

Supplement for Neuro-Biologically Plausible Probabilistic Inference Inspired by Topic Models

Anonymous Author(s)

Affiliation

Address

email

1 Full Network Derivation for LDA

Variational inference for LDA is summarized by iterative equations which take the form

$$p_{d,i,j} = \frac{\exp(\psi(\eta_{i,j}) - \psi(\bar{\eta}_j)) \exp \psi(\alpha_{d,j})}{\sum_k \exp(\psi(\eta_{i,k}) - \psi(\bar{\eta}_k)) \exp \psi(\alpha_{d,k})} \quad (1)$$

$$\alpha_{d,j} = \alpha_0 + \sum_i R_{d,i} p_{d,i,j} \quad (2)$$

$$\eta_{i,j} = \eta_0 + \sum_d R_{d,i} p_{d,i,j} \quad (3)$$

When the i th element of vector $\eta_j = \eta_{i,j}$ the approximate posterior is given by a product of marginals which take the form

$$q(\beta_j | \eta_j) = \text{Dir}(\eta_j) \quad (4)$$

$$q(\theta_d) = \text{Dir}(\alpha_d) \quad (5)$$

$$q(N_{d,i,:}) = \text{Mult}(p_{d,i,:}, R_i) \quad (6)$$

$$(7)$$

To obtain the network equations we define intermediate computation quantities:

$$w_{i,j} = \exp(\psi(\eta_{i,j})) \quad (8)$$

$$\bar{\eta}_j = \sum_i \eta_{i,j} \quad (9)$$

$$\bar{R}_{d,i} = \sum_j w_{i,j} \exp(\psi(\alpha_{d,j})) \quad (10)$$

$$\bar{N}_{d,i,j} = w_{i,j} \frac{R_{d,i}}{\bar{R}_{d,i}} \quad (11)$$

This allows us to rewrite the update equations for inference on single document as

$$\bar{N}_{d,i,j} = w_{i,j} \frac{R_{d,i}}{\bar{R}_{d,i}} \quad (12)$$

$$\alpha_{d,j} = \alpha_0 + \exp(\psi(\bar{\eta}_j)) \exp(\alpha_{d,j}) \sum_i \bar{N}_{d,i,j} \quad (13)$$

Clearly, for fixed $w_{i,j}$, iterating between these two equations converges to $\alpha_{d,j}$. The network which performs inference is then built by constructing a continuous time dynamical system which builds converges to the same fixed point as these two iterative equations, i.e. it has the property that

$$\frac{d}{dt}\alpha_{d,j} = \frac{d}{dt}\bar{N}_{d,i,j} = 0 \quad (14)$$

when the iterative equations above are satisfied. One way to accomplish this is to set

$$\frac{d}{dt}\bar{N}_{d,i,j} = w_{i,j}R_{d,i} - \bar{N}_{d,i,j}\bar{R}_i \quad (15)$$

$$\frac{d}{dt}\alpha_{d,j} = \exp(\psi(\bar{\eta}_j))(\alpha_0 - \alpha) + \exp(\alpha_{d,j}) \sum_i \bar{N}_{d,i,j} \quad (16)$$

To implement the learning rule. We utilize the approximation $\exp(\psi(x)) \approx x - 1/2$ for $x > 1$ so that $w_{i,j} = \eta_{i,j} - 1/2$. This allows us to write Hoffman's online learning rule

$$\tau_L \frac{d}{dt}\eta_{i,j} = \eta_0 - \eta_{i,j} + R_{t,i}p_{t,i,j} \quad (17)$$

as

$$\tau_L \frac{d}{dt}w_{i,j} = \eta_0 - 1/2 - w_{i,j} + \sum_d N_{d,i,j} \exp(\psi(\alpha_{d,j})) \quad (18)$$

Rescaling the right hand side by a positive quantity does not affect the location of the fixed point of the dynamics and results in

$$\tau_L \frac{d}{dt}\eta_{i,j} = \exp(\psi(\bar{\eta}_j))(\eta_0 - \eta_{i,j}) + \exp(\psi(\bar{\eta}_j)) \sum_d R_{d,i}p_{d,i,j} \quad (19)$$

$$= \exp(\psi(\bar{\eta}_j))(\eta_0 - \eta_{i,j}) + \exp(\alpha_{d,j})\bar{N}_{d,i,j} \quad (20)$$

as desired.

