Queues with Skill Based Routing and Resource Pooling under First Come First Served / Assign Longest Idle Server Discipline

EURANDOM 2nd Dutch-Israeli Conference, September 2010

Gideon Weiss University of Haifa joint work with and work by: Ivo Adan, Rene Caldentey, Cor Hurkens, Ed Kaplan, Rishi Talreja, Jeremy Visschers, Damon Wischik, Ward Whitt

Skill based routing and resource pooling:



Motivating Question 1: Infinite matching Kaplan '84,

In assigning housing to entitled families several housing projects are available (Kaplan, '84.'88).

"... Households applying for public housing are allowed to specify those housing projects in which they are willing to live; when a public housing unit becomes newly available, of those households willing to live in the associated housing project, the one that has been waiting the longest is offered the unit...."

FCFS discipline by law

OR question, game theoretic

How many choices should you mark?

- Mark only 1st choice: long wait
- Mark several: shorter wait but worse allotment

if everyone marks many choices: long wait and bad allotment

Q: How many profile i applicants go to project j?

Motivating Question 2: Overloaded system (Talreja Whitt)

Skill based routing in a call center type system: Overloaded systems, arrival rates $\Sigma \lambda_i > \Sigma \mu_j$ Stabilized by abandonment general arrival streams, general service times, patience distributions F_i



OR Question: How does it perform ?

How should we assign skills and staff them ?

Q: How many type i customer get served by server of type j?

Gideon Weiss, University of Haifa, FCFS Multi-Type, ©2010

Matching rates

When you use skill based routing, How many type c customers get served by type j servers?

EASY ?





Gideon Weiss, University of Haifa, FCFS Multi-Type, ©2010

Overloaded system with abandonment (Talreja Whitt, ManSci'08)

Call Centers

Overloaded $\lambda = \Sigma \lambda_i > \Sigma \mu_{i=} \mu$

Abandonments: patience Fi

FCFS

Fluid Arguments:

Uniform acceleration - many server scaling: servers are busy almost all the time successive customers wait almost same time on fluid scale we may get GLOBAL FCFS - everyone with patience <W abandons - everyone with patience >W waits exactly W.

$$\sum \lambda_i (1 - F_i(W)) = \sum \mu_j$$



6

Matching rates

When you use skill based routing,

How many type c customers get served by type j servers ?

NOT EASY !!



	β_1	β_2	β_3
α_a	Х		
α_b		Х	
α_c			Х

5

Gideon Weiss, University of Haifa, FCFS Multi-Type, ©2010

Service rates in overloaded system

 $\textbf{Question:}~\mbox{find}~\nu_{ij}~\mbox{the rate of services of i by } j$

Solve for W from: $\sum \lambda_i (1 - F_i(W)) = \sum \mu_i$

Then solve
$$\begin{split} \sum_{i \in C(j)} v_{ij} &= \mu_j, \qquad j = 1, \dots, J \\ \sum_{j \in S(i)} v_{ij} &= \lambda_i \overline{F_i}(W), \quad i = 1, \dots, I \end{split}$$

These may have no positive solution, unique solution or many solutions

- Can you solve it ?
- Does the stochastic system converge ?

Solved for trees, complete graphs, hybrids



The Housing problem: Infinite matching (Caldentey Kaplan '02, Caldentey, Kaplan, W '09)

This is not a customer server situation ---Housing units leave with the households and don't come back for many years



Can we work with this ?

Define: X_n the Markov chain of: ordered list of unmatched customers state space is words in alphabet of i = 1, ..., I

Conjecture: The Markov chain X is ergodic iff for all non-trivial S,C:

 $\alpha(C) < \beta(S(C)), \qquad \beta(S) < \alpha(C(S)), \qquad **$

Theorem: If the Markov chain is ergodic then $r_{i,j}^n \xrightarrow{n \to \infty a.s.} r_{i,j}$

This chain is intractable, we could only do some special cases

Example 1 Almost complete graph --- each node connected to all but one

$$\hat{\gamma}_{i,j} = \frac{\alpha_i \beta_j \left[(1 - \alpha_i)(1 - \beta_j) - \alpha_j \beta_i \right]}{(1 - \alpha_i - \beta_i)(1 - \alpha_j - \beta_j)} \left/ \left(1 + \sum_{i=1}^{I} \frac{\alpha_i \beta_i}{1 - \alpha_i - \beta_i} \right) \right|$$

FCFS Infinite matching - Markovian description

Fraction i->j matches in first n : r_{i}^{n}

 $r_{i,j}^n \xrightarrow{a.s.} r_{i,j}$

С

 α_3

Markovian description

Data: i.i.d. α.β

graph G

FCFS

Consider then s^{n+1} he will first look at those c left by $s^1 \dots s^n$ If he finds a match he will take the first, and leave one less unmatched behind. Else he will look for a match among those not considered previously, until he finds a match, adding a geometric number of c from $C(s^{n+1})$.

