

Advances in Network Tomography

Edoardo M. Airolidi
eairolidi@stat.cmu.edu

supervised by:

Christos N. Faloutsos
christos@cs.cmu.edu

Abstract

Knowledge about the origin-destination (OD) traffic matrix allows us to solve problems in design, routing, configuration debugging, monitoring and pricing. Direct measurement of these flows is usually not implemented because it is too expensive. A recent work provided a quick method to learn the OD traffic matrix from a set of available standard measurements, which correspond traffic flows observed on the link of a network every 5 minutes. Such a time span allows for more computationally expensive methods that in turn yield a better estimate of the OD traffic matrix.

In this work we are the first to explicitly introduce time in learning the OD traffic matrix. The second contribution is that we are the first to use realistic non-Gaussian marginals, specifically the Gamma and the successful log-Normal ones. We combine both these ideas in a novel, doubly stochastic and time-varying Bayesian dynamical system, and provide a simple and elegant solution to obtain informative prior distributions for the stochastic dynamical behavior. Our method out-performs existing solutions in a realistic setting.

Acknowledgments

I am grateful to my advisor Prof. Christos N. Faloutsos for presenting me with the problem of the reconstruction of the origin-destination flows from observable link loads, for valuable comments and suggestions, and for his enthusiasm and his continuous support during all the phases of this exciting project.

I wish to thank Prof. Srinivasan Seshan, Dr. Russel Yount and Dr. Frank Kietzke for providing me with their expertise, and for retrieving origin-destination traffic flows and link loads on Carnegie Mellon local area network, necessary to validate the methods we proposed. Dr. Frank Kietzke contributed during several stages of this study with comments and suggestions, and kindly provided me with detailed explanations on technical aspects of the problem.

Further I wish to thank Prof. Stephen E. Fienberg, and Prof. Christopher Genovese for their comments and suggestions at an early stage of this project, and Prof. Anthony Brockwell for helpful discussion of several aspects of the problem, and for pointing me towards relevant literature.

Contents

| | | |
|----------|---|-----------|
| 1 | Introduction | 5 |
| 1.1 | Problem Definition | 5 |
| 2 | Literature Review | 7 |
| 2.1 | Transportation Research | 10 |
| 2.2 | Statistical Research | 10 |
| 2.2.1 | A Recent Local Maximum Likelihood Approach | 11 |
| 3 | Proposed Methods | 14 |
| 3.1 | Explicit Dynamics for Gaussian Origin-Destination Flows | 14 |
| 3.1.1 | A State-Space Representation for the Model | 15 |
| 3.1.2 | Ad-Hoc M-Step for the EM Algorithm | 16 |
| 3.1.3 | Two-Stages Maximization of the Likelihood | 17 |
| 3.2 | 1-Time Non-Gaussian Origin-Destination Flows | 17 |
| 3.2.1 | Computing the Support | 18 |
| 3.2.2 | Irreducibility of the Chain | 18 |
| 3.2.3 | Gamma and log-Normal Models | 19 |
| 3.2.4 | Informative priors | 21 |
| 3.3 | Combining Dynamics and Non-Gaussianity | 21 |
| 3.3.1 | Informative Priors for Stochastic Dynamics | 21 |
| 3.3.2 | Bayesian Dynamical Systems | 22 |
| 3.3.3 | Particle Filter via SIR-Move Algorithm | 22 |
| 3.4 | Multivariate Integration | 23 |
| 4 | Experiments | 24 |
| 4.1 | Exploring Carnegie Mellon Network | 24 |
| 4.1.1 | Empirical Distributions of the OD Flows | 26 |
| 4.1.2 | Modeling the Coefficients of Constant Association α_{ij} | 27 |
| 4.2 | Exact Recovery Algorithms for Sparse Traffic Situations | 27 |
| 4.3 | Intense Traffic Sub-Networks at Carnegie Mellon | 27 |
| 4.3.1 | Naive SVD Solution for Strongly Correlated Flows | 27 |
| 4.3.2 | Local Dynamical Behavior | 28 |
| 4.3.3 | A Case Study: the Star Network Topology | 29 |
| 5 | Conclusions | 34 |
| A | Gaussian Dynamical System | 37 |
| A.1 | EM algorithm | 37 |
| A.2 | More computational efficiency | 38 |
| A.3 | KF posteriors | 38 |
| A.3.1 | One Y to one X | 38 |
| A.3.2 | One Y to many X s | 39 |
| B | The Key to fig n.10. | 41 |

1 Introduction

Knowledge about the origin-destination (OD) traffic matrix allows us to solve problems in design, routing, configuration debugging, monitoring and pricing; in fact the OD traffic matrix provides us with valuable information about who is communicating to whom in a local area network, at any given time. Most routers are not able to measure the OD traffic flows, though. Further the direct measurement of OD traffic flows via SNMP queries is never implemented on the few models that would technically allow it, because usually infeasible. Approximate methods have been recently proposed in order to infer the OD traffic matrix from a set of standard measurements, traffic loads on the links of the network produced every 5 minutes. Such a delay between successive measurements allows for computational methods that produce better estimates for the OD traffic matrix.

We improved the models present in the literature by introducing two realistic assumptions: (1) we modeled the marginal OD traffic flows with skewed distributions like Gamma and log-Normal, and (2) we introduced time dependence among the OD flows explicitly by means of a stochastic dynamical behavior, in our self-organizing Bayesian dynamical system. The Gamma and log-Normal models reduced by 25% and 38%, respectively, the estimation error yielded by recently proposed solutions; the introduction of explicit stochastic local dynamics reduced the estimation error up to 41% and 46%, respectively. The magnitude of the improvements entailed by the simple ideas we propose went far beyond that of state-of-the-art resampling schemes that could be used to refine any given set of estimates. Further the stochastic dynamics played an essential role in our models; it served as the right channel where to introduce prior information about the OD flows, and mitigate the problem of multiple modes in a certain probabilistic mapping to be defined below.

The outline of this report is as follows: in the remainder of this section we formulate the problem with X s and Y s; in section 2 we survey past approaches to the problem and point out their weaknesses; in section 3 we detail our proposed methods; in section 4 we explore real data from Carnegie-Mellon Local Area Network, apply our methods and models, and compare our estimates to those provided by current state-of-the-art solutions; and finally in section 5 we conclude with some remarks.

1.1 Problem Definition

We begin by giving a mathematical characterization of the problem.

Problem. *Given observations $Y(t) := [Y_1(t), \dots, Y_\ell(t)]'$ over times $t = 1, \dots, T$, and a matrix $A_{(\ell \times \kappa)}$ such that $Y(t) = AX(t) \forall t$, we want to estimate the unobservable inputs $X(t) := [X_1(t), \dots, X_\kappa(t)]'$ over times $t = 1, \dots, T$. It is always $\kappa \geq \ell$, and in general $\kappa = O(\ell^2)$.*

The formulation above is common to many applications; in this project we were interested in network traffic analysis. We observed traffic loads on the links of a network with terminal nodes, and routing nodes¹, and we were interested in estimating the traffic flows between all, or selected pairs, of terminal nodes. Hence in our problem the vector of observations $Y(t)$ would contain the set of measurements on the links of the network available at time t , and the vector $X(t)$ would contain non-observable origin-destination flows between pairs of terminal nodes at time t . The relevant information to characterize the structure of the network² would be contained in the matrix A , time independent, which tells us exactly how the flows $X(t)$ combine to form the observed link loads $Y(t)$ through the equations $Y(t) = AX(t)$ at each point in time.

¹Namely routers and switches, which did not create or absorb traffic, but merely filtered it and/or redirected it.

²We restricted our attention to the cases where A would entail a deterministic routing scheme.

A First Example

Consider figure 1, the central node represents a router, whereas the external nodes represent sub-networks. In the terminology above the router is a routing node, and the sub-networks are terminal nodes, and we wanted to estimate the flows between them. Throughout this report we referred to the non-observable flows between pairs of terminal nodes (the dashed arrows) as origin-destination (OD) flows, whereas we referred to the measurable loads on the directed links of the graph (the solid arrows) as link flows. We studied traffic flows in terms of Kbytes. More precisely we were interested in estimating the probability of observing a certain amount of traffic $X_i(t)$ over each origin-destination route i . In the sample network in figure 1 the information transits through the

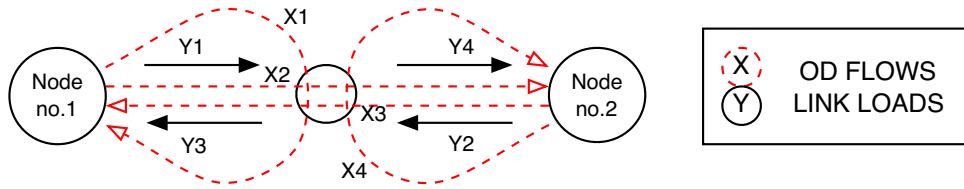


Figure 1: Basic routing scheme

routing node in the form of packets of different sizes. The four terminal nodes are connected to it through cables. The observation vector $Y(t)$ consists of four components at each time tick (5 minutes), namely the directed link loads we would measure at each of the four interfaces where the physical cables are connected; two measurements for the Kbytes going into the router, and two measurements for the Kbytes going out of the router. With two terminal nodes connected to the routing node the number of possible OD routes is 4, allowing for traffic from a node to itself³. Notice that since the router neither generates, nor absorbs traffic each one of the measurements can be recovered from the other three. In figure 2 we represented the mathematical problem in terms

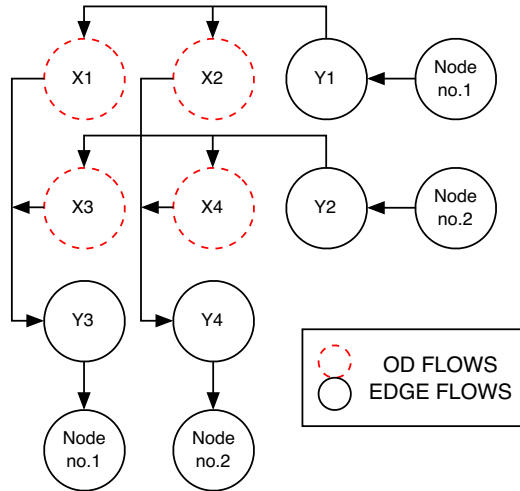


Figure 2: Basic mathematics with X s and Y s.

³The function of a routing node is to filter the packets that it receives, and redirect them. Traffic from a node i to itself amounts to the traffic that the router keeps local to that node i , by filtering it and sending it back.

of tables. The dashed circles, inner cells of the table, contain the non-observable OD flows, whereas solid circles, margins of the table, contain the available measurements, the link loads, at each time. In other words there is a sequence of tables over time; we are able to observe their margins and we want to make inference on the inner cells. In this report we denoted the observed link loads, in Kbytes, with Y s⁴, and we ordered them in a real valued random vector

$$Y = (Y_1 := Y_{1,in}, Y_2 := Y_{2,in}, Y_3 := Y_{1,out}, Y_4 := Y_{2,out})' \quad (1.1)$$

and we denoted the non-observable origin-destination flows, again in Kbytes, with X s and we ordered them in another real valued random vector

$$X = (X_1 := X_{\text{from } 1, \text{to } 1}, X_2 := X_{\text{from } 1, \text{to } 2}, X_3 := X_{\text{from } 2, \text{to } 1}, X_4 := X_{\text{from } 2, \text{to } 2})'. \quad (1.2)$$

For reasons we explain below we only needed to keep independent components of $Y(t)$ at each time t ; for example, in the matrix in equation 1.3 below we dropped a row before applying our methods, so that $\ell = 3$ and $\kappa = 4$ for this example. At every point in time the number of available origin-destination routes is $\kappa = O(\ell^2)$, where ℓ is the dimension of the vector of observed link loads. Origin-destination byte counts in X are related to the measurements in Y by means of the routing matrix A , that contains information on the deterministic routing scheme through the set of equations $Y = AX$. The matrix A summarizes the network structure in useful ways. For example $(AA')_{i,i}$ counts the OD routes passing through i , and $(AA')_{i,j}$ counts the number of OD routes passing through both i and j . The linear equations $Y = AX$ that represent the assignments in the network in figures 1 and 2 look like this:

$$\begin{bmatrix} Y_1 \\ Y_2 \\ Y_3 \\ Y_4 \end{bmatrix} = \begin{bmatrix} 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 \\ 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \end{bmatrix} \begin{bmatrix} X_1 \\ X_2 \\ X_3 \\ X_4 \end{bmatrix} \quad (1.3)$$

In table 1 below we summarize the symbols we used throughout this report.

2 Literature Review

The problem presented in section 1.1 is specific instance of a more general problem.

Classical Transportation Problem. *Given a directed graph $G(V, E)$, with K nodes and ℓ edges, identify a subset of $\kappa < K$ nodes in E and call them centroids. Let X be a vector with components $X_{od} := \{\text{the average number of trips going from node } o \text{ to node } d \text{ within a given time period}\}$. Each origin-destination flow X_{od} subdivides on the network into path flows ξ_{odk} , $k \in I_{od}$, where I_{od} is the subset of all path connecting the pair of centroids o and d .*

For a given edge $e \in E$, the sum of all path flows traversing this link is called the link load

$$Y_e = \sum_{od} \sum_{k \in I_{od}} \xi_{odk} I_{(e)}(k), \quad (2.1)$$

where $I_{(a)}(b)$ is 1 if $a = b$ and zero otherwise. More compactly equation 2.1 can be written as

$$Y = \Delta \xi = AX. \quad (2.2)$$

⁴We drop the time index to improve readability.

| Symbol | Meaning |
|----------------------------------|---|
| $Y(t), Y_t$ | $(\ell \times 1)$ column random vector containing the observable links loads at time t , as ordered according to equation 1.1. Y_t will be used when it is clear that we are considering vectors. |
| $Y_i(t), Y_t^i$ | Random number representing measurements of traffic loads at time t on the i^{th} link of the network, according to the ordering established in equation 1.1. |
| $\{Y\}$ | Set of random vectors as in $Y = \{Y(1), \dots, Y(T)\}$. |
| $X(t), X_t$ | $(\kappa \times 1)$ column random vector containing the unobservable origin-destination flows at time t , as ordered according to equation 1.2. We also call $X(t)$ the traffic matrix, following the interpretation of figure 2 in terms of tables. X_t will be used when it is clear that we are considering vectors. |
| $X_i(t), X_t^i$ | Random number representing unobservable flows at time t between the i^{th} origin-destination route in the network, according to the ordering established in equation 1.2. |
| $\{X\}$ | Set of random vectors as in $X = \{X(1), \dots, X(T)\}$. |
| A | $(\ell \times \kappa)$ routing matrix. Contains information about the deterministic routing scheme, through the set of equations $Y(t) = AX(t)$. It does not change over time. |
| ℓ | Number of links in the network for which we obtain measurements $Y_i(t)$ over time. |
| κ | Number of origin-destination flows $X(t)$ in the network which we are interested in estimating over time. |
| $\Theta, \Theta_t, \Psi, \Psi_t$ | Generic vectors of parameters to be defined. Subscript t indicates time dependence. |
| $N_D(\lambda, \Sigma)$ | Is a multivariate normal density on R^D , with mean column vector λ and variance covariance matrix Σ . |
| $\text{diag}(\lambda)$ | It is $(D \times D)$ diagonal matrix, for λ $(D \times 1)$ column vector, with diagonal elements $[\text{diag}(\lambda)]_{(i,i)} = \lambda_i, \quad i = 1, \dots, D$. |
| ϕ, ϕ_t | Scalar. Provides extra variability. Subscript t indicates time dependence. |
| t_W | Scalar. Amplitude of the time window for local likelihood methods. |
| OD | Origin-destination, refers to the X s. |
| IPFP | Iterative Proportional Fitting Procedure. |

Table 1: Summary of symbols and abbreviations.

In the definition of the problem above the centroids are the terminal nodes of our network, the assignment matrix Δ contains the same information as our routing matrix A , and X, Y are the vectors of origin-destination flows and link loads, respectively. An even more general version of the problem had been given, which would include random costs $C_e(Y)$ for traveling on a certain edge e , function on the observed traffic Y , random costs $U_k(Y)$ for choosing a certain path k , and proportional trip demand $P(C)$, possibly an implicit function of the travel cost vector C , that would yield a slightly different expression for the observed link loads, as in

$$Y = \Delta P(C) X, \quad (2.3)$$

and we would redefine the assignment map $A(X) := \Delta P(C) X$, possibly non linear.

Underlying Assumptions

The transportation problem in all its formulations is “old”, and several solutions had been proposed and rediscovered over the years. In order to keep track of the different characterizations we mention here some relevant dimensions: (L1) there is only one vector of flows X , that does not change over the sampling period; (L2) the link loads Y are measured with error; (L3) the assignment mapping $A(X)$ is endogenous $A(X) := \Delta P(C) X = \Delta P^*(Y) X = A(X, Y)$; (L4) the assignment mapping $A(X)$ is deterministic, as opposed to random; (L5) the assignment mapping $A(X)$ entails unique

paths, as opposed to multiple paths; (L6) the assignment mapping $A(X)$ is linear; (L7) prior information, in the form of a prior probability distribution on the unobservable flows X , is taken into account; (L7') starting values for X are needed; (L8) X and Y are integer random numbers. We present two examples below.

