

# Predicción de tráfico de Internet y aplicaciones

I.Bernal, J.Arakil, D.Morato, M.Izal, E.Magaña y L.A. Díez  
Departamento de Automática y Computación. Universidad Pública de Navarra  
Grupo de Redes, Sistemas y Servicios telemáticos  
Campus Arrosadía - 31006 Pamplona (Navarra)  
Teléfono: 948 168904 Fax: 948 168924  
E-mail: javier.aracil@unavarra.es

**Abstract** *In this paper we focus on traffic prediction as a means to achieve dynamic bandwidth allocation in a generic Internet link. Our findings show that coarse prediction (bytes per interval) proves advantageous to perform dynamic link dimensioning, even if we consider a part of the top traffic producers in the traffic predictor.*

## 1 Introducción

Hoy en día estamos asistiendo a un crecimiento imparable del tráfico de Internet. Ante tal demanda es un hecho que las operadoras desean dar calidad de servicio a sus usuarios y para ello es preciso dimensionar los enlaces. Los problemas del tráfico de Internet son muy diferentes de los de otros tipos de tráfico [1] y plantean un escenario de especial complejidad. En concreto tenemos que el tráfico de Internet presenta autosimilitud (*self-similarity*) y no estacionariedad.

Por el contrario, el tráfico telefónico es de *incrementos independientes* y por lo tanto aplican modelos de tipo  $/G/G/1$ . En el entorno de redes de banda ancha el ancho de banda efectivo se calcula con modelos de Markov en varias escalas de tiempo. Pero sin embargo, debido a la fuerte no estacionariedad y autosimilitud del tráfico de Internet, no existe hoy en día una teoría de dimensionamiento de enlaces de Internet.

### 1.1 Autosimilitud

Para entender bien lo que significa autosimilitud es necesario repasar conceptos básicos de independencia estadística. Sea  $X_1, X_2, \dots, X_n$  una muestra de  $n$  variables aleatorias independientes con media  $\mu$  y desviación estándar  $\sigma$ . Se cumple que:

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i \quad (1)$$

$$\text{var}(\bar{X}) = \sigma^2 n^{-1} \quad (2)$$

Por otro lado, sea el proceso de cuentas de paquetes  $X_i$  en intervalos de duración  $\delta$ , que mostramos en la figura 1.

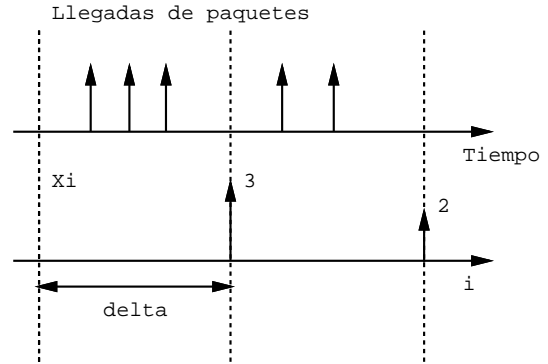


Figura 1: Cuentas por intervalos

Formamos ahora el proceso de agregación de intervalos (media muestral)

$$S_i^m = \frac{X_{im-m+1} + \dots + X_{im}}{m} \quad i = 1 \dots \lfloor \frac{n}{m} \rfloor \quad (3)$$

y observamos que no cumple el Teorema Central del límite en el caso de tráfico Internet, según el cual la varianza debe decaer con el número de muestras  $m$  en una proporción  $m^{-1}$ . La figura 2 muestra la varianza frente al nivel de agregación en coordenadas logarítmicas en ambos ejes para una traza real y para un caso de incrementos independientes. Observamos que la varianza decae *lentamente* en comparación con un proceso de incrementos independientes, con la forma:

$$\text{Var}(S_i^m) = \sigma^2 m^{-\beta} \quad (4)$$

con  $0 < \beta < 1$ .

