

doi: 10.3788/gzxb20144307.0706023

端到端网络流量的混合估计方法

蒋定德, 赵祖耀, 许宏伟, 王兴伟

(东北大学 信息科学与工程学院, 沈阳 110819)

摘 要: 利用主成分分析法获取端到端网络流量的主要特征分量并获得其初始估计结果. 为克服其初值敏感性将估计结果作为遗传算法的初始值、链路流量估计偏差作为遗传算法的适应度函数, 通过构建合适的交叉和变异概率函数来控制遗传算法的交叉和变异过程. 采用合适的约束迭代函数, 利用遗传算法通过迭代寻优获得端到端流量的估计结果, 仿真结果表明所提出的方法是可行的.

关键词: 端到端流量; 流量估计; 主成分分析; 迭代过程; 流量建模

中图分类号: TP393

文献标识码: A

文章编号: 1004-4213(2014)07-0706023-6

Mixed Estimation Approach to End-to-End Network Traffic

JIANG Ding-de, ZHAO Zu-yao, XU Hong-wei, WANG Xing-wei

(College of Information Science and Engineering, Northeastern University, Shenyang 110819, China)

Abstract: Principal component analysis was exploited to extract the principal features of end-to-end network traffic and to attain the initial estimation results. This results are taken as the prior value of genetic algorithm to overcome its sensitiveness to the prior value. The estimation biases of link traffic is regarded as the fitness function of genetic algorithm. The crossover and mutation probability functions are built to control its crossover and mutation processes. The appropriate iterative function with constraints is built. The genetic algorithm is used to attain the end-to-end traffic estimation results in the iterative way. Simulation results show that the proposed method is feasible.

Key words: End-to-end traffic; Traffic estimation; Principal component analysis; Iterative process; Traffic modeling

OCIS Codes: 060.4256; 060.4251; 060.4258; 060.1155

0 Introduction

With the rapid development of information technologies, communication networks are becoming more and more complex and their scales are increasing^[1-2]. The network traffic presents all kinds of new features, such as spatio-temporal correlation features, self-similarity nature, heavy-tailed distribution and so on^[3-5]. End-to-end traffic shows the global points of view about network behaviors and gives the network-wide traffic information to network operators. It is an important input parameter for

network management and traffic engineering. Due to the difficulty in directly attaining them, end-to-end traffic estimation has received more attention from researchers and operators around the world^[6-7].

The estimation methods in Refs. [7-9] is based on statistical theories to be put forward. They used the statistical models to set up networks' source-destination node flow model and exploited the statistical theory to estimate end-to-end traffic in a network. The approaches in Refs. [10-11] exploited the gravity model and link information to perform the inference of end-to-end traffic. And routing configuration and network

Foundation item: The National Natural Science Foundation of China (No. 61071124), the Specialized Research Fund for the Doctoral Program of Higher Education (No. 20100042120035), the Program for New Century Excellent Talents in University (No. NCET-11-0075) and the Fundamental Research Funds for the Central Universities (Nos. N120804004, N130504003)

First author: JIANG Ding-de, male, associate professor, Ph. D. degree, mainly focuses on network measurement and energy-efficient networks. Email: jiangdingde@ise.neu.edu.cn

Received: Nov. 4, 2013; **Accepted:** Jul. 4, 2013

<http://www.photon.ac.cn>

topology were utilized to build the effective estimation methods. Additionally, the part measurements of end-to-end traffic were used to help to attain the more accurate estimation results^[12]. However, these methods still hold the larger estimation errors.

A new mixed estimation approach to end-to-end traffic by combining genetic algorithm and principle component analysis are proposed in this paper. Genetic algorithm is used to attain the estimation of end-to-end traffic. Due to the sensitiveness of it to the initial value and highly dynamic, it is significantly difficult to directly use it to attain the accurate results. In this papers, the principal component analysis is exploited to overcome this problem. The principal feature of end-to-end traffic is extracted and the initial estimation is obtained. The genetic algorithm regards this initial solution as its initial value to begin to iterate step by step. Finally, the accurate estimations can be attained.

