

OPTIMAL LINEAR FEATURE TRANSFORMATIONS FOR SEMI-CONTINUOUS HIDDEN MARKOV MODELS

E. Günter Schukat-Talamazzini, Joachim Hornegger, Heinrich Niemann

Lehrstuhl für Mustererkennung (Informatik 5)
Friedrich-Alexander-Universität Erlangen-Nürnberg
Martensstraße 3
D-91058 Erlangen, F.R. of Germany

E-mail: {schukat,hornegger,niemann}@informatik.uni-erlangen.de

ABSTRACT

Linear discriminant or Karhunen-Loève transforms are established techniques for mapping features into a lower dimensional subspace. This paper introduces a uniform statistical framework, where the computation of the optimal feature reduction is formalized as a Maximum-Likelihood estimation problem. The experimental evaluation of this suggested extension of linear selection methods shows a slight improvement of the recognition accuracy.

1. INTRODUCTION

It is the ultimate goal of any probabilistic approach to speech recognition to capture the entire process of word production within *one single* homogeneous model. The statistical parameters of this model then can be optimized with respect to a large training sample of speech using standard parameter estimation techniques.

A first step in this direction was the introduction of the (discrete density) *hidden Markov model* (HMM, [1]), which constitutes a simple probabilistic description of acoustical word realization down to the level of labelled speech frames. Frame labelling was performed outside the scope of the word models by a phonetic classifier or a vector quantizer.

The undesirable exclusion of the feature level from word modelling stopped with the advent of *semi-continuous* models (SCHMM) [2] which incorporated the formerly external vector quantizer into the speech production model. Thus, in a SCHMM-based speech recognizer the entire processing sequence from the feature vectors to the word level is part of the global HMM, and its free statistical parameters can be jointly optimized by the Baum-Welch algorithm.

Unfortunately, the extraction of features from the speech wave now, as before, stays outside the probabilistic framework, and the development of the signal processing component is still subject to heuristics and intuition. In order to turn, at least a part of, the feature extraction stage into a parameterized and trainable building block of the recognizer, a decompo-

sition of the process is assumed: (1) initially, the signal is transformed into a sequence of high-dimensional short-time feature vectors $\mathbf{x}_t \in \mathbb{R}^D$; (2) the \mathbf{x}_t 's are mapped to vectors \mathbf{y}_t of drastically reduced dimension $d \ll D$. In the context of speech recognition, the first stage can be thought of as computation of the log power mel-spectral coefficients, followed by the enlargement of the resulting short-time parameter vector using first and second order temporal derivatives [3], or differences, or simply the spectral coefficients of an appropriate number of neighboring speech frames. The second stage is formed by a linear mapping involving standard feature extraction activities like, for instance, the cosine transform [4] and rotations of the coordinate system (Karhunen-Loève or linear discriminant transform, [5, 6]); finally, the features belonging to the d first coordinate axes are selected to be fed into the vector quantizer.

The basic idea presented in this paper is to incorporate the feature rotation step together with the subsequent dimensionality reduction ($D \rightarrow d$) into the SCHMM formalism. Given the model structure (including the target dimension d) and a representative training sample, the particular transformation matrix will be optimized with respect to the likelihood function of the overall model, resulting in the so-called *maximum likelihood rotation* (MLR).

The rest of the paper is organized as follows: the feature transform HMM is defined in Sect. 2. In Sect. 3, we will provide the derivation of Baum-Welch training formulae for the extended model; particularly, we show that the rotation-dependent part can be separated from the conventional SCHMM part of the Kullback-Leibler function. Sect. 4 is devoted to the of numerical optimization of the MLR expression as a function of the unknown rotation parameters.

2. FEATURE-TRANSFORM MARKOV MODELS

The basic idea underlying the *feature transform HMM* (FTHMM) is illustrated in Figure 1; we want to re-

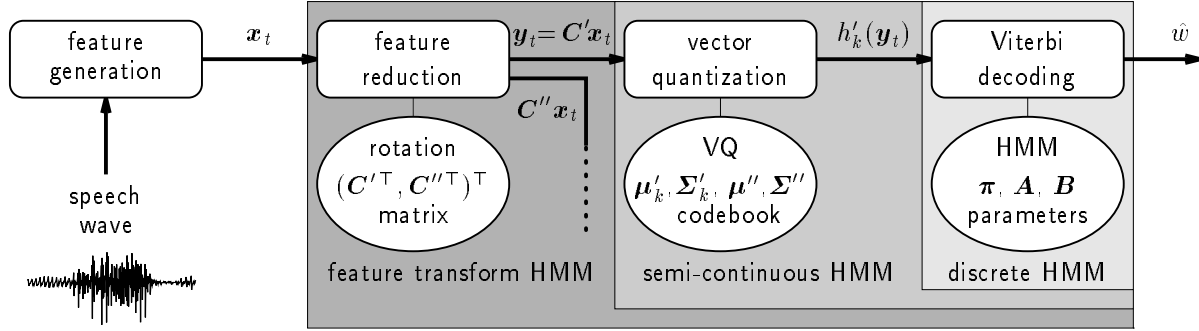


Figure 1: Different HMM-based speech recognizer architectures

place the ordinary semi-continuous HMM by the enlarged model

$$\boldsymbol{\lambda} = (\boldsymbol{\pi}, \mathbf{A}, \mathbf{B}, \mathbf{g}, \mathbf{C}) \quad ,$$

