

# **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

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Background</b>	<b>2</b>
2.1	Factor Analysis . . . . .	2
2.2	Bayesian Linear Regression . . . . .	3
2.2.1	Computing the Posterior . . . . .	3
2.3	Variational Inference . . . . .	5
<b>3</b>	<b>Related Work and Research Questions</b>	<b>6</b>
3.1	Related Work . . . . .	6
<b>4</b>	<b>Online Factor Analysis</b>	<b>8</b>
4.1	Online Stochastic Gradient Ascent . . . . .	8
4.1.1	Partial derivatives with respect to $\mathbf{F}$ . . . . .	10
4.1.2	Partial derivatives with respect to $\Psi$ . . . . .	10
4.1.3	Practical implementation . . . . .	12
4.2	Online Expectation-Maximisation . . . . .	13
4.2.1	E-step . . . . .	14
4.2.2	M-step . . . . .	14
4.2.3	Practical implementation . . . . .	15
4.3	Variational Inference . . . . .	16
4.3.1	Derivatives of the Variational Distribution . . . . .	17
4.3.2	Derivatives of the Prior . . . . .	18
4.3.3	Derivatives of the Likelihood . . . . .	18
4.3.4	Practical Implementation . . . . .	19
4.3.5	Comparison of VIFA and SLANG . . . . .	20
4.4	Factors Initialisation . . . . .	20

<b>5</b>	<b>Online Factor Analysis Experiments</b>	<b>23</b>
5.1	Methodology . . . . .	23
5.2	Results and Discussion . . . . .	25
<b>6</b>	<b>Linear Regression Experiments</b>	<b>29</b>
6.1	Posterior Estimation . . . . .	29
6.1.1	Methodology . . . . .	30
6.1.2	Results and Discussion . . . . .	31
6.2	Making Predictions . . . . .	35
6.2.1	Methodology . . . . .	37
6.2.2	Results and Discussion . . . . .	38
<b>7</b>	<b>Conclusions</b>	<b>40</b>
7.1	Final Reminder . . . . .	40
	<b>Bibliography</b>	<b>41</b>
<b>A</b>	<b>Supplementary Results</b>	<b>43</b>
A.1	Linear Regression Predictions . . . . .	43
A.2	Comparison of VIFA and SLANG . . . . .	43

# 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.

# Chapter 2

## Background

### 2.1 Factor Analysis

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  $\Psi$  is a  $D \times D$  diagonal matrix [1]. Putting all this together, an observation  $\theta \in \mathbb{R}^D$  is generated according to

$$\theta = \mathbf{F}\mathbf{h} + \mathbf{c} + \varepsilon. \quad (2.1)$$

In the context of this thesis, an observation  $\theta$  is the parameter vector of a neural network.

It follows that, given  $\mathbf{h}$ , the observations  $\theta$  are Gaussian distributed with mean  $\mathbf{F}\mathbf{h} + \mathbf{c}$  and covariance  $\Psi$  [1]. Formally,

$$p(\theta|\mathbf{h}) = \mathcal{N}(\mathbf{F}\mathbf{h} + \mathbf{c}, \Psi) = \frac{1}{\sqrt{(2\pi)^D |\Psi|}} \exp\left(-\frac{1}{2}(\theta - \mathbf{F}\mathbf{h} - \mathbf{c})^\top \Psi^{-1}(\theta - \mathbf{F}\mathbf{h} - \mathbf{c})\right), \quad (2.2)$$

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

$$p(\theta) = \mathcal{N}(\mathbf{c}, \mathbf{F}\mathbf{F}^\top + \Psi). \quad (2.3)$$

The parameters of the model are  $\mathbf{c}$ ,  $\mathbf{F}$  and  $\Psi$ . 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  $\theta_{\text{SWA}}$ . Having set the bias term, an expectation-maximisation (EM)

or singular value decomposition (SVD) algorithm can find the maximum likelihood estimates of  $\mathbf{F}$  and  $\Psi$  [1]. However, both methods require storing all the observations in memory, making them impractical for high-dimensional data, such as the parameter vectors of deep neural networks. Two alternative online algorithms are presented in Chapter 4.

## 2.2 Bayesian Linear Regression

A linear regression model is a mapping from vector inputs  $\mathbf{x} \in \mathbb{R}^D$  to scalar outputs  $y \in \mathbb{R}$  via an affine transformation. Given a set of observed input-output pairs,  $\mathcal{D} = \{(\mathbf{x}_n, y_n)\}_{n=1}^N$ , it is assumed that each output is generated according to

$$y_n = \theta^\top \mathbf{x}_n + \varepsilon \quad (2.4)$$

for some unknown  $\theta \in \mathbb{R}^D$ . The underlying signal,  $\theta^\top \mathbf{x}_n$ , is corrupted by additive noise,

$$\varepsilon \sim \mathcal{N}(0, \sigma^2), \quad (2.5)$$

for some  $\sigma > 0$  [1]. The model is often written with an explicit bias term, but this can be absorbed into  $\theta$  by adding a constant of one to the input, leading to the expression in Equation (2.4).

### 2.2.1 Computing the Posterior

Due to the additive noise, each  $y_n$  is a random variable, conditioned on  $\mathbf{x}_n$  and  $\theta$ . Since  $\varepsilon$  is Gaussian distributed, the conditional pdf of  $y_n$  is

$$p(y_n | \theta, \mathbf{x}_n) = \mathcal{N}(\theta^\top \mathbf{x}_n, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(y_n - \theta^\top \mathbf{x}_n)^2\right). \quad (2.6)$$

Assuming that the observations in  $\mathcal{D}$  are independent and identically distributed (iid), the log-likelihood of  $\theta$  having generated the data is

$$\begin{aligned} \log p(\mathcal{D} | \theta) &= \sum_{n=1}^N [\log p(y_n | \theta, \mathbf{x}_n) + \log p(\mathbf{x}_n)] \\ &= \sum_{n=1}^N \left[ -\frac{1}{2} \log 2\pi\sigma^2 - \frac{1}{2\sigma^2} (y_n - \theta^\top \mathbf{x}_n)^2 \right] + \sum_{n=1}^N \log p(\mathbf{x}_n) \\ &= -\frac{1}{2\sigma^2} \sum_{n=1}^N (y_n - \theta^\top \mathbf{x}_n)^2 - \frac{N}{2} \log \sigma^2 - \frac{N}{2} \log 2\pi + \sum_{n=1}^N \log p(\mathbf{x}_n) \\ &= -\frac{\beta}{2} \sum_{n=1}^N (y_n - \theta^\top \mathbf{x}_n)^2 + \frac{N}{2} \log \beta + \text{constant}, \end{aligned} \quad (2.7)$$



where  $\beta = \frac{1}{\sigma^2}$  [1]. In Bayesian linear regression, a prior distribution for  $\theta$  is also specified. A common choice is

$$p(\theta) = \mathcal{N}(\mathbf{0}, \alpha^{-1} \mathbf{I}) = \left( \frac{\alpha}{2\pi} \right)^{\frac{D}{2}} \exp \left( -\frac{\alpha}{2} \theta^\top \theta \right) \quad (2.8)$$

for some  $\alpha > 0$ , which is a hyperparameter known as the *precision* [1]. Applying the logarithm,

$$\begin{aligned} \log p(\theta) &= \frac{D}{2} \log \frac{\alpha}{2\pi} - \frac{\alpha}{2} \theta^\top \theta \\ &= -\frac{\alpha}{2} \theta^\top \theta + \frac{D}{2} \log \alpha + \text{constant}. \end{aligned} \quad (2.9)$$

Now, using Bayes' rule and Equation (2.7) and Equation (2.9), it follows that the log-posterior distribution of  $\theta$  for fixed  $\alpha$  and  $\beta$  is

$$\begin{aligned} \log p(\theta | \mathcal{D}) &= \log \frac{p(\mathcal{D} | \theta) p(\theta)}{p(\mathcal{D})} \\ &= \log p(\mathcal{D} | \theta) + \log p(\theta) - \log p(\mathcal{D}) \\ &= -\frac{\beta}{2} \sum_{n=1}^N (y_n - \theta^\top \mathbf{x}_n)^2 - \frac{\alpha}{2} \theta^\top \theta + \text{constant} \\ &= -\frac{\beta}{2} \sum_{n=1}^N (y_n^2 - 2y_n \theta^\top \mathbf{x}_n + (\theta^\top \mathbf{x}_n)^2) - \frac{\alpha}{2} \theta^\top \theta + \text{constant} \\ &= -\frac{\beta}{2} \sum_{n=1}^N y_n^2 + \beta \sum_{n=1}^N y_n \theta^\top \mathbf{x}_n - \frac{\beta}{2} \sum_{n=1}^N \theta^\top \mathbf{x}_n \mathbf{x}_n^\top \theta - \frac{\alpha}{2} \theta^\top \theta + \text{constant} \\ &= \left( \beta \sum_{n=1}^N y_n \mathbf{x}_n \right)^\top \theta - \frac{1}{2} \theta^\top \left( \alpha \mathbf{I} + \beta \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top \right) \theta + \text{constant} \\ &= \mathbf{b}^\top \theta - \frac{1}{2} \theta^\top \mathbf{A} \theta + \text{constant}, \end{aligned} \quad (2.10)$$

where

$$\mathbf{A} = \alpha \mathbf{I} + \beta \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top \quad \text{and} \quad \mathbf{b} = \beta \sum_{n=1}^N y_n \mathbf{x}_n. \quad (2.11)$$

Now, using a result from [1] to complete the square of Equation (2.10),

$$\begin{aligned} \log p(\theta | \mathcal{D}) &= -\frac{1}{2} (\theta - \mathbf{A}^{-1} \mathbf{b})^\top \mathbf{A} (\theta - \mathbf{A}^{-1} \mathbf{b}) + \frac{1}{2} \mathbf{b}^\top \mathbf{A}^{-1} \mathbf{b} + \text{constant} \\ &= -\frac{1}{2} (\theta - \mathbf{A}^{-1} \mathbf{b})^\top \mathbf{A} (\theta - \mathbf{A}^{-1} \mathbf{b}) + \text{constant} \\ &= -\frac{1}{2} (\theta - \mathbf{m})^\top \mathbf{S}^{-1} (\theta - \mathbf{m}) + \text{constant}, \end{aligned} \quad (2.12)$$

where

$$\mathbf{m} = \mathbf{A}^{-1} \mathbf{b} \quad \text{and} \quad \mathbf{S} = \mathbf{A}^{-1}. \quad (2.13)$$

Hence, the posterior distribution of  $\theta$  is

$$p(\theta | \mathcal{D}) = \mathcal{N}(\mathbf{m}, \mathbf{S}). \quad (2.14)$$

## 2.3 Variational Inference

For linear regression, computing the posterior in Equation (2.14) is possible because the log-likelihood,  $p(\mathcal{D}|\theta)$ , is a quadratic expression of  $\theta$ . However, when  $\theta$  is the parameter vector of a neural network with even a single non-linear hidden layer, the log-likelihood is not so simple and the posterior cannot be written in closed form. Instead, the marginal likelihood,  $p(\mathcal{D}) = \int_{\theta} p(\mathcal{D}|\theta)p(\theta)$ , must be evaluated, which is intractable for high-dimensional  $\theta$ .

An alternative strategy is to approximate the posterior directly with a Gaussian distribution  $q(\theta) = \mathcal{N}(\mu, \Sigma)$ , where the parameters  $\mu$  and  $\Sigma$  are chosen to minimise some measure of dissimilarity between  $q(\theta)$  and  $p(\theta|\mathcal{D})$ . Because  $q(\theta)$  is essentially a variable in an optimisation procedure, this approach is known as *variational inference* (VI) [1]. A common measure of dissimilarity between two distributions is the Kullback-Leibler (KL) divergence [1], which in this case is

$$\begin{aligned} \mathbb{D}_{KL}[q(\theta)||p(\theta|\mathcal{D})] &= \mathbb{E}_{q(\theta)} \left[ \log \frac{q(\theta)}{p(\theta|\mathcal{D})} \right] \\ &= \mathbb{E}_{q(\theta)} \left[ \log \frac{q(\theta)p(\mathcal{D})}{p(\mathcal{D}|\theta)p(\theta)} \right] \\ &= \mathbb{E}_{q(\theta)} [\log q(\theta) - \log p(\theta) - \log p(\mathcal{D}|\theta) + \log p(\mathcal{D})]. \end{aligned} \quad (2.15)$$

Since  $p(\mathcal{D})$  is a constant, it can be ignored and the optimal distribution can be obtained by solving

$$\min_{\mu, \Sigma} [\mathbb{E}_{q(\theta)} [\log q(\theta)] - \mathbb{E}_{q(\theta)} [\log p(\theta)] - \mathbb{E}_{q(\theta)} [\log p(\mathcal{D}|\theta)]]. \quad (2.16)$$

One way to do this is by using a stochastic gradient algorithm, which approximates the true gradient of the objective function with Monte Carlo estimates of the form

$$\frac{1}{ML} \sum_{m=1}^M \sum_{l=1}^L [\nabla_{\mu, \Sigma} \log q(\theta_l) - \nabla_{\mu, \Sigma} \log p(\theta_l) - \nabla_{\mu, \Sigma} \log p(y_m|\mathbf{x}_m, \theta_l)], \quad (2.17)$$

where  $\{(y_m, \mathbf{x}_m)\}_{m=1}^M$  is a batch of observed data sampled from  $\mathcal{D}$  and  $\theta_l \sim \mathcal{N}(\mu, \Sigma)$ . Normally it would not be possible to compute partial derivatives with respect to  $\mu$  and  $\Sigma$ , given that they only take part via the random sampling operation. However, this can in fact be done by using the *re-parameterisation trick* [3]. Instead of sampling from  $\mathcal{N}(\mu, \Sigma)$  directly, a random variable  $\mathbf{z}_l$  is sampled from  $\mathcal{N}(\mathbf{0}, \mathbf{I})$  and then transformed to  $\theta_l$  via a deterministic operation involving  $\mu$  and  $\Sigma$ . This simple trick means that Equation (2.17) can be evaluated and used in conjunction with SGD or any of its variants to optimise the objective in Equation (2.16). This is the idea behind the *auto-encoding variational Bayes* algorithm [8].

# Chapter 3

## Related Work and Research Questions

### 3.1 Related Work

Several methods have been suggested for optimising the variational objective in Equation (2.16). The Variational Online Gauss-Newton (VOGN) method [7] is one such approach. It learns the parameters of the variational distribution,  $\mu$  and  $\Sigma$ , via an approximate natural gradient algorithm. The online update rule is

$$\mu \leftarrow \mu - \eta_\mu \Sigma (\hat{\mathbf{g}}_\theta + \alpha \mu), \quad \Sigma^{-1} \leftarrow (1 - \eta_\Sigma) \Sigma^{-1} + \eta_\Sigma (\hat{\mathbf{G}}_\theta + \alpha \mathbf{I}), \quad (3.1)$$

where  $\eta_\mu, \eta_\Sigma > 0$  are learning rates,  $\alpha$  is the precision of the prior and  $\hat{\mathbf{g}}_\theta$  and  $\hat{\mathbf{G}}_\theta$  are, respectively, a Monte Carlo estimate of the gradient of the negative log-likelihood and an Empirical Fisher (EM) matrix, given  $\theta \sim q(\theta) = \mathcal{N}(\mu, \Sigma)$ . Formally,

$$\hat{\mathbf{g}}_\theta = -\frac{N}{M} \sum_{(\mathbf{x}_i, y_i) \in \mathcal{B}} \mathbf{g}_\theta^i, \quad \hat{\mathbf{G}}_\theta = -\frac{N}{M} \sum_{(\mathbf{x}_i, y_i) \in \mathcal{B}} \mathbf{g}_\theta^i (\mathbf{g}_\theta^i)^\top, \quad (3.2)$$

where  $\mathbf{g}_\theta^i = \nabla_\theta \log p(y_i | \mathbf{x}_i, \theta)$  and  $\mathcal{B}$  is a mini-batch of  $M$  training examples sampled from the full dataset  $\mathcal{D}$  of size  $N$ . When  $\theta$  is the parameter vector of a neural network, the gradients  $\mathbf{g}_\theta^i$  can be computed via back-propagation [17]. However, for large neural networks with  $D$  parameters, the update in Equation (3.1) is in general infeasible, as it involves inverting the  $D \times D$  covariance matrix  $\Sigma$ .