1 Problem statement and model

For a network with n nodes, specially a large-scale backbone network, network operators and researchers always hope to be able to know the activity situations of network traffic so that they can better learn network performance and make appropriate network plan. End-to-end traffic transfers more than one nodes and needs to multiple hops to arrive at the destination node. They are buried into the large link traffic, so it is not easy to find and extract them. The problem to discuss in this paper is how to estimate the end-to-end network traffic with the fairly good accuracy.

In network, end-to-end traffic goes through many links and nodes in accordance with the routing configuration information. They are aggregated into link traffic at a link. All end-to-end traffic and the whole link traffic meet some certain constrains by the routing information. So end-to-end traffic estimation can be modeled into the following equation

$$\mathbf{y}=\mathbf{A}\mathbf{x} \quad (1)$$

where $\mathbf{x}=(x_1, x_2, \dots, x_u)'$ and $\mathbf{y}=(y_1, y_2, \dots, y_v)'$ are a matrix denoting the link counts and end-to-end traffic in network, $u=n^2$ denotes the number of the end-to-end or Origin Destination(OD) flows, v stands for the number of links; \mathbf{A} denotes the route information matrix under the condition of a given topology and routing configuration, which is expressed by $\mathbf{A}=\{a_{ij}\}$, a_{ij} is the element of matrix \mathbf{A} . If OD flow j goes through link i , then $a_{ij}=1$, or $a_{ij}=0$.

Because end-to-end traffic represents the real information of networks, they should be nonnegative number. Eq. (1) can be converted as

$$\begin{cases} \mathbf{y}=\mathbf{A}\mathbf{x} \\ x_i \geq 0 \end{cases} \quad (2)$$

where $i=1, 2, \dots, u$.

Next, end-to-end traffic estimation is described by Eq. (2). It denotes how to find the correct and needed \mathbf{x} from \mathbf{y} and \mathbf{A} , meeting the nonnegative constraints in Eq. (2).

2 Initial value construction

This paper proposes an end-to-end traffic mixed estimation approach by exploiting genetic algorithm and principal component analysis. The approach takes advantage of two methods and avoids their shortcomings. However, due to the sensitiveness of genetic algorithm itself to a prior of the end-to-end traffic, this will lead to the considerable estimation errors. We exploit the principal component analysis to overcome this problem.

Principal component analysis is one of the best methods to find the most effective linear transform in the least mean square, namely high dimensional data to lower dimensions of space projection. Based on principal component analysis, we study how it can effectively capture the feature of end-to-end traffic in the large-scale backbone network as mentioned in Ref. [10].

By matrix theory, the z -time-slot sample data x_s of end-to-end traffic can be denoted into the following equation

$$\mathbf{x}_s=[x_1, x_2, \dots, x_z] \quad (3)$$

According to principal component analysis, the z -time-slot sample data denoted by x_s can be exploited to capture the inherent properties of end-to-end traffic. As a result, x_s can be decomposed as follows

$$\mathbf{x}_s=\mathbf{U}\mathbf{D}\mathbf{V}^T \quad (4)$$

where \mathbf{U} is a $u \times u$ matrix, \mathbf{D} is a $u \times z$ dialog matrix, \mathbf{V} is a $z \times z$ matrix.

The non-zero dialog elements of matrix \mathbf{D} reflect the energy spectrum of the end-to-end traffic \mathbf{x}_s , which correspond to the eigenvalue of $\mathbf{x}'_s \mathbf{x}_s$. If the sum of \mathbf{D}' 's non-zero dialog elements is s_D , we take $0.5s_D$ as the threshold. If the non-zero dialog element λ_i of matrix \mathbf{D} meets $\lambda_i > 0.5s_D$, the corresponding vectors of \mathbf{U} and \mathbf{V} can be picked out. Then we can select the appropriate k principal part of end-to-end traffic in such a way. By choosing the k top principal parts of the end-to-end sample traffic, \mathbf{x}_s can approximately be presented as follow

