

A Bayesian Variational Principle for Dynamic Self Organizing Maps

Anthony Fillion¹ and Thibaut Kulak¹ and François Blayo¹

1- NeoInstinct - NeoLab
3 rue Traversière, Lausanne - Switzerland

Abstract. We propose organisation conditions that yield a method for training SOM with adaptative neighborhood radius in a variational Bayesian framework. This method is validated on a non-stationary setting and compared in an high-dimensional setting with an other adaptative method.

1 Introduction

Self Organizing Maps (SOM, [1]) is a biologically inspired unsupervised vector quantization method for data modelization and visualization. A map is made of a collection of weights $\{w_z\}_z$ indexed by a discrete collection of points $z \in \{z_1, z_2, \dots, z_n\}$. The weights (or neurons) are meant to modelize some observation random variable \mathbf{x} . The points are usually regularly placed on a two dimensional grid. In this context, the SOM objective is to find representative weights with a well organized response to stimuli i.e. weights specialize their response to some kind of stimuli and weights close in the grid space are close in the observation space.

The original way to train such SOM is with Kohonen's iterations. Despite its simplicity, they require a time dependant neighborhood function whose neighborhood radius (or temperature) strictly decreases towards zero during training. In [2] and [3], heuristics are proposed in a non-probabilistic framework to adapt this temperature. They lead to "dynamic" SOM algorithms that are able to track non-stationary datasets. In [4], the variational EM framework is used to design a cost function for the SOM but the temperature parameter is still decreased as in Kohonen. In [5], the topologically preserving properties of the map stems from the basis function regularity. Hence those basis functions have to be carefully chosen and their number must typically grow exponentially with the observation space dimension.

We propose quantitative organization conditions that lead to an adaptative temperature choice in a variational Bayesian framework. High dimensional experiments show that our method is less sensitive to elasticity and outperforms that of [3].

2 Evidence Lower BOund

In [4], the neighborhood function is interpreted as the variational posterior in a Expectation-Minization (EM, [6]) procedure. Therefore, they propose to mini-

mize the following (opposite) ELBO to train SOMs

$$F(\rho, \theta) = \int \int \ln \left(\frac{q_{\mathbf{z}|\mathbf{x}}^\rho}{p_{\mathbf{x},\mathbf{z}}^\theta} \right) q_{\mathbf{z}|\mathbf{x}}^\rho d\mathbf{z} p_{\mathbf{x}} d\mathbf{x} \quad (1)$$

$$= \underbrace{\int \mathcal{D}(q_{\mathbf{z}|\mathbf{x}}^\rho, p_{\mathbf{z}|\mathbf{x}}^\theta) p_{\mathbf{x}} d\mathbf{x}}_{\text{organization}} + \underbrace{\mathcal{D}(p_{\mathbf{x}}, p_{\mathbf{x}}^\theta)}_{\text{modelization}}, \quad (2)$$

where, the data true density is $p_{\mathbf{x}}$ and the discrete latent variable is \mathbf{z} . The complete data model $p_{\mathbf{x},\mathbf{z}}^\theta \propto e^{-\frac{1}{2}|\frac{\mathbf{x}-w_{\mathbf{z}}}{\sigma}|^2}$ is chosen Gaussian and the variational (amortized, homoskedastic) posterior $q_{\mathbf{z}|\mathbf{x}}^\rho \propto e^{-\frac{1}{2}|\frac{\mathbf{z}-\mu_{\mathbf{x}}}{\lambda}|^2}$ as well. Their parameters are respectively $\theta = \{z \mapsto w_z, \sigma\}$, $\rho = \{x \mapsto \mu_x, \lambda\}$. The decomposition in the second equation reveals a modelization term that is the Kullback-Leibler (KL) divergence between the data density and the model marginal. As well as an other term that is the expected KL divergence between the variational posterior and the model posterior (or response to stimulus). The core idea is to choose the variational posterior organized so that this term will favor organized posterior models.