A computationally tractable version of VOGN can be obtained by using a low-rank plus diagonal approximation of the inverse of the covariance matrix. This is the approach taken by SLANG [13], which stands for stochastic, low-rank, approximate natural gradient. This algorithm sets  $\Sigma^{-1} = \mathbf{U}\mathbf{U}^\top + \mathbf{D}$  for some  $\mathbf{U} \in \mathbb{R}^{D \times K}$  and a  $D \times D$  diagonal matrix  $\mathbf{D}$ . Substituting this definition into Equation (3.1), the update can be

written as

$$\Sigma^{-1} \leftarrow [(1 - \eta_\Sigma)\mathbf{U}\mathbf{U}^\top + \eta_\Sigma\hat{\mathbf{G}}_\theta] + [(1 - \eta_\Sigma)\mathbf{D} + \eta_\Sigma\alpha\mathbf{I}]. \quad (3.3)$$

To make this tractable, the non-diagonal component of the right-hand side is approximated as

$$(1 - \eta_\Sigma)\mathbf{U}\mathbf{U}^\top + \eta_\Sigma\hat{\mathbf{G}}_\theta \approx \mathbf{Q}_{1:K}\Lambda_{1:K}\mathbf{Q}_{1:K}^\top, \quad (3.4)$$

where  $\Lambda_{1:K}$  is a  $K \times K$  diagonal matrix containing the  $K$  largest eigenvalues of  $(1 - \eta_\Sigma)\mathbf{U}\mathbf{U}^\top + \eta_\Sigma\hat{\mathbf{G}}_\theta$ , and the columns of  $\mathbf{Q}_{1:K} \in \mathbb{R}^{D \times K}$  are the corresponding eigenvectors. In order to compute this eigen-decomposition, the individual gradients  $\mathbf{g}_\theta^i$  for training examples  $(\mathbf{x}_i, y_i) \in \mathcal{B}$  are required. When using the standard back-propagation algorithm in deep learning frameworks such as PyTorch [14], this is only possible by looping over the mini-batch samples one-by-one and storing each  $\mathbf{g}_\theta^i$ , which is very inefficient. A faster approach - which was adopted by SLANG - is outlined in [4] for standard neural networks. However, the authors of SLANG admit that this speed-up is limited to feed-forward layers, and a separate implementation would be required to handle each different type of layer. This limits the applicability of SLANG in the case of more general neural networks, such as those with commonly used convolutional [9], residual [5], recurrent [6] or transformer layers [18].

# Chapter 4

## Online Factor Analysis

This chapter describes three algorithms which can be used to approximate the posterior distribution of the parameter vector of a neural network with a FA model. The first two algorithms fit a FA model to samples of the posterior obtained via some separate procedure, while the third algorithm estimates the FA posterior directly. Crucially, all three algorithms are executed online in an iterative fashion, making them suitable for high-dimensional deep neural networks.

### 4.1 Online Stochastic Gradient Ascent

Suppose first that samples from the posterior,  $\theta_1, \dots, \theta_T$ , are readily available. 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  $\Psi$  given the observed data is

$$L(\mathbf{F}, \Psi) = \frac{1}{T} \sum_{t=1}^T \log p(\theta_t | \mathbf{F}, \Psi). \quad (4.1)$$

The partial derivatives of the log-likelihood with respect to  $\mathbf{F}$  and  $\Psi$  are therefore

$$\nabla_{\mathbf{F}, \Psi} L(\mathbf{F}, \Psi) = \frac{1}{T} \sum_{t=1}^T \nabla_{\mathbf{F}, \Psi} \log p(\theta_t | \mathbf{F}, \Psi). \quad (4.2)$$

Computing  $\nabla_{\mathbf{F}, \Psi} L(\mathbf{F}, \Psi)$  in full would require a pass over all observations,  $\theta_1, \dots, \theta_T$ . However, a stochastic gradient algorithm can be used instead, as long as the expectation of the sample derivatives is proportional to  $\nabla_{\mathbf{F}, \Psi} L(\mathbf{F}, \Psi)$ . By using  $\nabla_{\mathbf{F}, \Psi} \log p(\theta_t | \mathbf{F}, \Psi)$ ,  $t = 1, \dots, T$ , as the sample derivatives, this condition clearly holds. Hence, as long as

the partial derivatives  $\nabla_{\mathbf{F}, \Psi} \log p(\theta_t | \mathbf{F}, \Psi)$  can be computed efficiently, they can be used in conjunction with SGD or any of its variants to optimise  $\mathbf{F}$  and  $\Psi$  online.

By adapting an argument for general latent variable models in [1] to FA, the required sample derivatives can be written as

$$\begin{aligned}
 \nabla_{\mathbf{F}, \Psi} \log p(\theta_t | \mathbf{F}, \Psi) &= \frac{1}{p(\theta_t | \mathbf{F}, \Psi)} \nabla_{\mathbf{F}, \Psi} p(\theta_t | \mathbf{F}, \Psi) \\
 &= \frac{1}{p(\theta_t | \mathbf{F}, \Psi)} \nabla_{\mathbf{F}, \Psi} \int_{\mathbf{h}_t} p(\theta_t, \mathbf{h}_t | \mathbf{F}, \Psi) \\
 &= \frac{1}{p(\theta_t | \mathbf{F}, \Psi)} \int_{\mathbf{h}_t} \nabla_{\mathbf{F}, \Psi} p(\theta_t, \mathbf{h}_t | \mathbf{F}, \Psi) \\
 &= \frac{1}{p(\theta_t | \mathbf{F}, \Psi)} \int_{\mathbf{h}_t} p(\theta_t, \mathbf{h}_t | \mathbf{F}, \Psi) \nabla_{\mathbf{F}, \Psi} \log p(\theta_t, \mathbf{h}_t | \mathbf{F}, \Psi) \quad (4.3) \\
 &= \int_{\mathbf{h}_t} \frac{p(\theta_t, \mathbf{h}_t | \mathbf{F}, \Psi)}{p(\theta_t | \mathbf{F}, \Psi)} \nabla_{\mathbf{F}, \Psi} \log p(\theta_t, \mathbf{h}_t | \mathbf{F}, \Psi) \\
 &= \int_{\mathbf{h}_t} p(\mathbf{h}_t | \theta_t, \mathbf{F}, \Psi) \nabla_{\mathbf{F}, \Psi} \log p(\theta_t, \mathbf{h}_t | \mathbf{F}, \Psi) \\
 &= \mathbb{E}_{p(\mathbf{h}_t | \theta_t, \mathbf{F}, \Psi)} [\nabla_{\mathbf{F}, \Psi} \log p(\theta_t, \mathbf{h}_t | \mathbf{F}, \Psi)].
 \end{aligned}$$

This is as far as the derivation in [1] goes. However, given the form of the FA model, it is possible to manipulate the sample derivatives further. In particular, using the fact that  $\mathbf{h}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  is independent from  $\mathbf{F}$  and  $\Psi$ ,

$$\begin{aligned}
 \nabla_{\mathbf{F}, \Psi} \log p(\theta_t, \mathbf{h}_t | \mathbf{F}, \Psi) &= \nabla_{\mathbf{F}, \Psi} \log (p(\theta_t | \mathbf{h}_t, \mathbf{F}, \Psi) p(\mathbf{h}_t | \mathbf{F}, \Psi)) \\
 &= \nabla_{\mathbf{F}, \Psi} \log (p(\theta_t | \mathbf{h}_t, \mathbf{F}, \Psi) p(\mathbf{h}_t)) \\
 &= \nabla_{\mathbf{F}, \Psi} (\log p(\theta_t | \mathbf{h}_t, \mathbf{F}, \Psi) + \log p(\mathbf{h}_t)) \quad (4.4) \\
 &= \nabla_{\mathbf{F}, \Psi} \log p(\theta_t | \mathbf{h}_t, \mathbf{F}, \Psi).
 \end{aligned}$$

Substituting Equation (4.4) into Equation (4.3),

$$\nabla_{\mathbf{F}, \Psi} \log p(\theta_t | \mathbf{F}, \Psi) = \mathbb{E}_{p(\mathbf{h}_t | \theta_t, \mathbf{F}, \Psi)} [\nabla_{\mathbf{F}, \Psi} \log p(\theta_t | \mathbf{h}_t, \mathbf{F}, \Psi)]. \quad (4.5)$$

Note that  $p(\theta_t | \mathbf{h}_t, \mathbf{F}, \Psi)$  is just the Gaussian distribution in Equation (2.2), given  $\mathbf{F}$  and  $\Psi$ . Hence, substituting  $\mathbf{c} = \theta_{\text{SWA}}$  into Equation (2.2) and applying the logarithm,

$$\begin{aligned}
 \log p(\theta_t | \mathbf{h}_t, \mathbf{F}, \Psi) &= -\frac{1}{2} (\theta_t - \mathbf{F} \mathbf{h}_t - \theta_{\text{SWA}})^\top \Psi^{-1} (\theta_t - \mathbf{F} \mathbf{h}_t - \theta_{\text{SWA}}) - \frac{1}{2} \log |\Psi| - \frac{D}{2} \log 2\pi \\
 &= -\frac{1}{2} (\mathbf{d}_t - \mathbf{F} \mathbf{h}_t)^\top \Psi^{-1} (\mathbf{d}_t - \mathbf{F} \mathbf{h}_t) - \frac{1}{2} \log |\Psi| - \frac{D}{2} \log 2\pi, \quad (4.6)
 \end{aligned}$$

where  $\mathbf{d}_t = \theta_t - \theta_{\text{SWA}}$ . This is convenient, since Equation (4.6) can be differentiated with respect to both  $\mathbf{F}$  and  $\Psi$ . Of course, this requires access to  $\theta_{\text{SWA}}$ , which is not

available during training. As a compromise - and following the approach of the SWAG covariance approximation -  $\theta_{\text{SWA}}$  can be replaced by the running average of the neural network's parameter vectors.

#### 4.1.1 Partial derivatives with respect to $\mathbf{F}$

From [16], for any symmetric matrix  $\mathbf{W}$ ,

$$\nabla_{\mathbf{A}}(\mathbf{x} - \mathbf{A}\mathbf{s})^T \mathbf{W}(\mathbf{x} - \mathbf{A}\mathbf{s}) = -2\mathbf{W}(\mathbf{x} - \mathbf{A}\mathbf{s})\mathbf{s}^T. \quad (4.7)$$

Hence, differentiating Equation (4.6) with respect to  $\mathbf{F}$  gives

$$\nabla_{\mathbf{F}} \log p(\theta_t | \mathbf{h}_t, \mathbf{F}, \Psi) = \Psi^{-1}(\mathbf{d}_t - \mathbf{F}\mathbf{h}_t)\mathbf{h}_t^T. \quad (4.8)$$

It then follows from Equation (4.5) that  $\nabla_{\mathbf{F}} \log p(\theta_t | \mathbf{F}, \Psi)$  is the expected value of  $\Psi^{-1}(\mathbf{d}_t - \mathbf{F}\mathbf{h}_t)\mathbf{h}_t^T$  over the distribution  $p(\mathbf{h}_t | \theta_t, \mathbf{F}, \Psi)$ . Letting  $\mathbb{E}[\cdot]$  denote  $\mathbb{E}_{p(\mathbf{h}_t | \theta_t, \mathbf{F}, \Psi)}[\cdot]$  to simplify the notation,

$$\begin{aligned} \nabla_{\mathbf{F}} \log p(\theta_t | \mathbf{F}, \Psi) &= \mathbb{E}[\Psi^{-1}(\mathbf{d}_t - \mathbf{F}\mathbf{h}_t)\mathbf{h}_t^T] \\ &= \Psi^{-1}(\mathbb{E}[\mathbf{d}_t\mathbf{h}_t^T] - \mathbb{E}[\mathbf{F}\mathbf{h}_t\mathbf{h}_t^T]) \\ &= \Psi^{-1}(\mathbf{d}_t\mathbb{E}[\mathbf{h}_t^T] - \mathbf{F}\mathbb{E}[\mathbf{h}_t\mathbf{h}_t^T]). \end{aligned} \quad (4.9)$$

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

$$\Sigma = (\mathbf{I} + \mathbf{F}^T \Psi^{-1} \mathbf{F})^{-1} \quad \text{and} \quad \mathbf{m}_t = \Sigma \mathbf{F}^T \Psi^{-1} \mathbf{d}_t. \quad (4.10)$$

Hence, using identities from [16],

$$\mathbb{E}[\mathbf{h}_t^T] = \mathbf{m}_t^T \quad \text{and} \quad \mathbb{E}[\mathbf{h}_t\mathbf{h}_t^T] = \Sigma + \mathbf{m}_t\mathbf{m}_t^T. \quad (4.11)$$

Finally, substituting Equation (4.11) into Equation (4.9),

$$\nabla_{\mathbf{F}} \log p(\theta_t | \mathbf{F}, \Psi) = \Psi^{-1}(\mathbf{d}_t\mathbf{m}_t^T - \mathbf{F}(\Sigma + \mathbf{m}_t\mathbf{m}_t^T)). \quad (4.12)$$

#### 4.1.2 Partial derivatives with respect to $\Psi$

In order to differentiate Equation (4.6) with respect to  $\Psi$ , it helps to use the fact that  $\Psi$  is diagonal. First consider  $\mathbf{X}^{-1} = \text{diag}(\frac{1}{x_1}, \dots, \frac{1}{x_D})$  and  $\mathbf{a} = (a_1, \dots, a_D)^T$ . Then

$$\mathbf{a}^T \mathbf{X}^{-1} \mathbf{a} = \sum_{d=1}^D \frac{a_d^2}{x_d}, \quad (4.13)$$

and so

$$\frac{\partial}{\partial x_d} \mathbf{a}^\top \mathbf{X}^{-1} \mathbf{a} = \frac{-a_d^2}{x_d^2} \quad (4.14)$$

for  $d = 1, \dots, D$ . Since the partial derivatives of Equation (4.13) with respect to the off-diagonal entries of  $\mathbf{X}$  are zero,

$$\begin{aligned} \nabla_{\mathbf{X}}(\mathbf{a}^\top \mathbf{X}^{-1} \mathbf{a}) &= \text{diag}\left(\frac{-a_1^2}{x_1^2}, \dots, \frac{-a_D^2}{x_D^2}\right) \\ &= -\text{diag}(\text{diag}(\mathbf{X}^{-2}) \odot \mathbf{a}^2), \end{aligned} \quad (4.15)$$

where  $\odot$  denotes the element-wise matrix product (with broadcasting, if applicable) and the square of a vector is applied element-wise. Also, 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. Substituting  $\mathbf{X} = \Psi$  and  $\mathbf{a} = \mathbf{d}_t - \mathbf{F}\mathbf{h}_t$ , into Equation (4.15),

$$\nabla_{\Psi}(\mathbf{d}_t - \mathbf{F}\mathbf{h}_t)^\top \Psi^{-1}(\mathbf{d}_t - \mathbf{F}\mathbf{h}_t) = -\text{diag}(\text{diag}(\Psi^{-2}) \odot (\mathbf{d}_t - \mathbf{F}\mathbf{h}_t)^2). \quad (4.16)$$