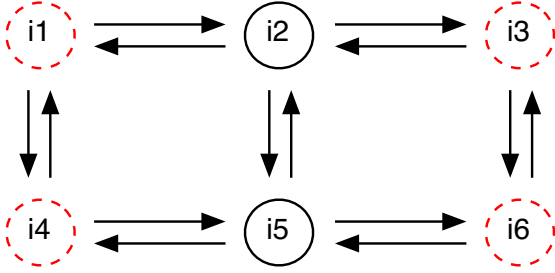


Figure 3: Proportional routing.

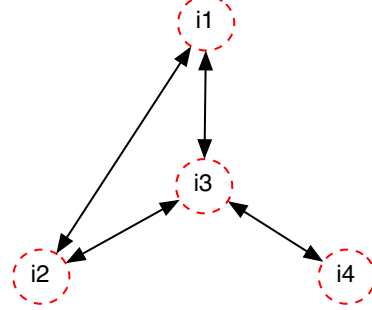


Figure 4: Unique path routing.

Example 1: (Cascetta and Nguyen) Consider the 6 node, 14 link and 4 centroid network structure in figure 3. Nodes $\{i_1, i_3, i_4, i_6\}$ are centroids (our terminal nodes), and we observe the loads Y_l on the links $\{(1, 4), (4, 1), (2, 5), (5, 2), (3, 6), (6, 3)\}$. The assignment mapping $A X$ is

| Link | X_{13} | X_{14} | X_{16} | X_{31} | X_{34} | X_{36} | X_{41} | X_{43} | X_{46} | X_{61} | X_{63} | X_{64} |
|-------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|
| (1,4) | 0.20 | 0.86 | 0.33 | | 0.33 | 0.02 | | | | | | 0.20 |
| (4,1) | | | | 0.20 | | | 0.86 | 0.33 | 0.20 | 0.33 | 0.02 | |
| (2,5) | 0.10 | 0.12 | 0.33 | 0.10 | 0.33 | 0.12 | | | 0.10 | | | 0.10 |
| (5,2) | 0.10 | | | 0.10 | 0.33 | | 0.12 | 0.33 | 0.10 | 0.33 | 0.12 | 0.10 |
| (3,6) | | 0.02 | 0.33 | 0.20 | 0.33 | 0.86 | | | 0.20 | | | |
| (6,3) | 0.20 | | | | | | 0.02 | 0.33 | | 0.33 | 0.86 | 0.20 |

Table 2: Routing scheme with proportional assignment.

Example 2: (Vardi) Consider the 4 node, 8 link and 4 centroid network structure in figure 4. All the nodes are centroids, and we observe the loads on the links Y_l on $\{(1, 2), (2, 1), (1, 3), (3, 1), (2, 3), (3, 2), (3, 4), (4, 3)\}$. Notice that in this case the zero sum constraint on the OD flows makes one of the observations redundant. The assignment mapping $A X$ is

| Link | X_{12} | X_{13} | X_{14} | X_{21} | X_{23} | X_{24} | X_{31} | X_{32} | X_{34} | X_{41} | X_{42} | X_{43} |
|-------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|
| (1,2) | 1 | | | | | | | | | | | |
| (1,3) | | 1 | 1 | | | 1 | | | | | | |
| (2,1) | | | | 1 | | 1 | 1 | | | 1 | | |
| (2,3) | | | | | 1 | | | | | | | |
| (3,2) | | | | | | | 1 | 1 | | 1 | 1 | |
| (3,4) | | | 1 | | | 1 | | | 1 | | | |
| (4,3) | | | | | | | | | | 1 | 1 | 1 |

Table 3: Routing scheme with unique paths.

The **estimation problem** boils down to finding reasonable origin-destination flows X that match the observed link loads. The assignment mapping is usually not invertible and there is more than one feasible solution.

2.1 Transportation Research

A relevant part of the literature dealt with transportation analysis (the number of travelers that commute, or the amount of freight shipped), and with traffic monitoring problems (the number of cars that move between different metropolitan areas of a city). If we exclude time (L1) and endogeneity of the assignment mapping (L3) the all proposed solutions to the estimation problem boiled down to the following problem.

Minimization problem:

$$\min_X p \cdot D_1(X, \{X_{obs}\}) + q \cdot D_2(Y(X), \{Y_{obs}\}) \quad \text{s.t.} \quad Y = f_A(X) \quad \text{and} \quad X \geq 0. \quad (2.4)$$

where D_1 and D_2 are distance measures, and the constraints make sure X, Y represent one of the positive, feasible solutions. Notice that the solution X is constant over time.

Intuition: in general there are many feasible origin-destination traffic matrices⁵ X matching the observed loads Y , in fact in section 1.1 we counted the number of unobservable flows $\kappa = O(\ell^2)$, ℓ being the number of available measurements. The minimization helps us choose among all the positive, feasible solutions the one OD traffic matrix (\hat{X}) that best matches the constraints on the measurements we have in terms of D_2 , and is closest according to D_1 to the prior information we are willing to consider, properly weighting these criteria through p and q . e.g. D_1 could be the entropy, D_2 the Euclidean distance, and $Y = f_A(X) = A X$.

The solutions given over the years were both deterministic and stochastic in nature, and found their rationale rooted in maximum likelihood estimation, generalized least squares, Bayesian inference, information theory and economics.

2.2 Statistical Research

In Statistics, early works dealt with inferences for the inner cells of a table, given values for the margins — again a constant solution would be obtained for the inner cells, starting from several sets of margins. A geometrical framework for the tables was given, then the focus shifted to different models for the counts, to slowly time-varying tables with constant parameters over time, and eventually to real valued cell entries and time-varying parameters.

Tables

Deming and Stephan (1940) first tackled the problem and gave their Iterative Proportional Fitting Procedure (IPFP) that returns a feasible solution for the inner cells, given any set of positive starting points that do not necessarily meet the constraints on the margins. Fienberg (1968) showed how any $(r \times c)$ two-way table can be represented by points within the $(rc - 1)$ dimensional simplex, and a complete account of the geometry of such objects is given in Fienberg (1968) and Fienberg (1970a). From Fienberg (1970b) we learned that any $(r \times c)$ two-way table can be identified by its margins and the $(r - 1) \cdot (c - 1)$ association coefficients $\alpha_{i,j}$. Further given any starting positive OD

⁵Recall that X is a vector that derives from a matrix through equation 1.2.

matrix — hence a set of $(r - 1) \cdot (c - 1)$ coefficients $\alpha_{i,j}$ — and a set of margins (Y) IPFP would guarantee a solution consistent with the margins, and with the same association coefficients $\alpha_{i,j}$. The path towards the solution lies on the manifold of constant interaction defined by the association coefficients of the initial OD matrix, and ends where the latter meets the $(r - 1)$ manifold defined by the row margins and the $(c - 1)$ manifold defined by the column margins — the three manifolds meet in exactly one point. Ireland and Kullback (1968) showed that IPFP solution minimizes the discrimination information, and not the variation of the χ^2 statistic due to Neyman, as originally conjectured by Deming and Stephan (1940). They showed the solution is BAN.

Origin-Destination Traffic Matrices

Vanderebei and Iannone (1994) assumed independent Poisson traffic counts for the entries of the OD matrix, developed three equivalent formulations of the EM algorithm and studied the fixed points of the EM operator. They were not able to prove that the log-likelihood function for their model was convex in general, but gave some partial results.

Vardi (1996) studied independent Poisson traffic counts, both under fixed and probabilistic routing schemes. He showed that the likelihood function may have an absolute maximum on the boundary whereas the unique solution of the maximum likelihood equations yields an internal, local maximum.

Tebaldi and West (1998) pointed out the need for informative priors in a conjugate and time independent setting. They studied independent Poisson traffic counts $P(X_t | \Theta_t)$ following Vardi (1996), and then choose independent conjugate priors for the parameters $P(\Theta_t)$. At any given time, they wanted to find the posterior distribution $P(X_t, \Theta_t | Y_t)$ using one set of measurements Y_t only. They found that the data was only able to limit the support of the Posterior distribution without adding information about likely values, and the priors would drive their inferences.

Cao et al. (2000) assumed independent Gaussian OD traffic flows, and used EM algorithm to derive estimates for the parameters, that depended on t . Their method involved finding a first approximation for the OD flows at time t (full details in section 2.2.1 below) and then using IPFP to get the final estimates $\hat{X}(t)$ perfectly matching the observed link loads. The main drawbacks of their approach were the assumptions: traffic IID over time, and Gaussian OD flows. Further, since in order to estimate Θ_t (time-dependent) they used a window of contiguous observations centered around time t , their method cannot provide estimates for $X(t)$ at the beginning and at the end of the time series. Cao et al. (2002) proposed the same model, and suggested a divide-and-conquer strategy to deal with larger scale problems as the number of origin-destination pairs grows.

2.2.1 A Recent Local Maximum Likelihood Approach

In this section we present the model proposed by Cao et al. (2000) and (2002), in order to point out some unsatisfactory aspects of it that we fully address in section 3. Denote the origin-destination flows by X s, the link loads by Y s, and the routing matrix by A , as discussed above. The following three equations

$$\begin{cases} X_t \sim N_I(\lambda, \Sigma) \text{ iid} \\ \Sigma = \phi \cdot \text{diag}(\lambda^\tau), \quad \tau \text{ is a known constant} \\ Y_t = AX_t \end{cases} \quad (2.5)$$

define a Gaussian process for the available measurements $Y_t \sim N_D(A\lambda, A\Sigma A')$, that sort of approximates a multivariate Poisson process with parameter vector λ . Sort of approximate since the parameter ϕ allows for extra Poisson variability.

We observe vectors Y_1, \dots, Y_T and we want to estimate the distribution of X_1, \dots, X_T . The local maximum likelihood approach consists in estimating $\Theta_t = (\Lambda_t, \Phi_t)'$, using measurements in a time window centered on t ⁶, to obtain an estimate $\hat{P} = P_{\hat{\Theta}_t}(X_t)$, or possibly $\hat{P} = P_{\hat{\Theta}_t}(X_t|X_t \geq 0)$, for the distribution of the OD flows at time t . Eventually we can compute a point estimate for the unobservable flows using the expected value of \hat{P} , as in $\hat{X}_t = E_{\hat{\Theta}}(X_t|X_t \geq 0) = E(X_t|Y_t, \hat{\Theta}, X_t \geq 0)$.

Four simplifying hypotheses are introduced in carrying out the calculations, namely: (C1) multivariate normality of the unobservable flows (vector X_t) at each time tick; (C2) independence of the flows between distinct pairs of terminal nodes ($X_i(t) \perp X_j(t)$, $i \neq j$, $\forall t$); (C3) independence of the vectors containing the flows over time (vector $X_t \perp X_s$, $t \neq s$); (C4) identical distribution of the flows (vector $X_t \stackrel{D}{=} X_s$, $t, s \in [t_a, t_b]$) at time ticks not too far apart in time ($|t_b - t_a| < \text{const}$).

C1: Multivariate Normality

The assumption (C1) is needed to keep the mathematics manageable. It is in contrast with both empirical evidence and state-of-the-art theoretical models, see for example Leland et al. (1993) or Carmona and Coutin (1998), that suggested that OD traffic flows in a network are bursty, self-similar and their sample paths resemble fractional Brownian motion or log-normal process paths.

C2: Independence of the Flows between Distinct Pairs of Terminal Nodes

The independence assumption (C2) is useful in proving the identifiability of all the parameters, and helps keep the dimensionality of the problem manageable. Intuitively, in the latent κ -dimensional space of the flows there are $\kappa+1$ parameters that need be estimated ($\lambda_1, \dots, \lambda_\kappa$, and ϕ). The linearity of the routing matrix A plus the multivariate normality of the flows X_t yield $\frac{\ell(\ell-1)}{2}$ parameters in the ℓ -dimensional space of the observations, and, since $\kappa = O(\ell^2)$ in our problem, it is possible to estimate the parameters underlying the distribution of the measurements, and invert the mapping provided by $A\Sigma A'$ exactly to recover the parameters of the distribution of the unobservable OD flows⁷.

Thus linearity of the routing matrix, and multivariate normality of the unobservable flows together make the work here. Independence is needed to keep the parameters in the latent space $O(\kappa) = O(\ell^2)$; a more complex variance-covariance matrix would push the parameters towards $O(\kappa^2)$, which would be too many. Assumption (C2) is a convenient modeling idea.

C3: Independence of the Flows over Time

Assumption (C3) is very much unsatisfactory. We are dealing with traffic flows over time, and terminal nodes can be as small as single PC stations⁸. It is somewhat unrealistic to model these flows as independent over time. It is true that the fitting procedure assumes independence of the unobservable flows X_t over time, and then re-introduces some time dependence from the back door, in the form of smoothing, but a more explicit model for the dynamic would be desirable.

In other words, estimates for Θ_t are obtained using the observations Y_t , $t = t - t_W, \dots, t + t_W$, the observations Y_t , $t = t + 1 - t_W, \dots, t + 1 + t_W$ would be used to estimate Θ_{t+1} , and so on. Any

⁶Hence *local* maximum likelihood.

⁷This argument can be made precise making good use of the fact that the routing nodes neither generate, nor absorbs traffic, but that is not the point here.

⁸For example this is the case in small wireless Local Area Networks that are present nowadays in most offices, private houses and public places. And these smaller networks are very much the ones that need to estimate the origin-destination traffic matrix, not being able to submit SNMP queries to the respective routers.

two successive estimates would have the observations Y_t , $t = t + 1 - t_W, \dots, t + t_W$ in common, hence introducing dependence between $\hat{\Theta}_t$ and $\hat{\Theta}_{t+1}$ and, through these, between the estimates \hat{X}_t and \hat{X}_{t+1} . Notice that X_t and X_{t+1} were assumed independent.

C4: Identical Distribution of the Flows

This assumption is questionable; Cao et al. assumed Θ constant for all of the observations in a window $[Y_{t-t_W}, Y_{t+t_W}]$ to be able to estimate it precisely. Eventually non-constant estimates $\hat{\Theta}_t$ are obtained for all of the observations in the window $[Y_{t-t_W}, Y_{t+t_W}]$.

T1: Multivariate Integration

The local likelihood procedure in section 2.2.1 requires to estimate Θ_t first, which completely specifies \hat{P}_t , and eventually computes point estimates \hat{X}_t for the origin-destination flows. Estimating \hat{X}_t requires the following computation

$$\begin{aligned}\hat{X}_t &= E(X_t|Y_t, \hat{\Theta}, X_t > 0) \\ &= \int_{R^+} x_t \cdot f(x_t|y_t, \hat{\theta}, x_t > 0) \cdot dx \\ &= \int_{R^+} x_t \cdot \frac{f(x_t|y_t, \hat{\theta}) I_{\{A^{-1}y_t\}}(x_t)}{\int_{R^+} f(x_t|y_t, \hat{\theta}) dx} \cdot dx\end{aligned}$$

that involves the multivariate integration $\int_{R^+ \cap \{A^{-1}y_t\}} f(x_t|y_t, \hat{\theta}) dx$ of the multivariate Normal distribution

$$N(\hat{\lambda} + \hat{\Sigma}A'(A\hat{\Sigma}A')^{-1}(Y_t - A\hat{\lambda}), \hat{\Sigma} - \hat{\Sigma}A'(A\hat{\Sigma}A')^{-1}A\hat{\Sigma}'), \quad \hat{\Sigma} = \hat{\phi} \cdot \text{diag}(\hat{\lambda}),$$

over the positive orthant, and $I_{\{A^{-1}y_t\}}(x_t)$ involves the computation of the support of $P(X_t|Y_t, \hat{\Theta})$. The alternative solution (T1) proposed by Cao et al. to avoid these issues consisted in making rough first guesses by estimating $\hat{X}_{t,i}$, $i = 1, \dots, I$ independently, as in

$$\hat{X}_i(t) = E(X_i(t)|Y_t, \hat{\Theta}, X_i(t) > 0).$$

These univariate integrations boiled down to the calculation of

$$E(Z|Z > 0) = \mu + \frac{\sigma}{\sqrt{2\pi}} \cdot \exp\left(-\frac{\mu}{2\sigma}\right) \cdot \Phi^{-1}\left(\frac{\mu}{\sigma}\right),$$

where $Z := \hat{X}_{t,i} \sim N(\mu, \sigma^2)$, and

$$\begin{aligned}\mu &= [\lambda + \Sigma A'(A\Sigma A')^{-1}(y_t - A\lambda)]_i \\ \sigma^2 &= [\Sigma - \Sigma A'(A\Sigma A')^{-1}A\Sigma']_{i,i}\end{aligned}$$

At this point the IPFP was used to get final estimates \hat{X}_t that satisfied the set of constraints $Y_t = A\hat{X}_t$ provided by the available set of measurements.

Some Remarks

We discussed above the assumption (C1) to (C4) underlying the model in Cao et al. and a particular solution (T1) they propose to deal with a nasty multivariate integration problem.

- Multivariate normality of the unobservable OD flows (C1) contrasts with both empirical and theoretical evidence, and is unsatisfactory. The independence of the unobservable flows (C3) could be relaxed introducing some explicit form of dynamics.
- The solution (T1) proposed to deal with the multivariate integration is unsatisfactory as we may be spoiling our efforts to obtain good estimates by integrating independently component by component.
- The local likelihood approach does not provide estimates that cover the whole sequence of times for which we have measurements available. This point is particularly painful if we consider that we can measure the link loads every 5 minutes, so that even using a pretty minimal window of 5 time ticks we would be ten minutes behind in terms of estimated OD traffic flows.

