

---

# Rectified Factor Networks

## Supplement

---

Djork-Arné Clevert, Andreas Mayr, Thomas Unterthiner and Sepp Hochreiter  
Institute of Bioinformatics, Johannes Kepler University, Linz, Austria  
{okko,mayr,unterthiner,hochreit}@bioinf.jku.at

### Contents

<b>1</b>	<b>Introduction</b>	<b>5</b>
<b>2</b>	<b>Rectified Factor Network (RFN) Algorithms</b>	<b>5</b>
<b>3</b>	<b>Convergence Proof for the RFN Learning Algorithm</b>	<b>9</b>
<b>4</b>	<b>Correctness Proofs for the RFN Learning Algorithms</b>	<b>11</b>
4.1	Diagonal Noise Covariance Update . . . . .	12
4.2	Full Noise Covariance Update . . . . .	14
<b>5</b>	<b>Maximum Likelihood Factor Analysis</b>	<b>15</b>
<b>6</b>	<b>The RFN Objective</b>	<b>18</b>
<b>7</b>	<b>Generalized Alternating Minimization</b>	<b>20</b>
<b>8</b>	<b>Gradient-based M-step</b>	<b>22</b>
8.1	Gradient Ascent . . . . .	22
8.2	Newton Update . . . . .	23
8.2.1	Newton Update of the Loading Matrix . . . . .	24
8.2.2	Newton Update of the Noise Covariance . . . . .	24
<b>9</b>	<b>Gradient-based E-Step</b>	<b>26</b>
9.1	Motivation for Rectifying and Normalization Constraints . . . . .	26
9.2	The Full E-step Objective . . . . .	27
9.3	E-step for Mean with Rectifying Constraints . . . . .	28
9.3.1	The E-Step Minimization Problem . . . . .	28
9.3.2	The Projection onto the Feasible Set . . . . .	29
9.4	E-step for Mean with Rectifying and Normalizing Constraints . . . . .	30

9.4.1	The E-Step Minimization Problem . . . . .	30
9.4.2	The Projection onto the Feasible Set . . . . .	31
9.5	Gradient and Scaled Gradient Projection and Projected Newton . . . . .	33
9.5.1	Gradient Projection Algorithm . . . . .	33
9.5.2	Scaled Gradient Projection and Projected Newton Method . . . . .	34
9.5.3	Combined Method . . . . .	35
<b>10</b>	<b>Alternative Gaussian Prior</b>	<b>36</b>
<b>11</b>	<b>Hyperparameters Selected for Method Assessment</b>	<b>38</b>
<b>12</b>	<b>Data Set I</b>	<b>39</b>
<b>13</b>	<b>Data Set II</b>	<b>43</b>
<b>14</b>	<b>RFN Pretraining for Convolution Nets</b>	<b>47</b>
<b>15</b>	<b>Running Times for RFN's Projected Newton Step</b>	<b>47</b>

## List of Theorems

1	Theorem (RFN Convergence) . . . . .	9
2	Theorem (RFN Correctness: Diagonal Noise Covariance Update) . . . . .	12
3	Theorem (RFN Correctness: Full Noise Covariance Update) . . . . .	14
4	Theorem (GAM Convergence Theorem) . . . . .	21
5	Theorem (Newton Update for Loading Matrix) . . . . .	24
6	Theorem (Newton Update for Noise Covariance) . . . . .	24
7	Theorem (Newton Update for Inverse Noise Covariance) . . . . .	25
8	Theorem (Projection: Rectifying) . . . . .	29
9	Theorem (Projection: Rectifying and Normalizing) . . . . .	31
10	Theorem (Theorem 5.4.5 in Kelley [1999]) . . . . .	34
11	Theorem (Lemma 5.5.1 in Kelley [1999]) . . . . .	35

## List of Algorithms

1	Rectified Factor Network . . . . .	6
2	Projection with E-Step Improvement . . . . .	7
3	Simple Projection: Rectifying . . . . .	7
4	Simple Projection: Rectifying and Normalization . . . . .	8
5	Scaled Newton Projection . . . . .	8
6	Scaled Projection With Reduced Matrix . . . . .	8
7	Weight Decay . . . . .	9
8	Dropout . . . . .	9

## 1 Introduction

This supplement contains additional information complementing the main manuscript and is structured as follows: First, the rectified factor network (RFN) learning algorithm with E- and M-step updates, weight decay and dropout regularization is given in Section 2. In Section 3, we proof that the (RFN) learning algorithm is a “generalized alternating minimization” (GAM) algorithm and converges to a solution that maximizes the RFN objective. The correctness of the RFN algorithm is proofed in Section 4. Section 5 describes the maximum likelihood factor analysis model and the model selection by the EM-algorithm. The RFN objective, which has to be maximized, is described in Section 6. Next, RFN’s GAM algorithm via gradient descent both in the M-step and the E-step is reported in the Section 7. The following sections 8 and 9 describe the gradient-based M- and E-step, respectively. In Section 10, we describe how the RFNs sparseness can be controlled by a Gaussian prior. Additional information on the selected hyperparameters of the benchmark methods is given in Section 11. The sections 12 and 13 describe the data generation of the benchmark datasets and report the results for three different experimental settings, namely for extracting 50 (undercomplete), 100 (complete) or 150 (overcomplete) factors / hidden units. In Section 14 describes experiments, that we have done to assess the performance of RFN *first layer* pretraining on *CIFAR-10* and *CIFAR-100* for three deep convolutional network architectures: (i) the AlexNet Ciresan et al. [2012], Krizhevsky et al. [2012], (ii) Deeply Supervised Networks (DSN) Lee et al. [2014], and (iii) our 5-Convolution-Network-In-Network (5C-NIN). Finally, Section 15 provides running times for RFN’s projected Newton step and for solving a quadratic program.

## 2 Rectified Factor Network (RFN) Algorithms

Algorithm 1 is the rectified factor network (RFN) learning algorithm. The RFN algorithm calls Algorithm 2 to project the posterior probability  $p_i$  onto the family of rectified and normalized variational distributions  $Q_i$ . Algorithm 2 guarantees an improvement of the E-step objective  $O = \frac{1}{n} \sum_{i=1}^n D_{\text{KL}}(Q_i \parallel p_i)$ . Projection Algorithm 2 relies on different projections, where a more complicated projection is tried if a simpler one failed to improve the E-step objective. If all following Newton-based gradient projection methods fail to decrease the E-step objective, then projection Algorithm 2 falls back to gradient projection methods. First the equality constraints are solved and inserted into the objective. Thereafter, the constraints are convex and gradient projection methods are applied. This approach is called “generalized reduced gradient method” Abadie and Carpentier [1969], which is our preferred alternative method. If this method fails, then Rosen’s gradient projection method Rosen [1961] is used. Finally, the method of Haug and Arora Haug and Arora [1979] is used.

First we consider Newton-based projection methods, which are used by Algorithm 2. Algorithm 4 performs a simple projection, which is the projected Newton method with learning rate set to one. This projection is very fast and ideally suited to be performed on GPUs for RFNs with many coding units. Algorithm 3 is the fast and simple projection without normalization even simpler than Algorithm 4. Algorithm 5 generalizes Algorithm 4 by introducing step sizes  $\lambda$  and  $\gamma$ . The step size  $\lambda$  scales the gradient step, while  $\gamma$  scales the difference between to old projection and the new projection. For both  $\lambda$  and  $\gamma$  annealing steps, that is, learning rate decay is used to find an appropriate update.

If these Newton-based update rules do not work, then Algorithm 6 is used. Algorithm 6 performs a scaled projection with a reduced Hessian matrix  $\mathbf{H}$  instead of the full Hessian  $\Sigma_p^{-1}$ . For computing  $\mathbf{H}$  an  $\epsilon$ -active set is determined, which consists of all  $j$  with  $\mu_j \leq \epsilon$ . The reduced matrix  $\mathbf{H}$  is the Hessian  $\Sigma_p^{-1}$  with  $\epsilon$ -active columns and rows  $j$  fixed to unit vector  $e_j$ .

The RFN algorithm allows regularization of the parameters  $\mathbf{W}$  and  $\Psi$  (off-diagonal elements) by weight decay. Priors on the parameters can be introduced. If the priors are convex functions, then convergence of the RFN algorithm is still ensured. The weight decay Algorithm 7 can optionally be used after the M-step of Algorithm 1. Coding units can be regularized by dropout. However dropout is not covered by the convergence proof for the RFN algorithm. The dropout Algorithm 8 is applied during the projection between rectifying and normalization. Methods like mini-batches or other stochastic gradient methods are not covered by the convergence proof for the RFN algorithm. However, in Gunawardana and Byrne [2005] it is shown how to generalize the GAM convergence

proof to mini-batches as it is shown for the incremental EM algorithm. Dropout and other stochastic gradient methods can be show to converge similar to mini-batches.

---

**Algorithm 1** Rectified Factor Network

---

**Input**

for  $1 \leq i \leq n$ :  $\mathbf{v}_i \in \mathbb{R}^m$ ,  
number of coding units  $l$

**Hyper-Parameters**

$\Psi_{\min}, W_{\max}, \eta_{\Psi}, \eta_W, \rho, \tau, 1 < \eta \leq 1$

**Initialization**

$\Psi = \tau \mathbf{I}$ ,  $\mathbf{W}$  element-wise random in  $[-\rho, \rho]$ ,

$\mathbf{C} = \frac{1}{n} \sum_{k=1}^n \mathbf{v}_k \mathbf{v}_k^T$ , STOP=false

**Main**

**while** STOP=false **do**

**—E-step1—**

**for all**  $1 \leq i \leq n$  **do**

$$(\mu_p)_i = (\mathbf{I} + \mathbf{W}^T \Psi^{-1} \mathbf{W})^{-1} \mathbf{W}^T \Psi^{-1} \mathbf{v}_i$$

**end for**

$$\Sigma = (\mathbf{I} + \mathbf{W}^T \Psi^{-1} \mathbf{W})^{-1}$$

**—Projection—**

    perform projection of  $(\mu_p)_i$  onto the feasible set by Algorithm 2 giving  $\mu_i$

**—E-step2—**

$$\mathbf{U} = \frac{1}{n} \sum_{i=1}^n \mathbf{v}_i \mu_i^T$$

$$\mathbf{S} = \frac{1}{n} \sum_{i=1}^n \mu_i \mu_i^T + \Sigma$$

**—M-step—**

$$\eta_W = \eta_{\Psi} = \eta$$

$$\mathbf{E} = \mathbf{C} - \mathbf{U} \mathbf{W}^T - \mathbf{W} \mathbf{U} + \mathbf{W} \mathbf{S} \mathbf{W}^T$$

**—W update—**

$$\mathbf{W} = \mathbf{W} + \eta_W (\mathbf{U} \mathbf{S}^{-1} - \mathbf{W})$$

**—diagonal  $\Psi$  update—**

$$\Psi_{kk} = \Psi_{kk} + \eta_{\Psi} (E_{kk} - \Psi_{kk})$$

**end for**

**—full  $\Psi$  update—**

$$\Psi = \Psi + \eta_{\Psi} (\mathbf{E} - \Psi)$$

**—bound parameters—**

$$\mathbf{W} = \text{median}\{-W_{\max}, \mathbf{W}, W_{\max}\}$$

$$\Psi = \text{median}\{\Psi_{\min}, \Psi, \max\{\mathbf{C}\}\}$$

    if stopping criterion is met: STOP=true

**end while**

---

---

**Algorithm 2** Projection with E-Step Improvement

---

**Goal**

obtain  $\mu_i^{\text{new}} = \mu_i$  that decrease the E-step objective

**Input**

$\Sigma^{\text{new}} = \Sigma_p, \Sigma^{\text{old}} = \Sigma_p^{\text{old}}$

for  $1 \leq i \leq n$ :  $(\mu_p)_i, \mu_i^{\text{old}}, p_i = \mathcal{N}((\mu_p)_i, \Sigma_p)$

simple projection P (rectified or rectified & normalized),

E-step objective:  $O = \frac{1}{n} \sum_{i=1}^n D_{\text{KL}}(Q_i \parallel p_i)$

$\gamma_{\min}, \lambda_{\min}, \rho_\gamma, \rho_\lambda, \epsilon$  (for  $\epsilon$ -active set)

**Main****—Simple Projection—**

perform `Newton Projection` by Algorithm 4 or Algorithm 3

**—Scaled Projection—**

**if**  $0 \leq \Delta O$  **then**

following loop for: (1)  $\gamma$ , (2)  $\lambda$ , or (3)  $\gamma$  and  $\lambda$  annealing

$\gamma = \lambda = 1$

**while**  $0 \leq \Delta O$  and  $\lambda > \lambda_{\min}$  and  $\gamma > \gamma_{\min}$  **do**

$\gamma = \rho_\gamma \gamma$  (skipped for  $\lambda$  annealing)

$\lambda = \rho_\lambda \lambda$  (skipped for  $\gamma$  annealing)

perform `Scaled Newton Projection` by Algorithm 5

**end while**

**end if**

**—Scaled Projection With Reduced Matrix—**

**if**  $0 \leq \Delta O$  **then**

determine  $\epsilon$ -active set as all  $j$  with  $\mu_j \leq \epsilon$

set  $H$  to  $\Sigma_p^{-1}$  with  $\epsilon$ -active columns and rows  $j$  fixed to  $e_j$

following loop for: (1)  $\gamma$ , (2)  $\lambda$ , or (3)  $\gamma$  and  $\lambda$  annealing

$\gamma = \lambda = 1$

**while**  $0 \leq \Delta O$  and  $\lambda > \lambda_{\min}$  and  $\gamma > \gamma_{\min}$  **do**

$\gamma = \rho_\gamma \gamma$  (skipped for  $\lambda$  annealing)

$\lambda = \rho_\lambda \lambda$  (skipped for  $\gamma$  annealing)

perform `Scaled Projection With Reduced Matrix` by Algorithm 6

**end while**

**end if**

**—General Gradient Projection—**

**while**  $0 \leq \Delta O$  **do**

use generalized reduced gradient Abadie and Carpentier [1969] OR

use Rosen's gradient projection Rosen [1961] OR

use method of Haug and Arora Haug and Arora [1979]

**end while**

---

---

**Algorithm 3** Simple Projection: Rectifying

---

**Goal**

for  $1 \leq i \leq n$ : project  $(\mu_p)_i$  onto feasible set giving  $\mu_i$

**Input**

$(\mu_p)_i$

**Main**

**for all**  $1 \leq j \leq l$  **do**

$\mu_{ij} = \max \left\{ 0, [(\mu_p)_i]_j \right\}$

**end for**

---

---

**Algorithm 4** Simple Projection: Rectifying and Normalization

---

**Goal**

for  $1 \leq i \leq n$ : project  $(\mu_p)_i$  onto feasible set giving  $\mu_i$

**Input**

for  $1 \leq i \leq n$ :  $(\mu_p)_i$

**Rectifier**

**for all**  $1 \leq i \leq n$  **do**

**for all**  $1 \leq j \leq l$  **do**

$$\hat{\mu}_{ij} = \max \left\{ 0, [(\mu_p)_i]_j \right\}$$

**end for**

**end for**

**Normalizer**

**for all**  $1 \leq i \leq n$  **do**

**if** at least one  $\hat{\mu}_{ij} > 0$  **then**

**for all**  $1 \leq j \leq l$  **do**

$$\mu_{ij} = \frac{\hat{\mu}_{ij}}{\sqrt{\frac{1}{n} \sum_{s=1}^n \hat{\mu}_{sj}^2}}$$

**end for**

**else**

**for all**  $1 \leq j \leq l$  **do**

$$\mu_{ij} = \begin{cases} \sqrt{n} & \text{for } j = \arg \max_j \{ [(\mu_p)_i]_j \} \\ 0 & \text{otherwise} \end{cases}$$

**end for**

**end if**

**end for**

---

---

**Algorithm 5** Scaled Newton Projection

---

**Goal**

perform a scaled Newton step with subsequent projection

**Input**

for  $1 \leq i \leq n$ :  $(\mu_p)_i$

for  $1 \leq i \leq n$ :  $\mu_i^{\text{old}}$

simple projection P (rectified or rectified & normalized),

$\lambda$  (gradient step size),  $\gamma$  (projection difference)

**Main**

$$d = P(\mu_i^{\text{old}} + \lambda((\mu_p)_i - \mu_i^{\text{old}}))$$

$$\mu_i^{\text{new}} = P(\mu_i^{\text{old}} + \gamma(d - \mu_i^{\text{old}}))$$

---

---

**Algorithm 6** Scaled Projection With Reduced Matrix

---

**Goal**

perform a scaled projection step with reduced matrix

**Input**

for  $1 \leq i \leq n$ :  $(\mu_p)_i$

for  $1 \leq i \leq n$ :  $\mu_i^{\text{old}}$

simple projection P (rectified or rectified & normalized),

$\lambda, \gamma, H, \Sigma_p^{-1}$

**Main**

$$d = P(\mu_i^{\text{old}} + \lambda H^{-1} \Sigma_p^{-1}((\mu_p)_i - \mu_i^{\text{old}}))$$

$$\mu_i^{\text{new}} = P(\mu_i^{\text{old}} + \gamma(d - \mu_i^{\text{old}}))$$

---



---

**Algorithm 7** Weight Decay

---

**Input**

Parameters  $\mathbf{W}$   
 Weight decay factors  $\gamma_G$  (Gaussian) and  $\gamma_L$  (Laplacian)

**Gaussian**

$$\mathbf{W} = \mathbf{W} - \gamma_G \mathbf{W}$$

**Laplacian**

$$\tilde{\mathbf{W}} = \text{median}\{-\gamma_L, \mathbf{W}, \gamma_L\}$$

$$\mathbf{W} = \mathbf{W} - \tilde{\mathbf{W}}$$


---

---

**Algorithm 8** Dropout

---

**Input**

for  $1 \leq i \leq n$ :  $\mu_i$   
 dropout probability  $d$

**Main**

**for all**  $1 \leq i \leq n$  **do**  
   **for all**  $1 \leq j \leq l$  **do**  
      $\text{Pr}(\delta = 0) = d$   
      $\mu_{ij} = \delta \mu_{ij}$   
   **end for**  
**end for**

---

### 3 Convergence Proof for the RFN Learning Algorithm

**Theorem 1** (RFN Convergence). *The rectified factor network (RFN) learning algorithm given in Algorithm 1 is a “generalized alternating minimization” (GAM) algorithm and converges to a solution that maximizes the objective  $\mathcal{F}$ .*

*Proof.* The factor analysis EM algorithm is given by Eq. (67) and Eq. (68) in Section 5. Algorithm 1 is the factor analysis EM algorithm with modified the E-step and the M-step. The E-step is modified by constraining the variational distribution  $Q$  to non-negative means and by normalizing its means across the samples. The M-step is modified to a Newton direction gradient step.

Like EM factor analysis, Algorithm 1 aims at maximizing the negative *free energy*  $\mathcal{F}$ , which is

$$\begin{aligned}
 \mathcal{F} &= \frac{1}{n} \sum_{i=1}^n \log p(\mathbf{v}_i) - \frac{1}{n} \sum_{i=1}^n D_{\text{KL}}(Q(\mathbf{h}_i) \parallel p(\mathbf{h}_i \mid \mathbf{v}_i)) \\
 &= \frac{1}{n} \sum_{i=1}^n \int Q(\mathbf{h}_i) \log p(\mathbf{v}_i) d\mathbf{h}_i - \frac{1}{n} \sum_{i=1}^n \int Q(\mathbf{h}_i) \log \frac{Q(\mathbf{h}_i)}{p(\mathbf{h}_i \mid \mathbf{v}_i)} d\mathbf{h}_i \\
 &= -\frac{1}{n} \sum_{i=1}^n \int Q(\mathbf{h}_i) \log \frac{Q(\mathbf{h}_i)}{p(\mathbf{h}_i, \mathbf{v}_i)} d\mathbf{h}_i \\
 &= -\frac{1}{n} \sum_{i=1}^n \int Q(\mathbf{h}_i) \log \frac{Q(\mathbf{h}_i)}{p(\mathbf{h}_i)} d\mathbf{h}_i + \frac{1}{n} \sum_{i=1}^n \int Q(\mathbf{h}_i) \log p(\mathbf{v}_i \mid \mathbf{h}_i) d\mathbf{h}_i \\
 &= \frac{1}{n} \sum_{i=1}^n \int Q(\mathbf{h}_i) \log p(\mathbf{v}_i \mid \mathbf{h}_i) d\mathbf{h}_i - \frac{1}{n} \sum_{i=1}^n D_{\text{KL}}(Q(\mathbf{h}_i) \parallel p(\mathbf{h}_i)) .
 \end{aligned} \tag{1}$$

$D_{\text{KL}}$  denotes the Kullback-Leibler (KL) divergence Kullback and Leibler [1951], which is larger than or equal to zero.

Algorithm 1 decreases  $\frac{1}{n} \sum_{i=1}^n D_{\text{KL}}(Q(\mathbf{h}_i) \parallel p(\mathbf{h}_i \mid \mathbf{v}_i))$  (the E-step objective) in its E-step under constraints for non-negative means and normalization. The constraint optimization problem from

Section 9.2 for the E-step is

$$\begin{aligned}
\min_{Q(\mathbf{h}_i)} \quad & \frac{1}{n} \sum_{i=1}^n D_{\text{KL}}(Q(\mathbf{h}_i) \parallel p(\mathbf{h}_i \mid \mathbf{v}_i)) \\
\text{s.t.} \quad & \forall_i : \boldsymbol{\mu}_i \geq \mathbf{0}, \\
& \forall_j : \frac{1}{n} \sum_{i=1}^n \mu_{ij}^2 = 1.
\end{aligned} \tag{2}$$

The M-step of Algorithm 1 aims at decreasing

$$\mathcal{E} = -\frac{1}{n} \sum_{i=1}^n \int_{\mathbb{R}^l} Q(\mathbf{h}_i) \log(p(\mathbf{v}_i \mid \mathbf{h}_i)) d\mathbf{h}_i. \tag{3}$$

Algorithm 1 performs one gradient descent step in the Newton direction to decrease  $\mathcal{E}$ , while EM factor analysis minimizes  $\mathcal{E}$ .

From the modification of the E-step and the M-step follows that Algorithm 1 is a *Generalized Alternating Minimization (GAM)* algorithm according to Gunawardana and Byrne [2005]. GAM is an EM algorithm that increases  $\mathcal{F}$  in the E-step and increases  $\mathcal{F}$  in the M-step (see also Section 7). The most important requirements for the convergence of the GAM algorithm according to Theorem 4 (Proposition 5 in Gunawardana and Byrne [2005]) are the increase of the objective  $\mathcal{F}$  in both the E-step and the M-step. Therefore we first show these two decreases before showing that all requirements of convergence Theorem 4 are met.

**Algorithm 1 ensures to decrease the M-step objective.** The M-step objective  $\mathcal{E}$  is convex in  $\mathbf{W}$  and  $\boldsymbol{\Psi}^{-1}$  according to Theorem 5 and Theorem 7. The update with  $\eta_W = \eta_\Psi = \eta = 1$  leads to the minimum of  $\mathcal{E}$  according to Theorem 5 and Theorem 7. The convexity of  $\mathcal{E}$  guarantees that each update with  $0 < \eta_W = \eta_\Psi = \eta \leq 1$  decreases the M-step objective  $\mathcal{E}$ , except the current  $\mathbf{W}$  and  $\boldsymbol{\Psi}^{-1}$  are already the minimizers.

