

网络流中反馈优先的策略

黄月芳, 史定华

(上海大学理学院数学系, 上海 200436)

摘要: 由于网络传输过程中部分数据包的丢失或错误, 终端用户一旦遇到这种情况, 在服务结束后会重新提出服务请求. 不同于传统的先进先出服务规则, 提出了反馈优先的服务策略. 通过构造拟生灭过程, 对网络流中的反馈优先策略的排队进行了分析. 运用谱展开方法得到了它的平稳队长分布, 进而求解了两种不同反馈策略排队的时延指标, 并证明了反馈优先的服务策略能提高服务质量.

关键词: 网络流; 服务质量; 反馈策略; 拟生灭过程; 谱展开方法

中图分类号: TP271 **文献标识码:** A **文章编号:** 0372-2112 (2004) 05-0802-03

The Strategy of Feedback Priority in Network Traffic

HUANG Yuefang, SHI Dinghua

(Department of Mathematics, Shanghai University, Shanghai 200436, China)

Abstract: Once the transmission error or lost data packets occurred in network flows, the terminal user might bring forward a new service requirement. Different with the traditional FCFS discipline, the strategy of feedback priority was proposed. The queue with this policy was analyzed by constructing a quasi birth and death process. Based on the spectral expansion method, the stationary distribution of queue length was obtained. Furthermore, the delay performances of the two queues with different feedback strategies were calculated, indicating that the strategy with feedback priority can improve the quality of service.

Key words: network traffic; quality of service; feedback strategy; quasi birth and death process; spectral expansion method

1 引言

随着当今信息时代的到来, 网上信息量迅猛增长, 多媒体的应用也在经历着前所未有的发展. 由于网络的数据流量以指数形式扩充, 在实际应用中仅靠单纯增加带宽已无法解决问题, 这对采用何种服务策略以提高网络服务质量(QoS)^[1]提出了严峻的挑战. 理想的服务策略应当对到达信元提供尽量低的时延, 以满足网络中高速通信流量的需求.

终端用户在接受服务结束以后, 由于种种原因, 例如传输过程中部分数据包的丢失或错误, 会以一定的反馈概率对中央处理器(Central processing unit)资源重新提出服务请求. 对此可以用最简单的 M/M/1 反馈排队来建模, 即假定: 信元(需求数据)以泊松分布到达 CPU, 中央处理器以指数分布的时延提供服务, 如果信息接收成功, 则需求被满足, 否则以一定的概率通过 USB 接口调入数据后重新经 CPU 提供服务, 这是网络流中常见的信息处理过程.

对待反馈的服务请求, 通常传统的做法是服务器按照先进先出(FIFO)的服务规则响应顾客请求, 即将反馈的顾客置于队长的尾部, 已有部分文章^[2,3]对上述反馈策略进行过建模分析. 我们在本文中提出一种新的策略, 即将反馈的客户请

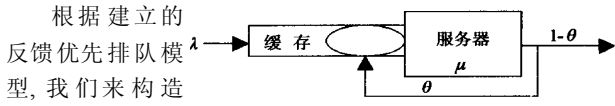
求赋予优先权. 如果有反馈的客户请求发生, 正在初次接受服务的信元立即停止服务, 直至反馈的客户请求服务结束后再继续接受服务; 若无反馈发生, 缓冲区中非反馈信元按照到达先后次序接受服务. 这种服务规则被称为反馈优先服务策略. 我们通过对反馈优先网络流建立排队模型, 构造了拟生灭(QBD)过程^[4], 并用谱展开方法^[5]得到了系统的平稳队长分布. 在此基础上, 计算了模型的平均队长、平均等待和逗留时间, 进而通过与传统反馈策略排队指标进行比较, 表明反馈优先的服务策略能提高网络服务质量. 相信本文结果将对网络通信具有一定的指导意义.

本文共分四部分. 第二节对反馈优先策略建立 QBD 过程并用谱展开方法得到队长的平稳分布. 在第三节中首先求解了模型的时延指标, 进而对两种服务策略的优劣进行了讨论. 最后给出了相应的数值模拟结果.

