

ISE 22/23 - Tema 5:

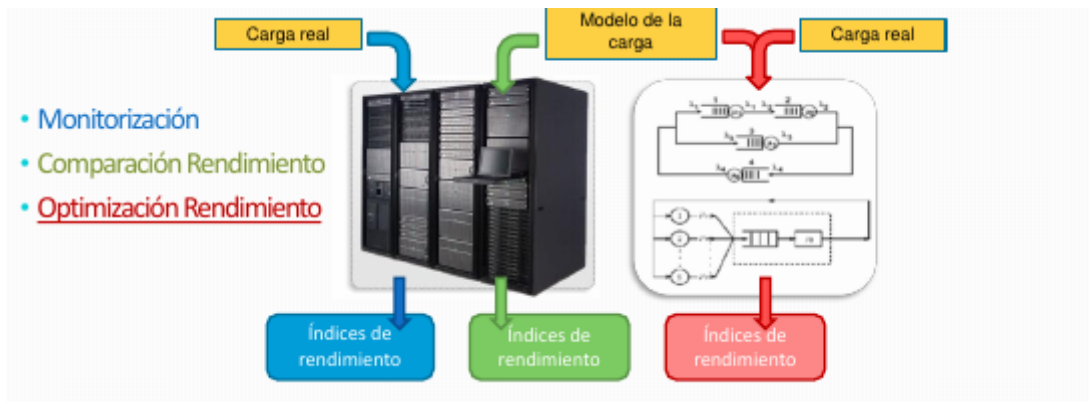
¿Cómo mejorar el rendimiento de mi servidor?

Objetivos:

- Proporcionar un modelo analítico de comportamiento de un sistema informático como punto de partida para obtener índices de rendimiento.
- Entender la importancia de los cuellos de botella como limitadores del rendimiento de los sistemas informáticos.
- Saber aplicar las leyes operacionales en ejemplos sencillos para obtener índices de rendimiento.
- Saber interpretar los límites optimistas del rendimiento que establece el análisis operacional.
- Saber evaluar de forma cuantitativa el efecto de diferentes terapias de mejora o estrategias de diseño sobre el rendimiento de un servidor.

5.1. Introducción - Redes de colas de espera

¿Cómo podemos mejorar el rendimiento de un servidor?



El modelo de un sistema informático

Abstracción del sistema informático real. Conjunto de dispositivos interrelacionados y trabajos que los usan (carga).

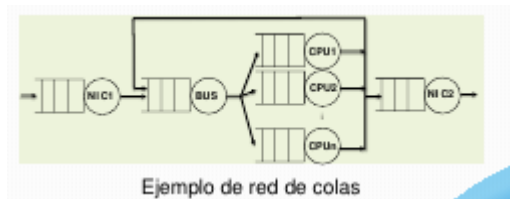
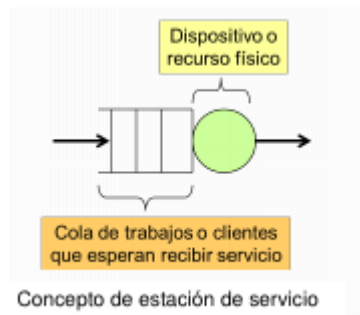
- Dispositivos (resources): núcleos lógicos, unidades de almacenamiento permanentes, tarjetas de red, etc.

- Trabajos (jobs): procesos, accesos, peticiones, etc.

Normalmente un dispositivo o recurso solo puede ser usado por un trabajo a la vez. El resto de trabajos tendrá que esperar.

Modelos basados en redes de colas (queueing networks):

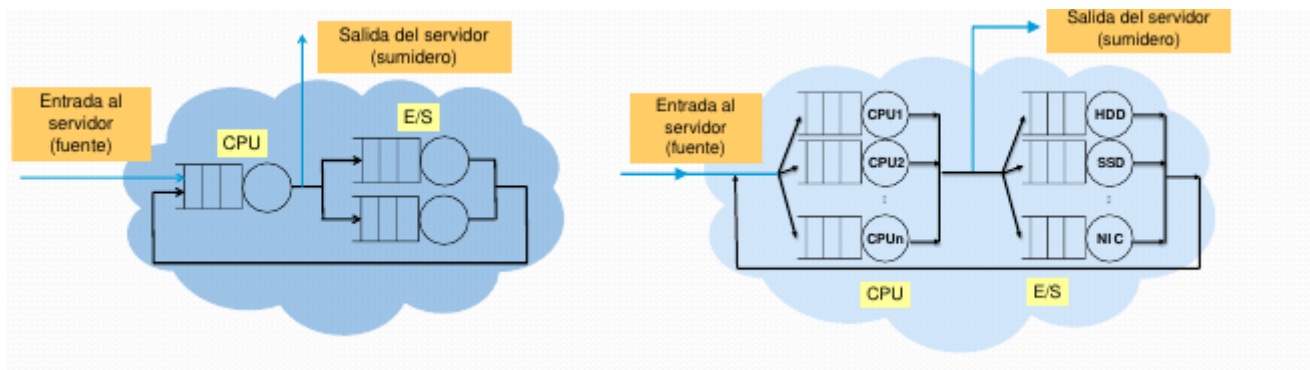
- Una red de colas está formada por un conjunto de estaciones de servicio conectadas entre sí.
 - Estación de servicio (service station): Objeto compuesto por un dispositivo (recurso físico) que presta un servicio y una cola de espera para los trabajos (clientes) que demandan un servicio de él.