$$\tilde{\mathbf{x}}_s=\mathbf{V}'\mathbf{D}'\mathbf{U}' \quad (5)$$

where \mathbf{V}' and \mathbf{D}' denote the main feature in the end-to-end sample traffic. That is, by the samples of the end-to-end traffic, we build the model to describe it, only using the k top principal components with accordance of principal component analysis theory. Accordingly, by this model, we can predict the future end-to-end traffic

according to the below equation

$$\mathbf{x}(t) = \mathbf{V}'\mathbf{D}'\mathbf{u}(t)' \quad (6)$$

where $t = z+1, z+2, \dots, T$ and T denotes the number of time slots to analyze; $\mathbf{u}(t)'$ is a vector denoting the feature flows at t moment. According to Eq. (1), the below equation can be attained

$$\mathbf{y}(t) = \mathbf{A}\mathbf{V}'\mathbf{D}'\mathbf{u}(t)' \quad (7)$$

where $t = z+1, z+2, \dots, T$.

By Eq. (7), we can get the estimation of $\mathbf{u}(t)'$, namely $\hat{\mathbf{u}}(t)'$. According to Eq. (6), the end-to-end traffic estimation at t moment is obtained

$$\hat{\mathbf{x}}(t) = \mathbf{V}'\mathbf{D}'\hat{\mathbf{u}}(t)' \quad (8)$$

where $t = z+1, z+2, \dots, T$.

Unfortunately, the estimation results by Eq. (8) have the larger estimation errors. Thus they can not directly be used as the final estimations of end-to-end traffic for network activities.

3 Iterative inference

By using the evolution and genetic mechanism, a given original solvable group will be able to converge to a feasible solution to this problem. According to genetic algorithm theory, because end-to-end traffic is satisfied with the constraints in Eq. (2), the subject function of genetic algorithm can be denoted as

$$\begin{cases} \min \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_{\mathbb{F}}^2 \\ \text{s. t. } \mathbf{x} \geq 0 \end{cases} \quad (9)$$

where \mathbf{A} indicates route matrix, \mathbf{y} indicates link matrix and \mathbf{x} indicates end-to-end traffic matrix. Thus on the condition of the minimum norm solution, the optimal solution to end-to-end traffic can be got.

For every feasible solution \mathbf{x} to end-to-end traffic, define the following difference

$$\Delta f = \mathbf{y} - \mathbf{A}\mathbf{x} \quad (10)$$

The crossover probability can be defined as

$$p_a = \frac{\Delta f - \Delta \bar{f}}{\Delta f_{\max} - \Delta \bar{f}} \quad (11)$$

The mutation probability can be defined as

$$p_b = \frac{\Delta f_{\max} - \Delta f}{\Delta f_{\max} - \Delta \bar{f}} \quad (12)$$

where $\Delta \bar{f}$ denotes the average difference in Eq. (10), and Δf_{\max} stands for the maximum difference.

Given the crossover and mutation thresholds \bar{p}_a and \bar{p}_b , if $p_a > \bar{p}_a$, then perform the crossover behavior, and if $p_b > \bar{p}_b$, then make the mutation behavior.

Then, according to Eq. (2) and (9), we use the genetic algorithm to construct the following iterative estimation about the end-to-end traffic

$$\begin{cases} \mathbf{x}^{r+1} = \mathbf{x}^r \Delta \mathbf{x}^r \\ \text{s. t. } \mathbf{x}^r \geq 0 \\ \min \|\mathbf{y} - \mathbf{A}\mathbf{x}^r\|_{\mathbb{F}}^2 \\ \mathbf{x}^0 = \hat{\mathbf{x}}(t), t = M+1, M+2, \dots, T \end{cases} \quad (13)$$

where $\Delta \mathbf{x}^r$ is the variable value of the end-to-end traffic at the r th iterative step according to the generic algorithm theory.

Accordingly, we can obtain the estimations of the end-to-end traffic by Eq. (13). The mixed algorithm proposed above can be described as follow

Step 1: Decompose the end-to-end sample traffic matrix \mathbf{x}_s according to Eq. (5). Then obtain the eigenvector matrix \mathbf{U} , the diagonal matrix \mathbf{D} , and the eigenflow matrix, \mathbf{V} , respectively.

