

Expectation Propagation for the Clutter Problem – Theory and Implementation

Matěj Korvas

June 3, 2013

Abstract

These lecture notes describe the algorithm of Expectation Propagation as applied to the Clutter problem, touching on the underlying theory. Important is also the practical part where I document my implementation of the algorithm.

1 Introduction

Expectation Propagation (EP for short) was introduced in [3] as an iterated version of the previously known Assumed-Density Filtering approximate inference algorithm. In the work [3], the author also shows how EP is applied to the clutter problem.

In the next section, we describe the EP algorithm in general, in Section 3, we formulate the clutter problem and derive formulas used in EP to solve it, and the final Section 4 discusses our implementation of EP applied to the clutter problem.

2 Expectation Propagation

Expectation Propagation is an approximate inference algorithm for graphical probabilistic models that factorise as follows:

$$p(\mathbf{z}, \mathbf{e}) = \prod_i f_i(\mathbf{z}, \mathbf{e}) \quad (1)$$

where \mathbf{z} is the vector of latent variables, \mathbf{e} is the vector of observed variables (evidence), and f_i are factors that depend on a non-empty subset of \mathbf{z} and a subset of \mathbf{e} . This factorisation naturally emerges in experiments with i.i.d. observations where f_0 is the prior on \mathbf{z} and f_i the posterior for the i -th observation for $i = 1, \dots$

TODO: Draw a figure of a graphical model that is typically used for doing EP.

EP approximates the factors f_i with factors \tilde{f}_i that belong to a convenient probability distribution family. The approximation aims to minimise the KL-divergence between a distribution computed using the exact factor f_i , and a distribution using the approximate factor \tilde{f}_i . If the approximating distribution family is chosen from the exponential family (which it typically is), minimising the KL-divergence is reduced to *matching moments*, i.e. setting a few moments of the estimating distribution (its sufficient statistics) to the values of corresponding moments of the distribution approximated. Choosing the family from the exponential family has also other benefits, including the fact that this family is closed under the operation of product (this property being assumed in the algorithm), and that Minka [3] proved the existence of a fixed point for the solution provided the family is exponential.

Choosing the approximating family is the first thing done in the algorithm. Next, approximated factors \tilde{f}_i and their product $Q = \prod_i \tilde{f}_i$ are initialised to uniform. The algorithm then proceeds in iterations, iteratively updating all the approximating factors in each of the outer iterations. When convergence is reached, the normalisation coefficient, an estimate of $p(\mathbf{e})$, is computed. A more detailed exposition of the algorithm follows.

1. Initialisation

All the approximate factors are initialised to uniform, meaning the initial approximation is non-informative. The product Q of the factors is computed accordingly. Typically, all the factors, as well as their product, are initialised to constant 1.

Factors that already belong to the chosen family can also be computed during initialisation, as such factors are always best approximated by themselves, not needing to be updated iteratively.

2. Outer loop

Following four steps are repeated until convergence.

2.1. Choose a factor \tilde{f}_i

Choose a factor to approximate.

2.2. Compute the cavity distribution $Q^{\setminus i}$

When updating the factor \tilde{f}_i , we would ideally want to minimise the KL-divergence between the true distribution and the resulting approximative distribution:

$$\arg \min_{\tilde{f}_i} \text{KL}(p \parallel \prod_i \tilde{f}_i). \quad (2)$$

However, there we would need to compute moments of p in order to optimise for this KL-divergence. If we were able to do that, we would not need to use approximate inference in the first place, so let us assume this is intractable. In that case, we have to substitute p with an approximation. The approximation used in EP is the following:

$$\hat{p} = \frac{1}{Z_i} f_i Q^{\setminus i} \quad (3)$$

where

$$Q^{\setminus i} \propto \prod_{j \neq i} \tilde{f}_j \quad (= Q / \tilde{f}_i). \quad (4)$$

Here, $Q^{\setminus i}$ is called the *cavity distribution*, as it is a distribution over \mathbf{z} obtained by multiplying all the approximate factors but the i -th one (thus creating the cavity in the distribution) and normalising (in order to make it a distribution). \hat{p} is defined as a product of the *exact* factor f_i with the rest of the factors *approximated*, normalised to 1, and the cavity distribution needs to be computed in order to express \hat{p} .

2.3. Compute the approximative distribution \mathcal{Q}_{new}

Whereas the previous step was concerned with computing the cavity distribution, computing the normalisation coefficient Z_i (as $\int_{\mathbf{Z}} f_i(\mathbf{z}) \mathcal{Q}^{\setminus i}(\mathbf{z}) d\mathbf{z}$) and the approximative posterior distribution \hat{p} is reserved for this step.

Having computed \hat{p} , we can minimise the KL-divergence to an updated \mathcal{Q}_{new} restricted to be in the approximating family \mathcal{F} :

$$\arg \min_{\mathcal{Q}_{\text{new}} \in \mathcal{F}} \text{KL}(\hat{p} \parallel \mathcal{Q}_{\text{new}}). \quad (5)$$

As mentioned earlier, this minimisation is achieved by matching moments of \mathcal{Q}_{new} to those of \hat{p} .

2.4. Update the factor

We can see the relation of the f_i , which we wish to approximate, to \mathcal{Q}_{new} by combining formulas (3) and (5):

$$\mathcal{Q}_{\text{new}} \approx \hat{p} = \frac{1}{Z_i} f_i \mathcal{Q}^{\setminus i}. \quad (6)$$

From here, we easily obtain the formula for the approximation of f_i :

$$f_i \approx \tilde{f}_i = Z_i \frac{\mathcal{Q}_{\text{new}}}{\mathcal{Q}^{\setminus i}}. \quad (7)$$

Thanks to the right hand side of Eq. (7) consisting of a division of distributions from the approximating family (and a coefficient), \tilde{f}_i will also be from that family (provided it is closed under division). Now, the approximate factor gets updated according to Eq. (7), and the outer loop is repeated.

