

Supplementary material: Fast interior-point inference in high-dimensional sparse, penalized state-space models

Eftychios A. Pnevmatikakis and Liam Paninski

1 Analysis of the Low Rank Approximation

We examine the number of singular values that are needed to capture a fraction θ of energy of U_t . If r is that number then the Singular Value Decomposition $L\Sigma L^T$ solves the following problem

$$\min \|U - L\Sigma L^T\|_F \text{ such that } \text{rank}(L\Sigma L^T) = r, \quad (1)$$

where $\|\cdot\|_F$ is the Frobenius norm and we have dropped the subscripts for simplicity. Suppose that each B_t is a d -dimensional gaussian vector with iid $\mathcal{N}(0, 1)$ entries and that each $\tilde{D}_t^{-1}E_t = \alpha I_d$ with $0 < \alpha < 1$. Then U is a random matrix with $[U]_{ij} \sim \mathcal{N}(0, \alpha^{2(i-1)})$. Let U_l be the matrix that consists of the first l rows of U and define k as the minimum number of rows required to capture a θ fraction of the energy,

$$k = \arg \min\{l : \mathbb{E}\|U_l\|_F^2 \geq \theta \mathbb{E}\|U\|_F^2\}. \quad (2)$$

We claim that with high probability $k \geq r$. To compute k we have

$$\begin{aligned} \mathbb{E}\|U_l\|_F^2 &= d \frac{1 - \alpha^{2l}}{1 - \alpha^2} \Rightarrow \\ \mathbb{E}\|U_l\|_F^2 \geq \theta \mathbb{E}\|U\|_F^2 &\Leftrightarrow (1 - \alpha^{2l}) \geq \theta(1 - \alpha^{2t}) \Rightarrow \\ k &= \left\lceil \frac{\log(1 - \theta(1 - \alpha^{2t}))}{2 \log(\alpha)} \right\rceil, \end{aligned} \quad (3)$$

where $\lceil \cdot \rceil$ is the ceil function. Note that k is independent of d . Therefore, we expect our low rank approximation to give substantial computational gains if

$$d \gg \left\lceil \frac{\log(1 - \theta(1 - \alpha^{2t}))}{2 \log(\alpha)} \right\rceil. \quad (4)$$

We can also compute a bound on the deviation of the effective rank of U from $k + c$ for some positive integer c , using large deviations theory. A weaker version of this is computing the deviation of $\|U_k\|_F^2$ from $\mathbb{E}(\|U_k\|_F^2)$ by estimating the probability

$$\mathbb{P}(\|U_{k+c}\|_F^2 \leq \mathbb{E}(\|U_k\|_F^2)). \quad (5)$$

This is the probability that more than $k + c$ rows are required to capture the θ fraction of the expected energy. Therefore this constitutes a bound on the probability that the effective rank

of U will be greater than $k + c$. $\|U_{k+c}\|_F^2$ can be considered as the sum of $k + c$ i.i.d. random variables Q_i , with

$$\|U_{k+c}\|_F^2 = \sum_{i=1}^{k+c} \alpha^{2(i-1)} Q_i, \quad (6)$$

where each Q_i is a chi-squared distribution with d degrees of freedom. Then from Cramer's theorem (Dembo and Zeitouni, 1993) we have that

$$\mathbb{P}(\|U_{k+c}\|_F^2 \leq \mathbb{E}(\|U_k\|_F^2)) \leq \exp(-d\kappa(\mathbb{E}(\|U_k\|_F^2))), \quad (7)$$

with

$$\kappa(\mathbb{E}(\|U_k\|_F^2)) := \sup_t \left(\mathbb{E}(t\|U_k\|_F^2) - \log(\mathbb{E}(e^{t\|U_{k+c}\|_F^2})) \right). \quad (8)$$

By using the moment generating function for a chi-squared random variable (which is defined on the interval $(-\infty, 0.5)$) we have

$$\kappa(\mathbb{E}(\|U_k\|_F^2)) := \sup_{t < 0.5} \underbrace{\left(t\mathbb{E}(\|U_k\|_F^2) + \frac{1}{2} \sum_{i=1}^{k+c} \log(1 - 2t\alpha^{2(i-1)}) \right)}_{f(t)}. \quad (9)$$

The maximizing t cannot be found in closed form. However, it can be shown that $f(t)$ is concave and that $f'(0) < 0$. As a result $\kappa(\|U_l\|_F^2) > f(0) > 0$. Therefore, the probability of a fixed deviation from the expected number of required rows k decays exponentially with the dimension d . Moreover, for a fixed d , numerical simulations show that the probability falls sharply with the order of the deviation. The exact rate will be pursued elsewhere.

