# Technical details of the identity/attribute model

Pietro Berkes, Richard E. Turner, and Maneesh Sahani

## 1 Notation

Indices $p, i$ and $q, j$ refer to identities and attributes and run over the interval $1, \ldots, d_b$ and $1, \ldots, d_a$, respectively. Indices $k$ and $l$ refer to input dimensions and run over $1, \ldots, d_y$. Capital letters stand for the set of all variables with the corresponding lowercase letter (e.g., $B = b_{1:T,1:d_c}$). For simplicity, we will omit hyperparameters from probability distributions and indicate random variables only. We define $\Theta$ to be the set of all parameters.

## 2 Generative model

The conditional dependencies in the model are defined by the directed graph in Figure 1.

### 2.1 Observations model

$$P(Y|B, A, \Theta) = \prod_t P(\mathbf{y}_t|\mathbf{b}_t, \mathbf{a}_t, W, \boldsymbol{\rho}) \tag{1}$$

$$= \prod_t \mathcal{N}_{\mathbf{y}_t}\left(\sum_p \mathbf{W}_p\, \mathbf{a}_{t,p}\, b_{t,p}, \boldsymbol{\Sigma}_y\right), \tag{2}$$

where $\boldsymbol{\Sigma}_y = \mathrm{diag}\left(\rho_k^{-1}\right)$ and

$$P(\rho_k|d, e) = \mathrm{Gamma}_{\rho_k}(d, e) = \frac{e^d}{\Gamma(d)}\rho_k^{d-1}\exp(-e\rho_k) \tag{3}$$

$$P(w_{kpq}|\gamma_{pq}) = \mathcal{N}_{w_{kpq}}(0, \gamma_{pq}^{-1}) \ . \tag{4}$$

### 2.2 Identity variables model

$$P(B|\mathbf{T}) = \prod_p \left(P(b_{1,p})\prod_{t>1}P(b_{t,p}|b_{t-1,p}, \mathbf{T})\right) \tag{5}$$

$$P(b_{1,p} = 1) = \pi_0 \tag{6}$$

$$P(b_{t,p} = \alpha|b_{t-1,p} = \beta) = T_{\beta,\alpha} \tag{7}$$

$$P(T) = \prod_{\beta=0}^{1} \mathrm{Dirichlet}\left(\{T_{\beta,0}, T_{\beta,1}\}\,|\,\mathbf{u}_\beta^{(T)}\right) \tag{8}$$
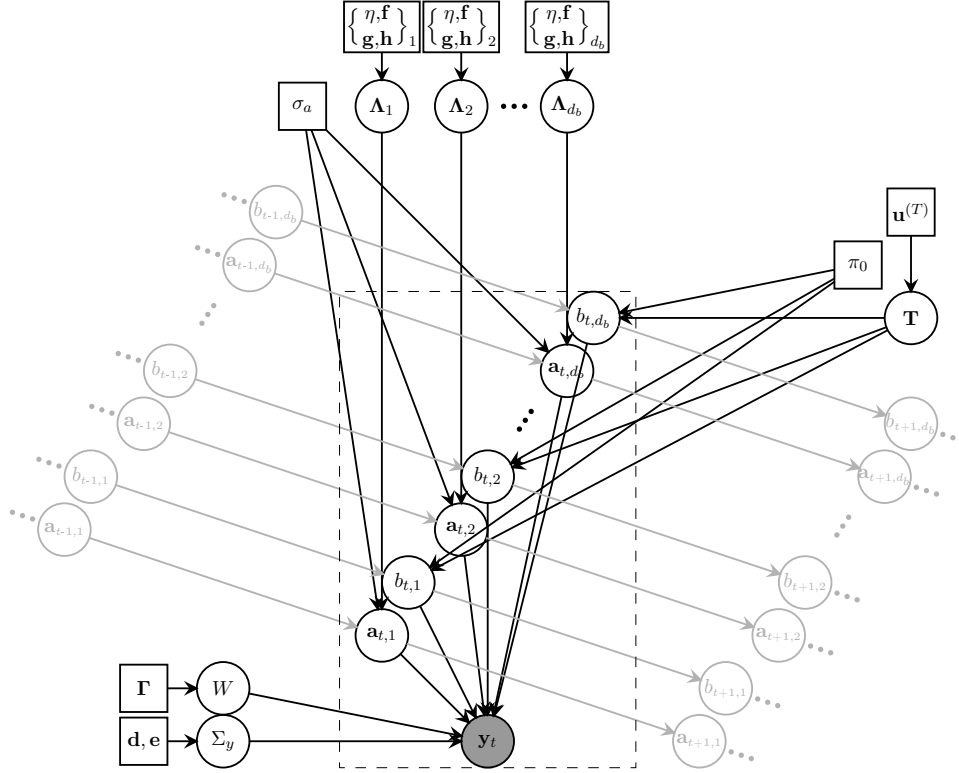
**Figure 1:** Directed graphical model representing the distribution of a single video frame. Circles represent random variables, and rectangles represent hyperparameters; the grey-filled circle represents the observed image; greyed symbols represent variables associated with neighbouring frames. The dashed box indicates that the variables within are replicated $T$ times (the length of an input sequence) in the complete model.

## 2.3 Attributes variables model

$$P(A|\mathbf{\Lambda}) = \prod_p \left( P(\mathbf{a}_{1,p}) \prod_{t>1} P(\mathbf{a}_{t,p}|\mathbf{a}_{t-1,p}, \mathbf{\Lambda}_p) \right) \tag{9}$$

$$P(\mathbf{a}_{1,p}) = \mathcal{N}_{\mathbf{a}_{1,p}}(\mathbf{0}, \sigma_a^2) \tag{10}$$

$$P(a_{t,pq}|a_{t-1,pq}, \lambda_{pq}) = \mathcal{N}_{a_{t,pq}}(\lambda_{pq}\, a_{t-1,pq}, 1 - \lambda_{pq}^2) \tag{11}$$

The prior over the dynamics parameters $\mathbf{\Lambda}$ is conjugate, but it has an unconventional form, because the mean and the standard deviation of $a_{t,pq}$ are coupled. However, it is still in the exponential family.

$$P(\mathbf{\Lambda}) = \prod_{p,q} P(\lambda_{pq}) \tag{12}$$

$$P(\lambda_{pq}) = \frac{1}{Z} \exp\left( -\eta_{pq} \log(1 - \lambda_{pq}^2) + \mathbf{\Phi}(\lambda_{pq})^T \mathbf{u}_{pq}^{(\lambda)} \right), \tag{13}$$

where,

$$\mathbf{u}_{pq}^{(\lambda)} := (f_{pq}, g_{pq}, h_{pq}) \tag{14}$$

$$\boldsymbol{\Phi}(\lambda) := \left((1 - \lambda^2)^{-1}, \lambda(1 - \lambda^2)^{-1}, \lambda^2(1 - \lambda^2)^{-1}\right)^T \tag{15}$$

are the hyperparameters and the sufficient statistics, respectively.