El modelo de servidor central

Es la red de colas que más se ha utilizado para representar el comportamiento básico de los programas en un servidor de cara a extraer información sobre su rendimiento.

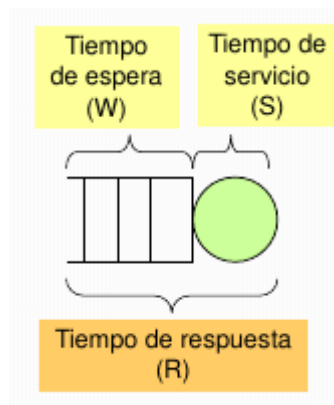
1. Un trabajo que “llega” al servidor comienza utilizando el procesador.
2. Después de “abandonar” el procesador, el trabajo puede:
 1. Terminar (sale del servidor), o bien
 2. Realizar un acceso a una unidad de entrada/salida (discos, red,...).
3. Después de una operación con una unidad de entrada/salida, el trabajo vuelve a “visitar” al procesador.



Algunas variables características a un trabajo, en una estación de servicio, en un instante concreto

- Tiempo de espera en cola (W, waiting time):
 - Tiempo transcurrido desde que el trabajo solicita hacer uso del recurso físico (se pone en la cola) hasta que realmente empieza a utilizarlo.
- Tiempo de servicio (S, service time):
 - Desde que el trabajo accede al recurso físico hasta que lo libera (tiempo que tarda el recurso físico en procesar el trabajo).
- Tiempo de respuesta (R, response time)
 - Suma de los dos tiempos anteriores.

$$R = W + S$$



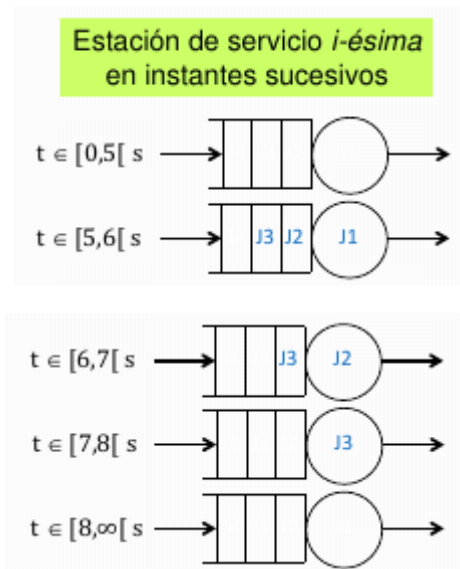
Recopilando estas medidas para múltiples trabajos, obtendremos distribuciones de probabilidad que caracterizan a esa estación de servicio.

Ejercicio

Suponga que la estación de servicio i -ésima de una red de colas tiene un tiempo de servicio constante $S_i = 1s$. Suponga que los trabajos (jobs) llegan con la siguiente distribución temporal:

- Durante los primeros 5 segundos no llega ningún trabajo.
- En $t=5s$ llegan 3 trabajos: J_1 , J_2 y J_3 (por ese orden).

Calcule los tiempos de espera en la cola y los tiempos de respuesta que experimentan cada uno de los tres trabajos. Calcule finalmente los valores medios de W y R .



