

**PRINCIPAL MANIFOLD BASED  
CORRELATION ANALYSIS**  
applied to  
**GENE REGULATION ANALYSIS OF  
GLIOBLASTOMA MULTIFORME**

Diplomarbeit

von

Patrick Michl

Betreuer: Prof. Dr. Willi Jäger

Universität Heidelberg

Fakultät für Mathematik

September 2017

## Acknowledgements

I would like to thank my professor, Dr. Willi Jäger, for the great opportunity to write this thesis. All of his support and advice have been greatly appreciated. Furthermore, I want to thank Dr. Rainer König for leading me into the subject. I owe my interest in the intriguing area of living cells to his teaching and influence. Also I want to thank Rebecca, my parents and wonderful and great friends to whom I owe that I came this far.

# Abstract

Gene regulation analysis is a challenging task, which requires the consideration of intricate dependency structures. These structures, however, frequently are only selectively understood in terms of parametric relationships, which also impedes the derivation of meaningful correlation measures. The present thesis addresses this issue by introducing a generalized correlation measure, which is based on principal manifolds. This is motivated by recent advances within the approximation of principal manifolds by deep structured Energy Base Models. Finally the application of this approach is demonstrated for gene regulation analysis of cDNA microarray data of Glioblastoma Multiforme.

# Contents

<b>Introduction</b>	<b>1</b>
<b>1 Principal Manifold based Correlation Analysis</b>	<b>4</b>
1.1 Correlation and Regression Dilution . . . . .	5
1.2 Linear Principal Manifold based Correlation . . . . .	9
1.3 Principal Manifold based Correlation . . . . .	17
<b>2 Estimation of Principal Manifolds by Energy Based Models</b>	<b>22</b>
2.1 Markov Random Fields . . . . .	23
2.2 Energy Based Models . . . . .	29
<b>3 Approximative tractable inference in Energy Based Models</b>	<b>44</b>
3.1 Computability and Sampling in EBMs . . . . .	45
3.2 Structured EBMs for Efficient Sampling . . . . .	62
<b>4 Application to Gene Regulation Analysis</b>	<b>73</b>
4.1 Statistical Modelling of Gene Regulatory Networks by EBMs . . . . .	74
4.2 Gene Regulation Analysis of Glioblastoma Multiforme . . . . .	82
<b>Discussion</b>	<b>87</b>

# List of Figures

1.1	Principal Curve for a 2-dimensional realization . . . . .	18
1.2	Elliptical $\mathcal{M}$ -distribution in 3 dimensions . . . . .	19
2.1	Examples for cliques in a Markov Random Field . . . . .	24
3.1	Graphical Model of a Restricted Boltzmann Machine . . . . .	66
3.2	Graphical Model of a Deep Boltzmann Machine . . . . .	70
3.3	Bipartite graph structure of a Deep Boltzmann Machine . . . . .	71
4.1	Modelling Gene Regulatory Networks (GRN) by conditioning EBMs .	78
4.2	Initialise Gauss-Bernoulli DBMs by conditioning EBMs . . . . .	79
4.3	Gene Regulatory Network of GAG degradation pathway . . . . .	83
4.4	Histogram of cDNA log-ratios of GBM dataset . . . . .	84
4.5	$L$ -Correlation for selected genes in GBM dataset . . . . .	85
4.6	Predicted Gene Regulation of GAG degradation pathway in GBM by $\mathcal{M}$ -Correlations . . . . .	86

# List of Algorithms

2.1	Steepest Gradient Descent in EBMs . . . . .	40
3.1	Mini-Batch Gradient Descent in EBMs . . . . .	48
3.2	Contrastive Divergence learning in EBMs . . . . .	60

# Introduction

Various diseases, like cancer, are characterized by modifications of gene expression that target specific cellular functions to support their growth or spreading. Accordingly identifying such modifications can be of tremendous use to understand the respective mechanisms of the diseases and eventually to reveal approaches to counteract them. Fortunately, great advances in high-throughput screening technologies increasingly support the identification of gene interactions by statistical association measures. This task can be quite challenging due to intricate systematic and random variations, that for finite sample sizes of gene expression profiles are only hardly distinguishable. It is therefore required to incorporate structural beliefs. Nevertheless in the domain of gene expression structural assumptions usually are subject to uncertainties and therefore in many cases have to be treated as weak constraints, rather than proven truths. As the assumption of an underlying structure, may have severe consequences for the tractability of statistical inference, such assumptions frequently are incorporated to simplify the statistical model. In the last several decades, the bar in computational statistics was continuously raised. Furthermore this development was underpinned by increased computational capabilities. Thereby amongst different modelling approaches, that benefited from this circumstance, a substantial part can be traced back to the incorporation of a differential structure within the statistical model. This structural assumption is biologically justified by the manifold hypothesis, which assumes, that the gene expression profiles of identical cell types are to be found scattered about a low dimensional smooth submanifold, which covers large parts of the embedding space.

## Previous Works

The statistical framework to address the manifold hypothesis was introduced in the late 80s by (Hastie et al. 1989) and generalizes principal component analysis to smooth principal manifolds. Since the approximation of non-linear principal manifolds, however, is accompanied by high computational efforts, it nearly took another twenty years for the development of an applicable method (Gorban et al. 2008). An even more acknowledged breakthrough in this direction, was the utilization of deep structured Energy Based Models (Salakhudinov et al. 2009). These provide many advantages in the approximation, foremost by the consideration of unobserved causal structures, which in particular motivated an application in gene expression analysis (Angermüller et al. 2016, Chen et al. 2016, Syafiandini et al. 2016). Nevertheless, in order to finally utilize principal manifolds for the quantification of gene interactions, there are several obstacles that have to be overcome.

## Contribution

A severe gap in the literature about principal manifolds regards the unavailability of a tractable association measure that considers its smooth manifold structure: Usually in such a case one would tend to use the mutual information, which - with respect to a density estimate - measures the expected deviation between the joint distribution and the product of marginal distributions over the realizations. In Energy Based Models, however, due to the intractability of the partition function, the probability of a realization can usually not be calculated directly. In order to close this gap a new association measure is introduced, by the  $\mathcal{M}$ -Correlation, which is shown to generalize the Pearson Correlation to principal manifolds. By its definition the application of the  $\mathcal{M}$ -Correlation in turn requires an approximation of the underlying principal manifold  $\mathcal{M}$ . Therefore the methods required for such an approximation by an Energy Based Model are comprehensively reconstituted. Thereby the focus in this thesis is put on the assumptions that are required for convergence and rapid mixing. Since the thesis chronologically started with the issue to model gene expression data, it is not surprising that these assumptions afterwards turn out to be suitable in this



very domain. Therefore the use of an  $\mathcal{M}$ -Correlation based analysis is empirically demonstrated in the gene regulation analysis of Glioblastoma Multiforme.

## Overview of the Thesis

The first chapter introduces the concepts that are required to generalize the Pearson correlation to smooth submanifolds of its embedding space. Thereby an approximation of the Pearson correlation is derived, which in a first step is extended to linear principal manifolds ( $L$ -Correlation) and finally to smooth principal manifolds ( $\mathcal{M}$ -Correlation). As the evaluation of the  $\mathcal{M}$ -Correlation requires an estimation of  $\mathcal{M}$ , the second chapter introduces Energy Based Models (EBM) as a statistical framework. Thereby it is shown that for some “mild” analytical assumptions a local ML-estimation can be derived by a respective gradient descent. Afterwards the third chapter is concerned with a tractable approximation of this estimation, as well as its application to a respective gradient descent. Thereby in a first step the use of Monte Carlo (MC) integration yield the Mini-batch stochastic gradient descent (MBGD), whose convergence is assured by the incorporation of a supplementary statistical assumption. Afterwards the MC approximation by itself is approximated by a Markov Chain Monte Carlo (MCMC) approximation, which applied to the MBGD algorithm then provides the Contrastive Divergence (CD) algorithm. In order to assure the convergence of the CD algorithm, some further analytical assumptions regarding the function space are introduced. However, to also support fast convergence, a sparsity assumption is introduced, that applied to specifically structured EBMs, given by Deep Boltzmann Machines (DBM), is shown to assure rapid mixing of the Markov Chain. In the fourth chapter, the assumptions previously introduced, are evaluated with respect to their applicability to gene expression data. Finally as all conditions are met, the  $\mathcal{M}$ -Correlation is empirically studied by comparison to other correlation measures with respect to their application on Glioblastoma Multiforme.

# Chapter 1

## Principal Manifold based Correlation Analysis

### Overview

*This chapter is intended to introduce the necessary concepts that are required to generalize the Pearson correlation to smooth submanifolds of its embedding space. In this purpose the first section reconstitutes elementary properties of the Pearson correlation to derive a representation with respect to the a regression line. The second section introduces principal components to derive this respective line, and thereupon to generalize it to linear subspaces. Thereby the theory is derived for generic elliptical distributions and the principal components are introduced as an orthogonal basis of the embedding space, that decorrelates elliptically distributed random vectors. In the subsequent section the spaces, that are spanned by principal components, are used to identify the tangent spaces of principal manifolds. As principal manifolds, however, are not assured to exist for arbitrary underlying densities, a class of smooth manifold based densities is introduced, that closes this gap. Finally with the  $\mathcal{M}$ -correlation an association measure is defined, which is shown to generalize the Pearson correlation to densities of this respective class.*

## 1.1 Correlation and Regression Dilution

A fundamental issue, that accompanies the analysis of multivariate data, concerns the quantification of statistically dependency structures by association measures. Many approaches in this direction can be traced back to the late 19<sup>th</sup> century, where the issue was closely related to the task, to extract laws of nature from two dimensional scatter plots. This in particular applies to the widespread Pearson correlation coefficient.

**Definition** (Pearson Correlation). *Let  $X: \Omega \rightarrow \mathbb{R}$  and  $Y: \Omega \rightarrow \mathbb{R}$  be random variables with finite variances  $\sigma_X^2$  and  $\sigma_Y^2$ . Then the Pearson correlation  $\rho_{X,Y}$  is defined by:*

$$\rho_{X,Y} := \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y} \quad (1.1)$$

Due to its popularity and simplicity the Pearson correlation has been generalized to a variety of different domains of application, including generic monotonous relationships, relationships between sets of random variables and asymmetric relationships (Zheng et al. 2010). In the purpose to provide a generalization to smooth curves and submanifolds, that allow an incorporation of structural assumptions, some elementary considerations have to be taken into account, that allow a separation between the pairwise quantification of dependencies and their global modelling. Pearson’s original motivation, was the regression of a straight line, that minimizes the averaged Euclidean distance to points, that are scattered about it (Pearson 1901, p561). Thereby his investigations were preceded by the observation, that for a measurement series the assumed “direction of causality” influences the estimate of the slope of the regression line. Thereby the direction of causality is implicated by the choice of an error model, that assumes one random variable to be error free and the other to account for the whole observed error. Pearson empirically observed, that for  $n \in \mathbb{N}$  points, given by i.i.d. realizations  $\mathbf{x} \in \mathbb{R}^n$  of  $X$  and  $\mathbf{y} \in \mathbb{R}^n$  of  $Y$ , the **least squares** regression line of  $\mathbf{y}$  on  $\mathbf{x}$  only equals the regression line of  $\mathbf{x}$  on  $\mathbf{y}$ , if all points perfectly fit on a straight line. In all other cases, however, the slopes of the respective regression lines turned out, not to be reciprocal and their product was found within the interval  $[0, 1)$ . This observation was decisive for Pearson’s definition of the correlation coefficient. Thereby  $\rho_{X,Y}$  is estimated by its empirical counterpart  $\rho_{x,y}$ , that replaces variances by sample variances and the covariance by the sample variance.

**Lemma 1.1.** *Let  $X: \Omega \rightarrow \mathbb{R}$  and  $Y: \Omega \rightarrow \mathbb{R}$  be random variables with  $n \in \mathbb{N}$  i.i.d. realizations  $\mathbf{x} \in \mathbb{R}^n$  and  $\mathbf{y} \in \mathbb{R}^n$ . Furthermore let  $\beta_x \in \mathbb{R}$  denote the slope of the linear regression of  $\mathbf{y}$  on  $\mathbf{x}$  and  $\beta_y \in \mathbb{R}$  the slope of the linear regression of  $\mathbf{x}$  on  $\mathbf{y}$ . Then:*

$$\rho_{x,y}^2 = \beta_x \beta_y \quad (1.2)$$

*Proof of Lemma 1.1.* The following proof is based on (Kenney et al. 1962). The least squares regression of  $\mathbf{y}$  on  $\mathbf{x}$  implicates, that for regression coefficients  $\alpha_x, \beta_x \in \mathbb{R}$  and a normal distributed random error  $\varepsilon := Y - (\beta_x X + \alpha_x)$  the log-likelihood of the realizations is maximized, if and only if the  $\ell^2$ -norm of the realizations of  $\varepsilon$  is minimized, such that:

$$SSE_y(\alpha_x, \beta_x) := \sum_{i=1}^n (y_i - (\beta_x x_i + \alpha_x))^2 \rightarrow \min \quad (1.3)$$

Since  $SSE_y$  is a quadratic function of  $\alpha_x$  and  $\beta_x$  and therefore convex, it has a unique global minimum at:

$$\frac{\partial}{\partial \alpha_x} SSE_y = 2 \sum_{i=1}^n (y_i - (\beta_x x_i + \alpha_x))(-x_i) = 0 \quad (1.4)$$

$$\frac{\partial}{\partial \beta_x} SSE_y = 2 \sum_{i=1}^n (y_i - (\beta_x x_i + \alpha_x))(-1) = 0 \quad (1.5)$$

By equating the coefficients, equations 1.4 and 1.5 can be rewritten as a system of linear equations of  $\alpha_x$  and  $\beta_x$ :

$$\alpha_x n + \beta_x \sum_{i=1}^n x_i = \sum_{i=1}^n y_i \quad (1.6)$$

$$\alpha_x \sum_{i=1}^n x_i + \beta_x \sum_{i=1}^n x_i^2 = \sum_{i=1}^n x_i y_i \quad (1.7)$$

Consequently in matrix notation the vector  $(\alpha_x, \beta_x)^T$  is determined by:

$$\begin{pmatrix} \alpha_x \\ \beta_x \end{pmatrix} = \begin{pmatrix} n & \sum_{i=1}^n x_i \\ \sum_{i=1}^n x_i & \sum_{i=1}^n x_i^2 \end{pmatrix}^{-1} \begin{pmatrix} \sum_{i=1}^n y_i \\ \sum_{i=1}^n x_i y_i \end{pmatrix} \quad (1.8)$$

Let  $\bar{x}$ ,  $\bar{y}$  respectively denote the *sample means*. Then by calculating the matrix inverse, the slope  $\beta_x$  equates to:

$$\beta_x = \left( \sum_{i=1}^n x_i y_i - n \bar{x} \bar{y} \right) \left( \sum_{i=1}^n x_i^2 - n \bar{x}^2 \right)^{-1} \quad (1.9)$$

Thereupon by substituting the *sample variance*:

$$\begin{aligned} \sigma_x^2 &:= \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 \\ &= \frac{1}{n} \left( \sum_{i=1}^n x_i^2 - n \bar{x}^2 \right) \end{aligned} \quad (1.10)$$

And the *sample covariance*:

$$\begin{aligned} \text{Cov}(\mathbf{x}, \mathbf{y}) &:= \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) \\ &= \frac{1}{n} \left( \sum_{i=1}^n x_i y_i - n \bar{x} \bar{y} \right) \end{aligned} \quad (1.11)$$

It follows from equation 1.9, that:

$$\beta_x = \frac{\text{Cov}(\mathbf{x}, \mathbf{y})}{\sigma_x^2} \quad (1.12)$$

Conversely the slope  $\beta_y$  of the linear regression of  $\mathbf{x}$  on  $\mathbf{y}$  mutatis mutandis equates to:

$$\beta_y = \frac{\text{Cov}(\mathbf{y}, \mathbf{x})}{\sigma_y^2} \quad (1.13)$$

By the symmetry of Cov it the follows, from equations 1.12 and 1.13 that:

$$\beta_x \beta_y = \frac{\text{Cov}(\mathbf{x}, \mathbf{y})^2}{\sigma_x^2 \sigma_y^2} = \rho_{xy}^2 \quad (1.14)$$

□

Lemma 1.1 shows, that  $\rho_{x,y}$  may be regarded as the geometric mean of the regression slopes  $\beta_x$  and  $\beta_y$ , where  $\beta_x$  and  $\beta_y$  respectively describe the causal relationships  $X \rightarrow Y$  and  $Y \rightarrow X$ . Thereby  $X$  and  $Y$  respectively are treated as error free regressor

variables to predict the corresponding response variable, that captures the overall error. The mutual linear relationship  $X \leftrightarrow Y$  is then described by a regression line, that equally treats errors in both variables. As an immediate consequence of this symmetry it follows, that this **total least squares** regression line is unique, and its slope  $\beta_x^*$ , that describes  $\mathbf{y}$  by  $\mathbf{x}$  is reciprocal to the slope  $\beta_y^*$ , that describes  $\mathbf{x}$  by  $\mathbf{y}$  such that  $\beta_x^* \beta_y^* = 1$ . In this sense  $\beta_x$  and  $\beta_y$  may be regarded as biased estimations of  $\beta_x^*$  and  $\beta_y^*$ . Thereby the bias generally is known as “regression dilution” or “regression attenuation”. For the case that both errors are independent and normal distributed, this bias can be corrected by a prefactor, that incorporates the error of the respective regressor variable. An application of this correction to lemma 1.1 then shows, that  $\rho_{X,Y}$  has a consistent estimations by the sample variances of  $\mathbf{x}$  and  $\mathbf{y}$  and the variances of their respective errors  $\varepsilon_X$  and  $\varepsilon_Y$ .

**Proposition 1.1.** *Let  $X: \Omega \rightarrow \mathbb{R}$  and  $Y: \Omega \rightarrow \mathbb{R}$  be random variables with  $n \in \mathbb{N}$  i.i.d. realizations  $\mathbf{x} \in \mathbb{R}^n$  and  $\mathbf{y} \in \mathbb{R}^n$  and random errors  $\varepsilon_X \sim \mathcal{N}(0, \eta_X^2)$  and  $\varepsilon_Y \sim \mathcal{N}(0, \eta_Y^2)$ . Then:*

$$\rho_{x,y}^2 \xrightarrow{P} \left(1 - \frac{\eta_X^2}{\sigma_x^2}\right) \left(1 - \frac{\eta_Y^2}{\sigma_y^2}\right), \text{ for } n \rightarrow \infty \quad (1.15)$$

*Proof of Proposition 1.1.* Let  $\beta_x \in \mathbb{R}$  be the slope of the ordinary least squares (OLS) regression line of  $\mathbf{y}$  on  $\mathbf{x}$ , where  $\mathbf{x}$  is assumed to realize  $X$  with a normal distributed random error  $\varepsilon_X \sim \mathcal{N}(0, \eta_X^2)$ . Then  $X$  decomposes into (i) an unobserved error free regressor variable  $X^*$  and (ii) the random error  $\varepsilon_X$ , such that:

$$X \sim X^* + \varepsilon_X \quad (1.16)$$

With respect to this decomposition, the slope  $\beta_x^*$  of the total least squares (TLS) regression line, that also considers  $\varepsilon_X$ , then is identified by the slope of the OLS regression of  $\mathbf{y}$  on  $\mathbf{x}^*$ , where  $\mathbf{x}^*$  realizes  $X^*$ . Thereupon let  $\sigma_{x^*}^2$  be the empirical variance of  $\mathbf{x}^*$ , then according to (Snedecor et al. 1967) it follows, that:

$$\beta_x \xrightarrow{P} \frac{\sigma_{x^*}^2}{\sigma_{x^*}^2 + \eta_X^2} \beta_x^*, \text{ for } n \rightarrow \infty \quad (1.17)$$

Since furthermore  $\varepsilon_X$  by definition is statistically independent from  $X^*$ , it can be

concluded, that:

$$\begin{aligned}\sigma_x^2 &= \text{Var}(X^* + \varepsilon_X) \\ &= \text{Var}(X^*) + \text{Var}(\varepsilon_X) = \sigma_{x^*}^2 + \eta_X^2\end{aligned}\tag{1.18}$$

Such that:

$$\frac{\sigma_{x^*}^2}{\sigma_{x^*}^2 + \eta_X^2} \stackrel{1.18}{=} \frac{\sigma_x^2 - \eta_X^2}{\sigma_x^2} = 1 - \frac{\eta_X^2}{\sigma_x^2}\tag{1.19}$$

And therefore by equation 1.17 that:

$$\beta_x \xrightarrow{P} \left(1 - \frac{\eta_X^2}{\sigma_x^2}\right) \beta_x^*, \text{ for } n \rightarrow \infty\tag{1.20}$$

Conversely let now  $\beta_y \in \mathbb{R}$  be the the slope of the OLS regression of  $\mathbf{x}$  on  $\mathbf{y}$ , where  $\mathbf{y}$  is assumed to realize  $Y$  with a random error  $\varepsilon_Y \sim \mathcal{N}(0, \eta_Y^2)$ . Then also the corrected slope  $\beta_y^*$  mutatis mutandis satisfies the relation given by equation 1.20 and by the representation of  $\rho_{x,y}$ , as given by lemma 1.1, it then can be concluded, that:

$$\rho_{x,y}^2 \stackrel{1.1}{=} \beta_x \beta_y \xrightarrow{P} \left(1 - \frac{\eta_X^2}{\sigma_x^2}\right) \left(1 - \frac{\eta_Y^2}{\sigma_y^2}\right) \beta_x^* \beta_y^*, \text{ for } n \rightarrow \infty\tag{1.21}$$

The proposition then follows by the uniqueness of the total least squares regression line for known variances  $\eta_X^2$  and  $\eta_Y^2$ , such that:

$$\beta_y^* = \frac{1}{\beta_x^*}$$

□

## 1.2 Linear Principal Manifold based Correlation

Within the same publication, in which Pearson introduced the correlation coefficient, he also developed a structured approach that determines the straight line, that minimizes the Euclidean distance (Pearson 1901, p563). His method, which later received attribution as the method of **Principal Component Analysis (PCA)**, however, even went further and allowed a canonical generalization of the problem in the following sense: For  $d \in \mathbb{N}$  let  $\mathbf{X}: \Omega \rightarrow \mathbb{R}^d$  be a multivariate random vector and for

$n \in \mathbb{N}$  let  $\mathbf{x} \in \mathbb{R}^{n \times d}$  be an i.i.d. realization of  $\mathbf{X}$ . Then for any given  $k \in \mathbb{N}$  with  $k \leq d$  the goal is, to determine an affine linear subspace  $L \subseteq \mathbb{R}^d$  of dimension  $k$ , that minimizes the summed Euclidean distance to  $\mathbf{x}$ . In order to solve this problem, the fundamental idea of Pearson was, to transfer the principal axis theorem from ellipsoids to multivariate Gaussian distributed random vectors. Thereupon, however, the method also can be formulated with respect to generic elliptical distributions.

**Definition** (Elliptical Distribution). *For  $d \in \mathbb{N}$  let  $\mathbf{X}: \Omega \rightarrow \mathbb{R}^d$  be a random vector. Then  $\mathbf{X}$  is elliptically distributed, iff there exists a random vector  $\mathbf{S}: \Omega \rightarrow \mathbb{R}^k$  with  $k \leq d$ , which distribution is invariant to rotations, a matrix  $A \in \mathbb{R}^{d \times k}$  of rank  $k$  and a vector  $\mathbf{b} \in \mathbb{R}^d$ , such that:*

$$\mathbf{X} \sim A\mathbf{S} + \mathbf{b} \quad (1.22)$$

Consequently a random vector  $\mathbf{X}$  is elliptically distributed, if it can be represented by an affine transformation of a radial symmetric distributed random vector  $\mathbf{S}$ . The decisive property, that underpins the choice of elliptical distributions, lies within their coincidence of linear and statistical dependencies, which allows to decompose  $\mathbf{X}$  in statistically independent components by a linear decomposition. This property allows, to substantiate the multidimensional “linear fitting problem” with respect to an orthogonal projection.

**Proposition 1.2.** *For  $d \in \mathbb{N}$  let  $\mathbf{X}: \Omega \rightarrow \mathbb{R}^d$  be an elliptically distributed random vector,  $L \subseteq \mathbb{R}^d$  an affine linear subspace of  $\mathbb{R}^d$  and  $\pi_L$  the orthogonal projection of  $\mathbb{R}^d$  onto  $L$ . Then the following statements are equivalent:*

- (i)  $L$  minimizes the Euclidean distance to  $\mathbf{X}$
- (ii)  $\mathbb{E}(\mathbf{X}) \in L$  and  $L$  maximizes the variance  $\text{Var}(\pi_L(\mathbf{X}))$

*Proof.* Let  $\mathbf{Y}_L := \mathbf{X} - \pi_L(\mathbf{X})$ , then the Euclidean distance between  $\mathbf{X}$  and  $L$  can be written:

$$d(\mathbf{X}, \pi_L(\mathbf{X}))^2 = \mathbb{E}(\|\mathbf{X} - \pi_L(\mathbf{X})\|_2^2) = \mathbb{E}(\mathbf{Y}_L^2) \quad (1.23)$$

This representation can furthermore be decomposed by using the algebraic formula for the variance:

$$\mathbb{E}(\mathbf{Y}_L^2) = \text{Var}(\mathbf{Y}_L) + \mathbb{E}(\mathbf{Y}_L)^2 \quad (1.24)$$



Let now be  $\mathbf{Y}_L^\perp := \pi_L(\mathbf{X})$ , then  $\mathbf{X} = \mathbf{Y}_L + \mathbf{Y}_L^\perp$  and  $\mathbf{Y}_L$  and  $\mathbf{Y}_L^\perp$  are uncorrelated, such that:

$$\text{Var}(\mathbf{X}) = \text{Var}(\mathbf{Y}_L + \mathbf{Y}_L^\perp) = \text{Var}(\mathbf{Y}_L) + \text{Var}(\mathbf{Y}_L^\perp) \quad (1.25)$$

From equations 1.23, 1.24 and 1.25 it follows, that:

$$d(\mathbf{X}, \pi_L(\mathbf{X}))^2 = \text{Var}(\mathbf{X}) - \text{Var}(\mathbf{Y}_L^\perp) + \mathbb{E}(\mathbf{Y}_L)^2 \quad (1.26)$$

Consequentially the Euclidean distance is minimized, if and only if the right side of equation 1.26 is minimized. The first term  $\text{Var}(\mathbf{X})$ , however, does not depend on  $L$  and since  $\mathbf{X}$  is elliptically distributed, the linear independence of  $\mathbf{Y}_L^\perp$  and  $\mathbf{Y}_L$  is sufficient for statistical independence. It follows, that the Euclidean distance is minimized, if and only if: (1) The term  $\mathbb{E}(\mathbf{Y}_L)^2$  is minimized and (2) the term  $\text{Var}(\mathbf{Y}_L^\perp)$  is maximized. Concerning (1) it follows, that:

$$\mathbb{E}(\mathbf{Y}_L)^2 = \mathbb{E}(\mathbf{X} - \pi_L(\mathbf{X}))^2 = (\mathbb{E}(\mathbf{X}) - \pi_L(\mathbb{E}(\mathbf{X})))^2$$

Therefore the term  $\mathbb{E}(\mathbf{Y}_L)^2$  is minimized, if and only if  $\pi_L(\mathbb{E}(\mathbf{X})) = \mathbb{E}(\mathbf{X})$ , which in turn means that  $\mathbb{E}(\mathbf{X}) \in L$ . Concerning (2), the proposition immediately follows by the definition of  $\mathbf{Y}_L^\perp$ .  $\square$

Let now be  $k \leq d$ . In order to derive an affine linear subspace  $L \subseteq \mathbb{R}^d$  that minimizes the Euclidean distance to  $\mathbf{X}$ , proposition ... states, that it suffices to provide an  $L$  which (1) is centred in  $\mathbf{X}$ , such that  $\mathbb{E}(\mathbf{X}) \in L$ , and (2) maximizes the variance of the projection. In order to maximize  $\text{Var}(\pi_L(\mathbf{X}))$ , however, it is beneficial to give a further representation.

**Lemma 1.2.** *For  $d \in \mathbb{N}$  let  $\mathbf{X}: \Omega \rightarrow \mathbb{R}^d$  be an elliptically distributed random vector,  $L \subseteq \mathbb{R}^d$  an affine linear subspace of  $\mathbb{R}^d$ , which for an  $k \leq d$ , a vector  $\mathbf{v} \in \mathbb{R}^d$  and an orthonormal basis  $\mathbf{u}_1, \dots, \mathbf{u}_k \in \mathbb{R}^d$  is given by:*

$$L = \mathbf{v} + \bigoplus_{i=1}^k \mathbb{R}\mathbf{u}_i$$

*Let further be  $\pi_L: \mathbb{R}^d \rightarrow L$  the orthogonal projection of  $\mathbb{R}^d$  onto  $L$ . Then the variance*

of the projection is given by:

$$\text{Var}(\pi_L(\mathbf{X})) = \sum_{i=1}^k \mathbf{u}_i^T \text{Cov}(\mathbf{X}) \mathbf{u}_i$$

*Proof.* Let  $L' := L + \mathbb{E}(\mathbf{X}) - \mathbf{v}$ , then the orthogonal projection  $\pi_{L'}(\mathbf{X})$  decomposes into individual orthogonal projections to the respective basis vectors, such that:

$$\pi_{L'}(\mathbf{X}) = \mathbb{E}(\mathbf{X}) + \sum_{i=1}^k \langle \mathbf{X} - \mathbb{E}(\mathbf{X}), \mathbf{u}_i \rangle \mathbf{u}_i \quad (1.27)$$

Let  $\hat{X}_i := \langle \mathbf{X}, \mathbf{u}_i \rangle \mathbf{u}_i$ , for  $i \in \{1, \dots, k\}$ . The total variance of this projection is then given by:

$$\text{Var}(\pi_{L'}(\mathbf{X})) \stackrel{1.27}{=} \text{Var} \left( \mathbb{E}(\mathbf{X}) + \sum_{i=1}^k \hat{X}_i - \sum_{i=1}^k \mathbb{E}(\hat{X}_i) \right) = \text{Var} \left( \sum_{i=1}^k \hat{X}_i \right) \quad (1.28)$$

Since the random variables  $\hat{X}_i$  by definition are uncorrelated, the algebraic formula for the variance can be used to decompose the variance:

$$\text{Var} \left( \sum_{i=1}^k \hat{X}_i \right) = \sum_{i=1}^k \text{Var}(\hat{X}_i) \quad (1.29)$$

By equating the term  $\text{Var}(\hat{X}_i)$ , for  $i \in \{1, \dots, k\}$  it follows, that:

$$\text{Var}(\hat{X}_i) = \text{Var}(\langle \mathbf{X}, \mathbf{u}_i \rangle \mathbf{u}_i) = \text{Var}(\mathbf{X}^T \mathbf{u}_i) \mathbf{u}_i^2 = \text{Var}(\mathbf{X}^T \mathbf{u}_i) \quad (1.30)$$

And furthermore by introducing the covariance matrix  $\text{Cov}(\mathbf{X})$ :

$$\text{Var}(\mathbf{X}^T \mathbf{u}_i) \stackrel{\text{def}}{=} \mathbb{E}((\mathbf{X}^T \mathbf{u}_i)^T (\mathbf{X}^T \mathbf{u}_i)) = \mathbf{u}_i^T \mathbb{E}(\mathbf{X}^T \mathbf{X}) \mathbf{u}_i \stackrel{\text{def}}{=} \mathbf{u}_i^T \text{Cov}(\mathbf{X}) \mathbf{u}_i \quad (1.31)$$

Summarized the equations 1.28, 1.29, 1.30 and 1.31 provide a representation for the variance of the projection to  $L'$ :

$$\text{Var}(\pi_{L'}(\mathbf{X})) = \sum_{i=1}^k \mathbf{u}_i^T \text{Cov}(\mathbf{X}) \mathbf{u}_i$$

Finally the total variance of the projection is invariant under translations of  $L$ , such that:

$$\begin{aligned} \text{Var}(\pi_L(\mathbf{X})) &= \text{Var}(\pi_L(\mathbf{X}) + \mathbb{E}(\mathbf{X}) - \mathbf{v}) = \text{Var}(\pi_{L'}(\mathbf{X})) \\ &\stackrel{1.29}{=} \sum_{i=1}^k \mathbf{u}_i^T \text{Cov}(\mathbf{X}) \mathbf{u}_i \end{aligned} \quad (1.32)$$

□

Lemma 1.2, shows, that for elliptically distributed random vectors  $\mathbf{X}$  the best fitting linear subspaces are completely determined by the expectation  $\mathbb{E}(\mathbf{X})$  and the covariance matrix  $\text{Cov}(\mathbf{X})$ . On this point it is important to notice, that the covariance matrix is symmetric, which allows its diagonalization with regard to real valued Eigenvalues.

