

Beta 在线匹配

黄晓宇¹, 曾青松², 杨磊¹

(1. 华南理工大学经济与贸易学院, 广东广州 510006; 2. 广州番禺职业技术学院信息工程学院, 广东广州 511483)

摘要: 二部图的在线匹配问题最早由 Karp 等人在 1990 年提出, 该问题在近年得到了广泛的关注, 在日常生活中的应用. 本文引入了 *Beta* 分布作为二部图节点间的邻接关系的统计先验, 提出了最大化节点的预留匹配能力准则作为在线匹配策略的评价度量, 设计了在线匹配算法 *BetaOM*, 并证明了该算法的正确性. 本文把 *BetaOM* 分别应用于基于人造数据和真实数据的在线匹配问题, 实验的结果显示该算法优于经典的 Greedy 算法和 Ranking 算法.

关键词: 二部图; 在线匹配; *Beta* 分布; 随机优化

中图分类号: O212.3 **文献标识码:** A **文章编号:** 0372-2112 (2017)05-1268-04

电子学报 URL: <http://www.ejournal.org.cn> **DOI:** 10.3969/j.issn.0372-2112.2017.05.033

Beta Online Matching

HUANG Xiao-yu¹, ZENG Qing-song², YANG Lei¹

(1. School of Economics and Commerce, South China University of Technology, Guangzhou, Guangdong 510006, China;

2. School of Information and Technology, Guangzhou Panyu Polytechnic, Guangzhou, Guangdong 511483, China)

Abstract: In recent years, the online bipartite graph matching problem has attracted substantial research attention as many real-life problems can be eventually reduced to it. In this work, we study the classic online matching problem that was initially investigated by Karp in 1990. We adopt the *Beta* distribution as the prior distribution of the adjacency relation among the nodes, and present a novel measure to evaluate the matching policy. We also design *BetaOM*, a *Beta* distribution based online matching algorithm, and mathematically prove its soundness. Experiments with *BetaOM* as well as the benchmark algorithms on both synthetic and real data demonstrate that the proposed *BetaOM* is superior to the others.

Key words: bipartite graph; online matching; *Beta* distribution; stochastic optimization

1 引言

本文研究如下的二部图在线匹配问题: 考察两个互不相交的节点集合 $S = \{s_1, s_2, \dots, s_n\}$ 和 $T = \{t_1, t_2, \dots, t_n\}$, 其中 S 为给定集合, 它的每个成员节点最多只能与一个 T 中的节点发生匹配; T 中的节点依下标递增顺序顺次到达, 每个节点 t_j 在抵达的同时还将揭示一个由 S 中所有和它相邻的节点构成的集合 N_j , 若 N_j 中存在不少于一个可用的节点, 则我们从这些节点中选择一个与 t_j 匹配, 并称与 t_j 的匹配是“成功”的, 反之, 则称该匹配为“失败”的. 我们的任务是要设计一个合适的匹配策略 π , 使之能为每个新来的节点 t_j , 选出合适的节点 s_j 与之相配, 使得在所有匹配完成后, 累计的成功匹配数目达到最大. 我们强调在上述设置中, π 从每个 N_j 中选择的是“可用的”节点, 原因在于, 在由 S, T 中的

所有节点构成的完整二部图中, 每个 s 节点都可以同时与 T 中多个不同的节点相邻, 因而相应的, 它也可以出现在多个不同的邻接集合 N 中. 但注意节点 s 的实际匹配次数不能大于 1, 所以, 对每个节点 t_j , $|N_j| > 0$ 并不等价于 N_j 中存在可用的 s 节点.

本文的研究属于经典的二部图在线匹配问题, 对该问题的研究最早可以追溯到 1990 年 Karp 等人的工作^[1]. 近 20 年来, 研究者陆续提出了该问题的多种变体, 如 *b-Matching*^[2]、*Adwords*^[3] 以及其它扩展^[4] 等. 但需要指出的是, 这些扩展的求解策略多数仍然植根于文献^[1] 中提出的 Ranking 算法. 事实上, 对文献^[2~4], 当它们研究的问题退化为与文献^[1] 相同的设置时, 相应的求解策略也都与文献^[1] 中提出的 Ranking 算法等价^[5].