$$W_i(J_1) = 0s$$

$$R_i(J_1) = 1s$$

$$W_i(J_2) = 1s$$

$$R_i(J_2) = 2s$$

$$W_i(J_3) = 2s$$

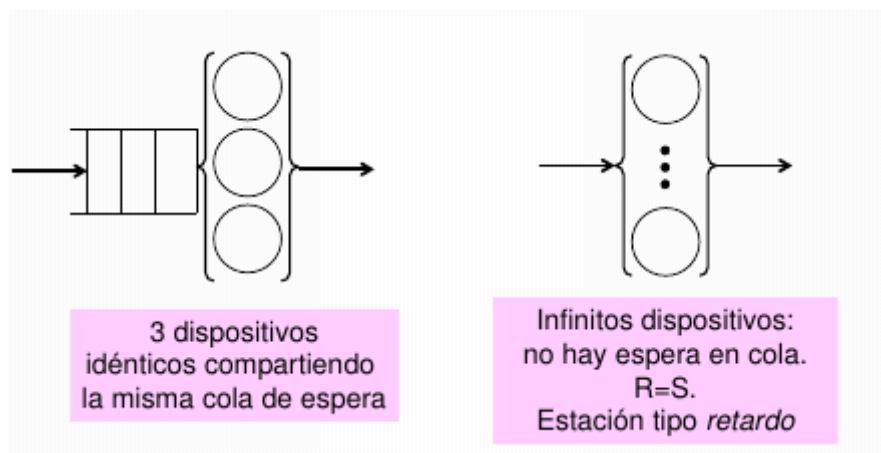
$$R_i(J_3) = 3s$$

$$\text{Valor medio } W_i = 1s$$

$$\text{Valor medio } R_i = 2s$$

Estaciones con más de un servidor

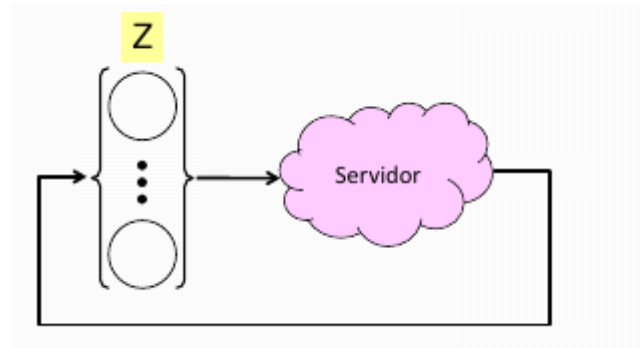
Son capaces de atender a más de un trabajo en paralelo:



El tiempo de reflexión Z (think time)

Es un parámetro (Z) que representa el tiempo que requiere el usuario antes de volver a lanzar una petición al servidor tras la respuesta de éste. Se suele modelar mediante una estación de servicio tipo retardo con un tiempo de servicio = Z .

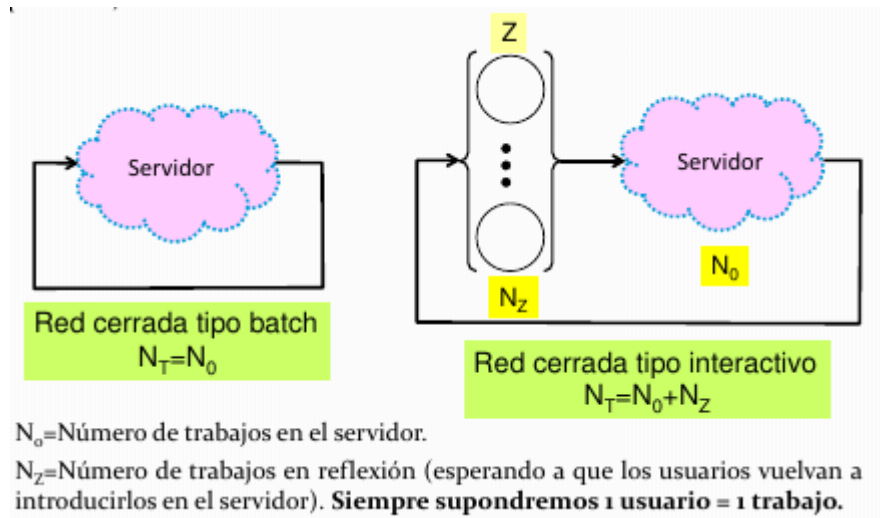
Para ello, realizamos una hipótesis adicional: cada usuario envía un único trabajo al servidor.



Redes de colas

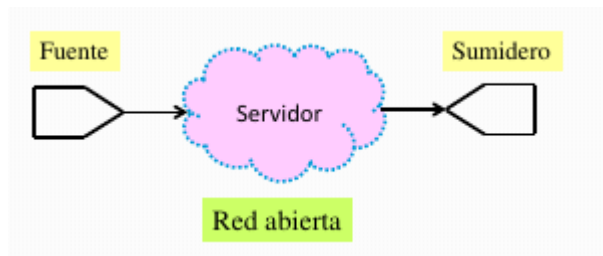
- Cerradas

- Presentan un número constante de trabajos que van recirculando por la red (NT). Dependiendo de si hay o no interacción con usuarios se distingue entre:
 - Redes de tipo batch (por lotes)
 - Redes interactivas

