

# Auto-Encoding with Stochastic Expectation Propagation in Latent Variable Models



**Vera Gangeskar Johne**

Department of Engineering

University of Cambridge

This dissertation is submitted for the degree of

*Master of Philosophy*

Fitzwilliam College

August 2016



## Declaration

I, Vera Gangeskar Johne of Fitzwilliam College, being a candidate for the M.Phil in Machine Learning, Speech and Language Technology, hereby declare that this report and the work described in it are my own work, unaided except as may be specified below, and that the report does not contain material that has already been used to any substantial extent for a comparable purpose. Total word count: 4702.



Vera Gangeskar Johne  
August 2016



## Acknowledgements

Professor Richard Turner is known to my fellow students and I for his clear and insightful explanations of challenging concepts. Coming from a computer science background into the degree, I decided to try to do my thesis under Richard to get a more mathematically complete understanding of statistical machine learning. I am so glad I did, thank you Richard!

Along with Rich, I was lucky to be supervised by Yingzhen Li. Yingzhen has that ability of never making the listener feel silly for asking a question. I have learnt a lot from you Yingzhen, thank you!

As I reflect on the past year, my 22 fellow course-mates and I have bonded through this unique experience as the first MLSALT students. Thank you for friendship and a stimulating environment: Both important ingredients for productivity.

I would also like to thank two professors from my undergraduate years whose mentorship has been invaluable on my academic journey thus far. First, Bachelor thesis supervisor Steve Hails whom encouraged me to apply, and who introduced me to modeling in the unknown. Secondly, Christophe Dessimoz who co-supervised me during an internship at the European Bioinformatics Institute where I got my first taste of ML.

Hans, my brother, whom introduced me to computer programming back in the day, and who challenges me to think individually.

Most of all I want to say thank you to my parents who have always encouraged free thinking, and exploration of knowledge.



# Abstract

Scalable Bayesian methods are becoming increasingly in demand. It is not uncommon to have datasets at peta- and exabyte scale. Moreover, datasets tend to be high dimensional and scattered with uncertainty. The Bayesian framework provides an elegant and principled system for reasoning under incomplete knowledge. We are interested in performing scalable inference and learning in directed probabilistic models where the latent space is both high dimensional and continuous, yielding an intractable posterior distribution. Seeing the success of variational methods in such settings through stochastic optimization and recognition models, a question arises naturally: Can similar techniques work in the Expectation Propagation framework?

This thesis proposes a new inference technique which utilizes Stochastic Expectation Propagation for model learning, and a recognition model (inference network) for posterior inference over local latent variables, resulting in constant space complexity. We derive a new local objective function that leverages from stochastic optimization. Auto-encoding-with-Stochastic-Expectation-Propagation (AESEP) demonstrates signs of successful learning for factor analysis.





# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Problem Scenario . . . . .	3
<b>2</b>	<b>Variational Inference</b>	<b>5</b>
2.1	Use of Reparameterization Trick . . . . .	6
2.2	Use of Recognition Models . . . . .	7
<b>3</b>	<b>Expectation Propagation</b>	<b>9</b>
3.1	EP . . . . .	9
3.2	SEP . . . . .	12
<b>4</b>	<b>Auto-encoding with SEP</b>	<b>14</b>
4.1	Approximating Factors and Parameter Tying . . . . .	14
4.2	Local KL-divergence . . . . .	15
4.3	Stochastic Objective . . . . .	16
4.4	Algorithm . . . . .	18
<b>5</b>	<b>Experiments</b>	<b>19</b>
5.1	Vanilla EP to Infer Latent Variables . . . . .	20
5.2	Auto-encoding with SEP . . . . .	20
<b>6</b>	<b>Conclusion</b>	<b>24</b>
	<b>References</b>	<b>25</b>
<b>A</b>	<b>Derivation of Variational Lowerbound</b>	<b>27</b>
<b>B</b>	<b>Moment Matching of EP</b>	<b>28</b>
<b>C</b>	<b>Exact EP on local latent variables of FA</b>	<b>29</b>



# 1 | Introduction

At the heart of many machine learning problems lies the task of pattern recognition: The task of unravelling hidden patterns in high-dimensional observations. One way to model this scenario is via writing a generative procedure. This involves assuming a latent space, and a set of generative rules relating the latent space to the observed space. For instance, a document (observed space) might be generated by first selecting a set of topics, then words from that topic (latent space) and combining them into sentences. In the literature such models are referred to as latent variable models (LVMs).

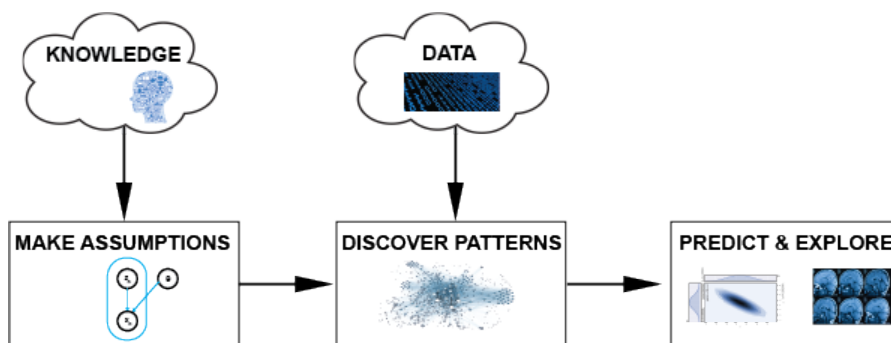


Figure 1.1: Probabilistic Pipeline.<sup>1</sup>

Often the goal in this scenario is inference, and the probabilistic framework offers a principled way for performing this task. In this pipeline (Figure 1.1), we use our prior knowledge to script a model, which is trained for discovering patterns with an available dataset. The trained model then finds its utility in making predictions about the unknown. The dynamics of this pipeline are governed by Bayesian probability, which has proven itself as the optimal language to reason under incomplete and uncertain information, [MacKay \(2003\)](#).

---

<sup>1</sup>Diagram inspired from Professor David Blei's talk *Scaling And Generalizing Approximate Bayesian Inference*, given 13/07/2016 in the department of Engineering, Cambridge University.

In practice, inference is often faced with the bottleneck of intractability. First is the analytical intractability, and second is the computational intractability. In the case of LVMs, these can stem from high-dimensional latent space, non-linear likelihoods, or large datasets.

In such circumstances, we must turn to approximate inference schemes. In fact, approximate inference techniques are often used in practice where exact inference is possible, to harvest their computational benefits. Over the past decade, the size of a typical dataset has dramatically escalated, and it is not uncommon to have streaming data (e.g. user clicks). This requires both quality analysis and computational efficiency. What is more, the goal is often decision making based on little information (e.g. Netflix movie suggestions to a new user). This has brought a wave of attention to scalable Bayesian learning, for example Stochastic Variational Inference, [Hoffman et al. \(2013\)](#), and Auto-encoding Variational Bayes (AEVB), [Kingma and Welling \(2013\)](#).

Approximate inference can be achieved by directly sampling from the intractable posterior (e.g. Markov Chain Monte Carlo), or through distributional approximations (e.g. Variational Inference and Expectation Propagation). Moreover, recent efforts have attempted to bridge the gap between the two schools of thought, [Salimans et al. \(2015\)](#). Variational Inference (VI) has generally accumulated more attention compared to Expectation Propagation (EP), [Minka \(2001c\)](#). Perhaps because VI casts Bayesian inference as an optimization problem yielding a well-defined lower bound to the model evidence. EP, on the other hand, iteratively approximates moments of factors that together form a global approximation through simpler local computations. It has been shown that variational methods are biased, and that EP outperforms VI in several applications, such as when likelihood functions are non-smooth (See [Turner and Sahani \(2008\)](#) and [Turner and Sahani \(2011\)](#)). This makes us eager to learn more about EP, and interesting research is already doing so (e.g. [Li et al. \(2016\)](#)).

There is however a hurdle with EP: each local approximation must be maintained in memory, leaving an  $O(N)$  memory footprint. As a result, Stochastic Expectation Propagation (SEP), [Li et al. \(2015\)](#), was proposed. SEP maintains a global approximation to the posterior through parameter tying of the EP factors, while updating local approximations in the EP spirit. This makes SEP a good choice for Bayesian learning in the large data setting with its  $O(1)$  space complexity.

Many notorious machine learning problems involve local latent variables (e.g. dimensionality reduction and topic modeling). AEVB successfully performed approximate bayesian variational inference using an auto-encoding map in latent variable models. However, in their

experiments they were not Bayesian about the global model parameters. This thesis combines SEP and AEVB into a new approximate inference scheme.

We propose doing Bayesian learning of the model parameters using SEP, and using an auto-encoding map (recognition model) to approximate the posterior over the latent variables.

The trajectory of this thesis is as follows. First the general problem scenario is mapped out by presenting a generic graphical model and a concrete set of goals. Chapter 2 highlights important theory that has been successful in the VI framework, which we braid into the proposed scheme. Chapter 3 summarizes traditional EP, and introduces SEP. Finally chapter 5 presents Auto-encoding-with-SEP (AESEP) and derives its objective function. In chapter 6 experiments and evaluation is done. We end with suggestions for the direction of future work.

## 1.1 Problem Scenario

We consider a generic class of models<sup>2</sup> with the dataset  $\mathbf{x}_{1:N}$ , where the  $N$  samples are independently and identically distributed (i.i.d.), and  $\mathbf{x}_n$  is either discrete or continuous. The observed variables are generated by some process with latent variable  $\mathbf{z}_n$  per datapoint, and a set of global parameters,  $\boldsymbol{\theta}$ . The local latent variables typically represents some hidden structure in our data, while the global parameters denote the model parameters. This structure is depicted in the graphical model in Figure 1.2.

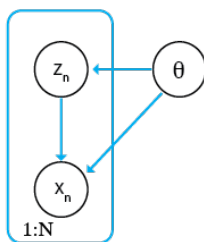


Figure 1.2: Graphical model for class of models under consideration. The model has observations  $\mathbf{x}_{1:N}$ , local latent variables  $\mathbf{z}_{1:N}$ , and global latent variables  $\boldsymbol{\theta}$ . The arrows represent the generative model.