3 Proposed Methods

We want to make inferences about non-observable origin-destination flows in a Local Area Network, starting from a set of measurable link loads. In such a situation a higher likelihood of the measurements may or may not yield better estimates of the non-observable flows. Good inferences would require a realistic model, able to capture relevant features of the data. Our best model contains explicit time dependence of the OD flows, a stochastic dynamical behavior, and skewed distributions for the OD flows; it is a step forward in comparison to the models present in the literature both in terms of degree of realism of the underlying assumptions, and goodness of the inferences. The Bayesian dynamical system we propose outperformed state-of-the-art models by reducing the estimation error of more than 45%, in ℓ_2 distance⁹, in a realistic setting.

The major problem we had to cope with was the existence of multiple modes in the filtering distributions $P(OD_t | Y_1, \dots, Y_t)$ at each time t . We solved it by introducing informative priors on certain parameters governing the dynamics of the system in order to identify a single, most likely posterior mode among possibly many of them.

In section 3.1 we introduce time dependence and make use of a simple linear Gaussian system to obtain preliminary estimates for the OD flows; in section 3.2 we introduce more realistic Gamma and log-Normal models for the OD flows, and we properly use the measurements at every time point t in order to make inferences on the OD flows at the same time t ; in section 3.3 we combine realistic non-Gaussian models for the traffic flows and explicit dynamical behavior by means of a Bayesian dynamical system, and obtain excellent estimates for the OD flows.

3.1 Explicit Dynamics for Gaussian Origin-Destination Flows

The dynamical model we propose here entails first-order Markov dependence between the non-observable flows at different time points, and its corresponding graphical representation is displayed in the right panel of figure 5 below. The left panel, instead, displays the graphical representation of the recent model proposed in Cao et al. (2000), where the absence of arrows between the non-observable (dashed) nodes indicates that the OD traffic flows at different times are independent. As the graphs suggest our model includes Cao’s model as a special case. Briefly, the time dependence we introduced among the non-observable OD flows with the graphical representation in figure 5,

⁹We performed several experiments using validation data, obtained by monitoring about 12321 OD traffic flows at Carnegie Mellon University.

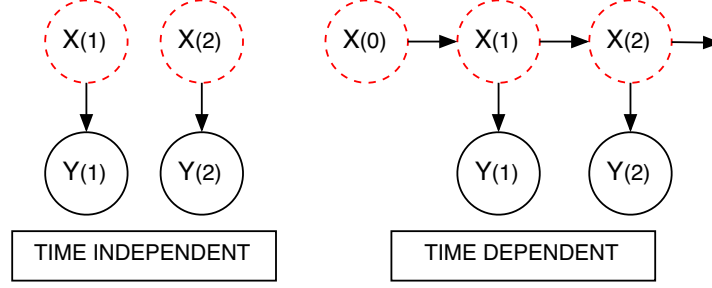


Figure 5: Left: graphical representation for a model with independent OD flows. Right: graphical representation for the model we propose in section 3.1 with first-order Markov dependence among the OD flows.

mathematically translates into a contribution of $P(X_{t+1}|X_t)$ in computing the posterior distribution $P(X_{t+1}|Y_{t+1}, \dots, Y_1)$ as the few lines below show for $t + 1 = 2$:

$$P(X_2|Y_2, Y_1) \propto \int P(X_2, X_1, Y_2|Y_1) dX_1$$

$$(\text{State-Space computation includes } P(X_2|X_1)) = \int P(X_2|X_1) P(X_1|Y_1) P(Y_2|X_2) dX_1$$

$$(\text{if the time points were IID } P(X_2|X_1) = P(X_2)) = P(X_2) P(Y_2|X_2) \cdot \int P(X_1|Y_1) dX_1$$

$$(\text{and the posterior would simplify into}) = P(X_2) P(Y_2|X_2)$$

3.1.1 A State-Space Representation for the Model

The model we used is defined as follows:

$$\begin{aligned} & \begin{cases} X_{t+1} = F_{t+1} \cdot X_t + Q_{t+1} \cdot \iota + \epsilon_{t+1} \\ Y_t = A \cdot X_t + \eta_t \end{cases} \\ &= \begin{cases} \begin{bmatrix} X_{t+1} \\ 1 \end{bmatrix} = \begin{bmatrix} F_{t+1} & Q_{t+1} \\ 0 & I \end{bmatrix} \cdot \begin{bmatrix} X_t \\ \iota \end{bmatrix} + \begin{bmatrix} \epsilon_{t+1} \\ 1 \end{bmatrix} \\ Y_t = [A|0] \cdot \begin{bmatrix} X_t \\ 1 \end{bmatrix} + \eta_t \end{cases} \quad (3.1) \\ &= \begin{cases} \tilde{X}_{t+1} = \tilde{F}_{t+1} \cdot \tilde{X}_t + \tilde{\epsilon}_{t+1} \\ Y_t = \tilde{A} \cdot \tilde{X}_t + \eta_t \end{cases} \end{aligned}$$

for $t \geq 1$, where $\iota = \iota \cdot (1, \dots, 1)'$ is a constant vector of the length κ , the parameter ϕ_t enters into the variance-covariance matrix of $\epsilon_t \sim N(0, \phi_t \cdot Q_t)$, $X_1 \sim N(0, V_1)$, $\eta_t \sim N(0, R_t)$, $X_1 \perp \epsilon_t$ and $X_1 \perp \eta_t$ for all $t \geq 1$, and finally Q_t is a diagonal matrix with elements $(q_{t,1}^\tau, \dots, q_{t,\kappa}^\tau)$, where τ is a known constant. In the model above whenever $F_t = 0$ there is a one-to-one mapping between $(q_{t,1}^\tau, \dots, q_{t,\kappa}^\tau, \phi_t)'$ and the unique elements in $E(Y_t), V(Y_t)$. Further the following lemma holds.

Lemma. *The linear Gaussian state-space model in equations 3.1 contains the model in Cao et al. (2000) defined by equations 2.5 in section 2.2.1 as a special case. Such a model can be obtained by simply setting $\iota = 1$ and $F_t = 0, \forall t$, hence imposing independence among the origin-destination flows X_t at different times.*

Local Linear Dynamics

In this report we considered a local linear dynamics for the OD flows. The local linearity assumption is not really a constraint as any non-linear dynamical behavior can be locally approximated to the first order by a linear behavior. Extensions of the Kalman recursions as the Extended Kalman Filter or the Unscented Kalman Filter deal more precisely with non-linear dynamics, but a local linear dynamical behavior will do for us without the need for further complications. In section 4.3.2 we provide some empirical evidence to support our choice for a diagonal F_t in our model.

3.1.2 Ad-Hoc M-Step for the EM Algorithm

In order to estimate the OD flows at successive times it was natural to compute the sequence of posterior distributions $X_1|Y_1$, $X_2|Y_1$ and $X_2|Y_1, Y_2$, $X_3|Y_1, Y_2$ and $X_3|Y_1, Y_2, Y_3$, and so on. We used the Kalman Filter to recursively compute filtered and smoothed posterior probability distributions for the OD flows; the celebrated Kalman Filter provides in fact recursions to find MS-optimal estimates of the hidden states once we know the other parameters of the state-space model — $F_t, Q_t, \phi_t, R_t, X_1, V_1$.

We needed estimates for these parameters as well, and a major concern was the estimation of F_t , because of its possibly enormous variance. In our case it was possible to write down the likelihood as:

$$f((Y_1, \dots, Y_T) | \Theta) = (2\pi)^{-T\ell/2} \left(\prod_{i=1}^T \det \Sigma_i \right)^{-1/2} \exp \left\{ -\frac{1}{2} \sum_{i=1}^T (Y_i - \hat{Y}_i)^T \Sigma_i^{-1} (Y_i - \hat{Y}_i) \right\},$$

and maximize it directly¹⁰, or alternatively use the EM algorithm, in the spirit of Ghahramani and Hinton (1996). Let's consider the simple case $F_t = F, Q_t = Q$ and $R_t = R$ for all $t \geq 1$. The parameters would be $\{X\}, A, F, Q, R, X_1, V_1$, and the EM algorithm would give the recursions to update the quantities involved, including the OD flows, $\{X\}$.

The E-step

The quantity of interest is the expected value of the log-likelihood of the complete data, given the observations $\mathcal{Q} = E[\log P(\{X\}, \{Y\}) | \{Y\}]$. We can write:

$$\begin{aligned} \log P(\{X\}, \{Y\}) &= \log P(X_1) \cdot \prod_{t=2}^T P(X_t | X_{t-1}) \cdot \prod_{t=1}^T P(Y_t | X_t) \\ &= -\frac{T(\ell+I)}{2} \log 2\pi - \frac{1}{2} X_1' V_1^{-1} X_1 - \frac{1}{2} \log |V_1| \\ &\quad - \frac{1}{2} \sum_{t=2}^T ([X_t - F X_{t-1}]' Q^{-1} [X_t - F X_{t-1}]) - \frac{T-1}{2} \log |Q| \\ &\quad - \frac{1}{2} \sum_{t=2}^T ([Y_t - A X_t]' R^{-1} [Y_t - A X_t]) - \frac{T}{2} \log |R| \end{aligned}$$

and making use of some properties of the trace operator it is easy to see that the expectation $\mathcal{Q} = E[\log P(\{X\}, \{Y\}) | \{Y\}]$ depends on the three quantities $E[X_t | \{Y\}]$, $E[X_t X_t' | \{Y\}]$ and $E[X_t X_{t-1}' | \{Y\}]$. The derivations of the recursions needed for the E-step follows that of Ghahramani and Hinton (1996), for the model in section 3.1.1 in terms of the new variables (tilde).

¹⁰In the equation above \hat{Y}_i and Σ_i are the means and variances of the one-step-ahead projections.

The M-step

Denote the three quantities $E[X_t|\{Y\}]$, $E[X_t X'_t|\{Y\}]$ and $E[X_t X'_{t-1}|\{Y\}]$ by \hat{X}_t , P_t and $P_{t,t-1}$ respectively, for convenience. The parameters which we needed to maximize \mathcal{Q} over were F , Q , ϕ , X_1 , V_1 , since A was given by the fixed routing scheme, and the variance-covariance matrix of the measurements R was zero¹¹. In our model $F=F(Q)$, in terms of the original variables (no tilde) for the model in section 3.1.1, and there is a new parameter ϕ so that we need to add updating equations for Q and ϕ as in:

Q : variance-covariance matrix of the OD flows.

$$\frac{\partial \mathcal{Q}}{\partial Q^{-1}} = 0 \quad \text{yields} \quad \frac{T-1}{2} (Q^{-1} - Q 11' Q) - \frac{1}{2} \left(\sum_{t=2}^T P_t - F^{new} \sum_{t=2}^T P_{t-1,t} \right) = 0 \quad (3.2)$$

ϕ : extra-Poisson variability for the OD flows.

$$\frac{\partial \mathcal{Q}}{\partial \phi} = 0 \quad \text{yields} \quad \phi^{new} = \frac{1}{T-1} \left(\sum_{t=2}^T P_t - F^{new} \sum_{t=2}^T P_{t-1,t} \right) \quad (3.3)$$

3.1.3 Two-Stages Maximization of the Likelihood

We also maximized the log-likelihood (\mathcal{L}) directly, in two stages. The first stage consisted of maximizing \mathcal{L} with respect to the all parameters, but F_t , using an implementation of the BFGS algorithm that allowed for box constraints. In the second stage we maximized \mathcal{L} with respect to F_t . Our experiments included multiple starting points, early stopping and regularization.

3.2 1-Time Non-Gaussian Origin-Destination Flows

In this section we leave the safe shores of Normality to introduce Gamma and log-Normal models for the OD flows in a time-independent Bayesian context; that is we properly aim to make inferences about the non-observable OD flows X_t by using the information contained in the available measurements at time t only, Y_t . Further we use informative priors to mitigate the problem of multiple posterior modes. Briefly, at each time t , given the link measurements Y_t , the information contained in the routing scheme $Y_t = A X_t$, and a model for the OD flows $P(X_t|\Theta_t)$ we computed the joint posterior distribution $P(X_t, \Theta_t|Y_t)$, and eventually found point estimates for the OD traffic flows \hat{X}_t .

The equations $Y_t = A X_t$ allow infinite solutions for X_t . A model for the flows $P(X_t|\Theta_t)$ induces a probabilistic mapping on $\{X_T : Y_t = A X_t\}$, but even so the posterior $P(X_t, \Theta_t|Y_t)$ may still have more than one mode. In fact as we used non-informative priors, the information contained in the data Y_t improved our knowledge about the OD flows by setting constraints on the support of $P(X_t, \Theta_t|Y_t)$ only. As a consequence, using the posterior mean as a point estimate for the OD flows would yield estimated flows either too high or too low, in the presence of multiple modes. Informative priors on the other hand were able to drive our inferences in the correct direction by making the posterior *more unimodal*.

The relevant issues here were how to sample efficiently in a highly constrained space via MCMC, and how to ensure the MCMC could explore the whole space where the OD flows lived (irreducibility of the chain). We checked convergence of the various chains using Raftery and Lewis, Geweke, and

¹¹We fixed $R = Id_D \cdot \epsilon$ as a practical solution to problem of computing R^{-1} in the likelihood.

Heidelberger and Welch convergence diagnostics. The results in this section extend to continuous, non-conjugate analysis the results in Tebaldi and West (1998). In the remainder of this section we drop the time subscripts, since the models we propose are to be fitted time by time.

3.2.1 Computing the Support

The routing matrix A is full row rank, by construction. In order to sample efficiently we took advantage of the following fact.

Fact. *There exists a permutation ρ of the columns of $A_{(\ell \times \kappa)}$ such that $[A]_{(i, \rho(j))} = [A_1 \mid A_2]$, where A_1 is $(\ell \times \ell)$ and has full rank, and A_2 is $(\ell \times (\kappa - \ell))$.*

As a consequence we were able to permute the components of X to get $[X]_{\rho(i)}' = [X_1 \mid X_2]'$, and $Y = AX = A_1 X_1 + A_2 X_2$, and we then ran chains in the lower dimensional space where X_2 lived, obtaining $X_1(X_2, Y)$ at each step, like so:

$$X_1 = A_1^{-1} \cdot (Y - A_2 X_2)$$

The space where X_2 lived was itself constrained by the available measurement Y , at any given time. Using a χ^2 proposal in the Metropolis steps, properly defining the full conditionals $P(X_{2,i} | X_{2,-i}, Y)$ to be zero if, say, a negative value for $X_{2,i}$ was sampled as a candidate, and accepting the candidate conditionally to a further check that all the corresponding $X_1(X_2, Y)$ were positive, allowed us to avoid a direct and expensive computation of the support.

3.2.2 Irreducibility of the Chain

We need to make sure that the chains we used were able to explore the entire space of OD traffic flows. It is possible to show that the *positivity condition* introduced by Besag (1974) holds, which is sufficient to ensure the irreducibility of the chain. Moreover we show that the support of the univariate full conditional distributions is convex.

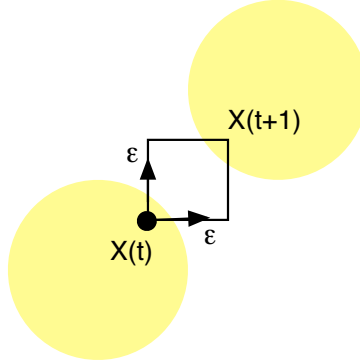


Figure 6: The chain is at $X(t)$, the circles represent the joint support. The possible moves according to the Gibbs scheme $X(t) + [0, \epsilon]'$ or $X(t) + [\epsilon, 0]'$ would yield outside of joint support, whereas a move $X(t) + [\epsilon, \epsilon]'$ would yield inside. We show that such a situation cannot happen, and whenever $X(t) + [\epsilon, \epsilon]'$ yields a *valid* point of the support of the joint distribution, the Gibbs moves also do.

Proof: (convexity of the support) We want to show that if the two $(\kappa - \ell)$ -dimensional vectors $x = (x_1, \dots, x_i, \dots, x_{\kappa-\ell})'$ and $y = (x_1, \dots, y_i, \dots, x_{\kappa-\ell})'$ differ in the i^{th} component, and if both

$A_1^{-1}(Y - A_2 x) \geq 0$ and $A_1^{-1}(Y - A_2 y) \geq 0$, then also $z = (x_1, \dots, \alpha x_i + (1 - \alpha)y_i, \dots, x_{\kappa-\ell})$ is such that $A_1^{-1}(Y - A_2 z) \geq 0$. This is easily proved true since from $z = \alpha x + (1 - \alpha)y$ follows that $A_1^{-1}(Y - A_2 z) = \alpha A_1^{-1}(Y - A_2 x) + (1 - \alpha) A_1^{-1}(Y - A_2 y)$, which is ≥ 0 being an average of two positive quantities. QED.

