

Performance Modeling of Computer Systems and Networks

Prof. Vittoria de Nitto Personè

The model for a service center:
analytical results

Università degli studi di Roma Tor Vergata
Department of Civil Engineering and Computer Science Engineering

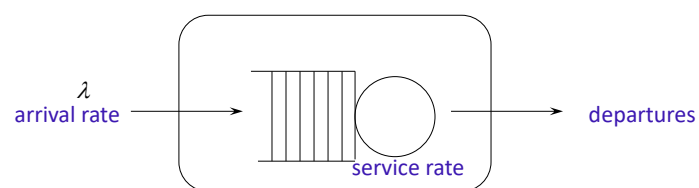
Copyright © Vittoria de Nitto Personè, 2021
<https://creativecommons.org/licenses/by-nc-nd/4.0/>



1

Analytical models

Server center



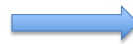
Little's law

$$E(T_s) = E(T_q) + E(S)$$

$$E(N_s) = \lambda E(T_s)$$

$$E(N_s) = E(N_q) + \rho$$

$$E(N_q) = \lambda E(T_q)$$



$$E(T_s) = \frac{E(N_s)}{\lambda}$$

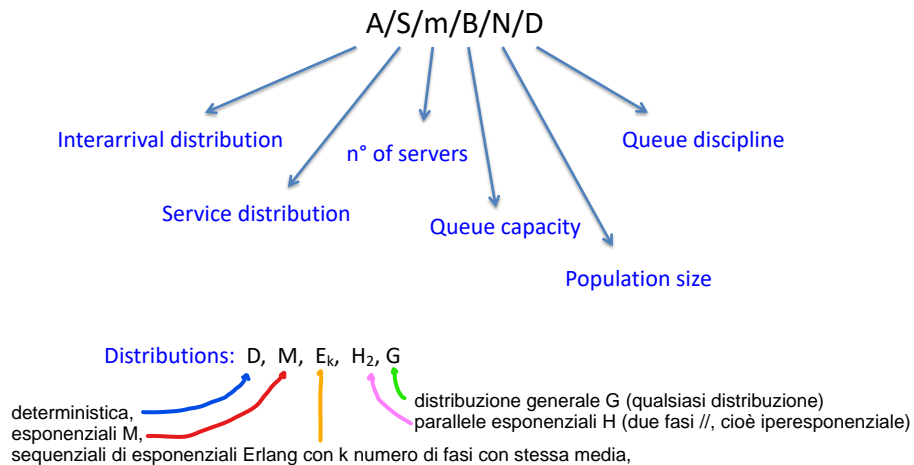
$$E(T_q) = \frac{E(N_q)}{\lambda}$$

Prof. Vittoria de Nitto Personè

2

2

The Kendall notation



Prof. Vittoria de Nitto Personè

3

3

Non-preemptive abstract scheduling

FIFO, LIFO-non-preemp, Random

It seems like

FIFO should have the best mean response time because jobs are serviced most closely to the time they arrive (rispetta ordine di arrivo)
 LIFO may make a job wait a very long time

all the above policies have **exactly the same mean response time.**

(hanno stessa media ma NON HANNO STESSA VARIANZA)

Prof. Vittoria de Nitto Personè

4

4

1930: The Khinchin Pollaczek equation (KP)

interrarivi esponenziali, distribuzione servente generale, 1 servente

M/G/1 abstract scheduling

$$E(N_Q) = \frac{\rho^2}{2(1-\rho)} \left[1 + \frac{\sigma^2(S)}{E(S)^2} \right]$$

$= C^2$
Squared coefficient of variation
Service time dispersion
 rapporto tra varianza e quadrato delle media

1. Any service time distribution
2. Poisson arrivals
3. Abstract discipline (FIFO, LIFO, RAND...)

Più il coefficiente cresce, più le prestazioni peggiorano!

