

Análisis del problema de control de congestión en redes utilizando argumentos de pasividad.

Análisis y Control de Sistemas No Lineales
Monografía Final

Andrés Ferragut

Marzo de 2009

1. Introducción

Un problema central en las redes de datos es el del reparto de recursos. Una red, desde un punto de vista simplificado, es un conjunto de enlaces que conectan diferentes nodos, y cada enlace puede cursar una cierta cantidad de tráfico máximo, medida por ejemplo en bits por segundo. Entre los nodos de la red se establecen rutas y por último, sobre estas rutas se abren diferentes conexiones.

En los primeros despliegues de Internet, cuando el poder de cómputo de los terminales y servidores no era grande y las conexiones que se abrían eran esporádicas, el uso del ancho de banda no era un problema y cada conexión podía intentar transmitir una cantidad arbitraria de datos, pues de todos modos esta cantidad no era grande debido al poco poder de cómputo y a su vez no obstruía otras comunicaciones, porque las mismas casi nunca eran simultáneas.

Sin embargo, a medida que creció la red en usuarios y capacidad de procesamiento, y aún antes de llegar a la etapa comercial de Internet, se observó el problema conocido como “congestion collapse” o colapso por congestión. Éste era debido a que las conexiones individuales no recibían realimentación sobre el estado de carga de los enlaces que utilizaba, y por lo tanto no reaccionaba, por ejemplo limitando la cantidad de tráfico que debía transmitir, ante una situación de congestión en su camino. Esto obviamente volvía todas las conexiones que compartían el recurso congestionado inutilizables, debido a que los buffers de los enlaces se llenaban y la consecuente pérdida de datos requería retransmisiones, es decir más paquetes. Este fenómeno de realimentación positiva es el “congestion collapse”.

Resultó evidente entonces que debía diseñarse un mecanismo de *control de congestión* que tuviera en cuenta este fenómeno y tomara acciones desde la fuente de datos para evitarlo. Esto llevó al diseño del mecanismo de control de congestión de TCP de V. Jacobson [1] (c.f. [9] para una explicación detallada del mismo). En él se introduce un mecanismo de ventana deslizante para controlar la cantidad de paquetes en circulación de una determinada conexión, y se modifica el tamaño de dicha ventana de acuerdo a señales que indiquen congestión. En el caso de TCP Tahoe (y actualizaciones posteriores como TCP Reno y Newreno), se utilizan las pérdidas de paquetes o un desorden en el arribo de reconocimientos (ACKs) como indicador de congestión y esto provoca una reducción drástica de la ventana, que luego crece lentamente si el flujo de paquetes es correcto.

Este diseño fue realizado utilizando la intuición de los ingenieros de redes de la época, pero carecía de un modelo claro de funcionamiento y del reparto de recursos que el mismo logra en una red concreta. Se trata sin embargo de un mecanismo de control que cumple con algunas premisas naturales al trabajar con un sistema distribuido y potencialmente enorme como la Internet. En primer lugar se trata de un control *descentralizado* en donde los usuarios reaccionan a señales de la red de manera independiente, cada uno reaccionando a pérdidas de sus propios paquetes. Asimismo, los enlaces descartan paquetes (generan la señal de realimentación) sin mirar de qué usuario provienen o el estado de los otros enlaces. Es deseable además que el mecanismo intente “llenar” la red, es decir, mantener el tráfico en los enlaces cerca de la capacidad máxima. Por último, también es deseable que el mecanismo sea estable y converja a un equilibrio razonable en términos de reparto de recursos, y que sea robusto a perturbaciones por tráfico no controlado o a los retardos inherentes a la propagación de los datos en la red.

Un punto de quiebre en el análisis de este tipo de sistemas se produce con el trabajo de Kelly et. al [5]. En dicho trabajo se modela el control de congestión implementado por los protocolos de capa de transporte en las redes actuales como una solución descentralizada a un problema de optimización convexa sujeto a restricciones de capacidad de la red.

Este descubrimiento llevó a que en la literatura comenzara a analizarse el control de congestión como un sistema de control no lineal distribuido en que los diferentes elementos de la red (usuarios, enlaces) siguen una dinámica en la que únicamente utilizan información disponible localmente. En el trabajo citado y subsiguientes [6, 7] se proponen nuevos mecanismos de control para lograr diferentes objetivos. El principal problema a analizar consistió en determinar la estabilidad de estos mecanismos, lo cual se realizó por métodos “ad-hoc”.

Sin embargo, en el trabajo de Wen y Arcak [10] se formula el problema desde un punto de vista más general y se aplican las herramientas de pasividad de sistemas de control para probar la estabilidad de una familia grande de dinámicas, que incluye a los controladores anteriormente propuestos.

A su vez, en un trabajo reciente [12] hemos aplicado esta misma técnica a la solución de un problema de optimización similar al de la capa de transporte, pero en el que intervienen las capas superiores, procurando establecer un punto de operación de la red que sea óptimo o justo desde un punto de vista del usuario y no de las conexiones individuales.

A lo largo de este trabajo iremos abordando estos puntos. En la sección 2 haremos una descripción breve de las dinámicas que modelan el control de congestión. A continuación en 3 resumiremos los resultados de la teoría de sistemas pasivos. En 4 se describe la aplicación de estos resultados a los problemas de control de congestión siguiendo a [10]. Luego, en 5 se describe el contexto de control del número de conexiones y se discute la aplicación de la pasividad a los problemas que allí surgen, que están conectados al control de congestión. Por último en 6 se presentan algunas conclusiones y posibles líneas de trabajo futuro.

2. Modelado del control de congestión

El control de congestión implementado por TCP para transferencias de archivos por Internet es un algoritmo de *ventana*. El transmisor mantiene un cierto valor $W(t)$ que consiste en la cantidad de paquetes que puede mantener “en tránsito” para una conexión dada. Se consideran en tránsito aquellos paquetes cuyo reconocimiento de recepción o ACK (de acknowledgement) aún no ha sido recibido.

El control comienza con una fase inicial denominada *slow start* que consiste en duplicar W , desde un valor inicial de un paquete, cada vez que una ventana completa ha sido reconocida. Esta fase permite descubrir rápidamente el orden de W alrededor del cual se debe trabajar y culmina cuando llega la primera señal de congestión, que en general es un paquete perdido, lo cual se detecta por que el reconocimiento ha tardado más de lo normal en llegar (timeout).

Aquí es que se pasa a una fase denominada *congestion avoidance* que en las transferencias largas de archivos es la que contiene la mayor parte de los datos, y es la que nos interesa modelar. Describamos entonces como funciona esta fase para el caso de TCP Reno, la versión más común en la actualidad.

En TCP Reno [3], la ventana se actualiza de la siguiente manera:

$$W \leftarrow W + \frac{1}{W} \quad \text{por cada ACK recibido} \quad (1a)$$

$$W \leftarrow W/2 \quad \text{por cada paquete perdido} \quad (1b)$$

aquí los paquetes perdidos se detectan por lo que se denomina “triple ACK duplicado” y consiste en dar por perdido un paquete cuando el receptor informa 3 veces consecutivas que se encuentra esperándolo.¹

Las conexiones en Internet además tienen normalmente un retardo o latencia entre los extremos comunicados. Se denomina Round Trip Time (*RTT*) al retardo de ida y vuelta de una conexión. Si suponemos que este retardo es significativamente mayor que el tiempo que se demora en transmitir los paquetes (esta hipótesis es válida para redes de alto ancho de banda y recorridos largos) podemos suponer que la fuente envía W paquetes cada *RTT* y la evolución de W será:

$$\begin{aligned} W_{n+1} &= W_n + \frac{1}{W_n} \cdot W_n = W_n + 1 \quad \text{si la ventana se recibe correctamente} \\ W_{n+1} &= W_n/2^k \quad \text{si se pierden } k \text{ paquetes} \end{aligned}$$

El objetivo de este control queda ahora más claro: la ventana sube linealmente en cada *RTT* mientras los paquetes son recibidos correctamente, lo que se interpreta como que no hay congestión. Cuando la ventana sube suficiente como para congestionar los enlaces que intervienen en la conexión, entonces se baja drásticamente el valor y por consiguiente la tasa de transmisión.

En [4], Mathis et. al modelaron el comportamiento del algoritmo (1) suponiendo que existe una probabilidad p de pérdida de paquetes (o mejor dicho, de que haya congestión) llegando a la siguiente relación entre

¹En TCP solo se informa el no. de secuencia del *próximo* paquete esperado, por lo que los paquetes siguientes a una pérdida generan estos ACKs repetidos

la tasa de transmisión de equilibrio (en paquetes/seg.) y los parámetros del sistema:

$$\hat{x} = \frac{1}{RTT} \sqrt{\frac{2}{3p}} \quad (2)$$

que denominaremos en adelante *fórmula de Mathis*.

Sin embargo, a los efectos de estudiar si el algoritmo de control converge a este valor de equilibrio es necesario un modelo dinámico del TCP. Para ello, Kelly et. al. en [5] propusieron el siguiente análisis en tiempo continuo.

Sea $x(t)$ la tasa de transmisión y $p(t)$ la probabilidad de pérdida de paquetes. Entonces $x(t)p(t)$ representa la proporción de paquetes que se pierden y $x(t)(1-p(t))$ la proporción que llega correctamente. De la ecuación (1) se tiene que si $W(t)$ es la ventana en tiempo t , ésta aumenta en $1/W(t)$ por cada paquete que llega correctamente y baja un factor $\beta = 1/2$ por cada paquete no recibido, por lo que podemos escribir que aproximadamente:

$$W(t) = \frac{x(t)(1-p(t))}{W(t)} - \beta W(t)x(t)p(t)$$

Observando además que si el RTT de la conexión es T entonces se tendrá que:

$$x(t) = W(t)/T \Rightarrow \dot{x}(t) = \dot{W}(t)/T$$

y sustituyendo en la ecuación anterior:

$$\dot{x} = \frac{1-p}{T^2} - \beta x^2 p \quad (3)$$

Observemos que en el equilibrio $\dot{x} = 0$ (y suponiendo que existe una probabilidad de pérdida \hat{p} de equilibrio) se cumple:

$$\hat{x} = \frac{1}{T} \sqrt{\frac{1-\hat{p}}{\beta \hat{p}}}$$

que para valores pequeños de probabilidad de pérdida recupera la esencia del análisis de Mathis.