Also, using the identity  $\nabla_{\mathbf{X}} \log |\mathbf{X}| = \mathbf{X}^{-\top}$  from [16] and the fact that  $\Psi^{-\top} = \Psi^{-1}$ ,

$$\nabla_{\Psi} \log |\Psi| = \Psi^{-1}. \quad (4.17)$$

Hence, the partial derivatives of Equation (4.6) with respect to  $\Psi$  are

$$\nabla_{\Psi} \log p(\theta_t | \mathbf{h}_t, \mathbf{F}, \Psi) = \frac{1}{2} \text{diag}(\text{diag}(\Psi^{-2}) \odot (\mathbf{d}_t - \mathbf{F}\mathbf{h}_t)^2) - \frac{1}{2} \Psi^{-1}. \quad (4.18)$$

Again, letting  $\mathbb{E}[\cdot]$  denote  $\mathbb{E}_{p(\mathbf{h}_t | \theta_t, \mathbf{F}, \Psi)}[\cdot]$ , it follows from Equation (4.5) that

$$\begin{aligned} 2 \cdot \nabla_{\Psi} \log p(\theta_t | \mathbf{F}, \Psi) &= \mathbb{E}[\text{diag}(\text{diag}(\Psi^{-2}) \odot (\mathbf{d}_t - \mathbf{F}\mathbf{h}_t)^2) - \Psi^{-1}] \\ &= \text{diag}(\mathbb{E}[\text{diag}(\Psi^{-2}) \odot (\mathbf{d}_t - \mathbf{F}\mathbf{h}_t)^2]) - \mathbb{E}[\Psi^{-1}] \\ &= \text{diag}(\text{diag}(\Psi^{-2}) \odot \mathbb{E}[(\mathbf{d}_t - \mathbf{F}\mathbf{h}_t)^2]) - \Psi^{-1}. \end{aligned} \quad (4.19)$$

Now, expanding the expectation inside Equation (4.19) and substituting in the expressions from Equation (4.11),

$$\begin{aligned} \mathbb{E}[(\mathbf{d}_t - \mathbf{F}\mathbf{h}_t)^2] &= \mathbb{E}[\mathbf{d}_t \odot \mathbf{d}_t] - 2\mathbb{E}[\mathbf{d}_t \odot (\mathbf{F}\mathbf{h}_t)] + \mathbb{E}[(\mathbf{F}\mathbf{h}_t) \odot (\mathbf{F}\mathbf{h}_t)] \\ &= \mathbf{d}_t \odot \mathbf{d}_t - 2\mathbf{d}_t \odot (\mathbf{F}\mathbb{E}[\mathbf{h}_t]) + \mathbb{E}[\text{diag}(\mathbf{F}\mathbf{h}_t \mathbf{h}_t^\top \mathbf{F}^\top)] \\ &= \mathbf{d}_t \odot \mathbf{d}_t - 2\mathbf{d}_t \odot (\mathbf{F}\mathbf{m}_t) + \text{diag}(\mathbb{E}[\mathbf{F}\mathbf{h}_t \mathbf{h}_t^\top \mathbf{F}^\top]) \\ &= \mathbf{d}_t \odot \mathbf{d}_t - 2\mathbf{d}_t \odot (\mathbf{F}\mathbf{m}_t) + \text{diag}(\mathbf{F}\mathbb{E}[\mathbf{h}_t \mathbf{h}_t^\top] \mathbf{F}^\top) \\ &= \mathbf{d}_t \odot \mathbf{d}_t - 2\mathbf{d}_t \odot (\mathbf{F}\mathbf{m}_t) + \text{diag}(\mathbf{F}(\Sigma + \mathbf{m}_t \mathbf{m}_t^\top) \mathbf{F}^\top). \end{aligned} \quad (4.20)$$



Finally, substituting Equation (4.20) into Equation (4.19) and rearranging,

$$\begin{aligned} \nabla_{\Psi} \log p(\theta_t | \mathbf{F}, \Psi) = & \text{diag} \left( \frac{1}{2} \text{diag}(\Psi^{-2}) \odot \left( \mathbf{d}_t \odot \mathbf{d}_t - 2\mathbf{d}_t \odot (\mathbf{F}\mathbf{m}_t) \right. \right. \\ & \left. \left. + \text{diag}(\mathbf{F}(\Sigma + \mathbf{m}_t\mathbf{m}_t^T)\mathbf{F}^T) \right) \right) - \frac{1}{2}\Psi^{-1}. \end{aligned} \quad (4.21)$$

### 4.1.3 Practical implementation

In practice, storing the full  $D \times D$  covariance matrix  $\Psi$  (or its inverse) would be infeasible for high-dimensional  $\theta_t \in \mathbb{R}^D$ . However, since  $\Psi$  is diagonal and the partial derivatives of the log-likelihood with respect to the off-diagonal are always zero, it suffices to work with the diagonal entries only. All occurrences of  $D \times D$  matrices can then be removed from the gradient calculations.

First, note that the partial derivatives with respect to both  $\mathbf{F}$  and  $\Psi$  depend on  $\mathbf{m}_t$  and  $\Sigma$  from Equation (4.10), which themselves depend on  $\Psi^{-1}$ . Let  $\psi = \text{diag}(\Psi)$  and  $\psi^{-n} = \text{diag}(\Psi^{-n})$  for  $n \in \mathbb{N}^+$ . Then

$$\mathbf{F}^T \Psi^{-1} = (\mathbf{F} \odot \psi^{-1})^T. \quad (4.22)$$

Now, setting  $\mathbf{C} = (\mathbf{F} \odot \psi^{-1})^T$  and substituting into Equation (4.10), it follows that

$$\Sigma = (\mathbf{I} + \mathbf{C}\mathbf{F})^{-1} \quad \text{and} \quad \mathbf{m}_t = \Sigma \mathbf{C} \mathbf{d}_t. \quad (4.23)$$

These more efficient values can then be used in Equation (4.12), which itself can be simplified to

$$\nabla_{\mathbf{F}} \log p(\theta_t | \mathbf{F}, \Psi) = \psi^{-1} \odot (\mathbf{d}_t \mathbf{m}_t^T - \mathbf{F}(\Sigma + \mathbf{m}_t \mathbf{m}_t^T)). \quad (4.24)$$

Also, the partial derivatives with respect to  $\Psi$  in Equation (4.21) can be re-written with respect to  $\psi$ , as

$$\begin{aligned} \nabla_{\Psi} \log p(\theta_t | \mathbf{F}, \Psi) = & \frac{1}{2} \psi^{-2} \odot \left( \mathbf{d}_t \odot \mathbf{d}_t - 2\mathbf{d}_t \odot (\mathbf{F}\mathbf{m}_t) \right. \\ & \left. + \text{sum}((\mathbf{F}(\Sigma + \mathbf{m}_t \mathbf{m}_t^T)) \odot \mathbf{F}, \text{dim} = 1) \right) - \frac{1}{2} \psi^{-1}, \end{aligned} \quad (4.25)$$

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

One final point is that the elements of  $\psi$  must be positive, since they represent variances. One way to achieve this is by using the re-parameterisation  $\psi = \exp \beta$  for

some unconstrained parameter  $\beta \in \mathbb{R}^D$ . Then the gradient updates can be performed on  $\beta$  instead of  $\psi$ . Using the chain rule from calculus,

$$\begin{aligned}\nabla_{\beta} \log p(\theta_t | \mathbf{F}, \Psi) &= \nabla_{\psi} \log p(\theta_t | \mathbf{F}, \Psi) \odot \nabla_{\beta} \exp \beta \\ &= \nabla_{\psi} \log p(\theta_t | \mathbf{F}, \Psi) \odot \exp \beta \\ &= \nabla_{\psi} \log p(\theta_t | \mathbf{F}, \Psi) \odot \psi,\end{aligned}\tag{4.26}$$

where  $\nabla_{\psi} \log p(\theta_t | \mathbf{F}, \Psi)$  is given in Equation (4.25). Pseudo code for the practical implementation is given in Algorithm 1.

---

**Algorithm 1** Online Stochastic Gradient Ascent for Factor Analysis

---

**Input:** Observation dimension  $D$ , latent dimension  $K$ , learning rate  $\alpha > 0$

- 1: Initialise  $\mathbf{F} \in \mathbb{R}^{D \times K}$ ,  $\psi = \mathbf{1}^D$ ,  $\bar{\theta}_0 = \mathbf{0}^D$
  - 2:  $\beta \leftarrow \log \psi$
  - 3: **for**  $t = 1, \dots, T$  **do**
  - 4:   Sample observation  $\theta_t \in \mathbb{R}^D$
  - 5:    $\bar{\theta}_t \leftarrow \bar{\theta}_{t-1} + \frac{1}{t} (\bar{\theta}_t - \bar{\theta}_{t-1})$
  - 6:    $\mathbf{d}_t \leftarrow \theta_t - \bar{\theta}_t$
  - 7:    $\mathbf{C} \leftarrow (\mathbf{F} \odot \psi^{-1})^\top$  (with broadcasting)
  - 8:    $\Sigma \leftarrow (\mathbf{I} + \mathbf{C}\mathbf{F})^{-1}$
  - 9:    $\mathbf{m}_t \leftarrow \Sigma \mathbf{C} \mathbf{d}_t$
  - 10:   Compute  $\nabla_{\mathbf{F}} \log p(\theta_t | \mathbf{F}, \Psi)$  according to Equation (4.24)
  - 11:   Compute  $\nabla_{\beta} \log p(\theta_t | \mathbf{F}, \Psi)$  according to Equation (4.26)
  - 12:    $\mathbf{F} \leftarrow \mathbf{F} + \alpha \nabla_{\mathbf{F}} \log p(\theta_t | \mathbf{F}, \Psi)$
  - 13:    $\beta \leftarrow \beta + \alpha \nabla_{\beta} \log p(\theta_t | \mathbf{F}, \Psi)$
  - 14:    $\psi \leftarrow \exp \beta$
  - 15: **end for**
  - 16:  $\theta_{\text{SWA}} \leftarrow \bar{\theta}_T$
  - 17: **return**  $\theta_{\text{SWA}}, \mathbf{F}, \psi$
- 

## 4.2 Online Expectation-Maximisation

Again, suppose that samples from the posterior,  $\theta_1, \dots, \theta_T$ , are available. The classic EM algorithm for FA iteratively optimises the log-likelihood of  $\mathbf{F}$  and  $\Psi$  by alternating “E” and “M” steps until convergence. Using properties of the Kullback-Leibler

divergence, it can be shown that

$$L(\mathbf{F}, \Psi) \geq - \sum_{t=1}^T \mathbb{E}_{q(\mathbf{h}_t|\theta_t)} [\log q(\mathbf{h}_t|\theta_t)] + \sum_{t=1}^T \mathbb{E}_{q(\mathbf{h}_t|\theta_t)} [\log p(\mathbf{h}_t, \theta_t|\mathbf{F}, \Psi)], \quad (4.27)$$

where the first and second terms on the right-hand side are called the *entropy* and *energy*, respectively, and  $q(\mathbf{h}_t|\theta_t)$ ,  $t = 1, \dots, T$ , are known as the *variational* distributions [1]. The EM algorithm optimises this lower bound on the log-likelihood with respect to  $\mathbf{F}$ ,  $\Psi$  and also  $q(\mathbf{h}_t|\theta_t)$ , hence the name “variational distributions”. The idea is that, by pushing up the lower bound, the log-likelihood  $L(\mathbf{F}, \Psi)$  will hopefully increase as well. In fact, it is guaranteed that each iteration of EM does not decrease  $L(\mathbf{F}, \Psi)$ , which again follows from the properties of the Kullback-Leibler divergence [1].

### 4.2.1 E-step

In the batch E-step,  $\mathbf{F}$  and  $\Psi$  are fixed and Equation (4.27) is maximised with respect to  $q(\mathbf{h}_t|\theta_t)$ ,  $t = 1, \dots, T$ . From [1], the optimal variational distributions are

$$q(\mathbf{h}_t|\theta_t) = p(\mathbf{h}_t|\theta_t, \mathbf{F}, \Psi) \propto \mathcal{N}(\mathbf{m}_t, \Sigma), \quad (4.28)$$

where  $\mathbf{m}_t$  and  $\Sigma$  are given in Equation (4.23). Note that this optimisation can be performed separately for each  $\theta_t$  as it is sampled, using the estimates of  $\mathbf{F}$  and  $\Psi$  on iteration  $t$ . However, in batch EM all  $q(\mathbf{h}_t|\theta_t)$  are updated on every iteration. This is clearly not possible in an online algorithm which discards  $\theta_t$  before sampling  $\theta_{t+1}$ . Therefore, as a compromise, in the online version each  $q(\mathbf{h}_t|\theta_t)$  will be computed once only, on iteration  $t$ , and held fixed thereafter. The only other detail is that the batch algorithm uses  $\theta_{\text{SWA}}$  to compute  $\mathbf{m}_t$ . As  $\theta_{\text{SWA}}$  is not available during training, it will be replaced by the running average of the neural network’s parameter vectors, as in the online SGA algorithm from Section 4.1.

### 4.2.2 M-step

In the batch M-step,  $q(\mathbf{h}_t|\theta_t)$ ,  $t = 1, \dots, T$ , are fixed and Equation (4.27) is maximised with respect to  $\mathbf{F}$  and  $\Psi$ . From [1], the optimal values are

$$\mathbf{F} = \mathbf{A}\mathbf{H}^{-1}, \quad (4.29)$$

where

$$\mathbf{A} = \frac{1}{T} \sum_{t=1}^T \mathbf{d}_t \mathbf{m}_t^\top \quad \text{and} \quad \mathbf{H} = \Sigma + \frac{1}{T} \sum_{t=1}^T \mathbf{m}_t \mathbf{m}_t^\top, \quad (4.30)$$

and

$$\Psi = \text{diag} \left( \text{diag} \left( \frac{1}{T} \sum_{i=1}^T \mathbf{d}_i \mathbf{d}_i^\top - 2\mathbf{F}\mathbf{A}^\top + \mathbf{F}\mathbf{H}\mathbf{F}^\top \right) \right). \quad (4.31)$$

Note that this optimisation involves summing over  $t = 1, \dots, T$ . Moreover, on each iteration all components of the sums in Equation (4.30) and Equation (4.31) are updated. As in the E-step, updating all components is not possible in an online algorithm. Therefore, in the online version these sums will be updated incrementally on each iteration.