**Algorithm 1 ensures to decrease the E-step objective.** The E-step decrease of Algorithm 1 is performed by Algorithm 2. According to Theorem 11 the scaled projection with reduced matrix ensures a decrease of the E-step objective for rectifying constraints (convex feasible set). According to Theorem 10 also gradient projection methods ensure a decrease of the E-step objective for rectifying constraints. For rectifying constraints and normalization, the feasible set is not convex because of the equality constraints. To optimize such problems, the generalized reduced gradient method Abadie and Carpentier [1969] solves each equality constraint for one variable and inserts it into the objective. For our problem Eq. (146) gives the solution and Eq. (147) the resulting convex constraints. Now scaled projection and gradient projection methods can be applied. For rectifying and normalizing constraints, also Rosen's Rosen [1961] and Haug & Arora's Haug and Arora [1979] gradient projection method ensures a decrease of the E-step objective since they can be applied to non-convex problems.

We show that the requirements as given in Section 7 for GAM convergence according to Theorem 4 (Proposition 5 in Gunawardana and Byrne [2005]) are fulfilled:

1. the learning rules, that is, the E-step and the M-step, are closed maps  $\longrightarrow$  ensured by continuous and continuous differentiable maps,
2. the parameter set is compact  $\longrightarrow$  ensured by bounding  $\boldsymbol{\Psi}$  and  $\mathbf{W}$ ,
3. the family of variational distributions is compact (often described by the feasible set of parameters of the variational distributions)  $\longrightarrow$  ensured by continuous and continuous differentiable functions for the constraints and by the bounds on the variational parameters  $\boldsymbol{\mu}$  and  $\boldsymbol{\Sigma}$  determined by bounds on the parameters and the data,
4. the support of the density models does not depend on the parameter  $\longrightarrow$  ensured by Gaussian models with full-rank covariance matrix,
5. the density models are continuous in the parameters  $\longrightarrow$  ensured by Gaussian models

6. the E-step has a unique maximizer  $\longrightarrow$  ensured by the convex, continuous, and continuous differentiable function that is minimized Dredze et al. [2008, 2012] together with compact feasible set for the variational parameters, the maximum may be local for non-convex feasible sets stemming from normalization,
7. the E-step increases the objective if not at the maximizer  $\longrightarrow$  ensured as shown above,
8. the M-step has a unique maximizer (this is not required)  $\longrightarrow$  ensured by minimizing a convex, continuous and continuous differentiable function in the model parameter and a convex feasible set, the maximum is a global maximum,
9. the M-step increases the objective if not at the maximizer  $\longrightarrow$  ensured as shown above.

□

Since this Proposition 5 in Gunawardana and Byrne [2005] is based on Zangwill’s generalized convergence theorem, updates of the RFN algorithm are viewed as point-to-set mappings Zangwill [1969]. Therefore the numerical precision, the choice of the methods in the E-step, and GPU implementations are covered by the proof. That the M-step has a unique maximizer is not required to proof Theorem 1 by Theorem 4. However we obtain an alternative proof by exchanging the variational distribution  $Q$  and the parameters  $(\mathbf{W}, \Psi)$ , that is, exchanging the E-step and the M-step. A theorem analog to Theorem 4 but with E-step and M-step conditions exchanged can be derived from Zangwill’s generalized convergence theorem Zangwill [1969].

The resulting model from the GAM procedure is at a local maximum of the objective given the model family and the family of variational distributions. *The solution minimizes the KL-distance between the family of full variational distributions and full model family.* “Full” means that both the observed and the hidden variables are taken into account, where for the variational distributions the probability of the observations is set to 1. The *desired family* is defined as the set of all probability distributions that assign probability one to the observation. In our case the family of variational distributions is not the desired family since some distributions are excluded by the constraints. Therefore the solution of the GAM optimization does not guarantee stationary points in likelihood Gunawardana and Byrne [2005]. This means that we do not maximize the likelihood but minimize

$$-\mathcal{F} \approx D_{\text{KL}}(Q(\mathbf{h}, \mathbf{v}) \parallel p(\mathbf{h}, \mathbf{v})) + c \quad (4)$$

according to Eq. (73), where  $c$  is a constant independent of  $Q$  and independent of the model parameters.

## 4 Correctness Proofs for the RFN Learning Algorithms

The RFN algorithm is correct if it has a low reconstruction error and explains the data covariance matrix by its parameters like factor analysis. We show in Theorem 2 and Theorem 3 that the RFN algorithm

1. minimizes the reconstruction error given  $\mu_i$  and  $\Sigma$  (the error is quadratic in  $\Psi$ );
2. explains the covariance matrix by its parameters  $\mathbf{W}$  and  $\Psi$  plus an estimate of the second moment of the coding units  $\mathbf{S}$ .

Since the minimization of the reconstruction error is based on  $\mu_i$ , the quality of reconstruction and covariance explanation depends on the correlation between  $\mu_i$  and  $\mathbf{v}_i$ . The larger the correlation between  $\mu_i$  and  $\mathbf{v}_i$ , the lower the reconstruction error and the better the explanation of the data covariance. We ensure maximal information in  $\mu_i$  on  $\mathbf{v}_i$  by the I-projection (the minimal Kullback-Leibler distance) of the posterior onto the family of rectified and normalized Gaussian distributions.

The reconstruction error for given mean values  $\mu_i$  is

$$\frac{1}{n} \sum_{i=1}^n \|\epsilon_i\|_2^2, \quad (5)$$

where

$$\boldsymbol{\epsilon}_i = \mathbf{v}_i - \mathbf{W} \boldsymbol{\mu}_i . \quad (6)$$

The reconstruction error for using the whole variational distribution  $Q(\mathbf{h}_i)$  instead of its means is  $\Psi$ . Below we will derive Eq. (17), which is

$$\Psi = \text{diag} \left( \frac{1}{n} \sum_{i=1}^n \boldsymbol{\epsilon}_i \boldsymbol{\epsilon}_i^T + \mathbf{W} \boldsymbol{\Sigma} \mathbf{W}^T \right) . \quad (7)$$

Therefore  $\Psi$  is the reconstruction error for given mean values plus the variance  $\mathbf{W} \boldsymbol{\Sigma} \mathbf{W}^T$  introduced by the hidden variables.

#### 4.1 Diagonal Noise Covariance Update

**Theorem 2** (RFN Correctness: Diagonal Noise Covariance Update). *The fixed point  $\mathbf{W}$  minimizes  $\text{Tr}(\Psi)$  given  $\boldsymbol{\mu}_i$  and  $\boldsymbol{\Sigma}$  by ridge regression with*

$$\text{Tr}(\Psi) = \frac{1}{n} \sum_{i=1}^n \|\boldsymbol{\epsilon}_i\|_2^2 + \left\| \mathbf{W} \boldsymbol{\Sigma}^{1/2} \right\|_F^2 , \quad (8)$$

where we used the error

$$\boldsymbol{\epsilon}_i = \mathbf{v}_i - \mathbf{W} \boldsymbol{\mu}_i \quad (9)$$

The model explains the data covariance matrix by

$$\mathbf{C} = \Psi + \mathbf{W} \mathbf{S} \mathbf{W}^T \quad (10)$$

up to an error, which is quadratic in  $\Psi$  for  $\Psi \ll \mathbf{W} \mathbf{W}^T$ . The reconstruction error

$$\frac{1}{n} \sum_{i=1}^n \|\boldsymbol{\epsilon}_i\|_2^2 \quad (11)$$

is quadratic in  $\Psi$  for  $\Psi \ll \mathbf{W} \mathbf{W}^T$ .

*Proof.* The fixed point equation for the  $\mathbf{W}$  update is

$$\Delta \mathbf{W} = \mathbf{U} \mathbf{S}^{-1} - \mathbf{W} = \mathbf{0} \Rightarrow \mathbf{W} = \mathbf{U} \mathbf{S}^{-1} . \quad (12)$$

Using the definition of  $\mathbf{U}$  and  $\mathbf{S}$ , the fixed point equation Eq. (12) gives

$$\mathbf{W} = \left( \frac{1}{n} \sum_{i=1}^n \mathbf{v}_i \boldsymbol{\mu}_i^T \right) \left( \frac{1}{n} \sum_{i=1}^n \boldsymbol{\mu}_i \boldsymbol{\mu}_i^T + \boldsymbol{\Sigma} \right)^{-1} \quad (13)$$

Therefore  $\mathbf{W}$  is a *ridge regression* estimate, also called *generalized Tikhonov regularization* estimate, which minimizes

$$\begin{aligned} & \frac{1}{n} \sum_{i=1}^n \|\mathbf{v}_i - \mathbf{W} \boldsymbol{\mu}_i\|_2^2 + \left\| \mathbf{W} \boldsymbol{\Sigma}^{1/2} \right\|_F^2 \\ &= \frac{1}{n} \sum_{i=1}^n \|\boldsymbol{\epsilon}_i\|_2^2 + \left\| \mathbf{W} \boldsymbol{\Sigma}^{1/2} \right\|_F^2 \\ &= \frac{1}{n} \sum_{i=1}^n \boldsymbol{\epsilon}_i^T \boldsymbol{\epsilon}_i + \text{Tr} \left( \mathbf{W} \boldsymbol{\Sigma}^{1/2} \boldsymbol{\Sigma}^{1/2} \mathbf{W}^T \right) \\ &= \text{Tr} \left( \frac{1}{n} \sum_{i=1}^n \boldsymbol{\epsilon}_i \boldsymbol{\epsilon}_i^T + \mathbf{W} \boldsymbol{\Sigma} \mathbf{W}^T \right) , \end{aligned} \quad (14)$$

where we used the reconstruction error

$$\boldsymbol{\epsilon}_i = \mathbf{v}_i - \mathbf{W} \boldsymbol{\mu}_i . \quad (15)$$

We obtain with this definition of the error

$$\begin{aligned}
& \frac{1}{n} \sum_{i=1}^n \epsilon_i \epsilon_i^T + \mathbf{W} \Sigma \mathbf{W}^T \\
&= \frac{1}{n} \sum_{i=1}^n \mathbf{v}_i \mathbf{v}_i^T - \frac{1}{n} \sum_{i=1}^n \mathbf{v}_i \mu_i^T \mathbf{W}^T - \frac{1}{n} \sum_{i=1}^n \mathbf{W} \mu_i \mathbf{v}_i^T \\
&\quad + \frac{1}{n} \sum_{i=1}^n \mathbf{W} \mu_i \mu_i^T \mathbf{W}^T + \mathbf{W} \Sigma \mathbf{W}^T \\
&= \mathbf{C} - \mathbf{U} \mathbf{W}^T - \mathbf{W} \mathbf{U}^T + \mathbf{W} \mathbf{S} \mathbf{W}^T.
\end{aligned} \tag{16}$$

Therefore from the fixed point equation for  $\Psi$  with the diagonal update rule follows

$$\Psi = \text{diag} \left( \frac{1}{n} \sum_{i=1}^n \epsilon_i \epsilon_i^T + \mathbf{W} \Sigma \mathbf{W}^T \right), \tag{17}$$

where “diag” projects a matrix to a diagonal matrix. From this follows that

$$\text{Tr}(\Psi) = \text{Tr} \left( \frac{1}{n} \sum_{i=1}^n \epsilon_i \epsilon_i^T + \mathbf{W} \Sigma \mathbf{W}^T \right). \tag{18}$$

Consequently, the fixed point  $\mathbf{W}$  minimizes  $\text{Tr}(\Psi)$  given  $\mu_i$  and  $\Sigma$ .

After convergence of the algorithm  $\Sigma = (\mathbf{I} + \mathbf{W}^T \Psi^{-1} \mathbf{W})^{-1}$  holds. The Woodbury identity (matrix inversion lemma) states

$$(\mathbf{W} \mathbf{W}^T + \Psi)^{-1} = \Psi^{-1} - \Psi^{-1} \mathbf{W} (\mathbf{I} + \mathbf{W}^T \Psi^{-1} \mathbf{W})^{-1} \mathbf{W}^T \Psi^{-1} \tag{19}$$

from which follows by multiplying the equation from right and left by  $\Psi$  that

$$\begin{aligned}
\mathbf{W} \Sigma \mathbf{W}^T &= \mathbf{W} (\mathbf{I} + \mathbf{W}^T \Psi^{-1} \mathbf{W})^{-1} \mathbf{W}^T \\
&= \Psi - \Psi (\mathbf{W} \mathbf{W}^T + \Psi)^{-1} \Psi
\end{aligned} \tag{20}$$

Inserting this equation Eq. (20) into Eq. (17) gives

$$\begin{aligned}
\Psi &= \text{diag} \left( \frac{1}{n} \sum_{i=1}^n \epsilon_i \epsilon_i^T + \Psi - \Psi (\mathbf{W} \mathbf{W}^T + \Psi)^{-1} \Psi \right) \\
&= \Psi + \text{diag} \left( \frac{1}{n} \sum_{i=1}^n \epsilon_i \epsilon_i^T - \Psi (\mathbf{W} \mathbf{W}^T + \Psi)^{-1} \Psi \right).
\end{aligned} \tag{21}$$

Therefore we have

$$\text{diag} \left( \frac{1}{n} \sum_{i=1}^n \epsilon_i \epsilon_i^T - \Psi (\mathbf{W} \mathbf{W}^T + \Psi)^{-1} \Psi \right) = \mathbf{0}. \tag{22}$$

It follows that

$$\begin{aligned}
\text{Tr} \left( \frac{1}{n} \sum_{i=1}^n \epsilon_i \epsilon_i^T \right) &= \text{Tr} \left( \Psi (\mathbf{W} \mathbf{W}^T + \Psi)^{-1} \Psi \right) \\
&\leq \text{Tr} \left( (\mathbf{W} \mathbf{W}^T + \Psi)^{-1} \right) \text{Tr}(\Psi)^2.
\end{aligned} \tag{23}$$

The inequality uses the fact that for positive definite matrices  $\mathbf{A}$  and  $\mathbf{B}$  inequality  $\text{Tr}(\mathbf{AB}) \leq \text{Tr}(\mathbf{A})\text{Tr}(\mathbf{B})$  holds Patel and Toda [1979]. Thus, for  $\Psi \ll \mathbf{W} \mathbf{W}^T$  the error  $\text{Tr} \left( \frac{1}{n} \sum_{i=1}^n \epsilon_i \epsilon_i^T \right) = \frac{1}{n} \sum_{i=1}^n \epsilon_i^T \epsilon_i$  is quadratic in  $\Psi$ .

Multiplying the fixed point equation Eq. (12) by  $\mathbf{S}$  gives  $\mathbf{U} = \mathbf{W} \mathbf{S}$ . Therefore we have:

$$\mathbf{W} \mathbf{U}^T = \mathbf{W} \mathbf{S} \mathbf{W}^T = \mathbf{U} \mathbf{W}^T. \tag{24}$$

Inserting Eq. (20) into the first line of Eq. (16) and Eq. (24) for simplifying the last line of Eq. (16) gives

$$\frac{1}{n} \sum_{i=1}^n \epsilon_i \epsilon_i^T - \Psi (W W^T + \Psi)^{-1} \Psi = C - \Psi - W S W^T. \quad (25)$$

Using the trace norm (nuclear norm or Ky-Fan n-norm) on matrices, Eq. (23) states that the left hand side is quadratic in  $\Psi$  for  $\Psi \ll W W^T$ . The trace norm of a positive semi-definite matrix is its trace and bounds the Frobenius norm Srebro [2004]. Furthermore, Eq. (22) states that the left hand side of this equation has zero diagonal entries. Therefore it follows that

$$C = \Psi + W S W^T \quad (26)$$

holds except an error, which is quadratic in  $\Psi$  for  $\Psi \ll W W^T$ . The diagonal is exactly modeled according to Eq. (22).  $\square$

Therefore the model corresponding to the fixed point explains the empirical matrix of second moments  $C$  by a noise part  $\Psi$  and a signal part  $W S W^T$ . Like factor analysis the data variance is explained by the model via the parameters  $\Psi$  (noise) and  $W$  (signal).

## 4.2 Full Noise Covariance Update

**Theorem 3** (RFN Correctness: Full Noise Covariance Update). *The fixed point  $W$  minimizes  $\text{Tr}(\Psi)$  given  $\mu_i$  and  $\Sigma$  by ridge regression with*

$$\text{Tr}(\Psi) = \frac{1}{n} \sum_{i=1}^n \|\epsilon_i\|_2^2 + \left\| W \Sigma^{1/2} \right\|_F^2, \quad (27)$$

where we used the error

$$\epsilon_i = v_i - W \mu_i \quad (28)$$

The model explains the data covariance matrix by

$$C = \Psi + W S W^T. \quad (29)$$

The reconstruction error

$$\frac{1}{n} \sum_{i=1}^n \|\epsilon_i\|_2^2 \quad (30)$$

is quadratic in  $\Psi$  for  $\Psi \ll W W^T$ .

*Proof.* The first part follows from previous Theorem 2. The fixed point equation for the  $\Psi$  update is

$$\Psi = C - U W^T - W U^T + W S W^T, \quad (31)$$

using Eq. (24) this leads to

$$C = \Psi + W S W^T. \quad (32)$$

From Eq. (16) follows for the fixed point of  $\Psi$  with the full update rule:

$$\Psi = \frac{1}{n} \sum_{i=1}^n \epsilon_i \epsilon_i^T + W \Sigma W^T. \quad (33)$$

Inserting Eq. (20) into Eq. (33) gives

$$\Psi = \frac{1}{n} \sum_{i=1}^n \epsilon_i \epsilon_i^T + \Psi - \Psi (W W^T + \Psi)^{-1} \Psi, \quad (34)$$

from which follows

$$\frac{1}{n} \sum_{i=1}^n \epsilon_i \epsilon_i^T = \Psi (W W^T + \Psi)^{-1} \Psi. \quad (35)$$

Thus, the error  $\text{Tr}(\frac{1}{n} \sum_{i=1}^n \epsilon_i \epsilon_i^T) = \frac{1}{n} \sum_{i=1}^n \epsilon_i^T \epsilon_i$  is quadratic in  $\Psi$ , for  $\Psi \ll W W^T$ .  $\square$

## 5 Maximum Likelihood Factor Analysis

We are given the data  $\{v\} = \{v_1, \dots, v_n\}$  which is assumed to be centered. Centering can be done by subtracting the mean  $\mu$  from the data. The model is

$$v = Wh + \epsilon, \quad (36)$$

where

$$h \sim \mathcal{N}(0, I) \quad \text{and} \quad \epsilon \sim \mathcal{N}(0, \Psi). \quad (37)$$

The model includes the *observations*  $v \in \mathbb{R}^m$ , the *noise*  $\epsilon \in \mathbb{R}^m$ , the *factors*  $h \in \mathbb{R}^l$ , the *factor loading matrix*  $W \in \mathbb{R}^{m \times l}$ , and the *noise covariance matrix*  $\Psi \in \mathbb{R}^{m \times m}$ . Typically we assume that  $\Psi$  is a diagonal matrix to explain data covariance by signal and not by noise. The data variance is explained through a signal part  $Wh$  and through a noise part  $\epsilon$ . The parameters of the model are  $W$  and  $\Psi$ . From the model assumption it follows that if  $h$  is given, then only the noise  $\epsilon$  is a random variable and we have

$$v | h \sim \mathcal{N}(Wh, \Psi). \quad (38)$$

We want to derive the *likelihood* of the data under the model, that is, the likelihood that the model has produced the data. Let  $E$  denote the expectation of the data including the prior distribution of the factors and the noise distribution. We obtain for the first two moments and the variance:

$$\begin{aligned} E(v) &= E(Wh + \epsilon) = WE(h) + E(\epsilon) = 0, \\ E(v v^T) &= E((Wh + \epsilon)(Wh + \epsilon)^T) = \\ &\quad WE(h h^T) W^T + WE(h) E(\epsilon^T) \\ &\quad + E(\epsilon) E(h^T) W^T + E(\epsilon \epsilon^T) = \\ &\quad W W^T + \Psi \\ \text{var}(v) &= E(v v^T) - (E(v))^2 = W W^T + \Psi. \end{aligned} \quad (39)$$

The observations are Gaussian distributed since their distribution is the product of two Gaussian densities divided by a normalizing constant. Therefore, the marginal distribution for  $v$  is

$$v \sim \mathcal{N}(0, W W^T + \Psi). \quad (41)$$

The log-likelihood  $\log \prod_{i=1}^n p(v_i)$  of the data  $\{v\}$  under the model  $(W, \Psi)$  is

$$\begin{aligned} \log \prod_{i=1}^n p(v_i) &= \log \prod_{i=1}^n (2\pi)^{-m/2} |W W^T + \Psi|^{-1/2} \\ &\quad \exp\left(-\frac{1}{2} \left(v_i^T (W W^T + \Psi)^{-1} v_i\right)\right) \\ &= -\frac{n m}{2} \log(2\pi) - \frac{n}{2} \log |W W^T + \Psi| \\ &\quad - \frac{1}{2} \sum_{i=1}^n v_i^T (W W^T + \Psi)^{-1} v_i, \end{aligned} \quad (42)$$

where  $|\cdot|$  denotes the absolute value of the determinant of a matrix.

To maximize the likelihood is difficult since a closed form for the maximum does not exist. Therefore, typically the expectation maximization (EM) algorithm is used to maximize the likelihood. For the EM algorithm a variational distribution  $Q$  is required which estimates the factors given the observations.