3. Evaluate the normalisation constant

After the algorithm has converged to a set of factors $\{\tilde{f}_i\}$, an approximate posterior $p(\mathbf{z}, \mathbf{e})$ can be computed as a product of the factors, according to the assumptions. If we are interested in $p(\mathbf{e})$, the model evidence, it can be computed now as

$$p(\mathbf{e}) = \int_{\mathbf{Z}} \prod_i \tilde{f}_i(\mathbf{z}) d\mathbf{z}. \quad (8)$$

3 The Clutter Problem

In the Clutter problem, we assume a sequence of d -dimensional i.i.d. observations being generated either from a normal distribution with an unknown mean with some probability, or from the “clutter” distribution. The model is specified by the following formulas:

$$W_i \stackrel{\text{i.i.d.}}{\sim} \text{Bern}(w_0) \quad (9)$$

$$\mathbf{x}_i \mid \boldsymbol{\mu} \stackrel{\text{ind.}}{\sim} W_i \mathcal{N}(\mathbf{0}_d, a\mathbf{I}_d) + (1 - W_i) \mathcal{N}(\boldsymbol{\mu}, \mathbf{I}_d) \quad (10)$$

The w_0 parameter determines the *proportion of clutter*, W_i select for each observation whether it was generated from the distribution of interest, $\mathcal{N}(\boldsymbol{\mu}, \mathbf{I}_d)$, or the clutter, and finally, $\boldsymbol{\mu}$ is the unknown mean of the

distribution we are trying to estimate. When learning the model, we will not learn W_i explicitly for each i , but rather treat the observations as identically distributed with the same proportion of clutter w_0 :

$$p(\mathbf{x} \mid \boldsymbol{\mu}) = w_0 \mathcal{N}(\mathbf{0}_d, a\mathbf{I}_d) + (1 - w_0) \mathcal{N}(\boldsymbol{\mu}, \mathbf{I}_d). \quad (11)$$

Finally, we adopt a broad Gaussian prior on $\boldsymbol{\mu}$:

$$\boldsymbol{\mu} \sim \mathcal{N}(\mathbf{0}_d, b\mathbf{I}_d). \quad (12)$$

TODO: draw the graphical model

This problem fits nicely the assumptions for EP:

1. It is intractable to do exact inference to find the value of $\boldsymbol{\mu}$. This is due to the fact that in the Bayesian network, the node for $\boldsymbol{\mu}$ has $(N + 1)$ independent parent nodes, a prior and the N likelihood factors, of which the N likelihood factors have 2 Gaussian components each. This results in the posterior for $\boldsymbol{\mu}$ consisting of 2^N N -dimensional Gaussians, corresponding to the 2^N subsets of observations that could have been generated from the true distribution (as opposed to the clutter).
2. The posterior is a product of factors that depend on a non-empty subset of the latent variables (which is $\{\boldsymbol{\mu}\}$ in this case) and a subset of the observed variables (either $\{\mathbf{x}_i\}$ for the likelihood factors, or \emptyset for the prior) – exactly as required.

Instantiating the general Eq. (1) for the Clutter problem, we get the following:

$$p((\boldsymbol{\mu}), (\mathbf{x}_1, \dots, \mathbf{x}_N)) = p(\boldsymbol{\mu}) \cdot \prod_{i=1}^N p(\mathbf{x}_i \mid \boldsymbol{\mu}). \quad (13)$$

In Eq. (13), the generic f_0 is instantiated as the prior $p(\boldsymbol{\mu})$, and the generic $f_i, i = 1, \dots$ as the likelihood $p(\mathbf{x}_i \mid \boldsymbol{\mu})$. In the following, we may use one or the other notation, whichever is more convenient.

We choose to approximate the factors, and hence also their product, by (unnormalised) spherical Gaussians, with one stipulation: the factors approximating the likelihoods may have their σ^2 parameter negative. This is an inherent property of the algorithm, and we discuss it later in Section 4. Still, each factor \tilde{f}_i can be represented by the triple $\langle \tilde{s}_i, \tilde{\mathbf{m}}_i, \tilde{v}_i \rangle$, describing its scale ($\int_{\mathbf{z}} \tilde{f}_i(\mathbf{z}) d\mathbf{z}$), mean, and variance, respectively:

$$\tilde{f}_i = \tilde{s}_i \mathcal{N}(\tilde{\mathbf{m}}_i, \tilde{v}_i \mathbf{I}_d). \quad (14)$$

Besides that, also the approximate posterior has the same form, and we shall denote its parameters as follows:

$$\mathcal{Q} = \mathcal{N}(\mathbf{m}, v\mathbf{I}_d). \quad (15)$$

Note that \mathcal{Q} is an approximating *distribution*, i.e. it is normalised to 1.

Since \tilde{f}_0 , the prior, already is a spherical Gaussian, its parameters can be set as part of initialisation:

$$\tilde{s}_0 = 1 \quad \tilde{\mathbf{m}}_0 = \mathbf{0}_d \quad \tilde{v}_0 = b. \quad (16)$$

This factor is exact and need not be updated anymore.

What remains is expressing the formulas (4), (5), (7) for a factor $\tilde{f}_i, i = 1, \dots$, and (8). The following paragraphs are devoted to this.

