

# Parameter estimation for the gestalt model

## 1 Generative model

A gestalt, a perceptual object, is characterised by a covariance component for the joint distribution of visual neural activity.

$$p(v \mid g) = \mathcal{N}(v; 0, C_v) \quad (1)$$

$$C_v = \sum_{k=1}^K g_k C_k \quad (2)$$

where  $K$  is the fixed number of possible gestalts in the visual scene and  $g_k$  is the strength of the gestalt number  $k$ , coming from a  $K$ -dimensional Gamma prior distribution with shape and scale parameters  $\alpha_g$  and  $\theta_g$  controlling the sparsity of the prior.

$$p(g) = \text{Gam}(g; \alpha_g, \theta_g) \quad (3)$$

The global contrast of the image patch is encoded by a scalar variable  $z$ , also coming from a Gamma prior

$$p(z) = \text{Gam}(z; \alpha_z, \theta_z) \quad (4)$$

The pixel intensities are generated from the neural activity through a set of linear projective field models, possibly Gabor filters,  $A$ , scaled by the contrast and adding some independent observational noise.

$$p(x \mid v, z) = \mathcal{N}(x; zAv, C_x) \quad (5)$$

$$C_x = \sigma_x I; \quad (6)$$

We might assume that a single composition of gestalts, characterised by the  $g$  vector, generates a batch of  $B$  independent images, described by cellular activities  $V = \{v_1 \dots v_B\}$  and observations  $X = \{x_1 \dots x_B\}$ . Then the following likelihood distributions hold

$$p(V | g) = \prod_{b=1}^B \mathcal{N}(v_b; 0, C_v) \quad (7)$$

$$p(X | V) = \prod_{b=1}^B \mathcal{N}(x_b; zAv_b, C_x) \quad (8)$$

## 2 Gibbs sampling as the E-step

An efficient way to collect samples from the joint posterior over all hidden variables is to employ a Gibbs sampling scheme, where we sample from the conditional posteriors. The first is over  $v$ , and can be defined as follows

$$p(v | x, g, z) = \frac{p(x | v, z, g)p(v | z, g)}{p(x | z, g)} = \frac{\mathcal{N}(x; zAv, \sigma_x I) \mathcal{N}(v; 0, C_v)}{\int_{-\infty}^{\infty} \mathcal{N}(x; zAv, \sigma_x I) \mathcal{N}(v; 0, C_v) dv} \quad (9)$$

The Gaussian over  $x$  can be rewritten to a Gaussian over  $v$  times a constant  $c_1$  in the following way

$$\begin{aligned} \mathcal{N}(x; zAv, \sigma_x I) &= c_1 \mathcal{N}(v; \frac{1}{z}(A^T A)^{-1} A^T x, \frac{\sigma_x}{z^2}(A^T A)^{-1}) = \\ &= c_1 \mathcal{N}(v; \frac{1}{z} A^+ x, \frac{\sigma_x}{z^2}(A^T A)^{-1}) \end{aligned} \quad (10)$$

where  $A^+$  is the Moore-Penrose pseudoinverse of  $A$ , with  $D_x = D_v \rightarrow A^+ = A^{-1}$ . Consequently, the product of two Gaussians in the numerator of Eq. 9 can also be written as a Gaussian over  $v$  introducing a new constant

$$\mathcal{N}(x; zAv, \sigma_x I) \mathcal{N}(v; 0, C_v) = c_1 c_2 \mathcal{N}(v; \mu_{post}, C_{post}) \quad (11)$$

The denominator of Eq. 9 is the integral of this formula, which evaluates to  $c_1 c_2$ , as the Gaussian integrates to one. This cancels the constants in the numerator, making the conditional posterior equal to the combined Gaussian over  $v$ , which, after expanding  $\mu_{post}$  and  $C_{post}$ , is

$$p(v | x, g, z) = \mathcal{N}\left(v; \frac{z}{\sigma_x} \left(\frac{z^2}{\sigma_x} A^T A + C_v^{-1}\right)^{-1} A^T x, \left(\frac{z^2}{\sigma_x} A^T A + C_v^{-1}\right)^{-1}\right) \quad (12)$$

and for a batch of size  $B$

$$p(V | X, g, z) = \prod_{b=1}^B \mathcal{N}(v_b; \mu_{post}(x_b), C_{post}) \quad (13)$$

which can be sampled directly from a Gaussian of dimension  $D_v \times B$ .

