

# IDIAP COMMUNICATION REPORT



## NOTES ON PROBABILISTIC LINEAR DISCRIMINANT ANALYSIS

Chris McCool      Laurent El Shafey

Idiap-Com-03-2013

JUNE 2013



# Notes on Probabilistic Linear Discriminant Analysis

Christopher Steven McCool and Laurent El Shafey

([csmccool79@gmail.com](mailto:csmccool79@gmail.com), [laurent.el-shafey@idiap.ch](mailto:laurent.el-shafey@idiap.ch))

May 15, 2012

Idiap Research Institute,  
Centre du Parc,  
PO Box 592,  
1920 Martigny,  
Switzerland

# Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
1.1	Acknowledgments . . . . .	3
<b>2</b>	<b>The PLDA Model</b>	<b>3</b>
2.1	Setting up the Basics . . . . .	3
2.2	Extending the notation for Multiple Samples of the Same Identity . . . . .	4
2.3	Completing the Formulation . . . . .	4
2.3.1	Deriving $p(\tilde{\mathbf{y}}_i \tilde{\mathbf{x}}_i, \Theta)$ and $p(\tilde{\mathbf{x}}_i)$ . . . . .	5
2.3.2	Some Sanity Checks on Array and Matrix Sizes . . . . .	5
2.4	Sufficient Statistics . . . . .	6
2.4.1	Complexity of the Problem . . . . .	6
2.4.2	Calculating $E[\tilde{\mathbf{y}}_i \tilde{\mathbf{x}}_i, \Theta]$ . . . . .	6
2.4.3	Calculating $E[\tilde{\mathbf{y}}_i\tilde{\mathbf{y}}_i^T \tilde{\mathbf{x}}_i, \Theta]$ . . . . .	7
2.5	Authentication Scores for PLDA . . . . .	7
<b>3</b>	<b>Learning the Parameters</b>	<b>8</b>
3.1	Differentiating with Respect to $\boldsymbol{\mu}$ . . . . .	9
3.2	Differentiating with Respect to $\mathbf{A}$ . . . . .	10
3.3	Differentiating with Respect to $\boldsymbol{\Sigma}$ . . . . .	10
<b>4</b>	<b>Implementation Details</b>	<b>12</b>
4.1	Log-likelihood of $\tilde{\mathbf{x}}$ . . . . .	13
4.1.1	Calculating $-\frac{1}{2}(\tilde{\mathbf{x}}_i - \tilde{\boldsymbol{\mu}})^T (\tilde{\boldsymbol{\Sigma}} + \tilde{\mathbf{A}}\tilde{\mathbf{A}}^T)^{-1} (\tilde{\mathbf{x}}_i - \tilde{\boldsymbol{\mu}})$ Efficiently: Method 1 . . . . .	13
4.1.2	Calculating $-\frac{1}{2}(\tilde{\mathbf{x}}_i - \tilde{\boldsymbol{\mu}})^T (\tilde{\boldsymbol{\Sigma}} + \tilde{\mathbf{A}}\tilde{\mathbf{A}}^T)^{-1} (\tilde{\mathbf{x}}_i - \tilde{\boldsymbol{\mu}})$ Efficiently: Method 2 . . . . .	15
4.1.3	Determinant of $(\tilde{\boldsymbol{\Sigma}} + \tilde{\mathbf{A}}\tilde{\mathbf{A}}^T)$ . . . . .	18
4.1.4	Final Solution . . . . .	20
4.2	Sufficient Statistics . . . . .	21
4.2.1	Solving for $E[\tilde{\mathbf{y}}_i \tilde{\mathbf{x}}_i, \Theta]$ . . . . .	21
4.2.2	Solving for $E[\tilde{\mathbf{y}}_i\tilde{\mathbf{y}}_i^T \tilde{\mathbf{x}}_i, \Theta]$ . . . . .	22
<b>A</b>	<b>Rules and Identities for Differentiation</b>	<b>24</b>
<b>B</b>	<b>Matrix Identities</b>	<b>24</b>
B.1	Inverse of a Product of Matrices . . . . .	24
B.2	Block Matrix Inversion . . . . .	24
B.3	Block LU Decomposition . . . . .	25
B.4	Block LDU Decomposition . . . . .	26
B.5	Square Root of a Matrix . . . . .	26
B.6	Log Determinant of a Diagonal Matrix . . . . .	26
B.7	Log Determinant of a Symmetric Real Matrix . . . . .	26
B.8	Inverse of a Block Diagonal Matrix . . . . .	27
<b>C</b>	<b>Intermediate Solutions for PLDA: Gaussian Priors</b>	<b>27</b>
C.1	Matrix Inversion: $(\mathbf{I} + \tilde{\mathbf{A}}^T\tilde{\boldsymbol{\Sigma}}^{-1}\tilde{\mathbf{A}})^{-1}$ . . . . .	27
C.1.1	Finding $(A - BD^{-1}C)^{-1}$ . . . . .	28
C.1.2	Finding $-D^{-1}C(A - BD^{-1}C)^{-1}$ . . . . .	28
C.1.3	Finding $D^{-1} + D^{-1}C(A - BD^{-1}C)^{-1}BD^{-1}$ . . . . .	28
C.2	LDU Decomposition for $(\mathbf{I} + \tilde{\mathbf{A}}^T\tilde{\boldsymbol{\Sigma}}^{-1}\tilde{\mathbf{A}})^{-1}$ . . . . .	30

# 1 Introduction

The following are a set of working notes on how to derive and implement probabilistic linear discriminant analysis (PLDA) using Gaussian priors as proposed in [4]. This includes:

- how to derive the update formulae using maximum likelihood (ML),
- how to derive the formulae for estimating the latent variables,
- how to derive the formulae to perform scoring, and
- finally simplifications for memory and computation time.

We start by introducing the PLDA model. Then we go on to describe how it can be derived so that we can train the parameters. After this we deal with how to derive an efficient way to implement the model to perform training and to compute the likelihood of some input samples given this model. For clarity we have provided a set of Appendices which contains several useful definitions and derivations used in the text.

## 1.1 Acknowledgments

We want to thank many people who helped us in writing down these derivations. First and foremost we would like to thank Sébastien Marcel who gave us the freedom to examine this model in more detail. Other people that we would like to thank are Carl Scheffler and the rest of the people at the Biometrics Group of the Idiap Research Institute.

## 2 The PLDA Model

The PLDA model of Prince and Elder [4] consists of a mean offset  $\boldsymbol{\mu}$ , a subspace describing the main directions of identity variation  $\mathbf{F}$ , a subspace describing the main directions of condition, or session variation,  $\mathbf{G}$  and a noise term  $\boldsymbol{\epsilon}_{ij}$ , where  $i$  is the  $i^{\text{th}}$  identity and  $j$  is their  $j^{\text{th}}$  observation. We can then describe PLDA as the process to represent the  $D_x$  dimensional feature vector  $\mathbf{x}_{ij}$  such that

$$\mathbf{x}_{ij} = \boldsymbol{\mu} + \mathbf{F}\mathbf{h}_i + \mathbf{G}\mathbf{w}_{ij} + \boldsymbol{\epsilon}_{ij}, \quad (1)$$

where  $\mathbf{h}_i$  is the latent variable representing identity variation (the weight of the directions for identity variation) and  $\mathbf{w}_{ij}$  is the latent variable representing condition or session variation (the weight of the directions for identity variation).

### 2.1 Setting up the Basics

To solve this set of equations we place certain restrictions upon it. First, we assume that  $p(\mathbf{h}_i)$  is a zero mean unit standard deviation multivariate Gaussian ( $\mathcal{N}_{\mathbf{h}}[\mathbf{0}, \mathbf{I}]$ ), that  $p(\mathbf{w}_{ij})$  is a zero mean unit standard deviation multivariate Gaussian ( $\mathcal{N}_{\mathbf{w}}[\mathbf{0}, \mathbf{I}]$ ) and that  $\boldsymbol{\epsilon}_{ij}$  is assumed to be Gaussian with a diagonal covariance  $\boldsymbol{\Sigma}$ . We can then write Equation 1 as a conditional probability  $Pr(\mathbf{x}_{ij}|\mathbf{h}_i, \mathbf{w}_{ij}, \boldsymbol{\Theta})$  where the parameters for our PLDA model are defined as  $\boldsymbol{\Theta} = [\boldsymbol{\mu}, \mathbf{F}, \mathbf{G}, \boldsymbol{\Sigma}]$ . We can now write the following,

$$p(\mathbf{h}_i) = \mathcal{N}_{\mathbf{h}}[\mathbf{0}, \mathbf{I}], \quad (2)$$

$$p(\mathbf{w}_{ij}) = \mathcal{N}_{\mathbf{w}}[\mathbf{0}, \mathbf{I}], \quad (3)$$

$$p(\mathbf{x}_{ij}|\mathbf{h}_i, \mathbf{w}_{ij}, \boldsymbol{\Theta}) = \mathcal{N}_{\mathbf{x}}[\boldsymbol{\mu} + \mathbf{F}\mathbf{h}_i + \mathbf{G}\mathbf{w}_{ij}, \boldsymbol{\Sigma}]. \quad (4)$$

## 2.2 Extending the notation for Multiple Samples of the Same Identity

The notation from the previous section can be extended to include the case of multiple observations for a single identity  $i$ . In this case  $\tilde{\mathbf{x}}_i = [\mathbf{x}_{i1}^T, \mathbf{x}_{i2}^T, \dots, \mathbf{x}_{iJ_i}^T]^T$ ,  $\tilde{\boldsymbol{\mu}} = [\boldsymbol{\mu}^T, \boldsymbol{\mu}^T, \dots, \boldsymbol{\mu}^T]^T$ ,  $\tilde{\mathbf{y}}_i = [\mathbf{h}_i^T, \mathbf{w}_{i1}^T, \mathbf{w}_{i2}^T, \dots, \mathbf{w}_{iJ_i}^T]^T$  and  $\tilde{\boldsymbol{\epsilon}}_i = [\boldsymbol{\epsilon}_{i1}^T, \boldsymbol{\epsilon}_{i2}^T, \dots, \boldsymbol{\epsilon}_{iJ_i}^T]$  where  $J_i$  is the number of observations for identity  $i$ ; these are vectors of size  $(J_i D_x, 1)$ ,  $(J_i D_x, 1)$ ,  $(D_F + J_i D_G, 1)$  and  $(J_i D_x, 1)$  respectively. Using this we can write Equation 1 as,

$$\tilde{\mathbf{x}}_i = \tilde{\boldsymbol{\mu}} + \tilde{\mathbf{A}}\tilde{\mathbf{y}}_i + \tilde{\boldsymbol{\epsilon}}_i. \quad (5)$$

This is the same as Equation 11 in [4]. We take the example of  $J_i = 3$  (we will keep the same value of  $J_i$  for the rest of the examples in these notes) and demonstrate the structure of the variables we have introduced.

$$\tilde{\mathbf{A}} = \begin{bmatrix} \mathbf{F} & \mathbf{G} & \mathbf{0} & \mathbf{0} \\ \mathbf{F} & \mathbf{0} & \mathbf{G} & \mathbf{0} \\ \mathbf{F} & \mathbf{0} & \mathbf{0} & \mathbf{G} \end{bmatrix}, \quad (6)$$

$$\tilde{\boldsymbol{\Sigma}} = \begin{bmatrix} \boldsymbol{\Sigma} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\Sigma} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \boldsymbol{\Sigma} \end{bmatrix}. \quad (7)$$

We can then write the probabilities as,

$$p(\tilde{\mathbf{x}}_i | \tilde{\mathbf{y}}_i, \boldsymbol{\Theta}) = \mathcal{N}_{\tilde{\mathbf{x}}_i} [\tilde{\boldsymbol{\mu}} + \tilde{\mathbf{A}}\tilde{\mathbf{y}}_i, \tilde{\boldsymbol{\Sigma}}], \quad (8)$$

$$p(\tilde{\mathbf{y}}_i) = \mathcal{N}_{\tilde{\mathbf{y}}_i} [\mathbf{0}, \mathbf{I}]. \quad (9)$$

It can be seen that in this more compact formulation the matrix  $\tilde{\mathbf{A}}$  consists of the matrices  $\mathbf{F}$  and  $\mathbf{G}$  which are repeated  $J_i$  times and the matrix  $\tilde{\boldsymbol{\Sigma}}$  consists of the matrix  $\boldsymbol{\Sigma}$  also repeated  $J_i$  times. For later use, we now define another latent variable  $\mathbf{y}_{ij} = [\mathbf{h}_i^T, \mathbf{w}_{ij}^T]^T$  which is the latent variable for the  $j^{\text{th}}$  observation of identity  $i$ .

## 2.3 Completing the Formulation

With the above equations we can now write down  $p(\tilde{\mathbf{x}}_i)$  and  $p(\tilde{\mathbf{y}}_i | \tilde{\mathbf{x}}_i, \boldsymbol{\Theta})$  by using the identities B42-B45 (from page 689) of [1]. We have rewritten these identities below.

$$p(\mathbf{y}_B) = \mathcal{N}_{\mathbf{y}_B} [\boldsymbol{\mu}_{\mathbf{y}_B}, \boldsymbol{\Sigma}_B], \quad (10)$$

$$p(\mathbf{x}_B | \mathbf{y}_B) = \mathcal{N}_{\mathbf{x}_B} [\mathbf{A}_B \mathbf{y}_B + \mathbf{b}_B, \mathbf{L}_B^{-1}], \quad (11)$$

$$\boldsymbol{\Lambda}_B = \boldsymbol{\Sigma}_B^{-1}. \quad (12)$$

With these we can define the following:

$$p(\mathbf{x}_B) = \mathcal{N}_{\mathbf{x}_B} [\mathbf{A}_B \boldsymbol{\mu}_{\mathbf{y}_B} + \mathbf{b}_B, \mathbf{L}_B^{-1} + \mathbf{A}_B \boldsymbol{\Sigma}_B \mathbf{A}_B^T], \quad (13)$$

$$p(\mathbf{y}_B | \mathbf{x}_B) = \mathcal{N}_{\mathbf{y}_B} [\boldsymbol{\Sigma}^* (\mathbf{A}_B^T \mathbf{L}_B (\mathbf{x}_B - \mathbf{b}_B) + \boldsymbol{\Sigma}_B^{-1} \boldsymbol{\mu}_{\mathbf{y}_B}), \boldsymbol{\Sigma}^*], \quad (14)$$

$$\boldsymbol{\Sigma}^* = (\boldsymbol{\Sigma}_B^{-1} + \mathbf{A}_B^T \mathbf{L}_B \mathbf{A}_B)^{-1}. \quad (15)$$

Where we have used the subscript  $B$  to clearly indicate that these are the formulae from B42-B45 but with  $\mathbf{x}$  and  $\mathbf{y}$  reversed (which has no effect); we do this because it makes it easier to see the relationship for our problem.

### 2.3.1 Deriving $p(\tilde{\mathbf{y}}_i|\tilde{\mathbf{x}}_i, \Theta)$ and $p(\tilde{\mathbf{x}}_i)$

To derive the equations for  $p(\tilde{\mathbf{y}}_i|\tilde{\mathbf{x}}_i, \Theta)$  and  $p(\tilde{\mathbf{x}}_i)$  we use the above identities and make the following substitutions.

$$\boldsymbol{\mu}_{\mathbf{y}_B} = 0, \quad (16)$$

$$\boldsymbol{\Sigma}_B = \mathbf{I}, \quad (17)$$

$$\mathbf{y}_B = \tilde{\mathbf{y}}_i, \quad (18)$$

$$\mathbf{x}_B = \tilde{\mathbf{x}}_i, \quad (19)$$

$$\mathbf{A}_B = \tilde{\mathbf{A}}, \quad (20)$$

$$\mathbf{b}_B = \tilde{\boldsymbol{\mu}}, \quad (21)$$

$$\boldsymbol{\Sigma}^* = \left( \mathbf{I}^{-1} + \tilde{\mathbf{A}}^T \tilde{\boldsymbol{\Sigma}}^{-1} \tilde{\mathbf{A}} \right)^{-1} = \left( \mathbf{I} + \tilde{\mathbf{A}}^T \tilde{\boldsymbol{\Sigma}}^{-1} \tilde{\mathbf{A}} \right)^{-1}, \quad (22)$$

$$\mathbf{L}_B^{-1} = \tilde{\boldsymbol{\Sigma}}. \quad (23)$$

Therefore, we get the following for  $p(\tilde{\mathbf{x}}_i)$ :

$$p(\tilde{\mathbf{x}}_i) = \mathcal{N}_{\tilde{\mathbf{x}}_i} \left[ \tilde{\mathbf{A}}\mathbf{0} + \tilde{\boldsymbol{\mu}}, \tilde{\boldsymbol{\Sigma}} + \tilde{\mathbf{A}}\mathbf{I}\tilde{\mathbf{A}}^T \right] = \mathcal{N}_{\tilde{\mathbf{x}}_i} \left[ \tilde{\boldsymbol{\mu}}, \tilde{\boldsymbol{\Sigma}} + \tilde{\mathbf{A}}\tilde{\mathbf{A}}^T \right]. \quad (24)$$

And we get the following for  $p(\tilde{\mathbf{y}}_i|\tilde{\mathbf{x}}_i, \Theta)$ :

$$p(\tilde{\mathbf{y}}_i|\tilde{\mathbf{x}}_i, \Theta) = \mathcal{N}_{\tilde{\mathbf{y}}_i} \left[ \boldsymbol{\Sigma}^* \left( \tilde{\mathbf{A}}^T \tilde{\boldsymbol{\Sigma}}^{-1} (\tilde{\mathbf{x}}_i - \tilde{\boldsymbol{\mu}}) + \mathbf{I}^{-1}\mathbf{0} \right), \boldsymbol{\Sigma}^* \right] = \mathcal{N}_{\tilde{\mathbf{y}}_i} \left[ \boldsymbol{\Sigma}^* \tilde{\mathbf{A}}^T \tilde{\boldsymbol{\Sigma}}^{-1} (\tilde{\mathbf{x}}_i - \tilde{\boldsymbol{\mu}}), \boldsymbol{\Sigma}^* \right], \quad (25)$$

$$p(\tilde{\mathbf{y}}_i|\tilde{\mathbf{x}}_i, \Theta) = \mathcal{N}_{\tilde{\mathbf{y}}_i} \left[ \left( \mathbf{I} + \tilde{\mathbf{A}}^T \tilde{\boldsymbol{\Sigma}}^{-1} \tilde{\mathbf{A}} \right)^{-1} \tilde{\mathbf{A}}^T \tilde{\boldsymbol{\Sigma}}^{-1} (\tilde{\mathbf{x}}_i - \tilde{\boldsymbol{\mu}}), \left( \mathbf{I} + \tilde{\mathbf{A}}^T \tilde{\boldsymbol{\Sigma}}^{-1} \tilde{\mathbf{A}} \right)^{-1} \right]. \quad (26)$$

### 2.3.2 Some Sanity Checks on Array and Matrix Sizes

We will make some statements so that if you go through these formulae things will hopefully make a bit more sense to you.

- $D_G$  is associated with the subspace  $\mathbf{G}$  (condition or session subspace) which is of size  $(D_x, D_G)$
- $D_F$  is associated with the subspace  $\mathbf{F}$  (identity subspace) which is of size  $(D_x, D_F)$
- $\tilde{\mathbf{y}}_i$  is the latent variables for the  $J_i$  observations of identity  $i$  and is a vector of size  $(D_F + J_i D_G, 1)$
- $\tilde{\mathbf{A}}$  is a matrix of size  $(J_i D_x, D_F + J_i D_G)$  and its transpose  $\tilde{\mathbf{A}}^T$  is a matrix of size  $(D_F + J_i D_G, J_i D_x)$
- $\tilde{\boldsymbol{\Sigma}}$  is a matrix of size  $(J_i D_x, J_i D_x)$
- $\tilde{\boldsymbol{\mu}}$  is a vector of size  $(J_i D_x, 1)$
- $\tilde{\mathbf{x}}_i$  is a vector of size  $(J_i D_x, 1)$

## 2.4 Sufficient Statistics

In this section we describe how to estimate the sufficient statistics to train the model. These statistics are the first order moment and second order moment of the latent variables. We find these quantities by making use of Equation 26. It will be shown that in both cases they depend upon the quantity  $(\mathbf{I} + \tilde{\mathbf{A}}^T \tilde{\Sigma}^{-1} \tilde{\mathbf{A}})^{-1}$ .

