络，其施加在电弧流上的下限和上限。图2示出了如算法1.1中所述构建过程1.1凸起成本流程问题可以通过转换为等同的最小成本流动问题来解决[40]。网络的最小成本流程对(EQ. (10))提供了最佳分配。

1. 网络中的每个弧与弧形和上束相关联，表示为[l, u]和电弧流f的凸起函数。

Algorithm 1. Optimization solving algorithm.

Input: The set of experts V ; the set of queries Q ; the matching score matrix $R_{M \times N}$; the COI matrix $U_{M \times N}$; Number of expertise level K ; m, n_1, n_2 as described above.
Output: An assignment of experts to queries maximizing objective function (10).

- 1.1 Create a network G with source node S and sink node T ;
- 1.2 **foreach** $q_j \in Q$ **do**
- 1.3 Create $K+1$ nodes, denoted as $Q_j, Q_{j1}, \dots, Q_{jk}$ respectively;
- 1.4 Add an arc from source node S to node Q_j , with zero cost and flow constraint $[m, m]$;
- 1.5 Add an arc from node Q_j to Q_{jk} , with square cost function $\mu_k f^2$ and flow constraint $[0, m]$;
- 1.6 **foreach** $v_i \in V$ **do**
- 1.7 Create a node V_i ;
- 1.8 Add an arc from V_i to sink node T , with square cost function βf^2 and flow constraint $[n_1, n_2]$;
- 1.9 **foreach** $v_i \in V, q_j \in Q$, s.t. $U_{ij} = 1$ **do**
- 1.10 $k =$ expert level of v_i ;
- 1.11 Add an arc from Q_{jk} to V_i , with linear cost function $-(R_{ij} - \lambda I_{ij})f$ and flow constraint $[0, 1]$;
- 1.12 Compute the minimum cost flow on G ;
- 1.13 **foreach** $v_i \in V, q_j \in Q$, s.t. $U_{ij} = 1$ **do**
- 1.14 $k =$ expert level of v_i ;
- 1.15 **if** flow $f(Q_{jk}, V_i) = 1$ **then** Assign query q_j to expert v_i ;

Theorem 1. Algorithm 1 based on minimum convex cost flow gives an optimal solution.

Proof. First the minimum convex cost flow problem (MCCF) can be formulated as the following optimization problem:

$$\begin{aligned} \text{Min } & \sum_{(a,b) \in E(G)} C_{ab}(f(a,b)) \\ \text{s.t. } & \forall a \in V(G), \sum_{b:(a,b) \in E(G)} f(a,b) = \sum_{b:(b,a) \in E(G)} f(b,a) \\ & \forall (a,b) \in E(G), l_{ab} \leq f(a,b) \leq u_{ab} \end{aligned} \quad (15)$$

The model is defined on a directed network $G = (V(G), E(G))$ with lower bound l_{ab} , upper bound u_{ab} and a convex cost function $C_{ab}(f(a,b))$ associated with every arc (a,b) .

Now we prove that minimizing (Eq. (15)) on the graph G constructed in Algorithm 1 is equivalent to maximizing (Eq. (10)). For simplicity, we use I_{ij} to denote $|T(q_j) \cap T(v_i)| / |T(q_j)|$. For the constructing process, we see a feasible flow on G is mapping to a query-expert assignment. The flow from S to Q_j indicates the number of experts assigned with query q_j , and the flow from V_i to T indicates the number of queries assigned to expert v_i . And the cost between V_i and T is corresponding to the load balance soft penalty function (Eq. (4)). The meaning of the flow from Q_j to V_i is the number of k -th-level experts assigned to q_j , thus we impose a square cost function $\mu_k f^2$ on the arcs which is equivalent to the negative of the authority balance penalty. The flow from Q_{jk} to V_i means we assign query q_j to expert v_i , it is easy to find that no query will be assigned to the same expert twice since we give an upper bound of 1 on the arc, while the cost is equivalent to the negative of matching score and topic average confidence score. Therefore, our problem can be reduced to an equivalent MCCF problem, where the objective

function of MCCF problem (Eq. (15)) is the negative form of (Eq. (10)).

In practice, it is not necessary to add all (Q_{jk}, V_i) arcs. To further reduce the complexity of the algorithm, we first greedily generate an assignment and preserve corresponding arcs, then keep only $c \cdot m$ arcs for Q_{jk} and $c \cdot n_2$ arcs for V_i which have highest matching score (c is a fixed constant). We call this process Arc-Reduction, which will reduce the number of arcs in the network without influencing the performance too much. To process large scale data, we can leverage the parallel implementation of convex cost flow [41].

4.5. Online matching

After an automatic expertise matching process, the user may provide feedbacks. Typically, there are two types of user feedbacks: (1) pointing out a false match; (2) specifying a new match. Online matching aims to adjust the matching result according to the user feedback. One important requirement is how to perform the adjustment at real time. In our framework, we provide online interactive adjustment without recalculating the whole cost flow. For both types of feedbacks, we can easily accomplish online adjustment by canceling some flows and augmenting new assignments in our framework. We give Algorithm 2 to consider the first type of feedback, which still produces an optimal solution.

Algorithm 2. Online matching algorithm.

Input: A minimum cost network flow f on G corresponding to the current assignment; an inappropriate match (v_i, q_j) .

Output: A new assignment.

- 2.1 $k =$ expert level of v_i ;
- 2.2 **if** $f(Q_{jk}, V_i) = 1$ **then**
- 2.3 Construct the residual network $G(f)$;
- 2.4 Compute the shortest path P_{back} from T to S on $G(f)$ which contains backward arc (V_i, Q_{jk}) ;
- 2.5 Cancel(rollback) 1 unit of flow along P_{back} and update $G(f)$;
- 2.6 Remove arc (Q_{jk}, V_i) from G and update $G(f)$;
- 2.7 Compute shortest augmenting path P_{aug} from S to T ;
- 2.8 Augment 1unit of flow along P_{aug} ;

Lemma 1 (Negative cycle optimality conditions). Ahuja et al. [40] A feasible solution f^* is an optimal solution of the minimum cost flow problem if and only if it satisfies the negative cycle optimality conditions: namely, the residual network $G(f^*)$ contains no negative cost cycle.

Theorem 2. Algorithm 2 produces an optimal solution in the network without assignment (q_j, v_i) .

Proof. According to Lemma 1, the residual network $G(f)$ contains no negative cost cycle since the given flow f has the minimum cost. In Algorithm 2, we remove the inappropriate match (v_i, q_j) and adjust the network flow in line (2.3)–(2.5). Denote the feasible flow in the network after line (2.5) as f' . According to the SAP (short augmenting path) algorithm of cost flow, if f' has the minimum cost (i.e., $G(f')$ contains no negative cycle), the algorithm will give the optimal solution. We show the optimality of f' by contradiction. Assume $G(f')$ contains a negative cycle C , C must intersect with the shortest path P_{back} computed online (2.3), since the original $G(f)$ contains no negative cycle. Thus merging C into path P_{back} will generate a shorter path, which contradicts with the assumption that P_{back} is shortest. Therefore, f' has the minimum

算法1.优化求解算法。

输入：专家 v ;查询 Q ;匹配得分矩阵 $R_{m \times n}$;COI矩阵 $U_{m \times n}$;专业知识级别 k ;如上所述的 $M, N1, N2$ 。产出：向询问专家的分配最大化客观函数 (10)。1.1 使用源节点 S 和宿节点 T 创建网络 G ;1.2 foreach $q_j \in Q$ do 1: 3 1: 4 1: 5

创建 $k+1$ 节点，表示为 $q_j, q_{j1}, \dots, q_{jk}$ ；从源节点 S 向节点 Q_j 添加一个弧，零成本和流量约束 m, m ；从节点 Q_j 到 Q_{jk} 的电弧，方形成本函数 MKF 和流量约束 $0, m$ ；

4.5. Online matching

在自动专业知识匹配过程之后，用户可以提供反馈。通常，有两种类型的用户反馈：(1)指出假匹配；(2)指定新匹配。在线匹配旨在根据用户反馈调整匹配结果。一个重要要求是如何实时进行调整。在我们的框架中，我们提供在线交互式调整，而无需重新计算整个成本流程。对于这两种类型的反馈，我们可以通过取消某些流程并在框架中增强新分配来轻松完成在线调整。我们给予

1.6 FOREACH VI AV DO 1: 7
1: 8 创建节点VI;将VI的电弧添加到沉积节点T, 方形成本函数 B F 2和流量约束 n_1, n_2 ；
1.9 FOREACH VI AV, Q_j AQ, S: T: UIJ 1做1:10 1:11
 K vi 的专家水平;从QJK向VI添加电弧，线性成本函数 rij lijj 和流量约束 $0, 1$ ；

1.12 计算 g 的最小成本流量;1.13 FOREACH VI AV, Q_j AQ, S: T: UIJ 1做1:14 1:15
 K vi 的专家水平;如果流 $f(Q_{jk}, V_i) = 1$, 则 v_i 1然后将查询 q_j 分配给专家 v_i ；

定理1.基于最小凸起成本流的算法1提供了最佳解决方案。

证明。首先，可以将最小凸起成本流量问题 (MCCF) 作为以下优化问题 (MCCF) 为：

$$\begin{aligned} \text{Min } & \sum_{(a,b) \in E(G)} C_{ab}(f(a,b)) \\ \text{s.t. } & \forall a \in V(G), \sum_{b:(a,b) \in E(G)} f(a,b) = \sum_{b:(b,a) \in E(G)} f(b,a) \\ & \forall (a,b) \in E(G), l_{ab} \leq f(a,b) \leq u_{ab} \end{aligned} \quad (15)$$