with initial probabilities $\boldsymbol{\pi}$, transition probabilities \mathbf{A} , mixture coefficients \mathbf{B} , a parameterized codebook \mathbf{g} of mixture components, and the feature transformation matrix \mathbf{C} . The crucial issue of the FTHMM is that $\boldsymbol{\lambda}$ is designed to model the production of sequences of the original feature vectors $\mathbf{x}_t \in \mathbb{R}^D$; this is opposite to the SCHMM which generates the reduced vectors $\mathbf{y}_t \in \mathbb{R}^d$ (see Figure 1).

If the mapping $\mathbf{x}_t \mapsto \mathbf{y}_t$ is assumed to be singular — which is true, of course, if the dimension of the feature space is actually reduced — no inverse exists, and no probabilistic model generating the \mathbf{x}_t can be found. This problem is circumvented by adopting a “soft” variant of feature space reduction.

$$b_j(\mathbf{x}) = \sum_{k=1}^K b_{jk} g_k(\mathbf{x}) = h''(\mathbf{C}''\mathbf{x}) \cdot \sum_{k=1}^K b_{jk} h'_k(\mathbf{C}'\mathbf{x})$$

The k -th semi-continuous mixture component g_k of our model is factorized into

- a *codebook-class dependent* distribution h'_k operating on the most discriminative part $\mathbf{C}'\mathbf{x} \in \mathbb{R}^d$ of \mathbf{x} and
- a *codebook-class independent* distribution h'' operating on the uninformative noise part $\mathbf{C}''\mathbf{x} \in \mathbb{R}^{D-d}$ of \mathbf{x} .

The rows of both \mathbf{C}' and \mathbf{C}'' together form the *orthonormal* $D \times D$ matrix $\mathbf{C} = (\mathbf{C}'^\top, \mathbf{C}''^\top)^\top$, i.e., the product $\mathbf{C}\mathbf{C}^\top$ equals the identity matrix \mathbf{I}_D . The initial components y_1, \dots, y_d of the rotated vector $\mathbf{C}\mathbf{x}$ should carry the relevant information for speech recognition. The remaining coefficients y_{d+1}, \dots, y_D will not contribute to word or sentence identification since their probability $h''(\mathbf{C}''\mathbf{x})$ may be factored out as a constant term in the above expression for $b_j(\mathbf{x})$. Consequently, the computation of $h''(\mathbf{C}''\mathbf{x})$ becomes obsolete, and only the class-conditional probabilities $h'_k(\mathbf{C}'\mathbf{x})$ — recall that d is usually much smaller than D — have to be provided during Viterbi decoding.

3. BAUM-WELCH REESTIMATION

Let $\mathbf{X} = \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T$ be a sequence of D -dimensional real valued feature vectors for estimating the parameters. In classical HMMs the EM algorithm is used for solving the incomplete data estimation problem: it is hidden which state sequence produces the observable output data. Let $\mathbf{q} = q_1, \dots, q_T$ be the sequence of state indices from an N -state model. Dealing with semi-continuous HMMs, we need furthermore a specification of the mixture components. Let $\mathbf{k} = k_1, k_2, \dots, k_T$ denote a sequence of mixture density components. The pair (\mathbf{q}, \mathbf{k}) represents that at time t the state s_{q_t} uses the density component g_{k_t} for generating an output vector.

The estimation of the FTHMM parameters is performed iteratively using the EM Algorithm [7]. For that purpose we have to maximize the Kullback-Leibler statistics

$$Q(\boldsymbol{\lambda}, \hat{\boldsymbol{\lambda}}) = \sum_{\mathbf{q}} \sum_{\mathbf{k}} P(\mathbf{X}, \mathbf{q}, \mathbf{k} | \boldsymbol{\lambda}) \cdot \log P(\mathbf{X}, \mathbf{q}, \mathbf{k} | \hat{\boldsymbol{\lambda}})$$

with respect to $\hat{\boldsymbol{\lambda}}$. After expanding the log likelihood expression of the right hand side of the above equation and reordering the summation $Q(\boldsymbol{\lambda}, \hat{\boldsymbol{\lambda}})$ separates into three well-known terms (see, for instance, [2, p.160]) depending exclusively on $\boldsymbol{\pi}$, \mathbf{A} , and \mathbf{B} , respectively. The remaining expression

