

Extending the Bayesian Deep Learning Method MultiSWAG

Scott Brownlie

Master of Science
School of Informatics
University of Edinburgh
2021

Abstract

This skeleton demonstrates how to use the `infthesis` style for MSc dissertations in Artificial Intelligence, Cognitive Science, Computer Science, Data Science, and Informatics. It also emphasises the page limit, and that you must not deviate from the required style. The file `skeleton.tex` generates this document and can be used as a starting point for your thesis. The abstract should summarise your report and fit in the space on the first page.

Acknowledgements

Any acknowledgements go here.

Table of Contents

1	Introduction	1
1.1	Using Sections	2
1.2	Citations	2
2	Your next chapter	3
3	Factor Analysis	4
3.1	Online Gradient Algorithm for Factor Analysis	5
3.1.1	Derivatives with respect to \mathbf{F}	6
3.1.2	Derivatives with respect to Ψ	6
3.2	Online EM for Factor Analysis	7
4	Conclusions	9
4.1	Final Reminder	9
	Bibliography	10

Chapter 1

Introduction

The preliminary material of your report should contain:

- The title page.
- An abstract page.
- Optionally an acknowledgements page.
- The table of contents.

As in this example `skeleton.tex`, the above material should be included between:

```
\begin{preliminary}  
...  
\end{preliminary}
```

This style file uses roman numeral page numbers for the preliminary material.

The main content of the dissertation, starting with the first chapter, starts with page 1. ***The main content must not go beyond page 40.***

The report then contains a bibliography and any appendices, which may go beyond page 40. The appendices are only for any supporting material that's important to go on record. However, you cannot assume markers of dissertations will read them.

You may not change the dissertation format (e.g., reduce the font size, change the margins, or reduce the line spacing from the default 1.5 spacing). Over length or incorrectly-formatted dissertations will not be accepted and you would have to modify your dissertation and resubmit. You cannot assume we will check your submission before the final deadline and if it requires resubmission after the deadline to conform to the page and style requirements you will be subject to the usual late penalties based on your final submission time.

1.1 Using Sections

Divide your chapters into sub-parts as appropriate.

1.2 Citations

Citations (such as [?] or [?]) can be generated using BibTeX. For more advanced usage, the `natbib` package is recommended. You could also consider the newer `biblatex` system.

These examples use a numerical citation style. You may also use (Author, Date) format if you prefer.

Chapter 2

Your next chapter

A dissertation usually contains several chapters.

Chapter 3

Factor Analysis

FA is a latent variable model which generates observations $\theta \in \mathbb{R}^d$ as follows. First, a latent vector $\mathbf{h} \in \mathbb{R}^K$, for some $K < d$, is sampled from $p(\mathbf{h}) = \mathcal{N}(\mathbf{0}, \mathbf{I})$. Next, \mathbf{h} is transformed onto a K -dimensional linear subspace of \mathbb{R}^d by left-multiplying it by a *factor loading* matrix $\mathbf{F} \in \mathbb{R}^{d \times K}$. The origin of this subspace is then shifted by adding a bias term $\mathbf{c} \in \mathbb{R}^d$. Finally, the data is perturbed by adding some zero mean Gaussian noise $\varepsilon \in \mathbb{R}^d$ sampled from $\mathcal{N}(\mathbf{0}, \Psi)$, where Ψ is a $d \times d$ diagonal matrix [1]. Putting all this together, an observation $\theta \in \mathbb{R}^d$ is generated according to

$$\theta = \mathbf{F}\mathbf{h} + \mathbf{c} + \varepsilon. \quad (3.1)$$

In this context, an observation θ is the parameter vector of a neural network.