Step 2: By the principal component analysis, extract the k top principal components, and then obtain the parameters of the model about the end-to-end traffic from the sample traffic, namely \mathbf{V}' and \mathbf{D}' .

Step 3: According to Eq. (7), get the estimations of the eigenflow $\hat{\mathbf{u}}(t)'$ at time t by principal component analysis.

Step 4: According to Eq. (8), get the estimations of the end-to-end traffic $\hat{\mathbf{x}}(t)$ at time t by principal component analysis.

Step 5: Set $\mathbf{x}^0 = \hat{\mathbf{x}}(t), t = z+1, z+2, \dots, T$, begin the genetic evolution process, and get the optimal solution \mathbf{x}^r .

Step 6: Find the optimal solution \mathbf{x}^r which is satisfied with $\mathbf{x}^r \geq 0$ and $\min \|\mathbf{y} - \mathbf{A}\mathbf{x}^r\|_{\mathbb{F}}^2$.

Step 7: Build the iterative equation $\mathbf{x}^{r+1} = \mathbf{x}^r + \Delta \mathbf{x}^r$, and get the resulting solution meeting with Eq. (13).

Step 8: If the estimation process is over, then save the results to file and exit, or go back to Step 3.

Through above iterative inference process, we can obtain the final estimation of end-to-end traffic.

4 Simulation results and analysis

To verify the performance of the mixed algorithm, the data from the real network is used to perform several simulation. In this paper, we will discuss this method (for short Mixed) and other two approaches (namely PCA and GA) in the estimation biases and relative errors. PCA were reported as a good estimation for end-to-end traffic in Ref. [10], while GA is implemented by us to directly estimate end-to-end traffic.

Fig. 1 denotes the relative biases of ODs 40, 80, 99 and 137 of our method. From Fig. 1(a) and (b), we can see that our method have the estimation relative biases less than 1 over most of the estimation time. Fig. 1(c) and (d) indicates that our method also hold have the similar estimation relative biases. Fig. 1 shows that our method can make the fairly accurate estimation for end-to-end traffic over the different moments.