当前, 在对在线匹配算法的研究中, 一个重要的特

点是算法的设计和分析普遍以对手攻击设置或完全随机设置为前提,对于这一现象,Bahmani 与 Kapralov 有非常精到的评价:“(对在线匹配模型)无论是对手攻击假设还是完全随机假设都过于严格了.事实上,在现实的场景中,对于具体的问题,我们经常有丰富的统计先验信息可供策略设计者使用.”^[6]受此启发,在本文的工作中,我们将研究先验知识在问题求解中的应用.具体的,我们将为节点集合 S 与 T 间的连接关系引入统计先验(Beta 分布),借助这一先验,我们得以充分刻画和利用问题的领域信息,从而有望获得更好的匹配结果.

2 算法

2.1 算法设计

我们沿用第 1 节的记号,此外,对每个节点 t_j ,我们假设节点 s_i 以概率 $p_i \sim \text{Beta}(a_i, b_i)$ 与之相邻,这里 a_i 是截至节点 t_j, s_i 累计被选中匹配的次數, $b_i = j - a_i$ 则为 s_i 未被选中的累计次数.由熟知结论, $p_i = \frac{a_i}{a_i + b_i}$.

我们首先给出如下定义:

定义 1 (节点的预留匹配能力) 对 S 元素的任一排列 $\pi = s_{z_1} s_{z_2} \dots s_{z_n}$, 其中 $\{z_1, z_2, \dots, z_n\} = \{1, 2, \dots, n\}$, 定义 π 中第 i 个位置上的元素 s_{z_i} 的预留匹配能力为:

$$F_{\pi}(i, s_{z_i}) = \sum_{k=i+1}^n \sum_{j=k}^n p_{z_j} \quad (1)$$

定义 1 有非常直观的解释:我们把 π 视为一个在线匹配策略,对每一个到达的元素 t_i, π 从 S 中选择元素 s_{z_i} 与之相配,并计算后续匹配(从第 $i+1$ 轮开始到最后)成功的概率.注意对第 $k(\geq i+1)$ 轮匹配, $s_{z_i} \sim s_{z_{i+1}}$ 已不能使用,所以该轮匹配成功的概率为 $\sum_{j=k}^n p_{z_j}$.我们对所有可能的 k 值求和,即得到了式(1).

对于第 i 轮匹配,记 S 中所有尚未匹配的节点集合为 U_i ,记 $A_i = U_i \cap N_i$,容易看出,若 $|A_i| > 1$,则以下式(2)的目标是一个较为合适的选择,其原因在于它为后续的匹配成功提供了最大的可能.

$$s_{z_i} = \arg \max_{s \in A_i} F_{\pi}(i, s) \quad (2)$$

下面计算 $F_{\pi}(i, s)$,我们有:

$$\begin{aligned} F_{\pi}(i, s) &= \sum_{k=i+1}^n \sum_{j=k}^n p_{z_j} = \sum_{j=i+1}^n \sum_{k=i+1}^j p_{z_j} \\ &= \sum_{j=i+1}^n (j-i)p_{z_j} \end{aligned}$$

根据熟知的排序不等式,当 $F_{\pi}(i, s)$ 达到最大时,如下论断成立:

结论 1 诸 $p_{z_j}(j = i+1, i+2, \dots, n)$ 与下标 j 以相同的递增序排列.

结论 2 与式(2)的解 s_{z_i} 对应的概率 p_{z_i} 必定是集合

A_i 中各节点对应的概率中的最小者.

根据上述讨论,我们得到算法 1 的设计.

算法 1 Beta 在线匹配算法(BetaOM)

输入:待匹配节点集 $S = \{s_1, s_2, \dots, s_n\}$.

1. 初始化 n 维计数向量 $c = [0, 0, \dots, 0]'$;
2. FOR $k = 1$ to n
3. 揭示第 k 个新来节点 t_k 在 S 中的邻居向量 v , 其中若 t_k 与 s_i 匹配,则置分量 $v_i = 1$,否则置 $v_i = 0$;
4. 更新 $c \leftarrow c + v$;
5. 根据 v , 从 S 中选择对应的 c 分量最小的可用节点作为 t_k 的匹配节点;
6. ENDFOR

2.2 算法分析

本节用于分析算法 1.