Para terminar de entender el funcionamiento del sistema hace falta describir cómo se genera la señal de realimentación $p(t)$, es decir, cómo se relacionan los rates de usuarios individuales con la probabilidad de pérdida. Además, sería deseable entender el equilibrio de este algoritmo que corre en cada usuario por separado. El análisis de [5] nos da la pista al reescribir la ecuación (3) de la siguiente manera:

$$\dot{x} = k(x) (U'(x) - p) \quad U(x) = -\frac{1}{\beta T^2 x} \quad (4)$$

en esta expresión $k(x)$ es positiva y $U(x)$ es una función cóncava y creciente de x y se ha despreciado el término $1-p$ de (3). El equilibrio de esta ecuación entonces puede interpretarse en términos económicos: si $U(x)$ es la utilidad del usuario por obtener una tasa x y p es el precio por unidad entonces el equilibrio es el que maximiza el excedente del usuario, el valor óptimo a "comprar".

Hagamos intervenir ahora a la topología de la red. Sea R una matriz de tamaño $L \times N$ siendo L la cantidad de enlaces y N la cantidad de conexiones presentes en la red, donde $R_{li} = 1$ si la conexión i utiliza el enlace l y 0 en otro caso. Si $x = (x_i)$, $i = 1, \dots, N$ es el vector de tasas de transmisión de cada conexión, el rate que llega a cada enlace será $y = (y_l)$ con $y_l = \sum_i R_{li} x_i$ o bien:

$$y = Rx \quad (5)$$

Cada enlace ahora tendrá una cierta probabilidad de descartar paquetes en su buffer, y esto se relaciona con la capacidad c_l del enlace y el rate de entrada y_l al mismo. Un posible modelo estático para este fenómeno es tomar:

$$p_l = \frac{1}{y_l} (y_l - c_l)^+ \quad (6)$$

es decir, la probabilidad de pérdida o precio del enlace l es igual a la proporción de paquetes por encima de la capacidad del mismo. Observemos que esta es una función monótona de y_l .

Si ahora suponemos que la probabilidad de pérdida es pequeña y que cada enlace se comporta de manera independiente, podemos aproximar a primer orden la probabilidad de pérdida que ve la conexión i , q_i como la suma de las probabilidades de pérdida en cada enlace que atraviesa la conexión, es decir $q_i = \sum_l R_{li} p_l$ o bien:

$$q = R^t p \quad (7)$$

Por lo cual, juntando las ecuaciones (4), (5), (6), (7) se tiene la siguiente dinámica:

$$\dot{x}_i = k(x_i) (U'_i(x_i) - q_i) \quad (8a)$$

$$y = Rx \quad (8b)$$

$$p_l = f(y_l) = \frac{1}{y_l} (y_l - c_l)^+ \quad (8c)$$

$$q = R^t p \quad (8d)$$

que es esquematizada en la figura 1.

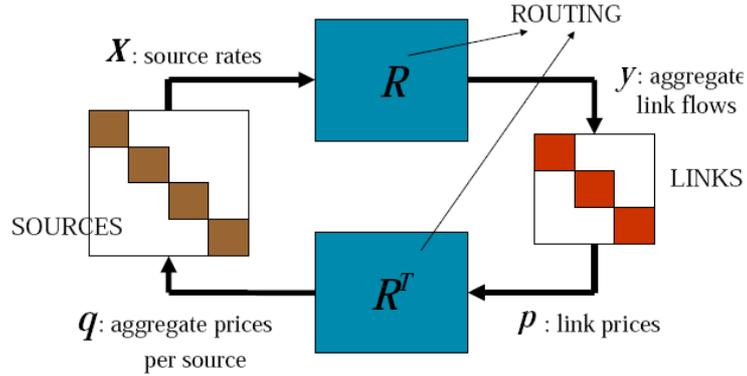


Figura 1: Diagrama de lazo cerrado del control de congestión

El equilibrio de este sistema verifica:

$$U'(\hat{x}) = \hat{q}, \hat{y} = R\hat{x}, \hat{p} = f(\hat{y}), \hat{q} = R^t \hat{p}$$

que es exactamente la condición de optimalidad del siguiente problema de optimización convexa:

$$\text{máx} \left\{ \sum_i U_i(x_i) - \sum_l \int_0^{y_l} f_l(s) ds : x \in \mathbb{R}^n, y = Rx \right\} \quad (9)$$

donde $f_l(y) = \frac{1}{y} (y - c_l)^+$ actúa como *función de barrera* o de penalidad, que en el fondo lo que hace es imponer aproximadamente la restricción:

$$Rx \leq c$$

donde $c = (c_l)$ es el vector que contiene las capacidades de los enlaces. En la figura 2 puede verse la función de barrera.

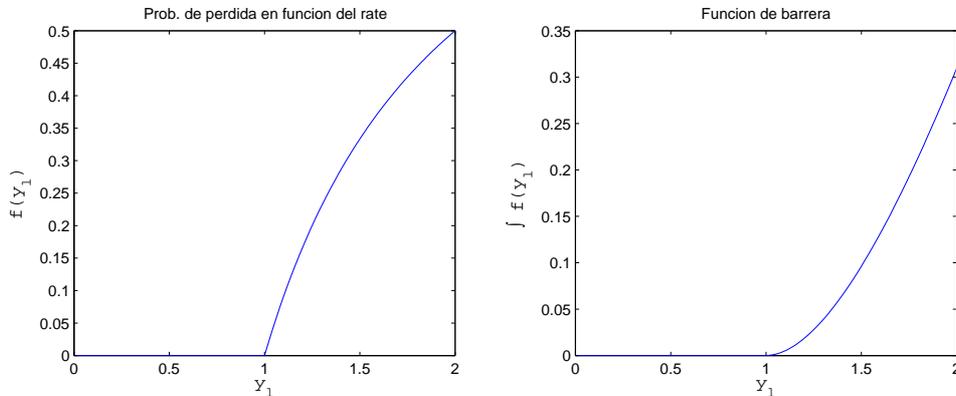


Figura 2: Función de barrera para el problema primal de la ecuación (9) con $c_l = 1$

Observación 2.1. El problema (9) es un problema convexo. Esto surge de observar que las U_i son cóncavas por hipótesis y la función de barrera es la integral de una función creciente, por lo que es convexa. Por lo tanto la función objetivo es diferencia de una función cóncava y una convexa, por lo que es cóncava. A la dinámica (8) se le denomina dinámica *primal* porque actúa sobre las variables primales del problema de optimización, las x_i .

El resultado de este análisis es que este modelo de TCP resuelve aproximadamente el siguiente problema de optimización convexa:

$$\begin{aligned} \text{máx} \quad & \sum_i U_i(x_i) \\ \text{s.t.} \quad & Rx \leq c \end{aligned} \quad (10)$$

El problema (10) es un problema de tipo economía de mercado, en el que se busca maximizar la utilidad global de los actores sujeto a una restricción de recursos disponibles, que en este caso son las capacidades de los enlaces. En la economía es habitual encontrar este tipo de modelos y de algoritmos descentralizados para converger al óptimo. Es aquí que cobra sentido el uso del término “precios de congestión” para las probabilidades de pérdida, en este caso dadas por $p_l = f(y_l)$.

Se tiene además la siguiente proposición, cuya prueba puede verse en [10], Apéndice III.

Proposición 2.1. *Si las funciones U_i son estrictamente cóncavas, f es estrictamente creciente y R es de rango completo entonces cada uno de los problemas (9) y (10) tiene su solución única \hat{x}, \hat{p} .*

Sería interesante disponer de un algoritmo que resuelva *exactamente* el problema (10). Para ello usamos el principio de dualidad de Lagrange [11]. El lagrangeano del problema (10), llamando p a los precios o multiplicadores, es:

$$L(x, p) = \sum_{i=1}^N U_i(x_i) - \sum_{l=1}^L p_l(y_l - c_l)$$

Como el problema es convexo, no hay gap de dualidad y el óptimo debe verificar:

$$(\hat{x}, \hat{p}) : \min_{p \geq 0} \max_{x \geq 0} L(x, p)$$

por lo que con la notación anterior, en el óptimo se debe verificar:

$$U'_i(\hat{x}_i) = \hat{q}_i$$

y el multiplicador debe verificar la “complementary slackness” [11]:

$$\hat{p}_l \begin{cases} = 0 & \hat{y}_l < c_l \\ \leq 0 & \hat{y}_l = c_l \end{cases}$$

Para tratar de alcanzar este equilibrio, en [6] los autores propusieron una dinámica *dual* en el que la dinámica del sistema está en la actualización de precios y las fuentes siempre reaccionan inmediatamente al precio actual q_i fijando x_i en el valor óptimo para estos precios. La dinámica completa es:

$$x_i = U_i^{-1}(q_i) \quad (11a)$$

$$y = Rx \quad (11b)$$

$$\dot{p}_l = \gamma(y_l - c_l)_p^+ \quad (11c)$$

$$\dot{q} = R^t p \quad (11d)$$

Aquí, $\gamma > 0$ y la función $(\cdot)_p^+$ corresponde a una saturación para que el precio no tome valores negativos. Más exactamente:

$$(x)_p^+ = \begin{cases} x & \text{si } p > 0 \\ \text{máx}(x, 0) & \text{si } p \leq 0 \end{cases}$$

En [6] se hace notar que si $\gamma = 1/c_l$ entonces este modelo se corresponde aproximadamente a tomar como precio de congestión el retardo de cola sufrido por los paquetes. Esto llevó a la implementación de la versión FAST de TCP por los autores.

