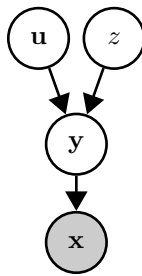


# Derivations for the GSM

Máté Lengyel

February 11, 2011

## 1 Definition



Following Wainwright and Simoncelli (2000), we define the GSM as the following generative model:<sup>1</sup>

$$\mathbf{u} \sim \mathcal{N}(\mathbf{u}; \mathbf{0}, \mathbf{C}) \quad (1)$$

$$z \sim \text{Gamma}(k, \theta) \quad (2)$$

$$\mathbf{y} = z \mathbf{u} \quad (3)$$

$$\mathbf{x}|\mathbf{y} \sim \mathcal{N}(\mathbf{A}\mathbf{y}, \sigma_{\mathbf{x}}^2 \mathbf{I}) \quad (4)$$

$$P(\mathbf{x}_{1:T}) = \prod_{t=1}^T P(\mathbf{x}_t) \quad (5)$$

---

<sup>1</sup>This is in fact a slight extension of the original model because Wainwright and Simoncelli (2000) only defined it down to the level of filter coefficient,  $\mathbf{y}$  (implicitly assuming that deterministic estimates at that level suffice), while we define the model all the way down to the level of pixels,  $\mathbf{x}$ .

## 2 Some preliminaries

$$P(\mathbf{x}|\mathbf{u}, z) = \mathcal{N}(\mathbf{x}; z \mathbf{A}\mathbf{u}, \sigma_x^2 \mathbf{I}) \quad (6)$$

$$= \mathcal{N}(z \mathbf{A}\mathbf{u}; \mathbf{x}, \sigma_x^2 \mathbf{I}) \quad (7)$$

$$= b \cdot \mathcal{N}(\mathbf{u}; \mathbf{m}, \mathbf{D}) \quad \text{for } N_u \leq N_x \quad (8)$$

with

$$b = \frac{\sqrt{|2\pi\mathbf{D}|}}{\sqrt{|2\pi\sigma_x^2\mathbf{I}|}} e^{-\frac{1}{2} \left( \frac{1}{\sigma_x^2} \mathbf{x}^\top \mathbf{x} - \mathbf{m}^\top \mathbf{D}^{-1} \mathbf{m} \right)} = z^{-N_u} \underbrace{\frac{\sqrt{|2\pi\sigma_x^2(\mathbf{A}^\top \mathbf{A})^{-1}|}}{\sqrt{|2\pi\sigma_x^2\mathbf{I}|}} e^{-\frac{1}{2\sigma_x^2} \mathbf{x}^\top [\mathbf{I} - \mathbf{A}(\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top] \mathbf{x}}}_{\text{does not depend on } \mathbf{u} \text{ or } z} \quad (9)$$

$$\mathbf{m} = \mathbf{D} (z \mathbf{A})^\top \frac{1}{\sigma_x^2} \mathbf{x} = \frac{1}{z} (\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top \mathbf{x} \quad (10)$$

$$\mathbf{D}^{-1} = z^2 \mathbf{A}^\top (\sigma_x^2 \mathbf{I})^{-1} \mathbf{A} = \frac{z^2}{\sigma_x^2} \mathbf{A}^\top \mathbf{A} \quad (11)$$

Technically, this only works for the (under)complete case (because  $\mathbf{A}^\top \mathbf{A}$  needs to be invertable), but as we shall see later this is not going to be a problem.

## 3 Inference

$$P(\mathbf{y}, \mathbf{u}, z|\mathbf{x}) = P(\mathbf{y}|\mathbf{u}, z, \mathbf{x}) P(\mathbf{u}|z, \mathbf{x}) P(z|\mathbf{x}) \quad (12)$$