我们首先指出,在算法 1 中,计数向量 c 具有如下性质:

性质 1 (逼近性) 对 $1 \leq i \leq n$, 记节点 s_i 与 T 中节点匹配的真实概率为 p_i , 则由 Beta 分布的性质,在经过 k 轮循环之后,对 p_i 的经验估计为 $\frac{c_i}{c}$. 根据 Hoeffding 不等式,对任意常数 $\varepsilon > 0$, 有:

$$\mathbb{P}(|p_i - \frac{c_i}{c}| > \varepsilon) < 2\exp(-2k\varepsilon^2) \quad (3)$$

上述性质显示,当 $k \rightarrow \infty$ (从而 $n \rightarrow \infty$), 由算法 1 得到的估计 $\frac{c_i}{c}$ 将逼近真正的 p_i .

性质 2 (保序性) 不失一般性,令 $p_1 < p_2 < \dots < p_n$, 记 $\delta = \min\{p_{i+1} - p_i | i = 1, 2, \dots, n\}$,

在不等式(3)中,我们取 $\varepsilon = \frac{\delta}{2}$, 并令不等号右端 $2\exp(-2k\varepsilon^2) < \tau$, 可得:

$$k > -\frac{2}{\delta^2} \ln \frac{\tau}{2} \quad (4)$$

结合式(3)、(4), 令 $\tau \rightarrow 0$, 则对 $k \rightarrow \infty$, 估计值 $\frac{c_i}{c}$ 将以概率 $(1 - \tau)$ 集中在以 p_i 为中心, $\frac{\delta}{2}$ 为半径的区域内,

注意根据 δ 的定义,上述事实意味着对 $1 \leq i \leq n-1, \frac{c_i}{c} < \frac{c_{i+1}}{c}$. 即对 $i \neq j, \frac{c_i}{c} \leq \frac{c_j}{c}$ 当且仅当 $p_i \leq p_j$.

根据性质(1)和(2),对每一个节点 s_i , 我们都可以使用估计值 $\frac{c_i}{c}$ 近似它被匹配的真实概率 p_i , 并且由所有这些近似值构成的序列和它们对应的真实值构成的序

列具有相同的序关系。

此外,还需指出的是,在算法 1 中,我们没有直接计算 $\frac{c_i}{c}$,这是由于,在算法的第 k 轮循环中,为获得式(2)的解 s_{z_i} ,根据前文讨论的结果,我们只需知道各 S 节点对应的匹配概率之间的序关系即可,而无需计算其精确值。

以下我们证明算法 1 的合理性,为此,我们首先引入下述假设 1。

假设 1(可替代节点假设) 对于任意两个节点 s_i 与 s_j ,若以下两个条件同时成立:

- (1) $p_i \leq p_j$;
- (2) 存在某个 k ,使得 $\{s_i, s_j\} \subset N_k$.

则对所有的 $k < l \leq n$,有: $s_i \in N_l \rightarrow s_j \in N_l$,即在所有可能的匹配中,节点 s_i 均可被 s_j 替代。

在下文中,若 s_i 可被 s_j 替代,则我们记为 $s_i \leq s_j$ 。

下面我们证明:若 S 中的节点满足上述的可替代节点假设,则由算法 1 得到的匹配序列在每一个节点上的预留匹配能力都是所有可能序列中的最大者,因而在最大化预留匹配能力值的意义下,算法 1 是最佳的匹配策略。

我们首先有如下引理:

引理 1 对任一匹配策略 π ,若存在步骤 k , π 从 A_k 选择的节点为 s_{z_k} ,且 $p_{z_k} > \min\{p_z | s_z \in A_k\}$ 。则必定存在另一策略 π' ,满足:

- (1) 存在 $s_{z_k} \in S$,使得 $F_{\pi'}(k, s_{z_k}) > F_{\pi}(k, s_{z_k})$;
- (2) 对所有其它步骤 $l \neq k$, $F_{\pi'}(l, s_{z_l}) \geq F_{\pi}(l, s_{z_l})$ 。

证明 由条件(1),令 $p_{z_k} = \min\{p_z | s_z \in A_k\}$,记 p_{z_k} 对应的节点为 s_{z_k} (若 A_k 中有多个节点对应的均值取最小值,则从中取任意一个为 s_{z_k})。记 π 得到的匹配序列为 $\pi = s_{z_1} s_{z_2} \dots s_{z_k} \dots s_{z_n}$,注意这里允许匹配失败的情况,即允许存在节点 t_l, A_l 为空。以下分两种情形讨论:

情形 1 π 中包含了节点 s_{z_k} ,由 $s_{z_k} \leq s_{z_k}$,注意 s_{z_k} 必出现在 s_{z_k} 之后,因此,根据假设 1, s_{z_k} 与 s_{z_k} 在 π 中的位置可以互换,容易验证,两者交换后得到的结果序列即为 π' 。