**First order moment of  $\tilde{\mathbf{y}}_i$**  From Equation 26 the expected value of  $\tilde{\mathbf{y}}_i$  is given by,

$$E[\tilde{\mathbf{y}}_i | \tilde{\mathbf{x}}_i, \Theta] = (\mathbf{I} + \tilde{\mathbf{A}}^T \tilde{\Sigma}^{-1} \tilde{\mathbf{A}})^{-1} \tilde{\mathbf{A}}^T \tilde{\Sigma}^{-1} (\tilde{\mathbf{x}}_i - \tilde{\boldsymbol{\mu}}). \quad (27)$$

This is the same solution as provided in Appendix 1, Equation 22 of [4].

**Second order moment of  $\tilde{\mathbf{y}}_i$**  From Equation 26 we know  $\text{Var}[\tilde{\mathbf{y}}_i | \tilde{\mathbf{x}}_i, \Theta]$  and so we have the following,

$$E[\tilde{\mathbf{y}}_i \tilde{\mathbf{y}}_i^T | \tilde{\mathbf{x}}_i, \Theta] = \text{Var}[\tilde{\mathbf{y}}_i | \tilde{\mathbf{x}}_i, \Theta] + E[\tilde{\mathbf{y}}_i | \tilde{\mathbf{x}}_i, \Theta] E[\tilde{\mathbf{y}}_i | \tilde{\mathbf{x}}_i, \Theta]^T. \quad (28)$$

This comes from the fact that  $\text{Var}[y] = E[y^2] - E[y]^2$ , which is the previous equation but reformulated in terms of  $E[y^2]$ . Continuing the derivation, we take  $E[\tilde{\mathbf{y}}_i | \tilde{\mathbf{x}}_i, \Theta]$  from Equation 27 and  $\text{Var}[\tilde{\mathbf{y}}_i | \tilde{\mathbf{x}}_i, \Theta]$  is taken directly from Equation 26. This gives us,

$$E[\tilde{\mathbf{y}}_i \tilde{\mathbf{y}}_i^T | \tilde{\mathbf{x}}_i, \Theta] = (\mathbf{I} + \tilde{\mathbf{A}}^T \tilde{\Sigma}^{-1} \tilde{\mathbf{A}})^{-1} + E[\tilde{\mathbf{y}}_i | \tilde{\mathbf{x}}_i, \Theta] E[\tilde{\mathbf{y}}_i | \tilde{\mathbf{x}}_i, \Theta]^T. \quad (29)$$

This is the same equation as provided in Appendix 1, Equation 23 of [4].

### 2.4.1 Complexity of the Problem

Given the above formulae, the next question is how to go about actually implementing this.

1. Calculating the first order moment,  $E[\tilde{\mathbf{y}}_i | \tilde{\mathbf{x}}_i, \Theta] = (\mathbf{I} + \tilde{\mathbf{A}}^T \tilde{\Sigma}^{-1} \tilde{\mathbf{A}})^{-1} \tilde{\mathbf{A}}^T \tilde{\Sigma}^{-1} (\tilde{\mathbf{x}}_i - \tilde{\boldsymbol{\mu}})$ , implies the calculation of the large matrix  $(\mathbf{I} + \tilde{\mathbf{A}}^T \tilde{\Sigma}^{-1} \tilde{\mathbf{A}})^{-1}$ .
2. Calculating the second order moment,  $E[\tilde{\mathbf{y}}_i \tilde{\mathbf{y}}_i^T | \tilde{\mathbf{x}}_i, \Theta] = (\mathbf{I} + \tilde{\mathbf{A}}^T \tilde{\Sigma}^{-1} \tilde{\mathbf{A}})^{-1} + E[\tilde{\mathbf{y}}_i | \tilde{\mathbf{x}}_i, \Theta] E[\tilde{\mathbf{y}}_i | \tilde{\mathbf{x}}_i, \Theta]^T$ , also implies the calculation of the same matrix.

So the main stumbling block is being able to calculate  $(\mathbf{I} + \tilde{\mathbf{A}}^T \tilde{\Sigma}^{-1} \tilde{\mathbf{A}})^{-1}$ . This matrix depends on the number of samples,  $J_i$ , for a particular client  $i$ . We will describe the derivation later, but below we provide the scalable and exact equations for this problem.

### 2.4.2 Calculating $E[\tilde{\mathbf{y}}_i | \tilde{\mathbf{x}}_i, \Theta]$

To calculate  $E[\tilde{\mathbf{y}}_i | \tilde{\mathbf{x}}_i, \Theta]$  the final scalable solution becomes, for  $J_i = 3$ ,

$$E[\tilde{\mathbf{y}}_i | \tilde{\mathbf{x}}_i, \Theta] = \begin{bmatrix} (\mathbf{I} + J_i \mathbf{F}^T \mathbf{Q} \mathbf{F})^{-1} \sum_{j=1}^{J_i} \mathbf{F}^T \mathbf{Q} \tilde{\mathbf{x}}_{ij} \\ - (\mathbf{I} + \mathbf{G}^T \Sigma^{-1} \mathbf{G})^{-1} \mathbf{G}^T \Sigma^{-1} \mathbf{F} (\mathbf{I} + J_i \mathbf{F}^T \mathbf{Q} \mathbf{F})^{-1} \sum_{j=1}^{J_i} \mathbf{F}^T \mathbf{Q} \tilde{\mathbf{x}}_{ij} + (\mathbf{I} + \mathbf{G}^T \Sigma^{-1} \mathbf{G})^{-1} \mathbf{G}^T \Sigma^{-1} \tilde{\mathbf{x}}_{i1} \\ - (\mathbf{I} + \mathbf{G}^T \Sigma^{-1} \mathbf{G})^{-1} \mathbf{G}^T \Sigma^{-1} \mathbf{F} (\mathbf{I} + J_i \mathbf{F}^T \mathbf{Q} \mathbf{F})^{-1} \sum_{j=1}^{J_i} \mathbf{F}^T \mathbf{Q} \tilde{\mathbf{x}}_{ij} + (\mathbf{I} + \mathbf{G}^T \Sigma^{-1} \mathbf{G})^{-1} \mathbf{G}^T \Sigma^{-1} \tilde{\mathbf{x}}_{i2} \\ - (\mathbf{I} + \mathbf{G}^T \Sigma^{-1} \mathbf{G})^{-1} \mathbf{G}^T \Sigma^{-1} \mathbf{F} (\mathbf{I} + J_i \mathbf{F}^T \mathbf{Q} \mathbf{F})^{-1} \sum_{j=1}^{J_i} \mathbf{F}^T \mathbf{Q} \tilde{\mathbf{x}}_{ij} + (\mathbf{I} + \mathbf{G}^T \Sigma^{-1} \mathbf{G})^{-1} \mathbf{G}^T \Sigma^{-1} \tilde{\mathbf{x}}_{i3} \end{bmatrix}. \quad (30)$$



This gives us that,

$$E[\mathbf{h}_i|\tilde{\mathbf{x}}_i, \Theta] = (\mathbf{I} + J_i \mathbf{F}^T \mathbf{Q} \mathbf{F})^{-1} \sum_{j=1}^{J_i} \mathbf{F}^T \mathbf{Q} \tilde{\mathbf{x}}_{ij}, \quad (31)$$

and

$$E[\mathbf{w}_{ij}|\mathbf{h}_i, \mathbf{x}_{ij}, \Theta] = (\mathbf{I} + \mathbf{G}^T \Sigma^{-1} \mathbf{G})^{-1} \mathbf{G}^T \Sigma^{-1} (\tilde{\mathbf{x}}_{ij} - \mathbf{F} E[\mathbf{h}_i|\tilde{\mathbf{x}}_i, \Theta]). \quad (32)$$

Where  $\mathbf{Q} = \Sigma^{-1} - \Sigma^{-1} \mathbf{G} (\mathbf{I} + \mathbf{G}^T \Sigma^{-1} \mathbf{G})^{-1} \mathbf{G}^T \Sigma^{-1} = (\Sigma + \mathbf{G} \mathbf{G}^T)^{-1}$  and  $\tilde{\mathbf{x}}_{ij} = \mathbf{x}_{ij} - \boldsymbol{\mu}$ . This solution will be discussed later in Section 4. Note that we have included  $\mathbf{h}_i$  on the right hand side of  $E[\mathbf{w}_{ij}|\mathbf{h}_i, \mathbf{x}_{ij}, \Theta]$  as it is assumed to have been pre-calculated, but this is not strictly necessary; this can be calculated each time but in this case we would write  $E[\mathbf{w}_{ij}|\tilde{\mathbf{x}}_i, \Theta]$  and imply a greater overhead.

### 2.4.3 Calculating $E[\tilde{\mathbf{y}}_i \tilde{\mathbf{y}}_i^T | \tilde{\mathbf{x}}_i, \Theta]$

To calculate  $E[\tilde{\mathbf{y}}_i \tilde{\mathbf{y}}_i^T | \tilde{\mathbf{x}}_i, \Theta]$  we need to calculate,

$$E[\tilde{\mathbf{y}}_i \tilde{\mathbf{y}}_i^T | \tilde{\mathbf{x}}_i, \Theta] = (\mathbf{I} + \tilde{\mathbf{A}}^T \tilde{\Sigma}^{-1} \tilde{\mathbf{A}})^{-1} + E[\tilde{\mathbf{y}}_i | \tilde{\mathbf{x}}_i, \Theta] E[\tilde{\mathbf{y}}_i | \tilde{\mathbf{x}}_i, \Theta]^T. \quad (33)$$

However, in practice we do not actually want to do this jointly for all of the  $J_i$  observations. We do need to have  $E[\tilde{\mathbf{y}}_i | \tilde{\mathbf{x}}_i, \Theta]$  estimated jointly but we then consider it one sample at a time and so we go to the notation  $E[\mathbf{y}_{ij} | \tilde{\mathbf{x}}_i, \Theta]$ ; we use the notation with  $\mathbf{y}_{ij}$  to be explicit. We also need an appropriate expression for  $(\mathbf{I} + \tilde{\mathbf{A}}^T \tilde{\Sigma}^{-1} \tilde{\mathbf{A}})^{-1}$ . This leads to the following scalable solution,

$$E[\mathbf{y}_{ij} \mathbf{y}_{ij}^T | \tilde{\mathbf{x}}_i, \Theta] = \begin{bmatrix} \mathbf{T}_{ul} & \mathbf{T}_{ur} \\ \mathbf{T}_{lr} & \mathbf{T}_{lr} \end{bmatrix} + E[\mathbf{y}_{ij} | \tilde{\mathbf{x}}_i, \Theta] E[\mathbf{y}_{ij} | \tilde{\mathbf{x}}_i, \Theta]^T, \quad (34)$$

where,

$$\mathbf{T}_{ul} = (\mathbf{I} + J_i \mathbf{F}^T \mathbf{Q} \mathbf{F})^{-1}, \quad (35)$$

$$\mathbf{T}_{ur} = \mathbf{T}_{ur}^T = -(\mathbf{I} + \mathbf{G}^T \Sigma^{-1} \mathbf{G})^{-1} \mathbf{G}^T \Sigma^{-1} \mathbf{F} (\mathbf{I} + J_i \mathbf{F}^T \mathbf{Q} \mathbf{F})^{-1}, \quad (36)$$

$$\mathbf{T}_{lr} = (\mathbf{I} - \mathbf{T}_{ul} \mathbf{F}^T \Sigma^{-1} \mathbf{G}) (\mathbf{I} + \mathbf{G}^T \Sigma^{-1} \mathbf{G})^{-1}. \quad (37)$$

This solution will be described in more detail in Section 4.

## 2.5 Authentication Scores for PLDA

A verification score for PLDA can be derived by examining Figure 2B in [4]. In this figure there are two models  $M_0$  and  $M_1$ . For  $M_0$  the latent identity variables are considered to be independent while for  $M_1$  the latent identity variables is shared by all of the observations. Given two observations one for enrolment  $\mathbf{x}_i$  of client  $i$  and another for testing  $\mathbf{x}_t$  of an unknown client  $t$ , this can then be expressed as follows.

For  $M_0$  each of the observations,  $\mathbf{x}_i$  and  $\mathbf{x}_t$ , are considered to have independent latent identity variables  $\mathbf{h}_i$  and  $\mathbf{h}_p$ . They are marginalised independently as follows,

$$p(\mathbf{x}_i, \mathbf{x}_p | M_0) = \int \int p(\mathbf{x}_i, \mathbf{h}_i, \mathbf{w}_{i1}) p(\mathbf{x}_t, \mathbf{h}_p, \mathbf{w}_{p1}) d\mathbf{h}_i d\mathbf{w}_{i1} d\mathbf{h}_p d\mathbf{w}_{p1}, \quad (38)$$

$$= \int \int p(\mathbf{x}_i, \mathbf{h}_i, \mathbf{w}_{i1}) d\mathbf{h}_i d\mathbf{w}_{i1} \int \int p(\mathbf{x}_t, \mathbf{h}_p, \mathbf{w}_{p1}) d\mathbf{h}_p d\mathbf{w}_{p1}. \quad (39)$$

For  $M_1$  the two observations are considered to have the same latent identity variable  $\mathbf{h}_i$  and so are jointly marginalised. However, the session variables are still considered independently,

$$p(\mathbf{x}_i, \mathbf{x}_p | M_1) = \int \int p(\mathbf{x}_i, \mathbf{x}_p, \mathbf{h}_i, \mathbf{w}_{i1}, \mathbf{w}_{it}) d\mathbf{h}_i d\mathbf{w}_{i1} d\mathbf{w}_{it}, \quad (40)$$

$$= \int \left[ \int p(\mathbf{x}_i, \mathbf{h}_i, \mathbf{w}_{i1}) d\mathbf{w}_{i1} \int p(\mathbf{x}_p, \mathbf{h}_i, \mathbf{w}_{ip}) d\mathbf{w}_{ip} \right] d\mathbf{h}_i. \quad (41)$$

We can then use these two models to form a discriminant function using the likelihood ratio ( $LR$ ),

$$LR = \frac{p(\mathbf{x}_i, \mathbf{x}_p | M_1)}{p(\mathbf{x}_i, \mathbf{x}_p | M_0)}. \quad (42)$$

The logarithm of this can also be taken to get the log-likelihood ratio ( $LLR$ ),

$$LLR = \ln [p(\mathbf{x}_i, \mathbf{x}_p | M_1)] - \ln [p(\mathbf{x}_i, \mathbf{x}_p | M_0)]. \quad (43)$$

### 3 Learning the Parameters

We now come to the problem of how to learn our parameters  $\Theta = [\boldsymbol{\mu}, \mathbf{F}, \mathbf{G}, \boldsymbol{\Sigma}]$ . For this purpose, an expectation-maximisation (EM) algorithm can be used and such an approach was described in [4] as indicated by Equation 26 in Appendix 1 of their work. We rewrite this equation below,

$$Q(\Theta_t, \Theta_{t-1}) = \sum_{i=1}^I \sum_{j=1}^{J_i} \int p(\mathbf{y}_{ij} | \mathbf{x}_{i1}, \dots, \mathbf{x}_{iJ_i}, \Theta_{t-1}) \ln [p(\mathbf{x}_{ij} | \mathbf{y}_{ij}, \Theta_t) p(\mathbf{y}_{ij})] d\mathbf{y}_{ij}. \quad (44)$$

The first term is obtained when performing the expectation part (E-step) by keeping the parameters fixed at  $t-1$  ( $\Theta_{t-1}$ ), this is solved by calculating the sufficient statistics as described in Section 2.4. The second term,  $\ln [p(\mathbf{x}_{ij} | \mathbf{y}_{ij}, \Theta) p(\mathbf{y}_{ij})]$ , is the maximisation part (M-step) of the EM derivation. We rewrite this equation below,

$$\ln [p(\mathbf{x}_{ij} | \mathbf{y}_{ij}, \Theta) p(\mathbf{y}_{ij})] = \ln [p(\mathbf{x}_{ij} | \mathbf{y}_{ij}, \Theta)] + \ln [p(\mathbf{y}_{ij})]. \quad (45)$$

To maximise Equation 45 we differentiate with respect to  $\Theta$  and set the left hand side to 0. So let us write it fully to understand what we will get for the different parts of  $\Theta$  that we need to maximise,

$$\ln [p(\mathbf{x}_{ij} | \mathbf{y}_{ij}, \Theta) p(\mathbf{y}_{ij})] = \ln [p(\mathbf{x}_{ij} | \mathbf{y}_{ij}, \Theta)] + \ln [p(\mathbf{y}_{ij})]. \quad (46)$$

For the first term we have,

$$\ln [p(\mathbf{x}_{ij} | \mathbf{y}_{ij}, \Theta)] = -\frac{D_x}{2} \ln(2\pi) - \frac{1}{2} \ln(\det(\boldsymbol{\Sigma})) - \frac{1}{2} (\mathbf{x}_{ij} - \boldsymbol{\mu} - \mathbf{A}\mathbf{y}_{ij})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_{ij} - \boldsymbol{\mu} - \mathbf{A}\mathbf{y}_{ij}), \quad (47)$$

$$\mathbf{A} = [\mathbf{F} \quad \mathbf{G}], \quad (48)$$

$$\mathbf{y}_{ij} = [\mathbf{h}_i^T, \mathbf{w}_{ij}^T]^T. \quad (49)$$

For the second term we have,

$$\ln [p(\mathbf{y}_{ij})] = \ln [\mathcal{N}_{\mathbf{y}}[\mathbf{0} | \mathbf{I}]], \quad (50)$$

$$= \ln \left[ (2\pi)^{-\frac{D_F + D_G}{2}} \det(\mathbf{I})^{-\frac{1}{2}} \exp \left( -\frac{1}{2} \mathbf{y}_{ij}^T \mathbf{I} \mathbf{y}_{ij} \right) \right], \quad (51)$$

$$= -\frac{D_F + D_G}{2} \ln(2\pi) - \frac{1}{2} \ln(\det(\mathbf{I})) - \frac{1}{2} \mathbf{y}_{ij}^T \mathbf{y}_{ij}. \quad (52)$$

From the above equations we can see that  $\ln[p(\mathbf{y}_{ij})]$  does not depend on the parameters  $\Theta$ , and so we can ignore this part in the derivation. Therefore, we concentrate on differentiating Equation 47 with respect to the parameters  $\mu$ ,  $\Sigma$  and  $\mathbf{A}$  (which includes  $\mathbf{F}$  and  $\mathbf{G}$ ). Using the definitions from Appendix A we go onto to differentiate Equation 47. We take up the derivations for each parameter below.