Por último, podemos mezclar la dinámica primal y dual en un solo algoritmo: las fuentes buscan adaptarse al precio actual y el precio se mueve según cuánto se violan las restricciones de capacidad. La dinámica

completa queda:

$$\dot{x}_i = k(U'_i(x_i) - q_i) \quad (12a)$$

$$y = Rx \quad (12b)$$

$$\dot{p}_i = \gamma(y_i - c_i)_{p_i}^+ \quad (12c)$$

$$q = R^t p \quad (12d)$$

que se denomina dinámica *primal-dual*. Aquí k, γ son constantes positivas.

En la sección 4 se estudia la estabilidad de estas dinámicas usando argumentos de pasividad. Para ello, pasemos primero al estudio de la noción de pasividad.

3. Pasividad en sistemas de control

El enfoque de *pasividad* en el análisis de sistemas realimentados de control es una alternativa para probar estabilidad analizando diferentes partes del lazo e identificando en ellas sistemas *pasivos* en el sentido de que no aportan energía al sistema (o directamente la disipan). El principal resultado es que la realimentación negativa de sistemas pasivos da lugar a un sistema pasivo y por ende estable.

El concepto de pasividad proviene del análisis de circuitos con componentes pasivas (resistencias, inductores, capacitores) y de estructuras mecánicas. En ellos, la noción de energía está bien definida y el sistema es pasivo si la energía *almacenada* en el mismo es menor o igual a la energía *entregada* al sistema.

Esto lleva a la siguiente definición. Dado el sistema dinámico:

$$\dot{x} = f(x, u) \quad (13a)$$

$$y = h(x, u) \quad (13b)$$

con $f : \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}^n$ localmente Lipschitz, $h : \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}^m$ continua, $f(0, 0) = 0, h(0, 0) = 0$. Obsérvese que el sistema tiene el mismo número de entradas y salidas.

Definición 3.1. *El sistema (13) se dice pasivo si existe una función semidefinida positiva, continuamente diferenciable $V(x)$, llamada función de almacenamiento tal que:*

$$u^t y \geq \frac{\partial V}{\partial x} f(x, u) + \epsilon u^t u + \delta y^t y + \rho \psi(x), \quad \forall (x, u) \in \mathbb{R}^n \times \mathbb{R}^m \quad (14)$$

con ϵ, δ y ρ constantes positivas y ψ una función semidefinida positiva de x tal que $\psi(x(t)) \equiv 0 \Rightarrow x(t) \equiv 0$ para toda solución $x(t)$ con cualquier u .

El término $\rho \psi(x)$ se denomina *tasa de disipación de estados*. A su vez, se clasifican los sistemas pasivos según los valores de las constantes positivas involucradas en:

- *Sin pérdida (lossless)* si en (14) se da la igualdad y a su vez $\epsilon = \delta = \rho = 0$.
- *Estrictamente pasivo en la entrada* si $\epsilon > 0$.
- *Estrictamente pasivo en la salida* si $\delta > 0$.
- *Estrictamente pasivo en el estado* si $\rho > 0$.

y estos nombres se combinan en el caso en que varias de las constantes sean positivas.

Observación 3.1. La condición sobre $\psi(x)$ se interpreta de la siguiente manera: si ψ es definida positiva, el término de disipación correspondiente al estado será positivo siempre que el sistema no esté en el origen. Pero aún si ψ es semidefinida, la condición adicional garantiza que no hay soluciones que no disipen energía que no sean la trayectoria de equilibrio (es decir, no hay soluciones no triviales en la variedad $\psi(x) = 0$).

A su vez, para el caso de sistemas estrictamente pasivos en la salida, es importante la siguiente propiedad adicional:

Definición 3.2. *El sistema (13) se dice de origen observable si ninguna solución no trivial de $\dot{x} = f(x, 0)$ puede quedarse en $S = \{x : h(x, 0) = 0\}$.*

Esta definición dice que no puede ocurrir que el sistema sin forzamiento ($u = 0$) produzca una salida nula. Esto básicamente garantiza que el término $y^t y$ sea positivo aún cuando no haya entrada y se relaciona con la observabilidad en el sentido de sistemas lineales. En el caso lineal el sistema no forzado será:

$$\begin{aligned}\dot{x} &= Ax \\ y &= Cx\end{aligned}$$

En este caso el par (A, C) es observable es equivalente a:

$$y(t) = Ce^{At}x(0) \equiv 0 \Leftrightarrow x(0) = 0 \Leftrightarrow x(t) \equiv 0$$

por lo que la definición 3.2 es equivalente a la observabilidad en este caso.

El siguiente lema relaciona la pasividad con las nociones de estabilidad de Lyapunov del sistema no forzado.

Lema 3.1. *Consideremos el sistema (13) entonces:*

1. *Si el sistema es pasivo con $V(x)$ definida positiva entonces el origen de $\dot{x} = f(x, 0)$ es estable.*
2. *Si el sistema es estrictamente pasivo en la salida con una función V definida positiva y observable en el origen entonces el origen de $\dot{x} = f(x, 0)$ es asintóticamente estable.*
3. *Si el sistema es estrictamente pasivo en el estado con V definida positiva entonces el origen es asintóticamente estable.*

Si además $V(x)$ es radialmente no acotada, en 2 y 3 entonces el origen es un atractor global.

Demostración:

Para el primer caso, $V(x)$ es apropiada como función de Lyapunov para el sistema por ser definida positiva y por ser pasivo se tiene que:

$$\dot{V} = \frac{\partial V}{\partial x} f(x, 0) \leq u^t y = 0$$

por lo que el origen es estable.

Para el segundo caso nuevamente tomamos V como función de Lyapunov y observamos que:

$$\dot{V} = \frac{\partial V}{\partial x} f(x, 0) \leq -\delta h^t(x, 0)h(x, 0) \leq 0$$

y además si $\dot{V} = 0$ se tiene que $h(x, 0) = 0$, es decir las soluciones que verifican $\dot{V} = 0$ deben permanecer en el conjunto S definido previamente, y por la observabilidad en el origen esto solo ocurre para la solución $x(t) \equiv 0$ por lo que, aplicando el teorema de LaSalle el origen es asintóticamente estable, y será globalmente asintóticamente estable si V es radialmente no acotada.

Por último, tomando nuevamente V como función de Lyapunov se verifica que:

$$\dot{V} = \frac{\partial V}{\partial x} f(x, 0) \leq -\rho \psi(x) \leq 0$$

y nuevamente si $\dot{V} = 0$ deberá cumplirse que $\psi(x(t)) \equiv 0$ y esto solo ocurre para la solución $x(t) \equiv 0$ por la condición de la función ψ . □

Para terminar esta parte, definiremos la pasividad de un sistema no lineal sin memoria:

Definición 3.3. *Sea $\Psi : \mathbb{R}^m \rightarrow \mathbb{R}^m$ con $y = \Psi(u)$ una no linealidad. Diremos que es pasiva si:*

$$u^t y \geq \epsilon u^t u + \delta y^t y$$

Si además $\epsilon > 0$ diremos que es pasiva en la entrada y si $\delta > 0$ diremos que es pasiva en la salida.

A continuación, vincularemos la definición de pasividad con una condición en frecuencia para el caso de sistemas lineales. Para ello recordaremos una definición y el lema de Kalman-Yakubovich-Popov.

Definición 3.4. *Una matriz de transferencias $p \times p$ $Z(s)$ se dice real positiva si:*

- *Los elementos de $Z(s)$ son analíticos en el semiplano derecho.*

- Cualquier polo imaginario puro de cualquier elemento de $Z(s)$ es simple y la matriz de residuos es semidefinida positiva y
- para cada ω tal que $j\omega$ no es polo de Z la matriz $Z(j\omega) + Z^T(-j\omega)$ es semidefinida positiva.

A su vez, una función de transferencia $Z(s)$ es estrictamente real positiva si $Z(s - \varepsilon)$ es real positiva para algún $\varepsilon > 0$.

Observemos que en el caso $p = 1$ de un sistema SISO la condición anterior se reduce a que $\Re(Z(j\omega)) \geq 0$ $\forall \omega$ o bien que el diagrama de Nyquist de Z esté completamente en el semiplano derecho.

