# Variational Inference for the PLDS

Lars Buesing

March 6, 2017

## 1 Model

$$
\begin{aligned}
p(\mathbf{x}) &= \mathcal{N}(\mathbf{x}|\mu, \Sigma) && (1) \\
p(\mathbf{y}|\mathbf{x}) &= \prod_n p(y_n|\eta_n), \quad \eta := W\mathbf{x} + \overline{\mathbf{d}} && (2)
\end{aligned}
$$

For latent dynamical system we know that $\Lambda := \Sigma^{-1}$ is tri-diagonal, and $W$ block-diagonal:

$$
\Lambda = \begin{pmatrix}
Q_0^{-1} + AQ^{-1}A^\top & A^\top Q^{-1} & & \\
Q^{-1}A & Q^{-1} + AQ^{-1}A^\top & A^\top Q^{-1} & \\
& \ddots & \ddots & \ddots
\end{pmatrix} \tag{3}
$$

$$
W = \text{blk-diag}(\underbrace{C, \ldots, C}_{T\text{-times}}) \tag{4}
$$

## 2 Inference problem

Gaussian variational approximation:

$$
q(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\mathbf{m}, V) \tag{5}
$$

Variational lower bound:

$$
\begin{aligned}
\mathcal{L}(\mathbf{m}, V) &\leq \log p(\mathbf{y}) && (6) \\
\mathcal{L}(\mathbf{m}, V) &= \frac{1}{2}\left(\log|V| - \text{tr}[\Sigma^{-1}V] - (\mathbf{m}-\mu)^\top \Sigma^{-1}(\mathbf{m}-\mu)\right) + \sum_n \mathbb{E}_{q(\mathbf{x})}[\log p(y_n|\eta_n)] && (7)
\end{aligned}
$$

For Poisson with exp link function we can compute $\mathbb{E}_{q(\mathbf{x})}[\log p(y_n|\eta_n)]$, otherwise (eg for Bernoulli observations) use a local variational lower bound on the integrated likelihood:

$$
\begin{aligned}
\mathbb{E}_{q(\mathbf{x})}[\log p(y_n|\eta_n)] &\geq -f_n(\overline{m}_n, \overline{v}_n) && (8) \\
f_n(\overline{m}_n, \overline{v}_n) &= -y_n\overline{m}_n + \exp(\overline{m}_n + \overline{v}_n/2) && (9) \\
\overline{\mathbf{m}} &:= W\mathbf{m} + \overline{\mathbf{d}} && (10) \\
\overline{\mathbf{v}} &:= \text{diag}(WVW^\top) && (11)
\end{aligned}
$$

The bound then reads:

$$
\mathcal{L}(\mathbf{m}, V) = \frac{1}{2}\left(\log|V| - \text{tr}[\Sigma^{-1}V] - (\mathbf{m}-\mu)^\top \Sigma^{-1}(\mathbf{m}-\mu)\right) - \sum_n f_n(\overline{m}_n, \overline{v}_n) \tag{12}
$$

For convex $f_n$ (true for exp-PLDS): strictly concave optimization in $\mathbf{m}, V$
Possible optimization strategies:

1. Direct optimization over $\mathbf{m}, V$: strictly concave, however $V$ dense; does not make use of Markovian structure of the model

2. Optimization over $\mathbf{m}, V^{-1}$: Opper et al show that optimal $V^* = (\Sigma^{-1} + W^\top \text{diag}(\lambda)W)^{-1}$; hence for tri-diagonal $\Sigma^{-1}$ and block-diagonal $W$ then $V^*$ is also tri-diagonal; however optimization over $\mathbf{m}, \lambda$ is not convex and converges slowly according to [Seeger et al. ICML2013]

3. Solve the dual optimization as proposed in [Seeger et al. ICML2013]: convex, makes use of Markovian structure of the model

