

Analysis of Robust Soft Learning Vector Quantization

J.J.G. de Vries

J.J.G.de.Vries@student.rug.nl

Institute of Mathematics and Computing Science,
University of Groningen.

Abstract. One of the popular methods for multiclass classification is Learning Vector Quantization (LVQ). There have been developed several variants of LVQ lately, among which Robust Soft Learning Vector Quantization, or RSLVQ for short. An introductory study showed that RSLVQ performs better than other LVQ algorithms, even very close to the optimal linear classifier, within a controlled environment. In order to study its performance in detail, we performed a mathematical analysis of the algorithm, in the form of a system of coupled Ordinary Differential Equations (ODE's), which might also help development of an optimal LVQ algorithm. Following from our analysis, we compare the potential performance of RSLVQ in relation to other LVQ variants and present a guideline for settings of the control parameter, i.e. the softness parameter.

1 Introduction

Learning Vector Quantization (LVQ), originally posed by Kohonen [4, 3] and known by the name LVQ1, is a method of online supervised competitive learning. Many variations on the basic scheme of LVQ1 have been suggested, among which LVQ2.1 and LVQ3 [3, 5], GLVQ [7] and RSLVQ [8, 9] with the aim of obtaining better generalization behavior.

During learning, data samples and their class labels are presented sequentially, or so called 'on-line'. From a set of prototype vectors, defined in the same (potentially high dimensional) space as the data, the closest (set of) prototype(s) is determined and updated such that if the class label coincides with the class label of the data sample, the prototype is attracted to the data, otherwise repelled. The data, carrying labels of different classes, is assumed to be distributed around a specified number of prototypes. Note that there can be more than one prototype per class, enabling a good fit of prototypes to data that contains highly complex class boundaries. Hopefully the prototypes represent the data well after they have settled and training is finished. Classification now can be done by determining the closest of all prototypes and returning the class label corresponding to this winning prototype. Therefore the decision boundaries between the prototypes can also be seen as the Voronoi tessellation of the feature space.

There are a lot of variations of LVQ-algorithms, which mainly differ in which specific prototypes are updated (for example only the closest conflicting prototypes or the closest prototype with corresponding and closest with conflicting label) and how these prototypes are updated. The generic structure of an LVQ algorithm can be expressed in the following way:

$$\begin{aligned}
 w_l^\mu &= w_l^{\mu-1} + \Delta w_l^\mu, \\
 &= w_l^{\mu-1} + \frac{\eta}{N} f(\{w_i^{\mu-1}\}, \xi^\mu, \sigma^\mu)(\xi^\mu - w_l^{\mu-1}) \quad (1) \\
 &\text{with } l, i = 1 \dots c, \mu = 1, 2, \dots
 \end{aligned}$$

Where w_l^μ is prototype w_l at time step μ , η is the so-called learning rate, N is the dimensionality of the system. The specific form of f is determined by the algorithm used to perform LVQ. The softness of RSLVQ determines the extent to which correctly classified example data (i.e. the closest prototype has a coinciding class label) causes an update of the prototypes. The hard or crisp variant of RSLVQ, in which the softness is taken in the limit to 0, called Learning From Mistakes (LFM) has been analyzed mathematically by Biehl, Ghosh and Hammer [1, 2] in a controlled environment. The study described in this paper is an extension of their study and analyzes truly RSLVQ using the same modelled environment and forms the mathematical background of findings in an introductory study on the performance of RSLVQ [10].

This document is organized as follows: section 2 describes the model in which RSLVQ is analyzed, after which section 3 provides a detailed description of the RSLVQ algorithm. Section 4 describes the analysis globally and section 5 shows the experiments and results of the analysis. Finally section 6 concludes the paper and section 7 gives an overview of future work. The detailed mathematical analysis is attached in appendices A and B.

2 Model

To be able to analyze RSLVQ, we need to restrict the model in which we observe RSLVQ. Biehl et al. [1] defined the following model consisting of high dimensional data originating from a mixture of two overlapping Gaussian clusters, of two classes. We assume that data vectors $\xi \in \mathbb{R}^N$ of class $\sigma \in \{\pm\}$ are drawn independently, with probability $P(\xi)$, according to the following distribution:

$$P(\xi) = \sum_{\sigma=\pm 1} p_\sigma P(\xi|\sigma) \quad (2)$$

with

$$P(\xi|\sigma) = \frac{1}{(2\pi v_\sigma)^{N/2}} \exp\left(-\frac{1}{2v_\sigma}(\xi - \lambda B_\sigma)^2\right) \quad (3)$$

The Gaussian clusters are centered around $\lambda \mathbf{B}_\sigma$ with variances v_σ . The prior probabilities p_σ of both classes ($\sigma \in \{+1, -1\}$ or $\{+, -\}$ for short) satisfy $p_+ + p_- = 1$. The vectors \mathbf{B}_σ are chosen to be orthonormal, i.e. $\mathbf{B}_+^2 = \mathbf{B}_-^2 = 1$ and $\mathbf{B}_+ \cdot \mathbf{B}_- = 0$, so λ specifies the distance (i.e. $\lambda\sqrt{2}$) between the cluster centers. Note that, since λ is chosen such that the clusters overlap, the classification task is clearly not linearly separable. The data points $\{\xi^\mu, \sigma^\mu\}$ are now presented sequentially so that at each time step $\mu = 1, 2, \dots$ a new uncorrelated vector ξ^μ , along with its label σ^μ , independently drawn according to the density (3), is presented.

There will be fit two prototypes to these clusters, each representing one of the two classes, i.e.:

$$\mathbf{w}_s^\mu \in \mathbb{R}^N \text{ with } s \in \{\pm 1\}, \mu = 1, 2, \dots \quad (4)$$

2.1 Characteristic Quantities

A set of suitable order parameters or characteristic quantities that describe the system has been found by Biehl et al. [1] to be the following:

$$\begin{aligned} R_{S\sigma}^\mu &= \mathbf{w}_S^\mu \cdot \mathbf{B}_\sigma \\ Q_{ST}^\mu &= \mathbf{w}_S^\mu \cdot \mathbf{w}_T^\mu \end{aligned} \quad (5)$$

with $\sigma, S, T \in \{\pm 1\}, \mu = 1, 2, \dots$

The self-overlaps Q_{++}^μ, Q_{--}^μ and the symmetric cross-overlap $Q_{+-}^\mu = Q_{-+}^\mu$ relate to the lengths and relative angle between the prototype vectors. The quantities $R_{+\sigma}^\mu$ and $R_{-\sigma}^\mu$ specify the projections of the prototype vectors into the plane spanned by the vectors \mathbf{B}_σ . These characteristic quantities have also been found to express the generalization error in the following way:

$$\begin{aligned} \epsilon_g &= p_+ \epsilon_+ + p_- \epsilon_- \\ \epsilon_\sigma &= \phi \left(\frac{Q_{\sigma\sigma} - Q_{-\sigma-\sigma} - 2\lambda(R_{\sigma\sigma} - R_{-\sigma\sigma})}{2\sqrt{v_\sigma} \sqrt{Q_{++} - 2Q_{+-} + Q_{--}}} \right) \end{aligned} \quad (6)$$

with $\phi(z) = \int_{-\infty}^z \frac{\exp(-x^2/2)}{\sqrt{2\pi}} dx$

Note that for all time steps $\mu = 1, 2, \dots$ these quantities can be determined, resulting in the learning curve or the typical generalization error $\epsilon_g(\alpha)$ after on-line training with $\mu = \alpha N$ random examples.

Details of the calculation can be found, as well as more information on the characteristic quantities, in appendix A.

3 Robust Soft Learning Vector Quantization

As proposed by Seo and Obermayer [9] RSLVQ is a generic algorithm in which different assumptions on the distribution of the data can be made. Note that the assumption about the data distribution to train the prototypes within the LVQ algorithm might differ from the true data distribution used in the model, just as with real world applications in which the true data distribution mostly is not known beforehand. RSLVQ is defined by the following extension of the general update step (1):

$$\begin{aligned} w_i^\mu &= w_i^{\mu-1} + \Delta w_i^\mu \\ &= w_i^{\mu-1} + \tilde{\eta} \begin{cases} (P_l(\bar{l}|\xi^\mu) - P(\bar{l}|\xi^\mu)) \frac{\partial f(\xi, w_i^\mu)}{\partial w_i^\mu} & \text{if } l = \sigma^\mu \\ -P(\bar{l}|\xi^\mu) \frac{\partial f(\xi, w_i^\mu)}{\partial w_i^\mu} & \text{if } l \neq \sigma^\mu \end{cases} \end{aligned} \quad (7)$$

where $P_l(\bar{l}|\xi^\mu)$ and $P(\bar{l}|\xi^\mu)$ are assignment probabilities:

$$\begin{aligned} P_l(\bar{l}|\xi^\mu) &= \frac{p(\bar{l}) \exp(f(\xi, w_i^\mu))}{\sum_{\bar{l}=\sigma^\mu} p(\bar{l}) \exp(f(\xi, w_i^\mu))} \\ P(\bar{l}|\xi^\mu) &= \frac{p(\bar{l}) \exp(f(\xi, w_i^\mu))}{\sum_{\bar{l}} p(\bar{l}) \exp(f(\xi, w_i^\mu))} \end{aligned} \quad (8)$$

$P_l(\bar{l}|\xi^\mu)$ describes the posterior probability that the data sample $\{\xi^\mu, \sigma^\mu\}$ is assigned to prototype w_i of class l , given that the data sample was generated by the correct class. $P(\bar{l}|\xi^\mu)$ describes the posterior probability that the data sample is assigned to prototype w_i of all prototypes of all classes. $f(\xi, w_i^\mu)$ describes the assumed distribution of the data around the prototypes in such a way that $K(\bar{l}) \exp(f(\xi, w_i^\mu))$ gives the probability that the data vector ξ^μ is assigned to prototype w_i .

In our settings we assume a Gaussian distribution, i.e. $K(\bar{l}) = (2\pi v_{\bar{l}})^{N/2}$ and $f(\xi, w_i^\mu) = -(\xi - w_i^\mu)^2 / 2v_{\bar{l}}$, implying $\frac{\partial f(\xi, w_i^\mu)}{\partial w_i^\mu} = \frac{\xi - w_i^\mu}{v_{\bar{l}}}$, and the prototypes w_S all have the same width and strength, i.e. the variances and priors are equal: $\forall S :: v_S = v_{soft}, p(S) = \frac{1}{\#w_S}$.

With these assumptions the update rule (7) becomes:

$$w_i^\mu = w_i^{\mu-1} + \frac{\tilde{\eta}}{v_{soft}} \begin{cases} (P_l(\bar{l}|\xi^\mu) - P(\bar{l}|\xi^\mu))(\xi - w_i^\mu) & \text{if } l = \sigma^\mu \\ -P(\bar{l}|\xi^\mu)(\xi - w_i^\mu) & \text{if } l \neq \sigma^\mu \end{cases} \quad (9)$$

where

$$P_l(\bar{l}|\xi^\mu) = \frac{\exp(-(\xi - w_i^\mu)^2 / 2v_{soft})}{\sum_{\bar{l}=\sigma^\mu} \exp(-(\xi - w_i^\mu)^2 / 2v_{soft})}$$

$$P(\bar{l}|\xi^\mu) = \frac{\exp(-(\xi - w_{\bar{l}}^\mu)^2/2v_{soft})}{\sum_{\bar{l}} \exp(-(\xi - w_{\bar{l}}^\mu)^2/2v_{soft})} \quad (10)$$

Since in our model only two prototypes, each representing a different class, are used, equation (10) can be written as:

$$\begin{aligned} P_l(l|\xi^\mu) &= \frac{\exp(-(\xi - w_l^\mu)^2/2v_{soft})}{\exp(-(\xi - w_l^\mu)^2/2v_{soft})} = 1 \\ P(l|\xi^\mu) &= \frac{\exp(-(\xi - w_l^\mu)^2/2v_{soft})}{\exp(-(\xi - w_+^\mu)^2/2v_{soft}) + \exp(-(\xi - w_-^\mu)^2/2v_{soft})} \\ &= \frac{1}{1 + \exp\left(\frac{((\xi - w_+^\mu)^2 - (\xi - w_-^\mu)^2)/2v_{soft}}{1}\right)} \end{aligned} \quad (11)$$

Putting this back into equation (9) yields:

$$\begin{aligned} w_l^\mu &= w_l^{\mu-1} + \frac{\bar{\eta}}{v_{soft}} \left(\delta_{l\sigma^\mu} - \frac{1}{1 + \exp\left(\frac{d_l^\mu - d_{-l}^\mu}{2v_{soft}}\right)} \right) (\xi^\mu - w_l^{\mu-1}) \\ &\text{with } d_l^\mu = (\xi - w_l^\mu)^2 \end{aligned} \quad (12)$$

It is however more convenient to rescale the learning rate with the dimensionality of the data (N), so we rewrite:

$$w_l^\mu = w_l^{\mu-1} + \frac{\eta}{Nv_{soft}} (\delta_{l\sigma^\mu} - \Omega_l) (\xi^\mu - w_l^{\mu-1}) \quad (13)$$

Where $\delta_{l\sigma^\mu}$ is the Kronecker delta and

$$\Omega_l = \frac{1}{1 + \exp\left(\frac{d_l^\mu - d_{-l}^\mu}{2v_{soft}}\right)} \quad (14)$$

However, in the form of equation (12), the update function cannot be integrated analytically, as we would like in the further analysis. As an alternative route we approximate the update (12) by an LVQ variant which facilitates further analytic treatment. We use the observation that $\frac{1}{1+\exp(x)}$ is very similar to $\Phi\left(\frac{-x}{c}\right)$, where $c \in \mathbb{R}$ is a constant which controls the slope of the Φ -function and we rewrite:

$$w_l^\mu = w_l^{\mu-1} + \frac{\eta}{Nv_{soft}} (\delta_{l\sigma^\mu} - \Phi_l) (\xi^\mu - w_l^{\mu-1}) \quad (15)$$

Where

$$\Phi_l = \Phi\left(\frac{d_{-l}^\mu - d_l^\mu}{2cv_{soft}}\right) \quad (16)$$

$$\Phi(z) = \int_{-\infty}^z \frac{e^{-x^2/2}}{\sqrt{2\pi}} dx \quad (17)$$

To obtain the value of c let us set equal the slope at $x = 0$, therefore observe the derivative of both activation functions:

$$\begin{aligned}\frac{d}{dx} \frac{1}{1+e^x} &= -\frac{e^x}{(1+e^x)^2} \\ \frac{d}{dx} \Phi\left(\frac{-x}{c}\right) &= -\frac{e^{-\frac{x^2}{2c^2}}}{c\sqrt{2\pi}}\end{aligned}\tag{18}$$

Now set these derivatives equal and plug in $x = 0$:

$$\begin{aligned}-\frac{e^0}{(1+e^0)^2} &= -\frac{e^{-\frac{0}{2c^2}}}{c\sqrt{2\pi}} \\ -\frac{1}{4} &= -\frac{1}{c\sqrt{2\pi}} \\ c &= \frac{4}{\sqrt{2\pi}}\end{aligned}\tag{19}$$

4 Analysis

The mathematical analysis of RSLVQ consists of the following steps:

- Describe development of the characteristic quantities in terms of recurrence relations
- Turning of recurrence relations into differential equations
- Performing averages on the differential equations

These three steps are described in appendix B and give us mathematical descriptions of the development of the system with learning time. There are multiple variants of the differential equations. First of all there are two variants of RSLVQ: original (equation (13)) and with Φ -approximation (equation(15)), both in the limit $\eta \rightarrow 0$. In this limit it is possible to analytically determine the ODE's for Φ -approximated RSLVQ while the ODE's for original RSLVQ contain numerical integrations.

The most interesting results are the stationary states of the system of coupled ODE's. These states can be obtained using large learning times or, more reliably, by searching for zeros of the right hand sides of the ODE's. To search for zeros of a seven-dimensional non-linear system is, however, difficult, therefore we search for zeros in the sum of squared right hand side terms, with the restriction that the solution is physical, i.e. the covariance matrix C_k , see equation (34), should be positive semi-definite, which is the case when all its eigenvalues are non-negative.