3 Organization conditions

Given an observation \mathbf{x} , let $w_{z_{\mathbf{x}}^*}$ be its closest weight. If the distance to this weight $|w_{z_{\mathbf{x}}^*} - w_{\mathbf{z}}|$ in observation space is increasing with the distance $|z_{\mathbf{x}}^* - \mathbf{z}|$ in point space then the map is organized. This mean that if the Bayesian posterior (or response) $p_{\mathbf{z}|\mathbf{x}}^\theta \propto e^{-\frac{1}{2}|\frac{\mathbf{x}-w_{\mathbf{z}}}{\sigma}|^2}$ is plotted as a function of \mathbf{z} , it would be strictly decreasing around its maximum. This is illustrated in Fig. 1. Hence we propose the following organization conditions:

1: The response $p_{\mathbf{z}|\mathbf{x}}^\theta$ has a dominating mode at $z_{\mathbf{x}}^* = \arg \min_{\mathbf{z}} |\mathbf{x} - w_{\mathbf{z}}|$.

2: $|w_{\mathbf{z}} - w_{z_{\mathbf{x}}^*}| = \frac{1}{\eta} |\mathbf{z} - z_{\mathbf{x}}^*|$ when \mathbf{z} is in the neighborhood of $z_{\mathbf{x}}^*$.

The first condition ensures that each weight specializes on some kind of stimulus. The second one ensures that the weight distance is proportional to the point distance by a scale η . In other words, the map preserves distances locally.

3.1 Variational posterior selection

We now choose ρ such that the variational posterior $q_{\mathbf{z}|\mathbf{x}}^\rho$ verifies the organization conditions. Because the variational posterior has a unique mode at $\mu_{\mathbf{x}}$, choosing

$$\mu_{\mathbf{x}} = z_{\mathbf{x}}^*$$

(like Kohonen does) favors reponses that verify the first order organization condition. We also have,

$$-\ln p_{\mathbf{z}|\mathbf{x}}^\theta = \frac{1}{\sigma^2} \langle \mathbf{x} - w_{z_{\mathbf{x}}^*} | w_{z_{\mathbf{x}}^*} - w_{\mathbf{z}} \rangle + \frac{1}{2} \left| \frac{w_{\mathbf{z}} - w_{z_{\mathbf{x}}^*}}{\sigma} \right|^2 + c,$$

with c being a constant w.r.t \mathbf{z} . Because $z_{\mathbf{x}}^*$ is a mode, the scalar product is almost null for \mathbf{z} in the neighborhood of $z_{\mathbf{x}}^*$. Therefore, if \mathbf{z} is in the neighborhood of $z_{\mathbf{x}}^*$, the organization term will favor responses such that $\left| \frac{w_z - w_{z^*}}{\sigma} \right| \simeq \left| \frac{z - z^*}{\lambda} \right|$. Thus choosing

$$\lambda = \eta\sigma,$$

favors responses satisfying the second organization condition.

3.2 Stochastic Gradient calculations

With the previous selection for ρ ($\mu_{\mathbf{x}} = z_{\mathbf{x}}^*, \lambda = \eta\sigma$), the objective F in Eq.(1) only depends on θ . Algorithm 1 computes a stochastic approximation over \mathbf{x} of its derivatives along $\theta = \{w, \sigma\}$ assuming that $z_{\mathbf{x}}^* = \arg \min_z |\mathbf{x} - w_z|$ has no gradient. Then any stochastic gradient descent method can be used to minimize F . A Python implementation is available at: github.com/anthony-Neo/VDSOM.

Algorithm 1: Stochastic gradient of F

Data: η : elasticity;
 $Z = \{z_i\}_{1 \leq i \leq n}$: grid points;
 $\{w_z\}_{z \in Z}, \sigma$: weights and variance;
 x, m : observation and observation space dimension
Result: $g_{\sigma}, \{g_{w_z}\}_{z \in Z}$ the stochastic gradients of F

```

1 for  $z, y \in Z \times Z$  do
2    $d_{z,y} := |z - y|^2$ ;
3    $f_z := |x - w_z|^2$ ;
4    $\ln p_z := -m \ln \sigma - \frac{f_z}{2\sigma^2}$ ;
5  $z^* := \arg \min_z f_z$ ;
6  $q_z := \text{softmax}(-\frac{d_{z,z^*}}{2\eta^2\sigma^2})$ ;
7  $d^* := \sum_z d_{z,z^*} \times q_z$ ;
8  $g_{\sigma} := \frac{m}{\sigma} + \frac{1}{\sigma^3} \sum_z [\eta(1 + \ln q_z - \ln p_z)(d_{z,z^*} - d^*) - f_z] q_z$ ;
9 for  $z \in Z$  do
10   $g_{w_z} := -\frac{q_z}{\sigma^2} (x - w_z)$ ;
```