### 4.2.3 Practical implementation

Let  $\bar{\mathbf{A}}_t$  and  $\bar{\mathbf{H}}_t$  be the estimates of  $\mathbf{A}$  and  $\mathbf{H}$  from Equation (4.30), respectively, after iteration  $t$  of online EM. That is,

$$\bar{\mathbf{A}}_t = \frac{1}{t} \sum_{i=1}^t \mathbf{d}_i \mathbf{m}_i^\top \quad \text{and} \quad \bar{\mathbf{H}}_t = \Sigma + \frac{1}{t} \sum_{i=1}^t \mathbf{m}_i \mathbf{m}_i^\top. \quad (4.32)$$

Then the estimates of  $\mathbf{F}$  and  $\Psi$  on iteration  $t$  become

$$\mathbf{F} = \bar{\mathbf{A}}_t \bar{\mathbf{H}}_t^{-1} \quad (4.33)$$

and

$$\Psi = \text{diag} \left( \text{diag} \left( \frac{1}{t} \sum_{i=1}^t \mathbf{d}_i \mathbf{d}_i^\top - 2\mathbf{F}\bar{\mathbf{A}}_t^\top + \mathbf{F}\bar{\mathbf{H}}_t \mathbf{F}^\top \right) \right). \quad (4.34)$$

As in Algorithm 1, it suffices to work with  $\psi = \text{diag}(\Psi)$  instead of  $\Psi$ . The diagonal entries of Equation (4.34) can be computed more efficiently as

$$\begin{aligned} \psi &= \frac{1}{t} \sum_{i=1}^t \mathbf{d}_i^2 - 2 \cdot \text{sum}(\mathbf{F} \odot \bar{\mathbf{A}}_t, \text{dim} = 1) + \text{sum}((\mathbf{F}\bar{\mathbf{H}}_t) \odot \mathbf{F}, \text{dim} = 1) \\ &= \frac{1}{t} \sum_{i=1}^t \mathbf{d}_i^2 + \text{sum}((\mathbf{F}\bar{\mathbf{H}}_t) \odot \mathbf{F} - 2\mathbf{F} \odot \bar{\mathbf{A}}_t, \text{dim} = 1). \end{aligned} \quad (4.35)$$

All that remains to complete the online algorithm is to define the incremental update rules for  $\bar{\mathbf{A}}_t$  and  $\bar{\mathbf{H}}_t$ . Both are similar, and for  $\bar{\mathbf{A}}_t$  the derivation is

$$\begin{aligned} \bar{\mathbf{A}}_t &= \frac{1}{t} \sum_{i=1}^t \mathbf{d}_i \mathbf{m}_i^\top \\ &= \frac{1}{t} \left( \frac{t-1}{t-1} \sum_{i=1}^{t-1} \mathbf{d}_i \mathbf{m}_i^\top + \mathbf{d}_t \mathbf{m}_t^\top \right) \\ &= \frac{1}{t} \left( t \cdot \frac{1}{t-1} \sum_{i=1}^{t-1} \mathbf{d}_i \mathbf{m}_i^\top - \frac{1}{t-1} \sum_{i=1}^{t-1} \mathbf{d}_i \mathbf{m}_i^\top + \mathbf{d}_t \mathbf{m}_t^\top \right) \\ &= \frac{1}{t} \left( t\bar{\mathbf{A}}_{t-1} - \bar{\mathbf{A}}_{t-1} + \mathbf{d}_t \mathbf{m}_t^\top \right) \\ &= \bar{\mathbf{A}}_{t-1} + \frac{1}{t} (\mathbf{d}_t \mathbf{m}_t^\top - \bar{\mathbf{A}}_{t-1}). \end{aligned} \quad (4.36)$$

Pseudo code for the practical implementation is given in Algorithm 2.

---

**Algorithm 2** Online Expectation-Maximisation for Factor Analysis
 

---

**Input:** Observation dimension  $D$ , latent dimension  $K$

- 1: Initialise  $\mathbf{F} \in \mathbb{R}^{D \times K}$ ,  $\Psi = \mathbf{1}^D$
  - 2: Initialise  $\bar{\theta}_0 = \mathbf{0}^D$ ,  $\bar{\mathbf{A}}_0 = \mathbf{0}^{D \times K}$ ,  $\bar{\mathbf{B}}_0 = \mathbf{0}^{K \times K}$ ,  $\bar{\mathbf{d}}^2_0 = \mathbf{0}^D$
  - 3: **for**  $t = 1, \dots, T$  **do**
  - 4:   Sample observation  $\theta_t \in \mathbb{R}^D$
  - 5:    $\bar{\theta}_t \leftarrow \bar{\theta}_{t-1} + \frac{1}{t}(\theta_t - \bar{\theta}_{t-1})$
  - 6:    $\mathbf{d}_t \leftarrow \theta_t - \bar{\theta}_t$
  - 7:    $\mathbf{C} \leftarrow (\mathbf{F} \odot \Psi^{-1})^\top$  (with broadcasting)
  - 8:    $\Sigma \leftarrow (\mathbf{I} + \mathbf{C}\mathbf{F})^{-1}$
  - 9:    $\mathbf{m}_t \leftarrow \Sigma \mathbf{C} \mathbf{d}_t$
  - 10:    $\bar{\mathbf{B}}_t \leftarrow \bar{\mathbf{B}}_{t-1} + \frac{1}{t}(\mathbf{m}_t \mathbf{m}_t^\top - \bar{\mathbf{B}}_{t-1})$
  - 11:    $\bar{\mathbf{H}}_t \leftarrow \Sigma + \bar{\mathbf{B}}_t$
  - 12:    $\bar{\mathbf{A}}_t \leftarrow \bar{\mathbf{A}}_{t-1} + \frac{1}{t}(\mathbf{d}_t \mathbf{m}_t^\top - \bar{\mathbf{A}}_{t-1})$
  - 13:    $\mathbf{F} \leftarrow \bar{\mathbf{A}}_t \bar{\mathbf{H}}_t^{-1}$
  - 14:    $\bar{\mathbf{d}}^2_t \leftarrow \bar{\mathbf{d}}^2_{t-1} + \frac{1}{t}(\mathbf{d}_t^2 - \bar{\mathbf{d}}^2_{t-1})$
  - 15:    $\Psi \leftarrow \bar{\mathbf{d}}^2_t + \text{sum}((\mathbf{F} \bar{\mathbf{H}}_t) \odot \mathbf{F} - 2\mathbf{F} \odot \bar{\mathbf{A}}_t, \text{dim} = 1)$
  - 16: **end for**
  - 17:  $\theta_{\text{SWA}} \leftarrow \bar{\theta}_T$
  - 18: **return**  $\theta_{\text{SWA}}, \mathbf{F}, \Psi$
- 

### 4.3 Variational Inference

Algorithms 1 and 2 both assume that samples from the true posterior can be obtained via some separate procedure. This section describes an alternative method based on VI which learns a FA posterior directly.

In Section 2.3, a general VI algorithm was outlined for learning an approximate Gaussian posterior of the parameter vector of a neural network. Suppose that the approximate posterior is a FA model. That is,  $q(\theta) = \mathcal{N}(\mathbf{c}, \mathbf{F}\mathbf{F}^\top + \Psi)$  for some  $\mathbf{c} \in \mathbb{R}^D$ ,  $\mathbf{F} \in \mathbb{R}^{D \times K}$  and  $D \times D$  diagonal matrix  $\Psi$ . Then the gradient of the VI objective in Equation (2.16) becomes

$$\nabla_{\mathbf{c}, \mathbf{F}, \Psi} \mathbb{E}_{q(\theta)} [\log q(\theta)] - \nabla_{\mathbf{c}, \mathbf{F}, \Psi} \mathbb{E}_{q(\theta)} [\log p(\theta)] - \nabla_{\mathbf{c}, \mathbf{F}, \Psi} \mathbb{E}_{q(\theta)} [\log p(\mathcal{D}|\theta)], \quad (4.37)$$

where the expectation is taken over samples  $\theta \sim \mathcal{N}(\mathbf{c}, \mathbf{F}\mathbf{F}^\top + \Psi)$ . These derivatives are derived in the following sections.

### 4.3.1 Derivatives of the Variational Distribution

The logarithm of the pdf of the FA posterior is

$$\log q(\theta) = -\frac{1}{2}(\theta - \mathbf{c})^\top (\mathbf{F}\mathbf{F}^\top + \Psi)^{-1}(\theta - \mathbf{c}) - \frac{1}{2} \log |\mathbf{F}\mathbf{F}^\top + \Psi| - \frac{D}{2} \log 2\pi. \quad (4.38)$$

It is given in [16] that, for any symmetric matrix  $\mathbf{A}$  and stochastic vector  $\mathbf{x}$  with mean  $\mathbf{m}$  and covariance  $\mathbf{M}$ ,

$$\mathbb{E}[\mathbf{x}^\top \mathbf{A} \mathbf{x}] = \text{Tr}(\mathbf{A}\mathbf{M}) + \mathbf{m}^\top \mathbf{A} \mathbf{m}. \quad (4.39)$$

Hence, setting  $\mathbf{S} = \mathbf{F}\mathbf{F}^\top + \Psi$ , it follows that

$$\begin{aligned} \mathbb{E}_{q(\theta)}[(\theta - \mathbf{c})^\top \mathbf{S}^{-1}(\theta - \mathbf{c})] &= \text{Tr}(\mathbf{S}^{-1} \mathbb{V}_{q(\theta)}[(\theta - \mathbf{c})]) + \mathbb{E}_{q(\theta)}[(\theta - \mathbf{c})]^\top \mathbf{S}^{-1} \mathbb{E}_{q(\theta)}[(\theta - \mathbf{c})] \\ &= \text{Tr}(\mathbf{S}^{-1}(\mathbf{F}\mathbf{F}^\top + \Psi)) + \mathbf{0} \mathbf{S}^{-1} \mathbf{0} \\ &= \text{Tr}(\mathbf{I}) \\ &= D. \end{aligned} \quad (4.40)$$

Therefore,

$$\mathbb{E}_{q(\theta)}[\log q(\theta)] = -\frac{D}{2} - \frac{1}{2} \log |\mathbf{S}| - \frac{D}{2} \log 2\pi. \quad (4.41)$$

Then

$$\nabla_{\mathbf{c}} \mathbb{E}_{q(\theta)}[\log q(\theta)] = \mathbf{0} \quad (4.42)$$

and, using the identity  $\nabla_{\mathbf{X}} \log |\mathbf{X}| = \mathbf{X}^{-\top}$  from [16] and the fact that  $\mathbf{S}$  is symmetric,

$$\nabla_{\mathbf{S}} \mathbb{E}_{q(\theta)}[\log q(\theta)] = -\frac{1}{2} \mathbf{S}^{-1}. \quad (4.43)$$

Hence, it follows from the chain rule of calculus that

$$\begin{aligned} \nabla_{\mathbf{F}} \mathbb{E}_{q(\theta)}[\log q(\theta)] &= -\frac{1}{2} \mathbf{S}^{-1} (\nabla_{\mathbf{F}} \mathbf{S}) \\ &= -\frac{1}{2} \mathbf{S}^{-1} (2\mathbf{F}) \\ &= -\mathbf{S}^{-1} \mathbf{F} \end{aligned} \quad (4.44)$$

and

$$\begin{aligned} \nabla_{\Psi} \mathbb{E}_{q(\theta)}[\log q(\theta)] &= \text{diag} \left( \text{diag} \left( -\frac{1}{2} \mathbf{S}^{-1} (\nabla_{\Psi} \mathbf{S}) \right) \right) \\ &= -\frac{1}{2} \text{diag} \left( \text{diag} (\mathbf{S}^{-1} \mathbf{I}) \right) \\ &= -\frac{1}{2} \text{diag} (\text{diag} (\mathbf{S}^{-1})). \end{aligned} \quad (4.45)$$

### 4.3.2 Derivatives of the Prior

Let the prior be  $p(\theta) = \mathcal{N}(\mathbf{0}, \alpha^{-1}\mathbf{I})$  for some precision  $\alpha > 0$ . Then

$$\begin{aligned}\mathbb{E}_{q(\theta)}[\log p(\theta)] &= \mathbb{E}_{q(\theta)}\left[-\frac{1}{2}\theta^\top \alpha \mathbf{I} \theta - \frac{1}{2} \log |\alpha^{-1}\mathbf{I}| - \frac{D}{2} \log 2\pi\right] \\ &= -\frac{\alpha}{2} \mathbb{E}_{q(\theta)}[\theta^\top \mathbf{I} \theta] + \text{constant} \\ &= -\frac{\alpha}{2} (\text{Tr}(\mathbf{F}\mathbf{F}^\top + \Psi) + \mathbf{c}^\top \mathbf{c}) + \text{constant},\end{aligned}\tag{4.46}$$

where the last line follows from Equation (4.39). Hence,

$$\nabla_{\mathbf{c}} \mathbb{E}_{q(\theta)}[\log p(\theta)] = -\alpha \mathbf{c},\tag{4.47}$$

$$\nabla_{\mathbf{F}} \mathbb{E}_{q(\theta)}[\log p(\theta)] = -\alpha \mathbf{F},\tag{4.48}$$

$$\nabla_{\Psi} \mathbb{E}_{q(\theta)}[\log p(\theta)] = -\frac{\alpha}{2} \mathbf{I},\tag{4.49}$$

which follow from the matrix calculus identities in [16].

### 4.3.3 Derivatives of the Likelihood

Let the log-likelihood,  $\log p(\mathcal{D}|\theta)$ , be parameterised by a neural network with weights  $\theta$ . Then the back-propagation algorithm [17] can compute  $\nabla_{\theta} \log p(\mathcal{D}|\theta)$ . From here, the re-parameterisation trick [3] can be used to obtain the partial derivatives with respect to  $\mathbf{c}$ ,  $\mathbf{F}$  and  $\Psi$ . First note that, by re-writing Equation (2.1), sampling  $\theta \sim q(\theta) = \mathcal{N}(\mathbf{c}, \mathbf{F}\mathbf{F}^\top + \Psi)$  is equivalent to setting

$$\theta = \mathbf{F}\mathbf{h} + \mathbf{c} + \Psi^{1/2}\mathbf{z}\tag{4.50}$$

from some  $\mathbf{h} \sim \mathcal{N}(\mathbf{0}^K, \mathbf{I}^{K \times K})$  and  $\mathbf{z} \sim \mathcal{N}(\mathbf{0}^D, \mathbf{I}^{D \times D})$ , where the exponent is applied to  $\Psi$  element-wise. Given  $\mathbf{h}$  and  $\mathbf{z}$ , this is a deterministic function of  $\mathbf{c}$ ,  $\mathbf{F}$  and  $\Psi$ . It then follows from the chain rule of calculus that

$$\nabla_{\mathbf{c}} \log p(\mathcal{D}|\theta) = \nabla_{\theta} \log p(\mathcal{D}|\theta),\tag{4.51}$$

$$\nabla_{\mathbf{F}} \log p(\mathcal{D}|\theta) = \nabla_{\theta} \log p(\mathcal{D}|\theta) \mathbf{h}^\top,\tag{4.52}$$

$$\nabla_{\Psi} \log p(\mathcal{D}|\theta) = \text{diag} \left( \text{diag} \left( \frac{1}{2} \nabla_{\theta} \log p(\mathcal{D}|\theta) \mathbf{z}^\top \Psi^{-1/2} \right) \right).\tag{4.53}$$

Due to the functional form of a non-linear neural network, the expectation of these derivatives over  $q(\theta)$  cannot be computed in closed-form and must be approximate by Monte Carlo estimates using samples  $\theta \sim q(\theta)$ .