It follows that, given \mathbf{h} , the observations θ are Gaussian distributed with mean $\mathbf{F}\mathbf{h} + \mathbf{c}$ and covariance Ψ . Formally,

$$p(\theta|\mathbf{h}) = \mathcal{N}(\mathbf{F}\mathbf{h} + \mathbf{c}, \Psi) = \frac{1}{\sqrt{(2\pi)^d |\Psi|}} \exp\left(-\frac{1}{2}(\theta - \mathbf{F}\mathbf{h} - \mathbf{c})^\top \Psi^{-1}(\theta - \mathbf{F}\mathbf{h} - \mathbf{c})\right), \quad (3.2)$$

where $|\Psi|$ is the *determinant* of Ψ . From [1], integrating $p(\theta|\mathbf{h})$ over \mathbf{h} gives the marginal distribution

$$p(\theta) = \mathcal{N}(\mathbf{c}, \mathbf{F}\mathbf{F}^\top + \Psi). \quad (3.3)$$

The parameters of the model are \mathbf{c} , \mathbf{F} and Ψ . The value of \mathbf{c} which maximises the likelihood of the observed data is the empirical mean of the observations [1], which in this case is θ_{SWA} . Having set the bias term, an EM or SVD algorithm can find the maximum likelihood estimates of \mathbf{F} and Ψ [1]. However, both methods require storing all the observations in memory, making them impractical for high-dimensional data with lots of observations. Two alternative online algorithms are presented in Section 3.1 and 3.2.

3.1 Online Gradient Algorithm for Factor Analysis

In situations where learning latent variable models with the classic EM algorithm is slow, [1] suggests optimising the log-likelihood of the model parameters via gradient methods. Since FA is a latent variable model, this approach can be applied here. In this case the log-likelihood of the parameters \mathbf{F} and Ψ is

$$L(\mathbf{F}, \Psi) = \frac{1}{T} \sum_{t=1}^T \log p(\theta_t | \mathbf{F}, \Psi). \quad (3.4)$$

By adapting the gradient derivation for general latent variable model in [1] to FA,

$$\begin{aligned} \nabla_{\mathbf{F}, \Psi} \log p(\theta | \mathbf{F}, \Psi) &= \frac{1}{p(\theta | \mathbf{F}, \Psi)} \nabla_{\mathbf{F}, \Psi} p(\theta | \mathbf{F}, \Psi) \\ &= \frac{1}{p(\theta | \mathbf{F}, \Psi)} \nabla_{\mathbf{F}, \Psi} \int_{\mathbf{h}} p(\theta, \mathbf{h} | \mathbf{F}, \Psi) \\ &= \frac{1}{p(\theta | \mathbf{F}, \Psi)} \int_{\mathbf{h}} \nabla_{\mathbf{F}, \Psi} p(\theta, \mathbf{h} | \mathbf{F}, \Psi) \\ &= \frac{1}{p(\theta | \mathbf{F}, \Psi)} \int_{\mathbf{h}} p(\theta, \mathbf{h} | \mathbf{F}, \Psi) \nabla_{\mathbf{F}, \Psi} \log p(\theta, \mathbf{h} | \mathbf{F}, \Psi) \quad (3.5) \\ &= \int_{\mathbf{h}} \frac{p(\theta, \mathbf{h} | \mathbf{F}, \Psi)}{p(\theta | \mathbf{F}, \Psi)} \nabla_{\mathbf{F}, \Psi} \log p(\theta, \mathbf{h} | \mathbf{F}, \Psi) \\ &= \int_{\mathbf{h}} p(\mathbf{h} | \theta, \mathbf{F}, \Psi) \nabla_{\mathbf{F}, \Psi} \log p(\theta, \mathbf{h} | \mathbf{F}, \Psi) \\ &= \mathbb{E}_{\mathbf{h} \sim p(\mathbf{h} | \theta, \mathbf{F}, \Psi)} [\nabla_{\mathbf{F}, \Psi} \log p(\theta, \mathbf{h} | \mathbf{F}, \Psi)]. \end{aligned}$$