该模型在具有下限实验室，上限 UAB 和凸起成本函数 $C_{ab}(f(a,b))$ 的定向网络 $G = (V(G), E(G))$ 上定义，与每个弧 (a, b) 相关联。现在我们证明，在算法1中构造的图表G上最小化 (EQ. (15)) 相当于最大化 (EQ. (10))。为简单起见，我们使用IJ表表示 $9t_qj \setminus t_{vi} = 9t_qj - 9t_{vi}$ 。对于构建过程，我们看到G上的可射到查询专家分配。从S到 Q_j 的流量表示查询 Q_j 分配的专家数量，并且从VI到T的流量表示分配给专家VI的查询数。并且VI和T之间的成本与负载平衡软罚函数 (Eq. (4)) 相对应。从 Q_j 到 Q_{jk} 的流量的含义是分配给 Q_j 的 k 级专家的数量，因此我们在弧上强加一个平方成本函数 MKF 。2. 相当于权威余额惩罚的负面。从 QJK 到VI的流程意味着我们将查询 Q_j 分配给Expert VI，很容易发现没有查询将被分配给同一专家两次，因为我们在弧上提供1的上限，而成本相当于匹配分数和主题平均置信度得分。因此，我们的问题可以减少到相同的MCCF问题，其中目标

构建残余网络 $G = F$ ；计算来自T到S的最短路径 P_{Back} ，其中包含后向弧 vi, q_{jk} ；CANCELINGRALLBACK 1沿逆向的流量单位，更新 $G = F$ ；从 g 中删除弧 q_{jk}, vi 和更新 $g = f$ ；计算到 t 的最短增强路径波动；沿着 $Phug$ 的流量增加1个；

引理1 (负周期最优性条件)。Ahuja等。一种可行的解决方案 F_n 是且仅当它满足负周期最优性条件时，最小成本流动问题的最佳解决方案是：即，残差网络 $G(F_n)$ 不包含负成本循环。

定理2.算法2在没有分配的情况下在网络中产生最佳解决方案 (Q_j, VI)。

证明。根据LEMMA 1，由于给定的流 F 具有最小成本，因此残余网络 $G(F)$ 不包含负成本循环。在算法2中，我们删除不当匹配 (VI, Q_j)，并在线 (2.3) – (2.5) 调整网络流量。在线 (2.5) 作为 F 。

0.根据成本流量的SAP (短增强路径) 算法，如果 F 0具有最小成本 (即， $G(F_0)$ 不包含负周期)，则该算法将提供最佳解决方案。我们通过矛盾显示 F 0的最优性。假设 $G(F_0)$ 包含负周期 C ， C 必须与在线计算的最短路径送 (2.3) 相交，因为原始 $g(f)$ 不包含负周期。因此，将 C 合并到路径送方案将产生较短的路径，这与假设 P_{BACK} 最短的假设相矛盾。因此， F_0 具有最小值。

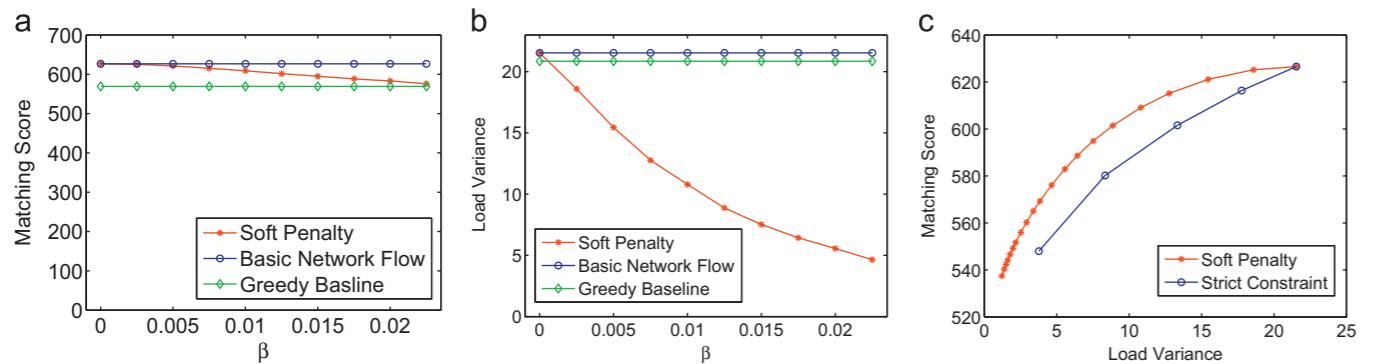


Fig. 3. (a) and (b) illustrate how soft penalty function influences the matching score (MS) and load variance with different β respectively. (c) gives a comparison between soft penalty function and strict constraint methods towards load balance.

cost. Accordingly, **Algorithm 2** gives the optimal solution after augmenting the new assignment. \square

5. Experimental results

The proposed approach for expertise matching is very general and can be applied to many application to align experts and queries. We evaluate the proposed framework on two different genres of expertise matching problems: paper-reviewer assignment and course-teacher assignment. Three experiments on different datasets are conducted to show the effectiveness of the proposed method. All datasets, code, and detailed results are publicly available.² All the experiments are carried out on a PC running Windows XP with Intel Core2 Quad CPU Q9550 (2.83 GHz), 4G RAM.

5.1. Paper-reviewer assignment experiment

Dataset. The paper-reviewer dataset consists of 338 papers and 354 reviewers. The reviewers are the program committee members of KDD'09 and the 338 papers are those published on KDD'08, KDD'09, and ICDM'09. For each reviewer, we collect her/his all publications from academic search system Arnetminer³ [42] to generate the expertise document. As for the COI problem, we generate the COI matrix U according to the coauthor relationship in the last five years and the organization they belong to. Finally, we set that a paper should be reviewed by $m=5$ experts, and an expert at most reviews $n_2=10$ papers.

Baseline methods and evaluation metrics. We employ a greedy algorithm as the baseline. The greedy algorithm assigns experts with highest matching score to each query, while keeping the load balance for each expert (i.e., $|Q(v_i)| \leq n_2$) and avoiding the conflict-of-interest.

As there are no standard answers, in order to quantitatively evaluate our method, we define the following metrics:

Matching score (MS): It is defined as the accumulative matching score.

$$MS = \sum_{v_i \in V} \sum_{q_j \in Q(v_i)} U_{ij} R_{ij}$$

Load variance (LV): It is defined as the variance of the number of papers assigned to different reviewers.

$$LV = \sum_{i=1}^M \left(|Q(v_i)| - \frac{\sum_{i=1}^M |Q(v_i)|}{M} \right)^2$$

Expertise variance (EV): It is defined as the variance of the number of top level reviewers assigned to different papers.

$$EV = \sum_{j=1}^N \left(|V(q_j) \cap V^1| - \frac{\sum_{j=1}^N |V(q_j) \cap V^1|}{N} \right)^2$$

Results. In this experiment, we tune different parameters to analyze the influence on the accumulative matching score. We also evaluate the efficiency of our proposed approach.

We first set $\mu=0$ and tune the parameter β to find out the effects of soft penalty function. Fig. 3(a) illustrates how soft penalty function influences the matching score with different β . We see that the matching score decreases slightly with β increasing. Fig. 3(b) shows the effects of load variance with β varied. We see that the load variance changes very fast toward balance.

In Fig. 3(c), we compare the two different methods to achieve load balance, namely, strict constraint and soft penalty. The two LV-MS curves are respectively generated by setting different minimum numbers n_1 for strict constraint and varying the weight parameter β for soft load balance penalty. The curves show that soft penalty outperforms strict constraint towards load balance.

Then we set β to 0 to test the effects of authority balance. Experts are divided into two levels base on their H -index, and we set $\mu_2=0$ to consider the balance of the senior reviewers only. Fig. 4 presents the accumulative matching score (a) and expertise variance (b) with μ_1 varied.

Further, we analyze the effects of different constraints. Specifically, we first remove all constraints (using Eq. (1) only), and then add the constraints one by one in the order (load balance, authority balance, and COI). In each step, we perform expertise matching using our approach. Table 2 lists the accumulative matching score obtained in each step. We see that the load balance constraint will reduce the expertise matching score, while the other constraints have little negative effect. This is because senior experts are often good at many aspects (the matching score between them and many queries are large), thus assigned with heavy load in traditional matching. In our method, we try to get a more reasonable assignment by adding load balance constraint, which will restrict the work load of those senior experts. As a result, the matching score decreases.

To clearly illustrate the effect of load-balance constraint, we present Fig. 5, in which we see that traditional information-retrieval based method assigns many papers to senior reviewers, while some reviewers do not get any work. The load-balance constraint is necessary to generate a reasonable matching.

Finally, we evaluate the efficiency performance of the proposed algorithm. We compare the CPU time of the original optimal algorithm and the version with Arc-Reduction. As shown in Fig. 6, the Arc-Reduction process can significantly reduce the time consumption.

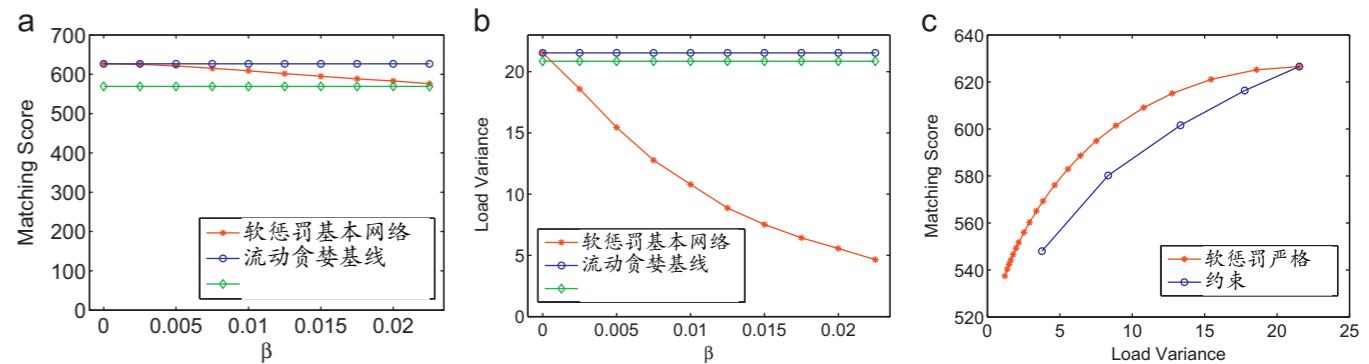


图3. (a) 和 (b) 示出了软惩罚功能如何影响匹配分数 (MS) 和与不同B的负载方差。

成本。因此，算法2在增强新分配后给出最佳解决方案。&

5. Experimental results

所提出的专业匹配方法非常一般，可以应用于许多应用程序来对准专家和查询。我们评估了两种不同类型的专业知识匹配问题的拟议框架：纸张-审阅者分配和课程-教师分配。进行了在不同数据集上进行的三个实验以显示所提出的方法的有效性。所有数据集，代码和详细结果都是公开可用的。