Proof: (irreducibility) The Gibbs sampler scheme involves iterative sampling from the full conditional distributions $P(Z_i | Z_{(-i)} = z_{(-i)})$, for $i = 1, \dots, N$ and Z vector. A sufficient condition to ensure the irreducibility of the chain, Besag (1974), requires that the support of the full conditional distributions is positive where that of the joint distribution of Z is positive, that is:

$$\text{if } P(Z_i = z_i, Z_{(-i)} = z_{(-i)}) > 0 \Rightarrow P(Z_i | Z_{(-i)} = z_{(-i)}) > 0. \quad (3.4)$$

2D case: we show that condition 3.4 holds. Specifically consider the situation displayed in figure 6 above, where there are $\kappa - \ell = 2$ components of X_2 that we need to sample from. The chain is at a point $X_2 > 0$ where the joint support is positive and $A_1^{-1}(Y - A_2 X_2) > 0$, and it moves by $(+\epsilon, +\epsilon)'$ to the point $X_2 + (\epsilon, \epsilon)'$ where the joint support is also positive and $A_1^{-1}(Y - A_2 (X_2 + (\epsilon, \epsilon)')) > 0$. We want to show that whenever both X_2 and $X_2 + (\epsilon, \epsilon)'$ are feasible, it is possible to pass from the former to the latter by means of component-wise moves, as we would with Gibbs moves; that is, the support of the full conditionals must be positive either at $A_1^{-1}(Y - A_2 (X_2 + (0, \epsilon)'))$ or at $A_1^{-1}(Y - A_2 (X_2 + (\epsilon, 0)'))$. In other words we want to show that

$$\{ A_1^{-1}(Y - A_2 X_2) \geq 0 \quad \wedge \quad A_1^{-1}(Y - A_2 (X_2 + (\epsilon, \epsilon)')) \geq 0 \} \quad (3.5)$$

implies

$$\{ A_1^{-1}(Y - A_2 (X_2 + (\epsilon, 0)')) \geq 0 \quad \vee \quad A_1^{-1}(Y - A_2 (X_2 + (0, \epsilon)')) \geq 0 \}. \quad (3.6)$$

Assume that 3.5 holds. Notice that $A_1^{-1}(Y - A_2 (X_2 + (\epsilon, \epsilon)')) = A_1^{-1}(Y - A_2 X_2 - \epsilon(A_2^{11}, A_2^{21})' - \epsilon(A_2^{12}, A_2^{22})') \geq 0$. Add $A_1^{-1}(Y - A_2 X_2) \geq 0$, non negative by assumption, and rearrange terms to get $A_1^{-1}(Y - A_2 X_2 - \epsilon(A_2^{11}, A_2^{21})') + A_1^{-1}(Y - A_2 X_2 - \epsilon(A_2^{12}, A_2^{22})') \geq 0$ which cannot be the sum of two negative quantities. QED.

Similar derivations show that whenever the joint support has positive probability at $A_1^{-1}(Y - A_2 (X_2 - (\epsilon, \epsilon)'))$ then it also possible for the chain to get there either through $A_1^{-1}(Y - A_2 (X_2 - (0, \epsilon)'))$ or through $A_1^{-1}(Y - A_2 (X_2 - (\epsilon, 0)'))$; and that the same condition holds as we consider the moves to the points $A_1^{-1}(Y - A_2 (X_2 + (\epsilon, -\epsilon)'))$ and $A_1^{-1}(Y - A_2 (X_2 + (-\epsilon, \epsilon)'))$.

General case: the proof is exactly the same as in the 2D case, but more tedious. Now X_2 and $(\epsilon, \dots, \epsilon)'$ are $\kappa - \ell = n$ -dimensional. Assume a $A_1^{-1}(Y - A_2 X_2) \geq 0$ and $A_1^{-1}(Y - A_2 (X_2 + (\epsilon, \dots, \epsilon)')) \geq 0$ hold true. Rewrite $A_1^{-1}(Y - A_2 (X_2 + (\epsilon, \dots, \epsilon)'))$ as $A_1^{-1}(Y - A_2 X_2 - \epsilon(A_2^{11}, A_2^{21}, \dots, A_2^{n1})' - \dots - \epsilon(A_2^{1n}, A_2^{2n}, \dots, A_2^{nn})')$ ≥ 0 . Add $(n - 1) \times A_1^{-1}(Y - A_2 X_2) \geq 0$, non negative by assumption, and rearrange terms to get $A_1^{-1}(Y - A_2 X_2 - \epsilon(A_2^{11}, A_2^{21}, \dots, A_2^{n1})') + \dots + A_1^{-1}(Y - A_2 X_2 - \epsilon(A_2^{1n}, A_2^{2n}, \dots, A_2^{nn})') \geq 0$, which cannot be the sum of n negative terms. QED.

Again similar derivations show that condition 3.4 holds as we consider moves to other points $X_2 + (\pm\epsilon, \dots, \pm\epsilon)'$.

3.2.3 Gamma and log-Normal Models

In this section we give some details about the models for the origin-destination flows X_t . We drop the time index for convenience of exposition, and subscripts refer to components of X_t . We introduced Gamma and log-Normal models with a parsimonious parameterization that related means

and variances. i.e. we assumed single OD flows $X_i \sim \text{Gamma}(\alpha_i, \beta)$ with distinct parameters α_i , and we imposed a common scale parameter β , as in

$$X_i \sim \text{Gamma}(\alpha_i, \beta) = \frac{X_i^{\alpha_i-1} e^{-\frac{X_i}{\beta}}}{\beta^{\alpha_i} \Gamma(\alpha_i)}, \quad X_i > 0, i = 1, \dots, \kappa, \quad (3.7)$$

or, alternatively, we assumed single OD flows $X_i \sim \text{log-Normal}(\mu_i, \phi \cdot \mu_i^k)$, for a fixed constant k , with distinct means and variances proportional to them times a common scaling factor ϕ , as in

$$X_i \sim \text{log-Normal}(\mu_i, \phi \mu_i^k) = \frac{1}{\sqrt{2\pi\phi\mu_i^k}} e^{-\frac{1}{2\phi\mu_i^k}(\log(X_i)-\mu_i)^2}, \quad X_i > 0, i = 1, \dots, \kappa, \quad (3.8)$$

Full Conditionals

Say $\Theta = (\alpha_1, \dots, \alpha_\kappa, \beta)'$ then $P(X, \Theta) = \prod_{i=1}^\kappa P(X_i|\Theta) P(\Theta) = \prod_{i=1}^\kappa P(X_i|\alpha_i, \beta) P(\alpha_i) P(\beta)$. We wanted $\alpha_i \in (0, \infty)$ and $\beta \in (0, \infty)$. We assumed improper priors for α_i and $1/\beta$. Then, noticing that $P(\Theta|X, Y) = P(\Theta|X) I_{\{A^{-1}Y\}}(X)$, we obtained the following full conditionals.

$$\begin{aligned} P(\tfrac{1}{\beta}|X, Y) &\propto \prod P(X_i|\alpha_i, \beta) \cdot P(\tfrac{1}{\beta}) \\ &\propto \exp\left\{\sum(\alpha_i - 1)\log(X_i) - \sum \frac{X_i}{\beta} + \sum \alpha_i \log(\tfrac{1}{\beta})\right\} \\ &\propto \text{Gamma}\left[\sum \alpha_i + 1, \frac{1}{\sum X_i}\right] \end{aligned}$$

$$\begin{aligned} P(\alpha_i|X, Y) &\propto P(X_i|\alpha_i, \beta) \cdot P(\alpha_i) \\ &\propto \exp\left\{(\alpha_i - 1)\log(X_i) - \frac{X_i}{\beta} - \alpha_i \log(\beta) - \log \Gamma(\alpha_i)\right\} \\ &\propto \frac{e^{-\alpha_i \log(\frac{\beta}{X_i})}}{\Gamma(\alpha_i)}, \quad \alpha_i > 0 \end{aligned}$$

whereas for $P(X|Y, \Theta)$ we used the fact in 3.2.1 to conclude that $P(X|Y, \Theta) = P(X_2|Y, \Theta) \times P(X_1(X_2)|Y, \Theta)$; hence for $X_i \in X_2$ and $X_j \in X_1$ it followed:

$$\begin{aligned} P(X_i|X_{(-i)}, Y, \Theta) &\propto P(X_i|\Theta) \cdot P(X_1|Y, \Theta) \\ &= \text{Gamma}_{X_i}(\alpha_i, \beta) \cdot \prod_j \text{Gamma}_{X_j}(\alpha_j, \beta) I_{\{A^{-1}Y\}}(X_j) \end{aligned} \quad (3.9)$$

We explored the posterior distributions using the Gibbs algorithm, with Metropolis steps to sample from $P(X_i|Y, \Theta)$ and $P(\alpha_i|X, Y)$, using χ^2 and Uniform proposals, an improper prior for alpha proportional to a constant, and several priors for β : one proportional to a constant, one proportional to $\frac{1}{\beta}$, and one proportional to $\frac{1}{\beta^2}$.

In the same fashion we obtained the full conditionals for the log-Normal model, with constant prior for μ_i and priors proportional to $\frac{1}{\phi^2}$ (used in the calculations below), and to $\frac{1}{\phi}$ for ϕ .

$$\begin{aligned} P(\mu_i|X, Y) &\propto \prod P(X_i|\mu_i, \phi) \cdot P(\mu_i) \\ &\propto \frac{1}{\mu_i^{\frac{k}{2}}} e^{-\frac{1}{2\phi} \left(\frac{\log(X_i)-\mu_i}{\mu_i^k}\right)^2} \end{aligned}$$

$$\begin{aligned} P(\phi|X, Y) &\propto P(X_i|\mu_i, \phi) \cdot P(\phi) \\ &\propto \frac{1}{\phi^{\frac{k}{2}+2}} e^{-\frac{1}{2\phi} \sum_i \left(\frac{\log(X_i)-\mu_i}{\mu_i^k}\right)^2} \end{aligned}$$

$$\begin{aligned} P(X_i|X_{(-i)}, Y, \Theta) &\propto P(X_i|\Theta) \cdot P(X_1|Y, \Theta) \\ &= \text{log-Normal}_{X_i}(\mu_i, \phi) \cdot \prod_j \text{log-Normal}_{X_j}(\mu_j, \phi) I_{\{A^{-1}Y\}}(X_j) \end{aligned}$$

3.2.4 Informative priors

In order to mitigate the loss of precision of the estimates \hat{X}_t , due to the many feasible solutions, we introduced informative priors for the parameters in $\Theta = (\alpha_i, \beta) = (\mu_i, \phi)$ governing the distribution of the OD flows. In doing so we basically imposed *soft* constraints on the parameters, and made the posterior distributions *more unimodal*.

In the Gamma model, the priors for the α_i s were truncated Normal distributions with a huge variance, centered around the preliminary estimates obtained with the linear dynamical system in section 3.1, and the prior for β was proportional to $\frac{1}{\beta^2}$. Similarly, for the log-Normal model we use truncated Normal priors for the μ_i s and a prior proportional to $\frac{1}{\phi^2}$ for ϕ . As a result the posterior mode closer to the preliminary estimates would become more likely than the other modes, and the bias introduced by these other modes would fade. Hypothetically, though, if the preliminary estimates were off by large from the true OD flows, and at the same time corresponded to a posterior mode, then the estimates would be driven towards the wrong feasible solution. This never happened in our experiments in section 4, where the estimates consistently improved. Further in 60% of the time points we obtained better estimates even with non-informative priors. The results we obtained were quite interesting.

3.3 Combining Dynamics and Non-Gaussianity

In this we section combine skewed distributions and explicit time dependence for the OD flows, and we show how we mitigated the loss of precision of the estimates due to the possibly multiple posterior modes in this dynamic setting. In the absence of a dynamical behavior suggested by some physical law, or known to some degree, we took leave from the classical analysis of time series and found the solution, again, in the use of informative priors. in a Bayesian dynamical system.

Briefly a Bayesian dynamical system provides a set of posterior distributions $P(X_t, \Theta_t | Y_t, \dots, Y_1)$ as t varies, whereas in section 3.2 we obtained a series of posterior distributions $P(X_t, \Theta_t | Y_t)$. In this setting we used informative priors on the parameters underlying the stochastic dynamics F_t of the Bayesian dynamical systems defined in section 3.3.2, instead of informative priors on X_t . We estimated the parameters using a Particle Filter; one of the problem of particle filters, partially addresses by the resample-move algorithm proposed by Gilks and Berzuini (2001), is to guarantee a set of particles *rich enough* to describe the state of the system at any given time. We used the sequence of posterior distributions $P(X_t, \Theta_t | Y_t)$, that contain information about the state of the system, to improve our inferences.

3.3.1 Informative Priors for Stochastic Dynamics

We propose the use of informative priors on the explicit stochastic dynamics F_t , in order to provide a set of particles $\Theta^{(t)}$, at every time t , such that their distribution $P(\Theta^{(t)}) \approx P(\Theta_t | Y_t)$ would correspond to the marginal posterior we obtained in the static Bayesian analysis of section 3.2. In other words we propose a stochastic F_t that solves the convolution problem $\Theta_{t+1} = F_t \Theta_t$, given that we have roughly good ideas about how Θ_t should look like at every time t , that is, according to $P(\Theta_t | Y_t)$.

In the Gamma model the convolution could not be solved exactly, since the product of Gamma distributions is no longer a Gamma. We took advantage of the fact that the ratio of Gamma distributions is Pearson Type VI, and we estimated the parameters underlying the random variables $F_t \sim \text{InverseGamma}$, such that $\Theta_{t+1} = F_t \Theta_t$, where $\Theta_t \sim \text{Gamma}$, and $\Theta_{t+1} \sim \text{PearsonTypeVI}$. The log-Normal distribution has the nice property that the product of log-Normals is log-Normal,

hence we were able to solve the convolution problem exactly at all times. We estimated the parameters underlying the random variables $F_t \sim \log\text{-Normal}$, such that $\Theta_{t+1} = F_t \Theta_t$, where $\Theta_t, \Theta_{t+1} \sim \log\text{-Normal}$.

3.3.2 Bayesian Dynamical Systems

For the Gamma model define $\Theta_t := (\alpha_t, \beta_t)'$ and $\Psi_t := (\gamma_t, \delta_t)'$. Then:

$$\begin{cases} \Theta_{t+1} &= F_t \cdot \Theta_t \\ P(X_t^i | \Theta_t) &\sim \text{Gamma}(\alpha_{t,i}, \beta_t) \\ Y_t &= A \cdot X_t, \quad t \geq 1 \end{cases} \quad i = 1, \dots, \kappa \quad (3.10)$$

where $P(\Theta_0^i) \sim \text{Gamma}_i$, and $P(F_t^{ii} | \Psi_t) \sim \text{Inv Gamma}(\gamma_{t,i}, \delta_{t,i})$, $i = 1, \dots, \kappa$. For the log-Normal model define $\Theta_t := (\mu_t, \phi_t)'$ and $\Psi_t := (\gamma_t, \delta_t)'$. Then:

$$\begin{cases} \Theta_{t+1} &= F_t \cdot \Theta_t \\ P(X_t^i | \Theta_t) &\sim \log\text{-Normal}(\mu_{t,i}, \phi_t \cdot \mu_{t,i}^\tau) \\ Y_t &= A \cdot X_t, \quad t \geq 1 \end{cases} \quad i = 1, \dots, \kappa \quad (3.11)$$

where τ is a fixed scalar, $P(\Theta_0^i) \sim \log\text{-Normal}_i$, and $P(F_t^{ii} | \Psi_t) \sim \log\text{-Normal}(\gamma_{t,i}, \delta_{t,i})$, $i = 1, \dots, \kappa$.

We estimated Ψ_t from the posterior distributions $P(\Theta_t | Y_t)$, that represented our best guesses about the evolution of Θ_t at each time t .

3.3.3 Particle Filter via SIR-Move Algorithm

In order to filter the posterior distributions of the origin-destination flows and estimate the parameters of the models above, we implemented the sample-resample-move algorithm of Gilks and Berzuini (2001), which we briefly outline below:

1. Initialization, $t = 0$.

For $i = 1, \dots, N$, sample $\tilde{x}_0^{(i)} \sim P_{X_0}$ and set $t=1$.

2. Importance sampling step

For $i = 1, \dots, N$, sample $\tilde{x}_0^{(i)} \sim P_{X_t | X_{t-1}^{(i)}}$ and set $\tilde{x}_{0:t}^{(i)} = (\tilde{x}_{0:t-1}^{(i)}, \tilde{x}_t^{(i)})$.

For $i = 1, \dots, N$, evaluate the importance weights $\tilde{\omega}_t^{(i)} = P_{Y_t | X_t = \tilde{x}_0^{(i)}}(y_t)$.

Normalize the importance weights.

3. Resampling step

Resample with replacement N particles $(x_{0:t}^{(i)} : i = 1, \dots, N)$ from the set $(\tilde{x}_{0:t}^{(i)} : i = 1, \dots, N)$ according to the importance weights

4. Move step

Move each of the N particles $(x_{0:t}^{(i)} : i = 1, \dots, N)$ a few steps according to the MCMC in section 3.2

5. Set $t \leftarrow t + 1$ and go to step 2.

The key point of this algorithm is to be able to sample from both P_{X_0} and $P_{X_{t+1} | X_t}$, and to be able to compute the importance weights.

3.4 Multivariate Integration

The local maximum likelihood approach of Cao et al. involves the estimation of Θ to obtain \hat{P} , and then the computation of the multivariate integral $E(X_t|Y_t, \hat{\Theta}, X_t > 0)$ to obtain point estimates for the OD flows \hat{X}_t . The multivariate integral is nasty and the solution (T1) proposed by Cao et al. is to estimate $\hat{X}_i(t)$ component-wise by means of a convenient closed form solution, and then to adjust these one-dimensional estimates using IPFP so that they satisfy the constraints imposed by the observed link loads Y_t . The approximate component-wise integrations ignore *joint* information and introduce an additional source of error in a non-controllable way.