### 3.1 Inferring $\mathbf{u}$

$$P(\mathbf{u}|z, \mathbf{x}) \propto P(\mathbf{u}) P(\mathbf{x}|\mathbf{u}, z) \quad (13)$$

$$\propto \mathcal{N}(\mathbf{u}; \mathbf{0}, \mathbf{C}) \mathcal{N}(\mathbf{u}; \mathbf{m}, \mathbf{D}) \quad (14)$$

$$= \mathcal{N}(\mathbf{u}; \boldsymbol{\mu}(z, \mathbf{x}), \boldsymbol{\Sigma}(z)) \quad (15)$$

with

$$\boldsymbol{\mu}(z, \mathbf{x}) = \boldsymbol{\Sigma}(z) (\mathbf{C}^{-1} \mathbf{0} + \mathbf{D}^{-1} \mathbf{m}) = \frac{z}{\sigma_x^2} \boldsymbol{\Sigma}(z) \mathbf{A}^\top \mathbf{x} \quad (16)$$

$$\boldsymbol{\Sigma}(z) = (\mathbf{C}^{-1} + \mathbf{D}^{-1})^{-1} = \left( \mathbf{C}^{-1} + \frac{z^2}{\sigma_x^2} \mathbf{A}^\top \mathbf{A} \right)^{-1} \quad (17)$$

So, as we see, both  $\boldsymbol{\mu}(z, \mathbf{x})$  and  $\boldsymbol{\Sigma}(z)$  are well-defined even in the overcomplete case.

We will also need the marginal posterior:

$$P(\mathbf{u}|\mathbf{x}) = \sum_z P(z|\mathbf{x}) P(\mathbf{u}|z, \mathbf{x}) \quad (18)$$

$$= \sum_z P(z|\mathbf{x}) \mathcal{N}(\mathbf{u}; \boldsymbol{\mu}(z, \mathbf{x}), \boldsymbol{\Sigma}(z)) \quad (19)$$

### 3.2 Inferring $\mathbf{y}$

It follows trivially from the above:

$$P(\mathbf{y}|\mathbf{u}, z, \mathbf{x}) = \delta(\mathbf{y} - z \mathbf{u}) \quad (20)$$

$$P(\mathbf{y}|z, \mathbf{x}) = \mathcal{N}(\mathbf{y}; z \boldsymbol{\mu}(z, \mathbf{x}), z^2 \boldsymbol{\Sigma}(z)) \quad (21)$$

$$P(\mathbf{y}|\mathbf{x}) = \sum_z P(z|\mathbf{x}) \mathcal{N}(\mathbf{y}; z \boldsymbol{\mu}(z, \mathbf{x}), z^2 \boldsymbol{\Sigma}(z)) \quad (22)$$

### 3.3 Inferring $z$

$$P(z|\mathbf{x}) \propto P(z) P(\mathbf{x}|z) \quad (23)$$

Rather than deriving  $P(\mathbf{x}|z)$ , the (marginal) likelihood of  $z$ , through lengthy algebraic manipulations, we build on the following intuition. One can think of  $\mathbf{x}$  (given  $z$ ) as a deterministically scaled (by  $z \mathbf{A}$ ) version of  $\mathbf{u}$  (a multivariate Gaussian random variable with known mean,  $\mathbf{0}$ , and covariance,  $\mathbf{C}$ ) plus a (multivariate) Gaussian noise term (with  $\mathbf{0}$  mean and  $\sigma_x^2 \mathbf{I}$  covariance). This insight yields a simple form for the probability of  $\mathbf{x}$  given  $z$ , which is just the likelihood we need (see also Bishop (2006, p. 93)):

$$P(\mathbf{x}|z) = \mathcal{N}(\mathbf{x}; \mathbf{0}, \sigma_x^2 \mathbf{I} + z^2 \mathbf{A} \mathbf{C} \mathbf{A}^\top) \quad (24)$$

## 4 Learning

The objective of learning is to maximise the likelihood of the parameters,  $\vartheta = \{\sigma_x^2, \mathbf{A}, \mathbf{C}, k, \theta\}$ :

$$P(\mathbf{x}_{1:T}|\vartheta) = \prod_t P(\mathbf{x}_t|\vartheta) = \prod_t \int dz P(z|k, \theta) P(\mathbf{x}_t|z, \sigma_x^2, \mathbf{A}, \mathbf{C}) \quad (25)$$

$$= \prod_t \int dz P(z|k, \theta) \mathcal{N}(\mathbf{x}_t; \mathbf{0}, \sigma_x^2 \mathbf{I} + z^2 \mathbf{A} \mathbf{C} \mathbf{A}^\top) \quad (26)$$

where in the last step we used Eq. 24.

