

以各种制约因素优化匹配

Wenbin Tang a, , Jie Tang a, Tao Lei a, Chenhao Tan b, Bo Gao a, Tian Li c

清华大学计算机科学与技术系, 北京10084, 中国康奈尔大学计算机科学系, 北京
北京北京市电科学科技科学科学科学科学科学科学科学科学科学科学

ARTICLE INFO

Available online 31 July 2011

关键词: 专业知识匹
配约束优化纸张 –
审阅者分配主题模型

ABSTRACT

本文研究了与各种限制的专业知识问题。专业匹配, 旨在找到专家和查询之间的对齐, 是许多应用程序中的常见问题, 如会议纸张 – 审阅者分配, 产品 – 评论员对齐和工具 – 内匹配。大多数现有方法将此问题正式化为信息 – 检索问题, 并专注于独立查找每个查询的一组专家。然而, 在现实世界的系统中, 通常需要考虑各种约束。例如, 为了审查一篇论文, 希望至少有一个高级审查员指导审查过程。一个重要问题是: “我们可以设计一个框架, 以便有效地找到各种约束下的专业知识匹配的最佳解决方案吗?” 本文通过在基于约束的优化框架中制定专业知识匹配问题来探讨这种方法。在拟议的框架中, 专业知识匹配问题与凸起成本流量问题相关联, 这保证了在各种约束下的最佳解决方案。我们还介绍了一个在线匹配算法, 以支持实时结合用户反馈。拟议的方法已经在两种不同类型的专业知识匹配问题上进行了评估, 即会议纸张 – 审阅者分配和教师课程任务。实验结果验证了拟议方法的有效性。根据提出的方法, 我们还开发了一个纸张审阅者建议的在线系统, 这些建议已被用于纸张 – 审阅者在大会上的审阅者分配, 从会议组织者的反馈是非常积极的。© 2011 Elsevier B.V. 保留所有权利。

1. Introduction

电脑技术和人类集体智能融合最近被出现为用户在互联网上查找和分享信息的流行方式。例如, Chacha.com 是最大的移动搜索引擎之一, 已经吸引了用户以回答超过3亿个问题; Epinions.com, 消费者审查网站, 已经收集了数千千万的产品评论。以其独特的人类智能的独特使用, 基于人的计算提供了新的方向;但是, 它也会带来一些全新的挑战。称为专业知识匹配的一个关键问题是如何用问题 (查询) 对齐人类专家。直截了当, 我们希望被分配回答问题的人体专家具有与问题相关的具体专业知识。但它显然是不够的。理想的匹配系统还应该考虑现实世界中的各种限制, 例如, 专家只能回答一定数量的问题 (负载余额); 作为不同的权威程度

专家可能很大程度上有所不同, 希望每个问题都可以通过至少一个高级专家 (权威BAL) 回答/审查;一个问题可能与多个不同方面 (主题) 有关, 因此预计所有指定专家的合并专业知识可能会涵盖问题的所有方面 (主题覆盖)。问题引起了不同域名的相当兴趣。例如, 已经使用诸如挖掘Web [1], 潜在语义索引 [2], 概率主题建模 [3,4], 整数线性编程 [5] 来对诸如挖掘纸张审阅者分配进行若干工程, 最小成本流动 [6] 和域知识和匹配模型的混合方法 [7]。还制定了一些诸如 [8–13] 的系统, 以帮助提案 – 审阅者和纸张审阅者分配。然而, 大多数现有方法主要关注提高测量查询和专家之间相关性的准确性, 即如何查找每个查询的 (或等级) 相关专家, 但忽略不同的约束或使用启发式解决约束。此外, 这些方法通常不考虑用户反馈。另一方面, 有一些方法专注于专家发现。例如, 方等人。 [14] 提出了一个专家查找的概率模型, 佩特科维等人。 [15] 雇用

* Corresponding author.

E-mail address: tangwb06@gmail.com (W. Tang).