4 Numerical Experiments

4.1 Non-stationary distributions

We used 8×7500 steps of the Adam stochastic optimizer [7] with a learning rate of $\alpha = 10^{-3}$. The variance parameter is initialized at $\sigma_0 = 5$ and the initial weights sample a normal Gaussian. The grid is a 15×15 node lattice with regularly spaced points between $-1, 1$ and the elasticity parameter is $\eta = 1$. During the first half of the iterations, the “moons” data set from sci-kit [8] learn

is sampled. Then, during the second half, the “circles” data set is sampled. Step 0 and each 7500 steps, data, weights and edges are plotted in Fig. 2 from the upper left corner to the lower right one. We see that the SOM tracks the changing data set, it correctly fits the observation density (rather than its support as in [3]) and it preserves the grid neighborhood. In Fig. 3, the log standard deviation σ and the log distortion (mean over the samples of the min squared distance with the weights, cf [3]) are plotted against time. Spikes appear half time which means that the method detects the change in data and adapts neighborhood radius accordingly.

4.2 High-dimensional distributions

The DSOM algorithm of [3] is compared to ours (VDSOM) on 20000 samples of the MNIST Fashion dataset on a 10×10 toroidal grid. DSOM uses a learning rate of $\alpha = 10^{-3}$ while VDSOM uses the same configuration as before. In Fig. 4, the distortion is plotted as a function of the elasticity η . VDSOM outperforms DSOM while being less sensitive. In Fig. 5, the weights of both methods trained with their respective optimal elasticity are displayed on a grid. Note that some of the DSOM weights are noise while this is not the case for VDSOM. VDSOM weights also seem less blurred.

5 Conclusion

We proposed organisation conditions that yield a method for training SOM with adaptative neighborhood radius in a variational Bayesian framework. This method has been validated in a non-stationary setting and compared in an high-dimensional setting with an other adaptative method.

References

- [1] T. Kohonen. *Self-Organizing Maps*. Springer Series in Information Sciences. Springer Berlin Heidelberg, 2012.
- [2] Erik Berglund and Joaquin Sitte. The parameterless self-organizing map algorithm. *IEEE transactions on neural networks / a publication of the IEEE Neural Networks Council*, 17:305–16, 04 2006.
- [3] Nicolas Rougier and Yann Boniface. Dynamic self-organising map. *Neurocomputing*, 74(11):1840–1847, 2011.
- [4] Jakob J Verbeek, Nikos Vlassis, and Ben JA Kröse. Self-organizing mixture models. *Neurocomputing*, 63:99–123, 2005.
- [5] Christopher M Bishop, Markus Svensén, and Christopher KI Williams. Gtm: The generative topographic mapping. *Neural computation*, 10(1):215–234, 1998.
- [6] C.M. Bishop. *Pattern Recognition and Machine Learning*. Information Science and Statistics. Springer, 2006.
- [7] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [8] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.

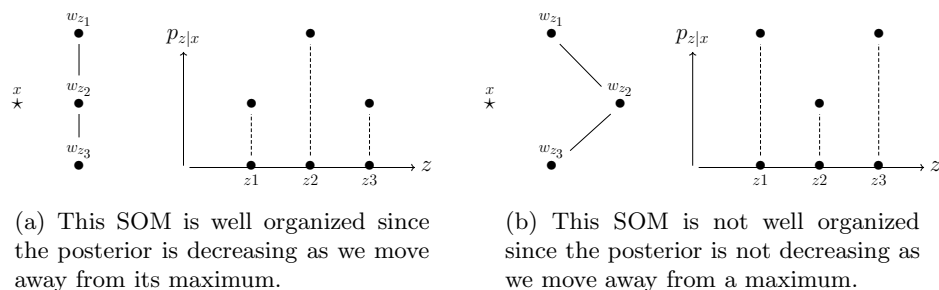


Fig. 1: On the left part of each subfigure, an observation x and a SOM with 1D neighborhood visualised in its 2D observation space. On the right part, the graph of its posterior.

