直流 信息组织 信息批判

35-39

计算机科学1997 Vol. 24№. 6

基于 Intranet 的集成化信息系统信息的组织和规划*)

Information Organization and Planning of Intranet Based Integrated Information System

陈卫东 杨建军 鲁东明 潘云鹤

G202

F2](.)

(浙江大学人工智能研究所 杭州 310027)

篇 要 In order to optimize the Intranet information distribution, this paper introduces "SPP" and "KNAPSACK" to solve the "Access Bottleneck Problem" and "Communication Cost Problem" effectively, moreover, it provides a new idea to the similar problem solving.

关键词 Internet.Intranet.Web, Network.Access bottleneck.Communication cost.SPP.KNAP-SACK

1. 问题的提出

一个企业的专用网即 Intranet 网,在 Internet 上具有两种功能。对外,主动发布信息,介绍其最新产品和技术,在公众面前为企业作宣传,以图丰厚的利润回报;对内,其自身也是一个 Internet 用户,必然要访问内部网以外的各种信息,以便充分了解市场,在商业竞争中保持有利的地位。面对如此浩瀚无垠的 Internet 信息和来自企业内外众多的访问请求,必须设法有效地组织和规划内部网的信息资源,以达到减轻服务器的负担,降低访问费用,提高 Intranet 网使用效率等目的。为了更好地理解问题本身,现归纳成如下两个方面的问题。

问题之一:一般地,在企业内部分布信息的时候,总是在逻辑上将相应的信息主题分成块结构。这样,内容上彼此独立的信息往往组成不同的信息块,分布在企业内部不同的 Web 服务器上。由于各信息块内容不同,其重要性及感兴趣的程度也不同,因此它们被访问的频率也必然不一样。如果这些信息块毫无规律地安置在各 Web 服务器上,必然导致各个服务器被访问的频率差别显著。最极端的情况,则是其中的一两台服务器被频繁访问。甚至有可能负担过重而拥塞,而其他服务器则显得"门前冷落鞍马

稀"。一旦出现这种情况,最直接的后果是资源浪费, 更进一步说,因为频繁地告诉用户"因访问者太多, 请您稍后再试"诸如此类的信息,会使网上用户对此 结点失去耐心,从而失去某些无形资产。所以有必要 对各信息块被访问的频率进行统计,在此基础上建 立数学模型、确定解决方案、使得各服务器所承受的 负担尽可能地均衡,防止造成访问瓶颈。

问题之二:面对丰富多彩的 Internet 信息资源,Intranet 用户必然同时也是 Internet 用户,必定频繁地访问 Internet 网络。由于企业对 Internet 信息访问是有针对性的,有其明显的目的,所以必然只对其中的一些信息资源感兴趣,访问也比较频繁,造成通信费用的增长。另一方面,必然有不同的用户对相同的信息主题产生兴趣,屡屡击人相同的网址调用这些信息,同样也造成通信费用的增长。为了有效地降低通信费用,可以将那些被频繁访问的 Internet 信息下载至 Intranet 内部的 Web 服务器上,使之成为内部信息,Intranet 用户在内部就可以对其进行访问。一旦成为内部信息,其通信费用和访问速度就不可同日而语了。所以也必须在统计的基础上,建立数学模型,确定解决方案,切实有效地降低通信费用,提高访问速度。

上述两上问题是基于 Intranet 的分布式信息系统信息组织和规划的关键问题。前者是针对内部信

^{*)}国家863高科技项目、国家科技攻关项目、国家自然科学基金等资助。陈卫东 博士生,主要研究方向;计算机网络,人工智能,杨雕军 博士生,主要研究方向;计算机网络,人工智能, CSCW,播云鹤 校长,教授,博士生导师,主要研究方向;计算机图形学,人工智能,认知科学,CAD等。

息的组织安放的,我们概括为"访问瓶颈问题";后者则针对外部通信,我们称之为"通信费用问题"。本文将就这两个问题分别提出其数学模型和求解方法。

2. 访问瓶颈问题的数学模型及求解

2.1 数学模型

首先,我们将"访问瓶颈"问题整理成如下描述: 将一定数量的信息块存放到规定数目的服务器上,由于各信息块被访问的次数不一样,存放的原则是 使各服务器被访问的次数(为存放在其上的信息块 被访问次数之和)尽可能相等。