Gideon Weiss, University of Haifa, FCFS Multi-Type, ©2010



Gideon Weiss, University of Haifa, FCFS Multi-Type, ©2010

9

A Lyapunov function of Fayolle, Malyshev and Menshikov.

To show positive recurrence we need to show that $E(\Delta f(x,y)) \le -h < 0$





Back to Queueing - Manufacturing system `N' model

Adan, Foley, and Mcdonald, '08 `N' system queueing model:

- Memoryless,
- FCFS.

If system is empty, and type a arrives, machine 1 will take him w.p η

Typical state: Machine m1 serves first customer X type b customers behind it Machine m2 serves next type a customer y customers behind second machine



Results: steady state exact asymptotics, for large (x,y)

$$\Pi(x,y) \sim c(x \mid y) \left(\frac{\lambda_a + \lambda_b}{\mu_1 + \mu_2}\right)^y \left(\frac{\lambda_b}{\mu_1}\right)$$

Gideon Weiss, University of Haifa, FCFS Multi-Type, ©2010



13

A reversible multi-type loss system (Adan, Hurkens, W '10)

Loss system:

State: set of idle machines $\{M_1, M_2, \dots, M_k\}$

 $\begin{array}{c} C \\ S \end{array} \begin{array}{c} \hline b \\ \hline a \\ \hline \end{array} \begin{array}{c} c \\ \hline a \\ \hline \end{array} \begin{array}{c} c \\ \hline \end{array} \end{array} \begin{array}{c} c \\ \hline \end{array} \begin{array}{c} c \\ \hline \end{array} \end{array} \begin{array}{c} c \\ \hline \end{array} \end{array} \begin{array}{c} c \\ \end{array} \end{array}$

out: $\{M_1, M_2, \dots, M_{k-1}\}$ --> $\{M_1, M_2, \dots, M_k\}$ rate μ_{Mk}

Assume MC reversible. Detailed balance:

In: $\{M_1, M_2, ..., M_k\} \longrightarrow \{M_1, M_2, ..., M_{k-1}\}$ rate $\lambda_{M_k}\{M_1, ..., M_k\}$

 $\lambda_{M_k} \{M_1, \dots, M_k\} = \sum_{i \in C(M_k)} \lambda_i P(i \text{ chooses server } M_k \mid \{M_1, \dots, M_k\} \text{ idle})$

$$\pi\{M_1, \dots, M_{k-1}\} \mu_{M_k} = \pi\{M_1, \dots, M_k\} \lambda_{M_k}\{M_1, \dots, M_k\}$$

Solution:

$$\pi\{M_1, \dots, M_k\} = \pi\{\emptyset\} \frac{\mu_{M_1} \ \mu_{M_2} \ \cdots \ \mu_{M_k}}{\lambda_{M_1}\{M_1\}\lambda_{M_2}\{M_1, M_2\}\cdots\lambda_{M_k}\{M_1, \dots, M_k\}}$$

Gideon Weiss, University of Haifa, FCFS Multi-Type, ©2010

Skill based queues - FCFS (Visschers '00, Visschers, Adan, W '10)

Service exponential μ_j Graph G Arrivals Poisson λ_i Poisson FCFS service, random assignment to idle servers. State $\mathbf{s} = (M_1, n_1, \cdots, M_k, n_k)$



Thm: Use loss system assignment, then: MC obeys partial balance, product form:

$$\pi_{Queue}(\mathbf{s}) = \frac{\lambda_{M_1} \{M_1 \cdots M_K\} \lambda_{M_2} \{M_2 \cdots M_K\} \cdots \lambda_{M_k} \{M_{k+1} \cdots M_K\}}{\mu_{M_1} (\mu_{M_1} + \mu_{M_2}) \cdots (\mu_{M_1} + \cdots + \mu_{M_k})} \rho_1^{n_1} \cdots \rho_k^{n_k}}$$
$$\rho_l = \frac{\lambda_{U\{M_1, \dots, M_l\}}}{\mu_{M_1} + \cdots + \mu_{M_l}}$$

16

Assignment condition

Solution:

$$\pi_{Loss}\{M_1,\ldots,M_k\} = \pi\{\emptyset\} \frac{\mu_{M_1} \ \mu_{M_2} \ \cdots \ \mu_{M_k}}{\lambda_{M_1}\{M_1\}\lambda_{M_2}\{M_1,M_2\}\cdots\lambda_{M_k}\{M_1,\ldots,M_k\}}$$