The generative procedure can be described in two steps:  $\mathbf{z}_n$  is generated from a prior dis-

<sup>2</sup>Many well-known machine learning problems have graphical models of this form, e.g. Bayesian mixture models, matrix factorization (Factor analysis, PCA), Dirichlet process mixtures, latent dirichlet allocation, multilevel regression etc.

tribution  $p(\mathbf{z}_n|\boldsymbol{\theta})$ , then  $\mathbf{x}_n$  is generated from a conditional distribution  $p(\mathbf{x}_n|\mathbf{z}_n, \boldsymbol{\theta})$ . The probabilistic model can be written as:

$$p(\boldsymbol{\theta}, \mathbf{z}, \mathbf{x}) = p(\boldsymbol{\theta}) \prod_{n=1}^N p(\mathbf{z}_n, \mathbf{x}_n|\boldsymbol{\theta}) = p(\boldsymbol{\theta}) \prod_{n=1}^N p(\mathbf{z}_n|\boldsymbol{\theta})p(\mathbf{x}_n|\boldsymbol{\theta}, \mathbf{z}_n) \quad (1.1)$$

Both the global and local latent variables are unknown. Bayes' rule is applied to do posterior inference over the latent space:

$$p(\mathbf{z}_{1:N}, \boldsymbol{\theta}|\mathbf{x}_{1:N}) = \frac{p(\boldsymbol{\theta}) \prod_{n=1}^N p(\mathbf{z}_n|\boldsymbol{\theta})p(\mathbf{x}_n|\boldsymbol{\theta}, \mathbf{z}_n)}{p(\mathbf{x}_{1:N})}$$

The computation of the marginal likelihood,  $p(\mathbf{x}_{1:N})$ , involves integrating out the latent space, and is commonly a source of intractability as discussed in the introduction. To perform this task, we propose an approximate inference scheme with the following goals:

1. Efficient Approximate Bayesian learning of the model parameters  $\boldsymbol{\theta}$ .
2. Efficient approximate posterior inference on the local latent variables  $\mathbf{z}$  given an observed value  $\mathbf{x}$ .

In AESEP, SEP is used to achieve goal 1., and recognition models for goal 2.

## 2 | Variational Inference

In this chapter exciting research that has led to improved approximate posterior inference in the VI framework is introduced. In particular we discuss the reparameterization trick, [Papaspiliopoulos et al. \(2007\)](#), for improved gradient estimates, and recognition models for efficient approximate posterior inference, [Kingma and Welling \(2013\)](#). Both techniques are employed in the proposed system.

VI phrases Bayesian inference as an optimization problem. An approximate posterior distribution,  $q_\phi(\mathbf{z}|\mathbf{x})$ , is fit to the exact posterior,  $p_\theta(\mathbf{z}|\mathbf{x})$ , by optimizing the approximation's parameters,  $\phi$ , such that the lower bound,  $\mathcal{L}$ , of the marginal likelihood is maximized. The marginal log-likelihood is composed of the individual datapoint likelihood terms,  $\log p(X) = \sum_{n=1}^N \log p(\mathbf{x}_n)$ . Each individual term can be expressed in the form:

$$\log p(\mathbf{x}_n) = KL(q_\phi(\mathbf{z}|\mathbf{x})||p_\theta(\mathbf{z}|\mathbf{x})) + \mathcal{L}(\boldsymbol{\theta}, \phi; \mathbf{x}_n)$$

Where the KL term is the Kullback-Leibler (KL) divergence between two distributions. The variational lower bound is derived in [appendix A](#), and the KL-divergence is also defined. The result is the following form of the variational lower bound (also known as the ELBO):

$$\log p(\mathbf{x}_n) \geq \mathcal{L}(\boldsymbol{\theta}, \phi; \mathbf{x}_n) = \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[ \log p_\theta(\mathbf{x}, \mathbf{z}) - q_\phi(\mathbf{z}|\mathbf{x}) \right]$$

This optimization requires gradients with respect to both the generative parameters,  $\boldsymbol{\theta}$ , and the variational parameters,  $\phi$ . The gradients with respect to the variational parameters becomes challenging where naive Monte Carlo (MC) estimates of the gradients yield high variance, [Paisley et al. \(2012\)](#). Similar issues are faced in AESEP when taking gradients with respect to the recognition model and SEP parameters.

## 2.1 Use of Reparameterization Trick

The reparameterization trick is a strategy for generating samples from a distribution, e.g.  $z \sim q(z|\mathbf{x}, \phi)$ . The random variable  $\mathbf{z}$  can be reparameterized using a differentiable transform:  $\tilde{\mathbf{z}} = g(\boldsymbol{\epsilon}, \mathbf{x}|\phi)$ . The noise variable,  $\boldsymbol{\epsilon}$ , is independent from  $\phi$ , and is sampled from some prior distribution:  $\boldsymbol{\epsilon} \sim p(\boldsymbol{\epsilon})$ . The idea is that the mapping,  $g(\cdot)$ , is differentiable with respect to both  $\phi$  and  $\boldsymbol{\epsilon}$ . For example, suppose  $q(z) \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ , and let  $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ , then the transformation becomes:

$$\mathbf{z} = g(\boldsymbol{\epsilon}, \mathbf{x}|\phi) = \boldsymbol{\mu} + \boldsymbol{\Sigma}^{1/2}\boldsymbol{\epsilon} \quad (2.1)$$

The differentiable transform depends on the nature of the approximate distribution  $q$ . Intuitively, this is simply a shift in the source of stochasticity to  $\boldsymbol{\epsilon}$ ; merely substitution in terms of calculus. However, some advantageous consequences appear when computing samples of the gradients. Suppose we want to sample from  $q$  with respect to a function of  $\mathbf{z}$ ,  $f(\mathbf{z})$ , then by the chain rule this becomes a derivative with respect to  $\boldsymbol{\epsilon}$ :

$$\nabla_{\phi} \mathbb{E}_{q(z|\mathbf{x}_n, \phi)}[f(\mathbf{z})] = \nabla_{\boldsymbol{\epsilon}} \mathbb{E}_{p(\boldsymbol{\epsilon})}[f(g(\boldsymbol{\epsilon}, \mathbf{x}_n|\phi))]$$

The idea of variance reduction is to modify a function of a random variable,  $\mathbf{z}_n$ , such that the variance decreases, but the expectations stays the same, Paisley et al. (2012). The reparameterization trick reduces the variance of the samples, especially if the differential transform has an isotropic covariance. The mechanics of the calculus above is a result of the deterministic mapping which explicitly results in the substitution:  $q(\mathbf{z}|\mathbf{x}, \phi)d\mathbf{z} = p(\boldsymbol{\epsilon})d\boldsymbol{\epsilon}$ . From this it follows that:

$$\int q(\mathbf{z}|\mathbf{x}, \phi)f(\mathbf{z})d\mathbf{z} = \int p(\boldsymbol{\epsilon})f(\mathbf{z})d\boldsymbol{\epsilon} = \int p(\boldsymbol{\epsilon})f(g(\boldsymbol{\epsilon}, \mathbf{x}|\phi))d\boldsymbol{\epsilon}$$

The reparameterization trick has been applied in the variational setting in a vast amount of published literature. AEVB shows explicitly how it is used to form a more accurate approximation of the ELBO, such that the MC estimate is differentiable with respect to the variational parameters. AESEP utilizes this trick in each local minimization for both the local and global latent variables.



## 2.2 Use of Recognition Models

One of the stated goals is efficient approximate inference of the local latent variables. The implicit dependency between an observation  $\mathbf{x}_n$  and corresponding  $\mathbf{z}_n$  (See Figure 1.2) invites the use of a recognition model,  $q(\mathbf{z}_{1:N}|\mathbf{x}_{1:N}, \phi)$ , to approximate the exact posterior  $p(\mathbf{z}_{1:N}|\mathbf{x}_{1:N}, \theta)$ .

This framework has the advantage of bringing flexibility to the approximation by putting few constraints on the recognition model parameters  $\phi$ . In contrast, the common mean field approximation<sup>1</sup>, which assumes a factorized posterior, not only constraints the posterior but also requires parameters for each factor to be stored in memory. Therefore, the recognition model can also be thought of as parameter tying.

To achieve a deterministic mapping from an observation  $\mathbf{x}_n$  to  $\mathbf{z}_n$ , AEVB makes use of the reparameterization trick (Figure 2.1b.). For example, if the recognition model is assumed Gaussian  $q(\mathbf{z}_n|\mathbf{x}_n, \phi) = \mathcal{N}(\boldsymbol{\mu}_n, \boldsymbol{\Sigma}_n)$ , then the output of the deterministic map given an observation,  $\mathbf{x}_n$ , is  $\boldsymbol{\mu}_n$  and  $\boldsymbol{\Sigma}_n$ . Given the Gaussian assumed structure, the reparameterization in Equation 2.1 can be utilized to deterministically map  $\mathbf{x}_n$  to  $\mathbf{z}_n$ .

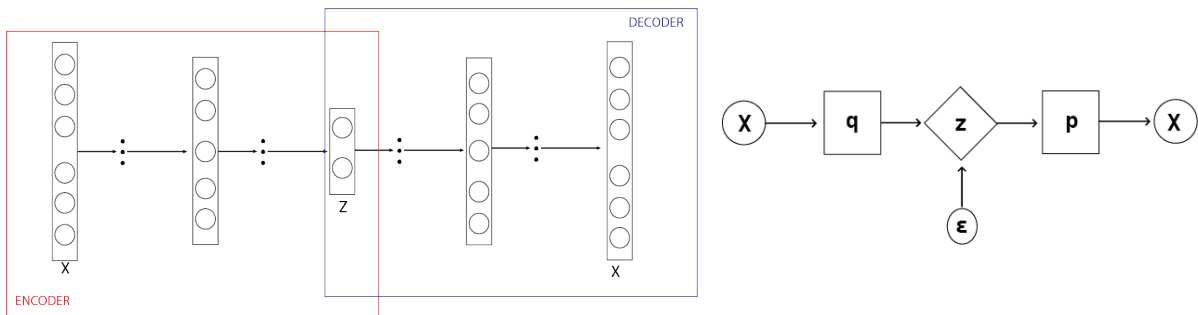


Figure 2.1: a. Autoencoder showing the mapping and terminology used. b. Reparameterization trick's role in the encoding/decoding process, where the kite-shaped square around  $\mathbf{z}$  implies the random variable is deterministic. Figure b. is inspired by D. Kingma's talk on AEVB, Kingma (2014).