### 4.3.4 Practical Implementation

In order to compute the gradients in Equations (4.44) and (4.45), the inverse of  $\mathbf{S} = \mathbf{F}\mathbf{F}^\top + \Psi$  is required. Using the Woodbury identity [16],

$$\begin{aligned} (\mathbf{F}\mathbf{F}^\top + \Psi)^{-1} &= \Psi^{-1} - \Psi^{-1}\mathbf{F}(\mathbf{I} + \mathbf{F}^\top\Psi^{-1}\mathbf{F})^{-1}\mathbf{F}^\top\Psi^{-1} \\ &= \Psi^{-1} - (\Psi^{-1} \odot \mathbf{F})(\mathbf{I} + \mathbf{F}^\top(\Psi^{-1} \odot \mathbf{F}))^{-1}(\Psi^{-1} \odot \mathbf{F})^\top \\ &= \Psi^{-1} - \mathbf{A}(\mathbf{I} + \mathbf{F}^\top\mathbf{A})^{-1}\mathbf{A}^\top, \end{aligned} \quad (4.54)$$

where  $\Psi = \text{diag}(\Psi)$  and  $\mathbf{A} = \Psi^{-1} \odot \mathbf{F}$  (with broadcasting). It follows that

$$\begin{aligned} \nabla_{\mathbf{F}} \mathbb{E}_{q(\theta)}[\log q(\theta)] &= -(\mathbf{F}\mathbf{F}^\top + \Psi)^{-1}\mathbf{F} \\ &= -\Psi^{-1}\mathbf{F} + \mathbf{A}(\mathbf{I} + \mathbf{F}^\top\mathbf{A})^{-1}(\mathbf{A}^\top\mathbf{F}) \\ &= -\mathbf{A} + \mathbf{A}(\mathbf{I} + \mathbf{F}^\top\mathbf{A})^{-1}(\mathbf{A}^\top\mathbf{F}) \\ &= -\mathbf{A} + \mathbf{A}(\mathbf{I} + \mathbf{B})^{-1}\mathbf{B}^\top \\ &= -\mathbf{A} + \mathbf{C}\mathbf{B}^\top, \end{aligned} \quad (4.55)$$

where  $\mathbf{B} = \mathbf{F}^\top\mathbf{A}$  and  $\mathbf{C} = \mathbf{A}(\mathbf{I} + \mathbf{B})^{-1}$ . Also,

$$\begin{aligned} \nabla_{\Psi} \mathbb{E}_{q(\theta)}[\log q(\theta)] &= -\frac{1}{2} \text{diag}((\mathbf{F}\mathbf{F}^\top + \Psi)^{-1}) \\ &= -\frac{1}{2} \text{diag}(\Psi^{-1} - \mathbf{A}(\mathbf{I} + \mathbf{F}^\top\mathbf{A})^{-1}\mathbf{A}^\top) \\ &= -\frac{1}{2}\Psi^{-1} + \frac{1}{2} \text{diag}(\mathbf{C}\mathbf{A}^\top) \\ &= -\frac{1}{2}\Psi^{-1} + \frac{1}{2} \text{sum}(\mathbf{C} \odot \mathbf{A}, \text{dim} = 1). \end{aligned} \quad (4.56)$$

Note that these calculations involve no explicit use of  $D \times D$  matrices. All other references to  $D \times D$  matrices can be removed by writing the partial derivatives in Equations (4.49) and (4.53) with respect to  $\Psi$ , as

$$\nabla_{\Psi} \mathbb{E}_{q(\theta)}[\log p(\theta)] = -\frac{\alpha}{2} \mathbf{1}^D, \quad (4.57)$$

$$\nabla_{\Psi} \log p(\mathcal{D}|\theta) = \frac{1}{2} \nabla_{\theta} \log p(\mathcal{D}|\theta) \odot (\Psi^{-1/2} \odot \mathbf{z}). \quad (4.58)$$

Finally, as in Algorithm 1, the re-parameterisation  $\Psi = \exp \beta$  can again be used to ensure that the variances remain positive. Since  $\nabla_{\beta} \Psi = \Psi$ , it follows from the chain rule that the corresponding partial derivatives with respect to  $\beta$  are

$$\nabla_{\beta} \mathbb{E}_{q(\theta)}[\log q(\theta)] = -\frac{1}{2} + \frac{1}{2} \text{sum}(\mathbf{C} \odot \mathbf{A}, \text{dim} = 1) \odot \Psi, \quad (4.59)$$

$$\nabla_{\beta} \mathbb{E}_{q(\theta)}[\log p(\theta)] = -\frac{\alpha}{2} \Psi, \quad (4.60)$$



$$\nabla_{\beta} \log p(\mathcal{D}|\theta) = \frac{1}{2} \nabla_{\theta} \log p(\mathcal{D}|\theta) \odot (\psi^{1/2} \odot \mathbf{z}). \quad (4.61)$$

Pseudo code for the practical implementation is given in Algorithm 3. Being a VI method for approximating a posterior with a FA model, the algorithm is named VIFA.

### 4.3.5 Comparison of VIFA and SLANG

There are some notable similarities between VIFA and SLANG. Both learn an approximate posterior of a neural network by minimising the VI objective in Equation (2.16). Moreover, both use a low-rank plus diagonal approximation of a  $D \times D$  matrix to make the computation tractable - VIFA approximates the posterior covariance, whereas SLANG approximates its inverse. Both are online gradient-based algorithms. However, while VIFA uses vanilla gradients, SLANG uses natural gradients. The authors of SLANG claim that natural gradient methods are preferable when optimising the parameters of a distribution [13], although the use of an EF matrix to approximate the Fisher information matrix in natural gradient algorithms has been called into question [10].

Here it is claimed that VIFA is a much more general and practical algorithm. It can be seen from Algorithm 3 that VIFA is implemented using only simple matrix operations and the standard back-propagation algorithm. In contrast, sampling each neural network parameter vector in SLANG requires computing a Cholesky decomposition, and each update to the approximate posterior involves an eigen-decomposition. Moreover, this eigen-decomposition requires the back-propagated gradients of individual training examples. As pointed out in Section 3.1, this is problematic, since a separate implementation of back-propagation is needed for each different type of neural network layer. Perhaps the biggest advantage of VIFA is that it only requires the gradient with respect to the mean negative log-likelihood of a mini-batch of training examples. This is easy to obtain using standard automatic differentiation in deep learning frameworks such as PyTorch [14]. Hence, VIFA is a general algorithm which can be applied to any neural network architecture with no extra effort.

## 4.4 Factors Initialisation

One detail missing from Algorithms 1, 2 and 3 is how  $\mathbf{F}$  is initialised. To encourage diverse factors, it makes sense to initialise  $\mathbf{F}$  with orthogonal columns. In practice, this

**Algorithm 3** VIFA

---

**Input:** Dataset  $\mathcal{D} = \{(\mathbf{x}_n, y_n)\}_{n=1}^N$ , observation dimension  $D$ , latent dimension  $K$ , precision  $\alpha > 0$ , mini-batch size  $M$ , Monte Carlo average size  $L$ , learning rates  $\eta_{\mathbf{c}}, \eta_{\mathbf{F}}, \eta_{\beta} > 0$

- 1: Initialise  $\mathbf{c} = \mathbf{0}^D$ ,  $\mathbf{F} \in \mathbb{R}^{D \times K}$ ,  $\psi = \mathbf{1}^D$
- 2:  $\beta \leftarrow \log \psi$
- 3: Initialise  $\hat{\mathbf{g}}_{\mathbf{c}} = \mathbf{0}^D$ ,  $\hat{\mathbf{G}}_{\mathbf{F}} = \mathbf{0}^{D \times K}$ ,  $\hat{\mathbf{g}}_{\beta} = \mathbf{0}^D$
- 4:  $\mathbf{A} \leftarrow \psi^{-1} \odot \mathbf{F}$  (with broadcasting)
- 5:  $\mathbf{B} \leftarrow \mathbf{F}^T \mathbf{A}$
- 6:  $\mathbf{C} \leftarrow \mathbf{A}(\mathbf{I} + \mathbf{B})^{-1}$
- 7:  $\mathbf{D} \leftarrow -\mathbf{A} + \mathbf{C}\mathbf{B}^T$
- 8:  $\mathbf{s} \leftarrow -\frac{1}{2} + \frac{1}{2} \text{sum}(\mathbf{C} \odot \mathbf{A}, \text{dim} = 1) \odot \psi$
- 9: **for**  $t = 1, \dots, T$  **do**
- 10:     Sample a mini-batch  $\mathcal{B}$  of  $M$  training examples from  $\mathcal{D}$
- 11:     Sample  $\mathbf{h} \sim \mathcal{N}(\mathbf{0}^K, \mathbf{I}^{K \times K})$
- 12:     Sample  $\mathbf{z} \sim \mathcal{N}(\mathbf{0}^D, \mathbf{I}^{D \times D})$
- 13:      $\boldsymbol{\theta} \leftarrow \mathbf{F}\mathbf{h} + \mathbf{c} + \psi^{1/2} \odot \mathbf{z}$
- 14:      $\hat{\mathbf{g}}_{\boldsymbol{\theta}} \leftarrow N \cdot \nabla_{\boldsymbol{\theta}} \left( -\frac{1}{M} \sum_{(\mathbf{x}, y) \in \mathcal{B}} \log p(y|\mathbf{x}, \boldsymbol{\theta}) \right)$  (via back-propagation)
- 15:      $\hat{\mathbf{g}}_{\mathbf{c}} \leftarrow \hat{\mathbf{g}}_{\mathbf{c}} + \alpha \mathbf{c} + \hat{\mathbf{g}}_{\boldsymbol{\theta}}$
- 16:      $\hat{\mathbf{G}}_{\mathbf{F}} \leftarrow \hat{\mathbf{G}}_{\mathbf{F}} + \mathbf{D} + \alpha \mathbf{F} + \hat{\mathbf{g}}_{\boldsymbol{\theta}} \mathbf{h}^T$
- 17:      $\hat{\mathbf{g}}_{\beta} \leftarrow \hat{\mathbf{g}}_{\beta} + \mathbf{s} + \frac{\alpha}{2} \psi + \frac{1}{2} \hat{\mathbf{g}}_{\boldsymbol{\theta}} \odot (\psi^{1/2} \odot \mathbf{z})$
- 18:     **if**  $t \bmod L = 0$  **then**
- 19:          $\mathbf{c} \leftarrow \mathbf{c} - \eta_{\mathbf{c}} \left( \frac{1}{L} \hat{\mathbf{g}}_{\mathbf{c}} \right)$
- 20:          $\mathbf{F} \leftarrow \mathbf{F} - \eta_{\mathbf{F}} \left( \frac{1}{L} \hat{\mathbf{G}}_{\mathbf{F}} \right)$
- 21:          $\beta \leftarrow \beta - \eta_{\beta} \left( \frac{1}{L} \hat{\mathbf{g}}_{\beta} \right)$
- 22:          $\psi \leftarrow \exp \beta$
- 23:         Reset  $\hat{\mathbf{g}}_{\mathbf{c}}, \hat{\mathbf{G}}_{\mathbf{F}}$  and  $\hat{\mathbf{g}}_{\beta}$  according to line 3
- 24:         Recalculate  $\mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{D}$  and  $\mathbf{s}$  according to lines 4-8
- 25:     **end if**
- 26: **end for**
- 27: **return**  $\mathbf{c}, \mathbf{F}, \psi$

---

can be achieved by first generating a matrix  $\mathbf{A} \in \mathbb{R}^{D \times K}$ , whose entries are independently sampled from a standard normal distribution, and then computing its reduced QR decomposition<sup>1</sup>. This decomposition is  $\mathbf{A} = \mathbf{Q}\mathbf{R}$ , where  $\mathbf{Q} \in \mathbb{R}^{D \times K}$  is orthogonal and  $\mathbf{R} \in \mathbb{R}^{K \times K}$  is upper triangular. Then set  $\mathbf{F} = \mathbf{Q}$ .

---

<sup>1</sup><https://pytorch.org/docs/1.9.0/generated/torch.linalg.qr.html>

# Chapter 5

## Online Factor Analysis Experiments

### 5.1 Methodology

The aim of these experiments is to test how well the online FA algorithms from Chapter 4 are able to fit observations sampled from actual FA models. The form of a FA model is given in Equation (2.3). In particular, it has parameters  $\mathbf{c}$ ,  $\mathbf{F}$  and  $\Psi$ . The maximum likelihood estimate of  $\mathbf{c}$ , given the data, is just the empirical mean of the sampled observations [1]. This is computed exactly in both online FA algorithms (Algorithms 1 and 2) by maintaining a running average of the observations. The other parameters,  $\mathbf{F}$  and  $\Psi$ , appear in the FA model as part of the full covariance matrix,  $\mathbf{F}\mathbf{F}^\top + \Psi$ . Note that right-multiplying  $\mathbf{F}$  by any orthogonal matrix  $\mathbf{A} \in \mathbb{R}^{K \times K}$  will result in the exact same covariance [1], since

$$(\mathbf{F}\mathbf{A})(\mathbf{F}\mathbf{A})^\top + \Psi = \mathbf{F}\mathbf{A}\mathbf{A}^\top\mathbf{F}^\top + \Psi = \mathbf{F}\mathbf{F}^\top + \Psi. \quad (5.1)$$

Therefore, the  $\mathbf{F}$  estimated by an online FA algorithm cannot be directly compared to the true  $\mathbf{F}$ . The comparison must be made between the full covariance matrices.

Since the covariance matrices have shape  $D \times D$ , memory requirements dictate that  $D$ , the observation dimension, cannot be too large. Therefore, values of  $D = 100$  and  $D = 1000$  were used in all experiments in this section. In all cases the latent dimension was set to  $K = 10$ , meaning that two different observation to latent ratios were tested:  $\frac{D}{K} = 10$  and  $\frac{D}{K} = 100$ . Since computing the maximum likelihood estimate of  $\mathbf{c} \in \mathbb{R}^D$  is trivial, in all experiments the entries of  $\mathbf{c}$  were simply sampled independently from a standard normal distribution. Each factor loading matrix  $\mathbf{F}$  was generated in such a way that its columns spanned the  $K$ -dimensional latent space and the conditioning number of the resultant covariance matrix could be controlled. This was achieved by finding

the first  $K$  eigenvectors of a symmetric positive semi-definite matrix  $\mathbf{M} \in \mathbb{R}^{D \times D}$ , and then setting the columns of  $\mathbf{F}$  equal to these eigenvectors multiplied by a vector  $\mathbf{s} > \mathbf{0} \in \mathbb{R}^D$  (element-wise). By construction, the  $K$  eigenvectors are linearly independent and therefore span the latent space, and scaling them by  $\mathbf{s}$  affects the conditioning number of  $\mathbf{F}\mathbf{F}^\top$ . The conditioning number of a matrix is equal to

$$\max_{i,j} \frac{|\lambda_i|}{|\lambda_j|}, \quad (5.2)$$

where the  $\lambda_i$  and  $\lambda_j$  are eigenvalues of the matrix [3]. This ratio can be controlled by specifying the range of  $\mathbf{s}$ , or alternatively, the range of  $\mathbf{s}^2$ , which is called the *spectrum*. The larger the ratio of the upper to lower bound of the spectrum, the larger the ratio of the eigenvalues. In each experiment the spectrum was sampled from a uniform distribution with one of the following ranges:  $[1, 10]$ ,  $[1, 100]$  or  $[1, 1000]$ . Finally, the diagonal entries of  $\Psi$  were sampled from a uniform distribution with upper bound equal to the maximum value of  $\mathbf{s}^2$ . This is consistent with the FA assumption that an observation,  $\boldsymbol{\theta} = \mathbf{F}\mathbf{h} + \mathbf{c} + \boldsymbol{\varepsilon}$ , is generated by corrupting the signal  $\mathbf{F}\mathbf{h} + \mathbf{c}$  with some random noise  $\boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \Psi)$ . The full details of how FA models were generated are given in Algorithm 4.

---

**Algorithm 4** Generate a Factor Analysis Model

---

**Input:** Observation dimension  $D$ , latent dimension  $K$ , spectrum range  $[a, b]$

- 1: Generate  $\mathbf{c} \in \mathbb{R}^D$  by sampling entries independently from  $\mathcal{N}(0, 1)$
  - 2: Generate  $\mathbf{A} \in \mathbb{R}^{D \times D}$  by sampling entries independently from  $\mathcal{N}(0, 1)$
  - 3:  $\mathbf{M} \leftarrow \mathbf{A}\mathbf{A}^\top$
  - 4: Compute the  $K$  eigenvectors,  $\mathbf{v}_1, \dots, \mathbf{v}_K \in \mathbb{R}^D$ , corresponding to the  $K$  largest eigenvalues of  $\mathbf{M}$
  - 5: Construct the matrix  $\mathbf{V}_K \in \mathbb{R}^{D \times K}$  with columns  $\mathbf{v}_1, \dots, \mathbf{v}_K$
  - 6: Generate  $\mathbf{s}^2 \in \mathbb{R}^D$  by sampling entries independently from  $\mathcal{U}(a, b)$
  - 7:  $\mathbf{s} \leftarrow \sqrt{\mathbf{s}^2}$  (square root is applied element-wise)
  - 8:  $\mathbf{F} \leftarrow \mathbf{V}_K \odot \mathbf{s}$  (with broadcasting)
  - 9:  $s_{\max} \leftarrow \max(\mathbf{s}^2)$  (largest element of  $\mathbf{s}^2$ )
  - 10: Generate  $\boldsymbol{\psi} \in \mathbb{R}^D$  by sampling entries independently from  $\mathcal{U}(0, s_{\max})$
  - 11:  $\Psi \leftarrow \text{diag}(\boldsymbol{\psi})$
  - 12: **return**  $\mathbf{c}, \mathbf{F}, \Psi$
- 

Having generated a FA model, observations were sampled according to Equation (2.1). Using the data, the parameters of the model were then estimated by three sep-

arate FA algorithms: online SGA (Algorithm 1), online EM (Algorithm 2) and the `scikit-learn` `FactorAnalysis` estimator [15]. The `scikit-learn` implementation is based on the batch SVD algorithm from [1]. The *randomised* version of SVD was used, which for most applications is “sufficiently precise while providing significant speed gains” [15]. Each algorithm was tested on samples of varying size, from 100 up to 100,000 observations. The latent dimension of each approximate model was set to the true latent dimension  $K$ . All other hyperparameters of the `scikit-learn` algorithm were set to their default values. For both online algorithms,  $\mathbf{F}$  was initialised as described in Section 4.4, and online SGA ran with a learning rate of 0.001. Additionally, both online algorithms were allowed a *warm-up* period of 100 iterations during which the running averages,  $\bar{\boldsymbol{\theta}}_t, \bar{\mathbf{A}}_t, \bar{\mathbf{B}}_t$  and  $\bar{\mathbf{d}}^2_t$ , were updated, while  $\mathbf{F}$  and  $\Psi$  remained fixed (TODO: perhaps include in an appendix the derivation that Michael sent, about how this is equivalent to specifying a prior). Each experiment was repeated ten times. In each trial a different random seed was used for generating the true FA model and also initialising the parameters of the online algorithms.