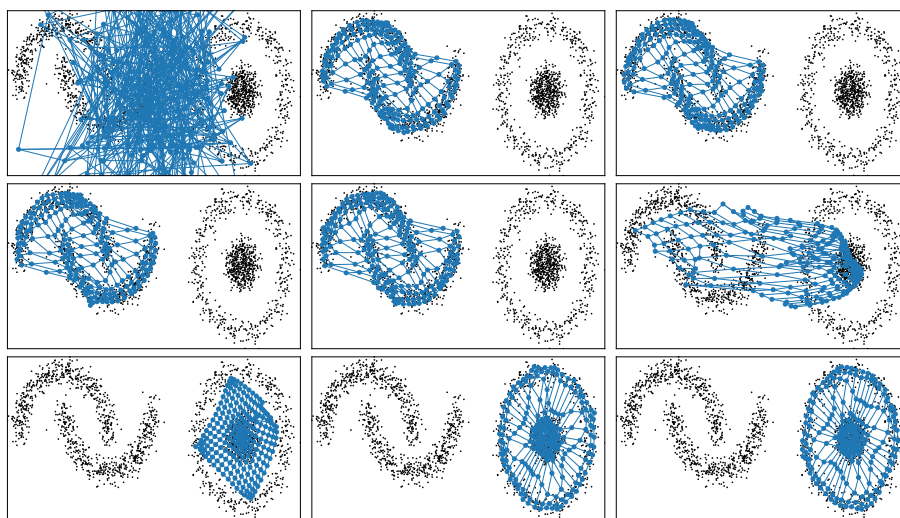


Fig. 2: Visualization of the map during the iterations, on a changing dataset

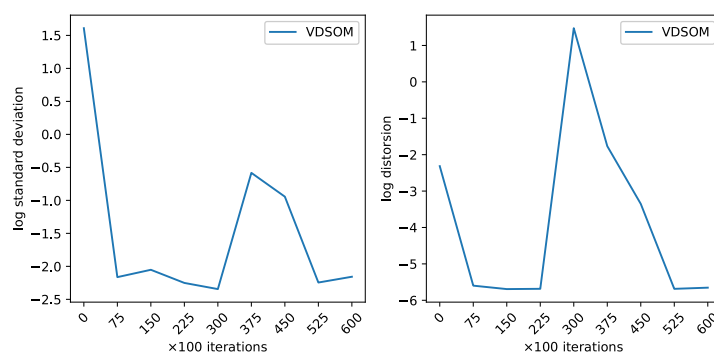


Fig. 3: Graphs of σ and distortion during the iterations, on a changing dataset

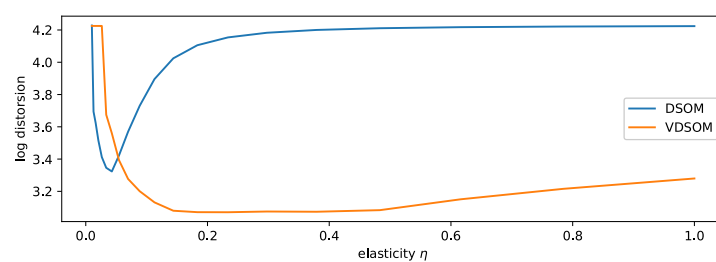
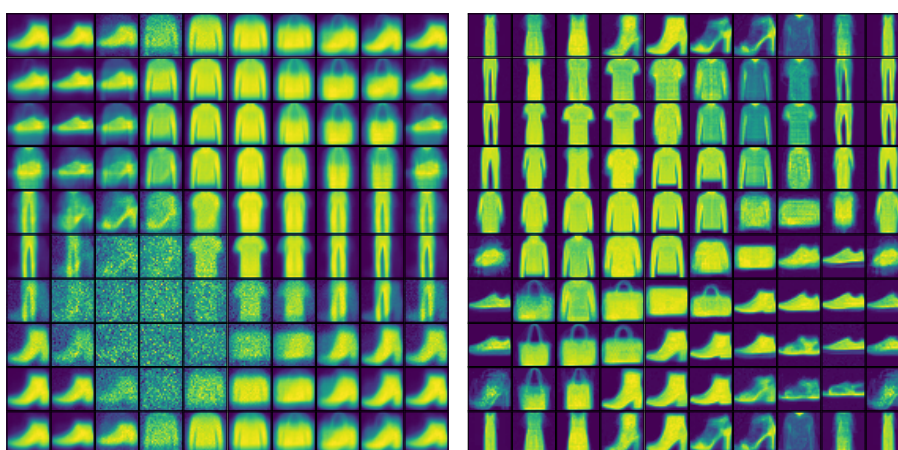


Fig. 4: sensitivity



(a) DSOM

(b) VDSOM

Fig. 5: Weights on the 2D grid