AEVB gives a clear picture of the system through the use of autoencoders (Figure 2.1a.). The deterministic encoding map (e.g. observation to mean and variance) is a neural network, and likewise is the decoding map. Therefore the only parameters necessary are the weights of the network. In the general framework, any regression map could be utilized. To do approximate

---

<sup>1</sup> $q(\mathbf{z}_{1:N}) = \prod_{n=1}^N q_n(\mathbf{z}_n)$

VI, AEVB utilizes this to evaluate the variational lower bound. Specifically, once an observation  $\mathbf{x}_n$  has been mapped to a latent variable  $\mathbf{z}_n$ , this is used to evaluate the likelihood of that particular observation which is necessary for the evaluation of the ELBO.

# 3 | Expectation Propagation

This chapter starts by outlining traditional EP, and explains how EP is able to achieve quality approximations of the posterior. Following last chapter's brief introduction to VI, the contrast between the two becomes clear. Once the EP framework is established, SEP is introduced which is how learning is done in the proposed system.

## 3.1 EP

Expectation propagation approximates the unnormalized posterior, or the joint distribution, and returns a tractable normalized posterior. For convenience, let  $*$  denote unnormalized densities. Suppose we have a dataset,  $X$ , and  $\theta$  is a vector containing all latent variables, both global and/or local, then the intractable exact posterior,  $p(\theta|X)$ , can be approximated with a tractable approximate distribution  $q(\theta)$ :

$$p(\theta, X) = p(\theta) \prod_{n=1}^N p(\mathbf{x}_n|\theta) \approx q^*(\theta) = p(\theta) \prod_{n=1}^N f_n(\theta)$$

Each likelihood term is approximated by a corresponding factor,  $f_n(\theta)$ . Often these factors are constrained to be members of a chosen exponential family. The result is a posterior in the same family<sup>1</sup>.

Like variational inference, EP too, is based on minimizing the KL- divergence. However, EP obtains its personality by taking the KL between the exact posterior and the approximate posterior: the reverse form compared with the KL used in VI.

$$KL(p(\theta|X)||q(\theta)) = \int p(\theta|X) \log \frac{p(\theta|X)}{q(\theta)} d\theta = \mathbb{E}_{p(\theta|X)} \left[ \log \frac{p(\theta|X)}{q(\theta)} \right]$$

---

<sup>1</sup>The product of distributions of the same exponential family is also a member of that family.

The above minimization is intractable since taking an expectation with respect to the intractable exact posterior is analytically impossible. EP gets around this by iteratively optimizing each factor in the context of all the other factors, using what is referred to as the cavity distribution: The product of all factors (including the prior), except for the current factor,  $n$ , being updated. The cavity distribution for factor  $n$  is:

$$q_{-n}(\boldsymbol{\theta}) = \frac{q(\boldsymbol{\theta})}{f_n(\boldsymbol{\theta})}$$

The local minimization for a factor  $n$ , becomes the following.  $p(\mathbf{x}_n|\boldsymbol{\theta})q_{-n}(\boldsymbol{\theta})$  is referred to as the tilted distribution.

$$KL\left(\frac{p(\mathbf{x}_n|\boldsymbol{\theta})q_{-n}(\boldsymbol{\theta})}{Z_n} \parallel f_n(\boldsymbol{\theta})q_{-n}(\boldsymbol{\theta})\right) \quad (3.1)$$

The tilted distribution's presence in the local minimizations is fundamental to EP's accurate global approximations. To offer clairvoyance, one can instead imagine an iterative procedure where the KL-divergence is computed directly between  $p(\mathbf{x}_n|\boldsymbol{\theta})$  and  $f_n(\mathbf{x})$ . The product of  $n$  such entirely independent minimizations would give a globally unreliable approximation. It might seem that one local minimization does not incorporate any information about the other  $N - 1$  minimizations, but this is not the case as the cavity distribution captures the current estimate from the  $N - 1$  observations. In other words, the advantage of the cavity distribution arises due to the fact that the cavity defines the areas where the posterior probability is high. Figure 3.1 visualizes the role played by the tilted and cavity distribution in each likelihood approximation.

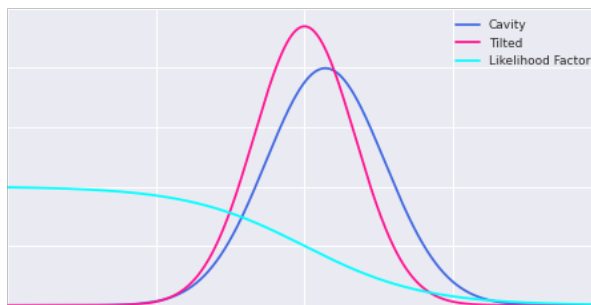


Figure 3.1: Example of EP update visualising the stability EP gets from the cavity distribution in the KL minimization. In this case the likelihood factors are assumed Gaussian. Observe that even though the likelihood terms do not resemble a Gaussian, the cavity distribution ensures that the approximation is more accurate where the posterior has a high value. Figure idea taken from Gelman et al. (2014).

Furthermore, when the approximating factors are members of the exponential family the minimization (Equation 3.1) reduces to moment matching of the moments of the tilted distribution. The proof of this result is provided in appendix B. The result in the EP minimization is the following projection of moments:

$$f_n(\boldsymbol{\theta})q_{-n}(\boldsymbol{\theta}) \leftarrow \text{proj}[p(\mathbf{x}_n|\boldsymbol{\theta})q_{-n}(\boldsymbol{\theta})]$$

Finally, the updated factor is found by dividing out the cavity distribution. In the rare case where the exact likelihood term and the approximating factor are in the same family, moment matching is trivial (For an example see chapter 5.1). Usually the likelihood takes a complicated form, and moments must be sampled from the exact likelihood term  $p(\mathbf{x}_n|\boldsymbol{\theta})$ . The algorithmic steps of EP can be summarized:

1. **Initialization.** All factors are initialized. Usually restricted to be members of an exponential family, e.g.  $f_n(\boldsymbol{\theta}) \leftarrow \mathcal{N}(\mathbf{0}, \mathbf{I})$ . Also initialize the unnormalized posterior:  $q^*(\boldsymbol{\theta}) = p(\boldsymbol{\theta}) \prod_{n=1}^N f_n(\boldsymbol{\theta})$
2. **Iteration.** Until convergence:
  - (a) **Cavity.** Compute the cavity distribution:
 
$$q_{-n}(\boldsymbol{\theta}) = \frac{q(\boldsymbol{\theta})}{f_n(\boldsymbol{\theta})}$$
  - (b) **Tilted.** Compute the tilted distribution:
 
$$p(\mathbf{x}_n|\boldsymbol{\theta})q_{-n}(\boldsymbol{\theta})$$
  - (c) **Moment matching.**

$$q^{new}(\boldsymbol{\theta}) = f_n(\boldsymbol{\theta})q_{-n}(\boldsymbol{\theta}) \leftarrow \text{proj}[p(\mathbf{x}_n|\boldsymbol{\theta})q_{-n}(\boldsymbol{\theta})]$$
  - (d) **Normalisation.** Compute the marginal likelihood of the tilted:
 
$$Z_n = \int p(\mathbf{x}_n|\boldsymbol{\theta})q_{-n}(\boldsymbol{\theta})d\boldsymbol{\theta}$$
  - (e) **Update.** Update the new factor and store
 
$$f_n(\boldsymbol{\theta}) = Z_n \frac{q^{new}(\boldsymbol{\theta})}{q_n(\boldsymbol{\theta})}$$

The local nature of EP makes it easy to parallelize and distribute. Moreover, the local computations are cheap compared with a more complex global minimization. Unfortunately, there are two main drawbacks to EP. First is the  $O(N)$  memory footprint as a result of the cavity's requirement to store the  $N$  approximating factors in memory. Second, traditional EP has no guarantee of convergence, [Minka \(2001a\)](#).

## 3.2 SEP

The beauty of EP's local computations has a critical limitation: The number of approximating factors linearly increase as the number of observations increase<sup>2</sup>. This becomes an issue when doing approximate bayesian parameter learning over large datasets. As a reference point, VI does not face this type of memory scaling because it maintains a global approximation. SEP also maintains a global approximation, while updating the global approximation in a local manner, resulting in  $O(1)$  space complexity.

The eventual goal of EP is to summarize the posterior through the likelihood factors. In its optimization it is assumed that the contribution of each likelihood term would be averaged out at convergence. SEP directly assumes this through parameter tying:

$$f(\boldsymbol{\theta})^N \triangleq \prod_{n=1}^N f_n(\boldsymbol{\theta}) \approx \prod_{n=1}^N p(\mathbf{x}_n|\boldsymbol{\theta})$$

Explicitly, only  $f(\boldsymbol{\theta})$  is maintained, and the unnormalized posterior becomes:

$$q(\boldsymbol{\theta}) \propto f(\boldsymbol{\theta})^N p(\boldsymbol{\theta})$$

SEP uses the same iterative local KL-minimization as EP to learn the current datapoint's approximating factor  $f_n(\boldsymbol{\theta})$ . However, the minimization is in the context of the averaged cavity,  $q_{-1}(\boldsymbol{\theta}) \propto \frac{q(\boldsymbol{\theta})}{f(\boldsymbol{\theta})}$ : summarizing the average contribution of each datapoint to the posterior.

$$KL\left(\frac{p(\mathbf{x}_n|\boldsymbol{\theta})q_{-1}(\boldsymbol{\theta})}{Z_n}\right) \parallel p(\boldsymbol{\theta})f(\boldsymbol{\theta})^N$$

To achieve a local approximation similarly perform moment matching and divide out the cavity:

$$f_n(\boldsymbol{\theta}) \leftarrow \text{proj}[q_{-1}(\boldsymbol{\theta})p(\mathbf{x}_n|\boldsymbol{\theta})]/q_{-1}(\boldsymbol{\theta})$$

Finally, the global factor is updated using damping:

$$f(\boldsymbol{\theta}) = f(\boldsymbol{\theta})^{1-\alpha} f_n(\boldsymbol{\theta})^\alpha$$

Setting the damping parameter  $\alpha = \frac{1}{N}$  is the most obvious choice, but the SEP authors also suggests other choices, such as an adaptive damping parameter, and decreasing  $\alpha$  according to the Robbins-Monroe condition.

---

<sup>2</sup>Prior to SEP, solving this required a visit to the hardware store or the cloud.