### 3.1 Differentiating with Respect to $\mu$

We need to differentiate Equation 47 with respect to  $\mu$ . We begin this differentiation below:

$$\begin{aligned} \frac{\partial}{\partial \mu} \ln[p(\mathbf{x}_{ij}|\mathbf{y}_{ij}, \Theta)] &= \\ \frac{\partial}{\partial \mu} \left( -\frac{D_x}{2} \ln(2\pi) - \frac{1}{2} \ln(\det(\Sigma)) - \frac{1}{2} (\mathbf{x}_{ij} - \mu - \mathbf{A}\mathbf{y}_{ij})^T \Sigma^{-1} (\mathbf{x}_{ij} - \mu - \mathbf{A}\mathbf{y}_{ij}) \right), \end{aligned} \quad (53)$$

$$= -0 - 0 - \frac{\partial}{\partial \mu} \left[ \frac{1}{2} (\mathbf{x}_{ij} - \mu - \mathbf{A}\mathbf{y}_{ij})^T \Sigma^{-1} (\mathbf{x}_{ij} - \mu - \mathbf{A}\mathbf{y}_{ij}) \right]. \quad (54)$$

We make the substitution  $\mathbf{a}_{ij} = (\mathbf{x}_{ij} - \mu - \mathbf{A}\mathbf{y}_{ij})$  and  $\frac{\partial \mathbf{a}_{ij}}{\partial \mu} = -1$  so that we can rewrite the equation above as,

$$= \frac{\partial}{\partial \mathbf{a}_{ij}} \frac{\partial \mathbf{a}_{ij}}{\partial \mu} \left( -\frac{1}{2} \mathbf{a}_{ij}^T \Sigma^{-1} \mathbf{a}_{ij} \right) = \frac{\partial}{\partial \mathbf{a}_{ij}} (-1) \left( -\frac{1}{2} \mathbf{a}_{ij}^T \Sigma^{-1} \mathbf{a}_{ij} \right). \quad (55)$$

We then apply the identity in Equation 193 to get,

$$\frac{\partial}{\partial \mu} \ln[p(\mathbf{x}_{ij}|\mathbf{y}_{ij}, \Theta)] = (-1) \left( -\frac{1}{2} \right) \left( \Sigma^{-1} + (\Sigma^{-1})^T \right) \mathbf{a}_{ij}, \quad (56)$$

$$= \left( \frac{1}{2} \right) (2\Sigma^{-1}) (\mathbf{x}_{ij} - \mu - \mathbf{A}\mathbf{y}_{ij}), \quad (57)$$

$$= \Sigma^{-1} (\mathbf{x}_{ij} - \mu - \mathbf{A}\mathbf{y}_{ij}). \quad (58)$$

Above we have used the fact that  $\Sigma^{-1} = (\Sigma^{-1})^T$ . We can now maximise this by setting the left hand side (LHS) to zero and reintroducing the sums to get the following,

$$0 = \sum_{i=1}^I \sum_{j=1}^{J_i} \Sigma^{-1} (\mathbf{x}_{ij} - \mu - \mathbf{A}\mathbf{y}_{ij}), \quad (59)$$

$$0 = \sum_{i=1}^I \sum_{j=1}^{J_i} \mathbf{x}_{ij} - \sum_{i=1}^I \sum_{j=1}^{J_i} \mu - \sum_{i=1}^I \sum_{j=1}^{J_i} \mathbf{A}\mathbf{y}_{ij}, \quad (60)$$

$$\sum_{i=1}^I \sum_{j=1}^{J_i} \mu = \sum_{i=1}^I \sum_{j=1}^{J_i} (\mathbf{x}_{ij} - \mathbf{A}\mathbf{y}_{ij}), \quad (61)$$

$$\mu = \frac{1}{I J_i} \sum_{i=1}^I \sum_{j=1}^{J_i} (\mathbf{x}_{ij} - \mathbf{A}\mathbf{y}_{ij}). \quad (62)$$

The result for  $\mu$ , Equation 62, does not match the one presented in [4] as they use a simplification that  $\mu$  is the mean of the data itself; as such it is not actually updated. The most likely explanation for this, which was never given, is that because the latent variables should be

zero mean then  $\sum_{i=1}^I \sum_{j=1}^{J_i} \mathbf{A} \mathbf{y}_{ij}$  should be 0. If we make this assumption then we get the same equation which is,

$$\boldsymbol{\mu} = \frac{1}{IJ_i} \sum_{i=1}^I \sum_{j=1}^{J_i} \mathbf{x}_{ij}. \quad (63)$$

### 3.2 Differentiating with Respect to $\mathbf{A}$

We need to differentiate Equation 47 with respect to  $\mathbf{A}$ . We begin this differentiation below:

$$\begin{aligned} \frac{\partial}{\partial \mathbf{A}} \ln [p(\mathbf{x}_{ij} | \mathbf{y}_{ij}, \boldsymbol{\Theta})] = \\ \frac{\partial}{\partial \mathbf{A}} \left( -\frac{D_x}{2} \ln(2\pi) - \frac{1}{2} \ln(\det(\boldsymbol{\Sigma})) - \frac{1}{2} (\mathbf{x}_{ij} - \boldsymbol{\mu} - \mathbf{A} \mathbf{y}_{ij})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_{ij} - \boldsymbol{\mu} - \mathbf{A} \mathbf{y}_{ij}) \right), \end{aligned} \quad (64)$$

$$= -0 - 0 - \frac{\partial}{\partial \mathbf{A}} \left[ \frac{1}{2} (\mathbf{x}_{ij} - \boldsymbol{\mu} - \mathbf{A} \mathbf{y}_{ij})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_{ij} - \boldsymbol{\mu} - \mathbf{A} \mathbf{y}_{ij}) \right]. \quad (65)$$

We now substitute  $\boldsymbol{\alpha}_{ij} = (\mathbf{x}_{ij} - \boldsymbol{\mu})$  into the above equation so that we can rewrite it as,

$$= -\frac{\partial}{\partial \mathbf{A}} \left[ \frac{1}{2} (\boldsymbol{\alpha}_{ij} - \mathbf{A} \mathbf{y}_{ij})^T \boldsymbol{\Sigma}^{-1} (\boldsymbol{\alpha}_{ij} - \mathbf{A} \mathbf{y}_{ij}) \right]. \quad (66)$$

This takes a form similar to that of the identity of Equation 198 and so we use this identity to get the following:

$$\frac{\partial}{\partial \mathbf{A}} \ln [p(\mathbf{x}_{ij} | \mathbf{y}_{ij}, \boldsymbol{\Theta})] = (-2) \left( -\frac{1}{2} \right) \boldsymbol{\Sigma}^{-1} (\boldsymbol{\alpha}_{ij} - \mathbf{A} \mathbf{y}_{ij}) \mathbf{y}_{ij}^T. \quad (67)$$

To maximise with respect to  $\mathbf{A}$ , we set the LHS of this equation to zero and reintroduce the sums. This gives us the following result,

$$0 = \sum_{i=1}^I \sum_{j=1}^{J_i} \boldsymbol{\Sigma}^{-1} (\boldsymbol{\alpha}_{ij} - \mathbf{A} \mathbf{y}_{ij}) \mathbf{y}_{ij}^T, \quad (68)$$

$$0 = \boldsymbol{\Sigma}^{-1} \left[ \sum_{i=1}^I \sum_{j=1}^{J_i} (\mathbf{x}_{ij} - \boldsymbol{\mu}) \mathbf{y}_{ij}^T - \sum_{i=1}^I \sum_{j=1}^{J_i} \mathbf{A} \mathbf{y}_{ij} \mathbf{y}_{ij}^T \right], \quad (69)$$

$$\mathbf{A} \sum_{i=1}^I \sum_{j=1}^{J_i} \mathbf{y}_{ij} \mathbf{y}_{ij}^T = \sum_{i=1}^I \sum_{j=1}^{J_i} (\mathbf{x}_{ij} - \boldsymbol{\mu}) \mathbf{y}_{ij}^T, \quad (70)$$

$$\mathbf{A} = \sum_{i=1}^I \sum_{j=1}^{J_i} (\mathbf{x}_{ij} - \boldsymbol{\mu}) E[\mathbf{y}_{ij}^T] \left( \sum_{i=1}^I \sum_{j=1}^{J_i} E[\mathbf{y}_{ij} \mathbf{y}_{ij}^T] \right)^{-1}. \quad (71)$$

This is the same as the update rule provided in [4].

### 3.3 Differentiating with Respect to $\boldsymbol{\Sigma}$

We need to differentiate Equation 47 with respect to  $\boldsymbol{\Sigma}$ . We begin this differentiation below:

$$\begin{aligned} \frac{\partial}{\partial \boldsymbol{\Sigma}} \ln [p(\mathbf{x}_{ij} | \mathbf{y}_{ij}, \boldsymbol{\Theta})] = \\ \frac{\partial}{\partial \boldsymbol{\Sigma}} \left[ -\frac{D_x}{2} \ln(2\pi) - \frac{1}{2} \ln(\det(\boldsymbol{\Sigma})) - \frac{1}{2} (\mathbf{x}_{ij} - \boldsymbol{\mu} - \mathbf{A} \mathbf{y}_{ij})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_{ij} - \boldsymbol{\mu} - \mathbf{A} \mathbf{y}_{ij}) \right], \end{aligned} \quad (72)$$

We substitute  $\mathbf{a}_{ij} = (\mathbf{x}_{ij} - \boldsymbol{\mu} - \mathbf{A}\mathbf{y}_{ij})$  so that we can rewrite the equation above as,

$$\frac{\partial}{\partial \boldsymbol{\Sigma}} \ln [p(\mathbf{x}_{ij} | \mathbf{y}_{ij}, \boldsymbol{\Theta})] = \frac{\partial}{\partial \boldsymbol{\Sigma}} \left[ -\frac{D_x}{2} \ln(2\pi) - \frac{1}{2} \ln(\det(\boldsymbol{\Sigma})) - \frac{1}{2} \mathbf{a}_{ij}^T \boldsymbol{\Sigma}^{-1} \mathbf{a}_{ij} \right]. \quad (73)$$

We can then use the identities in Equations 194 and 195. Using these identities we get the following:

$$\frac{\partial}{\partial \boldsymbol{\Sigma}} \ln [p(\mathbf{x}_{ij} | \mathbf{y}_{ij}, \boldsymbol{\Theta})] = \left( -0 - \frac{1}{2} \boldsymbol{\Sigma}^{-T} + \frac{1}{2} \boldsymbol{\Sigma}^{-T} \mathbf{a}_{ij} \mathbf{a}_{ij}^T \boldsymbol{\Sigma}^{-1} \right), \quad (74)$$

$$= -\frac{1}{2} \boldsymbol{\Sigma}^{-T} + \frac{1}{2} \boldsymbol{\Sigma}^{-T} (\mathbf{x}_{ij} - \boldsymbol{\mu} - \mathbf{A}\mathbf{y}_{ij}) (\mathbf{x}_{ij} - \boldsymbol{\mu} - \mathbf{A}\mathbf{y}_{ij})^T \boldsymbol{\Sigma}^{-1}. \quad (75)$$

We need to take the above result, set the LHS to zero and reintroduce the sum. This gives us the following result (remembering that we want to find an expression for  $\boldsymbol{\Sigma}$ ).

$$0 = -\frac{1}{2} \sum_{i=1}^I \sum_{j=1}^{J_i} \boldsymbol{\Sigma}^{-T} + \frac{1}{2} \sum_{i=1}^I \sum_{j=1}^{J_i} \boldsymbol{\Sigma}^{-T} (\mathbf{x}_{ij} - \boldsymbol{\mu} - \mathbf{A}\mathbf{y}_{ij}) (\mathbf{x}_{ij} - \boldsymbol{\mu} - \mathbf{A}\mathbf{y}_{ij})^T \boldsymbol{\Sigma}^{-1}, \quad (76)$$

$$\boldsymbol{\Sigma}^{-1} \sum_{i=1}^I \sum_{j=1}^{J_i} 1 = \boldsymbol{\Sigma}^{-1} \sum_{i=1}^I \sum_{j=1}^{J_i} (\mathbf{x}_{ij} - \boldsymbol{\mu} - \mathbf{A}\mathbf{y}_{ij}) (\mathbf{x}_{ij} - \boldsymbol{\mu} - \mathbf{A}\mathbf{y}_{ij})^T \boldsymbol{\Sigma}^{-1}, \quad (77)$$

$$\boldsymbol{\Sigma} \boldsymbol{\Sigma}^{-1} \sum_{i=1}^I \sum_{j=1}^{J_i} 1 = \boldsymbol{\Sigma} \boldsymbol{\Sigma}^{-1} \sum_{i=1}^I \sum_{j=1}^{J_i} (\mathbf{x}_{ij} - \boldsymbol{\mu} - \mathbf{A}\mathbf{y}_{ij}) (\mathbf{x}_{ij} - \boldsymbol{\mu} - \mathbf{A}\mathbf{y}_{ij})^T \boldsymbol{\Sigma}^{-1}, \quad (78)$$

$$I J_i = \sum_{i=1}^I \sum_{j=1}^{J_i} (\mathbf{x}_{ij} - \boldsymbol{\mu} - \mathbf{A}\mathbf{y}_{ij}) (\mathbf{x}_{ij} - \boldsymbol{\mu} - \mathbf{A}\mathbf{y}_{ij})^T \boldsymbol{\Sigma}^{-1}, \quad (79)$$

$$I J_i \boldsymbol{\Sigma} = \sum_{i=1}^I \sum_{j=1}^{J_i} (\mathbf{x}_{ij} - \boldsymbol{\mu} - \mathbf{A}\mathbf{y}_{ij}) (\mathbf{x}_{ij} - \boldsymbol{\mu} - \mathbf{A}\mathbf{y}_{ij})^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\Sigma}, \quad (80)$$

As we want  $\boldsymbol{\Sigma}$  to be diagonal, we finally get

$$\boldsymbol{\Sigma} = \frac{1}{I J_i} \sum_{i=1}^I \sum_{j=1}^{J_i} \text{diag} \left[ (\mathbf{x}_{ij} - \boldsymbol{\mu} - \mathbf{A}\mathbf{y}_{ij}) (\mathbf{x}_{ij} - \boldsymbol{\mu} - \mathbf{A}\mathbf{y}_{ij})^T \right]. \quad (81)$$

Everything else that follows is trying to get this expression to match up with what is presented in [4]. It should also provide a more stable estimate of  $\boldsymbol{\Sigma}$ .

Let's substitute  $\mathbf{b}_{ij} = \mathbf{x}_{ij} - \boldsymbol{\mu}$ ,

$$\boldsymbol{\Sigma} = \frac{1}{I J_i} \sum_{i=1}^I \sum_{j=1}^{J_i} \text{diag} \left[ (\mathbf{b}_{ij} - \mathbf{A}\mathbf{y}_{ij}) (\mathbf{b}_{ij} - \mathbf{A}\mathbf{y}_{ij})^T \right], \quad (82)$$

$$\boldsymbol{\Sigma} = \frac{1}{I J_i} \sum_{i=1}^I \sum_{j=1}^{J_i} \text{diag} \left[ \mathbf{b}_{ij} \mathbf{b}_{ij}^T - \mathbf{b}_{ij} \mathbf{y}_{ij}^T \mathbf{A}^T - \mathbf{A} \mathbf{y}_{ij} \mathbf{b}_{ij}^T + \mathbf{A} \mathbf{y}_{ij} \mathbf{y}_{ij}^T \mathbf{A}^T \right], \quad (83)$$

$$\boldsymbol{\Sigma} = \frac{1}{I J_i} \sum_{i=1}^I \sum_{j=1}^{J_i} \text{diag} \left[ \mathbf{b}_{ij} \mathbf{b}_{ij}^T - 2 \mathbf{A} \mathbf{y}_{ij} \mathbf{b}_{ij}^T + \mathbf{A} \mathbf{y}_{ij} \mathbf{y}_{ij}^T \mathbf{A}^T \right]. \quad (84)$$

We now make use of our previous definition of  $\mathbf{A}$ ,

$$\mathbf{A} = \sum_{i=1}^I \sum_{j=1}^{J_i} (\mathbf{x}_{ij} - \boldsymbol{\mu}) E[\mathbf{y}_{ij}]^T \left( \sum_{i=1}^I \sum_{j=1}^{J_i} E[\mathbf{y}_{ij} \mathbf{y}_{ij}^T] \right)^{-1}. \quad (85)$$

Taking this and substituting in,

$$\boldsymbol{\Sigma} = \frac{1}{IJ_i} \sum_{i=1}^I \sum_{j=1}^{J_i} \text{diag} [\mathbf{b}_{ij} \mathbf{b}_{ij}^T - 2\mathbf{A} \mathbf{y}_{ij} \mathbf{b}_{ij}^T + \mathbf{A} \mathbf{y}_{ij} \mathbf{y}_{ij}^T \mathbf{A}^T]. \quad (86)$$

$$= \frac{1}{IJ_i} \text{diag} \left[ \sum_{i=1}^I \sum_{j=1}^{J_i} \mathbf{b}_{ij} \mathbf{b}_{ij}^T - \sum_{i=1}^I \sum_{j=1}^{J_i} 2\mathbf{A} \mathbf{y}_{ij} \mathbf{b}_{ij}^T + \sum_{i=1}^I \sum_{j=1}^{J_i} \mathbf{A} \mathbf{y}_{ij} \mathbf{y}_{ij}^T \mathbf{A}^T \right], \quad (87)$$

As we will be taking the diagonal of  $\boldsymbol{\Sigma}$ , this can be further written as,

$$= \frac{1}{IJ_i} \text{diag} \left[ \sum_{i=1}^I \sum_{j=1}^{J_i} \mathbf{b}_{ij} \mathbf{b}_{ij}^T - 2\mathbf{A} \sum_{i=1}^I \sum_{j=1}^{J_i} \mathbf{y}_{ij} \mathbf{b}_{ij}^T + \mathbf{A} \sum_{i=1}^I \sum_{j=1}^{J_i} \mathbf{y}_{ij} \mathbf{y}_{ij}^T \mathbf{A}^T \right]. \quad (88)$$

This is the point at which we put in our definition for  $\mathbf{A}$ , for the last term only,

$$\boldsymbol{\Sigma} = \frac{1}{IJ_i} \text{diag} \left[ \dots + \left( \sum_{i=1}^I \sum_{j=1}^{J_i} \mathbf{y}_{ij} \mathbf{b}_{ij}^T \right) \left( \sum_{i=1}^I \sum_{j=1}^{J_i} \mathbf{y}_{ij} \mathbf{y}_{ij}^T \right)^{-1} \sum_{i=1}^I \sum_{j=1}^{J_i} \mathbf{y}_{ij} \mathbf{y}_{ij}^T \mathbf{A}^T \right], \quad (89)$$

$$= \frac{1}{IJ_i} \text{diag} \left[ \dots + \left( \sum_{i=1}^I \sum_{j=1}^{J_i} \mathbf{y}_{ij} \mathbf{b}_{ij}^T \right) \left( \sum_{i=1}^I \sum_{j=1}^{J_i} \mathbf{y}_{ij} \mathbf{y}_{ij}^T \right)^{-1} \left( \sum_{i=1}^I \sum_{j=1}^{J_i} \mathbf{y}_{ij} \mathbf{y}_{ij}^T \right) \mathbf{A}^T \right], \quad (90)$$

$$= \frac{1}{IJ_i} \text{diag} \left[ \dots + \left( \sum_{i=1}^I \sum_{j=1}^{J_i} \mathbf{y}_{ij} \mathbf{b}_{ij}^T \right) \mathbf{I} \mathbf{A}^T \right], \quad (91)$$

$$= \frac{1}{IJ_i} \text{diag} \left[ \sum_{i=1}^I \sum_{j=1}^{J_i} \mathbf{b}_{ij} \mathbf{b}_{ij}^T - 2\mathbf{A} \sum_{i=1}^I \sum_{j=1}^{J_i} \mathbf{y}_{ij} \mathbf{b}_{ij}^T + \left( \sum_{i=1}^I \sum_{j=1}^{J_i} \mathbf{y}_{ij} \mathbf{b}_{ij}^T \right) \mathbf{I} \mathbf{A}^T \right], \quad (92)$$

$$= \frac{1}{IJ_i} \text{diag} \left[ \sum_{i=1}^I \sum_{j=1}^{J_i} \mathbf{b}_{ij} \mathbf{b}_{ij}^T - 2\mathbf{A} \sum_{i=1}^I \sum_{j=1}^{J_i} \mathbf{y}_{ij} \mathbf{b}_{ij}^T + \mathbf{A} \sum_{i=1}^I \sum_{j=1}^{J_i} \mathbf{y}_{ij} \mathbf{b}_{ij}^T \right], \quad (93)$$

$$= \frac{1}{IJ_i} \text{diag} \left[ \sum_{i=1}^I \sum_{j=1}^{J_i} (\mathbf{x}_{ij} - \boldsymbol{\mu}) (\mathbf{x}_{ij} - \boldsymbol{\mu})^T - \mathbf{A} \sum_{i=1}^I \sum_{j=1}^{J_i} \mathbf{y}_{ij} (\mathbf{x}_{ij} - \boldsymbol{\mu})^T \right]. \quad (94)$$

This provides us with the same solution as Prince and Elder. We can now write,

$$\boldsymbol{\Sigma} = \frac{1}{IJ_i} \sum_{i=1}^I \sum_{j=1}^{J_i} \text{diag} \left[ (\mathbf{x}_{ij} - \boldsymbol{\mu}) (\mathbf{x}_{ij} - \boldsymbol{\mu})^T - \mathbf{A} \mathbf{y}_{ij} (\mathbf{x}_{ij} - \boldsymbol{\mu})^T \right]. \quad (95)$$

## 4 Implementation Details

The following is a quick summary of some of the implementation details for PLDA. This is not a guide on how to implement PLDA, although it could help considerably.