The conditional posterior over  $g$  is defined as follows

$$p(g | X, V, z) = \frac{p(X | g, V, z)p(g | V, z)}{p(X | V, z)} = \frac{p(V | g)p(g)}{p(V)} \quad (14)$$

which can be sampled by an MCMC sampling scheme with the following target

$$\log p(g | X, V) \sim -\frac{1}{2} \left[ B \log(\det(C_v)) + \sum_{b=1}^B v_b^T C_v^{-1} v_b \right] + (\alpha - 1) \sum_{k=1}^K \log(g_k) \quad (15)$$

We can cycle over elements of  $g$  with a Gibbs sampling scheme too. The unnormalised conditionals, assuming an independent (e.g. Gamma) prior look as follows

$$\begin{aligned} \log p(g_j | g_{\neg j}, X, V) &= \frac{p(V | g_j, g_{\neg j}, X)p(g_j | g_{\neg j}, X)}{p(V | g_{\neg j}, X)} = \\ &= \frac{p(V | g)p(g_j)}{p(V | g_{\neg j})} \sim p(V | g)p(g_j) \end{aligned} \quad (16)$$

These one-dimensional targets have to be sampled by MCMC too.

The conditional posterior over  $z$  will look the following

$$p(z | X, V, g) = \frac{p(X | g, z, V)p(z | V, g)}{p(X | V, g)} \sim p(X | z, V)p(z) \quad (17)$$

the log-posterior being

$$\log p(z | X, V) \sim -\frac{1}{2} \left[ B D_x \log(\sigma_x) + \frac{1}{\sigma_x} \sum_{b=1}^B (x_b - z A v_b)^T (x_b - z A v_b) \right] + \log p(z) \quad (18)$$

also to be sampled by MCMC.

### 3 Hamiltonian sampling as the M-step

The complete posterior:

$$\begin{aligned} p(v, g, z | x) &\sim p(x | v, g, z)p(v | g)p(g)p(z) \quad (19) \\ \log p(v, g, z | x) &\sim -\frac{1}{2} \left[ \frac{1}{\sigma_x} (x - z A v)^T (x - z A v) + \log(\det(C_v)) + v^T (C_v)^{-1} v \right] + \\ &+ \sum_{j=1}^K \left[ (sh_g - 1) \log(g_j) - \frac{g_j}{sc_g} \right] + (sh_z - 1) \log(z) - \frac{z}{sc_z} \end{aligned} \quad (20)$$

the gradient with respect to different variable types assuming a Gamma prior over  $g$  and  $z$ :

$$\frac{\partial \log p(v, g, z | x)}{\partial v} = -\frac{z}{\sigma_x} \left[ z A^T A + (C_v)^{-1} \right] v + \frac{z}{\sigma_x} A^T x \quad (21)$$

$$\frac{\partial \log p(v, g, z | x)}{\partial g_i} = -\frac{1}{2} \text{Tr} \left[ (C_v^{-1} - C_v^{-1} v v^T C_v^{-1}) C_i \right] + \frac{s h_g - 1}{g_i} - \frac{g_i}{s c_g} \quad (22)$$

$$\frac{\partial \log p(v, g, z | x)}{\partial z} = 2 v^T A^T A v z - x^T A v + \frac{s h_z - 1}{z} - \frac{z}{s c_z} \quad (23)$$

## 4 Gradient descent as the M-step

The complete-data likelihood with respect to a set of batch observations of size  $B$  is the following

$$p(\mathbf{V}, G, \mathbf{X} | C_{1..K}) = \prod_{n=1}^N p(X_n | V_n) p(V_n | g_n) p(g_n) = \prod_{n=1}^N p(g_n) \prod_{b=1}^B p(x_{nb} | v_{nb}) p(v_{nb} | g_n) \quad (24)$$