The algorithmic steps are very similar to the detailed outline given in the EP section, with the SEP updates as explained. The algorithm is summarized in Figure 3.2, and the contrast to EP can be seen. In particular take note of step 2, the computation of the cavity, clearly showing how SEP uses the averaged contribution of each datapoint. Theoretically, this is a compromise to the quality of the approximation, however the gain in terms of resources by changing to SEP from EP might often be more dramatic. The authors achieve impressive results on both synthetic and real datasets.

---

**Algorithm 1** EP
 

---

- 1: choose a factor  $f_n$  to refine:
  - 2: compute cavity distribution  
 $q_{-n}(\boldsymbol{\theta}) \propto q(\boldsymbol{\theta})/f_n(\boldsymbol{\theta})$
  - 3: compute tilted distribution  
 $\tilde{p}_n(\boldsymbol{\theta}) \propto p(\mathbf{x}_n|\boldsymbol{\theta})q_{-n}(\boldsymbol{\theta})$
  - 4: moment matching:  
 $f_n(\boldsymbol{\theta}) \leftarrow \text{proj}[\tilde{p}_n(\boldsymbol{\theta})]/q_{-n}(\boldsymbol{\theta})$
  - 5: inclusion:  
 $q(\boldsymbol{\theta}) \leftarrow q_{-n}(\boldsymbol{\theta})f_n(\boldsymbol{\theta})$
- 

---

**Algorithm 2** SEP
 

---

- 1: choose a datapoint  $\mathbf{x}_n \sim \mathcal{D}$ :
  - 2: compute cavity distribution  
 $q_{-1}(\boldsymbol{\theta}) \propto q(\boldsymbol{\theta})/f(\boldsymbol{\theta})$
  - 3: compute tilted distribution  
 $\tilde{p}_n(\boldsymbol{\theta}) \propto p(\mathbf{x}_n|\boldsymbol{\theta})q_{-1}(\boldsymbol{\theta})$
  - 4: moment matching:  
 $f_n(\boldsymbol{\theta}) \leftarrow \text{proj}[\tilde{p}_n(\boldsymbol{\theta})]/q_{-1}(\boldsymbol{\theta})$
  - 5: inclusion:  
 $q(\boldsymbol{\theta}) \leftarrow q_{-1}(\boldsymbol{\theta})f_n(\boldsymbol{\theta})$
  - 6: *implicit update*:  
 $f(\boldsymbol{\theta}) \leftarrow f(\boldsymbol{\theta})^{1-\frac{1}{N}}f_n(\boldsymbol{\theta})^{\frac{1}{N}}$
- 

Figure 3.2: Comparing EP and SEP algorithms. The procedure is done iteratively until convergence. In step one, EP chooses a factor to refine, while SEP chooses a datapoint (as local factors are not kept in memory) to compute the local minimization for updating the global approximation.

# 4 | Auto-encoding with SEP

This chapter introduces the proposed approximate inference scheme: Auto-encoding with SEP. First, we motivate the EP-like approximations. We incorporate parameter tying strategies to define a new set of approximations. Then we derive a local objective function by minimizing a KL-divergence in an EP-like fashion. Following this we propose the use of a sampling based stochastic optimization of the objective function.

## 4.1 Approximating Factors and Parameter Tying

The Problem Scenario introduced the generic probabilistic model (Equation 1.1) being considered. Mathematically, the joint likelihood  $p(\mathbf{z}_n, \mathbf{x}_n | \boldsymbol{\theta})$  is a function  $t_n(\mathbf{x}_n, \mathbf{z}_n, \boldsymbol{\theta})$  with parameters  $\mathbf{x}_n, \mathbf{z}_n$  and  $\boldsymbol{\theta}$ . In alignment with the EP philosophy, we approximate  $t_n(\mathbf{x}_n, \mathbf{z}_n, \boldsymbol{\theta})$  with  $\tilde{t}_n(\mathbf{x}_n, \mathbf{z}_n, \boldsymbol{\theta})$ .

$$p(\mathbf{z}_n | \boldsymbol{\theta})p(\mathbf{x}_n | \boldsymbol{\theta}, \mathbf{z}_n) = t_n(\mathbf{x}_n, \mathbf{z}_n, \boldsymbol{\theta}) \approx \tilde{t}_n(\mathbf{x}_n, \mathbf{z}_n, \boldsymbol{\theta}) = f_n(\mathbf{z}_n)\tilde{t}_n(\boldsymbol{\theta})$$

This can further be approximated using SEP parameter tying on the  $\tilde{t}_n(\boldsymbol{\theta})$  factors, such that  $\tilde{t}(\boldsymbol{\theta})^N = \prod_{n=1}^N \tilde{t}_n(\boldsymbol{\theta})$ , resulting in the following expression:

$$\tilde{t}_n(\mathbf{x}_n, \mathbf{z}_n, \boldsymbol{\theta}) = f_n(\mathbf{z}_n)\tilde{t}(\boldsymbol{\theta})$$

At this stage each  $f_n(\mathbf{z}_n)$  could be approximated using EP, yielding an SEP-EP model. We tie parameters at one more level. Intuitively there is a local dependency between the observed  $\mathbf{x}_n$  and the corresponding latent variable  $\mathbf{z}_n$ . Unlike the traditional EP approach, where each  $f_n$  factor would have separate parameters,  $\boldsymbol{\lambda}_n$ , we can exploit this intuitive relationship via an encoder map (like in chapter 2.2). Such that  $f_n(\mathbf{z}_n) = f(g_\phi(\mathbf{x}_n))$ . For convenience we write  $f(\mathbf{z}_n)$ , to represent the deterministic mapping from  $\mathbf{x}_n$  to  $\mathbf{z}_n$  via a recognition model utilizing the reparameterization trick. This makes the final approximation with both sources



of parameter tying:

$$p(\mathbf{z}_n|\boldsymbol{\theta})p(\mathbf{x}_n|\boldsymbol{\theta}, \mathbf{z}_n) = t_n(\mathbf{x}_n, \mathbf{z}_n, \boldsymbol{\theta}) \approx \tilde{t}_n(\mathbf{x}_n, \mathbf{z}_n, \boldsymbol{\theta}) = f(\mathbf{z}_n)\tilde{t}(\boldsymbol{\theta})$$

The complete approximation can be expressed as:

$$p(\mathbf{z}_{1:N}, \mathbf{x}_{1:N}, \boldsymbol{\theta}) = p(\boldsymbol{\theta}) \prod_{n=1}^N p(\mathbf{z}_n|\boldsymbol{\theta})p(\mathbf{x}_n|\boldsymbol{\theta}, \mathbf{z}_n) \approx p(\boldsymbol{\theta})\tilde{t}(\boldsymbol{\theta})^N \prod_{n=1}^N f(\mathbf{z}_n)$$

## 4.2 Local KL-divergence

The EP minimization step takes the KL-divergence between the tilted distribution, and the unnormalised approximate posterior. These densities are expressed respectively:

$$q(\boldsymbol{\theta}, \mathbf{z}_{1:N}) \propto p(\boldsymbol{\theta}) \prod_n \tilde{t}_n(\mathbf{x}_n, \mathbf{z}_n, \boldsymbol{\theta}) \quad (4.1)$$

$$q^{tilt}(\boldsymbol{\theta}, \mathbf{z}_{1:N}) = \frac{1}{Z} \frac{q(\boldsymbol{\theta}, \mathbf{z}_{1:N})}{\tilde{t}_n(\mathbf{x}_n, \mathbf{z}_n, \boldsymbol{\theta})} t_n(\mathbf{x}_n, \mathbf{z}_n, \boldsymbol{\theta}) \quad (4.2)$$

$Z$  is the normalization constant, which ensures that  $q^{tilt}(\boldsymbol{\theta}, \mathbf{z}_{1:N})$  integrates to unity. By taking the KL-divergence between Equation 4.1 and 4.2, a local objective,  $E$ , is derived.

$$\begin{aligned} E &= KL\left(q^{tilt}(\boldsymbol{\theta}, \mathbf{z}_{1:N}) \parallel q(\boldsymbol{\theta}, \mathbf{z}_{1:N})\right) \\ &= \int \dots \int \frac{1}{Z} \frac{q(\boldsymbol{\theta}, \mathbf{z}_{1:N})}{\tilde{t}_n(\mathbf{x}_n, \mathbf{z}_n, \boldsymbol{\theta})} t_n(\mathbf{x}_n, \mathbf{z}_n, \boldsymbol{\theta}) \log \frac{\frac{1}{Z} \frac{q(\boldsymbol{\theta}, \mathbf{z}_{1:N})}{\tilde{t}_n(\mathbf{x}_n, \mathbf{z}_n, \boldsymbol{\theta})} t_n(\mathbf{x}_n, \mathbf{z}_n, \boldsymbol{\theta})}{q(\boldsymbol{\theta}, \mathbf{z}_{1:N})} d\boldsymbol{\theta} d\mathbf{z}_{1:N} \\ &= \int d\mathbf{z}_{\neq n} \int \int \frac{1}{Z} \left[ \prod_{m \neq n} \tilde{t}_n(\mathbf{x}_m, \mathbf{z}_m, \boldsymbol{\theta}) \right] p(\boldsymbol{\theta}) t_n(\mathbf{x}_n, \mathbf{z}_n, \boldsymbol{\theta}) \log \frac{t_n(\mathbf{x}_n, \mathbf{z}_n, \boldsymbol{\theta})}{Z \tilde{t}_n(\mathbf{x}_n, \mathbf{z}_n, \boldsymbol{\theta})} d\boldsymbol{\theta}, \mathbf{z}_n \\ &= \int \int \frac{1}{Z} \tilde{t}(\boldsymbol{\theta})^{N-1} p(\boldsymbol{\theta}) p(\mathbf{z}_n, \mathbf{x}_n | \boldsymbol{\theta}) \log \frac{p(\mathbf{z}_n, \mathbf{x}_n | \boldsymbol{\theta})}{Z f_n(\mathbf{z}_n) \tilde{t}(\boldsymbol{\theta})} d\boldsymbol{\theta} d\mathbf{z}_n \\ &= \int \int \frac{1}{Z} \tilde{t}(\boldsymbol{\theta})^{N-1} p(\boldsymbol{\theta}) p(\mathbf{z}_n, \mathbf{x}_n | \boldsymbol{\theta}) \frac{\Gamma(\mathbf{x}_n, \mathbf{z}_n, \boldsymbol{\theta})}{\Gamma(\mathbf{x}_n, \mathbf{z}_n, \boldsymbol{\theta})} \log \Gamma(\mathbf{x}_n, \mathbf{z}_n, \boldsymbol{\theta}) d\boldsymbol{\theta} d\mathbf{z}_n \\ &= \int \int \frac{1}{Z} \tilde{t}(\boldsymbol{\theta})^{N-1} p(\boldsymbol{\theta}) p(\mathbf{z}_n, \mathbf{x}_n | \boldsymbol{\theta}) \frac{Z f_n(\mathbf{z}_n) \tilde{t}_n(\boldsymbol{\theta})}{p(\mathbf{x}_n, \mathbf{z}_n | \boldsymbol{\theta})} \Gamma(\mathbf{x}_n, \mathbf{z}_n, \boldsymbol{\theta}) \log \Gamma(\mathbf{x}_n, \mathbf{z}_n, \boldsymbol{\theta}) d\boldsymbol{\theta} d\mathbf{z}_n \\ &= \int \int \tilde{t}(\boldsymbol{\theta})^N p(\boldsymbol{\theta}) f_n(\mathbf{z}_n) \Gamma(\mathbf{x}_n, \mathbf{z}_n, \boldsymbol{\theta}) \log \Gamma(\mathbf{x}_n, \mathbf{z}_n, \boldsymbol{\theta}) d\boldsymbol{\theta} d\mathbf{z}_n \\ &= \int \int q(\boldsymbol{\theta}) f_n(\mathbf{z}_n) \Gamma(\mathbf{x}_n, \mathbf{z}_n, \boldsymbol{\theta}) \log \Gamma(\mathbf{x}_n, \mathbf{z}_n, \boldsymbol{\theta}) d\boldsymbol{\theta} d\mathbf{z}_n \end{aligned}$$