**Lemma 1.3.** *For  $d \in \mathbb{N}$  let  $\mathbf{X}: \Omega \rightarrow \mathbb{R}^d$  be an elliptically distributed random vector,  $L \subseteq \mathbb{R}^d$  an affine linear subspace of  $\mathbb{R}^d$ , which for an  $k \leq d$ , a vector  $\mathbf{v} \in \mathbb{R}^d$  and an orthonormal basis  $\mathbf{u}_1, \dots, \mathbf{u}_k \in \mathbb{R}^d$  is given by:*

$$L = \mathbf{v} + \bigoplus_{i=1}^k \mathbb{R} \mathbf{u}_i$$

*Let further be  $\pi_L: \mathbb{R}^d \rightarrow L$  the orthogonal projection of  $\mathbb{R}^d$  onto  $L$ , as well as  $\lambda_1, \dots, \lambda_d \in \mathbb{R}$  the eigenvalues of  $\text{Cov}(\mathbf{X})$ . Then there exist numbers  $a_1, \dots, a_d \in [0, 1]$  with  $\sum_{i=1}^d a_i = k$ , such that:*

$$\text{Var}(\pi_L(\mathbf{X})) = \sum_{i=1}^d \lambda_i a_i$$

*Proof.* From lemma 1.2 it follows, that:

$$\text{Var}(\pi_L(\mathbf{X})) = \sum_{j=1}^k \mathbf{u}_j^T \text{Cov}(\mathbf{X}) \mathbf{u}_j$$

Since the covariance matrix  $\text{Cov}(\mathbf{X})$  is a symmetric matrix, there exists an orthonormal basis transformation matrix  $S \in \mathbb{R}^{d \times d}$  and a diagonal matrix  $D \in \mathbb{R}^{d \times d}$ , such

that  $\text{Cov}(\mathbf{X}) = S^T D S$ . Then the variance  $\text{Var}(\pi_L(\mathbf{X}))$  has a decomposition, given by:

$$\text{Var}(\pi_L(\mathbf{X})) = \sum_{j=1}^k \mathbf{u}_j^T S^T D S \mathbf{u}_j = \sum_{j=1}^k (S \mathbf{u}_j)^T D S \mathbf{u}_j$$

For  $j \in \{1, \dots, k\}$  let now  $\mathbf{c}_j := S \mathbf{u}_j$  and for  $i \in \{1, \dots, n\}$  let the number  $a_i \in \mathbb{R}$  be defined by:

$$a_i := \sum_{j=1}^k (\mathbf{c}_{ji})^2$$

Then according to Lemma 1.2 the variance  $\text{Var}(\pi_L(\mathbf{X}))$  can be decomposed:

$$\begin{aligned} \sum_{j=1}^k \mathbf{c}_j^T D \mathbf{c}_j &= \sum_{j=1}^k \sum_{i=1}^d \mathbf{c}_{ji} \lambda_i \mathbf{c}_{ji} \\ &= \sum_{i=1}^d \lambda_i \sum_{j=1}^k (\mathbf{c}_{ji})^2 = \sum_{i=1}^d \lambda_i a_i \end{aligned}$$

Furthermore since  $\mathbf{u}_1, \dots, \mathbf{u}_k$  is an orthonormal basis and  $S$  an orthonormal matrix it follows that also  $\mathbf{c}_1, \dots, \mathbf{c}_k$  is an orthonormal basis. Consequentially for  $i \in \{1, \dots, d\}$  it holds, that:

$$a_i = \sum_{j=1}^k (\mathbf{c}_{ji})^2 \leq \sum_{j=1}^k \|\mathbf{c}_{ji}\|_2 \leq 1$$

And furthermore by its definition it follows, that  $a_i \geq 0$ , such that  $a_i \in [0, 1]$ . Besides this the sum over all  $a_i$  equates to:

$$\sum_{i=1}^d a_i = \sum_{i=1}^d \sum_{j=1}^k (\mathbf{c}_{ji})^2 = \sum_{j=1}^k \mathbf{c}_j^T \mathbf{c}_j = k$$

□

With reference to the principal axis transformation, the eigenvectors of the covariance matrix are then termed **principal components** and affine linear subspaces of the embedding space as **linear principal manifolds**.

**Definition** (Linear Principal Manifold). For  $d \in \mathbb{N}$  let  $\mathbf{X}: \Omega \rightarrow \mathbb{R}^d$  be a random vector. Then a vector  $\mathbf{c} \in \mathbb{R}^d$  with  $\mathbf{c} \neq 0$  is a principal component for  $\mathbf{X}$ , iff there exists an  $\lambda \in \mathbb{R}$ , such that :

$$\text{Cov}(\mathbf{X}) \cdot \mathbf{c} = \lambda \mathbf{c} \quad (1.33)$$

Furthermore let  $L \subseteq \mathbb{R}^d$  be an affine linear subspace of  $\mathbb{R}^d$  with dimension  $k \leq d$ . Then  $L$  is a linear  $k$ -principal manifold for  $\mathbf{X}$ , if there exists a set  $\mathbf{c}_1, \dots, \mathbf{c}_k$  of linear independent principal components for  $\mathbf{X}$ , such that:

$$L = \mathbb{E}(\mathbf{X}) + \bigoplus_{i=1}^k \mathbb{R}\mathbf{c}_i$$

Then  $L$  is termed maximal, iff the sum of the Eigenvalues  $\lambda_1, \dots, \lambda_k$ , that correspond to the principal components  $\mathbf{c}_1, \dots, \mathbf{c}_k$  is maximal.

**Proposition 1.3.** For  $d \in \mathbb{N}$  let  $\mathbf{X}: \Omega \rightarrow \mathbb{R}^d$  be an elliptically distributed random vector and  $L \subseteq \mathbb{R}^d$  an affine linear subspace. Then the following statements are equivalent:

- (i)  $L$  minimizes the Euclidean distance to  $\mathbf{X}$
- (ii)  $L$  is a maximal linear principal manifold for  $\mathbf{X}$

*Proof.* “ $\implies$ ” Let  $\pi_L: \mathbb{R}^d \hookrightarrow L$  denote the orthogonal projection of  $\mathbb{R}^d$  onto  $L$ . Then according to proposition 1.2  $L$  minimizes the averaged Euclidean distance to  $\mathbf{X}$ , if and only if (i)  $\mathbb{E}(\mathbf{X}) \in L$  and (ii)  $L$  maximizes the variance  $\text{Var}(\pi_L(\mathbf{X}))$ . In particular (i) is satisfied, if and only if an orthonormal basis  $\mathbf{u}_1, \dots, \mathbf{u}_k \in \mathbb{R}^d$  can be chosen, such that:

$$L = \mathbb{E}(\mathbf{X}) + \bigoplus_{i=1}^k \mathbb{R}\mathbf{u}_i$$

Then according to Lemma 1.3 there exist numbers  $a_1, \dots, a_d \in [0, 1]$  with  $\sum_{i=1}^d a_i = k$ , such that:

$$\text{Var}(\pi_L(\mathbf{X})) = \sum_{i=1}^d \lambda_i a_i$$

Thereupon (ii) is satisfied, if and only if the numbers  $a_i$  maximize this sum. Since the covariance matrix  $\text{Cov}(\mathbf{X})$  is positive semi-definite, the eigenvalues  $\lambda_i$  are not negative such that the sum is maximized for:

$$a_i = \begin{cases} 1 & \text{for } i \in \{1, \dots, k\} \\ 0 & \text{else} \end{cases}$$

Such that:

$$\begin{aligned} \sum_{j=1}^k \mathbf{u}_j^T \text{Cov}(\mathbf{X}) \mathbf{u}_j &= \sum_{i=1}^d \lambda_i a_i \\ &= \sum_{i=1}^k \lambda_i = \sum_{j=1}^k \mathbf{c}_j^T \text{Cov}(\mathbf{X}) \mathbf{c}_j \end{aligned}$$

Accordingly the choice  $\mathbf{u}_j = \mathbf{c}_j$  for  $i \in \{1, \dots, k\}$  maximizes  $\text{Var}(\pi_L(\mathbf{X}))$  and  $L$  has a representation, given by:

$$L = \mathbb{E}(\mathbf{X}) + \bigoplus_{i=1}^k \mathbb{R} \mathbf{c}_i$$

” $\Leftarrow$ ” Let  $L$  have a representation as given by (ii), then (1)  $\mathbb{E}(\mathbf{X}) \in L$  and (2) the variance  $\text{Var}(\pi_L(\mathbf{X}))$  is maximized. According to Proposition 1.2 it follows, that  $L$  minimizes the Euclidean distance to  $\mathbf{X}$ .  $\square$

**Definition** ( $L$ -Correlation). For  $d \in \mathbb{N}$  let  $\mathbf{X}: \Omega \rightarrow \mathbb{R}^d$  be a random vector and  $L$  a maximal linear principal manifold for  $\mathbf{X}$ . Then for any  $i, j \in \{1, \dots, d\}$  let the  $L$ -Correlation between  $X_i$  and  $X_j$  be defined by:

$$\rho_{X_i, X_j|L}^2 := R_i R_j \tag{1.34}$$

where with the orthogonal projection  $\pi_L: \mathbb{R}^d \rightarrow L$  for any  $i \in \{1, \dots, d\}$  the **reliability** of  $X_i$  with respect to  $L$  is given by:

$$R_i := 1 - \frac{\text{Var}_i(\mathbf{X} - \pi_L(\mathbf{X}))}{\text{Var}_i(\mathbf{X})} \tag{1.35}$$

**Proposition 1.4.** For an elliptically distributed random vector the  $L$ -Correlation generalizes the Pearson Correlation to maximal linear principal manifolds.

*Proof.* Let  $\mathbf{X}: \Omega \rightarrow \mathbb{R}^2$  be an elliptically distributed random vector and  $L$  a maximal

linear 1-principal manifold for  $\mathbf{X}$ . Then for  $i \in \{1, 2\}$  the random error of the variable  $X_i$  has a variance:

$$\eta_{X_i}^2 = \text{Var}_{X_i}(\mathbf{X} - \pi_L(\mathbf{X}))$$

Such that by the definition of the reliability it follows, that:

$$R_i \stackrel{1.35}{=} 1 - \frac{\eta_{X_i}^2}{\sigma_{X_i}^2}$$

Consequently:

$$\rho_{X_i, X_j|L}^2 = \left(1 - \frac{\eta_{X_i}^2}{\sigma_{X_i}^2}\right) \left(1 - \frac{\eta_{X_j}^2}{\sigma_{X_j}^2}\right)$$

With  $n \in \mathbb{N}$  i.i.d. realizations  $\mathbf{x} \in \mathbb{R}^{n \times 2}$  of  $\mathbf{X}$  an empirical  $L$ -Correlation  $\rho_{x_i, x_j|L}^2$  is then given by replacing the variances by the sample variances. Then by proposition 1.1 it follows, that:

$$\rho_{x_i, x_j}^2 \xrightarrow{P} \rho_{x_i, x_j|L}^2, \text{ for } n \rightarrow \infty$$

□

### 1.3 Principal Manifold based Correlation

Linear principal manifolds allow the projection of a random vector  $\mathbf{X}: \Omega \rightarrow \mathbb{R}^d$  onto a linear subspace  $L \subseteq \mathbb{R}^d$ , which maximally preserves the linear dependency structure of  $\mathbf{X}$  in terms of its covariances. Thereby for the orthogonal projection  $\pi_L: \mathbb{R}^d \hookrightarrow L$ , the variance on  $L$ , given by  $\text{Var}(\pi_L(\mathbf{X}))$ , is referred as the **explained variance** and the orthogonal deviation  $\text{Var}(\mathbf{X} - \pi_L(\mathbf{X}))$  as the **unexplained variance**. Thereupon by the assumption, that  $\mathbf{X}$  is elliptically distributed, it can be concluded, that linear independence coincides with statistical independence, that that  $\pi_L(\mathbf{X})$  and  $\mathbf{X} - \pi_L(\mathbf{X})$  are statistically independent and therefore allow the following decomposition:

$$\underbrace{\text{Var}(\mathbf{X})}_{\text{total variance}} = \underbrace{\text{Var}(\pi_L(\mathbf{X}))}_{\text{explained variance}} + \underbrace{\text{Var}(\mathbf{X} - \pi_L(\mathbf{X}))}_{\text{unexplained variance}}$$

This decomposition, as shown by theorem 1.4, is of fundamental importance for the correlation over linear Principal Manifolds, since it determines the reliabilities of the

respective random variable  $X$  by the ratio:

$$R = 1 - \frac{\text{explained variance}}{\text{total variance}}$$

On this point of the discussion it's just a small step to generalize the principal

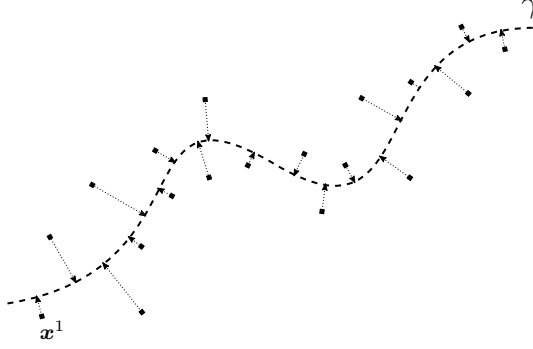


Figure 1.1: Principal Curve for a 2-dimensional realization

components, by a smooth curves  $\gamma: [a, b] \rightarrow \mathbb{R}^d$  (figure 1.1). This is particular appropriate, if the assumption of an elliptically distribution can only hardly be justified, like for **observed dynamical systems**. Thereby the evolution function generates a smooth submanifold  $\mathcal{M} \subseteq \mathbb{R}^d$  within the observation space  $\mathbb{R}^d$ , and an “error free” observation can be identified by a random vector  $\mathbf{X}^*$ , with outcomes on  $\mathcal{M}$ . Additionally, however, the observation function may be regarded to be subjected to a measurement error  $\varepsilon$ . By the assumption, that  $\varepsilon$  has an elliptical distribution, then the distribution of the observable random vector  $\mathbf{X}$  is represented by a elliptical  $\mathcal{M}$ -distribution.

**Definition** ( $\mathcal{M}$ -Distribution). For  $d \in \mathbb{N}$  let  $\mathbf{X}^*: \Omega \rightarrow \mathbb{R}^d$  be a random vector and  $\mathcal{M} \subseteq \mathbb{R}^d$  a smooth  $k$ -submanifold of  $\mathbb{R}^d$  with  $k \leq d$ . Then  $\mathbf{X}^*$  is  $\mathcal{M}$ -distributed, iff for the probability density  $P$ , which is induced by  $\mathbf{X}^*$ , it holds, that:

$$P(\mathbf{X}^* = \mathbf{x}) > 0 \Leftrightarrow \mathbf{x} \in \mathcal{M} \tag{1.36}$$

Thereupon a random vector  $\mathbf{X}: \Omega \rightarrow \mathbb{R}^d$  is elliptically  $\mathcal{M}$ -distributed, iff there exists an  $\mathcal{M}$ -distributed random vector  $\mathbf{X}^*: \Omega \rightarrow \mathbb{R}^d$  and an elliptically distributed random



error  $\boldsymbol{\varepsilon}: \Omega \rightarrow \mathbb{R}^d$ , such that:

$$\mathbf{X} \sim \mathbf{X}^* + \boldsymbol{\varepsilon} \quad (1.37)$$

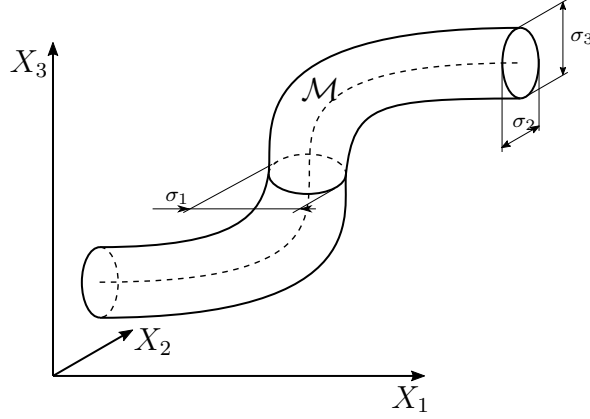


Figure 1.2: Elliptical  $\mathcal{M}$ -distribution in 3 dimensions

The assumption, that the observed random vector  $\mathbf{X}$ , is elliptically  $\mathcal{M}$ -distributed, is very general, but allows an estimation of  $\mathcal{M}$  by minimizing the averaged Euclidean distance to  $\mathbf{X}$ . Thereby the tangent spaces  $T_x\mathcal{M}$  have a basis, given by  $k$  principal components of local infinitesimal covariances, such that the remaining  $d - k$  principal components describe the normal space  $N_x\mathcal{M}$ , which is orthogonal to the tangent space  $T_x\mathcal{M}$ . Since  $T_x\mathcal{M}$  and  $N_x\mathcal{M}$  are equipped with an induced Riemannian metric, which is simply given by the standard scalar product, there exists a minimal orthogonal projection  $\pi_{\mathcal{M}}: \mathbb{R}^d \hookrightarrow \mathcal{M}$ , that maps any realization  $\mathbf{x}$  of  $\mathbf{X}$  to a closest point on  $\mathcal{M}$ . Then proposition 1.2 motivates properties for  $\mathcal{M}$  to minimize the averaged Euclidean distance to realizations of  $\mathbf{X}$ . This provides the definition of smooth  $k$ -principal manifolds (Hastie et al. 1989, p513).

**Definition** (Principal Manifold). For  $d \in \mathbb{N}$  let  $\mathbf{X}: \Omega \rightarrow \mathbb{R}^d$  be a random vector,  $\mathcal{M} \subseteq \mathbb{R}^d$  a (smooth)  $k$ -submanifold of  $\mathbb{R}^d$  with  $k \leq d$  and  $\pi_{\mathcal{M}}: \mathbb{R}^d \hookrightarrow \mathcal{M}$  a minimal orthogonal projection onto  $\mathcal{M}$ . Then  $\mathcal{M}$  is a (smooth)  $k$ -principal manifold for  $\mathbf{X}$ , iff  $\forall \mathbf{x} \in \mathcal{M}$  it holds, that:

$$\mathbb{E}(\mathbf{X} \in \pi_{\mathcal{M}}^{-1}(\mathbf{x})) = \mathbf{x} \quad (1.38)$$

Furthermore  $\mathcal{M}$  is termed maximal, iff  $\mathcal{M}$  maximizes the explained variance  $\text{Var}(\pi_{\mathcal{M}}(\mathbf{X}))$ .

By extending the local properties of the tangent spaces to the underlying manifold, by propositions 1.2 and 1.3 it can be concluded, that maximal principal manifolds minimize the Euclidean distance to  $\mathbf{X}$ . Intuitively this can be understood as follows: The principal manifold property assures, that:

$$\text{Var}(\mathbf{X}) = \text{Var}(\pi_{\mathcal{M}}(\mathbf{X})) + \text{Var}(\mathbf{X} - \pi_{\mathcal{M}}(\mathbf{X}))$$

Consequently the choice of  $\mathcal{M}$  maximizes  $\text{Var}(\pi_{\mathcal{M}}(\mathbf{X}))$  if and only if it minimizes  $\text{Var}(\mathbf{X} - \pi_{\mathcal{M}}(\mathbf{X}))$ , which equals the variance of the error and therefore the Euclidean distance. At closer inspection, however, it turns out, that in difference to linear principal manifolds, the maximization problem is ill-defined for arbitrary smooth principal manifolds, since for any finite number of realizations trivial solutions can be found by smooth principal manifolds, that interpolate the realizations and therefore provide a perfect explanation. In order to close this gap, further structural assumptions have to be incorporated, either by a parametric family  $\{\mathbf{f}_{\theta}\}_{\theta \in \Theta}$  that restricts the possible solutions - or by a regularization, as given in the elastic map algorithm that penalizes long distances and strong curvature (Gorban et al. 2008). Due to the complexity of this topic, however, it is left to the second chapter, where Energy based models are used to overcome this deficiency. In the following the generalization of the correlation to smooth principal manifolds for convenience is defined with respect to a principal manifold  $\mathcal{M}$ , which is maximal “with respect to appropriate restrictions”.

**Definition** ( $\mathcal{M}$ -Correlation). For  $d \in \mathbb{N}$  let  $\mathbf{X}: \Omega \rightarrow \mathbb{R}^d$  be a random vector,  $\mathcal{M}$  a smooth principal manifold for  $\mathbf{X}$ , which is maximal “with respect to appropriate restrictions” and  $\pi_{\mathcal{M}}: \mathbb{R}^d \rightarrow \mathcal{M}$  a minimal orthogonal projection. Then for any  $i, j \in \{1, \dots, d\}$  let the  $\mathcal{M}$ -Correlation between  $X_i$  and  $X_j$  be defined by:

$$\rho_{X_i, X_j | \mathcal{M}}^2 := R_i R_j \int_{\mathcal{M}} S_{i,j}(\mathbf{x}) S_{j,i}(\mathbf{x}) dP_{\mathcal{M}} \quad (1.39)$$

where for  $i \in \{1, \dots, d\}$  the **reliability** of  $X_i$  with respect to  $\mathcal{M}$  is given by:

$$R_i := 1 - \frac{\text{Var}_i(\mathbf{X} - \pi_{\mathcal{M}}(\mathbf{X}))}{\text{Var}_i(\mathbf{X})} \quad (1.40)$$

and for  $i, j \in \{1, \dots, d\}$  the **local sensitivity** of  $X_i$  with respect to  $X_j$  by:

$$S_{i,j}(\mathbf{x}) := \left. \frac{\partial}{\partial x_j} (\mathbf{x} - \pi_{\mathcal{M}}(\mathbf{x})) \right|_i \quad (1.41)$$

**Proposition 1.5.** *For an elliptical  $\mathcal{M}$ -distributed random vector  $\mathbf{X}$  the  $\mathcal{M}$ -Correlation generalizes the  $L$ -Correlation to smooth principal manifolds.*

*Proof.* Let  $L$  be a maximal linear principal manifold for  $\mathbf{X}: \Omega \rightarrow \mathbb{R}^d$ , and  $\mathbf{x}$  a realization of  $\mathbf{X}$  then there exists an  $\beta \in \mathbb{R}$  with:

$$S_{i,j}(\mathbf{x}) = \left. \frac{\partial}{\partial x_j} (\mathbf{x} - \pi_L(\mathbf{x})) \right|_i \equiv \beta$$

Furthermore for  $c \neq 0$ :

$$S_{j,i}(\mathbf{x}) = \left. \frac{\partial}{\partial x_i} (\mathbf{x} - \pi_L(\mathbf{x})) \right|_j \equiv \frac{1}{\beta}$$

Such that  $S_{i,j}(\mathbf{x})S_{j,i}(\mathbf{x}) = 1$ . Consequently for  $\mathcal{M} = L$  it follows, that:

$$\begin{aligned} \rho_{X_i, X_j | \mathcal{M}}^2 &= R_i R_j \int_{\mathcal{M}} S_{i,j}(\mathbf{x}) S_{j,i}(\mathbf{x}) dP_{\mathcal{M}} \\ &= R_i R_j \int_{\mathcal{M}} dP_{\mathcal{M}} \\ &= R_i R_j = \rho_{X_i, X_j | L}^2 \end{aligned}$$

□

## Chapter 2

# Estimation of Principal Manifolds by Energy Based Models

### Overview

*The evaluation of the generalized Pearson correlation requires the estimation of a principal manifold. This chapter is intended to introduce a statistical framework for this estimation, which allows the incorporation of meaningful structural assumptions. In this purpose, in the first section, random vectors with respect to their dependency structure are equipped with a topology, which provides the notation of Markov random field. Subsequently by the Hammersley-Clifford Theorem is shown, that this structural assumption allows a clique factorization of the joint probability distribution. Thereupon in the second section the positivity assumption of  $\mathcal{M}$ -distributions over their respective manifold  $\mathcal{M}$ , is incorporated within clique factorized distributions to define a statistical model by an Energy Based Model. Afterwards it is shown, that the resulting model has a canonical representation by a log-linear parametrisation, which allows an explicit representation of the log-likelihood gradient, which finally motivates a local maximum likelihood estimation of a principal manifold by a steepest gradient descent in the negative log-likelihood gradient. The chapter closes by the proof of its convergence under some mild analytical assumptions, regarding the cumulant function of the parametrisation.*

## 2.1 Markov Random Fields

An intuitive and generic statistical framework to model causal relationship structures is given by probabilistic **graphical models**. Thereby the joint probability distribution is modelled by a graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ , where the vertices  $v \in \mathcal{V}$  of the graph represent random variables  $X_v$  and the edges  $(u, v) \in \mathcal{E}$  conditional probabilities  $P(X_v \mid X_u)$  between the respective random variables. The strength of graphical modelling is revealed by the incorporation of latent variables, such that the modelled random vector  $\mathbf{X}$  decomposes into observables  $\mathbf{V}$  and latent random variables  $\mathbf{H}$  by  $\mathbf{X} = (\mathbf{V}, \mathbf{H})$ . Then the joint probability distribution  $P(\mathbf{V}, \mathbf{H})$  encodes structural uncertainties of indirect causalities, within the topology of the underlying graph  $\mathcal{G}$ , which in particular is highly desirable for observations in environments of strong structural uncertainty. Graphical modelling of principal manifolds requires that the errors of the variables are treated symmetrically and therefore, that the graph  $\mathcal{G}$  is undirected. This property induces a topological structure to the random vector  $\mathbf{X}$ , which is substantiated by the terminology of **random fields**.

**Definition** (Random Field). *Let  $(M, \tau)$  be a topological space and  $\mathbf{X}$  a random vector, such that  $M$  is an index set for the random variables in  $\mathbf{X}$ . Then  $\mathbf{X}$  is a random field over  $(M, \tau)$ .*

*Notation.* The terminology of random fields canonically applies to undirected graphs as follows: Let  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  be an undirected graph. Then any edge  $(v, u) \in \mathcal{E}$  can be identified with the unit interval in  $\mathbb{R}$  and consequently any path  $(v_i, u_i)_{i=1}^n \in \mathcal{E}^n$  with the length  $n \in \mathbb{R}$ . Let now for path connected vertices, the function  $d: \mathcal{V} \times \mathcal{V} \rightarrow \mathbb{R}_0^+ \cup \{\infty\}$  denote the length of shortest paths between them and infinity for unconnected vertices. Then  $d$  is a distance measure over  $\mathcal{V}$  and induces a canonical topology  $\mathcal{T}_{\mathcal{G}}$  from  $\mathbb{R}$  to  $\mathcal{V}$ , which makes  $(\mathcal{V}, \mathcal{T}_{\mathcal{G}})$  a topological space. Therefore in the following a “random field  $\mathbf{X}$  over an undirected graph  $\mathcal{G}$ ” denotes a random field over the respective topological space  $(\mathcal{V}, \mathcal{T}_{\mathcal{G}})$ .

Besides the fundamental structural assumption **(S1)**, that the the observables  $\mathbf{V}$  can be represented as the vertices of an undirected graph  $\mathcal{G}$ , the representation of the joint probability distribution  $P(\mathbf{X})$  by an **undirected graphical model** also requires, that the conditional probabilities can be represented by the edges of  $\mathcal{G}$ . This

makes a further structural assumption **(S2)** necessary, which requires that  $\mathbf{X}$  satisfies the local Markov property in  $\mathcal{G}$ . Then  $\mathbf{X}$  is known as a **Markov random field**.

**Definition** (Markov Random Field). *Let  $\mathbf{X}$  be a random field over an undirected graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  and  $P$  the joint probability distribution of  $\mathbf{X}$ . Then  $\mathbf{X}$  is a Markov Random Field (MRF) over  $\mathcal{G}$ , iff  $\forall v \in \mathcal{V}$ , the random variable  $X_v$  conditionally only depends on random variables, indexed within the local neighbourhood  $\mathcal{N}(v) \subseteq \mathcal{V}$ , such that:*

$$P(X_v \mid \{X_i\}_{i \in \mathcal{V} \setminus \{v\}}) = P(X_v \mid \{X_i\}_{i \in \mathcal{N}(v)}) \quad (2.1)$$

Although the structural assumptions **(S1)** and **(S2)** for  $\mathbf{X}$ , at first sight, seem very general and not restrictive, it turns out, that they have severe consequences for the joint probability distribution  $P(\mathbf{X})$ . As a preparatory step, it is necessary to reconstitute the notation cliques. Thereby for an undirected graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ , a clique  $c \subseteq \mathcal{V}$  in the graph denotes a completely connected subset of vertices  $\mathcal{V}$ , such that  $\forall i, j \in c \Rightarrow (i, j) \in \mathcal{E}$  (see figure 2.1). This allows to define a **clique factorization** of  $P$  over  $\mathcal{G}$ .