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？

图。图1示出了纸张审阅者匹配问题的示例（将审阅者分配给每份纸张）。在问题中，对应于审阅者（即，审阅者的专业知识）对应的主题可以是“机器学习”，“数据挖掘”，“计算-ORY”等。也，每篇论文也是如此在不同的主题上发行。有一些要求是一个好的任务。首先，对于每篇论文，指定的审稿人的专业知识应该涵盖论文的主题，所有审稿人都应该有负载余额（每个审阅者只能审查一定数量的论文）。此外，一些审稿人可能是高级，有些人可能是平均的。我们始终希望每份纸张的审查过程可以由至少一个高级评价者“监督”。另一个例子是患者-医生匹配案例，主题追加到医生的专业知识包括“儿科”，“Rheuma-Tological”，“神经病学”等。每位医生都有不同的专业知识学位主题，而患者的疾病也有关于主题的相关分布。理想情况下，在安排患者的咨询时，指定医生的主题应包含患者疾病的潜在原因（例如，SLE患者的咨询需要风湿病学家，肾病学家，心理学家和神经病学家），以及所有医生应该有一个负载余额，以便没有医生过剩。贡献。在本文中，我们正式定义了专业知识匹配问题并提出了基于约束的优化框架来解决问题。具体地，专业知识匹配问题作为凸起成本流问题，目标是在某些约束下找到具有最小成本的可行流量。理论上我们证明，所提出的框架可以在各种结构下实现最佳解决方案，并开发一种有效的算法来解决它。本文是我们之前的会议论文的延伸和改进[17]，与以前的工作有所不同。

（1）我们考虑匹配的“多主题覆盖”匹配来重新正式化我们的框架，这对纸张审阅者分配问题非常重要。我们展示了主题覆盖措施（例如，覆盖范围和信心）也可以纳入我们的框架。

（2）引入了一个额外的数据集以评估我们对多主题覆盖的方法的性能。实验结果证实了效果和效率

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头脑会议的论文。会议组织者的反馈确认了拟议方法的有效性。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性分配纸张。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MCCALLUM [4]通过提出一个新颖的主题模型作者-

Persons-pomplea-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

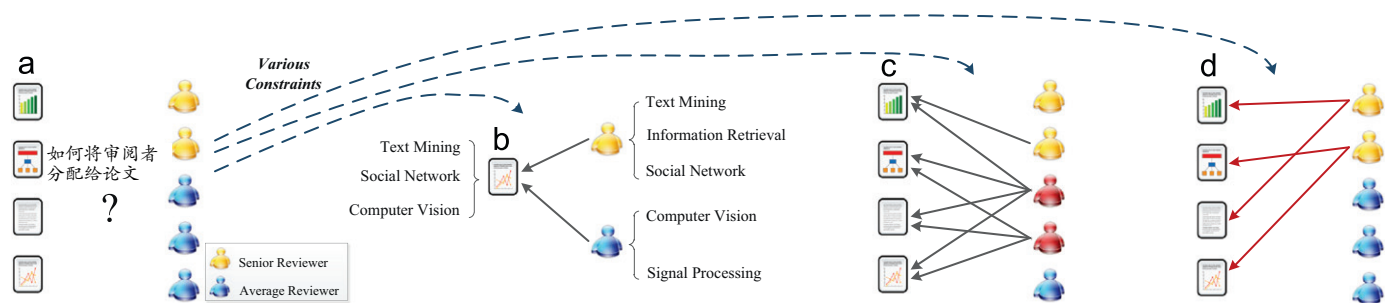


图1. (a) 纸张审阅者分配问题的图示。为了做出理想的作业，应考虑一些约束。例如：(b) 每篇论文与一个或多个主题有关，每个审稿人也具有不同主题的专业知识。对于一篇论文，所有指定审核人员的专业知识的组合应涵盖本文的所有主题。(c) 作业应该是“负载均衡”。显示了一个糟糕的例子，其中标有红色的两个评论员被分配了太多的作品，同时蓝色没有得到任何分配。(d) 审阅者的权威程度可能很大程度上，也有意义上，每篇论文至少有一个高级评论员分配，以便她/他可以指导审查流程。

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，Guervo/S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
t	题目
v	候选专家
q	集查询
vi	一个专家
qj	一个查询
yviz	主题
z	给定专题
vi	yqjz给定查询
qj	t的概率
z	给定查询
qj	t (vi) 关于专家
VI	T (Qj) 的主要相关主题集的一组主
T (Qj)	要相关主题
Qj	

还制定了[8–10,33,34,11–13]，以帮助提案 – 审阅者和纸张审阅者任务。但是，匹配问题仍然被视为信息检索问题，这显然无法导致最佳解决方案。在本文中，我们的目标是在基于约束的优化框架中模拟专业知识匹配问题，并提出了一种努力解决框架的算法。我们对现有工作的差异是：(a) 我们提供优化框架，将专业知识匹配和各种限制集合在一起；(b) 由于简单地定义新的（硬或软）约束，可以轻松扩展框架，因为新约束可以将新约束组合到优化框架中；(c) 框架可以保证最佳解决方案。

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 $V1/4[vi]$ ，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 T 方面的专业知识（称为主题），每个专家 vi 在不同主题上具有不同的专业知识学位。此外，鉴于第 $1/4$ 季度/ qj 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 FP $w9y$ 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）相关联的单词根据关于每个主题的词分布，即 P $w9y$ 。因此，分布中具有最高概率的单词将建议主题所代表的神学。