5.1. Paper-reviewer assignment experiment

数据集。 纸张审阅者数据集由338篇论文和354名审稿人组成。审阅者是KDD'09的计划委员会，338篇论文是KDD'08，KDD'09和ICDM'09发表的文件。对于每个审阅者，我们从学术搜索系统AR Netminers³ [42]中收集她/他的所有出版物，以产生专业知识文件。至于COI问题，我们根据过去五年的同志关系和他们所属的组织生成COI矩阵U。最后，我们设定了一篇论文，由M=5专家审查，以及大多数评论的专家N=10论文。基线方法和评估指标。我们使用贪婪的算法作为基线。贪婪算法将具有最高匹配分数的专家分配给每个查询，同时保持每个专家的负载余额（即， $9q_i/v_i - 9n_2$ ），并避免兴趣的混合。由于没有标准答案，为了定量评估我们的方法，我们定义了以下度量：匹配分数 (MS)：它被定义为累积匹配分数。

$$MS = \sum_{v_i \in V} \sum_{q_j \in Q(v_i)} U_{ij} R_{ij}$$

加载方差 (LV)： 它被定义为分配给不同审阅者的论文数量的方差。

$$LV = \sum_{i=1}^M \left(|Q(v_i)| - \frac{\sum_{i=1}^M |Q(v_i)|}{M} \right)^2$$

专业方案 (EV)： 它被定义为分配给不同论文的顶级审稿人数的差异。

$$EV = \sum_{j=1}^N \left(|V(q_j) \cap V^1| - \frac{\sum_{j=1}^N |V(q_j) \cap V^1|}{N} \right)^2$$

结果。在这个实验中，我们调整了不同的参数来分析对累积匹配分数的影响。我们还评估了我们提出的方法的效率。我们首先设置 $m=0$ 并调整参数 b 以找出软罚函数的效果。图。图3 (a)示出了软惩罚功能如何影响匹配分数与不同的 b 。我们看到匹配分数略微减少， B 增加。图。图3 (b)显示了与 B 变化的负载方差的影响。我们看到负载方差变化非常快，平衡。在图3 (c)中，我们比较两种不同的方法来实现负载平衡，即严格的限制和软罚性。两个LV-MS曲线分别通过设置不同的最小数量 N_1 来用于严格约束并改变重量参数 B 进行软载余额惩罚。曲线表明，软惩罚优于负载平衡的严格限制。然后我们将 b 设置为0以测试权威余额的影响。专家分为他们的H-Index的两个层面，我们设置了 $M=2$ ，以考虑仅考虑高级审阅者的余额。图4显示了 M_1 变化的累积匹配得分 (A) 和专业知识方差 (B)。此外，我们分析了不同约束的影响。具体而言，我们首先删除所有约束（仅使用Eq. (1)），然后按顺序（负载余额，权限余额和COI）逐个添加约束。在每一步中，我们使用我们的方法执行专业知识匹配。表2列出了每个步骤中获得的累积匹配分数。我们看到负载余额约束将减少专业匹配分数，而其他约束则具有很小的负面影响。这是因为在许多方面（它们之间的匹配分数和许多查询都是很大的），因此在传统匹配中分配了重负荷。在我们的方法中，我们尝试通过添加负载余额约束来获得更合理的分配，这将限制这些高级专家的工作负担。结果，匹配得分降低。为了清楚地说明负载平衡约束的影响，我们存在图5，我们看到传统信息检索的方法为高级审阅者分配许多文件，而一些审稿人则没有任何工作。负载余额约束是生成合理匹配的必要条件。最后，我们评估了所提出的算法的效率性能。我们比较原始最佳算法的CPU时间和带有弧度的版本。如图1所示，弧度减少过程可以显着降低时间消耗。

² <http://www.arnetminer.com/expertisematching>

³ <http://arnetminer.org>

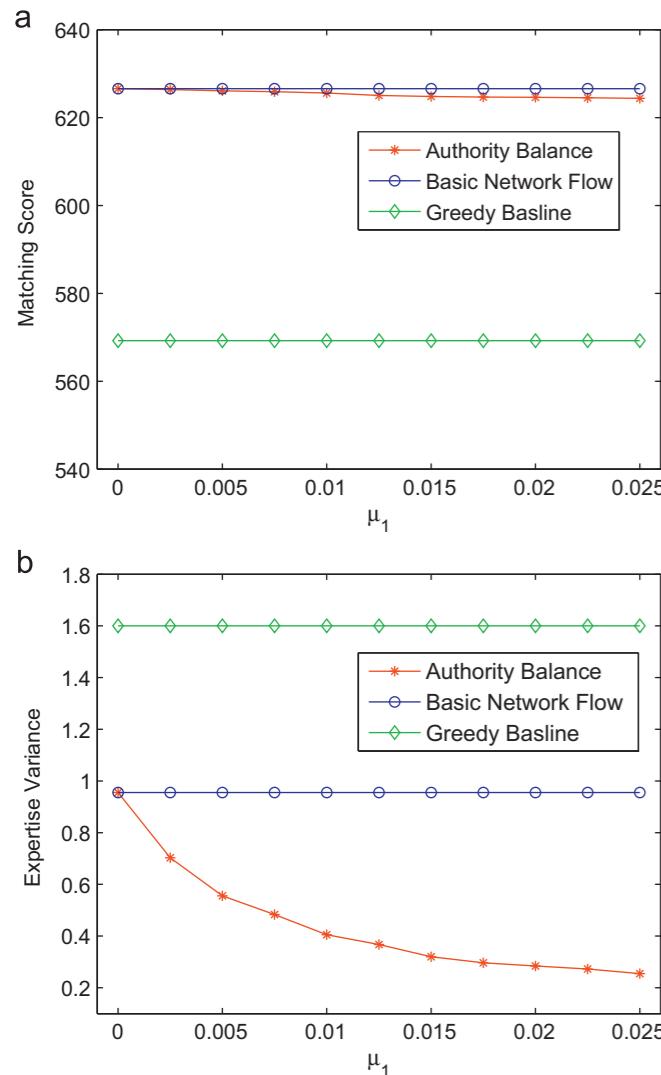


Fig. 4. Matching score (MS) and expertise variance (EV) with μ_1 varied.

Table 2
Effects of different constraints on matching score.

Constraint	Matching score
Basic objective function (Eq. (1))	635.51
+ Load balance soft penalty with $\beta = 0.02$	592.83
+ Authority balance with $\mu = (0.02, 0)^T$	599.37
+ COI	590.14

For example, when setting $c=12$ in this problem, we can achieve a $>3\times$ speedup without any loss in matching score.

We further use a case study (as shown in Tables 3 and 4) to demonstrate the effectiveness of our approach. We see that the result is reasonable. For example, Lise Getoor, whose research interests include relational learning, is assigned with a lot of papers about social network.

5.2. Multi-topic paper-reviewer assignment experiment

Dataset. We use another dataset (D2) to verify the performance on “topic coverage”. The dataset D2 is provided by [3],⁴ consisting

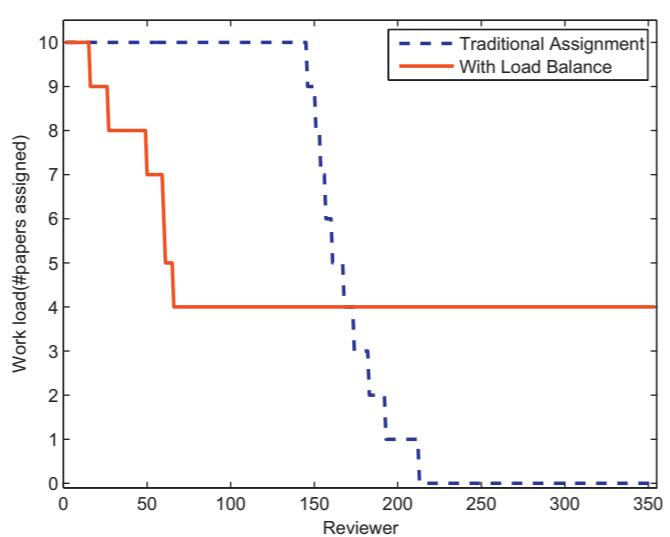


Fig. 5. The work load (number of paper assigned) of every reviewer.

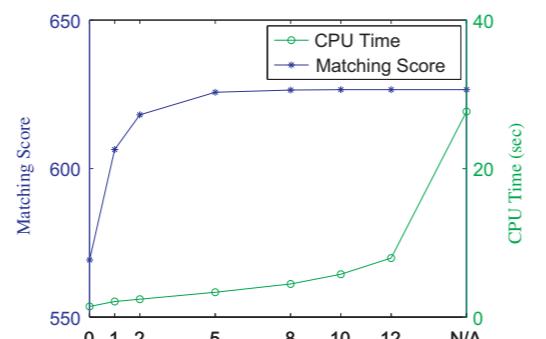


Fig. 6. Efficiency performance (s).

Table 3
Example assigned papers to three reviewers.