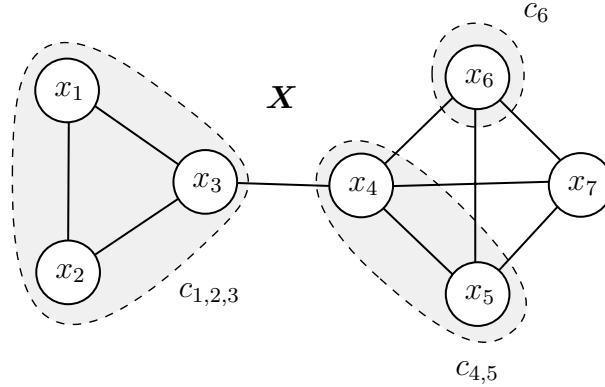


Figure 2.1: Examples for cliques in a Markov Random Field

**Definition** (Clique factorization). *Let  $\mathbf{X}$  be a random field over an undirected graph  $\mathcal{G}$  with cliques  $\mathcal{C}$  and  $P$  the probability distribution, which is induced by  $\mathbf{X}$ . Then  $P$  factorizes over  $\mathcal{G}$ , iff there exists a set of positive functions  $\{\Phi_c\}_{c \in \mathcal{C}}$ , termed **clique potentials**:*

$$\Phi_c: \prod_{i \in c} \text{img}(X_i) \rightarrow \mathbb{R}_0^+, \quad \forall c \in \mathcal{C} \quad (2.2)$$

and a normalization factor  $Z \in \mathbb{R}_0^+$ , termed **partition function**, such that:

$$P(\mathbf{X} = \mathbf{x}) = \frac{1}{Z} \prod_{c \in \mathcal{C}} \Phi_c(\mathbf{x}_c), \text{ with } Z := \int_{\mathcal{X}} \prod_{c \in \mathcal{C}} \Phi_c(\mathbf{x}_c) d\mathbf{x} \quad (2.3)$$

**Lemma 2.1.** *Let  $\mathbf{X}$  be a random field over an undirected graph  $\mathcal{G}$  and  $P$  the joint probability distribution of  $\mathbf{X}$ , such that  $P$  factorizes over  $\mathcal{G}$ . Then  $\mathbf{X}$  is a Markov Random Field over  $\mathcal{G}$ .*

*Proof of Lemma 2.1.* Let  $\mathcal{C}$  be the set of cliques in  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  and for any vertex  $i \in \mathcal{V}$  let  $\mathcal{N}(i) \subset V$  denote the local neighbourhood of  $i$  in  $\mathcal{G}$ . Since  $P$  factorizes over  $\mathcal{G}$ , there exists a set of positive functions  $\{\Phi_c\}_{c \in \mathcal{C}}$  with  $\Phi_c: \prod_{i \in c} \text{img}(X_i) \rightarrow \mathbb{R}_0^+$ ,  $\forall c \in \mathcal{C}$ , such that for  $\mathcal{V}_i^- := \mathcal{V} \setminus \mathcal{N}(i) \cup \{i\}$  it holds, that:

$$\begin{aligned} P(X_i \mid \{X_j\}_{j \in \mathcal{N}(i)}) &= \frac{P(X_i, \{X_j\}_{j \in \mathcal{N}(i)})}{P(\{X_j\}_{j \in \mathcal{N}(i)})} \\ &= \frac{\sum_{\mathcal{V}_i^-} \prod_{c \in \mathcal{C}} \Phi_c(\mathbf{x}_c)}{\int_{X_i} \sum_{\mathcal{V}_i^-} \prod_{c \in \mathcal{C}} \Phi_c(\mathbf{x}_c) d\mathbf{x}} \end{aligned} \quad (2.4)$$

Thereupon for  $i \in \mathcal{V}$  let  $\mathcal{C}_i^+ := \{c \in \mathcal{C} \mid i \in c\}$  denote the set of cliques in  $\mathcal{G}$ , that contain  $i$  and  $\mathcal{C}_i^- := \mathcal{C} \setminus \mathcal{C}_i^+$  the remaining set of cliques in  $\mathcal{G}$ , that do not contain  $i$ . Then for any  $i \in \mathcal{V}$  it holds, that:

$$\prod_{c \in \mathcal{C}} \Phi_c(\mathbf{x}_c) = \prod_{c \in \mathcal{C}_i^+} \Phi_c(\mathbf{x}_c) \prod_{c \in \mathcal{C}_i^-} \Phi_c(\mathbf{x}_c) \quad (2.5)$$

Since by the definition of  $\mathcal{V}_i^-$  and  $\mathcal{C}_i$  for any clique  $c \in \mathcal{C}_i$  it holds, that  $c \subseteq \mathcal{V}_i$  and therefore:

$$\begin{aligned} \sum_{\mathcal{V}_i^-} \prod_{c \in \mathcal{C}} \Phi_c(\mathbf{x}_c) &\stackrel{2.5}{=} \sum_{\mathcal{V}_i^-} \prod_{c \in \mathcal{C}_i^+} \Phi_c(\mathbf{x}_c) \prod_{c \in \mathcal{C}_i^-} \Phi_c(\mathbf{x}_c) \\ &= \prod_{c \in \mathcal{C}_i^+} \Phi_c(\mathbf{x}_c) \sum_{\mathcal{V}_i^-} \prod_{c \in \mathcal{C}_i^-} \Phi_c(\mathbf{x}_c) \end{aligned} \quad (2.6)$$

Also the factor  $\sum_{\mathcal{V}_i^-} \prod_{c \in \mathcal{C}_i^-} \Phi_c(\mathbf{x}_c)$  does not involve  $X_i$  and therefore can be factored

out, such that:

$$\sum_{\mathcal{V}_i^-} \prod_{c \in \mathcal{C}} \Phi_c(\mathbf{x}_c) \stackrel{2.6}{=} \left( \sum_{\mathcal{V} \setminus \mathcal{V}_i} \prod_{c \in \mathcal{C}_i^-} \Phi_c(\mathbf{x}_c) \right) \prod_{c \in \mathcal{C}_i^+} \Phi_c(\mathbf{x}_c) \quad (2.7)$$

This also applies to the denominator on the right side of equation ..., such that:

$$\int_{X_i} \sum_{\mathcal{V}_i^-} \prod_{c \in \mathcal{C}} \Phi_c(\mathbf{x}_c) d\mathbf{x} \stackrel{2.7}{=} \left( \sum_{\mathcal{V}_i^-} \prod_{c \in \mathcal{C}_i^-} \Phi_c(\mathbf{x}_c) \right) \int_{X_i} \prod_{c \in \mathcal{C}_i} \Phi_c(\mathbf{x}_c) d\mathbf{x} \quad (2.8)$$

From equations 2.4, 2.7 and 2.8 it then follows, that:

$$\begin{aligned} P(X_i \mid \{X_j\}_{j \in \mathcal{N}(i)}) &= \frac{\prod_{c \in \mathcal{C}_i^+} \Phi_c(\mathbf{x}_c)}{\int_{X_i} \prod_{c \in \mathcal{C}_i^+} \Phi_c(\mathbf{x}_c) d\mathbf{x}} \\ &= \frac{\prod_{c \in \mathcal{C}_i^-} \Phi_c(\mathbf{x}_c)}{\prod_{c \in \mathcal{C}_i^-} \Phi_c(\mathbf{x}_c)} \frac{\prod_{c \in \mathcal{C}_i^+} \Phi_c(\mathbf{x}_c)}{\int_{X_i} \prod_{c \in \mathcal{C}_i^+} \Phi_c(\mathbf{x}_c) d\mathbf{x}} \\ &= \frac{\prod_{c \in \mathcal{C}} \Phi_c(\mathbf{x}_c)}{\int_{X_i} \prod_{c \in \mathcal{C}} \Phi_c(\mathbf{x}_c) d\mathbf{x}} \\ &= \frac{P(\mathbf{X})}{P(\{X_j\}_{j \in \mathcal{V}_i^-})} = P(X_i \mid \{X_j\}_{j \in \mathcal{V} \setminus \{i\}}) \end{aligned}$$

□

**Theorem 2.1** (Hammersley-Clifford). *Let  $\mathbf{X}$  be a random field over an undirected graph  $\mathcal{G}$  and  $P$  the probability distribution, which is induced by  $\mathbf{X}$ . Then the following statements are equivalent:*

(1)  $\mathbf{X}$  is a Markov Random Field over  $\mathcal{G}$

(2)  $P$  factorizes over  $\mathcal{G}$

*Proof of Theorem 2.1.* Since the backward direction has already been proofed by lemma 2.1, it suffices to proof the forward direction. In the following for the undirected graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  a set of functions  $\{\Phi_s\}_{s \subseteq \mathcal{V}}$  is constructed, which is shown to satisfy the requirements for a clique factorization. Let  $P$  be the joint probability,



which is induced by  $\mathbf{X}$  and for any subset  $z \subseteq \mathcal{V}$ , let  $\pi_z$  be defined by:

$$\pi_z(\mathbf{x}_z) := P(\{X_i = x_i\}_{i \in z} \mid \{X_i\}_{i \in \mathcal{V} \setminus z}) P(\{X_i\}_{i \in \mathcal{V} \setminus z}) \quad (2.9)$$

Thereupon for any subset  $s \subseteq \mathcal{V}$  let the function  $\Phi_s: \prod_{i \in s} \text{img}(X_i) \rightarrow \mathbb{R}$  be defined by:

$$\Phi_s(\mathbf{x}_s) := \prod_{z \subseteq s} \pi_z(\mathbf{x}_z)^{\sigma(s, z)}, \text{ with } \sigma(s, z) := (-1)^{|s| - |z|} \quad (2.10)$$

Then the set of functions  $\{\Phi_s\}_{s \subseteq \mathcal{V}}$  satisfies the following requirements: (i) For any  $s \subseteq \mathcal{V}$  the function  $\Phi_s$  is positive, since any factor is positive (ii) The power  $\sigma(s, z)$  is 1 if  $|s| - |z|$  is even and  $-1$  if  $|s| - |z|$  is odd. Now let  $z \subseteq \mathcal{V}$ . Then  $\pi_z$  occurs as a factor within the function  $\Phi_z$ , with the power  $\sigma(z, z) = 1$ . Furthermore  $\pi_z$  occurs as a factor in functions over any subsets of  $V$  that contain  $z$  and one additional element  $a \in \mathcal{V}$  with the power  $\sigma(z \cup \{a\}, z) = -1$ . Continuing this process it can be seen by the binomial equation:

$$\sum_{k=0}^n \binom{n}{k} (-1)^k = (1 - 1)^n = 0, \quad \forall 0 \leq k \leq n \quad (2.11)$$

that for any  $z \subseteq \mathcal{V}$  the factors  $\pi_z$  are cancelled out within the product  $\phi_V$ , except for the case  $z = \mathcal{V}$ . In this case  $\pi_z = P(\mathbf{X})$ , such that:

$$P(\mathbf{X} = \mathbf{x}) = \prod_{s \subseteq V} \Phi_s(\mathbf{x}_s) \quad (2.12)$$

Let now be  $s \subseteq \mathcal{V}$ , such that  $s$  is not a clique in  $\mathcal{G}$ . Then by definition there exist two vertices  $i, j \in s$  within  $s$ , which are not connected, such that  $(i, j) \notin \mathcal{E}$ . Consequently the subsets  $z \subseteq s$  of  $s$  can be distinguished by their inclusion of  $i$  and  $j$ , such that for any  $w \subseteq s \setminus \{i, j\}$  there respectively exist the four subsets  $w$ ,  $w \cup \{i\}$ ,  $w \cup \{j\}$  and  $w \cup \{i, j\}$  within  $s$ . This allows to rewrite  $\Phi_s$  as follows:

$$\Phi_s \stackrel{2.5}{=} \prod_{z \subseteq s} \pi_z^{\sigma(s, z)} = \prod_{w \subseteq s \setminus \{i, j\}} \left( \frac{\pi_w}{\pi_{w \cup \{i\}}} \frac{\pi_{w \cup \{i, j\}}}{\pi_{w \cup \{j\}}} \right)^{\sigma(s, z)} \quad (2.13)$$

In the following it will be shown, that  $\Phi_s(\mathbf{x}_s) = 1$ . In the purpose of a shorter representation of the proof, the notations of random variables and their conditional

probability are slightly abbreviated. First, for a subset of vertices  $a \subseteq V$  let  $\mathbf{X}_a := \{X_i\}_{i \in a}$ . Then the conditional probability is extended by an underline to mark fixed random variables, such that for any  $a, b \subseteq V$  with  $a \cap b = \emptyset$ :

$$P(\mathbf{X}_a, \underline{\mathbf{X}}_b) := P(\mathbf{X}_a \mid \mathbf{X}_b)P(\mathbf{X}_b) \quad (2.14)$$

Then by using Bayes' Theorem, the factor  $\pi_w$  can be rewritten by the conditional dependency of  $X_i$ :

$$\begin{aligned} \pi_w &\stackrel{\text{def}}{=} P(\mathbf{X}_w, \underline{\mathbf{X}}_{V \setminus w}) = P(\mathbf{X}_w, \underline{\mathbf{X}}_{\{i\}}, \underline{\mathbf{X}}_{V \setminus (w \cup \{i\})}) \\ &\stackrel{\text{def}}{=} P(\mathbf{X}_w \mid \underline{\mathbf{X}}_{\{i\}}, \underline{\mathbf{X}}_{V \setminus (w \cup \{i\})})P(\underline{\mathbf{X}}_{\{i\}}, \underline{\mathbf{X}}_{V \setminus (w \cup \{i\})}) \\ &\stackrel{\text{Bayes}}{=} P(\underline{\mathbf{X}}_{\{i\}} \mid \mathbf{X}_w, \underline{\mathbf{X}}_{V \setminus (w \cup \{i\})})P(\mathbf{X}_w, \underline{\mathbf{X}}_{V \setminus (w \cup \{i\})}) \end{aligned} \quad (2.15)$$

This representation can also be applied to the factor  $\pi_{w \cup \{i\}}$ :

$$\begin{aligned} \pi_{w \cup \{i\}} &\stackrel{\text{def}}{=} P(\mathbf{X}_{w \cup \{i\}}, \underline{\mathbf{X}}_{V \setminus (w \cup \{i\})}) \\ &\stackrel{\text{def}}{=} P(\mathbf{X}_w, \mathbf{X}_{\{i\}} \mid \underline{\mathbf{X}}_{V \setminus (w \cup \{i\})})P(\underline{\mathbf{X}}_{V \setminus (w \cup \{i\})}) \\ &\stackrel{\text{Bayes}}{=} P(\mathbf{X}_{\{i\}} \mid \mathbf{X}_w, \underline{\mathbf{X}}_{V \setminus (w \cup \{i\})})P(\mathbf{X}_w, \underline{\mathbf{X}}_{V \setminus (w \cup \{i\})}) \end{aligned} \quad (2.16)$$

It appears, that the priors, given by the factors on the right sides of equations ... and ... are equal. Thereby the random variable  $X_j$  is fixed. Since  $i$  and  $j$ , however, are not connected, the Markov property, given by equation 2.1, requires that  $X_i$  and  $X_j$  are conditionally independent, if their neighbours are given. Consequently the random variables  $X_j$  in equations ... and ... can be chosen freely and it follows, that:

$$\begin{aligned} \frac{\pi_w}{\pi_{w \cup \{i\}}} &= \frac{P(\underline{\mathbf{X}}_{\{i\}} \mid \mathbf{X}_w, \underline{\mathbf{X}}_{V \setminus (w \cup \{i\})})P(\mathbf{X}_w, \underline{\mathbf{X}}_{V \setminus (w \cup \{i\})})}{P(\mathbf{X}_{\{i\}} \mid \mathbf{X}_w, \underline{\mathbf{X}}_{V \setminus (w \cup \{i\})})P(\mathbf{X}_w, \underline{\mathbf{X}}_{V \setminus (w \cup \{i\})})} \\ &= \frac{P(\underline{\mathbf{X}}_{\{i\}} \mid \mathbf{X}_{w \cup \{j\}}, \underline{\mathbf{X}}_{V \setminus (w \cup \{i, j\})})P(\mathbf{X}_{w \cup \{j\}}, \underline{\mathbf{X}}_{V \setminus (w \cup \{i, j\})})}{P(\mathbf{X}_{\{i\}} \mid \mathbf{X}_{w \cup \{j\}}, \underline{\mathbf{X}}_{V \setminus (w \cup \{i, j\})})P(\mathbf{X}_{w \cup \{j\}}, \underline{\mathbf{X}}_{V \setminus (w \cup \{i, j\})})} \\ &= \frac{P(\mathbf{X}_{w \cup \{j\}}, \underline{\mathbf{X}}_{V \setminus (w \cup \{j\})})}{P(\mathbf{X}_{w \cup \{j\}}, \underline{\mathbf{X}}_{V \setminus (w \cup \{i, j\})})} \stackrel{\text{def}}{=} \frac{\pi_{w \cup \{j\}}}{\pi_{w \cup \{i, j\}}} \end{aligned} \quad (2.17)$$

Thereupon it immediately follows that:

$$\frac{\pi_w}{\pi_{w \cup \{i\}}} \frac{\pi_{w \cup \{i, j\}}}{\pi_{w \cup \{j\}}} = 1 \quad (2.18)$$

And by equation 2.13 that  $\Phi_s(\mathbf{x}_s) = 1$ , if  $s$  is not a clique in  $\mathcal{G}$ . Let therefore  $\mathcal{C}$  be the set of cliques in  $\mathcal{G}$ , then:

$$P(\mathbf{X} = \mathbf{x}) \stackrel{2.12}{=} \prod_{s \subseteq \mathcal{V}} \Phi_s(\mathbf{x}_s) = \prod_{c \in \mathcal{C}} \Phi_c(\mathbf{x}_c)$$

□

## 2.2 Energy Based Models

The distribution assumption of an  $\mathcal{M}$ -distribution, as defined in section 1.3, requires, that the probability distribution is strictly positive over  $\mathcal{M}$ . If furthermore it is assumed, that all observed random variables are subjected to a Gaussian random error with non-vanishing variances, then any outcome  $\mathbf{x}$  in the sample space  $(S, \Sigma)$  of  $\mathbf{X}$  has a strictly positive probability. This motivates a further structural assumption **(S3)**:

$$P(\mathbf{X} = \mathbf{x}) > 0, \forall \mathbf{x} \in S \quad (2.19)$$

By assuming **(S1)** and **(S2)**, and therefore that  $\mathbf{X}$  is a Markov random field, the Hammersley-Clifford Theorem can be applied and it follows, that the probability distribution  $P$ , which is induced by  $\mathbf{X}$  factorizes over a family of positive functions, given by the clique potentials  $\{\Phi_c\}_{c \in \mathcal{C}}$ . It can be concluded that, with respect to **(S1)** and **(S2)** the additional assumption **(S3)** is equivalent to requirement that all clique potentials are strictly positive. The strictly positivity of a clique potential, however, in turn is equivalent to the property, that it can be represented by an exponential function  $\exp: \mathbb{R} \rightarrow (0, \infty)$  and therefore  $P$  by a **Boltzmann distribution**.

**Definition** (Boltzmann distribution). *Let  $(S, \Sigma)$  be a measurable space and  $P$  a probability density over  $(S, \Sigma)$ . Then  $P$  is a Boltzmann distribution, iff there exists a scalar function  $E: S \rightarrow \mathbb{R}$ , termed an **energy function**, and a normalization factor*

$Z \in \mathbb{R}^+$ , termed a **partition function**, such that:

$$P(\mathbf{X} = \mathbf{x}) = \frac{1}{Z} e^{-E(\mathbf{x})}, \text{ with } Z := \int_S e^{-E(\mathbf{x})} d\mathbf{x} \quad (2.20)$$

*Remark.* The sign convention within the definition of Boltzmann distributions serves no functional purpose, except the compatibility with the terminology, known from statistical physics, which also is referenced by the “partition function”  $Z$  and the “energy function”  $E$ .

**Corollary 2.1.** *Let  $\mathbf{X}$  be a random field over an undirected graph  $\mathcal{G}$  and  $P$  the induced probability distribution. Then the following statements are equivalent:*

- (1)  $\mathbf{X}$  is a Markov Random Field over  $\mathcal{G}$ , such that  $P$  is strictly positive
- (2)  $P$  is a Boltzmann distribution, that factorizes over  $\mathcal{G}$

*Proof of Corollary 2.1.* The Corollary immediately follows as a consequence of the Hammersley-Clifford Theorem:

“ $\implies$ ” Let  $\mathbf{X}$  be a Markov Random Field over  $\mathcal{G}$  and  $\mathcal{C}$  the set of cliques in  $\mathcal{G}$ . Then according to the Hammersley-Clifford Theorem it follows, that there exists a set of clique potentials  $\{\Phi_c\}_{c \in \mathcal{C}}$  with  $\Phi_c: \prod_{i \in c} \text{img}(X_i) \rightarrow \mathbb{R}_0^+$ ,  $\forall c \in \mathcal{C}$ , such that  $P$  factorizes over  $\mathcal{G}$ :

$$P(\mathbf{X} = \mathbf{x}) = \frac{1}{Z} \prod_{c \in \mathcal{C}} \Phi_c(\mathbf{x}_c) \quad (2.21)$$

If furthermore any outcome  $\mathbf{x}$  has a positive probability  $P(\mathbf{X} = \mathbf{x})$ , then the clique potentials  $\Phi_c$  have to be strictly positive. This allows to define an energy function as the sum of logarithmic clique potentials:

$$E(\mathbf{x}) := - \sum_{c \in \mathcal{C}} \log \Phi_c(\mathbf{x}_c) \quad (2.22)$$

Since  $P(\mathbf{X} = \mathbf{x}) = \exp(-E(\mathbf{x}))/Z$  it follows, that  $P$  is a Boltzmann distribution.

“ $\impliedby$ ” Let  $P$  factorize over  $\mathcal{G}$ , then according to the Hammersley-Clifford Theorem it follows, that  $\mathbf{X}$  is a Markov Random Field over  $\mathcal{G}$ . If furthermore  $P$  is a Boltzmann distribution, then by definition  $P$  is strictly positive.  $\square$

By the use Corollary 2.1, the positivity assumption **(S3)** can immediately be incorporated within a statistical model: Let  $\mathbf{X}$  be a MRF over an undirected graph  $\mathcal{G}$ , such that the induced probability  $P$  is strictly positive. Then according to corollary 2.1  $P$  is given by a Boltzmann distribution, that factorizes over  $\mathcal{G}$ . Consequently the model space  $\mathcal{M}$  of a statistical model for  $\mathbf{X}$ , can be described by a family of energy functions, which is therefore referred as an **Energy Based Model**.

**Definition** (Energy Based Model). *Let  $\mathbf{X}$  be an  $(S, \Sigma)$ -valued random field over a finite undirected graph  $\mathcal{G}$  and  $(S, \Sigma, \mathcal{P})$  a statistical model for  $\mathbf{X}$ . Then  $(S, \Sigma, \mathcal{P})$  is an Energy Based Model (EBM) for  $\mathbf{X}$  over  $\mathcal{G}$ , iff for all  $P \in \mathcal{P}$  it holds, that (i)  $P$  factorizes over  $\mathcal{G}$  and (ii)  $P$  is a Boltzmann distribution.*

**Example 2.1** (Boltzmann Machine). *In the contemporary literature about probabilistic graphical models the terminology of EBMs and Boltzmann Machines is frequently used equivalent to emphasize the underlying Boltzmann distribution. In the narrow sense of its original context, however, the **Boltzmann Machine** denotes a specific statistical model, that has been introduced to study the properties of stochastic Hopfield-Networks (Fahlman et al. 1983). In this context for a given undirected graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  with  $n$  vertices, any vertex  $i \in \mathcal{V}$  represents a thresholded binary neuron with an activation threshold  $b_i \in \mathbb{R}$  and any edge  $(i, j) \in \mathcal{E}$  a synapse with a synaptic weight  $w_{i,j} \in \mathbb{R}$ . With regard to the parameters  $W \in \mathbb{R}^{n \times n}$  and  $\mathbf{b} \in \mathbb{R}^n$ , the energy function is the given by :*

$$E_{W,\mathbf{b}}(\mathbf{x}) = - \left( \frac{1}{2} \mathbf{x}^T W \mathbf{x} + \mathbf{b}^T \mathbf{x} \right) \quad (2.23)$$

*Thereupon the conditional probabilities can be derived from the energy function. Thereby the contribution of a single neuron  $i \in \mathcal{V}$  to the total energy sums up to:*

$$\Delta E_i = \frac{1}{2} \sum_{j=1} w_{i,j} x_j + b_i \quad (2.24)$$

*Due to the underlying Boltzmann distribution it then follows, that:*

$$P_{W,\mathbf{b}}(X_i = 1 \mid \{X_j = x_j\}_{j \neq i}) = \frac{1}{1 + \exp(\Delta E_i)} =: \text{sigm}(\Delta E_i) \quad (2.25)$$

In comparison to more generic undirected graphical models, the structure of EBMs allows unconstrained statistical inference with the space of energy functions, which in difference to the clique potentials are not required to be positive. In order to take advantage of this circumstance, however, a parametrisation is required. This parametrisation is given by a **log-linear parametrisation**.

**Definition** (Log-linear parametrisation). *Let  $(S, \Sigma, \mathcal{P})$  be an EBM for  $\mathbf{X}$  over  $\mathcal{G}$  and  $P: \Theta \rightarrow \mathcal{P}$  a parametrisation for  $\mathcal{P}$ . Then  $P$  is a **log-linear parametrisation** for  $\mathcal{P}$ , iff there exists a sufficient statistic  $\phi$  for  $\mathbf{X}$ , termed **feature function**, such that  $\forall \theta \in \Theta$ :*

$$P_{\theta}(\mathbf{X} = \mathbf{x}) = \exp(\theta^T \phi(\mathbf{x}) - \Lambda(\theta)), \forall \mathbf{x} \in S \quad (2.26)$$

where the **cumulant function**  $\Lambda: \Theta \rightarrow \mathbb{R}$  is given by:

$$\Lambda(\theta) := \log \int_S \exp(\theta^T \phi(\mathbf{x})) \, d\mathbf{x} \quad (2.27)$$

The existence of the log-linear parametrization for EBMs over MRFs then immediately follows as a consequence of the Hammersley-Clifford theorem.

**Lemma 2.2.** *Let  $(S, \Sigma, \mathcal{P})$  be an EBM for  $\mathbf{X}$  over  $\mathcal{G}$  and let  $\mathcal{P}$  have a finite dimensional parametrisation  $P^*: \Theta \rightarrow \mathcal{M}$ . Then there exists a log-linear parametrisation for  $\mathcal{P}$ .*

*Proof of Lemma 2.2.* The proof is constructive: Let  $\mathcal{C}$  be the set of cliques in  $\mathcal{G}$ . According to equations 2.21 and 2.22 it follows, that any probability distribution  $Q \in \mathcal{P}$  is identified by a unique energy function  $E_Q$ , which can be written as a sum of logarithmic clique potentials, such that:

$$E_Q(\mathbf{x}) = - \sum_{c \in \mathcal{C}} \log \Phi_{Q,c}(\mathbf{x}_c) \quad (2.28)$$

Consequently a parametrisation for  $\mathcal{P}$  can be chosen, that decomposes in individual parametrisations for the cliques of the graph. Let now be  $c \in \mathcal{C}$  a fixed clique and  $n_{P^*} \in \mathbb{N}$  the dimension of the finite dimension parametrisation  $P^*$ . Then, since any  $Q \in \mathcal{P}$  factorizes over  $\mathcal{G}$ , it can be concluded, that any family of clique potentials  $\{\Phi_{Q,c}\}_{Q \in \mathcal{P}}$  has a minimal parametrisation of finite dimension  $n_c < n_{P^*}$ . Further-

more since any potential function is strictly positive, there exists a positive function  $\phi_c: \prod_{i \in c} \text{img}(X_i) \rightarrow \mathbb{R}^{n_c}$ , and a parameter space  $\Theta_c \subseteq \mathbb{R}^{n_c}$ , such that for all  $Q \in \mathcal{P}$ , there exists a vector  $\theta_c \in \Theta_c$  with:

$$\Phi_{Q,c}(\mathbf{x}_c) = \exp(\theta_c^T \phi_c(\mathbf{x}_c)) \quad (2.29)$$

Consequently there exist parameter spaces  $\Theta_c$  and functions  $\phi_c$ , such that for all  $Q \in \mathcal{P}$  the corresponding energy function  $E_Q$  can be identified by a parameter  $\theta \in \prod_{c \in \mathcal{C}} \Theta_c$  with:

$$E_Q(\mathbf{x}) = - \sum_{c \in \mathcal{C}} \log \Phi_{Q,c}(\mathbf{x}_c) = - \sum_{c \in \mathcal{C}} \theta_c^T \phi_c(\mathbf{x}_c) \quad (2.30)$$

Since the energy functions  $E_Q$  in turn identify all  $Q \in \mathcal{P}$ , it follows, that for any  $Q \in \mathcal{P}$ , there exist parameters  $\theta_c \in \Theta_c$ , termed **clique parameters**, such that  $\forall \mathbf{x} \in S$ :

$$Q(\mathbf{X} = \mathbf{x}) = \frac{1}{Z(\theta)} \exp \left( \sum_{c \in \mathcal{C}} \theta_c^T \phi_c(\mathbf{x}_c) \right) \quad (2.31)$$

Let therefore  $\Lambda(\theta) := \log Z(\theta)$ ,  $\theta := \prod_{c \in \mathcal{C}} \theta_c$  and  $\phi := \prod_{c \in \mathcal{C}} \phi_c$  and the parametrisation  $P: \Theta \rightarrow \mathcal{P}$  be given by:

$$P_\theta(\mathbf{X} = \mathbf{x}) := \exp(\theta^T \phi(\mathbf{x}) - \Lambda(\theta)), \forall \mathbf{x} \in S$$

Then (i)  $P$  is a parametrisation for  $\mathcal{M}$ , since  $\text{img}(P) = \mathcal{P}$  and (ii)  $P$  is a log-linear parametrisation.  $\square$

In the following the log-linear parametrisation will be used for observation based statistical inference in EBMs. Let  $\mathbf{X}$  be a MRF over a graph  $\mathcal{G}$  with  $d$  vertices and let  $\mathbf{x}$  be an i.i.d realization of  $\mathbf{X}$  of length  $n \in \mathbb{N}$ , such that  $\mathbf{x} \in \mathbb{R}^{n \times d}$ . Furthermore let  $(S, \Sigma, \mathcal{P})$  be an EBM for  $\mathbf{X}$  with a log-linear parametrisation  $P: \Theta \rightarrow \mathcal{P}$ . Then the question arises, which parameters  $\hat{\theta} \in \Theta$  “best possibly” explain the observation of  $\mathbf{x}$ , which provides an estimation for the “true parameter”. Thereby in the absence of further knowledge, given by a prior distribution, the estimation  $\hat{\theta}$  is required to

globally maximize the likelihood function:

$$L_{\mathbf{x}}(\boldsymbol{\theta}) := \prod_{i=1}^n P_{\boldsymbol{\theta}}(\mathbf{X} = \mathbf{x}^i)$$

Then the set of all maximum likelihood estimations constitutes a subset  $\Theta_{\text{ML}} \subseteq \Theta$ . Nevertheless, such a global search can be quite challenging if  $L_{\mathbf{x}}$  is non-concave. Then the determination of a global maximum would require to evaluate  $L_{\mathbf{x}}$  for any single parameter in  $\Theta$ . For the case, however, that  $L_{\mathbf{x}} \in C^2(\Theta)$ , then at least the local maximum likelihood estimations  $\hat{\boldsymbol{\theta}} \in \Theta_{\text{max}}$  are locally determined by the well-known criteria:

(ML1)  $L_{\mathbf{x}}$  has a critical point at  $\hat{\boldsymbol{\theta}}$ , such that  $\nabla L_{\mathbf{x}}(\hat{\boldsymbol{\theta}}) = 0$

(ML2)  $L_{\mathbf{x}}$  has a negative curvature at  $\hat{\boldsymbol{\theta}}$ , such that  $\nabla^2 L_{\mathbf{x}}(\hat{\boldsymbol{\theta}}) < 0$

These criteria, however, are more conveniently equated by the log-likelihood function  $\log L_{\mathbf{x}}$ , which has identical critical points and sign of curvature, but in difference to  $L_{\mathbf{x}}$ , decomposes into a sum of independent point realizations:

$$\begin{aligned} \log L_{\mathbf{x}}(\boldsymbol{\theta}) &\stackrel{\text{i.i.d.}}{=} \log \prod_{i=1}^d P_{\boldsymbol{\theta}}(\mathbf{X} = \mathbf{x}^i) \\ &= \sum_{i=1}^n \log L_{\mathbf{x}^i}(\boldsymbol{\theta}), \forall \boldsymbol{\theta} \in \Theta \end{aligned} \tag{2.32}$$

Let now be  $\mathcal{C}$  the set of cliques in  $\mathcal{G}$  and the log-linear parametrisation  $P$  be given by clique parameters  $\{\boldsymbol{\theta}_c\}_{c \in \mathcal{C}}$ , such that  $\Theta = \prod_{c \in \mathcal{C}} \Theta_c$ . Then the gradient  $\nabla \log L_{\mathbf{x}}$  decomposes into a sum of partial derivatives  $\partial_c \log L_{\mathbf{x}}$  with  $\partial_c := \partial / \partial \boldsymbol{\theta}_c$ , such that:

$$\nabla \log L_{\mathbf{x}}(\boldsymbol{\theta}) = \sum_{c \in \mathcal{C}} \partial_c \log L_{\mathbf{x}}(\boldsymbol{\theta}) \tag{2.33}$$

The log-linear parametrisation can then be used to derive an energy representation of the partial derivatives  $\partial_c \log L_{\mathbf{x}}$ .