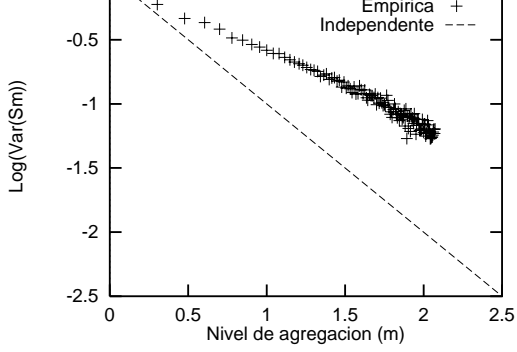


Figura 2: Varianza frente a nivel de agregación

La varianza decae de forma lenta simplemente porque no se cumplen las hipótesis del teorema central del límite. En concreto, las variables aleatorias que contribuyen a la media muestral (bytes por intervalo) no son independientes. Por el contrario, se cumple que el proceso de llegadas de tráfico es asintóticamente autosimilar de segundo orden. Sea  $\rho^{(m)}(j)$  la autocorrelación (*lag*  $j$ ,  $j > 1$ ) de  $S_i^m$ . Tenemos que:

$$\lim_{m \rightarrow \infty} \rho^{(m)}(j) = \frac{1}{2}((j+1)^{2H} - 2j^{2H} + (j-1)^{2H}) \quad (5)$$

Un proceso asintóticamente autosimilar de segundo orden sufre *dependencia a largo plazo*. Esto es, la autocorrelación del proceso decae lentamente y no es sumable, en contraste con procesos Poissonianos/Markovianos. Es importante observar que la dependencia a largo plazo es una propiedad asintótica: no importan los valores *absolutos* de la autocorrelación sino la *forma* de la misma.

En la figura 3 mostramos el efecto de la dependencia a largo plazo. En las gráficas de la izquierda tenemos tráfico de internet frente al tráfico de Poisson en las de la derecha, para varias escalas de tiempo de 10ms., 100ms. y 1 seg (ver ecuación 3). La caída lenta de la varianza provoca *ráfagas en cualquier escala de tiempo*, al contrario que un proceso de Poisson donde tenemos una suavización hacia la media conforme vamos agregando bytes por intervalo y formadon así el tráfico en escalas de tiempo mayores.

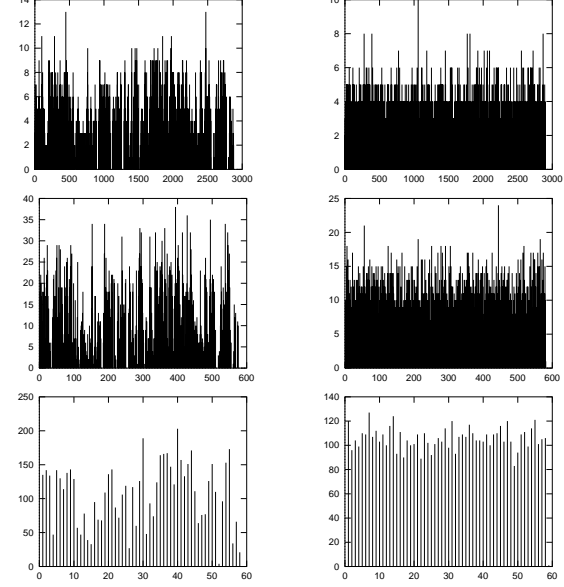


Figura 3: Trafico de Internet frente a Poisson

## 1.2 Causas de la dependencia a largo plazo

La dependencia a largo plazo se produce por el efecto del multiplex de fuentes on-off con varianza infinita [2], como se muestra en la figura 4. Cada una de estas ráfagas (conexiones TCP, por ejemplo) introduce correlación en la escala de tiempo de su duración. Estudios experimentales demuestran que las ráfagas que vienen de la transmisión de ficheros en Internet tienen varianza infinita [3].