Update formula for the cavity distribution

The general formula is as follows:

$$\mathcal{Q}^{\setminus i} \propto \mathcal{Q} / \tilde{f}_i. \quad (4 - \text{repeated})$$

After substituting the values of \mathcal{Q} and \tilde{f}_i , represented as shown in Eqs. (15) and (14), respectively, we obtain the following:

$$\mathcal{Q}^{\setminus i} \propto \frac{\mathcal{N}(\mathbf{m}, v\mathbf{I}_d)}{\tilde{s}_i \mathcal{N}(\tilde{\mathbf{m}}_i, \tilde{v}_i\mathbf{I}_d)}. \quad (17)$$

The parameters of $\mathcal{Q}^{\setminus i}$ can be computed using the formula for the ratio of Gaussians,

$$\mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1) / \mathcal{N}(\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2) = C \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \quad (18)$$

where

$$\begin{aligned} \boldsymbol{\Sigma} &= (\boldsymbol{\Sigma}_1^{-1} - \boldsymbol{\Sigma}_2^{-1})^{-1} \\ \boldsymbol{\mu} &= \boldsymbol{\Sigma} (\boldsymbol{\Sigma}_1^{-1} \boldsymbol{\mu}_1 - \boldsymbol{\Sigma}_2^{-1} \boldsymbol{\mu}_2) \\ C &= \sqrt{\frac{|\boldsymbol{\Sigma}| |\boldsymbol{\Sigma}_2|}{(2\pi)^d |\boldsymbol{\Sigma}_1|}} \exp \left\{ -\frac{1}{2} (\boldsymbol{\mu}_1^T \boldsymbol{\Sigma}_1^{-1} \boldsymbol{\mu}_1 - \boldsymbol{\mu}_2^T \boldsymbol{\Sigma}_2^{-1} \boldsymbol{\mu}_2 - \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}) \right\}. \end{aligned} \quad (19)$$

As the result, we can express $\mathcal{Q}^{\setminus i}$ in terms of its parameters $\mathbf{m}^{\setminus i}$ (mean) and $v^{\setminus i}$ ($v^{\setminus i}\mathbf{I}_d$ being the variance-covariance matrix) as follows:

$$\mathbf{m}^{\setminus i} = v^{\setminus i}(\mathbf{m}v^{-1} - \tilde{\mathbf{m}}_i\tilde{v}_i^{-1}) \quad v^{\setminus i} = (v^{-1} - \tilde{v}_i^{-1})^{-1}. \quad (20)$$

Update formula for \mathcal{Q}

In computing \mathcal{Q}_{new} according to Eq. (5), we have to compute \hat{p} and then its first and second moment in order to arrive at the spherical normal distribution minimising the KL-divergence to \hat{p} . In the definition of \hat{p} in Eq. (3), the quantity Z_i is yet to be computed. It is the normalisation constant of $f_i \mathcal{Q}^{\setminus i}$, i.e.:

$$Z_i = \int_{\mathbb{R}^d} f_i(\boldsymbol{\mu}) \mathcal{Q}^{\setminus i}(\boldsymbol{\mu}) d\boldsymbol{\mu}. \quad (21)$$

Parameters of $\mathcal{Q}^{\setminus i}$ were obtained in the previous step, and f_i was defined as the likelihood for \mathbf{x}_i (cf. Eq. (11)):

$$f_i(\boldsymbol{\mu}) = w_0 \mathcal{N}(\mathbf{x}_i; \mathbf{0}_d, a\mathbf{I}_d) + (1 - w_0) \mathcal{N}(\mathbf{x}_i; \boldsymbol{\mu}, \mathbf{I}_d). \quad (22)$$

Substituting into Eq. (21), we get

$$Z_i = \int_{\mathbb{R}^d} [w_0 \mathcal{N}(\mathbf{x}_i; \mathbf{0}_d, a\mathbf{I}_d) + (1 - w_0) \mathcal{N}(\mathbf{x}_i; \boldsymbol{\mu}, \mathbf{I}_d)] \mathcal{N}(\boldsymbol{\mu}; \mathbf{m}^{\setminus i}, v^{\setminus i}\mathbf{I}_d) d\boldsymbol{\mu} \quad (23)$$

$$= \int_{\mathbb{R}^d} w_0 \mathcal{N}(\mathbf{x}_i; \mathbf{0}_d, a\mathbf{I}_d) \mathcal{N}(\boldsymbol{\mu}; \mathbf{m}^{\setminus i}, v^{\setminus i}\mathbf{I}_d) d\boldsymbol{\mu} \quad (24)$$

$$+ \int_{\mathbb{R}^d} (1 - w_0) \mathcal{N}(\mathbf{x}_i; \boldsymbol{\mu}, \mathbf{I}_d) \mathcal{N}(\boldsymbol{\mu}; \mathbf{m}^{\setminus i}, v^{\setminus i}\mathbf{I}_d) d\boldsymbol{\mu}$$

$$= w_0 \mathcal{N}(\mathbf{x}_i; \mathbf{0}_d, a\mathbf{I}_d) + (1 - w_0) \int_{\mathbb{R}^d} \mathcal{N}(\mathbf{x}_i - \boldsymbol{\mu}; \mathbf{0}_d, \mathbf{I}_d) \mathcal{N}(\boldsymbol{\mu}; \mathbf{m}^{\setminus i}, v^{\setminus i}\mathbf{I}_d) d\boldsymbol{\mu} \quad (25)$$

$$= w_0 \mathcal{N}(\mathbf{x}_i; \mathbf{0}_d, a\mathbf{I}_d) + (1 - w_0) \mathcal{N}(\mathbf{x}_i; \mathbf{m}^{\setminus i}, (v^{\setminus i} + 1)\mathbf{I}_d) \quad (26)$$

where, going from (25) to (26), we used the result [1] about convolution of Gaussians.