## 5.2 Results and Discussion

Figure 5.1 shows the relative distance between the true covariance matrix of each FA model and the corresponding estimated covariance matrix of each learning algorithm, as a function of the number of samples used to fit the approximate models. The distance between two matrices is measured by the Frobenius norm (also sometimes called the Euclidean norm) of the difference between the two matrices. The distance is then divided by the Frobenius norm of the true covariance matrix to get the relative distance.

In the first row of the figure are the results corresponding to FA models whose factor loading matrices were generated with a spectrum range of  $[1, 10]$ . When  $D = 100$  (top-left), online EM approximates the true covariance matrix better than online SGA when the number of samples,  $T$ , is less than 20,000. For  $T \geq 20,000$ , the results are very close and the standard error bars overlap. Moreover, for  $T \geq 50,000$ , both online algorithms are on par with batch SVD. When  $D = 1000$  (top-right), online EM beats SGA for  $T \geq 200$  and matches batch SVD when  $T = 100,000$ . This is in contrast to SGA, which appears to get stuck after  $T = 20,000$  and does not decrease the distance any further. In the second row of the figure are the results corresponding to a spectrum range of  $[1, 100]$ . For both  $D = 100$  (middle-left) and  $D = 1000$  (middle-right), online EM outperforms SGA for  $T \geq 200$ . Both plots are similar in the sense that online EM

initially decreases the distance quite dramatically, even to a greater extent than batch SVD, then the rate of improvement slows before it again catches up with batch SVD at  $T = 100,000$ . In the third row of the figure are the results corresponding to a spectrum range of  $[1, 1000]$ . In this case, online EM is again superior to SGA. Up until  $T = 5000$  the results are similar to those in the second row of the figure. However, in this case online EM does not learn as efficiently as  $T$  becomes large and in the end it is not able to match batch SVD.

Figure 5.2 shows the 2-Wasserstein distance (TODO: add citation) between the Gaussian distribution defined by each estimated FA model and the Gaussian distribution defined by the true FA model, for the same combinations of observation dimension, latent dimension, spectrum range and sample size. Since all algorithms are able to compute the mean of the true Gaussian distribution exactly, the Wasserstein distance in this case is just a different measure of distance between the true and estimated covariance matrices. Therefore, the shape of the plots in Figure 5.2 is very similar to the shape of the plots in Figure 5.1, albeit the y-axis scales are different.

In summary, these results suggest that online EM is in general better at fitting FA models than online SGA, for all combinations of observation dimension, latent dimension, spectrum range and sample size. The SGA algorithm could possibly be improved by tuning its learning rate, or even using a learning rate scheduler. However, tuning hyperparameters is often a hassle and the fact that this is not required for online EM is another distinct advantage. Compared to batch SVD, the results achieved by online EM are almost identical when  $T = 100,000$ , except when the spectrum range is equal to  $[1, 1000]$ . This case would appear to be more challenging for online EM due to the larger conditioning number of the covariance matrix. It appears to get stuck in a local minimum at around  $T = 1000$  and the rate of improvement only starts to increase again after  $T = 20,000$ . Perhaps online EM would be able to catch batch SVD if more than 100,000 samples were used. However, this would be an unreasonably high number of neural network parameter vectors to sample for the use case in this thesis.

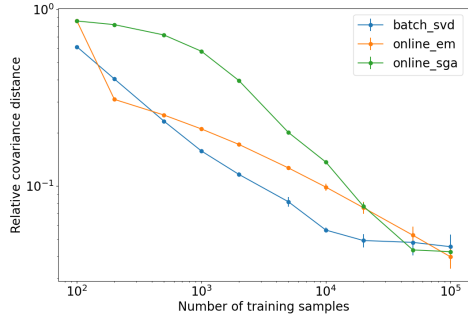
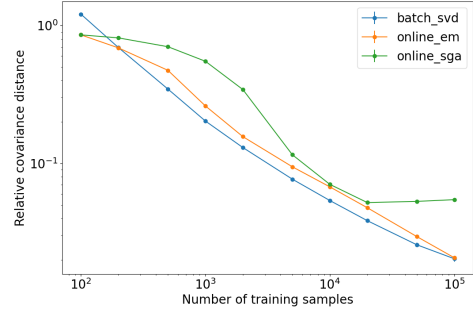
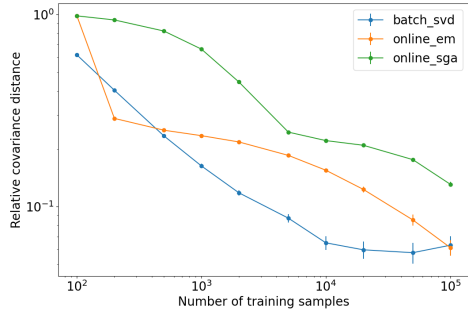
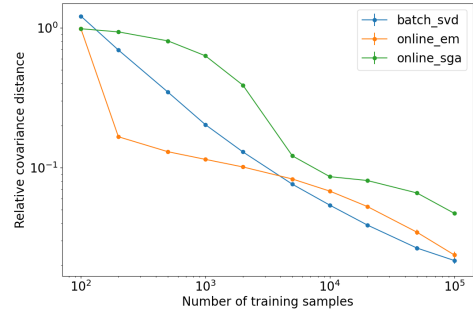
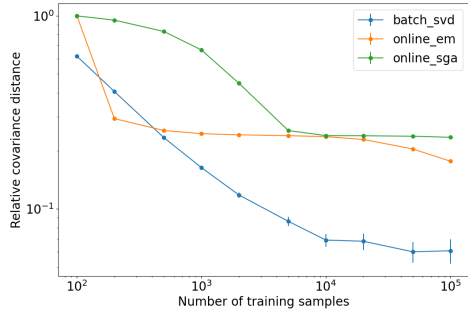
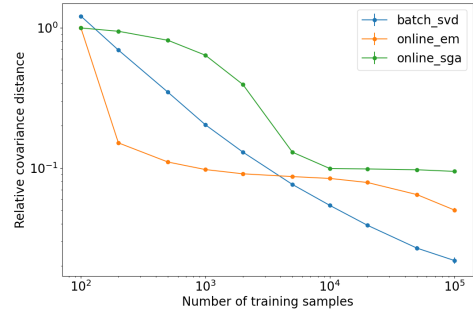
(a)  $D = 100, K = 10, \mathbf{s}^2 \sim \mathcal{U}(1, 10)$ (b)  $D = 1000, K = 10, \mathbf{s}^2 \sim \mathcal{U}(1, 10)$ (c)  $D = 100, K = 10, \mathbf{s}^2 \sim \mathcal{U}(1, 100)$ (d)  $D = 1000, K = 10, \mathbf{s}^2 \sim \mathcal{U}(1, 100)$ (e)  $D = 100, K = 10, \mathbf{s}^2 \sim \mathcal{U}(1, 1000)$ (f)  $D = 1000, K = 10, \mathbf{s}^2 \sim \mathcal{U}(1, 1000)$ 

Figure 5.1: The relative distance of the estimated FA covariance matrices from the true covariance matrix as a function of the number of samples used to learn the models. The blue, orange and green lines show, for batch SVD, online EM and online SGA, respectively, the Frobenius norm of the difference between the true covariance matrix and the estimated covariance matrix divided by the Frobenius norm of the true covariance matrix. Each data point shows the mean value over ten trials with different random seeds, and standard error bars are also plotted. The different plots correspond to different combinations of the observation dimension  $D$ , the latent dimension  $K$  and the range of the spectrum  $\mathbf{s}^2$ .



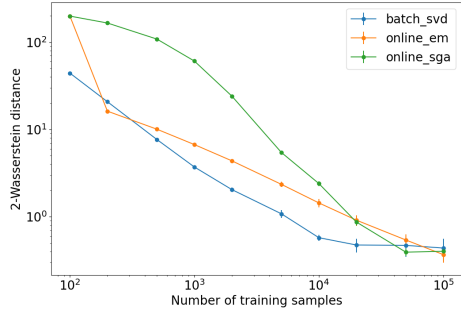
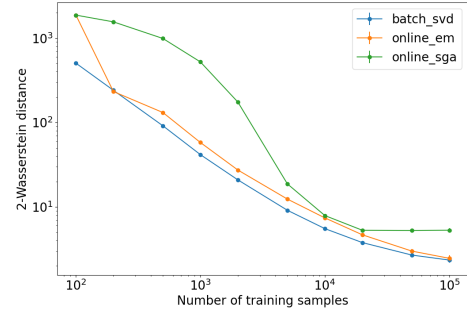
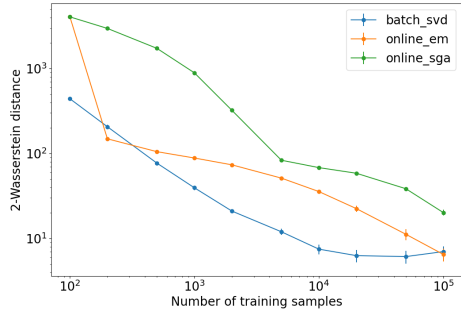
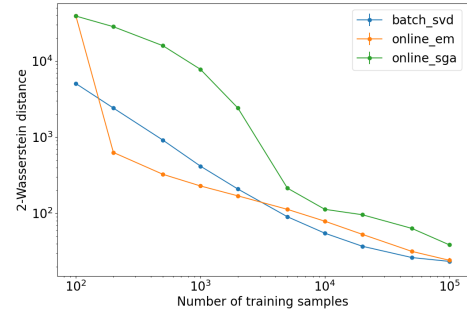
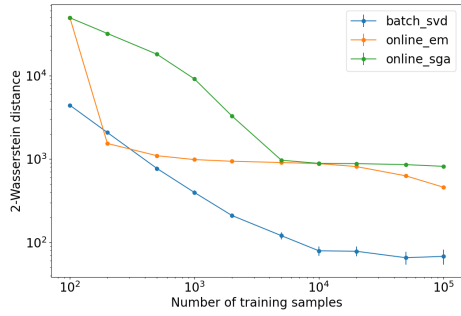
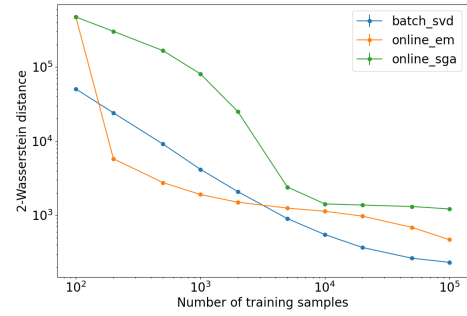
(a)  $D = 100, K = 10, \mathbf{s}^2 \sim \mathcal{U}(1, 10)$ (b)  $D = 1000, K = 10, \mathbf{s}^2 \sim \mathcal{U}(1, 10)$ (c)  $D = 100, K = 10, \mathbf{s}^2 \sim \mathcal{U}(1, 100)$ (d)  $D = 1000, K = 10, \mathbf{s}^2 \sim \mathcal{U}(1, 100)$ (e)  $D = 100, K = 10, \mathbf{s}^2 \sim \mathcal{U}(1, 1000)$ (f)  $D = 1000, K = 10, \mathbf{s}^2 \sim \mathcal{U}(1, 1000)$ 

Figure 5.2: The 2-Wasserstein distance between the Gaussian distribution defined by each estimated FA model and the Gaussian distribution defined by the true FA model. The blue, orange and green lines show the 2-Wasserstein distance corresponding to batch SVD, online EM and online SGA, respectively. Each data point shows the mean value over ten trials with different random seeds, and standard error bars are also plotted. The different plots correspond to different combinations of the observation dimension  $D$ , the latent dimension  $K$  and the range of the spectrum  $\mathbf{s}^2$ .

# Chapter 6

## Linear Regression Experiments

### 6.1 Posterior Estimation

The ability of Algorithms 1 and 2 to learn the posterior of a neural network's parameter vector hinges on obtaining samples of the parameter vector which are representative of the true posterior. In the SWAG approach, such samples are thought to be obtained from the iterates of the parameter vector which are encountered along the SGD trajectory while training the neural network with a high constant learning rate, following an initial pre-training phase [11]. In the case of linear regression - which is a neural network with no hidden layers - this assumption can be tested directly, since the true posterior of its parameter vector can be computed in closed form. Recall, however, that the posterior distribution in Equation (2.14) refers to specific values of  $\alpha$  and  $\beta$ . The parameter  $\beta$  is  $\frac{1}{\sigma^2}$ , where  $\sigma$  is the standard deviation of the outputs  $y_n$  in Equation (2.6). The precision  $\alpha$  is a hyperparameter which controls the width of the prior  $p(\theta)$ . In Equation (2.10), the log-posterior of  $\theta$  was written as

$$\log p(\theta|\mathcal{D}) = -\frac{\beta}{2} \sum_{n=1}^N (y_n - \theta^\top \mathbf{x}_n)^2 - \frac{\alpha}{2} \theta^\top \theta + \text{constant}. \quad (6.1)$$

Hence, the maximum *a posteriori* (MAP) estimate of  $\theta$  is that which maximises

$$-\frac{\beta}{2} \sum_{n=1}^N (y_n - \theta^\top \mathbf{x}_n)^2 - \frac{\alpha}{2} \theta^\top \theta. \quad (6.2)$$

Since  $\beta > 0$ , it is equivalent to minimise this expression scaled by  $-\frac{2}{\beta}$ , that is,

$$\sum_{n=1}^N (y_n - \theta^\top \mathbf{x}_n)^2 + \frac{\alpha}{\beta} \theta^\top \theta. \quad (6.3)$$

This objective function is analogous to the L2-regularised training loss often used in non-Bayesian linear regression [1], which is

$$\sum_{n=1}^N (y_n - \theta^\top \mathbf{x}_n)^2 + \lambda \theta^\top \theta \quad (6.4)$$

for some  $\lambda > 0$ . Setting  $\lambda = \frac{\alpha}{\beta}$ , the optimal  $\theta$  found by L2-regularised linear regression is the same as the MAP estimate of  $\theta$  in Bayesian linear regression with prior  $p(\theta) = \mathcal{N}(\mathbf{0}, \alpha^{-1} \mathbf{I})$ . This means that, given data  $\mathcal{D}$  and a particular value of  $\alpha$ , the SGD iterates sampled while minimising Equation (6.4) with  $\lambda = \frac{\alpha}{\beta}$  should in theory be representative of the true posterior  $p(\theta|\mathcal{D})$ . Also, for the same data  $\mathcal{D}$  and hyperparameters  $\alpha$  and  $\beta$ , the VI approach in Algorithm 3 should be able to approximate  $p(\theta|\mathcal{D})$  directly. The purpose of the experiments in this section is to test these hypotheses.