In a similar way, we can also compute a bound on the slightly more relevant probability. Assuming $T \rightarrow \infty$

$$\mathbb{P}(\|U_{k+c}\|_F^2 \leq \theta \|U\|_F^2) = \mathbb{P}\left(\|U_{k+c}\|_F^2 \leq \frac{\theta}{1-\theta} \|U_{\setminus(k+c)}\|_F^2\right) = \mathbb{P}\left(\|U_{k+c}\|_F^2 \leq \frac{\theta}{1-\theta} \alpha^{2(k+c)} \|V\|_F^2\right), \quad (10)$$

where $U_{\setminus l}$ is the matrix U without its first l rows and V is an independent copy of U . Following the same reasoning as before, and using that $\alpha^{2k} \approx 1 - \theta$

$$\mathbb{P}(\|U_{k+c}\|_F^2 \leq \theta \|U\|_F^2) \leq \exp\left(-\frac{d}{2} \sup_{-\frac{\alpha^{-2c}}{2\theta} < t < \frac{1}{2}} \left(\sum_{i=1}^{k+c} \log(1 - 2t\alpha^{2(i-1)}) + \sum_{i=1}^{\infty} \log(1 + 2t\theta\alpha^{2c}\alpha^{2(i-1)}) \right)\right) \quad (11)$$

It can again be shown that the supremum is greater than zero for all $c > 0$, and that it also increases with c , which establishes that the probability of the effective rank being greater than the bound of (4) falls exponentially with the dimension d and sharply with the order of the deviation c .

Remark 1.1. *The bound of (4) is in practice rather loose. A more detailed analysis shows that with the inclusion of the “noise term” $(F_t^{-1} + U_t \tilde{D}_t^{-1} U_t^T)^{-1/2}$, the effective rank drops, and (4) appears in the limiting situation where the measurement noise is infinite. Moreover, our analysis does not account for the recursive nature of the low rank approximations. Using these facts tighter bounds can be derived. A detailed analysis is presented in (Pneumatikakis et al., 2012).*

2 Proof of \tilde{H} being positive definite

We can write the forward-backward recursion of the Block-Thomas algorithm in matrix-vector form. The backward recursion

$$\begin{aligned} \mathbf{s}_T &= \mathbf{q}_T, \\ \mathbf{s}_t &= \mathbf{q}_t - \Gamma_t \mathbf{s}_{t+1}, \quad t = T-1, \dots, 1 \end{aligned} \tag{12}$$

can be written as

$$\begin{bmatrix} \mathbf{s}_1 \\ \mathbf{s}_2 \\ \vdots \\ \mathbf{s}_{T-1} \\ \mathbf{s}_T \end{bmatrix} = - \underbrace{\begin{bmatrix} 0 & \Gamma_1 & 0 & \dots & 0 \\ 0 & 0 & \Gamma_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \ddots & \vdots \\ 0 & 0 & \dots & 0 & \Gamma_{T-1} \\ 0 & 0 & \dots & 0 & 0 \end{bmatrix}}_{\Gamma} \begin{bmatrix} \mathbf{s}_1 \\ \mathbf{s}_2 \\ \vdots \\ \mathbf{s}_{T-1} \\ \mathbf{s}_T \end{bmatrix} + \begin{bmatrix} \mathbf{q}_1 \\ \mathbf{q}_2 \\ \vdots \\ \mathbf{q}_{T-1} \\ \mathbf{q}_T \end{bmatrix}. \tag{13}$$

Similarly, the forward recursion

$$\begin{aligned} \mathbf{q}_1 &= -M_1^{-1} \nabla_1, \\ \mathbf{q}_t &= -M_t^{-1} (\nabla_t + E_{t-1} \mathbf{q}_{t-1}), \quad t = 2, \dots, T \end{aligned} \tag{14}$$