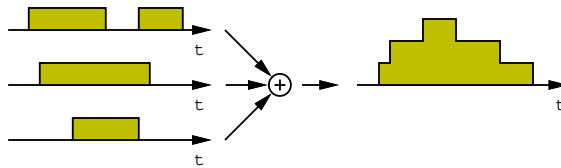


Figura 4: Superposición de fuentes on-off

Para que la varianza de la duración de la ráfaga sea infinita la distribución de la misma debe seguir la forma:

$$P(X > t) \sim Kt^{-\alpha} \quad 1 < \alpha < 2 \quad (6)$$

que se observa que es el caso para conexiones reales de Internet, como mostramos en la figura 5. En esta figura se muestra la distribución (función de supervivencia) de la duración de conexiones FTP.

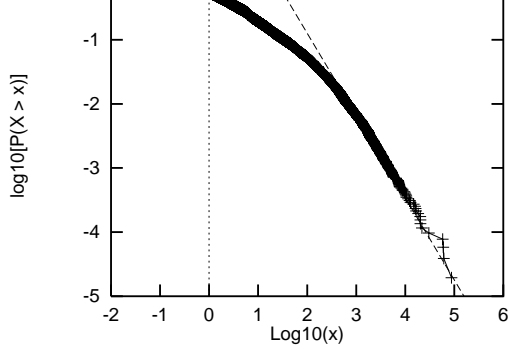


Figura 5: Duracion de conexiones FTP

### 1.3 No estacionariedad

Por otro lado, el tráfico de Internet adolece de una fuerte no estacionariedad, como se observa en la figura 6, que muestra varias escalas de tiempo de una traza de tráfico real. En conclusión, las características de alta intermitencia (dependencia a largo plazo) y no estacionariedad hacen que el dimensionamiento a-priori de enlaces de Internet sea difícil de realizar en la práctica.

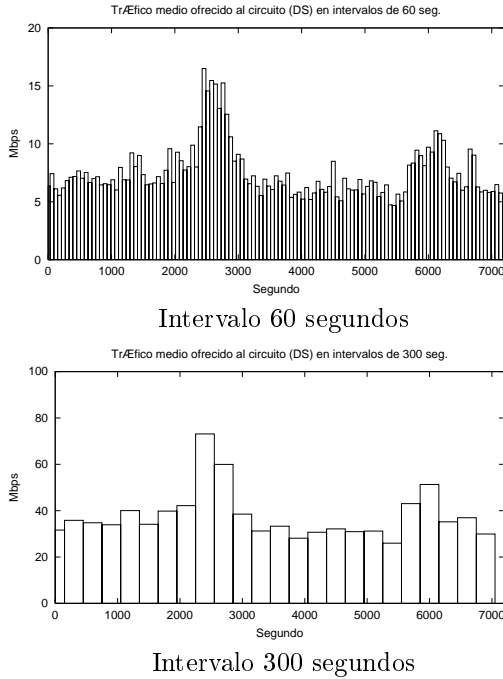


Figura 6: No estacionariedad

## 2 Planteamiento del problema

De los apartados anteriores observamos que el tráfico de Internet muestra dependencia y no estacionariedad. El modelado solo es posible en estadísticos de primer y segundo orden y eso no es suficiente para una correcta estimación de los recursos a asignar. Ante estas condiciones de tráfico fuertemente dinámico cabe pensar en otros métodos que hagan del problema de la correlación una ventaja. En concreto, los algoritmos de predicción

alta correlación.

Sea  $\delta$  el intervalo de tiempo de predicción y sean  $X_k$  y  $\hat{X}_k$  los tráficos (número de bytes) real y estimado en el intervalo. Probaremos una serie de estimadores de implementación sencilla para  $\hat{X}_k$ . En concreto, utilizaremos el *Método interpolador de Lagrange*, ya que es un estimador lineal simple que se puede usar con intervalos de longitud constante (como es nuestro caso), que sigue la expresión :

$$p(x) = \sum_{i=1}^{n+1} y_i \prod_{\substack{j=1 \\ j \neq i}}^{n+1} \frac{x - x_j}{x_i - x_j} \quad (7)$$

Particularizaremos para obtener polinomios de orden  $n$  con  $n < 3$  y obtenemos ecuaciones sencillas lineales, de coste computacional muy reducido para un hipotético asignador de ancho de banda localizado en un router o conmutador que gobierna un enlace.