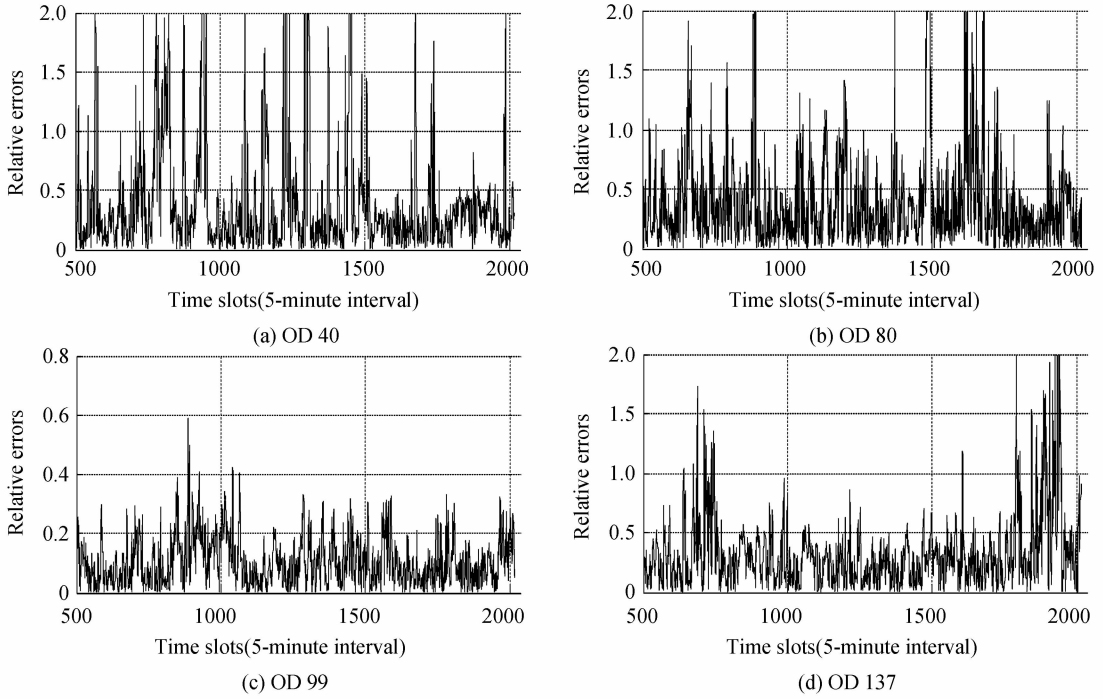


Fig. 1 Relative biases of ODs 40, 80, 125 and 137

To further analyze the estimation performance of our method, we compare our method with PCA and GA for the relative biases. Fig. 2 analyze the relative biases of three methods for ODs 40, 80, 99, and 137.

Fig. 2 indicate that in contrast to PCA and GA, our method exhibits the lower relative biases for end-to-end traffic. Thus our method can attain the better estimation results.

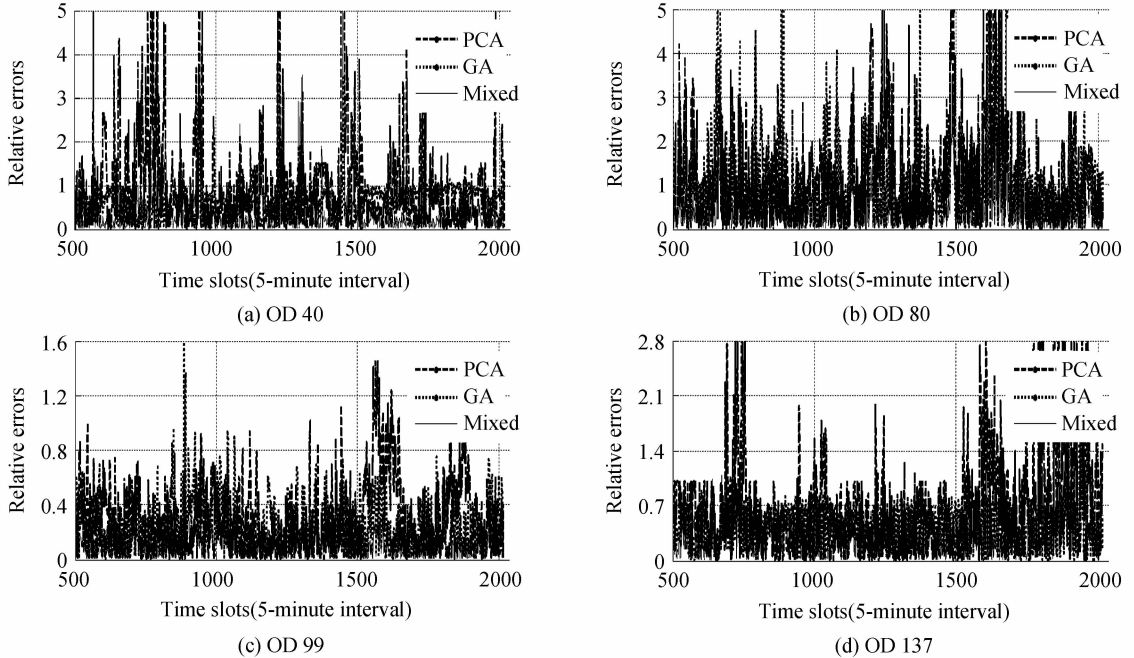


Fig. 2 Relative bias comparisones of ODs 40,80, 125 and 137

To further discuss the performance of our method, we make some analysis of the relative errors of three methods, including the relative spatial and temporal errors. The relative spatial errors are denoted as^[13-14]

$$\text{err}_{\text{sperr}}(k) = \frac{\|\hat{\mathbf{x}}_T(k) - \mathbf{x}_T(k)\|_2}{\|\mathbf{x}_T(k)\|_2} \quad (14)$$

where $k=1, \dots, u$, u is the number of the OD flows in

the large-scale network. T is total measurement time, $\|\cdot\|_2$ is the norm of L_2 . $\hat{\mathbf{x}}_T(k)$ indicates the estimation results of the end-to-end traffic for OD k during the whole times T and $\mathbf{x}_T(k)$ is the real value.