We show in section 4 that the need for inference arises in situations where the traffic is intense, and it is rare to observe no traffic between two nodes over the time window used to estimate Θ_t . As a consequence an unconstrained multivariate integration would rarely give negative results, and if it did that may be an indication of zero traffic¹², especially if the negative flows in \hat{X}_t were limited to few components. Further, in the convenient Gaussian setting it is possible to find the exact value for the integral with an iterative procedure. We propose two different solutions to deal with the multivariate integration (T1):

- (1) Estimate X_t using the mean of the unconstrained integral $E(X_t|Y_t, \hat{\Theta})$ and adjust the resulting OD pairs when negative, as in $\max\{\hat{X}_t, 0\}$.
- (2) Estimate X_t using the solution to the constrained minimization problem:

$$\left| \begin{array}{l} \hat{X}_t = \arg \min_X \frac{1}{2}(X - \hat{\lambda}_t)' \hat{\Sigma}^{-1} (X - \hat{\lambda}_t) + \frac{1}{2}(Y_t - AX)' (A\hat{\Sigma}_t A')^{-1} (Y_t - AX) \\ \text{subject to } X \geq \ell. \end{array} \right. \quad (3.12)$$

provided by Lagrange method.

The heuristic in (1) is not unknown in statistics; it underlies proposed corrections to many estimators derived using method of moments, as well as the famous Stein estimator for the mean of a multivariate normal distribution. Moreover we show that such a correction is more sensible to the extent of preserving some geometric properties of the estimates that we obtain. In fact once we estimate $\hat{\Theta}$, the statistical relation among OD flows X and observations Y is known to be

$$(X_t, Y_t)' | \hat{\Theta} \sim N \left(\begin{bmatrix} \hat{\lambda} \\ A\hat{\lambda} \end{bmatrix}, \begin{bmatrix} \hat{\Sigma} & \hat{\Sigma}A' \\ A\hat{\Sigma} & A\hat{\Sigma}A' \end{bmatrix} \right), \quad (3.13)$$

and estimating X_t component-wise would introduce a non-controllable source of error.

The solution in (2) is the exact solution to the problem of finding $E(X_t|Y_t, \hat{\Theta}, X_t > 0)$, the posterior means we would get using Normal prior and Normal likelihood. The closed form solution is hard to obtain, but we give an iterative algorithm that converges to the exact solution, and may be preferred in those cases where the correction in (1) involves several coordinates and an alternative method may be preferred. Suppressing the time subscript of X_t we can write the Lagrangian for the problem as

$$\mathcal{L}(X, \nu) = \frac{1}{2} \left[(X - \hat{\lambda}_t)' \hat{\Sigma}^{-1} (X - \hat{\lambda}_t) + (Y_t - AX)' (A\hat{\Sigma}_t A')^{-1} (Y_t - AX) \right] + \sum_{i=1}^{\kappa} \nu_i \cdot (X_i - \ell_i)$$

¹²Notice that multiple solutions may exist.

and solving for X we get

$$X = \left(\hat{\Sigma}^{-1} + A'(A\hat{\Sigma}_t A')^{-1}A \right)^{-1} \left(\hat{\Sigma}^{-1}\hat{\lambda} + A'(A\hat{\Sigma}_t A')^{-1}Y + \nu \right). \quad (3.14)$$

A reasonable algorithm is to start with the unconstrained multipliers by setting $\nu_i = 0$ for $i = 1, \dots, \kappa$, and then iterate until convergence: compute X from 3.14 and set

$$\nu_i = \begin{cases} \max \left\{ 0, \nu_i - \frac{(X_i - \ell_i)}{\frac{\partial X_i}{\partial \ell_i}} \right\} & \text{if } X_i > \ell_i \\ \nu_i - \frac{(X_i - \ell_i)}{\frac{\partial X_i}{\partial \ell_i}} & \text{if } X_i < \ell_i \end{cases}$$

for $i = 1, \dots, \kappa$, to get (X, ν) at each iteration.

4 Experiments

Here we present the results we obtained on star network topologies, using real network traffic data of the Carnegie Mellon LAN, and of a small LAN at Lucent Technologies.

4.1 Exploring Carnegie Mellon Network

The exploratory data analysis revealed an uneven distribution of the traffic over the possible routes, and a continuous flow of traffic to and from the external world. The origin-destination flows in the validation dataset we were able to retrieve were definitely not Gaussian, but rather skewed. Sub-problems where the OD traffic was sparse could be reduced and solved exactly when small, and solved by entropy related heuristics when large; to this extent we gave and employed two algorithms, `match` and `split-match`, whose underlying assumptions were validated by the data. The idea underlying the algorithms is that traffic one-to-many¹³ is more likely to happen when traffic volume is intense. As the traffic volume got intense we used the models we proposed in section 3 to recover the origin-destination flows.

Carnegie Mellon Data Collection

The data was gathered by Dr. Frank Kietzke and Dr. Russel Yount at the Carnegie Mellon Network Group, using the machines of Prof. Srinivasan Seshan at Information Networking Institute.

Carnegie Mellon network is very complex and even as our large CISCO routers implement all the sets of SNMP protocols we could not obtain validation data directly via SNMP queries to the routers; in this pilot study we obtained validation data via a network management application. Four datasets were collected, and we helped discover critical bugs in the network management program along the way, as some of the data we collected did not make sense. We experimented first hand how critical the collection of origin-destination flow is! A consequence of collecting data through the network management application, was that whenever the traffic became too intense, the program would trash OD flows in the queue before we could actually measure them, and this fact resulted in slightly lower traffic flows. There was enough traffic in the data, though, to actually see the expected daily patterns, though, and the overall volume was reasonable, so that we considered the validation data we collected as being the real origin-destination flows.

For more information, or to glance at the on-line link loads on the various routers, please visit the Carnegie Mellon University Network Group home page at <http://stats.net.cmu.edu/>.

¹³One source to many destinations in a 5-minute time interval.

Carnegie Mellon Data Overview

The table below summarizes the observed volume over a period of slightly less than two days.

Over about 40hr (480 measurements):

| | | | |
|------------------|------------------|--------|--------|
| . to ext world | 3513.963.250.400 | 3514GB | 49.55% |
| . from ext world | 2053.731.241.843 | 2054GB | 28.96% |
| . CMU (tcp/ip) | 1523.047.061.952 | 1523GB | 21.48% |
| . CMU (other) | 312.847.796 | 313MB | 0.01% |

We fitted a simple linear model to get an idea of how much traffic CMU generated on average over the collection period; it was about 14GB every 5 minutes.

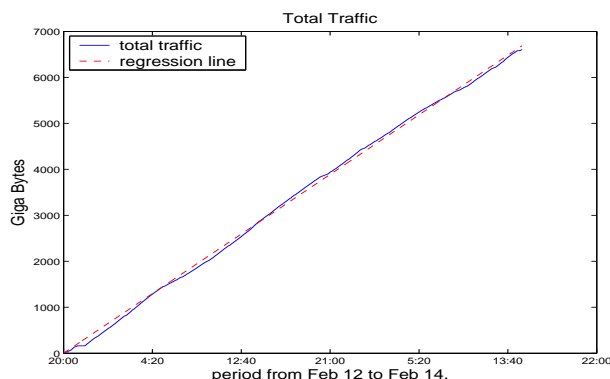


Figure 7: Total traffic in GB.

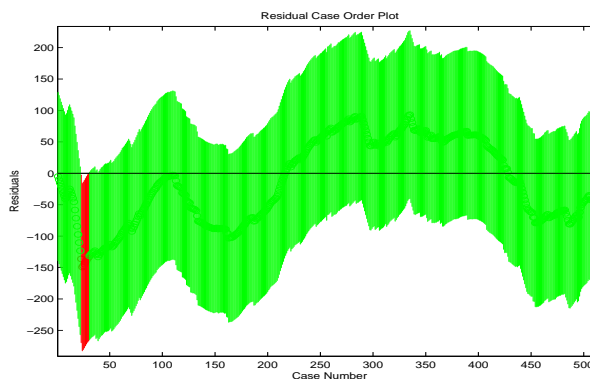


Figure 8: Case vs. residuals.

$$b = 13.9194 \quad b.\text{Int} (95\%) = [13.8981, 13.9406] \quad R^2 = 0.9720$$

There was a pattern in the residuals. It was not very clear in the differenced traffic, but it popped up looking at the percentages of traffic by destination. Carnegie Mellon network can be thought as a collection of local routers (mainly corresponding to physical buildings), all of them connected to two main router-switches (the Cores). The two Cores serve two different purposes. Core no.1 is the main switch, whereas Core no.2 is mainly a backup system, which in turn filters most of the traffic to and from the external world.

Total traffic:

| | CORE.1 | CORE.2 |
|-----------------------|--------------|--------------|
| . CMU internal | 1099 GB/5min | 423 GB/5min |
| . to external world | 1476 GB/5min | 2038 GB/5min |
| . from external world | 935 GB/5min | 1120 GB/5min |

A seasonal component would have been desirable in our models, but we did not have enough data to carry out a meaningful analysis; estimating the seasonal component using one day data was not a sensible thing to do. We notice, however, that our estimates \hat{X}_t provide a good starting point to fit, for example, a non-parametric seasonal component in the non-observable space where the OD flows live, should more data be available in the future.

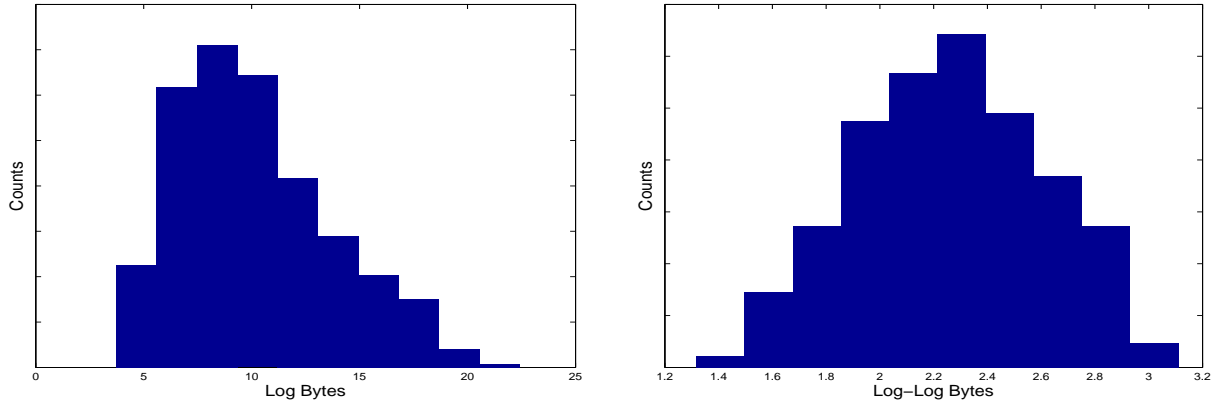


Figure 9: We considered all the 12321 OD flows at all times: we present histograms of the origin-destination flows transformed by $\log(X_t)$ in the left panel, and by $\log \log(X_t)$ in the right panel, omitting the zero counts.

4.1.1 Empirical Distributions of the OD Flows

We began our exploration by verifying on our data two empirical laws that are frequently associated with network traffic time series, namely: (1) the log-Normal distribution of the origin-destination flows, and (2) the 80%-20% distribution of the traffic over the available routes.

As far as the log-Normality of the OD flows is concerned, the plots in figure 9 above suggest that the data were far from Gaussian. Even a *log* transformation was not strong enough to obtain normality, a *log log* transformation worked better. About the distribution of the traffic over the possible origin-destination pairs, few routes definitely accounted for most of the traffic, as the plots in figure 10 below show, no matter what time of the day was.

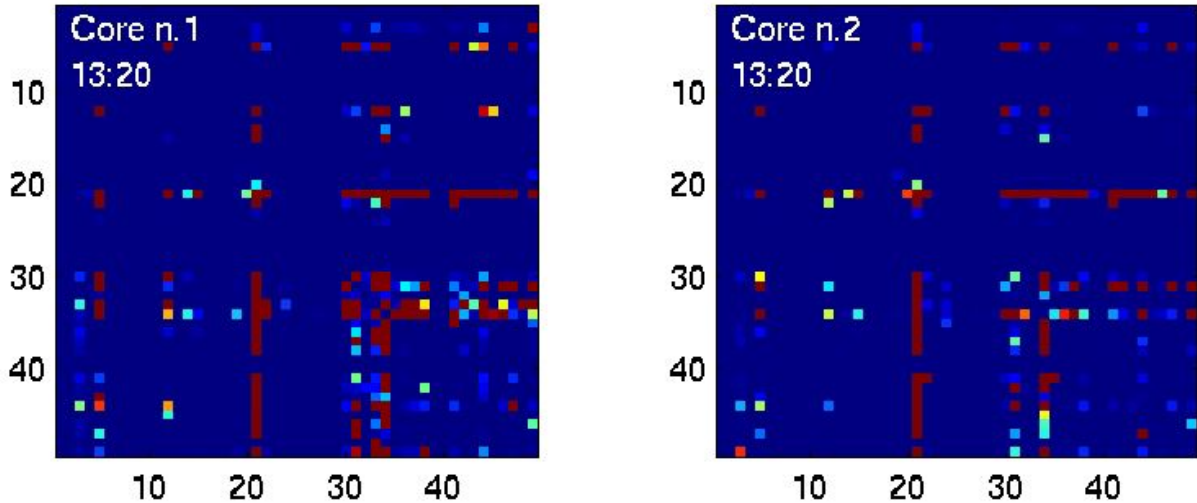


Figure 10: We present origin-destination traffic matrices corresponding to the two Cores at 13:20pm on day one of the collection period (the key to relate numbers to departments is presented in appendix B). The chromatic scale shows low traffic $< 1MB$ in blue, intense traffic $> 100MB$ in red, and various levels of medium traffic $\in [1MB, 100MB]$ in light blue, yellow and orange — a Black-White figure would show low traffic in black, medium traffic in white, and intense traffic in gray.

4.1.2 Modeling the Coefficients of Constant Association α_{ij}

We computed the coefficients of constant association α_{ij} corresponding to the sequence of OD traffic matrices looking for some patterns. The aim was to see whether it was possible to model these coefficients directly. The coefficients were as irregular as the data itself, in all the possible parameterizations, though. To model them would have required the same effort as to model the data, hence in this report we chose to model the data directly.

4.2 Exact Recovery Algorithms for Sparse Traffic Situations

Before studying inference methods for the OD flows, we devise two algorithms that are able to recover exactly the flows in sparse traffic situations. The algorithms `match` and `split-match` are based on the simple assumption that whenever some *in* and *out* link loads match exactly (say x bytes out of node B and x bytes into node A), and the traffic is low, we can safely attribute the traffic observed on the matching pair of links (B, A) to the corresponding origin-destination flow (traffic from B to A).

The 5-minute data are sparse. This fact alone allowed us to recover exactly the origin-destination flows for 91.97% of the time points, in example sub-networks like {CFA, Wean, Baker-Porter}. The algorithm `match` allowed us to recover exactly up to 98.08% of the data. The algorithm `split-match` augments `match` to deal with situations where only one link load is zero, and allowed us to recover exactly up to 98.51% of the data. An extension of the heuristic underlying these algorithms to cases where link loads do not match exactly entails arguments based on entropy.

4.3 Intense Traffic Sub-Networks at Carnegie Mellon

In this section we use the methods we presented in section 3: in section 4.3.2 we explore different dynamical behaviors; in section 4.3.3 we fully develop an example star network, and compare our methods, in terms of estimation error, to the state-of-the-art methods available today.

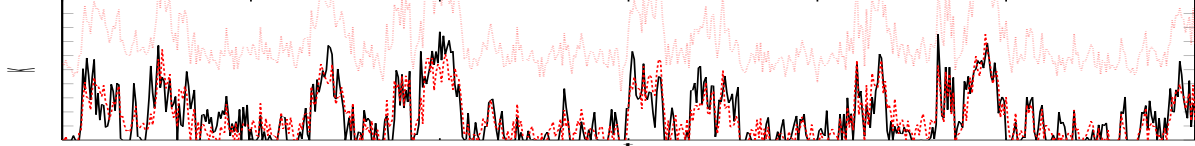
4.3.1 Naive SVD Solution for Strongly Correlated Flows

The nature of the problem suggested generalized inverses. We computed the solution corresponding to a Moore-Penrose generalized inverse, based on SVD decomposition. The fact that we had more non-observable OD flows than available measurements at each time point, made the results reliable only in presence of strong cross-correlations among the OD flows. The more the origin-destination flows got uncorrelated, the worse got the inferences based on the generalized inverse. Since we cannot observe the origin-destination flows, decisions on whether there is cross-correlation among them or not should be based on the observable link loads: this seemed impractical and we took a different direction.

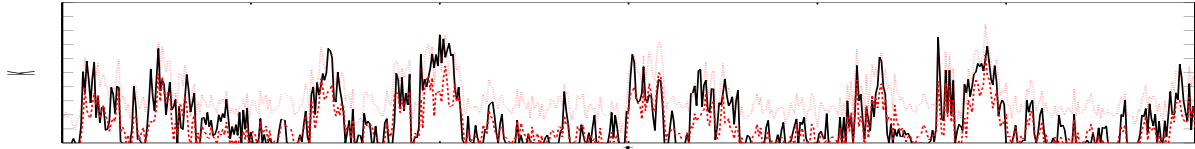
It is worth mentioning that classical time series analysis would find the solution to the under-determinacy of the problem in the presence of cross-correlations, and/or in partial knowledge about the dynamics underlying the traffic. Given the heterogeneity of the communication networks, assuming a specific form for these would harm the generality of our approach. This is the reason why we looked for a solution elsewhere: namely in the multiple use of data, and informative priors.