Rho è la vera misura del carico, non lambda, il quale pesa pochissimo come carico, in quanto
 $\rho = \lambda \cdot E[S]$

Prof. Vittoria de Nitto Personè

5

5

TIPO M/G/1, SI ANALIZZA UN SOLO JOB, QUINDI QUESTI SONO I PERCORSI CHE PUO' SEGUIRE UN JOB. IL SERVENTE E' SEMPRE SINGOLO

Phase-type distributions

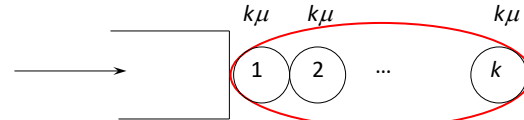
Exponential



unica fase di servizio
sto modellando tempo job con variabile
esponenziale

Servente singolo

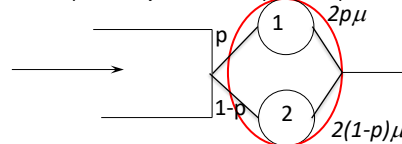
k-Erlang



E' SEMPRE UNICO SERVENTE,
CAMBIA COME MODELLO TEMPO
(comparo a parità di media)

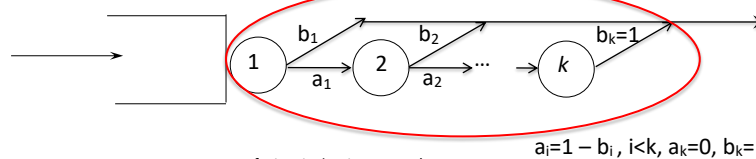
modello tempo unico job con k tempi tutti esponenziali e distribuiti allo stesso modo.

hyperexponential distribution



modello 1 tempo di servizio con 2 fasi
alternative esponenziali con due tassi diversi
media: $1/\mu$ (media del servizio uguale)

Cox distribution



$a_i = 1 - b_i, i < k, a_k = 0, b_k = 1$

Prof. Vittoria de Nitto Personè

che senso ha modellare
queste cose, se
mediamente sono uguali?

modello diverse
variabilità!!

6

serie di fasi esponenziali, nella erlang passo per tutti gli stati, qui il job potrebbe fare un solo stadio ed uscire, fare due stadi e uscire,...., farli tutti! Ma sempre un job c'è.

The Khinchin Pollaczek equation (KP)

$$E(N_Q) = \frac{\rho^2}{2(1-\rho)} [1 + C^2]$$

The mean queue population grows as C^2

$$D \longrightarrow C^2=0$$

$$E_k \longrightarrow C^2 = \frac{1}{k}, k \geq 1$$

$$M \longrightarrow C^2=1$$

$$H_2 \longrightarrow C^2 = g(p) = \frac{1}{2p(1-p)} - 1$$

variabilità cresce nel verso della freccia.

in funzione della probabilità,
dipende da p.
per p = 0.5, ottengo 1 come l'esponenziale.

$$p = 0.6 \quad C^2 = 1.08\bar{3}$$

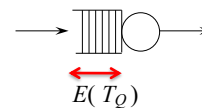
$$p = 0.7 \quad C^2 = 1.38095$$

$$p = 0.8 \quad C^2 = 2.125$$

$$p = 0.9 \quad C^2 = 4.\bar{5}$$

The Khinchin Pollaczek equation (KP)

M/G/1 abstract scheduling



$$E(T_Q) = \frac{E(N_Q)}{\lambda} = \frac{\rho^2}{\lambda 2(1-\rho)} [1 + C^2] = \frac{\rho}{1-\rho} \frac{C^2 + 1}{2} E(S)$$

$\rho = \lambda \cdot E[S]$

Normalmente il tempo di servizio deve proporzionale all'attesa, quindi se il coefficiente con C^2 fosse 30, non andrebbe tanto bene!

"mu" è lo stesso, quindi io confronto prestazioni supponendo stessa media, utilizzazione etc...
cambia solo la varianza. La media deve essere uguale, sennò confronto due cose diverse.

Analytical models
M/G/1

The Khinchin Pollaczek equation (KP)

$$g(\rho) = \frac{1}{2\rho(1-\rho)} - 1$$

$$E(N_Q) = \frac{\rho^2}{2(1-\rho)} [1 + C^2], \quad E(T_Q) = \frac{\rho}{1-\rho} \frac{C^2 + 1}{2} E(S)$$