For determining optimal settings of the control parameter v_{soft} we interpret the generalization error in the stationary states as a function of v_{soft} . This function can then be used to find the optimal setting by searching for the minimum.

5 Experiments and Results

Several experiments have been conducted, some of them containing a comparison of simulations and differential equations. Note that the differential equations are determined in the limit $N \rightarrow \infty$. The simulations were performed with $N = 100$, which is obviously sufficient to match the theory for $N \rightarrow \infty$, see [1] for a discussion of finite N corrections.

5.1 Simulations versus ODE

First we will show that the ODE's indeed describe the system by comparison of the development of the characteristic quantities of simulated training with those of the ODE's.

Original RSLVQ As one can see from figure 1, the ODE's match the simulations exactly, for original RSLVQ. Note that the learning rate of $\eta = 0.05$ already is small enough for the simulations to match the ODE's which are valid in the limit $\eta \rightarrow 0$. The softness v_{soft} is chosen such that it is not too small and not too large because it has a similar effect as the learning rate: too large makes the system taking too big steps resulting in poor convergence and too small causes the system to converge very slowly and besides that it shows the limiting behavior of LFM, i.e. poor performance and instability issues.

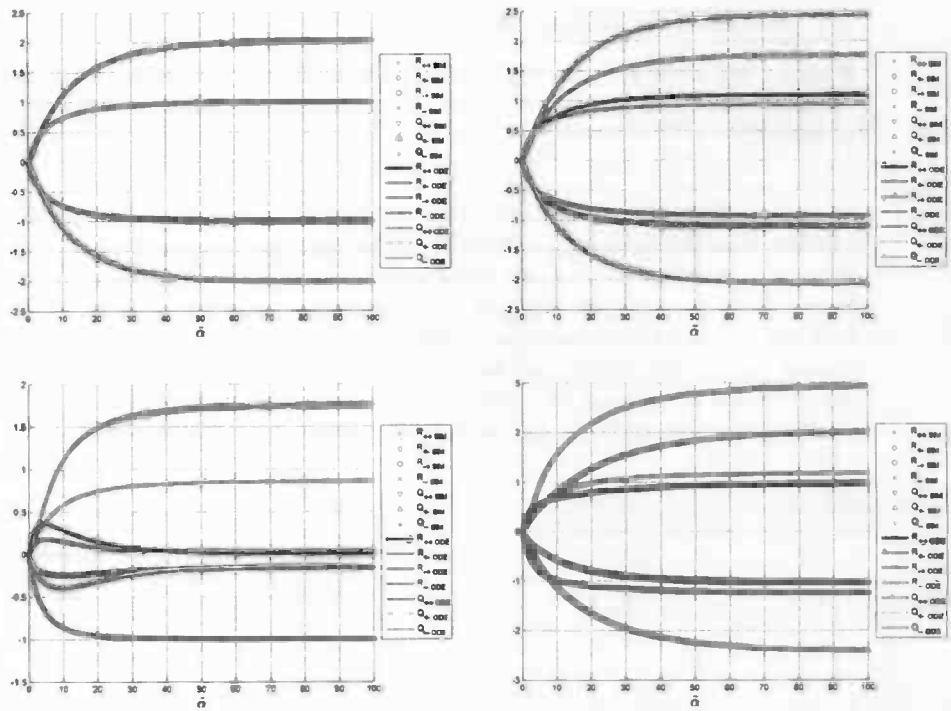


Fig. 1. Development of the characteristic quantities with the learning time $\bar{\alpha}$ of the original RSLVQ for both simulations (thick lines) and ODE's (thinner lines that lay on top of the thick lines) for different systems. Top left: $p_+ = p_- = 0.5, v_+ = v_- = 0.5$; top right: $p_+ = p_- = 0.5, v_+ = 0.25, v_- = 0.81$; bottom left: $p_+ = 0.7, v_+ = v_- = 1$; bottom right: $p_+ = 0.7, v_+ = 0.25, v_- = 0.81$. The softness v_{soft} has been set to 1 and the learning rate used in the simulations is $\eta = 0.05$.

RSLVQ with Φ -approximation The same experiment has also been conducted with RSLVQ with Φ -approximation for both simulations and ODE's.

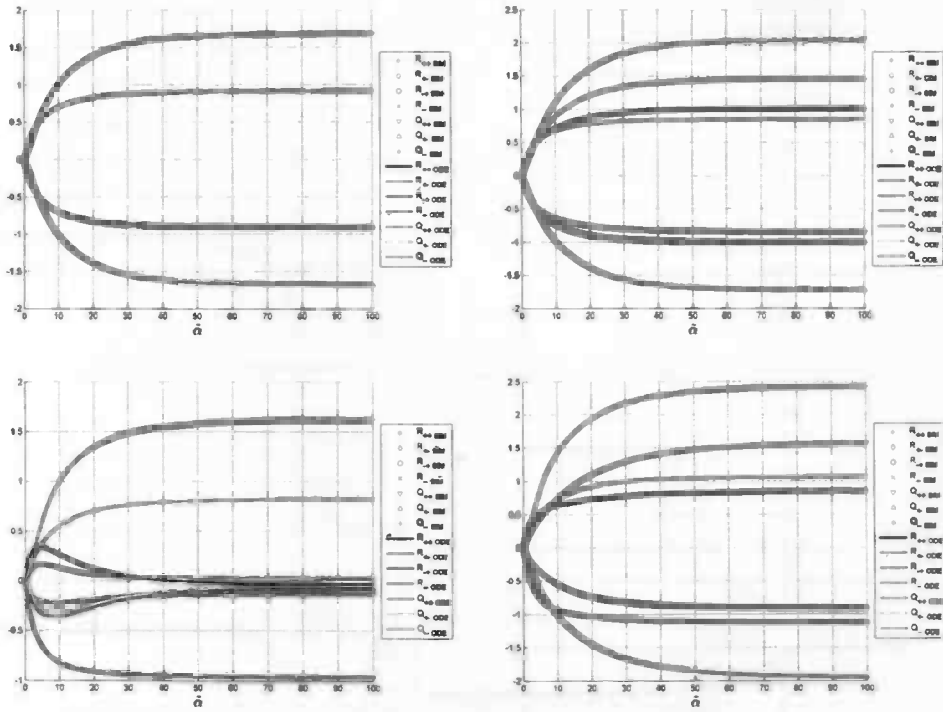


Fig. 2. Development of the characteristic quantities with the learning time $\tilde{\alpha}$ of RSLVQ with Φ -approximation for both simulations (thick lines) and ODE's (thinner lines that lay on top of the thick lines) for different systems. Top (bottom): (un)equal prior and left (right): (un)equal data variances; See figure 1 for detailed information on the settings.

As one can see from figure 2, the ODE's match the simulations exactly, for Φ -approximated RSLVQ as well.

5.2 Comparison

Let us now compare the original with the Φ -approximated version of RSLVQ:

As one can see from figure 3, the Φ -approximation does influence the development of the characteristic quantities, i.e. the ODE's (and therefore the simulations) differ from original RSLVQ and Φ -approximated RSLVQ, however the tendency of each of the quantities is the same and deviations are not too

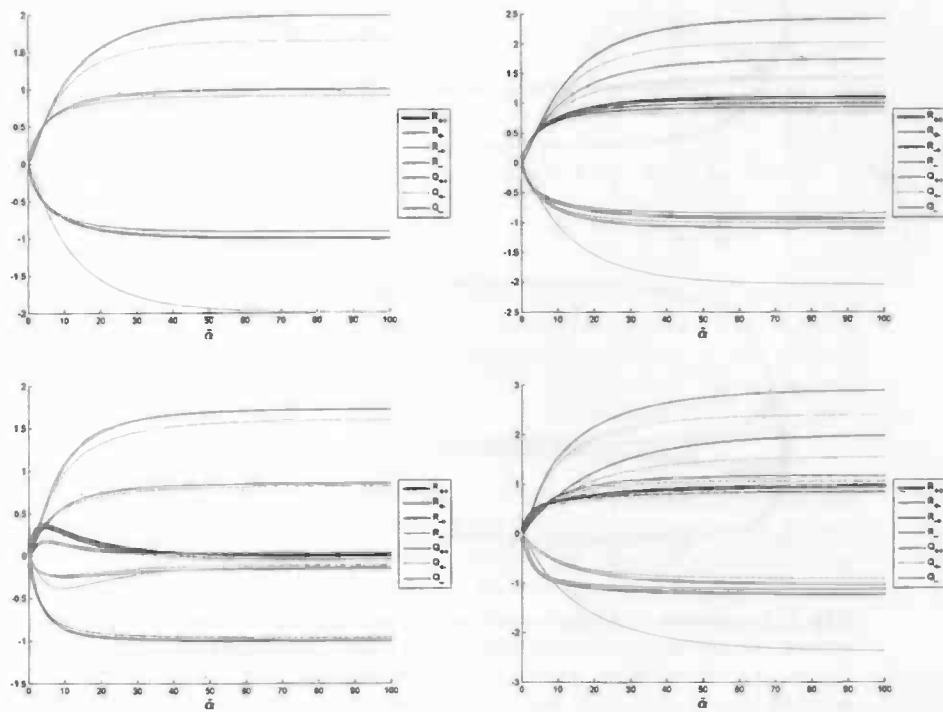


Fig. 3. Comparison of the characteristic quantities with the learning time $\bar{\alpha}$ of RSLVQ with (dotted) and without (solid) Φ -approximation for different systems. Top (bottom): (un)equal prior and left (right): (un)equal data variances; See figure 1 for detailed information on the settings.

large. Moreover the generalization ability is not affected as can be seen in figure 4, which shows a perfect match of the development of the generalization error during learning for original and Φ -approximated RSLVQ. Since the generalization ability is the main target of our study, we conclude that the ODE's for Φ -approximated RSLVQ describe the training process of the original RSLVQ algorithm well and can be used to study the algorithms performance. Note however that these ODE's are only valid for small η , i.e. in the limit $\eta \rightarrow 0$.

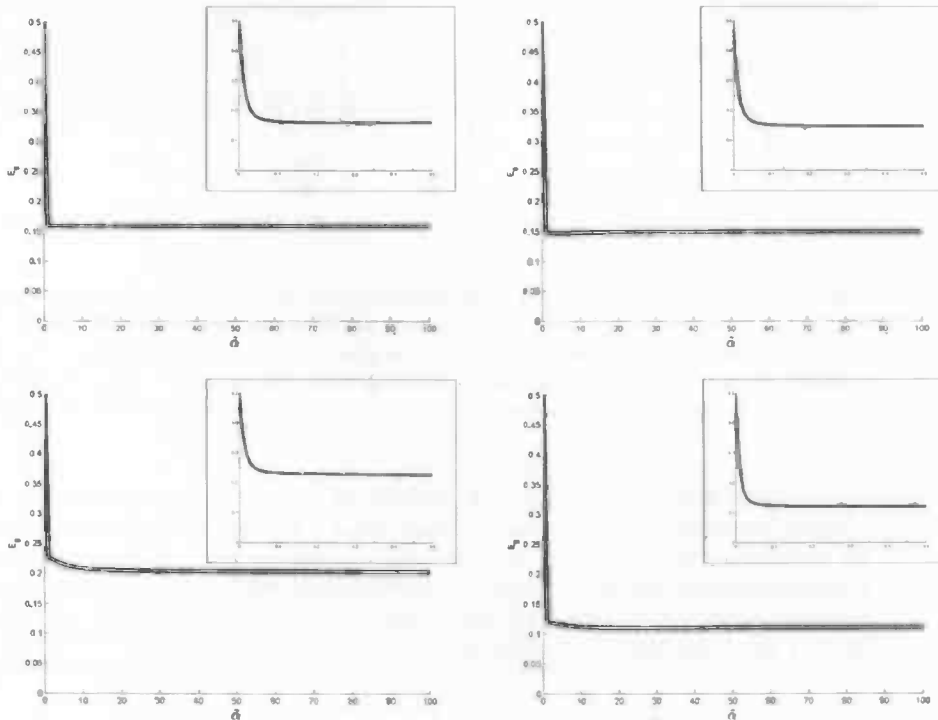


Fig. 4. Comparison of the generalization error (E_g) with the learning time $\tilde{\alpha}$ of RSLVQ with (small red) and without (thick black) Φ -approximation for different systems. Top (bottom): (un)equal prior and left (right): (un)equal data variances; See figure 1 for detailed information on the settings. The insets show a close up of the first part of the graphs, i.e. $\tilde{\alpha} \leq 0.5$.

5.3 Asymptotic Performance

The performance of the LVQ algorithms is measured in terms of the generalization error in the stationary states. These states correspond to zeros in the

derivatives of the ODE's. Because it is difficult to search for the zeros of 7 coupled ODE's, we search for the zeros of the sum of quadratic right hand side terms of the ODE's. The optimal generalization error can be determined for all settings of priors per setting of the data variances.

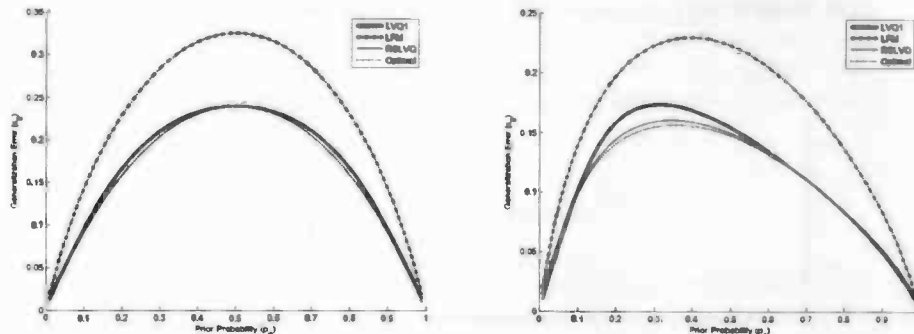


Fig. 5. Asymptotic performance, i.e. the generalization error in the stationary states, of RSLVQ in comparison with other LVQ algorithms. The dashed (lowest) line marks the best linear classifier, the continuous black line marks LVQ1, LFM is marked with the chained line and the red line represents RSLVQ, with the optimal or close to optimal choices of the softness parameter for left: $v_+ = v_- = 1$ and right: $v_+ = 0.25, v_- = 0.81$.

As figure 5 shows, RSLVQ with optimal or close to optimal choices of the softness parameter performs well beyond other LVQ algorithms, even optimal for equal data variances and very close to optimal for unequal data variances. Perhaps this last bit of performance can be gained by using multiple softness parameters, i.e. one per class or per prototype, to enable the system to fit fully to data with unequal data variances.

5.4 Optimal softness

One can write the limiting performance (generalization error) as a function of the softness parameter (v_{soft}) and use this function to numerically find the minimum and therefore the optimal setting of the softness parameter. This has been done for several settings, of which table 1 gives an overview.

As one can see from table 1, the optima found by numerical search are close to the maximum of data variances ($\max(v_+, v_-)$). However, to interpret these numbers correctly, let us look at how the function behaves in the surrounding area.

As one can see from figure 6, there is a large flat region in the generalization error for small values of v_{soft} for both equal and unequal priors. There is presumably a single optimal value, but the true minimum is apparently very flat, yielding very robust algorithm with respect to misestimation of the softness.

	$v_+ = 0,25$ $v_- = 0,25$	$v_+ = 0,5$ $v_- = 0,5$	$v_+ = 1$ $v_- = 1$	$v_+ = 0,25$ $v_- = 0,81$
$p_+ = 0,5$	0,3125	0,51094	1,0625	0,78722
$p_+ = 0,7$	0,25156	0,51328	0,98145	0,83025

Table 1. Optimal settings of v_{soft} , found by numerical search for minima of the generalization error.