The first point to note is that there are values that we will use throughout that are probably worth storing somewhere. A non-exhaustive list is:

- $\ln(2\pi)$
- $-\frac{D_{\mathbf{x}}}{2} \ln(2\pi)$
- $\ln[\det(\boldsymbol{\Sigma})]$

- $\Sigma^{-1}$  because we assume this to be diagonal this is actually a vector inversion.
- $\mathbf{Q} = \Sigma^{-1} - \Sigma^{-1}\mathbf{G}(\mathbf{I} + \mathbf{G}^T\Sigma^{-1}\mathbf{G})^{-1}\mathbf{G}^T\Sigma^{-1}$  because it is used for estimating the latent variables and scoring.
- $\mathbf{F}^T\Sigma^{-1}\mathbf{G}$  and  $\mathbf{F}^T\mathbf{Q}$
- $(\mathbf{I} + J_i\mathbf{F}^T\mathbf{Q}\mathbf{F})^{-1}$  because this is used for estimating the latent variables and for scoring.

We now take up scalable and exact solutions for the problems of evaluating the likelihood and estimating the sufficient statistics.

#### 4.1 Log-likelihood of $\tilde{\mathbf{x}}$

This involves calculating the likelihood of the following,

$$p(\tilde{\mathbf{x}}_i) = \mathcal{N}_{\tilde{\mathbf{x}}}[\tilde{\boldsymbol{\mu}}, \tilde{\boldsymbol{\Sigma}} + \tilde{\mathbf{A}}\tilde{\mathbf{A}}^T]. \quad (96)$$

We can simplify this by calculating the log-likelihood which gives us,

$$\ln[p(\tilde{\mathbf{x}}_i)] = -\frac{J_i D_x}{2} \ln(2\pi) - \frac{1}{2} \ln[\det(\tilde{\boldsymbol{\Sigma}} + \tilde{\mathbf{A}}\tilde{\mathbf{A}}^T)] - \frac{1}{2} (\tilde{\mathbf{x}}_i - \tilde{\boldsymbol{\mu}})^T (\tilde{\boldsymbol{\Sigma}} + \tilde{\mathbf{A}}\tilde{\mathbf{A}}^T)^{-1} (\tilde{\mathbf{x}}_i - \tilde{\boldsymbol{\mu}}).$$

For each of the three parts of this calculation we note the following:

1.  $-\frac{J_i D_x}{2} \ln(2\pi)$  is trivial to calculate as it depends only on the dimensionality of the feature space ( $D_x$ ) and the number of samples ( $J_i$ ),
2.  $-\frac{1}{2} \ln[\det(\tilde{\boldsymbol{\Sigma}} + \tilde{\mathbf{A}}\tilde{\mathbf{A}}^T)]$  will be difficult to calculate as it involves calculating the determinant of potentially very large matrices, therefore, simplifications need to be used, and
3.  $-\frac{1}{2} (\tilde{\mathbf{x}}_i - \tilde{\boldsymbol{\mu}})^T (\tilde{\boldsymbol{\Sigma}} + \tilde{\mathbf{A}}\tilde{\mathbf{A}}^T)^{-1} (\tilde{\mathbf{x}}_i - \tilde{\boldsymbol{\mu}})$  will be difficult to calculate as it involves multiplying a very large matrix, therefore, simplifications need to be used.

Because Point 1 can be relatively easily dealt with we only address Points 2 and 3 in the following sections.

##### 4.1.1 Calculating $-\frac{1}{2} (\tilde{\mathbf{x}}_i - \tilde{\boldsymbol{\mu}})^T (\tilde{\boldsymbol{\Sigma}} + \tilde{\mathbf{A}}\tilde{\mathbf{A}}^T)^{-1} (\tilde{\mathbf{x}}_i - \tilde{\boldsymbol{\mu}})$ Efficiently: Method 1

We will examine what the matrix  $(\tilde{\boldsymbol{\Sigma}} + \tilde{\mathbf{A}}\tilde{\mathbf{A}}^T)^{-1}$  is. First we derive fully what  $(\tilde{\boldsymbol{\Sigma}} + \tilde{\mathbf{A}}\tilde{\mathbf{A}}^T)$  is,

$$\tilde{\mathbf{A}} = \begin{bmatrix} \mathbf{F} & \mathbf{G} & \mathbf{0} & \mathbf{0} \\ \mathbf{F} & \mathbf{0} & \mathbf{G} & \mathbf{0} \\ \mathbf{F} & \mathbf{0} & \mathbf{0} & \mathbf{G} \end{bmatrix}, \quad (97)$$

$$\tilde{\mathbf{A}}\tilde{\mathbf{A}}^T = \begin{bmatrix} \mathbf{F} & \mathbf{G} & \mathbf{0} & \mathbf{0} \\ \mathbf{F} & \mathbf{0} & \mathbf{G} & \mathbf{0} \\ \mathbf{F} & \mathbf{0} & \mathbf{0} & \mathbf{G} \end{bmatrix} \begin{bmatrix} \mathbf{F}^T & \mathbf{F}^T & \mathbf{F}^T \\ \mathbf{G}^T & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{G}^T & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{G}^T \end{bmatrix}, \quad (98)$$

$$\tilde{\mathbf{A}}\tilde{\mathbf{A}}^T = \begin{bmatrix} \mathbf{F}\mathbf{F}^T + \mathbf{G}\mathbf{G}^T, & \mathbf{F}\mathbf{F}^T, & \mathbf{F}\mathbf{F}^T \\ \mathbf{F}\mathbf{F}^T, & \mathbf{F}\mathbf{F}^T + \mathbf{G}\mathbf{G}^T, & \mathbf{F}\mathbf{F}^T \\ \mathbf{F}\mathbf{F}^T, & \mathbf{F}\mathbf{F}^T, & \mathbf{F}\mathbf{F}^T + \mathbf{G}\mathbf{G}^T \end{bmatrix}, \quad (99)$$

$$\tilde{\boldsymbol{\Sigma}} + \tilde{\mathbf{A}}\tilde{\mathbf{A}}^T = \begin{bmatrix} \tilde{\boldsymbol{\Sigma}} + \mathbf{F}\mathbf{F}^T + \mathbf{G}\mathbf{G}^T, & \mathbf{F}\mathbf{F}^T, & \mathbf{F}\mathbf{F}^T \\ \mathbf{F}\mathbf{F}^T, & \tilde{\boldsymbol{\Sigma}} + \mathbf{F}\mathbf{F}^T + \mathbf{G}\mathbf{G}^T, & \mathbf{F}\mathbf{F}^T \\ \mathbf{F}\mathbf{F}^T, & \mathbf{F}\mathbf{F}^T, & \tilde{\boldsymbol{\Sigma}} + \mathbf{F}\mathbf{F}^T + \mathbf{G}\mathbf{G}^T \end{bmatrix}. \quad (100)$$

This matrix is actually very big ( $J_i D_x, J_i D_x$ ) and it would need to be inverted so calculating this efficiently is necessary if we want to implement scoring this way. If we take the Woodbury matrix identity,

$$(A + UCV)^{-1} = A^{-1} - A^{-1}U(C^{-1} + VA^{-1}U)^{-1}VA^{-1}, \quad (101)$$

and apply it in a similar way to [2] page 15 (Section 3.2.3) we would get the following,

$$\mathbf{\Lambda} = \left( \tilde{\Sigma} + \tilde{\mathbf{A}}\tilde{\mathbf{A}}^T \right)^{-1} = \tilde{\Sigma}^{-1} - \tilde{\Sigma}^{-1}\tilde{\mathbf{A}} \left( \mathbf{I} + \tilde{\mathbf{A}}^T\tilde{\Sigma}^{-1}\tilde{\mathbf{A}} \right)^{-1} \tilde{\mathbf{A}}^T\tilde{\Sigma}^{-1}, \quad (102)$$

$$\mathbf{\Lambda} = \tilde{\Sigma}^{-1} - \mathbf{\Gamma}\mathbf{\Gamma}^T. \quad (103)$$

Where,  $\mathbf{\Gamma} = \tilde{\Sigma}^{-1}\tilde{\mathbf{A}} \left( \mathbf{I} + \tilde{\mathbf{A}}^T\tilde{\Sigma}^{-1}\tilde{\mathbf{A}} \right)^{-\frac{1}{2}}$ .

Using the above we can now rewrite the quadratic term in the Gaussian that we need to evaluate as being,

$$(\tilde{\mathbf{x}}_i - \tilde{\boldsymbol{\mu}})^T \tilde{\Sigma}^{-1} (\tilde{\mathbf{x}}_i - \tilde{\boldsymbol{\mu}}) - (\tilde{\mathbf{x}}_i - \tilde{\boldsymbol{\mu}})^T \mathbf{\Gamma}\mathbf{\Gamma}^T (\tilde{\mathbf{x}}_i - \tilde{\boldsymbol{\mu}}). \quad (104)$$

The first part,  $(\tilde{\mathbf{x}}_i - \tilde{\boldsymbol{\mu}})^T \tilde{\Sigma}^{-1} (\tilde{\mathbf{x}}_i - \tilde{\boldsymbol{\mu}})$ , is efficient to compute as  $\tilde{\Sigma}^{-1}$  is diagonal. The second part,  $(\tilde{\mathbf{x}}_i - \tilde{\boldsymbol{\mu}})^T \mathbf{\Gamma}\mathbf{\Gamma}^T (\tilde{\mathbf{x}}_i - \tilde{\boldsymbol{\mu}})$ , is more complicated. But we can note that  $(\tilde{\mathbf{x}}_i - \tilde{\boldsymbol{\mu}})^T \mathbf{\Gamma} = \left( \mathbf{\Gamma}^T (\tilde{\mathbf{x}}_i - \tilde{\boldsymbol{\mu}}) \right)^T$  and so we only ever need to evaluate one half of the equation. And that in the end  $\mathbf{\Gamma}^T (\tilde{\mathbf{x}}_i - \tilde{\boldsymbol{\mu}})$  is a vector of dimensions  $(D_F + J_i D_G, 1)$ . Therefore, the final result of the second part is just the magnitude of this vector. This comes from [2].

The above solution is still problematic. This is because  $\tilde{\Sigma}^{-1}\tilde{\mathbf{A}}$  is a matrix of size  $(J_i D_x, D_F + J_i D_G)$  and is quadratic with the number of samples  $J_i$ . So let us expand this out and see if we can reduce the load on the memory somehow. We want to calculate  $\mathbf{\Gamma}^T (\tilde{\mathbf{x}}_i - \tilde{\boldsymbol{\mu}})$  which is given by,

$$\mathbf{\Gamma}^T (\tilde{\mathbf{x}}_i - \tilde{\boldsymbol{\mu}}) = \left[ \tilde{\Sigma}^{-1}\tilde{\mathbf{A}} \left( \mathbf{I} + \tilde{\mathbf{A}}^T\tilde{\Sigma}^{-1}\tilde{\mathbf{A}} \right)^{-\frac{1}{2}} \right]^T \begin{bmatrix} \bar{\mathbf{x}}_{i1} \\ \bar{\mathbf{x}}_{i2} \\ \bar{\mathbf{x}}_{i3} \end{bmatrix}, \quad (105)$$

$$= \left[ \left( \mathbf{I} + \tilde{\mathbf{A}}^T\tilde{\Sigma}^{-1}\tilde{\mathbf{A}} \right)^{-\frac{1}{2}} \right]^T \tilde{\mathbf{A}}^T\tilde{\Sigma}^{-1} \begin{bmatrix} \bar{\mathbf{x}}_{i1} \\ \bar{\mathbf{x}}_{i2} \\ \bar{\mathbf{x}}_{i3} \end{bmatrix}, \quad (106)$$

$$= \tilde{\mathbf{Z}}\tilde{\mathbf{A}}^T\tilde{\Sigma}^{-1} \begin{bmatrix} \bar{\mathbf{x}}_{i1} \\ \bar{\mathbf{x}}_{i2} \\ \bar{\mathbf{x}}_{i3} \end{bmatrix}, \quad (107)$$

$$= \tilde{\mathbf{Z}} \begin{bmatrix} \mathbf{F}^T\tilde{\Sigma}^{-1}, & \mathbf{F}^T\tilde{\Sigma}^{-1}, & \mathbf{F}^T\tilde{\Sigma}^{-1} \\ \mathbf{G}^T\tilde{\Sigma}^{-1} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{G}^T\tilde{\Sigma}^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{G}^T\tilde{\Sigma}^{-1} \end{bmatrix} \begin{bmatrix} \bar{\mathbf{x}}_{i1} \\ \bar{\mathbf{x}}_{i2} \\ \bar{\mathbf{x}}_{i3} \end{bmatrix}, \quad (108)$$

$$= \tilde{\mathbf{Z}} \begin{bmatrix} \sum_{j=1}^{J_i} \mathbf{F}^T\tilde{\Sigma}^{-1}\bar{\mathbf{x}}_{ij} \\ \mathbf{G}^T\tilde{\Sigma}^{-1}\bar{\mathbf{x}}_{i1} \\ \mathbf{G}^T\tilde{\Sigma}^{-1}\bar{\mathbf{x}}_{i2} \\ \mathbf{G}^T\tilde{\Sigma}^{-1}\bar{\mathbf{x}}_{i3} \end{bmatrix}. \quad (109)$$

Where  $\tilde{\mathbf{Z}} = \left[ \left( \mathbf{I} + \tilde{\mathbf{A}}^T\tilde{\Sigma}^{-1}\tilde{\mathbf{A}} \right)^{-\frac{1}{2}} \right]^T$  and  $\bar{\mathbf{x}}_{i1} = \mathbf{x}_{i1} - \boldsymbol{\mu}$ .

Obviously the above way of writing it is much more efficient in terms of memory as  $\tilde{\mathbf{Z}}$  is a matrix of size  $(D_F + J_i D_G, D_F + J_i D_G)$  and the second vector is of size  $(D_F + J_i D_G, 1)$ . The calculation of the matrix  $\tilde{\mathbf{Z}}$  will be a limiting factor in any application as it involves calculating the inverse and a square root. However, this only has to be derived once for a given number of observations  $J_i$ . Still, this is not the best way to solve this problem. We describe a better, more efficient way of doing this below.



#### 4.1.2 Calculating $-\frac{1}{2}(\tilde{\mathbf{x}}_i - \tilde{\boldsymbol{\mu}})^T (\tilde{\boldsymbol{\Sigma}} + \tilde{\mathbf{A}}\tilde{\mathbf{A}}^T)^{-1} (\tilde{\mathbf{x}}_i - \tilde{\boldsymbol{\mu}})$ Efficiently: Method 2

Starting from where we left off before what we actually want to calculate is,

$$\begin{aligned} & \begin{bmatrix} \sum_{j=1}^{J_i} \mathbf{F}^T \boldsymbol{\Sigma}^{-1} \bar{\mathbf{x}}_{ij} \\ \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \bar{\mathbf{x}}_{i1} \\ \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \bar{\mathbf{x}}_{i2} \\ \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \bar{\mathbf{x}}_{i3} \end{bmatrix}^T \tilde{\mathbf{Z}}^T \tilde{\mathbf{Z}} \begin{bmatrix} \sum_{j=1}^{J_i} \mathbf{F}^T \boldsymbol{\Sigma}^{-1} \bar{\mathbf{x}}_{ij} \\ \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \bar{\mathbf{x}}_{i1} \\ \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \bar{\mathbf{x}}_{i2} \\ \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \bar{\mathbf{x}}_{i3} \end{bmatrix}, \quad (110) \\ & = \left[ \sum_{j=1}^{J_i} \bar{\mathbf{x}}_{ij}^T \boldsymbol{\Sigma}^{-1} \mathbf{F}, \quad \bar{\mathbf{x}}_{i1}^T \boldsymbol{\Sigma}^{-1} \mathbf{G}, \quad \dots, \quad \bar{\mathbf{x}}_{i3}^T \boldsymbol{\Sigma}^{-1} \mathbf{G} \right] \left( \mathbf{I} + \tilde{\mathbf{A}}^T \tilde{\boldsymbol{\Sigma}}^{-1} \tilde{\mathbf{A}} \right)^{-1} \begin{bmatrix} \sum_{j=1}^{J_i} \mathbf{F}^T \boldsymbol{\Sigma}^{-1} \bar{\mathbf{x}}_{ij} \\ \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \bar{\mathbf{x}}_{i1} \\ \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \bar{\mathbf{x}}_{i2} \\ \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \bar{\mathbf{x}}_{i3} \end{bmatrix}. \quad (111) \end{aligned}$$

This comes from the definition of  $\tilde{\mathbf{Z}}$ . Now we can use the matrix inversion lemma, using an intermediate form of the Woodbury identity (coming from an LDU decomposition) again to simplify the inversion of  $(\mathbf{I} + \tilde{\mathbf{A}}^T \tilde{\boldsymbol{\Sigma}}^{-1} \tilde{\mathbf{A}})$ . This leads to the above equation becoming,

$$= \mathbf{a} \begin{bmatrix} I_{D_F} & 0 \\ -D^{-1}C & I_{J_i D_G} \end{bmatrix} \begin{bmatrix} (A - BD^{-1}C)^{-1} & 0 \\ 0 & D^{-1} \end{bmatrix} \begin{bmatrix} I_{D_F} & -BD^{-1} \\ 0 & I_{J_i D_G} \end{bmatrix} \mathbf{b}. \quad (112)$$

Where,

$$\mathbf{a} = \left[ \sum_{j=1}^{J_i} \bar{\mathbf{x}}_{ij}^T \boldsymbol{\Sigma}^{-1} \mathbf{F}, \quad \bar{\mathbf{x}}_{i1}^T \boldsymbol{\Sigma}^{-1} \mathbf{G}, \quad \dots, \quad \bar{\mathbf{x}}_{i3}^T \boldsymbol{\Sigma}^{-1} \mathbf{G} \right], \quad (113)$$

and

$$\mathbf{b} = \begin{bmatrix} \sum_{j=1}^{J_i} \mathbf{F}^T \boldsymbol{\Sigma}^{-1} \bar{\mathbf{x}}_{ij} \\ \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \bar{\mathbf{x}}_{i1} \\ \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \bar{\mathbf{x}}_{i2} \\ \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \bar{\mathbf{x}}_{i3} \end{bmatrix}. \quad (114)$$

We will now proceed to decompose this set of matrix multiplications to derive a more efficient solution. But first we define the elements that we will use below.