Reviewer	Assigned papers
Lise Getoor	Evaluating statistical tests for within-network classifiers of... Discovering organizational structure in dynamic social network Connections between the lines: augmenting social networks with text MetaFac: community discovery via relational hypergraph factorization Relational learning via latent social dimensions Influence and correlation in social networks
Wei Fan	Mining data streams with labeled and unlabeled training examples Vague one-class learning for data streams Active selection of sensor sites in remote sensing applications Name-ethnicity classification from open sources Consensus group stable feature selection Categorizing and mining concept drifting data streams
Jie Tang	Co-evolution of social and affiliation networks Influence and correlation in social networks Feedback effects between similarity and social influence... Mobile call graphs: beyond power-law and lognormal distributions Audience selection for online brand advertising: privacy-friendly...

of 73 queries and 189 reviewers. The 73 queries are paper abstracts from SIGIR'07, where each of them is related to at least two topics. The documents of reviewers are the combination of all

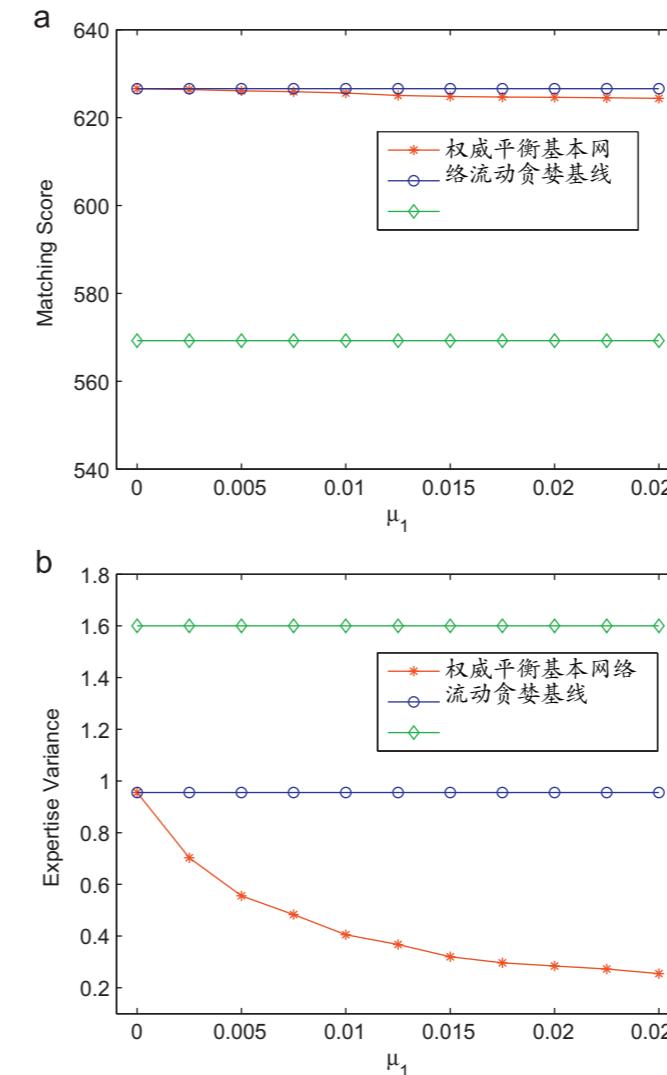


图4. 匹配分数 (MS) 和专业知识方差 (EV) 各种不同。

表2不同约束对匹配分数的影响。

Constraint	Matching score
基本目标函数 (eq. (1))	635.51
+ 负载平衡软惩罚与 $b_0:0.02$	592.83
+ COI	590.14

例如，当设置 $C=12$ 在这个问题中，我们可以在没有任何匹配分数的情况下实现 43% 的加速。我们进一步使用案例研究（如表3和4所示）来证明我们方法的有效性。我们看到结果是合理的。例如，Lise Getoor，其研究兴趣包括关系学习，被分配了很多关于社交网络的论文。

5.2. Multi-topic paper-reviewer assignment experiment

数据集。 我们使用另一个数据集 (D2) 来验证“主题覆盖”上的性能。数据集 D2 由 [3] 提供，其中 4 个查询由 73 个查询和 189 名审阅者组成。73 个查询是 Sigr'07 的纸质摘要，其中每个主题与至少两个主题有关。审稿人的文件是所有人的结合。

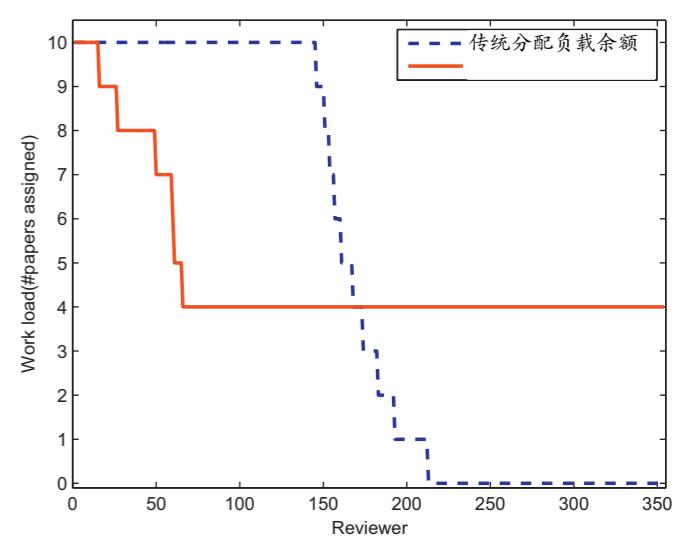


Fig. 5. The work load (number of paper assigned) of every reviewer.

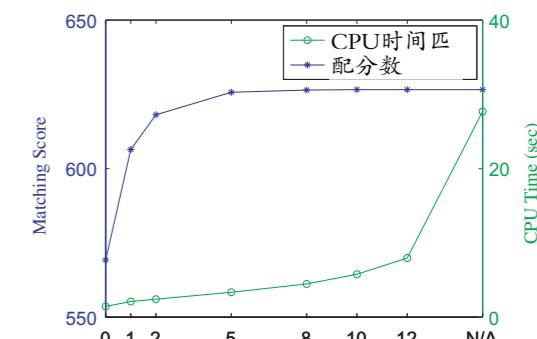


Fig. 6. Efficiency performance (s).

表3示例将文件分配给三位审阅者。

Reviewer Assigned papers

Lise 评估网络内的统计学测试，以便在线的动态社交网络连接中发现组织结构：通过文本 METAFAC 的增强社交网络；通过关系超照片分解关系学习通过潜在社会网络的影响和社交网络中的相关性

Wei 粉丝挖掘数据流具有标记和未标记的培训示例模糊单级学习数据流中的遥感应用中的传感器站点的活动选择来自开放消息来源的种族分类，共识组稳定特征选择分类和采矿概念漂移数据流

杰唐社会和隶属网络的共同演变在社交网络中的影响和相关性，相似性和社会影响力的反馈效应：超越幂律和逻辑正态分布版在线品牌广告的观众选择：隐私友好

⁴ <http://timan.cs.uiuc.edu/data/review.html>

score and the load balance constraint, we get the final configuration as shown Fig. 7.

Results. We exam our approach both with and without load balance constraint. The results of all methods on dataset D2 is presented in Table 5, from which we see that both our approach and Topic-KL outperform the two baseline methods on Coverage, Confidence, AverageConfidence and F_{score} . Interestingly, the baseline methods achieve better Confidence. This may be because the assigned reviewers of a paper are more likely to cover the same topic many times. Our approach without load balance achieves comparable performance with Topic-KL. It's noticeable when adding load balance constraint, the proposed method still produces good results, which validates the usefulness of our approach. In Table 5, we use our approach-L10, our approach-L5 to represent the proposed method with different load balance settings, i.e., $n_2=10$ and $n_2=5$.

Parameter sensitivity. Now we investigate the influences of the parameters in the framework. We first consider the multiplier σ , which controls the importance of Coverage score. When setting

Table 5
Comparison of all methods on four different measures: Coverage, Confidence, AverageCoverage and F_{score} . The number of topics is set to 20 when learning topic model.

Methods	Coverage (%)	Confidence (%)	AverageConfidence (%)	F_{score} (%)
Baseline-Pr	74	62	46	63
Baseline-KL	75	62	45	63
Topic-KL	87	58	53	67
Our approach	87	60	51	67
Our approach-L10	86	58	49	66
Our approach-L5	80	59	48	65

$\sigma=0$, the approach degenerates into the greedy baseline with poor performance. With a larger σ , both Coverage and Confidence increase and archive the best result at $\sigma=0.2$. It is also noticeable that the approach is not very sensitive to the parameter, as both Coverage and Confidence become stable with a relatively large σ (Fig. 8).

The number of topics also affects the performance. We employ the Gibbs sampling algorithm to learn topic models for different number of topics (e.g. 10, 30, 50 topics). In addition, we vary the number of related topics (i.e., $|T(q_j)|$ and $|T(v_i)|$) from 1 to 7, and $T(q_j), T(v_i)$ are determined by selecting top-k topics in θ_{q_j} and θ_{v_i} as we have discussed in Section 4.1. The sensitivity curves of Coverage and Confidence are plotted in Fig. 8 and Fig. 9. With the help of pairwise matching scores, Coverage and Confidence remain >80% and >54% even we set $T(q_j)=T(v_i)=1$ and topic number to be 10. Moreover, we see that too large or too small number of related topics do not produce good results. The appropriate number of related topics is about 3 or 4, which is fairly near the ground truth. Setting too few topics (e.g. 10) may hurt the final performance, but using enough topics (20, 30, 50) would not make big differences. One explanation is, small number of topics may limit the discriminative power of topics.