**Proposition 2.1.** *Let  $(S, \Sigma, \mathcal{P})$  be an EBM for  $\mathbf{X}$  over  $\mathcal{G}$  with cliques  $\mathcal{C}$  and  $P: \Theta \rightarrow \mathcal{M}$  a differentiable log-linear parametrization with clique parameters  $\{\boldsymbol{\theta}_c\}_{c \in \mathcal{C}}$  and energy functions  $\{E_{\boldsymbol{\theta}}\}_{\boldsymbol{\theta} \in \Theta}$ . Then for a realization  $\mathbf{x}$  of  $\mathbf{X}$  of length  $n \in \mathbb{N}$  and a clique*



$c \in \mathcal{C}$  the partial derivative  $\partial_c \log L_{\mathbf{x}}$  has an **energy representation**, which for  $\boldsymbol{\theta} \in \Theta$  is given by:

$$\partial_c \log L_{\mathbf{x}}(\boldsymbol{\theta}) = n \langle -\partial_c \log E_{\boldsymbol{\theta}} \rangle_{\text{d}} - n \langle -\partial_c \log E_{\boldsymbol{\theta}} \rangle_{\text{m}} \quad (2.34)$$

where  $\langle \cdot \rangle_{\text{d}}$  denotes an expectation over  $S$ , with respect to the empirical **data distribution**  $P_{\mathbf{x}}$ , given by:

$$P_{\mathbf{x}}(\mathbf{X}_c = \mathbf{y}_c) := \frac{1}{n} \sum_{i=1}^n \mathbb{I}(\mathbf{y}_c - \mathbf{x}_c^i), \text{ with } \mathbb{I}(\mathbf{x}) = \begin{cases} 1 & \mathbf{x} = \mathbf{0} \\ 0 & \text{else} \end{cases} \quad (2.35)$$

and  $\langle \cdot \rangle_{\text{m}}$  an expectation over  $S$ , with respect to the **model distribution**  $P_{\boldsymbol{\theta}}$

*Proof of Proposition 2.1.* The Proposition is proofed directly by calculation.

**Calculation of  $\log L_{\mathbf{x}}$**  For  $\boldsymbol{\theta} \in \Theta$  the log likelihood of  $\boldsymbol{\theta}$  with respect to  $\mathbf{x}$  is calculated by:

$$\log L_{\mathbf{x}}(\boldsymbol{\theta}) \stackrel{2.32}{=} \sum_{i=1}^n \log L_{\mathbf{x}^i}(\boldsymbol{\theta}) \quad (2.36)$$

Thereupon for any single realization  $\mathbf{x}^i$  the log-likelihood equates to:

$$\log L_{\mathbf{x}^i}(\boldsymbol{\theta}) \stackrel{\text{def}}{=} \log P_{\boldsymbol{\theta}}(\mathbf{X} = \mathbf{x}^i) \quad (2.37)$$

$$\begin{aligned} &\stackrel{\text{def}}{=} \log \exp \left( \sum_{c \in \mathcal{C}} \boldsymbol{\theta}_c^T \boldsymbol{\phi}_c(\mathbf{x}_c^i) - \Lambda(\boldsymbol{\theta}) \right) \\ &= \sum_{c \in \mathcal{C}} \boldsymbol{\theta}_c^T \boldsymbol{\phi}_c(\mathbf{x}_c^i) - \Lambda(\boldsymbol{\theta}) \end{aligned} \quad (2.38)$$

A substitution of equation 2.37 in 2.36 then gives:

$$\log L_{\mathbf{x}}(\boldsymbol{\theta}) = \sum_{i=1}^n \sum_{c \in \mathcal{C}} \boldsymbol{\theta}_c^T \boldsymbol{\phi}_c(\mathbf{x}_c^i) - \sum_{i=1}^n \Lambda(\boldsymbol{\theta}) \quad (2.39)$$

**Calculation of  $\partial_c \Lambda(\boldsymbol{\theta})$**  Let now be  $c \in \mathcal{C}$  a fixed clique in  $\mathcal{G}$ . Then  $\partial_c \Lambda(\boldsymbol{\theta})$  is calculated by the chain rule:

$$\partial_c \Lambda(\boldsymbol{\theta}) = \partial_c \log Z(\boldsymbol{\theta}) = \frac{1}{Z(\boldsymbol{\theta})} \partial_c Z(\boldsymbol{\theta}), \forall \boldsymbol{\theta} \in \Theta \quad (2.40)$$

With regard to the corresponding clique potential  $\Phi_c$ , given by:

$$\Phi_c(\mathbf{x}_c, \boldsymbol{\theta}_c) := \exp(\boldsymbol{\theta}_c^T \boldsymbol{\phi}_c(\mathbf{x}_c)), \forall \mathbf{x} \in S \quad (2.41)$$

it follows, that:

$$\begin{aligned} \prod_{z \in \mathcal{C} \setminus \{c\}} \Phi_z(\mathbf{x}_z, \boldsymbol{\theta}_z) &= \frac{Z(\boldsymbol{\theta})}{\phi_z(\mathbf{x}_z, \boldsymbol{\theta}_z)} \frac{1}{Z(\boldsymbol{\theta})} \prod_{z \in \mathcal{C}} \Phi_z(\mathbf{x}_z, \boldsymbol{\theta}_z) \\ &= \frac{Z(\boldsymbol{\theta})}{\Phi_z(\mathbf{x}_z, \boldsymbol{\theta}_z)} P_{\boldsymbol{\theta}}(\mathbf{X} = \mathbf{x}_s) \end{aligned} \quad (2.42)$$

This allows to calculate  $\partial_c Z(\boldsymbol{\theta})$  by:

$$\begin{aligned} \partial_c Z(\boldsymbol{\theta}) &= \int_S \partial_c \prod_{z \in \mathcal{C}} \Phi_z(\mathbf{x}_z, \boldsymbol{\theta}_z) d\mathbf{x} \\ &= \int_S (\partial_c \Phi_c(\mathbf{x}_c, \boldsymbol{\theta}_c)) \prod_{z \in \mathcal{C} \setminus \{c\}} \Phi_z(\mathbf{x}_z, \boldsymbol{\theta}_z) d\mathbf{x} \\ &= \int_S (\partial_c \Phi_c(\mathbf{x}_c, \boldsymbol{\theta}_c)) \frac{Z(\boldsymbol{\theta})}{\Phi_z(\mathbf{x}_z, \boldsymbol{\theta}_z)} P_{\boldsymbol{\theta}}(\mathbf{X} = \mathbf{x}_s) d\mathbf{x} \end{aligned} \quad (2.43)$$

This term can be simplified by a reverse application of the chain rule, such that:

$$\begin{aligned} \partial_c Z(\boldsymbol{\theta}) &= Z(\boldsymbol{\theta}) \int_S \partial_c \log \Phi_c(\mathbf{x}_c, \boldsymbol{\theta}_c) P_{\boldsymbol{\theta}}(\mathbf{X} = \mathbf{x}_s) d\mathbf{x} \\ &\stackrel{\text{def}}{=} Z(\boldsymbol{\theta}) \langle \partial_c \log \Phi_c(\boldsymbol{\theta}_c) \rangle_{\mathbf{m}} \\ &= Z(\boldsymbol{\theta}) \langle \boldsymbol{\phi}_c \rangle_{\mathbf{m}} \end{aligned} \quad (2.44)$$

A final substitution of equation 2.44 in equation 2.40 then gives:

$$\partial_c \Lambda(\boldsymbol{\theta}) = \frac{1}{Z(\boldsymbol{\theta})} Z(\boldsymbol{\theta}) \langle \boldsymbol{\phi}_c \rangle_{\mathbf{m}} = \langle \boldsymbol{\phi}_c \rangle_{\mathbf{m}} \quad (2.45)$$

**Calculation of  $\partial_c \log L_{\mathbf{x}}$**  From equation 2.39 it follows, that:

$$\begin{aligned} \partial_c \log L_{\mathbf{x}}(\boldsymbol{\theta}) &\stackrel{2.39}{=} \sum_{i=1}^n \phi_c(\mathbf{x}_c^i) - \sum_{i=1}^n \partial_c \Lambda(\boldsymbol{\theta}) \\ &\stackrel{2.45}{=} \sum_{i=1}^n \phi_c(\mathbf{x}_c^i) - n \langle \phi_c \rangle_{\mathbf{m}} \end{aligned} \quad (2.46)$$

The sum on the right side of equation 2.46 can be rewritten by the empirical distribution:

$$\begin{aligned} \sum_{i=1}^n \phi_c(\mathbf{x}_c^i) &= \int_S \phi_c(\mathbf{x}_c) \mathbb{I}(\mathbf{x}_c - \mathbf{x}_c^i) d\mathbf{x} \\ &\stackrel{\text{def}}{=} \int_S \phi_c(\mathbf{x}_c) n P_{\mathbf{x}}(\mathbf{X}_c = \mathbf{x}_c) d\mathbf{x} \\ &\stackrel{\text{def}}{=} n \langle \phi_c \rangle_{\mathbf{d}} \end{aligned} \quad (2.47)$$

Then a substitution of equation 2.47 in equation 2.46 yields:

$$\partial_c \log L_{\mathbf{x}}(\boldsymbol{\theta}) = n \langle \phi_c \rangle_{\mathbf{d}} - n \langle \phi_c \rangle_{\mathbf{m}} \quad (2.48)$$

Finally due to the identity:

$$\partial_c \log E_{\boldsymbol{\theta}}(\mathbf{x}) = -\phi_c(\mathbf{x}_c)$$

it follows, that:

$$\partial_c \log L_{\mathbf{x}}(\boldsymbol{\theta}) = n \langle -\partial_c \log E_{\boldsymbol{\theta}} \rangle_{\mathbf{d}} - n \langle -\partial_c \log E_{\boldsymbol{\theta}} \rangle_{\mathbf{m}} \quad (2.49)$$

□

**Example** (*Supplement to Example 2.1*). *With respect to its log-linear parametrisation it follows, that two kinds of cliques are to be distinguished in Boltzmann Machines: For any clique, that is given by a single vertex  $c = \{i\} \subseteq \mathcal{V}$ , the clique parameter is given by the activation threshold  $b_i$  and the feature function by  $\phi_i(x_i) = x_i$ , such that:*

$$-\frac{\partial}{\partial b_i} \log E_{W, \mathbf{b}}(\mathbf{x}) \stackrel{2.23}{=} x_i$$

Therefore according to Proposition 2.1 it follows, that:

$$\frac{\partial}{\partial b_i} \log L_{\mathbf{x}}(W, \mathbf{b}) = n \langle x_i \rangle_{\text{d}} - n \langle x_i \rangle_{\text{m}} \quad (2.50)$$

Furthermore for a clique, that is given by a pair  $c = \{i, j\} \subseteq \mathcal{V}$ , with  $(i, j) \in \mathcal{E}$  the clique parameter is given by the synaptic weight  $w_{ij}$  and the feature function by  $\phi_{ij}(x_i, x_j) = x_i x_j$ , such that:

$$-\frac{\partial}{\partial w_{ij}} \log E_{W, \mathbf{b}}(\mathbf{x}) \stackrel{2.23}{=} x_i x_j$$

And therefore:

$$\frac{\partial}{\partial w_{ij}} \log L_{\mathbf{x}}(W, \mathbf{b}) = n \langle x_i x_j \rangle_{\text{d}} - n \langle x_i x_j \rangle_{\text{m}} \quad (2.51)$$

Due to the definition of the empirical data distribution  $P_{\mathbf{x}}$  with respect to a realization  $\mathbf{x}$  of  $\mathbf{X}$ , the calculation of  $\partial_c \log L_{\mathbf{x}}$ , by its energy representation of Proposition 2.1 requires, that all random variables, which are represented by the clique  $c \in \mathcal{C}$  are realized. The mandatory next step therefore regards the generalization of this representation to latent random variables. In this purpose the underlying MRF  $\mathbf{X}$  is decomposed into an observable component  $\mathbf{V}$  and a latent component  $\mathbf{H}$ , such that  $\mathbf{X} = (\mathbf{V}, \mathbf{H})$ . Thereupon a straight forward approach to capture the latent component within the representation is to generalize the energy functions by **free energy functions**.

**Corollary 2.2.** *Let  $(S, \Sigma, \mathcal{P})$  be an EBM for  $\mathbf{X}$  over  $\mathcal{G}$  with cliques  $\mathcal{C}$  and a differentiable log-linear parametrization  $P: \Theta \rightarrow \mathcal{P}$  with clique parameters  $\{\theta_c\}_{c \in \mathcal{C}}$  and energy functions  $\{E_{\theta}\}_{\theta \in \Theta}$ . Furthermore let the  $\mathbf{X} = (\mathbf{V}, \mathbf{H})$  with  $\text{img} \mathbf{V} = (S_V, \Sigma_V)$  and  $\text{img} \mathbf{H} = (S_H, \Sigma_H)$ . Then for any  $\theta \in \Theta$  let the **free energy function**  $F_{\theta}: S_V \rightarrow \mathbb{R}$  be defined by:*

$$F_{\theta}(\mathbf{v}) := -\log \int_{S_H} e^{-E_{\theta}(\mathbf{v}, \mathbf{h})} d\mathbf{h}, \quad \forall \mathbf{v} \in S_V \quad (2.52)$$

Then for an i.i.d realization  $\mathbf{v}$  of  $\mathbf{V}$  of length  $n$  and for a clique  $c \in \mathcal{C}$  the partial derivative  $\partial_c \log L_{\mathbf{v}}$  has a **free energy representation**, which at  $\theta \in \Theta$  is given by:

$$\partial_c \log L_{\mathbf{v}}(\theta) = n \langle -\partial_c \log F_{\theta} \rangle_{\text{d}} - n \langle -\partial_c \log F_{\theta} \rangle_{\text{m}} \quad (2.53)$$

*Proof of Corollary 2.2.* The likelihood  $L_{\mathbf{v}}(\boldsymbol{\theta})$  is identified by the marginal probability  $P_{\boldsymbol{\theta}}(\mathbf{V} = \mathbf{v})$ , which by the law of total probability is given by:

$$\begin{aligned} P_{\boldsymbol{\theta}}(\mathbf{V} = \mathbf{v}) &= \int_{S_H} P_{\boldsymbol{\theta}}(\mathbf{V} = \mathbf{v}, \mathbf{H} = \mathbf{h}) d\mathbf{h} \\ &\stackrel{\text{def}}{=} \frac{1}{Z(\boldsymbol{\theta})} \int_{S_H} e^{-E_{\boldsymbol{\theta}}(\mathbf{v}, \mathbf{h})} d\mathbf{h}, \text{ with } Z(\boldsymbol{\theta}) := \int_S e^{-E_{\boldsymbol{\theta}}(\mathbf{x})} d\mathbf{x} \\ &\stackrel{2.52}{=} \frac{1}{Z'(\boldsymbol{\theta})} e^{-F_{\boldsymbol{\theta}}(\mathbf{v})}, \text{ with } Z'(\boldsymbol{\theta}) := \int_V e^{-F_{\boldsymbol{\theta}}(\mathbf{v})} d\mathbf{v} \end{aligned} \quad (2.54)$$

This shows, that the marginal distributions are Boltzmann distributions, where the energy functions are given by the free energy of the joint probability distribution  $P_{\boldsymbol{\theta}}(\mathbf{V}, \mathbf{H})$ . Since furthermore the parametrisation  $P_{\boldsymbol{\theta}}(\mathbf{V}, \mathbf{H})$  is assumed to be differentiable over  $\Theta$  it follows, that also  $P_{\boldsymbol{\theta}}(\mathbf{V})$  is differentiable over  $\Theta$  and Proposition 2.1 can be applied, such that:

$$\partial_c \log L_{\mathbf{v}}(\boldsymbol{\theta}) = n \langle -\partial_c \log F_{\boldsymbol{\theta}} \rangle_{\text{d}} - n \langle -\partial_c \log F_{\boldsymbol{\theta}} \rangle_{\text{m}} \quad (2.55)$$

□

The free energy representation for EBMs allows to derive the partial derivatives  $\partial_c \log L_{\mathbf{v}}$  with respect to all clique parameters  $\boldsymbol{\theta}_c$ . Let therefore  $P$  be a differentiable log-linear parametrisation with  $\Theta = \prod_{c \in \mathcal{C}} \Theta_c$ . Then according to equations 2.33 and 2.53 the gradient  $\nabla \log L_{\mathbf{v}}$  has a free energy representation, which is given by:

$$\nabla \log L_{\mathbf{v}}(\boldsymbol{\theta}) = n \sum_{c \in \mathcal{C}} (\langle -\partial_c \log F_{\boldsymbol{\theta}} \rangle_{\text{d}} - \langle -\partial_c \log F_{\boldsymbol{\theta}} \rangle_{\text{m}}) \quad (2.56)$$

This representation can be used to evaluate the criteria **(ML1)** and **(ML2)** for local maximum likelihood. Furthermore, however, the gradient  $\nabla \log L_{\mathbf{v}}$  identifies the direction that maximizes the increase of  $\log L_{\mathbf{v}}$ , which motivates a traversal of the model space  $\mathcal{P}$  along the gradient. For a discretization of the curve in a finite number of steps, which are proportional to the norm of the gradient, this strategy provides the **Steepest Gradient Descent** (see algorithm 2.1).

**Theorem 2.2** (Convergence of GD in EBMs). *Let  $(S, \Sigma, \mathcal{P})$  be an EBM for  $\mathbf{X}$  over*

---

**Algorithm 2.1** Steepest Gradient Descent in EBMs

---

**parameters**  
*prior parameter*  $\boldsymbol{\theta} \in \Theta$   
*step size*  $\eta > 0$   
*number of steps*  $N \in \mathbb{N}$

1: **procedure**  $\text{GD}(\boldsymbol{\theta}, \eta, N)$   
2:    $\boldsymbol{\theta}^{(1)} \leftarrow \boldsymbol{\theta}$   
3:   **for**  $i \in \{1, \dots, N-1\}$  **do**  
4:     **for**  $c \in \mathcal{C}$  **do**  
5:        $\boldsymbol{\theta}_c^{(i+1)} \leftarrow \boldsymbol{\theta}_c^{(i)} - \eta d(\langle -\partial_c \log F_{\boldsymbol{\theta}} \rangle_{\text{d}} - \langle -\partial_c \log F_{\boldsymbol{\theta}} \rangle_{\text{m}})$   
6:     **end for**  
7:   **end for**  
8:   **return**  $\boldsymbol{\theta}^{(N)}$   
9: **end procedure**

---

$\mathcal{G}$  with a log-linear parametrization  $P: \Theta \rightarrow \mathcal{M}$  with a cumulant function  $\Lambda$ , that satisfies:

(A1)  $\Lambda$  is two times continuously differentiable, i.e.  $\Lambda \in C^2(\Theta)$

(A2)  $\Lambda$  is pseudo concave

(A3)  $\nabla \Lambda$  is Lipschitz continuous

Furthermore for  $\mathbf{X} = (\mathbf{V}, \mathbf{H})$  let  $\mathbf{v}$  be an i.i.d. realization of  $\mathbf{V}$  and let  $\Theta_{\min}$  and  $\Theta_{\max}$  respectively denote the critical points of  $L_{\mathbf{v}}(\boldsymbol{\theta}) := P_{\boldsymbol{\theta}}(\mathbf{v})$ . Then for all prior parameters  $\boldsymbol{\theta} \in \Theta \setminus \Theta_{\min}$ , there exists a step size  $\eta > 0$  and a local ML estimation  $\hat{\boldsymbol{\theta}} \in \Theta_{\max}$ , such that:

$$\text{GD}(\boldsymbol{\theta}, \eta, N) \rightarrow \hat{\boldsymbol{\theta}}, \text{ for } N \rightarrow \infty$$

*Proof of Theorem 2.2.* According to (A1) the gradient  $\nabla \Lambda$  and the Hessian  $\nabla^2 \Lambda$  are well defined over  $\Theta$ , and therefore also  $\nabla \log L_{\mathbf{v}}$  and  $\nabla^2 \log L_{\mathbf{v}}$ . Let  $\boldsymbol{\theta} \in \Theta$  such that  $\boldsymbol{\theta}$  does not locally minimize  $L_{\mathbf{v}}$ . Furthermore for the set of cliques  $\mathcal{C}$  let the clique parameters be given by  $\{\boldsymbol{\theta}_c\}_{c \in \mathcal{C}}$ . Then for the case, that  $\boldsymbol{\theta}$  locally maximizes  $L_{\mathbf{v}}$ , it holds that  $\partial_c \log L_{\mathbf{v}}(\boldsymbol{\theta}) = 0, \forall c \in \mathcal{C}$ . Thereupon the Corollary is trivially satisfied,

since  $\forall i \in \mathbb{N}$ :

$$\begin{aligned}\boldsymbol{\theta}_c^{(i+1)} &\stackrel{\text{def}}{=} \boldsymbol{\theta}_c^{(i)} - \eta n (\langle -\partial_c \log F_{\boldsymbol{\theta}} \rangle_{\text{d}} - \langle -\partial_c \log F_{\boldsymbol{\theta}} \rangle_{\text{m}}) \\ &\stackrel{2.55}{=} \boldsymbol{\theta}_c^{(i)} + \eta \partial_c \log L_{\mathbf{v}} = \boldsymbol{\theta}_c^{(i)}\end{aligned}$$

Such that  $\boldsymbol{\theta}^{(i)} = \boldsymbol{\theta}$ ,  $\forall i \in \mathbb{N}$  and thus converges in  $\boldsymbol{\theta}$ , which locally maximizes  $L_{\mathbf{v}}$ . Let therefore  $\boldsymbol{\theta}$  be assumed, not to locally minimize or maximize the likelihood, such that  $\nabla \log L_{\mathbf{v}} \neq 0$ . Then according to **(A2)** for  $f := -\log L_{\mathbf{v}}$  it follows, that  $f$  is pseudo convex, such that a maximal  $U \subseteq \Theta$  with  $\boldsymbol{\theta} \in U$  can be chosen with  $f|_U$  is convex. Let  $\hat{\boldsymbol{\theta}} \in U$  minimize  $f$  in  $U$ , then  $\hat{\boldsymbol{\theta}}$  maximizes  $L_{\mathbf{v}}$  in  $U$ . In order to proof the Corollary, it therefore suffices to show that  $\lim_{i \rightarrow \infty} f(\boldsymbol{\theta}^{(i)}) = f(\hat{\boldsymbol{\theta}})$ , or equivalently:

$$\forall \varepsilon < 0 \exists \eta > 0, N \in \mathbb{N}: \left\| f(\boldsymbol{\theta}^{(N)}) - f(\hat{\boldsymbol{\theta}}) \right\| < \varepsilon$$

**Local Inequalities** The properties  $f$  inherits from  $\Lambda$  can be used to formulate inequalities in  $U$ . Since  $f|_U$  is convex, it follows that:

$$f(\boldsymbol{\theta}) \geq f(\boldsymbol{\vartheta}) + \nabla f(\boldsymbol{\vartheta})^T (\boldsymbol{\theta} - \boldsymbol{\vartheta}), \forall \boldsymbol{\theta}, \boldsymbol{\vartheta} \in U \quad (2.57)$$

Furthermore since  $L_{\mathbf{d}} \in C^2(\Theta)$  it follows, that also  $f \in C^2(\Theta)$  and according to **(A3)** since  $\nabla \Lambda$  is Lipschitz continuous also  $\nabla f$  is Lipschitz continuous, such that the largest Eigenvalue of the Hessian  $\nabla^2 f$  is upper bounded by a Lipschitz constant  $\delta > 0$ . Consequently:

$$f(\boldsymbol{\theta}) \leq f(\boldsymbol{\vartheta}) + \nabla f(\boldsymbol{\vartheta})^T (\boldsymbol{\theta} - \boldsymbol{\vartheta}) + \frac{\delta}{2} \|\boldsymbol{\theta} - \boldsymbol{\vartheta}\|^2, \forall \boldsymbol{\theta}, \boldsymbol{\vartheta} \in U \quad (2.58)$$

**Upper bound for  $f(\boldsymbol{\theta}^{(i+1)}) - f(\hat{\boldsymbol{\theta}})$**  The local inequalities can be applied to compare  $f(\boldsymbol{\theta}^{(i)})$  with  $f(\boldsymbol{\theta}^{(i+1)})$ . It follows, that:

$$\begin{aligned}
f(\boldsymbol{\theta}^{(i+1)}) &\stackrel{2.58}{\leq} f(\boldsymbol{\theta}^{(i)}) + \nabla f(\boldsymbol{\theta}^{(i)})^T (\boldsymbol{\theta}^{(i+1)} - \boldsymbol{\theta}^{(i)}) + \frac{\delta}{2} \|\boldsymbol{\theta}^{(i+1)} - \boldsymbol{\theta}^{(i)}\|^2 \\
&\stackrel{\text{def}}{=} f(\boldsymbol{\theta}^{(i)}) + \nabla f(\boldsymbol{\theta}^{(i)})^T (-\eta \nabla f(\boldsymbol{\theta}^{(i)})) + \frac{\delta}{2} \|\eta \nabla f(\boldsymbol{\theta}^{(i)})\|^2 \\
&= f(\boldsymbol{\theta}^{(i)}) - \eta \|\nabla f(\boldsymbol{\theta}^{(i)})\|^2 + \frac{\eta^2 \delta}{2} \|\nabla f(\boldsymbol{\theta}^{(i)})\|^2 \\
&= f(\boldsymbol{\theta}^{(i)}) - \eta \left(1 - \frac{\eta \delta}{2}\right) \|\nabla f(\boldsymbol{\theta}^{(i)})\|^2
\end{aligned} \tag{2.59}$$

Let now the step size  $\eta$  satisfy the inequality  $\eta \leq \frac{1}{\delta}$ , then according to equation 2.59 it can be concluded, that:

$$\begin{aligned}
f(\boldsymbol{\theta}^{(i+1)}) &\leq f(\boldsymbol{\theta}^{(i)}) - \frac{\eta}{2} \|\nabla f(\boldsymbol{\theta}^{(i)})\|^2 \\
&\stackrel{2.57}{\leq} f(\hat{\boldsymbol{\theta}}) + \nabla f(\boldsymbol{\theta}^{(i)})^T (\boldsymbol{\theta}^{(i)} - \hat{\boldsymbol{\theta}}) - \frac{\eta}{2} \|\nabla f(\boldsymbol{\theta}^{(i)})\|^2 \\
&= f(\hat{\boldsymbol{\theta}}) + \frac{1}{2\eta} \left( \|\boldsymbol{\theta}^{(i)} - \hat{\boldsymbol{\theta}}\|^2 - \|(\boldsymbol{\theta}^{(i)} - \eta \nabla f(\boldsymbol{\theta}^{(i)})) - \hat{\boldsymbol{\theta}}\|^2 \right) \\
&= f(\hat{\boldsymbol{\theta}}) + \frac{1}{2\eta} \left( \|\boldsymbol{\theta}^{(i)} - \hat{\boldsymbol{\theta}}\|^2 - \|\boldsymbol{\theta}^{(i+1)} - \hat{\boldsymbol{\theta}}\|^2 \right)
\end{aligned} \tag{2.60}$$

And therefore:

$$f(\boldsymbol{\theta}^{(i+1)}) - f(\hat{\boldsymbol{\theta}}) \stackrel{2.60}{\leq} \frac{1}{2\eta} \left( \|\boldsymbol{\theta}^{(i)} - \hat{\boldsymbol{\theta}}\|^2 - \|\boldsymbol{\theta}^{(i+1)} - \hat{\boldsymbol{\theta}}\|^2 \right) \tag{2.61}$$

**Upper bound for  $f(\boldsymbol{\theta}^{(N)}) - f(\hat{\boldsymbol{\theta}})$**  The upper bound for  $f(\boldsymbol{\theta}^{(i+1)}) - f(\hat{\boldsymbol{\theta}})$  can be used to determine an upper bound for the  $N$ -th step error  $f(\boldsymbol{\theta}^{(N)}) - f(\hat{\boldsymbol{\theta}})$ . For  $N \geq 2$



it follows, that:

$$\begin{aligned}
f(\boldsymbol{\theta}^{(N)}) - f(\hat{\boldsymbol{\theta}}) &\leq \frac{1}{N-1} \sum_{i=1}^{N-1} \left( f(\boldsymbol{\theta}^{(i)}) - f(\hat{\boldsymbol{\theta}}) \right) \\
&\leq \frac{1}{N-1} \sum_{i=1}^{N-1} \left( f(\boldsymbol{\theta}^{(i+1)}) - f(\hat{\boldsymbol{\theta}}) \right) \\
&\stackrel{??}{\leq} \frac{1}{2\eta(N-1)} \sum_{i=1}^{N-1} \left( \left\| \boldsymbol{\theta}^{(i)} - \hat{\boldsymbol{\theta}} \right\|^2 - \left\| \boldsymbol{\theta}^{(i+1)} - \hat{\boldsymbol{\theta}} \right\|^2 \right) \\
&= \frac{1}{2\eta(N-1)} \left( \left\| \boldsymbol{\theta} - \hat{\boldsymbol{\theta}} \right\|^2 - \left\| \boldsymbol{\theta}^{(N)} - \hat{\boldsymbol{\theta}} \right\|^2 \right) \\
&\leq \frac{1}{2\eta(N-1)} \left\| \boldsymbol{\theta} - \hat{\boldsymbol{\theta}} \right\|^2
\end{aligned} \tag{2.62}$$

Let now be  $\varepsilon > 0$ . Then for  $0 < \eta \leq \frac{1}{\delta}$  and

$$N > \frac{1}{2\eta\varepsilon} \left\| \boldsymbol{\theta} - \hat{\boldsymbol{\theta}} \right\|^2 + 1 \tag{2.63}$$

it follows, that:

$$\begin{aligned}
\left\| f(\boldsymbol{\theta}^{(N)}) - f(\hat{\boldsymbol{\theta}}) \right\| &= f(\boldsymbol{\theta}^{(N)}) - f(\hat{\boldsymbol{\theta}}) \\
&\stackrel{2.62}{\leq} \frac{1}{2\eta(N-1)} \left\| \boldsymbol{\theta} - \hat{\boldsymbol{\theta}} \right\|^2 \stackrel{2.63}{<} \varepsilon
\end{aligned}$$

□

## Chapter 3

# Approximative tractable inference in Energy Based Models

### Overview

*As the previous chapter closes by showing, that Energy based models allow a local maximum likelihood estimation of principal manifolds, the current chapter is concerned with a tractable approximation of this estimation. The first section starts by showing, that the log-likelihood gradient, as given by the free energy representation, for non-trivial cases, can not be solved analytically due to the underlying Boltzmann distribution. To address this issue, in a first step, a Monte Carlo (MC) approximation is applied to define a stochastic gradient descent, known as Mini-Batch Gradient Descent. Thereupon it is shown, that this algorithm under some additional statistical assumptions, almost surely converges in a local ML estimation. As the MC approximation, however, requires, that exact samples have to be drawn from the model it can not directly be applied to EBMs. Therefore, in a second step, the MC approximation by itself is approximated by a Markov Chain Monte Carlo (MCMC) approximation of the log-likelihood gradient which provides a corresponding stochastic gradient descent, known as Contrastive Divergence (CD) learning. The first section closes by showing, that CD learning under some supplementary analytical assumptions, regarding the function space, almost surely converge “close” to a local ML estimation, where the bias is quantified by an upper bound given by a respective KL -Divergence. Since the*

theoretical convergence, however, it not sufficient for tractable inference, the second section relates the tractability of CD learning to the mixing time of the Markov chain of the MCMC approximation and provides theoretical bounds. Finally with Restricted Boltzmann Machines (RBM) and Deep Boltzmann Machines (DBM) specifically structured EBM's are introduced that for an additional sparsity assumption assure the rapid mixing property and therefore tractable inference by CD learning.