在数学上,这是一个集合划分问题(Set Partitioning Problem,简称 SPP)。所谓集合划分问题,是指将一给定集合划分为指定个数的互不相交的子集,并使每个子集含有的元素大小之和尽可能一致。它的判定问题严格叙述为:

(SPP)实例:有穷集合 $A=\{a_1,a_2,\cdots,a_n\}$,以及每一个 $a\in A$ 的"大小" $w(a)\in R^+$;正整数 $m\in Z^+$ 。

问:是否存在一个关于 A 的划分 $\sigma=\{A_1,A_2,\cdots,A_m\}$,使 得 对 \forall ; == 1, 2, ..., m , 有 $\sum_{a\in A_j}w(a)=\frac{1}{m}\sum_{a\in A}w(a)$?

由此,"访问瓶颈"问题可以形式化描述为:将一 给定的信息主题集合 A=={a₁,a₂,···,a_n}(n 个信息主 题)划分为 m 个互不相交的子集(即分配至 m 台服 务器),并使每个子集(服务器)含有的元素大小(受 访问频率)之和尽可能一致。

在数学上,称使各子集中元素大小之和尽可能一致的 SPP 为 SPP 优化问题。对于它,人们通常从两个方面来考虑。其一是使最大的子集最小化(记为minimax-SPP),另一是使最小的子集最大化(记为maximin-SPP),下面给出其数学描述;

我们以正数集 $P = \{p_1, p_2, \cdots, p_n\}$ 代替 A,其中 p,就代表 p。的大小,P 划分成 $\sigma = \{M_1, M_2, \cdots, M_m\}$ m 个子集,同时 C。代表 M,中元素之和。那么,minimax-SPP 定义为:

 $\underset{\sigma}{\text{minimize}} \left\{ \underset{j}{\text{max}} C_{j} \right\};$

maximin-SPP 定义为:

maximize {minC_i},其中 o 为 P 的任一划分。

2.2 LPT 算法

人们从本世纪以来一直致力于构造一些好的性 • 36 • 能保证的多项式时间近似算法。我们采用其中最自然也是最有名的 LPT 算法 (Longest Processing Time)。LPT 算法思想是:1) 把元素按单调递减顺序排列,不妨设 $p_1 \ge p_2 \ge p_3 \ge \cdots \ge p_n$,构成序列 L;2) 将 L 中的当前元素放入当前和最小的子集中,然后从 L 中去掉它。

2.3 性能比较

记 $\sigma = \{M_1, M_2, \cdots, M_m\}$ 为由 LPT 得到的划分, $S^* = \{M_1^*, M_2^*, \cdots, M_m^*\}$ 为相应问题的**极**优划分, 并且分别记 M_1 和 M_1^* 的元素和为 C_1 和 C_1^* 。令 $M = \max_i C_1$, $W = \min_i C_i$; $M^* = \max_i C_i^*$, $W^* = \min_i C_i^*$, 则 有,

对于 minimax-SPP:
$$\frac{M}{M} \le \frac{4}{3} - \frac{1}{3m}$$
对于 maximin-SPP: $\frac{W}{W} \ge \frac{3m-1}{4m-2}$

2.4 推广到 NSPP

在企业专用网内,因为每个部门均有自己的 Web 服务器。在一般情况下,与本部门有关的信息 块将存放在该部门的 Web 服务器上,以便于使用和 管理。也就是说,有可能会指定某些信息块必须存放 在确定的服务器上,针对这种情况,相应地,我们可 将 SPP 进一步推广,这便是带核集划分问题,简称 为 NSPP。在数学上,可以描述为。