5.3. Course–teacher assignment experiment

Dataset. In the course–teacher assignment experiment, we manually crawled graduate courses from the department of Computer Science (CS) of four top universities, namely CMU, UIUC, Stanford, and MIT. In total, there are 609 graduate courses from the fall semester in 2008 to 2010 spring, and each course is instructed by 1–3 teachers. Our intuition is that teachers' research interest often match the graduate courses he/she is teaching. Thus we still use the teachers' recent (five years) publications as their

得分和负载余额约束, 我们得到了如图7的最终配置。7.结果。我们在没有负载余额约束的情况下考试我们的方法。DataSet D2上的所有方法的结果呈现在表5中, 我们看到我们的方法和主题-KL均优于覆盖, 信心, 平均值和FScore的两种基线方法。有趣的是, 基线方法实现更好的信心。这可能是因为纸张的指定审阅者更有可能多次涵盖相同的主题。我们没有负载余额的方法实现了主题-KL的可比性。在添加负载余额约束时, 它是明显的, 所提出的方法仍然可以提供良好的效果, 这验证了我们方法的有效性。在表5中, 我们使用我们的方法-L10, 我们的方法-15表示具有不同负载平衡设置的提出方法, 即N2_10和N2_5。参数灵敏度。现在我们研究了框架中参数的影响。我们首先考虑乘法器S, 控制覆盖率得分的重要性。设置时

$S \neq 0$, 方法退化为贪婪的基线, 性能差。具有较大的S, 覆盖率和置信度增加并归档S_0: 2的最佳结果。对于参数来说, 方法也不是非常敏感的, 因为覆盖率和置信度都具有相对较大的S(图8)。主题的数量也会影响性能。我们使用GIBBS采样算法来学习不同数量的主题(例如10,30,50主题)的主题模型。此外, 我们改变了从1到7的相关主题(即, 9t_qj_9和9t_vj_9)的数量, 而T(Qj), T(Vj)是通过选择YQJ和YVI中的Top-K主题来确定, 因为我们部分中讨论4.1。覆盖的遮盖力和置信度的灵敏度曲线在图4中绘制。参照图8和图9.在成对匹配分数的帮助下, 覆盖率和置信度仍然是480%和454%, 即使我们设置为T(Qj)_1(vi)_1和主题号是10。此外, 我们看到的相关主题数量太大或太多的相关主题不会产生良好的结果。适当数量的相关主题约为3或4, 这在实地真理附近。设置太少的主题(例如, 10)可能会损害最终表现, 但使用足够的主题(20,30,50)不会产生大差异。一个渐次进入是, 少数主题可能会限制主题的辨别力量。

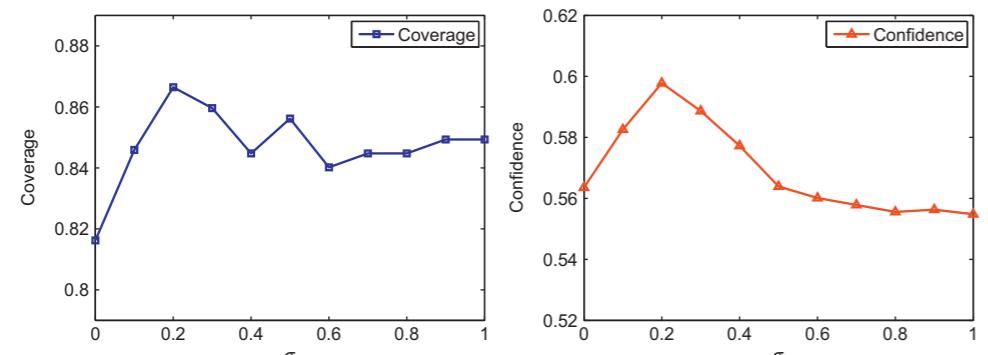


Fig. 8. Sensitivity of Coverage and Confidence to different Lagrangian multiplier σ .

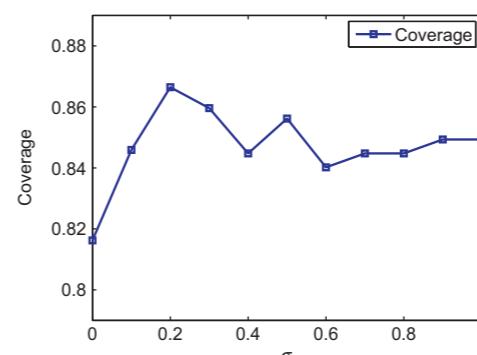


图8. 不同拉格朗日乘法器的覆盖和信心的敏感性。

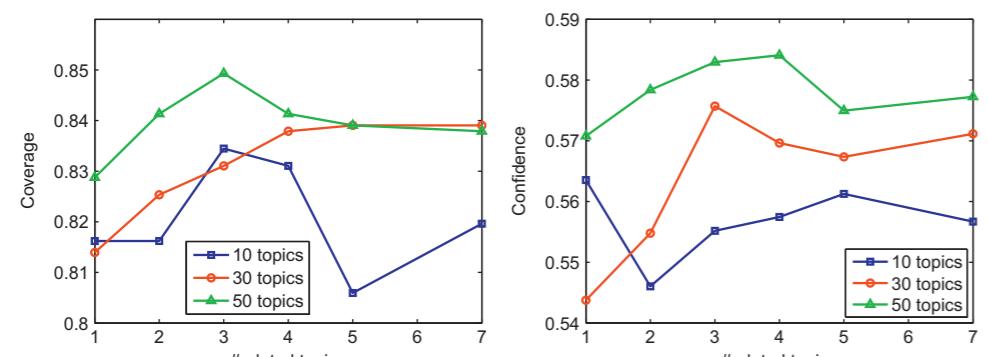


Fig. 9. Sensitivity of Coverage and Confidence to the number of related topics for different topic models, with $\sigma=0.3$.

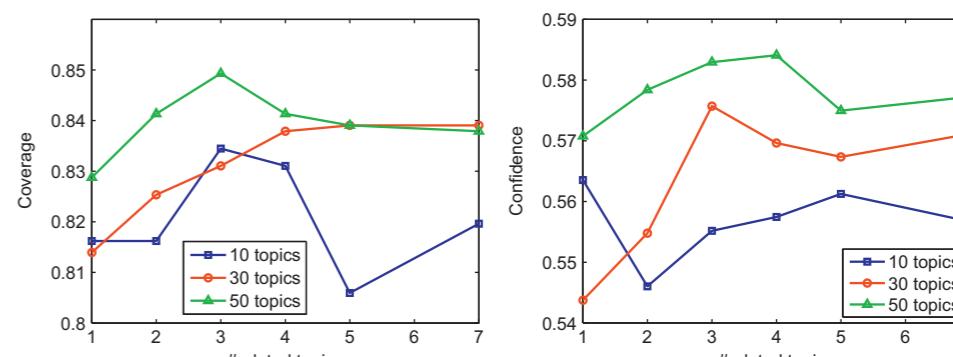


图9. 不同主题模型的覆盖率和相关主题次数的敏感性, $S=0:3$ 。

expertise documents, while the course description and course name are taken as the query.

Baseline methods and evaluation metrics. We employ the same greedy method used in experiment 5.1 as baseline. The real assignment is extracted as the ground truth. Thus, we perform the evaluation in terms of precision.

Results. Fig. 10(a) shows the assignment precision in the course–teacher assignment task by our approach and the baseline method, and (b) shows the effects of the parameter β on the

precision on UIUC data. The precision is defined as the ratio of the number of correct assignments (consistent with the ground truth data) over total number of assignments. As Fig. 10(a) shows, in all the datasets we collect from top universities, our algorithm outperforms the greedy method greatly. And in Fig. 10(b), as the β increases, the precision of our approach increases in general and decreases slowly after it exceeds the peak value. The peak value is more than 50% larger than the initial precision, which validates the effectiveness of the soft penalty approach.

We conduct a further analysis on the UIUC dataset. As Table 6 shows, some professors with publications in various domains, are likely to be assigned with many courses in the baseline algorithm. But in real situation, most professors, though with various background, want to focus on several directions. Thus some courses should be assigned to younger teachers. While in our algorithm, the situation is much better. And we can see that each teacher is assigned with a reasonable load as well as a centralized interest.

5.4. Online system

Based on the proposed method, we have developed an online system for paper–reviewer suggestions, which is available at <http://review.arnetminer.org/>. Fig. 11 shows a screenshot of the system. The input is a list of papers (with titles, abstracts, authors, and organization of each author) and a list of conference program committee (PC) members. We use the academic information stored in ArnetMiner to find the topic distribution for each paper and each PC member [42]. With the two input lists and the topic distribution, the system automatically finds the match between papers and authors. As shown in Fig. 11, there are 5–7 papers assigned to each PC member and the number of reviewers for each paper is set as 3. The system will also avoid the conflict-of-interest (COI) according to the coauthorship and co-organization relationship. In addition, users can provide feedbacks for online adjustment, by removing or confirm (fix) an assignment.

6. Conclusion and future work

In this paper, we study the problem of expertise matching in a constraint-based framework. We formalize the problem as a minimum convex cost flow problem. We theoretically prove that the proposed approach can achieve an optimal solution and develop an efficient algorithm to solve it. Experimental results on two different types of datasets demonstrate that the proposed approach can effectively and efficiently match experts with the queries. Also we provide an algorithm to consider user feedbacks in real time. We are now applying the proposed method to several real-world applications. Feedbacks from the users are very positive.

The general problem of expertise matching represents a new and an interesting research direction. There are many potential

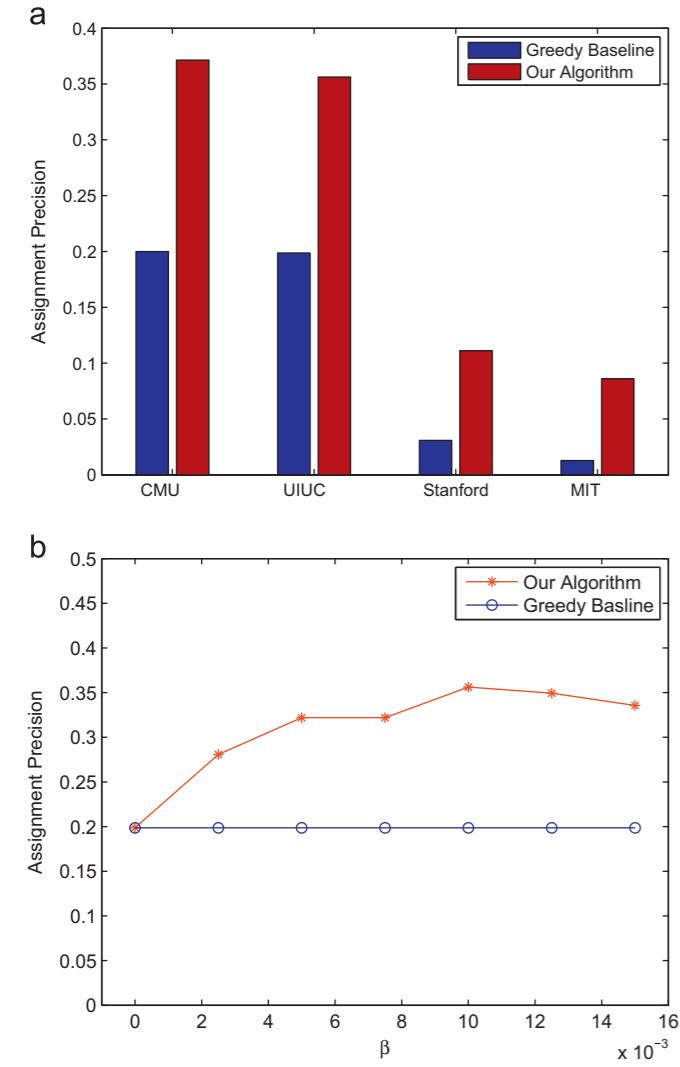


Fig. 10. Course-teacher assignment performance (%). (a) Course assignment results and (b) precision vs. β on UIUC data.