Likewise, the temporal relative errors are denotes as^[13-14]

$$\text{err}_{\text{terr}}(t) = \frac{\|\hat{\mathbf{x}}_u(t) - \mathbf{x}_u(t)\|_2}{\|\mathbf{x}_u(t)\|_2} \quad (15)$$

where $t = 1, \dots, T$, $\hat{\mathbf{x}}_u(t)$ indicates the estimation results of all the OD flows over in moment t and $\mathbf{x}_u(t)$ indicates the real value of all the OD flows.

Fig. 3 plots the spatial and temporal relative errors of three methods. From Fig. 3, we can find that compared with other two method, our approach holds the lower spatial and temporal relative errors. This illustrates that our method really can make the effective and accurate estimation for the end-to-end traffic. To verify the estimation ability of our method, we observe the Cumulative Distribution Function (CDF) of the relative errors of three approaches. Fig. 4 indicates their CDF. We can see easily that the CDF curves of our method are far on top of those of other two methods. This further validates that our method

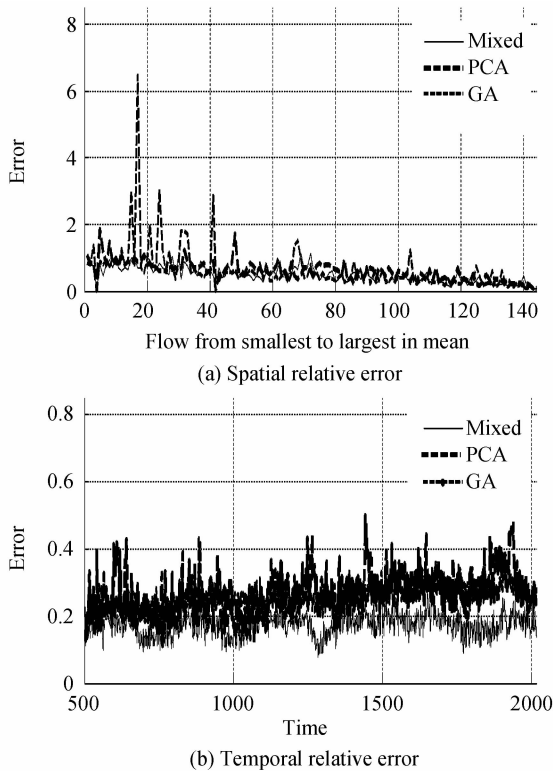


Fig. 3 Relative error analysis

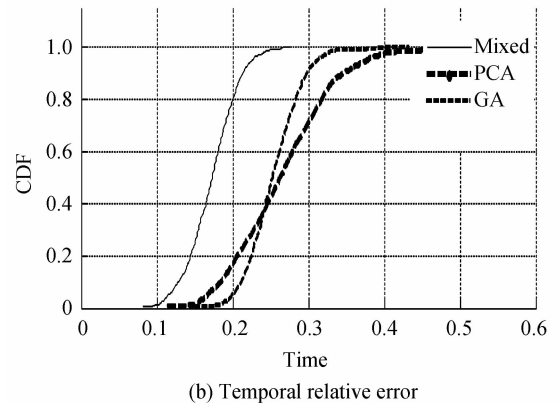
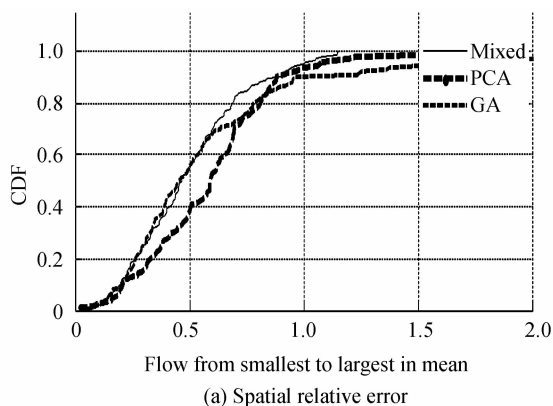


Fig. 4 CDF analysis of relative errors

hold the better estimation ability for end-to-end traffic.