Now, using the fact that $\mathbf{h} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ is independent from \mathbf{F} and Ψ ,

$$\begin{aligned} \nabla_{\mathbf{F}, \Psi} \log p(\theta, \mathbf{h} | \mathbf{F}, \Psi) &= \nabla_{\mathbf{F}, \Psi} \log (p(\theta | \mathbf{h}, \mathbf{F}, \Psi) p(\mathbf{h} | \mathbf{F}, \Psi)) \\ &= \nabla_{\mathbf{F}, \Psi} \log (p(\theta | \mathbf{h}, \mathbf{F}, \Psi) p(\mathbf{h})) \quad (3.6) \\ &= \nabla_{\mathbf{F}, \Psi} (\log p(\theta | \mathbf{h}, \mathbf{F}, \Psi) + \log p(\mathbf{h})) \\ &= \nabla_{\mathbf{F}, \Psi} \log p(\theta | \mathbf{h}, \mathbf{F}, \Psi). \end{aligned}$$

Substituting Equation (3.6) into Equation (3.5),

$$\nabla_{\mathbf{F}, \Psi} \log p(\theta | \mathbf{F}, \Psi) = \mathbb{E}_{\mathbf{h} \sim p(\mathbf{h} | \theta, \mathbf{F}, \Psi)} [\nabla_{\mathbf{F}, \Psi} \log p(\theta | \mathbf{h}, \mathbf{F}, \Psi)]. \quad (3.7)$$

Note that $p(\theta | \mathbf{h}, \mathbf{F}, \Psi)$ is just the Gaussian distribution in Equation (3.2) for the given values of \mathbf{F} and Ψ . Hence, substituting $\mathbf{c} = \theta_{\text{SWA}}$ into this Gaussian and applying the logarithm,

$$\begin{aligned} \log p(\theta | \mathbf{h}, \mathbf{F}, \Psi) &= -\frac{1}{2} (\theta - \mathbf{F}\mathbf{h} - \theta_{\text{SWA}})^\top \Psi^{-1} (\theta - \mathbf{F}\mathbf{h} - \theta_{\text{SWA}}) - \frac{1}{2} \log |\Psi| - \frac{d}{2} \log 2\pi \\ &= -\frac{1}{2} (\mathbf{d} - \mathbf{F}\mathbf{h})^\top \Psi^{-1} (\mathbf{d} - \mathbf{F}\mathbf{h}) - \frac{1}{2} \log |\Psi| - \frac{d}{2} \log 2\pi, \quad (3.8) \end{aligned}$$

where $\mathbf{d} = \boldsymbol{\theta} - \boldsymbol{\theta}_{\text{SWA}}$.

3.1.1 Derivatives with respect to \mathbf{F}

Differentiating Equation (3.8) with respect to \mathbf{F} ,

$$\nabla_{\mathbf{F}} \log p(\boldsymbol{\theta}|\mathbf{h}, \mathbf{F}, \boldsymbol{\Psi}) = \boldsymbol{\Psi}^{-1}(\mathbf{d} - \mathbf{F}\mathbf{h})\mathbf{h}^\top. \quad (3.9)$$

It follows that $\nabla_{\mathbf{F}} \log p(\boldsymbol{\theta}|\mathbf{F}, \boldsymbol{\Psi})$ is the expected value of $\boldsymbol{\Psi}^{-1}(\mathbf{d} - \mathbf{F}\mathbf{h})\mathbf{h}^\top$ over the distribution $p(\mathbf{h}|\boldsymbol{\theta}, \mathbf{F}, \boldsymbol{\Psi})$. Letting $\mathbb{E}[\cdot]$ denote $\mathbb{E}_{\mathbf{h} \sim p(\mathbf{h}|\boldsymbol{\theta}, \mathbf{F}, \boldsymbol{\Psi})}[\cdot]$ to simplify the notation,