情形 2 π 中不包含节点 s_{z_k} ,此时我们使用 s_{z_k} 代替 π 中的 s_{z_k} ,则替换后得到的序列即为 π' 。

根据引理 1,对匹配策略 π ,当且仅当 π 在每个候选集 A_k 中选取的节点都对应了与 A_k 相应的诸 p_z 参数中的最小者时, π 在每一个节点上的预留匹配能力都能取得最大值,因此,在最大化节点的预留匹配能力的意义下,算法 1 是最优的匹配策略。

3 实验

3.1 数据描述

我们分别使用 5 个不同的数据集 DS-1 ~ 5 来检验

本文算法的表现,其中 DS-1 和 DS-2 是我们根据特定分布生成的仿真数据,DS-3 ~ 5 则来源于真实的业务系统。

在 DS-1 中,我们设定 $|S| = |T| = 1,000$,且 S 和 T 两者间的邻接关系的分布较为均匀,每个 s_i 节点都以相应的概率 $p_i \in [0.02, 0.1]$ 和 T 中各节点相邻;在 DS-2 中,我们同样令 $|S| = |T| = 1,000$,但 S 和 T 两者间的邻接关系的分布则表现为一种高度不平衡的形式,我们把 S, T 间的 80% 的邻接关系都限制分布在 S 中前 20% 节点与 T 之间。

DS-3 来源于某婚恋网站,节点集 S, T 分别对应网站两个不相交的会员集合,这里 $|S| = |T| = 10,000$,仅当一对会员 $\langle s, t \rangle \in S \times T$ 在网站上匹配成功时,与他们对应的节点间才有邻接关系。

DS-4 出自电影评分数据集 MovieLens10M^①,我们根据一致分布从中随机选取了约 15,000 部电影构成集合 S ,选了 10,000 个用户构成集合 T ,当且仅当 T 中用户 t 对 S 中电影 s 给了评分时,我们才在它们对应的节点间建立邻接关系。

DS-5 来自众包点评网站豆瓣^②,包含了 100,000 个用户(T)对 80,000 本图书(S)的评分记录,与 DS-4 类似,若用户 t 对图书 s 给了评分,则我们在它们对应的节点间以边相连。

在所有实验中,我们都采用文献[1]中提出的 Greedy 和 Ranking 作为对比算法,我们把它们在 DS-1 ~ 5 上的匹配结果分别与 BetaOM 在相同数据上的匹配结果相比较,以评估 BetaOM 的表现。

表 1 随机设置实验结果(匹配成功均值)

	DS-1	DS-2	DS-3	DS-4	DS-5
Greedy	979.6	946.5	4,828.4	6,902.3	52,624
Ranking	977.9	945.5	4,755.2	7,323.2	51,154
BetaOM	978.4	993.1	4,945.0	8,410.5	54,468

3.2 随机实验

我们固定 DS-1 ~ 5 中各 s 节点的排序不变,各 t 节点则以随机排列顺序到达,对每一新来的 t 节点,我们分别以 Greedy、Ranking 和 BetaOM 把它匹配给 S 中节点。我们在每个数据集上都重复 10 次实验,并把各算法在每个数据集上的平均成功次数汇总在表 1 中。

从表 1 可以看出,在 DS-1 上, BetaOM 与基准算法的平均匹配成功次数均相当接近,但在 DS-2 ~ 5 上, BetaOM 的表现都显著优于对比算法,这显示了 BetaOM 的优越性。

① <http://grouplens.org/datasets/movielens/>

② <https://www.douban.com/>

3.3 对手攻击实验

注意 *BetaOM* 的实质是“把最容易匹配的节点留到最后”,所以一种自然的攻击思路是要努力误导 *BetaOM* 尽早消耗掉“最容易匹配”的节点.为此,我们把 DS-1 ~ 5 中各 s 节点根据其邻接的 t 节点数量作降序排列,对该排列中前 20% 的节点,我们统计 T 中各节点和它们间的邻接数量,并根据统计的结果对各 t 节点作升序排列.

我们让 T 中各节点按排序顺序依次到达,对每一新来的 t 节点,我们分别以 Greedy、Ranking 和 *BetaOM* 把它匹配给 S 中节点,实验的结果的统计汇总在表 2 中.