Sea ahora el sistema lineal:

$$\dot{x} = Ax + Bu \quad (15a)$$

$$y = Cx + Du \quad (15b)$$

entonces su matriz de transferencias asociada es $G(s) = C(sI - A)^{-1}B + D$. El lema de Kalman-Yakubovich-Popov establece la siguiente relación entre la condición de ser real positiva y la representación de estados del sistema:

Lema 3.2 (Kalman-Yakubovich-Popov). *Sea $G(s) = C(sI - A)^{-1}B + D$ una matriz de transferencias real positiva donde A es Hurwitz, (A, B) es controlable y (A, C) es observable, entonces $Z(s)$ es estrictamente real positiva si y solo si existe una matriz $P = P^t$ definida positiva, matrices W y L y una constante positiva ε tales que:*

$$PA + A^tP = -L^tL - \varepsilon P \quad (16a)$$

$$PB = C^t - L^tW \quad (16b)$$

$$W^tW = D + D^t \quad (16c)$$

Luego de estos preliminares pasamos a la siguiente proposición:

Proposición 3.1. *Un sistema lineal invariante en el tiempo con una realización mínima A, B, C, D y matriz de transferencias $G(s) = C(sI - A)^{-1}B + D$ real positiva es estrictamente pasivo en el estado.*

Demostración:

Sean P, W, L y ε como en (16) y tomemos $V(x) = \frac{1}{2}x^tPx$ como candidato a función de almacenamiento, entonces:

$$\begin{aligned} u^ty - \dot{V} &= u^t(Cx + Du) - x^tP(Ax + Bu) \\ &= u^tCx + u^tDu - x^tPAx - x^tPBu \\ &= u^tCx + \frac{1}{2}u^t(D + D^t)u - \frac{1}{2}x^t(PA + A^tP)x - x^tPBu \\ &= u^t(B^tP + W^tL)x + \frac{1}{2}u^tW^tWu + \frac{1}{2}x^t(L^tL)x + \frac{1}{2}\varepsilon x^tPx - x^tPBu \\ &= u^tW^tLx + \frac{1}{2}u^tW^tWu + \frac{1}{2}x^t(L^tL)x + \frac{1}{2}\varepsilon x^tPx \\ &= \frac{1}{2}(Lx + Wu)^t(Lx + Wu) + \frac{1}{2}\varepsilon x^tPx \\ &\geq \frac{1}{2}\varepsilon x^tPx \end{aligned}$$

por lo que la desigualdad de la definición (14) se cumple con $\rho\psi(x) = \frac{1}{2}\varepsilon x^tPx$ y por ser P definida positiva el sistema es estrictamente pasivo. \square

Corolario 3.1. *Dado un sistema lineal SISO con función de transferencia $G(s)$ tal que el diagrama de Nyquist de G verifica $\Re[G(j\omega)] > 0$ entonces el sistema es estrictamente pasivo.*

Usaremos este corolario más adelante para probar la estabilidad local de un sistema de control de congestión.

Antes de pasar al estudio de sistemas realimentados analizaremos el siguiente ejemplo en donde quedan en evidencia las diferentes formas de pasividad y se aplica el corolario anterior a su vez para decidir la pasividad del sistema en términos de su diagrama de Nyquist.

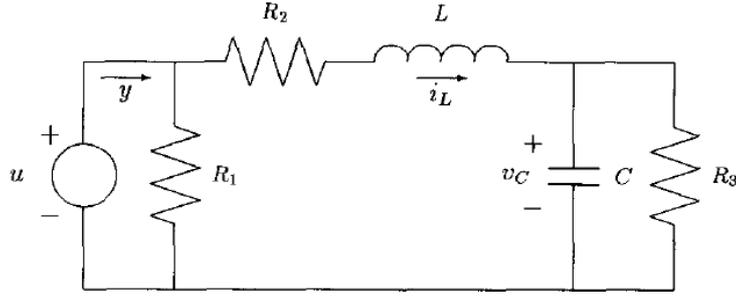


Figura 3: Circuito RLC de ejemplo

Ejemplo 3.1 (Khalil, ejemplo 10.10). Se considera el circuito de la figura 3, compuesto de una fuente de tensión y diferentes componentes pasivas.

Tomamos como entrada u la tensión de la fuente y como salida y la corriente por la misma. En este caso el producto uy representa la potencia instantánea entregada por la fuente al sistema y es de esperar que dicha potencia se disipe en las componentes del circuito.

Los estados serán x_1 la corriente por el inductor y x_2 la tensión en bornes del capacitor, lo que lleva a la siguiente representación de estados del sistema:

$$\begin{aligned} L\dot{x}_1 &= u - R_2x_1 - x_2 \\ C\dot{x}_2 &= x_1 - \frac{1}{R_3}x_2 \\ y &= x_1 + \frac{1}{R_1}u \end{aligned}$$

La energía en este circuito será almacenada en el inductor y el capacitor por lo que es de esperar que podamos tomar como función de almacenamiento a:

$$V(x) = \frac{1}{2}Lx_1^2 + \frac{1}{2}Cx_2^2$$

Se tiene entonces que:

$$\begin{aligned} \dot{V} &= Lx_1\dot{x}_1 + Cx_2\dot{x}_2 \\ &= x_1(u - R_2x_1 - x_2) + x_2(x_1 - \frac{1}{R_3}x_2) \\ &= x_1u - R_2x_1^2 - \frac{1}{R_3}x_2^2 \end{aligned}$$

Sustituyendo $x_1 = y - \frac{1}{R_1}u$ en el primer término se tiene que:

$$\dot{V} = uy - \frac{1}{R_1}u^2 - R_2x_1^2 - \frac{1}{R_3}x_2^2$$

de donde:

$$uy = \dot{V} + \frac{1}{R_1}u^2 + R_2x_1^2 + \frac{1}{R_3}x_2^2 \geq \dot{V}$$

por lo que el sistema es pasivo y precisamente los términos que involucran a las resistencias determinan la tasa de disipación.

Observemos por ejemplo que si $R_1 < \infty$ el sistema es estrictamente pasivo en la entrada, con $\varepsilon = 1/R_1$. Esto se condice con el hecho de que la resistencia R_1 en paralelo podría modelar pérdidas debidas a la fuente de tensión. El caso $R_2 > 0$ y $R_3 < \infty$ modela pérdidas en los integradores de este sistema de segundo orden, y corresponden al caso de pasividad en el estado.

Para ver la aplicación del corolario 3.1, calculemos la transferencia de este sistema para un valor de las componentes involucradas. Tomemos por ejemplo:

$$R_1 = R_3 = 1k\Omega, R_2 = 1\Omega, L = 1mHy, C = 1nF$$

En ese caso la función de transferencia del sistema queda:

$$G(s) = \frac{0,001s^2 + 2001s + 2,001 \cdot 10^9}{s^2 + 1,001 \cdot 10^6s + 1,001 \cdot 10^{12}}$$

cuyos diagramas de Bode y Nyquist pueden verse en la figura 4. Como se observa, el retardo de fase introducido por el sistema (estable) no supera los 90 grados y esto hace que la transferencia sea real positiva.

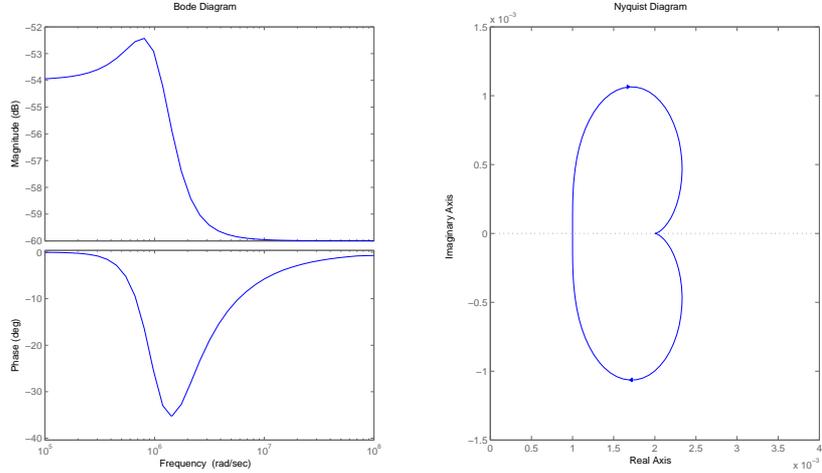


Figura 4: Diagramas de Bode y Nyquist de la transferencia del circuito 3

Pasemos ahora al análisis de sistemas realimentados. Para ello consideramos el sistema realimentado de la figura 5.

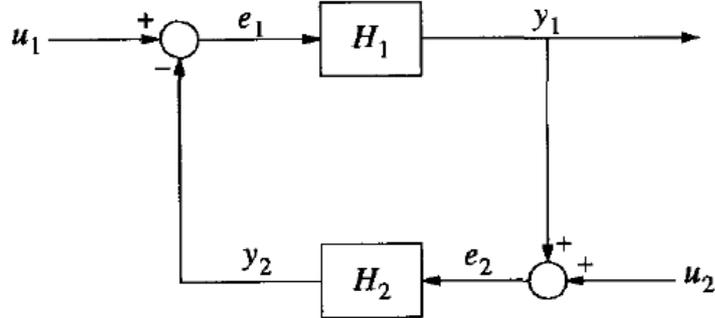


Figura 5: Realimentación negativa de sistemas de control.

Supongamos que H_1 y H_2 son dos sistemas dinámicos de la forma (13), es decir:

$$\begin{aligned} \dot{x}_i &= f_i(x_i, e_i) \\ y_i &= h_i(x_i, e_i) \end{aligned}$$

para $i = 1, 2$. Supongamos además que podemos escribir el modelo en espacio de estados del sistema realimentado completo en la forma de la (13) tomando:

$$x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}, u = \begin{bmatrix} u_1 \\ u_2 \end{bmatrix}, y = \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}$$

Se tendrá entonces el siguiente teorema:

Teorema 3.1. Si H_1 y H_2 descritos por el modelo anterior son sistemas pasivos con funciones de almacenamiento $V_1(x)$ y $V_2(x)$, es decir verifican:

$$e_i^t y_i \geq \frac{\partial V_i}{\partial x_i} f_i(x_i, e_i) + \epsilon_i e_i^t e_i + \delta_i y_i^t y_i + \rho_i \psi_i(x_i) \quad i = 1, 2$$

entonces el origen del sistema realimentado no forzado:

$$\dot{x} = f(x, 0)$$

es estable y todas las soluciones que comienzan suficientemente cerca del origen permanecen acotadas para todo $t \geq 0$. Si V_1 y V_2 son radialmente no acotadas entonces todas las soluciones son acotadas. Además, el origen es asintóticamente estable en cualquiera de los siguientes casos:

1. $\rho_1 > 0$ y $\rho_2 > 0$.
2. $\rho_1 > 0$, $\epsilon_1 + \delta_2 > 0$ y H_2 es de origen observable.
3. $\rho_2 > 0$, $\epsilon_2 + \delta_1 > 0$ y H_1 es de origen observable.
4. $\epsilon_1 + \delta_2 > 0$, $\epsilon_2 + \delta_1 > 0$ y ambos sistemas de origen observable.