2 模型介绍

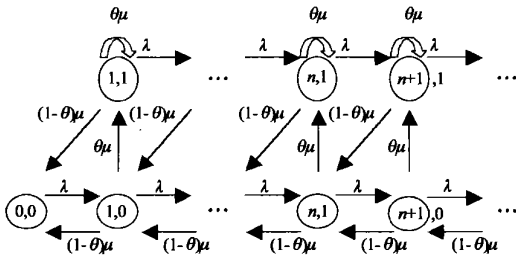
信元按照率为 $\lambda > 0$ 的泊松过程到达缓存, 服务器按照 FIFO 的排队规则逐个处理初次到达缓冲区的客户请求. 我们假设缓冲区的容量是无限的, 单个客户请求需要的服务时间是参数为 $\mu > 0$ 的指数分布. 当信元接受完服务以后, 或以概

率(1-H)立刻离开系统,或以反馈概率 $0 < H < 1$ 重新提出服务请求.我们规定反馈者优先的协议,即反馈的请求将抢占正在服务的首次到达系统的请求并使其进入缓存.具体的排队模型示意图如下:



根据建立的反馈优先排队模型,我们来构造 QBD 过程.令 N

(t) E_0 表示 t E_0 时刻缓存中顾客的个数加上正在服务的顾客,即队长.令 $S(t) = 0, 1$ 表示 t 时刻服务器所处的状态,其中 $S(t) = 0$ 表示 t 时刻服务器正在对首次接受服务的信元服务,这时系统内没有反馈的请求; $S(t) = 1$ 表示 t 时刻反馈的请求正在接受服务,这时系统中有一个反馈的请求, $(N(t) - 1)$ 个未接受过服务的信元.根据指数分布的无记忆性,显然 $\{N(t), S(t); t \in E_0\}$ 构成一个二维马氏链(Markov chain).将其状态坐标按照字典序排列,可得到该二维马氏链的状态转移图如下.



由转移图可以发现该二维马氏链是一个 QBD 过程^[4],其转移率矩阵可表示如下:

$$Q = \begin{bmatrix} B_0 & A_0 & & & \\ C_0 & B & A & & \\ & C & B & w & \\ & & w & w & \end{bmatrix}$$

其中

$$A_0 = [K \ 0] \quad B_0 = -K \quad C_0 = \begin{bmatrix} (1-H)L \\ (1-H)L \end{bmatrix} \quad A = \begin{bmatrix} K & 0 \\ 0 & K \end{bmatrix}$$

$$B = \begin{bmatrix} -(K+L) & HL \\ 0 & -(K+L)+HL \end{bmatrix} \quad C = \begin{bmatrix} (1-H)L & 0 \\ (1-H)L & 0 \end{bmatrix}$$

根据转移率矩阵 Q,我们可以得到以下两个定理:

定理 1 反馈优先排队系统平稳的充要条件是 $K < (1-H)L$.

证明 充分条件可根据漂移性条件^[3]推出.首先求解矩阵方程 $x(A+B+C) = 0$,得到向量 $x = c((1-H), H)$, c 为任意正实数.进而由漂移不等式 $xAe < xCe$ 可以算出 $K < (1-H)L$,其中 $e = (1, 1)c$.在平稳条件下, $[0, t]$ 时间段内系统平均输入 Kt 应小于平均输出 $L(1-H)t$,于是必要条件显然.

注:这里得到的平稳条件与众所周知的反馈非优先的模型结论相同.

定理 2 在平稳条件下,反馈优先排队的队长分布的衰减率为 $K/(1-H)L$,其平稳分布是

$$P_n = \frac{P_0}{(K+HL)} \left[\left(\frac{K}{K+L} \right)^n \#(HL, -HL) + \left(\frac{K}{(1-H)L} \right)^n \#(K, HL) \right],$$

$n = 1, 2, \dots$

而

$$P_0 = [(1-H)L - K] / [(1-H)L]$$

证明 我们利用谱展开方法^[5]来证明本定理.