假设我们有 T 主题，主题 $ZAF1$ TG 专业 vi 的专业知识学位 $zaf1$ 示为具有 p $zyviz$ 1 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 Qj 都可以表示为单词序列，即 dqj 。要代表每个专家 VI ，没有普遍性，我们也将视为一系列单词，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。

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 $[n1, n2]$ ，其中 $n1$ $m2$ ；（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 $N1$ 和 $N2$ 之间的范围内。在现有的作品中，Du mais等人。[2]和mimno等人。[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。

4.受限制的优化框架

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。

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

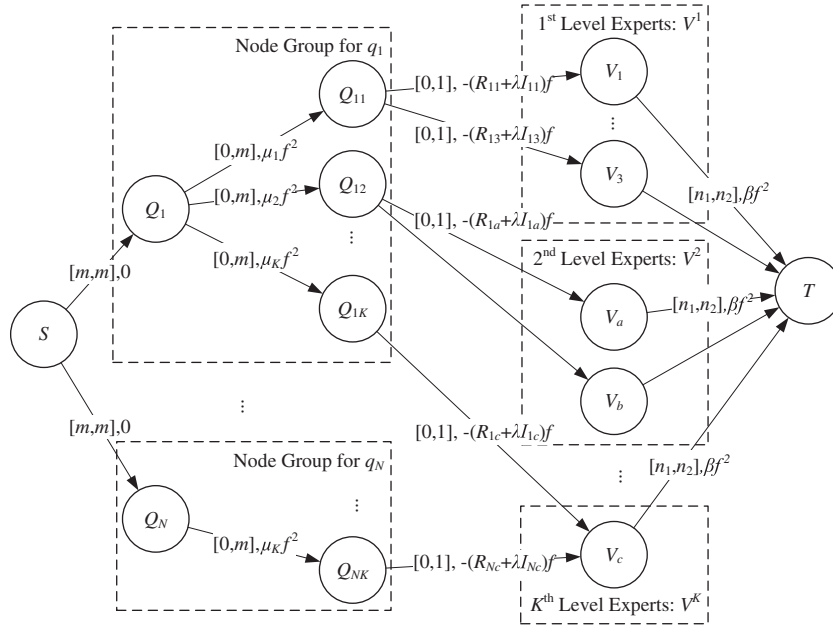


图2.根据客观函数 (10) 的凸起成本网络流的构造。

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 R_{IJ} , 以及如何优化框架。

4.2. Modeling multiple topics

主题建模的目标是将每个专家 V_i 与 T 维主题分发的矢量 Y_{VI} 艺术相关联, 也将每个查询 Q_j 与矢量 Y_{QJ} 艺术相关联。该主题分布可以以多种不同的方式获得。例如, 在纸张审阅者分配问题中, 每个审阅者都可以从预定义的类别中选择他们的专业知识。此外, 我们可以使用统计主题建模[35,36]来自动从输入数据中提取主题。在本文中, 我们使用主题建模方法初始化每个专家和每个查询的主题分布。要提取主题分发, 我们可以考虑我们有一组 M 个专家文档和 N 查询文档 (每个代表专家或查询)。通过累积与专家有关的内容信息, 可以获得专家的文件。例如, 我们将所有出版物文件与审阅者的专家文件相结合, 因此专家 V_i 的文件可以代表 D_{VI} FWIJG。每个查询也可以被视为文档。然后, 我们可以使用诸如LDA

[36]的主题模型来了解这些 T 主题方面。具体而言, 让 $D = D_{V1}, \dots, D_{VM}$ 是一组专家文件。根据LDA的整个系列的日志可能性是

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

其中 $C(w, d)$ 是文档 D 中的单词 W 的计数, $p(w|z, \phi_z)$ 是主题 z 成字 w 的概率, 并且 $p(z|d, \theta_d)$ 是包含主题 z 的文档 d 的概率。我们使用GIBBS采样算法[37,38]以了解每个专家的主题分发 Y_{VI} 。查询 Y_{QJ} 的主题分布可以从生产的 Y_{VI} 推断出来。

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

V_i 和查询 Q_j 之间的匹配分数 R_{IJ} 被解释为概率 R_{LM}

$IJ \quad P \quad qj9di \quad qwaqjp \quad w9di$

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

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, N 是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

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

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. 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价的最低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。

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 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 qj aq做1: 3 1: 4 1: 5

创建 $k-1$ 节点，表示为 $qj, qj1, \dots, QJK$ 分别；从源节点 S 向节点 Qj 添加一个弧，零成本和流量约束 m, m ；从节点 Qj 到 QJK 的电弧，方形成本函数 MKF 2和流量约束 $0, m$ ；