$$\begin{aligned} \nabla_{\mathbf{F}} \log p(\boldsymbol{\theta}|\mathbf{F}, \boldsymbol{\Psi}) &= \mathbb{E}[\boldsymbol{\Psi}^{-1}(\mathbf{d} - \mathbf{F}\mathbf{h})\mathbf{h}^\top] \\ &= \mathbb{E}[\boldsymbol{\Psi}^{-1}\mathbf{d}\mathbf{h}^\top] - \mathbb{E}[\boldsymbol{\Psi}^{-1}\mathbf{F}\mathbf{h}\mathbf{h}^\top] \\ &= \boldsymbol{\Psi}^{-1}\mathbf{d}\mathbb{E}[\mathbf{h}^\top] - \boldsymbol{\Psi}^{-1}\mathbf{F}\mathbb{E}[\mathbf{h}\mathbf{h}^\top]. \end{aligned} \quad (3.10)$$

From the E-step of the EM algorithm in [1], $p(\mathbf{h}|\boldsymbol{\theta}, \mathbf{F}, \boldsymbol{\Psi}) \propto \mathcal{N}(\mathbf{m}, \Sigma)$, where

$$\mathbf{m} = (\mathbf{I} + \mathbf{F}^\top \boldsymbol{\Psi}^{-1} \mathbf{F})^{-1} \mathbf{F}^\top \boldsymbol{\Psi}^{-1} \mathbf{d} \quad \text{and} \quad \Sigma = (\mathbf{I} + \mathbf{F}^\top \boldsymbol{\Psi}^{-1} \mathbf{F})^{-1}. \quad (3.11)$$

Hence, substituting $\mathbb{E}[\mathbf{h}^\top] = \mathbf{m}^\top$ and $\mathbb{E}[\mathbf{h}\mathbf{h}^\top] = \Sigma + \mathbf{m}\mathbf{m}^\top$ into Equation (3.10),

$$\nabla_{\mathbf{F}} \log p(\boldsymbol{\theta}|\mathbf{F}, \boldsymbol{\Psi}) = \boldsymbol{\Psi}^{-1}\mathbf{d}\mathbf{m}^\top - \boldsymbol{\Psi}^{-1}\mathbf{F}(\Sigma + \mathbf{m}\mathbf{m}^\top). \quad (3.12)$$

3.1.2 Derivatives with respect to $\boldsymbol{\Psi}$

In order to differentiate Equation (3.8) with respect to $\boldsymbol{\Psi}$, it helps to use the fact that $\boldsymbol{\Psi}$ is a diagonal matrix. First consider $\mathbf{X}^{-1} = \text{diag}(\frac{1}{x_1}, \dots, \frac{1}{x_d})$ and $\mathbf{a} = (a_1, \dots, a_d)^\top$. Then

$$\mathbf{a}^\top \mathbf{X}^{-1} \mathbf{a} = \sum_{i=1}^d \frac{a_i^2}{x_i}, \quad (3.13)$$

and so

$$\frac{\partial}{\partial x_i} \mathbf{a}^\top \mathbf{X}^{-1} \mathbf{a} = \frac{-a_i^2}{x_i^2} \quad (3.14)$$

for $i = 1, \dots, d$. Since the partial derivatives of Equation (3.13) with respect to the off-diagonal entries of \mathbf{X} are zero,

$$\nabla_{\mathbf{X}}(\mathbf{a}^\top \mathbf{X}^{-1} \mathbf{a}) = \text{diag}\left(\frac{-a_1^2}{x_1^2}, \dots, \frac{-a_d^2}{x_d^2}\right) = -\text{diag}\left(\text{diag}(\mathbf{X}^{-2}) \odot (\mathbf{a} \odot \mathbf{a})\right), \quad (3.15)$$

where \odot denotes the element-wise matrix product. Note that when applied to a d -length vector, $\text{diag}(\cdot)$ represents the $d \times d$ diagonal matrix with the vector on its diagonal, and when applied to a $d \times d$ matrix, $\text{diag}(\cdot)$ represents the d -length vector consisting of the diagonal entries of the matrix.