We consider a single data vector  $v_i$ . The posterior is also Gaussian with mean  $(\mu_p)_i$  and covariance matrix  $\Sigma_p$ :

$$\begin{aligned} h_i | v_i &\sim \mathcal{N}((\mu_p)_i, \Sigma_p) \\ (\mu_p)_i &= W^T (W W^T + \Psi)^{-1} v_i \\ \Sigma_p &= I - W^T (W W^T + \Psi)^{-1} W, \end{aligned} \quad (43)$$

where we used the fact that

$$\begin{aligned} \mathbf{a} &\sim \mathcal{N}(\boldsymbol{\mu}_a, \Sigma_{aa}), \quad \mathbf{u} \sim \mathcal{N}(\boldsymbol{\mu}_u, \Sigma_{uu}), \\ \Sigma_{ua} &= \text{Cov}(\mathbf{u}, \mathbf{a}) \text{ and } \Sigma_{au} = \text{Cov}(\mathbf{a}, \mathbf{u}) : \\ \mathbf{a} | \mathbf{u} &\sim \mathcal{N}(\boldsymbol{\mu}_a + \Sigma_{au} \Sigma_{uu}^{-1} (\mathbf{u} - \boldsymbol{\mu}_u), \Sigma_{aa} - \Sigma_{au} \Sigma_{uu}^{-1} \Sigma_{ua}) \end{aligned} \quad (44)$$

and

$$\mathbf{E}(\mathbf{h}\mathbf{v}) = \mathbf{W} \mathbf{E}(\mathbf{h} \mathbf{h}^T) = \mathbf{W}. \quad (45)$$

The EM algorithm sets  $Q$  to the posterior distribution for data vector  $\mathbf{v}_i$ :

$$Q_i(\mathbf{h}_i) = p(\mathbf{h}_i | \mathbf{v}_i; \mathbf{W}, \boldsymbol{\Psi}) = \mathcal{N}((\boldsymbol{\mu}_p)_i, \boldsymbol{\Sigma}_p), \quad (46)$$

therefore we obtain for standard EM

$$\boldsymbol{\mu}_i = (\boldsymbol{\mu}_q)_i = (\boldsymbol{\mu}_p)_i \quad (47)$$

$$\boldsymbol{\Sigma} = \boldsymbol{\Sigma}_q = \boldsymbol{\Sigma}_p. \quad (48)$$

The matrix inversion lemma (Woodbury identity) can be used to compute  $\boldsymbol{\mu}_i$  and  $\boldsymbol{\Sigma}$ :

$$(\mathbf{W} \mathbf{W}^T + \boldsymbol{\Psi})^{-1} = \boldsymbol{\Psi}^{-1} - \boldsymbol{\Psi}^{-1} \mathbf{W} (\mathbf{I} + \mathbf{W}^T \boldsymbol{\Psi}^{-1} \mathbf{W})^{-1} \mathbf{W}^T \boldsymbol{\Psi}^{-1}. \quad (49)$$

Using this identity, the mean and the covariance matrix can be computed as:

$$\begin{aligned} \boldsymbol{\mu}_i &= \mathbf{W}^T (\mathbf{W} \mathbf{W}^T + \boldsymbol{\Psi})^{-1} \mathbf{v}_i = (\mathbf{I} + \mathbf{W}^T \boldsymbol{\Psi}^{-1} \mathbf{W})^{-1} \mathbf{W}^T \boldsymbol{\Psi}^{-1} \mathbf{v}_i, \\ \boldsymbol{\Sigma} &= \mathbf{I} - \mathbf{W}^T (\mathbf{W} \mathbf{W}^T + \boldsymbol{\Psi})^{-1} \mathbf{W} = (\mathbf{I} + \mathbf{W}^T \boldsymbol{\Psi}^{-1} \mathbf{W})^{-1}. \end{aligned} \quad (50)$$

The EM algorithm maximizes a lower bound  $\mathcal{F}$  on the log-likelihood:

$$\begin{aligned} \mathcal{F} &= \log p(\mathbf{v}_i) - D_{\text{KL}}(Q(\mathbf{h}_i) \| p(\mathbf{h}_i | \mathbf{v}_i)) \\ &= \int Q(\mathbf{h}_i) \log p(\mathbf{v}_i) d\mathbf{h}_i - \int Q(\mathbf{h}_i) \log \frac{Q(\mathbf{h}_i)}{p(\mathbf{h}_i | \mathbf{v}_i)} d\mathbf{h}_i \\ &= - \int Q(\mathbf{h}_i) \log \frac{Q(\mathbf{h}_i)}{p(\mathbf{h}_i, \mathbf{v}_i)} d\mathbf{h}_i \\ &= - \int Q(\mathbf{h}_i) \log \frac{Q(\mathbf{h}_i)}{p(\mathbf{h}_i)} d\mathbf{h}_i + \int Q(\mathbf{h}_i) \log p(\mathbf{v}_i | \mathbf{h}_i) d\mathbf{h}_i \\ &= \int Q(\mathbf{h}_i) \log p(\mathbf{v}_i | \mathbf{h}_i) d\mathbf{h}_i - D_{\text{KL}}(Q(\mathbf{h}_i) \| p(\mathbf{h}_i)). \end{aligned} \quad (51)$$

$D_{\text{KL}}$  denotes the Kullback-Leibler (KL) divergence Kullback and Leibler [1951] which is larger than zero.

$\mathcal{F}$  is the EM objective which has to be maximized in order to maximize the likelihood. The **E-step** maximizes  $\mathcal{F}$  with respect to the variational distribution  $Q$ , therefore the E-step minimizes  $D_{\text{KL}}(Q(\mathbf{h}_i) \| p(\mathbf{h}_i | \mathbf{v}_i))$ . After the standard unconstrained E-step, the variational distribution is equal to the posterior, i.e.  $Q(\mathbf{h}_i) = p(\mathbf{h}_i | \mathbf{v}_i)$ . Therefore the KL divergence

$$D_{\text{KL}}(Q(\mathbf{h}_i) \| p(\mathbf{h}_i | \mathbf{v}_i)) = 0 \quad (52)$$

is zero, thus  $\mathcal{F}$  is equal to the log-likelihood  $\log p(\mathbf{v}_i)$  ( $\mathcal{F} = \log p(\mathbf{v}_i)$ ). The **M-step** maximizes  $\mathcal{F}$  with respect to the parameters  $(\mathbf{W}, \boldsymbol{\Psi})$ , therefore the M-step maximizes  $\int Q(\mathbf{h}_i) \log p(\mathbf{v}_i | \mathbf{h}_i) d\mathbf{h}_i$ .

We next consider again all  $n$  samples  $\{\mathbf{v}\} = \{\mathbf{v}_1, \dots, \mathbf{v}_n\}$ . The *expected reconstruction error*  $\mathcal{E}$  for these  $n$  data samples is

$$\mathcal{E} = -\frac{1}{n} \sum_{i=1}^n \int_{\mathbb{R}^l} Q(\mathbf{h}_i) \log (p(\mathbf{v}_i | \mathbf{h}_i)) d\mathbf{h}_i = \frac{1}{n} \sum_{i=1}^n \mathbb{E}_Q (\log (p(\mathbf{v}_i | \mathbf{h}_i))) \quad (53)$$

and objective to maximize becomes

$$\mathcal{F} = -\mathcal{E} - \frac{1}{n} \sum_{i=1}^n D_{\text{KL}}(Q(\mathbf{h}_i) \| p(\mathbf{h}_i)). \quad (54)$$



The M-step requires to minimize  $\mathcal{E}$ :

$$\mathcal{E} = \frac{m}{2} \log(2\pi) + \frac{1}{2} \log |\Psi| + \quad (55)$$

$$\begin{aligned} & \frac{1}{2n} \sum_{i=1}^n \mathbb{E}_Q \left( (v_i - \mathbf{W} h_i)^T \Psi^{-1} (v_i - \mathbf{W} h_i) \right) \\ &= \frac{m}{2} \log(2\pi) + \frac{1}{2} \log |\Psi| + \end{aligned} \quad (56)$$

$$\begin{aligned} & \frac{1}{2n} \sum_{i=1}^n \mathbb{E}_Q (v_i^T \Psi^{-1} v_i - 2 v_i^T \Psi^{-1} \mathbf{W} h_i + h_i^T \mathbf{W}^T \Psi^{-1} \mathbf{W} h_i) \\ &= \frac{m}{2} \log(2\pi) + \frac{1}{2} \log |\Psi| + \frac{1}{2n} \sum_{i=1}^n v_i^T \Psi^{-1} v_i \end{aligned} \quad (57)$$

$$\begin{aligned} & - \text{Tr} \left( \Psi^{-1} \mathbf{W} \sum_{i=1}^n \mathbb{E}_Q (h_i) v_i^T \right) + \frac{1}{2} \text{Tr} \left( \mathbf{W}^T \Psi^{-1} \mathbf{W} \sum_{i=1}^n \mathbb{E}_Q (h_i h_i^T) \right) \\ &= \frac{m}{2} \log(2\pi) + \frac{1}{2} \log |\Psi| + \frac{1}{2} \text{Tr} \left( \Psi^{-1} \frac{1}{n} \sum_{i=1}^n v_i v_i^T \right) \end{aligned} \quad (58)$$

$$\begin{aligned} & - \text{Tr} \left( \Psi^{-1} \mathbf{W} \frac{1}{n} \sum_{i=1}^n \mu_i v_i^T \right) + \frac{1}{2} \text{Tr} \left( \mathbf{W}^T \Psi^{-1} \mathbf{W} \frac{1}{n} \sum_{i=1}^n (\Sigma + \mu_i \mu_i^T) \right) \\ &= \frac{1}{2} (m \log(2\pi) + \log |\Psi| + \text{Tr}(\Psi^{-1} \mathbf{C})) \\ & - 2 \text{Tr}(\Psi^{-1} \mathbf{W} \mathbf{U}^T) + \text{Tr}(\mathbf{W}^T \Psi^{-1} \mathbf{W} \mathbf{S}), \end{aligned} \quad (59)$$

where  $\text{Tr}$  gives the trace of a matrix.

The derivatives with respect to the parameters are set to zero for the optimal parameters:

$$\nabla_{\mathbf{W}} \mathcal{E} = -\frac{1}{2n} \sum_{i=1}^n \Psi^{-1} \mathbf{W} \mathbb{E}_Q (h_i h_i^T) + \frac{1}{2n} \sum_{i=1}^n \Psi^{-1} v_i \mathbb{E}_Q^T (h_i) = \mathbf{0} \quad (60)$$

and

$$\begin{aligned} \nabla_{\Psi} \mathcal{E} &= -\frac{1}{2} \Psi^{-1} + \\ & \frac{1}{2n} \sum_{i=1}^n \mathbb{E}_Q \left( \Psi^{-1} (v_i - \mathbf{W} h_i) (v_i - \mathbf{W} h_i)^T \Psi^{-1} \right) = \mathbf{0}. \end{aligned} \quad (61)$$

Solving above equations gives:

$$\mathbf{W}^{\text{new}} = \left( \frac{1}{n} \sum_{i=1}^n v_i \mathbb{E}_{h_i|v_i}^T (h_i) \right) \left( \frac{1}{n} \sum_{i=1}^n \mathbb{E}_Q (h_i h_i^T) \right)^{-1} \quad (62)$$

and

$$\begin{aligned} \Psi^{\text{new}} &= \frac{1}{n} \sum_{i=1}^n \mathbb{E}_Q \left( (v_i - \mathbf{W}^{\text{new}} h_i) (v_i - \mathbf{W}^{\text{new}} h_i)^T \right) = \\ & \frac{1}{n} \sum_{i=1}^n v_i v_i^T - \frac{1}{n} \sum_{i=1}^n v_i \mathbb{E}_Q^T (h_i) (\mathbf{W}^{\text{new}})^T - \\ & \frac{1}{n} \sum_{i=1}^n \mathbf{W}^{\text{new}} \mathbb{E}_Q (h_i) v_i^T + \mathbf{W}^{\text{new}} \frac{1}{n} \sum_{i=1}^n \mathbb{E}_Q (h_i h_i^T) (\mathbf{W}^{\text{new}})^T. \end{aligned} \quad (63)$$

We obtain the following EM updates:

$$\mathbf{E}\text{-step:} \tag{64}$$

$$\begin{aligned} \boldsymbol{\mu}_i &= (\mathbf{I} + \mathbf{W}^T \boldsymbol{\Psi}^{-1} \mathbf{W})^{-1} \mathbf{W}^T \boldsymbol{\Psi}^{-1} \mathbf{v}_i, \\ \boldsymbol{\Sigma} &= (\mathbf{I} + \mathbf{W}^T \boldsymbol{\Psi}^{-1} \mathbf{W})^{-1}, \\ \mathbb{E}_Q(\mathbf{h}_i) &= \boldsymbol{\mu}_i \\ \mathbb{E}_Q(\mathbf{h}_i \mathbf{h}_i^T) &= \boldsymbol{\mu}_i \boldsymbol{\mu}_i^T + \boldsymbol{\Sigma} \end{aligned}$$

$$\mathbf{M}\text{-step:} \tag{65}$$

$$\begin{aligned} \mathbf{W}^{\text{new}} &= \left( \frac{1}{n} \sum_{i=1}^n \mathbf{v}_i \mathbb{E}_{\mathbf{h}_i|\mathbf{v}_i}^T(\mathbf{h}_i) \right) \left( \frac{1}{n} \sum_{i=1}^n \mathbb{E}_Q(\mathbf{h}_i \mathbf{h}_i^T) \right)^{-1} \\ \boldsymbol{\Psi}^{\text{new}} &= \frac{1}{n} \sum_{i=1}^n \mathbf{v}_i \mathbf{v}_i^T - \frac{1}{n} \sum_{i=1}^n \mathbf{v}_i \mathbb{E}_Q^T(\mathbf{h}_i) (\mathbf{W}^{\text{new}})^T - \\ &\quad \frac{1}{n} \sum_{i=1}^n \mathbf{W}^{\text{new}} \mathbb{E}_Q(\mathbf{h}_i) \mathbf{v}_i^T + \mathbf{W}^{\text{new}} \frac{1}{n} \sum_{i=1}^n \mathbb{E}_Q(\mathbf{h}_i \mathbf{h}_i^T) (\mathbf{W}^{\text{new}})^T. \end{aligned} \tag{66}$$

The EM algorithms can be reformulated as:

$$\mathbf{E}\text{-step:} \tag{67}$$

$$\begin{aligned} \boldsymbol{\mu}_i &= (\mathbf{I} + \mathbf{W}^T \boldsymbol{\Psi}^{-1} \mathbf{W})^{-1} \mathbf{W}^T \boldsymbol{\Psi}^{-1} \mathbf{v}_i, \\ \boldsymbol{\Sigma} &= (\mathbf{I} + \mathbf{W}^T \boldsymbol{\Psi}^{-1} \mathbf{W})^{-1}, \\ \mathbb{E}_Q(\mathbf{h}_i) &= \boldsymbol{\mu}_i \\ \mathbb{E}_Q(\mathbf{h}_i \mathbf{h}_i^T) &= \boldsymbol{\mu}_i \boldsymbol{\mu}_i^T + \boldsymbol{\Sigma} \end{aligned}$$

$$\mathbf{M}\text{-step:} \tag{68}$$

$$\begin{aligned} \mathbf{C} &= \frac{1}{n} \sum_{i=1}^n \mathbf{v}_i \mathbf{v}_i^T \\ \mathbf{U} &= \frac{1}{n} \sum_{i=1}^n \mathbf{v}_i \mathbb{E}_Q^T(\mathbf{h}_i) \end{aligned} \tag{69}$$

$$\mathbf{S} = \frac{1}{n} \sum_{i=1}^n \mathbb{E}_Q(\mathbf{h}_i \mathbf{h}_i^T) \tag{70}$$

$$\mathbf{W}^{\text{new}} = \mathbf{U} \mathbf{S}^{-1} \tag{71}$$

$$\boldsymbol{\Psi}^{\text{new}} = \mathbf{C} - \mathbf{U} \mathbf{W}^T - \mathbf{W} \mathbf{U}^T + \mathbf{W} \mathbf{S} \mathbf{W}^T. \tag{72}$$

## 6 The RFN Objective

*Our goal is to find a sparse, non-negative representation of the input which extracts structure from the input.* A sparse, non-negative representation is desired to code only events or objects that have caused the input. We assume that only few events or objects caused the input, therefore, we aim at sparseness. Furthermore, we do not want to code the degree of absence of events or objects. As the vast majority of events and objects is supposed to be absent, to code for their degree of absence would introduce a high level of random fluctuations.

We aim at extracting structures from the input, therefore generative models are used as they explicitly model input structures. For example factor analysis models the covariance structure of the data. However a generative model cannot enforce sparse, non-negative representation of the input.

The input representation of a generative model is the posterior’s mean, median, or mode. Generative models with rectified priors (zero probability for negative values) lead to rectified posteriors. However these posteriors do not have sparse means (they must be positive), that is, they do not yield sparse codes Frey and Hinton [1999]. For example, rectified factor analysis, which rectifies Gaussian priors and selects models using a variational Bayesian learning procedure, does not yield posteriors with sparse means Harva and Kaban [2005, 2007]. A generative model with hidden units  $\mathbf{h}$  and data  $\mathbf{v}$  is defined by its prior  $p(\mathbf{h})$  and its likelihood  $p(\mathbf{v} | \mathbf{h})$ . The posterior  $p(\mathbf{h} | \mathbf{v})$  supplies the input representation of a model by the posterior’s mean, median, or mode. However, the posterior depends on the data  $\mathbf{v}$ , therefore sparseness and non-negativity of its means cannot be guaranteed independent of the data. Problem at coding the input by generative models is the data-dependency of the posterior means.

Therefore we use the *posterior regularization method (posterior constraint method)* Ganchev et al. [2010], Graca et al. [2009, 2007]. The posterior regularization framework separates model characteristics from data dependent characteristics like the likelihood or posterior constraints. Posterior regularization incorporates data-dependent characteristics as constraints on model posteriors given the observed data, which are difficult to encode via model parameters by Bayesian priors.

A generative model with prior  $p(\mathbf{h})$  and likelihood  $p(\mathbf{v} | \mathbf{h})$  has the full model distribution  $p(\mathbf{h}, \mathbf{v}) = p(\mathbf{v} | \mathbf{h})p(\mathbf{h})$ . It can be written as  $p(\mathbf{h}, \mathbf{v}) = p(\mathbf{h} | \mathbf{v})p(\mathbf{v})$ , where  $p(\mathbf{h} | \mathbf{v})$  is the model posterior of the hidden variables and  $p(\mathbf{v})$  is the evidence, that is, the likelihood of the data to be produced by the model. The model family and its parametrization determines which structures are extracted from the data. Typically the model parameters enter the likelihood  $p(\mathbf{v} | \mathbf{h})$  and are adjusted to the observed data. For the posterior regularization method, a family  $\mathcal{Q}$  of allowed posterior distributions is introduced.  $\mathcal{Q}$  is defined by the expectations of constraint features. In our case the posterior means have to be non-negative. Distributions  $Q \in \mathcal{Q}$  are called *variational distributions* (see later for using this term). The full variational distribution is  $Q(\mathbf{h}, \mathbf{v}) = Q(\mathbf{h} | \mathbf{v})p_v(\mathbf{v})$  with  $Q(\mathbf{h} | \mathbf{v}) \in \mathcal{Q}$ . The distribution  $p_v(\mathbf{v})$  is the unknown distribution of observations as determined by the world or the data generation process. This distribution is approximated by samples drawn from the world, namely the training samples.  $p(\mathbf{h}, \mathbf{v})$  contains all model assumptions like the structures used to model the data, while  $Q(\mathbf{h}, \mathbf{v})$  contains all data dependent characteristics including data dependent constraints on the posterior.

The goal is to achieve  $Q(\mathbf{h}, \mathbf{v}) = p(\mathbf{h}, \mathbf{v})$ , to obtain (1) a desired structure that is extracted from the data and (2) desired code properties. However in general it is to achieve this identity, therefore we want to minimize the distance between these distributions. We use the Kullback-Leibler (KL) divergence Kullback and Leibler [1951]  $D_{\text{KL}}$  to measure the distance between these distributions. Therefore our objective is  $D_{\text{KL}}(Q(\mathbf{h}, \mathbf{v}) || p(\mathbf{h}, \mathbf{v}))$ . Minimizing this KL divergence (1) extracts the desired structure from the data by increasing the likelihood, that is,  $p_v(\mathbf{v}) \approx p(\mathbf{v})$ , and (2) enforces desired code properties by  $Q(\mathbf{h} | \mathbf{v}) \approx p(\mathbf{h} | \mathbf{v})$ . Thus, the code derived from  $Q(\mathbf{h} | \mathbf{v})$  has the desired properties and t extracts the desired input data structures.

We now approximate the KL divergence by approximating the expectation over  $p_v(\mathbf{v})$  by the empirical mean of samples  $\{\mathbf{v}\} = \{\mathbf{v}_1, \dots, \mathbf{v}_n\}$  drawn from  $p_v(\mathbf{v})$ :

$$\begin{aligned} D_{\text{KL}}(Q(\mathbf{h}, \mathbf{v}) || p(\mathbf{h}, \mathbf{v})) &= \int Q(\mathbf{h}, \mathbf{v}) \log \frac{Q(\mathbf{h}, \mathbf{v})}{p(\mathbf{h}, \mathbf{v})} d\mathbf{h} d\mathbf{v} \\ &= \int_V p_v(\mathbf{v}) \int_H Q(\mathbf{h} | \mathbf{v}) \log \frac{Q(\mathbf{h}, \mathbf{v})}{p(\mathbf{h}, \mathbf{v})} d\mathbf{h} d\mathbf{v} \\ &\approx \frac{1}{n} \sum_{i=1}^n \int_H Q(\mathbf{h} | \mathbf{v}_i) \log \frac{Q(\mathbf{h}, \mathbf{v}_i)}{p(\mathbf{h}, \mathbf{v}_i)} d\mathbf{h} \\ &= \frac{1}{n} \sum_{i=1}^n \int_H Q(\mathbf{h} | \mathbf{v}_i) \log \frac{Q(\mathbf{h} | \mathbf{v}_i)}{p(\mathbf{h}, \mathbf{v}_i)} d\mathbf{h} + \frac{1}{n} \sum_{i=1}^n \log p_v(\mathbf{v}_i). \end{aligned} \tag{73}$$

The last term  $\frac{1}{n} \sum_{i=1}^n \log p_v(\mathbf{v}_i)$  neither depends on  $Q$  nor on the model, therefore we will neglect it. In the following, we often abbreviate  $Q(\mathbf{h} | \mathbf{v}_i)$  by  $Q(\mathbf{h}_i)$  or write  $Q(\mathbf{h}_i | \mathbf{v}_i)$ , since the hidden variable is based on the observation  $\mathbf{v}_i$ . Similarly we often write  $p(\mathbf{h}_i, \mathbf{v}_i)$  instead of  $p(\mathbf{h}, \mathbf{v}_i)$  and even more often  $p(\mathbf{h}_i | \mathbf{v}_i)$  instead of  $p(\mathbf{h} | \mathbf{v}_i)$ .

We obtain the objective  $\mathcal{F}$  (to be maximized) of the *posterior constraint method* Ganchev et al. [2010], Graca et al. [2009, 2007]:

$$\mathcal{F} = \frac{1}{n} \sum_{i=1}^n \log p(\mathbf{v}_i) - \frac{1}{n} \sum_{i=1}^n D_{\text{KL}}(Q(\mathbf{h}_i) \parallel p(\mathbf{h}_i | \mathbf{v}_i)) \quad (74)$$

$$\begin{aligned} &= \frac{1}{n} \sum_{i=1}^n \int Q(\mathbf{h}_i) \log p(\mathbf{v}_i) d\mathbf{h}_i - \frac{1}{n} \sum_{i=1}^n \int Q(\mathbf{h}_i) \log \frac{Q(\mathbf{h}_i)}{p(\mathbf{h}_i | \mathbf{v}_i)} d\mathbf{h}_i \\ &= -\frac{1}{n} \sum_{i=1}^n \int Q(\mathbf{h}_i) \log \frac{Q(\mathbf{h}_i)}{p(\mathbf{h}_i, \mathbf{v}_i)} d\mathbf{h}_i \\ \text{nonumber} &= -\frac{1}{n} \sum_{i=1}^n \int Q(\mathbf{h}_i) \log \frac{Q(\mathbf{h}_i)}{p(\mathbf{h}_i)} d\mathbf{h}_i + \frac{1}{n} \sum_{i=1}^n \int Q(\mathbf{h}_i) \log p(\mathbf{v}_i | \mathbf{h}_i) d\mathbf{h}_i \quad (75) \\ &= \frac{1}{n} \sum_{i=1}^n \int Q(\mathbf{h}_i) \log p(\mathbf{v}_i | \mathbf{h}_i) d\mathbf{h}_i - \frac{1}{n} \sum_{i=1}^n D_{\text{KL}}(Q(\mathbf{h}_i) \parallel p(\mathbf{h}_i)) . \end{aligned}$$

The first line is the negative objective of the posterior constraint method while the third line is the negative Eq. (73) without the term  $\frac{1}{n} \sum_{i=1}^n \log p_v(\mathbf{v}_i)$ .