有限集 $A == \{p_1, p_2, \cdots, p_n\} \cup \{g_1, g_2, \cdots, g_m\}, 其$ 中 $P == \{p_1, p_2, \cdots, p_n\}$ 是非核元集,我们称对应的信息主题为非核信息主题集, $G == \{g_1, g_2, \cdots, g_m\}$ 为核元集,我们称对应信息主题为核信息主题集。它与SPP 区别在于在 A 的划分 $\sigma == \{M_1, M_2, \cdots, M_m\}$ 中,必须有 \forall j $g_i \in M_j$,其中对核元 g_i ,要求 g_i 非负,于是类似地得到两个问题 maximin-NSPP 和 minimax-NSPP。

当 LPT 应用于它们时,首先我们放置 g_i 于 M_i ,然后再按 LPT 放置非核元。我们称这样的经过修正的 LPT 为 MLPT。

同样,maximin-NSPP和 minimax-NSPP也有一个性能问题。再次引用上一节中的标记 M,M*,W和 W*:

对于 minimax-NSPP,
$$\frac{M}{M} \le \frac{3}{2} - \frac{1}{2m}$$

对于 maximin-NSPP,
$$\frac{\mathbf{W}}{\mathbf{W}} > \frac{2m-1}{3m-2}$$

在数学上,算法的分析必须考虑到最坏情况。但是,在工程上,这种极端的情况是很少出现的。所以,

引人 LPT 和 MLPT 算法就能有效地解决"访问瓶颈"问题。我们可以利用工具软件,首先对各带核和不带核的信息主题进行统计和分析,在此基础上对各信息主题进行优化组合,使得各 Web 服务器信息量分布合理,负担均衡,并使资源合理配置,达到 Intranet 稳定有效运行之目的。

3. 诵信费用问题的数学模型及求解

3.1 数学模型

首先,我们将"通信费用"问题整理成如下描述: 现企业内部网中有一节点专门用来存放从外部调入 的信息。因该节点容量有限,问怎样有目的地选择信 息放人节点,在不超过最大容量的情况下,使节点的 使用效率最高。在数学上,这是一个背包问题,其数 学描述如下:

设有一位旅行家,在出发前准备一只背包,限制背包内物体的总重量不超过 b。现有 n 类物体 p_1 , p_2 ,…, p_a , p_i 类中每个物体的重量为 w_i ,价值为 v_i 。 $1 \le i$ $\le n$,物体不能拆开,问怎样把物体放入背包内,在不超重的条件下使背包内物体的总价值最大?

设 $x_i = \{1,0$ 表示不装人背包,1表示装人背包。 则背包问题可归纳为下列数学问题:

目标函数;
$$\max z = \sum_{i=1}^{n} v_i x_i$$
;

求满足约束条件的并使目标函数值最大的解 (x_1,x_2,\cdots,x_n) 。

由此,"通信费用"问题可以形式化描述为;设在内部网 Intranet 之外有 n 个信息块,其信息量大小为 Si,每调用一次的费用为 Fi,对每个信息量为 Si 的信息块在一定时间内的访问次数为 Ti。现在要选择部分信息单元调入内部网某一站点 a,a 所允许的最大容量为 M。则在信息内容不能全部调入的情况下,问怎样把一些信息单元调入 a,在不超过最大容量的情况下,a 内信息的实用价值最大。

设 x₁={{}·0表示不调人内部网·1表示调人内部网。则通信费用问题可归纳为下列数学问题:

目标函数;
$$\max_{z=\sum_{i=1}^n v_i x_{i-1}}$$

求满足约束条件的并使目标函数值最大的解

 (x_1, x_2, \dots, x_n)

这里, $v_i = F_i * T_i$,表示所调人 a 的信息,若费用 越大,访问次数越多,则实用价值越大。

3.2 诵信费用问题的算法

在数学上, 背包问题可以用很多算法求解, 其中最直观的算法是"贪婪算法"。它是依照一种最优度量法即每步都取局部最优值最后求得解的过程。对应"通信费用"问题, 其算法思想可以简单描述如下: 先把要调人内部网的各信息块按 v./S. 的比值排序、比值高的信息块排在前面, 比值较低的排在后面, 足值、成序列。然后, 尽量挑选比值大的信息块先调入 a. 到最后, 可能还剩余部分容量小于按顺序该调入 a 的信息块的容量。此时, 依次比较其后的每一个信息块, 直到找到一个可放人的为止。若所有的信息块容量均大于剩余容量,则前面调入 a 的信息块即为解。