Para medir la bondad del estimador podemos usar la distribución de probabilidad del error  $\hat{X}_k - X_k$ . Pero es todavía más interesante el retardo en cola para un servidor con capacidad variable  $\frac{X_k}{\delta}$ . Esta última medida no sólo tiene en cuenta el error instantáneo sino también el acumulado y modela mejor un escenario real de predicción.

Por otro lado, es interesante estudiar no solo el caso de predicción con el total del tráfico sino predicción basada en un subconjunto de usuarios. Esto puede ser muy flexible en el caso de topologías de red donde todo el tráfico no pasa por un solo punto. Las fuentes más activas pueden informar a los routers en el camino extremo a extremo del tráfico que van a enviar, en sintonía con estándares recientes de conmutación por etiquetas [5]. De este modo la predicción no sólo es útil en enlaces de acceso sino en topologías genéricas de red. De hecho la fuerte no homogeneidad de los usuarios ayuda a predecir en base a *parte* del tráfico. La figura 7 muestra el porcentaje de tráfico del enlace frente al porcentaje de usuarios que lo producen.

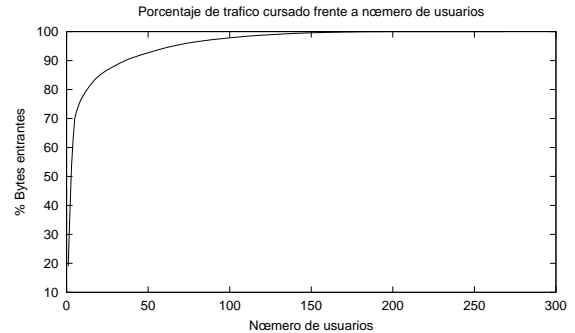


Figura 7: Porcentaje de tráfico frente a porcentaje usuarios que lo generan

Se observa claramente que es posible predecir con un porcentaje pequeño de usuarios y no con

regulares. La figura 8 muestra el usuario más activo de la muestra. Prácticamente transmite a tasa constante.

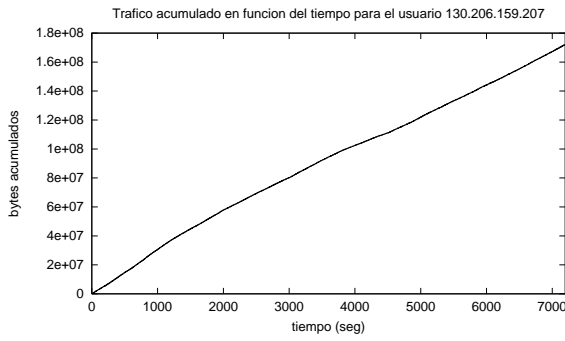
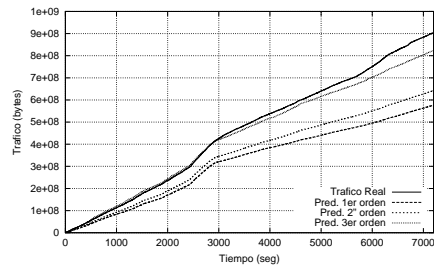