$$A = \left[ \mathbf{I} + J_i \mathbf{F}^T \boldsymbol{\Sigma}^{-1} \mathbf{F} \right], \quad (115)$$

$$B = \left[ \mathbf{F}^T \boldsymbol{\Sigma}^{-1} \mathbf{G}, \quad \mathbf{F}^T \boldsymbol{\Sigma}^{-1} \mathbf{G}, \quad \mathbf{F}^T \boldsymbol{\Sigma}^{-1} \mathbf{G} \right], \quad (116)$$

$$C = \begin{bmatrix} \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{F} \\ \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{F} \\ \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{F} \end{bmatrix} = B^T, \quad (117)$$

$$D = \begin{bmatrix} \mathbf{I} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{I} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G} \end{bmatrix}. \quad (118)$$

First we note that,

$$\mathbf{a} \begin{bmatrix} I_{D_F} & 0 \\ -D^{-1}C & I_{J_i D_G} \end{bmatrix} = \left( \begin{bmatrix} I_{D_F} & -BD^{-1} \\ 0 & I_{J_i D_G} \end{bmatrix} \mathbf{b} \right)^T. \quad (119)$$

This means that solving one provides us with the solution for the other, provided we take the transpose. For convenience we choose to solve the second one. We note that,

$$-BD^{-1} = - \left[ \mathbf{F}^T \boldsymbol{\Sigma}^{-1} \mathbf{G}, \quad \dots, \quad \mathbf{F}^T \boldsymbol{\Sigma}^{-1} \mathbf{G} \right] \begin{bmatrix} (\mathbf{I} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G})^{-1} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \ddots & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & (\mathbf{I} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G})^{-1} \end{bmatrix}, \quad (120)$$

$$-BD^{-1} = - \left[ \mathbf{F}^T \boldsymbol{\Sigma}^{-1} \mathbf{G} (\mathbf{I} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G})^{-1}, \dots, \mathbf{F}^T \boldsymbol{\Sigma}^{-1} \mathbf{G} (\mathbf{I} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G})^{-1} \right]. \quad (121)$$

The matrix  $-BD^{-1}$  is of size  $(D_F, J_i D_G)$  and  $\mathbf{b}$  is of size  $(D_F + J_i D_G, J_i D_x)$ . Thus, the multiplication of these two terms is given by,

$$\begin{bmatrix} I & -BD^{-1} \\ 0 & I \end{bmatrix} \mathbf{b} = \begin{bmatrix} I_{D_F}, & -BD^{-1} \\ 0, & I_{J_i D_G} \end{bmatrix} \begin{bmatrix} \sum_{j=1}^{J_i} \mathbf{F}^T \boldsymbol{\Sigma}^{-1} \bar{\mathbf{x}}_{ij} \\ \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \bar{\mathbf{x}}_{i1} \\ \vdots \\ \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \bar{\mathbf{x}}_{i3} \end{bmatrix}, \quad (122)$$

$$= \begin{bmatrix} -BD^{-1} \begin{bmatrix} \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \bar{\mathbf{x}}_{i1} \\ \vdots \\ \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \bar{\mathbf{x}}_{i3} \end{bmatrix} + \sum_{j=1}^{J_i} \mathbf{F}^T \boldsymbol{\Sigma}^{-1} \bar{\mathbf{x}}_{ij} \\ \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \bar{\mathbf{x}}_{i1} \\ \vdots \\ \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \bar{\mathbf{x}}_{i3} \end{bmatrix}, \quad (123)$$

$$= \begin{bmatrix} \left( \sum_{j=1}^{J_i} \mathbf{F}^T \boldsymbol{\Sigma}^{-1} \bar{\mathbf{x}}_{ij} - \sum_{j=1}^{J_i} \mathbf{F}^T \boldsymbol{\Sigma}^{-1} \mathbf{G} (\mathbf{I} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G})^{-1} \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \bar{\mathbf{x}}_{ij} \right) \\ \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \bar{\mathbf{x}}_{i1} \\ \vdots \\ \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \bar{\mathbf{x}}_{i3} \end{bmatrix}, \quad (124)$$

$$= \begin{bmatrix} \mathbf{F}^T \left( \sum_{j=1}^{J_i} \boldsymbol{\Sigma}^{-1} \bar{\mathbf{x}}_{ij} - \sum_{j=1}^{J_i} \boldsymbol{\Sigma}^{-1} \mathbf{G} (\mathbf{I} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G})^{-1} \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \bar{\mathbf{x}}_{ij} \right) \\ \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \bar{\mathbf{x}}_{i1} \\ \vdots \\ \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \bar{\mathbf{x}}_{i3} \end{bmatrix}, \quad (125)$$

$$= \begin{bmatrix} \mathbf{F}^T \left( \sum_{j=1}^{J_i} \left( \boldsymbol{\Sigma}^{-1} - \boldsymbol{\Sigma}^{-1} \mathbf{G} (\mathbf{I} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G})^{-1} \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \right) \bar{\mathbf{x}}_{ij} \right) \\ \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \bar{\mathbf{x}}_{i1} \\ \vdots \\ \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \bar{\mathbf{x}}_{i3} \end{bmatrix}, \quad (126)$$

$$= \begin{bmatrix} \sum_{j=1}^{J_i} \mathbf{F}^T \mathbf{Q} \bar{\mathbf{x}}_{ij} \\ \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \bar{\mathbf{x}}_{i1} \\ \vdots \\ \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \bar{\mathbf{x}}_{i3} \end{bmatrix}. \quad (127)$$

Where  $\mathbf{Q} = \boldsymbol{\Sigma}^{-1} - \boldsymbol{\Sigma}^{-1} \mathbf{G} (\mathbf{I} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G})^{-1} \mathbf{G}^T \boldsymbol{\Sigma}^{-1} = (\boldsymbol{\Sigma} + \mathbf{G} \mathbf{G}^T)^{-1}$ . We also denote another term for use later,  $\mathbf{L} = \boldsymbol{\Sigma}^{-1} \mathbf{G} (\mathbf{I} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G})^{-1} \mathbf{G}^T \boldsymbol{\Sigma}^{-1}$ . Therefore,

$$\mathbf{a} \begin{bmatrix} I_{D_F} & 0 \\ -D^{-1} C & I_{J_i D_G} \end{bmatrix} = \begin{bmatrix} \sum_{j=1}^{J_i} \mathbf{F}^T \mathbf{Q} \bar{\mathbf{x}}_{ij} \\ \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \bar{\mathbf{x}}_{i1} \\ \vdots \\ \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \bar{\mathbf{x}}_{i3} \end{bmatrix}^T, \quad (128)$$

$$= \left[ \sum_{j=1}^{J_i} \bar{\mathbf{x}}_{ij}^T \mathbf{Q}^T \mathbf{F}, \bar{\mathbf{x}}_{i1}^T \boldsymbol{\Sigma}^{-1} \mathbf{G}, \dots, \bar{\mathbf{x}}_{i3}^T \boldsymbol{\Sigma}^{-1} \mathbf{G} \right]. \quad (129)$$

Now let us work on the central term,

$$\begin{bmatrix} (A - BD^{-1}C)^{-1} & 0 \\ 0 & D^{-1} \end{bmatrix}. \quad (130)$$

First we note that we need to solve  $BD^{-1}C = BD^{-1}B^T$ . We first use the definition of  $BD^{-1}$  from before.

$$BD^{-1} = \left[ \mathbf{F}^T \boldsymbol{\Sigma}^{-1} \mathbf{G} (\mathbf{I} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G})^{-1}, \dots, \mathbf{F}^T \boldsymbol{\Sigma}^{-1} \mathbf{G} (\mathbf{I} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G})^{-1} \right]. \quad (131)$$

Then to find  $BD^{-1}B^T$  we get that,

$$BD^{-1}B^T = BD^{-1} \begin{bmatrix} \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{F} \\ \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{F} \\ \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{F} \end{bmatrix}, \quad (132)$$

$$= \left[ J_i \left( \mathbf{F}^T \boldsymbol{\Sigma}^{-1} \mathbf{G} (\mathbf{I} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G})^{-1} \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{F} \right) \right] = BD^{-1}C. \quad (133)$$

We use this to continue to solve for  $A - BD^{-1}C$ ,

$$A - BD^{-1}C = \mathbf{I} + J_i \mathbf{F}^T \boldsymbol{\Sigma}^{-1} \mathbf{F} - J_i \left( \mathbf{F}^T \boldsymbol{\Sigma}^{-1} \mathbf{G} (\mathbf{I} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G})^{-1} \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{F} \right), \quad (134)$$

$$A - BD^{-1}C = \mathbf{I} + J_i \mathbf{F}^T \boldsymbol{\Sigma}^{-1} \mathbf{F} - J_i \left( \mathbf{F}^T \mathbf{L} \mathbf{F} \right), \quad (135)$$

$$A - BD^{-1}C = \mathbf{I} + J_i \mathbf{F}^T (\boldsymbol{\Sigma}^{-1} \mathbf{F} - \mathbf{L} \mathbf{F}), \quad (136)$$

$$A - BD^{-1}C = \mathbf{I} + J_i \mathbf{F}^T (\boldsymbol{\Sigma}^{-1} - \mathbf{L}) \mathbf{F}, \quad (137)$$

$$A - BD^{-1}C = \mathbf{I} + J_i \mathbf{F}^T \mathbf{Q} \mathbf{F}. \quad (138)$$

Therefore,

$$\begin{bmatrix} (A - BD^{-1}C)^{-1} & 0 \\ 0 & D^{-1} \end{bmatrix} = \begin{bmatrix} (\mathbf{I} + J_i \mathbf{F}^T \mathbf{Q} \mathbf{F})^{-1} & 0 \\ 0 & D^{-1} \end{bmatrix}. \quad (139)$$

Now, let us keep going with this. Work out the right hand side which is,

$$\begin{bmatrix} (A - BD^{-1}C)^{-1} & 0 \\ 0 & D^{-1} \end{bmatrix} \begin{bmatrix} I_{D_F} & -BD^{-1} \\ 0 & I_{J_i D_F} \end{bmatrix} \mathbf{b} = \begin{bmatrix} (\mathbf{I} + J_i \mathbf{F}^T \mathbf{Q} \mathbf{F})^{-1} & 0 \\ 0 & D^{-1} \end{bmatrix} \begin{bmatrix} \sum_{j=1}^{J_i} \mathbf{F}^T \mathbf{Q} \bar{\mathbf{x}}_{ij} \\ \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \bar{\mathbf{x}}_{i1} \\ \vdots \\ \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \bar{\mathbf{x}}_{i3} \end{bmatrix}, \quad (140)$$

$$= \begin{bmatrix} (\mathbf{I} + J_i \mathbf{F}^T \mathbf{Q} \mathbf{F})^{-1} & 0 \\ 0 & \begin{bmatrix} (\mathbf{I} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G})^{-1} & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & (\mathbf{I} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G})^{-1} \end{bmatrix} \end{bmatrix} \begin{bmatrix} \sum_{j=1}^{J_i} \mathbf{F}^T \mathbf{Q} \bar{\mathbf{x}}_{ij} \\ \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \bar{\mathbf{x}}_{i1} \\ \vdots \\ \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \bar{\mathbf{x}}_{i3} \end{bmatrix}, \quad (141)$$

$$= \begin{bmatrix} (\mathbf{I} + J_i \mathbf{F}^T \mathbf{Q} \mathbf{F})^{-1} \sum_{j=1}^{J_i} \mathbf{F}^T \mathbf{Q} \bar{\mathbf{x}}_{ij} \\ (\mathbf{I} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G})^{-1} \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \bar{\mathbf{x}}_{i1} \\ \vdots \\ (\mathbf{I} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G})^{-1} \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \bar{\mathbf{x}}_{i3} \end{bmatrix}. \quad (142)$$

Then we incorporate the left hand side and we get,

$$\begin{aligned}
& \left[ \sum_{j=1}^{J_i} \bar{\mathbf{x}}_{ij}^T \mathbf{Q}^T \mathbf{F}, \quad \bar{\mathbf{x}}_{i1}^T \boldsymbol{\Sigma}^{-1} \mathbf{G}, \quad \dots, \quad \bar{\mathbf{x}}_{i3}^T \boldsymbol{\Sigma}^{-1} \mathbf{G} \right] \begin{bmatrix} \left( \mathbf{I} + J_i \mathbf{F}^T \mathbf{Q} \mathbf{F} \right)^{-1} \sum_{j=1}^{J_i} \mathbf{F}^T \mathbf{Q} \bar{\mathbf{x}}_{ij} \\ \left( \mathbf{I} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G} \right)^{-1} \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \bar{\mathbf{x}}_{i1} \\ \vdots \\ \left( \mathbf{I} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G} \right)^{-1} \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \bar{\mathbf{x}}_{i3} \end{bmatrix} \quad (143) \\
& = \left[ \sum_{j=1}^{J_i} \bar{\mathbf{x}}_{ij}^T \mathbf{Q}^T \mathbf{F} \left( \mathbf{I} + J_i \mathbf{F}^T \mathbf{Q} \mathbf{F} \right)^{-1} \sum_{j=1}^{J_i} \mathbf{F}^T \mathbf{Q} \bar{\mathbf{x}}_{ij} + \sum_{j=1}^{J_i} \bar{\mathbf{x}}_{ij}^T \boldsymbol{\Sigma}^{-1} \mathbf{G} \left( \mathbf{I} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G} \right)^{-1} \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \bar{\mathbf{x}}_{ij} \right]. \quad (144)
\end{aligned}$$

This is the efficient formula to do full integration log-likelihood scoring for PLDA. It is efficient for the following reasons.

1. Computing  $\left( \mathbf{I} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G} \right)^{-1}$  will be efficient because it is a matrix of size  $(D_G, D_G)$ . In fact it can be pre-computed and does not depend at all upon the number of samples  $J_i$ .
2.  $\mathbf{Q} = \boldsymbol{\Sigma}^{-1} - \boldsymbol{\Sigma}^{-1} \mathbf{G} \left( \mathbf{I} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G} \right)^{-1} \mathbf{G}^T \boldsymbol{\Sigma}^{-1}$  is efficient to compute because  $\left( \mathbf{I} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G} \right)^{-1}$  is efficient to compute as is  $\mathbf{G}^T \boldsymbol{\Sigma}^{-1}$  and  $\boldsymbol{\Sigma}^{-1} \mathbf{G}$  because  $\boldsymbol{\Sigma}^{-1}$  is diagonal. In fact it can be pre-computed and does not depend at all upon the number of samples  $J_i$ .
3. Computing  $\left( \mathbf{I} + J_i \mathbf{F}^T \mathbf{Q} \mathbf{F} \right)^{-1}$  is efficient because it is a matrix of size  $(D_F, D_F)$  and  $\mathbf{Q}$  is efficient to compute (and can even be pre-computed). Also, the matrix does not increase in size depending upon the number of samples  $J_i$ . However, it does depend upon the number of samples.
4.  $\mathbf{F}^T \mathbf{Q} \bar{\mathbf{x}}_{ij}$  is efficient to compute and is only of size  $(D_F, 1)$ , also  $\mathbf{G}^T \boldsymbol{\Sigma}^{-1} \bar{\mathbf{x}}_{ij}$  is efficient to compute and is only of size  $(D_G, 1)$ . This also indicates that enrolment should consist of obtaining these quantities and not on keeping the full feature vector.

#### 4.1.3 Determinant of $\left( \tilde{\boldsymbol{\Sigma}} + \tilde{\mathbf{A}} \tilde{\mathbf{A}}^T \right)$

How do we easily find the determinant of this quite large matrix? First we note that,

$$\det \left( \tilde{\boldsymbol{\Sigma}} + \tilde{\mathbf{A}} \tilde{\mathbf{A}}^T \right) = \det \left( \tilde{\boldsymbol{\Sigma}} \right) \det \left( \mathbf{I} + \tilde{\mathbf{A}}^T \tilde{\boldsymbol{\Sigma}}^{-1} \tilde{\mathbf{A}} \right). \quad (145)$$

This particular trick can be found in wikipedia (<http://en.wikipedia.org/wiki/Determinant>, Sylvester's determinant theorem). This works well because  $\tilde{\boldsymbol{\Sigma}}$  is diagonal and so  $\det \left( \tilde{\boldsymbol{\Sigma}} \right)$  is trivial to compute,

$$\det \left( \tilde{\boldsymbol{\Sigma}} \right) = \det \left( \boldsymbol{\Sigma} \right)^{J_i}$$

Taking up the second term, we note that  $\tilde{\mathbf{A}}^T \tilde{\boldsymbol{\Sigma}}^{-1} \tilde{\mathbf{A}}$  is of size  $[(D_F + J_i D_G), (D_F + J_i D_G)]$ . However, there should be a structure to  $\tilde{\mathbf{A}}^T \tilde{\boldsymbol{\Sigma}}^{-1} \tilde{\mathbf{A}}$ , so let us try to further simplify this,

$$\tilde{\mathbf{A}}^T \tilde{\boldsymbol{\Sigma}}^{-1} \tilde{\mathbf{A}} = \begin{bmatrix} \mathbf{F}^T & \mathbf{F}^T & \mathbf{F}^T \\ \mathbf{G}^T & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{G}^T & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{G}^T \end{bmatrix} \begin{bmatrix} \boldsymbol{\Sigma}^{-1} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\Sigma}^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \boldsymbol{\Sigma}^{-1} \end{bmatrix} \begin{bmatrix} \mathbf{F} & \mathbf{G} & \mathbf{0} & \mathbf{0} \\ \mathbf{F} & \mathbf{0} & \mathbf{G} & \mathbf{0} \\ \mathbf{F} & \mathbf{0} & \mathbf{0} & \mathbf{G} \end{bmatrix}, \quad (146)$$

$$= \begin{bmatrix} \mathbf{F}^T \boldsymbol{\Sigma}^{-1}, & \dots, & \mathbf{F}^T \boldsymbol{\Sigma}^{-1} \\ \mathbf{G}^T \boldsymbol{\Sigma}^{-1} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \ddots & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \end{bmatrix} \begin{bmatrix} \mathbf{F} & \mathbf{G} & \mathbf{0} & \mathbf{0} \\ \mathbf{F} & \mathbf{0} & \mathbf{G} & \mathbf{0} \\ \mathbf{F} & \mathbf{0} & \mathbf{0} & \mathbf{G} \end{bmatrix}, \quad (147)$$

$$= \begin{bmatrix} J_i \mathbf{F}^T \boldsymbol{\Sigma}^{-1} \mathbf{F}, & \mathbf{F}^T \boldsymbol{\Sigma}^{-1} \mathbf{G}, & \dots, & \mathbf{F}^T \boldsymbol{\Sigma}^{-1} \mathbf{G} \\ \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{F} & \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G} & \mathbf{0} & \mathbf{0} \\ \vdots & \mathbf{0} & \ddots & \mathbf{0} \\ \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{F} & \mathbf{0} & \mathbf{0} & \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G} \end{bmatrix}, \quad (148)$$

$$\mathbf{I} + \tilde{\mathbf{A}}^T \tilde{\boldsymbol{\Sigma}}^{-1} \tilde{\mathbf{A}} = \begin{bmatrix} \mathbf{I} + J_i \mathbf{F}^T \boldsymbol{\Sigma}^{-1} \mathbf{F}, & \mathbf{F}^T \boldsymbol{\Sigma}^{-1} \mathbf{G}, & \dots, & \mathbf{F}^T \boldsymbol{\Sigma}^{-1} \mathbf{G} \\ \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{F} & \mathbf{I} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G} & \mathbf{0} & \mathbf{0} \\ \vdots & \mathbf{0} & \ddots & \mathbf{0} \\ \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{F} & \mathbf{0} & \mathbf{0} & \mathbf{I} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G} \end{bmatrix}. \quad (149)$$

If need be we could further subdivide the problem. By decomposing the determinant of  $\mathbf{I} + \tilde{\mathbf{A}}^T \tilde{\boldsymbol{\Sigma}}^{-1} \tilde{\mathbf{A}}$  into a block matrix we can use the identity,

$$\det \begin{pmatrix} A & B \\ C & D \end{pmatrix} = \det(A) \det(D - CA^{-1}B) = \det(D) \det(A - BD^{-1}C). \quad (150)$$

The terms  $A$ ,  $B$ ,  $C$  and  $D$  are the same as before. So an efficient way to compute this is to note that,

$$\det(D) = \begin{bmatrix} \mathbf{I} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \ddots & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{I} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G} \end{bmatrix}, \quad (151)$$

$$= \det(\mathbf{I} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G})^{J_i}. \quad (152)$$

This is an  $(D_G, D_G)$  matrix whose determinant would be efficient to compute. Next, we already have an efficient expression for the second term and that is,

$$A - BD^{-1}C = \mathbf{I} + J_i \mathbf{F}^T \mathbf{Q} \mathbf{F}. \quad (153)$$

Consequently we need to determine,

$$\det(\mathbf{I} + J_i \mathbf{F}^T \mathbf{Q} \mathbf{F}). \quad (154)$$

This is an  $(D_F, D_F)$  matrix and could be computed efficiently.