**$\mathcal{F}$  is the objective in our framework which has to be maximized.** Maximizing  $\mathcal{F}$  (1) increases the model likelihood  $\frac{1}{n} \sum_{i=1}^n \log p(\mathbf{v}_i)$ , (2) finds a proper input representation by small  $D_{\text{KL}}(Q(\mathbf{h}_i) \parallel p(\mathbf{h}_i | \mathbf{v}_i))$ . Thus, the data representation (1) extracts structures from the data as imposed by the generative model while (2) ensuring desired code properties via  $Q \in \mathcal{Q}$ .

In the variational framework,  $Q$  is the variational distribution and  $\mathcal{F}$  is called the *negative free energy* Neal and Hinton [1998]. This physical term is used since variational methods were introduced for quantum physics by Richard Feynman Feynman [1972]. The hidden variables can be considered as the fictive causes or explanations of environmental fluctuations Friston [2012].

If  $p(\mathbf{h} | \mathbf{v}) \in \mathcal{Q}$ , then  $Q(\mathbf{h} | \mathbf{v}) = p(\mathbf{h} | \mathbf{v})$  and we obtain the classical EM algorithm. The EM algorithm maximizes the lower bound  $\mathcal{F}$  on the log-likelihood as seen at the first line of Eq. (74) and ensures in its E-step  $Q(\mathbf{h} | \mathbf{v}) = p(\mathbf{h} | \mathbf{v})$ .

## 7 Generalized Alternating Minimization

Instead of the EM algorithm we use the *Generalized Alternating Minimization (GAM)* algorithm Gunawardana and Byrne [2005] to allow for gradient descent both in the M-step and the E-step. The representation of an input by a generative model is the vector of the mean values of the posterior, that is, the most likely hidden variables that produced the observed data. We have to modify the E-step to enforce variational distributions which lead to sparse codes via zero values of the components of its mean vector. Sparse codes, that is, many components of the mean vector are zero, are obtained by enforcing non-negative means. This rectification is analog to rectified linear units for neural networks, which have enabled sparse codes for neural networks. Therefore the variational distributions are restricted to stem from a family with non-negative constraints on the means. To impose constraints on the posterior is known as the *posterior constraint method* Ganchev et al. [2010], Graca et al. [2009, 2007]. The posterior constraint method maximizes the objective both in the E-step and the M-step. The posterior constraint method is computationally infeasible for our approach, since we assume a large number of hidden units. For models with many hidden units, the maximization in the E-step would take too much time. The posterior constraint method does not support fast implementations on GPUs and stochastic gradients, which we want to allow in order to use mini-batches and dropout regularization.

Therefore we perform only one gradient descent step both in the E-step and in the M-step. Unfortunately, the convergence proofs of the EM algorithm are no longer valid. However we show that our algorithm is a generalized alternating minimization (GAM) method. Gunawardana and Byrne showed that the GAM converges Gunawardana and Byrne [2005] (see also Wu [1983]).

The following GAM convergence Theorem 4 is Proposition 5 in Gunawardana and Byrne [2005] and proves the convergence of the GAM algorithm to a solution that minimizes  $-\mathcal{F}$ .

**Theorem 4** (GAM Convergence Theorem). *Let the point-to-set map  $FB$  be the composition  $B \circ F$  of point-to-set maps  $F : \mathcal{D} \times \Theta \rightarrow \mathcal{D} \times \Theta$  and  $B : \mathcal{D} \times \Theta \rightarrow \mathcal{D} \times \Theta$ . Suppose that the point-to-set maps  $F$  and  $B$  are defined so that*

- (1)  $F$  and  $B$  are closed on  $\mathcal{D}' \times \Theta$
- (2)  $F(\mathcal{D}' \times \Theta) \subseteq \mathcal{D} \times \Theta$  and  $B(\mathcal{D}' \times \Theta) \subseteq \mathcal{D} \times \Theta$

*Suppose also that  $F$  is such that all  $(Q'_X, \theta') \in F(Q_X, \theta)$  have  $\theta' = \theta$  and satisfy*

$$(GAM.F): \quad D_{KL}(Q'_X \parallel p_{X;\theta}) \leq D_{KL}(Q_X \parallel p_{X;\theta})$$

*with equality only if*

$$(EQ.F): \quad Q_X = \arg \min_{Q'_X \in \mathcal{D}} D_{KL}(Q'_X \parallel p_{X;\theta}),$$

*with  $Q_X$  being the unique minimizer. Suppose also that the point-to-set map  $B$  is such that all  $(Q'_X, \theta') \in B(Q_X, \theta)$  have  $Q'_X = Q_X$  and satisfy*

$$(GAM.B): \quad D_{KL}(Q_X \parallel p_{X;\theta'}) \leq D_{KL}(Q_X \parallel p_{X;\theta})$$

*with equality only if*

$$(EQ.B): \quad \theta \in \arg \min_{\xi \in \Theta} D_{KL}(Q_X \parallel p_{X;\xi}).$$

*Then,*

- (1) *the point-to-set map  $FB$  is closed on  $\mathcal{D}' \times \Theta$*
- (2)  *$FB(\mathcal{D}' \times \Theta) \subseteq \mathcal{D} \times \Theta$*

*and  $FB$  satisfies the GAM and EQ conditions of the GAM convergence theorem, that is, Theorem 3 in Gunawardana and Byrne [2005].*

*Proof.* See Proposition 5 in Gunawardana and Byrne [2005]. □

The point-to-set mappings allow extended E-step and M-steps without unique iterates. Therefore, Theorem 4 holds for different implementations, different hardware, different precisions of the algorithm under consideration.

For a GAM method to converge, we have to ensure that the objective increases in both the E-step and the M-step.  $Q$  is from a constrained family of variational distributions, while the posterior and the full distribution (observation and hidden units) are both derived from a model family. The model family is a parametrized family. For our models (i) the support of the density models does not depend on the parameter and (ii) the density models are continuous in their parameters. GAM convergence requires both (i) and (ii). Furthermore, both the E-step and the M-step must have unique maximizers and they increase the objective if they are not at a maximum point.

The learning rules, that is, the E-step and the M-step are closed maps as they are continuous functions. The objective for the E-step is strict convex in all its parameters for the variational distributions, simultaneously Dredze et al. [2008, 2012]. It is quadratic for the mean vectors on which constraints are imposed. The objective for the M-step is convex in both parameters  $\mathbf{W}$  and  $\Psi^{-1}$  (we sometimes estimate  $\Psi$  instead of  $\Psi^{-1}$ ). The objective is quadratic in the loading matrix  $\mathbf{W}$ . For rectifying only, we guarantee unique global maximizers by convex and compact sets for both the family of desired distributions and the set of possible parameters. For this convex optimization problem with one *global* maximum. For rectifying and normalizing, the family of desired distributions is *not* convex due to equality constraints introduced by the normalization. However we can guarantee *local* unique maximizers.

Summary of the requirements for GAM convergence Theorem 4:

1. the learning rules, that is, the E-step and the M-step, are closed maps,
2. the parameter set is compact,

3. the family of variational distributions is compact (often described by the feasible set of parameters of the variational distributions),
4. the support of the density models does not depend on the parameter,
5. the density models are continuous in the parameters,
6. the E-step has a unique maximizer,
7. the E-step increases the objective if not at the maximizer,
8. the M-step has a unique maximizer (not required by Theorem 4),
9. the M-step increases the objective if not at the maximizer.

The resulting model from the GAM procedure is at a local maximum of the objective given the model family and the family of variational distributions. *The solution minimizes the KL-distance between the family of full variational distributions and full model family.* “Full” means that both the observed and the hidden variables are taken into account, where for the variational distributions the probability of the observations is set to 1. The *desired family* is defined as the set of all probability distributions that assign probability one to the observation. In our case the family of variational distributions is not the desired family since some distributions are excluded by the constraints. Therefore the solution of the GAM optimization does not guarantee stationary points in likelihood Gunawardana and Byrne [2005]. This means that we do not maximize the likelihood but minimize the KL-distance between variational distributions and model.

## 8 Gradient-based M-step

### 8.1 Gradient Ascent

The gradients in the M-step are:

$$\nabla_{\mathbf{W}} \mathcal{E} = \frac{1}{2n} \sum_{i=1}^n \Psi^{-1} \mathbf{v}_i \mathbf{E}_Q^T(\mathbf{h}_i) - \frac{1}{2n} \sum_{i=1}^n \Psi^{-1} \mathbf{W} \mathbf{E}_Q(\mathbf{h}_i \mathbf{h}_i^T)$$

and

$$\nabla_{\Psi} \mathcal{E} = -\frac{1}{2} \Psi^{-1} + \frac{1}{2n} \sum_{i=1}^n \mathbf{E}_Q \left( \Psi^{-1} (\mathbf{v}_i - \mathbf{W} \mathbf{h}_i) (\mathbf{v}_i - \mathbf{W} \mathbf{h}_i)^T \Psi^{-1} \right). \quad (76)$$

Alternatively, we can estimate  $\Psi^{-1}$  which leads to the derivatives:

$$\nabla_{\Psi^{-1}} \mathcal{E} = \frac{1}{2} \Psi - \frac{1}{2n} \sum_{i=1}^n \mathbf{E}_Q \left( (\mathbf{v}_i - \mathbf{W} \mathbf{h}_i) (\mathbf{v}_i - \mathbf{W} \mathbf{h}_i)^T \right). \quad (77)$$

Scaling the gradients leads to:

$$2 \nabla_{\mathbf{W}} \mathcal{E} = \Psi^{-1} \frac{1}{n} \sum_{i=1}^n \mathbf{v}_i \mathbf{E}_Q^T(\mathbf{h}_i) - \Psi^{-1} \mathbf{W} \frac{1}{n} \sum_{i=1}^n \mathbf{E}_Q(\mathbf{h}_i \mathbf{h}_i^T) \quad (78)$$

and

$$\begin{aligned} 2 \nabla_{\Psi} \mathcal{E} = & \quad (79) \\ & - \Psi^{-1} + \Psi^{-1} \left( \frac{1}{n} \sum_{i=1}^n \mathbf{v}_i \mathbf{v}_i^T - \frac{1}{n} \sum_{i=1}^n \mathbf{v}_i \mathbf{E}_Q^T(\mathbf{h}_i) \mathbf{W}^T \right. \\ & \left. - \frac{1}{n} \sum_{i=1}^n \mathbf{W} \mathbf{E}_Q(\mathbf{h}_i) \mathbf{v}_i^T + \mathbf{W} \frac{1}{n} \sum_{i=1}^n \mathbf{E}_Q(\mathbf{h}_i \mathbf{h}_i^T) \mathbf{W}^T \right) \Psi^{-1}. \end{aligned}$$

or

$$2 \nabla_{\Psi^{-1}} \mathcal{E} = \Psi - \left( \frac{1}{n} \sum_{i=1}^n \mathbf{v}_i \mathbf{v}_i^T - \frac{1}{n} \sum_{i=1}^n \mathbf{v}_i \mathbb{E}_Q^T(\mathbf{h}_i) \mathbf{W}^T - \frac{1}{n} \sum_{i=1}^n \mathbf{W} \mathbb{E}_Q(\mathbf{h}_i) \mathbf{v}_i^T + \mathbf{W} \frac{1}{n} \sum_{i=1}^n \mathbb{E}_Q(\mathbf{h}_i \mathbf{h}_i^T) \mathbf{W}^T \right). \quad (80)$$

Only the sums

$$\mathbf{U} = \frac{1}{n} \sum_{i=1}^n \mathbf{v}_i \mathbb{E}_Q^T(\mathbf{h}_i) \quad (81)$$

and

$$\mathbf{S} = \frac{1}{n} \sum_{i=1}^n \mathbb{E}_Q(\mathbf{h}_i \mathbf{h}_i^T) \quad (82)$$

must be computed for both gradients.

$$\mathbf{C} = \frac{1}{n} \sum_{i=1}^n \mathbf{v}_i \mathbf{v}_i^T \quad (83)$$

is the estimated covariance matrix (matrix of second moments for zero mean).

**The generalized EM algorithm update rules are:**

$$\mathbf{C} = \frac{1}{n} \sum_{i=1}^n \mathbf{v}_i \mathbf{v}_i^T \quad (84)$$

**E-step:**

$$\begin{aligned} \boldsymbol{\mu}_i &= \mathbf{W}^T (\mathbf{W} \mathbf{W}^T + \Psi)^{-1} \mathbf{v}_i = (\mathbf{I} + \mathbf{W}^T \Psi^{-1} \mathbf{W})^{-1} \mathbf{W}^T \Psi^{-1} \mathbf{v}_i, \\ \boldsymbol{\Sigma} &= \mathbf{I} - \mathbf{W}^T (\mathbf{W} \mathbf{W}^T + \Psi)^{-1} \mathbf{W} = (\mathbf{I} + \mathbf{W}^T \Psi^{-1} \mathbf{W})^{-1}, \\ \mathbb{E}_Q(\mathbf{h}_i) &= \boldsymbol{\mu}_i \\ \mathbb{E}_Q(\mathbf{h}_i \mathbf{h}_i^T) &= \boldsymbol{\mu}_i \boldsymbol{\mu}_i^T + \boldsymbol{\Sigma} \\ \mathbf{U} &= \frac{1}{n} \sum_{i=1}^n \mathbf{v}_i \mathbb{E}_Q^T(\mathbf{h}_i) \\ \mathbf{S} &= \frac{1}{n} \sum_{i=1}^n \mathbb{E}_Q(\mathbf{h}_i \mathbf{h}_i^T) \end{aligned}$$

**M-step:** (85)

$$\begin{aligned} \Delta \mathbf{W} &= \Psi^{-1} \mathbf{U} - \Psi^{-1} \mathbf{W} \mathbf{S} \\ \Delta \Psi &= -\Psi^{-1} + \Psi^{-1} (\mathbf{C} - \mathbf{U} \mathbf{W}^T - \mathbf{W} \mathbf{U} + \mathbf{W} \mathbf{S} \mathbf{W}^T) \Psi^{-1}. \end{aligned}$$

## 8.2 Newton Update

Instead of gradient ascent, we now consider a Newton update step. The Newton update for finding the roots of  $\frac{\partial f}{\partial \mathbf{v}}$  is

$$\mathbf{v}_{n+1} = \mathbf{v}_n - \eta \mathbf{H}^{-1} \nabla_{\mathbf{v}} f(\mathbf{v}_n), \quad (86)$$

where  $\eta$  is a small step size and  $\mathbf{H}$  is the Hessian of  $f$  with respect to  $\mathbf{v}$  evaluated at  $\mathbf{v}_n$ . We denote the update direction by

$$\Delta \mathbf{v} = -\mathbf{H}^{-1} \nabla_{\mathbf{v}} f(\mathbf{v}_n). \quad (87)$$

### 8.2.1 Newton Update of the Loading Matrix

**Theorem 5** (Newton Update for Loading Matrix). *The M-step objective  $\mathcal{E}$  is quadratic in  $\mathbf{W}$ , thus convex in  $\mathbf{W}$ . The Newton update direction for  $\mathbf{W}$  in the M-step is*

$$\Delta \mathbf{W} = \mathbf{U} \mathbf{S}^{-1} - \mathbf{W}. \quad (88)$$

*Proof.* The M-step objective is the *expected reconstruction error*  $\mathcal{E}$ , which is according to Eq. (55)

$$\begin{aligned} \mathcal{E} = & -\frac{1}{n} \sum_{i=1}^n \int_{\mathbb{R}^l} Q(\mathbf{h}_i) \log(p(\mathbf{v}_i | \mathbf{h}_i)) d\mathbf{h}_i = \frac{1}{2} \left( m \log(2\pi) + \log|\Psi| \right. \\ & \left. + \text{Tr}(\Psi^{-1} \mathbf{C}) - 2 \text{Tr}(\Psi^{-1} \mathbf{W} \mathbf{U}^T) + \text{Tr}(\mathbf{W}^T \Psi^{-1} \mathbf{W} \mathbf{S}) \right), \end{aligned} \quad (89)$$

where  $\text{Tr}$  gives the trace of a matrix. This is a quadratic function in  $\mathbf{W}$ , as stated in the theorem.

The Hessian  $\mathbf{H}_{\mathbf{W}}$  of (2 $\mathcal{E}$ ) with respect to  $\mathbf{W}$  as a vector is:

$$\begin{aligned} \mathbf{H}_{\mathbf{W}} &= \frac{\partial \text{vec}(2 \nabla_{\mathbf{W}} \mathcal{E})}{\partial \text{vec}(\mathbf{W})^T} = \frac{\partial \text{vec}(-\Psi^{-1} \mathbf{U} + \Psi^{-1} \mathbf{W} \mathbf{S})}{\partial \text{vec}(\mathbf{W})^T} \\ &= \mathbf{S} \otimes \Psi^{-1}, \end{aligned} \quad (90)$$

where  $\otimes$  is the Kronecker product of matrices.  $\mathbf{H}_{\mathbf{W}}$  is positive definite, thus the problem is convex in  $\mathbf{W}$ . The inverse of  $\mathbf{H}_{\mathbf{W}}$  is

$$\mathbf{H}_{\mathbf{W}}^{-1} = \mathbf{S}^{-1} \otimes \Psi. \quad (91)$$

For the product of the inverse Hessian with the gradient we have:

$$\begin{aligned} \mathbf{H}_{\mathbf{W}}^{-1} \text{vec}(-\Psi^{-1} \mathbf{U} + \Psi^{-1} \mathbf{W} \mathbf{S}) &= \text{vec}(\Psi(-\Psi^{-1} \mathbf{U} + \Psi^{-1} \mathbf{W} \mathbf{S}) \mathbf{S}^{-1}) \\ &= \text{vec}(-\mathbf{U} \mathbf{S}^{-1} + \mathbf{W}). \end{aligned} \quad (92)$$

If we apply a Newton update, then the update direction for  $\mathbf{W}$  in the M-step is

$$\Delta \mathbf{W} = \mathbf{U} \mathbf{S}^{-1} - \mathbf{W}. \quad (93)$$

□

This is the exact EM update if the step-size  $\eta$  is 1. Since the objective is a quadratic function in  $\mathbf{W}$ , one Newton update would lead to the exact solution.

### 8.2.2 Newton Update of the Noise Covariance

We define the expected approximation error by

$$\begin{aligned} \mathbf{E} &= \mathbf{C} - \mathbf{U} \mathbf{W}^T - \mathbf{W} \mathbf{U} + \mathbf{W} \mathbf{S} \mathbf{W}^T \\ &= \frac{1}{n} \sum_{i=1}^n \mathbb{E}_Q \left( (\mathbf{v}_i - \mathbf{W} \mathbf{h}_i) (\mathbf{v}_i - \mathbf{W} \mathbf{h}_i)^T \right). \end{aligned} \quad (94)$$

**$\Psi$  as parameter.**

**Theorem 6** (Newton Update for Noise Covariance). *The Newton update direction for  $\Psi$  as parameter in the M-step is*

$$\Delta \Psi = \mathbf{E} - \Psi. \quad (95)$$

An update with  $\Delta \Psi$  ( $\eta = 1$ ) leads to the minimum of the M-step objective  $\mathcal{E}$ .

*Proof.* The M-step objective is the *expected reconstruction error*  $\mathcal{E}$ , which is according to Eq. (55)

$$\begin{aligned} \mathcal{E} = & -\frac{1}{n} \sum_{i=1}^n \int_{\mathbb{R}^l} Q(\mathbf{h}_i) \log(p(\mathbf{v}_i | \mathbf{h}_i)) d\mathbf{h}_i = \frac{1}{2} \left( m \log(2\pi) + \log|\Psi| \right. \\ & \left. + \text{Tr}(\Psi^{-1} \mathbf{C}) - 2 \text{Tr}(\Psi^{-1} \mathbf{W} \mathbf{U}^T) + \text{Tr}(\mathbf{W}^T \Psi^{-1} \mathbf{W} \mathbf{S}) \right), \end{aligned} \quad (96)$$



where  $\text{Tr}$  gives the trace of a matrix.

Since

$$2 \nabla_{\Psi} \mathcal{E} = \Psi^{-1} - \Psi^{-1} E \Psi^{-1}, \quad (97)$$

is

$$\Psi = E \quad (98)$$

the minimum of  $\mathcal{E}$  with respect to  $\Psi$ . Therefore an update with  $\Delta \Psi = E - \Psi$  leads to the minimum.

The Hessian  $H_{\Psi}$  of  $(2\mathcal{E})$  with respect to  $\Psi$  as a vector is:

$$\begin{aligned} H_{\Psi} &= \frac{\partial \text{vec}(2 \nabla_{\Psi} \mathcal{E})}{\partial \text{vec}(\Psi)^T} = \frac{\partial \text{vec}(\Psi^{-1} - \Psi^{-1} E \Psi^{-1})}{\partial \text{vec}(\Psi)^T} \\ &= -\Psi^{-1} \otimes \Psi^{-1} + \Psi^{-1} \otimes (\Psi^{-1} E \Psi^{-1}) + (\Psi^{-1} E \Psi^{-1}) \otimes \Psi^{-1}. \end{aligned} \quad (99)$$

The expected approximation error  $E$  is a sample estimate for  $\Psi$ , therefore we have  $\Psi \approx E$ . The Hessian may not be positive definite for some values of  $E$ , like for small values of  $E$ . In order to guarantee a positive definite Hessian, more precisely an approximation to it, for minimization, we set

$$E = \Psi \quad (100)$$

and obtain

$$H_{\Psi} = \Psi^{-1} \otimes \Psi^{-1}. \quad (101)$$

We derive an approximate Newton update that is very close to the Newton update.

The inverse of the approximated  $H_{\Psi}$  is

$$H_{\Psi}^{-1} = \Psi \otimes \Psi. \quad (102)$$

For the product of the inverse Hessian with the gradient we have:

$$\begin{aligned} H_{\Psi}^{-1} \text{vec}(\Psi^{-1} - \Psi^{-1} E \Psi^{-1}) &= \text{vec}(\Psi (\Psi^{-1} - \Psi^{-1} E \Psi^{-1}) \Psi) \\ &= \text{vec}(\Psi - E). \end{aligned} \quad (103)$$

If we apply a Newton update, then the update direction for  $\Psi$  in the M-step is

$$\Delta \Psi = E - \Psi. \quad (104)$$

This is the exact EM update if the step-size  $\eta$  is 1.  $\square$

**$\Psi^{-1}$  as parameter.**

**Theorem 7** (Newton Update for Inverse Noise Covariance). *The M-step objective  $\mathcal{E}$  is convex in  $\Psi^{-1}$ . The Newton update direction for  $\Psi^{-1}$  as parameter in the M-step is*

$$\Delta \Psi^{-1} = \Psi^{-1} - \Psi^{-1} E \Psi^{-1}. \quad (105)$$

*A first order approximation of this Newton direction for  $\Psi$  in the M-step is*

$$\Delta \Psi = E - \Psi. \quad (106)$$

*An update with  $\Delta \Psi$  ( $\eta = 1$ ) leads to the minimum of the M-step objective  $\mathcal{E}$ .*

*Proof.* The M-step objective is the *expected reconstruction error*  $\mathcal{E}$ , which is according to Eq. (55)

$$\begin{aligned} \mathcal{E} &= -\frac{1}{n} \sum_{i=1}^n \int_{\mathbb{R}^l} Q(\mathbf{h}_i) \log(p(\mathbf{v}_i | \mathbf{h}_i)) d\mathbf{h}_i = \frac{1}{2} \left( m \log(2\pi) + \log |\Psi| \right. \\ &\quad \left. + \text{Tr}(\Psi^{-1} C) - 2 \text{Tr}(\Psi^{-1} W U^T) + \text{Tr}(W^T \Psi^{-1} W S) \right), \end{aligned} \quad (107)$$

where  $\text{Tr}$  gives the trace of a matrix.