2 Full VB Derivation for Online DDM inference

2.1 Network Implementation

We assume that the quantity of interest is the current distribution of latent causes, $p(\theta(t)|R(\tau), \tau = 0..T)$, where $R_i(t)$ represents the spikes from neuron i as a sum of delta functions. If no spikes occur then no evidence is presented and posterior inference over $\theta(t)$ is simply given by an undriven Kalman filter with parameters (Λ, Σ_D) . A recurrent neural network which uses a PPC to encode a posterior which evolves according to a Kalman filter has the property that neural responses are linearly related to the inverse covariance matrix of the posterior and inverse covariance times the posterior mean. In the absence of evidence, it is easy to show that these quantities must evolve according to quadratic recurrent dynamics which implement divisive normalization. Thus, the patterns of neural activity which linearly encode them must do so as well. When a new spike arrives, optimal inference is no longer possible and a variational approximation must be utilized. As is shown in the supplement this variational approximation is similar to the variational approximation used for LDA. As a result, a network which can divisively inhibit its synapses is once again sufficient for the purposes of implementing approximate Bayesian inference.

108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161

$$\begin{aligned} \frac{d}{dt} \Sigma^{-1} &= \text{Diag} \left(\exp \left(\boldsymbol{\mu} + \frac{1}{2} \text{Diag}(\Sigma) \right) \right) + \mathbf{\Lambda}^T \Sigma^{-1} + \Sigma^{-1} \mathbf{\Lambda} - \Sigma^{-1} \Sigma_D \Sigma^{-1} \\ \frac{d}{dt} (\Sigma^{-1} \boldsymbol{\mu})_j &= \sum_i R_i^t \frac{w_{ij} \exp(\mu_{t,j})}{\sum_k w_{ik} \exp(\mu_{t,k})} - \exp \left(\mu_j + \frac{1}{2} \text{Diag}(\Sigma)_j \right) + (\mathbf{\Lambda}^T \Sigma^{-1} \boldsymbol{\mu} - \Sigma^{-1} \Sigma_D \Sigma^{-1} \boldsymbol{\mu})_j \end{aligned} \quad (21)$$

2.2 Full Derivation

Just as the Dirichlet-Multinomial formulation of LDA has an equivalent Gamma-Poisson formulation, so too can we consider a Log Normal-Poisson formulation of the correlated topic model. The generative model takes the form

$$\begin{aligned} \mathbf{z}(t) &\sim N(\mathbf{z}(t - \Delta t) + \Delta t \mathbf{\Lambda} (\mu_0 - \mathbf{z}(t - \Delta t)), \Delta t \Sigma_D) \\ c_j(t) &= \exp(z_j(t)) \\ g(t) &\sim N(g(t - \Delta t) + \Delta t \lambda_g (\mu_{g_0} - g(t - \Delta t)), \Delta t \sigma_g^2) \\ \mathbf{w}_i &\sim \text{Dir}(\boldsymbol{\alpha}) \\ N_{ij}(t) &\sim \text{Poi}(w_{ij} c_j(t) \Delta t) \\ R_i(t) &= \sum_j N_{ij}(t) \sim \text{Poi} \left(\Delta t \sum_j w_{ij} c_j(t) \right) \end{aligned} \quad (22)$$

Since we are doing online inference we are only interested in updating the posterior for a single time step in the forward direction using the posterior from the previous time step as the posterior for the next. For a time step of size Δt .

$$\begin{aligned} \log(\mathbf{R}^t, \mathbf{N}^t, \mathbf{z}^t) &\sim \log \left(\delta \left(R_i^t - \sum_j N_{ij}^t \right) \right) + \sum_{ij} N_{ij}^t (\log w_{ij} + z_j^t + \log \Delta t) - \Delta t \sum_j \exp(z_j^t) \\ &\quad - \frac{1}{2} (\mathbf{z}^t - \mu_{t-\Delta t} - \Delta t \mathbf{\Lambda} (\mu_0 - \mu_{t-\Delta t})) \cdot \\ &\quad \left((\mathbf{I} - \Delta t \mathbf{\Lambda}) \Sigma_{t-\Delta t} (\mathbf{I} - \Delta t \mathbf{\Lambda})^T + \Delta t \Sigma_D \right)^{-1} (\mathbf{z}^t - \mu_{t-\Delta t} - \Delta t \mathbf{\Lambda} (\mu_0 - \mu_{t-\Delta t})) \end{aligned} \quad (23)$$

A variational posterior of the form

$$q(\mathbf{N}, \mathbf{z}) = q(\mathbf{N}_{i,:} | R_i^t, p_{i,:}^t) q(\mathbf{z} | \boldsymbol{\mu}^t, \Sigma^t) \quad (24)$$

This immediately results in an update rule for the Multinomial distribution which takes the familiar form,

$$\langle N_{ij}^t \rangle = R_i \frac{w_{ij} \exp(\mu_j^t)}{\sum_k w_{ik} \exp(\mu_k^t)} \quad (25)$$