$\Gamma$  is the following expression:

$$\Gamma(\mathbf{x}_n, \mathbf{z}_n, \boldsymbol{\theta}) = \frac{p(\mathbf{x}_n, \mathbf{z}_n | \boldsymbol{\theta})}{Z f_n(\mathbf{z}_n) \tilde{t}_n(\boldsymbol{\theta})} \quad (4.3)$$

### 4.3 Stochastic Objective

The objective can be approximated by employing sampling schemes. The unobserved variables are now sampled,  $\{\boldsymbol{\theta}^{(m)}, \mathbf{z}_n^{(m)}\}_{m=1}^M \sim q(\boldsymbol{\theta})f(\mathbf{z}_n)$ , such that  $\boldsymbol{\theta}$  is sampled from the SEP posterior approximation, and  $\mathbf{z}_n$  from the recognition model approximation. The samples are taken via the reparameterization trick, to achieve benefits (see chapter 2.1) when computing the gradients with respect to both the SEP and recognition model parameters. Approximating  $E$  via sampling results in:

$$E \approx \tilde{E} = \frac{1}{M} \sum_{m=1}^M \Gamma(\mathbf{x}_n, \mathbf{z}_n^{(m)}, \boldsymbol{\theta}^{(m)}) \log \Gamma(\mathbf{x}_n, \mathbf{z}_n^{(m)}, \boldsymbol{\theta}^{(m)}) \quad (4.4)$$

In a similar manner, the normalization constant in Equation 4.3 can be approximated from the unnormalized  $\Gamma$ :

$$Z \approx \frac{1}{M} \sum_{m=1}^M \Gamma^*(\mathbf{x}_n, \mathbf{z}_n^{(m)}, \boldsymbol{\theta}^{(m)}) \quad (4.5)$$

At this stage, the tradeoff between a biased and an unbiased estimate becomes relevant. Note that an unbiased estimate of  $Z$ , does not mean that  $\frac{1}{Z}$  is unbiased. Theoretical work on unbiased estimates of  $\frac{1}{Z}$  exists. To achieve an unbiased estimate for the final approximation, separate sets of samples must be used in the approximation of  $Z$  and  $\Gamma^*$ . Unfortunately, this leads to a resulting approximation with high variance. Intuitively, the two sets invite two sources of randomness, which leads to the observed increased variance. Self normalized important sampling reduces the variance by using one set of samples for both approximations.

Following Equation 4.4, we utilize this fact which results in the following biased estimate:

$$\begin{aligned}\tilde{E} &= \frac{1}{M} \sum_{m=1}^M \frac{1}{Z} \Gamma^*(\mathbf{x}_n, \mathbf{z}_n^{(m)}, \boldsymbol{\theta}^{(m)}) \log \frac{\Gamma^*(\mathbf{x}_n, \mathbf{z}_n^{(m)}, \boldsymbol{\theta}^{(m)})}{Z} \\ &= \frac{1}{M} \sum_{m=1}^M \frac{\Gamma^*(\mathbf{x}_n, \mathbf{z}_n^{(m)}, \boldsymbol{\theta}^{(m)})}{Z} \left[ \log \Gamma^*(\mathbf{x}_n, \mathbf{z}_n^{(m)}, \boldsymbol{\theta}^{(m)}) - \log Z \right]\end{aligned}$$

Plugging in for  $Z$  using Equation 4.5:

$$\begin{aligned}&= \frac{1}{M} \sum_{m=1}^M \frac{\Gamma^*(\mathbf{x}_n, \mathbf{z}_n^{(m)}, \boldsymbol{\theta}^{(m)})}{\frac{1}{M} \sum_{m=1}^M \Gamma^*(\mathbf{x}_n, \mathbf{z}_n^{(m)}, \boldsymbol{\theta}^{(m)})} \left[ \log \Gamma^*(\mathbf{x}_n, \mathbf{z}_n^{(m)}, \boldsymbol{\theta}^{(m)}) - \log \left( \frac{1}{M} \sum_{m=1}^M \Gamma^*(\mathbf{x}_n, \mathbf{z}_n^{(m)}, \boldsymbol{\theta}^{(m)}) \right) \right] \\ &= \sum_{m=1}^M \frac{\Gamma^*(\mathbf{x}_n, \mathbf{z}_n^{(m)}, \boldsymbol{\theta}^{(m)})}{\sum_{m=1}^M \Gamma^*(\mathbf{x}_n, \mathbf{z}_n^{(m)}, \boldsymbol{\theta}^{(m)})} \left[ \log \Gamma^*(\mathbf{x}_n, \mathbf{z}_n^{(m)}, \boldsymbol{\theta}^{(m)}) - \log \left( \sum_{m=1}^M \Gamma^*(\mathbf{x}_n, \mathbf{z}_n^{(m)}, \boldsymbol{\theta}^{(m)}) \right) - \log \frac{1}{M} \right]\end{aligned}$$

Next the self normalising weights  $\tilde{w}$  are defined:

$$\begin{aligned}\tilde{w}^{(m)} &= \frac{\Gamma^*(\mathbf{x}_n, \mathbf{z}_n^{(m)}, \boldsymbol{\theta}^{(m)})}{\sum_{m=1}^M \Gamma^*(\mathbf{x}_n, \mathbf{z}_n^{(m)}, \boldsymbol{\theta}^{(m)})} \\ \log \tilde{w}^{(m)} &= \log \Gamma^*(\mathbf{x}_n, \mathbf{z}_n^{(m)}, \boldsymbol{\theta}^{(m)}) - \log \sum_{m=1}^M \Gamma^*(\mathbf{x}_n, \mathbf{z}_n^{(m)}, \boldsymbol{\theta}^{(m)})\end{aligned}$$

Combining the expression, we get:

$$\tilde{E} = \sum_{m=1}^M \tilde{w}^{(m)} \log \tilde{w}^{(m)} + M \log(M)$$

Since the objective is being minimized the constant can be ignored, giving:

$$\tilde{E} = \sum_{m=1}^M \tilde{w}^{(m)} \log \tilde{w}^{(m)}$$

## 4.4 Algorithm

We denote  $\psi$  as the parameters of the global SEP factor (i.e. the moments of the chosen parametric distribution), and  $\phi$  the parameters of the recognition model (e.g. the weights of a neural network).  $\alpha$  and  $\beta$  are the damping constants for the SEP and recognition model, respectively, and  $g(\cdot)$  is a differentiable transform via the reparameterization trick.

1. **Initialization.** Initialize the global SEP factor and the Recognition model parameters:  
 $\psi, \phi \leftarrow \text{Initialize}$
2. **Iteration.** Until convergence:
  - **Datapoint.** Choose a datapoint  $\mathbf{x}_n$ .
  - **Sample.** Random samples for the reparameterization of  $\theta$  and  $\mathbf{z}_n$   
 $\epsilon_\theta \sim p(\epsilon)$   
 $\epsilon_z \sim p(\epsilon)$
  - **Gradients.** Compute gradients of local objective  $\tilde{E}$   
 $\mathbf{g} \leftarrow \nabla_{\psi, \phi} \tilde{E}(\mathbf{x}_n, g_\phi(\mathbf{x}_n, \epsilon_z), g_\psi(\mathbf{x}_n, \epsilon_\theta))$
  - **Optimize.** Optimize parameters using gradient based method to achieve local optimal set of parameters.  
 $\psi_n, \phi_n$
  - **SEP Update.** Update global SEP factor  
 $\tilde{t}(\theta) = \tilde{t}(\theta)^{1-\alpha} \tilde{t}_n(\theta)^\alpha$
  - **Update Recognition Model.** Update recognition model using damping  
 $\phi \leftarrow \phi^{1-\beta} \phi_n^\beta$

## 5 | Experiments

We use Factor Analysis (FA) as an instance of the generic latent variable model (Figure 1.2) to test AESEP. FA is a linear model where each observation,  $\mathbf{x}_n \in \mathbb{R}^p$ , is generated by some lower dimensional latent variable,  $\mathbf{z}_n \in \mathbb{R}^q$ , through a lineal transformation,  $\mathbf{W} \in \mathbb{R}^{p \times q}$ , with added isotropic noise,  $\sigma_x$ . The transform can for example be:  $\mathbf{x}_n = \mathbf{W}\mathbf{z}_n + \sigma_x \mathbf{I}$ . This has conditional density<sup>1</sup>:

$$p(\mathbf{x}_n | \mathbf{z}_n, \mathbf{W}, \sigma_x^2) \propto \exp\{-0.5(\mathbf{x}_n - \mathbf{W}\mathbf{z}_n)^T (\sigma_x^2 \mathbf{I})^{-1} (\mathbf{x}_n - \mathbf{W}\mathbf{z}_n)\}$$

