

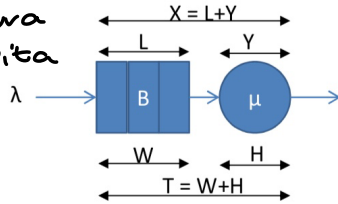
→ Sistemi di servizio

→ un'astrazione utile per studiare le prestazioni di molti sistemi reali.

è caratterizzato da un processo di arrivo, da una coda di attesa e da (uno o più) server.

$B=0$ → sistema a pura perdita

$B=\infty$ → sistema a pura attesa



λ = frequenza di arrivo	Y = # in servizio
$1/\mu$ = tempo medio di servizio	X = # nel sistema
B = dimensione del buffer	W = tempo speso in coda
N = numero di server	H = tempo di servizio
L = # in coda	T = tempo speso nel sistema

→ Formula di Little

→ vale per sotto condizioni molto generali riguardanti i sistemi di servizi a pura attesa:

sistema stabile, Δ
sistema work-conserving

→ il lavoro fatto dal server non si perde.

$$\bar{X} = \lambda \bar{T}$$

• Vale anche per i sottosistemi:

$$\bar{L} = \lambda \bar{W} \quad \text{e} \quad \bar{Y} = \lambda \bar{H} = \frac{\lambda}{\mu} = \rho$$

carico del sistema

→ Distribuzioni degli arrivi dei servizi

Supponiamo di osservare il flusso di auto su un'autostrada e la coda ad un casello. Se il flusso è di 10 auto al minuto e la coda è stabile, è lunga in media 100 auto.

Formula di Little?

• Il tempo medio $\bar{T} = \frac{100}{10} = 10 \text{ min}$
 λ = frequenza di arrivo

• $\mu > 10 \text{ auto/min}$

→ Modelli Markoviani

↳ la distribuzione dei tempi di interarrivo e di servizio è esponenziale

$$P\{t_i \leq t\} = 1 - e^{-\lambda t}$$

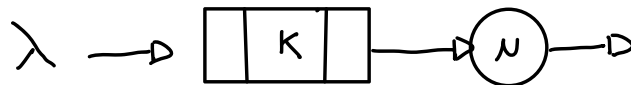
$$\overline{A(t)} = \lambda t$$

→ Sistemi senza memoria

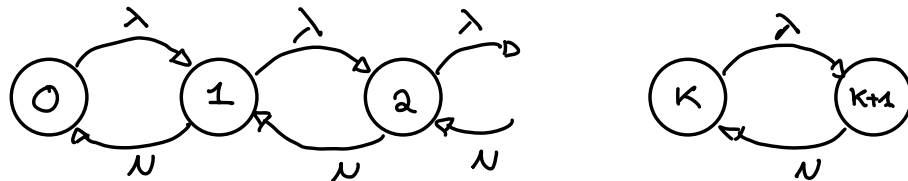
↳ il comportamento del sistema a partire da un momento in cui viene osservato, non dipende da quanto successo prima.

→ Coda M/M/1/K

- Una coda con un server, interarrivo e servizio esponenziali e buffer con K posizioni si chiama coda M/M/1/K



- Il sistema si risolve con una catena di Markov in cui ogni stato rappresenta il numero di elementi nel sistema;



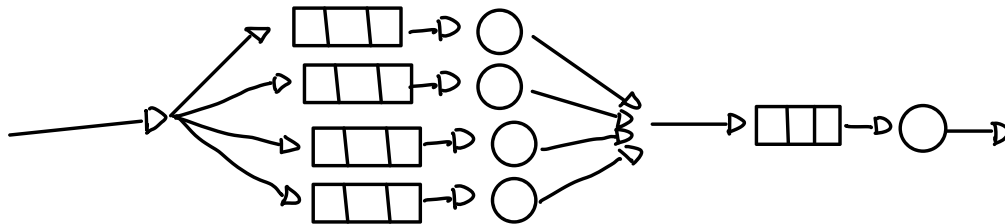
- Se P_i è la probabilità dello stato i-esimo

$$\begin{cases} \lambda P_0 = \mu P_1 \\ \lambda P_1 = \mu P_2 \\ \lambda P_K = \mu P_{K+1} \\ \sum_{i=0}^{K+1} P_i = 1 \end{cases}$$

$$P_0 = \frac{1-\rho}{1-\rho^{K+2}} \quad P_j = \frac{1-\rho}{1-\rho^{K+2}} \rho^j$$

Esempio

- Un sistema informativo Web-based sia costituita da quattro server di front-end operanti a divisione di carico e un server di backoffice che gestisce tutte le richieste. Si supponga che ciascun server di front-end elabori una richiesta in media 20 ms e il server 10 ms. Se al sistema arrivano 80 richieste al secondo e il tempo di risposta è di 95 ms. Quanti server di front-end stanno probabilmente funzionando.



- Abbiamo due sottosistemi in cascata e supponiamo che tutto il traffico offerto al primo venga offerto al primo venga offerto anche al secondo.

$$T = \frac{1}{(N-\lambda)} = \frac{1}{(100-80)} = 50 \text{ msec}$$

Otteniamo:

$$- 1 \text{ server } T = \frac{1}{(N-\lambda_{\text{server}})} = \frac{1}{(50-80)} = \text{congestione}$$

$$- 2 \text{ server } T = \frac{1}{(50-40)} = 100 \text{ ms}$$

$$- 3 \text{ server } T = \frac{1}{(50-26,6)} = 48 \text{ ms}$$

$$- 4 \text{ server } T = \frac{1}{(50-20)} = 33,3 \text{ ms}$$

La somma dei ritardi con 3 server è di 92,8 ms
(~ 95 ms)

-D Sistemi a perdita e principio PASTA

PASTA

(Poisson Arrivals
see Time Averages)

Gli arrivi «campionano» lo
stato del sistema e la
distribuzione di probabilità
dello stato

-D Se un sistema a coda
ha un buffer finito e il
processo degli arrivi per-
mette che giungano al
sistema più job di quanti
ne può contenere.

$$p_{k+1} = \frac{1-p}{1-p^{k+2}} p^{k+1}$$

-D Coda M/D/1

Una coda con un server,
tempo di intervallo esponenziale
e tempo di servizio deterministico

-D L'analisi del sistema è
matematicamente meno
semplice di quella relativa
ad una coda M/D/1

$$\bar{X} = \left(1 - \frac{\rho}{2}\right) \frac{\rho}{1-\rho}$$

$$\bar{T} = \left(1 - \frac{\rho}{2}\right) \frac{1}{\mu - \lambda}$$

$$\bar{L} = \frac{\rho^2}{2(1-\rho)}$$

$$\bar{Y} = \rho$$

$$\bar{W} = \frac{1}{\mu} \frac{\rho}{2(1-\rho)}$$

$$\bar{H} = \frac{1}{\mu}$$