Setting $\mathbf{X} = \Psi$ and $\mathbf{a} = \mathbf{d} - \mathbf{Fh}$, it follows that

$$\nabla_{\Psi}(\mathbf{d} - \mathbf{Fh})^{\top} \Psi^{-1}(\mathbf{d} - \mathbf{Fh}) = -\text{diag} \left(\text{diag}(\Psi^{-2}) \odot ((\mathbf{d} - \mathbf{Fh}) \odot (\mathbf{d} - \mathbf{Fh})) \right). \quad (3.16)$$

Also, using the identity $\nabla_{\mathbf{X}} \log |\mathbf{X}| = \mathbf{X}^{-\top}$ and the fact that $\Psi^{-\top} = \Psi^{-1}$,

$$\nabla_{\Psi} \log |\Psi| = \Psi^{-1}. \quad (3.17)$$

Hence, using Equation (3.16) and Equation (3.17), the partial derivatives of (3.8) with respect to Ψ are

$$\nabla_{\Psi} \log p(\theta | \mathbf{h}, \mathbf{F}, \Psi) = \frac{1}{2} \text{diag} \left(\text{diag}(\Psi^{-2}) \odot ((\mathbf{d} - \mathbf{Fh}) \odot (\mathbf{d} - \mathbf{Fh})) \right) - \frac{1}{2} \Psi^{-1}. \quad (3.18)$$

Again, letting $\mathbb{E}[\cdot]$ denote $\mathbb{E}_{\mathbf{h} \sim p(\mathbf{h} | \theta, \mathbf{F}, \Psi)}[\cdot]$, it follows that

$$\begin{aligned} 2 \cdot \nabla_{\Psi} \log p(\theta | \mathbf{F}, \Psi) &= \mathbb{E} \left[\text{diag} \left(\text{diag}(\Psi^{-2}) \odot ((\mathbf{d} - \mathbf{Fh}) \odot (\mathbf{d} - \mathbf{Fh})) \right) - \Psi^{-1} \right] \\ &= \text{diag} \left(\mathbb{E} \left[\text{diag}(\Psi^{-2}) \odot ((\mathbf{d} - \mathbf{Fh}) \odot (\mathbf{d} - \mathbf{Fh})) \right] \right) - \mathbb{E}[\Psi^{-1}] \\ &= \text{diag} \left(\text{diag}(\Psi^{-2}) \odot \mathbb{E}[(\mathbf{d} - \mathbf{Fh}) \odot (\mathbf{d} - \mathbf{Fh})] \right) - \Psi^{-1}. \end{aligned} \quad (3.19)$$

The expectation inside Equation (3.19) is

$$\begin{aligned} \mathbb{E}[(\mathbf{d} - \mathbf{Fh}) \odot (\mathbf{d} - \mathbf{Fh})] &= \mathbb{E}[\mathbf{d} \odot \mathbf{d}] - 2\mathbb{E}[\mathbf{d} \odot \mathbf{Fh}] + \mathbb{E}[\mathbf{Fh} \odot \mathbf{Fh}] \\ &= \mathbf{d} \odot \mathbf{d} - 2\mathbf{d} \odot \mathbf{F}\mathbb{E}[\mathbf{h}] + \mathbb{E}[\text{diag}(\mathbf{Fh}\mathbf{h}^{\top} \mathbf{F})] \\ &= \mathbf{d} \odot \mathbf{d} - 2\mathbf{d} \odot \mathbf{Fm} + \text{diag}(\mathbb{E}[\mathbf{Fh}\mathbf{h}^{\top} \mathbf{F}]) \\ &= \mathbf{d} \odot \mathbf{d} - 2\mathbf{d} \odot \mathbf{Fm} + \text{diag}(\mathbf{F}\mathbb{E}[\mathbf{h}\mathbf{h}^{\top}] \mathbf{F}^{\top}) \\ &= \mathbf{d} \odot \mathbf{d} - 2\mathbf{d} \odot \mathbf{Fm} + \text{diag}(\mathbf{F}(\Sigma + \mathbf{m}\mathbf{m}^{\top}) \mathbf{F}^{\top}) \\ &= \mathbf{d} \odot \mathbf{d} - 2\mathbf{d} \odot \mathbf{Fm} + \text{sum}(\mathbf{F}(\Sigma + \mathbf{m}\mathbf{m}^{\top}) \odot \mathbf{F}, \text{dim} = 1), \end{aligned} \quad (3.20)$$

where $\text{sum}(\cdot, \text{dim} = 1)$ denotes the operation of summing along the rows of a matrix.