Update rules for the Gaussian posterior on $\mathbf{z}(t)$ are obtained for two different cases. Case 1: no spikes are present. In this case we minimize the residual KL divergence by gradient descent, i.e. by taking the derivative with respect to the parameters of interest of

$$\begin{aligned} KL(q(\mathbf{z} | \boldsymbol{\mu}_t, \Sigma_t) | p) &\sim \sum_{ij} \langle N_{ij}^t \rangle \mu_{t,j} - \Delta t \sum_j \exp(\mu_{t,j} + \frac{1}{2} \sigma_{t,j}^2) \\ &\quad - \frac{1}{2} \text{Tr} \left((\Sigma_t + \mu_t \mu_t^T) \left((\mathbf{I} - \Delta t \mathbf{\Lambda}) \Sigma_{t-\Delta t} (\mathbf{I} - \Delta t \mathbf{\Lambda})^T + \Delta t \Sigma_D \right)^{-1} \right) \\ &\quad + \mu_t \cdot \left((\mathbf{I} - \Delta t \mathbf{\Lambda}) \Sigma_{t-\Delta t} (\mathbf{I} - \Delta t \mathbf{\Lambda})^T + \Delta t \Sigma_D \right)^{-1} (\mu_{t-\Delta t} + \Delta t \mathbf{\Lambda} (\mu_0 - \mu_{t-\Delta t})) \\ &\quad + \frac{1}{2} \log |\Sigma_t| \end{aligned} \quad (26)$$

This results in

$$\begin{aligned} \frac{d}{dt} \Sigma^{-1} &= \text{Diag} \left(\exp \left(\mu + \frac{1}{2} \text{diag}(\Sigma) \right) \right) + \Lambda^T \Sigma^{-1} + \Sigma^{-1} \Lambda - \Sigma^{-1} \Sigma_D \Sigma^{-1} \\ \frac{d}{dt} \Sigma^{-1} \mu &= - \exp \left(\mu + \frac{1}{2} \text{diag}(\Sigma) \right) + \Lambda^T \Sigma^{-1} \mu - \Sigma^{-1} \Lambda \mu_0 - \Sigma^{-1} \Sigma_D \Sigma^{-1} \mu \end{aligned} \quad (27)$$

In the second case, we suppose that a spike in the input layer has occurred. Since the underlying spike generation process is Poisson, this event occurs in a time window of size zero. This allows us to greatly simplify inference by considering only those terms in the log probability which are present in the limit as Δt goes to zero.

This results in update equation:

$$\begin{aligned} \langle N_{ij}^t \rangle &= R_i^t \frac{w_{ij} \exp(\mu_{t,j})}{\sum_k w_{ik} \exp(\mu_{t,k})} \\ \Sigma_t^{-1} &= \Sigma_{t-}^{-1} \\ \Sigma_t^{-1} \mu_t &= \Sigma_{t-}^{-1} \mu_{t-} + \sum_i \langle \mathbf{N}_i^t \rangle \end{aligned} \quad (28)$$

Now, full VB inference would require that we evolve the dynamics of equation 27 until a spike occurs then iterate the equations 28 until convergence before continuing. This was the procedure used to generate the x-axis values of the blue dots in Figure 6 right in the main text. A more biologically plausible scheme would only iterate 28 only once. This results in the DDM inference equation in the main text. This approximation results in almost zero information loss due in large part to the fact that prior provided by knowledge of the dynamics is actually quite informative.

3 Simulations

The points in Figure 6 were generated by sampling from the generative models for LDA and the DDM and then applying the network implementation of VBEM or VBEM itself to those samples. The output of both algorithms is a set of natural parameters associated with each olfactory scene or in the case of the DDM each time step.

3.1 LDA

For LDA, we only plotted the natural parameters of the gamma distributed posterior over concentration, $q_c(c|\alpha)$. Here, the natural parameter is the shape parameter of the gamma distributed posterior which is proportional to the network's estimate of the concentration. We considered a bank of 100 odors and 20 olfactory receptor neurons. Odors were drawn randomly from a symmetric Dirichlet distribution with parameter $\alpha = 1/2$. Each olfactory scene contained on average 5 odors with concentration significantly different from zero. Average concentration level was chosen so that each olfactory receptor neuron fired on average 25 spikes per trial.

3.2 DDM

For the DDM the natural parameters of the posterior distribution over log concentration are the precision matrix (inverse covariance matrix) and the precision matrix times them posterior mean. In this case we also considered 100 odors and 20 sensory input neurons and generated odors from a symmetric Dirichlet with $\alpha = 1/2$. The dynamics of the OU process which governed the behaviour of the odor concentration were given by $\Lambda = 2\mathbf{I}$ and $\Sigma_D = 4\log(21)\mathbf{I}$. This ensures that, once again, there are about 5 odors which are present at a significant level of concentration at any given time.