Since

$$2 \nabla_{\Psi^{-1}} \mathcal{E} = -\Psi + E \quad (108)$$

is

$$\Psi = E \quad (109)$$

the minimum of  $\mathcal{E}$  with respect to  $\Psi^{-1}$ . Therefore an update with  $\Delta\Psi = E - \Psi$  leads to the minimum.

The Hessian  $H_{\Psi^{-1}}$  of  $(2\mathcal{E})$  with respect to  $\Psi^{-1}$  as a vector is:

$$H_{\Psi^{-1}} = \frac{\partial \text{vec}(2 \nabla_{\Psi^{-1}} \mathcal{E})}{\partial \text{vec}(\Psi^{-1})^T} = \frac{\partial \text{vec}(-\Psi + E)}{\partial \text{vec}(\Psi^{-1})^T} = \Psi \otimes \Psi. \quad (110)$$

Since the Hessian is positive definite, the E-step objective  $\mathcal{E}$  is convex in  $\Psi^{-1}$ , which is the first statement of the theorem.

The inverse of  $H_{\Psi^{-1}}$  is

$$H_{\Psi^{-1}}^{-1} = \Psi^{-1} \otimes \Psi^{-1}. \quad (111)$$

For the product of the inverse Hessian with the gradient we have:

$$\begin{aligned} H_{\Psi^{-1}}^{-1} \text{vec}(-\Psi + E) &= \text{vec}(\Psi^{-1} (-\Psi + E) \Psi^{-1}) \\ &= \text{vec}(-\Psi^{-1} + \Psi^{-1} E \Psi^{-1}). \end{aligned} \quad (112)$$

If we apply a Newton update, then the update direction for  $\Psi^{-1}$  in the M-step is

$$\Delta\Psi^{-1} = \Psi^{-1} - \Psi^{-1} E \Psi^{-1}. \quad (113)$$

We now can approximate the update for  $\Psi$  by the first terms of the Taylor expansion:

$$\Psi + \Delta\Psi = (\Psi^{-1} + \Delta\Psi^{-1})^{-1} \approx \Psi - \Psi \Delta\Psi^{-1} \Psi. \quad (114)$$

We obtain for the update of  $\Psi$

$$\Delta\Psi = -\Psi \Delta\Psi^{-1} \Psi = E - \Psi. \quad (115)$$

This is the exact EM update if the step-size  $\eta$  is 1.  $\square$

The Newton update derived from  $\Psi^{-1}$  as parameter is the Newton update for  $\Psi$ . Consequently, the Newton direction for both  $\Psi$  and  $\Psi^{-1}$  is in the M-step

$$\Delta\Psi = E - \Psi. \quad (116)$$

## 9 Gradient-based E-Step

### 9.1 Motivation for Rectifying and Normalization Constraints

The representation of data vector  $v$  by the model is the variational mean vector  $\mu_q$ . In order to obtain sparse codes we want to have non-negative  $\mu_q$ . We enforce non-negative mean values by constraints and optimize by projected Newton methods and by gradient projection methods. Non-negative constraints correspond to rectifying in the neural network field. Therefore we aim to construct sparse codes in analogy to the rectified linear units used for neural networks.

We constrain the variational distributions to the family of normal distributions with non-negative mean components. Consequently we introduce non-negative or **rectifying constraints**:

$$\mu \geq 0, \quad (117)$$

where the inequality “ $\geq$ ” holds component-wise.

However generative models with many coding units face a problem. They tend to *explain away small and rare signals by noise*. For many coding units, model selection algorithms prefer models with coding units which do not have variation and, therefore, are removed from the model. Other coding units hardly contribute to explain the observations. The likelihood is larger if small and rare signals are explained by noise, than the likelihood if coding units are used to explain such signals. Coding units without variance are kept on their default values, where they have maximal contribution to the likelihood. If they are used for coding, they deviate from their maximal values for each sample. In accumulation these deviations decrease the likelihood more than it is increased by explaining small or rare signals. For our RFN models the problem can become severe, since we aim at models with up to several tens of thousands of coding units. To avoid the explaining away problem, we enforce the selected models to use all their coding units on an equal level. We do that by keeping the variation of each noise-free coding unit across the training set at one. Consequently, we introduce a **normalization constraint** for each coding unit  $1 \leq j \leq l$ :

$$\frac{1}{n} \sum_{i=1}^n \mu_{ij}^2 = 1. \quad (118)$$

This constraint means that the noise-free part of each coding unit has variance one across samples.

We will derive methods to increase the objective in the E-step both for only rectifying constraints and for rectifying and normalization constraints. These methods ensure to reduce the objective in the E-step to guarantee convergence via the GAM theory. The resulting model from the GAM procedure is at a local maximum of the objective given the model family and the family of variational distributions. *The solution minimizes the KL-distance between the family of full variational distributions and full model family.* “Full” means that both the observed and the hidden variables are taken into account.

## 9.2 The Full E-step Objective

The E-step maximizes  $\mathcal{F}$  with respect to the variational distribution  $Q$ , therefore the E-step minimizes the Kullback-Leibler divergence (KL-divergence) Kullback and Leibler [1951]  $D_{\text{KL}}(Q(\mathbf{h}) \parallel p(\mathbf{h} \mid \mathbf{v}))$ . The KL-divergence between  $Q$  and  $p$  is

$$D_{\text{KL}}(Q \parallel p) = \int Q(\mathbf{h}) \log \frac{Q(\mathbf{h})}{p(\mathbf{h} \mid \mathbf{v})} d\mathbf{h}. \quad (119)$$

*Rectifying constraints* introduce non-negative constraints. The minimization with respect to  $Q(\mathbf{h}_i)$  gives the constraint minimization problem:

$$\begin{aligned} \min_{Q(\mathbf{h}_i)} \quad & \frac{1}{n} \sum_{i=1}^n D_{\text{KL}}(Q(\mathbf{h}_i) \parallel p(\mathbf{h}_i \mid \mathbf{v}_i)) \\ \text{s.t.} \quad & \forall_i : \boldsymbol{\mu}_i \geq \mathbf{0}, \end{aligned} \quad (120)$$

where  $\boldsymbol{\mu}_i$  is the mean vector of  $Q(\mathbf{h}_i)$ .

*Rectifying and normalizing constraints* introduce non-negative constraints and equality constraints. The minimization with respect to  $Q(\mathbf{h}_i)$  gives the constraint minimization problem:

$$\begin{aligned} \min_{Q(\mathbf{h}_i)} \quad & \frac{1}{n} \sum_{i=1}^n D_{\text{KL}}(Q(\mathbf{h}_i) \parallel p(\mathbf{h}_i \mid \mathbf{v}_i)) \\ \text{s.t.} \quad & \forall_i : \boldsymbol{\mu}_i \geq \mathbf{0}, \\ & \forall_j : \frac{1}{n} \sum_{i=1}^n \mu_{ij}^2 = 1, \end{aligned} \quad (121)$$

where  $\boldsymbol{\mu}_i$  is the mean vector of  $Q(\mathbf{h}_i)$ .

First we consider the families from which the model and from which the variational distributions stem. The posterior of the model with Gaussian prior  $p(\mathbf{h})$  is Gaussian (see Section 5):

$$p(\mathbf{h} \mid \mathbf{v}) \sim (2\pi)^{-\frac{l}{2}} |\boldsymbol{\Sigma}_p|^{-\frac{1}{2}} \exp \left( -\frac{1}{2} (\mathbf{h} - \boldsymbol{\mu}_p)^T \boldsymbol{\Sigma}_p^{-1} (\mathbf{h} - \boldsymbol{\mu}_p) \right). \quad (122)$$

To be as close as possible to the posterior distribution, we restrict  $Q$  to be from a Gaussian family:

$$Q(\mathbf{h}) \sim (2\pi)^{-\frac{l}{2}} |\Sigma_q|^{-\frac{1}{2}} \exp\left(-\frac{1}{2} (\mathbf{h} - \boldsymbol{\mu}_q)^T \Sigma_q^{-1} (\mathbf{h} - \boldsymbol{\mu}_q)\right). \quad (123)$$

For Gaussians, the Kullback-Leibler divergence between  $Q$  and  $p$  is

$$D_{\text{KL}}(Q \parallel p) = \frac{1}{2} \left\{ \text{Tr}(\Sigma_p^{-1} \Sigma_q) + (\boldsymbol{\mu}_p - \boldsymbol{\mu}_q)^T \Sigma_p^{-1} (\boldsymbol{\mu}_p - \boldsymbol{\mu}_q) - l - \ln \frac{|\Sigma_q|}{|\Sigma_p|} \right\}. \quad (124)$$

This Kullback-Leibler divergence is convex in the mean vector  $\boldsymbol{\mu}_q$  and the covariance matrix  $\Sigma_q$  of  $Q$ , simultaneously Dredze et al. [2008, 2012].

We now minimize Eq. (124) with respect to  $Q$ . For the moment we do not care about the constraints introduced by non-negativity and by normalization. Eq. (124) has a quadratic form in  $\boldsymbol{\mu}_q$ , where  $\Sigma_q$  does not enter, and terms in  $\Sigma_q$ , where  $\boldsymbol{\mu}_q$  does not enter. Therefore we can separately minimize for  $\Sigma_q$  and for  $\boldsymbol{\mu}_q$ .

For the minimization with respect to  $\Sigma_q$ , we require

$$\frac{\partial}{\partial \Sigma_q} \text{Tr}(\Sigma_p^{-1} \Sigma_q) = \Sigma_p^{-T} \quad (125)$$

and

$$\frac{\partial}{\partial \Sigma_q} \ln |\Sigma_q| = \Sigma_q^{-T}. \quad (126)$$

For optimality the derivative of the objective  $D_{\text{KL}}(Q \parallel p)$  with respect to  $\Sigma_q$  must be zero:

$$\frac{\partial}{\partial \Sigma_q} D_{\text{KL}}(Q \parallel p) = \frac{1}{2} \Sigma_p^{-T} - \frac{1}{2} \Sigma_q^{-T} = \mathbf{0}. \quad (127)$$

This gives

$$\Sigma = \Sigma_q = \Sigma_p. \quad (128)$$

We often drop the index  $q$  since for  $1 \leq i \leq n$  all covariance matrices  $\Sigma_q$  are equal to  $\Sigma_p$ .

The mean vector  $\boldsymbol{\mu}_q$  of  $Q$  is the solution of the minimization problem:

$$\min_{\boldsymbol{\mu}} \frac{1}{2} (\boldsymbol{\mu}_p - \boldsymbol{\mu})^T \Sigma_p^{-1} (\boldsymbol{\mu}_p - \boldsymbol{\mu}) \quad (129)$$

which is equivalent to

$$\min_{\boldsymbol{\mu}} \frac{1}{2} \boldsymbol{\mu}^T \Sigma_p^{-1} \boldsymbol{\mu} - \boldsymbol{\mu}_p^T \Sigma_p^{-1} \boldsymbol{\mu}. \quad (130)$$

The derivative and the Hessian of this objective is:

$$\frac{\partial}{\partial \boldsymbol{\mu}} D_{\text{KL}}(Q \parallel p) = \Sigma_p^{-1} (\boldsymbol{\mu} - \boldsymbol{\mu}_p), \quad (131)$$

$$\frac{\partial^2}{\partial^2 \boldsymbol{\mu}} D_{\text{KL}}(Q \parallel p) = \Sigma_p^{-1}. \quad (132)$$

### 9.3 E-step for Mean with Rectifying Constraints

#### 9.3.1 The E-Step Minimization Problem

Rectifying is realized by non-negative constraints. The mean vector  $\boldsymbol{\mu}_q$  of  $Q$  is the solution of the minimization problem:

$$\begin{aligned} \min_{\boldsymbol{\mu}} \quad & \frac{1}{2} (\boldsymbol{\mu} - \boldsymbol{\mu}_p)^T \Sigma_p^{-1} (\boldsymbol{\mu} - \boldsymbol{\mu}_p) \\ \text{s.t.} \quad & \boldsymbol{\mu} \geq \mathbf{0}. \end{aligned} \quad (133)$$

This is a convex quadratic minimization problem with non-negativity constraints (convex feasible set).

If  $\lambda$  is the Lagrange multiplier for the constraints, then the dual is

$$\begin{aligned} \min_{\lambda} \quad & \frac{1}{2} \lambda^T \Sigma_p \lambda + \mu_p^T \lambda \\ \text{s.t.} \quad & \lambda \geq \mathbf{0} . \end{aligned} \quad (134)$$

The Karush-Kuhn-Tucker conditions require for the optimal solution for each component  $1 \leq j \leq l$ :

$$\lambda_j \mu_j = 0 . \quad (135)$$

Further the derivative of the Lagrangian with respect to  $\mu$  gives

$$\Sigma_p^{-1} \mu - \Sigma_p^{-1} \mu_p - \lambda = \mathbf{0} \quad (136)$$

which can be written as

$$\mu - \mu_p - \Sigma_p \lambda = \mathbf{0} . \quad (137)$$

This minimization problem cannot be solved directly. Therefore we perform a gradient projection or projected Newton step to decrease the objective.

### 9.3.2 The Projection onto the Feasible Set

To decrease the objective, we perform a gradient projection or a projected Newton step. We will base our algorithms on *Euclidean least distance projections*. If projected onto convex sets, these projections do not increase distances. The Euclidean projection onto the feasible set is denoted by  $P$ , that is, the map that takes  $\mu_p$  to its nearest point  $\mu$  (in the  $L^2$ -norm) in the feasible set.

For rectifying constraints, the projection  $P$  (Euclidean least distance projection) of  $\mu_p$  onto the convex feasible set is given by the solution of the convex optimization problem:

$$\begin{aligned} \min_{\mu} \quad & \frac{1}{2} (\mu - \mu_p)^T (\mu - \mu_p) \\ \text{s.t.} \quad & \mu \geq \mathbf{0} . \end{aligned} \quad (138)$$

The following Theorem 8 shows that update Eq. (139) is the projection  $P$  defined by optimization problem Eq. (138).

**Theorem 8** (Projection: Rectifying). *The solution to optimization problem Eq. (138), which defines the Euclidean least distance projection, is*

$$\mu_j = [P(\mu_p)]_j = \begin{cases} 0 & \text{for } (\mu_p)_j \leq 0 \\ (\mu_p)_j & \text{for } (\mu_p)_j > 0 \end{cases} \quad (139)$$

*Proof.* For the projection we have the minimization problem:

$$\begin{aligned} \min_{\mu} \quad & \frac{1}{2} (\mu - \mu_p)^T (\mu - \mu_p) \\ \text{s.t.} \quad & \mu \geq \mathbf{0} . \end{aligned} \quad (140)$$

The Lagrangian  $L$  with multiplier  $\lambda \geq \mathbf{0}$  is

$$L = \frac{1}{2} (\mu - \mu_p)^T (\mu - \mu_p) - \lambda^T \mu . \quad (141)$$

The derivative with respect to  $\mu$  is

$$\frac{\partial L}{\partial \mu} = \mu - \mu_p - \lambda = \mathbf{0} . \quad (142)$$

The Karush-Kuhn-Tucker (KKT) conditions require for the optimal solution that for each constraint  $j$ :

$$\lambda_j \mu_j = 0 . \quad (143)$$

If  $0 < (\mu_p)_j$  then Eq. (142) requires  $0 < \mu_j$  because the Lagrangian  $\lambda_j$  is larger than or equal to zero:  $0 \leq \lambda_j$ . From the KKT conditions Eq. (143) follows that  $\lambda_j = 0$  and, therefore,  $0 < \mu_j = (\mu_p)_j$ . If  $(\mu_p)_j < 0$  then  $0 < \mu_j - (\mu_p)_j$ , because the constraints of the primal problem require  $0 \leq \mu_j$ . From Eq. (142) follows that  $0 < \lambda_j$ . From the KKT conditions Eq. (143) follows that  $(\mu_p)_j = 0$  and  $0 < \lambda_j = -(\mu_p)_j$ . If  $(\mu_p)_j = 0$ , then Eq. (142) and the KKT conditions Eq. (143) lead to  $(\mu_p)_j = \mu_j = \lambda_j = 0$ .

Therefore the solution of problem Eq. (138) is

$$\mu_j = \begin{cases} (\mu_p)_j & \text{for } (\mu_p)_j > 0 \text{ and } \lambda_j = 0 \\ 0 & \text{for } (\mu_p)_j \leq 0 \text{ and } \lambda_j = -(\mu_p)_j \end{cases} . \quad (144)$$

This finishes the proof.  $\square$

## 9.4 E-step for Mean with Rectifying and Normalizing Constraints

### 9.4.1 The E-Step Minimization Problem

If we also consider normalizing constraints, then we have to minimize all KL-divergences simultaneously. The normalizing constraints connect the single optimization problems for each sample  $\mathbf{v}_i$ . For the E-step, we obtain the minimization problem:

$$\begin{aligned} \min_{\boldsymbol{\mu}_i} \quad & \frac{1}{n} \sum_{i=1}^n (\boldsymbol{\mu}_i - (\boldsymbol{\mu}_p)_i)^T \boldsymbol{\Sigma}_p^{-1} (\boldsymbol{\mu}_i - (\boldsymbol{\mu}_p)_i) \\ \text{s.t.} \quad & \forall_i : \boldsymbol{\mu}_i \geq \mathbf{0} \quad , \quad \forall_j : \frac{1}{n} \sum_{i=1}^n \mu_{ij}^2 = 1 . \end{aligned} \quad (145)$$

The “ $\geq$ ”-sign is meant component-wise. The  $l$  equality constraints lead to non-convex feasible sets. The solution to this optimization problem are the means vectors  $\boldsymbol{\mu}_i$  of  $Q(\mathbf{h}_i)$ .

**Generalized Reduced Gradient.** The equality constraints can be solved for one variable which is then inserted into the objective. The equality constraint gives for each  $1 \leq j \leq l$ :

$$\mu_{1j}^2 = n - \sum_{i=2}^n \mu_{ij}^2 \quad \text{or} \quad \mu_{1j} = \sqrt{n - \sum_{i=2}^n \mu_{ij}^2} . \quad (146)$$

These equations can be inserted into the objective and, thereby, we remove the variables  $\mu_{1j}$ . We have to ensure that the  $\mu_{1j}$  exist by

$$\sum_{i=2}^n \mu_{ij}^2 \leq n . \quad (147)$$

These constraints define a convex set feasible set. To solve the each equality constraints for a variable and insert it into the objective is called *generalized reduced gradient* method Abadie and Carpentier [1969]. For solving the reduced problem, we can use methods for constraint optimization were we now ensure a convex feasible set. These methods solve the original problem Eq. (145). We only require an improvement of the objective with a feasible value. For the reduced problem, we perform one step of a *gradient projection method*.

**Gradient Projection Methods.** Also for the original problem Eq. (145), *gradient projection methods* can be used. The gradient projection method has been generalized by Rosen to *non-linear constraints* Rosen [1961] and was later improved by Haug and Arora [1979]. The gradient projection algorithm of Rosen works for *non-convex feasible sets*. The idea is to linearize the nonlinear constraints and solve the problem. Subsequently a restoration move brings the solution back to the constraint boundaries.

#### 9.4.2 The Projection onto the Feasible Set

To decrease the objective, we perform a gradient projection, a projected Newton step, or a step of the generalized reduced method. We will base our algorithms on *Euclidean least distance projections*. If projected onto convex sets, these projections do not increase distances. The Euclidean projection onto the feasible set is denoted by  $P$ , that is, the map that simultaneously takes  $\{(\mu_p)_i\}$  to the nearest points  $\{\mu_i\}$  (in the  $L^2$ -norm) in the feasible set.

For rectifying and normalizing constraints the projection (Euclidean least distance projection) of  $\{(\mu_p)_i\}$  onto the **non-convex** feasible set leads to the optimization problem

$$\begin{aligned} \min_{\mu_i} \quad & \frac{1}{n} \sum_{i=1}^n (\mu_i - (\mu_p)_i)^T (\mu_i - (\mu_p)_i) \\ \text{s.t.} \quad & \forall_i : \mu_i \geq \mathbf{0}, \\ & \forall_j : \frac{1}{n} \sum_{i=1}^n \mu_{ij}^2 = 1. \end{aligned} \quad (148)$$

By using  $(\mu_i - (\mu_p)_i)^T (\mu_i - (\mu_p)_i) = \mu_i^T \mu_i - 2\mu_i^T (\mu_p)_i + (\mu_p)_i^T (\mu_p)_i$ , we see that the objective contains the sum  $\sum_{ij} \mu_{ij}^2$ . The constraints enforce this sum to be constant. Therefore inserting the equality constraints into the objective, optimization problem Eq. (148) is equivalent to

$$\begin{aligned} \min_{\mu_i} \quad & -\frac{1}{n} \sum_{i=1}^n \mu_i^T (\mu_p)_i \\ \text{s.t.} \quad & \forall_i : \mu_i \geq \mathbf{0}, \\ & \forall_j : \frac{1}{n} \sum_{i=1}^n \mu_{ij}^2 = 1. \end{aligned} \quad (149)$$

The following Theorem 9 shows that updates Eq. (150) and Eq. (151) form the projection defined by optimization problem Eq. (148).