4.3.2 Local Dynamical Behavior

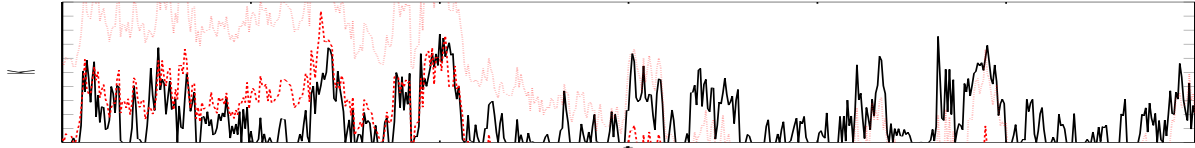
Here we provide some evidence to suggest that using a diagonal matrix F_t in the model defined by equations 3.1 is a reasonable choice. Considerations about the variability and bias of the inferences guided our choice. We played with the matrix F_t and tried several structures for it. We present below a plot for each F_t structure; the black (solid) lines represent the true origin-destination flows, and we superimpose the sequence of posterior modes obtained with the Kalman recursions (red, dashed lines), along with the upper bands at $+3$ posterior standard deviations (red, dotted lines). First we considered independent OD flows, by setting $F_t = 0$.



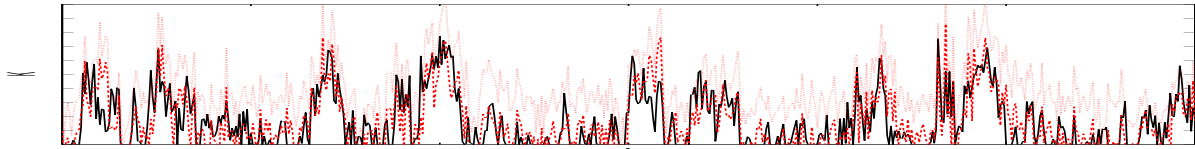
Forcing the absence of dynamics had an explosive effect on the variance of the posterior distributions, especially if compared to a diagonal $F_t = F$, independent AR(1) processes, below.



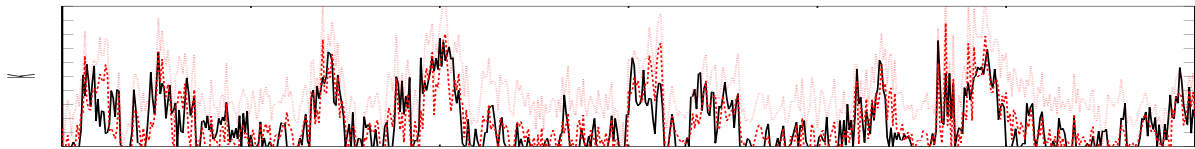
Using a fully unconstrained $F_t = F$, fixed in time, gave bad estimates with a strong drift.



Finally we had to decide between a diagonal time-varying F_t , which yielded



and a full time-varying F_t , which yielded



As we added more parameters and let F_t change fully unconstrained over time, the filtered OD flows got better, but the variability of the entries of F_t exploded. The plots below show the distributions of the coefficients of variation $\frac{\sigma_t}{\mu_t}$ of elements of F_t , for the two cases of unconstrained and diagonal F_t using a windows of 25 time points, that give a sense of the entailed variabilities.

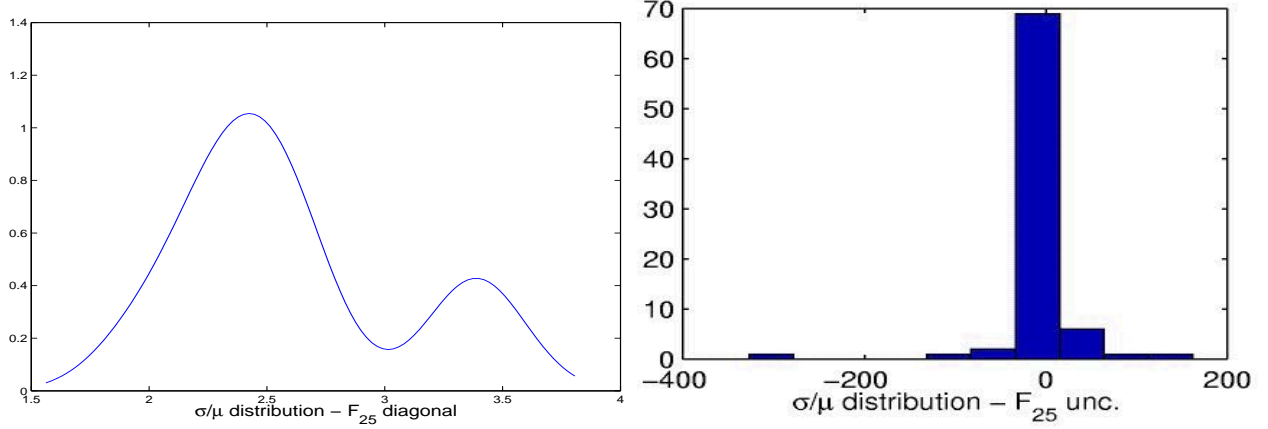
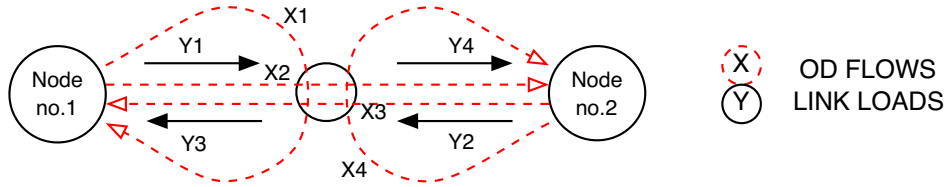


Figure 11: As we consider a fully unconstrained matrix F_t an index of the pure variability of its entries ($\frac{\sigma}{\mu}$) explodes. Left: F_t diagonal. Right: F_t unconstrained.

We concluded that a diagonal dynamic matrix F_t , that can vary quickly over time, was a reasonably good choice.

4.3.3 A Case Study: the Star Network Topology

In order to test the performance of our methods, we sampled the non-observable origin-destination flows from our validation dataset, and computed the corresponding link loads for a simple star network topology composed of one router and two nodes connected to it. Then we attempted to



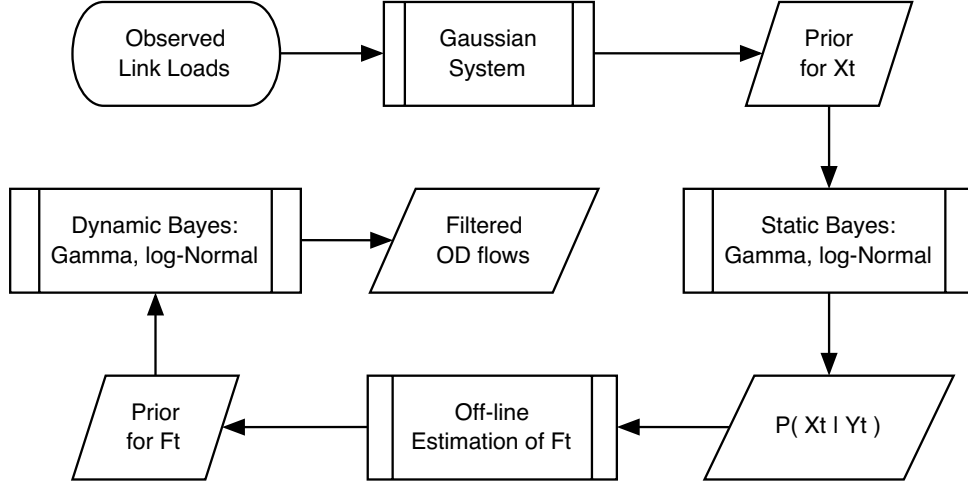
reconstruct the origin-destination flows, starting from the available measurements, the link loads. Throughout this section we denote the vector of origin-destination flows at time t by X_t , and the vector of observed link loads at time t by Y_t .

Overview

The primary object of interest in the experiments in this section was the sequence of filtering distributions $P(X_t, \Theta_t | Y_t, \dots, Y_1)$. In particular we were interested in the mean of the marginal distribution $\hat{X}_t = E(X_t | Y_t, \dots, Y_1)$ which we would use as point estimate for the OD flows at time t .

Given the heterogeneity of the communication networks out there, and without universal recipes for the cross-correlation structure and the dynamics, we took leave from classical time series modeling “a la Box and Jenkins” that would rely on these objects, and on the decomposition of the OD flows using trend, seasonal component, and ARMA process, to overcome the under-determinacy of the problem, and moved towards approaches that make use of data augmentation, like simulated annealing. We propose a *new way* to augment the data and generate extra information: we used different models, sequentially, on copies on the data, whereas simulated annealing would use the

same model on copies of the data, simultaneously. In other words instead of maximizing a power



of the likelihood of the data $P^r(Y_t|X_t, \Theta_t)$, we translated preliminary estimates into informative priors and then used Bayes theorem to compute the posteriors and obtain new, refined estimates; the use of different models allowed us to take advantage of a larger set of useful properties, which could not be packed into one single model.

In order to compare inferences obtained with different methods we computed the estimation errors; the ℓ_2 distance between the true OD flows in the validation set, and the estimates. We first compared our linear Gaussian dynamical system with the local likelihood approach in Cao et al. The Gaussian system in section 3.1.1 yielded better estimates, as we expected, since we

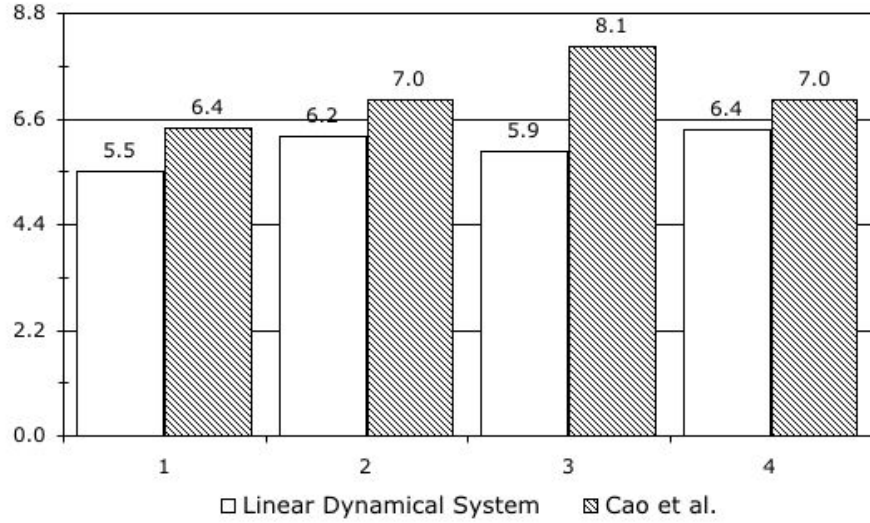


Figure 12: The bars represent the estimation errors (ℓ_2 distance between the true OD flows in a validation set and the estimates) obtained with our linear Gaussian dynamical system, and with the method proposed by Cao et al. In the example star topology we considered, there were 4 OD flows that needed be estimated.

showed the model by Cao et al. is a special case of it. At this stage we tried also to learn the structure of F_t by means of the EM algorithm in Ghahramani and Hinton (1996) and explored some

more classical time series modeling without exciting results, mainly because of the few observations about the link loads available. Before passing to the next method is worth noticing that the contribution of the Gaussian system, fitted using a two-stage maximization of the likelihood, is that its parameterization entails a one-to-one correspondence between the parameters that govern means and variances of the non-observable OD flows, and the parameters that govern means and variances of the observable link loads. This feature allowed us to roughly *identify* a most likely area where to look for the correct solution, among the many feasible ones. Asymptotic properties of the estimates based on local likelihood methods were discussed in Loader (1999).

Next we compared the static Bayesian method in section 3.2 with our Gaussian system, and the method by Cao et al. We used informative priors as soft thresholds in order to make *more*

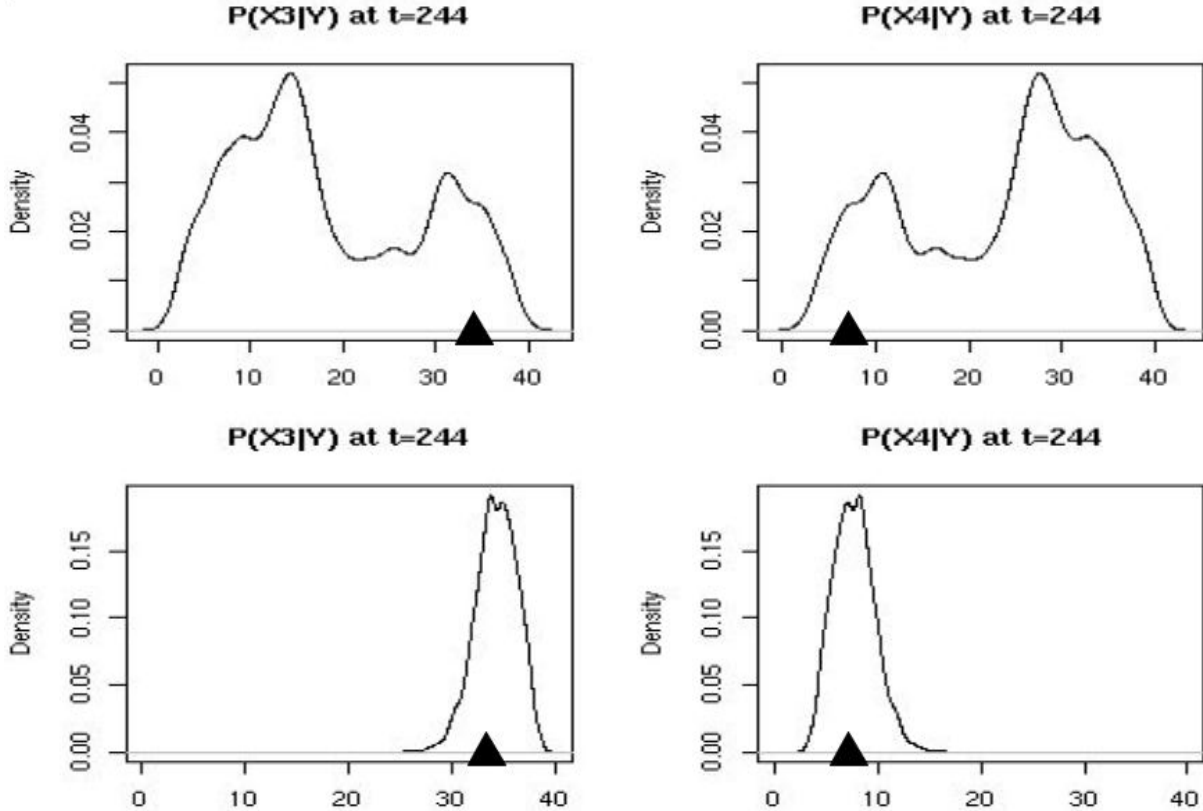


Figure 13: Example posterior distributions for the OD flows: X_3 and X_4 at time $t = 244$. The traffic on the X axes is measured in Kbytes, and the figures show the posterior distributions we obtained with non-informative priors (top panel) and with informative priors (bottom panel) based on our Gaussian system. The triangles represent the true OD Flows, whereas our point estimates for the OD flows would correspond to the means of the posterior distributions. Making the posterior distributions *more unimodal* improved the inferences, reducing the bias entailed by extra modes.

unimodal the posterior distribution of the OD flows; these distributions would represent our best probabilistic guesses about the amount of traffic on each OD route at time t , given the observed link loads at time t only. The informative priors we used had a huge variance, but were centered on the point estimates obtained with the Gaussian system. Modeling the OD flows with realistic, non-Gaussian distributions improved a lot our point estimates, means of the posterior distributions above. Further informative priors helped reduce the bias entailed by multiple modes.

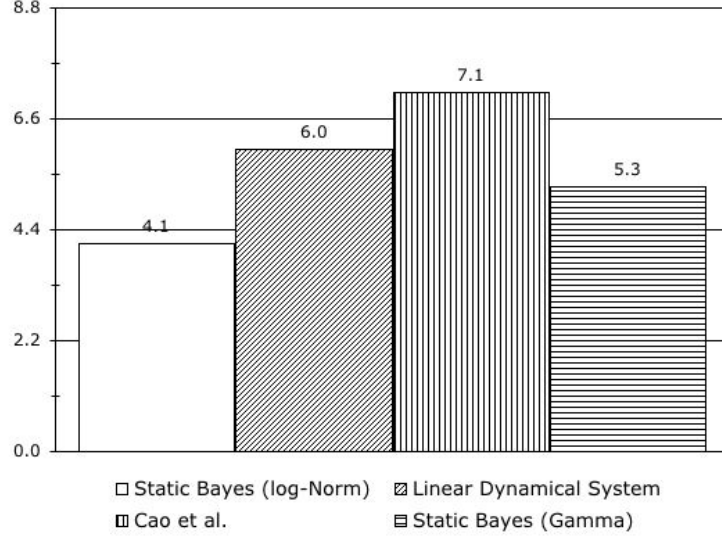


Figure 14: The bars represent the estimation errors (average ℓ_2 distance between the true OD flows in a validation set and the estimates) obtained with different methods: our static Bayesian method based on Gamma models and log-Normal models for the OD flows, our linear Gaussian dynamical system, and the method proposed by Cao et al.