### 3.1 Computability and Sampling in EBM's

Due to the structure of EBM's, the number of potential states increases exponentially with the number of vertices in the underlying graph. This property has significant consequences for the computability of the Steepest Gradient Descent in EBM's: By reconstituting its definition in Corollary 2.2, it appears, that for any clique  $c \in \mathcal{C}$  the update for the clique parameter  $\theta_c$  is proportional to the partial derivative  $\partial_c \log L_{\mathbf{v}}$ , which according to equation 2.53 decomposes in two parts. These are respectively referred to as the **positive phase** and the **negative phase**:

$$\partial_c \log L_{\mathbf{v}}(\theta) = n \underbrace{\langle -\partial_c \log F_{\theta} \rangle_{\mathbf{d}}}_{\text{=:positive phase}} - n \underbrace{\langle -\partial_c \log F_{\theta} \rangle_{\mathbf{m}}}_{\text{=:negative phase}} \quad (3.1)$$

Thereby the naming reflects the respective effects for statistical inference: Whereas the positive phase increases the probability for a realization  $\mathbf{v}$  of  $\mathbf{V}$ , the negative phase decreases the probability of predictions, generated by the model distribution. With respect to the empirical data distribution  $P_{\mathbf{v}}$  the positive phase thereby equates to:

$$\begin{aligned} \langle -\partial_c \log F_{\theta} \rangle_{\mathbf{d}} &\stackrel{\text{def}}{=} - \int_S P_{\mathbf{v}}(\mathbf{V} = \mathbf{v}) \partial_c \log F_{\theta}(\mathbf{v}) \, \mathrm{d}\mathbf{v} \\ &= -\frac{1}{n} \sum_{i=1}^n \partial_c \log F_{\theta}(\mathbf{v}^i) \end{aligned} \quad (3.2)$$

With regard to a log-linear parametrisation, the partial derivative of the free energy, then can be rewritten as an expectation of the feature  $\phi_c$  with respect to the condi-

tional model distribution  $P_{\theta}(\mathbf{H} \mid \mathbf{V})$ . In order to derive a complete representation of the positive phase as an integral over  $\phi_c$ , the empirical distribution  $P_v(\mathbf{V})$  can be used to define the **extended data distribution**:

$$P_{v,\theta}(\mathbf{V}, \mathbf{H}) := P_v(\mathbf{V})P_{\theta}(\mathbf{H} \mid \mathbf{V}) \quad (3.3)$$

Then the log-likelihood gradient has a feature representation, which is given by:

$$\partial_c \log L_v(\theta) = \underbrace{n \mathbb{E}_{\mathbf{X} \sim P_{v,\theta}}(\phi_c(\mathbf{X}_c))}_{\text{positive phase}} - \underbrace{n \mathbb{E}_{\mathbf{X} \sim P_{\theta}}(\phi_c(\mathbf{X}_c))}_{\text{negative phase}} \quad (3.4)$$

In order to calculate the positive phase, it has to be distinguished, if the clique  $c$  represents only observables, or also latent variables. In the first case the feature  $\phi_c$ , by definition, does not depend on any latent random variables, such that:

$$\begin{aligned} \mathbb{E}_{\mathbf{X} \sim P_{v,\theta}}(\phi_c(\mathbf{X}_c)) &= - \int_S \phi_c(\mathbf{x}_c) P_{v,\theta}(\mathbf{X}_c = \mathbf{x}_c) d\mathbf{x} \\ &= - \frac{1}{n} \sum_{i=1}^n \int_{S_H} \phi_c(\mathbf{v}_c^i, \mathbf{h}_c) P_{\theta}(\mathbf{H} = \mathbf{h} \mid \mathbf{V} = \mathbf{v}^i) d\mathbf{h} \\ &= - \frac{1}{n} \sum_{i=1}^n \phi_c(\mathbf{v}_c^i) \int_{S_H} P_{\theta}(\mathbf{H} = \mathbf{h} \mid \mathbf{V} = \mathbf{v}^i) d\mathbf{h} \\ &= - \frac{1}{n} \sum_{i=1}^n \phi_c(\mathbf{v}_c^i) \end{aligned} \quad (3.5)$$

If, however, the clique  $c$  also comprises latent random variables, then the corresponding feature  $\phi_c$  also depends on latent random variables and therefore can not be pulled out of the integral over  $S_H$ . Then due to the underlying Boltzmann distribution, this integral, even for quite simple feature functions - like quadratic terms - is not analytically solvable. This motivates an approximation by **Monte Carlo** (MC) integration: Let  $\mathbf{v}$  by an i.i.d. realization of  $\mathbf{V}$  and  $\mathbf{X}_d \sim P_{v,\theta}$  an  $S$ -valued i.i.d random sample of length  $m$ . Then for any  $c \in \mathcal{C}$  also  $\phi_c(\mathbf{X}_{d,c})$  is a random sample of length  $m$  and describes data distributed realizations of the feature functions  $\phi_c$ . Furthermore let:

$$\hat{\mathbf{MC}}_{d,c}^{(m)} := \frac{1}{m} \sum_{i=1}^m \phi_c(\mathbf{X}_{d,c})^i, \forall m \in \mathbb{N} \quad (3.6)$$

Since the expectation of the random sample  $\boldsymbol{\mu}_c := \mathbb{E}(\boldsymbol{\phi}_c(\mathbf{X}_{d,c}))$  by its definition equals the positive phase it can be concluded from the law of large numbers, that:

$$\hat{\mathbf{M}}\mathbf{C}_{d,c}^{(m)} \xrightarrow{\text{a.s.}} \mathbb{E}_{\mathbf{X} \sim P_{v,\boldsymbol{\theta}}}(\boldsymbol{\phi}_c(\mathbf{X})), \text{ for } m \rightarrow \infty \quad (3.7)$$

Consequently  $\hat{\mathbf{M}}\mathbf{C}_{d,c}^{(m)}$  is an unbiased estimator for  $\mathbb{E}_{\mathbf{X} \sim P_{v,\boldsymbol{\theta}}}(\boldsymbol{\phi}_c)$ . Let now be  $\mathbf{X}_m$  an  $S$ -valued i.i.d random sample with  $\mathbf{X}_m \sim P_{\boldsymbol{\theta}}$  and let:

$$\hat{\mathbf{M}}\mathbf{C}_{m,c}^{(m)} := \frac{1}{m} \sum_{i=1}^m \boldsymbol{\phi}_c(\mathbf{X}_{m,c})^i, \forall m \in \mathbb{N} \quad (3.8)$$

Then analogues to the positive phase, the sequence  $(\hat{\mathbf{M}}\mathbf{C}_{m,c}^{(m)})_{m \in \mathbb{N}}$  almost surely converges towards the negative phase:

$$\hat{\mathbf{M}}\mathbf{C}_{m,c}^{(m)} \xrightarrow{\text{a.s.}} \mathbb{E}_{\mathbf{X} \sim P_{\boldsymbol{\theta}}}(\boldsymbol{\phi}_c(\mathbf{X})), \text{ for } m \rightarrow \infty \quad (3.9)$$

This motivates the definition of an MC approximation for  $\partial_c \log L_v(\boldsymbol{\theta})$  by:

$$\hat{\mathbf{M}}\mathbf{C}_c^{(m)} := n \underbrace{\frac{1}{m} \sum_{i=1}^m \boldsymbol{\phi}_c(\mathbf{X}_{d,c})^i}_{\hat{\mathbf{M}}\mathbf{C}_{d,c}^{(m)}} - n \underbrace{\frac{1}{m} \sum_{i=1}^m \boldsymbol{\phi}_c(\mathbf{X}_{m,c})^i}_{\hat{\mathbf{M}}\mathbf{C}_{m,c}^{(m)}}, \forall m \in \mathbb{N} \quad (3.10)$$

Then according to equations 3.1, 3.7 and 3.9 it follows, that for all  $c \in \mathcal{C}$ :

$$\hat{\mathbf{M}}\mathbf{C}_c^{(m)} \xrightarrow{\text{a.s.}} \partial_c \log L_v(\boldsymbol{\theta}), \text{ for } m \rightarrow \infty \quad (3.11)$$

such that  $\hat{\mathbf{M}}\mathbf{C}_c^{(m)}$  yields an unbiased estimator for  $\partial_c \log L_v(\boldsymbol{\theta})$ . Let now  $m \in \mathbb{N}$  be fixed. Then the convergence motivates to use the MC approximation:

$$\partial_c \log L_v(\boldsymbol{\theta}) \approx \hat{\mathbf{M}}\mathbf{C}_c^{(m)}$$

to define an update rule for a stochastic gradient descent. This is known as a **Mini-Batch Gradient Descent** (see Algorithm 3.1).

The practical usability of the MC approximation, however, depends on its efficiency in terms of its convergence rate. For the positive phase let  $\boldsymbol{\phi}_c(\mathbf{X}_{d,c})$  have a

---

**Algorithm 3.1** Mini-Batch Gradient Descent in EBM
 

---

**parameters**

*prior parameter*  $\boldsymbol{\theta} \in \Theta$

*step size*  $\eta > 0$

*mini-batch size*  $m \in \mathbb{N}$

*number of steps*  $N \in \mathbb{N}$

1: **procedure** MBGD( $\boldsymbol{\theta}, \eta, m, N$ )

2:    $\boldsymbol{\theta}^{(1)} \leftarrow \boldsymbol{\theta}$

3:   **for**  $i \in \{1, \dots, N-1\}$  **do**

4:     Sample  $\mathbf{x}_d$  from  $P_{\mathbf{v}, \boldsymbol{\theta}}$  of length  $m$

5:     Sample  $\mathbf{x}_m$  from  $P_{\boldsymbol{\theta}}$  of length  $m$

6:     **for**  $c \in \mathcal{C}$  **do**

7:

$$\boldsymbol{\theta}_c^{(i+1)} \leftarrow \boldsymbol{\theta}_c^{(i)} - \eta \frac{N_d}{m} \sum_{j=1}^m (\phi_c(\mathbf{x}_{d,c}^j) - \phi_c(\mathbf{x}_{m,c}^j))$$

8:     **end for**

9:   **end for**

10:   **return**  $\boldsymbol{\theta}^{(N)}$

11: **end procedure**

---

finite variance  $\boldsymbol{\sigma}_c^2$ . Then the Central Limit Theorem postulates, that the sequence

$$(\sqrt{m}(\hat{\mathbf{M}}\mathbf{C}_{d,c}^{(m)} - \boldsymbol{\mu}_c))_{m \in \mathbb{N}}$$

has a stationary distribution, which is given by a normal distribution  $\mathcal{N}(0, \boldsymbol{\sigma}_c^2)$ , such that:

$$m\text{Var}(\hat{\mathbf{M}}\mathbf{C}_{d,c}^{(m)}) \xrightarrow{\text{a.s.}} \boldsymbol{\sigma}_c^2, \text{ for } m \rightarrow \infty \quad (3.12)$$

It therefore can be concluded, that for  $m \rightarrow \infty$  the variance decreases by the factor  $1/m$  and the convergence rate is given by  $1/\sqrt{m}$ . For the case, that  $\phi_c(\mathbf{X}_{d,c})$  also has a finite third absolute moment  $\boldsymbol{\rho}_c$ , the Berry-Esseen Theorem allows to extend this result to the individual distributions, that approximate  $\mathcal{N}(0, \boldsymbol{\sigma}_c^2)$ . In this case the Monte Carlo integration  $\hat{\mathbf{M}}\mathbf{C}_{d,c}^{(m)}$  efficiently approximates the positive phase with a convergence rate  $1/\sqrt{m}$ . And if furthermore  $\phi_c(\mathbf{X}_{m,c})$  has a finite second and a finite third absolute moment, then also  $\hat{\mathbf{M}}\mathbf{C}_{d,c}^{(m)}$  and therefore  $\hat{\mathbf{M}}\mathbf{C}_c^{(m)}$  has a convergence rate  $1/\sqrt{m}$ . Since the standard error of this approximation is furthermore proportional to  $\|\partial_c \log L_{\mathbf{v}}(\boldsymbol{\theta})\|$ , the intuition suggests, that in “sufficiently harmless” likelihood

landscapes for some  $m \in \mathbb{N}$  the Mini-Batch Gradient Descent should converge towards a parameter, that locally maximizes  $L_{\mathbf{v}}$ . Thereby, however, it has to be considered, that the approximation can mislead by a consistent realization of very improbable examples, such that - in difference to the Steepest Gradient Descent  $\mathbf{GD}$  - the best possible result, that can be achieved by  $\mathbf{MBGD}$  is almost sure convergence.

**Theorem 3.1** (Convergence of  $\mathbf{MBGD}$  in EBMs). *Let  $(S, \Sigma, \mathcal{M})$  be an EBM for  $\mathbf{X} = (\mathbf{V}, \mathbf{H})$  over  $\mathcal{G}$  with cliques  $\mathcal{C}$  and a log-linear parametrization  $P: \Theta \rightarrow \mathcal{M}$ , that additionally to (A1), (A2) and (A3) satisfies:*

(A4) *for  $k \in \{2, 3, 4\}$  there exist  $a_k, b_k \geq 0$ , such that:*

$$\int_S \|\nabla P_{\boldsymbol{\theta}}(\mathbf{x})\|^k d\mathbf{x} < a_k + b_k \|\boldsymbol{\theta}\|^k, \forall \boldsymbol{\theta} \in \Theta$$

*Furthermore for a realization  $\mathbf{v}$  of  $\mathbf{V}$  let the critical points of  $L_{\mathbf{v}}: \boldsymbol{\theta} \mapsto P_{\boldsymbol{\theta}}(\mathbf{v})$  respectively be given by  $\Theta_{\min}$  and  $\Theta_{\max}$ . Then for all mini-batch sizes  $m \in \mathbb{N}$  and all prior parameters  $\boldsymbol{\theta} \in \Theta \setminus \Theta_{\min}$ , there exists a step size  $\eta > 0$  and a local maximum likelihood estimation  $\hat{\boldsymbol{\theta}}_{\text{ML}} \in \Theta_{\max}$ , such that:*

$$\mathbf{MBGD}(\boldsymbol{\theta}, \eta, m, N) \xrightarrow{\text{a.s.}} \hat{\boldsymbol{\theta}}_{\text{ML}}, \text{ for } N \rightarrow \infty$$

*Notes on the proof of Theorem 3.1.* The Theorem has been proven by (Bottou 1998, p26 ff.) and uses the framework of stochastic approximation theory as introduced by (Kushner et al. 1978) and (Ljung et al. 1983). Thereby the proof essentially utilizes the assumption (A4) to provide an upper bound for the standard error of the approximation  $\hat{\mathbf{m}}_{\mathcal{C}}^{(m)}$  for the case  $m = 1$ . Since the standard error, however, is proportional to  $1/\sqrt{m}$ , the case  $m = 1$  already describes the “worst case”. Consequently the a.s. convergence is in particular assured for an  $m \in \mathbb{N}$  with  $m > 1$ .  $\square$

The Mini-Batch Gradient Descent efficiently solves the computability issues of the Steepest Gradient Descent. Nevertheless, it can not directly be applied to EBMs, due to the intractable task to sample  $\mathbf{X}_{\text{d}}$  and  $\mathbf{X}_{\text{m}}$  from the extended data distribution and the model distribution: This can be understood as follows: Usually one would try to start with a set of random variables, that respectively are statistically independent from all other variables. For an underlying acyclic directed graph, this would

be given by those random variables that are represented by the roots of the graph. Then in a series of subsequent steps all further random variables could successively be sampled from previous realizations, until all random variables are realized. For EBMs, however, the situation is much more complicated, since due to the undirected edges no random variable can be chosen, that are independent from the remaining, such that exact inference is generally intractable. This motivates approximative inference in EBMs, where the two main classes of algorithms can be distinguished (Wainwright et al. 2007): The first class comprises variational methods, such as mean-field approximations. Thereby the strategy of the methods is, to replace the model space by a “tractable” family of probability distributions, like fully-factorized distributions, which typically allows a fast and accurate approximation. Nevertheless this approximation is restricted to the “tractable” family, such that with respect to the original model space the estimates can be correspondingly worse. A very different strategy is applied by **Markov chain Monte Carlo** (MCMC) based algorithms, which introduce a **Markov chain** within the sampling procedure, which stationary distribution equals the desired target distribution. Thereby the idea is to define an initial  $S$ -valued random vector  $\mathbf{X}^{(0)}$  with a known probability distribution  $P^{(0)}$ , e.g. a uniform distribution, and thereupon to repeatedly apply a **transition function**  $T: (S, \Sigma) \rightarrow (S, \Sigma)$ , which generates a sequence of random vectors  $(\mathbf{X}^{(t)})_{t \in \mathbb{N}}$  by:

$$\mathbf{X}^{(0)} \xrightarrow{T} \mathbf{X}^{(1)} \xrightarrow{T} \mathbf{X}^{(2)} \xrightarrow{T} \dots \quad (3.13)$$

Then the sequence, by its definition is a Markov chain, since every random vector  $\mathbf{X}^{(t+1)}$  conditionally only depends on its predecessor  $\mathbf{X}^{(t)}$ . Such Markov chains have an intuitive representation within the function space  $\mathcal{P} \subseteq L^1(S, \Sigma)$  of probability distributions over  $(S, \Sigma)$ . By defining the **transition operator**  $\tau: \mathcal{P} \rightarrow \mathcal{P}$  as the push forward operator  $\tau(P) := T_*P$ ,  $\forall P \in \mathcal{P}$ , the Markov chain 3.13 is accompanied by a sequence of probability distributions  $(P^{(t)})_{t \in \mathbb{N}}$  in  $\mathcal{P}$ :

$$P^{(0)} \xrightarrow{\tau} P^{(1)} \xrightarrow{\tau} P^{(2)} \xrightarrow{\tau} \dots$$

Due to the  $L_1$ -norm of the embedding space  $L^1(S, \Sigma)$ , it follows, that  $\mathcal{P}$  is a normed space with respect to the induced norm, and therefore in particular a topological



space. Furthermore, with respect to the equivalence of measures in  $L^1$ ,  $\mathcal{P}$  is also a vector space and therefore a topological vector space. These properties can be used to show, that the transition operator is a continuous linear operator on  $\mathcal{P}$  (Zucca 2002), which motivates to regard  $(P^{(t)})_{t \in \mathbb{N}}$  as a trajectory within a dynamical system, where the **evolution function**  $\tau: \mathbb{N} \times \mathcal{P} \rightarrow \mathcal{P}$  is given by::

$$\tau^t(P) := \underbrace{(\tau \circ \dots \circ \tau)}_t(P), \forall P \in \mathcal{P} \quad (3.14)$$

The challenge for sampling in EBMs can therefore be described by the task to determine an evolution function  $\tau$  and an initial point  $P^{(0)}$ , such that:

**(MC1)** The Markov chain  $(\mathbf{X}^{(t)})_{t \in \mathbb{N}}$  has a stationary distribution  $\pi$

**(MC2)** The stationary distribution  $\pi$  equals the desired target distribution

In the first step, it has is shown, that the Markov chain has a stationary distribution.

**Proposition 3.1.** *Let  $(S, \Sigma)$  be a measurable space and  $\mathbf{X}$  an  $S$ -valued MRF over an undirected connected graph  $\mathcal{G}$ . Then for all measurable functions  $T: (S, \Sigma) \rightarrow (S, \Sigma)$  the transition operator  $\tau$  which is induced by  $T$  has a unique fixed point.*

*Proof of Proposition 3.1.* For the discrete case, that considers a finite state space, the proof is given by (Gelfand et al. 1990) and comprises two steps: In the first step the existence of a fixed point for  $\tau$  is shown as a consequence of a fixed point theorem. In the second step the property of the underlying graph to be connected is used to show its uniqueness. Thereupon according to (Casella et al. 1992, p170ff) the continuous case follows by a generalization of the respective discrete concepts by their continuous pendants.

**Existence** Let  $P$  be a probability distribution over a finite measurable space  $(S, \Sigma)$  with  $S = \{1, \dots, n\}$  and  $\Sigma = \mathcal{P}(S)$ . Then any state  $i \in S$  has a probability:

$$p_i := P(\mathbf{X} = i) \in [0, 1] \quad (3.15)$$

such that  $P$  can be represented by a vector  $\mathbf{p} \in [0, 1]^n$ . Let now be  $\mathcal{P} \subseteq L^1(S, \Sigma)$  the set of all probability distributions over  $(S, \Sigma)$ . Then by the normalization condition

$$\sum_{i=1}^n p_i = 1 \quad (3.16)$$

the set  $\mathcal{P}$  forms a convex and compact subset of  $\mathbb{R}^n$ , which is termed a **probability simplex**. Since furthermore the transition operator  $\tau$ , that is induced by the transition function  $T$  is a continuous linear operator on  $\mathcal{P}$ , the Brouwer Fixed Point Theorem can be applied and it follows that  $\tau$  at least has one fixed point in  $\mathcal{P}$ .

**Uniqueness** For any states  $i, j \in S$ , the transition from state  $i$  to  $j$  has a unique transition probability:

$$p_{i,j} := P(T \circ \mathbf{X} = j \mid \mathbf{X} = i) \in [0, 1] \quad (3.17)$$

such that, the transition operator  $\tau$  can be represented by a transition matrix  $A := (p_{i,j})_{i,j \in S}$ . Then the trajectory of an initial distribution  $P^{(0)} \in \mathcal{P}$  with a representation  $\mathbf{p}^{(0)} \in \mathbb{R}^n$ , has a representation, given by:

$$\mathbf{p}^{(t)} = A^t \mathbf{p}^{(0)}, \forall t \in \mathbb{N} \quad (3.18)$$

Therefore by assuming that  $\mathbf{p}$  represents a fixed point of  $\tau$ , then  $\mathbf{p}$  has to satisfy the relationship  $\mathbf{p} = A\mathbf{p}$ , such that  $\mathbf{p}$  is an Eigenvector of  $A$ . Thereby the matrix  $A$  by definition is non-negative. Thereupon, however, by the requirement, that the underlying undirected graph  $\mathcal{G}$  is complete, it follows that  $A$  is indeed strictly positive, since all vertices within  $\mathcal{G}$  are path connected such that  $p_{i,j} > 0$  for all states  $i, j \in S$ . This allows an application of the Perron-Frobenius theorem, which states, that in this case the largest Eigenvalue  $\lambda_1$  has an algebraic multiplicity of 1 and equals the value 1. Consequently there exists a unique Eigenvector  $\mathbf{v}_1$  to the Eigenvalue 1 and for all other Eigenvectors  $\mathbf{v}_2, \dots, \mathbf{v}_n$  it holds, that their respective eigenvalues  $\lambda_2 \dots, \lambda_n$  are strictly smaller than 1, such that their respective Eigenspaces decay to 0 for  $t \rightarrow \infty$ . This shows that  $\tau$  has a unique fixed point.  $\square$

According to Proposition 3.1 for an underlying connected graph any Markov chain  $(\mathbf{X}^{(t)})_{t \in \mathbb{N}}$ , which is induced by a measurable function  $T$ , has a unique stationary distribution  $\pi$ . Furthermore, however, it shows that the stationary distribution does not depend on the the initial distribution  $P^{(0)}$ , such that requirement **(MC1)** is always satisfied. This allows e.g. to draw the initial random sample from a uniform distributions by:

$$X^{(0)} \sim P^{(0)} := \mathcal{U}_V \otimes \mathcal{U}_H$$

Thereupon in order to also satisfy requirement **(MC2)** in the next step a transition function  $T$  has to be determined, such that the stationary distribution  $\pi$  equals the desired target distribution. For a probability distribution, which is completely known by its conditional probabilities, such a transition function is given by the **Gibbs Sampler** (Geman et al. 1984).

**Definition** (Gibbs Sampler). *Let  $(S, \Sigma)$  be a measurable space with  $S \subseteq \mathbb{R}^d$  and  $\mathbf{X}$  an  $S$ -valued random vector. Then a measurable function  $T: (S, \Sigma) \rightarrow (S, \Sigma)$  is a systematic sweep Gibbs sampler for a probability distribution  $P$  over  $(S, \Sigma)$ , iff for  $\mathbf{X}' := T \circ \mathbf{X}$  it holds, that:*

$$X'_i \sim P(X_i \mid \{X'_j\}_{1 \leq j < i}, \{X_j\}_{i < j \leq d}), \forall i \in \{1, \dots, d\} \quad (3.19)$$

The recursive definition of the systematic sweep Gibbs sampler can easily be entangled by iterating over the indices: For a given  $S$ -valued random vector  $\mathbf{X}$  the application of the measurable function  $T$  provides a further  $S$ -valued random vector  $\mathbf{X}'$ . Thereby according to its definition, the random variable  $X'_1$ , conditionally only depends on the random variables  $\{X_j\}_{1 < j \leq d}$  and therefore can be sampled at first. Since any subsequent steps, do not depend on  $X_1$ , but only on  $X'_1$ , the random variable  $X'_1$  can be thought as an update of  $X_1$ . Thereupon iteratively, the respective  $i$ -th random variable  $X'_i$  is sampled from the previously updated random variables  $\{X'_j\}_{1 \leq j < i}$ , and the remaining initial random variables  $\{X_j\}_{i < j \leq d}$  until the last random variable  $X'_d$  is sampled from the previously updated random variables  $\{X'_j\}_{1 \leq j < d}$ . Let now be  $\tau$  the transition operator, which is induced by a Gibbs sampler  $T$  for a probability distribution  $P$ . Then the above iteration can be used to show that  $P$  is a fixed point for  $\tau$ .