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 2和流量约束 $n1, n2$ ；

1.9 FOREACH VI AV, QJ 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 $liij$ 和流量约束 $0, 1$ ；

1.12计算 g 的最小成本流量；1.13 FOREACH VI AV, QJ

AQ, $S: T: UIJ$ 1做1:14 1:15

K_{vi} 的专家水平；如果流 f_{qjk} ，则 vi 1然后将查询 qj 分配给专

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：

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

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 表示 $9t_{qj} \setminus t_{vi} 9=9t_{qj} 9$ 。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 QJ 的流量表示查询 Qj 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 之间的成本与负载均衡软罚函数（Eq.（4））相对应。从 QJ 到 QJK 的流量的含义是分配给 QJ 的 k th级专家的数量，因此我们在弧上强加一个平方成本函数 $mk(f)$ 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 QJ 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

MCCF问题的功能（EQ.（15））是（EQ.（10））的负形式。在实践中，没有必要添加所有（ QJK, VI ）弧。为了进一步降低算法的复杂性，我们首先贪婪地生成分配并保留相应的弧，然后只能保持 QJK 的 $C-M$ 弧和用于 VI 的弧度，具有最高匹配得分（ C 是固定常数）。我们称此过程还原，这将减少网络中的弧数而不会影响性能太大。为了处理大规模数据，我们可以利用凸起成本流的并行实现[41]。

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

算法2考虑第一种类型的反馈，它仍然产生最佳解决方案。

算法2.在线匹配算法。

输入： G 上的最小成本网络流量 F 对应于当前分配；不当匹配 (vi, qj) 。输出：新分配。2.1k $pertvi$ 水平；2.2如果 f_{qjk}, vi 1那么2: 3 2: 4

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, qjk ；CANCELINGRALLBACK 1沿迎战的流量单位，更新 $G-F$ ；从 g 中删除 arc_{qjk}, 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；

2.5
2.6
2.7
2.8

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（ QJ, 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, QJ) ，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

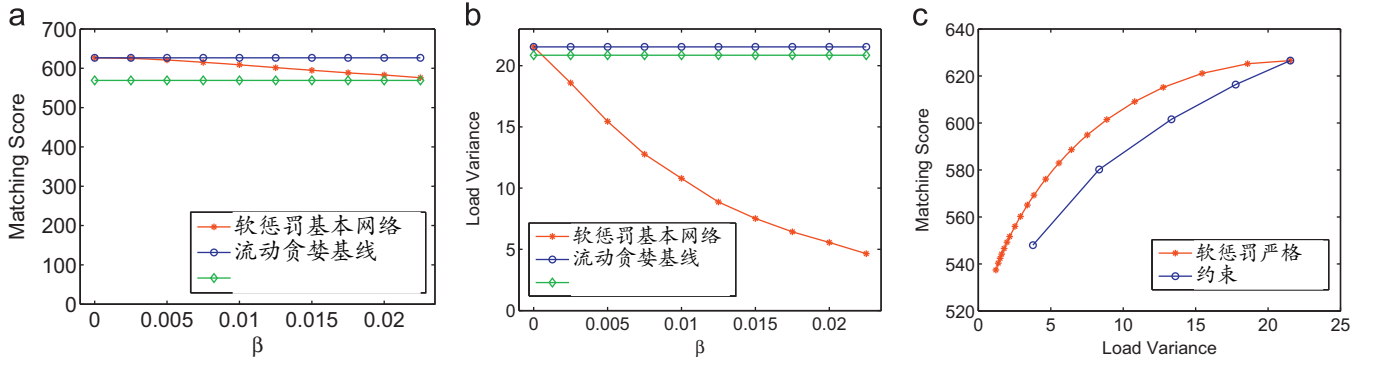


图3. (a) 和 (b) 示出了软惩罚功能如何影响匹配分数 (MS) 和与不同B的负载方差。 (c) 在软惩罚函数和严格的负载平衡方面进行了比较。

成本。因此，算法2在增强新分配后给出最佳解决方案。 &

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。

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 Netminer3 [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N2 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， $9q \quad vi \quad 9rn2$ ），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。

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

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

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 $m \in 0$ 并调整参数 b 以找出软惩罚函数的效果。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，即严格的限制和软惩罚性。两个LV-MS曲线分别通过设置不同的最小数量 $N1$ 来产生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 $M2 \quad 0$ ，以考虑仅考虑高级审阅者的余额。图4显示了 $M1$ 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 $E_{q_i} \quad (1)$ ），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荷。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。如图6所示，弧度减少过程可以显著降低时间消耗。