From the likelihood, observe that the posterior is the following:

$$p(\mathbf{z}_n | \mathbf{x}_n, \mathbf{W}, \sigma_x^2) \sim \mathcal{N}([\mathbf{I} + \mathbf{W}^T (\sigma_x^2 \mathbf{I})^{-1} \mathbf{W}]^{-1} \mathbf{W}^T (\sigma_x^2 \mathbf{I})^{-1} \mathbf{x}_n, [\mathbf{I} + \mathbf{W}^T (\sigma_x^2 \mathbf{I})^{-1} \mathbf{W}]^{-1}) \quad (5.1)$$

Equation 5.1 is compared to the posterior inference approximation produced by the recognition model to verify the approximation. Because SEP's convergence properties are not well understood, Li et al. (2015), evaluation of model learning by SEP is more challenging. This is in contrast to VI with a well-defined lower bound on marginal likelihood. As an alternative, we use the predictive marginal likelihood as our metric. Another option is to plot the fixed points of the EP energy function (See Minka (2001b)). However, this is also questionable as the direct relationship between EP and SEP is an open question, Li et al. (2015). The exact generative model used to generate synthetic data is:

$$p(\mathbf{z}) = \mathcal{N}(\mathbf{0}, \mathbf{I})$$

$$p(\mathbf{x} | \mathbf{z}, \theta) = \mathcal{N}(\mathbf{W}\mathbf{z}, \sigma_x^2 \mathbf{I})$$

$$w_{i,j} \sim \mathcal{N}(0, \sigma_w^2)$$

$$\sigma_{x,w} \sim \mathcal{N}(0, 1)$$

---

<sup>1</sup>See Fokoué (2009) for a detailed introduction to Factor Analysis.

Implementation details can be found on [https://github.com/verajohne/Autoencoding\\_SEP](https://github.com/verajohne/Autoencoding_SEP)

## 5.1 Vanilla EP to Infer Latent Variables

We show that EP performs exact inference in the case of factor analysis. This also serves as a clear demonstration of moment matching. The likelihood terms of FA are Gaussian, therefore when the approximated factors are assumed Gaussian, moments are easily computed through the exact projection:

$$f_n(\mathbf{z}_n) = p(\mathbf{x}_n | \mathbf{z}_n, \boldsymbol{\theta})$$

This derivation is included in appendix C. Figure 5.1 shows that EP successfully performs exact inference over the latent space in FA.

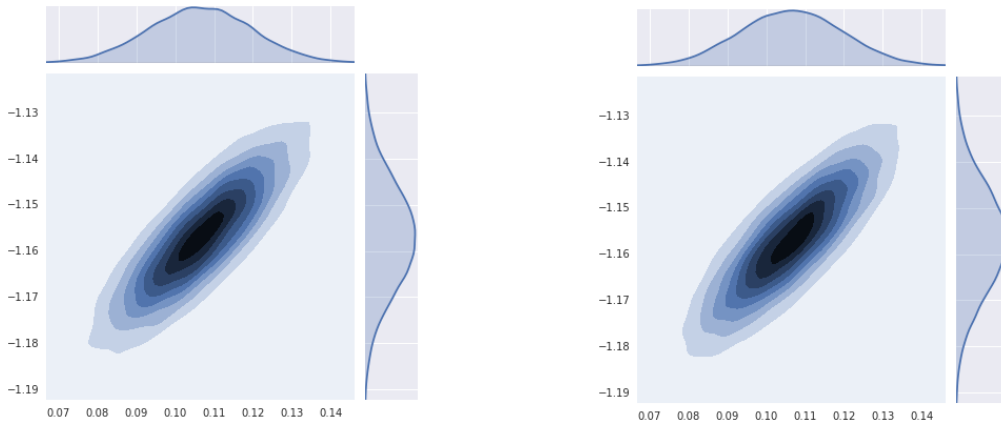


Figure 5.1: Posterior over local latent variable  $\mathbf{z}_n$  where the latent space is 2 dimensional. samples are generated using a. Equation 5.1 b. Exact EP inference on  $\mathbf{z}_n$ . The small differences in the covariance is due to the stochasticity in the sampler.

## 5.2 Auto-encoding with SEP

The global SEP factor,  $t(\boldsymbol{\theta})$ , is assumed Gaussian with mean,  $\mathbf{u}$ , and covariance,  $\mathbf{V}$ , such that  $t(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta}; \mathbf{u}, \mathbf{V})$ . Since factor analysis is a linear model, the recognition model is utilized to learn the following encoding:

$$f(\mathbf{z}) = \mathcal{N}(\mathbf{z}; \mathbf{R}\mathbf{x}, \mathbf{S})$$

$\mathbf{R}$  and  $\mathbf{S}$  are shared across all  $N$  data-points. Thus, the only variables maintained in memory are  $\mathbf{R}, \mathbf{S}, \mathbf{u}, \mathbf{V}$ ; independent of dataset size. Indeed, an increase in dimensionality of either latent or observed space will result in larger memory occupancy of each individual variable.

A number of observations are worth noting for future efforts. Initial testing used separate MC samples for estimates of  $\Gamma^*$  and the normalization constant  $Z$ . This resulted in numerical issues due to the two sources of variance (see chapter 4.3). As a result, we turned to self normalizing importance sampling. This solved the numerical issues, but noisy gradients were still observed, despite the presence of the reparameterization trick as a control variate. One strategy that reduced observed noise was fixing the random number generator seed in each local minimization. This constrained the variance between each batch of MC samples.

Moreover, since the objective function is a MC estimate, stochastic gradient based optimizers, such as Adam Kingma and Ba (2014), can be leveraged. The implementation of Adam dramatically improved computation time (See Figure 5.2).

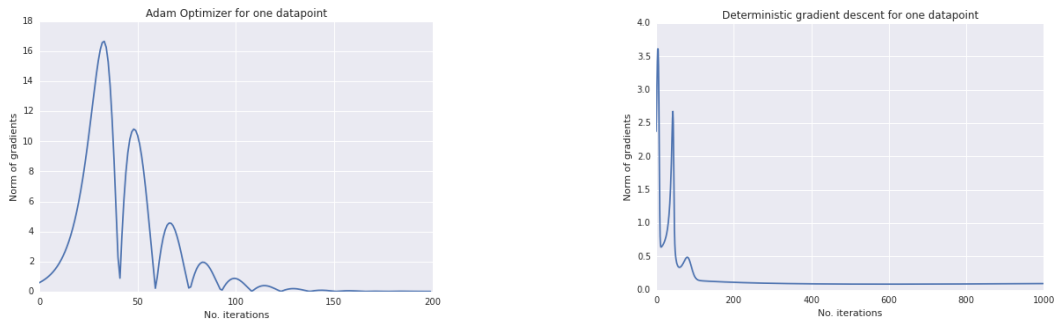


Figure 5.2: The norm of the gradients of  $\mathbf{u}, \mathbf{R}, \mathbf{S}$ . Taken during the minimization of the first datapoint during the first iteration, using a. Adam b. deterministic gradient descent (GD). Observing the range of the axes shows that (i) Adam stabilized much faster than deterministic GD (x-axis) (ii) Adam adds momentum to the search (y-axis). The parameters of Adam were initialized as advised by the authors

The initial sanity check set  $\mathbf{R}, \mathbf{S}, \mathbf{u}, \mathbf{V}$  to their optimal values, and it was observed that marginal likelihood stayed maximized, as well as the parameters remaining at their optimal values; indicative of a stable local objective function. Prior to the proposed joint optimization of AESEP, both SEP and the recognition model were tested separately. Figure 5.3 shows the results after fixing the SEP factor to optimal values, and learning the encoding map. Comparing the exact and approximated posterior, we can conclude that the recognition model learns the required encoding.

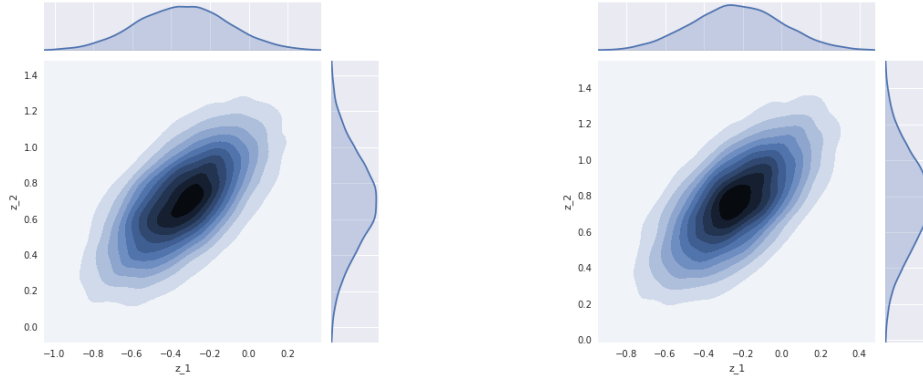


Figure 5.3: Learning  $\mathbf{R}$ , while keeping  $\mathbf{u}$  and  $\mathbf{V}$  fixed. The two plots show the posterior over local latent variable  $\mathbf{z}_n$  where the latent space is 2 dimensional. a. Exact posterior (Equation 5.1) b. Approximated posterior.

For the next experiment the setup was flipped; fixing the recognition model parameters,  $\mathbf{R}$  and  $\mathbf{S}$ , and training SEP (See Figure 5.4). In theory the predictive marginal likelihood should gradually increase towards the true marginal likelihood, and this is observed when large datasets were used; indicating that SEP requires more observations in order to generalize. When small datasets are used AESEP stabilizes at a local minima. Albeit in a noisy-fashion, predictive marginal likelihood is slowly increasing. However, observe that, unlike in VI, where we expect to see a monotonically increase marginal likelihood, this is not expected in EP-like algorithms.

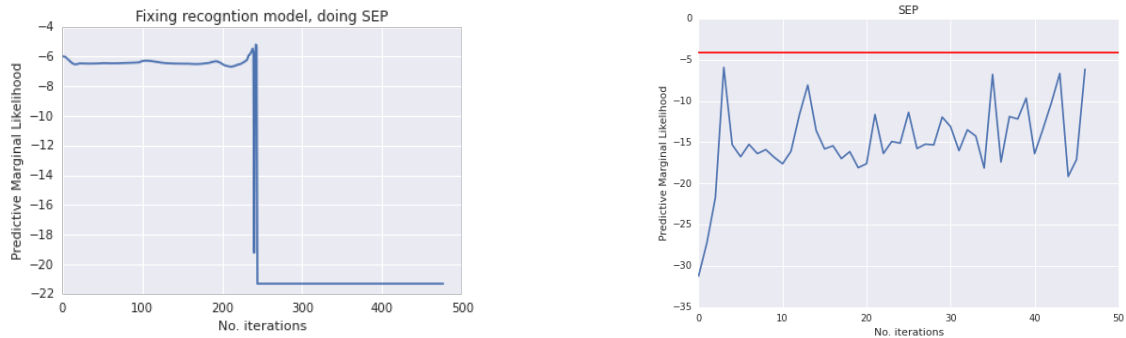


Figure 5.4: Fixing the recognition model, and learning the generative parameters using SEP. a. Small dataset b. Large dataset. The red line indicates true marginal likelihood.