Let's denote the logarithm of this by  $\mathcal{L} = \log p(\mathbf{V}, G, \mathbf{X} | C_{1..K})$ . We can approximate the integral of this logarithm over the joint posterior by averaging over  $L$  samples from it, separately for each observation  $x_n$ . As we will seek the values of the precision components  $C_{1..K}$  that maximise this integral, we can discard each term not depending on these parameters. This way we arrive to the following expression

$$\begin{aligned} \mathcal{L} &\sim \sum_{n=1}^N \frac{1}{L} \sum_{l=1}^L -\frac{1}{2} \left[ B \log \left( \det \left( C_v^{(l,n)} \right) \right) + \sum_{b=1}^B v^{(l,n,b)T} (C_v^{(l,n)})^{-1} v^{l,n,b} \right] = \\ &= -\frac{1}{2L} \sum_{m=1}^{NL} \left[ B \log \left( \det (C_v^m) \right) + \sum_{b=1}^B v^{(m,b)T} (C_v^m)^{-1} v^{m,b} \right] \end{aligned} \quad (25)$$

noting that the double summation over  $L$  samples over all  $N$  observations always happens on the same terms, so we can substitute it with a single sum that iterates over the full sample set.

To ensure that the optimisation procedure does not produce precision matrices that are not positive definite, we can optimise for the Cholesky upper triangle matrix instead of the precision matrix, as this also specifies the Gaussian completely.

$$C_k = U_k^T U_k \quad (26)$$

$$C_v = \sum_{k=1}^K g_k U_k^T U_k \quad (27)$$

$$\frac{\partial C_v^m}{\partial [U_k]_{i,j}} = g_k^m \frac{\partial (U_k^T U_k)}{\partial [U_k]_{i,j}} = g_k^m (U_k^T J^{ij} + J^{ji} U_k) \equiv g_k^m \hat{U}_k^{ij} \quad (28)$$

where  $J^{ij}$  is the single-entry matrix so that its element at index  $(i, j)$  is 1, and 0 everywhere else. Then by the chain rule, the derivative of  $\mathcal{L}$  according to an element of  $U_k$  looks like this

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial [U_k]_{i,j}} &= -\frac{1}{2L} \sum_{m=1}^{NL} \text{Tr} \left[ \frac{\partial \mathcal{L}^m}{\partial C_v^m} \frac{\partial C_v^m}{\partial [U_k]_{i,j}} \right] = \\ &= -\frac{1}{2L} \sum_{m=1}^{LN} \text{Tr} \left[ \left[ B (C_v^m)^{-1} - \sum_{b=1}^B (C_v^m)^{-1} v^{m,b} v^{(m,b)T} (C_v^m)^{-1} \right] g_k^m \hat{U}_k^{ij} \right] = \\ &= -\frac{1}{2L} \text{Tr} \left[ \sum_{m=1}^{LN} g_k^m \left[ B (C_v^m)^{-1} - (C_v^m)^{-1} \left( \sum_{b=1}^B v^{m,b} v^{(m,b)T} \right) (C_v^m)^{-1} \right] \hat{U}_k^{ij} \right] \end{aligned} \quad (29)$$

The regularities of the  $\hat{U}_k$  matrices allow us to replace the trace with a much more efficient computation:

$$M = -\frac{1}{2L} \left[ \sum_{m=1}^{LN} g_k^m \left[ B (C_v^m)^{-1} - (C_v^m)^{-1} \left( \sum_{b=1}^B v^{m,b} v^{(m,b)T} \right) (C_v^m)^{-1} \right] \right] \quad (30)$$

$$\frac{\partial \mathcal{L}}{\partial [U_k]_{i,j}} = \sum_{a=1}^{D_v} [M]_{j,a} [U_k]_{i,a} + [M]_{a,j} [U_k]_{i,a} \quad (31)$$

As a generalised M-step, we can move the parameters in the direction of the gradient scaled by a learning rate

$$[U_k]_{i,j}^{new} = [U_k]_{i,j}^{old} + \epsilon \frac{\partial \mathcal{L}}{\partial [U_k]_{i,j}} \quad (32)$$