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

³ <http://arnetminer.org>

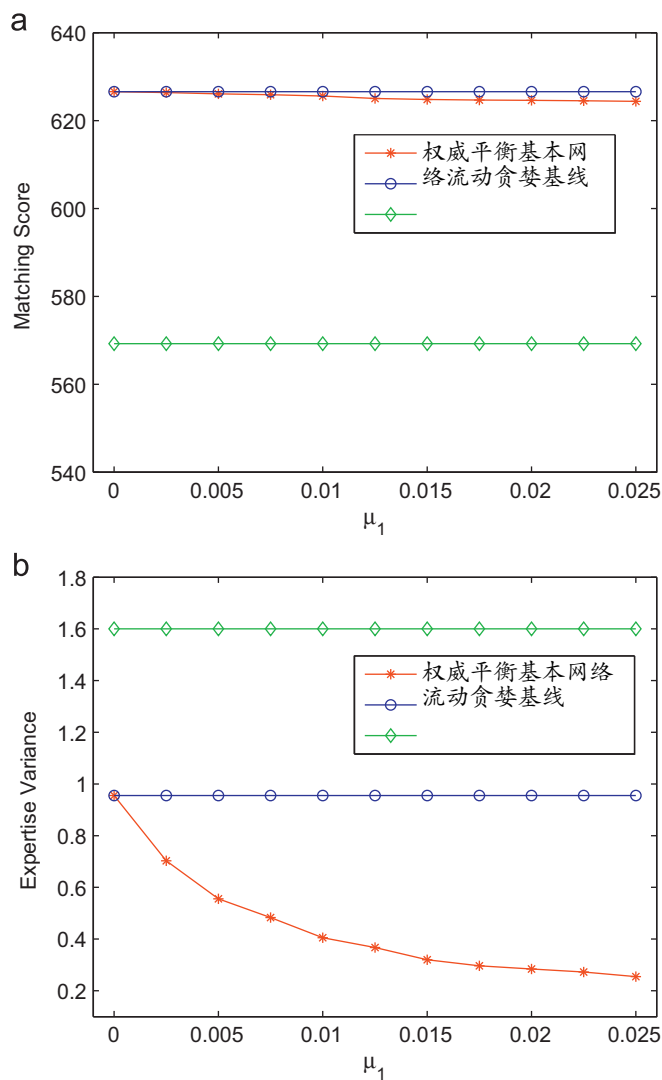


图4.匹配分数（MS）和专业知识方差（EV）各种不同。

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	
0:02	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，被分配有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]，其中4个由73个查询和189名审阅者组成。73个查询是Sigr07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

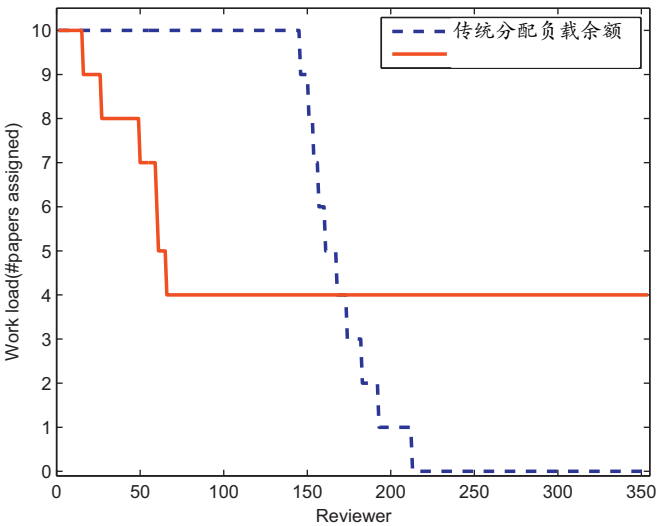


Fig. 5. The work load (number of paper assigned) of every reviewer.

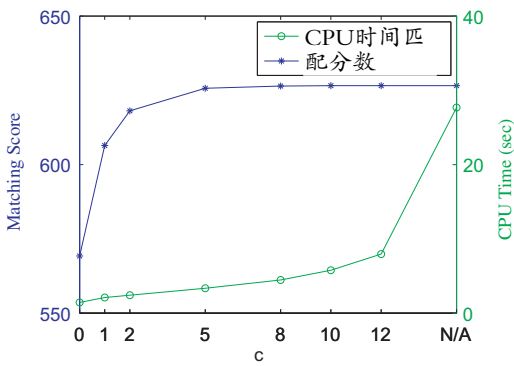


Fig. 6. Efficiency performance (s).