## 3   Variational approximation

In addition to the approximation required by VBEM between parameters and latent variables, another factorization is introduced to make the posterior tractable. This is between different identity variables at different times. Moreover, we factorize the distributions over the weights and the input noise to learn independently the signal and noise distribution. All other factorizations arise naturally.

$$Q(A, B, W, \mathbf{T}, \boldsymbol{\Lambda}, \boldsymbol{\rho}) \tag{16}$$

$$= Q(A, B)\, Q(W, \mathbf{T}, \boldsymbol{\Lambda}, \boldsymbol{\rho}) \tag{17}$$

$$= \prod_t \prod_p Q_{t,p}(b_{t,p}, \mathbf{a}_{t,p}) \prod_{pq} Q(\mathbf{w}_{pq})Q(\mathbf{T})Q(\boldsymbol{\Lambda})Q(\boldsymbol{\rho})\,. \tag{18}$$

An alternative factorization of the latent variables could be between identity and attribute variables, keeping the temporal correlations intact:

$$Q(A, B) = \prod_p Q(b_{1:T,p})Q(\mathbf{a}_{1:T,p})\,. \tag{19}$$

The distributions of $b_{1:T,a}$ and $\mathbf{a}_{1:T,a}$ can then be inferred using the forward-backward algorithm and Kalman smoothing, respectively. This approximation leads to smoother inferred signals. However, it is unclear which approximation is to be preferred for learning, or whether it is possible to combine the two approximations to obtain a more accurate estimation of the joint distribution. We experimented with initializing the distribution for the second approximation using the result of the first, obtaining encouraging results. For extensive runs, however, we used the first approximation alone because the computations are much faster.

## 4   VBE-Step

$$Q(A, B) = \prod_t \prod_p Q_{t,p}(b_{t,p}, \mathbf{a}_{t,p}) \tag{20}$$

From the VBEM theory [1] we know that the optimal approximation to the real

posterior is given by

$$Q_{t,p}(b_{t,p}, \mathbf{a}_{t,p}) \propto \exp\left\langle \log P(Y, B, A|\Theta) \right\rangle_{Q(\Theta)Q_{s,q \neq t,p}} \tag{21}$$

$$= \exp\left[\left\langle \log \prod_p P(b_{1,p}) \prod_{t>1} P(b_{t,p}|b_{t-1,p}, \mathbf{T}) \right\rangle_{Q(\mathbf{T})Q_{s,q\neq t,p}}\right] \tag{22}$$

$$\cdot \exp\left[\left\langle \log \prod_p P(\mathbf{a}_{1,p}) \prod_{t>1} P(\mathbf{a}_{t,p}|\mathbf{a}_{t-1,p}, \mathbf{\Lambda}_p) \right\rangle_{Q(\mathbf{\Lambda}_p)Q_{s,q\neq t,p}}\right]$$

$$\cdot \exp\left[\left\langle \log \prod_t P(\mathbf{y}_t|\mathbf{b}_t, \mathbf{a}_t, W, \boldsymbol{\rho}) \right\rangle_{Q(W)Q(\boldsymbol{\rho})Q_{s,q\neq t,p}}\right]$$

$$\propto Q_{t,p}^1(B) \ Q_{t,p}^2(A) \ Q_{t,p}^3(B, A) \tag{23}$$

Consider the first two terms, $Q_{t,p}^1(B)$ and $Q_{t,p}^2(A)$. For $1 < t < T$,

$$Q_{t,p}^1(b_{t,p}) = \frac{1}{Z^1} \exp\left[ \sum_{b_{t-1,p}} Q_{t-1,p}(b_{t-1,p}) \left\langle \log T_{b_{t-1,p}, b_{t,p}} \right\rangle \right.$$

$$\left. + \sum_{b_{t+1,p}} Q_{t+1,p}(b_{t+1,p}) \left\langle \log T_{b_{t,p}, b_{t+1,p}} \right\rangle \right]. \tag{24}$$

For $t = 1$

$$Q_{1,p}^1(b_{1,p}) = \frac{1}{Z^1} \exp\left[ \log P(b_{1,p}) + \sum_{b_{2,p}} Q_{2,p}(b_{2,p}) \left\langle \log T_{b_{1,p}, b_{2,p}} \right\rangle \right], \tag{25}$$

and for $t = T$

$$Q_{T,p}^1(b_{T,p}) = \frac{1}{Z^1} \exp\left[ \sum_{b_{T-1,p}} Q_{T-1,p}(b_{T-1,p}) \left\langle \log T_{b_{T-1,p}, b_{T,p}} \right\rangle \right]. \tag{26}$$

It is important to normalize this distribution (i.e., to compute $\frac{1}{Z^1}$), as we are going to see later.