# 3 Variational inference via dual optimization

## 3.1 Optimization to solve

Dual problem:

$$\min_{\lambda} \quad \frac{1}{2}(\lambda - \mathbf{y})^\top W\Sigma W^\top (\lambda - \mathbf{y}) - (W\mu + \overline{\mathbf{d}})^\top (\lambda - \mathbf{y}) - \frac{1}{2}\log|A_\lambda| + \sum_n f^*(\lambda_n) \tag{13}$$

$$\text{subject to} \quad \lambda_i > 0$$

where

$$f^*(\lambda_i) \quad := \quad \lambda_i(\log \lambda_i - 1) \tag{14}$$

$$A_\lambda \quad := \quad \Sigma^{-1} + W^\top \operatorname{diag}(\lambda)W \tag{15}$$

The optimal variational parameters for $q(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\mathbf{m}^*, V^*)$ are given by:

$$\mathbf{m}^* \quad = \quad \mu - \Sigma W^\top(\lambda^* - \mathbf{y}) \tag{16}$$

$$V^* \quad = \quad (\Sigma^{-1} + W^\top \operatorname{diag}(\lambda^*)W)^{-1} = A_{\lambda^*}^{-1} \tag{17}$$

## 3.2 How to optimize?

Use gradient based methods:

$$\nabla_\lambda \quad = \quad W\Sigma W^\top(\lambda - \mathbf{y}) - W\mu - \overline{\mathbf{d}} + \log \lambda - \frac{1}{2}\operatorname{diag}(WA_\lambda^{-1}W^\top)$$

$$= \quad \underbrace{W\Sigma W\lambda}_{O(N)} + \underbrace{\log \lambda}_{O(N)} - \frac{1}{2}\operatorname{diag}(\underbrace{WA_\lambda^{-1}W^\top}_{O(N)}) - \underbrace{W(\Sigma W^\top \mathbf{y} + \mu)}_{\text{pre-compute}}$$

Hessian:

$$H_\lambda \quad = \quad \operatorname{diag}(\lambda)^{-1} + W\Sigma W^\top + (WA_\lambda^{-1}W^\top) \circ (WA_\lambda^{-1}W^\top)$$

Iterate:

$$\mathbf{m}^k \quad = \quad \mu + \Sigma W^\top \mathbf{y} - \Sigma W^\top \lambda^k$$

$$A^k \quad = \quad \Sigma^{-1} + W^\top \operatorname{diag}(\lambda^k)W$$

$$\nabla^k \quad = \quad \log \lambda^k - W\mathbf{m}^k - \overline{\mathbf{d}} - \frac{1}{2}\operatorname{diag}(W(A^k)^{-1}W^\top)$$

$$\lambda^{k+1} \quad = \quad \lambda^k - \nu\nabla^k$$

Computing the block-diagonal elements of $A^k$ is equivalent to Kalman smoothing and requires a forward-backward pass through the data which costs $O(Td^3)$.
What's the relation to Laplace approximation?

$$\nabla^k \quad = \quad -\Sigma^{-1}(\mathbf{x} - \mu) + W^\top(\mathbf{y} - \exp(W\mathbf{x} + \overline{\mathbf{d}}))$$

$$H^k \quad = \quad -(\Sigma^{-1} + W^\top \operatorname{diag}(\exp(W\mathbf{x} + \overline{\mathbf{d}}))W)$$

## 3.3 Kalman smoothing

The matrix $A_\lambda$ equals exactly the precision matrix of a LDS with dynamics given by $A, Q$ and observations sampled from $\mathcal{N}(C\mathbf{x}_t, \operatorname{diag}(\lambda_t))$. Hence, calculating the block-diagonal of $A_\lambda^{-1}$ is exactly equivalent to calculating the smoothed posterior covariance of this LDS. Let $P_{t|t}$ denote the filtered covariance, $P_{t+1|t}$ the one-step-ahead covariance and $P_{t|T}$ the smoothed covariance of this model.

$$(A_\lambda^{-1})_{(t-1)d+1:td,(t-1)d+1:td} \quad \overset{!}{=} \quad P_{t|T} \tag{18}$$

We use the Kalman smoother recursions. The forward pass reads:

$$P_{t+1|t} = AP_{t|t}A^\top + Q \tag{19}$$

$$P_{t+1|t+1} = \left(P_{t+1|t}^{-1} + C^\top \operatorname{diag}(\lambda_t)C\right)^{-1} \tag{20}$$

$$= \left(I + P_{t+1|t}C^\top \operatorname{diag}(\lambda_t)C\right) \backslash P_{t+1|t} \tag{21}$$

$$P_{0|0} = Q_0 \tag{22}$$

The backward pass is given by:

$$C_t = P_{t|t}A^\top / P_{t+1|t} \tag{23}$$

$$P_{t|T} = P_{t|t} + C_t \left(P_{t+1|T} - P_{t+1|t}\right) C_t^\top \tag{24}$$

The initialization for the backward pass $P_{T|T}$ is calculated the last step of the forward pass.