Equation 26 reveals important invariances in the model. Namely, the likelihood (or predictive density) only depends on  $\mathbf{A} \mathbf{C} \mathbf{A}^\top$  and is thus invariant under the following reparametrisation:  $\mathbf{C} \rightarrow \mathbf{C}'$ ,  $\mathbf{A} \rightarrow \mathbf{A} \mathbf{C}'^{\frac{1}{2}} \mathbf{C}'^{-\frac{1}{2}}$ . Thus, it does not make sense to learn *both*  $\mathbf{C}$  and  $\mathbf{A}$ . Likewise, the scale of  $z$  (parametrised by  $\theta$  in the case of the Gamma prior we are using) can also be folded into either  $\mathbf{A}$  or  $\mathbf{C}$ , since  $z$  simply multiplies  $\mathbf{A} \mathbf{C} \mathbf{A}^\top$  and so the predictive density is invariant under  $\theta \rightarrow \theta'$ ,  $\mathbf{C} \rightarrow \frac{\theta^2}{\theta'^2} \mathbf{C}$  (or  $\mathbf{A} \rightarrow \frac{\theta}{\theta'} \mathbf{A}$ ).

Thus, a sensible combination of parameters to be learned would be to fix  $\mathbf{C}$ , eg. to be the identity, and  $\theta$ , eg. to be 1, and learn the rest of the parameters. However, in the context of modelling V1 it may be more instructive to fix  $\mathbf{A}$ , to be a filter bank of Gabor filters, reminiscent of V1 receptive fields (though  $\mathbf{A}$  really defines *projective* not receptive fields), and  $\theta$ , and learn the rest of the parameters (including  $\mathbf{C}$ ). Note that by fixing  $\mathbf{A}$  we do lose some expressive power in statistical terms but may obtain more interpretable results in biological terms. Nevertheless, for completeness, we derive the learning rules both for  $\mathbf{A}$  and  $\mathbf{C}$ . (And we don't derive the learning rule for  $k$  since there are no image data sets at our disposal that would be controlled for having natural global luminance / contrast statistics.)

We use EM (Dempster et al., 1977) to perform maximum likelihood iteratively. In each E-step we compute the (sufficient statistics of the) posterior over latent variables,  $\mathbf{u}$ ,  $\mathbf{y}$ , and  $z$ , and in each M-step we compute the values of the new parameters,  $\vartheta^*$  such that they maximise the negative free energy (Neal and Hinton, 1998) given the (sufficient statistics of the) posterior computed in the

previous E-step: <sup>2</sup>

$$\vartheta^* = \operatorname{argmax}_{\vartheta'} \sum_t \int dz d\mathbf{u} P(\mathbf{y}, \mathbf{u}, z | \mathbf{x}_t; \vartheta) \ln P(\mathbf{x}_t, \mathbf{y}, \mathbf{u}, z; \vartheta') \quad (27)$$

$$= \operatorname{argmax}_{\vartheta'} \sum_t \int dz d\mathbf{u} P(\mathbf{y}, \mathbf{u}, z | \mathbf{x}_t; \vartheta) \cdot \left[ \ln P(\mathbf{u}; \mathbf{C}') + \ln P(\mathbf{x}_t | \mathbf{y}; \sigma_x^{2'}, \mathbf{A}') + \dots \right] \quad (28)$$

$$= \operatorname{argmax}_{\vartheta'} \sum_t \int d\mathbf{u} P(\mathbf{u} | \mathbf{x}_t; \vartheta) \ln P(\mathbf{u}; \mathbf{C}') + \sum_t \int d\mathbf{y} P(\mathbf{y} | \mathbf{x}_t; \vartheta) \ln P(\mathbf{x}_t | \mathbf{y}; \sigma_x^{2'}, \mathbf{A}') + \dots \quad (29)$$

## 4.1 M-step for C

Preliminaries:

$$\ln P(\mathbf{u}; \mathbf{C}) \propto \ln |\mathbf{C}^{-1}| - \mathbf{u}^\top \mathbf{C}^{-1} \mathbf{u} + \dots \quad (30)$$

$$\frac{\partial}{\partial C_{ij}^{-1}} \ln P(\mathbf{u}; \mathbf{C}) \propto \operatorname{Tr} \left( \mathbf{C} \frac{\partial \mathbf{C}^{-1}}{\partial C_{ij}^{-1}} \right) - u_i u_j = C_{ij} - u_i u_j \quad (31)$$

$$\nabla_{\mathbf{C}^{-1}} \ln P(\mathbf{u}; \mathbf{C}) \propto \mathbf{C} - \mathbf{u} \mathbf{u}^\top \quad (32)$$