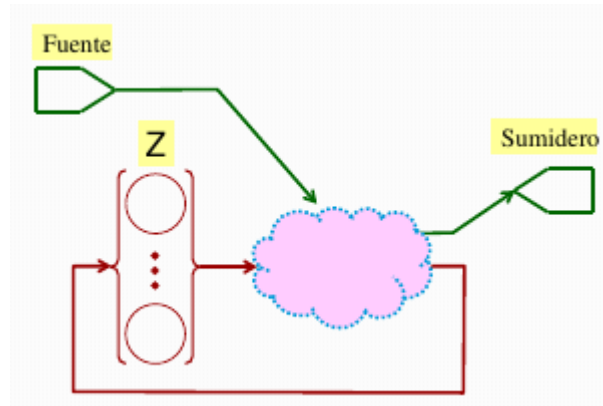


- Abiertas

- Los trabajos llegan a la red a través de una fuente externa que no controlamos. Tras ser procesados, salen de ella a través de uno o más sumideros. No existe realimentación entre sumidero y fuente.
- El número de trabajos en el servidor (N_0) puede variar con el tiempo.



- **Mixtas**
 - Cuando el modelo no corresponde a ninguno de los dos anteriores.



5.2. Variables y leyes operacionales

El análisis operacional

Técnica de análisis de redes de colas basada en valores medios de diferentes variables medibles (variables operacionales) del servidor.

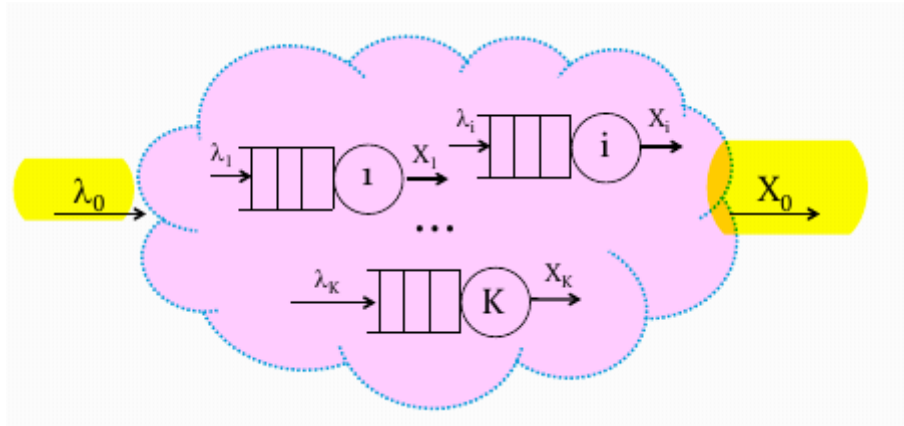
Tiene sentido usar valores medios porque son servidores, muchas veces ejecutamos el mismo proceso todo el rato. Si fueran computadores personales NO.

En general, se va a omitir el "medio" en las variables, es decir, quitar la barrita que indica que estamos hablando de media.

- Nos proporcionará relaciones generales entre las variables operacionales (leyes operacionales).
- Nos permitirá calcular las prestaciones del servidor para los casos de baja y alta carga por medio de cálculos muy sencillos.
- Nos permitirá evaluar los efectos en el rendimiento de diferentes modificaciones en los recursos del servidor.

Variables del servidor y de cada estación de servicio

- El servidor contiene K estaciones de servicio (recursos o dispositivos).
- A todo el servidor en su globalidad lo denotamos como dispositivo “cero”.

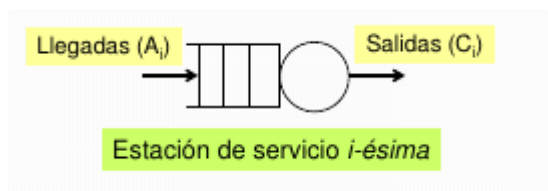


λ_0 -- peticiones por segundo

X_0 -- peticiones completadas por segundo

Variables operacionales básicas de cada estación de servicio

- Variable global temporal
 - T -- Duración del periodo de medida para el que se extrae el modelo.
- Variables operacionales básicas de la estación de servicio i-ésima, medidas durante el tiempo T:
 - A_i -- Número de trabajos solicitados a la estación (llegadas, ARRIVALS).
 - B_i -- Tiempo que el dispositivo ha estado en uso (=ocupado) (BUSY time).
 - C_i -- Número de trabajos completados por la estación (salidas, COMPLETIONS).