It is however possible that for some settings of the softness parameter the (close to) optimal generalization error is reached earlier during training than for others. Either a too small or too large value of v_{soft} results in slow or poor convergence. Furthermore a too small softness (in the limit $v_{soft} \rightarrow 0$) will reach the limiting behavior of LFM, which has stability issues.

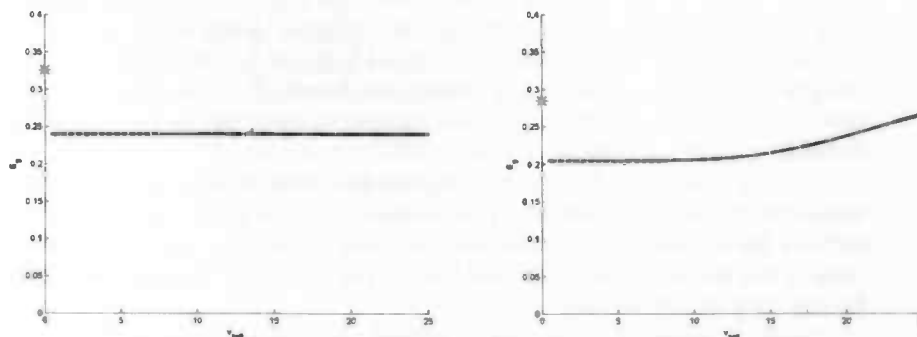


Fig. 6. Dependence of the generalization error for large $\tilde{\alpha}$ on v_{soft} with left: $p_+ = 0.5$ and right: $p_+ = 0.7$, $v_+ = v_- = 1$, showing robustness with respect to settings of the softness parameter. The red star marks the generalization error for LFM, however with finite η because for LFM the limit $\eta \rightarrow 0$ causes instability.

Note that different settings of the data variances showed similar graphs, however for smaller data variances the extension of the close to optimal range of softness decreases. Taking this into consideration, the expectation is that a setting of $1 \leq v_{soft} \leq 2$ will give good performance within reasonable training time for most applications.

6 Discussion

We showed that we can mathematically describe the learning behavior of RSLVQ by using a system of 7 coupled ordinary differential equations. In the limit $\eta \rightarrow 0$ we can even calculate them analytically for RSLVQ with Φ -approximation. Without Φ -approximation we encounter numerical integrations. The ODE's describe

the system well since they fit the to the simulations, for both original and Φ -approximated RSLVQ. The Φ -approximation does influence the behavior of the system, however the quantity of most interest, the generalization error, is not affected by the approximation.

The performance of RSLVQ is well beyond other LVQ algorithms, confirming the findings of the introductory study that was based on simulations only, and even optimal for equal data variances and close to optimal for unequal data variances. Finally we saw that there is presumably a single optimal setting of the softness parameter, but there is also a flat region of close to optimal softness, enabling an easy choice for setting this control parameter for practical applications.

7 Future work

We were able to describe the learning behavior of RSLVQ mathematically in the limit $\eta \rightarrow 0$. ODE's for finite η contain large numerical integrations, causing large time needed for solving the ODE's. It would however be interesting to compare the results for finite η with the results we found. Perhaps some simplifying assumptions enable to get rid of some numerical integrations in the ODE's for finite η .

RSLVQ turned out to perform slightly less then optimal for unequal data variances. Perhaps this last bit of performance can be gained by using multiple softness parameters, for example one per class or one per prototype (this coincides in our model). This could well lead to the optimal LVQ algorithm, at least for the data model we used.

Finally it might be interesting to extend the analysis to more than two prototypes since RSLVQ cannot show its full potential in our model yet, because of the choice for using only two prototypes. By using more than two prototypes per class, the algorithm is able to fit to nonlinear decision boundaries. This however would imply an elaborate extension of the calculations because our analysis is based on using two prototypes.

Acknowledgements

First and foremost I would like to thank my supervisor, and leader of the *LVQ research group*, Michael Biehl for his guidance throughout my master thesis project. He was always available to answer my questions. I am also very grateful to Anarta Ghosh, who has helped me a lot with the calculations.

I would like to thank Aree Witoelar, member of the *LVQ research group*, for his help on debugging the calculations and code for solving the differential equations. The biweekly meetings of the *LVQ research group* have also provided me more insight in various problems and worked as inspiration for my project.

Finally I would like to express my gratitude for the support from my family during the project.

References

1. M. Biehl, A. Ghosh and B. Hammer, *Dynamics and generalization ability of LVQ algorithms*, in *Journal of Machine Learning Research* 8 (Feb):323-360, 2007.
2. A. Ghosh, M. Biehl and B. Hammer, *Dynamical analysis of LVQ type learning rules*, *Workshop on the Self-Organizing-Map, WSOM'05*, 2005.
3. T. Kohonen, *Improved versions of learning vector quantization*, *Proceedings of the International Joint conference on Neural Networks (Sand Diego, 1990)*, 1:545-550, 1990.
4. T. Kohonen, *Learning vector quantization*, M. Arbib, editor, *The handbook of brain theory and neural networks*, 537-540, MIT Press, Cambridge, MA, 1995.
5. T. Kohonen, *Self-Organizing Maps*, Springer, Berlin, 1997.
6. G. Reents and R. Urbanczik, *Self-Averaging and On-Line Learning*, *Physical review letters*. Vol. 80. No. 24, pp. 5445-5448, 1998.
7. A.S. Sato and K. Yamada, *Generalized learning vector quantization*, In G. Tesauero, D. Touretzky and T. Leen, editors, *Advances in Neural Information Processing Systems*, volume 7, 423-429, 1995.
8. S. Seo, M. Bode and K. Obermayer, *Soft nearest prototype classification*, *IEEE Transactions on Neural Networks*, 13(2):390-398, 2003.
9. S. Seo and K. Obermayer, *Soft learning vector quantization*, *Neural computation*, 15: 1589-1603, 2003.
10. J.J.G. de Vries, *The behaviour of RSLVQ in a controlled environment*, *Rijksuniversiteit Groningen*, internal document, 2006.

A Statistics of the model

A.1 Notations

Let us first introduce the following notations:

For any $x \in \mathbb{R}^N$, $x^2 \equiv x \cdot x$; \cdot denotes scalar product. $\langle \cdot \rangle$ denotes the average (expectation) over $p(\xi)$ and can be expressed in the following form:

$$\langle \cdot \rangle = \sum_{\sigma=\pm 1} p_{\sigma} \langle \cdot \rangle_{\sigma} \quad (20)$$

Where $\langle \cdot \rangle_{\sigma}$ is the conditional average for class σ .

A.2 Statistics of the data

From equation (3) follows that the components ξ_j are statistically independent, Gaussian distributed, quantities with variance v_{σ} and mean $\langle \xi_j \rangle_{\sigma} = \lambda(B_{\sigma})_j$. Furthermore, note that for a statistical quantity $X \sim N(\mu, \sigma)$ it holds that $\sigma^2 = \langle X^2 \rangle - \langle X \rangle^2$, so $\langle X^2 \rangle = \langle X \rangle^2 + \sigma^2$. It follows that:

$$\begin{aligned} \langle \xi^2 \rangle_{\sigma} &= \sum_{j=1}^N \langle \xi_j^2 \rangle_{\sigma} \\ &= \sum_{j=1}^N (v_{\sigma} + \langle \xi_j \rangle_{\sigma}^2) \\ &= \sum_{j=1}^N (v_{\sigma} + (\lambda(B_{\sigma})_j)^2) \\ &= v_{\sigma} N + \lambda^2 \sum_{j=1}^N (B_{\sigma})_j^2 \\ &= v_{\sigma} N + \lambda^2 \end{aligned} \quad (21)$$

Note that in the last step it is used that $\sum_{j=1}^N (B_{\sigma})_j^2 = 1$.

Thus we obtain:

$$\begin{aligned} \langle \xi \cdot \xi \rangle &= \sum_{\sigma=\pm 1} p_{\sigma} \langle \xi^2 \rangle_{\sigma} \\ &= p_{+1}(Nv_{+1} + \lambda^2) + p_{-1}(Nv_{-1} + \lambda^2) \\ &= N(p_{+1}v_{+1} + p_{-1}v_{-1}) + \lambda^2 \\ &\approx N(p_{+1}v_{+1} + p_{-1}v_{-1}) \\ &\quad \left[\because N \gg \lambda \right] \end{aligned} \quad (22)$$

A.3 Order Parameters and Projections

Define the order parameters (R_{lm}, Q_{lm}) and the projections (b_l, h_l) as follows:

$$\begin{aligned} R_{lm} &= w_l \cdot B_m \\ Q_{lm} &= w_l \cdot w_m \\ b_l &= \xi \cdot B_l \\ h_l &= \xi \cdot w_l \end{aligned} \quad (23)$$

Define,

$$x = (h_1, h_{-1}, b_1, b_{-1}) \quad (24)$$

A.4 Statistics of the Projections

Given that each training vector is independent of all previous ones, the statistical properties of the projections are well defined for large N . The central limit theorem yields that their joint density, $p(h_{+1}, h_{-1}, b_{+1}, b_{-1}) = p(x)$, is normally distributed and fully specified by the corresponding conditional averages and covariances.

First Order Statistics of h :

$$\begin{aligned} \langle h_l \rangle_k &= \langle w_l \cdot \xi \rangle_k \\ &= w_l \cdot \langle \xi \rangle_k \\ &= w_l \cdot \lambda B_k \\ &= \lambda R_{lk} \end{aligned} \quad (25)$$

First Order Statistics of b :

$$\begin{aligned} \langle b_l \rangle_k &= \langle B_l \cdot \xi \rangle_k \\ &= B_l \cdot \langle \xi \rangle_k \\ &= B_l \cdot \lambda B_k \\ &= \begin{cases} \lambda & \text{if } l = k, \text{ note } B_l^2 = 1 \\ 0 & \text{if } l \neq k, \text{ note } B_l \cdot B_k = 0 \end{cases} \\ &= \lambda \delta_{lk} \end{aligned} \quad (26)$$

Where δ_{lk} is the Kronecker delta. Hence the conditional means of x for two classes can be expressed in the following way:

$$\mu_{+1} = \begin{pmatrix} \lambda R_{+1,+1} \\ \lambda R_{-1,+1} \\ 1 \\ 0 \end{pmatrix} \text{ and } \mu_{-1} = \begin{pmatrix} \lambda R_{+1,-1} \\ \lambda R_{-1,-1} \\ 0 \\ 1 \end{pmatrix} \quad (27)$$

Second Order Statistics of h : $\langle h_l h_m \rangle_k - \langle h_l \rangle_k \langle h_m \rangle_k$

To compute the conditional variance let us first look at the average,

$$\begin{aligned}
\langle h_l h_m \rangle_k &= \langle (\mathbf{w}_l \cdot \boldsymbol{\xi})(\mathbf{w}_m \cdot \boldsymbol{\xi}) \rangle_k \\
&= \left\langle \left(\sum_{i=1}^N (\mathbf{w}_l)_i (\boldsymbol{\xi})_i \right) \left(\sum_{j=1}^N (\mathbf{w}_m)_j (\boldsymbol{\xi})_j \right) \right\rangle_k \\
&= \sum_{i=1}^N (\mathbf{w}_l)_i (\mathbf{w}_m)_i \langle (\boldsymbol{\xi})_i^2 \rangle_k + \sum_{i=1}^N \sum_{j=1, j \neq i}^N (\mathbf{w}_l)_i (\mathbf{w}_m)_j \langle (\boldsymbol{\xi})_i (\boldsymbol{\xi})_j \rangle_k \\
&= \sum_{i=1}^N (\mathbf{w}_l)_i (\mathbf{w}_m)_i [v_k + \lambda^2 (\mathbf{B}_k)_i (\mathbf{B}_k)_i] \\
&\quad + \sum_{i=1}^N \sum_{j=1, j \neq i}^N (\mathbf{w}_l)_i (\mathbf{w}_m)_j \lambda^2 (\mathbf{B}_k)_i (\mathbf{B}_k)_j \\
&\quad \left[\begin{array}{l} \cdot \text{ components of } \boldsymbol{\xi} \text{ have variance } v_k, \text{ also see equation (21),} \\ \cdot \text{ and are independent} \end{array} \right] \\
&= v_k \sum_{i=1}^N (\mathbf{w}_l)_i (\mathbf{w}_m)_i + \lambda^2 \sum_{i=1}^N (\mathbf{w}_l)_i (\mathbf{w}_m)_i (\mathbf{B}_k)_i (\mathbf{B}_k)_i \\
&\quad + \lambda^2 \sum_{i=1}^N \sum_{j=1, j \neq i}^N (\mathbf{w}_l)_i (\mathbf{w}_m)_j (\mathbf{B}_k)_i (\mathbf{B}_k)_j \\
&= v_k \mathbf{w}_l \cdot \mathbf{w}_m + \lambda^2 \sum_{i=1}^N \sum_{j=1}^N (\mathbf{w}_l)_i (\mathbf{w}_m)_j (\mathbf{B}_k)_i (\mathbf{B}_k)_j \\
&= v_k \mathbf{w}_l \cdot \mathbf{w}_m + \lambda^2 (\mathbf{w}_l \cdot \mathbf{B}_k) (\mathbf{w}_m \cdot \mathbf{B}_k) \\
&= v_k Q_{lm} + \lambda^2 R_{lk} R_{mk} \tag{28}
\end{aligned}$$

Hence we have,

$$\begin{aligned}
\langle h_l, h_m \rangle_k - \langle h_l \rangle_k \langle h_m \rangle_k &= v_k Q_{lm} + \lambda^2 R_{lk} R_{mk} - \lambda^2 R_{lk} R_{mk} \\
&= v_k Q_{lm} \tag{29}
\end{aligned}$$

Second Order Statistics of b : Similar to equation (28) we get the second order statistics for b as follows:

$$\begin{aligned}
\langle b_l b_m \rangle_k &= \langle (\mathbf{B}_l \cdot \boldsymbol{\xi})(\mathbf{B}_m \cdot \boldsymbol{\xi}) \rangle_k \\
&= \left\langle \left(\sum_{i=1}^N (\mathbf{B}_l)_i (\boldsymbol{\xi})_i \right) \left(\sum_{j=1}^N (\mathbf{B}_m)_j (\boldsymbol{\xi})_j \right) \right\rangle_k \\
&= \sum_{i=1}^N (\mathbf{B}_l)_i (\mathbf{B}_m)_i \langle (\boldsymbol{\xi})_i^2 \rangle_k + \sum_{i=1}^N \sum_{j=1, j \neq i}^N (\mathbf{B}_l)_i (\mathbf{B}_m)_j \langle (\boldsymbol{\xi})_i (\boldsymbol{\xi})_j \rangle_k
\end{aligned}$$

$$\begin{aligned}
&= \sum_{i=1}^N (\mathbf{B}_l)_i (\mathbf{B}_m)_i [v_k + \lambda^2 (\mathbf{B}_k)_i (\mathbf{B}_k)_i] \\
&\quad + \sum_{i=1}^N \sum_{j=1, j \neq i}^N (\mathbf{B}_l)_i (\mathbf{B}_m)_j \lambda^2 (\mathbf{B}_k)_i (\mathbf{B}_k)_j \\
&\quad \left[\begin{array}{l} \cdot \cdot \text{ components of } \xi \text{ have variance } v_k, \text{ also see equation (21),} \\ \cdot \cdot \text{ and are independent} \end{array} \right] \\
&= v_k \sum_{i=1}^N (\mathbf{B}_l)_i (\mathbf{B}_m)_i + \lambda^2 \sum_{i=1}^N (\mathbf{B}_l)_i (\mathbf{B}_m)_i (\mathbf{B}_k)_i (\mathbf{B}_k)_i \\
&\quad + \lambda^2 \sum_{i=1}^N \sum_{j=1, j \neq i}^N (\mathbf{B}_l)_i (\mathbf{B}_m)_j (\mathbf{B}_k)_i (\mathbf{B}_k)_j \\
&= v_k \mathbf{B}_l \cdot \mathbf{B}_m + \lambda^2 \sum_{i=1}^N \sum_{j=1}^N (\mathbf{B}_l)_i (\mathbf{B}_m)_j (\mathbf{B}_k)_i (\mathbf{B}_k)_j \\
&= v_k \mathbf{B}_l \cdot \mathbf{B}_m + \lambda^2 (\mathbf{B}_l \cdot \mathbf{B}_k) (\mathbf{B}_m \cdot \mathbf{B}_k) \\
&= \delta_{lm} v_k + \lambda^2 \delta_{lk} \delta_{mk} \\
&\quad \left[\begin{array}{l} \cdot \cdot \mathbf{B}_l \cdot \mathbf{B}_k = \begin{cases} 1 & \text{if } l = k \\ 0 & \text{if } l \neq k \end{cases} \end{array} \right] \\
&= \delta_{lm} (v_k + \lambda^2 \delta_{lk}) \tag{30}
\end{aligned}$$