Table 6
Case study: professors with many courses assigned in UIUC (2008, fall – 2010, spring).

Professor	Pub papers	Courses assigned (baseline)	Courses assigned (our approach)
Jose Meseguer	237	23 courses Database systems (2008, spring) Programming languages and compilers (2008, spring) Iterative and multigrid methods (2009, spring) Programming languages and compilers (2009, spring)	7 courses Programming languages and compilers (2008, spring) Programming language semantics (2008, spring) Programming languages and compilers (2008, fall) Programming languages and compilers (2009, spring)
ChengXiang Zhai	117	18 courses Computer vision (2009, spring) Text information systems (2009, spring) Stochastic processes and applic (2009, fall) Computer vision (2008, spring)	7 courses Text information systems (2008, spring) Stochastic processes and applic (2008, fall) Text information systems (2009, spring) Stochastic processes and applic (2009, fall)

专业文件，而课程描述和课程名称被视为查询。基线方法和评估指标。我们使用实验5.1中使用的相同的贪婪方法作为基线。真正的分配是作为地面真理提取的。因此，我们在精度方面进行评估。结果。图10 (a) 通过我们的方法和基线方法示出了课程—教师分配任务中的分配精度，并且 (b) 显示了参数B对的效果

UIUC数据的精度。精度被定义为在总分配总数上的正确分配数量（与地面真实数据一致）的比率。如图10 (a) 所示，我们从顶部大学收集的所有数据集中，我们的算法大大地表现了贪婪的方法。并且在图1中。如图10 (b) 所示，随着B的增加，我们的方法的精度一般增加并且在超过峰值后缓慢降低。峰值比初始精度大超过50%，验证了软罚球方法的有效性。我们对UIUC数据集进行进一步的分析。作为表6。

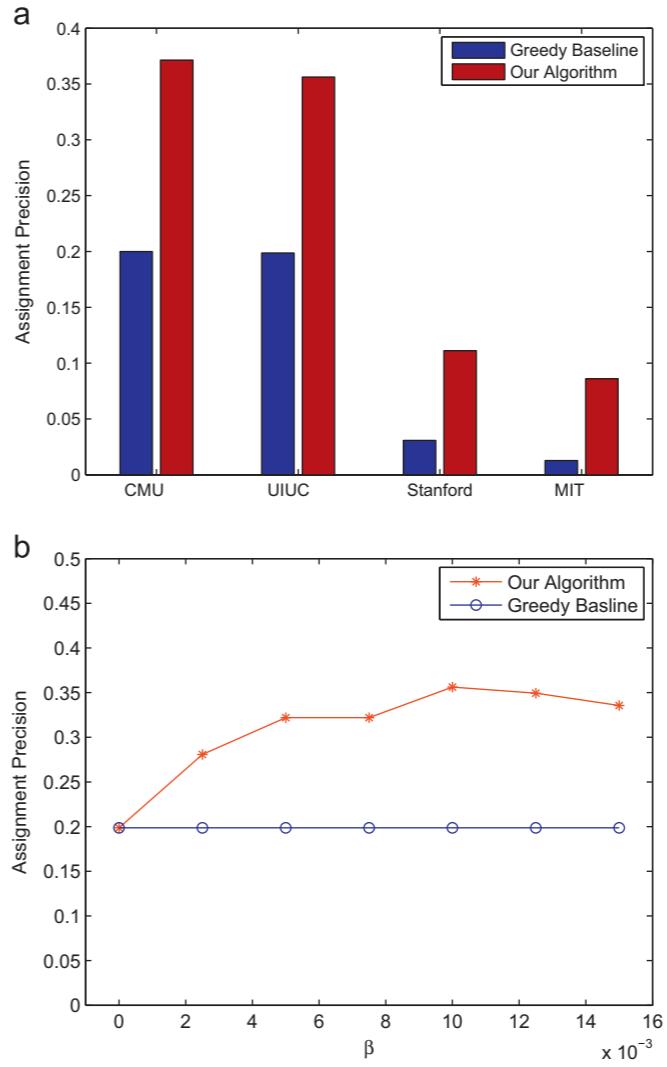


图10.课程教师分配性能 (%). (a) 课程分配结果和 (b) 对UIUC数据的精度与B.

表6案例研究：UIUC中分配了许多课程的教授（2008年，2010年秋季，春季）。

Pub Papers教授分配 (基线) 课程 (我们的方法)	
Jose Meseguer 237	23课程7课程数据库系统 (2008年，春季) 编程语言和编译器 (2008, Spring) 编程语言和编译器 (2008, Spring) 编程语言语言语言语言语言 (2008, Spring) 迭代和多重资源方法 (2009, Spring) 编程语言和编译器 (2008年，秋季) 编程语言和编译器 (2009年，Spring) 编程语言和编译器 (2009年，春季)

成翔翟117课程7课程电脑愿景 (2009年，春季) 文本信息系统 (2008年，春季) 文本信息系统 (2009年，春季) 随机流程和亚申请 (2008, 秋季) 随机流程和亚申请 (2009, 秋季) 文本信息系统 (2009年, Spring) 计算机视觉 (2008, Spring) 随机流程和申请 (2009年, 秋季)

shows, some professors with publications in various domains, are likely to be assigned with many courses in the baseline algorithm. But in real situation, most professors, though with various background, want to focus on several directions. Thus some courses should be assigned to younger teachers. While in our algorithm, the situation is much better. And we can see that each teacher is assigned with a reasonable load as well as a centralized interest.

5.4. Online system

Based on the proposed method, we have developed an online system for paper–reviewer suggestions, which is available at <http://review.arnetminer.org/>. Fig. 11 shows a screenshot of the system. The input is a list of papers (with titles, abstracts, authors, and organization of each author) and a list of conference program committee (PC) members. We use the academic information stored in ArnetMiner to find the topic distribution for each paper and each PC member [42]. With the two input lists and the topic distribution, the system automatically finds the match between papers and authors. As shown in Fig. 11, there are 5–7 papers assigned to each PC member and the number of reviewers for each paper is set as 3. The system will also avoid the conflict-of-interest (COI) according to the coauthorship and co-organization relationship. In addition, users can provide feedbacks for online adjustment, by removing or confirm (fix) an assignment.

6. 结论和未来的工作

在本文中，我们研究了基于约束框架的专业知识问题。我们将问题形式化为最小凸起成本流动问题。理论上，我们证明了所提出的方法可以实现最佳解决方案并开发一种有效的算法来解决它。两种不同类型的数据集上的实验结果表明，所提出的方法可以有效且有效地将专家与查询匹配。我们还提供了一种算法，可以实时考虑用户反馈。我们现在将提议的方法应用于几个现实世界的应用程序。来自用户的反馈非常积极。专业知识匹配的一般问题代表了一种新的和有趣的研究方向。有很多潜力

Assign Result Grouped By Reviewers:

#Reviewer/Paper: 3 #Papers/Reviewer: Lower bound: 5 Upper bound: 7 Beta: 0.0 Reassign ReassignWithFeedbacks

Home | Paper List | Reviewer List | Relevance View | Save | Export(.xls)

Expand All | Collapse All

[+] Ah-Hwee Tan (See Details)

- SIGMA: MPI for Large Scale Machine Learning
- Privacy Preserving Frequency-based Learning Algorithms in 2-Part Fully Distributed Setting
- Learning from simple to complex
- Active exploration for link-based preference learning using Gaussian processes
- The Refinement of Chartist Knowledge for Stock Price Index Forecasting Using Feature Extraction Neural Networks (FENNs)
- Active Learning via Generalized Queries with Minimum Cost

[+] Alexandre V. Evfimievski (See Details)

- Applying Multidimensional Association Rule Mining to Feedback-based Recommendation Systems
- k-Support Anonymity based on Pseudo Taxonomy for Outsourcing Frequent Itemset Mining
- Mining complex periodic behaviors for moving objects
- Malware Detection Based on Objective-Oriented Association Mining
- Fast mining for epistatic interactions
- Versatile Publishing for Privacy Preservation

[+] Alexandros Ntoulas (See Details)

[+] Amol Ghoting (See Details)

Fig. 11. Screenshot of the online system.

future directions of this work. One interesting issue is to apply the proposed framework to question answer (e.g., Yahoo! Answer), where one of the most important issues is how to identify who can answer a new question. Another interesting issue is to incorporate some supervised information into our framework to further improve the performance of expertise matching. Finally, it is important to consider the influence between users when extending expertise matching to the social network.

Acknowledgments

The work is supported by the Natural Science Foundation of China (Nos. 61073073, 60703059, and 60973102), Chinese National Key Foundation Research (No. 60933013), and National High-tech R&D Program (No. 2009AA01Z138). The work was done when the last author was visiting Tsinghua University.