表 2 对手攻击设置实验结果(匹配成功次数)

	DS-1	DS-2	DS-3	DS-4	DS-5
Greedy	994	930	4,952	7,510	54,879
Ranking	994	939	5,170	7,840	54,467
<i>BetaOM</i>	993	990	5,187	8,191	55,355

表 2 的结果显示在攻击设置下,*BetaOM* 的表现仍然优于基准算法.此外,对比表 1 和 2,我们的一个吃惊的发现是表 2 展现的结果完全超出预期:在 DS-1 ~ 5 这五个数据集上,攻击策略仅在 DS-2 上取得了一定效果,但该策略对 *BetaOM* 的影响远远低于它在基准算法上造成的影响;在 DS-1 和 DS-3 ~ 5 上,攻击策略反而进一步提升了所有算法的匹配成功率.为探究其原因,我们对实验数据作了进一步的统计分析,结果显示,与 S 集合中“最好”的那部分节点连接最少的 t 节点亦即 T 中与 S 全集连接最少的节点,因此,上述攻击策略在客观上使得“最难”匹配的 t 节点集中出现在了其它“容易”的节点之前,因而提高了这些“最难”的节点的匹配成功率,而对于其它“容易”的 t 节点,由于可以和它们匹配的 s 节点分布较广,因而即使人为的把它们出现的顺序推后,也不会对匹配的成功率造成太大的影响.

由此可见,*BetaOM* 具有较强的抵御攻击能力.

4 结论

二部图的在线匹配模型在近年得到了广泛的应用,本文研究了经典的二部图在线匹配问题.在本文的工作中,我们为二部图节点的邻接关系引入了 *Beta* 先验,以此为依据,我们提出了评价二部图的边匹配的量化度量—节点的预留匹配能力,基于该度量,我们设计了二部图在线匹配算法 *BetaOM*,并证明了该算法的正确性.我们把提出的算法分别应用于人造数据和真实数据,实验的结果显示本文的算法优于经典的 Greedy 算法和 Ranking 算法.

对于未来的工作,我们认为,本文的研究还可以从如下两方面作进一步的拓展:(1)在传统的在线匹配研究中,最大化算法的竞争比一直是相关工作关注的中心议

题;与之相比,本文的研究却围绕文中提出的最大化节点的预留匹配能力准则而展开.观察本文的实验结果,我们猜测以上两个准则间应具有一致性,但具体的分析还有待在后续的工作中开展.(2)本文算法的成功得益于统计先验 *Beta* 分布的引入,对此,一个自然的考虑是,能否通过引入表达能力更强的分布,设计更好的在线匹配算法?甚至更进一步的,是否存在一般性的准则,能指导我们根据具体问题的领域特点,选择合适的先验分布?对这些问题的研究,也具有非常重要的意义.

参考文献

- [1] Richard M Karp, Umesh V Vazirani, et al. An optimal algorithm for on-line bipartite matching [A]. Proceedings of the Twenty-second Annual ACM Symposium on Theory of Computing [C]. New York: ACM, 1990. 352 – 358.
- [2] Bala Kalyanasundaram, Kirk R Pruhs. An optimal deterministic algorithm for online b-matching [J]. Theoretical Computer Science, 2000, 233(1): 319 – 325.
- [3] Aranyak Mehta, Amin Saberi, Umesh Vazirani, Vijay Vazirani. Adwords and generalized online matching [J]. Journal of the ACM, 2007, 54(5): 22 – 40.
- [4] Gagan Aggarwal, Gagan Goel, et al. Online vertex weighted bipartite matching and single-bid budgeted allocations [A]. Proceedings of the Twenty-second Annual ACM-SIAM Symposium on Discrete Algorithms [C]. Philadelphia, USA: SIAM, 2011. 1253 – 1264.
- [5] Aranyak Mehta. Online matching and ad allocation [J]. Theoretical Computer Science, 2012, 8(4): 265 – 368.
- [6] Bahman Bahmani, Michael Kapralov. Improved bounds for online stochastic matching [A]. European Symposium on Algorithms [C]. Springer, 2010. 170 – 181.

作者简介



黄晓宇 男,1977 年生,广东茂名,博士、讲师,研究方向为统计与优化理论、机器学习.
E-mail: echxy@scut.edu.cn



曾青松 男,1976 年生,湖南邵东人,博士、副教授,研究方向为模式识别与数据挖掘.