**Theorem 9** (Projection: Rectifying and Normalizing). *If at least one  $(\mu_p)_{ij}$  is positive for  $1 \leq j \leq l$ , then the solution to optimization problem Eq. (148), which defines the Euclidean least distance projection, is*

$$\begin{aligned} \hat{\mu}_{ij} &= \begin{cases} 0 & \text{for } (\mu_p)_{ij} \leq 0 \\ (\mu_p)_{ij} & \text{for } (\mu_p)_{ij} > 0 \end{cases} \\ \mu_{ij} &= [P((\mu_p)_i)]_j = \frac{\hat{\mu}_{ij}}{\sqrt{\frac{1}{n} \sum_{i=1}^n \hat{\mu}_{ij}^2}}. \end{aligned} \quad (150)$$

*If all  $(\mu_p)_{ij}$  are non-positive for  $1 \leq j \leq l$ , then the optimization problem Eq. (148) has the solution*

$$\mu_{ij} = \begin{cases} \sqrt{n} & \text{for } j = \arg \max_j \{(\mu_p)_{ij}\} \\ 0 & \text{otherwise} \end{cases}. \quad (151)$$

*Proof.* In the following we show that updates Eq. (150) and Eq. (151) are the projection onto the feasible set. For the projection of  $\{(\mu_p)_i\}$  onto the feasible set, we have the minimization problem:

$$\begin{aligned} \min_{\mu_i} \quad & \frac{1}{n} \sum_{i=1}^n (\mu_i - (\mu_p)_i)^T (\mu_i - (\mu_p)_i) \\ \text{s.t.} \quad & \forall_i : \mu_i \geq \mathbf{0}, \\ & \forall_j : \frac{1}{n} \sum_{i=1}^n \mu_{ij}^2 = 1. \end{aligned} \quad (152)$$

The feasible set is non-convex because of the quadratic equality constraint. The Lagrangian with multiplier  $\lambda \geq 0$  is

$$L = \frac{1}{n} \sum_{i=1}^n (\mu_i - (\mu_p)_i)^T (\mu_i - (\mu_p)_i) - \sum_{i=1}^n \lambda_i^T \mu_i + \sum_j \tau_j \left( \frac{1}{n} \sum_{i=1}^n \mu_{ij}^2 - 1 \right). \quad (153)$$

The Karush-Kuhn-Tucker (KKT) conditions require for the optimal solution:

$$\lambda_{ij} \mu_{ij} = 0 \quad \text{and} \quad \tau_j \left( \frac{1}{n} \sum_{i=1}^n \mu_{ij}^2 - 1 \right) = 0. \quad (154)$$

The derivative of  $L$  with respect to  $\mu_{ij}$  is

$$\frac{\partial L}{\partial \mu_{ij}} = \frac{2}{n} (\mu_{ij} - (\mu_p)_{ij}) - \lambda_{ij} + \frac{2}{n} \tau_j \mu_{ij} = 0. \quad (155)$$

We multiply this equation by  $\mu_{ij}$  and obtain:

$$\frac{2}{n} (\mu_{ij}^2 - (\mu_p)_{ij} \mu_{ij}) - \lambda_{ij} \mu_{ij} + \frac{2}{n} \tau_j \mu_{ij}^2 = 0. \quad (156)$$

The KKT conditions give  $\lambda_{ij} \mu_{ij} = 0$ , therefore this term can be removed from the equation. Next we sum over  $i$ :

$$\frac{2}{n} \sum_{i=1}^n (\mu_{ij}^2 - (\mu_p)_{ij} \mu_{ij}) + \frac{2}{n} \sum_{i=1}^n \tau_j \mu_{ij}^2 = 0. \quad (157)$$

Using the equality constraint  $1/n \sum_{i=1}^n \mu_{ij}^2 = 1$  and dividing by 2 and gives:

$$1 - \frac{1}{n} \sum_{i=1}^n (\mu_p)_{ij} \mu_{ij} + \tau_j = 0. \quad (158)$$

Solving for  $\tau_j$  leads to:

$$\tau_j = \frac{1}{n} \sum_{i=1}^n (\mu_p)_{ij} \mu_{ij} - 1. \quad (159)$$

We insert  $\tau_j$  into Eq. (155)

$$-(\mu_p)_{ij} - \frac{n}{2} \lambda_{ij} + \left( \frac{1}{n} \sum_{s=1}^n (\mu_p)_{sj} \mu_{sj} \right) \mu_{ij} = 0. \quad (160)$$

We immediately see, that if  $\mu_{ij} = 0$  then  $(\mu_p)_{ij} = -\frac{n}{2} \lambda_{ij} < 0$ . Therefore we can assume  $\mu_{ij} > 0$ . Multiplying Eq. (160) with  $\mu_{ij}$  and using the KKT conditions gives

$$-(\mu_p)_{ij} \mu_{ij} + \left( \frac{1}{n} \sum_{s=1}^n (\mu_p)_{sj} \mu_{sj} \right) \mu_{ij}^2 = 0. \quad (161)$$

Therefore  $(\mu_p)_{ij} \mu_{ij}$  and  $\frac{1}{n} \sum_{s=1}^n (\mu_p)_{sj} \mu_{sj}$  have the same sign or  $\mu_{ij} = 0$ . Since  $0 \leq \mu_{ij}$ , we deduce that  $(\mu_p)_{ij}$  and  $\frac{1}{n} \sum_{s=1}^n (\mu_p)_{sj} \mu_{sj}$  have the same sign or  $\mu_{ij} = 0$ . Since the sum is independent of  $i$ , all  $(\mu_p)_{ij}$  with  $\mu_{ij} > 0$  have the same sign for  $1 \leq i \leq n$ . Solving Eq. (160) for  $\mu_{ij}$  gives

$$\mu_{ij} = \frac{(\mu_p)_{ij} + \frac{n}{2} \lambda_{ij}}{\frac{1}{n} \sum_{s=1}^n (\mu_p)_{sj} \mu_{sj}}. \quad (162)$$



**I.** If all  $(\mu_p)_{ij}$  are non-positive for  $1 \leq j \leq l$ , then the sum  $\frac{1}{n} \sum_{s=1}^n (\mu_p)_{sj} \mu_{sj}$  is negative. From the first order derivative of the Lagrangian in Eq. (155), we can compute the second order derivative

$$\frac{\partial^2 L}{\partial \mu_{ij} \partial \mu_{ij}} = \frac{2}{n} + \frac{2}{n} \tau_j = 2 \sum_{i=1}^n (\mu_p)_{ij} \mu_{ij} < 0. \quad (163)$$

We inserted the expression of Eq. (159) for  $\tau_j$ . Since all mixed second order derivatives are zero, the (projected) Hessian of the Lagrangian is diagonal with negative entries. Therefore it is strict negative definite. Thus, the second order necessary conditions cannot be fulfilled. The minimum is a border point of the constraints.

For each  $j$  for which all  $(\mu_p)_{ij}$  are non-positive for  $1 \leq j \leq l$ , optimization problem Eq. (149) defines a plane that has a normal vector in the positive orthant (hyperoctant). For such a  $j$  the corresponding equality constraint defines a hypersphere. Minimization means that the plane containing the solution is parallel to the original plane and should be as close to the origin as possible. If we move the plane parallel from the origin into the positive orthant, then the first intersection with the hypersphere is

$$\mu_{ij} = \begin{cases} \sqrt{n} & \text{for } j = \arg \max_j \{(\mu_p)_{ij}\} \\ 0 & \text{otherwise} \end{cases}. \quad (164)$$

This is the solution for  $\mu_{ij}$  with  $1 \leq j \leq l$  to our minimization problem.

**II.** If one  $(\mu_p)_{ij}$  is positive, then from Eq. (160) with this  $(\mu_p)_{ij}$  follows that  $\frac{1}{n} \sum_{s=1}^n (\mu_p)_{sj} \mu_{sj}$  is positive, otherwise Eq. (160) has only negative terms on the left hand side. In particular, the second order necessary conditions are always fulfilled as Eq. (163) is positive. For  $(\mu_p)_{ij} < 0$  it follows from Eq. (160) that  $\lambda_{ij} > 0$  and from the KKT conditions that  $\mu_{ij} = 0$ . For  $(\mu_p)_{ij} > 0$  it follows from Eq. (160) that  $\mu_{ij} > 0$  and from the KKT conditions that  $\lambda_{ij} = 0$ . Therefore we define:

$$\hat{\mu}_{ij} = \begin{cases} 0 & \text{for } (\mu_p)_{ij} \leq 0 \\ (\mu_p)_{ij} & \text{for } (\mu_p)_{ij} > 0 \end{cases}, \quad (165)$$

We write the solution as

$$\mu_{ij} = \frac{\hat{\mu}_{ij}}{\frac{1}{n} \sum_{s=1}^n (\mu_p)_{sj} \mu_{sj}} = \alpha_j \hat{\mu}_{ij}. \quad (166)$$

We now use the equality constraint:

$$\frac{1}{n} \sum_{i=1}^n \mu_{ij}^2 = \alpha_j^2 \frac{1}{n} \sum_{i=1}^n \hat{\mu}_{ij}^2 = 1. \quad (167)$$

Solving for  $\alpha_j$  gives:

$$\alpha_j = \frac{1}{\sqrt{\frac{1}{n} \sum_{i=1}^n \hat{\mu}_{ij}^2}}. \quad (168)$$

Therefore the solution is

$$\mu_{ij} = \frac{\hat{\mu}_{ij}}{\sqrt{\frac{1}{n} \sum_{i=1}^n \hat{\mu}_{ij}^2}}. \quad (169)$$

This finishes the proof. □

## 9.5 Gradient and Scaled Gradient Projection and Projected Newton

### 9.5.1 Gradient Projection Algorithm

The *projected gradient descent* or *gradient projection algorithm* Bertsekas [1976], Kelley [1999] performs first a gradient step and then projects the result to the *feasible set*. The projection onto the

feasible set is denoted by  $P$ , that is, the map that takes  $\mu$  into the nearest point (in the  $L^2$ -norm) in the feasible set to  $\mu$ . The feasible set must be convex, however later we will introduce gradient projection methods for non-convex feasible sets.

The gradient projection method is in our case

$$\mu_{k+1} = P(\mu_k + \lambda \Sigma_p^{-1}(\mu_p - \mu_k)) . \quad (170)$$

The Lipschitz constant for the gradient is  $\|\Sigma_p^{-1}\|_s = e_{\max}(\Sigma_p^{-1})$ , the largest eigenvalue of  $\Sigma_p^{-1}$ . The following statement is Theorem 5.4.5 in Kelley [1999].

**Theorem 10** (Theorem 5.4.5 in Kelley [1999]). *The sufficient decrease condition*

$$D_{\text{KL}}(Q(\mu_{k+1}) \parallel p) - D_{\text{KL}}(Q(\mu_k) \parallel p) \leq \frac{-\alpha}{\lambda} \|\mu_k - \mu_{k+1}\|^2 \quad (171)$$

(e.g. with  $\alpha = 10^{-4}$ ) holds for all  $\lambda$  such that

$$0 < \lambda \leq \frac{2(1 - \alpha)}{e_{\max}(\Sigma_p^{-1})} . \quad (172)$$

*Proof.* See Kelley [1999]. □

*Theorem 10 guarantees that we can increase the objective by gradient projection in the E-step, except the case where we already reached the maximum.*

For a fast upper bound on the maximal eigenvalue we use

$$e_{\max}(\Sigma_p^{-1}) \leq \text{Tr}(\Sigma_p^{-1}) \quad (173)$$

and

$$e_{\max}(\Sigma_p^{-1}) \leq \|\mathbf{W}\|_s^2 \|\Psi^{-1}\|_s - 1 , \quad (174)$$

where the latter follows from

$$\Sigma_p^{-1} = \mathbf{I} + \mathbf{W}^T \Psi^{-1} \mathbf{W} . \quad (175)$$

Improved methods for finding an appropriate  $\lambda$  by line search methods have been proposed Birgin et al. [2000], Serafini et al. [2005]. We use a search with  $\lambda = \beta^t$  with  $t = 0, 1, 2, \dots$  and  $\beta = 2^{-1}$  or  $\beta = 10^{-1}$ .

A special version of the gradient projection method is the *generalized reduced method* Abadie and Carpentier [1969]. This method is able to solve our optimization problem with equality constraints. The gradient projection method has been generalized by Rosen to non-linear constraints Rosen [1961]. The gradient projection algorithm of Rosen can also be used for a region which is not convex. The idea is to linearize the nonlinear constraints and solve the problem. Subsequently a restoration move brings the solution back to the constraint boundaries. Rosen's gradient projection method was improved by Haug and Arora [1979]. *These methods guarantee that we can increase the objective in the E-step for non-convex feasible sets, except the case where we already reached the maximum.* These algorithms for non-convex feasible sets will only give a local maximum. Also the GAM algorithm will only find a local maximum.

### 9.5.2 Scaled Gradient Projection and Projected Newton Method

Both the *scaled gradient projection algorithm* and the *projected Newton method* were proposed in Bertsekas [1982]. We follow Kelley [1999].

The idea is to use a Newton update instead of the a gradient update:

$$\mu_{k+1} = P(\mu_k + \lambda \mathbf{H}^{-1} \Sigma_p^{-1}(\mu_p - \mu_k)) . \quad (176)$$

$\mathbf{H}^{-1}$  can be an arbitrary strict positive definite matrix. If we set  $\mathbf{H}^{-1} = \Sigma_p$ , then we have a Newton update of the *projected Newton method* Bertsekas [1982]. For  $\lambda = 1$  we obtain

$$\mu_{k+1} = P(\mu_p) . \quad (177)$$

otherwise

$$\boldsymbol{\mu}_{k+1} = \mathbf{P}((1 - \lambda)\boldsymbol{\mu}_k + \lambda\boldsymbol{\mu}_p) . \quad (178)$$

The search direction for the unconstrained problem can be rotated by  $\mathbf{H}^{-1}$  to be orthogonal to the direction of decrease in the inactive directions for the constrained problem.

To escape this possible problem, an  $\epsilon$ -active set is introduced which contains all  $j$  with  $\mu_j \leq \epsilon$ . All columns and rows of the Hessian having an index in the  $\epsilon$ -active set are fixed to  $\mathbf{e}_j$ . After sorting the indices of the  $\epsilon$ -active set together, they form a block which is the sub-identity matrix.  $\mathbf{H}$  is set to the Hessian  $\boldsymbol{\Sigma}_p$  where the  $\epsilon$ -active set columns and rows are replaced by unit vectors.

The following Theorem 11 is Lemma 5.5.1 in Kelley [1999]. *Theorem 11 states that the objective decreases using the reduced Hessian in the projected Newton method for convex feasible sets.*

**Theorem 11** (Lemma 5.5.1 in Kelley [1999]). *The sufficient decrease condition*

$$D_{\text{KL}}(Q(\boldsymbol{\mu}_{k+1}) \parallel p) - D_{\text{KL}}(Q(\boldsymbol{\mu}_k) \parallel p) \leq -\alpha (\boldsymbol{\mu}_k - \boldsymbol{\mu}_p)^T \boldsymbol{\Sigma}_p^{-1} (\boldsymbol{\mu}_k - \boldsymbol{\mu}_{k+1}) \quad (179)$$

holds for all  $\lambda$  smaller than a bound depending on  $\mathbf{H}$  and  $\epsilon$ .

*Proof.* See Kelley [1999]. □

In practical applications, a proper  $\lambda$  is found by line search. The *projected Newton method* uses  $\lambda = 1$  to set  $\epsilon$  Bertsekas [1982]:

$$\epsilon = \|\boldsymbol{\mu}_k - \mathbf{P}(\boldsymbol{\mu}_p)\| . \quad (180)$$

### 9.5.3 Combined Method

Following Kim et al. [2006], Serafini et al. [2005] we use the following very general update rule, which includes the gradient projection algorithm, the scaled gradient projection algorithm, and the projected Newton method.

We use following update for the E-step:

$$\begin{aligned} \mathbf{d}_{k+1} &= \mathbf{P}(\boldsymbol{\mu}_k + \lambda \mathbf{H}^{-1} \boldsymbol{\Sigma}_p^{-1} (\boldsymbol{\mu}_p - \boldsymbol{\mu}_k)) , \\ \boldsymbol{\mu}_{k+1} &= \mathbf{P}(\boldsymbol{\mu}_k + \gamma (\mathbf{d}_{k+1} - \boldsymbol{\mu}_k)) . \end{aligned} \quad (181)$$

We have to project twice since the equality constraint produces a manifold in the parameter space.

We iterate this update until we see a decrease of the objective in the E-step:

$$D_{\text{KL}}(Q_{k+1} \parallel p) - D_{\text{KL}}(Q_k \parallel p) < 0 . \quad (182)$$

For the constraints we have only to optimize the mean vector  $\boldsymbol{\mu}$  to ensure

$$D_{\text{KL}}(Q(\boldsymbol{\mu}_{k+1}) \parallel p) - D_{\text{KL}}(Q(\boldsymbol{\mu}_k) \parallel p) < 0 . \quad (183)$$

Even

$$D_{\text{KL}}(Q(\boldsymbol{\mu}_{k+1}) \parallel p) = D_{\text{KL}}(Q(\boldsymbol{\mu}_k) \parallel p) \quad (184)$$

can be sufficient if minimizing  $\boldsymbol{\Sigma}_{k+1} = \boldsymbol{\Sigma}_p$  ensures

$$D_{\text{KL}}(Q_{k+1} \parallel p) < D_{\text{KL}}(Q_k \parallel p) . \quad (185)$$

We use following schedule:

1.
  - $\mathbf{H}^{-1} = \boldsymbol{\Sigma}_p$
  - $\lambda = 1$
  - $\gamma = 1$

That is

$$\boldsymbol{\mu}_{k+1} = \mathbf{P}(\boldsymbol{\mu}_p) . \quad (186)$$

2.
  - $\mathbf{H}^{-1} = \Sigma_p$
  - $\lambda = 1$
  - $\gamma \in (0, 1]$

That is

$$\boldsymbol{\mu}_{k+1} = \mathbf{P}((1 - \gamma) \boldsymbol{\mu}_k + \gamma \mathbf{P}(\boldsymbol{\mu}_p)) . \quad (187)$$

3.
  - $\mathbf{H}^{-1} = \Sigma_p$
  - $\lambda \in (0, 1]$
  - $\gamma = 1$

That is

$$\boldsymbol{\mu}_{k+1} = \mathbf{P}((1 - \lambda) \boldsymbol{\mu}_k + \lambda \boldsymbol{\mu}_p) . \quad (188)$$

4.
  - $\mathbf{H}^{-1} = \Sigma_p$
  - $\lambda \in (0, 1]$
  - $\gamma \in (0, 1]$

That is

$$\boldsymbol{\mu}_{k+1} = \mathbf{P}((1 - \gamma) \boldsymbol{\mu}_k + \gamma \mathbf{P}((1 - \lambda) \boldsymbol{\mu}_k + \lambda \boldsymbol{\mu}_p)) . \quad (189)$$

5.
  - $\mathbf{H}^{-1} = \mathbf{R}(\Sigma_p)$
  - $\lambda \in (0, 1]$
  - $\gamma \in (0, 1]$

$\mathbf{R}(\Sigma_p)$  denotes the reduced matrix (Hessian or a positive definite) according to the projected Newton method or the scaled gradient projection algorithm. For convex feasible sets we can guarantee at this level already an increase of the objective at the E-step.

6.
  - $\mathbf{H}^{-1} = \mathbf{I}$
  - $\lambda \in (0, 1]$
  - $\gamma \in (0, 1]$

This is the gradient projection algorithm. In particular we include the generalized reduced method and Rosen's gradient projection method. At this step we guarantee an increase of the objective at the E-step even for non-convex feasible sets because we also use complex methods for constraint optimization.

Step 5. ensures an improvement if only using rectifying constraints according to the theory of projected Newton methods Kelley [1999]. Step 6. ensures an improvement if using both rectifying constraints and normalizing constraints, because we use known methods for constraint optimization. To set  $\boldsymbol{\mu}_{k+1} = \boldsymbol{\mu}_k$  is sufficient to increase the objective at the E-step if  $\Sigma_{k+1} = \Sigma_p$  decreases the KL divergence. However we will not always set  $\boldsymbol{\mu}_{k+1} = \boldsymbol{\mu}_k$  to avoid accumulation points outside the solution set.

## 10 Alternative Gaussian Prior

We assume  $\mathbf{h}$  is Gaussian with covariance  $\mathbf{M}$  and mean  $\boldsymbol{\xi}$

$$\mathbf{h} \sim \mathcal{N}(\boldsymbol{\xi}, \mathbf{M}) . \quad (190)$$

We derive the posterior for this prior.

The likelihood is Gaussian since a affine transformation of a Gaussian random variable is again a Gaussian random variable and the convolution of two Gaussians is Gaussian, too. Thus,  $\mathbf{v} = \mathbf{W}\mathbf{h} + \boldsymbol{\epsilon}$  is Gaussian if  $\mathbf{h}$  and  $\boldsymbol{\epsilon}$  are both Gaussian. For the prior moments we have

$$\mathbf{E}(\mathbf{h}) = \boldsymbol{\xi} , \quad (191)$$

$$\mathbf{E}(\mathbf{h}\mathbf{h}^T) = \mathbf{M} + \boldsymbol{\xi} \boldsymbol{\xi}^T , \quad (192)$$

$$\text{var}(\mathbf{h}) = \mathbf{M} \quad (193)$$

and for the likelihood of  $\mathbf{v}$  we obtain the moments

$$\mathbb{E}(\mathbf{v}) = \mathbf{W}\boldsymbol{\xi}, \quad (194)$$

$$\begin{aligned} \mathbb{E}(\mathbf{v}\mathbf{v}^T) &= \mathbf{W} \mathbb{E}(\mathbf{h}\mathbf{h}^T) \mathbf{W}^T + \boldsymbol{\Psi} \\ &= \mathbf{W} \mathbf{M} \mathbf{W}^T + \boldsymbol{\Psi} + \mathbf{W} \boldsymbol{\xi} \boldsymbol{\xi}^T \mathbf{W}^T, \end{aligned} \quad (195)$$

$$\text{var}(\mathbf{v}) = \mathbf{W} \mathbf{M} \mathbf{W}^T + \boldsymbol{\Psi}. \quad (196)$$

We need some algebraic identities to derive the posterior. The Woodbury matrix identity gives

$$\mathbf{M} - \mathbf{M} \mathbf{W}^T (\mathbf{W} \mathbf{M} \mathbf{W}^T + \boldsymbol{\Psi})^{-1} \mathbf{W} \mathbf{M} = (\mathbf{M}^{-1} + \mathbf{W}^T \boldsymbol{\Psi}^{-1} \mathbf{W})^{-1}. \quad (197)$$

Multiplying this equation from the left hand side with  $\boldsymbol{\Psi}^{-1} \mathbf{W}$  gives

$$\begin{aligned} &\boldsymbol{\Psi}^{-1} \mathbf{W} (\mathbf{M}^{-1} + \mathbf{W}^T \boldsymbol{\Psi}^{-1} \mathbf{W})^{-1} \\ &= \boldsymbol{\Psi}^{-1} \mathbf{W} \mathbf{M} - \boldsymbol{\Psi}^{-1} \mathbf{W} \mathbf{M} \mathbf{W}^T (\mathbf{W} \mathbf{M} \mathbf{W}^T + \boldsymbol{\Psi})^{-1} \mathbf{W} \mathbf{M} \\ &= \boldsymbol{\Psi}^{-1} (\mathbf{W} \mathbf{M} \mathbf{W}^T + \boldsymbol{\Psi}) (\mathbf{W} \mathbf{M} \mathbf{W}^T + \boldsymbol{\Psi})^{-1} \mathbf{W} \mathbf{M} - \\ &\quad \boldsymbol{\Psi}^{-1} \mathbf{W} \mathbf{M} \mathbf{W}^T (\mathbf{W} \mathbf{M} \mathbf{W}^T + \boldsymbol{\Psi})^{-1} \mathbf{W} \mathbf{M} \\ &= (\boldsymbol{\Psi}^{-1} (\mathbf{W} \mathbf{M} \mathbf{W}^T + \boldsymbol{\Psi}) - \boldsymbol{\Psi}^{-1} \mathbf{W} \mathbf{M} \mathbf{W}^T) (\mathbf{W} \mathbf{M} \mathbf{W}^T + \boldsymbol{\Psi})^{-1} \mathbf{W} \mathbf{M} \\ &= (\mathbf{W} \mathbf{M} \mathbf{W}^T + \boldsymbol{\Psi})^{-1} \mathbf{W} \mathbf{M}. \end{aligned} \quad (198)$$

It follows that

$$\mathbf{M} \mathbf{W}^T (\mathbf{W} \mathbf{W}^T + \boldsymbol{\Psi})^{-1} \mathbf{a} = (\mathbf{M}^{-1} + \mathbf{W}^T \boldsymbol{\Psi}^{-1} \mathbf{W})^{-1} \mathbf{W}^T \boldsymbol{\Psi}^{-1} \mathbf{a}. \quad (199)$$