Finally, substituting Equation (3.20) into Equation (3.19) and rearranging,

$$\begin{aligned} \nabla_{\Psi} \log p(\theta | \mathbf{F}, \Psi) &= \frac{1}{2} \text{diag} \left(\text{diag}(\Psi^{-2}) \odot (\mathbf{d} \odot \mathbf{d} - 2\mathbf{d} \odot \mathbf{Fm} \right. \\ &\quad \left. + \text{sum}(\mathbf{F}(\Sigma + \mathbf{m}\mathbf{m}^{\top}) \odot \mathbf{F}, \text{dim} = 1)) \right) - \frac{1}{2} \Psi^{-1}. \end{aligned} \quad (3.21)$$

3.2 Online EM for Factor Analysis

The batch EM algorithm for FA in [1] can be adapted to an online version. The E-step of the batch algorithm sets the variational distribution $q(\mathbf{h} | \theta_t, \mathbf{F}, \Psi) \propto \mathcal{N}(\mathbf{m}_t, \Sigma)$

for each θ_t , where \mathbf{m}_t and Σ are the parameters in Equation (3.11) with \mathbf{d} replaced by $\mathbf{d}_t = \theta_t - \theta_{\text{SWA}}$. This can be done separately for each θ_t as it is sampled, using the current estimates of \mathbf{F} and Ψ . The only other detail is that θ_{SWA} , which is not available during training, must be replaced by the running average $\bar{\theta}_t$.

Modifying the M-step requires a bit more thought, as it involves summing over all θ_t . The M-step sets

$$\mathbf{F} = \mathbf{A}\mathbf{H}^{-1}, \quad (3.22)$$

where

$$\mathbf{A} = \frac{1}{T} \sum_{t=1}^T \mathbf{d}_t \mathbf{m}_t^\top \quad \text{and} \quad \mathbf{H} = \Sigma + \frac{1}{T} \sum_{t=1}^T \mathbf{m}_t \mathbf{m}_t^\top, \quad (3.23)$$

and

$$\Psi = \text{diag} \left(\frac{1}{T} \sum_{t=1}^T \mathbf{d}_t \mathbf{d}_t^\top - 2\mathbf{F}\mathbf{A}^\top + \mathbf{F}\mathbf{H}\mathbf{F}^\top \right). \quad (3.24)$$

Batch EM iterates the E and M-steps above until \mathbf{F} and Ψ converge. On each iteration of the M-step, all components of the sums in Equation (3.23) and Equation (3.24) are updated. Clearly, this is not possible in an online algorithm which only holds a single θ_t in memory at any one time. A compromise is to update the sums incrementally on epoch t , with \mathbf{d}_t and \mathbf{m}_t derived from θ_t and $\bar{\theta}_t$, and then fix these components of the sums for the remainder of the algorithm. This is the approach that will be adopted in the project.

Chapter 4

Conclusions

4.1 Final Reminder

The body of your dissertation, before the references and any appendices, *must* finish by page 40. The introduction, after preliminary material, should have started on page 1.

You may not change the dissertation format (e.g., reduce the font size, change the margins, or reduce the line spacing from the default 1.5 spacing). Over length or incorrectly-formatted dissertations will not be accepted and you would have to modify your dissertation and resubmit. You cannot assume we will check your submission before the final deadline and if it requires resubmission after the deadline to conform to the page and style requirements you will be subject to the usual late penalties based on your final submission time.

Bibliography

- [1] David Barber. *Bayesian Reasoning and Machine Learning*. Cambridge University Press, 2012.