En todos estos casos si V_1 y V_2 son radialmente no acotadas, el origen es atractor global.

Demostración:

Si no hay forzamiento, $u_1 = u_2 = 0$ por lo que $e_1 = -y_2$ y $e_2 = y_1$. Usaremos entonces $V(x) = V_1(x_1) + V_2(x_2)$ como función de Lyapunov del sistema lazo cerrado, se tiene que:

$$\begin{aligned} \dot{V} &= \frac{\partial V_1}{\partial x_1} f_1(x_1, e_1) + \frac{\partial V_2}{\partial x_2} f_2(x_2, e_2) \\ &\leq e_1^t y_1 - \epsilon_1 e_1^t e_1 - \delta_1 y_1^t y_1 - \rho_1 \psi_1(x_1) + e_2^t y_2 - \epsilon_2 e_2^t e_2 - \delta_2 y_2^t y_2 - \rho_2 \psi_2(x_2) \\ &= -y_2^t y_1 - \epsilon_1 y_2^t y_2 - \delta_1 y_1^t y_1 - \rho_1 \psi_1(x_1) + y_1^t y_2 - \epsilon_2 y_1^t y_1 - \delta_2 y_2^t y_2 - \rho_2 \psi_2(x_2) \\ &= -\rho_1 \psi_1(x_1) - \rho_2 \psi_2(x_2) - (\epsilon_2 + \delta_1) y_1^t y_1 - (\epsilon_1 + \delta_2) y_2^t y_2 \leq 0 \end{aligned}$$

Si no hacemos hipótesis sobre las constantes, solo tenemos que $\dot{V} \leq 0$ y por lo tanto el origen es estable y los conjuntos de la forma $\{V(x) \leq c\}$ son positivamente invariantes, por lo que las trayectorias que comienzan suficientemente cerca del origen permanecen acotadas.

Para la estabilidad asintótica aplicamos el Teorema de LaSalle, por lo que hay que probar que $\dot{V}(x) \equiv 0$ implica $x(t) \equiv 0$. En el primer caso, si $\rho_i > 0$ $i = 1, 2$, $\dot{V} = 0$ implica $\psi_1(x_1) = \psi_2(x_2) = 0$ por lo que $x(t) \equiv 0$ gracias a las hipótesis sobre las funciones ψ .

El segundo y tercer caso son simétricos. Por ejemplo, en el segundo caso se tiene que si $\dot{V} = 0$:

$$\begin{aligned} \rho_1 > 0 &\Rightarrow \psi_1(x_1(t)) \equiv 0 \Rightarrow x_1(t) \equiv 0 \\ \epsilon_1 + \delta_2 > 0 &\Rightarrow y_2(t) \equiv 0 \Rightarrow e_1(t) \equiv 0 \end{aligned}$$

De estas ecuaciones se concluye que $y_1(t) = e_2(t) \equiv 0$ por lo que el sistema 2 no tiene entrada externa y además $y_2 \equiv 0$ de donde, por la observabilidad en el origen de H_2 se concluye que $x_2(t) \equiv 0$. Por lo tanto si $\dot{V} \equiv 0$ se tiene que $x(t) \equiv 0$ que es lo que se quería probar.

El último caso es similar, observando que:

$$\begin{aligned} \epsilon_1 + \delta_2 > 0 &\Rightarrow y_2(t) = e_1(t) \equiv 0 \\ \epsilon_2 + \delta_1 > 0 &\Rightarrow y_1(t) = e_2(t) \equiv 0 \end{aligned}$$

y de la observabilidad en el origen de H_1 y H_2 se concluye que $x_1(t) \equiv 0$ y $x_2(t) \equiv 0$.

En el caso de que V_1 y V_2 son radialmente no acotadas, V también lo será, por lo que la estabilidad asintótica del origen será global. \square

Por último enunciamos el siguiente teorema cuya demostración es similar al anterior y puede verse en [2]

Teorema 3.2. Sea H_1 como antes, pasivo con función de almacenamiento V_1 y H_2 una función sin memoria $y_2 = \Psi(e_2)$, también pasiva. Entonces el origen del sistema realimentado no forzado:

$$\dot{x} = f(x, 0)$$

con $x = x_1$ es estable y todas las soluciones que comienzan suficientemente cerca del origen permanecen acotadas para todo $t \geq 0$. Si V_1 es radialmente no acotada entonces todas las soluciones son acotadas. Además, el origen es asintóticamente estable en cualquiera de los siguientes casos:

1. H_1 es pasivo en el estado ($\rho_1 > 0$).
2. $\epsilon_2 + \delta_1 > 0$ y H_1 es de origen observable.

En todos estos casos si V_1 es radialmente no acotada, el origen es atractor global.

A continuación aplicaremos estos resultados para probar estabilidad global del control de congestión en Internet.

4. El control de congestión como un sistema pasivo y estabilidad

Estudiemos entonces las dinámicas introducidas en la sección 2 bajo la óptica de la pasividad. Comencemos por la dinámica primal dada por las ecuaciones (8), cuyo diagrama de bloques puede verse en la figura 6.

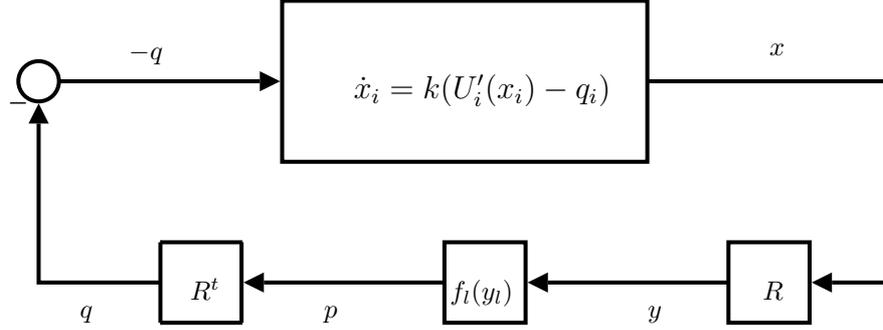


Figura 6: Dinámica primal como un sistema realimentado.

Aquí, H_1 es un sistema diagonal dado por $\dot{x}_i = k_i(U'_i(x_i) - q_i)$, $i = 1, \dots, N$ de entrada $-q$ y salida x . H_2 es una no linealidad dada por R y las funciones de barrera $f_i(p_i)$.

A los efectos de estudiar la pasividad, es conveniente realizar el cambio de variables:

$$\begin{aligned} x &\mapsto x - \hat{x} \\ y &\mapsto y - \hat{y} \\ p &\mapsto p - \hat{p} \\ q &\mapsto q - \hat{q} \end{aligned}$$

para centrar el sistema alrededor del equilibrio. Se tienen entonces las siguientes proposiciones:

Proposición 4.1. *El sistema H_1 de entrada $-(q - \hat{q})$ y salida $x - \hat{x}$ es pasivo en el estado.*

Demostración:

Tomemos como función de almacenamiento:

$$V_1(x - \hat{x}) = \frac{1}{2k} \sum_{i=1}^N (x_i - \hat{x}_i)^2$$

Entonces:

$$\begin{aligned} \dot{V} &= \sum_{i=1}^N (x_i - \hat{x}_i)(U'_i(x_i) - q_i) \\ &= \sum_{i=1}^N (x_i - \hat{x}_i)(U'_i(x_i) - U'_i(\hat{x}_i) + \hat{q}_i - q_i) \\ &= \sum_{i=1}^N (x_i - \hat{x}_i)(U'_i(x_i) - U'_i(\hat{x}_i)) + (x - \hat{x})^t(\hat{q} - q) \\ &\leq (x - \hat{x})^t(\hat{q} - q) \end{aligned}$$

donde en la segunda ecuación se suma y resta $U'_i(\hat{x}_i) = \hat{q}_i$ y el primer término de la tercera ecuación es negativo por ser U_i estrictamente cóncava y por lo tanto U'_i una función decreciente.

Tomando $\rho\psi(x) = -\sum_{i=1}^N (x_i - \hat{x}_i)(U'_i(x_i) - U'_i(\hat{x}_i)) \geq 0$ que es definida positiva (por ser las U'_i estrictamente decrecientes) se tiene la pasividad en el estado. \square

Proposición 4.2. *El sistema H_2 de entrada $x - \hat{x}$ y salida $q - \hat{q}$ es pasivo.*

Demostración:

En este caso, $q - \hat{q} = \Psi(x - \hat{x})$ dada por:

$$q - \hat{q} = R^t \text{diag}(h(R(x - \hat{x})))$$

con $h_l(y - \hat{y}) = f_l(y_l) - \hat{p}$. Entonces:

$$\begin{aligned} (x - \hat{x})^t (q - \hat{q}) &= (x - \hat{x})^t R^t (p - \hat{p}) \\ &= (y - \hat{y})^t (p - \hat{p}) \\ &= \sum_{l=1}^L (y_l - \hat{y}_l) (f_l(y_l) - f_l(\hat{y}_l)) \\ &\geq 0 \end{aligned}$$

donde la desigualdad se debe a que f_l es una función creciente para cada l . \square

Aplicando los resultados de la sección 3 se tiene el siguiente teorema, cuya demostración es es directa del teorema 3.2 y las proposiciones anteriores.