The posterior  $p(\mathbf{h} \mid \mathbf{v})$  is derived from Gaussian conditioning because both the likelihood  $p(\mathbf{v})$  and the prior  $p(\mathbf{h})$  are Gaussian distributed. The conditional distribution  $p(\mathbf{a} \mid \mathbf{b})$  of two random variables  $\mathbf{a}$  and  $\mathbf{b}$  that both follow a Gaussian distribution is a Gaussian:

$$\mathbf{a} \sim \mathcal{N}(\boldsymbol{\mu}_a, \boldsymbol{\Sigma}_{aa}), \quad (200)$$

$$\mathbf{b} \sim \mathcal{N}(\boldsymbol{\mu}_b, \boldsymbol{\Sigma}_{bb}), \quad (201)$$

$$\boldsymbol{\Sigma}_{ba} = \text{Cov}(\mathbf{b}, \mathbf{a}), \quad (202)$$

$$\boldsymbol{\Sigma}_{ab} = \text{Cov}(\mathbf{a}, \mathbf{b}), \quad (203)$$

$$\mathbf{a} \mid \mathbf{b} \sim \mathcal{N}(\boldsymbol{\mu}_a + \boldsymbol{\Sigma}_{ab} \boldsymbol{\Sigma}_{bb}^{-1} (\mathbf{b} - \boldsymbol{\mu}_b), \boldsymbol{\Sigma}_{aa} - \boldsymbol{\Sigma}_{ab} \boldsymbol{\Sigma}_{bb}^{-1} \boldsymbol{\Sigma}_{ba}). \quad (204)$$

Therefore we need the second moments between  $\mathbf{v}$  and  $\mathbf{h}$ :

$$\mathbb{E}(\mathbf{v}\mathbf{h}^T) = \mathbb{E}(\mathbf{W}\mathbf{h}\mathbf{h}^T) + \mathbb{E}(\boldsymbol{\epsilon}\mathbf{h}^T) = \mathbf{W} (\mathbf{M} + \boldsymbol{\xi} \boldsymbol{\xi}^T). \quad (205)$$

The covariances between  $\mathbf{v}$  and  $\mathbf{h}$  are

$$\text{Cov}(\mathbf{v}, \mathbf{h}) = \mathbb{E}(\mathbf{v}\mathbf{h}^T) - \mathbb{E}(\mathbf{v})\mathbb{E}(\mathbf{h}^T) \quad (206)$$

$$= \mathbf{W} \mathbf{M} + \mathbf{W} \boldsymbol{\xi} \boldsymbol{\xi}^T - \mathbf{W} \boldsymbol{\xi} \boldsymbol{\xi}^T = \mathbf{W} \mathbf{M},$$

$$\text{Cov}(\mathbf{h}, \mathbf{v}) = \mathbb{E}(\mathbf{h}\mathbf{v}^T) - \mathbb{E}(\mathbf{h})\mathbb{E}(\mathbf{v}^T) = \mathbf{M} \mathbf{W}^T. \quad (207)$$

Thus, the mean of  $p(\mathbf{h} \mid \mathbf{v})$  is

$$\begin{aligned} \boldsymbol{\mu}_{h|\mathbf{v}} &= \boldsymbol{\xi} + \mathbf{M} \mathbf{W}^T (\mathbf{W} \mathbf{M} \mathbf{W}^T + \boldsymbol{\Psi})^{-1} (\mathbf{v} - \mathbf{W} \boldsymbol{\xi}) \\ &= \boldsymbol{\xi} + (\mathbf{M}^{-1} + \mathbf{W}^T \boldsymbol{\Psi}^{-1} \mathbf{W})^{-1} \mathbf{W}^T \boldsymbol{\Psi}^{-1} (\mathbf{v} - \mathbf{W} \boldsymbol{\xi}) \\ &= (\mathbf{M}^{-1} + \mathbf{W}^T \boldsymbol{\Psi}^{-1} \mathbf{W})^{-1} (\mathbf{M}^{-1} + \mathbf{W}^T \boldsymbol{\Psi}^{-1} \mathbf{W}) \boldsymbol{\xi} \\ &\quad + (\mathbf{M}^{-1} + \mathbf{W}^T \boldsymbol{\Psi}^{-1} \mathbf{W})^{-1} \mathbf{W}^T \boldsymbol{\Psi}^{-1} (\mathbf{v} - \mathbf{W} \boldsymbol{\xi}) \\ &= (\mathbf{M}^{-1} + \mathbf{W}^T \boldsymbol{\Psi}^{-1} \mathbf{W})^{-1} \\ &\quad (\mathbf{M}^{-1} \boldsymbol{\xi} + \mathbf{W}^T \boldsymbol{\Psi}^{-1} \mathbf{W} \boldsymbol{\xi} + \mathbf{W}^T \boldsymbol{\Psi}^{-1} \mathbf{v} - \mathbf{W}^T \boldsymbol{\Psi}^{-1} \mathbf{W} \boldsymbol{\xi}) \\ &= (\mathbf{M}^{-1} + \mathbf{W}^T \boldsymbol{\Psi}^{-1} \mathbf{W})^{-1} (\mathbf{W}^T \boldsymbol{\Psi}^{-1} \mathbf{v} + \mathbf{M}^{-1} \boldsymbol{\xi}). \end{aligned} \quad (208)$$

The covariance matrix of  $p(\mathbf{h} \mid \mathbf{v})$  is

$$\begin{aligned}\Sigma_{\mathbf{h}|\mathbf{v}} &= \mathbf{M} - \mathbf{M} \mathbf{W}^T (\mathbf{W} \mathbf{M} \mathbf{W}^T + \Psi)^{-1} \mathbf{W} \mathbf{M} \\ &= (\mathbf{M}^{-1} + \mathbf{W}^T \Psi^{-1} \mathbf{W})^{-1} .\end{aligned}\tag{209}$$

In particular, the variable  $\xi$  may be used to enforce more sparseness by setting its components to negative values. Since the covariance matrix  $\Sigma_{\mathbf{h}|\mathbf{v}}$  is positive semi-definite, we ensure that

$$\xi^T (\mathbf{M}^{-1} + \mathbf{W}^T \Psi^{-1} \mathbf{W})^{-1} \xi \geq 0 .\tag{210}$$

If  $\xi = -\rho \mathbf{1}$  ( $\mathbf{1}$  is the vector with all components being one), then the largest absolute components of  $\Sigma_{\mathbf{h}|\mathbf{v}} \xi$  must be negative. Thus,  $\xi = -\rho \mathbf{1}$  leads to sparser solutions.

## 11 Hyperparameters Selected for Method Assessment

The performance of rectified factor networks (RFNs) as unsupervised methods for data representation was compared with:

- (1) **RFN**: rectified factor networks,
- (2) **RFNn**: RFNs without normalization,
- (3) **DAE**: denoising autoencoders with rectified linear units,
- (4) **RBM**: restricted Boltzmann machines with Gaussian visible units and hidden binary units,
- (5) **FAsp**: factor analysis with Jeffrey's prior ( $p(z) \propto 1/z$ ) on the hidden units which is sparser than a Laplace prior,
- (6) **FAlap**: factor analysis with Laplace prior on the hidden units,
- (7) **ICA**: independent component analysis by FastICA Hyvärinen and Oja [1999],
- (8) **SFA**: sparse factor analysis with a Laplace prior on the parameters,
- (9) **FA**: standard factor analysis,
- (10) **PCA**: principal component analysis.

The number of components are fixed to 50, 100, or 150 for each method. The used hyperparameters are listed in Tab. 1.

Table 1: Hyperparameters of all methods that were used to assess the performance of rectified factor networks (RFNs) as unsupervised methods for data representation.

Method	Used hyperparameters
RFN	{learning rate=0.1, iterations=1000}
RFNn	{learning rate=0.1, iterations=1000}
DAE	{corruption level=0.2, learning rate=1e-04, iterations=1000}
RBM	{learning rate=0.01, iterations=1000}
FAsp	{iterations=500}
FAlap	{iterations=500}
SFA	{Laplace weight decay factor=5e-05, iterations=500}

## 12 Data Set I

The number of components are fixed to 50, 100 or 150.

We generated nine different benchmark data sets (D1 to D9), where each data set consists of 100 instances for averaging the results. Each instance consists of 100 samples and 100 features resulting in a  $100 \times 100$  data matrix. Into these data matrices, structures are implanted as biclusters. A bicluster is a pattern consisting of a particular number of features which is found in a particular number of samples. The size of the bicluster is given by the number of features that form the pattern and by the number of samples in which the pattern is found. The data sets had different noise levels and different bicluster sizes. We considered large and small bicluster sizes, where large biclusters have 20–30 samples and 20–30 features, while small biclusters have 3–8 samples and 3–8 features. The signal strength (scaling factor) of a pattern in a sample was randomly chosen according to the Gaussian  $\mathcal{N}(1, 1)$ . Finally, to each data matrix background noise was added, where the noise is distributed according to a zero-mean Gaussian with standard deviation 1, 5, or 10. The data sets are described in Tab. 2. The remaining components of the spanning outer product vectors were drawn by  $\mathcal{N}(0, 0.01)$ .

Table 2: Overview over the datasets. Shown is the background noise (“noise”), the number of large biclusters ( $n_1$ ), and the number of small biclusters ( $n_2$ ).

	D1	D2	D3	D4	D5	D6	D7	D8	D9
noise	1	5	10	1	5	10	1	5	10
$n_1$	10	10	10	15	15	15	5	5	5
$n_2$	10	10	10	5	5	5	15	15	15

Table 3: Comparison for 50 factors / hidden units extracted by RFN, RFN without normalization (RFNn), denoising autoencoder (DAE), restricted Boltzmann machines (RBM), factor analysis with a very sparse prior (FAsp), factor analysis with a Laplace prior (FAlap), independent component analysis (ICA), sparse factor analysis (SFA), factor analysis (FA), and principal component analysis (PCA) on nine data sets. Criteria are: sparseness of the coding units (SP), reconstruction error (ER), and the difference between the empirical and the model covariance matrix (CO). The lower right column block gives the average SP (%), ER and CO. Results reported here, are the mean together with the standard deviation of 100 instances. The maximal value in the table and the maximal standard deviation was set to 999 and to 99, respectively.

	D1			D2			D3			D4			D5		
	SP	ER	CO	SP	ER	CO	SP	ER	CO	SP	ER	CO	SP	ER	CO
RFN	74±0	58±1	5±0	75±0	233±3	66±1	75±0	456±5	253±6	74±0	63±1	6±1	75±0	236±3	68±2
RFNn	73±0	85±3	13±2	75±0	272±3	85±2	75±0	531±6	321±7	72±0	95±4	17±2	74±0	276±4	89±3
DAE	65±0	65±2	—	66±0	233±2	—	66±0	456±4	—	65±1	71±2	—	66±0	237±2	—
RBM	25±2	86±3	—	11±1	287±3	—	10±1	558±5	—	25±2	94±3	—	11±1	292±3	—
FAsp	39±1	232±31	654±99	40±1	999±41	999±99	41±1	999±99	999±99	38±1	318±33	999±99	40±1	999±48	999±99
FAlap	4±0	53±2	144±36	4±0	224±5	185±5	5±0	439±9	692±16	4±0	55±2	180±39	4±0	226±5	192±6
ICA	2±0	34±0	—	2±0	164±2	—	2±0	324±4	—	2±0	35±0	—	2±0	166±2	—
SFA	1±0	42±1	11±2	1±0	206±4	56±2	1±0	406±9	215±7	1±0	42±1	13±2	1±0	208±4	58±2
FA	1±0	42±1	6±1	1±0	206±4	54±2	1±0	407±8	210±6	1±0	42±1	8±1	1±0	208±4	56±2
PCA	1±0	34±0	—	0±0	164±2	—	0±0	324±4	—	1±0	35±0	—	0±0	166±2	—
	D6			D7			D8			D9			average		
	SP	ER	CO	SP	ER	CO	SP	ER	CO	SP	ER	CO	SP	ER	CO
RFN	75±0	458±5	256±6	75±0	53±1	4±1	75±0	230±3	64±1	75±0	454±5	251±5	75±0	249±3	108±3
RFNn	75±0	532±6	323±7	73±0	73±3	10±2	75±0	268±3	82±2	75±0	528±6	317±7	74±0	295±4	140±4
DAE	66±0	458±4	—	65±0	58±1	—	66±0	230±2	—	66±0	453±5	—	66±0	251±3	—
RBM	10±1	561±5	—	23±2	76±2	—	11±1	282±3	—	10±1	555±5	—	15±1	310±4	—
FAsp	40±2	999±99	999±99	39±1	152±26	345±99	40±1	999±31	999±99	41±1	999±99	999±99	40±1	999±63	999±99
FAlap	5±0	443±9	701±15	4±0	50±2	110±37	4±0	221±5	177±4	5±0	439±10	686±15	4±0	239±6	341±19
ICA	2±0	325±4	—	2±0	34±0	—	2±0	163±2	—	2±0	322±4	—	2±0	174±2	—
SFA	1±0	408±9	217±7	1±0	42±1	8±2	1±0	204±4	54±2	1±0	405±9	213±7	1±0	218±5	94±3
FA	1±0	409±9	212±7	1±0	42±1	4±1	1±0	205±4	53±2	1±0	405±8	208±6	1±0	218±4	90±3
PCA	0±0	325±4	—	1±0	34±0	—	0±0	163±2	—	0±0	322±4	—	0±0	174±2	—



Table 4: Comparison for 100 factors / hidden units extracted by RFN, RFN without normalization (RFNn), denoising autoencoder (DAE), restricted Boltzmann machines (RBM), factor analysis with a very sparse prior (FAsp), factor analysis with a Laplace prior (FAlap), independent component analysis (ICA), sparse factor analysis (SFA), factor analysis (FA), and principal component analysis (PCA) on nine data sets. Criteria are: sparseness of the coding units (SP), reconstruction error (ER), and the difference between the empirical and the model covariance matrix (CO). The lower right column block gives the average SP (%), ER and CO. Results reported here, are the mean together with the standard deviation of 100 instances. The maximal value in the table and the maximal standard deviation was set to 999 and to 99, respectively.

	D1			D2			D3			D4			D5		
	SP	ER	CO	SP	ER	CO	SP	ER	CO	SP	ER	CO	SP	ER	CO
RFN	79±1	23±3	2±0	82±1	63±9	16±3	82±1	120±17	61±15	78±1	27±3	2±1	82±1	62±7	16±3
RFNn	77±0	61±4	6±1	80±0	169±4	36±2	80±0	326±8	135±6	76±1	73±4	9±2	79±0	171±5	37±2
DAE	67±0	48±2	—	70±0	134±1	—	70±0	260±2	—	67±0	54±2	—	70±0	137±1	—
RBM	14±1	81±3	—	4±0	266±3	—	4±0	514±6	—	15±1	88±2	—	4±0	270±3	—
FAsp	72±0	233±32	499±99	62±0	999±43	999±99	56±0	999±99	999±99	71±0	320±34	878±99	62±0	999±49	999±99
FAlap	6±0	27±3	202±17	6±0	38±3	756±33	6±0	74±5	999±83	6±0	31±3	274±23	6±0	39±3	778±34
ICA	3±2	0±0	—	3±1	0±0	—	3±1	0±0	—	3±2	0±0	—	3±1	0±0	—
SFA	1±0	6±0	30±5	1±0	14±0	68±3	1±0	28±1	243±8	1±0	8±0	38±5	1±0	15±0	72±3
FA	1±0	6±0	18±3	1±0	14±0	50±2	1±0	28±1	182±7	1±0	8±0	24±4	1±0	15±0	52±2
PCA	4±0	0±0	—	2±0	0±0	—	1±0	0±0	—	4±0	0±0	—	2±0	0±0	—
	D6			D7			D8			D9			average		
	SP	ER	CO	SP	ER	CO	SP	ER	CO	SP	ER	CO	SP	ER	CO
RFN	82±1	120±16	60±13	80±1	18±2	1±0	82±1	61±7	15±3	82±1	122±13	60±11	81±1	68±9	26±6
RFNn	80±0	329±7	137±6	78±0	49±3	4±1	80±0	165±4	34±1	80±0	325±7	134±6	79±0	185±5	59±3
DAE	70±0	261±2	—	68±0	39±2	—	70±0	132±1	—	70±0	259±2	—	69±0	147±2	—
RBM	4±0	517±6	—	12±1	71±2	—	4±0	261±3	—	4±0	512±5	—	7±1	287±4	—
FAsp	56±1	999±99	999±99	73±0	149±28	237±62	62±0	999±34	999±99	56±0	999±99	999±99	63±0	999±65	999±99
FAlap	6±0	74±6	999±91	6±0	22±3	134±14	6±0	37±2	733±28	6±0	73±6	999±84	6±0	46±4	985±45
ICA	3±1	0±0	—	3±2	0±0	—	3±1	0±0	—	3±1	0±0	—	3±1	0±0	—
SFA	1±0	28±1	247±8	1±0	5±0	21±5	1±0	14±0	64±2	1±0	27±1	240±7	1±0	16±1	114±5
FA	1±0	28±1	184±8	1±0	5±0	11±3	1±0	14±0	47±2	1±0	27±1	179±7	1±0	16±1	83±4
PCA	1±0	0±0	—	4±0	0±0	—	2±0	0±0	—	1±0	0±0	—	2±0	0±0	—

Table 5: Comparison for 150 factors / hidden units extracted by RFN, RFN without normalization (RFNn), denoising autoencoder (DAE), restricted Boltzmann machines (RBM), factor analysis with a very sparse prior (FAsp), factor analysis with a Laplace prior (FAlap), independent component analysis (ICA), sparse factor analysis (SFA), factor analysis (FA), and principal component analysis (PCA) on nine data sets. Criteria are: sparseness of the coding units (SP), reconstruction error (ER), and the difference between the empirical and the model covariance matrix (CO). The lower right column block gives the average SP (%), ER and CO. The lower right column block gives the average SP, ER and CO. Results reported here, are the mean together with the standard deviation of 100 instances. The maximal value in the table and the maximal standard deviation was set to 999 and to 99, respectively.

	D1				D2				D3				D4				D5			
	SP	ER	CO		SP	ER	CO		SP	ER	CO		SP	ER	CO		SP	ER	CO	
RFN	83±1	7±2	0±1		86±0	15±1	3±1		86±2	33±20	18±23		83±1	9±2	1±0		86±1	15±3	4±1	
RFNn	79±0	48±3	4±1		81±0	129±3	21±1		81±0	250±7	80±4		78±0	60±4	6±1		81±0	131±3	22±1	
DAE	68±0	44±2	—		72±0	118±1	—		72±0	229±2	—		68±0	50±2	—		72±0	120±2	—	
RBM	10±1	81±3	—		3±0	265±3	—		3±0	514±6	—		10±1	88±2	—		3±0	270±4	—	
FAsp	83±1	233±32	340±71		79±0	999±43	999±99		77±0	999±99	999±99		81±1	320±34	574±99		79±1	999±49	999±99	
FAlap	4±0	27±3	295±25		4±0	38±3	791±41		3±0	74±5	999±91		4±0	31±3	394±31		4±0	39±3	817±39	
ICA	3±2	0±0	—		3±1	0±0	—		3±1	0±0	—		3±2	0±0	—		3±1	0±0	—	
SFA	1±0	6±0	49±7		1±0	14±0	173±4		1±0	28±1	632±10		1±0	8±0	61±7		1±0	15±0	181±5	
FA	1±0	6±0	40±5		1±0	14±0	160±4		1±0	28±1	590±10		1±0	8±0	51±6		1±0	15±0	168±4	
PCA	4±0	0±0	—		2±0	0±0	—		1±0	0±0	—		4±0	0±0	—		2±0	0±0	—	
	D6				D7				D8				D9				average			
	SP	ER	CO		SP	ER	CO		SP	ER	CO		SP	ER	CO		SP	ER	CO	
RFN	86±1	30±13	15±16		84±2	5±3	0±1		86±0	14±1	3±1		86±1	30±8	15±9		85±1	17±6	7±6	
RFNn	81±0	251±6	81±3		80±0	37±3	2±0		81±0	126±3	20±1		81±0	248±6	79±3		80±0	142±4	35±2	
DAE	72±0	230±2	—		70±0	36±2	—		72±0	116±1	—		72±0	227±2	—		71±0	130±2	—	
RBM	3±0	516±6	—		8±1	71±2	—		3±0	260±4	—		3±0	511±5	—		5±0	286±4	—	
FAsp	77±0	999±99	999±99		84±0	149±28	168±55		80±0	999±34	999±99		77±1	999±99	999±99		80±0	999±65	999±99	
FAlap	3±0	74±6	999±97		4±0	22±3	198±17		4±0	37±2	768±40		3±0	73±6	999±93		4±0	46±4	976±53	
ICA	3±1	0±0	—		3±2	0±0	—		3±1	0±0	—		3±1	0±0	—		3±1	0±0	—	
SFA	1±0	28±1	640±11		1±0	5±0	34±6		1±0	14±0	164±3		1±0	27±1	625±9		1±0	16±1	285±7	
FA	1±0	28±1	596±10		1±0	5±0	27±5		1±0	14±0	153±3		1±0	27±1	583±9		1±0	16±1	263±6	
PCA	1±0	0±0	—		4±0	0±0	—		2±0	0±0	—		1±0	0±0	—		2±0	0±0	—	

## 13 Data Set II

This data sets was generate as described in Section 12, but instead of drawing the remaining components of the spanning outer product vectors from  $\mathcal{N}(0, 0.01)$ , they were now drawn from  $\mathcal{N}(0, 0.5)$ .

Table 6: Comparison for 50 factors / hidden units extracted by RFN, RFN without normalization (RFNn), denoising autoencoder (DAE), restricted Boltzmann machines (RBM), factor analysis with a very sparse prior (FAsp), factor analysis with a Laplace prior (FAlap), independent component analysis (ICA), sparse factor analysis (SFA), factor analysis (FA), and principal component analysis (PCA) on nine data sets. Criteria are: sparseness of the coding units (SP), reconstruction error (ER), and the difference between the empirical and the model covariance matrix (CO). The lower right column block gives the average SP (%), ER and CO. Results reported here, are the mean together with the standard deviation of 100 instances. The maximal value in the table and the maximal standard deviation was set to 999 and to 99, respectively.