5 Conclusion

This paper studies the estimation problem of end-to-end traffic. By combining genetic algorithm and principal component analysis, a mixed estimation approach is proposed to make the accurate estimation for end-to-end traffic in the large-scale networks.

Principle component analysis is exploited to extract the principal feature components of end-to-end traffic and to attain the initial estimation results. Then genetic algorithm adopts the initial results to perform the further estimation and attain the final estimation results. Simulation results indicates that the proposed approach is feasible and promising.

References

- [1] WANG X, ZHANG D. A novel method to estimate IP traffic matrix[C]. Proceedings of WiCOM'10, 2010; 1-4.
- [2] CAI Ting, HUANG Shan-guo, LI Xin, *et al.* Dynamic survivable mapping algorithm based on ant colony optimization in IP over WDM networks[J]. *Acta Photonica Sinica*, 2012, **41**(12): 1400-1404.
- [3] JIANG Ding-de, XU Zheng-zheng, NIE Lai-sen, *et al.* An approximate approach to end-to-end traffic in communication networks[J]. *Chinese Journal of Electronics*, 2012, **21**(4): 705-710.
- [4] GUAN Ai-hong, WANG Bo-yun, FU Hong-liang, *et al.* A deflection routing mechanism based on priority and burst segmentation in optical burst switching networks[J]. *Acta Photonica Sinica*, 2012, **41**(2): 127-132.
- [5] AKGUL T, BAYKUT S, KANTARCI M, *et al.* Periodicity-based anomalies in self-similar network traffic flow measurements[C]. Proceedings of TIM'11, 2011, **60**(4): 1358-1366.
- [6] JIANG Ding-de, QIN Wen-da, TANG Qing-yi, *et al.* An estimation approach to traffic matrix in optical networks based on network tomography[J/OL]. [2013-11-04]. <http://www.photon.ac.cn/CN/abstract/abstract20262.shtml>.
- [7] TEBALDI C, WEST M. Bayesian inference on network traffic using link count data[J]. *Journal of American Statistics Association*, 1998, **93**(442): 557-576.
- [8] JUVA I, KUUSELA P, VIRTAMO J. A case study on traffic matrix estimation under Gaussian distribution [C]. In Proceedings of NTS'04, 2004: 49-60.

-
- [9] VATON S, BEDO J. Network traffic matrix; How can one learn the prior distributions from the link counts only[C]. In Proceedings of ICC'04, 2004; 2138-2142.
- [10] SOULE A, LAKHINA A, TAFT N, *et al.* Traffic matrices: balancing measurements, inference and modeling [C]. In Proceedings of SIGMETRICS'05, 2005, **33**(1): 362-373.
- [11] ZHANG Y, ROUGHAN M, DUFFIELD N, *et al.* Fast accurate computation of large-scale IP traffic matrices from link loads[C]. In Proceedings of SIGMETRICS'03, 2003, **31**(3): 206-217.
- [12] TAKEDA T, SHIONOTO K. Traffic matrix estimation in large-scale IP networks[C]. In Proceedings of LANMAN'10, 2010; 1-6.
- [13] JIANG Ding-de, XU Zheng-zheng, CHEN Zhen-hua, *et al.* Joint time-frequency sparse estimation of large-scale network traffic[J]. *Computer Networks*, 2011, **55**(10): 3533-3547.
- [14] JIANG Ding-de, XU Zheng-zheng, XU Hong-wei, *et al.* An approximation method of origin-destination flow traffic from link load counts[J]. *Computers and Electrical Engineering*, 2011, **37**(6): 1106-1121.