### 6.1.1 Methodology

Synthetic data was generated for these experiments as follows. First, 1000 inputs  $\mathbf{x} \in \mathbb{R}^2$  were sampled from a multivariate zero mean Gaussian distribution with covariance matrix

$$\begin{bmatrix} 1 & 0.5 \\ 0.5 & 1 \end{bmatrix}. \quad (6.5)$$

Next, the parameter vector  $\theta \in \mathbb{R}^2$  was sampled from  $\mathcal{N}(\mathbf{0}, \alpha^{-1} \mathbf{I})$  with  $\alpha = 0.01$ . Then the outputs  $y \in \mathbb{R}$  were generated according to Equation (2.4) with  $\beta = 0.1$ . Using this data, the true posterior in Equation (2.14) was evaluated.

For the purpose of obtaining samples from the posterior, the linear regression model was pre-trained for 500 epochs of SGD with a learning rate of 0.001 and a batch size of 100. Then, starting from the pre-trained model, a further 100 epochs of SGD were executed, during which the parameter vector sampled after each mini-batch gradient update was collected. During the collection period, the batch size was again set to 100 but different values of the learning rate were tested, namely 0.1 and 0.01. As for regularisation, as well as the theoretically motivated value of  $\lambda = \frac{\alpha}{\beta} = 0.1$ , a value of  $\lambda = 0.001$  was also tested. Since there were  $\frac{1000}{100} = 10$  batch updates per epoch, a total of 1000 parameter vectors were collected over the final 100 epochs. The empirical mean and covariance of these samples was then computed and compared to the ground truth.

For the VI experiments, Algorithm 3 was executed with the same synthetic data and precision  $\alpha$ . In this case, the training process ran for 5000 epochs with a batch size of

$M = 100$  and a Monte Carlo average size of  $L = 10$ . The learning rates  $\eta_{\mathbf{c}}$ ,  $\eta_{\mathbf{F}}$  and  $\eta_{\beta}$  were set to 0.01, 0.0001 and 0.01, respectively. The reason for using a smaller learning rate for  $\mathbf{F}$  is that its contribution to the full covariance matrix is  $\mathbf{F}\mathbf{F}^T$ . Two different latent dimensions were tested for the approximate FA posterior, namely  $K = 1$  and  $K = 2$ . When computing the log-likelihood term involved in the gradient calculations, the variance of the outputs  $y$  was set to the true value,  $\sigma^2 = \frac{1}{\beta} = 10$ . Finally, to improve numerical stability, any gradients with Frobenius norm greater than 10 were rescaled to have norm of exactly 10.

All experiments were repeated ten times. In each trial a different random seed was used for generating the data, the true parameter vector and the initial parameters of the linear regression model.

### 6.1.2 Results and Discussion

The results of the posterior sampling experiments are shown in Table 6.1. It is apparent from the third column that the empirical mean of the sampled parameter vectors does not match the true posterior mean when the theoretically motivated value of  $\lambda = 0.1$  is used. However, if the regularisation strength is reduced to  $\lambda = 0.001$ , the distance between the means is much smaller. Unfortunately, the distance between the covariance matrices of the two distributions is relatively large for all combinations of the hyperparameters, as can be seen in the fourth column. The best result is achieved when using a learning rate of 0.1 with  $\lambda = 0.1$ , but even then the average relative distance is 0.62.

Figure 6.1 shows, for each combination of the hyperparameters, the sampled parameter vectors for one of the ten trials in each posterior sampling experiment. These are plotted on top of the contours of the pdf of the true posterior distribution. Note that in each case the pre-trained parameter vector is very close to the true posterior mean. When  $\lambda = 0.1$ , SGD moves away from the pre-trained solution to a different region of the parameter space (subplots (a) and (c)). When  $\lambda = 0.001$ , SGD bounces around the pre-trained solution and the empirical mean of the parameter vectors,  $\theta_{\text{SWA}}$ , is right on the true posterior mean. However, when the learning rate is 0.1 the empirical covariance is in the wrong direction (subplot (b)), and when the learning rate is 0.01 the empirical covariance is far too narrow (subplot (d)).

These results suggest that, in general, SGD does not sample parameter vectors from the true posterior distribution. This is a very simple example of a linear model in two dimensions. Even when tuning the learning rate and regularisation strength, it is not

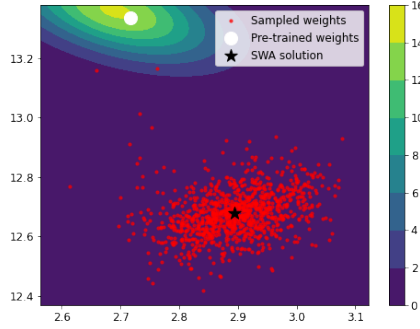
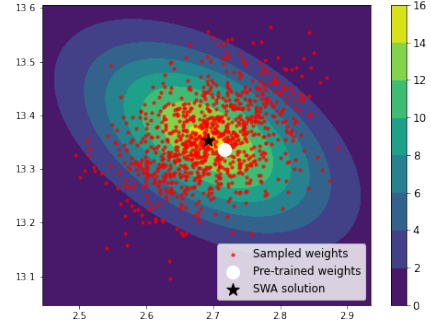
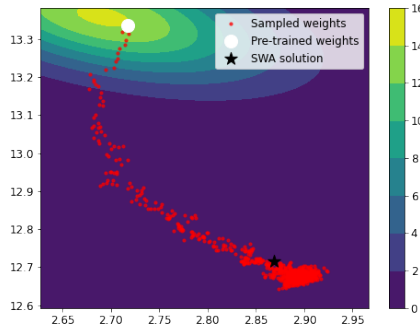
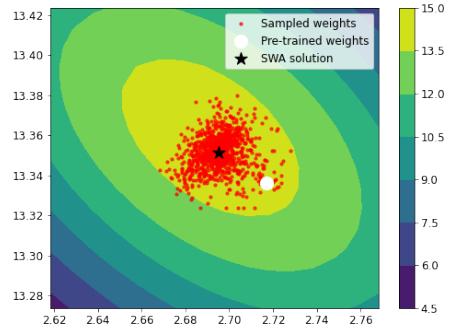
(a) Learning rate of 0.1 and  $\lambda = 0.1$ (b) Learning rate of 0.1 and  $\lambda = 0.001$ (c) Learning rate of 0.01 and  $\lambda = 0.1$ (d) Learning rate of 0.01 and  $\lambda = 0.001$ 

Figure 6.1: Parameter vectors sampled while training a linear regression model via SGD, for different combinations of learning rate and regularisation strength. These are shown as red dots, plotted on top of the contours of the pdf of the true posterior distribution. In each plot, the white dot denotes the pre-trained parameter vector and the black star is the empirical mean of the red dots.

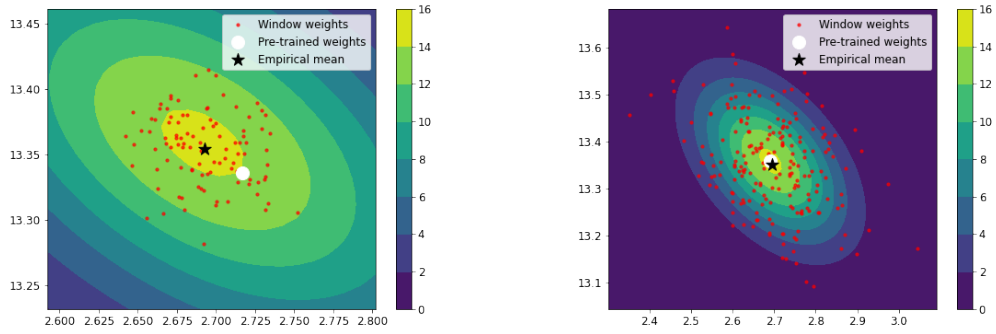
Learning Rate	$\lambda$	Relative Distance from Mean	Relative Distance from Covariance	Wasserstein Distance
0.1	0.1	$0.0666 \pm 0.0078$	$0.6200 \pm 0.0750$	$1.3480 \pm 0.5087$
0.1	0.001	$0.0006 \pm 0.0001$	$0.8065 \pm 0.0101$	$0.0078 \pm 0.0004$
0.01	0.1	$0.0620 \pm 0.0070$	$1.7110 \pm 0.6187$	$1.1776 \pm 0.4434$
0.01	0.001	$0.0010 \pm 0.0002$	$0.9870 \pm 0.0042$	$0.0219 \pm 0.0006$

Table 6.1: Distances between the true posterior distribution of the parameter vector of a linear regression model and the empirical distribution of the parameter vectors sampled along the SGD trajectory while training the model. Relative distances between the true and empirical means and covariances are shown. Each relative distance is the Frobenius norm of the difference between the true parameter and the empirical parameter divided by the Frobenius norm of the true parameter. Also shown is the 2-Wasserstein distance between the true posterior Gaussian distribution and the empirical Gaussian distribution. Each distance is the mean value over ten trials with different random seeds, and standard errors are also given.

possible to obtain samples which are representative of the true posterior. For certain values of  $\lambda$  (but not the theoretically motivated one), the mean solution  $\theta_{\text{SWA}}$  is a good approximation of the true posterior mean, but the same cannot be said for the empirical covariance. The most that can be expected of the online FA algorithms is that they fit the sampled parameter vectors well. Therefore, given the data observed in Figure 6.1, it would be unreasonable to expect them to learn the true posterior distribution. In the case of a high-dimensional multi-layer neural network, approximating the true posterior would be even more unlikely. Since the ground truth cannot be compute in closed form, hyperparameter tuning would be far more difficult. Therefore, contrary to what is claimed in [11], this suggest that SWAG does not actually do approximate Bayesian inference in general.

It was also observed in [12] that the empirical covariance of the linear regression SGD iterates does not match the true posterior covariance. As such, the authors proposed a different technique for posterior sampling called *finite window iterate averaging*. Instead of using the individual parameter vectors sampled after each mini-batch update of SGD, this algorithm averages the parameter vectors in fixed-length disjoint windows. These average parameter vectors are then considered samples from the pos-

terior. Figure 6.2 shows examples of such samples for the same learning rate and regularisation strength as that in Figure 6.1 (b). The theoretical work in [12] suggests that the window size should be set to the number of training examples divided by the batch size, so that all training examples are used to generate a single posterior sample. In these experiments this value is  $1000/100 = 10$ . However, this setting does not produce samples which are representative of the true posterior, as observed in Figure 6.2 (a). Better samples can be obtained by setting the batch size to 10 and the window size to 50, as shown in Figure 6.2 (b). Unfortunately, this result is not theoretically motivated and was only obtained after significant tuning of both the batch size and window size. In the case of non-linear neural networks for which the true posterior is not known, this would not be possible. Therefore, this algorithm does not seem to be a viable option in practice.



(a) Batch size of 100, window size of 10      (b) Batch size of 10, window size of 50

Figure 6.2: Samples constructed via finite window iterate averaging for different combinations of batch size and window size. Each red dot is the mean of a set of linear regression parameter vectors sampled during consecutive mini-batch updates of SGD, with a learning rate of 0.1 and the regularisation strength set to 0.001. The contours show the pdf of the true posterior distribution. The white dot denotes the pre-trained parameter vector and the black star is the empirical mean of the red dots.

Moving on to VI, the results of applying Algorithm 3 to the data are shown in Table 6.2. Compared to using the SGD iterates, the more direct VI approach does a much better job of approximating the true covariance matrix. For the FA model with a single latent dimension, the average distance of the estimated covariance from the ground truth is 0.0983, which is much lower than the best result of 0.6200 from Table 6.1. The much improved approximation of the covariance can also be observed in Figure

6.3 for one of the ten VI trials. In this example, both the diagonal and off-diagonal entries of the covariance matrices of the two FA models are very similar to the ground truth. Moreover, in these experiments the FA model with a single latent dimension achieved better results than the FA with two latent dimensions. This is encouraging, since applying this algorithm to estimate the posterior of a deep neural network is only feasible if the latent dimension can be set to a value much less than the number of learnable parameters in the model.

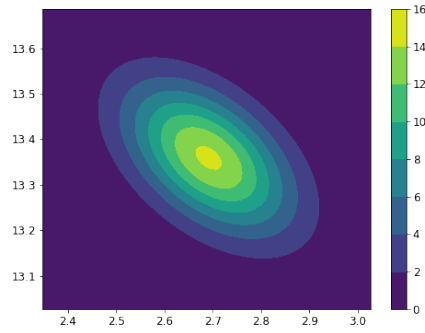
Latent Dimension of FA Model	Relative Distance from Mean	Relative Distance from Covariance	Wasserstein Distance
1	$0.0031 \pm 0.0005$	$0.0983 \pm 0.0129$	$0.0017 \pm 0.0004$
2	$0.0040 \pm 0.0008$	$0.1053 \pm 0.0199$	$0.0027 \pm 0.0006$

Table 6.2: Distances between the true posterior distribution of the parameter vector of a linear regression model and the approximate FA posterior estimated using VIFA. Relative distances between the true and approximate means and covariances are shown. Each relative distance is the Frobenius norm of the difference between the true parameter and the empirical parameter divided by the Frobenius norm of the true parameter. Also shown is the 2-Wasserstein distance between the true posterior Gaussian distribution and the empirical Gaussian distribution. Each distance is the mean value over ten trials with different random seeds, and standard errors are also given.

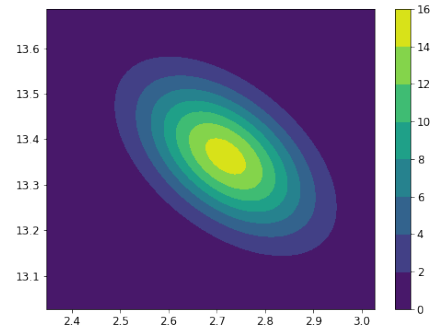
## 6.2 Making Predictions

Regardless of whether SWAG does Bayesian inference or not, it was shown in [11] that this method can lead to more accurate predictions. Once a Gaussian distribution has been learned from the SGD iterates, parameter vectors of the model can easily be sampled. For each sampled parameter vector, the resultant model can be used to make a different set of predictions for the test data. Final predictions are then obtained by averaging the predictions from all models. The aim of the experiments in this section is to test, for linear regression, whether this approach leads to more accurate predictions than standard SGD, when the Gaussian distribution of the parameter vectors is a FA model learned by either online SGA or online EM.

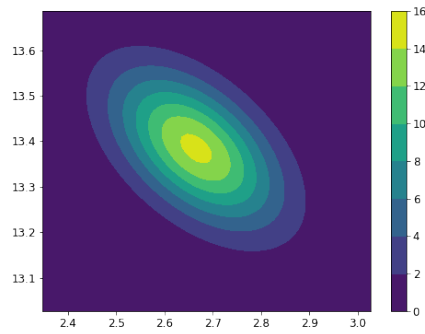




(a) Ground truth



(b) VIFA: latent dim = 1



(c) VIFA: latent dim = 2

Figure 6.3: The true posterior pdf of a linear regression model with two learnable parameters, plus the pdfs of two FA models fit to the same data using VIFA. Both approximations are very similar to the ground truth, including the FA model with only a single latent dimension.

### 6.2.1 Methodology

For these experiments, four real regression datasets from the UCI Machine Learning Repository [2] were used. Namely, the Energy Efficiency<sup>1</sup>, Boston Housing<sup>2</sup>, Concrete Compressive Strength<sup>3</sup> and Yacht Hydrodynamics<sup>4</sup> datasets. More details about these datasets are given in Table 6.3. Note that the Energy Efficiency dataset has two target variables, heating load and cooling load, but in these experiments only the heating load was used.