Hence we have,

$$\begin{aligned}
\langle b_l b_m \rangle_k - \langle b_l \rangle_k \langle b_m \rangle_k &= \delta_{lm} (v_k + \lambda^2 \delta_{lk}) - \lambda^2 \delta_{lk} \delta_{mk} \\
&= \delta_{lm} v_k \tag{31}
\end{aligned}$$

Covariance of h and b : To compute the conditional variance, $\langle h_l b_m \rangle_k - \langle h_l \rangle_k \langle b_m \rangle_k$, let us first look at the average,

$$\begin{aligned}
\langle h_l b_m \rangle_k &= \langle (\mathbf{w}_l \cdot \boldsymbol{\xi})(\mathbf{B}_m \cdot \boldsymbol{\xi}) \rangle_k \\
&= \left\langle \left(\sum_{i=1}^N (\mathbf{w}_l)_i (\boldsymbol{\xi})_i \right) \left(\sum_{j=1}^N (\mathbf{B}_m)_j (\boldsymbol{\xi})_j \right) \right\rangle_k \\
&= \left\langle \sum_{i=1}^N (\mathbf{w}_l)_i (\mathbf{B}_m)_i (\boldsymbol{\xi})_i (\boldsymbol{\xi})_i + \sum_{i=1}^N \sum_{j=1, j \neq i}^N (\mathbf{w}_l)_i (\mathbf{B}_m)_j (\boldsymbol{\xi})_i (\boldsymbol{\xi})_j \right\rangle_k \\
&= \sum_{i=1}^N (\mathbf{w}_l)_i (\mathbf{B}_m)_i \langle (\boldsymbol{\xi})_i (\boldsymbol{\xi})_i \rangle_k \\
&\quad + \sum_{i=1}^N \sum_{j=1, j \neq i}^N (\mathbf{w}_l)_i (\mathbf{B}_m)_j \langle (\boldsymbol{\xi})_i (\boldsymbol{\xi})_j \rangle_k
\end{aligned}$$

$$\begin{aligned}
&= \sum_{i=1}^N (\mathbf{w}_l)_i (\mathbf{B}_m)_i [v_k + \lambda^2 (\mathbf{B}_k)_i (\mathbf{B}_k)_i] \\
&+ \sum_{i=1}^N \sum_{j=1, j \neq i}^N (\mathbf{w}_l)_i (\mathbf{B}_m)_j \lambda^2 (\mathbf{B}_k)_i (\mathbf{B}_k)_j \\
&= v_k \mathbf{w}_l \cdot \mathbf{B}_m + \lambda^2 \sum_{i=1}^N \sum_{j=1}^N (\mathbf{w}_l)_i (\mathbf{B}_m)_j (\mathbf{B}_k)_i (\mathbf{B}_k)_j \\
&= v_k \mathbf{w}_l \cdot \mathbf{B}_m + \lambda^2 (\mathbf{w}_l \cdot \mathbf{B}_k) (\mathbf{B}_m \cdot \mathbf{B}_k) \\
&\quad \left[\because \mathbf{B}_m \cdot \mathbf{B}_k = \begin{cases} 1 & \text{if } m = k \\ 0 & \text{if } m \neq k \end{cases} \right] \\
&= v_k R_{lm} + \lambda^2 R_{lk} \delta_{mk}
\end{aligned} \tag{32}$$

Hence we have,

$$\begin{aligned}
\langle h_l, b_m \rangle_k - \langle h_l \rangle_k \langle b_m \rangle_k &= v_k R_{lm} + \lambda^2 R_{lk} \delta_{mk} - \lambda R_{lk} \lambda \delta_{mk} \\
&= v_k R_{lm}
\end{aligned} \tag{33}$$

The conditional density of x for class k is $N(\mu_k, C_k)$, where, μ_k is the conditional mean vector for class k and C_k is the conditional covariance matrix for class k . In our model the conditional density of x is a 4-dimensional Gaussian. Following from equations (29, 31, 33), the covariance matrix by C_k and can be expressed as follows:

$$C_k = v_k \begin{pmatrix} Q_{+1,+1} & Q_{+1,-1} & R_{+1,+1} & R_{+1,-1} \\ Q_{+1,-1} & Q_{-1,-1} & R_{-1,+1} & R_{-1,-1} \\ R_{+1,+1} & R_{-1,+1} & 1 & 0 \\ R_{+1,-1} & R_{-1,-1} & 0 & 1 \end{pmatrix} \tag{34}$$

B Detailed derivation of system of coupled ODE's

Because each training sample that is presented is independent of its predecessors we can conceptually think of training as a stochastic process, to be precise a Markov process (i.e. the future states do depend on the current state but are independent of previous states). If the underlying distribution of the training data is simple enough, the whole dynamics of the system can be analyzed using a few characteristic quantities $\{R_{lm}, Q_{lm}\}$. These order parameters are self-averaging [6] in the thermodynamic limit ($N \rightarrow \infty$), allowing to analyze the stochastic evolution of the system in terms of deterministic evolution of the characteristic quantities. The evolution of characteristic quantities is described by a system of coupled differential equations.

B.1 Recurrence relation for the characteristic quantities

Recurrence relation for **R**: From equation (15) follows:

$$\begin{aligned}
R_{lm}^\mu &= w_l^\mu \cdot B_m \\
&= \left(w_l^{\mu-1} + \Delta w_l^\mu \right) \cdot B_m \\
&= R_{lm}^{\mu-1} + \Delta w_l^\mu \cdot B_m \\
&= R_{lm}^{\mu-1} + \frac{\eta}{Nv_{soft}} (\delta_{l\sigma^\mu} - \Phi_l) (\xi^\mu - w_l^{\mu-1}) \cdot B_m \\
&= R_{lm}^{\mu-1} + \frac{\eta}{Nv_{soft}} \left(\delta_{l\sigma^\mu} \xi^\mu - \delta_{l\sigma^\mu} w_l^{\mu-1} - \Phi_l \xi^\mu \right. \\
&\quad \left. + \Phi_l w_l^{\mu-1} \right) \cdot B_m \\
&= R_{lm}^{\mu-1} + \frac{\eta}{Nv_{soft}} \left(\delta_{l\sigma^\mu} b_m^\mu - \delta_{l\sigma^\mu} R_{lm}^{\mu-1} - \Phi_l b_m^\mu + \Phi_l R_{lm}^{\mu-1} \right) \quad (35)
\end{aligned}$$

Recurrence relation for **Q**:

$$\begin{aligned}
Q_{lm}^\mu &= w_l^\mu \cdot w_m^\mu \\
&= \left(w_l^{\mu-1} + \Delta w_l^\mu \right) \cdot \left(w_m^{\mu-1} + \Delta w_m^\mu \right) \\
&= w_l^{\mu-1} \cdot w_m^{\mu-1} + w_l^{\mu-1} \cdot \Delta w_m^\mu + \Delta w_l^\mu \cdot w_m^{\mu-1} + \Delta w_l^\mu \cdot \Delta w_m^\mu \\
&= Q_{lm}^{\mu-1} + w_l^{\mu-1} \cdot \left(\frac{\eta}{Nv_{soft}} (\delta_{m\sigma^\mu} - \Phi_m) (\xi^\mu - w_m^{\mu-1}) \right) \\
&\quad + w_m^{\mu-1} \cdot \left(\frac{\eta}{Nv_{soft}} (\delta_{l\sigma^\mu} - \Phi_l) (\xi^\mu - w_l^{\mu-1}) \right) \\
&\quad + \left(\frac{\eta}{Nv_{soft}} (\delta_{m\sigma^\mu} - \Phi_m) (\xi^\mu - w_m^{\mu-1}) \right) \\
&\quad \cdot \left(\frac{\eta}{Nv_{soft}} (\delta_{l\sigma^\mu} - \Phi_l) (\xi^\mu - w_l^{\mu-1}) \right) \\
&= Q_{lm}^{\mu-1} + \frac{\eta}{Nv_{soft}} w_l^{\mu-1} \cdot \left(\delta_{m\sigma^\mu} \xi^\mu - \delta_{m\sigma^\mu} w_m^{\mu-1} - \Phi_m \xi^\mu + \Phi_m w_m^{\mu-1} \right) \\
&\quad + \frac{\eta}{Nv_{soft}} w_m^{\mu-1} \cdot \left(\delta_{l\sigma^\mu} \xi^\mu - \delta_{l\sigma^\mu} w_l^{\mu-1} - \Phi_l \xi^\mu + \Phi_l w_l^{\mu-1} \right) \\
&\quad + \left(\frac{\eta}{Nv_{soft}} \right)^2 \left(\delta_{m\sigma^\mu} \xi^\mu - \delta_{m\sigma^\mu} w_m^{\mu-1} - \Phi_m \xi^\mu + \Phi_m w_m^{\mu-1} \right) \\
&\quad \cdot \left(\delta_{l\sigma^\mu} \xi^\mu - \delta_{l\sigma^\mu} w_l^{\mu-1} - \Phi_l \xi^\mu + \Phi_l w_l^{\mu-1} \right) \\
&= Q_{lm}^{\mu-1} + \frac{\eta}{Nv_{soft}} \left(\delta_{m\sigma^\mu} h_l^\mu - \delta_{m\sigma^\mu} Q_{lm}^{\mu-1} - \Phi_m h_l^\mu \right)
\end{aligned}$$

$$\begin{aligned}
& +\Phi_m Q_{lm}^{\mu-1} + \delta_{l\sigma^\mu} h_m^\mu - \delta_{l\sigma^\mu} Q_{lm}^{\mu-1} - \Phi_l h_m^\mu + \Phi_l Q_{lm}^{\mu-1} \Big) \\
& + \left(\frac{\eta}{N v_{soft}} \right)^2 \left(\delta_{l\sigma^\mu} \delta_{m\sigma^\mu} (\xi^\mu)^2 - \delta_{l\sigma^\mu} \delta_{m\sigma^\mu} w_l^{\mu-1} \xi^\mu - \delta_{m\sigma^\mu} \Phi_l (\xi^\mu)^2 \right. \\
& \quad + \delta_{m\sigma^\mu} w_l^{\mu-1} \Phi_l \xi^\mu - \delta_{l\sigma^\mu} \delta_{m\sigma^\mu} w_m^{\mu-1} \xi^\mu + \delta_{l\sigma^\mu} \delta_{m\sigma^\mu} w_l^{\mu-1} w_m^{\mu-1} \\
& \quad + \delta_{m\sigma^\mu} w_m^{\mu-1} \Phi_l \xi^\mu - \delta_{m\sigma^\mu} w_l^{\mu-1} w_m^{\mu-1} \Phi_l - \delta_{l\sigma^\mu} \Phi_m (\xi^\mu)^2 \\
& \quad + \delta_{l\sigma^\mu} w_l^{\mu-1} \Phi_m \xi^\mu + \Phi_l \Phi_m (\xi^\mu)^2 - \Phi_l \Phi_m w_l^{\mu-1} \xi^\mu \\
& \quad + \delta_{l\sigma^\mu} w_m^{\mu-1} \Phi_m \xi^\mu - \delta_{l\sigma^\mu} w_l^{\mu-1} w_m^{\mu-1} \Phi_m - w_m^{\mu-1} \Phi_l \Phi_m \xi^\mu \\
& \quad \left. + w_l^{\mu-1} w_m^{\mu-1} \Phi_l \Phi_m \right) \\
& = Q_{lm}^{\mu-1} + Q_{lm}^{\mu-1} \frac{\eta}{N v_{soft}} \left(\Phi_m + \Phi_l - \delta_{m\sigma^\mu} - \delta_{l\sigma^\mu} \right) \\
& \quad + h_l^\mu \frac{\eta}{N v_{soft}} \left(\delta_{m\sigma^\mu} - \Phi_m \right) + h_m^\mu \frac{\eta}{N v_{soft}} \left(\delta_{l\sigma^\mu} - \Phi_l \right) \\
& \quad + \left(\frac{\eta}{v_{soft}} \right)^2 \left((\xi^\mu)^2 \left(\delta_{l\sigma^\mu} \delta_{m\sigma^\mu} - \delta_{m\sigma^\mu} \Phi_l - \delta_{l\sigma^\mu} \Phi_m + \Phi_l \Phi_m \right) \right. \\
& \quad \quad + \delta_{l\sigma^\mu} \delta_{m\sigma^\mu} \left(Q_{lm} - h_l - h_m \right) + \delta_{m\sigma^\mu} \left(h_l \Phi_l + h_m \Phi_l - Q_{lm} \Phi_l \right) \\
& \quad \quad \left. + \delta_{l\sigma^\mu} \left(h_l \Phi_m + h_m \Phi_m - Q_{lm} \Phi_m \right) + \Phi_l \Phi_m \left(Q_{lm} - h_l - h_m \right) \right) \\
& = Q_{lm}^{\mu-1} + \frac{\eta}{N v_{soft}} Q_{lm}^{\mu-1} \left(\Phi_m + \Phi_l - \delta_{m\sigma^\mu} - \delta_{l\sigma^\mu} \right) \\
& \quad + h_l^\mu \frac{\eta}{N v_{soft}} \left(\delta_{m\sigma^\mu} - \Phi_m \right) + h_m^\mu \frac{\eta}{N v_{soft}} \left(\delta_{l\sigma^\mu} - \Phi_l \right) \\
& \quad + \left(\frac{\eta}{v_{soft}} \right)^2 \left((\xi^\mu)^2 \left(\delta_{l\sigma^\mu} \delta_{m\sigma^\mu} - \delta_{m\sigma^\mu} \Phi_l - \delta_{l\sigma^\mu} \Phi_m + \Phi_l \Phi_m \right) \right. \\
& \quad \quad + \delta_{l\sigma^\mu} \delta_{m\sigma^\mu} \left(Q_{lm} - h_l - h_m \right) + \delta_{m\sigma^\mu} \left(h_l \Phi_l + h_m \Phi_l - Q_{lm} \Phi_l \right) \\
& \quad \quad \left. + \delta_{l\sigma^\mu} \left(h_l \Phi_m + h_m \Phi_m - Q_{lm} \Phi_m \right) + \Phi_l \Phi_m \left(Q_{lm} - h_l - h_m \right) \right) \\
& \approx Q_{lm}^{\mu-1} + \frac{\eta}{N v_{soft}} Q_{lm}^{\mu-1} \left(\Phi_m + \Phi_l - \delta_{m\sigma^\mu} - \delta_{l\sigma^\mu} \right) \\
& \quad + h_l^\mu \frac{\eta}{N v_{soft}} \left(\delta_{m\sigma^\mu} - \Phi_m \right) + h_m^\mu \frac{\eta}{N v_{soft}} \left(\delta_{l\sigma^\mu} - \Phi_l \right) \\
& \quad + \left(\frac{\eta}{N v_{soft}} \right)^2 (\xi^\mu)^2 \left(\delta_{l\sigma^\mu} \delta_{m\sigma^\mu} - \delta_{m\sigma^\mu} \Phi_l - \delta_{l\sigma^\mu} \Phi_m + \Phi_l \Phi_m \right) \tag{36}
\end{aligned}$$

In the last step we neglect the terms of $O(\frac{1}{N^2})$ in (36), note that $(\xi^\mu)^2$ is in the order of $O(N)$, see equation (22), so this term remains.