---

<sup>2</sup>Thus we are performing *complete* M steps to reach the minimum corresponding to the posterior in each iteration as opposed to incomplete M-steps which would just move us slightly towards that minimum – for a discussion of these issues see Neal and Hinton, 1998

And so:

$$\mathbf{C}^* = \operatorname{argmax}_{\mathbf{C}'} \sum_t \int d\mathbf{u} P(\mathbf{u}|\mathbf{x}_t) \ln P(\mathbf{u}; \mathbf{C}') \quad (33)$$

$$\mathbf{0} = \left. \nabla_{\mathbf{C}'^{-1}} \right|_{\mathbf{C}'^{-1}=\mathbf{C}^{*-1}} \sum_t \int d\mathbf{u} P(\mathbf{u}|\mathbf{x}_t) \ln P(\mathbf{u}; \mathbf{C}') \quad (34)$$

$$= \sum_t \int d\mathbf{u} P(\mathbf{u}|\mathbf{x}_t) \left. \nabla_{\mathbf{C}'^{-1}} \right|_{\mathbf{C}'^{-1}=\mathbf{C}^{*-1}} \ln P(\mathbf{u}; \mathbf{C}') \quad (35)$$

$$= \sum_t \sum_z P(z|\mathbf{x}_t) \int d\mathbf{u} P(\mathbf{u}|z, \mathbf{x}_t) (\mathbf{C}^* - \mathbf{u} \mathbf{u}^\top) \quad (36)$$

$$= \sum_t \sum_z P(z|\mathbf{x}_t) [\mathbf{C}^* - (\langle \mathbf{u}|z, \mathbf{x}_t \rangle \langle \mathbf{u}^\top|z, \mathbf{x}_t \rangle + \operatorname{Cov}[\mathbf{u}|z, \mathbf{x}_t])] \quad (37)$$

$$= \sum_t \sum_z P(z|\mathbf{x}_t) [\mathbf{C}^* - (\boldsymbol{\mu}(z, \mathbf{x}_t) \boldsymbol{\mu}^\top(z, \mathbf{x}_t) + \boldsymbol{\Sigma}(z))] \quad (38)$$

$$= T \mathbf{C}^* - \boldsymbol{\Xi}_{\mathbf{u}} \quad (39)$$

with

$$\boldsymbol{\Xi}_{\mathbf{u}} = \sum_z \boldsymbol{\lambda}(z) \quad (40)$$

$$\boldsymbol{\lambda}(z) = \boldsymbol{\Sigma}(z) l(z) + \sum_t P(z|\mathbf{x}_t) \boldsymbol{\mu}(z, \mathbf{x}_t) \boldsymbol{\mu}^\top(z, \mathbf{x}_t) \quad (41)$$

$$l(z) = \sum_t P(z|\mathbf{x}_t) \quad (42)$$

$$\mathbf{C}^* = \frac{1}{T} \boldsymbol{\Xi}_{\mathbf{u}} \quad (43)$$

## 4.2 M-step for A

Preliminaries:

$$\ln P(\mathbf{x}|\mathbf{y}; \sigma_{\mathbf{x}}^2, \mathbf{A}) \propto -(\mathbf{x} - \mathbf{A}\mathbf{y})^\top (\mathbf{x} - \mathbf{A}\mathbf{y}) + \dots = \sum_i \left( x_i - \sum_j A_{ij} y_j \right)^2 \quad (44)$$

$$\frac{\partial}{\partial A_{ij}} \ln P(\mathbf{x}|\mathbf{y}; \sigma_{\mathbf{x}}^2, \mathbf{A}) \propto - \left( x_i - \sum_{j'} A_{ij'} y_{j'} \right) y_j \quad (45)$$

$$\nabla_{\mathbf{A}} \ln P(\mathbf{x}|\mathbf{y}; \sigma_{\mathbf{x}}^2, \mathbf{A}) \propto -(\mathbf{x} - \mathbf{A}\mathbf{y}) \mathbf{y}^\top \quad (46)$$