can be written in matrix-vector form as

$$\begin{bmatrix} \mathbf{q}_1 \\ \mathbf{q}_2 \\ \vdots \\ \mathbf{q}_{T-1} \\ \mathbf{q}_T \end{bmatrix} = - \underbrace{\begin{bmatrix} 0 & 0 & \dots & 0 & 0 \\ M_2^{-1} E_1 & 0 & \dots & 0 & 0 \\ \vdots & \ddots & \ddots & \vdots & \vdots \\ 0 & \dots & M_{T-1}^{-1} E_{T-2} & 0 & 0 \\ 0 & \dots & 0 & M_T^{-1} E_{T-1} & 0 \end{bmatrix}}_E \begin{bmatrix} \mathbf{q}_1 \\ \mathbf{q}_2 \\ \vdots \\ \mathbf{q}_{T-1} \\ \mathbf{q}_T \end{bmatrix} - \underbrace{\begin{bmatrix} M_1^{-1} & 0 & \dots & 0 & 0 \\ 0 & M_2^{-1} & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \dots & M_{T-1}^{-1} & 0 \\ 0 & 0 & \dots & 0 & M_T^{-1} \end{bmatrix}}_{M^{-1}} \begin{bmatrix} \nabla_1 \\ \nabla_2 \\ \vdots \\ \nabla_{T-1} \\ \nabla_T \end{bmatrix} \tag{15}$$

Combining (13) and (15) we have

$$\mathbf{s} = -(I + \Gamma)^{-1} (I + E)^{-1} M^{-1} \nabla, \tag{16}$$

where Γ, E, M are matrices defined in (13) and (15). Since $\mathbf{s} = -H^{-1} \nabla$ it follows that the Hessian is equal to

$$H = M(I + E)(I + \Gamma). \tag{17}$$

In the case of the LRBT algorithm, if we define $\tilde{M}_t^{-1} = \tilde{D}_t^{-1} - L_t \Sigma_t L_t^T$ and $\tilde{\Gamma}_t = \tilde{M}_t^{-1} E_t^T$, we have that

$$\begin{aligned} \tilde{\mathbf{q}}_t &= -\tilde{M}_t^{-1} (\nabla_t + E_{t-1} \tilde{\mathbf{q}}_{t-1}) \\ \tilde{\mathbf{s}}_t &= \tilde{\mathbf{q}}_t - \tilde{\Gamma}_t \tilde{\mathbf{s}}_{t+1}. \end{aligned} \tag{18}$$

Therefore, an equivalent representation holds in the sense that

$$\tilde{\mathbf{s}} = -\tilde{H}^{-1} \nabla, \quad \text{with } \tilde{H} = \tilde{M}(I + \tilde{E})(I + \tilde{\Gamma}), \tag{19}$$

where the block matrices $\tilde{M}, \tilde{E}, \tilde{\Gamma}$ are defined in the same way as their exact counterparts M, E, Γ . We can rewrite \tilde{H} as

$$\tilde{H} = \tilde{M}(I + \tilde{E})\tilde{M}^{-1}\tilde{M}(I + \tilde{\Gamma}) \quad (20)$$

A straight calculation shows that

$$\tilde{M}(I + \tilde{\Gamma}) = (\tilde{M}(I + \tilde{E}))^T = \begin{bmatrix} \tilde{M}_1 & E_1^T & 0 & \dots & 0 \\ 0 & \tilde{M}_2 & E_2^T & \dots & 0 \\ \vdots & \vdots & \ddots & \ddots & \vdots \\ 0 & 0 & \dots & \tilde{M}_{T-1} & E_{T-1}^T \\ 0 & 0 & \dots & 0 & \tilde{M}_t \end{bmatrix}, \quad (21)$$

and the approximate Hessian can be written as

$$\tilde{H} = \begin{bmatrix} \tilde{M}_1 & E_1^T & 0 & \dots & 0 \\ 0 & \tilde{M}_2 & E_2^T & \dots & 0 \\ \vdots & \vdots & \ddots & \ddots & \vdots \\ 0 & 0 & \dots & \tilde{M}_{T-1} & E_{T-1}^T \\ 0 & 0 & \dots & 0 & \tilde{M}_t \end{bmatrix}^T \begin{bmatrix} \tilde{M}_1^{-1} & 0 & 0 & \dots & 0 \\ 0 & \tilde{M}_2^{-1} & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \ddots & \vdots \\ 0 & 0 & \dots & \tilde{M}_{T-1}^{-1} & 0 \\ 0 & 0 & \dots & 0 & \tilde{M}_t^{-1} \end{bmatrix} \begin{bmatrix} \tilde{M}_1 & E_1^T & 0 & \dots & 0 \\ 0 & \tilde{M}_2 & E_2^T & \dots & 0 \\ \vdots & \vdots & \ddots & \ddots & \vdots \\ 0 & 0 & \dots & \tilde{M}_{T-1} & E_{T-1}^T \\ 0 & 0 & \dots & 0 & \tilde{M}_t \end{bmatrix}, \quad (22)$$