The mean value and variance of \hat{p} can be derived for a general form of the factor f_i . Hence, we will simplify the next derivations by rewriting \hat{p} in the following form:

$$\hat{p}(\boldsymbol{\mu}) = \frac{1}{Z(\mathbf{m}, \boldsymbol{\Sigma})} f(\boldsymbol{\mu}) \mathcal{N}(\boldsymbol{\mu}; \mathbf{m}, \boldsymbol{\Sigma}) \quad (27)$$

where

$$Z(\mathbf{m}, \boldsymbol{\Sigma}) = \int_{\mathbb{R}^d} f(\boldsymbol{\mu}) \mathcal{N}(\boldsymbol{\mu}; \mathbf{m}, \boldsymbol{\Sigma}) d\boldsymbol{\mu}. \quad (28)$$

The two moments will be found from derivatives of Z :

$$\frac{dZ(\mathbf{m}, \boldsymbol{\Sigma})}{d\mathbf{m}} = \int_{\mathbb{R}^d} \frac{d}{d\mathbf{m}} (f(\boldsymbol{\mu}) \mathcal{N}(\boldsymbol{\mu}; \mathbf{m}, \boldsymbol{\Sigma})) d\boldsymbol{\mu} \quad (29)$$

$$= \int_{\mathbb{R}^d} \frac{f(\boldsymbol{\mu})}{\sqrt{(2\pi)^d |\boldsymbol{\Sigma}|}} \exp \left\{ -\frac{1}{2} (\boldsymbol{\mu} - \mathbf{m})^T \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu} - \mathbf{m}) \right\} ((\boldsymbol{\mu} - \mathbf{m})^T \boldsymbol{\Sigma}^{-1}) d\boldsymbol{\mu} \quad (30)$$

$$= \int_{\mathbb{R}^d} Z(\mathbf{m}, \boldsymbol{\Sigma}) \hat{p}(\boldsymbol{\mu}) ((\boldsymbol{\mu} - \mathbf{m})^T \boldsymbol{\Sigma}^{-1}) d\boldsymbol{\mu} \quad (31)$$

$$= Z(\mathbf{m}, \boldsymbol{\Sigma}) \left(\int_{\mathbb{R}^d} \boldsymbol{\mu} \hat{p}(\boldsymbol{\mu}) d\boldsymbol{\mu} - \int_{\mathbb{R}^d} \mathbf{m} \hat{p}(\boldsymbol{\mu}) d\boldsymbol{\mu} \right)^T \boldsymbol{\Sigma}^{-1} \quad (32)$$

$$= Z(\mathbf{m}, \boldsymbol{\Sigma}) \cdot (\mathbb{E}_{\hat{p}}[\boldsymbol{\mu}] - \mathbf{m})^T \boldsymbol{\Sigma}^{-1}, \quad (33)$$

$$\frac{dZ(\mathbf{m}, \boldsymbol{\Sigma})}{d\boldsymbol{\Sigma}} = \int_{\mathbb{R}^d} \frac{d}{d\boldsymbol{\Sigma}} (f(\boldsymbol{\mu}) \mathcal{N}(\boldsymbol{\mu}; \mathbf{m}, \boldsymbol{\Sigma})) d\boldsymbol{\mu} \quad (34)$$

$$= \int_{\mathbb{R}^d} \frac{f(\boldsymbol{\mu})}{\sqrt{(2\pi)^d |\boldsymbol{\Sigma}|}} \exp \left\{ -\frac{1}{2} (\boldsymbol{\mu} - \mathbf{m})^T \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu} - \mathbf{m}) \right\} \left(\frac{1}{2} \boldsymbol{\Sigma}^{-T} (\boldsymbol{\mu} - \mathbf{m}) (\boldsymbol{\mu} - \mathbf{m})^T \boldsymbol{\Sigma}^{-T} \right) \\ - \frac{1}{2} \frac{f(\boldsymbol{\mu})}{\sqrt{(2\pi)^d |\boldsymbol{\Sigma}|}} \boldsymbol{\Sigma}^{-T} \exp \left\{ -\frac{1}{2} (\boldsymbol{\mu} - \mathbf{m})^T \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu} - \mathbf{m}) \right\} d\boldsymbol{\mu} \quad (35)$$

$$= \frac{1}{2} \boldsymbol{\Sigma}^{-1} \left[\int_{\mathbb{R}^d} \boldsymbol{\mu} \boldsymbol{\mu}^T Z(\mathbf{m}, \boldsymbol{\Sigma}) \hat{p}(\boldsymbol{\mu}) d\boldsymbol{\mu} - \int_{\mathbb{R}^d} \boldsymbol{\mu} \mathbf{m}^T Z(\mathbf{m}, \boldsymbol{\Sigma}) \hat{p}(\boldsymbol{\mu}) d\boldsymbol{\mu} \right. \\ \left. - \int_{\mathbb{R}^d} \mathbf{m} \boldsymbol{\mu}^T Z(\mathbf{m}, \boldsymbol{\Sigma}) \hat{p}(\boldsymbol{\mu}) d\boldsymbol{\mu} + \int_{\mathbb{R}^d} \mathbf{m} \mathbf{m}^T Z(\mathbf{m}, \boldsymbol{\Sigma}) \hat{p}(\boldsymbol{\mu}) d\boldsymbol{\mu} \right] \boldsymbol{\Sigma}^{-1} \\ - \frac{1}{2} \int_{\mathbb{R}^d} \boldsymbol{\Sigma}^{-1} Z(\mathbf{m}, \boldsymbol{\Sigma}) \hat{p}(\boldsymbol{\mu}) d\boldsymbol{\mu} \quad (36)$$

$$= Z(\mathbf{m}, \boldsymbol{\Sigma}) \cdot \left\{ \frac{1}{2} \boldsymbol{\Sigma}^{-1} [\mathbb{E}_{\hat{p}}[\boldsymbol{\mu} \boldsymbol{\mu}^T] - \mathbb{E}_{\hat{p}}[\boldsymbol{\mu}] \mathbf{m}^T - \mathbf{m} \mathbb{E}_{\hat{p}}[\boldsymbol{\mu}]^T + \mathbf{m} \mathbf{m}^T] \boldsymbol{\Sigma}^{-1} - \frac{1}{2} \boldsymbol{\Sigma}^{-1} \right\} \quad (37)$$

where \mathbf{X}^{-T} is a shorthand for $(\mathbf{X}^{-1})^T (= (\mathbf{X}^T)^{-1})$ and we applied matrix calculus results from [2].