Teorema 4.1. *El equilibrio del controlador (8) es globalmente asintóticamente estable.*

Pasemos ahora a la dinámica dual dada por las ecuaciones (11) y que puede verse en la figura 7

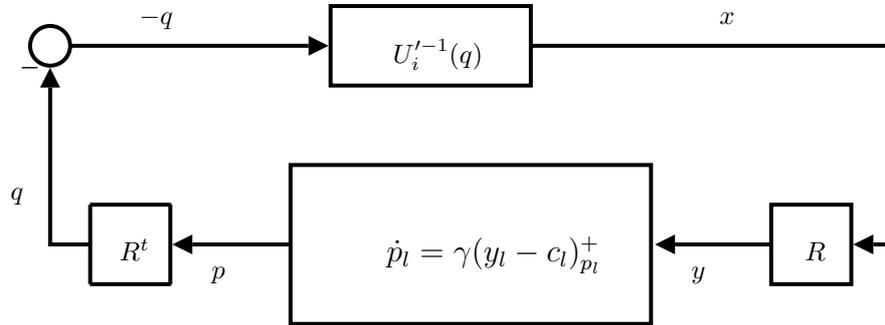


Figura 7: Dinámica dual como un sistema realimentado.

El análisis es similar al anterior, nuevamente hacemos el cambio de variable para centrar el sistema en el equilibrio. Ahora H_1 es una no linealidad y H_2 contiene el ruteo y la dinámica de los enlaces.

Proposición 4.3. *El sistema H_1 de entrada $-(q - \hat{q})$ y salida $x - \hat{x}$ es pasivo.*

Demostración:

Se tendrá que:

$$-(q - \hat{q})^t (x - \hat{x}) = -\sum_{i=1}^N (q_i - \hat{q}_i) (U'_i^{-1}(q_i) - U'_i^{-1}(\hat{q}_i)) \leq 0$$

donde hemos usado que U' es monótona decreciente (por ser U cóncava) y por lo tanto su inversa U'^{-1} también lo es. \square

Proposición 4.4. *El sistema H_2 de entrada $x - \hat{x}$ y salida $q - \hat{q}$ es pasivo.*

Demostración:

Tomemos como función de almacenamiento:

$$V_2(p) = \frac{1}{2\gamma} \sum_{l=1}^L (p_l - \hat{p}_l)^2$$

Entonces:

$$\dot{V} = \sum_{l=1}^L (p_l - \hat{p}_l)(y_l - c_l)_{p_l}^+$$

La demostración reposa entonces en dos desigualdades. La primera, debida a la saturación:

$$(p_l - \hat{p}_l)(y_l - c_l)_{p_l}^+ \leq (p_l - \hat{p}_l)(y_l - c_l)$$

ya que si la saturación no está activa, son iguales, y si está activa entonces el término de la izquierda es cero y el de la derecha es $-\hat{p}_l(y_l - c_l) \geq 0$ ya que $\hat{p}_l \geq 0$ y si la saturación actúa es porque $y_l < c_l$, por lo que el producto de ambos es positivo.

La segunda desigualdad es:

$$(p_l - \hat{p}_l)(y_l - c_l) \leq (p_l - \hat{p}_l)(y_l - \hat{y}_l)$$

la cual surge de la “complementary slackness” de p e y . Si $\hat{y}_l = c_l$ entonces se da la igualdad. Si $\hat{y}_l < c_l$ entonces $\hat{p}_l = 0$ por lo que el término de la derecha es mayor ya que $p_l \geq 0$.

Juntando estas desigualdades se tiene que:

$$\begin{aligned} \dot{V} &= \sum_{l=1}^L (p_l - \hat{p}_l)(y_l - c_l)_{p_l}^+ \\ &\leq \sum_{l=1}^L (p_l - \hat{p}_l)(y_l - \hat{y}_l) \\ &= (p - \hat{p})^t (y - \hat{y}) \\ &= (p - \hat{p})^t R(x - \hat{x}) \\ &\quad (q - \hat{q})^t (x - \hat{x}) \end{aligned}$$

lo que muestra que es pasivo. □

Observemos que no se puede establecer que el sistema sea estrictamente pasivo en el estado por lo que para probar estabilidad debemos refinar el argumento anterior.

Se tendrá entonces el siguiente teorema.

Teorema 4.2. *El equilibrio del controlador (11) es globalmente asintóticamente estable.*

Demostración:

Que el equilibrio es estable y las trayectorias son acotadas es consecuencia directa del teorema 3.2 y las proposiciones anteriores, ya que es realimentación de dos sistemas pasivos.

Para la estabilidad asintótica del equilibrio debemos utilizar el Teorema de LaSalle. Damos aquí solamente la idea de la demostración, que puede verse en [8]. Tomando como función de Lyapunov la función de almacenamiento de la dinámica dual se tiene que $\dot{V} = 0 \Leftrightarrow p_l = \hat{p}_l$ o $y_l = c_l$. Además $x_i = U_i^{\prime-1}(q_i)$. Observemos que estas son las condiciones KKT de optimalidad del problema (10) que caracteriza (\hat{x}, \hat{p}) , el único (c.f. 2.1) punto de equilibrio del sistema. □

Por último, combinando las proposiciones 4.1 y 4.4 se tiene el siguiente teorema.

Teorema 4.3. *El equilibrio de la dinámica primal-dual (12a) es globalmente asintóticamente estable.*

Demostración:

La dinámica primal dual es la realimentación negativa de los sistemas H_1 de 4.1 y H_2 de 4.4, que son pasivos, por lo que el sistema lazo cerrado es estable. Para la estabilidad asintótica basta tomar como función de Lyapunov la suma de las funciones de almacenamiento:

$$V(x, p) = \frac{1}{2k} \sum_{i=1}^N (x_i - \hat{x}_i)^2 + \frac{1}{2\gamma} \sum_{l=1}^l (p_l - \hat{p}_l)^2$$

y observar que $\dot{V} = 0 \Leftrightarrow x_i = \hat{x}_i, p_l = \hat{p}_l \forall i, l$ que es el equilibrio deseado. □

5. Control del número de conexiones. Estabilidad y justicia.

Desde el trabajo ya citado de Kelly [5], la noción de justicia en el reparto de recursos que se realiza en una red ha quedado asociada al modelo propuesto en el problema de optimización (10). Este modelo, eligiendo funciones de utilidad apropiadas permite tener en cuenta diferentes nociones de justicia y el análisis de las secciones anteriores muestra cómo es posible diseñar controladores que hagan operar a la red en el óptimo de (10).

Sin embargo, si bien desde el punto de vista teórico el modelo resulta satisfactorio, en la práctica existen algunas limitaciones que hacen que la red no se comporte exactamente según este modelo.

La primera limitante es que el protocolo TCP más utilizado actualmente, TCP Newreno, utiliza un mecanismo de ventana como el que se describió en la sección 2, y este mecanismo determina completamente la función de utilidad. Por lo tanto, no es posible obtener diferentes nociones de justicia sin cambiar las implementaciones de TCP, algo difícil de hacer a escala global. A su vez, la interacción entre TCP diferentes puede llevar a injusticias que no son deseables en la red. Esto ha limitado el despliegue incremental de nuevos mecanismos de control de congestión.

La segunda limitación es una consecuencia de la implementación actual de TCP: como vimos, la función de utilidad que modela el TCP Newreno está dada por:

$$U(x) = -\frac{1}{T^2x}$$

siendo T el retardo de ida y vuelta de la conexión. Esto implica que conexiones de más retardo pesarán menos en la función de utilidad y por consiguiente obtendrán menos recursos. Para mitigar este efecto, surgieron diferentes mecanismos en los que los clientes utilizan *múltiples* conexiones para aumentar el ancho de banda total obtenido.

5.1. Modelo de múltiples conexiones.

Incluyamos el problema de las múltiples conexiones en el modelo. La asignación deseada corresponde al siguiente problema de optimización:

$$\begin{aligned} \text{máx} \quad & \sum_i U_i(\rho_i) \\ \text{s.t.} \quad & R\rho \leq c \end{aligned} \tag{17}$$

donde ρ_i es el *rate total* obtenido por el usuario i , y $U_i(\rho)$ es la función de utilidad asociada a dicho usuario. U debe ser una función cóncava y en general se utiliza la siguiente familia, denominada α -utilidades:

$$\begin{aligned} U_i(\rho) &= k_i \frac{\rho^{1-\alpha}}{(1-\alpha)}, \quad \alpha > 0, \quad \alpha \neq 1 \\ U_i(\rho) &= k_i \log(\rho), \quad \alpha = 1 \end{aligned}$$

La ventaja de esta familia es que permite modelar diferentes nociones de justicia. Por ejemplo, el caso $\alpha \rightarrow 0$ corresponde a maximizar el throughput de la red, mientras que el caso $\alpha \rightarrow \infty$ corresponde a la asignación denominada *max-min fairness*, que es la asignación más justa un sentido económico. El caso $\alpha = 1$ se encuentra a mitad de camino y se denomina *proportional fairness*. Observemos que el modelo de TCP corresponde a tomar $\alpha = 2$.

Ahora bien, si los usuarios pueden abrir múltiples conexiones el rate total ρ_i corresponderá al agregado de las mismas. Si suponemos que cada conexión tiene la misma función de utilidad $U_{i,TCP}(x)$ modelando el TCP, x_i es el rate de cada conexión individual y n_i es el número de conexiones del usuario i , entonces la asignación de recursos que hará la red será tal que:

$$\begin{aligned} \text{máx} \quad & \sum_i n_i U_{i,TCP}(x_i) \\ \text{s.t.} \quad & \sum_i R_{li} n_i x_i \leq c \end{aligned} \tag{18}$$

y el rate obtenido por el usuario será $\rho_i = n_i x_i$.

El problema (18) es simplemente la adaptación del problema (10) al caso en que un usuario abre múltiples conexiones con la misma función de utilidad.

Nos encontramos entonces en una situación en que la red impone su propia noción de justicia, dada por las $U_{i,TCP}$. A su vez, existe una presión de parte de los usuarios a abrir más conexiones, ya que si el usuario

i abre una conexión extra, su porción de la torta en el reparto se vuelve mayor. Esto lleva a una situación del tipo “tragedia de los comunes” en el que no hay incentivos para que los usuarios mantengan individualmente acotados su número de conexiones.