4 Table of properties for LDA, ICA, and L1-based networks

Neural network implementations of LDA, ICA, and L1 make very different predictions for neural computation. ICA is a feedforward algorithm. Network implementations of L1 typically use linear recurrent inhibition to compute a prediction of the input image which is then subtracted from the actual image. Sparsity is then achieved via a competitive race to explain mechanism. The network

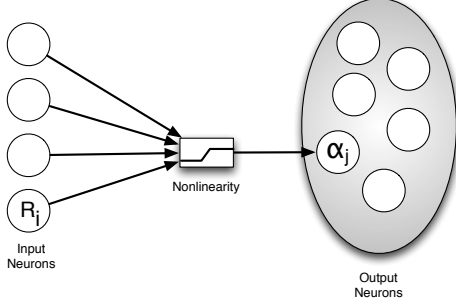


Figure 1: ICA network model

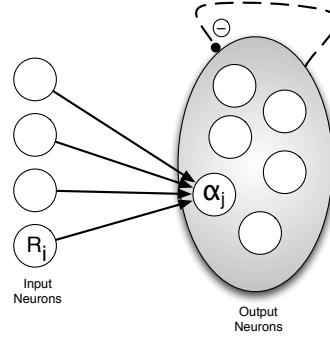


Figure 2: L1 network model

implementation of LDA, on the other hand, divides the image by the prediction to generate the residual signal which drives the refinement of the representation. Other differences are summarized in the table below.

Table 1: Network Model Comparison

Attribute	ICA	L1	LDA
Reccurently Modulated	N	Y	Y
Divisive Normalization	N	N	Y
Subtractive Inhibition	Y	Y	Y
Only Positive Weights	N	N	Y
Contrast Invariance	N	N	Y

5 Comparison of LDA, ICA and L1 regression on Natural Images

Like ICA and L1 penalized regression, LDA has at its core a linear mixing model capable of implementing a sparse prior. In the Gamma-Poisson formulation, this is achieved by setting the α parameter of the gamma distribution on the un-normalized topic distribution to a value less than one. This results in greater excess kurtosis than the L1 prior. The primary differences between the LDA and these other methods are that (1) LDA naturally functions in the domain of count data, and (2) LDA naturally enforces a non-negativity constraint on the mixture components. In this sense it is similar to Non-negative Matrix Factorization (NMF). In the visual domain NMF applied to a data set of whitened natural images has been shown to be capable of reproducing oriented receptive fields similar to those obtained by Olshausen and Field using L1 penalized regression and van Hateren and van der Schaaf using ICA [1, 2, 3]. Here we extend this result by showing that NMF implemented via variational inference applied to LDA also generates oriented filters when applied to *un-whitened* natural images. We used ICA code and natural images from [3] and L1 code from [2]. This results in a marked difference between LDA and ICA and L1 penalized regression, the latter two of which fail miserably at reproducing any kind of a sensible representation when applied to un-whitened data. See Figure 3.

References

- [1] Patrik O Hoyer. Non-negative Matrix Factorization with Sparseness Constraints. *Journal of Machine Learning Research*, 5:1457–1469, 2004.
- [2] B.A. Olshausen and D.J. Field. Emergence of simple-cell receptive field properties by learning a sparse code for natural images. *Nature*, 381(6583):607–609, 1996.
- [3] J.H. van Hateren and A. van der Schaaf. Independent component filters of natural images compared with simple cells in primary visual cortex. *Proceedings of the Royal Society of London. Series B: Biological Sciences*, 265(1394):359, 1998.

270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323

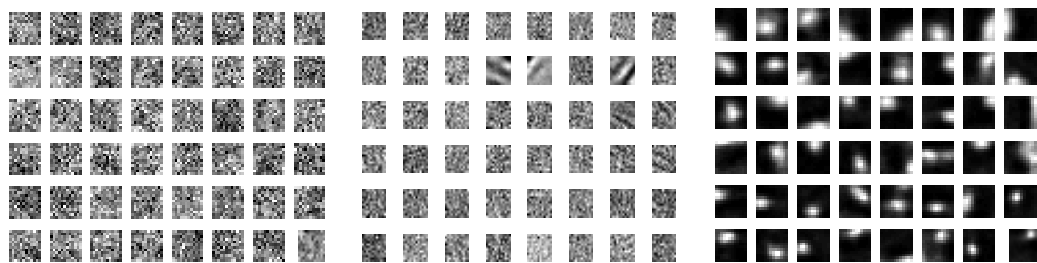


Figure 3: Filters found with (Left) L1 Olshausen&Field (Middle) ICA (Right) LDA when trained with *un-whitened* natural image data.