Eqs. (33) and (37) give us formulas for the moments we are interested in. However, they include the term $Z(\mathbf{m}, \boldsymbol{\Sigma})$, which is an inconvenient integral to compute. Taking the derivative of the log instead will

get us rid of this term:

$$\frac{d \log Z(\mathbf{m}, \Sigma)}{d\mathbf{m}} = \frac{1}{Z(\mathbf{m}, \Sigma)} \frac{dZ(\mathbf{m}, \Sigma)}{d\mathbf{m}} = (\mathbb{E}_{\hat{p}}[\boldsymbol{\mu}] - \mathbf{m})^T \Sigma^{-1} \quad (38)$$

$$\begin{aligned} \frac{d \log Z(\mathbf{m}, \Sigma)}{d\Sigma} &= \frac{1}{Z(\mathbf{m}, \Sigma)} \frac{d \log Z(\mathbf{m}, \Sigma)}{d\Sigma} \\ &= \frac{1}{2} \Sigma^{-1} [\mathbb{E}_{\hat{p}}[\boldsymbol{\mu}\boldsymbol{\mu}^T] - \mathbb{E}_{\hat{p}}[\boldsymbol{\mu}]\mathbf{m}^T - \mathbf{m}\mathbb{E}_{\hat{p}}[\boldsymbol{\mu}]^T + \mathbf{m}\mathbf{m}^T] \Sigma^{-1} - \frac{1}{2} \Sigma^{-1} \end{aligned} \quad (39)$$

The first and second moment are now obtained easily from Eqs. (38) and (39) by shuffling them a bit. We get:

$$\mathbb{E}_{\hat{p}}[\boldsymbol{\mu}] = \mathbf{m} + \Sigma \left(\frac{d \log Z(\mathbf{m}, \Sigma)}{d\mathbf{m}} \right)^T \quad (40)$$

$$\begin{aligned} \mathbb{E}_{\hat{p}}[\boldsymbol{\mu}\boldsymbol{\mu}^T] - \mathbb{E}_{\hat{p}}[\boldsymbol{\mu}]\mathbb{E}_{\hat{p}}[\boldsymbol{\mu}]^T &= \Sigma \left(2 \frac{d \log Z(\mathbf{m}, \Sigma)}{d\Sigma} + \Sigma^{-1} \right) \Sigma - \left[-\mathbb{E}_{\hat{p}}[\boldsymbol{\mu}]\mathbf{m}^T - \mathbf{m}\mathbb{E}_{\hat{p}}[\boldsymbol{\mu}]^T + \mathbf{m}\mathbf{m} \right] \\ &\quad - \mathbb{E}_{\hat{p}}[\boldsymbol{\mu}]\mathbb{E}_{\hat{p}}[\boldsymbol{\mu}]^T \end{aligned} \quad (41)$$

$$\begin{aligned} &= 2\Sigma \frac{d \log Z(\mathbf{m}, \Sigma)}{d\Sigma} \Sigma + \Sigma \\ &\quad - \left[(\mathbb{E}_{\hat{p}}[\boldsymbol{\mu}] - \mathbf{m})(\mathbb{E}_{\hat{p}}[\boldsymbol{\mu}] - \mathbf{m})^T - \mathbb{E}_{\hat{p}}[\boldsymbol{\mu}]\mathbb{E}_{\hat{p}}[\boldsymbol{\mu}]^T \right] - \mathbb{E}_{\hat{p}}[\boldsymbol{\mu}]\mathbb{E}_{\hat{p}}[\boldsymbol{\mu}]^T \end{aligned} \quad (42)$$

$$= 2\Sigma \frac{d \log Z(\mathbf{m}, \Sigma)}{d\Sigma} \Sigma + \Sigma - \Sigma \left(\frac{d \log Z(\mathbf{m}, \Sigma)}{d\mathbf{m}} \right)^T \frac{d \log Z(\mathbf{m}, \Sigma)}{d\mathbf{m}} \Sigma^T \quad (43)$$

$$= \Sigma \left[2 \frac{d \log Z(\mathbf{m}, \Sigma)}{d\Sigma} - \left(\frac{d \log Z(\mathbf{m}, \Sigma)}{d\mathbf{m}} \right)^T \frac{d \log Z(\mathbf{m}, \Sigma)}{d\mathbf{m}} \right] \Sigma + \Sigma \quad (44)$$