	D1			D2			D3			D4			D5		
	SP	ER	CO	SP	ER	CO	SP	ER	CO	SP	ER	CO	SP	ER	CO
RFN	72±1	74±2	11±1	75±0	240±3	72±2	75±0	462±5	260±6	72±1	79±2	12±1	75±0	244±3	75±2
RFNn	68±1	122±5	32±4	74±0	285±4	97±3	74±0	537±7	331±8	65±1	144±6	48±6	74±0	290±4	102±4
DAE	61±0	82±2	—	66±0	243±2	—	66±0	461±4	—	60±0	88±2	—	66±0	247±3	—
RBM	22±1	106±3	—	11±1	301±3	—	10±1	566±6	—	22±1	113±3	—	11±1	308±4	—
FAsp	37±1	469±38	999±99	40±1	999±50	999±99	40±2	999±99	999±99	37±1	610±44	999±99	40±1	999±58	999±99
FAlap	4±0	50±1	392±66	4±0	228±5	135±13	5±0	443±9	406±18	4±0	51±1	477±63	4±0	230±6	147±18
ICA	2±0	35±0	—	2±0	168±2	—	2±0	327±4	—	2±0	35±0	—	2±0	170±2	—
SFA	1±0	42±1	26±3	1±0	210±5	61±2	1±0	409±8	220±6	1±0	41±1	32±4	1±0	211±5	63±2
FA	1±0	42±1	13±2	1±0	210±4	58±2	1±0	409±8	214±6	1±0	41±1	17±2	1±0	212±5	60±2
PCA	0±0	35±0	—	0±0	168±2	—	0±0	327±4	—	0±0	35±0	—	0±0	170±2	—
	D6			D7			D8			D9			average		
	SP	ER	CO	SP	ER	CO	SP	ER	CO	SP	ER	CO	SP	ER	CO
RFN	75±0	464±5	264±6	73±0	68±2	9±1	75±0	237±3	69±1	75±0	459±5	257±6	74±0	259±3	114±3
RFNn	74±0	541±6	336±8	71±1	106±4	23±3	74±0	279±3	91±2	75±0	533±6	325±8	72±1	315±5	154±5
DAE	66±0	465±4	—	62±0	75±2	—	66±0	238±2	—	66±0	458±4	—	64±0	262±3	—
RBM	10±1	570±6	—	20±1	97±3	—	11±1	294±3	—	10±1	562±5	—	14±1	324±4	—
FAsp	41±1	999±99	999±99	38±1	335±32	999±99	41±1	999±40	999±99	41±1	999±99	999±99	39±1	999±69	999±99
FAlap	5±0	447±9	413±19	4±0	49±1	292±57	4±0	227±5	123±11	5±0	443±9	401±17	4±0	241±5	310±31
ICA	2±0	329±4	—	2±0	35±0	—	2±0	167±2	—	2±0	325±4	—	2±0	177±2	—
SFA	1±0	412±8	223±7	1±0	42±1	19±3	1±0	209±4	59±2	1±0	408±9	218±7	1±0	221±5	102±4
FA	1±0	412±8	217±7	1±0	42±1	10±1	1±0	209±4	57±2	1±0	409±9	213±7	1±0	221±5	95±3
PCA	0±0	329±4	—	0±0	35±0	—	0±0	167±2	—	0±0	325±4	—	0±0	177±2	—

Table 7: Comparison for 100 factors / hidden units extracted by RFN, RFN without normalization (RFNn), denoising autoencoder (DAE), restricted Boltzmann machines (RBM), factor analysis with a very sparse prior (FAsp), factor analysis with a Laplace prior (FAlap), independent component analysis (ICA), sparse factor analysis (SFA), factor analysis (FA), and principal component analysis (PCA) on nine data sets. Criteria are: sparseness of the coding units (SP), reconstruction error (ER), and the difference between the empirical and the model covariance matrix (CO). The lower right column block gives the average SP (%), ER and CO. Results reported here, are the mean together with the standard deviation of 100 instances. The maximal value in the table and the maximal standard deviation was set to 999 and to 99, respectively.

	D1			D2			D3			D4			D5		
	SP	ER	CO	SP	ER	CO	SP	ER	CO	SP	ER	CO	SP	ER	CO
RFN	76±1	34±3	4±1	82±1	67±8	18±3	82±1	124±16	63±12	75±1	38±3	5±1	82±1	69±10	19±5
RFNn	71±1	110±7	25±4	79±0	180±5	42±2	80±0	331±8	139±7	65±2	143±9	47±8	79±0	185±5	45±3
DAE	63±0	66±2	—	70±0	142±2	—	70±0	264±3	—	62±0	73±2	—	70±0	146±2	—
RBM	12±1	100±3	—	5±0	282±4	—	4±0	522±6	—	12±1	106±3	—	5±1	288±4	—
FAsp	71±0	474±38	999±99	62±0	999±53	999±99	56±1	999±99	999±99	70±0	616±44	999±99	62±0	999±60	999±99
FAlap	6±0	21±2	425±28	6±0	40±2	827±35	6±0	75±6	999±99	6±0	23±2	523±32	6±0	42±3	865±43
ICA	3±2	0±0	—	3±1	0±0	—	3±1	0±0	—	3±2	0±0	—	3±1	0±0	—
SFA	1±0	10±0	71±7	1±0	15±0	84±4	1±0	28±1	254±8	1±0	12±0	87±8	1±0	16±0	92±5
FA	1±0	10±0	48±5	1±0	15±0	59±3	1±0	28±1	189±7	1±0	12±1	61±6	1±0	16±0	64±3
PCA	4±0	0±0	—	2±0	0±0	—	1±0	0±0	—	3±0	0±0	—	2±0	0±0	—
	D6			D7			D8			D9			average		
	SP	ER	CO	SP	ER	CO	SP	ER	CO	SP	ER	CO	SP	ER	CO
RFN	82±1	127±17	65±14	77±1	30±3	3±1	82±1	64±8	17±4	82±1	123±15	62±13	80±1	75±9	28±6
RFNn	80±0	334±8	141±7	74±1	86±4	14±2	79±0	174±4	39±2	80±0	329±7	137±6	76±1	208±6	70±5
DAE	70±0	266±2	—	64±0	57±2	—	70±0	138±1	—	70±0	262±2	—	68±0	157±2	—
RBM	4±0	527±6	—	11±1	92±2	—	4±0	274±4	—	4±0	518±6	—	7±1	301±4	—
FAsp	56±0	999±99	999±99	71±0	338±33	999±99	62±1	999±42	999±99	56±1	999±99	999±99	63±0	999±74	999±99
FAlap	6±0	75±6	999±89	6±0	18±2	337±24	6±0	40±3	793±37	6±0	74±6	999±89	6±0	45±3	999±53
ICA	3±1	0±0	—	3±1	0±0	—	3±1	0±0	—	3±1	0±0	—	3±1	0±0	—
SFA	1±0	28±1	260±9	1±0	8±0	52±7	1±0	15±0	76±3	1±0	28±1	248±7	1±0	18±1	136±6
FA	1±0	28±1	193±8	1±0	8±0	33±5	1±0	15±0	54±2	1±0	28±1	185±6	1±0	18±1	99±5
PCA	1±0	0±0	—	4±0	0±0	—	2±0	0±0	—	1±0	0±0	—	2±0	0±0	—

Table 8: Comparison for 150 factors / hidden units extracted by RFN, RFN without normalization (RFNn), denoising autoencoder (DAE), restricted Boltzmann machines (RBM), factor analysis with a very sparse prior (FAsp), factor analysis with a Laplace prior (FAlap), independent component analysis (ICA), sparse factor analysis (SFA), factor analysis (FA), and principal component analysis (PCA) on nine data sets. Criteria are: sparseness of the factors (SP) reported in %, reconstruction error (ER), and the difference between the empirical and the model covariance matrix (CO). The lower right column block gives the average SP (%), ER and CO. Results reported here, are the mean together with the standard deviation of 100 instances. The maximal value in the table and the maximal standard deviation was set to 999 and to 99, respectively.

	D1			D2			D3			D4			D5		
	SP	ER	CO	SP	ER	CO	SP	ER	CO	SP	ER	CO	SP	ER	CO
RFN	81±1	12±2	1±1	86±0	16±1	4±1	86±0	29±4	15±5	80±1	15±5	2±2	86±1	17±5	5±3
RFNn	72±1	100±8	19±4	80±0	137±4	24±1	81±0	254±6	83±4	66±0	113±3	52±5	80±0	141±4	26±2
DAE	64±0	62±2	—	71±0	125±2	—	72±0	232±2	—	63±0	69±2	—	71±0	129±2	—
RBM	8±0	101±3	—	4±0	282±4	—	3±0	521±6	—	8±0	106±3	—	4±0	289±4	—
FAsp	81±1	474±38	999±99	79±0	999±53	999±99	77±1	999±99	999±99	80±1	616±44	999±99	79±1	999±60	999±99
FAlap	4±0	21±2	607±34	4±0	40±2	879±40	3±0	75±6	999±96	4±0	23±2	749±42	4±0	42±3	926±45
ICA	3±2	0±0	—	3±1	0±0	—	3±1	0±0	—	3±2	0±0	—	3±1	0±0	—
SFA	1±0	10±0	103±9	1±0	15±0	204±7	1±0	28±1	656±12	1±0	12±0	126±10	1±0	16±0	220±8
FA	1±0	10±0	87±8	1±0	15±0	187±5	1±0	28±1	611±11	1±0	12±1	108±9	1±0	16±0	200±6
PCA	4±0	0±0	—	2±0	0±0	—	1±0	0±0	—	3±0	0±0	—	2±0	0±0	—
	D6			D7			D8			D9			average		
	SP	ER	CO	SP	ER	CO	SP	ER	CO	SP	ER	CO	SP	ER	CO
RFN	86±1	29±7	15±6	82±1	10±3	1±1	86±1	17±10	5±9	86±1	31±19	16±13	84±1	20±6	7±4
RFNn	81±0	255±6	84±3	76±1	74±5	9±2	81±0	133±3	23±1	81±0	250±7	81±4	77±0	162±5	45±3
DAE	72±0	234±2	—	65±0	53±2	—	72±0	122±1	—	72±0	230±2	—	69±0	140±2	—
RBM	3±0	525±6	—	8±0	93±3	—	3±0	273±4	—	3±0	517±6	—	5±0	301±4	—
FAsp	77±1	999±99	999±99	81±1	338±33	673±99	79±0	999±42	999±99	77±1	999±99	999±99	79±1	999±74	999±99
FAlap	3±0	75±6	999±94	4±0	18±2	479±31	4±0	40±3	831±43	3±0	74±6	999±95	4±0	45±3	999±88
ICA	3±1	0±0	—	3±1	0±0	—	3±1	0±0	—	3±1	0±0	—	3±1	0±0	—
SFA	1±0	28±1	668±12	1±0	8±0	78±8	1±0	15±0	188±5	1±0	28±1	644±9	1±0	18±1	321±9
FA	1±0	28±1	622±11	1±0	8±0	64±7	1±0	15±0	173±4	1±0	28±1	599±9	1±0	18±1	294±8
PCA	1±0	0±0	—	4±0	0±0	—	2±0	0±0	—	1±0	0±0	—	2±0	0±0	—

## 14 RFN Pretraining for Convolution Nets

We assess the performance of RFN *first layer* pretraining on *CIFAR-10* and *CIFAR-100* for three deep convolutional network architectures: (i) the AlexNet Krizhevsky et al. [2012], (ii) Deeply Supervised Networks (DSN) Lee et al. [2014], and (iii) our 5-Convolution-Network-In-Network (5C-NIN).

Both CIFAR datasets contain 60k 32x32 RGB-color images, which were divided into 50k train and 10k test sets, split between 10 (CIFAR10) and 100 (CIFAR100) categories. Both datasets are preprocessed as described in Goodfellow et al. [2013] by global contrast normalization and ZCA whitening. Additionally, the datasets were augmented by padding the images with four zero pixels at all borders. For data augmentation, at the beginning of every epoch, images in the training set were distorted by random translation and random flipping in horizontal and vertical directions. For the AlexNet, we neither preprocessed nor augmented the datasets.

Inspired by Lin et al. [2013]’s Network In Network, we constructed a 5-Convolution-Network-In-Network (5C-NIN) architecture with five convolutional layers, each followed by a 2x2 max-pooling layer (stride 1) and a multilayer perceptron (MLP) convolutional layer. ReLUs were used for the convolutional layers and dropout for regularization. We followed Krizhevsky [2009] for weight initialization, learning rates, and learning policies. The networks were trained using mini-batches of size 100 and 128 for 5C-NIN and AlexNet, respectively.

For RFN pretraining, we randomly extracted 5x5 patches from the training data to construct 192 filters for DSN and 5C-NIN while 32 for AlexNet. These filters constitute the first convolutional layer of each network which is then trained using default setting. For assessing the improvement by RFNs, we repeated training with randomly initialized weights in the first layer. The results are presented in Tab. 9. For comparison, the lower panel of the table reports the performance of the currently top performing networks: Network In Network (NIN, Lin et al. [2013]), Maxout Networks (MN, Goodfellow et al. [2013]) and DeepCNiN Graham [2014]. *In all cases pretraining with RFNs decreases the test error rate.*

Table 9: The upper panel shows results of convolutional deep networks with first layer pretrained by RFN (“RFN”) and with first layer randomly initialized (“org”). The first column gives the network architecture, namely, AlexNet, Deeply Supervised Networks (DSN), and our 5-Convolution-Network-In-Network (5C-NIN). The test error rates are reported (for CIFAR-100 DSN model was missing). Currently best performing networks Network In Network (NIN), Maxout Networks (MN), and DeepCNiN are reported in the lower panel. In all cases pretraining with RFNs decreased the test error rate.

Dataset	CIFAR-10		CIFAR-100		augmented
	org	RFN	org	RFN	
AlexNet	18.21	18.04	46.18	45.80	
DSN	7.97	7.74	34.57	-	✓
5C-NIN	7.81	7.63	29.96	29.75	✓
NIN	8.81	-	35.68	-	✓
MN	9.38	-	38.57	-	✓
DeepCNiN	6.28	-	24.30	-	✓

## 15 Running Times for RFN’s Projected Newton Step

In this section, we report the running times for RFN’s projected Newton step and for solving a quadratic program using NumPy Python and CVXOPT (Python Software for Convex Optimization), respectively. Both benchmarks were profiled with the same hardware using only the CPU. Fig. 1 shows the run times for various problem sizes in [s] both approaches. The projected Newton step complexity per iteration is  $O(nl)$ , see Fig. 2. In contrast, a quadratic program solver typically requires for the  $(nl)$  variables (the means of the hidden units for all samples)  $O(n^4 l^4)$  steps to find the minimum Ben-Tal and Nemirovski [2001].

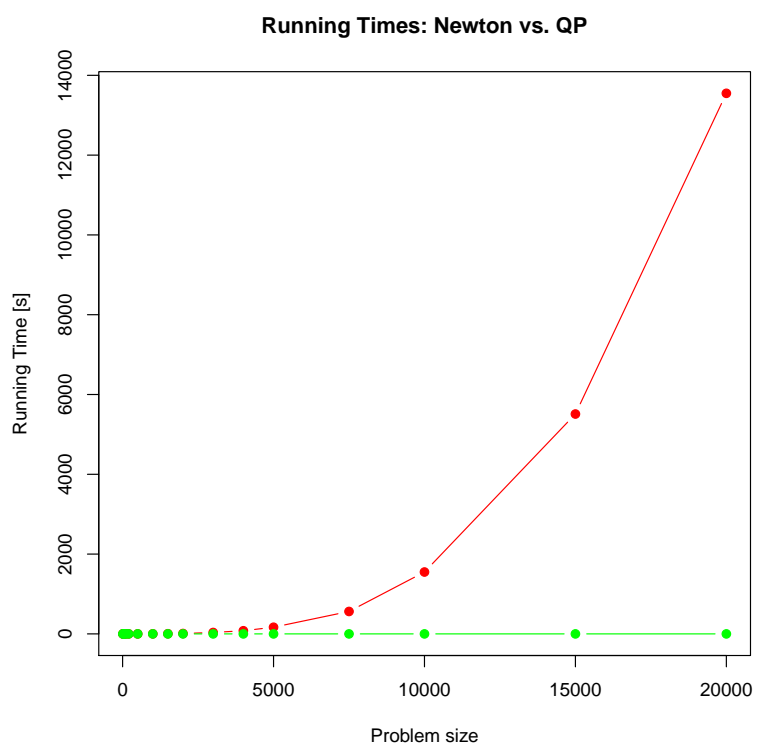


Figure 1: Running times for various problem sizes in [s] of RFN's projected Newton step and of quadratic program solver.



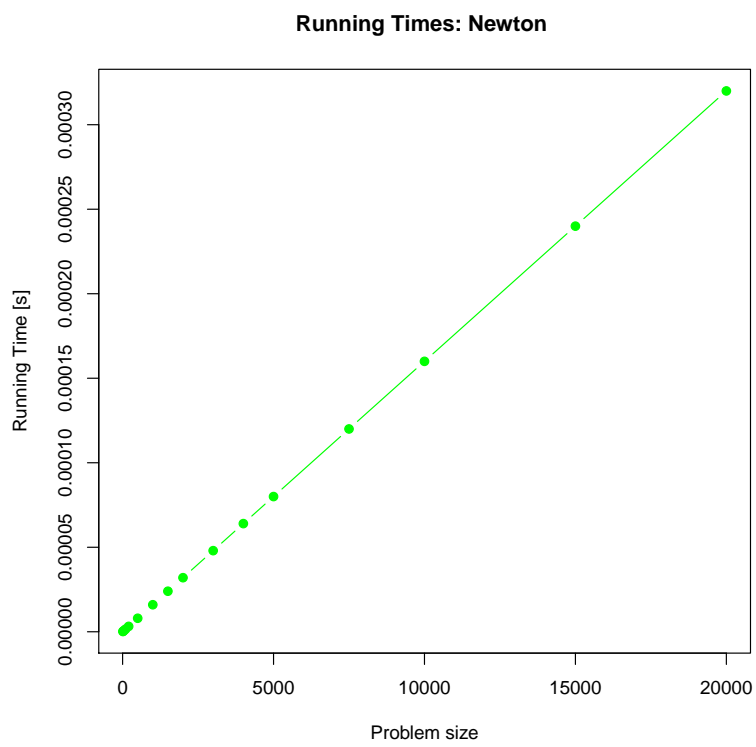


Figure 2: Running times for various problem sizes in [s] of RFN's projected Newton step.

## References

- J. Abadie and J. Carpentier. *Optimization*, chapter Generalization of the Wolfe Reduced Gradient Method to the Case of Nonlinear Constraints. Academic Press, 1969.
- A. Ben-Tal and A. Nemirovski. *Interior Point Polynomial Time Methods for Linear Programming, Conic Quadratic Programming, and Semidefinite Programming*, chapter 6, pages 377–442. Society for Industrial and Applied Mathematics, 2001. doi: 10.1137/1.9780898718829.ch6.
- D. P. Bertsekas. On the Goldstein-Levitin-Polyak gradient projection method. *IEEE Trans. Automat. Control*, 21:174–184, 1976.
- D. P. Bertsekas. Projected Newton methods for optimization problems with simple constraints. *SIAM J. Control Optim.*, 20:221–246, 1982.
- E. G. Birgin, J. M. Martínez, and M. Raydan. Nonmonotone spectral projected gradient methods on convex sets. *Siam Journal on Optimization*, 10(4):1196–1211, 2000. doi: 10.1137/S1052623497330963.
- D. C. Ciresan, U. Meier, and J. Schmidhuber. Multi-column deep neural networks for image classification. In *IEEE Conference on Computer Vision and Pattern Recognition CVPR 2012*, 2012. Long preprint arXiv:1202.2745v1 [cs.CV].
- M. Dredze, K. Crammer, and F. Pereira. Confidence-weighted linear classification. In *Proceedings of the 25th international conference on Machine learning (ICML08)*, volume 25, pages 264–271. ACM New York, 2008. doi: 10.1145/1390156.1390190.
- M. Dredze, K. Crammer, and F. Pereira. Confidence-weighted linear classification for text categorization. *Journal of Machine Learning Research*, 13(1):1891–1926, 2012.
- R. P. Feynman. *Statistical Mechanics*. Benjamin, Reading, MA, 1972.
- B. J. Frey and G. E. Hinton. Variational learning in nonlinear Gaussian belief networks. *Neural Computation*, 11(1):193–214, 1999.
- K. Friston. A free energy principle for biological systems. *Entropy*, 14:2100–2121, 2012.
- K. Ganchev, J. Graca, J. Gillenwater, and B. Taskar. Posterior regularization for structured latent variable models. *Journal of Machine Learning Research*, 11:2001–2049, 2010.
- I. J. Goodfellow, D. Warde-Farley, M. Mirza, A. Courville, and Y. Bengio. Maxout networks. *ArXiv e-prints*, 2013.
- J. V. Graca, K. Ganchev, and B. Taskar. Expectation maximization and posterior constraints. In J.C. Platt, D. Koller, Y. Singer, and S.T. Roweis, editors, *Advances in Neural Information Processing Systems*, volume 20, pages 569–576, 2007.
- J. V. Graca, K. Ganchev, B. Taskar, and F. Pereira. Posterior vs. parameter sparsity in latent variable models. In Y. Bengio, D. Schuurmans, J. D. Lafferty, C. K. I. Williams, and A. Culotta, editors, *Advances in Neural Information Processing Systems*, volume 22, pages 664–672, 2009.
- Benjamin Graham. Fractional max-pooling. *CoRR*, abs/1412.6071, 2014. URL <http://arxiv.org/abs/1412.6071>.
- A. Gunawardana and W. Byrne. Convergence theorems for generalized alternating minimization procedures. *Journal of Machine Learning Research*, 6:2049–2073, 2005.
- M. Harva and A. Kaban. A variational bayesian method for rectified factor analysis. In *Proc. Int. Joint Conf. on Neural Networks (IJCNN’05)*, pages 185–190, 2005.
- M. Harva and A. Kaban. Variational learning for rectified factor analysis. *Signal Processing*, 87(3): 509–527, 2007. doi: 10.1016/j.sigpro.2006.06.006.
- E. J. Haug and J. S. Arora. *Applied optimal design*. J. Wiley & Sons, New York, 1979.
- A. Hyvärinen and E. Oja. A fast fixed-point algorithm for independent component analysis. *Neural Comput.*, 9(7):1483–1492, 1999.
- C. T. Kelley. *Iterative Methods for Optimization*. Society for Industrial and Applied Mathematics (SIAM), Philadelphia, 1999.
- D. Kim, S. Sra, and I. S. Dhillon. A new projected quasi-Newton approach for the nonnegative least squares problem. Technical Report TR-06-54, Department of Computer Sciences, University of Texas at Austin, 2006.

- A. Krizhevsky. Learning multiple layers of features from tiny images. Master's thesis, Department of Computer Science, University of Toronto, 2009.
- A. Krizhevsky, I. Sutskever, and G. E. Hinton. ImageNet classification with deep convolutional neural networks. In F. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 25*, pages 1097–1105. Curran Associates, Inc., 2012.
- S. Kullback and R. A. Leibler. On information and sufficiency. *Annals of Mathematical Statistics*, 22:79–86, 1951.
- C.-Y. Lee, S. Xie, P. Gallagher, Z. Zhang, and Z. Tu. Deeply-Supervised Nets. *ArXiv e-prints*, 2014.
- Min Lin, Qiang Chen, and Shuicheng Yan. Network in network. *CoRR*, abs/1312.4400, 2013. URL <http://arxiv.org/abs/1312.4400>.
- R. Neal and G. E. Hinton. A view of the EM algorithm that justifies incremental, sparse, and other variants. In M. I. Jordan, editor, *Learning in Graphical Models*, pages 355–368. MIT Press, Cambridge, MA, 1998.
- R. Patel and M. Toda. Trace inequalities involving hermitian matrices. *Linear Algebra and its Applications*, 23:13–20, 1979. doi: 10.1016/0024-3795(79)90089-2.
- J. B. Rosen. The gradient projection method for nonlinear programming. part ii. nonlinear constraints. *Journal of the Society for Industrial and Applied Mathematics*, 9(4):514–532, 1961. doi: 10.1137/0109044.
- T. Serafini, G. Zanghirati, and L. Zanni. Gradient projection methods for quadratic programs and applications in training support vector machines. *Optimization Methods and Software*, 20(2-3): 353–378, 2005. doi: 10.1080/10556780512331318182.
- N. Srebro. *Learning with Matrix Factorizations*. PhD thesis, Department of Electrical Engineering and Computer Science, Massachusetts Institute of Technology, 2004.
- C. F. J. Wu. On the convergence properties of the EM algorithm. *Annals of Statistics*, 11(1):95–103, 1983.
- W. I. Zangwill. *Nonlinear Programming: A Unified Approach*. Prentice Hall, Englewood Cliffs, N.J., 1969.