B.2 Differential equations

From equation (35) and (36) follows:

$$\begin{aligned}
\frac{R_{lm}^\mu - R_{lm}^{\mu-1}}{1/N} &= \frac{\eta}{v_{soft}} \left(\delta_{l\sigma^\mu} b_m^\mu - \delta_{l\sigma^\mu} R_{lm}^{\mu-1} - \Phi_l b_m^\mu + \Phi_l R_{lm}^{\mu-1} \right) \quad (37) \\
\frac{Q_{lm}^\mu - Q_{lm}^{\mu-1}}{1/N} &= \frac{\eta}{v_{soft}} Q_{lm}^{\mu-1} \left(\Phi_m + \Phi_l - \delta_{m\sigma^\mu} - \delta_{l\sigma^\mu} \right) \\
&\quad + h_l^\mu \frac{\eta}{v_{soft}} \left(\delta_{m\sigma^\mu} - \Phi_m \right) + h_m^\mu \frac{\eta}{v_{soft}} \left(\delta_{l\sigma^\mu} - \Phi_l \right) \\
&\quad + \frac{1}{N} \left(\frac{\eta}{v_{soft}} \right)^2 (\xi^\mu)^2 \left(\delta_{l\sigma^\mu} \delta_{m\sigma^\mu} - \delta_{m\sigma^\mu} \Phi_l - \delta_{l\sigma^\mu} \Phi_m + \Phi_l \Phi_m \right) \quad (38)
\end{aligned}$$

We are interested in the mean values of these characteristic quantities, therefore performing averages over the sequence of data. Since the data samples are independent of all previous data samples, the system, including R_{lm} and Q_{lm} , is independent of data sample ξ^μ . Besides this we define:

$$\alpha = \mu/N \quad (39)$$

When considering the thermodynamic limit ($N \rightarrow \infty$), α can be considered as a continuous variable with $\Delta\alpha = 1/N$. Using equation (39) and exploiting the independence we get:

$$\begin{aligned}
\frac{dR_{lm}}{d\alpha} &= \frac{\eta}{v_{soft}} \left(\langle \delta_{l\sigma} b_m \rangle - \langle \delta_{l\sigma} \rangle R_{lm} - \langle \Phi_l b_m \rangle + \langle \Phi_l \rangle R_{lm} \right) \\
&\quad \left[\begin{aligned} \langle \delta_{l\sigma} \rangle &= \sum_{\sigma=\pm 1} p_\sigma \langle \delta_{l\sigma} \rangle_\sigma = p_l \\ \langle \delta_{l\sigma} b_m \rangle &= \sum_{\sigma=\pm 1} p_\sigma \langle \delta_{l\sigma} b_m \rangle_\sigma = p_l \langle b_m \rangle_l = p_l \delta_{lm} \lambda \end{aligned} \right] \\
&= \frac{\eta}{v_{soft}} \left(p_l (\delta_{lm} \lambda - R_{lm}) - \langle \Phi_l b_m \rangle + \langle \Phi_l \rangle R_{lm} \right) \quad (40) \\
\frac{dQ_{lm}}{d\alpha} &= \frac{\eta}{v_{soft}} Q_{lm} \left(\langle \Phi_m \rangle + \langle \Phi_l \rangle - \langle \delta_{m\sigma} \rangle - \langle \delta_{l\sigma} \rangle \right) \\
&\quad + \frac{\eta}{v_{soft}} \left(\langle h_l \delta_{m\sigma} \rangle - \langle h_l \Phi_m \rangle + \langle h_m \delta_{l\sigma} \rangle - \langle h_m \Phi_l \rangle \right) \\
&\quad + \frac{1}{N} \left(\frac{\eta}{v_{soft}} \right)^2 \left(\langle \delta_{l\sigma^\mu} \delta_{m\sigma^\mu} \xi^2 \rangle - \langle \delta_{m\sigma^\mu} \Phi_l \xi^2 \rangle - \langle \delta_{l\sigma^\mu} \Phi_m \xi^2 \rangle + \langle \Phi_l \Phi_m \xi^2 \rangle \right) \\
&\quad \left[\begin{aligned} \langle X \xi^2 \rangle &\approx \sum_{\sigma=\pm 1} p_\sigma \langle \xi^2 \rangle_\sigma \langle X \rangle_\sigma = \sum_{\sigma=\pm 1} p_\sigma v_\sigma N \langle X \rangle_\sigma \\ \langle \delta_{m\sigma} X \rangle &= \sum_{\sigma=\pm 1} p_\sigma \langle \delta_{m\sigma} X \rangle_\sigma = p_m \langle X \rangle_m \\ &\text{and see equations (22, 27)} \end{aligned} \right] \\
&\approx \frac{\eta}{v_{soft}} Q_{lm} \left(\langle \Phi_m \rangle + \langle \Phi_l \rangle - p_m - p_l \right) \\
&\quad + \frac{\eta}{v_{soft}} \left(p_m \lambda R_{lm} - \langle h_l \Phi_m \rangle + p_l \lambda R_{ml} - \langle h_m \Phi_l \rangle \right)
\end{aligned}$$

$$\begin{aligned}
& + \frac{1}{N} \left(\frac{\eta}{v_{soft}} \right)^2 \left(N p_l v_l \delta_{lm} - p_m v_m N \langle \Phi_l \rangle_m \right. \\
& \quad \left. - p_l v_l N \langle \Phi_m \rangle_l + \sum_{\sigma=\pm 1} p_\sigma v_\sigma N \langle \Phi_l \Phi_m \rangle_\sigma \right) \\
& = \frac{\eta}{v_{soft}} Q_{lm} \left(\langle \Phi_m \rangle + \langle \Phi_l \rangle - p_m - p_l \right) \\
& + \frac{\eta}{v_{soft}} \left(p_m \lambda R_{lm} - \langle h_l \Phi_m \rangle + p_l \lambda R_{ml} - \langle h_m \Phi_l \rangle \right) \\
& + \left(\frac{\eta}{v_{soft}} \right)^2 \left(\left(\sum_{\sigma=\pm 1} p_\sigma v_\sigma \langle \Phi_l \Phi_m \rangle_\sigma \right) + \delta_{lm} p_l v_l - p_m v_m \langle \Phi_l \rangle_m - p_l v_l \langle \Phi_m \rangle_l \right) \\
& \approx \frac{\eta}{v_{soft}} Q_{lm} \left(\langle \Phi_m \rangle + \langle \Phi_l \rangle - p_m - p_l \right) \\
& + \frac{\eta}{v_{soft}} \left(p_m \lambda R_{lm} - \langle h_l \Phi_m \rangle + p_l \lambda R_{ml} - \langle h_m \Phi_l \rangle \right) \tag{41}
\end{aligned}$$

In the final step we neglect the terms of $O(\eta^2)$, which is correct in the limit $\eta \rightarrow 0$. To compute the remaining averages in the differential equations above (40,41) let us look at the function Φ_σ .

$$\begin{aligned}
\Phi_\sigma & = \Phi \left(\frac{d_{-\sigma} - d_\sigma}{2c v_{soft}} \right) \\
& = \Phi \left(\frac{1}{2c v_{soft}} \left((\xi - w_{-\sigma})^2 - (\xi - w_\sigma)^2 \right) \right) \\
& = \Phi \left(\frac{1}{2c v_{soft}} \left(2\xi \cdot w_\sigma - 2\xi \cdot w_{-\sigma} + w_{-\sigma}^2 - w_\sigma^2 \right) \right) \\
& = \Phi \left(\frac{1}{2c v_{soft}} \left(2h_\sigma - 2h_{-\sigma} + Q_{-\sigma-\sigma} - Q_{\sigma\sigma} \right) \right) \\
& = \Phi \left(\left(\frac{1}{c v_{soft}}, \frac{-1}{c v_{soft}}, 0, 0 \right) \cdot (h_{+\sigma}, h_{-\sigma}, b_{+\sigma}, b_{-\sigma}) + \frac{Q_{-\sigma-\sigma} - Q_{\sigma\sigma}}{2c v_{soft}} \right) \tag{42}
\end{aligned}$$

Hence we have,

$$\Phi_\sigma = \Phi(\alpha_\sigma \cdot \mathbf{x} - \beta_\sigma) \tag{43}$$

Where, $\alpha_\sigma = \left(\frac{\sigma}{c v_{soft}}, \frac{-\sigma}{c v_{soft}}, 0, 0 \right)$ and $\beta_\sigma = -\left(\frac{Q_{-\sigma, -\sigma} - Q_{\sigma, \sigma}}{2c v_{soft}} \right)$.

Following from equation (20), we can calculate $\langle \Phi(\alpha_\sigma \cdot \mathbf{x} - \beta_\sigma) \rangle$ and $\langle (\mathbf{x})_n \Phi(\alpha_\sigma \cdot \mathbf{x} - \beta_\sigma) \rangle$, where $(\mathbf{x})_n$ is the n^{th} component of \mathbf{x} ; $n \in \{1, 2, 3, 4\}$, see equation (24), from their conditional averages $\langle \dots \rangle_\sigma$.

B.3 Averaging Φ

As can be seen in equation (43), we encounter Φ -functions of the following generic form:

$$\Phi_s = \Phi(\alpha_s \cdot x - \beta_s)$$

$$\text{with } \Phi(x) = \int_{-\infty}^x \frac{e^{-\frac{z^2}{2}}}{\sqrt{2\pi}} dz \quad (44)$$

Now consider averages $\langle (x)_n \Phi_s \rangle_k$ and $\langle \Phi_s \rangle_k$:

$$\begin{aligned} \langle (x)_n \Phi_s \rangle_k &= \frac{1}{(2\pi)^{\frac{4}{2}} (\det(C_k))^{\frac{1}{2}}} \int_{\mathbb{R}^4} (x)_n \Phi(\alpha_s \cdot x - \beta_s) \\ &\quad \exp\left(-\frac{1}{2}(x - \mu_k)^T C_k^{-1} (x - \mu_k)\right) dx \\ &= \frac{1}{(2\pi)^2 (\det(C_k))^{\frac{1}{2}}} \int_{\mathbb{R}^4} (x' + \mu_k)_n \Phi(\alpha_s \cdot x' + \alpha_s \cdot \mu_k - \beta_s) \\ &\quad \exp\left(-\frac{1}{2}x'^T C_k^{-1} x'\right) dx' \end{aligned} \quad (45)$$

with $x' = x - \mu_k$. Note that this shift does not change the integral since the limits are $-\infty$ and ∞ in all \mathbb{R}^4 dimensions. Now decompose C_k into $C_k = C_k^{\frac{1}{2}} C_k^{\frac{1}{2}}$. This is possible since C_k is a covariance matrix and therefore positive semidefinite, hence $C_k^{\frac{1}{2}}$ exists. Now let $x' = C_k^{\frac{1}{2}} y$, resulting in $dx' = \det(C_k^{\frac{1}{2}}) dy = (\det(C_k))^{\frac{1}{2}} dy$:

$$\begin{aligned} \langle (x)_n \Phi_s \rangle_k &= \frac{1}{(2\pi)^2 (\det(C_k))^{\frac{1}{2}}} \int_{\mathbb{R}^4} (C_k^{\frac{1}{2}} y + \mu_k)_n \Phi(\alpha_s C_k^{\frac{1}{2}} y + \alpha_s \cdot \mu_k - \beta_s) \\ &\quad \exp\left(-\frac{1}{2}(C_k^{\frac{1}{2}} y)^T C_k^{-1} C_k^{\frac{1}{2}} y\right) (\det(C_k))^{\frac{1}{2}} dy \\ &= \frac{1}{(2\pi)^2} \int_{\mathbb{R}^4} (C_k^{\frac{1}{2}} y + \mu_k)_n \Phi(\alpha_s C_k^{\frac{1}{2}} y + \alpha_s \cdot \mu_k - \beta_s) \\ &\quad \exp\left(-\frac{1}{2}y^2\right) dy \\ &= \frac{1}{(2\pi)^2} \int_{\mathbb{R}^4} (C_k^{\frac{1}{2}} y)_n \Phi(\alpha_s C_k^{\frac{1}{2}} y + \alpha_s \cdot \mu_k - \beta_s) \exp\left(-\frac{1}{2}y^2\right) dy \\ &\quad + \frac{1}{(2\pi)^2} \int_{\mathbb{R}^4} (\mu_k)_n \Phi(\alpha_s \cdot C_k^{\frac{1}{2}} y + \alpha_s \cdot \mu_k - \beta_s) \exp\left(-\frac{1}{2}y^2\right) dy \\ &= I + \frac{1}{(2\pi)^2} \int_{\mathbb{R}^4} (\mu_k)_n \Phi(\alpha_s \cdot x' + \alpha_s \cdot \mu_k - \beta_s) \\ &\quad \exp\left(-\frac{1}{2}(C_k^{-\frac{1}{2}} x')^2\right) (\det C_k)^{-\frac{1}{2}} dx' \\ &= I + (\mu_k)_n \frac{1}{(2\pi)^2 (\det C_k)^{\frac{1}{2}}} \int_{\mathbb{R}^4} \Phi(\alpha_s \cdot (x - \mu_k) + \alpha_s \cdot \mu_k - \beta_s) \end{aligned}$$

$$\begin{aligned}
& \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_k) \mathbf{C}_k^{-\frac{1}{2}} (\mathbf{x} - \boldsymbol{\mu}_k)\right) d\mathbf{x} \\
&= I + (\boldsymbol{\mu}_k)_n \frac{1}{(2\pi)^2 (\det \mathbf{C}_k)^{\frac{1}{2}}} \int_{\mathbb{R}^4} \Phi(\boldsymbol{\alpha}_s \cdot \mathbf{x} - \beta_s) \\
& \quad \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_k) \mathbf{C}_k^{-\frac{1}{2}} (\mathbf{x} - \boldsymbol{\mu}_k)\right) d\mathbf{x} \\
&= I + (\boldsymbol{\mu}_k)_n \langle \Phi_s \rangle_k
\end{aligned} \tag{46}$$

Where

$$\begin{aligned}
I &= \frac{1}{(2\pi)^2} \int_{\mathbb{R}^4} (\mathbf{C}_k^{\frac{1}{2}} \mathbf{y})_n \Phi(\boldsymbol{\alpha}_s \mathbf{C}_k^{\frac{1}{2}} \mathbf{y} + \boldsymbol{\alpha}_s \cdot \boldsymbol{\mu}_k - \beta_s) \exp\left(-\frac{1}{2} \mathbf{y}^2\right) d\mathbf{y} \\
&= \frac{1}{(2\pi)^2} \int_{\mathbb{R}^4} \sum_{j=1}^4 \left((\mathbf{C}_k^{\frac{1}{2}})_{nj}(\mathbf{y})_j \right) \Phi(\boldsymbol{\alpha}_s \mathbf{C}_k^{\frac{1}{2}} \mathbf{y} + \boldsymbol{\alpha}_s \cdot \boldsymbol{\mu}_k - \beta_s) \exp\left(-\frac{1}{2} \mathbf{y}^2\right) d\mathbf{y} \\
&= \frac{1}{(2\pi)^2} \sum_{j=1}^4 \left(\int_{\mathbb{R}} (\mathbf{C}_k^{\frac{1}{2}})_{nj}(\mathbf{y})_j \Phi(\boldsymbol{\alpha}_s \mathbf{C}_k^{\frac{1}{2}} \mathbf{y} + \boldsymbol{\alpha}_s \cdot \boldsymbol{\mu}_k - \beta_s) \exp\left(-\frac{1}{2} (\mathbf{y})_j^2\right) d(\mathbf{y})_j \right)
\end{aligned} \tag{47}$$

Now define, for $j \in \{1, 2, 3, 4\}$:

$$I_j = \int_{\mathbb{R}} (\mathbf{C}_k^{\frac{1}{2}})_{nj}(\mathbf{y})_j \Phi(\boldsymbol{\alpha}_s \cdot \mathbf{C}_k^{\frac{1}{2}} \mathbf{y} + \boldsymbol{\alpha}_s \cdot \boldsymbol{\mu}_k - \beta_s) \exp\left(-\frac{1}{2} (\mathbf{y})_j^2\right) d(\mathbf{y})_j \tag{48}$$