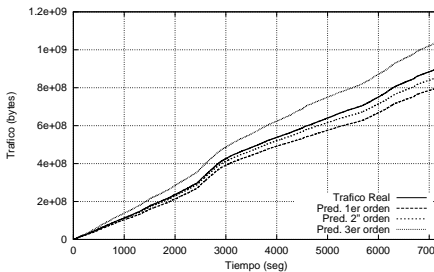
Figura 8: Usuario más activo

### 3 Resultados

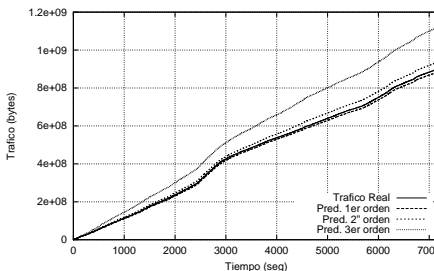
Los resultados preliminares se obtienen con Polinomios interpoladores de Lagrange de primer,segundo y tercer orden en intervalos de 1 segundo. Este intervalo de un segundo es un tiempo superior a RTTs típicos en la Internet y permite que el asignador de ancho de banda (un conmutador ATM con ABR por ejemplo) tenga tiempo suficiente para adecuar las condiciones del circuito a la nueva carga de tráfico. En primer lugar la figura 9 muestra una comparación visual de tráfico real frente a tráfico obtenido mediante predicción.



Predicción con 4 usuarios



Predicción con 30 usuarios



Predicción con 100 usuarios

Se observa que la predicción con un número reducido de usuarios (30 a 100) obtiene buenos resultados. El número total de usuarios en la muestra es de 300. Este resultado se relaciona perfectamente con el resultado observado en la figura 8, en la cual podemos observar como los 30 usuarios más activos, generan más del 80% del tráfico total. De hecho se puede observar como la distribución de probabilidad error de predicción, se estabiliza bastante al utilizar un número de usuarios para la predicción superior a 30, como se muestra en la figura 10

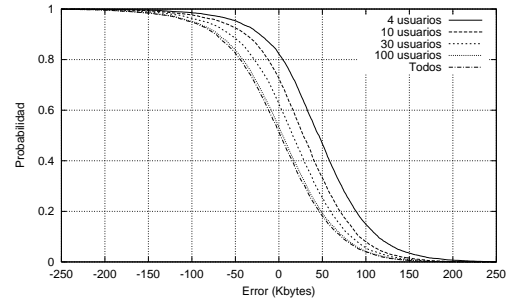
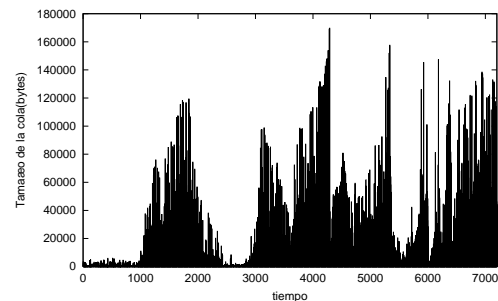


Figura 10: Función de densidad del error de predicción

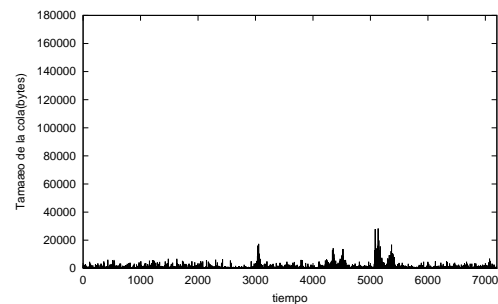
Más interesante todavía es el número de usuarios en cola, obtenido mediante simulación de la ecuación de Lindley:

$$Q_{n+1} = \max \{Q_n + A_{n+1} - C_{n+1}, 0\} \quad (8)$$

donde  $Q_n$  es el número de bytes en cola intervalo  $n$ ,  $A_n$  es el número de bytes que llegan en intervalo  $n$  y  $C_n$  es la capacidad del servidor en el instante  $n$ .



N=10 usuarios



N=30 usuarios

Figura 11: Número de usuarios en cola

En nuestro caso  $C_n = A_n$ . La figura 11 muestra el número de bytes en cola  $Q_n$  en el intervalo de medida. Se observa que con 30 usuarios se consigue estabilizar la cola en torno a valores muy bajos (menores de 30 KBytes), mostrando la viabilidad práctica de la idea.

En la figura 12 comparamos la predicción del tráfico real con un FBM (Fractional Brownian Motion) [4], que es un proceso gaussiano asintóticamente autosimilar de segundo orden, muy utilizado para modelar tráfico de Internet. La figura 12 muestra en este caso la función de supervivencia  $P(X > x)$  del retardo en cola.