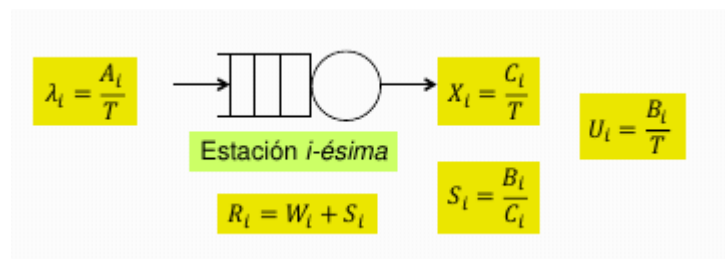


Variables operacionales deducidas

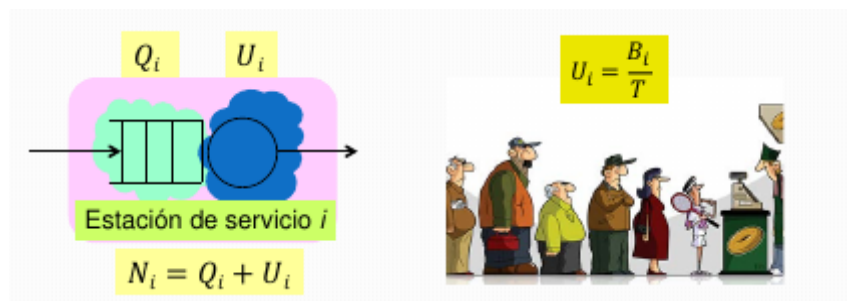
Se deben poder estimar a partir de las variables básicas:

- λ_i Tasa media de llegada (arrival rate) trabajos/segundo

- Arrivals/Tiempo -- A_i/T
- X_i Productividad media (throughput) trabajos/segundo
 - Completions/Tiempo -- C_i/T
- S_i Tiempo medio de servicio (service time) segundos [/trabajo]
 - Busy/Completions -- B_i/C_i
- W_i Tiempo medio de espera en cola (waiting time) segundos [/trabajo]
- R_i Tiempo medio de respuesta (response time) segundos [/trabajo]
- U_i Utilización media (utilization) sin unidades
 - Busy/Tiempo -- U_i/T
 - $U_i = X_i * S_i$



- N_i : Número medio de trabajos en la estación de servicio (cola más recurso).
- Q_i : Número medio de trabajos en cola de espera (jobs in queue). U_i : Número medio de trabajos siendo servidos por el dispositivo,
 - $U_i = N_i - Q_i$.
 - Coincide numéricamente con la utilización media = proporción de tiempo que el dispositivo ha estado en uso (busy) con respecto al intervalo total de medida (T) (como máximo 1 si $B_i=T$).



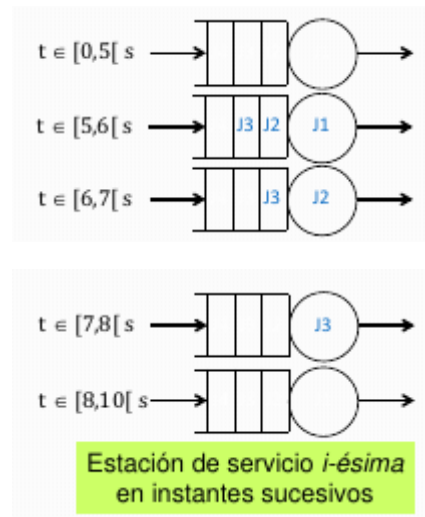
Ejercicio

Suponga que la estación de servicio i -ésima de una red de colas tiene un tiempo de servicio constante $S_i=1s$. Suponga que los trabajos llegan con la siguiente distribución temporal:

- Durante los primeros 5 segundos no llega ningún trabajo.

- En $t=5s$ llegan 3 trabajos: J_1 , J_2 y J_3 (por ese orden).

Para el intervalo de medida $[0, 10[s$, calcule A_i , B_i , C_i , λ_i , X_i , U_i , Q_i , N_i



A_i - Arrivals, trabajos en total: $A_i = 3$ trabajos

B_i - Busy, tiempo ocupado: $B_i = 3s$ (porque el tiempo de servicio S es siempre 1 en este caso)

C_i - Completions, trabajos finalizados: $C_i = 3$ trabajos

λ_i - Tasa media de llegada: $A_i/T = 3/10 = 0.3$ trabajos por segundo

X_i - Productividad media: $C_i/T = 3/10 = 0.3$ trabajos por segundo

U_i - Utilización media: $B_i/T = 3/10 = 0.3$

Q_i - Número medio de trabajos en cola de espera:

- Durante el intervalo 0-5 había 0 trabajos a la cola
- En el intervalo 5-6 habían 2 trabajos a la cola
- En el intervalo 6-7 había 1 trabajo a la cola
- En el intervalo 7-10 había 0 trabajos a la cola
 - $(0 \cdot 5s + 2 \cdot 1s + 1 \cdot 1s + 0 \cdot 3s) / 10s = 0,3$ trabajos

N_i - Número medio de trabajos en la estación de servicio (cola más recurso).