So one can write:

$$I = \frac{1}{(2\pi)^2} \sum_{j=1}^4 I_j \tag{49}$$

Now apply partial integration to I_j . The rule of partial integration says: $\int f(x)g'(x)dx = f(x)g(x) - \int g(x)f'(x)dx$. Applied to I_j this gives:

$$\begin{aligned}
f((\mathbf{y})_j) &= \Phi(\boldsymbol{\alpha}_s \mathbf{C}_k^{\frac{1}{2}} \mathbf{y} + \boldsymbol{\alpha}_s \cdot \boldsymbol{\mu}_k - \beta_s) \\
g'((\mathbf{y})_j) &= (\mathbf{C}_k^{\frac{1}{2}})_{nj}(\mathbf{y})_j \exp\left(-\frac{1}{2} (\mathbf{y})_j^2\right) \\
g((\mathbf{y})_j) &= -(\mathbf{C}_k^{\frac{1}{2}})_{nj} \exp\left(-\frac{1}{2} (\mathbf{y})_j^2\right) \\
f'((\mathbf{y})_j) &= \frac{\partial}{\partial (\mathbf{y})_j} \Phi(\boldsymbol{\alpha}_s \mathbf{C}_k^{\frac{1}{2}} \mathbf{y} + \boldsymbol{\alpha}_s \cdot \boldsymbol{\mu}_k - \beta_s)
\end{aligned} \tag{50}$$

Applying the partial integration yields:

$$I_j = \underbrace{\left[-(\mathbf{C}_k^{\frac{1}{2}})_{nj} \exp\left(-\frac{1}{2} (\mathbf{y})_j^2\right) \Phi(\boldsymbol{\alpha}_s \mathbf{C}_k^{\frac{1}{2}} \mathbf{y} + \boldsymbol{\alpha}_s \cdot \boldsymbol{\mu}_k - \beta_s) \right]_{-\infty}^{\infty}}_{=0}$$

$$\begin{aligned}
& + \int_{\mathbb{R}} (C_k^{\frac{1}{2}})_{nj} \exp\left(-\frac{1}{2}(\mathbf{y})_j^2\right) \frac{\partial}{\partial(\mathbf{y})_j} \Phi(\alpha_s C_k^{\frac{1}{2}} \mathbf{y} + \alpha_s \cdot \boldsymbol{\mu}_k - \beta_s) d(\mathbf{y})_j \\
& = \int_{\mathbb{R}} (C_k^{\frac{1}{2}})_{nj} \exp\left(-\frac{1}{2}(\mathbf{y})_j^2\right) \frac{\partial}{\partial(\mathbf{y})_j} \Phi(\alpha_s C_k^{\frac{1}{2}} \mathbf{y} + \alpha_s \cdot \boldsymbol{\mu}_k - \beta_s) d(\mathbf{y})_j \quad (51)
\end{aligned}$$

Filling this in into equation (49) gives:

$$\begin{aligned}
I & = \frac{1}{(2\pi)^2} \sum_{j=1}^4 \left(\int_{\mathbb{R}} (C_k^{\frac{1}{2}})_{nj} \exp\left(-\frac{1}{2}(\mathbf{y})_j^2\right) \frac{\partial}{\partial(\mathbf{y})_j} \Phi(\alpha_s C_k^{\frac{1}{2}} \mathbf{y} + \alpha_s \cdot \boldsymbol{\mu}_k - \beta_s) d(\mathbf{y})_j \right) \\
& = \frac{1}{(2\pi)^2} \sum_{j=1}^4 \left((C_k^{\frac{1}{2}})_{nj} \right) \int_{\mathbb{R}^4} \frac{\partial}{\partial(\mathbf{y})_j} \left(\Phi(\alpha_s C_k^{\frac{1}{2}} \mathbf{y} + \alpha_s \cdot \boldsymbol{\mu}_k - \beta_s) \right) \exp\left(-\frac{1}{2}\mathbf{y}^2\right) d\mathbf{y} \quad (52)
\end{aligned}$$

Now consider:

$$\begin{aligned}
\frac{\partial \left(\Phi(\alpha_s C_k^{\frac{1}{2}} \mathbf{y} + \alpha_s \cdot \boldsymbol{\mu}_k - \beta_s) \right)}{\partial(\mathbf{y})_j} & = \phi(\alpha_s C_k^{\frac{1}{2}} \mathbf{y} + \alpha_s \cdot \boldsymbol{\mu}_k - \beta_s) \frac{\partial}{\partial(\mathbf{y})_j} (\alpha_s C_k^{\frac{1}{2}} \mathbf{y}) \\
& = \sum_{i=1}^4 (\alpha_s)_i (C_k^{\frac{1}{2}})_{ij} \phi(\alpha_s C_k^{\frac{1}{2}} \mathbf{y} + \alpha_s \cdot \boldsymbol{\mu}_k - \beta_s) \quad (53)
\end{aligned}$$

Where $\phi(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}$ is the standard normal probability density function and in the last step is used that for the components $(\mathbf{y})_i$ with $i \neq j$ the derivative to $(\mathbf{y})_j$ is zero, resulting in: $\frac{\partial}{\partial(\mathbf{y})_j} (\alpha_s C_k^{\frac{1}{2}} \mathbf{y}) = \sum_{i=1}^4 (\alpha_s)_i (C_k^{\frac{1}{2}})_{ij}$. Hence,

$$\begin{aligned}
I & = \frac{1}{(2\pi)^2} \sum_{j=1}^4 \left((C_k^{\frac{1}{2}})_{nj} \sum_{i=1}^4 (\alpha_s)_i (C_k^{\frac{1}{2}})_{ij} \right) \\
& \quad \int_{\mathbb{R}^4} \phi(\alpha_s C_k^{\frac{1}{2}} \mathbf{y} + \alpha_s \cdot \boldsymbol{\mu}_k - \beta_s) \exp\left(-\frac{1}{2}\mathbf{y}^2\right) d\mathbf{y} \\
& = \frac{1}{(2\pi)^2} (C_k \alpha_s)_n \int_{\mathbb{R}^4} \phi(\alpha_s C_k^{\frac{1}{2}} \mathbf{y} + \alpha_s \cdot \boldsymbol{\mu}_k - \beta_s) \exp\left(-\frac{1}{2}\mathbf{y}^2\right) d\mathbf{y} \quad (54)
\end{aligned}$$

Note that in the last equation it is used that $\sum_{j=1}^4 \left((C_k^{\frac{1}{2}})_{nj} \sum_{i=1}^4 (\alpha_s)_i (C_k^{\frac{1}{2}})_{ij} \right) = (C_k \alpha_s)_n$, which is true for symmetrical matrices $C_k^{\frac{1}{2}}$. Since C_k is a covariance matrix hence it is symmetrical, see equation (34), and positive definite, there exists at least one decomposition into $C_k^{\frac{1}{2}}$ also being symmetrical.

Note also that $\exp(-\frac{1}{2}\mathbf{y}^2) d\mathbf{y}$ is a measure which is invariant under rotation of the coordinate axes. Now rotate the system in such a way that one of the axes, say \tilde{y} , is aligned with the vector $C_k^{\frac{1}{2}} \alpha_s$. Using that $\frac{1}{\sqrt{2\pi}} \int_{\mathbb{R}} \exp(-\frac{1}{2}z^2) dz = 1$, the remaining three coordinates can be integrated over, yielding:

$$I = \frac{1}{\sqrt{2\pi}} (\mathbf{C}_k \boldsymbol{\alpha}_s)_n \int_{\mathbb{R}} \phi(\|\mathbf{C}_k^{\frac{1}{2}} \boldsymbol{\alpha}_s\| \tilde{y} + \boldsymbol{\alpha}_s \cdot \boldsymbol{\mu}_k - \beta_s) \exp\left(-\frac{1}{2} \tilde{y}^2\right) d\tilde{y} \quad (55)$$

Now define:

$$\tilde{\alpha}_{s,k} = \|\mathbf{C}_k^{\frac{1}{2}} \boldsymbol{\alpha}_s\| = \sqrt{\boldsymbol{\alpha}_s \mathbf{C}_k \boldsymbol{\alpha}_s} \quad \text{and} \quad \tilde{\beta}_{s,k} = \boldsymbol{\alpha}_s \cdot \boldsymbol{\mu}_k - \beta_s \quad (56)$$

and rename $z = \tilde{\alpha}_{s,k} \tilde{y}$, so $d\tilde{y} = \frac{1}{\tilde{\alpha}_{s,k}} dz$, to obtain:

$$\begin{aligned} I &= \frac{(\mathbf{C}_k \boldsymbol{\alpha}_s)_n}{\sqrt{2\pi}} \int_{\mathbb{R}} \phi(\tilde{\alpha}_{s,k} \tilde{y} + \tilde{\beta}_{s,k}) \exp\left(-\frac{1}{2} \tilde{y}^2\right) d\tilde{y} \\ &= \frac{(\mathbf{C}_k \boldsymbol{\alpha}_s)_n}{\sqrt{2\pi} \tilde{\alpha}_{s,k}} \int_{\mathbb{R}} \phi(z + \tilde{\beta}_{s,k}) \exp\left(-\frac{1}{2} \left(\frac{z}{\tilde{\alpha}_{s,k}}\right)^2\right) dz \\ &= \frac{(\mathbf{C}_k \boldsymbol{\alpha}_s)_n}{\sqrt{2\pi} \tilde{\alpha}_{s,k}} \int_{\mathbb{R}} \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(z + \tilde{\beta}_{s,k})^2}{2}\right) \exp\left(-\frac{1}{2} \left(\frac{z}{\tilde{\alpha}_{s,k}}\right)^2\right) dz \\ &= \frac{(\mathbf{C}_k \boldsymbol{\alpha}_s)_n}{\sqrt{2\pi} \tilde{\alpha}_{s,k}} \int_{\mathbb{R}} \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2} z^2 \left(1 + \frac{1}{\tilde{\alpha}_{s,k}^2}\right) - z \tilde{\beta}_{s,k} - \frac{1}{2} \tilde{\beta}_{s,k}^2\right) dz \\ &= \frac{(\mathbf{C}_k \boldsymbol{\alpha}_s)_n}{\sqrt{2\pi} \tilde{\alpha}_{s,k}} \exp\left(-\frac{1}{2} \tilde{\beta}_{s,k}^2\right) \int_{\mathbb{R}} \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2} z^2 \left(1 + \frac{1}{\tilde{\alpha}_{s,k}^2}\right) - z \tilde{\beta}_{s,k}\right) dz \quad (57) \end{aligned}$$

Now use that $\int_{-\infty}^{\infty} \frac{dx}{\sqrt{2\pi}} \exp\left(-\frac{1}{2} ax^2 + bx\right) = \frac{1}{\sqrt{a}} \exp\left(\frac{b^2}{2a}\right)$:

$$\begin{aligned} I &= \frac{(\mathbf{C}_k \boldsymbol{\alpha}_s)_n}{\sqrt{2\pi} \tilde{\alpha}_{s,k}} \exp\left(-\frac{1}{2} \tilde{\beta}_{s,k}^2\right) \frac{1}{\sqrt{1 + \frac{1}{\tilde{\alpha}_{s,k}^2}}} \exp\left(\frac{\tilde{\beta}_{s,k}^2}{2(1 + \frac{1}{\tilde{\alpha}_{s,k}^2})}\right) \\ &= \frac{(\mathbf{C}_k \boldsymbol{\alpha}_s)_n}{\sqrt{2\pi} \tilde{\alpha}_{s,k}} \frac{1}{\sqrt{1 + \frac{1}{\tilde{\alpha}_{s,k}^2}}} \exp\left(-\frac{1}{2} \frac{\tilde{\beta}_{s,k}^2}{(1 + \tilde{\alpha}_{s,k}^2)}\right) \quad (58) \end{aligned}$$

Summarizing equations (47) to (58) gives, from equation (45):

$$\langle (\mathbf{x})_n \Phi_s \rangle_k = \frac{(\mathbf{C}_k \boldsymbol{\alpha}_s)_n}{\sqrt{2\pi} \tilde{\alpha}_{s,k}} \frac{1}{\sqrt{1 + \frac{1}{\tilde{\alpha}_{s,k}^2}}} \exp\left(-\frac{1}{2} \frac{\tilde{\beta}_{s,k}^2}{(1 + \tilde{\alpha}_{s,k}^2)}\right) + (\boldsymbol{\mu}_k)_n \langle \Phi_s \rangle_k \quad (59)$$

Next consider the average $\langle \Phi_s \rangle_k$. Similar calculations as for $\langle (\mathbf{x})_n \Phi_s \rangle_k$ give:

$$\begin{aligned} \langle \Phi_s \rangle_k &= \frac{1}{(2\pi)^{\frac{1}{2}} (\det(\mathbf{C}_k))^{\frac{1}{2}}} \int_{\mathbb{R}^4} \Phi(\boldsymbol{\alpha}_s \cdot \mathbf{x} - \beta_s) \\ &\quad \exp\left(-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_k)^T \mathbf{C}_k^{-1} (\mathbf{x} - \boldsymbol{\mu}_k)\right) d\mathbf{x} \end{aligned}$$

$$\begin{aligned}
&= \frac{1}{(2\pi)^2(\det(\mathbf{C}_k))^{\frac{1}{2}}} \int_{\mathbb{R}^4} \Phi(\alpha_s \cdot \mathbf{x}' + \alpha_s \cdot \mu_k - \beta_s) \\
&\quad \exp\left(-\frac{1}{2}\mathbf{x}'^T \mathbf{C}_k^{-1} \mathbf{x}'\right) d\mathbf{x}' \\
&= \frac{1}{(2\pi)^2(\det(\mathbf{C}_k))^{\frac{1}{2}}} \int_{\mathbb{R}^4} \Phi(\alpha_s \mathbf{C}_k^{\frac{1}{2}} \mathbf{y} + \alpha_s \cdot \mu_k - \beta_s) \\
&\quad \exp\left(-\frac{1}{2}(\mathbf{C}_k^{\frac{1}{2}} \mathbf{y})^T \mathbf{C}_k^{-1} \mathbf{C}_k^{\frac{1}{2}} \mathbf{y}\right) (\det(\mathbf{C}_k))^{\frac{1}{2}} d\mathbf{y} \\
&= \frac{1}{(2\pi)^2} \int_{\mathbb{R}^4} \Phi(\alpha_s \mathbf{C}_k^{\frac{1}{2}} \mathbf{y} + \alpha_s \cdot \mu_k - \beta_s) \exp\left(-\frac{1}{2}\mathbf{y}^2\right) d\mathbf{y} \\
&= \frac{1}{\sqrt{2\pi}} \int_{\mathbb{R}} \Phi(\|\alpha_s \mathbf{C}_k^{\frac{1}{2}}\| \tilde{y} + \alpha_s \cdot \mu_k - \beta_s) \exp\left(-\frac{1}{2}\tilde{y}^2\right) d\tilde{y} \\
&= \frac{1}{\sqrt{2\pi}} \int_{\mathbb{R}} \Phi(\tilde{\alpha}_{s,k} \tilde{y} + \tilde{\beta}_{s,k}) \exp\left(-\frac{1}{2}\tilde{y}^2\right) d\tilde{y} \tag{60}
\end{aligned}$$

Now apply $\int_{-\infty}^{\infty} \frac{dx}{\sqrt{2\pi}} \exp(-\frac{1}{2}ax^2 + bx)\Phi(cx + d) = \frac{1}{\sqrt{a}} \exp\left(\frac{b^2}{2a}\right)\Phi\left(\frac{bc+ad}{\sqrt{a^2+ac^2}}\right)$, which holds for ($a > 0$):

$$\langle \Phi_s \rangle_k = \Phi\left(\frac{\tilde{\beta}_{s,k}}{\sqrt{1 + \tilde{\alpha}_{s,k}^2}}\right) \tag{61}$$