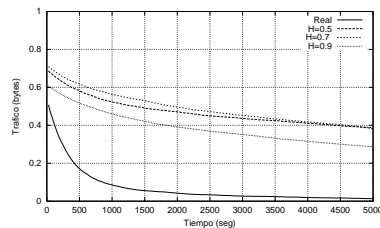


Figura 12: Comparación con un FBM

Observamos que los FBM modelan bien la dependencia a largo plazo *pero no a corto plazo*. Este resultado apunta en la dirección de modelado con fuentes de dependencia a corto plazo, tipo cadenas de Markov para aquellos escenarios donde el parámetro de relevancia sea la dependencia a corto y no a largo plazo.

Los resultados anteriores muestran el caso en que se predice el número de bytes que llegan en un intervalo. Es el caso más sencillo y queda por estudiar que ocurre con las características del tráfico dentro del intervalo, que puede presentar escalado multifractal [6]. Posiblemente estas características dentro del intervalo hacen que la predicción sea muy grosera al no tenerlas en cuenta y es necesario introducir más parámetros aparte de los bytes en bruto. Sin embargo existen múltiples escenarios donde *si la red dispone de una estimación de los bytes por intervalo es suficiente para mejorar en gran medida las prestaciones*. Una posible aplicación son los entornos de “Burst Switching” donde, gracias a la predicción, es posible paralelizar el tiempo de paquetización con la reserva de recursos. Al conocer de antemano los bytes a transmitir, gracias a la predicción, se envía el mensaje de reserva de recursos antes de comenzar la paquetización, como se muestra en la figura 13.

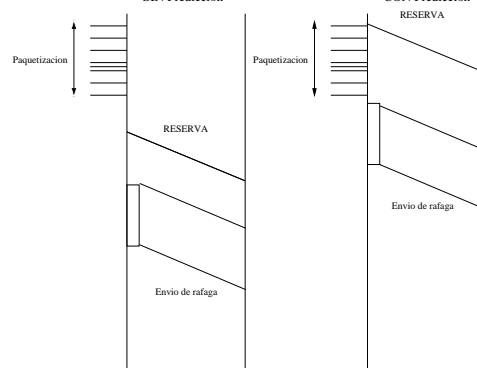


Figura 13: Optical Burst Switching

## 4 Conclusiones y trabajos futuros

En este artículo hemos presentado métodos de predicción de coste computacional muy bajo que pueden ser utilizados para dimensionar dinámicamente enlaces de la Internet que requieren una estimación de bytes por intervalo. Queda como trabajo futuro el análisis de prestaciones con distintos modelos de tráfico dentro de cada intervalo y la selección de algoritmos de predicción óptimos en este último caso.

## Referencias

- [1] K. Park and W. Willinger (Editors). *Self-similar Network Traffic and Performance Evaluation*. Wiley Interscience, 2000.
- [2] W. Willinger, M. S. Taqqu, R. Sherman, and D. V. Wilson. Self-similarity through high-Variability: Statistical analysis of ethernet LAN traffic at the source level. *IEEE/ACM Transactions on Networking*, 5(1), Febrero 1997.
- [3] M. E. Crovella and A. Bestavros. Self-similarity in world wide web traffic: Evidence and possible causes. *IEEE/ACM Transactions on Networking*, 5(6):835–846, Diciembre 1997.
- [4] I. Norros. On the Use of Fractional Brownian Motion in the Theory of Connectionless Networks. *IEEE Journal on Selected Areas in Communications*, 13(6):953–962, Agosto 1995.
- [5] *IEEE Communications Magazine* Special Issue on MPLS, Diciembre 1999
- [6] A. Erramilli, O. Narayan and A. Neidhart Performance Impacts of Multi-Scaling in Wide Area TCP/IP Traffic *INFOCOM 2000*, Tel Aviv, Israel