定义矩阵 $Q(z) = A + Bz + Cz^2$,其行列式记为 $\det(Q(z))$.解方程 $\det(Q(z)) = 0$,我们得到两个满足 $0 < z < 1$ 的实根

$$z_1 = \frac{K}{K+L}, z_2 = \frac{K}{(1-H)L}$$

由方程 $u_j Q(z_j) = 0$ 得相应的特征向量

$$u_1 = (1, -1), \quad u_2 = (K, HL)$$

根据谱展开方法,队长的衰减率为 $\max(z_1, z_2)$,即 $K/(1-H)L$.这时

$$P_n = k_1 \left(\frac{K}{K+L} \right)^n \cdot \#(1, -1) + k_2 \left(\frac{K}{(1-H)L} \right)^n \cdot \#(K, HL), \quad n = 1, 2, \dots$$

通过求解联立方程 $P_0 B_0 + P_1 C_0 = 0$ 和 $P_0 A_0 + P_1 B_0 + P_2 C = 0$,可以得到 $k_2 = \frac{P_0}{K+HL} \# \frac{K}{(1-H)L}$, $k_1 = \frac{HL^2(1-H)}{K+L} k_2 = \frac{P_0}{K+HL}$

$\# \frac{K}{K+L} \# HL$.再对 $\{P_n\}_{n \in E_0}$ 进行归一化,即 $P_0 + \sum_{n=1}^{\infty} P_n e = 1$,得到

$$P_0 = [(1-H)L - K] / [(1-H)L]$$

从而我们推导出了队长的平稳分布.

3 最优策略比较

在系统平稳的条件下,我们来比较反馈优先策略与传统的反馈策略的优劣.根据第二节的结果,可以推导出以下结果.

定理 3 在平稳条件下,反馈优先模型的排队性能指标是:

$$\text{平均队长 } L = \sum_{n=1}^{\infty} n P_n e = \frac{K}{(1-H)L - K}$$

信元从开始接受服务到离开系统的广义平均服务时间

$$T = \sum_{i=1}^{\infty} i (1-H) H^i \# L^{-1} = \frac{1}{(1-H)L}$$

信元的平均等待时间(到达至开始接受服务) $W = L \# T$

信元的平均逗留时间(到达至离开系统)

$$S = (L+1) \# T = \frac{1}{(1-H)L - K}$$

当 $H=0$ 时上述结果与 M/M/1 排队结果相同.

当反馈的客户请求并不具有优先权时,它的广义平均服务时间比较复杂.我们经过仔细分析得到下面的定理.

定理 4 在平稳条件下,传统的反馈系统的排队性能指标是:

$$\text{平均队长 } L = L$$

$$\text{信元的平均等待时间 } W = L \# L^{-1}$$

信元从开始接受服务到离开系统的广义平均服务时间

$$T = T + (1-H) H Q \sum_{i=0}^{\infty} [\#(Q+H)]^i L^{-1}$$

$$= T\# \frac{1- H^2- 2Qf^2+ Qf^3}{1- H^2- Qf}$$

其中 $Q= K/L$ 表示在单个信元的服务时间内到达的平均信元数.

信元的平均逗留时间(到达至离开系统) $S= (L+ 1)\#T$

证明 由于优先权的存在并不改变队长的分布, 从而对于传统的反馈策略, 其平均队长与反馈优先模型的平均队长相同. 现在反馈信元没有优先权, 所以到达信元的平均等待时间只需要等待系统中存在的平均信元服务完毕, 每个服务时间为 L^{-1} , 即可开始服务. 前两个指标得证.

困难是证明第三个指标. 因为这时反馈信元与新到达系统的信元一样按照 FIFO 的排队规则逐个接受服务, 对这种传统的反馈策略该反馈(标记)信元从开始接受服务到离开系统的广义平均服务时间可分为两部分: 第一部分是标记信元按反馈优先策略的广义平均服务时间; 第二部分是标记信元在等待服务期间因新增信元所需要的服务时间. 容易知道, 标记信元反馈的次数服从几何分布: $(1- H)H^i, i= 0, 1, 2, \dots$, 在标记信元等待服务期间, 第 1 次新增加的信元数目为 Q , 第 2 次新增加的信元数目为 $Q(1+ (Q+ H))$; 第 i 次新增加的信元数目为 $Q(1+ (Q+ H)^i)$, 依次类推. 而每个新增信元的平均服务时间都是 L^{-1} , 于是, 第三个指标得证. 第四个指标显然.