For the second term we get

$$Q_{t,p}^2(\mathbf{a}_{t,p}) = \mathcal{N}_{\mathbf{a}_{t,p}}\left(\boldsymbol{\mu}_{t,p}^2, \boldsymbol{\Sigma}_{t,p}^2\right), \tag{27}$$

where for $1 < t < T$

$$\boldsymbol{\Sigma}_{t,p}^{2\ -1} = \left\langle \boldsymbol{\Sigma}_p^{-1} \right\rangle + \left\langle \boldsymbol{\Lambda}_p^T \boldsymbol{\Sigma}_p^{-1} \boldsymbol{\Lambda}_p \right\rangle \tag{28}$$

$$\boldsymbol{\mu}_{t,p}^{2\ T} \boldsymbol{\Sigma}_{t,p}^{2\ -1} = \left\langle \mathbf{a}_{t-1,p} \right\rangle^T \left\langle \boldsymbol{\Lambda}_p^T \boldsymbol{\Sigma}_p^{-1} \right\rangle + \left\langle \mathbf{a}_{t+1,p} \right\rangle^T \left\langle \boldsymbol{\Sigma}_p^{-1} \boldsymbol{\Lambda}_p^T \right\rangle, \tag{29}$$

and $\boldsymbol{\Lambda}_p = \mathrm{diag}\left(\lambda_{pq}\right)$ and $\boldsymbol{\Sigma}_p = \mathbf{I} - \boldsymbol{\Lambda}_p^2$. For $t = 1$ we have

$$\boldsymbol{\Sigma}_{1p}^{2\ -1} = \frac{1}{\sigma_a^2}\mathbf{I} + \left\langle \boldsymbol{\Lambda}_p^T \boldsymbol{\Sigma}_p^{-1} \boldsymbol{\Lambda}_p \right\rangle \tag{30}$$

$$\boldsymbol{\mu}_{1p}^{2\ T} \boldsymbol{\Sigma}_{1p}^{2\ -1} = \left\langle \mathbf{a}_{2p} \right\rangle^T \left\langle \boldsymbol{\Sigma}_p^{-1} \boldsymbol{\Lambda}_p^T \right\rangle, \tag{31}$$

and for $t = T$

$$\mathbf{\Sigma}_{T,p}^{2}{}^{-1} = \langle \mathbf{\Sigma}_p^{-1} \rangle \tag{32}$$

$$\boldsymbol{\mu}_{T,p}^{2}{}^{T} \mathbf{\Sigma}_{Tp}^{2}{}^{-1} = \langle \mathbf{a}_{T-1,p} \rangle^{T} \langle \mathbf{\Lambda}_p^{T} \mathbf{\Sigma}_p^{-1} \rangle \ . \tag{33}$$

Putting everything together we obtain

$$Q_{t,p}(b_{t,p}, \mathbf{a}_{t,p}) = \frac{1}{Z} Q_{t,p}^{1}(b_{t,p}) Q_{t,p}^{2}(a_{t,p})$$
$$\cdot \exp \Big\langle -\frac{d_y}{2} \log 2\pi - \frac{1}{2} \sum_k \log \rho_k^{-1}$$
$$- \frac{1}{2} \big( \mathbf{y}_t - \sum_i \mathbf{W}_i a_{t,i} b_{t,i} \big)^{T} \mathbf{\Sigma}_y^{-1} \big( \mathbf{y}_t - \sum_i \mathbf{W}_i a_{t,i} b_{t,i} \big) \Big\rangle \ . \tag{34}$$

For $b_{t,p} = 0$

$$Q_{t,p}(b_{t,p} = 0, \mathbf{a}_{t,p}) = \frac{1}{Z} Q_{t,p}^{1}(b_{t,p} = 0) Q_{t,p}^{2}(a_{t,p})$$
$$\cdot \exp \Big\langle -\frac{d_y}{2} \log 2\pi - \frac{1}{2} \sum_k \log \rho_k^{-1}$$
$$- \frac{1}{2} \big( \mathbf{y}_t - \sum_{i \neq p} \mathbf{W}_i a_{t,i} b_{t,i} \big)^{T} \mathbf{\Sigma}_y^{-1} \big( \mathbf{y}_t - \sum_{i \neq p} \mathbf{W}_i a_{t,i} b_{t,i} \big) \Big\rangle \tag{35}$$

$$=: \frac{1}{Z} Q_{t,p}^{1}(b_{t,p} = 0) Q_{t,p}^{2}(\mathbf{a}_{t,p}) \cdot C \ . \tag{36}$$

Since the last term in Eq. 35 does not depend on $b_{t,p}$ or $\mathbf{a}_{t,p}$, we abbreviated it as a constant $C$. For $b_{t,p} = 1$ we get the same expression, plus an additional term

$$Q_{t,p}(b_{t,p} = 1, \mathbf{a}_{t,p}) = \frac{1}{Z} C \cdot Q_{t,p}^{1}(b_{t,p} = 1) Q_{t,p}^{2}(\mathbf{a}_{t,p})$$
$$\cdot \exp \Big\{ -\frac{1}{2} \Big[ -2 \Big( \mathbf{y}_t - \sum_{i \neq p} \langle \mathbf{W}_i \rangle \langle \mathbf{a}_{t,i} b_{t,i} \rangle \Big)^{T} \langle \mathbf{\Sigma}_y^{-1} \rangle \langle \mathbf{W}_p \rangle \mathbf{a}_{t,p}$$
$$+ \mathbf{a}_{t,p}^{T} \langle \mathbf{W}_p^{T} \mathbf{\Sigma}_y^{-1} \mathbf{W}_p \rangle \mathbf{a}_{t,p} \Big] \Big\} \ . \tag{37}$$

Expanding $Q_{t,p}^{2}(\mathbf{a}_{t,p})$ as in Eq. 27 and collecting the terms that contain $\mathbf{a}_{t,p}$

$$Q_{t,p}(b_{t,p} = 1, \mathbf{a}_{t,p}) = \frac{1}{Z} C C' \cdot Q_{t,p}^{1}(b_{t,p} = 1) \, \mathcal{N}_{\mathbf{a}_{t,p}}(\boldsymbol{\mu}_{t,p}, \mathbf{\Sigma}_{t,p}) \tag{38}$$

where

$$\mathbf{\Sigma}_{t,p}^{-1} = \mathbf{\Sigma}_{t,p}^{2}{}^{-1} + \langle \mathbf{W}_p^{T} \mathbf{\Sigma}_y^{-1} \mathbf{W}_p \rangle \tag{39}$$