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METAFACT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：超越幂律和逻辑正式发行版在线品牌广告的观众选择：隐私友好

摘要由相应的审稿人在Sigr1971–2006发表。在这个数据集中，25个主题，如“文本挖掘”，“群集”和“语言模型”的人员专家，以及每个查询的相关主题和每个评论员都被手动标记为金标准进行评估。我们使用Lemur Toolkit [43]进行预处理以授权每个查询和文档，删除常见的停止单词。设置。由于数据集中的作者和审核人员是匿名的，因此我们在本实验中不考虑COI和权威余额。为了与[3]中的设置一致，我们设定了由M=3专家审查的纸张，并使用覆盖范围（EQ。（6）），信心（EQ。（7）），AverageConfidence（EQ（8））和FScore作为措施。在所有方法中，分别选择具有最大Y值的顶主题作为每个查询和专家的相关主题。比较方法。我们在[3]中实施三种方法作为比较方法。第一个基线方法使用语言模型来检索每个查询的审阅者文档。具体地，给定查询QJ，排名审阅者的一种方式是使用查询QJ给定审阅者文档 d_i 的概率，即，

$$p(q_j|d_{v_i}) = \prod_{w \in q_j} \frac{tf(w, d_{v_i}) + \mu p(w|D)}{N_{d_{v_i}} + \mu} \quad (16)$$

如果NDVI是DVI的文档大小，TF_W是DVI中WORD W的发生时间的数量，P_{w9d}是整个评论者文档收集的UNIGRAM语言模型，M是一个虚拟设置到1000的Dirichlet平滑因子。在另外，我们使用P_{w9d}的平滑版本，所以当Word W在整个文档中看不见时，它不会返回零

表4五个随机纸的审稿人列表。

Paper	Assigned reviewers
C. Lee Giles, Jie Tang, Matthew Richardson, Hady Wiawan Lauw, Elena Zheleva	垃圾邮件过滤后勤回归Rong Jin, Chengxiang Zhai, Saharon Rosset, Masashi Sugiyama, Annalisa Appice
Surtucture学习非平滑排名亏损xian-sheng hua, 刘铁燕刘, 云博曹, 洛伦加塞纳无监督重复数据删除采用跨场依赖性诚信霍, 深斋瓦尔·瓦尔, Max康策, 唐纳德·梅兹勒, 奥伦库兰兰社交网络中的信息途径结构。李吉尔德, 沃尔夫冈Nejdl, Melanie Gnasa, Michalis Faloutsos, Cameron Marlow	

collection, i.e.,

$$p(w|D) = \frac{tf(w, D) + 1}{N_D + V} \quad (17)$$

其中V是估计的词汇大小（即，不同词语的数量）。第二个基线方法使用KL分歧进行审阅者检索[26]：

$$\sum_{w \in q_j, tf(w, d_i) > 0} p(w|q_j) \log \left(1 + \frac{tf(w, d_i)}{\mu p(w|D)} \right) + \log \frac{\mu}{\mu + N_{d_i}} \quad (18)$$

我们选择[3]中提出的最佳方法作为第三比较方法。这种方法的直觉是选择一组审阅者，该审阅者共同努力实现对查询主题的最相似的主题分布。具体地，给定纸质查询，该方法通过选择具有最小kl分歧值的新审阅者选择审阅者：

$$D(\theta_q || \theta_{r_1, \dots, r_{k-1}}^{r_k}) \quad (19)$$

其中YQ是查询和YRK

R1的主题分布，……，RK-1是先前选择的评论者R1，Y，RK，即，即，

$$p(z|\theta_{r_1, \dots, r_{k-1}}^{r_k}) = \frac{\sigma}{k-1} \sum_{i=1}^{k-1} p(z|\theta_{r_i}) + (1-\sigma)p(z|\theta_{r_k}) \quad (20)$$

其中S表示依赖先前选择的审稿人R1，Y，RK-1。在下面的讨论中，这三种方法将被称为基线-PR，基线-KL和主题-KL。与负责任主题相匹配。最大化全球覆盖率实际上是一个np-hard set涵盖问题。因此，我们添加了一个假设使其易于，其中对于特定的查询QJ，每个分配的专家 v_i 只能选择一个负责主题 t_{q_j} ，并介绍了查询的本主题。假设的优化问题提供了原始的优化问题，并且可以轻松地将我们的框架中。如图7所示，我们为查询QJ创建9TQJ9-1节点，其中QJ对应于查询，QJK表示QJ的第k相关主题。我们将两个弧从QJ添加到QJK，每k，分别具有 σ_{yjk} 和0的成本。弧的成本正在映射到主题覆盖范围，其中S是用于调整重要性的乘数。然后，对于相关主题包括k的每个专家 V_i ，我们将QJK添加到节点 v_j 的弧，这意味着专家能够负责主题。VJ连接到专家节点 V_i ，其中节点 v_j 用于避免查询分配给同一专家两次。结合匹配