**Corollary 3.1.** *Let  $(S, \Sigma, \mathcal{P})$  be an EBM for  $\mathbf{X}$  over an undirected connected graph  $\mathcal{G}$  and  $P: \Theta \rightarrow \mathcal{P}$  a parametrization of  $\mathcal{P}$ . Then for all  $\boldsymbol{\theta} \in \Theta$ , the Gibbs sampler  $T_{\boldsymbol{\theta}}: (S, \Sigma) \rightarrow (S, \Sigma)$  for  $P_{\boldsymbol{\theta}}$  induces a Markov chain that converges in probability towards  $P_{\boldsymbol{\theta}}$ .*

*Proof of Corollary 3.1.* Let  $\tau = T_{\boldsymbol{\theta}*}$  be the transition operator, induced by  $T_{\boldsymbol{\theta}}$ . Then for a discrete sample space according to Proposition 3.1  $\tau$  has a unique fixed point. It therefore suffices to show, that  $P_{\boldsymbol{\theta}}$  is a fixed point of  $\tau$ . This can be proofed by complete induction. Let  $\mathbf{X}' := T_{\boldsymbol{\theta}} \circ \mathbf{X}$ . Then:

**Induction Start** For  $i = 1$ , it holds that:

$$\begin{aligned} \tau(P_{\boldsymbol{\theta}})(\mathbf{X})_1 &\stackrel{\text{def}}{=} (T_{\boldsymbol{\theta}*}P_{\boldsymbol{\theta}})(\mathbf{X})_1 \\ &\stackrel{\text{def}}{=} \int_{\mathbf{x} \in S} P_{\boldsymbol{\theta}}(X_1 \mid \{X_j\}_{1 < j \leq d}) P_{\boldsymbol{\theta}}(\{X_j = x_j\}_{1 < j \leq d}) d\mathbf{x} \\ &\stackrel{\text{tp}}{=} P_{\boldsymbol{\theta}}(X_1) \\ &\Rightarrow X'_1 \sim X_1 \end{aligned} \tag{3.20}$$

**Induction Step** Let  $X'_j \sim X_j, \forall j \leq i$ . Then for  $i \mapsto i + 1$ , it holds, that:

$$\begin{aligned} \tau(P_{\boldsymbol{\theta}})(\mathbf{X})_{i+1} &\stackrel{\text{def}}{=} \int_S P_{\boldsymbol{\theta}}(X_{i+1} \mid \{Y_j\}_{1 \leq j \leq i}, \{X_j\}_{i+1 < j \leq d}) \\ &\quad P_{\boldsymbol{\theta}}(\{Y_j = x_j\}_{1 \leq j \leq i}, \{X_j = x_j\}_{i+1 < j \leq d}) d\mathbf{x} \\ &= \int_S P_{\boldsymbol{\theta}}(X_{i+1} \mid \{X_j = x_j\}_{j \neq i+1}) \\ &\quad P_{\boldsymbol{\theta}}(\{X_j = x_j\}_{j \neq i+1}) d\mathbf{x} \\ &\stackrel{\text{tp}}{=} P_{\boldsymbol{\theta}}(X_{i+1}) \\ &\Rightarrow X'_{i+1} \sim X_{i+1} \end{aligned} \tag{3.21}$$

Equations 3.20 and 3.21 show, that  $\tau(P_{\boldsymbol{\theta}})(\mathbf{X})_i = P_{\boldsymbol{\theta}}(X_i)$  for all  $i \in \{1, \dots, d\}$ , such that  $\tau(P_{\boldsymbol{\theta}}) = P_{\boldsymbol{\theta}}$ .  $\square$

In the following Corollary 3.1 is to be applied to the desired target distributions, which are given by the extended data distribution  $P_{\mathbf{v}, \boldsymbol{\theta}}$  and the model distribution  $P_{\boldsymbol{\theta}}$ . In the case of the model distribution for an initial distribution  $P_{\boldsymbol{\theta}}^{(0)}$ , which is chosen

by a uniform distribution  $\mathcal{U}_V \otimes \mathcal{U}_H$ , it follows that the Markov Chain  $(\mathbf{X}_m^{(k)})_{k \in \mathbb{N}}$  with:

$$\mathbf{X}_m^{(k)} \sim T_{\boldsymbol{\theta}^*}^k(\mathcal{U}_V \otimes \mathcal{U}_H), \forall i \in \mathbb{N} \quad (3.22)$$

has a stationary distribution  $P_{\boldsymbol{\theta}}$ . For the extended data distribution  $P_{\mathbf{v}, \boldsymbol{\theta}}$ , however, some further considerations have to be taken into account. Due to its definition by:

$$P_{\mathbf{v}, \boldsymbol{\theta}}(\mathbf{V}, \mathbf{H}) := P_{\mathbf{v}}(\mathbf{V})P_{\boldsymbol{\theta}}(\mathbf{H} | \mathbf{V})$$

it follows, that the observables are determined by the empirical data distribution  $P_{\mathbf{v}}(\mathbf{V})$  and therefore independent from the latent random variables, such that the conditional probabilities  $P_{\mathbf{v}, \boldsymbol{\theta}}(\mathbf{V} | \mathbf{H})$  are given by uniform distributions. If therefore the initial distribution  $P_{\mathbf{v}, \boldsymbol{\theta}}^{(0)}$  of the Markov chain is chosen according to:

$$P_{\mathbf{v}, \boldsymbol{\theta}}^{(0)}(\mathbf{V}, \mathbf{H}) = P_{\mathbf{v}}(\mathbf{V}) \otimes \mathcal{U}_H(\mathbf{H})$$

then the Gibbs sampler for  $P_{\mathbf{v}, \boldsymbol{\theta}}$  does not alter the marginal distribution of the observables, given by  $P_{\mathbf{v}}$ . It therefore suffices to only update the realizations of the latent random variables. Since the conditional probabilities  $P_{\mathbf{v}, \boldsymbol{\theta}}(\mathbf{H} | \mathbf{V})$  in turn are given by  $P_{\boldsymbol{\theta}}(\mathbf{H} | \mathbf{V})$  these updates can be applied by the Gibbs sampler  $T_{\boldsymbol{\theta}}$  for the model distribution  $P_{\boldsymbol{\theta}}$ , where the observables are **clamped** to their initial values, which is denoted by  $T_{\boldsymbol{\theta}, H}$ . Thereupon for  $i \in \mathbb{N}$  let  $T_{\boldsymbol{\theta}, H}^k$  denote the pushforward measure, which is induced by  $k$  applications of  $T_{\boldsymbol{\theta}, H}$ . Then the Markov chain  $(\mathbf{X}_d^{(k)})_{k \in \mathbb{N}}$  with:

$$\mathbf{X}_d^{(k)} \sim T_{\boldsymbol{\theta}, H}^k(P_{\mathbf{v}} \otimes \mathcal{U}_H), \forall k \in \mathbb{N} \quad (3.23)$$

converges in probability towards  $P_{\mathbf{v}, \boldsymbol{\theta}}$  for  $k \rightarrow \infty$ . It therefore can be concluded that samples  $\mathbf{X}_d \sim P_{\mathbf{v}, \boldsymbol{\theta}}$  and  $\mathbf{X}_m \sim P_{\boldsymbol{\theta}}$  can be approximated by the respective Markov Chains  $(\mathbf{X}_d^{(k)})_{k \in \mathbb{N}}$  and  $(\mathbf{X}_m^{(k)})_{k \in \mathbb{N}}$ . The next step therefore is to apply this approximation to the Monte Carlo approximations of the positive and the negative phase. In this purpose let the MCMC approximation for the positive phase be given by:

$$\text{MCMC}_{d,c}^{(m,k)} := \frac{1}{m} \sum_{i=1}^m \phi_c(\mathbf{X}_{d,c}^{(k)})^i, \forall m, k \in \mathbb{N} \quad (3.24)$$

Then according to Corollary 3.1 and equation 3.7 it follows, that:

$$\mathbf{MCMC}_{d,c}^{(m,k)} \xrightarrow{\text{a.s.}} \mathbb{E}_{\mathbf{X} \sim P_{v,\theta}}(\phi_c), \text{ for } m, k \rightarrow \infty \quad (3.25)$$

Furthermore let:

$$\mathbf{MCMC}_{m,c}^{(m,k)} := \frac{1}{m} \sum_{i=1}^m \phi_c(\mathbf{X}_{m,c}^{(k)})^i, \forall m, k \in \mathbb{N} \quad (3.26)$$

Then according to Corollary 3.1 and equation 3.9 it follows, that:

$$\mathbf{MCMC}_{m,c}^{(m,k)} \xrightarrow{\text{a.s.}} \mathbb{E}_{\mathbf{X} \sim P_\theta}(\phi_c), \text{ for } m, k \rightarrow \infty \quad (3.27)$$

Summarized by equation 3.4 it then follows, that for the definition:

$$\mathbf{MCMC}_c^{(m,k_d,k_m)} := n \underbrace{\frac{1}{m} \sum_{i=1}^m \phi_c(\mathbf{X}_{d,c}^{(k_d)})^i}_{\mathbf{MCMC}_{d,c}^{(m,k_d)}} - n \underbrace{\frac{1}{m} \sum_{i=1}^m \phi_c(\mathbf{X}_{m,c}^{(k_m)})^i}_{\mathbf{MCMC}_{m,c}^{(m,k_m)}}, \forall m, k_d, k_m \in \mathbb{N} \quad (3.28)$$

According to equations 3.25 and 3.27 it holds, that:

$$\mathbf{MCMC}_c^{(m,k_d,k_m)} \xrightarrow{\text{a.s.}} \partial_c \log L_v(\theta), \text{ for } m, k_d, k_m \rightarrow \infty \quad (3.29)$$

The almost sure convergence towards  $\partial_c \log L_v(\theta)$  motivates an approximation:

$$\partial_c \log L_v(\theta) \approx \mathbf{MCMC}_c^{(m,k_d,k_m)}$$

within the update rule of a stochastic gradient descent. Thereby, in difference to the Monte Carlo approximation  $\hat{\mathbf{M}}\mathbf{C}_c^{(m)}$ , it has to be noted, that for any fixed  $k_d, k_m \in \mathbb{N}$  the probability distributions  $P_{v,\theta}^{(k_d)}$  and  $P_\theta^{(k_m)}$ , which are generated by the Markov Chains, are generally different from their respective equilibrium distributions  $P_{v,\theta}$  and  $P_\theta$ , such that  $\mathbf{MCMC}_c^{(m,k_d,k_m)}$  is a biased estimator for  $\partial_c \log L_v(\theta)$ . One approach to make this bias tangible is to identify the gradient by the Kullback-Leibler divergence.

**Lemma 3.1.** *For any i.i.d realization  $\mathbf{v}$  of  $\mathbf{V}$  and any parameter  $\theta \in \Theta$  the following statements are equivalent:*

(1) The likelihood  $L_{\mathbf{v}}(\boldsymbol{\theta})$  is locally maximized by  $\boldsymbol{\theta}$

(2) The KL-Divergence  $D_{\text{KL}}(P_{\mathbf{v},\boldsymbol{\theta}}\|P_{\boldsymbol{\theta}})$  is locally minimized by  $\boldsymbol{\theta}$

*Proof of Lemma 3.1.* The KL-divergence is defined by:

$$D_{\text{KL}}(P_{\mathbf{v},\boldsymbol{\theta}}\|P_{\boldsymbol{\theta}}) := \int_S P_{\mathbf{v},\boldsymbol{\theta}}(\mathbf{x}) \log \frac{P_{\mathbf{v},\boldsymbol{\theta}}(\mathbf{x})}{P_{\boldsymbol{\theta}}(\mathbf{x})} d\mathbf{x} \quad (3.30)$$

Such that:

$$\begin{aligned} D_{\text{KL}}(P_{\mathbf{v},\boldsymbol{\theta}}\|P_{\boldsymbol{\theta}}) &\stackrel{3.30}{=} \int_S P_{\mathbf{v},\boldsymbol{\theta}}(\mathbf{v}, \mathbf{h}) \log \frac{P_{\mathbf{v}}(\mathbf{v})P_{\boldsymbol{\theta}}(\mathbf{h} | \mathbf{v})}{P_{\boldsymbol{\theta}}(\mathbf{v})P_{\boldsymbol{\theta}}(\mathbf{h} | \mathbf{v})} d(\mathbf{v}, \mathbf{h}) \\ &= \int_{S_V} \left( \int_{S_H} P_{\mathbf{v},\boldsymbol{\theta}}(\mathbf{v}, \mathbf{h}) d\mathbf{h} \right) \log \frac{P_{\mathbf{v}}(\mathbf{v})}{P_{\boldsymbol{\theta}}(\mathbf{v})} d\mathbf{v} \\ &\stackrel{\text{tp}}{=} \int_{S_V} P_{\mathbf{v}}(\mathbf{v}) \log P_{\mathbf{v}}(\mathbf{v}) d\mathbf{v} - \int_{S_V} P_{\mathbf{v}}(\mathbf{v}) \log P_{\boldsymbol{\theta}}(\mathbf{v}) d\mathbf{v} \end{aligned} \quad (3.31)$$

Since the first term on the right side does not depend on  $\boldsymbol{\theta}_c$ , it follows, that:

$$\begin{aligned} \partial_c D_{\text{KL}}(P_{\mathbf{v},\boldsymbol{\theta}}\|P_{\boldsymbol{\theta}}) &\stackrel{3.31}{=} -\partial_c \int_{S_V} P_{\mathbf{v}}(\mathbf{v}) \log P_{\boldsymbol{\theta}}(\mathbf{v}) d\mathbf{v} \\ &= -\partial_c \sum_{i=1}^n \log P_{\boldsymbol{\theta}}(\mathbf{v}^i) \\ &\stackrel{\text{def}}{=} -\partial_c \log L_{\mathbf{v}}(\boldsymbol{\theta}) \end{aligned} \quad (3.32)$$

It follows, that **(ML1)** has an equivalent formulation by:

$$\nabla L_{\mathbf{v}}(\boldsymbol{\theta}) = 0 \stackrel{3.32}{\Leftrightarrow} \nabla D_{\text{KL}}(P_{\mathbf{v},\boldsymbol{\theta}}\|P_{\boldsymbol{\theta}}) = 0$$

and furthermore **(ML2)** by:

$$\nabla^2 L_{\mathbf{v}}(\boldsymbol{\theta}) < 0 \stackrel{3.32}{\Leftrightarrow} \nabla^2 D_{\text{KL}}(P_{\mathbf{v},\boldsymbol{\theta}}\|P_{\boldsymbol{\theta}}) > 0$$

□

The KL-divergence  $D_{\text{KL}}(P\|Q)$  measures the discriminative information of  $P$  with respect to  $Q$  and may intuitively be thought as the distance of  $Q$  from  $P$  from the perspective of  $P$ . This allows an intuitive interpretation of Lemma 3.1: The likelihood

function  $L_{\mathbf{v}}$  is locally maximized, if and only if the distance between the observed data  $P_{\mathbf{v},\boldsymbol{\theta}}$  and the model prediction  $P_{\boldsymbol{\theta}}$  is locally minimized. By an MCMC approximation, for given  $k_d, k_m \in \mathbb{N}$ , the probability distributions  $P_{\mathbf{v},\boldsymbol{\theta}}$  and  $P_{\boldsymbol{\theta}}$  are then respectively substituted by their approximations  $P_{\mathbf{v},\boldsymbol{\theta}}^{(k_d)}$  and  $P_{\boldsymbol{\theta}}^{(k_m)}$ . It follows, that:

$$\lim_{m \rightarrow \infty} \text{MCMC}_c^{(m, k_d, k_m)} \stackrel{\text{a.s.}}{=} -\partial_c D_{\text{KL}}(P_{\mathbf{v},\boldsymbol{\theta}}^{(k_d)} \| P_{\boldsymbol{\theta}}^{(k_m)}), \forall c \in \mathcal{C} \quad (3.33)$$

**Lemma 3.2.** *The bias of a stochastic gradient descent, that uses the MCMC approximation for  $m \rightarrow \infty$  is upper bounded by:*

$$\text{bias}_{\text{MCMC}} \leq 2D_{\text{KL}}(P_{\boldsymbol{\theta}}^{(k_m)} \| P_{\boldsymbol{\theta}})$$

*Proof of Lemma 3.2.* By Lemma 3.1 and equation 3.33 it can be concluded, that in a stochastic gradient descent, that uses the MCMC approximation, for  $m \rightarrow \infty$  the fixed points of the gradient descent are given by parameters  $\boldsymbol{\theta}_{\text{CD}} \in \Theta$ , that satisfy the requirement:

$$\partial_c D_{\text{KL}}(P_{\mathbf{v},\boldsymbol{\theta}}^{(k_d)} \| P_{\boldsymbol{\theta}}^{(k_m)}) = 0, \forall c \in \mathcal{C} \quad (3.34)$$

Thereby the bias of an estimation  $\boldsymbol{\theta}_{\text{MCMC}}$  in the function space is upper bounded with respect to the corresponding local ML estimation  $\boldsymbol{\theta}_{\text{ML}}$ . This follows as a consequence of the Pythagorean theorem for dually flat manifolds (Amari 2016). It then follows, that:

$$\begin{aligned} \text{bias}_{\text{MCMC}} &= D_{\text{KL}}(P_{\mathbf{v},\boldsymbol{\theta}}^{(k_d)} \| P_{\boldsymbol{\theta}}^{(k_m)}) - D_{\text{KL}}(P_{\mathbf{v},\boldsymbol{\theta}} \| P_{\boldsymbol{\theta}}) \\ &\leq \underbrace{D_{\text{KL}}(P_{\mathbf{v},\boldsymbol{\theta}} \| P_{\mathbf{v},\boldsymbol{\theta}}^{(k_d)})}_{=:\text{bias}_d} + \underbrace{D_{\text{KL}}(P_{\boldsymbol{\theta}}^{(k_m)} \| P_{\boldsymbol{\theta}})}_{=:\text{bias}_m} \end{aligned} \quad (3.35)$$

By the definition of the extended data distribution, however, it follows that sufficiently close to the ML estimate  $\hat{\boldsymbol{\theta}}_{\text{ML}}$  it holds, that  $\text{bias}_d \leq \text{bias}_m$ , such that:

$$\text{bias}_{\text{MCMC}} \leq 2D_{\text{KL}}(P_{\boldsymbol{\theta}}^{(k_m)} \| P_{\boldsymbol{\theta}}) \quad (3.36)$$

□

Consequently in order to minimize the bias of the MCMC based stochastic gradi-



ent descent, the discriminative information between  $P_{\theta}^{(k_m)}$  and  $P_{\theta}$  is to be minimized. A very efficient approach is, to initialize the Markov Chain by a prior  $P_{\theta}^{(0)}$ , that is already “close” to its equilibrium distribution  $P_{\theta}$ . Then there is justified hope, that even a small number of sampling steps suffices to generate a posterior  $P_{\theta}^{(k_m)}$  that sufficiently approximates  $P_{\theta}$ . A key idea in this direction was, that near the fixed points of a gradient descent, the model distribution  $P_{\theta}$  will be found “close” to the extended data distribution  $P_{\mathbf{v},\theta}$ . Consequently in this case the posterior data distribution  $P_{\mathbf{v},\theta}^{(k_a)}$ , which approximates  $P_{\mathbf{v},\theta}$  will presumably provide a “good prior” for the model distribution. This consideration motivated an MCMC based stochastic gradient descent, which is referred as **contrastive divergence** learning (see algorithm 3.2, (author?) 2002).

Since its first publication (Hinton 2002) CD learning gained a lot of attention due to its empirically success. Thereby in many domains of application the CD algorithm appeared to generate sequences  $(\theta^{(t)})_{t \in \mathbb{N}}$  that converge “close” to local ML estimates  $\hat{\theta}_{\text{ML}}$ . Due to its high applicability the pragmatic consensus prevailed, to start the density estimation by CD learning with very short Markov chains and, in the case of convergence, to increase the chain length to get closer to a local ML estimate. This pragmatic approach was accompanied by a great enthusiasm in the machine learning community to further investigations on “good priors” which provided many variants of CD learning, including *Persistent CD* (Tieleman 2008), *Fast Persistent CD* (Tieleman et al. 2009) and *Tempered MCMC* (Desjardins et al. 2010). The fundamental question, if the algorithm converges, however, has only been answered recently and still requires assumptions, that are difficult to verify. Thereby analogous to the Mini-batch Gradient Descent, the convergence can be traced back to the standard error of the approximation of the gradient.

**Theorem 3.2** (Convergence of CD in EBMs). *Let  $(S, \Sigma, \mathcal{P})$  be an EBM for  $\mathbf{X} = (\mathbf{V}, \mathbf{H})$  over  $\mathcal{G}$  with cliques  $\mathcal{C}$  and a log-linear parametrization  $P: \Theta \rightarrow \mathcal{M}$ , that satisfies (A1), (A2), (A3), (A4) and:*

(A5) *The graph  $\mathcal{G}$  is connected*

(A6) *For any Gibbs sampler  $T_{\theta}$ , which is defined by  $P_{\theta}$ , let the induced transition operator  $\tau$  have an  $\mathcal{L}_2$ -spectral gap  $1 - \lambda_2 > 0$*

---

**Algorithm 3.2** Contrastive Divergence learning in EBMs

---

**parameters***prior distribution parameter  $\boldsymbol{\theta} \in \Theta$* *step size  $\eta > 0$* *mini-batch size  $m \in \mathbb{N}$* *positive phase chain length  $k_d$* *negative phase chain length  $k_m$* *number of steps  $N \in \mathbb{N}$* 

```
1: procedure CD( $\boldsymbol{\theta}, \eta, m, k_d, k_m, N$ )  
2:    $\boldsymbol{\theta}^{(1)} \leftarrow \boldsymbol{\theta}$   
3:   for  $i \in \{1, \dots, N-1\}$  do  
4:     Sample  $\mathbf{x}_d$  from  $P_v(\mathbf{V}) \otimes \mathcal{U}(\text{img} \mathbf{H})$  of length  $m$   
5:     for  $j \in \{1, \dots, k_d\}$  do  
6:       Sample  $\mathbf{x}_d$  from  $T_{\boldsymbol{\theta}, \mathbf{H}}(\mathbf{x}_d)$   
7:     end for  
8:     Sample  $\mathbf{x}_m \sim \mathbf{x}_d$  of length  $m$   
9:     for  $j \in \{1, \dots, k_m\}$  do  
10:      Sample  $\mathbf{x}_m \sim T_{\boldsymbol{\theta}}(\mathbf{x}_m)$   
11:    end for  
12:    for  $c \in \mathcal{C}$  do  
13:      
$$\boldsymbol{\theta}_c^{(i+1)} \leftarrow \boldsymbol{\theta}_c^{(i)} - \eta \frac{d}{m} \sum_{j=1}^d (\phi_c(\mathbf{x}_{d,c}^j) - \phi_c(\mathbf{x}_{m,c}^j))$$
  
14:    end for  
15:  end for  
16:  return  $\boldsymbol{\theta}^{(N)}$   
17: end procedure
```

---

(A7) Define a metric on the set of transition operators by:

$$\rho(\tau_{\boldsymbol{\theta}}, \tau_{\boldsymbol{\vartheta}}) := \sup_{\pi: |\pi| \leq 1} \sup_{\mathbf{x} \in S} |\tau_{\boldsymbol{\theta}}(\pi)(\mathbf{x}) - \tau_{\boldsymbol{\vartheta}}(\pi)(\mathbf{x})|$$

and assume, that  $\exists \xi > 0$ , such that:

$$\rho(\tau_{\boldsymbol{\theta}}, \tau_{\boldsymbol{\vartheta}}) \leq \xi \|\boldsymbol{\theta} - \boldsymbol{\vartheta}\|, \forall \boldsymbol{\theta}, \boldsymbol{\vartheta} \in \Theta$$

Then for all realizations  $\mathbf{v}$  of  $\mathbf{V}$ , prior parameters  $\boldsymbol{\theta} \in \Theta$  and mini-batch sizes  $m \in \mathbb{N}$ , there exists a step size  $\eta > 0$ , Markov Chain lengths  $k_d, k_m \in \mathbb{N}$  and a parameter  $\hat{\boldsymbol{\theta}}_{\text{CD}} \in \Theta$ , such that the sequence  $(\boldsymbol{\theta}^{(t)})_{t \in \mathbb{N}}$  which is generated by CD satisfies that:

$$\limsup_{t \rightarrow \infty} \left\| \frac{1}{t} \sum_{s=1}^t \boldsymbol{\theta}^{(s)} - \hat{\boldsymbol{\theta}}_{\text{CD}} \right\| \xrightarrow{P} 0$$

Notes on the proof of Theorem 3.2. The Theorem has been proven by (Wu et al. 2016) and uses the terminology of stochastic approximation theory. Thereby the enumeration of the assumptions as given in (Wu et al. 2016) differs from the above. This is because, in the above case, the assumptions are gradually expanded from the algorithms GD to MBGD and from MBGD to CD. Also the proof in (Wu et al. 2016) is given for generic exponential families, where in particular the above assumption (A5) is formulated as a positivity constraint of the transition density function. However, by reconstituting Proposition 3.1, it can be concluded that the given assumption (A5) of a connected graph  $\mathcal{G}$ , is sufficient in EBMs for the transition operator  $\tau_{\boldsymbol{\theta}}$  of a Gibbs Sampler  $T_{\boldsymbol{\theta}}$ , to have a unique fixed point. In this case the Markov Chain converges as the components of all eigenvalues of  $\tau_{\boldsymbol{\theta}}$  decay to zero, except for the largest Eigenvalue, which according to the Perron-Frobenius Theorem is identical to 1. Consequently the convergence rate is then determined by the **spectral gap**  $1 - \lambda_2$ , which denotes the difference between the largest Eigenvalue 1 and the second largest Eigenvalue  $\lambda_2$ . Therefore in order to enforce a strictly positive convergence rate, it has to be assumed, that (A6)  $1 - \lambda_2(\boldsymbol{\theta}) > 0, \forall \boldsymbol{\theta} \in \Theta$ . Finally it has to be avoided, that the sequence  $(\boldsymbol{\theta}_t)_{t \in \mathbb{N}}$ , which is generated by the gradient descent, causes arbitrary large jumps in the function space, which in (A7) is ensured by a corresponding Lipschitz continuity with respect to the function space. The proof given by (Wu et

al. 2016) treats the case of an infinite mini-batch size  $m \rightarrow \infty$  and an exact extended data distribution, such that  $P_{\mathbf{v}, \boldsymbol{\theta}}^{(k_d)} = P_{\mathbf{v}, \boldsymbol{\theta}}$ . This result, however, immediately is generalized, by an incorporation of assumption **(A4)** and the corresponding results from Theorem 3.1 as well as an additional Markov chain for the positive phase, defined by equation 3.23.  $\square$

## 3.2 Structured EBMs for Efficient Sampling

In the previous section, it has been shown, that with CD learning a stochastic gradient descent in EBMs, which uses an MCMC approximation of the gradient, under some generic assumptions is assured to converge. Nevertheless the theoretical existence of such a limit does not implicate, that its approximation is practically tractable. Thereby the decisive factor, which determines this behaviour, is the number of iterations which is required for the Markov chain to get “close” to its stationary distribution. This number has a vivid meaning in the respective  $L^1$ -space: Since the stationary Markov chain may be regarded as a dynamical system with a unique fixed point, this quantity equals the minimum number of steps, which is required for an arbitrary probability distribution to enter an  $\varepsilon$ -ball about the stationary distribution. This provides the notation of the **mixing time** of a Markov chain.

**Definition** (Mixing Time). *Let  $(\mathbf{X}^{(t)})_{t \in \mathbb{N}}$  be an irreducible aperiodic Markov chain with a state space  $(S, \Sigma)$ , a stationary distribution  $\pi$  and a transition operator  $\tau$ . Furthermore let  $\mathcal{P} \subset L^1(S, \Sigma)$  be the set of all probability distributions over  $(S, \Sigma)$ . Then for all  $t \in \mathbb{N}$  the maximal variation distance to the stationary distribution is defined by:*

$$\delta_{\max}(t) := \max_{\mu \in \mathcal{P}} \frac{1}{2} \|\tau^t(\mu) - \pi\|_1 \quad (3.37)$$

*Thereupon the for an  $\varepsilon > 0$  the  $\varepsilon$ -mixing time of  $(\mathbf{X}^{(t)})_{t \in \mathbb{N}}$  is defined by the smallest number of steps  $t$ , that is required for  $\delta_{\max}(t) \leq \varepsilon$ :*

$$t_{\text{mix}}(\varepsilon) := \arg \min_{t \in \mathbb{N}} (\delta_{\max}(t) \leq \varepsilon) \quad (3.38)$$

A Markov chain is then referred to be **torpid mixing**, if its mixing time increases exponential with the state space and **rapid mixing** if the mixing time maximally

has a polynomial growth. Since by an application to arbitrary EBMs the mixing time apparently can vary between these cases (Liu et al. 2014) it got all the more important to study this behaviour for specific classes of EBMs, or the other way around: To find EBMs where the Gibbs sampler is rapid mixing. It is therefore important to provide bounds for the mixing time, that depend on the underlying structure of the respective EBM. An important assumption, that has been incorporated within the previous section was **(A5)**, which requires, that the underlying graph is connected. This assumption assures, that the Markov chain has a unique stationary distribution, since for an EBM over a complete graph any state transition has a strictly positive probability. Then according to the Perron-Frobenius Theorem the largest eigenvalue  $\lambda_1$  of the transition operator  $\tau$  equals 1 and has an algebraic multiplicity of 1, such that its corresponding eigenvector comprises the only component, that resists a repeated application of  $\tau$ , whereas all others sooner or later decay to zero. At this point it is getting apparent, that the second (absolutely) largest eigenvalue  $\lambda_2$  and therefore the spectral gap  $1 - \lambda_2$  has some considerable influence on the mixing time. This property has been exploited in Theorem 3.2 by assumption **(A6)**, which requires that  $1 - \lambda_2 > 0$ , to ensure a positive convergence rate and therefore a finite mixing time. Moreover, however, the spectral gap can also be used to quantitatively bound the mixing time.

**Proposition 3.2.** *Let  $(\mathbf{X}^{(t)})_{t \in \mathbb{N}}$  be an irreducible aperiodic Markov chain with a state space  $(S, \Sigma)$ , an absolutely positive stationary distribution  $\pi$  and a transition operator  $\tau$ . Furthermore let  $\tau$  have a spectral gap  $1 - \lambda_2$ . Then for all  $0 < \varepsilon < 1$  it holds, that:*

$$(t_{\text{rel}} - 1) \ln \frac{1}{2\varepsilon} \leq t_{\text{mix}}(\varepsilon) \leq \left\lceil t_{\text{rel}} \left( \frac{1}{2} \ln \frac{1}{\pi_{\min}} + \ln \frac{1}{2\varepsilon} \right) \right\rceil \quad (3.39)$$

where:

$$t_{\text{rel}} := \frac{1}{1 - \lambda_2}, \text{ and } \pi_{\min} := \min_{\mathbf{x} \in S} \pi(\mathbf{X} = \mathbf{x})$$

*Notes on the proof of Proposition 3.2.* The Proposition has originally been proven for continuous-time reversible Markov processes by (Aldous 1982). On this basis the upper bound for time-discrete Markov chains has been proven by (Diaconis et al. 1991) and the lower bound by (Levon et al. 2009, Theorem 12.4).  $\square$

A fundamental disadvantage of the spectral bounds, as given by Proposition 3.2

is, that the spectrum of the transition operator is usually not tractable. Thereupon, however, the spectral gap also provides no intuitive interpretation with respect to the underlying structure. The next step is therefore to find an invariant quantity of the underlying Markov chain with an intuitive meaning for its structure. For this purpose, it is closer to the intuition, to consider the Markov chain as a random walk within a graph over the state space. Thereby the states can be regarded as the vertices of a graph and the transitions between the states as the edges. The mixing time may then identifies the minimum number of steps, which are required, until the probability of any position is sufficiently close to its stationary probability. It is clear, that this property requires, that the random walk is able to traverse the whole state space, such that its representation by a graph requires that the graph is connected. Thereupon, however, it also follows intuitively, that a “bottleneck” in the graph locally impedes a traversal and therefore delays the mixture. In order to transfer this idea to the whole state space, a measure is required, that quantifies its “global connectivity”. This quantity is given by the **conductance** (King 2003, p1).