$$\boldsymbol{\mu}_{t,p}^{T} \mathbf{\Sigma}_{t,p}^{-1} = \boldsymbol{\mu}_{t,p}^{2}{}^{T} \mathbf{\Sigma}_{t,p}^{2}{}^{-1} + \Big( \mathbf{y}_t - \sum_{i \neq p} \langle \mathbf{W}_i \rangle \langle \mathbf{a}_{t,i} b_{t,i} \rangle \Big)^{T} \langle \mathbf{\Sigma}_y^{-1} \rangle \langle \mathbf{W}_p \rangle \tag{40}$$

$$C' = \exp \Big[ -\frac{1}{2} \log |\mathbf{\Sigma}_{t,p}^{2}| + \frac{1}{2} \log |\mathbf{\Sigma}_{t,p}|$$
$$- \frac{1}{2} \boldsymbol{\mu}_{t,p}^{2}{}^{T} \mathbf{\Sigma}_{t,p}^{2}{}^{-1} \boldsymbol{\mu}_{t,p}^{2} + \frac{1}{2} \boldsymbol{\mu}_{t,p}^{T} \mathbf{\Sigma}_{t,p}^{-1} \boldsymbol{\mu}_{t,p} \Big] \tag{41}$$

We can obtain the normalization constant using the identity

$$1 = \int Q_{t,p}(b_{t,p} = 0, \mathbf{a}_{t,p}) \, \mathrm{d}\mathbf{a}_{t,p} + \int Q_{t,p}(b_{t,p} = 1, \mathbf{a}_{t,p}) \, \mathrm{d}\mathbf{a}_{t,p} \tag{42}$$

$$= \frac{1}{Z} C \cdot Q_{t,p}^1(b_{t,p} = 0) \int Q_{t,p}^2(\mathbf{a}_{t,p}) \, \mathrm{d}\mathbf{a}_{t,p}$$

$$+ \frac{1}{Z} C C' \cdot Q_{t,p}^1(b_{t,p} = 1) \int \mathcal{N}_{\mathbf{a}_{t,p}}(\boldsymbol{\mu}_{t,p}, \boldsymbol{\Sigma}_{t,p}) \, \mathrm{d}\mathbf{a}_{t,p} \tag{43}$$

$$= \frac{1}{Z} C \cdot \left( Q_{t,p}^1(b_{t,p} = 0) + C' \, Q_{t,p}^1(b_{t,p} = 1) \right), \tag{44}$$

and thus

$$\frac{1}{Z} C = \left( Q_{t,p}^1(b_{t,p} = 0) + C' \, Q_{t,p}^1(b_{t,p} = 1) \right)^{-1}. \tag{45}$$

Notice that $C$ and $C'$ depend on $t$ and $p$, although the indices have been dropped to avoid clutter.

The distribution $Q(B, A)$ can thus be obtained by first computing and normalizing $Q_{t,p}^1(b_{t,p})$ (Eq. 24–26); computing $\boldsymbol{\mu}_{t,p}^2$ and $\boldsymbol{\Sigma}_{t,p}^2$ (Eq. 28–33); computing $\boldsymbol{\mu}_{t,p}$, $\boldsymbol{\Sigma}_{t,p}$, and $C'$ (Eq. 39–41); computing the normalization constant (Eq. 45); and finally computing the sufficient statistics as follows:

$$Q_{t,p}(b_{t,p} = 1) = \frac{1}{Z} C C' \, Q_{t,p}^1(b_{t,p} = 1) \tag{46}$$

$$\langle \mathbf{a}_{t,p} \rangle_{Q_{t,p}} = Q_{t,p}(b_{t,p} = 1) \, \boldsymbol{\mu}_{t,p} + Q_{t,p}(b_{t,p} = 0) \, \boldsymbol{\mu}_{t,p}^2 \tag{47}$$

$$\langle \mathbf{a}_{t,p} b_{t,p} \rangle_{Q_{t,p}} = Q_{t,p}(b_{t,p} = 1) \, \boldsymbol{\mu}_{t,p} \tag{48}$$

$$\left\langle (b_{t,p} \mathbf{a}_{t,p})(b_{t,p} \mathbf{a}_{t,p})^T \right\rangle_{Q_{t,p}} = \left( \boldsymbol{\Sigma}_{t,p} + \boldsymbol{\mu}_{t,p} \boldsymbol{\mu}_{t,p}^T \right) Q_{t,p}(b_{t,p} = 1) \tag{49}$$

$$i \neq p : \quad \left\langle \mathbf{a}_{t,i} \, \mathbf{a}_{t,p}^{\,T} \right\rangle_{Q_{t,i} Q_{t,p}} = \langle \mathbf{a}_{t,i} \rangle_{Q_{t,i}} \langle \mathbf{a}_{t,p} \rangle_{Q_{t,p}}^T \tag{50}$$

$$\left\langle \mathbf{a}_{t,i} \, \mathbf{a}_{t,i}^{\,T} \right\rangle_{Q_{t,i}} = Q_{t,p}(b_{t,p} = 0) \left( \boldsymbol{\Sigma}_{t,p}^2 + \boldsymbol{\mu}_{t,p}^2 \boldsymbol{\mu}_{t,p}^{2\,T} \right)$$

$$+ Q_{t,p}(b_{t,p} = 1) \left( \boldsymbol{\Sigma}_{t,p} + \boldsymbol{\mu}_{t,p} \boldsymbol{\mu}_{t,p}^T \right). \tag{51}$$

# 5 VBM-Step

In the VBM-Step, functional maximization of the free energy for parameters $\Theta$ gives

$$Q(\Theta) = \frac{1}{Z} P(\Theta) \exp \left\langle \log P(Y, B, A | \Theta) \right\rangle_{Q(B,A)}. \tag{52}$$