Finally this gives us the following solution.

$$\det(\tilde{\boldsymbol{\Sigma}} + \tilde{\mathbf{A}} \tilde{\mathbf{A}}^T) = \det(\tilde{\boldsymbol{\Sigma}}) \det(\mathbf{I} + \tilde{\mathbf{A}}^T \tilde{\boldsymbol{\Sigma}}^{-1} \tilde{\mathbf{A}}), \quad (155)$$

$$= [\det(\boldsymbol{\Sigma})]^{J_i} [\det(\mathbf{I} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G})]^{J_i} [\det(\mathbf{I} + J_i \mathbf{F}^T \mathbf{Q} \mathbf{F})]. \quad (156)$$

Taking the logarithm we would get,

$$\ln(\det(\tilde{\boldsymbol{\Sigma}} + \tilde{\mathbf{A}} \tilde{\mathbf{A}}^T)) = \ln([\det(\boldsymbol{\Sigma})]^{J_i} [\det(\mathbf{I} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G})]^{J_i} [\det(\mathbf{I} + J_i \mathbf{F}^T \mathbf{Q} \mathbf{F})]), \quad (157)$$

$$= \ln([\det(\boldsymbol{\Sigma})]^{J_i}) + \ln([\det(\mathbf{I} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G})]^{J_i}) + \ln([\det(\mathbf{I} + J_i \mathbf{F}^T \mathbf{Q} \mathbf{F})]), \quad (158)$$

$$= J_i \ln(\det(\boldsymbol{\Sigma})) + J_i \ln(\det(\mathbf{I} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G})) + \ln(\det(\mathbf{I} + J_i \mathbf{F}^T \mathbf{Q} \mathbf{F})). \quad (159)$$

Providing a more useful final solution this gives us that,

$$-\frac{1}{2} \ln \left[ \det \left( \tilde{\Sigma} + \tilde{\mathbf{A}} \tilde{\mathbf{A}}^T \right) \right] = -\frac{1}{2} \ln \left( [\det(\Sigma)]^{J_i} [\det(\mathbf{I} + \mathbf{G}^T \Sigma^{-1} \mathbf{G})]^{J_i} \det(\mathbf{I} + J_i \mathbf{F}^T \mathbf{Q} \mathbf{F}) \right), \quad (160)$$

$$= -\frac{1}{2} \left[ J_i \ln(\det(\Sigma)) + J_i \ln(\det(\mathbf{I} + \mathbf{G}^T \Sigma^{-1} \mathbf{G})) + \ln(\det(\mathbf{I} + J_i \mathbf{F}^T \mathbf{Q} \mathbf{F})) \right], \quad (161)$$

$$= -\frac{J_i}{2} \ln(\det(\Sigma)) - \frac{J_i}{2} \ln(\det(\mathbf{I} + \mathbf{G}^T \Sigma^{-1} \mathbf{G})) - \frac{1}{2} \ln(\det(\mathbf{I} + J_i \mathbf{F}^T \mathbf{Q} \mathbf{F})). \quad (162)$$

#### 4.1.4 Final Solution

Now using the above solution we can extend it and find an efficient way to calculate the sufficient statistics. We take this up below. But before we do that, what is the final equation to calculate the log-likelihood?

$$\ln [p(\tilde{\mathbf{x}}_i)] = -\frac{J_i D_x}{2} \ln(2\pi) - \frac{1}{2} \ln \left[ \det \left( \tilde{\Sigma} + \tilde{\mathbf{A}} \tilde{\mathbf{A}}^T \right) \right] - \frac{1}{2} (\tilde{\mathbf{x}}_i - \tilde{\boldsymbol{\mu}})^T \left( \tilde{\Sigma} + \tilde{\mathbf{A}} \tilde{\mathbf{A}}^T \right)^{-1} (\tilde{\mathbf{x}}_i - \tilde{\boldsymbol{\mu}}), \quad (163)$$

$$\begin{aligned} \ln [p(\tilde{\mathbf{x}}_i)] = & -\frac{J_i D_x}{2} \ln(2\pi) - \frac{1}{2} \ln \left[ \det \left( \tilde{\Sigma} + \tilde{\mathbf{A}} \tilde{\mathbf{A}}^T \right) \right] \\ & - \frac{1}{2} (\tilde{\mathbf{x}}_i - \tilde{\boldsymbol{\mu}})^T \tilde{\Sigma}^{-1} (\tilde{\mathbf{x}}_i - \tilde{\boldsymbol{\mu}}) + \frac{1}{2} (\tilde{\mathbf{x}}_i - \tilde{\boldsymbol{\mu}})^T \boldsymbol{\Gamma} \boldsymbol{\Gamma}^T (\tilde{\mathbf{x}}_i - \tilde{\boldsymbol{\mu}}), \end{aligned} \quad (164)$$

$$\begin{aligned} \ln [p(\tilde{\mathbf{x}}_i)] = & -\frac{J_i D_x}{2} \ln(2\pi) - \frac{1}{2} \ln \left[ \det \left( \tilde{\Sigma} + \tilde{\mathbf{A}} \tilde{\mathbf{A}}^T \right) \right] - \frac{1}{2} (\tilde{\mathbf{x}}_i - \tilde{\boldsymbol{\mu}})^T \tilde{\Sigma}^{-1} (\tilde{\mathbf{x}}_i - \tilde{\boldsymbol{\mu}}) \\ & + \frac{1}{2} \sum_{j=1}^{J_i} \tilde{\mathbf{x}}_{ij}^T \mathbf{Q}^T \mathbf{F} \left( \mathbf{I} + J_i \mathbf{F}^T \mathbf{Q} \mathbf{F} \right)^{-1} \sum_{j=1}^{J_i} \mathbf{F}^T \mathbf{Q} \tilde{\mathbf{x}}_{ij} + \frac{1}{2} \sum_{j=1}^{J_i} \tilde{\mathbf{x}}_{ij}^T \Sigma^{-1} \mathbf{G} \left( \mathbf{I} + \mathbf{G}^T \Sigma^{-1} \mathbf{G} \right)^{-1} \mathbf{G}^T \Sigma^{-1} \tilde{\mathbf{x}}_{ij}. \end{aligned} \quad (165)$$

We further note the following,

$$-\frac{1}{2} \ln \left[ \det \left( \tilde{\Sigma} + \tilde{\mathbf{A}} \tilde{\mathbf{A}}^T \right) \right] = -\frac{J_i}{2} \ln(\det(\Sigma)) - \frac{J_i}{2} \ln(\det(\mathbf{I} + \mathbf{G}^T \Sigma^{-1} \mathbf{G})) - \frac{1}{2} \ln(\det(\mathbf{I} + J_i \mathbf{F}^T \mathbf{Q} \mathbf{F})), \quad (166)$$

and

$$-\frac{1}{2} (\tilde{\mathbf{x}}_i - \tilde{\boldsymbol{\mu}})^T \tilde{\Sigma}^{-1} (\tilde{\mathbf{x}}_i - \tilde{\boldsymbol{\mu}}) = -\frac{1}{2} \begin{bmatrix} \tilde{\mathbf{x}}_{i1}^T & \cdots & \tilde{\mathbf{x}}_{i3}^T \end{bmatrix} \begin{bmatrix} \Sigma^{-1} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \ddots & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \Sigma^{-1} \end{bmatrix} \begin{bmatrix} \tilde{\mathbf{x}}_{i1} \\ \vdots \\ \tilde{\mathbf{x}}_{i3} \end{bmatrix}, \quad (167)$$

$$= -\frac{1}{2} \begin{bmatrix} \tilde{\mathbf{x}}_{i1}^T \Sigma^{-1} & \cdots & \tilde{\mathbf{x}}_{i3}^T \Sigma^{-1} \end{bmatrix} \begin{bmatrix} \tilde{\mathbf{x}}_{i1} \\ \vdots \\ \tilde{\mathbf{x}}_{i3} \end{bmatrix}, \quad (168)$$

$$= -\frac{1}{2} \sum_{j=1}^{J_i} \tilde{\mathbf{x}}_{ij}^T \Sigma^{-1} \tilde{\mathbf{x}}_{ij}. \quad (169)$$

Thus, this can all be efficiently computed and even pre-computed for a model. The very final formulation that we have is,

$$\begin{aligned} \ln [p(\tilde{\mathbf{x}}_i)] &= -\frac{J_i D_x}{2} \ln(2\pi) - \frac{J_i}{2} \ln(\det(\boldsymbol{\Sigma})) - \frac{J_i}{2} \ln(\det(\mathbf{I} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G})) \\ &\quad - \frac{1}{2} \ln\left(\det\left(\mathbf{I} + J_i \mathbf{F}^T \mathbf{Q} \mathbf{F}\right)\right) + \frac{1}{2} \sum_{j=1}^{J_i} \tilde{\mathbf{x}}_{ij}^T \mathbf{Q}^T \mathbf{F} \left(\mathbf{I} + J_i \mathbf{F}^T \mathbf{Q} \mathbf{F}\right)^{-1} \sum_{j=1}^{J_i} \mathbf{F}^T \mathbf{Q} \tilde{\mathbf{x}}_{ij} \\ &\quad + \frac{1}{2} \sum_{j=1}^{J_i} \tilde{\mathbf{x}}_{ij}^T \boldsymbol{\Sigma}^{-1} \mathbf{G} \left(\mathbf{I} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G}\right)^{-1} \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \tilde{\mathbf{x}}_{ij} - \frac{1}{2} \sum_{j=1}^{J_i} \tilde{\mathbf{x}}_{ij}^T \boldsymbol{\Sigma}^{-1} \tilde{\mathbf{x}}_{ij}. \end{aligned} \quad (170)$$

## 4.2 Sufficient Statistics

The sufficient statistics for the PLDA model are  $E[\tilde{\mathbf{y}}_i | \tilde{\mathbf{x}}_i, \boldsymbol{\Theta}]$  and  $E[\tilde{\mathbf{y}}_i \tilde{\mathbf{y}}_i^T | \tilde{\mathbf{x}}_i, \boldsymbol{\Theta}]$ . These both involve calculations with the matrix  $(\mathbf{I} + \tilde{\mathbf{A}}^T \tilde{\boldsymbol{\Sigma}}^{-1} \tilde{\mathbf{A}})^{-1}$ . Above, we showed that for the log-likelihood we can exploit the structure of this matrix to find a scalable and exact solution. We will now show how to do the same thing for the sufficient statistics.

### 4.2.1 Solving for $E[\tilde{\mathbf{y}}_i | \tilde{\mathbf{x}}_i, \boldsymbol{\Theta}]$

First, we reuse the matrix  $-BD^{-1}$ ,

$$-BD^{-1} = - \left[ \mathbf{F}^T \boldsymbol{\Sigma}^{-1} \mathbf{G} \left(\mathbf{I} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G}\right)^{-1}, \dots, \mathbf{F}^T \boldsymbol{\Sigma}^{-1} \mathbf{G} \left(\mathbf{I} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G}\right)^{-1} \right], \quad (171)$$

and its transpose,

$$-D^{-1}C = - \begin{bmatrix} \left(\mathbf{I} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G}\right)^{-1} \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{F} \\ \vdots \\ \left(\mathbf{I} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G}\right)^{-1} \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{F} \end{bmatrix}. \quad (172)$$

We then note that,

$$\left(\mathbf{I} + \tilde{\mathbf{A}}^T \tilde{\boldsymbol{\Sigma}}^{-1} \tilde{\mathbf{A}}\right)^{-1} \tilde{\mathbf{A}}^T \tilde{\boldsymbol{\Sigma}}^{-1} (\tilde{\mathbf{x}}_i - \tilde{\boldsymbol{\mu}}) = \left(\mathbf{I} + \tilde{\mathbf{A}}^T \tilde{\boldsymbol{\Sigma}}^{-1} \tilde{\mathbf{A}}\right)^{-1} \begin{bmatrix} \sum_{j=1}^{J_i} \mathbf{F}^T \boldsymbol{\Sigma}^{-1} \tilde{\mathbf{x}}_{ij} \\ \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \tilde{\mathbf{x}}_{i1} \\ \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \tilde{\mathbf{x}}_{i2} \\ \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \tilde{\mathbf{x}}_{i3} \end{bmatrix}, \quad (173)$$

$$= \begin{bmatrix} I_{D_F} & 0 \\ -D^{-1}C & I_{J_i D_G} \end{bmatrix} \begin{bmatrix} (A - BD^{-1}C)^{-1} & 0 \\ 0 & D^{-1} \end{bmatrix} \begin{bmatrix} I_{D_F} & -BD^{-1} \\ 0 & I_{J_i D_F} \end{bmatrix} \mathbf{b}. \quad (174)$$

We have already solved most of this so we take up the solution halfway through by noting that we need to find,

$$= \begin{bmatrix} I_{D_F} & 0 \\ -D^{-1}C & I_{J_i D_G} \end{bmatrix} \begin{bmatrix} \left(\mathbf{I} + J_i \mathbf{F}^T \mathbf{Q} \mathbf{F}\right)^{-1} \sum_{j=1}^{J_i} \mathbf{F}^T \mathbf{Q} \tilde{\mathbf{x}}_{ij} \\ \left(\mathbf{I} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G}\right)^{-1} \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \tilde{\mathbf{x}}_{i1} \\ \vdots \\ \left(\mathbf{I} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G}\right)^{-1} \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \tilde{\mathbf{x}}_{i3} \end{bmatrix}, \quad (175)$$

$$= \begin{bmatrix} I_{D_F} & & 0 \\ - \left[ \begin{array}{c} (\mathbf{I} + \mathbf{G}^T \Sigma^{-1} \mathbf{G})^{-1} \mathbf{G}^T \Sigma^{-1} \mathbf{F} \\ \vdots \\ (\mathbf{I} + \mathbf{G}^T \Sigma^{-1} \mathbf{G})^{-1} \mathbf{G}^T \Sigma^{-1} \mathbf{F} \end{array} \right] & \begin{bmatrix} I_{D_G} & 0 & 0 \\ 0 & I_{D_G} & 0 \\ 0 & 0 & I_{D_G} \end{bmatrix} & \begin{bmatrix} (\mathbf{I} + J_i \mathbf{F}^T \mathbf{Q} \mathbf{F})^{-1} \sum_{j=1}^{J_i} \mathbf{F}^T \mathbf{Q} \bar{\mathbf{x}}_{ij} \\ (\mathbf{I} + \mathbf{G}^T \Sigma^{-1} \mathbf{G})^{-1} \mathbf{G}^T \Sigma^{-1} \bar{\mathbf{x}}_{i1} \\ \vdots \\ (\mathbf{I} + \mathbf{G}^T \Sigma^{-1} \mathbf{G})^{-1} \mathbf{G}^T \Sigma^{-1} \bar{\mathbf{x}}_{i3} \end{bmatrix} \end{bmatrix}, \quad (176)$$

$$= \begin{bmatrix} (\mathbf{I} + J_i \mathbf{F}^T \mathbf{Q} \mathbf{F})^{-1} \sum_{j=1}^{J_i} \mathbf{F}^T \mathbf{Q} \bar{\mathbf{x}}_{ij} \\ - (\mathbf{I} + \mathbf{G}^T \Sigma^{-1} \mathbf{G})^{-1} \mathbf{G}^T \Sigma^{-1} \mathbf{F} (\mathbf{I} + J_i \mathbf{F}^T \mathbf{Q} \mathbf{F})^{-1} \sum_{j=1}^{J_i} \mathbf{F}^T \mathbf{Q} \bar{\mathbf{x}}_{ij} + (\mathbf{I} + \mathbf{G}^T \Sigma^{-1} \mathbf{G})^{-1} \mathbf{G}^T \Sigma^{-1} \bar{\mathbf{x}}_{i1} \\ \vdots \\ - (\mathbf{I} + \mathbf{G}^T \Sigma^{-1} \mathbf{G})^{-1} \mathbf{G}^T \Sigma^{-1} \mathbf{F} (\mathbf{I} + J_i \mathbf{F}^T \mathbf{Q} \mathbf{F})^{-1} \sum_{j=1}^{J_i} \mathbf{F}^T \mathbf{Q} \bar{\mathbf{x}}_{ij} + (\mathbf{I} + \mathbf{G}^T \Sigma^{-1} \mathbf{G})^{-1} \mathbf{G}^T \Sigma^{-1} \bar{\mathbf{x}}_{i3} \end{bmatrix}. \quad (177)$$

This provides us with a very interesting solution because we have now separated the factors and all the matrices we have to invert are small. The final solution is obviously,

$$E [\tilde{\mathbf{y}}_i | \tilde{\mathbf{x}}_i, \Theta] = \begin{bmatrix} (\mathbf{I} + J_i \mathbf{F}^T \mathbf{Q} \mathbf{F})^{-1} \sum_{j=1}^{J_i} \mathbf{F}^T \mathbf{Q} \bar{\mathbf{x}}_{ij} \\ (\mathbf{I} + \mathbf{G}^T \Sigma^{-1} \mathbf{G})^{-1} \mathbf{G}^T \Sigma^{-1} \left[ \bar{\mathbf{x}}_{i1} - \mathbf{F} (\mathbf{I} + J_i \mathbf{F}^T \mathbf{Q} \mathbf{F})^{-1} \sum_{j=1}^{J_i} \mathbf{F}^T \mathbf{Q} \bar{\mathbf{x}}_{ij} \right] \\ \vdots \\ (\mathbf{I} + \mathbf{G}^T \Sigma^{-1} \mathbf{G})^{-1} \mathbf{G}^T \Sigma^{-1} \left[ \bar{\mathbf{x}}_{i3} - \mathbf{F} (\mathbf{I} + J_i \mathbf{F}^T \mathbf{Q} \mathbf{F})^{-1} \sum_{j=1}^{J_i} \mathbf{F}^T \mathbf{Q} \bar{\mathbf{x}}_{ij} \right] \end{bmatrix}. \quad (178)$$

This leads to,

$$E [\mathbf{h}_i | \tilde{\mathbf{x}}_i, \Theta] = (\mathbf{I} + J_i \mathbf{F}^T \mathbf{Q} \mathbf{F})^{-1} \sum_{j=1}^{J_i} \mathbf{F}^T \mathbf{Q} \bar{\mathbf{x}}_{ij}, \quad (179)$$

$$E [\mathbf{w}_{ij} | \tilde{\mathbf{x}}_i, \Theta] = (\mathbf{I} + \mathbf{G}^T \Sigma^{-1} \mathbf{G})^{-1} \mathbf{G}^T \Sigma^{-1} \left[ \bar{\mathbf{x}}_{ij} - \mathbf{F} (\mathbf{I} + J_i \mathbf{F}^T \mathbf{Q} \mathbf{F})^{-1} \sum_{j=1}^{J_i} \mathbf{F}^T \mathbf{Q} \bar{\mathbf{x}}_{ij} \right], \quad (180)$$

$$= (\mathbf{I} + \mathbf{G}^T \Sigma^{-1} \mathbf{G})^{-1} \mathbf{G}^T \Sigma^{-1} [\bar{\mathbf{x}}_{ij} - \mathbf{F} E [\mathbf{h}_i | \tilde{\mathbf{x}}_i, \Theta]]. \quad (181)$$

#### 4.2.2 Solving for $E [\tilde{\mathbf{y}}_i \tilde{\mathbf{y}}_i^T | \tilde{\mathbf{x}}_i, \Theta]$

To solve this we use the equation,

$$E [\tilde{\mathbf{y}}_i \tilde{\mathbf{y}}_i^T | \tilde{\mathbf{x}}_i, \Theta] = (\mathbf{I} + \tilde{\mathbf{A}}^T \tilde{\Sigma}^{-1} \tilde{\mathbf{A}})^{-1} + E [\tilde{\mathbf{y}}_i | \tilde{\mathbf{x}}_i, \Theta] E [\tilde{\mathbf{y}}_i | \tilde{\mathbf{x}}_i, \Theta]^T. \quad (182)$$

This implies that we have to store the full matrix  $(\mathbf{I} + \tilde{\mathbf{A}}^T \tilde{\Sigma}^{-1} \tilde{\mathbf{A}})^{-1}$  whose size depends on the number of samples, however, in fact we use this on a per latent variable basis. That is, we find  $\mathbf{y}_{ij} = [\mathbf{h}_i, \mathbf{w}_{ij}]$  for each  $j = [1, 2, \dots, J_i]$  and they are treated separately. This is because it is only used in Equation 71. This suggests that it might be possible to do this efficiently if we can find an expression for  $(\mathbf{I} + \tilde{\mathbf{A}}^T \tilde{\Sigma}^{-1} \tilde{\mathbf{A}})^{-1}$ . We take this problem up below.