- Podemos calcularlo como la Q_i pero contando también el trabajo en ejecución o $N_i = Q_i + U_i = 0.6$ trabajos

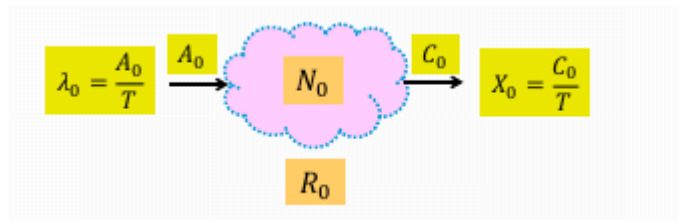
Variables operacionales de un servidor

Variables operacionales básicas de un servidor:

- A0 Número de trabajos solicitados al servidor (arrivals).
- C0 Número de trabajos completados por el servidor (completions).

Variables operacionales deducidas de un servidor:

- λ_0 Tasa media de llegada al servidor (arrival rate).
- X0 Productividad media del servidor (throughput).
- N0 Número medio de trabajos en el servidor (#jobs) = $N_1 + N_2 + \dots + N_K$.
- R0 Tiempo medio de respuesta del servidor (response time) \equiv tiempo que tarda, de media, el servidor en procesar una petición.



Razón de visita y demanda de servicio

- Razón media de visita V_i (VISIT ratio).
 - Representa la proporción entre el número de trabajos completados por el servidor y el número de trabajos completados por la estación de servicio i -ésima.
 - Nos indica el número de veces que, de media, un trabajo “visita” la estación de servicio i -ésima antes de abandonar el servidor.

$$V_i = \frac{C_i}{C_0}$$

- Demanda media de servicio D_i (service DEMAND).
 - **Cantidad de tiempo** que, por término medio, el dispositivo de la estación de servicio i -ésima le ha dedicado a cada trabajo que abandona el servidor (= que ha sido procesado por completo por el servidor).

$$D_i = \frac{B_i}{C_0} = V_i \times S_i$$

- Nótese que la demanda de servicio de una estación no tiene en cuenta la posible espera de un trabajo en su cola.

Ejercicio

Después de monitorizar el disco duro de un servidor web durante un periodo de 24 horas, se sabe que ha estado en uso (=ocupado) un total de 6 horas.

Asimismo, se han contabilizado durante ese periodo un total de 98590 peticiones de lectura/escritura al disco duro y un total de 98591 peticiones completadas. Se ha estimado que cada petición atendida por el servidor web ha requerido una media de 9,5 peticiones de lectura/escritura al disco duro. Calcule, para ese periodo de monitorización:

- a) La tasa media de llegada y la productividad media del disco duro.
- b) La utilización media del disco duro.
- c) El tiempo medio de servicio y la demanda media de servicio del disco duro.
- d) ¿Cuál es la productividad media del servidor web?

Nota: Todas las variables que se usan en este tema son valores medios por lo que, de aquí en adelante, normalmente no se indicará de forma explícita la palabra “medio” al referirnos a ellas.

- $T = 24h$
- $Bdd = 6h$
 - $dd ==$ disco duro
- $Add = 98590 \text{ tr}$
- $Cdd = 98591 \text{ tr}$
 - $tr ==$ trabajos, peticiones RW
 - $C > A$, es posible completar más trabajos de los que llegan porque puede que hubiera cosas pendientes antes de comenzar el monitoreo
- "Se ha estimado que cada petición atendida por el servidor web ha requerido una media de 9,5 peticiones de lectura/escritura al disco duro."

9.5 relaciona C_o y C_{dd}

- Si $C_o = 10 \gg C_{dd} = 9.5 * 10$
- Si $C_o = 100 \gg C_{dd} = 950$

$$9.5 = C_{dd}/C_o = V_{dd}$$

Es la razón de visita: por cada trabajo atendido, cuántas veces visita el dd.

a)

a.1) Tasa media de llegada: $\lambda_{dd} = Add/T$

$$\lambda_{dd} = 4108 \text{ tr}/24\text{h} = //\text{matemagia, factor de conversión y tal...} = 1.14 \text{ tr/s}$$

a.2) Productividad/Throughput del dd: $X_{dd} = C_{dd}/t$

Prácticamente igual a lo anterior, diferencia nimia porque A_{dd} y C_{dd} solo se llevan 1 de diferencia... así que digamos que $X_{dd} = 1.14 \text{ tr/s}$ también

b) Utilización media del dd: $U_{dd} = B_{dd}/T = 0.25$, 25% de utilización

c)

c.1) Tiempo medio de servicio: $S_{dd} = B_{dd}/C_{dd} = 0.22\text{s}$

c.2) Demanda media del servicio: $D_{dd} = B_{dd}/C_o = V_{dd} * S_{dd} = 9.5 * 0.22 = 2.1\text{s}$

d) Productividad del SERVER ojo:

- $X_o = C_o/T$
 - $V_{dd} = 9.5 = C_{dd}/C_o$;
 - $C_o = C_{dd}/V_{dd} = 98591/9.5 = 10378\text{tr}$

$$X_o = 10378/24\text{h} = 432\text{tr/h} = 0.12\text{tr/s}$$

Una buena práctica es volver a traducirlo todo al lenguaje dado en el enunciado, y no dejarlo con el genérico "trabajos":

1. Tasa media de llegada, λ_{dd} : 1.14 peticiones de RW al dd por segundo
Throughput del dd, X_{dd} : 1.14 peticiones de RW al dd por segundo
2. Utilización media del dd, U_{dd} : 25% de utilización
3. Tiempo medio de servicio, S_{dd} : 0.22s
Demanda media del servicio del disco duro, D_{dd} : 2.1s
4. Productividad del servidor, X_o : 0.12 peticiones de RW al dd por segundo

Leyes operacionales

El valor de todas las variables utilizadas en el análisis operacional dependerá del intervalo de observación T . **Existen, sin embargo, una serie de relaciones entre algunas variables operacionales que se mantienen válidas para cualquier intervalo de observación** y que no dependen de suposiciones sobre la distribución de los tiempos de servicio o de la forma en la que llegan los trabajos.

Estas relaciones se denominan leyes operacionales.

Estas leyes son tanto más útiles cuando se cumple la denominada **hipótesis del equilibrio de flujo**, que establece que si se escoge un intervalo de observación T suficientemente largo, se cumple que:

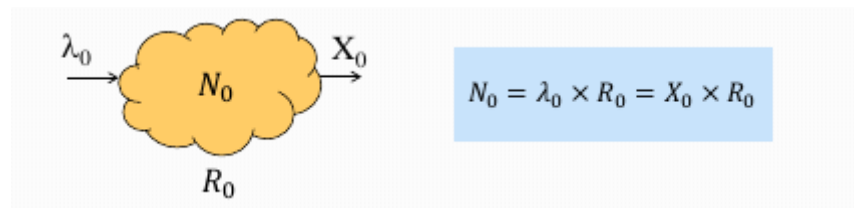
- El número de trabajos que completa el servidor coincide aproximadamente con los solicitados ($C_0 \approx A_0$). Dicho de otra manera, la productividad media coincide aproximadamente con la tasa media de llegada ($X_0 \approx \lambda_0$).
- El número de trabajos que completa cada estación de servicio coincide aproximadamente con los que se solicitan: ($C_i \approx A_i \rightarrow X_i \approx \lambda_i, \forall i=1...K$).

Las gallinas que entran por las que salen

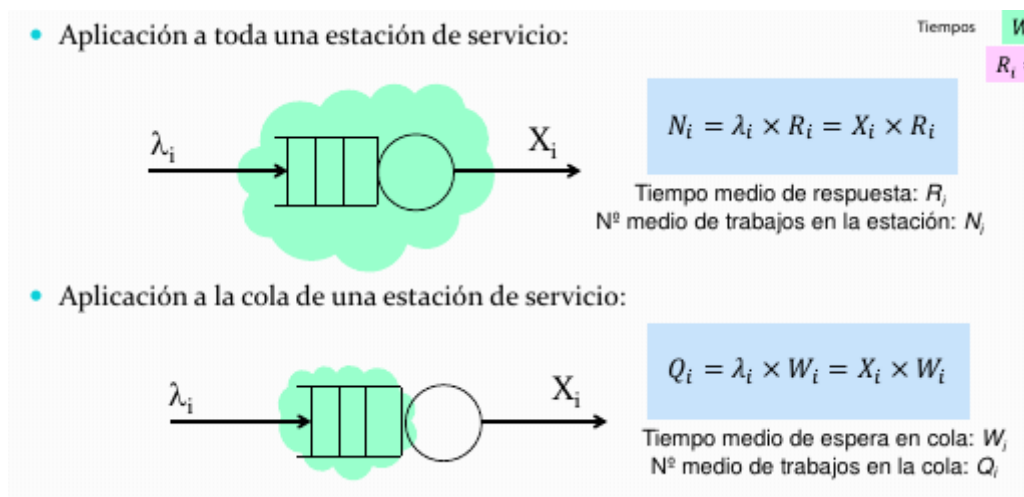
Ley de Little -- $N_0 = \lambda_0 \cdot R_0 = X_0 \cdot R_0$

La más importante - Simulación en SWAD

Sólo válida si está en equilibrio de flujo. Relaciona la productividad media X_0 con el tiempo medio de respuesta R_0 .