Finally we computed the filtering distributions of the OD flows at each time t ; our best probabilistic guesses about the amount of traffic on the OD routes at time t , given the observed link loads from time 1 to time t . In order to do so we used our best Bayesian dynamical system with stochastic

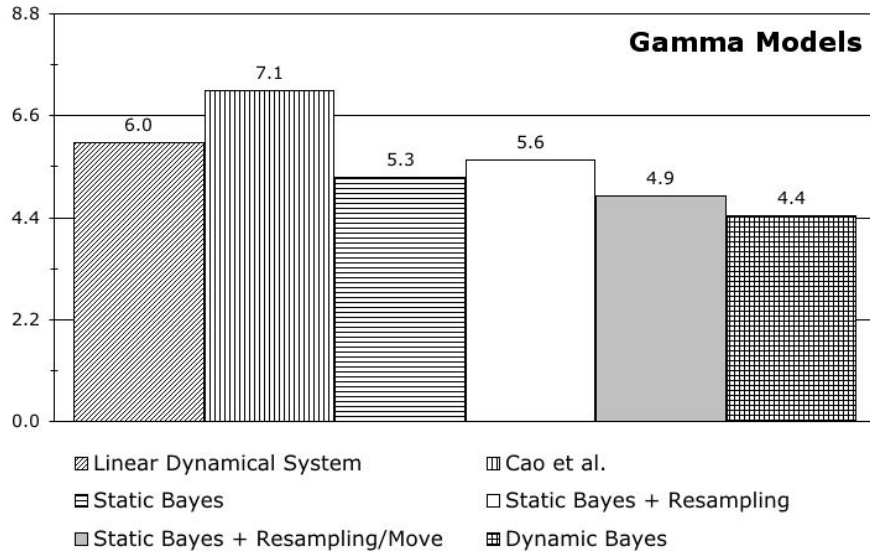


Figure 15: The bars represent the estimation errors (average ℓ_2 distance between the true OD flows in a validation set and the estimates) obtained with various methods. Our dynamic Bayesian method is a clear winner. We also included the estimation errors we obtained applying two methods not discussed in this report, which constitute valid competitors in our framework: an importance sampling scheme, and an importance sampling - move scheme in the spirit of Gilks and Berzuini (2001).

dynamics, and we introduced informative priors on the parameters governing such dynamical behavior, again, in order to identify a more likely area where to look for the correct solution, among the many feasible ones. The resulting inferences are better than the previous ones, both using

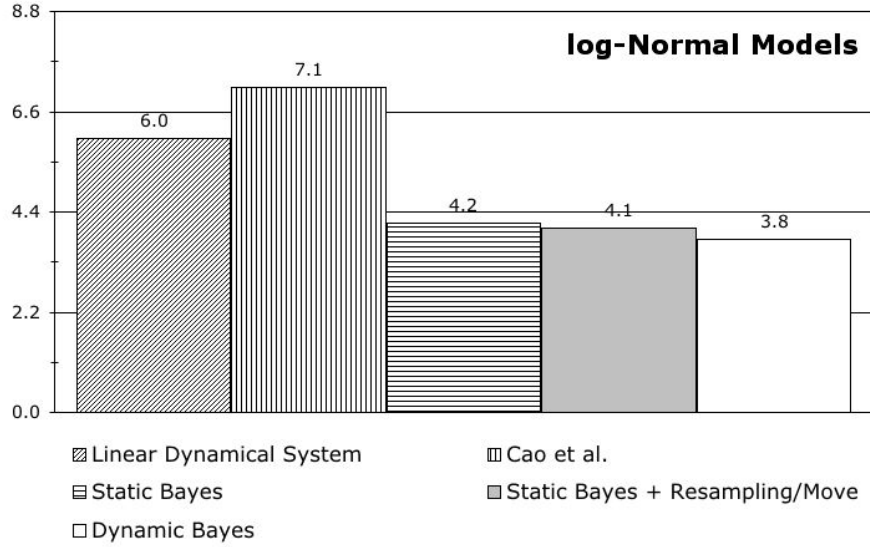


Figure 16: The bars represent the estimation errors (average ℓ_2 distance between the true OD flows in a validation set and the estimates) obtained with various methods. Our dynamic Bayesian method is a clear winner.

Gamma models and using log-Normal models for the OD flows. The informative priors on the

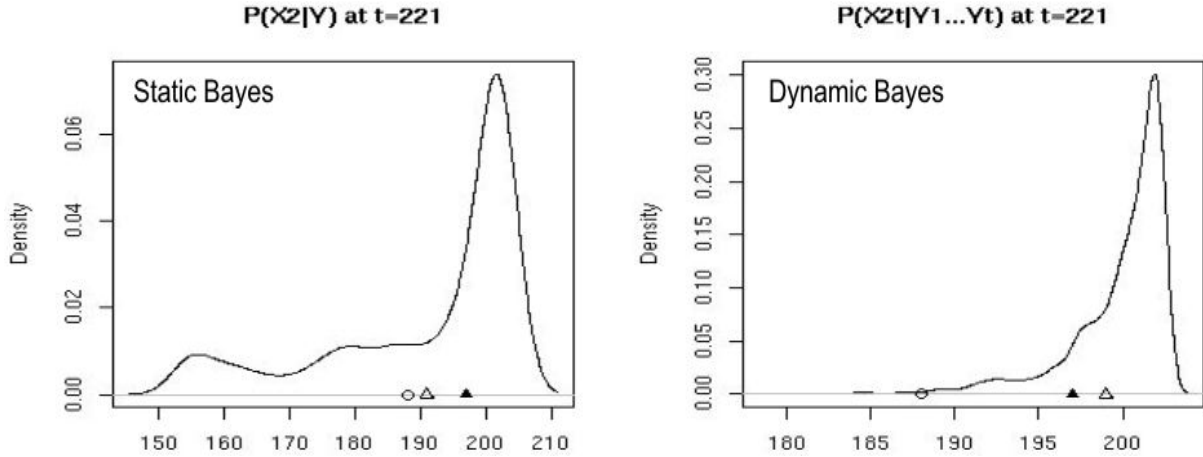


Figure 17: Example filtering distributions for the OD flows: X_2 at time $t = 221$. The traffic on the X axes is measured in Kbytes. The figures show the posterior distributions we obtained with informative priors in a static Bayesian setting (left panel), and with informative priors in a dynamic Bayesian setting (right panel). The triangles represent the true OD Flows, whereas our point estimates for the OD flows would correspond to the means of the posterior distributions. Using the information in the measurements Y_1, \dots, Y_T made the posterior distributions *less variable* and improved our inferences — notice the ranges.

parameters underlying the dynamical behavior of the system at time t were obtained using the pair of marginal posterior distributions $P(\Theta_t|Y_t)$ and $P(\Theta_{t+1}|Y_{t+1})$, that entailed our best guesses on

the distributions of the parameters from time t to time $t+1$, in the static Bayesian case. The reason for this improvement, which goes beyond the contribution of state-of-the-art resampling schemes, is that in the dynamic setting the point estimates we used for the OD flows were the means of the filtering distributions $P(X_t|Y_t, \dots, Y_1)$, that depended on all the observed link loads from time 1 to time t , whereas in a static setting the point estimates we used for the OD flows were the means of the posterior distributions $P(X_t|Y_t)$, that depended on the observed link loads at time t only: using more information reduced the variability of our estimates.

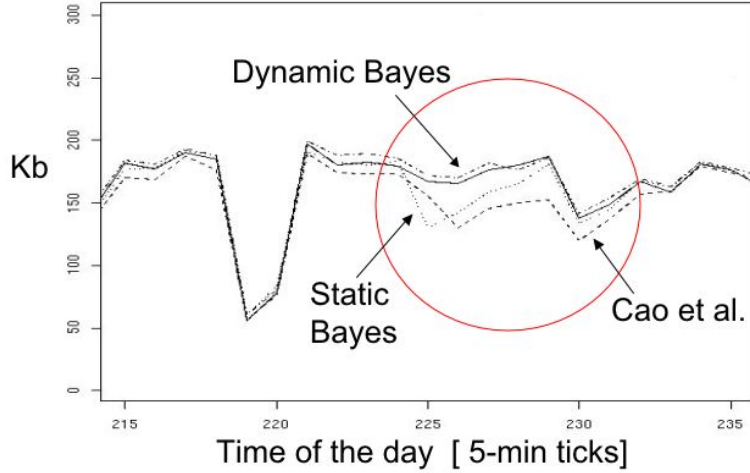


Figure 18: Example fit: the performance of the Bayesian dynamical system based on a log-Normal model for traffic flows, stochastic dynamics, and informative priors is quite amazing. The goodness of the estimates can be appreciated after sharp changes in the non-observable traffic flows. It definitely improves on the static Bayesian solution, and on the solution by Cao et al.

5 Conclusions

The problem of estimating the non-observable origin-destination flows in a network, starting from observations on the link loads, boils down to estimating κ numbers from ℓ numbers at each epoch, where $\kappa = O(\ell^2)$. It is clear that some extra information is needed. We found this extra information in the multiple use of the data¹⁴, and proposed a new methodology that was able to reduce by more than 45% the estimation error achieved by state-of-the-art solutions, in a realistic setting.

Introducing statistical models for the non-observable OD flows induced a probabilistic mapping on the observable link loads, the likelihood. The under-determinacy of the problem was still present, though, and showed up in the form of multiple modes in the likelihood. In such a situation maximizing the likelihood may or may not yield better inferences; good inferences would require a realistic model able to capture relevant features of the data. To this extent we introduced in our models: **1.** skewed distributions for the OD flows (Gamma and log-Normal), and **2.** explicit stochastic time dependence for the traffic.

Our best model is the Bayesian dynamical system in section 3.3.2. We introduced explicit time dependence by means of independent AR(1) processes, with stochastic (log-Normal and Inverse

¹⁴Classical time series analysis would look for this information in cross-correlation patterns and in physical laws for the dynamics of the OD flows. However, given the heterogeneity of the communication networks out there, any method that based its strength on such objects would lack the general applicability that we are able to claim.

Gamma) coefficients, for the means of OD flows, in order to allow for more flexibility on the dynamics of the OD flows themselves. Further the AR(1) processes were *local*, in the sense that the constants underlying the distribution of their coefficients were allowed to vary over time. It is here that we introduced the extra information, in the form of informative priors for these underlying constants.

The new methodology we proposed is based on multiple use of the data, and directly addressed the problem of multiple modes in the likelihood in two steps: **1.** in a first stage we identified the area where the correct solution was likely to be, by means of a model that entailed a one-

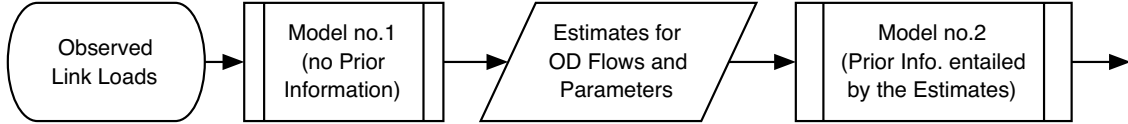


Figure 19: We proposed multiple use of the data to overcome the under-determinacy of the problem. Instead of relying on the same model for copies of the dataset, as in simulated annealing, we used different models on copies of the data, and Bayes theorem to properly update our beliefs at each stage.

to-one mapping between parameters underlying the OD flows and the link loads, **2.** in a second stage we translated the information entailed by the first-stage estimates into informative priors for the constants underlying the dynamics of the system, at each epoch. Our methodology fits in with recently proposed non-parametric Empirical Bayes approaches, that suggest how to fix the values of the constants underlying prior distributions in a Bayesian framework; it is also related to simulated annealing, in that we used different models on copies of the data, in a sequential fashion, instead of using the same model on copies of the data simultaneously.

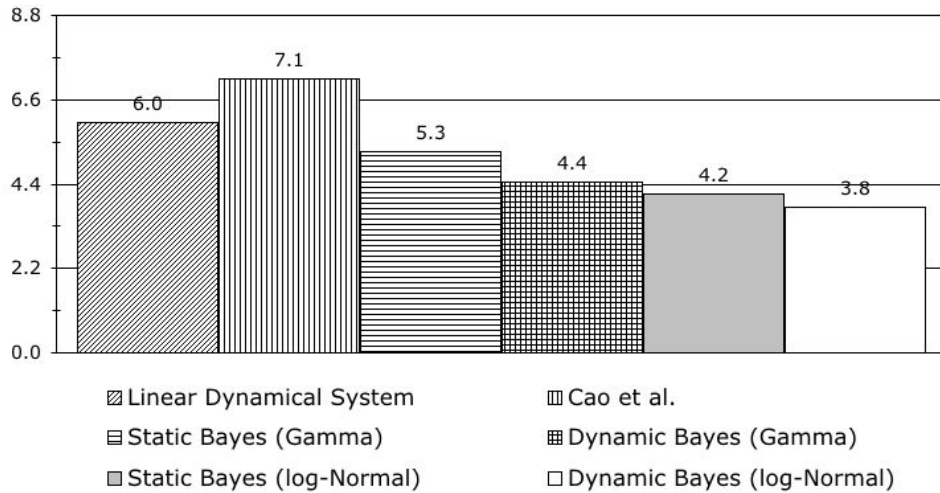


Figure 20: The bars represent the estimation errors (average ℓ_2 distance between the true OD flows in a validation set and the estimates) obtained with the methods we proposed; comparing the estimates obtained with our models to that obtained with that of Cao et al., the current state-of-the-art, in a realistic setting, we found that the static Bayesian models based on Gamma and log-Normal models for OD traffic flows reduced the errors by 25% and 38%, respectively, and that the introduction of a stochastic local dynamical behavior along with informative priors on it, reduced the errors by 41% and 46%, respectively.

A further advantage of our approach is that it produced estimates that were less variable than

those obtained with state-of-the-art solutions. We solved the problem in a dynamic setting, where the primary object of interest was the sequence of filtering distributions $P(X_t, \Theta_t | Y_t, \dots, Y_1)$. In particular the point estimates we used for the OD flows were the means of the marginal distributions $\hat{X}_t = E(X_t | Y_t, \dots, Y_1)$, which depended on the whole sequence of observed link loads $\{Y_t, \dots, Y_1\}$, and hence were less variable than the solutions given in a static setting present in the literature, based on a small set of observations around time t .

In conclusion, estimating origin-destination flows from link loads in a dynamic setting, using skewed models for the traffic, is the best available option at this time. We encourage the use of informative priors based on the data, in a novel spirit, to overcome the under-determinacy of the problem. Our Bayesian dynamical system based on a log-Normal model and informative priors reduced by more than 45% the estimation error achieved by state-of-the-art solutions in a realistic setting.

A Gaussian Dynamical System

A.1 EM algorithm

In order to estimate the parameters of the Gaussian system in the first-stage of our maximization procedure, we used both brute-force maximization of the log-likelihood as well as the EM algorithm. Here we detail the algebra of the EM algorithm, in the next subsection we show how to use EM along with one-step Newton-Raphson to speed up the computation.

The EM algorithm goes as follows: [idea]

1. start anywhere in the parametric space, say at $\Theta_0 \in \Omega$
2. write down $l(\text{complete data} | \Theta_0) = l((\mathbf{X}_1, \dots, \mathbf{X}_T) | \Theta_0)$
3. compute $Q(\Theta, \Theta_0) = E(l((\mathbf{X}_1, \dots, \mathbf{X}_T) | \Theta) | (\mathbf{Y}_1, \dots, \mathbf{Y}_T), \Theta_0)$ (E-step)
4. $\Theta_1 = \arg \max_{\Theta \in \Omega} Q(\Theta, \Theta_0)$ (M-step)
5. loop.