The matrix  $(\mathbf{I} + \tilde{\mathbf{A}}^T \tilde{\Sigma}^{-1} \tilde{\mathbf{A}})^{-1}$  can be defined, and calculated efficiently, by using the matrix inversion lemma. That is,



$$\begin{pmatrix} A & B \\ C & D \end{pmatrix}^{-1} = \begin{pmatrix} (A - BD^{-1}C)^{-1} & -(A - BD^{-1}C)^{-1}BD^{-1} \\ -D^{-1}C(A - BD^{-1}C)^{-1} & D^{-1} + D^{-1}C(A - BD^{-1}C)^{-1}BD^{-1} \end{pmatrix}, \quad (183)$$

and

$$\left(\mathbf{I} + \tilde{\mathbf{A}}^T \tilde{\Sigma}^{-1} \tilde{\mathbf{A}}\right) = \begin{bmatrix} \mathbf{I} + J_i \mathbf{F}^T \Sigma^{-1} \mathbf{F}, & \mathbf{F}^T \Sigma^{-1} \mathbf{G}, & \dots, & \mathbf{F}^T \Sigma^{-1} \mathbf{G} \\ \mathbf{G}^T \Sigma^{-1} \mathbf{F} & \mathbf{I} + \mathbf{G}^T \Sigma^{-1} \mathbf{G} & \mathbf{0} & \mathbf{0} \\ \vdots & \mathbf{0} & \ddots & \mathbf{0} \\ \mathbf{G}^T \Sigma^{-1} \mathbf{F} & \mathbf{0} & \mathbf{0} & \mathbf{I} + \mathbf{G}^T \Sigma^{-1} \mathbf{G} \end{bmatrix}. \quad (184)$$

We take some definitions from the appendix and reproduce them here.

$$(A - BD^{-1}C)^{-1} = (\mathbf{I}_{D_F} + J_i \mathbf{F}^T \mathbf{Q} \mathbf{F})^{-1}, \quad (185)$$

$$-D^{-1}C(A - BD^{-1}C)^{-1} = - \begin{bmatrix} (\mathbf{I}_{D_G} + \mathbf{G}^T \Sigma^{-1} \mathbf{G})^{-1} \mathbf{G}^T \Sigma^{-1} \mathbf{F} (\mathbf{I}_{D_F} + J_i \mathbf{F}^T \mathbf{Q} \mathbf{F})^{-1} \\ (\mathbf{I}_{D_G} + \mathbf{G}^T \Sigma^{-1} \mathbf{G})^{-1} \mathbf{G}^T \Sigma^{-1} \mathbf{F} (\mathbf{I}_{D_F} + J_i \mathbf{F}^T \mathbf{Q} \mathbf{F})^{-1} \\ (\mathbf{I}_{D_G} + \mathbf{G}^T \Sigma^{-1} \mathbf{G})^{-1} \mathbf{G}^T \Sigma^{-1} \mathbf{F} (\mathbf{I}_{D_F} + J_i \mathbf{F}^T \mathbf{Q} \mathbf{F})^{-1} \end{bmatrix}, \quad (186)$$

$$-(A - BD^{-1}C)^{-1} BD^{-1} = - \left[ D^{-1}C(A - BD^{-1}C)^{-1} \right]^T, \quad (187)$$

and

$$\begin{aligned} & D^{-1} + D^{-1}C(A - BD^{-1}C)^{-1}BD^{-1} \\ &= \begin{bmatrix} (\mathbf{I}_{D_G} + \mathbf{G}^T \Sigma^{-1} \mathbf{G})^{-1} + \mathbf{R} \mathbf{P}, & \mathbf{R} \mathbf{P} & \mathbf{R} \mathbf{P} \\ \mathbf{R} \mathbf{P} & \ddots & \mathbf{R} \mathbf{P} \\ \mathbf{R} \mathbf{P} & \mathbf{R} \mathbf{P} & (\mathbf{I}_{D_G} + \mathbf{G}^T \Sigma^{-1} \mathbf{G})^{-1} + \mathbf{R} \mathbf{P} \end{bmatrix}, \end{aligned} \quad (188)$$

where

$$\mathbf{P} = \mathbf{F}^T \Sigma^{-1} \mathbf{G} (\mathbf{I}_{D_G} + \mathbf{G}^T \Sigma^{-1} \mathbf{G})^{-1}, \quad (189)$$

and

$$\mathbf{R} = (\mathbf{I}_{D_G} + \mathbf{G}^T \Sigma^{-1} \mathbf{G})^{-1} \mathbf{G}^T \Sigma^{-1} \mathbf{F} (\mathbf{I}_{D_F} + J_i \mathbf{F}^T \mathbf{Q} \mathbf{F})^{-1}. \quad (190)$$

From the above definitions we can now define our matrix  $\left(\mathbf{I} + \tilde{\mathbf{A}}^T \tilde{\Sigma}^{-1} \tilde{\mathbf{A}}\right)^{-1}$  as being,

$$\left[ \begin{array}{c} (\mathbf{I}_{D_F} + J_i \mathbf{F}^T \mathbf{Q} \mathbf{F})^{-1} \\ - \begin{bmatrix} \mathbf{R} \\ \mathbf{R} \\ \mathbf{R} \end{bmatrix} \end{array} \begin{array}{c} - \begin{bmatrix} \mathbf{R}^T & \mathbf{R}^T & \mathbf{R}^T \end{bmatrix} \\ (\mathbf{I}_{D_G} + \mathbf{G}^T \Sigma^{-1} \mathbf{G})^{-1} + \mathbf{R} \mathbf{P}, & \mathbf{R} \mathbf{P} & \mathbf{R} \mathbf{P} \\ \mathbf{R} \mathbf{P} & \ddots & \mathbf{R} \mathbf{P} \\ \mathbf{R} \mathbf{P} & \mathbf{R} \mathbf{P} & (\mathbf{I}_{D_G} + \mathbf{G}^T \Sigma^{-1} \mathbf{G})^{-1} + \mathbf{R} \mathbf{P} \end{array} \right]. \quad (191)$$

Using the above definition it is easy to note that if we want to work out  $E[\mathbf{y}_{ij} \mathbf{y}_{ij}^T | \tilde{\mathbf{x}}_i, \Theta]$  we would use the following subrepresentation of  $\left(\mathbf{I} + \tilde{\mathbf{A}}^T \tilde{\Sigma}^{-1} \tilde{\mathbf{A}}\right)^{-1}$ ,

$$\begin{bmatrix} \mathbf{T}_{ul} & \mathbf{T}_{ur} \\ \mathbf{T}_{ll} & \mathbf{T}_{lr} \end{bmatrix} = \begin{bmatrix} (\mathbf{I}_{D_F} + J_i \mathbf{F}^T \mathbf{Q} \mathbf{F})^{-1} & -\mathbf{R}^T \\ -\mathbf{R} & (\mathbf{I}_{D_G} + \mathbf{G}^T \Sigma^{-1} \mathbf{G})^{-1} + \mathbf{R} \mathbf{P} \end{bmatrix}. \quad (192)$$

## A Rules and Identities for Differentiation

In this section we provide some of the useful identities for differentiation.

$$\frac{\partial \mathbf{x}^T \mathbf{B} \mathbf{x}}{\partial \mathbf{x}} = (\mathbf{B} + \mathbf{B}^T) \mathbf{x}, \quad (193)$$

$$\frac{\partial \mathbf{a}^T \mathbf{X}^{-1} \mathbf{b}}{\partial \mathbf{X}} = -\mathbf{X}^{-T} \mathbf{a} \mathbf{b}^T \mathbf{X}^{-T}, \quad (194)$$

$$\frac{\partial (\det(\mathbf{X}^k))}{\partial \mathbf{X}} = k \det(\mathbf{X}^k) \mathbf{X}^{-T}, \quad (195)$$

$$\frac{\partial \ln |\det(\mathbf{X})|}{\partial \mathbf{X}} = (\mathbf{X}^{-1})^T = (\mathbf{X}^T)^{-1}, \quad (196)$$

$$\frac{\partial (\mathbf{x} - \mathbf{A} \mathbf{s})^T \mathbf{W} (\mathbf{x} - \mathbf{A} \mathbf{s})}{\partial \mathbf{s}} = -2 \mathbf{A}^T \mathbf{W} (\mathbf{x} - \mathbf{A} \mathbf{s}), \quad (197)$$

$$\frac{\partial (\mathbf{x} - \mathbf{A} \mathbf{s})^T \mathbf{W} (\mathbf{x} - \mathbf{A} \mathbf{s})}{\partial \mathbf{A}} = -2 \mathbf{W} (\mathbf{x} - \mathbf{A} \mathbf{s}) \mathbf{s}^T, \quad (198)$$

$$\frac{\partial (\mathbf{x} - \mathbf{s})^T \mathbf{W} (\mathbf{x} - \mathbf{s})}{\partial \mathbf{s}} = -2 \mathbf{W} (\mathbf{x} - \mathbf{s}). \quad (199)$$

Equations 197 and 198 are valid if the matrix  $W$  is symmetric. These identities come from the matrix cookbook [3].

## B Matrix Identities

### B.1 Inverse of a Product of Matrices

$$(\mathbf{A} \mathbf{B})^{-1} = \mathbf{B}^{-1} \mathbf{A}^{-1}, \quad (200)$$

### B.2 Block Matrix Inversion

The matrix inversion lemma which is an instance of using the Schur complement can be derived in several ways. First we present the solution that we are interested in below,

$$\begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix}^{-1} = \begin{bmatrix} (\mathbf{A} - \mathbf{B} \mathbf{D}^{-1} \mathbf{C})^{-1} & -(\mathbf{A} - \mathbf{B} \mathbf{D}^{-1} \mathbf{C})^{-1} \mathbf{B} \mathbf{D}^{-1} \\ -\mathbf{D}^{-1} \mathbf{C} (\mathbf{A} - \mathbf{B} \mathbf{D}^{-1} \mathbf{C})^{-1} & \mathbf{D}^{-1} + \mathbf{D}^{-1} \mathbf{C} (\mathbf{A} - \mathbf{B} \mathbf{D}^{-1} \mathbf{C})^{-1} \mathbf{B} \mathbf{D}^{-1} \end{bmatrix}. \quad (201)$$

This identity consists of several parts which occur frequently. These parts are,

$$(\mathbf{A} - \mathbf{B} \mathbf{D}^{-1} \mathbf{C})^{-1}, \quad (202)$$

$$(\mathbf{A} - \mathbf{B} \mathbf{D}^{-1} \mathbf{C})^{-1} \mathbf{B} \mathbf{D}^{-1}. \quad (203)$$

In addition to the above identity there is another useful intermediate form of this inversion identity. This intermediate form is obtained by using the LDU decomposition (B.4) and taking the inverse. We then use the identity given in Equation 200 to get,

$$\begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix}^{-1} = \left( \begin{bmatrix} \mathbf{I} & \mathbf{B} \mathbf{D}^{-1} \\ \mathbf{0} & \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{A} - \mathbf{B} \mathbf{D}^{-1} \mathbf{C} & \mathbf{0} \\ \mathbf{0} & \mathbf{D} \end{bmatrix} \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{D}^{-1} \mathbf{C} & \mathbf{I} \end{bmatrix} \right)^{-1}, \quad (204)$$

$$= \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{D}^{-1} \mathbf{C} & \mathbf{I} \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{A} - \mathbf{B} \mathbf{D}^{-1} \mathbf{C} & \mathbf{0} \\ \mathbf{0} & \mathbf{D} \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{I} & \mathbf{B} \mathbf{D}^{-1} \\ \mathbf{0} & \mathbf{I} \end{bmatrix}^{-1}. \quad (205)$$

Tackling each term separately we derive their inverse by solving the equation  $AA^{-1} = I$ , where  $A^{-1}$  is the unknown. Using this we get the following,

$$\begin{bmatrix} I, & BD^{-1} \\ 0, & I \end{bmatrix}^{-1} = \begin{bmatrix} I, & -BD^{-1} \\ 0, & I \end{bmatrix}, \quad (206)$$

$$\begin{bmatrix} A - BD^{-1}C, & 0 \\ 0, & D \end{bmatrix}^{-1} = \begin{bmatrix} (A - BD^{-1}C)^{-1}, & 0 \\ 0, & D^{-1} \end{bmatrix}, \quad (207)$$

and

$$\begin{bmatrix} I & 0 \\ D^{-1}C & I \end{bmatrix}^{-1} = \begin{bmatrix} I & 0 \\ -D^{-1}C & I \end{bmatrix}. \quad (208)$$

Finally this yields the intermediate form of,

$$\begin{bmatrix} A & B \\ C & D \end{bmatrix}^{-1} = \begin{bmatrix} I & 0 \\ -D^{-1}C & I \end{bmatrix} \begin{bmatrix} (A - BD^{-1}C)^{-1}, & 0 \\ 0, & D^{-1} \end{bmatrix} \begin{bmatrix} I, & -BD^{-1} \\ 0, & I \end{bmatrix}. \quad (209)$$

### B.3 Block LU Decomposition

For the LU decomposition we begin as follows. For the  $L$  part we want,

$$\begin{bmatrix} A & B \\ C & D \end{bmatrix} \begin{bmatrix} I & 0 \\ x_1 & x_2 \end{bmatrix} = \begin{bmatrix} A + Bx_1 & Bx_2 \\ C + Dx_1 & Dx_2 \end{bmatrix}, \quad (210)$$

with  $C + Dx_1 = 0$  and  $Dx_2 = I$ . To achieve this,

$$Dx_1 = -C, \quad (211)$$

$$x_1 = -D^{-1}C, \quad (212)$$

and

$$Dx_2 = I, \quad (213)$$

$$x_2 = D^{-1}. \quad (214)$$

First, we substitute for  $x_1$ , this would give us,

$$\begin{bmatrix} A & B \\ C & D \end{bmatrix} \begin{bmatrix} I, & 0 \\ -D^{-1}C, & x_2 \end{bmatrix} = \begin{bmatrix} A - BD^{-1}C, & Bx_2 \\ C - DD^{-1}C, & Dx_2 \end{bmatrix} = \begin{bmatrix} A - BD^{-1}C, & Bx_2 \\ 0, & Dx_2 \end{bmatrix}, \quad (215)$$

Second, we substitute for  $x_2$ , this would give us,

$$\begin{bmatrix} A & B \\ C & D \end{bmatrix} \begin{bmatrix} I, & 0 \\ -D^{-1}C, & D^{-1} \end{bmatrix} = \begin{bmatrix} A - BD^{-1}C, & BD^{-1} \\ 0, & DD^{-1} \end{bmatrix} = \begin{bmatrix} A - BD^{-1}C, & BD^{-1} \\ 0, & I \end{bmatrix}. \quad (216)$$

This would finally give us the LU decomposition.

## B.4 Block LDU Decomposition

For the LDU decomposition, the following identity holds,

$$\begin{bmatrix} A & B \\ C & D \end{bmatrix} = \begin{bmatrix} I, & BD^{-1} \\ 0, & I \end{bmatrix} \begin{bmatrix} A - BD^{-1}C, & 0 \\ 0, & D \end{bmatrix} \begin{bmatrix} I, & 0 \\ D^{-1}C, & I \end{bmatrix}. \quad (217)$$

This can be proven by direct computation, noting that,

$$\begin{bmatrix} I, & BD^{-1} \\ 0, & I \end{bmatrix} \begin{bmatrix} A - BD^{-1}C, & 0 \\ 0, & D \end{bmatrix} \begin{bmatrix} I, & 0 \\ D^{-1}C, & I \end{bmatrix} = \begin{bmatrix} I, & BD^{-1} \\ 0, & I \end{bmatrix} \begin{bmatrix} A - BD^{-1}C, & 0 \\ C, & D \end{bmatrix}, \quad (218)$$

$$\begin{bmatrix} I, & BD^{-1} \\ 0, & I \end{bmatrix} \begin{bmatrix} A - BD^{-1}C, & 0 \\ C, & D \end{bmatrix} = \begin{bmatrix} A - BD^{-1}C + BD^{-1}C, & B \\ C, & D \end{bmatrix} = \begin{bmatrix} A, & B \\ C, & D \end{bmatrix}. \quad (219)$$

## B.5 Square Root of a Matrix

In general, a matrix can have several square root. However, a positive-definite matrix has precisely one positive-definite square root. The square root of a matrix can be found by first diagonalising it (when this is possible). The process is described below, but can also be found in more detail here [[http://en.wikipedia.org/wiki/Square\\_root\\_of\\_a\\_matrix](http://en.wikipedia.org/wiki/Square_root_of_a_matrix) (By diagonalization)]. If we consider a matrix  $M$  that has the following decomposition,  $D$  being diagonal,

$$M = ZDZ^{-1}, \quad (220)$$

and if we find the square root  $D^{\frac{1}{2}}$  of the diagonal matrix  $D$ , which means that  $D = D^{\frac{1}{2}}D^{\frac{1}{2}}$ , then we can compute a square root  $M^{\frac{1}{2}}$  of  $M$  as follows,

$$M^{\frac{1}{2}} = ZD^{\frac{1}{2}}Z^{-1}. \quad (221)$$

## B.6 Log Determinant of a Diagonal Matrix

The log determinant of a diagonal matrix is easy to compute. Assuming a diagonal matrix  $\Sigma$ , we have,

$$\ln |\det(\Sigma)| = \ln \left| \prod_{i=1}^I \sigma_{ii} \right|. \quad (222)$$

We can change this product around to be the sum of a log. This gives us,

$$\ln |\det(\Sigma)| = \sum_{i=1}^I \ln |\sigma_{ii}|. \quad (223)$$

## B.7 Log Determinant of a Symmetric Real Matrix

A symmetric real matrix is diagonalisable by orthogonal matrix, i.e., given a real symmetric matrix  $A$ ,  $Q^T A Q$  is diagonal for some orthogonal matrix  $Q$ . Furthermore, the log determinant of a symmetric real matrix can be found using its eigenvalue decomposition. Assuming such a symmetric real matrix  $A$ , its eigenvalue decomposition is given by,

$$A = Q\Lambda Q^T, \quad (224)$$

where  $\Lambda$  is the diagonal matrix containing the eigenvalue of  $A$ , and  $Q$  is orthogonal, which means that  $Q Q^T = I$ . Futhermore,

$$\ln |\det(A)| = \ln |\det(Q\Lambda Q^T)| = \ln |\det(\Lambda) \det(Q Q^T)| = \ln |\det(\Lambda) \det(I)| = \ln |\det(\Lambda)| \quad (225)$$

There is a built-in function called `slogdet()`, shipped with recent versions of SciPy/NumPy, which performs such a computation.