Hence the required averages are as follows:

$$\begin{aligned}
\langle \Phi_s \rangle_k &= \Phi\left(\frac{\tilde{\beta}_{s,k}}{\sqrt{1 + \tilde{\alpha}_{s,k}^2}}\right) \\
\langle (x)_n \Phi_s \rangle_k &= \frac{(\mathbf{C}_k \alpha_s)_n}{\sqrt{2\pi} \tilde{\alpha}_{s,k}} \frac{1}{\sqrt{1 + \frac{1}{\tilde{\alpha}_{s,k}^2}}} \exp\left(-\frac{1}{2} \frac{\tilde{\beta}_{s,k}^2}{(1 + \tilde{\alpha}_{s,k}^2)}\right) \\
&\quad + (\mu_k)_n \Phi\left(\frac{\tilde{\beta}_{s,k}}{\sqrt{1 + \tilde{\alpha}_{s,k}^2}}\right) \tag{62}
\end{aligned}$$

Where,

$$\tilde{\alpha}_{s,k} = \sqrt{\alpha_s \mathbf{C}_k \alpha_s} \quad \text{and} \quad \tilde{\beta}_{s,k} = \alpha_s \cdot \mu_k - \beta_s \tag{63}$$

B.4 Final form of the differential equations

Filling in the conditional averages of equation (62) into equations (40) and (41), yields the following system of coupled differential equations:

$$\begin{aligned}
\frac{dR_{lm}}{d\alpha} &= \frac{\eta}{v_{soft}} \left(p_l(\delta_{lm}\lambda - R_{lm}) - \sum_{\sigma=\pm 1} p_\sigma \langle b_m \Phi(\alpha_l \cdot x - \beta_l) \rangle_\sigma \right. \\
&\quad \left. + R_{lm} \sum_{\sigma=\pm 1} p_\sigma \langle \Phi(\alpha_l \cdot x - \beta_l) \rangle_\sigma \right) \\
&= \frac{\eta}{v_{soft}} \left(p_l(\delta_{lm}\lambda - R_{lm}) \right. \\
&\quad \left. - \left[\sum_{\sigma=\pm 1} p_\sigma \left(\frac{(C_\sigma \alpha_l)_{n_{bm}}}{\sqrt{2\pi} \tilde{\alpha}_{l,\sigma}} \frac{1}{\sqrt{1 + \frac{1}{\tilde{\alpha}_{l,\sigma}^2}}} \exp\left(-\frac{1}{2} \frac{\tilde{\beta}_{l,\sigma}^2}{(1 + \tilde{\alpha}_{l,\sigma}^2)}\right) \right. \right. \right. \\
&\quad \left. \left. \left. + (\mu_\sigma)_{n_{bm}} \Phi\left(\frac{\tilde{\beta}_{l,\sigma}}{\sqrt{1 + \tilde{\alpha}_{l,\sigma}^2}}\right) \right) \right] + R_{lm} \sum_{\sigma=\pm 1} p_\sigma \Phi\left(\frac{\tilde{\beta}_{l,\sigma}}{\sqrt{1 + \tilde{\alpha}_{l,\sigma}^2}}\right) \right) \quad (64)
\end{aligned}$$

$$\begin{aligned}
\frac{dQ_{lm}}{d\alpha} &= \frac{\eta}{v_{soft}} Q_{lm} \left(\sum_{\sigma=\pm 1} p_\sigma \langle \Phi(\alpha_m \cdot x - \beta_m) \rangle_\sigma + \sum_{\sigma=\pm 1} p_\sigma \langle \Phi(\alpha_l \cdot x - \beta_l) \rangle_\sigma - p_m - p_l \right) \\
&\quad + \frac{\eta}{v_{soft}} \left(p_m \lambda R_{lm} - \sum_{\sigma=\pm 1} p_\sigma \langle h_l \Phi(\alpha_m \cdot x - \beta_m) \rangle_\sigma \right. \\
&\quad \left. + p_l \lambda R_{ml} - \sum_{\sigma=\pm 1} p_\sigma \langle h_m \Phi(\alpha_l \cdot x - \beta_l) \rangle_\sigma \right) \\
&= \frac{\eta}{v_{soft}} Q_{lm} \left(\sum_{\sigma=\pm 1} p_\sigma \Phi\left(\frac{\tilde{\beta}_{m,\sigma}}{\sqrt{1 + \tilde{\alpha}_{m,\sigma}^2}}\right) + \sum_{\sigma=\pm 1} p_\sigma \Phi\left(\frac{\tilde{\beta}_{l,\sigma}}{\sqrt{1 + \tilde{\alpha}_{l,\sigma}^2}}\right) - p_m - p_l \right) \\
&\quad + \frac{\eta}{v_{soft}} \left(p_m \lambda R_{lm} + p_l \lambda R_{ml} \right. \\
&\quad \left. - \left[\sum_{\sigma=\pm 1} p_\sigma \left(\frac{(C_\sigma \alpha_m)_{n_{hl}}}{\sqrt{2\pi} \tilde{\alpha}_{m,\sigma}} \frac{1}{\sqrt{1 + \frac{1}{\tilde{\alpha}_{m,\sigma}^2}}} \exp\left(-\frac{1}{2} \frac{\tilde{\beta}_{m,\sigma}^2}{(1 + \tilde{\alpha}_{m,\sigma}^2)}\right) + (\mu_\sigma)_{n_{hl}} \Phi\left(\frac{\tilde{\beta}_{m,\sigma}}{\sqrt{1 + \tilde{\alpha}_{m,\sigma}^2}}\right) \right) \right] \right. \\
&\quad \left. - \left[\sum_{\sigma=\pm 1} p_\sigma \left(\frac{(C_\sigma \alpha_l)_{n_{hm}}}{\sqrt{2\pi} \tilde{\alpha}_{l,\sigma}} \frac{1}{\sqrt{1 + \frac{1}{\tilde{\alpha}_{l,\sigma}^2}}} \exp\left(-\frac{1}{2} \frac{\tilde{\beta}_{l,\sigma}^2}{(1 + \tilde{\alpha}_{l,\sigma}^2)}\right) + (\mu_\sigma)_{n_{hm}} \Phi\left(\frac{\tilde{\beta}_{l,\sigma}}{\sqrt{1 + \tilde{\alpha}_{l,\sigma}^2}}\right) \right) \right] \right) \quad (65)
\end{aligned}$$

Where,

$$\begin{aligned}
n_{bk} &= \begin{cases} 3 & \text{if } k = 1 \\ 4 & \text{if } k = -1 \end{cases} \\
n_{hk} &= \begin{cases} 1 & \text{if } k = 1 \\ 2 & \text{if } k = -1 \end{cases}
\end{aligned}$$

$$\begin{aligned}\tilde{\alpha}_{s,k} &= \sqrt{\alpha_s C_k \alpha_s}, \quad \tilde{\beta}_{s,k} = \alpha_s \cdot \mu_k - \beta_s \\ \alpha_\sigma &= \left(\frac{\sigma}{cv_{soft}}, \frac{-\sigma}{cv_{soft}}, 0, 0 \right), \quad \beta_\sigma = - \left(\frac{Q_{-\sigma, -\sigma} - Q_{\sigma, \sigma}}{2cv_{soft}} \right)\end{aligned}\quad (66)$$

In equations (64) and (65) one can see that the ODE's are linear in η . The learning rate η is therefore, in the limit $\eta \rightarrow 0$, no more than a scaling factor, which can be taken out by rescaling α to $\tilde{\alpha} = \alpha\eta$:

$$\begin{aligned}\frac{dR_{lm}}{d\tilde{\alpha}} &= \frac{dR_{lm}}{d\alpha\eta} \\ &= \frac{1}{v_{soft}} \left(p_l (\delta_{lm} \lambda - R_{lm}) \right. \\ &\quad \left. - \left[\sum_{\sigma=\pm 1} p_\sigma \left(\frac{(C_\sigma \alpha_l)_{nbm}}{\sqrt{2\pi} \tilde{\alpha}_{l,\sigma}} \frac{1}{\sqrt{1 + \frac{1}{\tilde{\alpha}_{l,\sigma}^2}}} \exp\left(-\frac{1}{2} \frac{\tilde{\beta}_{l,\sigma}^2}{(1 + \tilde{\alpha}_{l,\sigma}^2)}\right) \right. \right. \right. \\ &\quad \left. \left. \left. + (\mu_\sigma)_{nbm} \Phi\left(\frac{\tilde{\beta}_{l,\sigma}}{\sqrt{1 + \tilde{\alpha}_{l,\sigma}^2}}\right) \right) \right] + R_{lm} \sum_{\sigma=\pm 1} p_\sigma \Phi\left(\frac{\tilde{\beta}_{l,\sigma}}{\sqrt{1 + \tilde{\alpha}_{l,\sigma}^2}}\right) \right)\end{aligned}\quad (67)$$

$$\begin{aligned}\frac{dQ_{lm}}{d\tilde{\alpha}} &= \frac{dQ_{lm}}{d\alpha\eta} \\ &= \frac{1}{v_{soft}} Q_{lm} \left(\sum_{\sigma=\pm 1} p_\sigma \Phi\left(\frac{\tilde{\beta}_{m,\sigma}}{\sqrt{1 + \tilde{\alpha}_{m,\sigma}^2}}\right) + \sum_{\sigma=\pm 1} p_\sigma \Phi\left(\frac{\tilde{\beta}_{l,\sigma}}{\sqrt{1 + \tilde{\alpha}_{l,\sigma}^2}}\right) - p_m - p_l \right) \\ &\quad + \frac{1}{v_{soft}} \left(p_m \lambda R_{lm} + p_l \lambda R_{ml} \right. \\ &\quad \left. - \left[\sum_{\sigma=\pm 1} p_\sigma \left(\frac{(C_\sigma \alpha_m)_{nbl}}{\sqrt{2\pi} \tilde{\alpha}_{m,\sigma}} \frac{1}{\sqrt{1 + \frac{1}{\tilde{\alpha}_{m,\sigma}^2}}} \exp\left(-\frac{1}{2} \frac{\tilde{\beta}_{m,\sigma}^2}{(1 + \tilde{\alpha}_{m,\sigma}^2)}\right) + (\mu_\sigma)_{nbl} \Phi\left(\frac{\tilde{\beta}_{m,\sigma}}{\sqrt{1 + \tilde{\alpha}_{m,\sigma}^2}}\right) \right) \right] \right. \\ &\quad \left. - \left[\sum_{\sigma=\pm 1} p_\sigma \left(\frac{(C_\sigma \alpha_l)_{nhm}}{\sqrt{2\pi} \tilde{\alpha}_{l,\sigma}} \frac{1}{\sqrt{1 + \frac{1}{\tilde{\alpha}_{l,\sigma}^2}}} \exp\left(-\frac{1}{2} \frac{\tilde{\beta}_{l,\sigma}^2}{(1 + \tilde{\alpha}_{l,\sigma}^2)}\right) + (\mu_\sigma)_{nhm} \Phi\left(\frac{\tilde{\beta}_{l,\sigma}}{\sqrt{1 + \tilde{\alpha}_{l,\sigma}^2}}\right) \right) \right] \right)\end{aligned}\quad (68)$$

B.5 Original RSLVQ

Now let us look at the original definition of RSLVQ by Obermayer and Seo. The recurrence relations are the same, except we need to replace Φ_l by $\frac{1}{1 + \exp(\frac{d_l^\mu - d_{-l}^\mu}{2v_{soft}})}$ or Ω_l for shorthand.

Recurrence relation for \mathbf{R} : From equation (35) follows:

$$R_{lm}^\mu = R_{lm}^{\mu-1} + \frac{\eta}{N v_{soft}} \left(\delta_{l\sigma^\mu} b_m^\mu - \delta_{l\sigma^\mu} R_{lm}^{\mu-1} - \Omega_l b_m^\mu + \Omega_l R_{lm}^{\mu-1} \right) \quad (69)$$

Recurrence relation for Q: From equation (36) follows:

$$\begin{aligned}
Q_{lm}^\mu &\approx Q_{lm}^{\mu-1} + \frac{\eta}{Nv_{soft}} Q_{lm}^{\mu-1} \left(\Omega_m + \Omega_l - \delta_{m\sigma^\mu} - \delta_{l\sigma^\mu} \right) \\
&+ h_l^\mu \frac{\eta}{Nv_{soft}} \left(\delta_{m\sigma^\mu} - \Omega_m \right) + h_m^\mu \frac{\eta}{Nv_{soft}} \left(\delta_{l\sigma^\mu} - \Omega_l \right) \\
&+ \left(\frac{\eta}{Nv_{soft}} \right)^2 (\xi^\mu)^2 \left(\delta_{lm} - \delta_{m\sigma^\mu} \Omega_l - \delta_{l\sigma^\mu} \Omega_m + \Omega_l \Omega_m \right) \quad (70)
\end{aligned}$$

Differential equations: From equation (69) and (70) follows, similar to equations (40) and (41):

$$\frac{dR_{lm}}{d\alpha} = \frac{\eta}{v_{soft}} \left(p_l (\delta_{lm} \lambda - R_{lm}) - \langle \Omega_l b_m \rangle + \langle \Omega_l \rangle R_{lm} \right) \quad (71)$$

$$\begin{aligned}
\frac{dQ_{lm}}{d\alpha} &\approx \frac{\eta}{v_{soft}} Q_{lm} \left(\langle \Omega_m \rangle + \langle \Omega_l \rangle - p_m - p_l \right) \\
&+ \frac{\eta}{v_{soft}} \left(p_m \lambda R_l - \langle h_l \Omega_m \rangle + p_l \lambda R_{ml} - \langle h_m \Omega_l \rangle \right) \\
&+ \left(\frac{\eta}{v_{soft}} \right)^2 \left(\left(\sum_{\sigma=\pm 1} p_\sigma v_\sigma (\langle \Omega_l \Omega_m \rangle_\sigma + \delta_{lm}) \right) - p_m v_m \langle \Omega_l \rangle_m - p_l v_l \langle \Omega_m \rangle_l \right) \\
&\approx \frac{\eta}{v_{soft}} Q_{lm} \left(\langle \Omega_m \rangle + \langle \Omega_l \rangle - p_m - p_l \right) \\
&+ \frac{\eta}{v_{soft}} \left(p_m \lambda R_l - \langle h_l \Omega_m \rangle + p_l \lambda R_{ml} - \langle h_m \Omega_l \rangle \right) \quad (72)
\end{aligned}$$

Note again that in the final step we neglect the terms of $O(\eta^2)$, which is correct in the limit $\eta \rightarrow 0$. Similar to equation (42) we can generalize Ω_σ .

$$\begin{aligned}
\Omega_\sigma &= \frac{1}{1 + \exp\left(\frac{d_\sigma - d_{-\sigma}}{2v_{soft}}\right)} \\
&= \frac{1}{1 + \exp\left(\frac{1}{2v_{soft}} \left((\xi - w_\sigma)^2 - (\xi - w_{-\sigma})^2 \right) \right)} \\
&= \frac{1}{1 + \exp\left(\frac{1}{2v_{soft}} (2\xi \cdot w_{-\sigma} - 2\xi \cdot w_\sigma + w_\sigma^2 - w_{-\sigma}^2) \right)} \\
&= \frac{1}{1 + \exp\left(\frac{1}{2v_{soft}} (2h_{-\sigma} - 2h_\sigma + Q_{\sigma\sigma} - Q_{-\sigma-\sigma}) \right)} \\
&= \frac{1}{1 + \exp\left(\left(\frac{-1}{v_{soft}}, \frac{1}{v_{soft}}, 0, 0\right) \cdot (h_{+\sigma}, h_{-\sigma}, b_{+\sigma}, b_{-\sigma}) + \frac{Q_{\sigma\sigma} - Q_{-\sigma-\sigma}}{2v_{soft}}\right)} \\
&= \frac{1}{1 + \exp(\alpha_\sigma \cdot x + \beta_\sigma)} \quad (73)
\end{aligned}$$

Where, $\alpha_\sigma = \left(\frac{-\sigma}{v_{\text{soft}}}, \frac{\sigma}{v_{\text{soft}}}, 0, 0\right)$ and $\beta_\sigma = -\left(\frac{Q_{\sigma\sigma} - Q_{-\sigma-\sigma}}{2v_{\text{soft}}}\right)$.