Now for step 2 write

$$l((\mathbf{X}_1, \dots, \mathbf{X}_T) | \Theta = \theta) = -\frac{T}{2} \log(2\pi) - \frac{T}{2} \log |\Sigma| - \frac{1}{2} \sum_{t=1}^T (X_t - \lambda)' \Sigma^{-1} (X_t - \lambda).$$

For step 3 $\mathbf{X}_t | \mathbf{Y}_t, \Theta_0 \sim N(m^{(0)}, S^{(0)})$, yields

$$\begin{aligned} Q(\Theta, \Theta_0) &= -\frac{T}{2} \log(2\pi) - \frac{T}{2} \log |\Sigma| - \frac{1}{2} \sum_{t=1}^T E[(X_t - \lambda)' \Sigma^{-1} (X_t - \lambda) | \mathbf{Y}_t, \Theta_0] \\ &= -\frac{T}{2} (\log(2\pi) + \log |\Sigma| + \text{tr}(\Sigma^{-1} S^{(0)})) - \frac{1}{2} \sum_{t=1}^T (m_t^{(0)} - \lambda)' \Sigma^{-1} (m_t^{(0)} - \lambda) \end{aligned}$$

where

$$\begin{aligned} m_t^{(0)} &= \lambda_0 + \Sigma_0 A' (A \Sigma_0 A')^{-1} (Y_t - A \lambda_0), \quad \text{and} \\ S^{(0)} &= \Sigma_0 - \Sigma_0 A' (A \Sigma_0 A')^{-1} A \Sigma_0' \end{aligned}$$

since

$$(\mathbf{X}_t, \mathbf{Y}_t)' | \Theta_0 \sim N \left(\begin{bmatrix} \lambda_0 \\ A \lambda_0 \end{bmatrix}, \begin{bmatrix} \Sigma_0 & \Sigma_0 A' \\ A \Sigma_0 & A \Sigma_0 A' \end{bmatrix} \right).$$

For step 4 we want to find

$$\Theta_1 = \arg \max_{\Theta \in \Omega} Q(\Theta, \Theta_0).$$

Since

$$\begin{aligned} \Theta &= (\phi, \lambda_1, \dots, \lambda_I) \\ \Sigma &= \phi \cdot \text{diag}(\lambda_1, \dots, \lambda_I) \\ \Sigma^{-1} &= \frac{1}{\phi} \cdot \text{diag}(\frac{1}{\lambda_1}, \dots, \frac{1}{\lambda_I}) \\ |\Sigma| &= \phi^I \cdot \lambda_1 \lambda_2 \dots \lambda_I \\ \text{tr}(\Sigma^{-1} S^{(0)}) &= \frac{1}{\phi} \cdot \left(\frac{S_{1,1}^{(0)}}{\lambda_1} + \dots + \frac{S_{I,I}^{(0)}}{\lambda_I} \right) \\ \sum_{t=1}^T (m_t^{(0)} - \lambda)' \Sigma^{-1} (m_t^{(0)} - \lambda) &= \frac{1}{\phi} \sum_{t=1}^T \sum_{i=1}^I \left(\frac{(m_{t,i}^{(0)})^2}{\lambda_i} + \lambda_i - 2m_{t,i}^{(0)} \right) \end{aligned}$$

we can rewrite Q as

$$Q = -\frac{TI}{2} \log \phi - \frac{T}{2} \sum_{i=1}^I \log \lambda_i - \frac{T}{2\phi} \left(\frac{S_{1,1}^{(0)}}{\lambda_1} + \dots + \frac{S_{I,I}^{(0)}}{\lambda_I} \right) - \frac{1}{2\phi} \sum_{t=1}^T \sum_{i=1}^I \left(\frac{(m_{t,i}^{(0)})^2}{\lambda_i} + \lambda_i - 2m_{t,i}^{(0)} \right).$$

Next set the gradient $\nabla Q = \partial Q / \partial \Theta = 0$ to get

$$\begin{cases} \frac{\partial Q}{\partial \lambda_i} : & \phi \lambda_i - S_{i,i}^{(0)} - \frac{1}{T} \sum_{t=1}^T (m_{t,i}^{(0)})^2 + \lambda_i^2 = 0 & i = 1, \dots, I \\ \frac{\partial Q}{\partial \phi} : & \sum_{i=1}^I \left(\lambda_i - \frac{1}{T} \sum_{t=1}^T m_{t,i}^{(0)} \right) = 0. \end{cases} \quad (\text{A.1})$$

We can find an analytic solution for $\lambda(\phi)$, and then solve for ϕ . We keep ϕ and the λ s positive on their path to the solution, using fractional steps as necessary.

A.2 More computational efficiency

Let $f(\Theta) = (f_1(\Theta), \dots, f_I(\Theta), f_{I+1}(\Theta))'$ be the left hand side of (A.1) above; we used a one-step Newton-Raphson algorithm to speed up computations needed to update $\Theta^{(it)}$, whose convergence properties to the same local maximum are studied in Lange (1995), according to

$$\Theta^{(it+1)} = \Theta^{(it)} - \left[\dot{F}(\Theta^{(it)}) \right]^{-1} \cdot f(\Theta^{(it)}),$$

where \dot{F} is the Jacobian of $f(\Theta)$ with respect to Θ . We get

$$\begin{cases} \frac{\partial f_i}{\partial \lambda_j} : & = \phi \cdot I_{\{i\}}(j) + 2 \cdot \lambda_j \\ \frac{\partial f_{I+1}}{\partial \lambda_j} : & = 1 \\ \frac{\partial f_i}{\partial \phi} : & = \lambda_i \\ \frac{\partial f_{I+1}}{\partial \phi} : & = 0 \end{cases} \quad (\text{A.2})$$

we either computed the inverse analytically or numerically (penalizing the eigenvalues in the latter case to avoid singular matrices to numerical precision), and we used fractional steps as needed to keep parameter values positive on their path to the solution.

A.3 KF posteriors

A.3.1 One Y to one X

Kalman recursions in the scalar-scalar case. Start from

$$X_1^0 = \mu, \quad V_1^0 = V_1$$

compute the gain and the posteriors

$$\begin{aligned} K_1 &= \frac{V_1}{R+V_1} \\ X_1^1 &= \left(\frac{R}{R+V_1} \mu + \frac{V_1}{R+V_1} Y_1 \right) \\ V_1^1 &= \frac{RV_1}{R+V_1} \end{aligned}$$

and then again, project and correct

$$\begin{aligned}
X_2^1 &= F \left(\frac{R}{R+V_1} \mu + \frac{V_1}{R+V_1} Y_1 \right) \\
V_2^1 &= F^2 \left(\frac{R V_1}{R+V_1} \right) + Q \\
K_2 &= \frac{\left[F^2 \left(\frac{R V_1}{R+V_1} \right) + Q \right]}{R + \left[F^2 \left(\frac{R V_1}{R+V_1} \right) + Q \right]} \\
X_2^2 &= \frac{R}{R + \left[F^2 \left(\frac{R V_1}{R+V_1} \right) + Q \right]} \left[F \left(\frac{R}{R+V_1} \mu + \frac{V_1}{R+V_1} Y_1 \right) \right] + \frac{\left[F^2 \left(\frac{R V_1}{R+V_1} \right) + Q \right]}{R + \left[F^2 \left(\frac{R V_1}{R+V_1} \right) + Q \right]} Y_2 \\
V_2^2 &= \frac{R \left[F^2 \left(\frac{R V_1}{R+V_1} \right) + Q \right]}{R + \left[F^2 \left(\frac{R V_1}{R+V_1} \right) + Q \right]}
\end{aligned}$$

to obtain the general recursive formulas

$$\begin{aligned}
V_t^{t-1} &= F^2 (V_{t-1}^{t-1}) + Q \\
X_t^{t-1} &= F X_{t-1}^{t-1} \\
K_t &= \frac{[V_t^{t-1}]}{R + [V_t^{t-1}]} \\
X_t^t &= \frac{R}{R + [F^2 (V_{t-1}^{t-1}) + Q]} [F X_{t-1}^{t-1}] + \frac{[F^2 (V_{t-1}^{t-1}) + Q]}{R + [F^2 (V_{t-1}^{t-1}) + Q]} Y_t \\
&= (1 - K_t) [F X_{t-1}^{t-1}] + K_t Y_t \\
&= K_t Y_t + (1 - K_t) K_{t-1} F Y_{t-1} + (1 - K_t)(1 - K_{t-1}) K_{t-2} F^2 Y_{t-2} + \dots \\
&\quad \dots + (1 - K_t) \times \dots \times K_1 F^{t-1} Y_1 + (1 - K_t) \times \dots \times K_1 F^{t-1} \mu \\
V_t^t &= \frac{R [F^2 (V_{t-1}^{t-1}) + Q]}{R + [F^2 (V_{t-1}^{t-1}) + Q]} \\
&= \frac{R Q + R F^2 V_{t-1}^{t-1}}{R + Q + F^2 V_{t-1}^{t-1}}
\end{aligned}$$

A.3.2 One Y to many X s

Here we show how the Kalman recursions work, when the vector of observations Y has a lower dimension than the vector of states X : in this example Y is a scalar and X is a 2-dimensional vector. The *information* coming from Y is spread over the states in X ; the attribution takes into account the initial variance-covariance matrix of X and that of the corresponding error, and the variance of Y , and outlines how the states in X are projected, in the sense of L_2 , on the space spanned by the observation Y . Start from

$$X_1^0 = \begin{bmatrix} m_1 \\ m_2 \end{bmatrix}, \quad V_1^0 = \begin{bmatrix} s_1 & 0 \\ 0 & s_2 \end{bmatrix}$$

compute the gain and the posteriors

$$\begin{aligned}
K_1 &= \begin{bmatrix} \frac{s_1}{R+s_1+s_2} \\ \frac{s_2}{R+s_1+s_2} \end{bmatrix} \\
X_1^1 &= \begin{bmatrix} \frac{R+s_2}{R+s_1+s_2} m_1 - \frac{s_1}{R+s_1+s_2} m_2 + \frac{s_1}{R+s_1+s_2} Y_1 \\ \frac{R+s_1}{R+s_1+s_2} m_2 - \frac{s_2}{R+s_1+s_2} m_1 + \frac{s_2}{R+s_1+s_2} Y_1 \end{bmatrix} \\
V_1^1 &= \begin{bmatrix} \frac{s_1(R+s_2)}{R+s_1+s_2} & -\frac{s_1 s_2}{R+s_1+s_2} \\ -\frac{s_1 s_2}{R+s_1+s_2} & \frac{s_2(R+s_1)}{R+s_1+s_2} \end{bmatrix}
\end{aligned}$$

then project and correct as before

$$\begin{aligned}
X_2^1 &= \begin{bmatrix} f_1 \left(\frac{R+s_2}{R+s_1+s_2} m_1 - \frac{s_1}{R+s_1+s_2} m_2 + \frac{s_1}{R+s_1+s_2} Y_1 \right) \\ f_2 \left(\frac{R+s_1}{R+s_1+s_2} m_2 - \frac{s_2}{R+s_1+s_2} m_1 + \frac{s_2}{R+s_1+s_2} Y_1 \right) \end{bmatrix} \\
&=: \begin{bmatrix} f_1 \mu_1^1 \\ f_2 \mu_2^1 \end{bmatrix} \\
V_2^1 &= \begin{bmatrix} q_1 + f_1^2 \left(\frac{s_1(R+s_2)}{R+s_1+s_2} \right) & -f_1 f_2 \left(\frac{s_1 s_2}{R+s_1+s_2} \right) \\ -f_1 f_2 \left(\frac{s_1 s_2}{R+s_1+s_2} \right) & q_2 + f_2^2 \left(\frac{s_2(R+s_1)}{R+s_1+s_2} \right) \end{bmatrix} \\
&=: \begin{bmatrix} q_1 + f_1^2 \sigma_{1,1}^1 & -f_1 f_2 \sigma_{1,2}^1 \\ -f_1 f_2 \sigma_{1,2}^1 & q_2 + f_2^2 \sigma_{2,2}^1 \end{bmatrix}
\end{aligned}$$

to obtain the recursions

$$\begin{aligned}
K_{t+1} &= \begin{bmatrix} \frac{q_1 + f_1^2 \sigma_{1,1}^t - f_1 f_2 \sigma_{1,2}^t}{R + q_1 + q_2 + f_1^2 \sigma_{1,1}^t + f_2^2 \sigma_{2,2}^t - 2 f_1 f_2 \sigma_{1,2}^t} \\ \frac{q_2 + f_2^2 \sigma_{2,2}^t - f_1 f_2 \sigma_{1,2}^t}{R + q_1 + q_2 + f_1^2 \sigma_{1,1}^t + f_2^2 \sigma_{2,2}^t - 2 f_1 f_2 \sigma_{1,2}^t} \end{bmatrix} \\
V_{t+1}^{t+1} &= \begin{bmatrix} \frac{(q_1 + f_1^2 \sigma_{1,1}^t) [R + (q_2 + f_2^2 \sigma_{2,2}^t)] - (f_1 f_2 \sigma_{1,2}^t)^2}{R + q_1 + q_2 + f_1^2 \sigma_{1,1}^t + f_2^2 \sigma_{2,2}^t - 2 f_1 f_2 \sigma_{1,2}^t} & -\frac{R f_1 f_2 \sigma_{1,2}^t + (q_1 + f_1^2 \sigma_{1,1}^t) (q_2 + f_2^2 \sigma_{2,2}^t) - (f_1 f_2 \sigma_{1,2}^t)^2}{R + q_1 + q_2 + f_1^2 \sigma_{1,1}^t + f_2^2 \sigma_{2,2}^t - 2 f_1 f_2 \sigma_{1,2}^t} \\ -\frac{R f_1 f_2 \sigma_{1,2}^t + (q_1 + f_1^2 \sigma_{1,1}^t) (q_2 + f_2^2 \sigma_{2,2}^t) - (f_1 f_2 \sigma_{1,2}^t)^2}{R + q_1 + q_2 + f_1^2 \sigma_{1,1}^t + f_2^2 \sigma_{2,2}^t - 2 f_1 f_2 \sigma_{1,2}^t} & \frac{(q_2 + f_2^2 \sigma_{2,2}^t) [R + (q_1 + f_1^2 \sigma_{1,1}^t)] - (f_1 f_2 \sigma_{1,2}^t)^2}{R + q_1 + q_2 + f_1^2 \sigma_{1,1}^t + f_2^2 \sigma_{2,2}^t - 2 f_1 f_2 \sigma_{1,2}^t} \end{bmatrix} \\
X_{t+1}^{t+1} &= \begin{bmatrix} \frac{[R + (q_2 + f_2^2 \sigma_{2,2}^t) - f_1 f_2 \sigma_{1,2}^t] \cdot f_1 \mu_1^t - [(q_1 + f_1^2 \sigma_{1,1}^t) - f_1 f_2 \sigma_{1,2}^t] \cdot f_2 \mu_2^t + [(q_1 + f_1^2 \sigma_{1,1}^t) - f_1 f_2 \sigma_{1,2}^t] \cdot Y_2}{R + q_1 + q_2 + f_1^2 \sigma_{1,1}^t + f_2^2 \sigma_{2,2}^t - 2 f_1 f_2 \sigma_{1,2}^t} \\ \frac{[R + (q_1 + f_1^2 \sigma_{1,1}^t) - f_1 f_2 \sigma_{1,2}^t] \cdot f_2 \mu_2^t - [(q_2 + f_2^2 \sigma_{2,2}^t) - f_1 f_2 \sigma_{1,2}^t] \cdot f_1 \mu_1^t + [(q_2 + f_2^2 \sigma_{2,2}^t) - f_1 f_2 \sigma_{1,2}^t] \cdot Y_2}{R + q_1 + q_2 + f_1^2 \sigma_{1,1}^t + f_2^2 \sigma_{2,2}^t - 2 f_1 f_2 \sigma_{1,2}^t} \end{bmatrix}
\end{aligned}$$

B The Key to fig n.10.

- 1 Alcoa Users*
- 2 Alumni Network*
- 3 Anycast Addresses*
- 4 CERT*
- 5 CS*
- 6 Cyert Lab 128*
- 7 Electrical Computer Engineering*
- 8 GSIA QuickReg*
- 9 Groats*
- 10 Gruel*
- 11 INI Projects*
- 12 LAB-TEST 128.237*
- 13 Loopback (New)*
- 14 Loopback Router*
- 15 Math - Parallel Cluster*
- 16 Off Campus*
- 17 PSC Private*
- 18 Reserved 100*
- 19 S29-Reserved*
- 20 S30-Alcoa (Alcoa-side ISDN)*
- 21 S30-BuyersMart (T1)*
- 22 SEI 2*
- 23 Spirit House (dorm)*
- 24 Telerama CMU Housing*
- 25 Test Subnet*
- 26 VPN Local IP Space*
- 27 West-Net*
- 28 _offcampus - theplanet.com*
- 29 authbridge
- 30 bp
- 31 campus
- 32 cfa
- 33 cyh
- 34 cyh-a100
- 35 dh
- 36 gsia
- 37 hbh
- 38 hl
- 39 hyper
- 40 macosxlabs.com*
- 41 mi
- 42 mmch
- 43 ptc
- 44 res
- 45 sei-private*
- 46 sysdev
- 47 uc
- 48 voip
- 49 weh

References

- [1] Z. Bi, C.N. Faloutsos, and F. Korn. The “DGX” distribution for mining massive, skewed data. In *Seventh International ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, San Francisco, 2001.
- [2] J. Cao, D. Davis, S. VanDer Wiel, and B. Yu. Time-varying network tomography: router link data. *Journal of the American Statistical Association*, 95:1063–1075, 2000.
- [3] A.P. Dempster, N.M. Laird, and D.B. Rubin. Maximum likelihood from incomplete data via the EM algorithm, with discussion. *Journal of the Royal Statistical Society, Series B, Methodological*, 39:1–38, 1977.
- [4] I.S. Duff, A.M. Erisman, and J.K. Reid. *Direct methods for sparse matrices*. Clarendon Press, 1986.
- [5] S.E. Fienberg. An iterative procedure for estimation in contingency tables. *Annals of Mathematical Statistics*, 41:907–917, 1970a.
- [6] W.R. Gilks and C. Berzuini. Following a moving target — Monte Carlo inference for dynamic Bayesian models. *Journal of the Royal Statistical Society, Series B, Methodological*, 63:127–146, 2001.
- [7] T. Higuchi. Self-organizing time series model. In A. Doucet, N. de Freitas, and N. Gordon, editors, *Sequential Monte Carlo Methods in Practice*, pages 429–444. Springer-Verlag, 2001.
- [8] K. Lange. A gradient algorithm locally equivalent to the EM algorithm. *Journal of the Royal Statistical Society, Series B, Methodological*, 57:425–437, 1995.
- [9] C. Robert and G. Casella. *Monte Carlo Statistical Methods*. Springer-Verlag, 1999.
- [10] C. Tebaldi and M. West. Bayesian inference on network traffic using link count data. *Journal of the American Statistical Association*, 93:557–576, 1998.
- [11] R.J. Vanderbei and J. Iannone. An EM approach to od matrix estimation. Technical Report SOR 94-04, Princeton University, 1994.
- [12] Y. Vardi. Network tomography: estimating source-destination traffic intensities from link data. *Journal of the American Statistical Association*, 91:365–377, 1996.
- [13] M. West and J. Harrison. *Bayesian forecasting and dynamic models*. Springer-Verlag, 1997.