5.2. Control propuesto.

Frente a este problema, en el artículo [12] avanzamos hacia la proposición de un mecanismo de *control del número de conexiones* que permita imponer la noción de justicia de (17) pero manteniendo el funcionamiento actual de TCP. A continuación se describe brevemente la propuesta.

La idea es utilizar el número de conexiones como variable de control, y comparar el rate actual ρ_i que obtiene el usuario con lo que le correspondería según el precio actual de congestión q_i que ve el mismo usuario, dado por $U_i^{\prime-1}(q_i)$. Si este rate está por debajo de lo que corresponde, se le permite al usuario incrementar su número de conexiones. En caso contrario se disminuye, como forma de disminuir ρ_i . Supondremos para simplificar que TCP sigue la dinámica dual, es decir, los precios se actualizan según la ecuación (11) y el rate de TCP se adapta de manera instantánea a dicho precio, tomando $x_i = U_{i,TCP}^{\prime-1}(q_i) = f_{i,TCP}(q)$. La dinámica completa resulta:

$$\dot{n}_i = \beta (U_i^{\prime-1}(q_i) - \rho_i), \quad \beta > 0 \quad (19a)$$

$$x_i = f_{i,TCP}(q_i) \quad (19b)$$

$$\rho_i = n_i x_i \quad (19c)$$

El diagrama de bloques de este nuevo lazo puede verse en la figura 8.

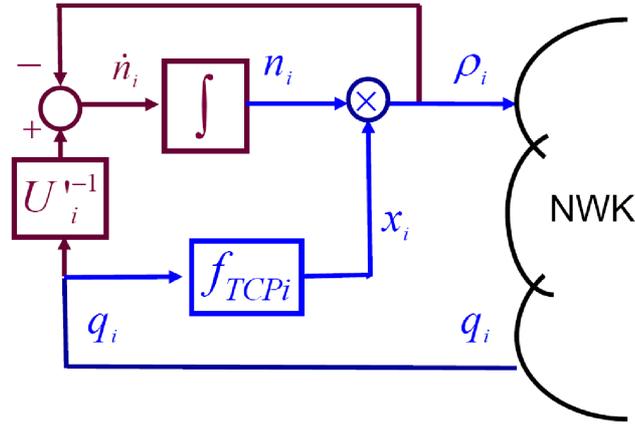


Figura 8: Diagrama de bloques del control del número de conexiones.

Aquí, la parte en azul corresponde al lazo de realimentación de TCP, y su entrada es el número de flujos. La propuesta es agregar un nuevo lazo (dibujado en violeta) en la escala del usuario para adaptar el número de flujos a la situación de congestión, de manera de mover la asignación de recursos de la red al punto justo.

Un punto de equilibrio $(n^*, x^*, \rho^*, y^*, p^*, q^*)$ del sistema completo satisface:

$$y^* = R\rho^*, q^* = R^T p^*, \rho_i^* = n_i^* x_i^*$$

y también:

$$U_i^{\prime}(\rho_i^*) = q_i^*, \quad (20)$$

$$p_i^*(c_i - y_i^*) = 0, \quad (21)$$

$$U_{TCPi}^{\prime}(x_i^*) = q_i^*. \quad (22)$$

(20) y (21) implican las condiciones KKT de optimalidad de (17). A su vez, (21) y (22) implican las condiciones KKT de (18), lo cual sale de considerar el Lagrangeano correspondiente (c.f. [12] para los detalles). Es decir, el equilibrio de esta dinámica es el óptimo del problema de reparto (17), pero a su vez se logra utilizando la asignación de recursos (18) que produce TCP, encontrando el número de flujos por usuario que hace coincidir ambos problemas. Una ventaja es que esto puede hacerse utilizando *el mismo precio de congestión* que proporciona la red, por lo cual no hay necesidad de nuevas señales de realimentación.

5.3. Análisis de estabilidad por pasividad.

Resta analizar la estabilidad de esta nueva dinámica y para ello utilizaremos argumentos de pasividad. Se tiene el siguiente teorema:

Teorema 5.1. *El equilibrio del control de número de flujos dado por (19) complementado por la dinámica dual en los enlaces es localmente asintóticamente estable*

Demostración:

La idea de la prueba reposa en observar que la porción de red del sistema $\rho - \rho^* \mapsto q - q^*$ no fue modificada, por lo que aplicando la proposición 4.4, este sistema es pasivo. Probaremos que la linealización del sistema $-(q - q^*) \mapsto \rho - \rho^*$ dado por (19) es pasiva. Observemos que como el sistema trabaja de forma descentralizada, basta con probar que el control de cada usuario es pasivo, ya que en ese caso existirá una función de almacenamiento $V_i(n_i)$ para cada usuario y la suma de estas funciones será función de almacenamiento del sistema completo.

Sean δn , δq , $\delta \rho$ y δx las desviaciones de las variables de interés respecto al equilibrio, se tendrá entonces que:

$$\begin{aligned}\delta \rho &= n^* \delta x + x^* \delta n \\ \delta x &= f'_{i,TCP}(q^*) \delta q = -b \delta q \\ s \delta n &= \beta((U'^{-1})'(q^*) \delta q - \delta \rho) = \beta(-a \delta q - \delta \rho)\end{aligned}$$

con $a, b > 0$ debido a que las funciones involucradas son decrecientes. Al despejar se obtiene:

$$\frac{\delta \rho}{\delta q} = \frac{n^* b s + \beta x^* a}{s + \beta x^*} = -G(s)$$

siendo $G(s)$ la transferencia de un compensador de adelanto atraso con cero y polo en el semiplano izquierdo. Por lo tanto, $\Re(G(j\omega)) > 0$ y aplicando el criterio del corolario 3.1 se tiene que el sistema linealizado de $-\delta q \mapsto \delta \rho$ es estrictamente pasivo, por lo que el lazo cerrado será asintóticamente estable. \square

Observación 5.1. En [12] puede verse también la generalización de este resultado al caso en que se cambia la dinámica de TCP por una dinámica no instantánea del tipo primal-dual. En este caso $G(s)$ ya no es un lead-lag pero sigue valiendo que $\Re(G(j\omega)) > 0$ por lo que se tiene el mismo resultado.

Este resultado induce a pensar que, utilizando argumentos de pasividad, sería posible probar que el sistema sin linealizar es pasivo, y por lo tanto el lazo cerrado es globalmente asintóticamente estable. Sin embargo, hasta el momento de escribir este trabajo no ha sido posible encontrar una función de almacenamiento apropiada para probar el resultado global.

5.4. Simulaciones.

Para finalizar, se presentan a continuación un par de escenarios de simulación donde se observa que el control propuesto logra cumplir con los objetivos planteados.

En primer lugar consideramos el caso de un enlace cuello de botella de capacidad 1 en el que convergen cuatro usuarios diferentes. Cada uno de estos usuarios abre conexiones de diferentes retardos (RTTs), y por lo tanto sus TCP logran diferentes velocidades. Más exactamente las utilidades TCP de los usuarios $i = 1, \dots, 4$ son de la forma:

$$U_{i,TCP}(x_i) = -\frac{1}{2T_i^2 x_i}$$

con $T_1 = 200ms$, $T_2 = 300ms$, $T_3 = 400ms$ y $T_4 = 800ms$. Si cada usuario abriera un único flujo la asignación resultante del problema (18) sería:

$$x_1^* = 0,41, \quad x_2^* = 0,28, \quad x_3^* = 0,21, \quad x_4^* = 0,10$$

y esta es la evolución que se observa en la figura 9.

La dinámica (19) permite corregir este sesgo hacia los usuarios de menor retardo. Para ello, elegimos como utilidad del usuario $U_i(\rho_i) = \log \rho_i$, la misma para todos los usuarios. La asignación óptima del problema (17) resulta entonces:

$$\rho_1^* = \rho_2^* = \rho_3^* = \rho_4^* = 0,25$$

y para lograr esto se debe compensar al usuario lento permitiéndole abrir más flujos. El resultado puede verse en la figura 10.

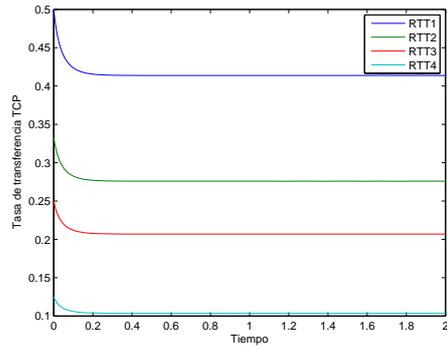


Figura 9: Reparto de recursos realizado por TCP.

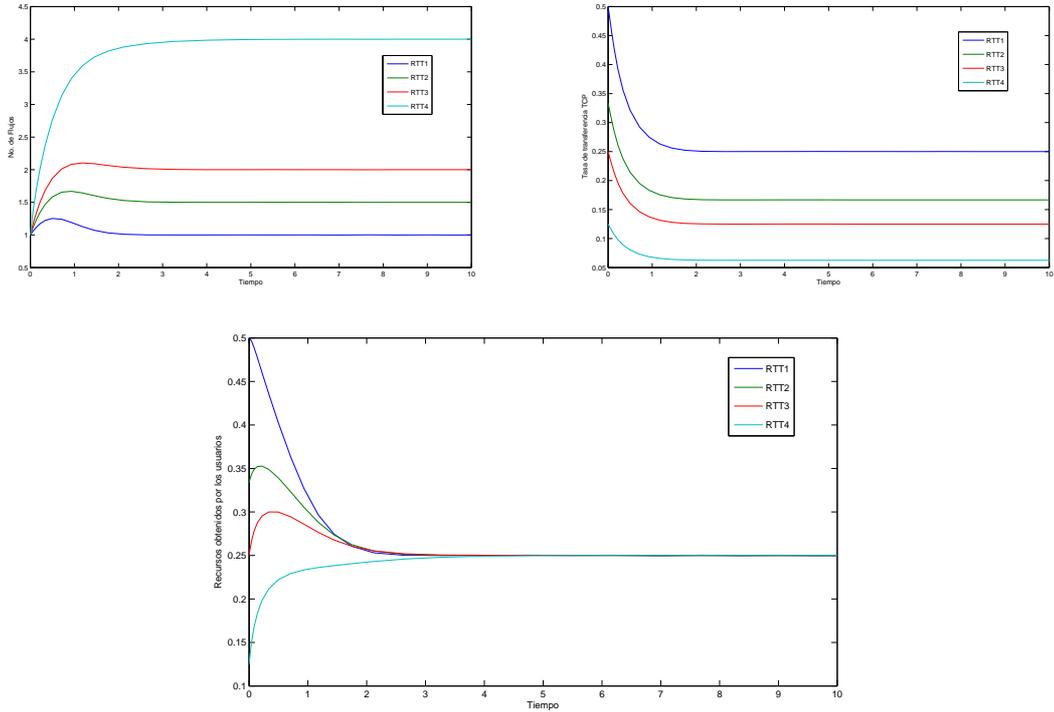


Figura 10: Reparto de recursos obtenido mediante el control de n . Las gráficas representan la evolución de n , x y ρ .