The recognition model successfully trained with SEP fixed, and the previous experiment indicate signs of convergence while fixing the recognition model. However, more experiments are mandated for reliable conclusions. Figure 5.5 shows the predictive marginal likelihood



using full AESEP. Comparing Figure 5.4 and 5.5 similar patterns are observed, except that in the case of full AESEP the stabilization at a local minima happens much sooner when small datasets are used.

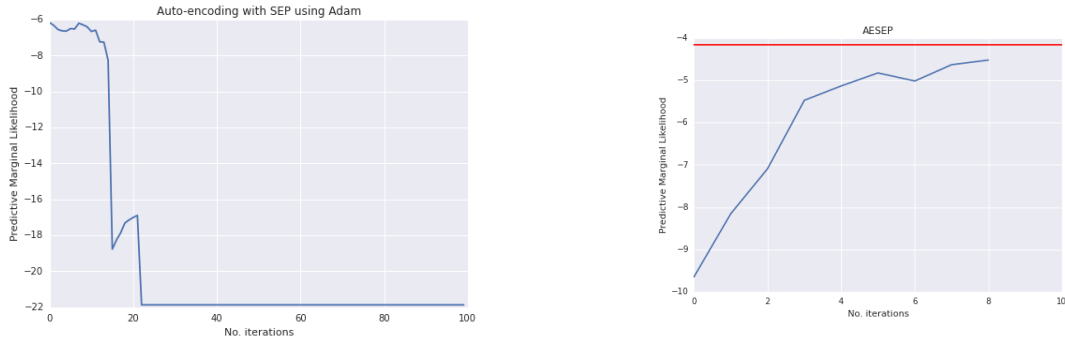


Figure 5.5: Predictive marginal likelihood for full AESEP with random initialization using a. Small Dataset b. Large dataset. The red line indicated true true marginal likelihood.

Experiments revealed AESEP to be difficult to tame; a behavior generally shared across EP-like algorithms. In the light of observed results, in training of AESEP, we propose the following starting points for further investigation. First, it is hypothesized that the two model constraints might be interfering: Does SEP and the recognition model have to learn at a similar rate (damping hyper-parameters)? Second, the theoretical claims of AESEP ought to be verified for larger and established datasets (e.g. MNIST). Third, lack of a global objective is restrictive in sensible visualizations of convergence. Predictive marginal likelihood are the recommended metrics. However, Energy function formulations are expected to ascribe an additional support for the theoretical claims. Fourth, profiling of the code reveals gradient computations as the bottleneck in time complexity, and hence alternative optimization strategies should be explored.

## 6 | Conclusion

This thesis has introduced AESEP for jointly performing Bayesian learning and inference in large latent variable models. We theoretically proved the existence of a local objective function that is minimized in the EP-fashion. With recent advances in stochastic optimization, we further developed a Monte Carlo approximation of the theoretical objective, and incorporated control variate schemes. Experiments indicate learning and inference on a factor analysis model. Applicability of AESEP to real-life scenarios is an open-ground and is left to future investigations.

# References

- Christopher M Bishop. Pattern recognition and machine learning, 2006.
- Ernest Fokoué. Bayesian computation of the intrinsic structure of factor analytic models. *Copyright© 2002-2009 Journal of Data Science*, 2009. (Page 19)
- Andrew Gelman, Aki Vehtari, Pasi Jylänki, Christian Robert, Nicolas Chopin, and John P Cunningham. Expectation propagation as a way of life. *arXiv preprint arXiv:1412.4869*, 2014. (Page 10)
- Matthew D Hoffman, David M Blei, Chong Wang, and John William Paisley. Stochastic variational inference. *Journal of Machine Learning Research*, 14(1):1303–1347, 2013. (Page 2)
- Diederik Kingma. Stochastic gradient variational bayes. [http://dpkingma.com/wordpress/wp-content/uploads/2014/05/2014-01\\_talk\\_ias.pdf](http://dpkingma.com/wordpress/wp-content/uploads/2014/05/2014-01_talk_ias.pdf), 2014. (Page 7)
- Diederik Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. (Page 21)
- Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013. (Pages 2 and 5)
- Yingzhen Li, José Miguel Hernández-Lobato, and Richard E Turner. Stochastic expectation propagation. In *Advances in Neural Information Processing Systems*, pages 2323–2331, 2015. (Pages 2 and 19)
- Yingzhen Li, Mark Rowland, Thang Bui, Daniel Hernandez-Lobato, Richard Turner, et al. Black-box alpha divergence minimization. In *Proceedings of The 33rd International Conference on Machine Learning*, pages 1511–1520, 2016. (Page 2)
- David JC MacKay. *Information theory, inference and learning algorithms*. Cambridge university press, 2003. (Page 1)

- 
- Thomas P Minka. The ep energy function and minimization schemes. *See [www.stat.cmu.edu/~minka/papers/learning.html](http://www.stat.cmu.edu/~minka/papers/learning.html)*, 2001a. (Page 11)
- Thomas P Minka. Expectation propagation for approximate bayesian inference. In *Proceedings of the Seventeenth conference on Uncertainty in artificial intelligence*, pages 362–369. Morgan Kaufmann Publishers Inc., 2001b. (Page 19)
- Thomas P Minka. *A family of algorithms for approximate Bayesian inference*. PhD thesis, Massachusetts Institute of Technology, 2001c. (Page 2)
- John Paisley, David Blei, and Michael Jordan. Variational bayesian inference with stochastic search. *arXiv preprint arXiv:1206.6430*, 2012. (Pages 5 and 6)
- Omiros Papaspiliopoulos, Gareth O Roberts, and Martin Sköld. A general framework for the parametrization of hierarchical models. *Statistical Science*, pages 59–73, 2007. (Page 5)
- Tim Salimans, Diederik P Kingma, Max Welling, et al. Markov chain monte carlo and variational inference: Bridging the gap. In *International Conference on Machine Learning*, pages 1218–1226, 2015. (Page 2)
- Richard Turner and Maneesh Sahani. Probabilistic amplitude and frequency demodulation. In *Advances in Neural Information Processing Systems*, pages 981–989, 2011. (Page 2)
- Richard E Turner and Maneesh Sahani. Two problems with variational expectation maximisation for time-series models. In *Workshop on Inference and Estimation in Probabilistic Time-Series Models*, volume 2, 2008. (Page 2)

# A | Derivation of Variational Lowerbound

$$\begin{aligned}\log p(\mathbf{x}) &= \log \int p_\theta(\mathbf{x}, \mathbf{z}) d\mathbf{z} \\ &= \log \int p_\theta(\mathbf{x}, \mathbf{z}) \frac{q_\phi(\mathbf{z}|\mathbf{x})}{q_\phi(\mathbf{z}|\mathbf{x})} d\mathbf{z} \\ &= \log \left( \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[ \frac{p_\theta(\mathbf{x}, \mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} \right] \right) \\ &\geq \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[ \log p_\theta(\mathbf{x}, \mathbf{z}) - \log q_\phi(\mathbf{z}|\mathbf{x}) \right] = \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi}; \mathbf{x})\end{aligned}$$

The last step uses Jensen's inequality. The final RHS is the variational lowerbound. VI considers the KL-divergence between the approximate and the exact posterior:

$$\begin{aligned}KL(q_\phi(\mathbf{z}|\mathbf{x})||p_\theta(\mathbf{z}|\mathbf{x})) &= \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \log \frac{q_\phi(\mathbf{z}|\mathbf{x})}{p_\theta(\mathbf{z}|\mathbf{x})} \\ &= \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[ -\log p(\mathbf{x}, \mathbf{z}) + \log p(\mathbf{x}) + \log q_\phi(\mathbf{z}|\mathbf{x}) \right] \\ &= -\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[ \log p_\theta(\mathbf{x}, \mathbf{z}) - \log q_\phi(\mathbf{z}|\mathbf{x}) \right] + \log p(\mathbf{x}) \\ &= -\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi}; \mathbf{x}) + \log p(\mathbf{x})\end{aligned}$$

# B | Moment Matching of EP

A distribution in the exponential family can be expressed in the following form:

$$p(\mathbf{z}|\boldsymbol{\eta}) = h(\mathbf{z})g(\boldsymbol{\eta})\exp\{\boldsymbol{\eta}^T u(\mathbf{z})\}$$

$\boldsymbol{\eta}$  are the natural parameters,  $g(\boldsymbol{\eta})$  ensures integration to unity:

$$g(\boldsymbol{\eta}) \int h(\mathbf{z})\exp\{\boldsymbol{\eta}^T u(\mathbf{z})\}d\mathbf{z} = 1$$

Differentiating this with respect to  $\boldsymbol{\eta}$  gives:

$$-\nabla \ln g(\boldsymbol{\eta}) = \mathbb{E}_{q(\boldsymbol{\eta})}[u(\mathbf{z})] \tag{B.1}$$

Taking the KL divergence between the exact,  $p$  and approximate distribution,  $q$ , gives:

$$\begin{aligned} KL(\mathbf{p}||\mathbf{q}) &= \int p(\mathbf{z})\log\frac{p(\mathbf{z})}{q(\mathbf{z})}d\boldsymbol{\theta} \\ &= \mathbb{E}_{p(\mathbf{z})}[\log p(\mathbf{z}) - \log q(\mathbf{z})] \\ &= \mathbb{E}_{p(\mathbf{z})}[\log q(\mathbf{z})] + \log p(\mathbf{z}) \\ &= \mathbb{E}_{p(\mathbf{z})}[-\ln(h(\mathbf{z})\exp\{\boldsymbol{\eta}^T u(\mathbf{z})\})] + \log p(\mathbf{z}) \\ &= \mathbb{E}_{p(\mathbf{z})}[-\ln g(\boldsymbol{\eta} - \boldsymbol{\eta}^T u(\mathbf{z}))] - \ln h(\mathbf{z}) + \log p(\mathbf{z}) \\ &= -\ln g(\boldsymbol{\eta}) - \boldsymbol{\eta}^T \mathbb{E}_{p(\mathbf{z})}[u(\mathbf{z})] + Constant \end{aligned}$$