## B.8 Inverse of a Block Diagonal Matrix

The inverse of a block diagonal matrix is a block diagonal matrix of the inverse blocks. That does not make a lot of sense so let us write it down,

$$\begin{bmatrix} A_1 & 0 & 0 & 0 & 0 \\ 0 & A_2 & 0 & 0 & 0 \\ 0 & 0 & A_3 & 0 & 0 \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ 0 & 0 & 0 & 0 & A_n \end{bmatrix}^{-1} = \begin{bmatrix} A_1^{-1} & 0 & 0 & 0 & 0 \\ 0 & A_2^{-1} & 0 & 0 & 0 \\ 0 & 0 & A_3^{-1} & 0 & 0 \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ 0 & 0 & 0 & 0 & A_n^{-1} \end{bmatrix}, \quad (226)$$

where  $[A_1, A_2, A_3, \dots, A_n]$  are a set of block matrices.

## C Intermediate Solutions for PLDA: Gaussian Priors

In this appendix we provide several intermediate solutions and representations that we use for PLDA when we have Gaussian priors.

### C.1 Matrix Inversion: $(\mathbf{I} + \tilde{\mathbf{A}}^T \tilde{\Sigma}^{-1} \tilde{\mathbf{A}})^{-1}$

The matrix  $(\mathbf{I} + \tilde{\mathbf{A}}^T \tilde{\Sigma}^{-1} \tilde{\mathbf{A}})^{-1}$  is central to many calculations for PLDA with a Gaussian prior. We provide a series of simplifications that can be used. We will tackle all of these in blocks or parts. To derive simplifications using this matrix we use its explicit form given by the block matrix inversion identity which we covered in the previous section.

For PLDA with a Gaussian prior the block matrix to invert is defined by the following,

$$A = [\mathbf{I}_{D_F} + J_i \mathbf{F}^T \Sigma^{-1} \mathbf{F}], \quad (227)$$

$$B = [ \mathbf{F}^T \Sigma^{-1} \mathbf{G}, \quad \mathbf{F}^T \Sigma^{-1} \mathbf{G}, \quad \mathbf{F}^T \Sigma^{-1} \mathbf{G} ], \quad (228)$$

$$C = \begin{bmatrix} \mathbf{G}^T \Sigma^{-1} \mathbf{F} \\ \mathbf{G}^T \Sigma^{-1} \mathbf{F} \\ \mathbf{G}^T \Sigma^{-1} \mathbf{F} \end{bmatrix} = B^T, \quad (229)$$

$$D = \begin{bmatrix} (\mathbf{I}_{D_G} + \mathbf{G}^T \Sigma^{-1} \mathbf{G}) & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & (\mathbf{I}_{D_G} + \mathbf{G}^T \Sigma^{-1} \mathbf{G}) & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & (\mathbf{I}_{D_G} + \mathbf{G}^T \Sigma^{-1} \mathbf{G}) \end{bmatrix}, \quad (230)$$

$$D^{-1} = \begin{bmatrix} (\mathbf{I}_{D_G} + \mathbf{G}^T \Sigma^{-1} \mathbf{G})^{-1} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & (\mathbf{I}_{D_G} + \mathbf{G}^T \Sigma^{-1} \mathbf{G})^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & (\mathbf{I}_{D_G} + \mathbf{G}^T \Sigma^{-1} \mathbf{G})^{-1} \end{bmatrix}. \quad (231)$$

$$BD^{-1} = [ \mathbf{F}^T \Sigma^{-1} \mathbf{G} (\mathbf{I}_{D_G} + \mathbf{G}^T \Sigma^{-1} \mathbf{G})^{-1}, \quad \dots, \quad \mathbf{F}^T \Sigma^{-1} \mathbf{G} (\mathbf{I}_{D_G} + \mathbf{G}^T \Sigma^{-1} \mathbf{G})^{-1} ]. \quad (232)$$

$$BD^{-1} B^T = BD^{-1} \begin{bmatrix} \mathbf{G}^T \Sigma^{-1} \mathbf{F} \\ \mathbf{G}^T \Sigma^{-1} \mathbf{F} \\ \mathbf{G}^T \Sigma^{-1} \mathbf{F} \end{bmatrix}, \quad (233)$$

$$= [ J_i ( \mathbf{F}^T \Sigma^{-1} \mathbf{G} (\mathbf{I}_{D_G} + \mathbf{G}^T \Sigma^{-1} \mathbf{G})^{-1} \mathbf{G}^T \Sigma^{-1} \mathbf{F} ) ] = BD^{-1} C. \quad (234)$$

### C.1.1 Finding $(A - BD^{-1}C)^{-1}$

We now define the inverted matrix  $(A - BD^{-1}C)^{-1}$  as this is used in several places.

$$(A - BD^{-1}C)^{-1} = \left[ \mathbf{I}_{D_F} + J_i \mathbf{F}^T \boldsymbol{\Sigma}^{-1} \mathbf{F} - J_i \left( \mathbf{F}^T \boldsymbol{\Sigma}^{-1} \mathbf{G} (\mathbf{I}_{D_G} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G})^{-1} \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{F} \right) \right]^{-1}, \quad (235)$$

$$= \left[ \mathbf{I}_{D_F} + J_i \left( \mathbf{F}^T \boldsymbol{\Sigma}^{-1} \mathbf{F} - \mathbf{F}^T \boldsymbol{\Sigma}^{-1} \mathbf{G} (\mathbf{I}_{D_G} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G})^{-1} \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{F} \right) \right]^{-1}, \quad (236)$$

$$= \left[ \mathbf{I}_{D_F} + J_i \mathbf{F}^T \left( \boldsymbol{\Sigma}^{-1} - \boldsymbol{\Sigma}^{-1} \mathbf{G} (\mathbf{I}_{D_G} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G})^{-1} \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \right) \mathbf{F} \right]^{-1}, \quad (237)$$

$$= (\mathbf{I}_{D_F} + J_i \mathbf{F}^T \mathbf{Q} \mathbf{F})^{-1}. \quad (238)$$

Where,

$$\mathbf{Q} = \boldsymbol{\Sigma}^{-1} - \boldsymbol{\Sigma}^{-1} \mathbf{G} (\mathbf{I}_{D_G} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G})^{-1} \mathbf{G}^T \boldsymbol{\Sigma}^{-1}. \quad (239)$$

### C.1.2 Finding $-D^{-1}C(A - BD^{-1}C)^{-1}$

We now need to find the final form for the matrix  $-D^{-1}C(A - BD^{-1}C)^{-1}$ . This is given by,

$$-D^{-1}C(A - BD^{-1}C)^{-1} = -D^{-1}C(\mathbf{I}_{D_F} + J_i \mathbf{F}^T \mathbf{Q} \mathbf{F})^{-1}, \quad (240)$$

$$= -D^{-1} \begin{bmatrix} \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{F} \\ \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{F} \\ \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{F} \end{bmatrix} (\mathbf{I}_{D_F} + J_i \mathbf{F}^T \mathbf{Q} \mathbf{F})^{-1}, \quad (241)$$

$$= -D^{-1} \begin{bmatrix} \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{F} (\mathbf{I}_{D_F} + J_i \mathbf{F}^T \mathbf{Q} \mathbf{F})^{-1} \\ \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{F} (\mathbf{I}_{D_F} + J_i \mathbf{F}^T \mathbf{Q} \mathbf{F})^{-1} \\ \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{F} (\mathbf{I}_{D_F} + J_i \mathbf{F}^T \mathbf{Q} \mathbf{F})^{-1} \end{bmatrix}, \quad (242)$$

$$= - \begin{bmatrix} (\mathbf{I}_{D_G} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G})^{-1} & \dots & \mathbf{0} \\ \vdots & \ddots & \vdots \\ \mathbf{0} & \dots & (\mathbf{I}_{D_G} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G})^{-1} \end{bmatrix} \begin{bmatrix} \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{F} (\mathbf{I}_{D_F} + J_i \mathbf{F}^T \mathbf{Q} \mathbf{F})^{-1} \\ \vdots \\ \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{F} (\mathbf{I}_{D_F} + J_i \mathbf{F}^T \mathbf{Q} \mathbf{F})^{-1} \end{bmatrix}, \quad (243)$$

$$= - \begin{bmatrix} (\mathbf{I}_{D_G} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G})^{-1} \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{F} (\mathbf{I}_{D_F} + J_i \mathbf{F}^T \mathbf{Q} \mathbf{F})^{-1} \\ (\mathbf{I}_{D_G} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G})^{-1} \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{F} (\mathbf{I}_{D_F} + J_i \mathbf{F}^T \mathbf{Q} \mathbf{F})^{-1} \\ (\mathbf{I}_{D_G} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G})^{-1} \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{F} (\mathbf{I}_{D_F} + J_i \mathbf{F}^T \mathbf{Q} \mathbf{F})^{-1} \end{bmatrix}. \quad (244)$$

### C.1.3 Finding $D^{-1} + D^{-1}C(A - BD^{-1}C)^{-1}BD^{-1}$

The last term that we are interested in is in the bottom right hand corner which is given by  $D^{-1} + D^{-1}C(A - BD^{-1}C)^{-1}BD^{-1}$ . We have the following expression for this matrix,

$$D^{-1} + D^{-1}C(A - BD^{-1}C)^{-1}BD^{-1} = D^{-1} - D^{-1}C(A - BD^{-1}C)^{-1}(-BD^{-1}), \quad (245)$$

$$= D^{-1} + D^{-1}C(A - BD^{-1}C)^{-1}BD^{-1}, \quad (246)$$

$$= D^{-1} + \begin{bmatrix} (\mathbf{I}_{D_G} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G})^{-1} \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{F} (\mathbf{I}_{D_F} + J_i \mathbf{F}^T \mathbf{Q} \mathbf{F})^{-1} \\ (\mathbf{I}_{D_G} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G})^{-1} \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{F} (\mathbf{I}_{D_F} + J_i \mathbf{F}^T \mathbf{Q} \mathbf{F})^{-1} \\ (\mathbf{I}_{D_G} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G})^{-1} \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{F} (\mathbf{I}_{D_F} + J_i \mathbf{F}^T \mathbf{Q} \mathbf{F})^{-1} \end{bmatrix} B D^{-1}, \quad (247)$$

$$= D^{-1} + \begin{bmatrix} (\mathbf{I}_{D_G} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G})^{-1} \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{F} (\mathbf{I}_{D_F} + J_i \mathbf{F}^T \mathbf{Q} \mathbf{F})^{-1} \\ (\mathbf{I}_{D_G} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G})^{-1} \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{F} (\mathbf{I}_{D_F} + J_i \mathbf{F}^T \mathbf{Q} \mathbf{F})^{-1} \\ (\mathbf{I}_{D_G} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G})^{-1} \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{F} (\mathbf{I}_{D_F} + J_i \mathbf{F}^T \mathbf{Q} \mathbf{F})^{-1} \end{bmatrix} \\ \left[ \mathbf{F}^T \boldsymbol{\Sigma}^{-1} \mathbf{G} (\mathbf{I}_{D_G} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G})^{-1}, \dots, \mathbf{F}^T \boldsymbol{\Sigma}^{-1} \mathbf{G} (\mathbf{I}_{D_G} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G})^{-1} \right], \quad (248)$$

We now need to make some substitutions. Let,

$$\mathbf{P} = \mathbf{F}^T \boldsymbol{\Sigma}^{-1} \mathbf{G} (\mathbf{I}_{D_G} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G})^{-1}, \quad (249)$$

$$\mathbf{R} = (\mathbf{I}_{D_G} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G})^{-1} \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{F} (\mathbf{I}_{D_F} + J_i \mathbf{F}^T \mathbf{Q} \mathbf{F})^{-1}, \quad (250)$$

Substituting these back in gives us,

$$D^{-1} + D^{-1} C (A - B D^{-1} C)^{-1} B D^{-1} = D^{-1} + \begin{bmatrix} \mathbf{R} \\ \mathbf{R} \\ \mathbf{R} \end{bmatrix} \begin{bmatrix} \mathbf{P} & \mathbf{P} & \mathbf{P} \end{bmatrix}, \quad (251)$$

$$= D^{-1} + \begin{bmatrix} \mathbf{R} \mathbf{P} & \mathbf{R} \mathbf{P} & \mathbf{R} \mathbf{P} \\ \mathbf{R} \mathbf{P} & \mathbf{R} \mathbf{P} & \mathbf{R} \mathbf{P} \\ \mathbf{R} \mathbf{P} & \mathbf{R} \mathbf{P} & \mathbf{R} \mathbf{P} \end{bmatrix}, \quad (252)$$

$$= \begin{bmatrix} (\mathbf{I}_{D_G} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G})^{-1} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & (\mathbf{I}_{D_G} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G})^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & (\mathbf{I}_{D_G} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G})^{-1} \end{bmatrix} + \begin{bmatrix} \mathbf{R} \mathbf{P} & \mathbf{R} \mathbf{P} & \mathbf{R} \mathbf{P} \\ \mathbf{R} \mathbf{P} & \mathbf{R} \mathbf{P} & \mathbf{R} \mathbf{P} \\ \mathbf{R} \mathbf{P} & \mathbf{R} \mathbf{P} & \mathbf{R} \mathbf{P} \end{bmatrix}, \quad (253)$$

$$= \begin{bmatrix} (\mathbf{I}_{D_G} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G})^{-1} + \mathbf{R} \mathbf{P} & \mathbf{R} \mathbf{P} & \mathbf{R} \mathbf{P} \\ \mathbf{R} \mathbf{P} & (\mathbf{I}_{D_G} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G})^{-1} + \mathbf{R} \mathbf{P} & \mathbf{R} \mathbf{P} \\ \mathbf{R} \mathbf{P} & \mathbf{R} \mathbf{P} & (\mathbf{I}_{D_G} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G})^{-1} + \mathbf{R} \mathbf{P} \end{bmatrix}. \quad (254)$$

We now define each entry by noting that,

$$\mathbf{R} \mathbf{P} = (\mathbf{I}_{D_G} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G})^{-1} \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{F} (\mathbf{I}_{D_F} + J_i \mathbf{F}^T \mathbf{Q} \mathbf{F})^{-1} \mathbf{F}^T \boldsymbol{\Sigma}^{-1} \mathbf{G} (\mathbf{I}_{D_G} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G})^{-1}, \quad (255)$$

and

$$(\mathbf{I}_{D_G} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G})^{-1} + \mathbf{R} \mathbf{P} = \\ (\mathbf{I}_{D_G} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G})^{-1} + (\mathbf{I}_{D_G} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G})^{-1} \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{F} (\mathbf{I}_{D_F} + J_i \mathbf{F}^T \mathbf{Q} \mathbf{F})^{-1} \mathbf{F}^T \boldsymbol{\Sigma}^{-1} \mathbf{G} (\mathbf{I}_{D_G} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G})^{-1}, \quad (256)$$

$$= (\mathbf{I}_{D_G} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G})^{-1} \left[ \mathbf{I}_{D_F} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{F} (\mathbf{I}_{D_F} + J_i \mathbf{F}^T \mathbf{Q} \mathbf{F})^{-1} \mathbf{F}^T \boldsymbol{\Sigma}^{-1} \mathbf{G} (\mathbf{I}_{D_G} + \mathbf{G}^T \boldsymbol{\Sigma}^{-1} \mathbf{G})^{-1} \right]. \quad (257)$$

## C.2 LDU Decomposition for $(\mathbf{I} + \tilde{\mathbf{A}}^T \tilde{\Sigma}^{-1} \tilde{\mathbf{A}})^{-1}$

If we consider the LDU decomposition for the matrix inversion we need to define the subparts for three block matrices. These matrices are given in Equation 209,

$$\begin{bmatrix} A & B \\ C & D \end{bmatrix}^{-1} = \begin{bmatrix} I & 0 \\ -D^{-1}C & I \end{bmatrix} \begin{bmatrix} (A - BD^{-1}C)^{-1} & 0 \\ 0 & D^{-1} \end{bmatrix} \begin{bmatrix} I & -BD^{-1} \\ 0 & I \end{bmatrix}. \quad (258)$$

We now define each of three matrices below,

$$\begin{bmatrix} I & 0 \\ -D^{-1}C & I \end{bmatrix} = \begin{bmatrix} \mathbf{I}_{D_F}, & \mathbf{0}_{(D_F, J_i D_G)} \\ - \begin{bmatrix} (\mathbf{I}_{D_G} + \mathbf{G}^T \Sigma^{-1} \mathbf{G})^{-1} \mathbf{G}^T \Sigma^{-1} \mathbf{F} \\ (\mathbf{I}_{D_G} + \mathbf{G}^T \Sigma^{-1} \mathbf{G})^{-1} \mathbf{G}^T \Sigma^{-1} \mathbf{F} \\ (\mathbf{I}_{D_G} + \mathbf{G}^T \Sigma^{-1} \mathbf{G})^{-1} \mathbf{G}^T \Sigma^{-1} \mathbf{F} \end{bmatrix} & \mathbf{I}_{J_i D_G} \end{bmatrix}, \quad (259)$$

$$\begin{bmatrix} (A - BD^{-1}C)^{-1} & 0 \\ 0 & D^{-1} \end{bmatrix} = \begin{bmatrix} (\mathbf{I}_{D_F} + J_i \mathbf{F}^T \mathbf{Q} \mathbf{F})^{-1} & \mathbf{0}_{(D_F, J_i D_G)} \\ \mathbf{0}_{(J_i D_G, D_F)}, & \begin{bmatrix} (\mathbf{I}_{D_G} + \mathbf{G}^T \Sigma^{-1} \mathbf{G})^{-1} & \dots & \mathbf{0} \\ \vdots & \ddots & \vdots \\ \mathbf{0} & \dots & (\mathbf{I}_{D_G} + \mathbf{G}^T \Sigma^{-1} \mathbf{G})^{-1} \end{bmatrix} \end{bmatrix}, \quad (260)$$

and

$$\begin{bmatrix} I & -BD^{-1} \\ 0 & I \end{bmatrix} = \begin{bmatrix} \mathbf{I}_{D_F}, & - \begin{bmatrix} \mathbf{F}^T \Sigma^{-1} \mathbf{G} (\mathbf{I}_{D_G} + \mathbf{G}^T \Sigma^{-1} \mathbf{G})^{-1}, & \dots, & \mathbf{F}^T \Sigma^{-1} \mathbf{G} (\mathbf{I}_{D_G} + \mathbf{G}^T \Sigma^{-1} \mathbf{G})^{-1} \end{bmatrix} \\ \mathbf{0}_{(J_i D_G, D_F)}, & \mathbf{I}_{J_i D_G} \end{bmatrix}. \quad (261)$$

## References

- [1] C. M. Bishop. *Pattern Recognition and Machine Learning*. Springer, 2006.
- [2] P. Li and S. J. D. Prince. *Advance in Face Image Analysis: Techniques and Technologies*, chapter Probabilistic Methods for Face Registration and Recognition (In Press). Idea Group Publishing, 2010.
- [3] K. B. Petersen and M. S. Pedersen. *The Matrix Cookbook, Version: November 14, 2008*. Technical University of Denmark, 2008.
- [4] S. J. D. Prince and J. H. Elder. Probabilistic linear discriminant analysis for inferences about identity. In *IEEE International Conference on Computer Vision*, pages 1–8, 2007.