References

- [1] C.B. Haym, H. Hirsh, W.W. Cohen, C. Nevill-manning, Recommending papers by mining the web, in: Proceedings of the 20th International Joint Conference on Artificial Intelligence (IJCAI'99), 1999, pp. 1–11.
- [2] S.T. Dumais, J. Nielsen, Automating the assignment of submitted manuscripts to reviewers, in: SIGIR'92: Proceedings of the 15th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, ACM, New York, NY, USA, 1992, pp. 233–244.
- [3] M. Karimzadehgan, C. Zhai, G. Belford, Multi-aspect expertise matching for review assignment, in: Proceedings of the 17th ACM International Conference on Information and Knowledge Management (CIKM'08), 2008, pp. 1113–1122.
- [4] D. Mimno, A. McCallum, Expertise modeling for matching papers with reviewers, in: Proceedings of the 13th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (SIGKDD'07), 2007, pp. 500–509.
- [5] M. Karimzadehgan, C. Zhai, Constrained multi-aspect expertise matching for committee review assignment, in: Proceedings of the 17th ACM International Conference on Information and Knowledge Management (CIKM'09), 2009, pp. 1697–1700.
- [6] D. Hartvigsen, J.C. Wei, R. Czuchlewski, The conference paper-reviewer assignment problem, Decision Sciences 30 (3) (1999) 865–876.
- [7] Y.-H. Sun, J. Ma, Z.-P. Fan, J. Wang, A hybrid knowledge and model approach for reviewer assignment, in: Proceedings of the 40th Hawaii International Conference on Systems Science (HICSS-40 2007), 2007, p. 47.
- [8] S. Hettich, M.J. Pazzani, Mining for proposal reviewers: lessons learned at the national science foundation, in: Proceeding of the 15th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD'06), 2006, pp. 862–871.
- [9] D. Conry, Y. Koren, N. Ramakrishnan, Recommender systems for the conference paper assignment problem, in: RecSys'09: Proceedings of the Third ACM Conference on Recommender Systems, ACM, New York, NY, USA, 2009, pp. 357–360.
- [10] N.D. Mauro, T.M.A. Basile, S. Ferilli, Grape: an expert review assignment component for scientific conference management systems, in: IEA/AIE'05: Proceedings of the 18th International Conference on Industrial and Engineering Applications of Artificial Intelligence and Expert Systems, 2005, pp. 789–798.
- [11] R. Van De Stadt, Cyberchair: A Web-Based Groupware Application To Facilitate The Paper Reviewing Process, URL <<http://borbala.com/cyberchair/wbgafrp.pdf>>.
- [12] Microsoft Conference Management Toolkit (cmt), URL <<http://cmt.research.microsoft.com/cmt/>>.
- [13] The EasyChair Software, URL <<http://www.easychair.org/>>.
- [14] H. Fang, C. Zhai, Probabilistic models for expert finding, in: Proceedings of the 29th European Conference on Information Retrieval Research ECIR'07, 2007, pp. 418–430.
- [15] D. Petkova, W.B. Croft, Hierarchical language models for expert finding in enterprise corpora, International Journal on Artificial Intelligence Tools (2008) 5–18.
- [16] K. Balog, L. Azzopardi, M. de Rijke, Formal models for expert finding in enterprise corpora, in: Proceedings of the 29th ACM SIGIR International Conference on Information Retrieval (SIGIR'06), 2006, pp. 43–55.
- [17] W. Tang, J. Tang, C. Tan, Expertise matching via cosstraint-based optimization, in: Proceedings of 2010 IEEE/WIC/ACM International Conferences on Web Intelligence (WI'2010), 2010.
- [18] D. Yimam, A. Kobsa, Demoir: a hybrid architecture for expertise modeling and recommender systems, in: Proceedings of the Ninth IEEE International Workshops on Enabling Technologies: Infrastructure for Collaborative Enterprises, 2000, pp. 67–74.
- [19] Y. Cao, J. Liu, S. Bao, H. Li, Research on expert search at enterprise track of trec 2005, in: TREC, 2005.
- [20] Y. Fu, W. Yu, Y. Li, Y. Liu, M. Zhang, S. Ma, Thuir At Trec 2005: Enterprise Track, 2005.
- [21] C. Basu, H. Hirsh, W.W. Cohen, C. Nevill-Manning, Technical paper recommendation: a study in combining multiple information sources, Journal of Artificial Intelligence Research 14 (2001) 231–252.
- [22] C. Basu, H. Hirsh, W. Cohen, Recommendation as classification: using social and content-based information in recommendation, in: Proceedings of the Fifteenth National Conference on Artificial Intelligence, AAAI Press, 1998, pp. 714–720.
- [23] D. Yarowsky, R. Florian, Taking the Load Off The Conference Chairs: Towards A Digital Paper-Routing Assistant, 1999.
- [24] J. Zhang, J. Tang, J. Li, Expert finding in a social network, Advances in Databases: Concepts Systems and Applications 23 (2010) 1066–1069.
- [25] D. Yimam-Seid, A. Kobsa, Expert finding systems for organizations: problem and domain analysis and the demoir approach, Sharing Expertise: Beyond Knowledge Management 23 (2003) 327–358.
- [26] C. Zhai, J. Lafferty, A study of smoothing methods for language models applied to ad hoc information retrieval, in: Proceedings of the 24th ACM SIGIR International Conference on Information Retrieval (SIGIR'01), 2001, pp. 334–342.

Assign Result Grouped By Reviewers:

#Reviewer/Paper: 3 #Papers/Reviewer: Lower bound: 5 Upper bound: 7 Beta: 0.0 Reassign ReassignWithFeedbacks

Home | Paper List | Reviewer List | Relevance View | Save | Export(.xls)

Expand All | Collapse All

[+] Ah-Hwee Tan (See Details)

- SIGMA: MPI for Large Scale Machine Learning
- Privacy Preserving Frequency-based Learning Algorithms in 2-Part Fully Distributed Setting
- Learning from simple to complex
- Active exploration for link-based preference learning using Gaussian processes
- The Refinement of Chartist Knowledge for Stock Price Index Forecasting Using Feature Extraction Neural Networks (FENNs)
- Active Learning via Generalized Queries with Minimum Cost

[+] Alexandre V. Evfimievski (See Details)

- Applying Multidimensional Association Rule Mining to Feedback-based Recommendation Systems
- k-Support Anonymity based on Pseudo Taxonomy for Outsourcing Frequent Itemset Mining
- Mining complex periodic behaviors for moving objects
- Malware Detection Based on Objective-Oriented Association Mining
- Fast mining for epistatic interactions
- Versatile Publishing for Privacy Preservation

[+] Alexandros Ntoulas (See Details)

[+] Amol Ghoting (See Details)

图11. 在线系统的屏幕截图。

这项工作的未来方向。一个有趣的问题是将拟议的框架应用于质疑答案（例如，雅虎！答案），其中最重要的问题之一是如何识别谁可以回答一个新问题。另一个有趣的问题是将一些受监管信息纳入我们的框架，以进一步提高专业匹配的绩效。最后，重要的是考虑用户在与社交网络匹配的专业知识时的影响。

Acknowledgments

该工作得到了中国自然科学基金（No.61073073,60703059和60973102）的支持，中国国家重点基础研究（60933013号）和国家高科技研发计划（No.2009AA01Z138）。最后一位作者访问清华大学时，这项工作是完成的。

References

- [1] C.B. Haym, H. Hirsh, W.W. Cohen, C. Nevill-Manning, Through开采网络推荐论文, 在: 第20届国际人工智能联席会议 (IJCAI'99), 1999, 第1–11页的第20届国际联席会议。[2] S.T. Dumais, J.Nielsen, 自动向审稿人员分配提交的稿件, 在: 第15届年度国际ACM SIGIR会议的关于信息检索, ACM, 纽约, 纽约, 美国, 1992, PP.的第15届国际ACM SIGIR会议, 233–244。[3] M. Karimzadehgan, C. Zhai, G. Belford, 多个方面专业匹配审查任务, 为: 第17届ACM信息和知识管理国际会议 (CIKM'08), 2008, PP. 1113–1122。[4] D. Mimmo, A. McCallum, 与审稿人员匹配论文的专业知识建模, 包括: 第13届ACM SIGKDD国际知识发现和数据挖掘国际会议 (SIGKDD'07), 2007, PP. 500–509。[5] M. Karimzadehgan, C. Zhai, 委员会审查任务的受限多方面专业匹配, 适用于: 第17届ACM信息和知识管理国际会议 (CIKM'09), 2009, PP. 1697–1700。[6] D. Hartvigsen, J.C. Wei, R. Czuchlewski, 会议论文 – 审阅者分配问题, 决策科学版30 (3) (1999) 865–876。[7] Y.-H. Sun, J. Ma, Z.-P. Fan, J. Wang, A hybrid knowledge and model approach for reviewer assignment, Decision Sciences 30 (3) (1999) 865–876。[8] S. Hettich, M.J. Pazzani, 提案审查员挖掘; 在国家科学基金会的经验教训, 在: 第15届ACM SIGKDD国际知识发现和数据挖掘会议上 (KDD'06), 2006年, pp.862–871。
- [9] D. Conry, Y. Koren, N. Ramakrishnan, Constence纸张分配问题的推荐系统, IN: Recsys'09: 第三个ACM会议的会议员工, ACM会议推荐者系统, ACM, 纽约, NY, 美国, 2009年, 第357–360页。[10] N.D.Mauro, T.M.A.巴斯蒂尔, S.Ferilli, 葡萄: 科学会议管理系统的专家审查任务组成部分, IEA / AIE'05: 第18届工业和工程应用系统的贸易委员会的人工智能和专家系统, 2005年, PP.的诉讼程序。789–798。[11] R. Van de Stadt, CyberChair: 基于Web的群件应用程序, 用于遵守纸质评论过程, URL / [HTTP://borbala.com/cyberchair/wbgafrp.pdf](http://borbala.com/cyberchair/wbgafrp.pdf)。[12] Microsoft会议管理工具包 (CMT), URL / <http://cmt.Research.microsoft.com/cmt/>。[13] EasyChair软件, URL / <http://www.easychair.org/>。[14] H. Fang, C. Zhai, 专家查找的概率模型, 在: 第29届欧洲信息检索研究会议核发会议核发组织07,2007, PP. 418–430。[15] D. Petkova, W.B. Croft, Enterprise Corpora中专家查找的分层语言模型, 国际人工智能工具 (2008) 5–18。[16] K. Balog, L. Azzopardi, M. de Rijke, 企业集团专家的正式模型, 包括: 第29届ACM SIGIR国际信息检索国际会议 (Sigir'2006), 2006, PP. 43–55。[17] W. Tang, J. Tang, C. Zhai, 通过基于COSNTRAINT的优化匹配的专业知识, IN: 2010年IEE E / WIC / ACM国际会议关于WEB Intelligence (WI'2010), 2010年的课程。[18] D. YIMAM, A. Kobsa, Demoir: 用专业知识建模和推荐系统的混合架构, 在: 第九届IEE EE国际研讨会上的培训技术的程序: 叠层进入的基础设施, 2000, PP. 67–74。[19] Y. Cao, J. Liu, S. Bao, H. Li, TREC 2005企业轨道专家搜索研究, IN: TREC, 2005。[20] Y. Fu, W. Yu, Y. Li, Y. Zhang, S. Ma, Thuir在TREC 2005: Enterprise Track, 2005。[21] C. Basu, H. Hirsh, WW科恩, C.尼维曼宁, 技术论文建议: 结合多种信息来源的研究, 人工智能研究杂志14 (2001) 231–252。[22] C. Basu, H. Hirsh, W. Cohen, 分类建议: 使用基于社会和基于内容的信息的建议书, 在: 第五十届全国人工智能大会上的诉讼程序, Aaai Press, 1998, PP. 714–720。[23] D. Yarowsky, R. Florian, 将负担从会议椅上取消: 走向数字纸张路由助理, 1999年。[24] J. Zhang, J. Tang, J. Li, 社交网络中的专家查找, 数据库, 23 (2010) 1066–1069。[25] D. Yimam-Seid, A. Kobsa, 组织的专家查找系统: 问题和域分析以及Demoir方法, 共享专业知识: 越超知识管理23 (2003) 327–358。[26] C. Zhai, J. Lafferty, 用于适用于临时信息检索的语言模型的平滑方法的研究, IN: 24th ACM SIGIR国际信息检索国际会议 (Sigir'01), 2001, PP. 334–342。