or

$$\tilde{H} = \begin{bmatrix} \tilde{M}_1 & E_1^T & 0 & \dots & 0 \\ E_1 & \tilde{M}_2 + E_1\tilde{M}_1^{-1}E_1^T & E_2^T & \dots & 0 \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ 0 & 0 & \dots & \tilde{M}_{T-1} + E_{T-2}\tilde{M}_{T-2}^{-1}E_{T-1}^T & E_{T-1}^T \\ 0 & 0 & \dots & E_{T-1} & \tilde{M}_t + E_{T-1}\tilde{M}_{T-1}^{-1}E_{T-1} \end{bmatrix}. \quad (23)$$

From (22) it follows that \tilde{H} is positive definite (PD), if the matrices \tilde{M}_t are also PD.

Lemma 2.1. *The matrices $\tilde{D}_t, t = 1, \dots, T$ are PD.*

Proof. In the case where A and V commute and A is stable, the matrix \tilde{D}_t is equal to

$$\tilde{D}_t = V^{-1}(I - (A^T A)^t)^{-1}(I - (A^T A)^{t+1}),$$

which is PD, by stability of A . The result holds also in the case where A and V do not commute, although the formulas are more complicated. \square

Lemma 2.2. *The matrices $\tilde{M}_t, t = 1, \dots, T$ are PD for any choice of the threshold θ .*

Proof. We introduce the matrices \hat{M}_t , defined as follows:

$$\begin{aligned} \hat{M}_1 &= M_1 \\ \hat{M}_t &= D_t + B_t^T W_t^{-1} B_t - E_{t-1} \tilde{M}_{t-1}^{-1} E_{t-1}^T. \end{aligned} \quad (24)$$

These matrices are the matrices obtained from the exact BT recursion $M_t = D_t + B_t^T W_t^{-1} B_t - E_{t-1} M_{t-1}^{-1} E_{t-1}^T$, applied to the approximate matrices \tilde{M}_{t-1}^{-1} . By using the relations

$$\begin{aligned} \tilde{M}_t^{-1} &= \tilde{D}_t^{-1} - L_t \Sigma_t L_t^T \\ \tilde{D}_t &= D_t - E_{t-1} \tilde{D}_{t-1}^{-1} E_{t-1}^T, \end{aligned} \quad (25)$$

we can rewrite \hat{M}_t as

$$\hat{M}_t = \tilde{D}_t + B_t^T W_t^{-1} B_t + E_{t-1} L_{t-1} \Sigma_{t-1} L_{t-1}^T E_{t-1}^T = \tilde{D}_t + O_t Q_t O_t^T. \quad (26)$$

Using (26) we see that \hat{M}_t is the sum of a PD matrix (\tilde{D}_t), and two semipositive definite (SPD) matrices (Σ_t is always PD by definition). Therefore, \hat{M}_t^{-1} is also PD and equals

$$\hat{M}_t^{-1} = \tilde{D}_t^{-1} - \underbrace{\tilde{D}_t^{-1} O_t (Q_t^{-1} + O_t^T \tilde{D}_t^{-1} O_t)^{-1} O_t^T \tilde{D}_t^{-1}}_{G_t}. \quad (27)$$

Now \tilde{M}_t^{-1} is obtained by the low rank approximation of G_t . We can write the singular value decomposition of G_t as

$$G_t = \begin{bmatrix} L_t & R_t \end{bmatrix} \begin{bmatrix} \Sigma_t & 0 \\ 0 & S_t \end{bmatrix} \begin{bmatrix} L_t^T \\ R_t^T \end{bmatrix}, \quad (28)$$

and have that

$$\tilde{M}_t^{-1} - \hat{M}_t^{-1} = R_t S_t R_t^T \quad (29)$$

Therefore $\tilde{M}_t^{-1} - \hat{M}_t^{-1}$ is SPD. Consequently \tilde{M}_t is the sum a PD and a SPD matrix and thus is PD. \square

A detailed proof of Theorem 3.4 will be presented in (Pnevmatikakis et al., 2012).

References

- Dembo, A. and Zeitouni, O. (1993). *Large deviations techniques and applications*. Springer, New York.
- Pnevmatikakis, E. A., Paninski, L., Rad, K. R., and Huggins, J. (2012). Fast Kalman filtering and forward-backward smoothing via a low-rank perturbative approach. In preparation.