**Definition** (Conductance). *Let  $(\mathbf{X}^{(t)})_{t \in \mathbb{N}}$  be an irreducible aperiodic Markov chain with a state space  $(S, \Sigma)$ , a strictly positive stationary distribution  $\pi$  and a transition operator  $\tau$  with a density  $d\tau$ . Then for any  $s \subset S$  and  $\bar{s} := S \setminus s$  the conductance of the cut  $(s, \bar{s})$  is defined by:*

$$\Phi_s := \frac{\text{vol}(\partial s)}{\min(\text{vol}(s), \text{vol}(\bar{s}))} \quad (3.40)$$

where:

$$\text{vol}(s) := \sum_{(i,j) \in s \times \mathcal{V}} d\tau(i, j), \text{ and } \text{vol}(\partial s) := \sum_{(i,j) \in s \times \bar{s}} d\tau(i, j)$$

Thereupon let  $\mathcal{S}^*$  be the set of all  $s \subset S$  with  $\pi(s) \leq \frac{1}{2}$ , then the conductance  $\Phi$  of  $\tau$  is defined by:

$$\Phi := \min_{s \in \mathcal{S}^*} \Phi_s \quad (3.41)$$

For a Markov chain, the conductance  $\Phi_s$  of a cut  $(s, \bar{s})$  can be thought as the conditional probability, that in its stationary distribution the Markov chain will leave the region  $s$  within a single application of the transition operator. Consequently, if for some region  $s \subset S$  the conductance  $\Phi_s$  of the cut  $(s, \bar{s})$  is very small, then

there necessarily exists an eigenvalue  $\lambda_2$  of the transition operator  $\tau$ , such that the components in the corresponding eigenspace only very slow decay to zero. This, however, is only the case, if  $\lambda_2$  is close to 1 such that the transition operator  $\tau$  has a small spectral gap  $1 - \lambda_2$ . This consideration is formalized by Cheeger's Inequality for Markov chains.

**Theorem 3.3** (Cheeger's Inequality for Markov chains). *Let  $(\mathbf{X}^{(t)})_{t \in \mathbb{N}}$  be an irreducible aperiodic and reversible Markov chain with a state space  $(S, \Sigma)$ , a strictly positive stationary distribution  $\pi$  and a transition operator  $\tau$ . Furthermore let  $\Phi$  be the conductance and  $\lambda_2$  the second largest eigenvalue of  $\tau$ . Then the spectral gap  $1 - \lambda_2$  is bounded by the conductance:*

$$\frac{1}{2}\Phi^2 \leq 1 - \lambda_2 \leq 2\Phi \quad (3.42)$$

*Notes on the proof of Theorem 3.3.* The Theorem has independently been proven by (Lawler et al. 1988) and (Jerrum et al. 1989). In both cases the proof essentially generalizes Cheeger's inequality, known from differential geometry.  $\square$

Since the goal of these considerations is, to define an upper bound for the mixing time, which is associated with an intuitive meaning, the previous results can be combined as follows:

**Corollary 3.2.** *Let  $(\mathbf{X}^{(t)})_{t \in \mathbb{N}}$  be an irreducible aperiodic Markov chain with a state space  $(S, \Sigma)$ , a stationary distribution  $\pi$  and a transition operator  $\tau$ . Furthermore let  $\Phi$  be the conductance of  $\tau$ . Then for all  $\varepsilon > 0$  it holds, that:*

$$t_{\text{mix}}(\varepsilon) \leq \left\lceil \frac{2}{\Phi^2} \ln \frac{1}{\varepsilon \pi_{\min}} \right\rceil \quad (3.43)$$

*Notes on the proof of Corollary 3.2.* The Corollary has been proven by (Jerison 2013, p6). The proof essentially combines the upper bound of Theorem 3.3 with the lower bound of Proposition 3.2.  $\square$

The application of the upper bound for the mixing time as given by Corollary 3.2, however, for each class of underlying EBM requires a separate and usually sophisticated mathematical analysis, that relates the conductance with the respective

parametrisation (Koller et al. 2009, p520). Consequently this effort has only been carried out for specific classes of EBMs which prove themselves in practice. One of the first classes of that kind, which attracted a broad attention, were **Restricted Boltzmann Machines (RBM)**.

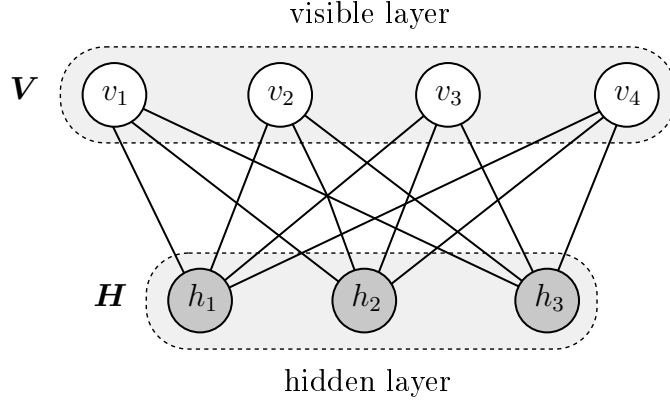


Figure 3.1: Graphical Model of a Restricted Boltzmann Machine

**Example 3.1** (Restricted Boltzmann Machine). *Soon after Boltzmann Machines have been proposed, with the “Harmonium” a more restrictive version was introduced as a design study for an artificial cognitive system (Smolensky 1986). The fundamental idea was, to distinguish a layer of observables from a layer of representational features, and to restrict the graph to edges between those two layers (see figure 3.1). Thereby for a **visible layer**  $\mathbf{V}$ , consisting of  $d$  observables and a **hidden layer**  $\mathbf{H}$ , consisting of  $m$  latent random variables the parametrisation comprises activation thresholds  $\mathbf{a} \in \mathbb{R}^d$  and  $\mathbf{b} \in \mathbb{R}^l$  and synaptic weights  $\mathbf{W} \in \mathbb{R}^{d \times l}$  between the layers, such that:*

$$E_{\mathbf{W}, \mathbf{a}, \mathbf{b}}(\mathbf{v}, \mathbf{h}) = -(\mathbf{v}^T \mathbf{W} \mathbf{h} + \mathbf{a}^T \mathbf{v} + \mathbf{b}^T \mathbf{h}) \quad (3.44)$$

*Then a Gibbs Sampler for the model distribution  $P_{\mathbf{W}, \mathbf{a}, \mathbf{b}}$  is determined by its conditional probabilities, which due to the bipartite graph structure for an observable  $V_i$  is given by:*

$$P_{\mathbf{W}, \mathbf{a}, \mathbf{b}}(V_i = 1 \mid \mathbf{H} = \mathbf{h}) \stackrel{2.25}{=} \text{sigm} \sum_{j=1}^l (w_{i,j} h_j + b_j) \quad (3.45)$$



and for a latent random variables  $H_j$  by:

$$P_{W,a,b}(H_j = 1 \mid \mathbf{V} = \mathbf{v}) \stackrel{2,25}{=} \text{sigm} \sum_{i=1}^d (w_{i,j}v_i + a_i) \quad (3.46)$$

Furthermore all observables  $V_i, V_j$  with  $i \neq j$  are pairwise conditionally independent, such that:

$$P_{W,a,b}(\mathbf{V} \mid \mathbf{H}) \stackrel{3,45}{=} \prod_{i=1}^d P_{W,a,b}(V_i \mid \mathbf{H}) \quad (3.47)$$

The same conclusion can be drawn for the latent random variables:

$$P_{W,a,b}(\mathbf{H} \mid \mathbf{V}) \stackrel{3,46}{=} \prod_{i=1}^l P_{W,a,b}(H_i \mid \mathbf{V}) \quad (3.48)$$

CD learning mainly owes its great popularity to its empirically successful, when applied to RBMs. The advantages compared to general BMs were particularly evident with respect to its mixing behaviour. For this purpose, however, it is necessary to reconstitute the sampling scheme in RBMs. For a partition of sets, which respectively comprise conditionally independent random variables, a Gibbs sampler can be formulated to parallelly update all random variables within a set “at once”. This provides a **Block Gibbs Sampler**.

**Definition** (Block Gibbs Sampler). *Let  $\mathbf{X}$  be an  $(S, \Sigma)$ -valued MRF over an undirected graph  $G = (\mathcal{V}, \mathcal{E})$  and  $(\mathcal{V}_i)_{i \in I}$  a partition of  $\mathcal{V}$ , such that for any  $i \in I$  the random variables in  $\mathcal{V}_i$  are pairwise conditionally independent. Then a measurable function  $T: (S, \Sigma) \rightarrow (S, \Sigma)$  is a block Gibbs sampler for a probability distribution  $P$  over  $(S, \Sigma)$ , iff for  $\mathbf{X}' := T \circ \mathbf{X}$  it holds, that:*

$$\mathbf{X}'_{\mathcal{V}_i} \sim P(\mathbf{X}_{\mathcal{V}_i} \mid \{\mathbf{X}_{\mathcal{V}_j}\}_{1 \leq j < i}, \{\mathbf{X}_{\mathcal{V}_j}\}_{i < j \leq |I|}), \forall i \in I \quad (3.49)$$

By applying the block Gibbs Sampler to RBMs the sampling scheme as used in CD learning then obtains a remarkably simple structure: The sampling starts by a Markov chain for the extended data distribution  $P_{\mathbf{v},\theta}(\mathbf{V}, \mathbf{H})$ . Thereby the observables are initially drawn from the marginal distribution  $P_{\mathbf{v},\theta}(\mathbf{V})$ , which by definition equals

the data distribution  $P_{\mathbf{v}}(\mathbf{V})$ , such that:

$$\mathbf{v}_{\text{d}}^{(0)} \sim P_{\mathbf{v}}$$

Thereupon, the latent random variables are not required to be drawn from a uniform distribution, since the conditional probability  $P_{\mathbf{v},\boldsymbol{\theta}}(\mathbf{H} \mid \mathbf{v}_{\text{d}}^{(0)})$ , which equals  $P_{\boldsymbol{\theta}}(\mathbf{H} \mid \mathbf{v}_{\text{d}}^{(0)})$ , is already completely known by equation 3.48, such that:

$$\mathbf{h}_{\text{d}}^{(0)} \sim P_{\boldsymbol{\theta}}(\mathbf{H} \mid \mathbf{v}_{\text{d}}^{(0)})$$

At this point, it is important to notice, that no further iterations in the Markov chain for  $P_{\mathbf{v},\boldsymbol{\theta}}$  are required, since the drawn sample  $(\mathbf{v}_{\text{d}}^{(0)}, \mathbf{h}_{\text{d}}^{(0)})$  already is an exact sample of  $P_{\mathbf{v},\boldsymbol{\theta}}$ , such that  $P_{\mathbf{v},\boldsymbol{\theta}} = P_{\mathbf{v},\boldsymbol{\theta}}^{(0)}$ . Thereupon the Markov chain for the model distribution  $P_{\boldsymbol{\theta}}$  is initialised by the last drawn sample for the data distribution, such that  $\mathbf{v}_{\text{m}}^{(0)} := \mathbf{v}_{\text{d}}^{(0)}$  and  $\mathbf{h}_{\text{m}}^{(0)} := \mathbf{h}_{\text{d}}^{(0)}$ . In any following iteration step the  $(i + 1)^{\text{th}}$  observables are updated from the  $i^{\text{th}}$  latent random variables:

$$\mathbf{v}_{\text{m}}^{(i+1)} \sim P_{\boldsymbol{\theta}}(\mathbf{V} \mid \mathbf{h}_{\text{m}}^{(i)}), \forall i \in \mathbb{N}$$

and then the  $(i + 1)^{\text{th}}$  latent random variables from the  $(i + 1)^{\text{th}}$  observables:

$$\mathbf{h}_{\text{m}}^{(i+1)} \sim P_{\boldsymbol{\theta}}(\mathbf{H} \mid \mathbf{v}_{\text{m}}^{(i+1)}), \forall i \in \mathbb{N}$$

Consequently when applied to RBMs, the sampling scheme of CD learning reduces to a single Markov chain, which alternates between the observables and latent random variables and therefore is referred as **alternating Gibbs sampling**. Thereby for any alternation, the respective random variables can be computed in parallel, which in particular for large models allows considerable time savings by the use of parallel computing. Nevertheless, there is a further and maybe even more important aspect, that justifies the success of CD learning, when applied to RBMs: In the first decade of the century, CD learning on RBMs has extensively been tested in numerous domains, where in many cases even for a single step approximation, the algorithm succeeded to converge in appropriate time. This empirical success of CD learning on RBMs received a theoretical explanation by evaluating the conductance for RBMs. Then

the weight matrix determines a tangible upper bound for the mixing time.

**Theorem 3.4** (Mixing time in RBMs). *Let  $(S, \Sigma, \mathcal{P})$  be an RBM with a visible layer  $\mathbf{V}$ , that comprises  $d$  observables and a hidden layer  $\mathbf{H}$ , consisting of  $l$  latent random variables, such that the model distribution  $P_{W,\mathbf{a},\mathbf{b}}$  is given by parameters  $\mathbf{a} \in \mathbb{R}^d$ ,  $\mathbf{b} \in \mathbb{R}^l$  and  $W \in \mathbb{R}^{d \times l}$ . Furthermore let:*

$$\|W\|_1 \|W^T\|_1 < 4 \quad (3.50)$$

*Then the mixing time  $t_{\text{mix}}$  of the alternating Gibbs Sampler for  $P_{W,\mathbf{a},\mathbf{b}}$  satisfies, that:*

$$t_{\text{mix}}(\varepsilon) \leq \frac{1}{\ln(4) - \ln(\|W\|_1 \|W^T\|_1)} \ln \left( \frac{\min(d, l)}{\varepsilon} \right) \quad (3.51)$$

*where the vectorial induced  $\ell_1$ -norm of the matrix is given by:*

$$\|W\|_1 := \max_{1 \leq j \leq l} \sum_{i=1}^d |W_{i,j}| \quad (3.52)$$

*Notes on the proof of Theorem 3.4.* The Theorem has been proven by (Tosh 2016, p4, Lemma 4). Thereby the proof utilizes Markovian couplings to provide an upper bound for the total variation distance and therefore to the mixing time (Tosh 2016, p3, Lemma 2). A subsequent application of the bounds to the individual Markov chains of (i) the observables and (ii) the latent random variables of an alternating Gibbs sampler then provides a generic upper bound for the mixing time (Tosh 2016, p4, Theorem 3). This upper bound, however, requires a restriction (Tosh 2016, p4) of the respective Markov chains which in the case of an RBMs is shown to be satisfied by the assumption of equation 3.50 (Tosh 2016, p4, Lemma 4) and (Tosh 2016, p18, Theorem 17) with reference to (Jerrum et al. 1993). This allows a derivation of the upper bound as given by equation 3.51 (Tosh 2016, p4, Corollary 5).  $\square$

Consequently for the case, that a the weight matrix of a given RBM satisfies the condition, given by equation 3.50, then the Markov chain of the alternative Gibbs sampling is polynomially bounded and therefore rapid mixing. In order to assure, that CD learning profits from rapid mixing, however, it has to be taken into account, that during each step of CD learning the weight matrix is updated. A possible approach

to enforce the rapid mixing property is therefore to regularize the  $\ell_1$ -norm within the update rule, such that the weights are bounded by the above criterion. At this point, however, it is important to notice, that an application of the  $\ell_1$ -regularization is only justified for the structural assumption of a respective sparsity, which of course is not satisfied for arbitrary observed data. Nevertheless, it can be concluded that RBMs in many cases allow tractable inference. Much before these theoretical results, however, the empirical success of RBMs already motivated their use as building blocks for more complicated “stacked” graphical models (Hinton et al. 2006). Thereby in particular the **Deep Boltzmann Machines (DBM)** (Salakhutdinov et al. 2009) may be regarded as a generalization of RBMs.

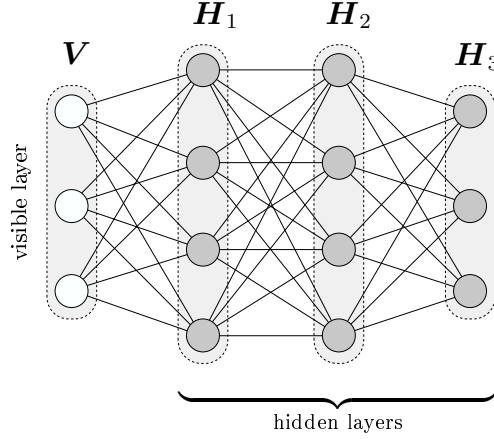


Figure 3.2: Graphical Model of a Deep Boltzmann Machine

**Example 3.2** (Deep Boltzmann Machine). *Deep Boltzmann Machines (DBM) generalize RBMs by decomposing the hidden component  $\mathbf{H}$  in  $L \in \mathbb{N}$  hidden layers  $\mathbf{H}_1, \mathbf{H}_2, \mathbf{H}_3 \dots \mathbf{H}_L$ , that form a hierarchical deep structure (see figure 3.2). Thereby the visible component  $\mathbf{V}$  and the first hidden layer  $\mathbf{H}_1$  are connected by weights  $W_1 \in \mathbb{R}^{d \times l_1}$  and any hidden layer  $\mathbf{H}_i$  with the subsequent layer  $\mathbf{H}_{i+1}$  by corresponding weights  $W_i \in \mathbb{R}^{l_i \times l_{i+1}}$ . With respective thresholds  $\mathbf{a} \in \mathbb{R}^d$  and  $\mathbf{b}_i \in \mathbb{R}^{l_i}$  the energy function of the DBM is therefore:*

$$E(\mathbf{v}, \mathbf{h}) = \underbrace{E(\mathbf{v}, \mathbf{h}_1)}_{\text{RBM}} + \sum_{i=1}^{L-1} E(\mathbf{h}_i, \mathbf{h}_{i+1}) \quad (3.53)$$

where the energy term  $E(\mathbf{v}, \mathbf{h}_1)$  equals that of an RBM:

$$E(\mathbf{v}, \mathbf{h}_1) := -(\mathbf{v}^\top W_1 \mathbf{h}_1 + \mathbf{a}^\top \mathbf{v} + \mathbf{b}_1^\top \mathbf{h}_1) \quad (3.54)$$

and the following terms respectively add contributions from additional connections and vertices of the  $i + 1$ th hidden layer:

$$E(\mathbf{h}_i, \mathbf{h}_{i+1}) := -(\mathbf{h}_i^\top W_{i+1} \mathbf{h}_{i+1} + \mathbf{b}_{i+1}^\top \mathbf{h}_{i+1}), \forall i \in \{1, \dots, L - 1\} \quad (3.55)$$

The structure of DBMs then forms a bipartite graph (see figure 3.3) which allows an application of an alternating Gibbs sampler.

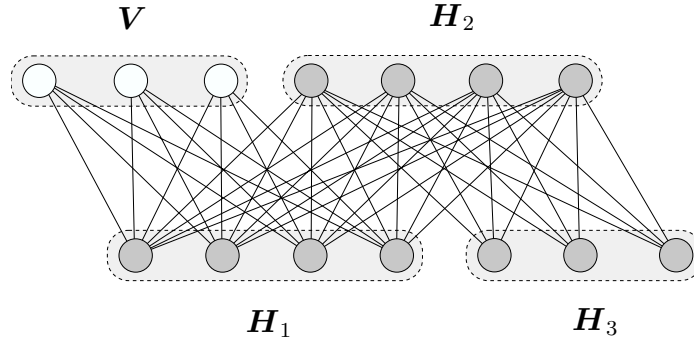


Figure 3.3: Bipartite graph structure of a Deep Boltzmann Machine

**Corollary 3.3** (Mixing time in DBMs). *Let  $(S, \Sigma, \mathcal{P})$  be an DBM with a visible layer  $\mathbf{V}$ , that comprises  $d$  observables, and  $l$  hidden layers  $\mathbf{H}_i, i \in \{2, \dots, L + 1\}$ , that respectively comprise  $l_i$  latent random variables. Then the model distribution  $P_{\boldsymbol{\theta}}$  is parametrised by (i) activation thresholds  $\mathbf{a} \in \mathbb{R}^d$  and  $\mathbf{b}_i \in \mathbb{R}^{l_i}$  and (ii) weight matrices  $W_1 \in \mathbb{R}^{d \times l_1}$  and  $W_i \in \mathbb{R}^{l_i \times l_{i+1}}$ . Let  $W$  be given by:*

$$W := \begin{pmatrix} W_1 & 0 & 0 & 0 \\ W_2^\top & W_3 & 0 & 0 \\ 0 & W_4^\top & \ddots & 0 \\ 0 & 0 & \ddots & \ddots \end{pmatrix}, \text{ with } \|W\|_1 \|W^\top\|_1 < 4$$

Then the mixing time  $t_{\text{mix}}$  of the alternating Gibbs Sampler for  $P_{\theta}$  satisfies, that:

$$t_{\text{mix}}(\varepsilon) \leq \frac{1}{\ln(4) - \ln(\|W\|_1 \|W^T\|_1)} \ln \left( \frac{\min(m_{\text{even}}, m_{\text{odd}})}{\varepsilon} \right) \quad (3.56)$$

where  $m_{\text{even}} := d + l_2 + l_4 + \dots$  is the total number of vertices in the even layers and  $m_{\text{odd}} := l_1 + l_3 + \dots$  the total number of vertices in the odd layers.

*Notes to the proof of Corollary 3.3.* The Corollary has been proven by (Tosh 2016, p5, Corollary 6). Thereby due to the bipartite graph structure (see figure 3.3), the Corollary follows as an immediate consequence of Theorem 3.4.  $\square$

## Chapter 4

# Application to Gene Regulation Analysis

### Overview

*Energy based models, as introduced in chapters 2 and 3, prove themselves to provide a suitable class of statistical models for the approximation of principal manifolds. Thereby, however, various structural assumptions had to be incorporated, which in the first section are evaluated in the domain of gene expression data. Afterwards biological structural knowledge is incorporated to define a deep structured EBM for the approximation of principal manifolds in gene expression data. This model is used in the second section to model a certain gene regulatory network upon gene expression profiles taken from Glioblastoma Multiforme (GBM) cells as well as of normal cells of identical type. Thereby the modelling comprises three steps: In the first step the assumed gene regulatory network is used to estimate parameters for the visible layer, the first hidden layer and their connecting weights. The estimated parameters are used as a prior to approximate a principal manifold, that by stratification equally considers GBM and normal cells. This joint model is then used as a prior to “finetune” the approximations in one case for GBM cells and in the other case for normal cells. Thereby the prior joint model is important to assure comparable results with respect to local ML estimates. Finally the differences in gene regulation are quantified by differences within the respective  $\mathcal{M}$ -Correlations. These results are compared to other*

*approaches.*

## 4.1 Statistical Modelling of Gene Regulatory Networks by EBMs

With respect to an application to gene regulation, maybe the foremost question that arises, is if an undirected probabilistic graphical model, that treats errors in variables equally is suitable. This question in particular appears by considering, that biochemical reactions as the building blocks of gene regulation, indeed have some considerate direction, which at certain conditions is given by the sign of the reaction rate. Nevertheless on the scale of gene regulation, the causality usually acts in both directions, to allow cells to preserve homoeostasis: Gene regulatory networks usually stabilize and adapt the production of proteins as an intracellular resource, or as a cellular response to changing environmental conditions, like the available nutrition. A well studied example for such an adaptive behaviour is the Lac operon model.

**Example 4.1** (Lac operon). *Glucose is an important energy source for many bacteria, like *E. coli*. If, however, glucose is not available within the cellular environment, then those bacteria also import the more complex lactose molecules. An increasing intracellular concentration of lactose in turn increases the probability for lactose repressor genes to be transcribed from the bacterial DNA, which, ceteris paribus, causes an increase in the concentration of lactose repressors. Thereupon an increasing concentration of lactose repressors, increases the probability for the degradation of lactose to glucose, such that, ceteris paribus, the concentration of lactose decreases and the concentration of glucose increases. Both mechanisms taken together, can be summarized by an overall effect, in which intracellular lactose induces its own degradation to glucose.*

The Lac operon gives an example for a causal relationship between concentration levels, i.e. lactose and lactose repressors, which has no designated direction. Although this example may not be regarded as universal, it at least motivates to represent the concentration levels of individual genes by an undirected graph, which is the structural assumption (S1). Let therefore  $\mathbf{V}$  denote observable concentration levels of genes in



a gene regulatory network. Thereupon the regulatory network may also involve interactions by further unobserved genes, which concentration levels are comprised by a latent random vector  $\mathbf{H}$ . Consequently the assumption that  $\mathbf{X} = (\mathbf{V}, \mathbf{H})$  covers all gene regulatory interactions of a given regulatory network, equals the structural assumption **(S2)**, that any  $X_i$  satisfies the local Markov property in an undirected graph  $\mathcal{G}$ . Let now  $\mathbf{V}$  be realized by a finite number of observations. Then for any outcome  $\mathbf{v}$ , which is NOT realised, it can not be distinguished, if the marginal probability  $P(\mathbf{V} = \mathbf{v})$  is vanishing or not. As this consideration virtually extends to the latent random vector, without loss of generality, it can be assumed, that  $P(\mathbf{X} = \mathbf{x}) > 0$  for all outcomes  $\mathbf{x}$ , which is the structural assumption **(S3)**. This shows, that the choice of EBMs for the statistical modelling of gene expression is considerate arguable. Thereupon the next step is to implement such a statistical model.

For this purpose a good starting point is to substantiate the random variables  $X_i$ . Gene expression data appears in multifarious variants due to their underlying screening technologies. Thereby in particular high-throughput technologies like DNA Microarray or Next Generation Sequencing allow a statistical treatment by quantitatively comparable variances. In the following the focus is given to gene expression profiles, given by cDNA microarrays. These microarrays consist of thousands of DNA sequences, referred as probes, that are printed in a high density array on a glass slide. Thereupon the sample, that is to be quantified, as well as a standardized reference sample, respectively are reverse-transcribed into cDNA and labelled with red and green fluorescent dyes. These target samples are jointly hybridized with the arrayed probes on the glass slide, which afterwards is discharged from the unbound residues. Finally a laser beam of defined wavelength is used to locally determine the intensities  $I_R(i)$  and  $I_G(i)$  of each spot  $i$  in the array. This allows to use the reference intensity  $I_G(i)$  to quantify the relative abundance of the respective DNA sequence in the sample by the ratio:

$$R_i := \frac{I_R(i)}{I_G(i)} \quad (4.1)$$

Thereupon a normal distributed additive Gaussian error is obtained by applying a logarithm to  $R_i$ . Thereby the most common used logarithm is given by the basis 2.

The cDNA **log-ratios** are then given by:

$$X_i := \log_2 R_i = \log_2 I_R(i) - \log_2 I_G(i) \quad (4.2)$$

By its definition it immediately follows, that the standardized reference sample adds an offset to the log-ratio  $X_i$ . Thereby it is common to choose this offset, with respect to the resolution of experiment, such that negative log-ratios may be regarded as measurement errors (see figure 4.4). The logarithm transformation allows to decompose any log-ratio  $X_i$  into a systematic component  $X_i^*$  which is determined by interactions and normal distributed random component, such that:

$$X_i = X_i^* + \varepsilon_i, \text{ with } \varepsilon_i \sim \mathcal{N}(0, \sigma_i^2) \quad (4.3)$$

The next step is therefore to describe the systematic component  $X_i^*$  by **direct causal interactions**. By equation ... and ... an error free ratio can be defined by  $R_i^* := R_i \exp(-\varepsilon_i)$ , such that:

$$X_i^* = \log_2(R_i^*) \quad (4.4)$$

Let now  $R_i^*$  represent the relative abundance of a gene  $G_i$ , then  $R_i^*$  by definition is proportional to the concentration level  $[G_i]$ . It therefore can be concluded that  $\exists \lambda_i > 0$  with  $R_i^* = \lambda_i [G_i]$  and by  $c_i := \log_2(\lambda_i)$ :

$$X_i^* \stackrel{4.4}{=} \log_2(\lambda_i [G_i]) = \log_2([G_i]) + c_i \quad (4.5)$$

Gene regulatory interactions occur in multifarious manifestations, which in their entirety can only hardly be considered in a tractable family of functions. It is therefore required to incorporate simplifying assumptions. A popular approach to provide a commensurate simplification is to describe direct causal interactions by **Hill-type interactions**. Thereby the concentration levels of a ligand L and an enzymes E under the under the assumption of cooperative binding is parametrised by a microscopic *dissociation constant*  $K_A \in \mathbb{R}_0^+$  and a *Hill coefficient*  $n \in \mathbb{N}$ , such that:

$$\ln([E]) = \frac{1}{1 + \exp(nK_A - n \ln([L]))} \quad (4.6)$$

If it is therefore assumed, that all direct gene interactions may be described by Hill-type interactions, it can be concluded, that for interacting genes  $G_i$  and  $G_j$  with respective log-ratios  $X_i^*$  and  $X_j^*$  there exist  $\alpha, \beta \in \mathbb{R}$ , such that:

$$X_i^* = \frac{1}{1 + \exp(\alpha + \beta X_j^*)} = \text{sigm}(\alpha + \beta X_j^*) \quad (4.7)$$

The relationship given by equation 4.7 shall now be used to describe conditional distributions in the EBM. For observables  $V_i, V_j$  this can be obtained by **conditioning** over an additional latent random variables  $H$ :

**Definition** (Conditioning). *Let  $\mathbf{V}$  be an MRF over an undirected graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ . Then for adjacent vertices  $i, j \in \mathcal{V}$  the corresponding random variables  $V_i, V_j$  are said to be conditioned on a random variable  $H$ , iff their joint probability  $P$  satisfies:*

$$P(V_i, V_j | H) = P(V_i | H)P(V_j | H) \quad (4.8)$$

Since conversely the marginal probability of the observables  $P(\mathbf{V})$  is given by a marginalisation of  $P$  over  $\mathbf{H}$ , conditioning may be regarded as an inverse operation to marginalisation. In the following, conditioning is used to locally derive a representation for Hill-type interactions by a n EBM with appropriate feature functions.

**Lemma 4.1.** *Let  $G_i$  and  $G_j$  denote interacting genes and  $V_i$  and  $V_j$  their respective log-ratios. Then  $V_i$  and  $V_j$  can be conditioned on a latent random variable  $H$ , such that  $P(V_i, | V_j)$  and  $P(V_j, | V_i)$  represent Hill-type interactions.*

*Proof of Lemma 4.1.* The proof is constructive: Let  $H$  be a binary random variable, such that the joint probability  $P(V_i, V_j, H)$  satisfies: (i)  $V_i$  is Gaussian distributed, such that for some  $w_i, a_i \in \mathbb{R}$ :

$$V_i \sim \mathcal{N}(V_i^*, \sigma_i^2), \text{ with } V_i^* = w_i H + a_i$$

(ii)  $V_j$  is Gaussian distributed, such that for some  $w_j, a_j \in \mathbb{R}$ :

$$V_j \sim \mathcal{N}(V_j^*, \sigma_j^2), \text{ with } V_j^* = w_j H + a_j$$

(iii)  $H$  is Bernoulli distributed, such that for some  $b \in \mathbb{R}$ :

$$H \sim \mathcal{B} \left( \text{sigm} \left( \frac{w_i}{\sigma_i} V_i + \frac{w_j}{\sigma_j} V_j + b \right) \right)$$

Then  $P(V_i, | V_j)$  and  $P(V_j, | V_i)$  represent Hill-type interactions. Furthermore it holds, that  $P(V_i | H, V_j) = P(V_i | H)$  and  $P(V_j | H, V_i) = P(V_j | H)$  and therefore  $P(V_i, V_j | H) = P(V_i | H)P(V_j | H)$ .  $\square$

The modelling of individual Hill-type interactions, by conditioning on latent random variables, can immediately be extended to gene regulatory networks (GRN). Thereby the model, could follow the design principle to respectively assign every edge of the GRN with a single latent variable, that represents the respective type of interaction by an appropriate feature function. In this way, the believed dependency structures can be applied to the observation, by a respective Maximum likelihood estimation.