## 5.1 Transition matrix T

We start with the distribution of the transition matrix $\mathbf{T}$

$$Q(\mathbf{T}) = \frac{1}{Z} P(\mathbf{T}) \cdot \exp \left\langle \log \prod_i \left( P(b_{1,i}) \prod_{t>1} P(b_{t,i} | b_{t-1,i}, \mathbf{T}) \right) \right\rangle \tag{53}$$

$$= \frac{1}{Z} \prod_{m=0}^{1} \mathrm{Dirichlet}(\{T_{m0}, T_{m1}\} | \mathbf{u}_m^{(T)}) \cdot \exp \left\langle \sum_i \sum_{t>1} \log T_{b_{t-1,i}, b_{t,i}} \right\rangle, \tag{54}$$

and since

$$\left\langle \log T_{b_{t-1,i},b_{t,i}} \right\rangle_{Q(B,A)} = \sum_{m,n=0}^{1} \left\langle \delta(b_{t-1,i}=m) \right\rangle \left\langle \delta(b_{t,i}=n) \right\rangle \log T_{mn} \qquad (55)$$

$$= \sum_{m,n=0}^{1} Q(b_{t-1,i}=m)Q(b_{t,i}=n) \log T_{mn} \qquad (56)$$

we obtain

$$Q(T) = \frac{1}{Z} \prod_{m=0}^{1} \frac{1}{\text{Beta}(\mathbf{u}_m^{(T)})} \prod_{n=0}^{1} T_{mn}^{\left[ u_{mn}^{(T)} - 1 + \sum_i \sum_{t>1} Q(b_{t-1,i}=m)Q(b_{t,i}=n) \right]} . \qquad (57)$$

Comparing Eq. 57 with the prior Dirichlet distribution over $T$, we see that the approximate posterior is also Dirichlet with parameters

$$Q(T) = \prod_{m=0}^{1} \text{Dirichlet}\left( \left\{ T_{m0}, T_{m1} \right\} \,|\, \tilde{\mathbf{u}}_m^{(T)} \right) \qquad (58)$$

$$\tilde{u}_{mn}^{(T)} = u_{mn}^{(T)} + \sum_i \sum_{t>1} Q(b_{t-1,i}=m)Q(b_{t,i}=n) . \qquad (59)$$

## 5.2   Attributes dynamics $\mathbf{\Lambda}$

$$Q(\lambda_{pq}) = \frac{1}{Z} P(\lambda_{pq}) \exp \left\langle \sum_{t>1} \log P(\mathbf{a}_{t,p}|\mathbf{a}_{t-1,p}, \mathbf{\Lambda}_p) \right\rangle_{Q(B,A)} \qquad (60)$$

$$= \frac{1}{Z} P(\lambda_{pq}) \exp \sum_{t>1} \left( -\frac{1}{2} \log(1-\lambda_{pq}^2) - \frac{1}{2(1-\lambda_{pq}^2)} \left\langle (\mathbf{a}_{t,pq} - \lambda_{pq}\mathbf{a}_{t-1,pq})^2 \right\rangle \right) \qquad (61)$$

$$= \frac{1}{Z} \exp \left[ -\eta_{ab} \log(1-\lambda_{pq}^2) + \frac{1}{1-\lambda_{pq}^2}(f_{pq} + \lambda_{pq}g_{pq} + \lambda_{pq}^2 h_{pq}) \right]$$

$$\exp \left[ -\frac{T-1}{2} \log(1-\lambda_{pq}^2) \right.$$

$$\left. -\frac{1}{1-\lambda_{pq}^2} \left( \frac{1}{2} \sum_{t>1} \left\langle a_{t,pq}^2 \right\rangle - \lambda_{pq} \sum_{t>1} \left\langle a_{t-1,pq} \right\rangle \left\langle a_{t,pq} \right\rangle + \frac{1}{2} \sum_{t>1} \lambda_{pq}^2 \left\langle a_{t-1,pq}^2 \right\rangle \right) \right] . \qquad (62)$$

The approximate posterior has the same functional form as the prior, with parameters

$$\tilde{\eta}_{pq} = \eta_{pq} + \frac{T-1}{2} \qquad (63)$$

$$\tilde{f}_{pq} = f_{pq} - \frac{1}{2} \sum_{t>1} \left\langle a_{t,pq}^2 \right\rangle \qquad (64)$$

$$\tilde{g}_{pq} = g_{pq} + \sum_{t>1} \left\langle a_{t-1,pq} \right\rangle \left\langle a_{t,pq} \right\rangle \qquad (65)$$

$$\tilde{h}_{pq} = h_{pq} - \frac{1}{2} \sum_{t>1} \left\langle a_{t-1,pq}^2 \right\rangle . \qquad (66)$$

7

## 5.3   Basis vectors and precision of observation noise $W, \rho_k$

$$Q(w_{kpq}) = \frac{1}{Z} P(w_{kpq}) \exp \left[ \sum_t \left( -\frac{1}{2} \sum_l \langle \log 2\pi \rho_l^{-1} \rangle \right. \right.$$

$$\left. \left. -\frac{1}{2} \sum_l \langle \rho_l \rangle \left\langle (y_{tl} - \sum_{ij} w_{lij} a_{t,ij} b_{t,i})^2 \right\rangle \right) \right] . \quad (67)$$

In the following, we will write $x_{t,ij}$ for the product $a_{t,ij} b_{t,i}$. Eliminating all terms that do not contain $w_{kpq}$, the exponent becomes

$$-\frac{T}{2} \langle \log 2\pi \rho_k^{-1} \rangle - \frac{1}{2} \langle \rho_k \rangle \sum_t \left( 2 \sum_{ij \neq pq} \langle x_{t,ij} x_{t,pq} \rangle \langle w_{kij} \rangle w_{kpq} \right.$$

$$\left. + \langle x_{t,pq}^2 \rangle w_{t,pq}^2 - 2 y_{tk} \langle x_{t,pq} \rangle w_{kpq} \right) \quad (68)$$

$$= -\frac{T}{2} \langle \log 2\pi \rho_k^{-1} \rangle - \frac{1}{2} \left[ \langle \rho_k \rangle \left( 2 \sum_{ij \neq pq} \sum_t \langle x_{t,ij} x_{t,pq} \rangle \langle w_{kij} \rangle - 2 \sum_t y_{tk} \langle x_{t,pq} \rangle \right) w_{kpq} \right.$$

$$\left. + \langle \rho_k \rangle \sum_t \langle x_{t,pq}^2 \rangle w_{kpq}^2 \right] . \quad (69)$$