"贪婪算法"虽然很直观,但是它的解并不一定是最忧解,有时甚至与最忧解相去甚远,误差超过可以忍受的范围。所以,在精度较高的场合下,一般使用分枝定界算法,该算法对一切可能的状态逐个搜索比较,将出现指数型的复杂性。因为我们主要是利用这种算法解决工程问题,必须综合考虑精度和复杂性。在兼顾两方面因素的前提下,本文采用一种近似算法,其算法思想主要是;

1) 问题的解由向量 (x_1,x_2,\cdots,x_n) 表示,取部分向量 (x_1,x_2,\cdots,x_k) 表示部分解, $k \le n$ 。此时 a 内信息总量为 $\sum_{i=1}^k S_i x_i$ 。存在以下2种可能情况:

①存在某个 j < k,使得 $M - \sum_{i=1}^{K} S_i x_i \ge S_i$,这表示在 a 内还可以放入某个信息块而容量不至于超过 M。则依次比较最后一个下载信息块之后的 v_i/S_i ,直 到找到一个放得下的信息块为止。若还有剩余容量,则做循环,直到所有的任何信息块均比较完。

②对于所有 j>k, $M-\sum_{i=1}^{k}S_ix_i \leqslant S_i$,这说明在 a 上已不能放入任何一个信息块,否则会超过 M。则 x_i =0, 当 j>k 时。

2)取部分向量(x₁,x₂,····,x_k,x_{k+1})表示部分解。 转以上步骤。

图中,h 为小于 n 的整数,IM 存放最佳结果的 I。

当 k≠0时,k 个元素的子集数目为(t)<nk,(3)

=1. $\sum_{k=0}^{k} {n \choose k} \le \sum_{k=0}^{k} n^k = \frac{n^{k+1}-1}{n-1} = O(n^k)$,故近似算 法的复杂性为 $O(n^{k+1})$ 。

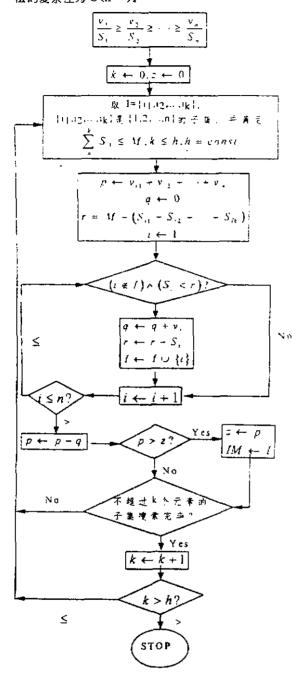


图1 通信费用问题算法框图

3.3 性能比较

设通信费用算法的最优解为 Z*,上述算法的解

为 2,可以证明:

$$\frac{Z'-Z}{Z} < \frac{1}{h+1}$$

该算法的性能与 h 的取值有关。当 n≥h*时、别 Z=Z*,即近似算法得出的解是最优解。当 h<h*时,此时近似算法不能保证一定获得最优解。

通过引人背包问题及其解法,我们成功地解决 了通信费用问题,它具有精度高,容易用程序实现等 优点,对于企业内部网具有重要的实际意义。

4 Web 使用情况统计技术

解决"访问瓶颈"问题和"通信费用"问题首先必须对 Web 使用情况进行统计。这两种问题的统计基础有所不同。对第一个问题的统计是基于 Web 服务器的一组日志文件,包括:访问日志、代理日志、错误日志。

所有的服务器访问信息都记录在访问日志(access log)中。访问日志的格式包括请求系统的主机名或 IP 地址、时间和日期、路径和被请求文档的文档名、请求成功信息或失败代码,以及传输的字节数。

关于读者阅读页面的信息越详细。Web 管理员就越能更好地跟踪 Web 的使用效果。

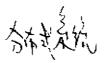
代理日志(agent log)包括来访问的代理的名称和版本号(通常是一个 Weh 浏览器)。Web 管理员可以通过此日志文件来查看读者通常使用哪一种浏览器来阅读他们的信息。