Makes sense only if for all permutations:

$$\lambda_{M_1}\{M_1\}\lambda_{M_2}\{M_1,M_2\}\cdots\lambda_{M_k}\{M_1,\dots,M_k\} = \lambda_{\bar{M}_1}\{\bar{M}_1\}\lambda_{\bar{M}_2}\{\bar{M}_1,\bar{M}_2\}\cdots\lambda_{\bar{M}_k}\{\bar{M}_1,\dots,\bar{M}_k\}$$

This uniquely determines

$$\lambda_M(S) = \sum_{i \in C(S)} \lambda_i \left/ \left(1 + \sum_{j \in S} \frac{\lambda_j(S \setminus M)}{\lambda_M(S \setminus j)} \right) \right.$$

Question: Can you find P(i, M | S)



Waiting time in queue

Using distributional form of Little's law we get the waiting time from arrival till start of service for a customer of type ${\bf c}\,$:



 W_c is a mixture of sums of exponentials, each equal to an M/M/1 waiting time, as if customer is going through a tandem queue till he finds a server,

$$W_{c} \sim \pi_{Queue} \{ M_{1}, \cdot, M_{2}, \cdot, \dots, M_{k}, \cdot \} \exp(\lambda_{U\{M_{1}\cdots M_{k}\}} - \mu_{\{M_{1}\cdots M_{k}\}}) * \\ * \exp(\lambda_{U\{M_{1}\cdots M_{k-1}\}} - \mu_{\{M_{1}\cdots M_{k-1}\}}) * \dots$$

Gideon Weiss, University of Haifa, FCFS Multi-Type, ©2010

Back to FCFS ∞ – Matching : A new Markov chain

Server s^n sees the last match of customers to each type of server, and counts unmatched customers inbetween them:



Symmetry lost ? Symmetry regained:

The Matching model is symmetric between customers and servers. The Loss and Queueing models are not !

Consider now the loss system, replace random assignment with ALIS: Assign arrivals to Longest Idle Servers

New system is not reversible, but is tractable:

State: (M_1, M_2, \ldots, M_k) idle servers, ordered, M_k has been idle longest

$$\pi_{Loss}(M_1,\ldots,M_k) = B \frac{\mu_{M_1} \ \mu_{M_2} \ \cdots \ \mu_{M_k}}{\lambda_{C\{M_1,M_2\cdots M_k\}} \cdots \lambda_{C\{M_{k-1},M_k\}} \lambda_{C\{M_k\}}}$$

Adding over permutations we get the same as for random assignment

$$\pi_{Loss}\{M_1, \dots, M_k\} = \pi\{\emptyset\} \frac{\mu_{M_1} \ \mu_{M_2} \ \cdots \ \mu_{M_k}}{\lambda_{M_1}\{M_1, \dots, M_k\} \cdots \lambda_{M_{k-1}}\{M_{k-1}, M_k\}\lambda_{M_k}\{M_k\}}$$

Gideon Weiss, University of Haifa, FCFS Multi-Type, ©2010

I tried to get expressions for special cases - just as horrible I programmed it and calculated for I=J=7 SUSPICION: These quantities are #P-hard to evaluate Gideon Weiss, University of Halfa, FCFS Multi-Type, ©2010



Gideon Weiss, University of Haifa, FCFS Multi-Type, ©2010

21

22

Resource pooling

Stability of Queue:

$$\lambda(C) < \mu(S(C)) \text{ for all } C$$
$$\lambda = \sum \lambda_i, \quad \mu = \sum \mu_j, \quad \rho = \frac{\lambda}{\mu}, \quad \alpha_i = \frac{\lambda_i}{\lambda}, \quad \beta_j = \frac{\mu_j}{\mu}$$

Queue is stable if $\,\rho\,$ is small enough

complete resource pooling



Overloaded system:

 $\label{eq:conjecture: under complete resource pooling, when $\rho > 1$ all the servers stay together, just the last queue blows up:$

$$\begin{split} &\lim_{t \to \infty} P(M_1, n_1, \cdots, n_{K-1}, M_K) = \pi_{Match} \left(M_1, n_1, \cdots, n_{K-1}, M_K \right) \\ &\lim_{t \to \infty} P(n_K < 100000) = 0 \end{split}$$

Overloaded system with abandonments:

Solve for global waiting time $\sum \lambda_i (1 - F_i(W)) = \sum \mu_i = \mu$

$$\alpha_i = \frac{\lambda_i (1 - F(W))}{\mu}, \quad \beta_j = \frac{\mu_j}{\mu}$$

Conjecture: if overloaded system with abandonment has complete resource pooling then under uniform acceleration scaled state is distributed like

$$\pi_{Match}(M_1, n_1, \cdots, n_{K-1}, M_K)$$

Gideon Weiss, University of Haifa, FCFS Multi-Type, ©2010