And so:

$$\mathbf{A}^* = \underset{\mathbf{A}'}{\operatorname{argmax}} \sum_t \int d\mathbf{y} P(\mathbf{y}|\mathbf{x}_t) \ln P(\mathbf{x}_t|\mathbf{y}; \sigma_x^{2'}, \mathbf{A}') \quad (47)$$

$$\mathbf{0} = \left. \nabla_{\mathbf{A}'} \right|_{\mathbf{A}'=\mathbf{A}^*} \sum_t \int d\mathbf{y} P(\mathbf{y}|\mathbf{x}_t) \ln P(\mathbf{x}_t|\mathbf{y}; \sigma_x^{2'}, \mathbf{A}') \quad (48)$$

$$= \sum_t \int d\mathbf{y} P(\mathbf{y}|\mathbf{x}_t) \left. \nabla_{\mathbf{A}'} \right|_{\mathbf{A}'=\mathbf{A}^*} \ln P(\mathbf{x}_t|\mathbf{y}; \sigma_x^{2'}, \mathbf{A}') \quad (49)$$

$$= \sum_t \sum_z P(z|\mathbf{x}_t) \int d\mathbf{y} P(\mathbf{y}|z, \mathbf{x}_t) (\mathbf{x}_t - \mathbf{A}^* \mathbf{y}) \mathbf{y}^\top \quad (50)$$

$$= \sum_t \sum_z P(z|\mathbf{x}_t) \left[ \mathbf{x}_t \int d\mathbf{y} P(\mathbf{y}|z, \mathbf{x}_t) \mathbf{y}^\top - \mathbf{A}^* \int d\mathbf{y} P(\mathbf{y}|z, \mathbf{x}_t) \mathbf{y} \mathbf{y}^\top \right] \quad (51)$$

$$= \sum_t \sum_z P(z|\mathbf{x}_t) [\mathbf{x}_t \langle \mathbf{y}^\top | z, \mathbf{x}_t \rangle - \mathbf{A}^* (\langle \mathbf{y} | z, \mathbf{x}_t \rangle \langle \mathbf{y}^\top | z, \mathbf{x}_t \rangle + \operatorname{Cov}[\mathbf{y} | z, \mathbf{x}_t])] \quad (52)$$

$$= \sum_t \sum_z P(z|\mathbf{x}_t) [\mathbf{x}_t z \boldsymbol{\mu}^\top(z, \mathbf{x}_t) - \mathbf{A}^* z^2 (\boldsymbol{\mu}(z, \mathbf{x}_t) \boldsymbol{\mu}^\top(z, \mathbf{x}_t) + \boldsymbol{\Sigma}(z))] \quad (53)$$

$$= \boldsymbol{\Psi}_y - \mathbf{A}^* \boldsymbol{\Xi}_y \quad (54)$$

with

$$\boldsymbol{\Psi}_y = \sum_t \mathbf{x}_t \sum_z P(z|\mathbf{x}_t) z \boldsymbol{\mu}^\top(z, \mathbf{x}_t) \quad (55)$$

$$\boldsymbol{\Xi}_y = \sum_z z^2 \boldsymbol{\lambda}(z) \quad (56)$$

and  $\boldsymbol{\lambda}(z)$  as defined in Eq. 41

$$\mathbf{A}^* = \boldsymbol{\Psi}_y \boldsymbol{\Xi}_y^{-1} \quad (57)$$

### 4.3 M-step for $\sigma_x^2$

Preliminaries:

$$\ln P(\mathbf{x}|\mathbf{y}; \sigma_x^{2'}, \mathbf{A}) = -\frac{1}{2\sigma_x^2} (\mathbf{x} - \mathbf{A}\mathbf{y})^\top (\mathbf{x} - \mathbf{A}\mathbf{y}) - \frac{N}{2} \ln \sigma_x^2 + \dots \quad (58)$$

$$\frac{\partial}{\partial \sigma_x^2} \ln P(\mathbf{x}|\mathbf{y}; \sigma_x^2, \mathbf{A}) \propto \frac{1}{\sigma_x^4} (\mathbf{x} - \mathbf{A}\mathbf{y})^\top (\mathbf{x} - \mathbf{A}\mathbf{y}) - \frac{N}{\sigma_x^2} \quad (59)$$

with  $N = \dim(\mathbf{x})$ .