Defining $R_{ijpq} := \sum_t \langle x_{t,ij} x_{t,pq} \rangle$ and $r_{kpq} := \sum_t y_{tk} \langle x_{t,pq} \rangle$, and because $P(w_{kpq}) = \mathcal{N}_{w_{kpq}}(0, \gamma_{pq}^{-1})$, the posterior reduces to a normal distribution

$$Q(w_{kpq}) = \mathcal{N}_{w_{kpq}}\left( \mu_{kpq}^{(w)}, \sigma_{kpq}^{(w)2} \right) \quad (70)$$

$$\sigma_{kpq}^{(w)2} = (\gamma_{pq} + \langle \rho_k \rangle R_{pqpq})^{-1} \quad (71)$$

$$\mu_{kpq}^{(w)} = \sigma_{kpq}^{(w)2} \langle \rho_k \rangle \left( r_{kpq} - \sum_{ij \neq pq} R_{ijpq} \langle w_{kij} \rangle \right) . \quad (72)$$

As for the distribution of the precision of the observations,

$$Q(\rho_k) = \frac{1}{Z} P(\rho_k) \exp \left\langle \sum_t \left( -\frac{1}{2} \sum_l \log 2\pi \rho_l^{-1} - \frac{1}{2} \sum_l \rho_l (y_{tl} - \sum_{ij} w_{lij} x_{t,ij})^2) \right) \right\rangle \quad (73)$$

$$= \frac{1}{Z} \exp \left[ (d_k - 1) \log \rho_k - e_k \rho_k \right.$$

$$+ \frac{T}{2} \log \rho_k - \frac{1}{2} \rho_k \left( \sum_t y_{tk}^2 - 2 \sum_{ij} \langle w_{kij} \rangle \sum_t y_{tk} \langle x_{t,ij} \rangle \right.$$

$$\left. \left. \textcolor{red}{+ \sum_{ijmn} \langle w_{kij} w_{kmn} \rangle \sum_t \langle x_{t,ij} x_{tmn} \rangle} \right) \right] . \quad (74)$$

The term in red expands as

$$\sum_{ijmn} \mu_{kij}^{(w)} \mu_{kmn}^{(w)} R_{ijmn} + \sum_{ij} \rho_k^{-1} (\gamma_{ij} + R_{ijij})^{-1} R_{ijij} . \quad (75)$$

8

When this expression is multiplied by $\rho_k$, the second term does not depend on $\rho_k$ and can be discarded. The posterior is thus a Gamma distribution with parameters

$$Q(\rho_k) = \text{Gamma}(d'_k, e'_k) \tag{76}$$

$$d'_k = d_k + \frac{T}{2} \tag{77}$$

$$e'_k = e_k + \frac{1}{2}\Big( \sum_t y_{tk}^2 - 2 \sum_{ij} \mu_{kij}^{(w)} r_{kij} + \sum_{ijmn} \mu_{kij}^{(w)} \mu_{kmn}^{(w)} R_{ijmn} \Big) \ . \tag{78}$$

# 6 Learning the ARD precision parameter

After an initial learning phase, we start learning the precision hyperparameter $\gamma_{pq}$ by maximizing the free energy,

$$\frac{\partial}{\partial \gamma_{pq}} \mathcal{F}(\gamma, Q(W), Q(\boldsymbol{\rho}), Q(\boldsymbol{\Lambda}), Q(\mathbf{T}), Q(B, A)) = \frac{\partial}{\partial \gamma_{pq}} \langle \log P(\mathbf{W}) \rangle_{Q(\mathbf{W})} \tag{79}$$

$$= \frac{\partial}{\partial \gamma_{pq}} \Big( \frac{d_y}{2} \log \gamma_{pq} - \frac{1}{2} \sum_k \gamma_{pq} \langle w_{kpq}^2 \rangle \Big) \tag{80}$$

$$= \frac{d_y}{2} \gamma_{pq}^{-1} - \frac{1}{2} \sum_k \langle w_{kpq}^2 \rangle \ . \tag{81}$$

Setting the derivative to zero we obtain

$$\gamma_{pq}^{-1} = \frac{1}{d_y} \sum_k \langle w_{kpq}^2 \rangle \tag{82}$$

$$= \frac{1}{d_y} \sum_k \Big( \mu_{kij}^{(w)2} + \sigma_{kpq}^{(w)2} \Big) \ . \tag{83}$$

# 7 Free energy

$$\mathcal{F}(Q(W), Q(\boldsymbol{\rho}), Q(\boldsymbol{\Lambda}), Q(\mathbf{T}), Q(B, A), \Theta) \tag{84}$$

$$= \Big\langle \log \frac{P(Y, B, A, W, \boldsymbol{\rho}, \mathbf{T}, \boldsymbol{\Lambda} | \Theta)}{Q(B, A) Q(W) Q(\boldsymbol{\rho}) Q(\boldsymbol{\Lambda}) Q(\mathbf{T})} \Big\rangle_{Q(B, A) Q(W) Q(\boldsymbol{\rho}) Q(\boldsymbol{\Lambda}) Q(\mathbf{T})} \tag{85}$$

$$= -\big\langle \log Q(B, A) \big\rangle_{Q(B, A)}$$

$$\quad - \big\langle \log Q(W) \big\rangle_{Q(W)} - \big\langle \log Q(\boldsymbol{\rho}) \big\rangle_{Q(\boldsymbol{\rho})}$$

$$\quad - \big\langle \log Q(\mathbf{T}) \big\rangle_{Q(\mathbf{T})} - \big\langle \log Q(\boldsymbol{\Lambda}) \big\rangle_{Q(\boldsymbol{\Lambda})}$$

$$\quad + \big\langle \log(P(B|\mathbf{T}) \big\rangle_{Q(B, A) Q(\mathbf{T})} + \big\langle \log(P(A|\boldsymbol{\Lambda}) \big\rangle_{Q(B, A) Q(\boldsymbol{\Lambda})}$$

$$\quad + \big\langle \log P(W) \big\rangle_{Q(W)} + \big\langle \log P(\boldsymbol{\rho}) \big\rangle_{Q(\boldsymbol{\rho})}$$

$$\quad + \big\langle \log P(\mathbf{T}) \big\rangle_{Q(\mathbf{T})} + \big\langle \log P(\boldsymbol{\Lambda}) \big\rangle_{Q(\boldsymbol{\Lambda})}$$

$$\quad + \big\langle \log P(Y|B, A, W, \boldsymbol{\rho}) \big\rangle_{Q(C, S) Q(W, \boldsymbol{\rho})} \ . \tag{86}$$