Puede aplicarse a cualquier estación de servicio, no sólo al servidor entero. La podemos aplicar a la estación completa, o a la cola de espera:



Ley de la utilización -- $U_i = X_i \cdot S_i = \lambda_i \cdot S_i$

Relaciona la utilización media de un dispositivo λ_i con el número de trabajos que completa por unidad de tiempo (=su productividad media, X_i) y el tiempo que le dedica, de media, a cada uno de ellos (=su tiempo de servicio medio, S_i).

Si equilibrio de flujo

$$U_i = X_i \times S_i = \lambda_i \times S_i$$

Ley del flujo forzado -- $X_i = X_0 \cdot V_i = \lambda_0 \cdot V_i = \lambda_i$

Las productividades X_0 , X_i (=flujos de salida) de cada estación de servicio tienen que ser proporcionales a la productividad global del servidor. La ley del flujo forzado relaciona la productividad del servidor con la de cada uno de los dispositivos que integran el mismo:

Si equilibrio de flujo

$$X_i = X_0 \times V_i = \lambda_0 \times V_i = \lambda_i$$

Recordemos que $V_i = C_i / C_0 = X_i / X_0$. V_i es el número de accesos del servidor a la estación de servicio por petición.

Relación Utilización-Demanda -- $U_i = X_0 \cdot D_i = \lambda_0 \cdot D_i$

Consecuencia de la ley de flujo forzado.

Las utilizaciones de cada dispositivo son proporcionales a las demandas de servicio del mismo, siendo la constante de proporcionalidad precisamente la productividad global del servidor (relación Utilización-Demanda de servicio):

Si equilibrio de flujo

$$U_i = X_0 \times D_i = \lambda_0 \times D_i \quad \text{Demo}$$

Recordemos que $D_i = B_i / C_0 = U_i / X_0$. D_i es la cantidad de tiempo que gasta la estación por cada trabajo que abandona el servidor.

Ejemplo

Un servidor de base de datos en equilibrio de flujo recibe una media de 120 consultas por minuto. Sabemos que su disco duro tarda, de media, 30ms en atender cada petición de E/S que le llega (48ms si incluimos la espera en la cola) y que su productividad es 25 peticiones de E/S completadas por segundo.

$$\lambda_0 = 120 \text{ tr/min} = 2 \text{ tr/s}$$

$$S_{dd} = 30 \text{ ms} = 0.03 \text{ s}$$

$$W_{dd} = 48 \text{ ms} - 30 \text{ ms} = 18 \text{ ms} = 0.018 \text{ s}$$

$$X_{dd} = 25 \text{ tr/s}$$

Calcule:

a) El número medio de peticiones de E/S en la cola de espera del disco duro.

b) ¿Cuánto tiempo, de media, consumen los accesos al disco duro por cada consulta que se realiza al servidor?

a) Q_{dd} (trabajos medios en cola del dd)

LEY DE LITTLE aplicada a la cola de espera de una estación:

$$Q_{dd} = \lambda_0 * W_{dd} = X_{dd} * W_{dd}$$

$$Q_{dd} = 25 * 0.018 = 0.45 \text{ tr en cola}$$

b) D_{dd} (tiempo que pasa el dd en cada tr que le manda el servidor)

LEY DE UTILIZACIÓN

$$U_{dd} = X_0 * D_{dd} = \lambda_0 * D_{dd}$$

$$U_{dd} = X_{dd} * S_{dd} = 0.03 * 25 = 0.75$$

$$0.75 = 2 * D_{dd}; D_{dd} = 0.375 \text{ s/tr}$$

Ley general del tiempo de respuesta

El tiempo medio de respuesta que experimenta, de media, una petición a un servidor en equilibrio de flujo se puede calcular teniendo en cuenta que cada una de ellas ha tenido que “visitar” V_i veces al dispositivo i -ésimo, requiriendo cada visita una media de R_i segundos:

$$R_0 = V_1 * R_1 + V_2 * R_2 \dots$$

Ley del tiempo de respuesta interactivo

Un servidor en una red cerrada siempre está en equilibrio de flujo (siempre supondremos que el tamaño de las colas es suficientemente grande, en este caso, $\geq NT$).

Al ser una red cerrada, el número total de trabajos (=clientes) en la red ($NT = N_Z + N_0$), es constante.

Aplicamos la ley de Little a diversas partes de la red de colas:

- Ley de Little aplicada a los clientes en reflexión: $N_Z = X_0 \times Z$, donde N_Z = Número **medio** de clientes (=trabajos) en reflexión.
- Ley de Little aplicada al servidor: $N_0 = X_0 \times R_0$



$$R_0 = \frac{N_T}{X_0} - Z$$

que se conoce como la *Ley del tiempo de respuesta interactivo*.

5.3. Límites optimistas del rendimiento

5.4. Técnicas de mejora del rendimiento

5.5. Algoritmos de resolución de redes de colas