Now, what remains to be computed in order to arrive at the KL-divergence minimiser are the derivatives

of $\log Z$:

$$\frac{d \log Z_i}{d \mathbf{m}^{\setminus i}} = \frac{d \log [w_0 \mathcal{N}(\mathbf{x}_i; \mathbf{0}_d, a \mathbf{I}_d) + (1 - w_0) \mathcal{N}(\mathbf{x}_i; \mathbf{m}^{\setminus i}, (v^{\setminus i} + 1) \mathbf{I}_d)]}{d \mathbf{m}^{\setminus i}} \quad (45)$$

$$= \frac{1}{w_0 \mathcal{N}(\mathbf{x}_i; \mathbf{0}_d, a \mathbf{I}_d) + (1 - w_0) \mathcal{N}(\mathbf{x}_i; \mathbf{m}^{\setminus i}, (v^{\setminus i} + 1) \mathbf{I}_d)} \cdot \frac{(1 - w_0) d \mathcal{N}(\mathbf{x}_i; \mathbf{m}^{\setminus i}, (v^{\setminus i} + 1) \mathbf{I}_d)}{d \mathbf{m}^{\setminus i}} \quad (46)$$

$$= \frac{1}{w_0 \mathcal{N}(\mathbf{x}_i; \mathbf{0}_d, a \mathbf{I}_d) + (1 - w_0) \mathcal{N}(\mathbf{x}_i; \mathbf{m}^{\setminus i}, (v^{\setminus i} + 1) \mathbf{I}_d)} \cdot \frac{\frac{1-w_0}{\sqrt{(2\pi)^d |(v^{\setminus i} + 1) \mathbf{I}_d|}} d \exp \left(-\frac{1}{2} (\mathbf{x}_i - \mathbf{m}^{\setminus i})^T ((v^{\setminus i} + 1) \mathbf{I}_d)^{-1} (\mathbf{x}_i - \mathbf{m}^{\setminus i}) \right)}{d \mathbf{m}^{\setminus i}} \quad (47)$$

$$= \frac{\frac{1-w_0}{\sqrt{(2\pi)^d |(v^{\setminus i} + 1) \mathbf{I}_d|}} \exp \left(-\frac{1}{2} (\mathbf{x}_i - \mathbf{m}^{\setminus i})^T ((v^{\setminus i} + 1) \mathbf{I}_d)^{-1} (\mathbf{x}_i - \mathbf{m}^{\setminus i}) \right)}{w_0 \mathcal{N}(\mathbf{x}_i; \mathbf{0}_d, a \mathbf{I}_d) + (1 - w_0) \mathcal{N}(\mathbf{x}_i; \mathbf{m}^{\setminus i}, (v^{\setminus i} + 1) \mathbf{I}_d)} \cdot \frac{d \left(-\frac{1}{2} (\mathbf{x}_i - \mathbf{m}^{\setminus i})^T ((v^{\setminus i} + 1) \mathbf{I}_d)^{-1} (\mathbf{x}_i - \mathbf{m}^{\setminus i}) \right)}{d \mathbf{m}^{\setminus i}} \quad (48)$$

$$= \frac{(1 - w_0) \mathcal{N}(\mathbf{x}_i; \mathbf{m}^{\setminus i}, (v^{\setminus i} + 1) \mathbf{I}_d)}{w_0 \mathcal{N}(\mathbf{x}_i; \mathbf{0}_d, a \mathbf{I}_d) + (1 - w_0) \mathcal{N}(\mathbf{x}_i; \mathbf{m}^{\setminus i}, (v^{\setminus i} + 1) \mathbf{I}_d)} \cdot \left(\mathbf{x}_i - \mathbf{m}^{\setminus i} \right)^T \left((v^{\setminus i} + 1) \mathbf{I}_d \right)^{-1} \quad (49)$$

Let us simplify the expression by introducing r as the probability of \mathbf{x}_i not being generated from the clutter, and realising that multiplication by the last term is equivalent to division by $(v^{\setminus i} + 1)$:

$$r := \frac{(1 - w_0) \mathcal{N}(\mathbf{x}_i; \mathbf{m}^{\setminus i}, (v^{\setminus i} + 1) \mathbf{I}_d)}{w_0 \mathcal{N}(\mathbf{x}_i; \mathbf{0}_d, a \mathbf{I}_d) + (1 - w_0) \mathcal{N}(\mathbf{x}_i; \mathbf{m}^{\setminus i}, (v^{\setminus i} + 1) \mathbf{I}_d)} \quad (50)$$

$$\frac{d \log Z_i}{d \mathbf{m}^{\setminus i}} = r \frac{(\mathbf{x}_i - \mathbf{m}^{\setminus i})^T}{v^{\setminus i} + 1} \quad (51)$$

The derivative of $\log Z_i$ by the variance parameter is obtained similarly (let Σ denote the second param-

eter of Z_i , which has the value $v^{\setminus i} \mathbf{I}_d$:

$$\frac{d \log Z_i}{d \Sigma} = \frac{d \log [w_0 \mathcal{N}(\mathbf{x}_i; \mathbf{0}_d, a \mathbf{I}_d) + (1 - w_0) \mathcal{N}(\mathbf{x}_i; \mathbf{m}^{\setminus i}, \Sigma + \mathbf{I}_d)]}{d \Sigma} \quad (52)$$

$$= \frac{1}{w_0 \mathcal{N}(\mathbf{x}_i; \mathbf{0}_d, a \mathbf{I}_d) + (1 - w_0) \mathcal{N}(\mathbf{x}_i; \mathbf{m}^{\setminus i}, \Sigma + \mathbf{I}_d)} \cdot \frac{(1 - w_0) d \mathcal{N}(\mathbf{x}_i; \mathbf{m}^{\setminus i}, \Sigma + \mathbf{I}_d)}{d \Sigma} \quad (53)$$

$$= \frac{(1 - w_0)(2\pi)^{-d/2}}{w_0 \mathcal{N}(\mathbf{x}_i; \mathbf{0}_d, a \mathbf{I}_d) + (1 - w_0) \mathcal{N}(\mathbf{x}_i; \mathbf{m}^{\setminus i}, \Sigma + \mathbf{I}_d)} \cdot \frac{d \left[|\Sigma + \mathbf{I}_d|^{-1/2} \exp \left(-\frac{1}{2} (\mathbf{x}_i - \mathbf{m}^{\setminus i})^T (\Sigma + \mathbf{I}_d)^{-1} (\mathbf{x}_i - \mathbf{m}^{\setminus i}) \right) \right]}{d \Sigma} \quad (54)$$