Now consider the averages $\langle (\mathbf{x})_n \Omega_s \rangle_k$ and $\langle \Omega_s \rangle_k$:

$$\begin{aligned}
\langle (\mathbf{x})_n \Omega_s \rangle_k &= \frac{1}{(2\pi)^{\frac{1}{2}} (\det(\mathbf{C}_k))^{\frac{1}{2}}} \int_{\mathbb{R}^4} \frac{(\mathbf{x})_n}{1 + \exp(\alpha_s \cdot \mathbf{x} + \beta_s)} \\
&\quad \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_k)^T \mathbf{C}_k^{-1}(\mathbf{x} - \boldsymbol{\mu}_k)\right) d\mathbf{x} \\
&= \frac{1}{(2\pi)^2 (\det(\mathbf{C}_k))^{\frac{1}{2}}} \int_{\mathbb{R}^4} \frac{(\mathbf{x}' + \boldsymbol{\mu}_k)_n}{1 + \exp(\alpha_s \cdot \mathbf{x}' + \alpha_s \cdot \boldsymbol{\mu}_k - \beta_s)} \\
&\quad \exp\left(-\frac{1}{2}\mathbf{x}'^T \mathbf{C}_k^{-1} \mathbf{x}'\right) d\mathbf{x}' \\
&= \frac{1}{(2\pi)^2 (\det(\mathbf{C}_k))^{\frac{1}{2}}} \int_{\mathbb{R}^4} \frac{(\mathbf{C}_k^{\frac{1}{2}} \mathbf{y} + \boldsymbol{\mu}_k)_n}{1 + \exp(\alpha_s \mathbf{C}_k^{\frac{1}{2}} \mathbf{y} + \alpha_s \cdot \boldsymbol{\mu}_k - \beta_s)} \\
&\quad \exp\left(-\frac{1}{2}(\mathbf{C}_k^{\frac{1}{2}} \mathbf{y})^T \mathbf{C}_k^{-1} \mathbf{C}_k^{\frac{1}{2}} \mathbf{y}\right) (\det(\mathbf{C}_k))^{\frac{1}{2}} d\mathbf{y} \\
&= \frac{1}{(2\pi)^2} \int_{\mathbb{R}^4} \frac{(\mathbf{C}_k^{\frac{1}{2}} \mathbf{y})_n}{1 + \exp(\alpha_s \mathbf{C}_k^{\frac{1}{2}} \mathbf{y} + \alpha_s \cdot \boldsymbol{\mu}_k - \beta_s)} \exp\left(-\frac{1}{2}\mathbf{y}^2\right) d\mathbf{y} \\
&\quad + \frac{1}{(2\pi)^2} \int_{\mathbb{R}^4} \frac{(\boldsymbol{\mu}_k)_n}{1 + \exp(\alpha_s \mathbf{C}_k^{\frac{1}{2}} \mathbf{y} + \alpha_s \cdot \boldsymbol{\mu}_k - \beta_s)} \exp\left(-\frac{1}{2}\mathbf{y}^2\right) d\mathbf{y} \\
&= I + (\boldsymbol{\mu}_k)_n \langle \Omega_s \rangle_k
\end{aligned} \tag{74}$$

Where

$$\begin{aligned}
I &= \frac{1}{(2\pi)^2} \int_{\mathbb{R}^4} \frac{(\mathbf{C}_k^{\frac{1}{2}} \mathbf{y})_n}{1 + \exp(\alpha_s \mathbf{C}_k^{\frac{1}{2}} \mathbf{y} + \alpha_s \cdot \boldsymbol{\mu}_k - \beta_s)} \exp\left(-\frac{1}{2}\mathbf{y}^2\right) d\mathbf{y} \\
&= \frac{1}{(2\pi)^2} \sum_{j=1}^4 \left(\int_{\mathbb{R}} \frac{(\mathbf{C}_k^{\frac{1}{2}})_{nj}(\mathbf{y})_j}{1 + \exp(\alpha_s \mathbf{C}_k^{\frac{1}{2}} \mathbf{y} + \alpha_s \cdot \boldsymbol{\mu}_k - \beta_s)} \exp\left(-\frac{1}{2}(\mathbf{y})_j^2\right) d(\mathbf{y})_j \right) \\
&= \frac{1}{(2\pi)^2} \sum_{j=1}^4 I_j
\end{aligned} \tag{75}$$

With

$$I_j = \int_{\mathbb{R}} \frac{(\mathbf{C}_k^{\frac{1}{2}})_{nj}(\mathbf{y})_j}{1 + \exp(\alpha_s \mathbf{C}_k^{\frac{1}{2}} \mathbf{y} + \alpha_s \cdot \boldsymbol{\mu}_k - \beta_s)} \exp\left(-\frac{1}{2}(\mathbf{y})_j^2\right) d(\mathbf{y})_j \tag{76}$$

Now apply partial integration to I_j . The rule of partial integration says:
 $\int f(x)g'(x)dx = f(x)g(x) - \int g(x)f'(x)dx$. Applied to I_j this gives:

$$\begin{aligned} f((\mathbf{y})_j) &= \frac{1}{1 + \exp(\alpha_s C_k^{\frac{1}{2}} \mathbf{y} + \alpha_s \cdot \mu_k - \beta_s)} \\ g'((\mathbf{y})_j) &= (C_k^{\frac{1}{2}})_{nj} (\mathbf{y})_j \exp\left(-\frac{1}{2}(\mathbf{y})_j^2\right) \\ g((\mathbf{y})_j) &= -(C_k^{\frac{1}{2}})_{nj} \exp\left(-\frac{1}{2}(\mathbf{y})_j^2\right) \\ f'((\mathbf{y})_j) &= \frac{\partial}{\partial(\mathbf{y})_j} \frac{1}{1 + \exp(\alpha_s C_k^{\frac{1}{2}} \mathbf{y} + \alpha_s \cdot \mu_k - \beta_s)} \end{aligned} \quad (77)$$

Applying the partial integration yields:

$$\begin{aligned} I_j &= \underbrace{\left[-(C_k^{\frac{1}{2}})_{nj} \exp\left(-\frac{1}{2}(\mathbf{y})_j^2\right) \frac{1}{1 + \exp(\alpha_s C_k^{\frac{1}{2}} \mathbf{y} + \alpha_s \cdot \mu_k - \beta_s)} \right]_{-\infty}^{\infty}}_{=0} \\ &+ \int_{\mathbb{R}} (C_k^{\frac{1}{2}})_{nj} \exp\left(-\frac{1}{2}(\mathbf{y})_j^2\right) \frac{\partial}{\partial(\mathbf{y})_j} \frac{1}{1 + \exp(\alpha_s C_k^{\frac{1}{2}} \mathbf{y} + \alpha_s \cdot \mu_k - \beta_s)} d(\mathbf{y})_j \\ &= \int_{\mathbb{R}} (C_k^{\frac{1}{2}})_{nj} \exp\left(-\frac{1}{2}(\mathbf{y})_j^2\right) \frac{\partial}{\partial(\mathbf{y})_j} \frac{1}{1 + \exp(\alpha_s C_k^{\frac{1}{2}} \mathbf{y} + \alpha_s \cdot \mu_k - \beta_s)} d(\mathbf{y})_j \end{aligned} \quad (78)$$

Filling this in into equation (75) gives:

$$\begin{aligned} I &= \frac{1}{(2\pi)^2} \sum_{j=1}^4 \left(\int_{\mathbb{R}} (C_k^{\frac{1}{2}})_{nj} \exp\left(-\frac{1}{2}(\mathbf{y})_j^2\right) \frac{\partial}{\partial(\mathbf{y})_j} \frac{1}{1 + \exp(\alpha_s C_k^{\frac{1}{2}} \mathbf{y} + \alpha_s \cdot \mu_k - \beta_s)} d(\mathbf{y})_j \right) \\ &= \frac{1}{(2\pi)^2} \sum_{j=1}^4 \left((C_k^{\frac{1}{2}})_{nj} \int_{\mathbb{R}^4} \frac{\partial}{\partial(\mathbf{y})_j} \left(\frac{1}{1 + \exp(\alpha_s C_k^{\frac{1}{2}} \mathbf{y} + \alpha_s \cdot \mu_k - \beta_s)} \right) \exp\left(-\frac{1}{2}\mathbf{y}^2\right) d\mathbf{y} \right) \end{aligned} \quad (79)$$

Now consider:

$$\begin{aligned} \frac{\partial}{\partial(\mathbf{y})_j} \frac{1}{1 + \exp(\alpha_s C_k^{\frac{1}{2}} \mathbf{y} + \alpha_s \cdot \mu_k - \beta_s)} &= - \frac{\exp(\alpha_s C_k^{\frac{1}{2}} \mathbf{y} + \alpha_s \cdot \mu_k - \beta_s)}{(1 + \exp(\alpha_s C_k^{\frac{1}{2}} \mathbf{y} + \alpha_s \cdot \mu_k - \beta_s))^2} \frac{\partial}{\partial(\mathbf{y})_j} (\alpha_s C_k^{\frac{1}{2}} \mathbf{y}) \\ &= - \sum_{i=1}^4 (\alpha_s)_i (C_k^{\frac{1}{2}})_{ij} \frac{\exp(\alpha_s C_k^{\frac{1}{2}} \mathbf{y} + \alpha_s \cdot \mu_k - \beta_s)}{(1 + \exp(\alpha_s C_k^{\frac{1}{2}} \mathbf{y} + \alpha_s \cdot \mu_k - \beta_s))^2} \end{aligned} \quad (80)$$

$$\begin{aligned}
I &= -\frac{1}{(2\pi)^2} (\mathbf{C}_k \boldsymbol{\alpha}_s)_n \int_{\mathbb{R}^4} \frac{\exp(\boldsymbol{\alpha}_s \mathbf{C}_k^{\frac{1}{2}} \mathbf{y} + \boldsymbol{\alpha}_s \cdot \boldsymbol{\mu}_k - \beta_s)}{\left(1 + \exp(\boldsymbol{\alpha}_s \mathbf{C}_k^{\frac{1}{2}} \mathbf{y} + \boldsymbol{\alpha}_s \cdot \boldsymbol{\mu}_k - \beta_s)\right)^2} \exp\left(-\frac{1}{2} \mathbf{y}^2\right) d\mathbf{y} \\
&= -\frac{1}{\sqrt{2\pi}} (\mathbf{C}_k \boldsymbol{\alpha}_s)_n \int_{\mathbb{R}} \frac{\exp(\|\mathbf{C}_k^{\frac{1}{2}} \boldsymbol{\alpha}_s\| \tilde{\mathbf{y}} + \boldsymbol{\alpha}_s \cdot \boldsymbol{\mu}_k - \beta_s)}{\left(1 + \exp(\|\mathbf{C}_k^{\frac{1}{2}} \boldsymbol{\alpha}_s\| \tilde{\mathbf{y}} + \boldsymbol{\alpha}_s \cdot \boldsymbol{\mu}_k - \beta_s)\right)^2} \exp\left(-\frac{1}{2} \tilde{\mathbf{y}}^2\right) d\tilde{\mathbf{y}} \\
&= -\frac{1}{\sqrt{2\pi}} (\mathbf{C}_k \boldsymbol{\alpha}_s)_n \int_{\mathbb{R}} \frac{\exp(\tilde{\alpha}_{s,k} \tilde{\mathbf{y}} + \tilde{\beta}_{s,k})}{\left(1 + \exp(\tilde{\alpha}_{s,k} \tilde{\mathbf{y}} + \tilde{\beta}_{s,k})\right)^2} \exp\left(-\frac{1}{2} \tilde{\mathbf{y}}^2\right) d\tilde{\mathbf{y}} \quad (81)
\end{aligned}$$

Filling this in into equation (74) gives:

$$\langle (\mathbf{x})_n \Omega_s \rangle_k = -\frac{1}{\sqrt{2\pi}} (\mathbf{C}_k \boldsymbol{\alpha}_s)_n \int_{\mathbb{R}} \frac{\exp(\tilde{\alpha}_{s,k} \tilde{\mathbf{y}} + \tilde{\beta}_{s,k})}{\left(1 + \exp(\tilde{\alpha}_{s,k} \tilde{\mathbf{y}} + \tilde{\beta}_{s,k})\right)^2} \exp\left(-\frac{1}{2} \tilde{\mathbf{y}}^2\right) d\tilde{\mathbf{y}} + (\boldsymbol{\mu}_k)_n \langle \Omega_s \rangle_k \quad (82)$$

Next consider the average $\langle \Omega_s \rangle_k$:

$$\begin{aligned}
\langle \Omega_s \rangle_k &= \frac{1}{(2\pi)^{\frac{4}{2}} (\det(\mathbf{C}_k))^{\frac{1}{2}}} \int_{\mathbb{R}^4} \frac{1}{1 + \exp(\boldsymbol{\alpha}_s \cdot \mathbf{x} + \beta_s)} \\
&\quad \exp\left(-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_k)^T \mathbf{C}_k^{-1} (\mathbf{x} - \boldsymbol{\mu}_k)\right) d\mathbf{x} \\
&= \frac{1}{(2\pi)^2 (\det(\mathbf{C}_k))^{\frac{1}{2}}} \int_{\mathbb{R}^4} \frac{1}{1 + \exp(\boldsymbol{\alpha}_s \cdot \mathbf{x}' + \boldsymbol{\alpha}_s \cdot \boldsymbol{\mu}_k - \beta_s)} \\
&\quad \exp\left(-\frac{1}{2} \mathbf{x}'^T \mathbf{C}_k^{-1} \mathbf{x}'\right) d\mathbf{x}' \\
&= \frac{1}{(2\pi)^2 (\det(\mathbf{C}_k))^{\frac{1}{2}}} \int_{\mathbb{R}^4} \frac{1}{1 + \exp(\boldsymbol{\alpha}_s \mathbf{C}_k^{\frac{1}{2}} \mathbf{y} + \boldsymbol{\alpha}_s \cdot \boldsymbol{\mu}_k - \beta_s)} \\
&\quad \exp\left(-\frac{1}{2} (\mathbf{C}_k^{\frac{1}{2}} \mathbf{y})^T \mathbf{C}_k^{-1} \mathbf{C}_k^{\frac{1}{2}} \mathbf{y}\right) (\det(\mathbf{C}_k))^{\frac{1}{2}} d\mathbf{y} \\
&= \frac{1}{(2\pi)^2} \int_{\mathbb{R}^4} \frac{1}{1 + \exp(\boldsymbol{\alpha}_s \mathbf{C}_k^{\frac{1}{2}} \mathbf{y} + \boldsymbol{\alpha}_s \cdot \boldsymbol{\mu}_k - \beta_s)} \exp\left(-\frac{1}{2} \mathbf{y}^2\right) d\mathbf{y} \\
&= \frac{1}{\sqrt{2\pi}} \int_{\mathbb{R}} \frac{1}{1 + \exp(\|\boldsymbol{\alpha}_s \mathbf{C}_k^{\frac{1}{2}}\| \tilde{\mathbf{y}} + \boldsymbol{\alpha}_s \cdot \boldsymbol{\mu}_k - \beta_s)} \exp\left(-\frac{1}{2} \tilde{\mathbf{y}}^2\right) d\tilde{\mathbf{y}} \\
&= \frac{1}{\sqrt{2\pi}} \int_{\mathbb{R}} \frac{1}{1 + \exp(\tilde{\alpha}_{s,k} \tilde{\mathbf{y}} + \tilde{\beta}_{s,k})} \exp\left(-\frac{1}{2} \tilde{\mathbf{y}}^2\right) d\tilde{\mathbf{y}} \quad (83)
\end{aligned}$$

where the remaining integration has to be performed numerically.