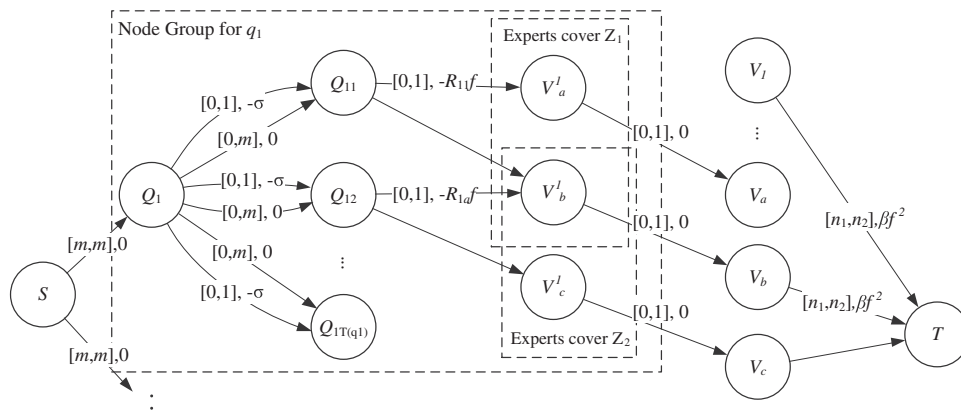


图7.与负责主题匹配的网络配置。

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用性。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， $9t_{qj}$ 和 $9t_{vi}$ ）的数量，而 $T(Qj)$ ， $T(Vi)$ 是通过选择YQJ和YVI中的Top-K主题来确定，因为我们在部分中讨论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 $T(Qj) = 1(vi)$ 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。

表5对四种不同措施的所有方法的比较：覆盖，信心，verageCoverage和fscore。学习主题模型时，主题的数量被设置为20。

Methods	Coverage (%)	Confidence (%)	AverageConfidence (%)	F _{score} (%)
Baseline-Pr	74	62	46	63
Baseline-KL	75	62	45	63
Topic-KL	87	58	53	67
我们的方法87 60 51 67我们的方法-L10 86 58 49 66我们的方法-15 80 59 48 65				

5.3. Course–teacher assignment experiment

数据集。在课程作业实验中，我们手动爬出了四所顶级大学计算机科学（CS）的研究生课程，即CMU，UIUC，Stanford和MIT。总共有2008年秋季学期的609个研究生课程春季，每门课程都被1–3名教师指示。我们的直觉是，教师的研究兴趣往往与他/她正在教学的研究生课程相匹配。因此，我们仍然使用教师最近（五年）出版物作为他们的出版物

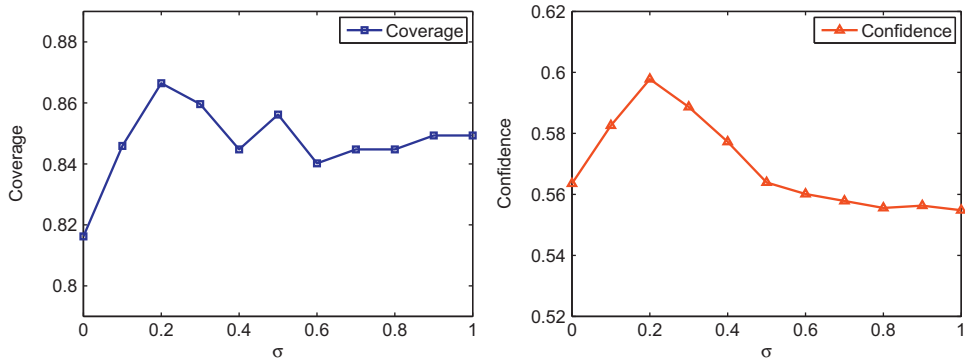


图8.不同拉格朗日乘法器的覆盖和信心的敏感性。

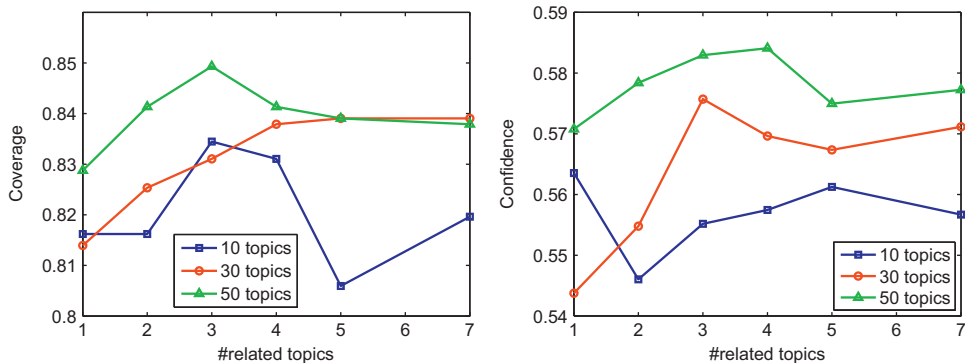


图9.不同主题模型的覆盖率和相关主题次数的敏感性，S = 0：3。

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。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。