$$= \frac{(1 - w_0)(2\pi)^{-d/2}}{w_0 \mathcal{N}(\mathbf{x}_i; \mathbf{0}_d, a \mathbf{I}_d) + (1 - w_0) \mathcal{N}(\mathbf{x}_i; \mathbf{m}^{\setminus i}, \Sigma + \mathbf{I}_d)} \cdot \left\{ |\Sigma + \mathbf{I}_d|^{-1/2} \exp \left[-\frac{1}{2} (\mathbf{x}_i - \mathbf{m}^{\setminus i})^T (\Sigma + \mathbf{I}_d)^{-1} (\mathbf{x}_i - \mathbf{m}^{\setminus i}) \right] \cdot \left(\frac{1}{2} (\Sigma + \mathbf{I}_d)^{-T} (\mathbf{x}_i - \mathbf{m}^{\setminus i}) (\mathbf{x}_i - \mathbf{m}^{\setminus i})^T (\Sigma + \mathbf{I}_d)^{-T} \right) - \frac{1}{2} |\Sigma + \mathbf{I}_d|^{-1/2} (\Sigma + \mathbf{I}_d)^{-T} \exp \left[-\frac{1}{2} (\mathbf{x}_i - \mathbf{m}^{\setminus i})^T (\Sigma + \mathbf{I}_d)^{-1} (\mathbf{x}_i - \mathbf{m}^{\setminus i}) \right] \right\} \quad (55)$$

$$= \frac{(1 - w_0)}{w_0 \mathcal{N}(\mathbf{x}_i; \mathbf{0}_d, a \mathbf{I}_d) + (1 - w_0) \mathcal{N}(\mathbf{x}_i; \mathbf{m}^{\setminus i}, \Sigma + \mathbf{I}_d)} \cdot \left[\mathcal{N}(\mathbf{x}_i; \mathbf{m}^{\setminus i}, \Sigma + \mathbf{I}_d) \cdot \frac{(\mathbf{x}_i - \mathbf{m}^{\setminus i})(\mathbf{x}_i - \mathbf{m}^{\setminus i})^T}{2(v^{\setminus i} + 1)^2} - \mathcal{N}(\mathbf{x}_i; \mathbf{m}^{\setminus i}, \Sigma + \mathbf{I}_d) \cdot \frac{\mathbf{I}_d}{2(v^{\setminus i} + 1)} \right] \quad (56)$$

$$= \frac{r}{2(v^{\setminus i} + 1)^2} \cdot \left[(\mathbf{x}_i - \mathbf{m}^{\setminus i})(\mathbf{x}_i - \mathbf{m}^{\setminus i})^T - (v^{\setminus i} + 1) \mathbf{I}_d \right] \quad (57)$$

Substituting into Eqs. (40) and (44), we finally arrive at the new parameters of \mathcal{Q} , \mathbf{m}_{new} (mean) and Σ_{new} (variance):

$$\mathbf{m}_{\text{new}} = \mathbf{m}^{\setminus i} + \Sigma \left(\frac{d \log Z_i}{d \mathbf{m}^{\setminus i}} \right)^T = \mathbf{m}^{\setminus i} + \Sigma r \frac{\mathbf{x}_i - \mathbf{m}^{\setminus i}}{v^{\setminus i} + 1} = \mathbf{m}^{\setminus i} + r \frac{v^{\setminus i}}{v^{\setminus i} + 1} (\mathbf{x}_i - \mathbf{m}^{\setminus i}) \quad (58)$$

$$\Sigma_{\text{new}} = \Sigma \left\{ 2 \frac{d \log Z_i}{d \Sigma} - \left(\frac{d \log Z_i}{d \mathbf{m}^{\setminus i}} \right)^T \frac{d \log Z_i}{d \mathbf{m}^{\setminus i}} \right\} \Sigma + \Sigma \quad (59)$$

$$= \Sigma \left\{ \frac{r}{(v^{\setminus i} + 1)^2} \cdot \left[(\mathbf{x}_i - \mathbf{m}^{\setminus i})(\mathbf{x}_i - \mathbf{m}^{\setminus i})^T - (v^{\setminus i} + 1) \mathbf{I}_d \right] - \left(r \frac{(\mathbf{x}_i - \mathbf{m}^{\setminus i})^T}{v^{\setminus i} + 1} \right)^T r \frac{(\mathbf{x}_i - \mathbf{m}^{\setminus i})^T}{v^{\setminus i} + 1} \right\} \Sigma + \Sigma \quad (60)$$

$$= r \left(\frac{v^{\setminus i}}{v^{\setminus i} + 1} \right)^2 \cdot \left[(1 - r)(\mathbf{x}_i - \mathbf{m}^{\setminus i})(\mathbf{x}_i - \mathbf{m}^{\setminus i})^T - (v^{\setminus i} + 1) \mathbf{I}_d \right] + \Sigma \quad (61)$$

where we again used the symbol Σ to denote $v^{-1}\mathbf{I}_d$.