The individual terms:

$$\left\langle \log Q(B, A) \right\rangle_{Q(B,A)} \tag{87}$$

$$= \sum_{p,t} \left\langle \log Q(b_{t,p}, \mathbf{a}_{t,p}) \right\rangle_{Q(b_{t,p}, \mathbf{a}_{t,p})} \tag{88}$$

$$= \sum_{p,t} \left[ \int \mathrm{d}\mathbf{a}_{t,p} \, Q(b_{t,p} = 1, \mathbf{a}_{t,p}) \log Q(b_{t,p} = 1, \mathbf{a}_{t,p}) \right.$$

$$\left. + \int \mathrm{d}\mathbf{a}_{t,p} \, Q(b_{t,p} = 0, \mathbf{a}_{t,p}) \log Q(b_{t,p} = 0, \mathbf{a}_{t,p}) \right] \tag{89}$$

Since

$$Q(b_{t,p} = 1, \mathbf{a}_{t,p}) = Q(b_{t,p} = 1) \, \mathcal{N}_{\mathbf{a}_{t,p}}(\boldsymbol{\mu}_{t,p}, \boldsymbol{\Sigma}_{t,p}) \tag{90}$$

$$Q(b_{t,p} = 0, \mathbf{a}_{t,p}) = Q(b_{t,p} = 0) \, \mathcal{N}_{\mathbf{a}_{t,p}}(\boldsymbol{\mu}_{t,p}^2, \boldsymbol{\Sigma}_{t,p}^2) \tag{91}$$

we get

$$\left\langle \log Q(B, A) \right\rangle_{Q(B,A)} \tag{92}$$

$$= \sum_{p,t} \left[ Q(b_{t,p} = 1) \Big[ \log Q(b_{t,p} = 1) - \mathrm{H}\left( \mathcal{N}_{\mathbf{a}_{t,p}}(\boldsymbol{\mu}_{t,p}, \boldsymbol{\Sigma}_{t,p}) \right) \Big] \right.$$

$$\left. + Q(b_{t,p} = 0) \Big[ \log Q(b_{t,p} = 0) - \mathrm{H}\left( \mathcal{N}_{\mathbf{a}_{t,p}}(\boldsymbol{\mu}_{t,p}^2, \boldsymbol{\Sigma}_{t,p}^2) \right) \Big] \right] \tag{93}$$

$$= \sum_{p,t} \left[ Q(b_{t,p} = 1) \Big[ \log Q(b_{t,p} = 1) - \frac{1}{2} \log |2\pi \boldsymbol{\Sigma}_{t,p}| - \frac{d_a}{2} \Big] \right.$$

$$\left. + Q(b_{t,p} = 0) \Big[ \log Q(b_{t,p} = 0) - \frac{1}{2} \log |2\pi \boldsymbol{\Sigma}_{t,p}^2| - \frac{d_a}{2} \Big] \right] \tag{94}$$

$$\left\langle \log Q(W) \right\rangle_{Q(W)} = \sum_{pq} \left\langle \log Q(\mathbf{w}_{pq}) \right\rangle_{Q(\mathbf{w}_{pq})} \tag{95}$$

$$= \sum_{pq} \left( -\frac{d_y}{2} - \frac{d_y}{2} \log(2\pi \sigma_{kpq}^2) \right) \tag{96}$$

$$\left\langle \log Q(\mathbf{T}) \right\rangle_{Q(\mathbf{T})} = \sum_{n=0}^{1} \left\langle \log Q(T_{n:}) \right\rangle_{Q(T_{n:})} \tag{97}$$

$$= \sum_{n=0}^{1} \left[ -\log \mathrm{Beta}(\tilde{\mathbf{u}}_n^{(T)}) + (\tilde{u}_{n0}^{(T)} - 1) \langle \log T_{n0} \rangle + (\tilde{u}_{n1}^{(T)} - 1) \langle \log T_{n1} \rangle \right] \tag{98}$$

$$\big\langle \log Q(\boldsymbol{\Lambda}) \big\rangle_{Q(\boldsymbol{\Lambda})} = \sum_{pq} \big\langle \log Q(\lambda_{pq}) \big\rangle \tag{99}$$

$$= \sum_{pq} \Bigg[ -\log \tilde{Z}_{pq} - \tilde{\eta}_{pq} \big\langle \log(1 - \lambda_{pq}^2) \big\rangle$$

$$+ \tilde{f}_{pq} \big\langle (1 - \lambda_{pq}^2)^{-1} \big\rangle + \tilde{g}_{pq} \big\langle \lambda_{pq}(1 - \lambda_{pq}^2)^{-1} \big\rangle + \tilde{h}_{pq} \big\langle \lambda_{pq}^2(1 - \lambda_{pq}^2)^{-1} \big\rangle \Bigg] \tag{100}$$

$$\big\langle \log P(B|\mathbf{T}) \big\rangle = \sum_i \big\langle \log P(b_{1,i}) \big\rangle + \sum_{t>1} \big\langle \log P(b_{t,i}|b_{t-1,i}, \mathbf{T}) \big\rangle \tag{101}$$

$$= \sum_i \Bigg[ \sum_{b_{1,i}} Q(b_{1,i}) \log P(b_{1,i}) + \sum_{t>1} \sum_{b_{t,i},b_{t-1,i}} Q(b_{t-1,i}, b_{t,i}) \big\langle \log T_{b_{t-1,i}b_{t,i}} \big\rangle \Bigg] \tag{102}$$

$$\big\langle \log P(A|\Lambda) \big\rangle = \sum_i \Bigg[ \big\langle \log P(\mathbf{a}_{1,i}) \big\rangle + \sum_{t>1} \big\langle \log P(\mathbf{a}_{t,i}|\mathbf{a}_{t-1,i}, \boldsymbol{\Lambda}) \big\rangle \Bigg] \tag{103}$$