- [27] K. Balog, L. Azzopardi, M. de Rijke, A Language Modeling Framework For Expert Finding, 2008.
- [28] X. Wei, W.B. Croft, Lda-based document models for ad-hoc retrieval, in: Proceedings of the 29th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval SIGIR'06, ACM, New York, NY, USA, 2006, pp. 178–185.
- [29] J.J.M. Guervós, P.A.C. Valdviés, Conference paper assignment using a combined greedy/evolutionary algorithm, in: PPSN, 2004, pp. 602–611.
- [30] E. Zitzler, L. Thiele, Multiobjective Evolutionary Algorithms: A Comparative Case Study and the Strength Pareto Approach, 1999.
- [31] M. Cámaras, J. Ortega, F. de Toro, A single front genetic algorithm for parallel multi-objective optimization in dynamic environments, Neurocomputing 72 (2009) 3570–3579.
- [32] C.J. Taylor, On the Optimal Assignment of Conference Papers to Reviewers, Technical Report, MS-CIS-08-30, Computer and Information Science Department, University of Pennsylvania, 2008.
- [33] M.A. Rodriguez, J. Bollen, An algorithm to determine peer-reviewers, in: Proceedings of the Conference on Information and Knowledge Management, ACM Press, Napa, California, 2008, pp. 319–328 doi: 10.1145/1458082.1458127.
- [34] S. Benferhat, J. Lang, Conference paper assignment, International Journal of Computational Intelligence Systems 16 (10) (2001) 1183–1192.
- [35] T. Hofmann, Probabilistic latent semantic indexing, in: Proceedings of the 31st Annual International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR'99), 1999, pp. 50–57.
- [36] D.M. Blei, A.Y. Ng, M.I. Jordan, Latent dirichlet allocation, Journal of Machine Learning Research 3 (2003) 993–1022.
- [37] T.L. Griffiths, M. Steyvers, Finding scientific topics, Proceedings of the National Academy of Sciences PNAS'04 (2004) 5228–5235.
- [38] J. Cao, T. Xia, J. Li, Y. Zhang, S. Tang, A density-based method for adaptive lda model selection, Neurocomputing 72 (2009) 1775–1781.
- [39] J. Tang, R. Jin, J. Zhang, A topic modeling approach and its integration into the random walk framework for academic search, in: Proceedings of 2008 IEEE International Conference on Data Mining (ICDM'08), 2008, pp. 1055–1060.
- [40] R.K. Ahuja, T.L. Magnanti, J.B. Orlin, Network Flows: Theory, Algorithms, and Applications, Prentice Hall, 1993.
- [41] P. Beraldi, F. Guerriero, R. Musmanno, Parallel algorithms for solving the convex minimum cost flow problem, Computational Optimization and Applications 18 (2) (2001) 175–190.
- [42] J. Tang, J. Zhang, L. Yao, J. Li, L. Zhang, Z. Su, Arnetminer: Extraction and mining of academic social networks, in: Proceedings of the 14th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (SIGKDD'08), 2008, pp. 990–998.
- [43] The Lemur Toolkit For Language Modeling And Information Retrieval, URL <<http://lemurproject.org/>>.



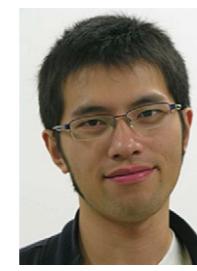
Wenbin Tang is a master student in Tsinghua University, supervised by Prof. Jie Tang. His research interests are text mining, social network and computer vision.



Jie Tang is an associate professor in Tsinghua University. His research interests are machine learning and text mining.



Tao Lei is currently a research assistant in KEG Group, Tsinghua University. He obtained a BS degree in Peking University. His research interests focus on machine learning and text mining.



Chenhao Tan is a Ph.D. student in Cornell University. His research interests are machine learning and social network.



Bo Gao is a software engineer in KEG Group, Tsinghua University. He is currently in charge of the development and maintenance of the academic social network ArnetMiner.



Tian Li is an undergraduate student in Beijing University of Aeronautics and Astronautics.

- [27] K. Balog, L. Azzopardi, M. de Rijke, 专家查找的语言建模框架, 2008年。[28] X. Wei, W.B.克罗夫特, 基于LDA的文档模型, 适用于Ad-hoc检索, In: 第29届国际ACM Sigir会议的关于信息检索和开发的ACM Sigir'06, ACM, 纽约, 纽约, 美国, 2006, PP. 178–185。[29] J.J.M. guervo, 使用组合贪婪/进化算法的会议纸张分配, IN: PPSN, 2004, PP. 602–611。[30] E. Zitzler, L. Thiele, 多目标进化算法: 一种比较案例研究和强度Pareto方法, 1999。[31] M. Ca'Mara, J. Ortega, F. de Toro, 单一前传递算法用于动态环境中的并行多目标优化, 神经会计机72 (2009) 3570–3579。[32] C.J. Taylor, 在审查员, 技术报告, MS-CIS-08-30, 密歇根大学, 2008年培养师, 技术报告, MS-CIS-08-30, 计算机和信息科学事件的最佳分配, 一种确定同行评审员的算法, IN: 信息和知识管理会议, ACM Press, Napa, California, 2008, PP. 319–328 Doi: 10.1145 / 1458082.1458127。[34] BENFERHAT, J. LANG, 会议论文任务, 国际计算智能系统16 (10) (2001) 1183–1192。[35] T. Hofmann, 概率潜在语义索引, 在: 31年度国际ACM Sigir会议上的研究和开发信息检索 (Sigir'99), 1999, PP. 50–57。[36] D.M. Blei, A.Y. Ng, m.i.约旦, 潜在的Dirichlet分配, 机器学习研究学报3 (2003) 993–1022。[37] T.L. Griffiths, M. Steyvers, 寻找科学主题, 国家科学院的诉讼程序PNAS'04 (2004) 5228–5235。[38] J. Cao, T. Xia, J. Li, Y. Zhang, S. Tang, 一种基于密度的自适应LDA模型选择方法, 神经计算器72 (2009) 1775–1781。[39] J. Tang, R. Jin, J. Zhang, 一个主题建模方法及其进入学术搜索随机步行框架的融合, 共度: 2008年I EEE数据挖掘国际会议 (ICDM'08), 2008年, PP. 1055–1060。[40] R.K. Ahuja, T.L. Magnanti, J.B. Orlin, 网络流量: 理论, 算法和应用, Prentice Hall, 1993。[41] P. Beraldi, F. Guerriero, R. Musmanno, 用于解决凸起最小成本流量问题的并行算法, 计算优化和应用 – 阳离子18 (2) (2001) 175–190。[42] J. Tang, J. Zhang, L. Yao, J. Li, L. Zhang, Z. Su, Arnetminer: 学术社交网络的提取和采矿, 在: 第14届ACM SIGKDD国际知识发现国际会议上和数据挖掘 (SIGKDD'08), 2008, PP. 990–998。[43] 语言建模和信息检索的狐猴工具包, URL <http://lemurproject.org/>。



Tao Lei目前是清华大学柯格集团的一名研究助理。他在北京大学获得了BS学位。他的研究兴趣专注于机器学习和文本挖掘。



Chenhao Tan是一个博士学位。康奈尔大学的学生。他的研究兴趣是机器学习和社交网络。



BO高是清华大学柯格集团的软件工程师。他目前负责学术社交网络ARNetminer的发展和维护。



温斌唐是清华大学的一名硕士学位, 由杰唐教授监督。他的研究兴趣是文本挖掘, 社交网络和计算机愿景。



杰唐是清华大学的副教授。他的研究兴趣是机器学习和文本挖掘。



天丽是北京航空航天大学的本科生。