Taking the gradient

$$-\nabla \ln g(\boldsymbol{\eta}) = \mathbb{E}_{p(\mathbf{z})}[u(\mathbf{z})] \tag{B.2}$$

Equation B.1 and B.2 imply moment matching:

$$\mathbb{E}_{q(\mathbf{z})}[u(\mathbf{z})] = \mathbb{E}_{p(\mathbf{z})}[u(\mathbf{z})]$$

# C | Exact EP on local latent variables of FA

The goal is to infer the latent variables  $\mathbf{z}_n$  with the following approximate distribution:

$$q^*(\mathbf{z}_{1:N}) = \prod_n f(\mathbf{z}_n)p(\mathbf{z}_n) = \prod_n \alpha_n \mathcal{N}(\mathbf{z}_n; \boldsymbol{\mu}_n, \boldsymbol{\Sigma}) \mathcal{N}(\mathbf{z}_n; \mathbf{0}, \mathbf{I}) = \prod_n q(\mathbf{z}_n)$$

EP steps

1. Cavity:

$$\frac{q^*(\mathbf{z}_{1:N})}{f_n(\mathbf{z}_n)} = q^{\setminus n}(\mathbf{z}_{1:N})$$

2. Tilted:

$$q^{\setminus n}(\mathbf{z}_{1:N})p(\mathbf{x}_n|\mathbf{z}_n, \boldsymbol{\theta})$$

3. Moment matching:

$$\begin{aligned} & KL\left(q^{\setminus n}(\mathbf{z}_{1:N})p(\mathbf{x}_n|\mathbf{z}_n, \boldsymbol{\theta}) \parallel q^{\setminus n}(\mathbf{z}_{1:N})f(\mathbf{z}_n)\right) \\ &= \int q^{\setminus n}(\mathbf{z}_{1:N})p(\mathbf{x}_n|\mathbf{z}_n, \boldsymbol{\theta}) \log \frac{q^{\setminus n}(\mathbf{z}_{1:N})p(\mathbf{x}_n|\mathbf{z}_n, \boldsymbol{\theta})}{q^{\setminus n}(\mathbf{z}_{1:N})f(\mathbf{z}_n)} \\ &= \int q^{\setminus n}(\mathbf{z}_{1:N})p(\mathbf{x}_n|\mathbf{z}_n, \boldsymbol{\theta}) \log \frac{p(\mathbf{x}_n|\mathbf{z}_n, \boldsymbol{\theta})}{f(\mathbf{z}_n)} \\ &\implies f_n(\mathbf{z}_n) = p(\mathbf{x}_n|\mathbf{z}_n, \boldsymbol{\theta}) \end{aligned}$$

From the generative model of factor analysis we know the form of the data likelihood, and as a result of exact moment matching we get the following result:

$$\begin{aligned}
f_n(\mathbf{z}_n) &= \mathcal{N}(\mathbf{x}_n; \mathbf{W}\mathbf{z}_n, \sigma_x^2 \mathbf{I}) \\
&= \frac{1}{\det(2\pi\sigma_x^2 \mathbf{I})} e^{-\frac{1}{2\sigma_x^2} (\mathbf{x}_n - \mathbf{W}\mathbf{z}_n)^T (\mathbf{x}_n - \mathbf{W}\mathbf{z}_n)} \\
&= \frac{1}{\det(2\pi\sigma_x^2 \mathbf{I})} e^{-\frac{1}{2\sigma_x^2} \mathbf{x}_n^T \mathbf{x}_n} e^{-\frac{1}{2\sigma_x^2} \mathbf{z}_n^T \mathbf{W}^T \mathbf{W} \mathbf{z}_n + \frac{1}{\sigma_x^2} \mathbf{z}_n^T \mathbf{W}^T \mathbf{x}_n} \\
&= \alpha_n \mathcal{N}(\mathbf{z}_n; \boldsymbol{\mu}_n, \boldsymbol{\Sigma})
\end{aligned}$$

By solving for  $\boldsymbol{\mu}_n$ ,  $\boldsymbol{\Sigma}$  and the likelihood scaling factor  $\alpha_n$ , a distribution over  $\mathbf{z}_n$  can be formed. By inspection, the covariance becomes:

$$\boldsymbol{\Sigma} = \left( \frac{\mathbf{W}^T \mathbf{W}}{\sigma_x^2} \right)^{-1}$$

In a similar manner, the mean is calculated:

$$\begin{aligned}
\mathbf{z}_n^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} &= \frac{1}{\sigma_x^2} \mathbf{z}_n^T \mathbf{W}^T \mathbf{x}_n \\
\frac{\mathbf{z}_n^T \mathbf{W}^T \mathbf{W} \boldsymbol{\mu}}{\sigma_x^2} &= \frac{\mathbf{z}_n^T \mathbf{W}^T \mathbf{x}_n}{\sigma_x^2} \\
\implies \boldsymbol{\mu} &= (\mathbf{W}^T \mathbf{W})^{-1} \mathbf{W}^T \mathbf{x}_n
\end{aligned}$$



Lastly, the likelihood scaling factor becomes:

$$\begin{aligned}
f_n(\mathbf{z}_n) &= \frac{1}{\det(2\pi\sigma_x^2\mathbf{I})} e^{-\frac{1}{2\sigma_x^2}\mathbf{x}_n^T\mathbf{x}_n} e^{-\frac{1}{2\sigma_x^2}\mathbf{z}_n^T\mathbf{W}^T\mathbf{W}\mathbf{z}_n + \frac{1}{\sigma_x^2}\mathbf{z}_n^T\mathbf{W}^T\mathbf{x}_n} \\
&\text{by defining the following variables} \\
a &= \frac{1}{\det(2\pi\sigma_x^2\mathbf{I})} e^{-\frac{1}{2\sigma_x^2}\mathbf{x}_n^T\mathbf{x}_n} \\
b &= -\frac{1}{2\sigma_x^2}\mathbf{z}_n^T\mathbf{W}^T\mathbf{W}\mathbf{z}_n + \frac{1}{\sigma_x^2}\mathbf{z}_n^T\mathbf{W}^T\mathbf{x}_n \\
&= a \times e^b \times \frac{e^{-\frac{1}{2}\boldsymbol{\mu}^T\boldsymbol{\Sigma}^{-1}\boldsymbol{\mu}}}{(2\pi)^{\frac{d}{2}}|\boldsymbol{\Sigma}|^{\frac{1}{2}}} \times \frac{(2\pi)^{\frac{d}{2}}|\boldsymbol{\Sigma}|^{\frac{1}{2}}}{e^{-\frac{1}{2}\boldsymbol{\mu}^T\boldsymbol{\Sigma}^{-1}\boldsymbol{\mu}}} \\
f_n(\mathbf{z}_n) &= \mathcal{N}(\mathbf{z}_n, \boldsymbol{\mu}, \boldsymbol{\Sigma}) \times \frac{a(2\pi)^{\frac{d}{2}}|\boldsymbol{\Sigma}|^{\frac{1}{2}}}{e^{-\frac{1}{2}\boldsymbol{\mu}^T\boldsymbol{\Sigma}^{-1}\boldsymbol{\mu}}} \\
\implies \alpha_n &= \frac{a(2\pi)^{\frac{d}{2}}|\boldsymbol{\Sigma}|^{\frac{1}{2}}}{e^{-\frac{1}{2}\boldsymbol{\mu}^T\boldsymbol{\Sigma}^{-1}\boldsymbol{\mu}}}
\end{aligned}$$

Substituting in the calculated mean and covariance and assuming the dimensions of  $\mathbf{z}$  is 2, the final likelihood factor takes the following tedious form.

$$\begin{aligned}
\alpha_n &= \frac{1}{\det(2\pi\sigma_x^2\mathbf{I})} e^{-\frac{1}{2\sigma_x^2}\mathbf{x}_n^T\mathbf{x}_n} \times \frac{(2\pi)^{\frac{d}{2}}|\boldsymbol{\Sigma}|^{\frac{1}{2}}}{e^{-\frac{1}{2}\boldsymbol{\mu}^T\boldsymbol{\Sigma}^{-1}\boldsymbol{\mu}}} \\
&= \frac{(2\pi)^{\frac{d}{2}}|\boldsymbol{\Sigma}|^{\frac{1}{2}}}{\det(2\pi\sigma_x^2\mathbf{I})} \times e^{-\frac{1}{2\sigma_x^2}\mathbf{x}_n^T\mathbf{x}_n + \frac{1}{2}\boldsymbol{\mu}^T\boldsymbol{\Sigma}^{-1}\boldsymbol{\mu}} \\
\boldsymbol{\mu}^T\boldsymbol{\Sigma}^{-1}\boldsymbol{\mu} &= \left( (\mathbf{W}^T\mathbf{W})^{-1}\mathbf{W}^T\mathbf{x}_n \right)^T \frac{\mathbf{W}^T\mathbf{W}}{\sigma_x^2} \left( (\mathbf{W}^T\mathbf{W})^{-1}\mathbf{W}^T\mathbf{x}_n \right) \\
&= \frac{1}{\sigma_x^2} \left( (\mathbf{W}^T\mathbf{W})^{-1}\mathbf{W}^T\mathbf{x}_n \right)^T (\mathbf{W}^T\mathbf{W})(\mathbf{W}^T\mathbf{W})^{-1}\mathbf{W}^T\mathbf{x}_n \\
&= \frac{1}{\sigma_x^2} \mathbf{x}_n^T \mathbf{W} (\mathbf{W}^T\mathbf{W})^{-1} \mathbf{W}^T \mathbf{x}_n \\
\alpha_n &= \frac{(2\pi)^{\frac{d}{2}}|\boldsymbol{\Sigma}|^{\frac{1}{2}}}{|2\pi\sigma_x^2\mathbf{I}|} e^{-\frac{1}{2\sigma_x^2}\mathbf{x}_n^T\mathbf{x}_n + \frac{1}{2\sigma_x^2}\mathbf{x}_n^T\mathbf{W}(\mathbf{W}^T\mathbf{W})^{-1}\mathbf{W}^T\mathbf{x}_n}
\end{aligned}$$