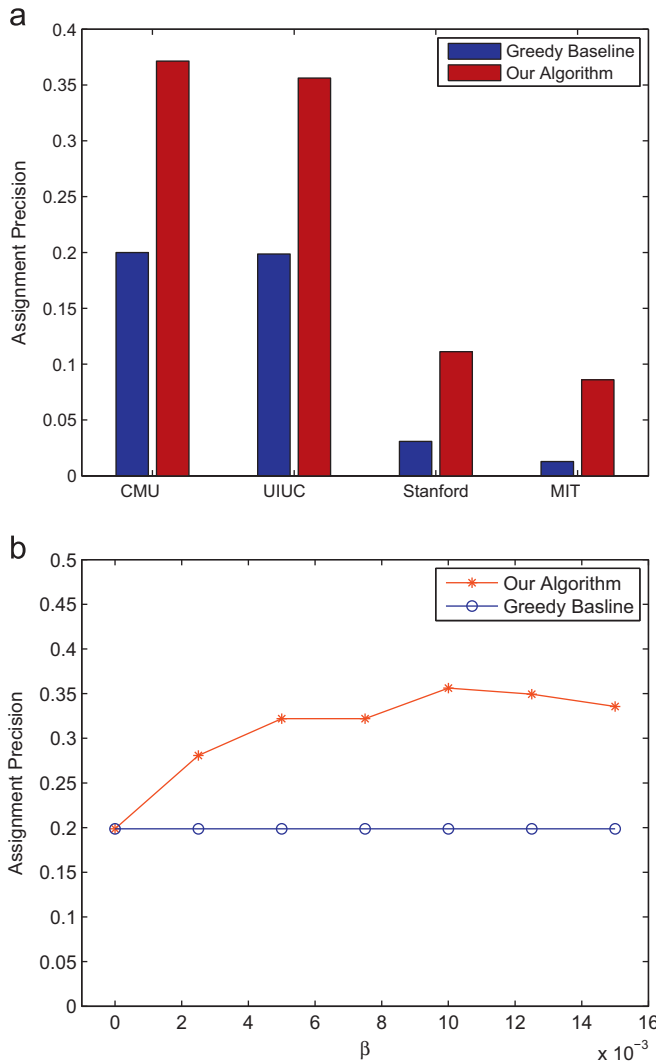


图10.课程教师分配性能 (%)。(a) 课程分配结果和 (b) 对UIUC数据的精度与B。

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，春天）

成翔程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: # Papers/Reviewer: Lower bound: Upper bound: Beta:

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

[-] Alexandre V. Evfimievski (See Details)

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

[+] 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.2009AA8Z138）。最后一位作者访问清华大学时，这项工作完成的。

References

- [1] C.B.Haym, H. Hirsh, W.W. Cohen, C. Nevill-Manning, 通过开采网络推荐论文，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，Sigir'92：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.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. MIMNO, 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.太阳, J.Ma, Z.-P. FAN, J. WANG, 审查员作业的混合知识和模型方法，在：40夏威夷系统科学国际会议（HICSS-40 2007），2007，p. 47。[8] S. Hettich, M.J. Pazzani, 提案审查员挖掘：在国家自然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, Constance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届工业和工程应用程序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, Cyber Chair：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/wbgafpr.pdf。[12] Microsoft会议管理工具包（CMT），URL / http://cmt.Research.microsoft.com/cmt/s。[13] EasyChair软件，URL / http://www.easychair.org/s。[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生组织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. TAN, 通过基于COSNTRAIINT的优化匹配的专业知识，IN: 2010年IEEE E / WIC / ACM国际会议关于WEB Intelligence（WT2010），2010年的课程。[18] D. YIMAM, A. Kobsa, Demoir：用于专业知识建模和推荐系统的混合架构，在：第九IEE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李, TREC 2005企业轨道专家搜索研究，IN: TREC, 2005。[20] Y.Fu, W. Yu, Y. Li, Y.刘, M. 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。数据库的进步：概念系统和应用程序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：超越知识管理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, 专家查找的语言建模框架, 2008年. [28] X. Wei, W.B. 克罗斯特, 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, p.a.c. Valdivieso, 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, 单一前遗传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, 宾夕法尼亚大学, 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, 神经计算机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IEEE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. Bersdi, 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/s>.



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。



天丽是北京航空航天大学的本科学士。



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。