由上述两个定理可以看出, 优先权的存在对信元的时延产生了巨大的影响. 与传统的反馈策略相比较, 虽然在反馈优先策略下信元的平均等待时间较长, 但由于 $2Qf^2- Qf^3 < Qf$, 其广义平均服务时间则较短, 故平均逗留时间也较短. 通常比较两种策略的优劣是根据信元从进入到最终离开系统的早晚, 即平均逗留时间长短来确定的, 因此我们提出的反馈优先策略优于传统的反馈协议.

4 数值模拟

本节对两种反馈策略的排队指标进行了数值模拟, 用图表的方式展现了第三节的理论结果. 通过改变反馈概率, 我们可以观测到各个平均排队指标的变化以及两种协议下策略优劣程度的差别.

图 1 反映了等待时间随反馈概率变化的情况, 图 2 反映了逗留时间随反馈概率变化的情况. 其中 $K= 0.4, L= 0.18$,

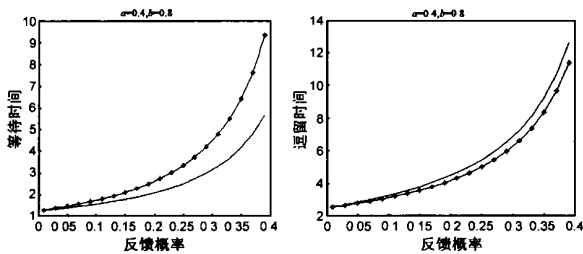


图 1

图 2

图中 m 表示为反馈优先策略的结果. 图 3 给出了 $K= 2, L= 8$ 的时延数值模拟结果, 图 4 给出了队长的模拟结果.

由图我们可以看出, 不同的策略协议将对网络流 QoS 产生较大的差别. 从满足客户请求最终时延的角度出发, 反馈优先策略比传统的反馈协议具有一定的优越性.

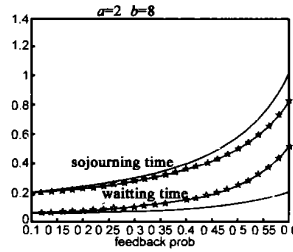


图 3

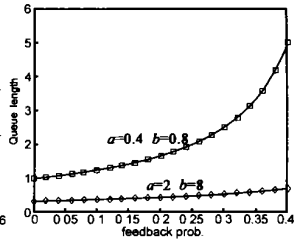


图 4

参考文献:

- [1] 黄文胜, 等. 提高 IP 网络服务质量的策略[J]. 信息工程大学学报, 2000, 1(4): 5- 8.
- [2] K Laevens, H Bruneel. Discrete-time queueing models with feedback for input-buffered ATM switches [J]. Performance Evaluation, 1996, 28 (4): 71- 87.
- [3] F Ishizaki, et al. Online sensitivity analysis of feedback controlled queueing systems with respect to buffer capacity [A]. Proc. of IEEE CDC 99[C]. USA, 1999. 40- 50.
- [4] M F Neuts. Matrix-geometric Solutions in Stochastic Models [M]. Baltimore: Johns Hopkins Univ. Press, 1981.
- [5] R Chakka. Performance and Reliability Modeling of Computing System Using Spectral Expansion [D]. Tyne: University of Newcastle, 1995.

作者简介:



黄月芳 女, 1975 年生于河南洛阳, 上海大学理学院在读博士, 研究方向: 自相似网络, 随机模型. Email: HYF-sh@sohu.com.



史定华 男, 1941 年生于江西, 教授, 博士生导师. 主要研究方向: 系统理论、随机模型、智能算法、生物信息等, 至今在国内外学术期刊已发表论文 80 余篇, 其中有 30 余篇被国际三大检索收录, 出版学术著作四部, 校对翻译著作一部, 现任中国运筹学会下可靠性学会副理事长, 排队论专业委员会副主任. Email: shidh2001@263.net.