**Corollary 4.1.** *Let the random vector  $\mathbf{V}$  denote observed log-ratios of a gene regulatory network, where all interactions are given by Hill-type interactions. Then there exists a canonical EBM, that models  $\mathbf{V}$  by conditioning (see figure 4.1).*

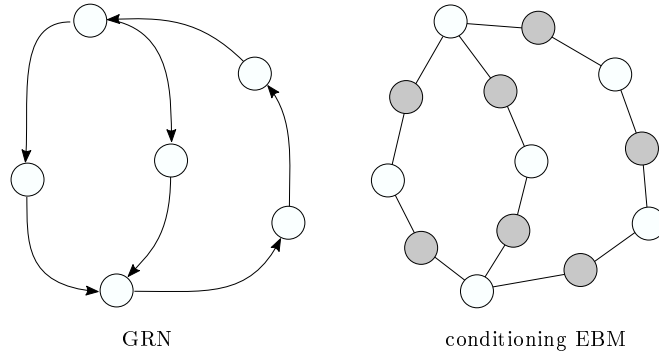


Figure 4.1: Modelling Gene Regulatory Networks (GRN) by conditioning EBMs

*Proof of Corollary 4.1.* For any edge  $(i, j) \in \mathcal{E}$  of the GRN let  $V_i$  and  $V_j$  be conditioned on a random variable  $H_{ij}$  as introduced in lemma 4.1. Then the joint probability can be represented by a Boltzmann distribution, where the energy function for

any outcome  $(\mathbf{v}, \mathbf{h})$  is given by:

$$E(\mathbf{v}, \mathbf{h}) = - \left( \sum_{i=1}^d \frac{(v_i - a_i)^2}{2\sigma_i^2} + \sum_{i=1}^d \sum_{(i,j) \in \mathcal{E}} w_{ij} \frac{v_i}{\sigma_i} h_{ij} + \sum_{(i,j) \in \mathcal{E}} h_{ij} b_{ij} \right) \quad (4.9)$$

□

It is important to notice, that the graph layout of an EBM that models random vector  $\mathbf{V}$  is not an intrinsic property of the observed random vector, but a matter of choice with respect to additional latent random variables, that are introduced within the model, to address structural uncertainties. In the context of traditional graphical models latent variables are typically used with some intended semantic, that models the dependency structure between observable variables (see figure 4.1). Nevertheless, this approach only accounts for uncertainties within believed dependencies, but not within the connection structure itself. A possible approach to overcome this restriction is to relax the connection structure of a conditioning EBM by a deep structured EBM. In this sense, the local ML estimate, which is obtained from the conditioning EBM is used as a prior for the deep structured EBM. This approach, however, requires a natural embedding in the parameter space. With respect to the conditioning of Hill-type interactions, the observables are represented by Gaussian distributions and the latent random variables by Bernoulli distributions. The parametrisation then can naturally be embedded in a Gauss-Bernoulli DBM (see figure 4.2).

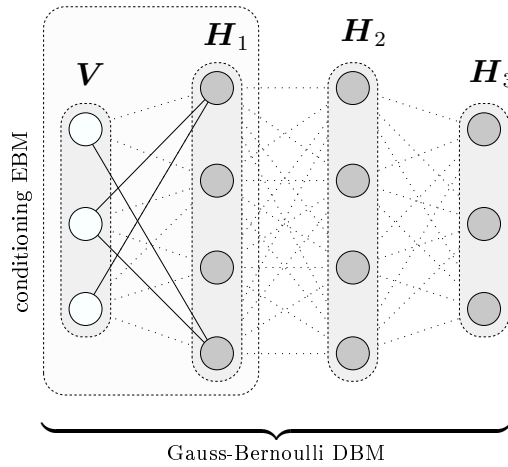


Figure 4.2: Initialise Gauss-Bernoulli DBMs by conditioning EBMs

**Example 4.2** (Gauss-Bernoulli Deep Boltzmann Machine). *By identifying the observables of a DBM with Gaussian distributions and the latent random variables with Bernoulli distribution the resulting deep structured EBM is referred as a Gauss-Bernoulli DBM (GBDBM). Thereby the observables  $\mathbf{V} = (V_1, V_2, \dots, V_l)$  are assumed to contain Gauss distributed errors  $\boldsymbol{\varepsilon} = (\varepsilon_1, \varepsilon_2, \dots, \varepsilon_l)$  with:*

$$\varepsilon_i \sim \mathcal{N}(0, \sigma_i^2)$$

*Thereupon  $\mathbf{V}$  and the first hidden layer  $\mathbf{H}_1$  are connected by weights  $W_1 \in \mathbb{R}^{d \times l_1}$  and any hidden layer  $\mathbf{H}_i$  with the subsequent layer  $\mathbf{H}_{i+1}$  by corresponding weights  $W_i \in \mathbb{R}^{l_i \times l_{i+1}}$ . With respective thresholds  $\mathbf{a} \in \mathbb{R}^d$  and  $\mathbf{b}_i \in \mathbb{R}^{l_i}$  the energy function of the GBDBM is then given by:*

$$E(\mathbf{v}, \mathbf{h}) = E(\mathbf{v}, \mathbf{h}_1) + \sum_{i=1}^{L-1} E(\mathbf{h}_i, \mathbf{h}_{i+1}) \quad (4.10)$$

*where the energy term  $E(\mathbf{v}, \mathbf{h}_1)$  equals:*

$$E(\mathbf{v}, \mathbf{h}_1) := - \left( \sum_{i=1}^d \frac{(v_i - a_i)^2}{2\sigma_i^2} + \sum_{i=1}^d \sum_{j=1}^{l_1} W_{1ij} \frac{v_i}{\sigma_i^2} h_{1j} + \sum_{j=1}^{l_1} b_{1j} h_{1j} \right) \quad (4.11)$$

*and the following terms respectively add contributions from additional connections and vertices of the  $i + 1$ th hidden layer:*

$$E(\mathbf{h}_i, \mathbf{h}_{i+1}) := -(\mathbf{h}_i^T W_{i+1} \mathbf{h}_{i+1} + \mathbf{b}_{i+1}^T \mathbf{h}_{i+1}), \forall i \in \{1, \dots, L-1\} \quad (4.12)$$

*Consequently the conditional distribution of  $V_i$  is given by:*

$$P(V_i | \mathbf{H}_1) = \mathcal{N} \left( \sum_{j=1}^{l_1} w_{ij} H_{1j} + a_i, \sigma_i^2 \right)$$

*The latent random variables  $\mathbf{H} = (\mathbf{H}_1, \mathbf{H}_2, \dots, \mathbf{H}_L)$  with  $\mathbf{H}_k = (H_{k1}, H_{k2}, \dots, H_{kl_k})$  are assumed to be Bernoulli distributed over a sigmoid activation function. Then the conditional distribution of  $H_{1j}$  is given by:*

$$\begin{aligned}
P(H_{1j} \mid \mathbf{V}, \mathbf{H}_2) &= \mathcal{B} \left( \text{sigm} \left( \sum_{i=1}^d W_{1ij} \frac{V_i}{\sigma_i^2} + \sum_{k=1}^{l_2} W_{2jk} H_k^{(2)} + b_j^{(1)} \right) \right) \\
P(H_{kj} \mid \mathbf{H}_{k-1}, \mathbf{H}_{k+1}) &= \mathcal{B} \left( \text{sigm} \left( \sum_{i=1}^{l_{k-1}} W_{kij} H_{(k-1)i} + \sum_{i=1}^{l_{k+1}} W_{(k+1)ji} H_{(k+1)i} \right) \right) \\
&\text{for } k \in \{2, \dots, L-1\} \\
P(H_{Lj} \mid \mathbf{H}_{L-1}) &= \mathcal{B} \left( \text{sigm} \left( \sum_{i=1}^{l_{L-1}} W_{(L-1)ij} H_{(L-1)i} \right) \right)
\end{aligned}$$

In contrast to traditional graphical models, deep learning approaches like GBDBMs usually do not prescribe this connection structure, but allow the model to adapt itself to the underlying observations. Thereby complicated non-linear dependencies between the observables are accomplished by indirect connections, that flow through multiply stacked latent variables. Accordingly it is all the more important to incorporate network characteristics, that restrict the model space to reasonable structured models. A particular characteristic of gene regulatory networks is a very small connectivity between 1.5 and 2.5 (Leclerc 2008, p1). This structural knowledge can be incorporated by enforcing this connectivity by an  $\ell_1$ -regularization in the update rule of the CD learning algorithm. This approach, also has positive effects on the tractability in of the local ML estimates: The weight matrix of a GBDBM is comprehensively given by:

$$W := \begin{pmatrix} W_1 & 0 & 0 & 0 \\ W_2^T & W_3 & 0 & 0 \\ 0 & W_4^T & \ddots & 0 \\ 0 & 0 & \ddots & \ddots \end{pmatrix}$$

Thereupon a respective  $\ell_1$ -regularization assures that  $\|W\|_1 \|W^T\|_1 < 4$ . By Corollary 3.3 it then can be concluded, that the Markov Chains in the MCMC approximation is rapid mixing.

## 4.2 Gene Regulation Analysis of Glioblastoma Multiforme

Glioblastoma multiforme is the most common malignant brain tumor in adults. At the scale of gene regulation Glioblastoma is characterised by gene deletions which in most cases affect the tumor suppressor gene TP53, the retinoblastoma suppressor gene RB-1 and further deletions in chromosome 22 as well as chromosome 10. Thereby the genetic damages often occur in combination. In newly developed primary glioblastomas more frequent losses of the PTEN gene occur by amplification of the EGFR gene (Ohgaki et al. 2004). In the secondary glioblastomas mutations of the TP53 gene often occur. In addition, point mutations in IDH1 and IDH2 genes coding for an isocitrate dehydrogenase are more common in this group, particularly the R132H mutation in the IDH1 gene (Watanabe et al. 2009). A further potential key player in the is given by Glycosaminoglycans (GAGs) (Afratis et al. 2012). Thereby GAGs play an intricate role in the extracellular matrix by specific interactions with growth factors and other transient components. The accumulated evidence regarding an altered structure of GAG in Glioblastoma indicates their importance disease progression. It may therefore be assumed, that the gene regulation of the GAG degradation pathway is altered by Glioblastoma. To investigate this question from the KEGG database a gene regulatory network has been derived (see figure 4.3).

Thereby the GAG gene regulatory network (see figure 4.3) comprises a large list of further candidates, that are assumed to indirectly interact with the GAG degradation pathway. (For convenience these candidates are represented by dots in the figure). Thereupon the challenge is to analyse the gene regulation with respect to cDNA log-ratios of a GBM dataset (see figure 4.4). This sample decomposes into three classes (i) “normal”, comprising a small subgroup of 10 patients, which are diagnosed with no GBM, (ii) “wildtype”, with a subgroup of 487 patients, that are diagnosed with GBM and finally (iii) “mutated” with a subgroup of 19 patients, that are diagnosed with GBM, where the TP53 gene occurs to be mutated.

In a first step the set of candidates has been reduced to genes with high correlation to the GAG degradation pathway. These have been identified by strong low dimensional  $L$ -correlations (see figure 4.5).



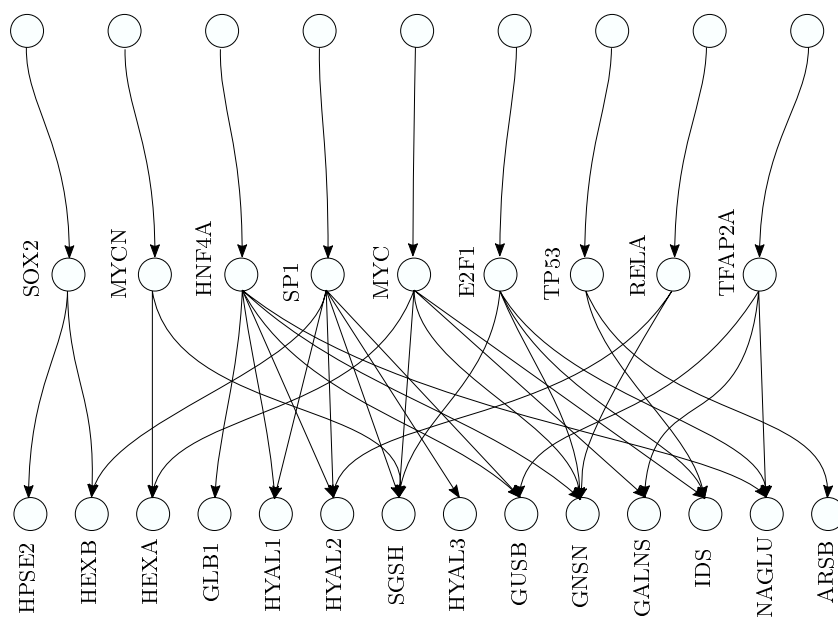


Figure 4.3: Gene Regulatory Network of GAG degradation pathway

Afterwards the GRN was reduced to these respective high correlating genes and embedded within a conditioning EBM (see figure 4.1) and a local ML estimate was derived for all samples by CD learning. The respective posterior in turn was used to initialise a GBDBM (see figure 4.2) which for reasons of tractability has been defined with only three hidden layers, where any subsequent hidden layer only comprised the half number of vertices than its predecessor. Thereupon a local ML estimate for this GBDBM was derived for all samples by CD learning. In order to calculate the  $\mathcal{M}$ -correlation upon the approximated local ML estimate, the reliabilities  $R_i$  have been estimated by the estimated variations  $\hat{\sigma}_i^2$ , such that for the finite realization  $\mathbf{v}$  of  $\mathbf{V}$ , the estimated reliability of  $X_i$  is given by:

$$\hat{R}_i = 1 - \frac{\hat{\sigma}_i^2}{\text{Var}_i(\mathbf{v})}$$

In order to approximate the integral of  $\mathcal{M}$ , a MCMC approximation has been used. Thereby a multiplicity of Markov Chains have been initialized by the empirical distribution and repeatedly pushed forward by an alternating Gibbs sampler. The gen-

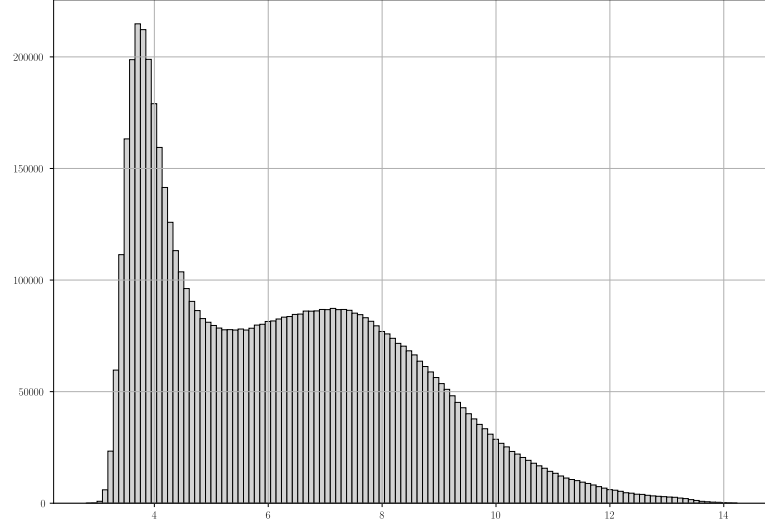


Figure 4.4: Histogram of cDNA log-ratios of GBM dataset

erated samples, that approximate the model distribution afterwards have been used to sample  $\mathbf{v}_{\mathcal{M}}$  from  $\mathcal{M}$  by deterministic sampling of the observables. Then by the law of large numbers it follows, that:

$$\frac{1}{m} \sum_{i=1}^m S_{i,j}(\mathbf{v}_{\mathcal{M}}^i) S_{j,i}(\mathbf{v}_{\mathcal{M}}^i) \xrightarrow{a.s.} \int_{\mathcal{M}} S_{i,j}(\mathbf{v}) S_{j,i}(\mathbf{v}) dP_{\mathcal{M}}, \text{ for } m \rightarrow \infty$$

Furthermore the local sensitivities  $S_{i,j}$  have been approximated by the “bruteforce” approximation:

$$\hat{S}_{i,j} = \partial_j \mathbb{E}(V_i \mid \mathbb{E}(\mathbf{H} \mid V_j))$$

This allowed an approximation of the  $\mathcal{M}$ -Correlation  $\rho_{V_i, V_j | \mathcal{M}}^2$ . The result has been used to define a graph structure, between the genes where strong  $\mathcal{M}$ -Correlations are represented by short distances of the respective vertices (see figure 4.6).

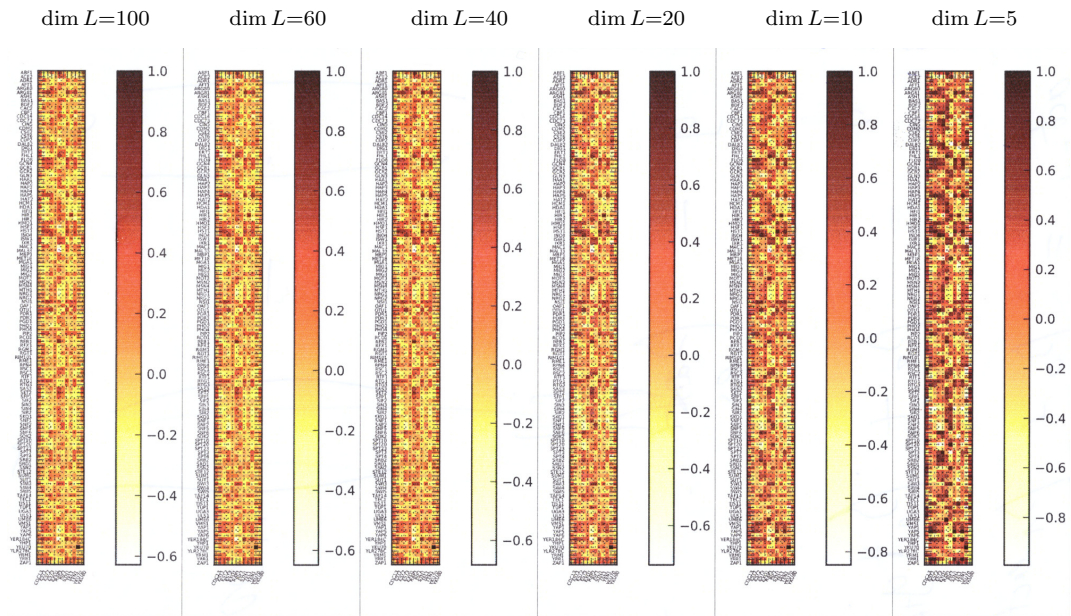


Figure 4.5:  $L$ -Correlation for selected genes in GBM dataset

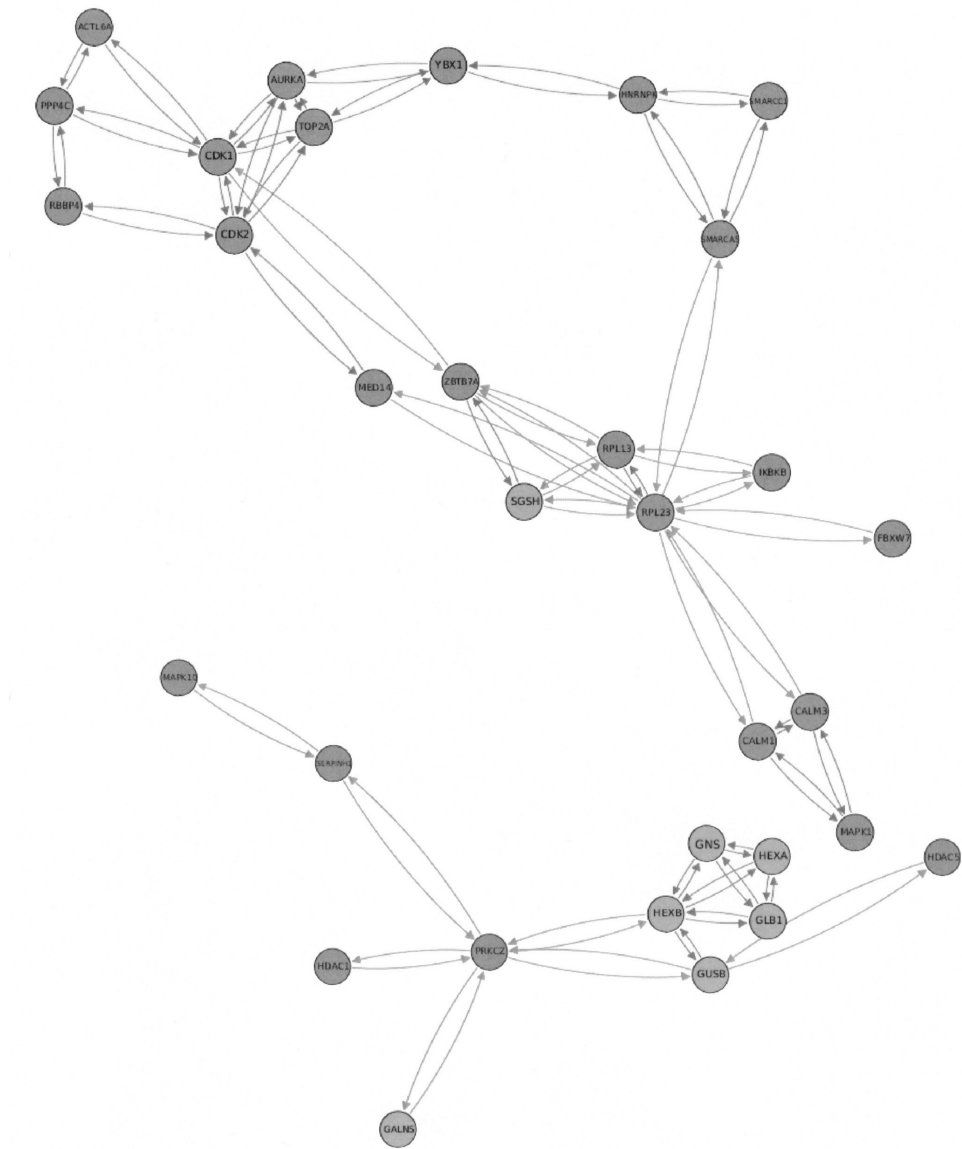


Figure 4.6: Predicted Gene Regulation of GAG degradation pathway in GBM by  $\mathcal{M}$ -Correlations

# Discussion

During the last decade deep structured graphical models motivated many different applications, ranging from classification, regression and dimensionality reduction tasks. The  $\mathcal{M}$ -Correlation shows, that the scope of their applicability also comprises the quantification of association structures. A great obstacle within the thesis, however, was that an application of the  $\mathcal{M}$ -Correlation to gene expression data can only hardly be evaluated with respect to its performance. Accordingly a commensurate quantification and comparison to other methods like mutual information based estimates is still pending. Since the opportunities in this direction are multifarious, further research in this area which considers differential and topological proximity structures is required.

# References

- [1] Nikos Afratis, Chrisostomi Gialeli, Dragana Nikitovic, Theodore Tsegenidis, Evgenia Karousou, Achilleas D. Theocharis, Mauro S. Pavão, George N. Tzanakakis, and Nikos K. Karamanos. Glycosaminoglycans: Key players in cancer cell biology and treatment. *FEBS Journal*, 279(7):1177–1197, 2012.
- [2] D J Aldous. Some inequalities for reversible Markov chains. *Journal of the London Mathematical Society*, 25(2):564–576, 1982.
- [3] Shun-Ichi Amari. *Information Geometry and Its Applications*. Springer, Tokyo, 1 edition, 2016.
- [4] Christof Angermueller, Tanel Pärnamaa, Leopold Parts, and Oliver Stegle. Deep learning for computational biology. *Molecular Systems Biology*, 12(7):878, 2016.
- [5] Leon Bottou. Online Learning and Stochastic Approximations. pages 1–34, 1998.
- [6] George Casella and Edward I. George. Explaining the Gibbs Sampler. *The American Statistician*, 46(3):167–174, 1992.
- [7] Yifei Chen, Yi Li, Rajiv Narayan, Aravind Subramanian, and Xiaohui Xie. Gene expression inference with deep learning. *Bioinformatics*, 32(12):1832–1839, 2016.
- [8] Guillaume Desjardins, Aaron Courville, Yoshua Bengio, Pascal Vincent, and Olivier Delalleau. Parallel Tempering for Training of Restricted Boltzmann Machines. *International Conference on Artificial Intelligence and Statistics*, 9:145–152, 2010.
- [9] Persi Diaconis and Daniel Stroock. Geometric Bounds for Eigenvalues of Markov Chains. *The Annals of Applied Probability*, 1(1):36, 1991.
- [10] Scott E Fahlman, Scott E Fahlman, Geoffrey E Hinton, Geoffrey E Hinton, Terrence J Sejnowski, and Terrence J Sejnowski. Massively Parallel Architectures for AI: NETL, THISTLE and Boltzmann Machines. *Proceedings of the 8th International Joint Conference on Artificial Intelligence*, pages 109–113, 1983.

- [11] Alan E. Gelfand and Adrian F M Smith. Sampling-based approaches to calculating marginal densities. *Journal of the American Statistical Association*, 85(410):398–409, 1990.
- [12] Stuart Geman and Donald Geman. Stochastic Relaxation, Gibbs Distributions, and the Bayesian Restoration of Images. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, PAMI-6(6):721–741, 1984.
- [13] a N Gorban and a Y Zinovyev. Principal Graphs and Manifolds. *Arxiv preprint arXiv08090490*, page 36, 2008.
- [14] Alexander N Gorban and Andrei Y Zinovyev. Elastic Maps and Nets and Their Application to Microarray Data Visualization. pages 1–35.
- [15] Trevor Hastie and Werner Stuetzle. Principal Curves. *Journal of the American Statistical Association*, 84(406):502–516, 1989.
- [16] G. Hinton and R. Salakhutdinov. Reducing the Dimensionality of Data with Neural Networks. *Science*, 313(July):504–507, 2006.
- [17] Geoffrey E Hinton. Training products of experts by minimizing contrastive divergence. *Neural Computation*, 14(8):1771–1800, 2002.
- [18] Daniel Jerison. General mixing time bounds for finite Markov chains via the absolute spectral gap. page 20, 2013.
- [19] Marc Jerrum and Alistair Sinclair. Approximating the Permanent. *SIAM Journal on Computing*, 18(6):1149–1178, 1989.
- [20] Mark Jerrum and Alistair Sinclair. Polynomial-time approximation algorithms for the ising model. *SIAM Journal on Computing*, 22(5):1087–1116, 1993.
- [21] JF Kenney and ES Keeping. Linear regression and correlation. *Mathematics of statistics*, 1:252–285, 1962.
- [22] Jamie King. Conductance and Rapidly Mixing Markov Chains. 2003.
- [23] Daphne Koller and Nir Friedman. *Probabilistic Graphical Models*. The MIT Press, 2009.
- [24] H. J. Kushner and D. S. Clark. Stochastic Approximation for Constrained and Unconstrained Systems. *Applied Math. Sci.*, 26, 1978.
- [25] Gregory F. Lawler and Alan D. Sokal. Bounds on the L2 Spectrum for Markov Chains and Markov Processes : A Generalization of Cheeger ’ s Inequality. *Transactions of the American Mathematical Society*, 309(2):557–580, 1988.

- [26] Robert D Leclerc. Survival of the sparsest: robust gene networks are parsimonious. *Molecular Systems Biology*, 4(213):2004–2009, 2008.
- [27] David a. Levin, Yuval Peres, and Elizabeth L. Wilmer. *Markov Chains and Mixing Times*. 2009.
- [28] Xianghang Liu and Justin Domke. Projecting Markov Random Field Parameters for Fast Mixing. *Neural Information Processing Systems*, pages 1–9, 2014.
- [29] L. Ljung and T. Söderström. *Theory and Practice of recursive identification*. MIT Press, Cambridge, 1983.
- [30] H. Ohgaki and P. Dessen. Genetic pathways to glioblastoma: a population-based study. *Cancer Research*, 64(19):6892–6899, 2004.
- [31] Karl Pearson. On lines and planes of closest fit to systems of points in space. *The London, Edinburgh, and Dublin Philosophical Magazine and Journal of Science*, 2(1):559–572, 1901.
- [32] Ruslan Salakhutdinov and Geoffrey Hinton. Deep Boltzmann Machines. *Proceedings of the 12th International Conference on Artificial Intelligence and Statistics*, 1(3):448–455, 2009.
- [33] Paul Smolensky. Information processing in dynamical systems: Foundations of harmony theory. In *Parallel Distributed Processing: Explorations in the Microstructure of Cognition, Volume 1: Foundations*, pages 194–281. MIT Press, 1986.
- [34] G. W. Snedecor and W. G. Cochran. *Statistical Methods*. Iowa State University Press, Ames, 6 edition, 1967.
- [35] Arida Ferti Syafiandini, Ito Wasito, Setiadi Yazid, Aries Fitriawan, and Mukhlis Amien. Multimodal Deep Boltzmann Machines for Feature Selection on Gene Expression Data.
- [36] Tijmen Tieleman. Training Restricted Boltzmann Machines using Approximations to the Likelihood Gradient. *Proceedings of the 25th International Conference on Machine Learning*, 307:7, 2008.
- [37] Tijmen Tieleman and Geoffrey Hinton. Using Fast Weights to Improve Persistent Contrastive Divergence. *Proceedings of the 26th Annual International Conference on Machine Learning (ICML)*, (8):1033–1040, 2009.
- [38] Christopher Tosh. Mixing Rates for the Alternating Gibbs Sampler over Restricted Boltzmann Machines and Friends. *International Conference on Machine Learning*, pages 1–7, 2016.
- [39] Martin J Wainwright and Michael I Jordan. Graphical Models, Exponential Families, and Variational Inference. *Foundations and Trends in Machine Learning*, 1(1-2):1–305, 2007.



- [40] T. Watanabe and S. Nobusawa. IDH1 mutations are early events in the development of astrocytomas and oligodendrogliomas. *The American Journal of Pathology*, 174(4):1149–1153, 2009.
- [41] Tung-Yu Wu, Bai Jiang, Yifan Jin, and Wing H. Wong. Convergence of Contrastive Divergence Algorithm in Exponential Family. pages 1–26, 2016.
- [42] Shurong Zheng, Ning-zhong Shi, and Zhengjun Zhang. Generalized Measures of Correlation. *Manuscript*, pages 1–45, 2010.
- [43] Fabio Zucca. On Some Properties of Transition Operators. *Extracta Mathematicae*, 17(2):201–209, 2002.

## **Erklärung**

Hiermit erkläre ich, dass ich meine Arbeit selbstständig unter Anleitung verfasst habe, dass ich keine anderen als die angegebenen Quellen und Hilfsmittel benutzt habe, und dass ich alle Stellen, die dem Wortlaut oder dem Sinne nach anderen Werken entlehnt sind, durch die Angabe der Quellen als Entlehnungen kenntlich gemacht habe.

Heidelberg, den 25. September 2017

Patrick Michl