Dataset	No. of Instances	No. of Input Variables
Energy Efficiency	768	8
Boston Housing	506	13
Concrete Compressive Strength	1030	8
Yacht Hydrodynamics	308	6

Table 6.3: UCI regression datasets used in experiments with linear regression.

For each dataset, the model was pre-trained for 500 epochs of SGD with a learning rate of 0.001, as in the synthetic experiments in Section 6.1. Then, starting from the pre-trained model, a further 100 epochs of SGD were executed. During this time, the parameter vectors sampled after each mini-batch gradient update were used to fit two FA models, one via online SGA and another via online EM. For these final 100 epochs, a learning rate of 0.1 and weight decay of 0.001 were used. These hyperparameters were chosen since it was observed in Figure 6.1 that this setting produced samples most representative of the true posterior for the synthetic data. During the pre-training phase, the weight decay was also set to 0.001 to make for a fairer comparison with the pre-trained model. For all epochs, the batch size was set to one tenth of the size of the full dataset. Both FA models had a latent dimension of three. Following the approach in the experiments in Chapter 5, the first 100 sampled parameter vectors were only used to update the running averages,  $\bar{\theta}_t$ ,  $\bar{\mathbf{A}}_t$ ,  $\bar{\mathbf{B}}_t$  and  $\bar{\mathbf{d}}^2_t$ , while  $\mathbf{F}$  and  $\Psi$  remained fixed. Again, the FA model updated via online SGA used a learning rate of 0.001. After training, test predictions were made using the pre-trained parameter vector, the empirical mean

<sup>1</sup><https://archive.ics.uci.edu/ml/datasets/energy+efficiency>

<sup>2</sup><https://archive.ics.uci.edu/ml/machine-learning-databases/housing>

<sup>3</sup><https://archive.ics.uci.edu/ml/datasets/concrete+compressive+strength>

<sup>4</sup><http://archive.ics.uci.edu/ml/datasets/yacht+hydrodynamics>

of the sampled parameter vectors,  $\theta_{\text{SWA}}$ , and two ensembles, one for each FA model. Each ensemble was constructed by sampling 30 parameter vectors from the FA model and averaging the resultant predictions, as in the UCI regression experiments in [11]. This process was repeated for 10 folds of cross-validation. Within each fold, the input features were standardised by subtracting their mean and dividing by their standard deviation, where these statistics were calculated using the fold's training data only.

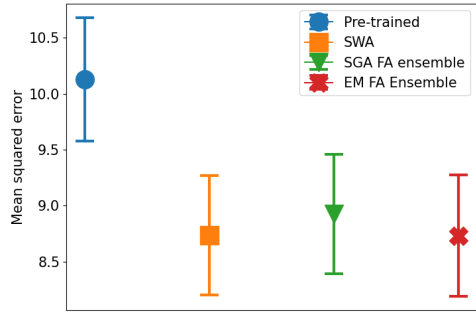
## 6.2.2 Results and Discussion

Figure 6.4 shows the average test mean squared error (MSE) for each dataset and prediction method. The average errors were computed over ten folds of cross-validation, and standard error bars are also plotted. Across all datasets, the SWA solution and the FA ensemble learned via online EM are the best prediction methods in general in terms of the average error. These two methods achieved almost identical results on all datasets. On three out of the four datasets they achieve lower average error than the pre-trained solution. The difference is relatively large in the case of the Energy Efficiency dataset, but for all other datasets the standard error bars have significant overlap. Recall that, in the experiments in Chapter 5, online EM was better at fitting FA models than online SGA. This appears to translate into more accurate linear regression ensembles, as the average MSE of the online EM ensemble is less than that of the online SGA ensemble for all datasets. The numerical MSE results are also given in Table A.1 of the appendix.

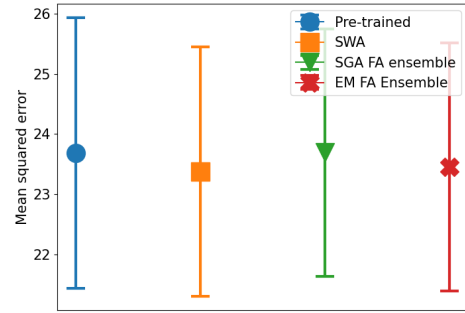
These results suggest that, in the case of linear regression, fitting a FA model to the SGD iterates via online EM and then using it to construct an ensemble can lead to more accurate predictions compared to standard SGD. However, this approach does not seem to improve on the much simpler SWA solution. Note that, in the case of linear regression, an ensemble of  $M$  parameter vectors sampled from a FA model can be written as

$$\frac{1}{M} \sum_{m=1}^M \theta_m^T \mathbf{x} = \left( \frac{1}{M} \sum_{m=1}^M \theta_m^T \right) \mathbf{x} = \theta_{\text{FA}}^T \mathbf{x}. \quad (6.6)$$

Hence, if the learned FA model is a good approximation of the empirical Gaussian distribution of the SGD iterates,  $\theta_{\text{FA}}$  will be similar to  $\theta_{\text{SWA}}$  for large  $M$ . This observation and the results in Figure 6.4 suggests that, if the accuracy of linear regression predictions is all that matters, the best approach is to use  $\theta_{\text{SWA}}$ .



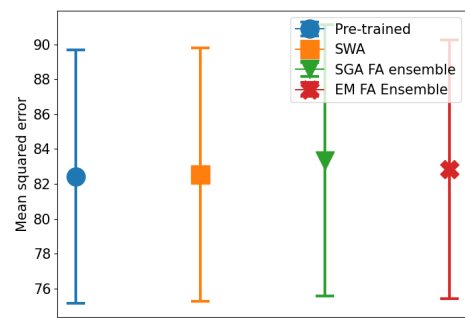
(a) Energy Efficiency



(b) Boston Housing



(c) Concrete Compressive Strength



(d) Yacht Hydrodynamics

Figure 6.4: Average test mean squared error of linear regression models on four UCI regression datasets. Results are shown for four different prediction methods: using the pre-trained parameter vector (blue), using the SWA solution (orange) and two ensembles formed by sampling parameter vectors from FA models fit to the SGD iterates, one via online SGA (green) and another learned via online EM (red). The markers show the average test error over ten folds of cross-validation, and standard error bars are also plotted.

# Chapter 7

## Conclusions

### 7.1 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.
- [2] Dheeru Dua and Casey Graff. UCI machine learning repository, 2017.
- [3] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT Press, 2016.
- [4] Ian J Goodfellow. Efficient per-example gradient computations. *arXiv preprint arXiv:1510.01799*, 2015.
- [5] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 770–778, 2016.
- [6] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9:1735–1780, 1997.
- [7] Mohammad E Khan, Didrik Nielsen, Voot Tangkaratt, Wu Lin, Yarin Gal, and Akash Srivastava. Fast and scalable Bayesian deep learning by weight-perturbation in Adam. *Proceedings of the International Conference on Machine Learning*, 35:2611—2620, 2018.
- [8] Diederik P Kingma and Max Welling. Auto-encoding variational Bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- [9] Alex Krizhevsky. Learning multiple layers of features from tiny images. *University of Toronto*, 2012.
- [10] Frederik Kunstner, Philipp Hennig, and Lukas Balles. Limitations of the empirical fisher approximation for natural gradient descent. *Proceedings of Advances in Neural Information Processing Systems*, 32:4158–4169, 2019.

- [11] Wesley Maddox, Timur Garipov, Pavel Izmailov, Dmitry Vetrov, and Andrew G Wilson. A simple baseline for bayesian uncertainty in deep learning. *Proceedings of Advances in Neural Information Processing Systems*, 32, 2019.
- [12] Stephan Mandt, Matthew D Hoffman, and David M Blei. Stochastic gradient descent as approximate Bayesian inference. *Journal of Machine Learning Research*, 18:4873–4907, 2017.
- [13] Aaron Mishkin, Frederik Kunstner, Didrik Nielsen, Mark Schmidt, and Mohammad E Khan. SLANG: Fast structured covariance approximations for Bayesian deep learning with natural gradient. *Proceedings of Advances in Neural Information Processing Systems*, 31:6248–6258, 2018.
- [14] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Köpf, Edward Yang, Zach DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. *Proceedings of Advances in Neural Information Processing Systems*, 32:8024–8035, 2019.
- [15] Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel, Andreas Müller, Joel Nothman, Gilles Louppe, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, Jake Vanderplas, Alexandre Passos, David Cournapeau, Matthieu Brucher, Matthieu Perrot, and Édouard Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- [16] Kaare B Petersen and Michael S Pedersen. *The Matrix Cookbook*. Technical University of Denmark, 2012.
- [17] David E Rumelhart, Geoffrey E Hinton, and Ronald J Williams. Learning representations by back-propagating errors. *Nature*, 323:533–536, 1986.
- [18] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Proceedings of Advances in Neural Information Processing Systems*, 30:5998–6008, 2017.

# Appendix A

## Supplementary Results

### A.1 Linear Regression Predictions

Table A.1 contains the same results as Figure 6.4, but in tabular format.

Dataset	Pre-trained	SWA	SGA-FA	EM-FA
Energy Efficiency	$10.13 \pm 0.55$	$8.74 \pm 0.54$	$8.93 \pm 0.53$	<b><math>8.73 \pm 0.54</math></b>
Boston Housing	$23.69 \pm 2.25$	<b><math>23.27 \pm 2.07</math></b>	$23.69 \pm 2.06$	$23.46 \pm 2.06$
Concrete Strength	$112.44 \pm 6.93$	$109.83 \pm 6.44$	$109.84 \pm 6.30$	<b><math>109.79 \pm 6.41</math></b>
Yacht Hydro.	<b><math>82.43 \pm 7.28</math></b>	$82.54 \pm 7.27$	$83.35 \pm 7.78$	$82.83 \pm 7.42$

Table A.1: Average test mean squared error of linear regression models on four UCI regression datasets. Results are show for four different prediction methods: using the pre-trained parameter vector (second column), using the SWA solution (third column) and two ensembles formed by sampling parameter vectors from FA models fit to the SGD iterates, one via online SGA (fourth column) and another learned via online EM (fifth column). Results are average over ten folds of cross-validation, and standard errors are also given. Bold text highlights the lowest average error for each dataset, without taking into account the standard errors.

### A.2 Comparison of VIFA and SLANG

$$\mathbf{c} \leftarrow \mathbf{c} - \eta(\alpha \mathbf{c} + \mathbf{g}_0) \quad (\text{A.1})$$



$$\mathbf{F} \leftarrow \mathbf{F} - \eta(-\mathbf{S}^{-1}\mathbf{F} + \alpha\mathbf{F} + \mathbf{g}_\theta\mathbf{h}^\top) \quad (\text{A.2})$$

$$\Psi \leftarrow \Psi - \frac{\eta}{2} \text{diag} \left( \text{diag}(\mathbf{S}^{-1} - \mathbf{g}_\theta\mathbf{z}^\top\Psi^{-1/2}) \right) + \frac{\eta\alpha}{2}\mathbf{I} \quad (\text{A.3})$$

$$\begin{aligned} \mathbf{F}\mathbf{F}^\top &\leftarrow \left( \mathbf{F} - \eta(-\mathbf{S}^{-1}\mathbf{F} + \alpha\mathbf{F} + \mathbf{g}_\theta\mathbf{h}^\top) \right) \left( \mathbf{F} - \eta(-\mathbf{S}^{-1}\mathbf{F} + \alpha\mathbf{F} + \mathbf{g}_\theta\mathbf{h}^\top) \right)^\top \\ &= \mathbf{F}\mathbf{F}^\top - \eta\mathbf{F}(-\mathbf{S}^{-1}\mathbf{F} + \alpha\mathbf{F} + \mathbf{g}_\theta\mathbf{h}^\top)^\top - \eta(-\mathbf{S}^{-1}\mathbf{F} + \alpha\mathbf{F} + \mathbf{g}_\theta\mathbf{h}^\top)\mathbf{F}^\top \\ &\quad + \eta^2(-\mathbf{S}^{-1}\mathbf{F} + \alpha\mathbf{F} + \mathbf{g}_\theta\mathbf{h}^\top)(-\mathbf{S}^{-1}\mathbf{F} + \alpha\mathbf{F} + \mathbf{g}_\theta\mathbf{h}^\top)^\top \\ &= \mathbf{F}\mathbf{F}^\top - \eta\mathbf{F}(-\mathbf{F}^\top\mathbf{S}^{-1} + \alpha\mathbf{F}^\top + \mathbf{h}\mathbf{g}_\theta^\top) - \eta(-\mathbf{S}^{-1}\mathbf{F} + \alpha\mathbf{F} + \mathbf{g}_\theta\mathbf{h}^\top)\mathbf{F}^\top \\ &\quad + \eta^2(-\mathbf{S}^{-1}\mathbf{F} + \alpha\mathbf{F} + \mathbf{g}_\theta\mathbf{h}^\top)(-\mathbf{F}^\top\mathbf{S}^{-1} + \alpha\mathbf{F}^\top + \mathbf{h}\mathbf{g}_\theta^\top) \\ &= \mathbf{F}\mathbf{F}^\top + \eta\mathbf{F}\mathbf{F}^\top\mathbf{S}^{-1} - \eta\alpha\mathbf{F}\mathbf{F}^\top - \eta\mathbf{F}\mathbf{h}\mathbf{g}_\theta^\top + \eta\mathbf{S}^{-1}\mathbf{F}\mathbf{F}^\top - \eta\alpha\mathbf{F}\mathbf{F}^\top - \eta\mathbf{g}_\theta\mathbf{h}^\top\mathbf{F}^\top \\ &\quad + \eta^2\mathbf{S}^{-1}\mathbf{F}\mathbf{F}^\top\mathbf{S}^{-1} - \eta^2\alpha\mathbf{S}^{-1}\mathbf{F}\mathbf{F}^\top - \eta^2\mathbf{S}^{-1}\mathbf{F}\mathbf{h}\mathbf{g}_\theta^\top \\ &\quad - \eta^2\alpha\mathbf{F}\mathbf{F}^\top\mathbf{S}^{-1} + \eta^2\alpha^2\mathbf{F}\mathbf{F}^\top + \eta^2\alpha\mathbf{F}\mathbf{h}\mathbf{g}_\theta^\top \\ &\quad - \eta^2\mathbf{g}_\theta\mathbf{h}^\top\mathbf{F}^\top\mathbf{S}^{-1} + \eta^2\alpha\mathbf{g}_\theta\mathbf{h}^\top\mathbf{F}^\top + \eta^2\mathbf{g}_\theta\mathbf{h}^\top\mathbf{h}\mathbf{g}_\theta^\top \\ &= \mathbf{F}\mathbf{F}^\top(1 - 2\eta\alpha + \eta^2\alpha^2) + \mathbf{F}\mathbf{F}^\top\mathbf{S}^{-1}(\eta - \eta^2\alpha) + \mathbf{S}^{-1}\mathbf{F}\mathbf{F}^\top(\eta - \eta^2\alpha) \\ &\quad + \eta^2\mathbf{S}^{-1}\mathbf{F}\mathbf{F}^\top\mathbf{S}^{-1} - \eta^2\mathbf{S}^{-1}\mathbf{F}\mathbf{h}\mathbf{g}_\theta^\top - \eta^2\mathbf{g}_\theta\mathbf{h}^\top\mathbf{F}^\top\mathbf{S}^{-1} \\ &\quad + \mathbf{F}\mathbf{h}\mathbf{g}_\theta^\top(\eta^2\alpha - \eta) + \mathbf{g}_\theta\mathbf{h}^\top\mathbf{F}(\eta^2\alpha - \eta) + \eta^2\mathbf{h}^\top\mathbf{h}\mathbf{g}_\theta\mathbf{g}_\theta^\top \end{aligned} \quad (\text{A.4})$$