And so:

$$\sigma_x^{*2} = \operatorname{argmax}_{\sigma_x^{2'}} \sum_t \int d\mathbf{y} P(\mathbf{y}|\mathbf{x}_t) \ln P(\mathbf{x}_t|\mathbf{y}; \sigma_x^{2'}, \mathbf{A}') \quad (60)$$

$$\mathbf{0} = \left. \frac{\partial}{\partial \sigma_x^{2'}} \right|_{\sigma_x^{2'} = \sigma_x^{*2}} \sum_t \int d\mathbf{y} P(\mathbf{y}|\mathbf{x}_t) \ln P(\mathbf{x}_t|\mathbf{y}; \sigma_x^{2'}, \mathbf{A}') \quad (61)$$

$$= \sum_t \int d\mathbf{y} P(\mathbf{y}|\mathbf{x}_t) \left. \frac{\partial}{\partial \sigma_x^{2'}} \right|_{\sigma_x^{2'} = \sigma_x^{*2}} \ln P(\mathbf{x}_t|\mathbf{y}; \sigma_x^{2'}, \mathbf{A}') \quad (62)$$

$$= \sum_t \sum_z P(z|\mathbf{x}_t) \int d\mathbf{y} P(\mathbf{y}|z, \mathbf{x}_t) \left[ \frac{1}{\sigma_x^{*4}} (\mathbf{x}_t - \mathbf{A}\mathbf{y})^\top (\mathbf{x}_t - \mathbf{A}\mathbf{y}) - \frac{N}{\sigma_x^{*2}} \right] \quad (63)$$

$$= \frac{T N}{\sigma_x^{*2}} - \frac{1}{\sigma_x^{*4}} \sum_t \sum_z P(z|\mathbf{x}_t) \int d\mathbf{y} P(\mathbf{y}|z, \mathbf{x}_t) \left[ \mathbf{x}_t^\top \mathbf{x}_t + (\mathbf{A}\mathbf{y})^\top (\mathbf{A}\mathbf{y}) - 2 \mathbf{y}^\top \mathbf{A}^\top \mathbf{x}_t \right] \quad (64)$$

$$= \frac{T N}{\sigma_x^{*2}} - \frac{1}{\sigma_x^{*4}} \left[ \sum_t \mathbf{x}_t^\top \mathbf{x}_t + \sum_z P(z|\mathbf{x}_t) \int d\mathbf{y} P(\mathbf{y}|z, \mathbf{x}_t) \left[ \mathbf{y}^\top \mathbf{A}^\top \mathbf{A} \mathbf{y} - 2 \mathbf{y}^\top \mathbf{A}^\top \mathbf{x}_t \right] \right] \quad (65)$$

$$= \frac{T N}{\sigma_x^{*2}} - \frac{1}{\sigma_x^{*4}} \sigma_{\|x\|}^2 \quad (66)$$

with

$$\begin{aligned} \sigma_{\|x\|}^2 &= \sum_t \mathbf{x}_t^\top \mathbf{x}_t + \sum_z z^2 \operatorname{Tr}(\mathbf{A}^\top \mathbf{A} \boldsymbol{\Sigma}(z)) l(z) + \\ &+ \sum_t \sum_z P(z|\mathbf{x}_t) \left[ z^2 \boldsymbol{\mu}^\top(z, \mathbf{x}_t) \mathbf{A}^\top \mathbf{A} \boldsymbol{\mu}(z, \mathbf{x}_t) - 2 z \boldsymbol{\mu}^\top(z, \mathbf{x}_t) \mathbf{A}^\top \mathbf{x}_t \right] \end{aligned} \quad (67)$$

and  $l(z)$  as defined in Eq. 42

$$\sigma_x^{*2} = \frac{\sigma_{\|x\|}^2}{T N} \quad (68)$$

## References

- Wainwright MJ, Simoncelli EP (2000) Scale mixtures of Gaussians and the statistics of natural images. In: *Adv. Neur. Inf. Proc. Syst. 12*. MIT Press, pp 855–861.
- Bishop CM (2006) . Pattern recognition and machine learning. New York: Springer.
- Dempster A, Laird N, Rubin D, et al. (1977) Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society. Series B (Methodological)* 39:1–38.
- Neal RM, Hinton GE (1998) A view of the em algorithm that justifies incremental, sparse, and other variants. In: *Learning in graphical models* (Jordan MI, ed.), Dordrecht, The Netherlands: Kluwer, pp 355–368.