However, this Σ_{new} is generally not a variance matrix of a spherical normal, which is the form we assume for the posterior distribution. Hence, we need to find the KL-divergence minimiser of a spherical normal from a normal with the general covariance Σ .¹ Let us solve this problem now, denoting the general multivariate normal with \mathcal{Q} and the spherical one with \mathcal{S} , and assuming the mean $\mathbf{0}$ for both, WLOG:

$$\arg \min_v \text{KL}(\mathcal{Q} \parallel \mathcal{S}) = \arg \min_v \int_{\mathbb{R}^d} \mathcal{N}(\mathbf{x}; \mathbf{0}, \Sigma) \log \frac{\mathcal{N}(\mathbf{x}; \mathbf{0}, \Sigma)}{\mathcal{N}(\mathbf{x}; \mathbf{0}, v\mathbf{I}_d)} d\mathbf{x} \quad (62)$$

$$= \arg \min_v \int_{\mathbb{R}^d} \mathcal{N}(\mathbf{x}; \mathbf{0}, \Sigma) \log \frac{\sqrt{|v\mathbf{I}_d|} \exp\left\{-\frac{1}{2}\mathbf{x}^T \Sigma^{-1} \mathbf{x}\right\}}{\sqrt{|\Sigma|} \exp\left\{-\frac{1}{2}\mathbf{x}^T v^{-1} \mathbf{I}_d \mathbf{x}\right\}} d\mathbf{x} \quad (63)$$

$$= \arg \min_v \int_{\mathbb{R}^d} \mathcal{N}(\mathbf{x}; \mathbf{0}, \Sigma) \left[\frac{d}{2} \log v - \frac{1}{2} \mathbf{x}^T \Sigma^{-1} \mathbf{x} - \frac{1}{2} \log |\Sigma| + \frac{1}{2} \mathbf{x}^T v^{-1} \mathbf{I}_d \mathbf{x} \right] d\mathbf{x} \quad (64)$$

$$= \arg \min_v \int_{\mathbb{R}^d} \mathcal{N}(\mathbf{x}; \mathbf{0}, \Sigma) (d \log v + \mathbf{x}^T \mathbf{x} / v) d\mathbf{x} \quad (65)$$

$$=: \arg \min_v G(v) \quad (66)$$

Because the function we minimise here is smooth for $v > 0$, we shall find the minimum by setting the derivative equal to zero:

$$0 = \frac{dG}{dv}(v^*) \quad (67)$$

$$= \left(\frac{d}{dv} \int_{\mathbb{R}^d} \mathcal{N}(\mathbf{x}; \mathbf{0}, \Sigma) (d \log v + \mathbf{x}^T \mathbf{x} / v) d\mathbf{x} \right) (v^*) \quad (68)$$

$$= \left(\int_{\mathbb{R}^d} \frac{d}{dv} \mathcal{N}(\mathbf{x}; \mathbf{0}, \Sigma) (d \log v + \mathbf{x}^T \mathbf{x} / v) d\mathbf{x} \right) (v^*) \quad (69)$$

$$= \int_{\mathbb{R}^d} \mathcal{N}(\mathbf{x}; \mathbf{0}, \Sigma) (d/v - \mathbf{x}^T \mathbf{x} / v^2) d\mathbf{x} \Big|_{v=v^*} \quad (70)$$

$$= \frac{d}{v} - \frac{1}{v^2} \mathbb{E}_{\mathcal{Q}}[\mathbf{x}^T \mathbf{x}] \Big|_{v=v^*} \quad (71)$$

Using the following identity for the product of a quadratic form with a Gaussian density function,

$$\int_{\mathbb{R}^d} (\mathbf{x} - \mathbf{x}_0)^T \mathbf{F}^{-1} (\mathbf{x} - \mathbf{x}_0) \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \Sigma) d\mathbf{x} = (\mathbf{x}_0 - \boldsymbol{\mu})^T \mathbf{F}^{-1} (\mathbf{x}_0 - \boldsymbol{\mu}) + \text{Tr}[\mathbf{F}^{-1} \Sigma], \quad (72)$$

we can express the minimiser v^* from Eq. (71):

$$v^* = \text{Tr}[\Sigma] / d. \quad (73)$$

This result can be combined with Eq. (61) to give us the updated variance of the spherical Gaussian

¹This Σ will be the Σ_{new} as given by Eq. (61).

posterior:

$$v_{\text{new}} = \text{Tr} \left[r \left(\frac{v^{i}}{v^{i} + 1} \right)^2 \cdot \left[(1 - r)(\mathbf{x}_i - \mathbf{m}^i)(\mathbf{x}_i - \mathbf{m}^i)^T - (v^{i} + 1)\mathbf{I}_d \right] + v^{i}\mathbf{I}_d \right] / d \quad (74)$$

$$= r \left(\frac{v^{i}}{v^{i} + 1} \right)^2 \cdot \left[(1 - r)\text{Tr} \left[(\mathbf{x}_i - \mathbf{m}^i)(\mathbf{x}_i - \mathbf{m}^i)^T \right] / d - (v^{i} + 1) \right] + v^{i} \quad (75)$$

$$= v^{i} - r \frac{(v^{i})^2}{v^{i} + 1} + \frac{r(1 - r)}{d} \left(\frac{v^{i}}{v^{i} + 1} \right)^2 \|\mathbf{x}_i - \mathbf{m}^i\|^2 \quad (76)$$

Eqs. (26), (58) and (76) give us the updated parameters for \mathcal{Q} , which was the objective of this step.

Update formula for \tilde{f}_i

Formula for the normalisation constant

TODO: how the results are read off

4 Implementation

TODO: how some parameters were instantiated, how uniform factors were expressed as Gaussians (infinite variance)

TODO: general properties of the implementation: Python, uses numpy, can be configured inside the source code, can be asked to draw plots interactively

TODO: example pictures from the algorithm

TODO: problems encountered (negative variance, infinities, unsensible normalisation constants)...but the algorithm converges well. Include some statistics of convergence properties (or, say, distance of the estimated mean from the true one)

References

- [1] BROMILEY, P. Products and convolutions of Gaussian distributions. Internal Report 2003-003, TINA Vision, 2003.
- [2] FACKLER, P. L. Notes on matrix calculus, 2005.
- [3] MINKA, T. P. Expectation Propagation for approximate Bayesian inference. In *Proceedings of the 17th Conference in Uncertainty in Artificial Intelligence* (2001), Morgan Kaufmann Publishers Inc., pp. 362–369.