在错误日志(error log)文件中,服务器记录了 所有的错误信息,它包括请求那些并不存在的页面, 以及那些不被允许的页面,试图运行并不存在的 CGI文件,以及用户验证错误等等。错误日志文件包 括日期和时间,发出请求的主机名和域名,错误的性 质及错误的原因。

由于 Web 服务器的日志文件是一种普通格式的文件,所以任何人都可以写一个程序来对它进行分析处理。如果一个服务器的日志文件是普通格式的,那么为分析这个服务器日志文件而写的程序,对每一个服务器都是适用的。

第二个问题的统计可以利用防火墙。代理型防火墙的代理服务软件对进出内部网的通信均加以控制。因此,可以利用代理服务软件进行 Intranet 内部用户对外访问的站点、主题以及传输字节的统计。

39-43



组与组通信

计算机通信

维普资讯 http://www.cqvip.com

计算机科学 1997 Vol. 24 № 6

分布式系统中组与组通信机制的研究*`

Research on Group and Group Communication Mechanisms in Distributed Systems

王兴伟 张应辉 刘积仁 李华天

TP3 3

(东土大学软件中心研究部 沈阳1100061

摘 要 Group and group communication is an important research field in distributed systems. In this paper, the basic concepts on group and group communication are presented at first, and then how to classify group and its communication mechanisms is described. The problems which should be solved when designing group and group communication systems are discussed. The new challenges posed by distributed multimedia group applications are introduced. Finally some conclusions are given.

关键词 Distributed systems Group Group communication Distributed multimedia

近年来,随着分布式系统的发展,组与组通信机制不仅在传统的分布式组应用领域(如复制文件系统、分布式名字服务系统)中得到进一步发展,而且在新型的分布式多媒体组应用领域(如计算机会议系统、远程学习系统、远程会诊系统等)也日益受到人们的重视。本文研究组与组通信的有关概念和机制以及多媒体对组机制提出的新的需求,以期推动

通过对内和对外两方面使用情况的统计,为我们解决"访问瓶颈"问题和"通信费用"问题提供了基础素材。

5. 展望

分布式信息系统信息组织的优化问题包含诸多因素,本文分析并成功地解决了其中两个主要问题。但事实并非如此简单。随着企业内部网的发展和完善,有可能碰到以下问题:内部网的各个 Web 服务器容量太小,不能容纳与日俱增的信息内容。同时,各个 Web 服务器上每块信息的容量和被访问次数都不相同,问如何将这些信息分布到不同的 Web 服务器上,在不超过各个服务器容量的情况下,使得每个服务器被访问次数尽可能相等?

企业的不断发展,必然会提出新的问题,新的问

组机制的研究与发展。

1 组与组通信机制基本概念

开放分布处理参考模型 RM-ODP 将组定义为一些对象的集合。这些对象由于结构上的原因或者是由于其行为具有共性而集结成组^[1]。〈x〉组就是由具有特征关系〈x〉的多个对象组成的集合。〈x〉描述

题需要新的方法来解决。本文提出的数学模型为企业信息规划奠定了解决各种问题的基础,并为解决问题提供了有关思路。

参考文献

- [1] Rick Stout 著、個下兵、卢荧泽, World Wide Web 参 考大全、海洋出版社, 1996
- [2] 郁松年、邱伟、组合数学、国防工业出版社、1995
- [3] Michael R. Gatey, David S. Johnson, Computers and Intractability; A Guide to the NP-Completeness, W. H. Freeman and Company, San Francisco, 1979
- [4] R. L. Graham Bounds on Multiprocessing Timing Anomalies (SIAM, Appl. Math., 17(2)
- [5] B. L. Deuemayer et al. Scheduling to Maximize the Minimum Processor Fininsh Time in a Multiprocessor System. SIAM. J. Alg. Disc. Math. Vol. 3, 190-196, 1982

^{*)&}quot;九五"国家科技攻关项目96-B08资助。