# 4 Appendix

## 4.1 Derivation of dual optimization

Original primal problem:

$$\max_{\mathbf{m},V} \quad \frac{1}{2}\left(\log|V| - \operatorname{tr}[\Sigma^{-1}V] - \|\mathbf{m}-\mu\|_{\Sigma^{-1}}^2\right) - \sum_n f_n(\overline{m}_n, \overline{v}_n)$$
$$\text{s.t.} \qquad\qquad V \in S^{++}$$

Expanded primal problem:

$$\operatorname{argmax}_{\mathbf{m},V,\rho,h} \quad \frac{1}{2}\left(\log|V| - \operatorname{tr}[\Sigma^{-1}V] - \|\mathbf{m}-\mu\|_{\Sigma^{-1}}^2\right) - \sum_n f_n(h_n, \rho_n)$$
$$\text{s.t.} \quad V \in S^{++}$$
$$h = W\mathbf{m} + \overline{\mathbf{d}}$$
$$\rho = \operatorname{diag}(WVW^\top)$$

Lagrangian:

$$\mathcal{L}(\mathbf{m},V,h,\rho,\alpha,\lambda) \;:=\; \frac{1}{2}\left(\log|V| - \operatorname{tr}[\Sigma^{-1}V] - \|\mathbf{m}-\mu\|_{\Sigma^{-1}}^2\right) - \sum_n f_n(h_n, \rho_n)$$
$$+\alpha^\top(h - W\mathbf{m} + \overline{\mathbf{d}}) + \frac{1}{2}\lambda^\top(\rho - \operatorname{diag}(WVW^\top))$$

Dual

$$D(\alpha,\lambda) \;:=\; \min_{\mathbf{m},V,h,\rho} L(\mathbf{m},V,h,\rho,\alpha,\lambda)$$
$$V^* \;=\; (\Sigma^{-1} + W^\top \operatorname{diag}(\lambda)W)^{-1} =: A_\lambda^{-1}$$
$$\mathbf{m}^* \;=\; \mu - \Sigma W^\top \alpha$$
$$\alpha \;=\; \lambda - y$$

Final reduced dual problem:

$$\operatorname{argmin}_\lambda \quad \frac{1}{2}(\lambda-\mathbf{y})^\top W\Sigma W(\lambda-\mathbf{y}) - (\overline{\mathbf{d}}+W\mu)^\top(\lambda-\mathbf{y}) - \frac{1}{2}\log|A_\lambda| + \sum_n f^*(\lambda_n)$$
$$\text{s.t.} \quad \lambda_i > 0$$

## 4.2 Duality basics

Primal problem with optimal values $p^*$:

$$\min \quad f(x)$$
$$\text{s.t.} \quad f_i(x) \le 0$$
$$h_i(x) = 0$$

Lagrange function with $\lambda_i \ge 0$:

$$L(x,\lambda,\nu) \;:=\; f(x) + \sum_i \lambda_i f_i(x) + \sum_i \nu_i h_i(x)$$

Dual:

$$g(\lambda,\nu) \;:=\; \inf_x L(x,\lambda,\nu)$$

Dual is a lower bound:

$$g(\lambda,\nu) \;\le\; p^*$$

This can be shown by bounding the constraint functions with linear functions from below. Dual problem:

$$\min \quad -g(\lambda,\nu)$$
$$\text{s.t.} \quad \lambda_i \ge 0$$

Dual is always convex!

**Slater's conditions** We have stong duality iff:

$$g(\lambda, \nu) \quad = \quad p^*$$

A sufficient condition for strong duality is: $f$ convex, no inequality constraints, primal feasible

**Dual function** The dual $f^*$ of a function $f$ is defined as:

$$f^*(y) \quad := \quad \sup_x \ y^\top x - f(x)$$