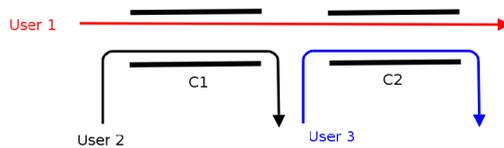


Figura 11: Red lineal del segundo ejemplo.

Exhaustivas pruebas de simulación muestran además que el equilibrio resulta globalmente estable, aunque como ya comentamos, no se tiene una prueba completa de este resultado.

Se considera ahora un segundo caso correspondiente a una red lineal como la de la figura 11.

Aquí tomamos $C_1 = 2$ y $C_2 = 1$ y el retardo de la ruta larga se tomó como el doble de los usuarios de rutas cortas, para simular una situación cercana a la realidad. Como antes, si cada fuente abre un único flujo

el óptimo del problema (18) resulta sesgado hacia las fuentes cortas y es:

$$x_1^* = 0,31, \quad x_2^* = 1,68, \quad x_3^* = 0,68$$

y la evolución del sistema hacia este punto puede verse en la primera gráfica de la figura 12.

Para corregir este sesgo, elegimos las utilidades de los usuarios $U_i(\rho_i)$ en la familia de las α -utilidades con $\alpha = 5$ lo que permite aproximar el comportamiento de la asignación “max-min”. Esta asignación es la más justa y resulta:

$$\rho_1^* = \rho_3^* = 0,5, \quad \rho_2^* = 1,5$$

es decir, el enlace más restrictivo que es el 2 es compartido de manera equitativa por los usuarios 1 y 3. Como el usuario 1 queda limitado por este enlace, se le permite utilizar al usuario 2 el ancho de banda remanente en el enlace 1. El resultado de aplicar el control propuesto en (19) puede verse en la segunda y tercera gráficas de la figura 12.

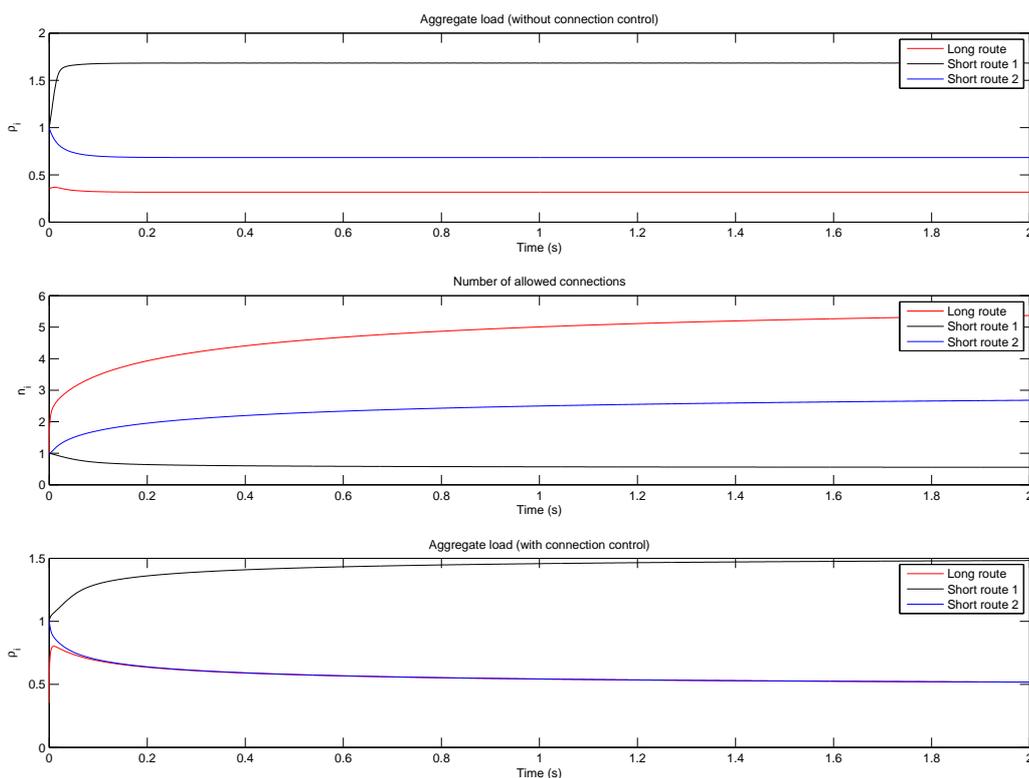


Figura 12: Control de número de flujos para lograr justicia max-min.

Como se observa, mediante el mecanismo propuesto se logra de manera práctica restablecer la justicia en la red. En la sección de conclusiones comentaremos algunos aspectos de la solución propuesta y posibles aproximaciones a la implementación.

6. Conclusiones y trabajo futuro.

En este trabajo se han presentado diferentes mecanismos de control de congestión, un problema de vital importancia para el buen funcionamiento de las redes de datos modernas. El crecimiento de la demanda por ancho de banda hace que sea necesario revisar los algoritmos de reparto de recursos diseñados para redes pequeñas y con poco tráfico, modelarlos, comprender el reparto de recursos que logran y estudiar su estabilidad.

La aproximación por pasividad ha permitido dar un marco común a todas las demostraciones que habían aparecido en la literatura del control de congestión. En este trabajo se han presentado las nociones de pasividad involucradas y algunos de los resultados de estabilidad de control de congestión relevantes asociados.

Por último, se presentó un modelo de control de número de conexiones que permite lograr el reparto de recursos propuesto en [5] pero sin tocar las implementaciones actuales del TCP. En dicho modelo se consideró al número de conexiones como fluido, es decir, se ha permitido que un usuario abra un número no entero de conexiones con el objetivo de llevar al sistema al equilibrio óptimo. Obviamente esto no es aplicable en la práctica, pero tanto en [12] como en trabajos posteriores los autores proponen un algoritmo de *control de admisión* de las conexiones que aproxima la ley fluida de manera apropiada y permite hacer funcionar al sistema alrededor del equilibrio deseado.

Sobre el trabajo futuro, una línea clara a seguir es intentar demostrar la pasividad del sistema (19) sin linealizar, para lograr probar estabilidad global. Con ciertas hipótesis adicionales, los autores en [12] probaron estabilidad global en el caso de un enlace, pero la prueba general permanece abierta. Es de destacar que el elemento más complejo en las demostraciones parece ser la ecuación $\rho = nx$. Esta no linealidad en forma de producto no permite utilizar la mayoría de los enfoques tradicionales de la teoría de control.

A su vez, resta trabajo por hacer en la extensión de este tipo de resultados y algoritmos a conexiones multi-camino, que pueden ser de interés para controlar el funcionamiento de redes peer-to-peer, cuya demanda de ancho de banda es la dominante en la Internet actual.

Referencias

- [1] V. Jacobson, *Congestion avoidance and control*, in Proc. of ACM/SIGCOMM '88, pp 314-329.
- [2] H. K. Khalil, *Nonlinear systems*, 2nd. Ed., Prentice-Hall, 1996.
- [3] IETF Request for Comments 2001, *TCP slow start, congestion avoidance, fast retransmit and fast recovery algorithms*, 1997.
- [4] M. Mathis, J. Semke, J. Mahdavi, T. Ott, *The macroscopic behavior of the TCP Congestion avoidance algorithm*, Computer Communications Review, Vol. 27 (3), julio 1997.
- [5] F.P. Kelly, A. Maulloo, D. Tan, *Rate control in communication networks: shadow prices, proportional fairness and stability*, Journal of the Operational Research Society, 49 (1998), 237-252.
- [6] S.H. Low, D.E. Lapsley, *Optimization flow control I: basic algorithm and convergence*, IEEE/ACM Transactions on Networking (1999), 861-875.
- [7] S.H. Low, F. Paganini, J. Doyle, *Internet Congestion Control: An Analytical Perspective*, IEEE Control Systems Magazine (2002).
- [8] F. Paganini, *A global stability result in network flow control*, Systems and Control Letters, vol. 46, pp. 165-172, 2002.
- [9] A. Tanenbaum, *Computer Networks*, 4a. edición, Prentice-Hall, 2002.
- [10] J. Wen, M. Arcak, *A Unifying Passivity Framework for Network Flow Control*, IEEE Transactions on Automatic Control, 49 (2), 2004.
- [11] S. Boyd, L. Vanderberghe, *Convex optimization*, Cambridge University Press, 2004.
- [12] A. Ferragut, F. Paganini, *Achieving network stability and user level fairness through admission control of TCP connections*, Conference on Information Systems and Sciences, Princeton, NJ, March 2008.