Service time	$E(N_Q)$ popolazione media in coda	$E(T_Q)$ tempo medio attesa in coda
Deterministic, M/D/1	$\frac{\rho^2}{2(1-\rho)}$	$\frac{\rho E(S)}{2(1-\rho)}$
Markovian, M/M/1	$\frac{\rho^2}{1-\rho}$	$\frac{\rho E(S)}{1-\rho}$
K-Erlang, M/E _k /1 $\sigma^2(S) = \frac{E(S)^2}{k}$	$\frac{\rho^2}{2(1-\rho)} \left(1 + \frac{1}{k}\right)$	$\frac{\rho E(S)}{2(1-\rho)} \left(1 + \frac{1}{k}\right)$
Hyperexpo, M/H ₂ /1 $\sigma^2(S) = E(S)^2 g(\rho)$	$\frac{\rho^2}{2(1-\rho)} (1 + g(\rho))$	$\frac{\rho E(S)}{2(1-\rho)} (1 + g(\rho))$

Prof. Vittoria de Nitto Personè

9

9

Nella colonna della popolazione media in coda, viene già incluso il C^2 con valore associato al caso in questione. Ciò ci fa capire che abbiamo sempre una dipendenza da rho^2 e non da C^2.
Se un provider deve fronteggiare arrivi raddoppiati, devo considerare anche il tasso di servizio raddoppiato, perchè rho in questo caso raddoppiato deve comunque essere uguale al rho prima di questo raddoppio, sennò non ha senso confrontare. Poichè mu = 1/E[S] allora il tempo di servizio si dimezza etc...

Analytical models
M/G/1

Service time Sensitivity

sensibilità rispetto a variazioni di tempi di servizio a parità di media (sennò non ha senso).

$$E(N_Q)_D \leq E(N_Q)_{E_k} \leq E(N_Q)_M \leq E(N_Q)_{H_2}$$

$$\sigma^2(N_Q)_D \leq \sigma^2(N_Q)_{E_k} \leq \sigma^2(N_Q)_M \leq \sigma^2(N_Q)_{H_2} \quad \text{non lo dimostriamo}$$

By considering $E(N_S) = E(N_Q) + \rho$, the same order holds for the variable N_S

By considering the Little's equation, the same order can be derived for the mean times $E(T_S)$ and $E(T_Q)$, but just for the 1° order moment, not for the variance

cioè in media, una distribuzione deterministica comporta tempo risposta medio minore dell'erlang etc..
ma ciò non si applica alle varianze tramite Little.

Prof. Vittoria de Nitto Personè

10

10

Discipline Sensitivity

non posso però basarmi su quello che i job chiedono, sarebbe size based
By definition, KP holds for any abstract service discipline, so

$$E(N_Q)_{\text{FIFO}} = E(N_Q)_{\text{LIFO}} = E(N_Q)_{\text{RAND}} = E(N_Q)_{\text{abstract}} \quad \text{lo scheduling non influenza KP}$$

$$\sigma^2(N_Q)_{\text{FIFO}} = \sigma^2(N_Q)_{\text{LIFO}} = \sigma^2(N_Q)_{\text{RAND}} = \sigma^2(N_Q)_{\text{abstract}}$$

By considering $E(N_S) = E(N_Q) + \rho$, the same equalities hold for the variable N_S

By considering the Little's equation, the same holds for $E(T_S)$ and $E(T_Q)$,

$$E(T_Q)_{\text{FIFO}} = E(T_Q)_{\text{LIFO}} = E(T_Q)_{\text{RAND}} = E(T_Q)_{\text{abstract}}$$

Is $\sigma^2(T_Q)$ the same for all these policies?

lo scheduling, in questo caso, conta!

Prof. Vittoria de Nitto Personè

11

11

Discipline Sensitivity

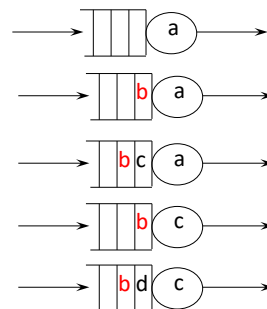
No! alcuni job potrebbero essere svantaggiati, incidendo su variabilità, MEDIA UGUALE.

LIFO can generate some extremely high response times because we have to wait for system to become empty to take care of that first arrival

$$\sigma^2(T_Q)_{\text{FIFO}} \leq \sigma^2(T_Q)_{\text{RAND}} \leq \sigma^2(T_Q)_{\text{LIFO}}$$

FIFO rispetta ordine arrivo, variabilità minima.
LIFO caso peggiore.

STIAMO PARLANDO DI VARIANZA,
NON DI MEDIA!



in questo esempio LIFO, "b" è svantaggiato.

Prof. Vittoria de Nitto Personè

12

12