$$= \sum_i \Bigg[ -\frac{d_a}{2} \log 2\pi\sigma_a^2 - \frac{1}{2\sigma_a^2} \big\langle \mathbf{a}_{1,i}^T \mathbf{a}_{1,i} \big\rangle_{Q(b_{1,i}, \mathbf{a}_{1,i})} \Bigg]$$

$$+ \sum_i \sum_{t>1} \Bigg[ -\frac{d_a}{2} \log 2\pi - \frac{1}{2} \sum_j \big\langle \log(1 - \lambda_{ij}^2) \big\rangle$$

$$- \frac{1}{2} \sum_j \Big( \big\langle (1 - \lambda_{ij}^2)^{-1} \big\rangle \big\langle a_{t,ij}^2 \big\rangle$$

$$- 2 \big\langle \lambda_{ij}(1 - \lambda_{ij}^2)^{-1} \big\rangle \big\langle a_{t,ij}\, a_{t-1,ij} \big\rangle$$

$$+ \big\langle \lambda_{ij}^2(1 - \lambda_{ij}^2)^{-1} \big\rangle \big\langle a_{t-1,ij}^2 \big\rangle \Big) \Bigg] \tag{104}$$

$$\big\langle \log P(W) \big\rangle_{Q(W)} = \sum_{pq} \big\langle \log P(\mathbf{w}_{pq}) \big\rangle_{Q(\mathbf{w}_{pq})} \tag{105}$$

$$= \sum_{pq} \Big( -\frac{d_y}{2} \log(2\pi\gamma_{pq}^{-1}) - \frac{1}{2}\gamma_{pq} \big\langle \mathbf{w}_{pq}^T \mathbf{w}_{pq} \big\rangle \Big) \tag{106}$$

$$\big\langle \log P(\mathbf{T}) \big\rangle_{Q(\mathbf{T})} = \sum_i \big\langle \log P(T_{i:}) \big\rangle_{Q(T_{i:})} \tag{107}$$

$$= \sum_i \Big[ -\log \mathrm{Beta}(\mathbf{u}_i^{(T)}) + (u_{i0}^{(T)} - 1)\big\langle \log T_{i0} \big\rangle + (u_{i1}^{(T)} - 1)\big\langle \log T_{i1} \big\rangle \Big] \tag{108}$$

$$\big\langle \log P(\mathbf{\Lambda}) \big\rangle_{Q(\mathbf{\Lambda})} = \sum_{ij} \big\langle \log P(\lambda_{ij}) \big\rangle \tag{109}$$

$$= \sum_{ij} \left[ -\log Z_{ij} - \eta_{ij} \big\langle \log(1 - \lambda_{ij}^2) \big\rangle \right.$$

$$\left. + f_{ij} \big\langle (1 - \lambda_{ij}^2)^{-1} \big\rangle + g_{ij} \big\langle \lambda_{ij}(1 - \lambda_{ij}^2)^{-1} \big\rangle + h_{ij} \big\langle \lambda_{ij}^2(1 - \lambda_{ij}^2)^{-1} \big\rangle \right] \tag{110}$$

$$-\big\langle \log Q(\boldsymbol{\rho}) \big\rangle_{Q(\boldsymbol{\rho})} + \big\langle \log P(\boldsymbol{\rho}) \big\rangle_{Q(\boldsymbol{\rho})} = -KL(Q(\boldsymbol{\rho})\|P(\boldsymbol{\rho})) \tag{111}$$

$$= -\sum_k KL(Q(\rho_k)\|P(\rho_k)) \tag{112}$$

$$= -\sum_w \left[ d_k' \log e_k' - d_k \log e_k - \log \frac{\Gamma(d_k')}{\Gamma(d_k)} \right.$$

$$\left. + (d_k' - d_k)(\Psi(d_k') - \log e_k') - d_k'\Big(1 - \frac{e_k}{e_k'}\Big) \right], \tag{113}$$

where $\Psi(\cdot)$ is the Digamma function.

$$\big\langle \log P(Y|B, A, W, \boldsymbol{\rho}) \big\rangle_{Q(B,A)Q(W)Q(\boldsymbol{\rho})} \tag{114}$$

$$= \sum_t \left[ -\frac{d_y}{2}\log 2\pi + \frac{1}{2}\sum_k \langle \log \rho_k \rangle - \frac{1}{2}\Big\langle \mathbf{y}_t^T \mathbf{\Sigma}_y^{-1} \mathbf{y}_t - 2\mathbf{y}_t^T \mathbf{\Sigma}_y^{-1} \sum_i \mathbf{W}_i \mathbf{a}_{t,i} b_{t,i} \right.$$

$$\left. + \sum_{i,p} (\mathbf{a}_{t,i} b_{t,i})^T \mathbf{W}_i^T \mathbf{\Sigma}_y^{-1} \mathbf{W}_p (\mathbf{a}_{t,p} b_{t,p}) \Big\rangle \right] \tag{115}$$

$$= -\frac{T d_y}{2}\log 2\pi + \frac{T}{2}\sum_k \langle \log \rho_k \rangle$$

$$- \frac{1}{2}\left[ \sum_t \mathbf{y}_t^T \big\langle \mathbf{\Sigma}_y^{-1} \big\rangle \mathbf{y}_t - 2\sum_t \mathbf{y}_t^T \big\langle \mathbf{\Sigma}_y^{-1} \big\rangle \sum_i \langle \mathbf{W}_i \rangle \langle \mathbf{a}_{t,i} b_{t,i} \rangle \right.$$

$$\left. + \sum_{i,p} \operatorname{trace}\Big( \big\langle \mathbf{W}_i^T \mathbf{\Sigma}_y^{-1} \mathbf{W}_p \big\rangle \sum_t \big\langle (\mathbf{a}_{t,p} b_{t,p})^T (\mathbf{a}_{t,i} b_{t,i}) \big\rangle \Big) \right] \tag{116}$$

# References

[1] Beal M (2003) Variational Algorithms for Approximate Bayesian Inference. Ph.D. thesis, Gatsby Computational Neuroscience Unit, University College London.