$$Qg(\boldsymbol{\lambda}, \hat{\boldsymbol{\lambda}}) = \sum_{k=1}^K \sum_{t=1}^T \left(\sum_{j=1}^N \zeta_t(j, k) \right) \cdot \log \hat{g}_k(\mathbf{x}_t)$$

has to be maximized in order to get reestimates of the codebook parameters; the *a posteriori* expectations

$$\zeta_t(j, k) = P(q_t = j, k_t = k | \mathbf{X}, \boldsymbol{\lambda})$$

are obtained through the classical forward-backward computations [2].

For a standard SCHMM, multivariate Gaussian mixture components $g_k(\mathbf{x}) = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$ with mean vectors $\boldsymbol{\mu}_k$ and covariance matrices $\boldsymbol{\Sigma}_k$ are assumed, and maximization of $Qg(\cdot, \cdot)$ results in the well-known ML estimates

$$\hat{\boldsymbol{\mu}}_k = \sum_{t=1}^T \frac{z_t(k)}{Z_k} \cdot \mathbf{x}_t, \quad \hat{\boldsymbol{\Sigma}}_k = \sum_{t=1}^T \frac{z_t(k)}{Z_k} \cdot (\mathbf{x}_t - \hat{\boldsymbol{\mu}}_k)(\mathbf{x}_t - \hat{\boldsymbol{\mu}}_k)^\top$$

for the semi-continuous codebook; we used the abbreviations $z_t(k) = \sum_j \zeta_t(j, k)$ and $Z_k = \sum_t z_t(k)$.

However, in the case of FTHMM mixtures

$$g_k(\mathbf{x}) = \mathcal{N}(\mathbf{C}'\mathbf{x} \mid \hat{\boldsymbol{\mu}}'_k, \hat{\boldsymbol{\Sigma}}'_k) \cdot \mathcal{N}(\mathbf{C}''\mathbf{x} \mid \hat{\boldsymbol{\mu}}'', \hat{\boldsymbol{\Sigma}}'')$$

with the means and covariances $\hat{\boldsymbol{\mu}}'_k, \hat{\boldsymbol{\Sigma}}'_k, \hat{\boldsymbol{\mu}}'',$ and $\hat{\boldsymbol{\Sigma}}''$ of the rotated Gaussians h'_k and h'' the separation of $Qg(\cdot, \cdot)$ becomes much more intricate because of the linear coupling between the partial matrices \mathbf{C}' and \mathbf{C}'' of the underlying feature transform.

Fortunately, the Kullback-Leibler function easily decomposes into $K + 1$ independent subexpressions

$$Qg(\boldsymbol{\lambda}, \hat{\boldsymbol{\lambda}}) = \sum_{k=1}^K Q_{h'_k}(\boldsymbol{\lambda}, \hat{\boldsymbol{\lambda}}) + Q_{h''}(\boldsymbol{\lambda}, \hat{\boldsymbol{\lambda}})$$

if we assume \mathbf{C}' and \mathbf{C}'' fixed for a moment. Setting the partial derivatives of the terms $Q_{h'_k}(\boldsymbol{\lambda}, \hat{\boldsymbol{\lambda}})$ to zero and solving for the unknown distribution parameters, we obtain the reestimates

$$\hat{\boldsymbol{\mu}}'_k = \mathbf{C}'\hat{\boldsymbol{\mu}}_k \quad \text{and} \quad \hat{\boldsymbol{\Sigma}}'_k = \mathbf{C}'\hat{\boldsymbol{\Sigma}}_k\mathbf{C}'^\top, \quad k = 1, \dots, K,$$

where $\hat{\boldsymbol{\mu}}_k$ and $\hat{\boldsymbol{\Sigma}}_k$ are the means and covariances of the standard SCHMM codebook distributions. Following analogous arguments, maximization of the likelihood function for the global FTHMM distribution h'' leads to the equations

$$\hat{\boldsymbol{\mu}}'' = \mathbf{C}''\hat{\boldsymbol{\mu}}, \quad \hat{\boldsymbol{\mu}} = \frac{1}{T} \sum_{t=1}^T \mathbf{x}_t$$

as well as

$$\hat{\boldsymbol{\Sigma}}'' = \mathbf{C}''\hat{\boldsymbol{\Sigma}}\mathbf{C}''^\top, \quad \hat{\boldsymbol{\Sigma}} = \frac{1}{T} \sum_{t=1}^T (\mathbf{x}_t - \hat{\boldsymbol{\mu}})(\mathbf{x}_t - \hat{\boldsymbol{\mu}})^\top$$

In order to solve the estimation problem for the rotation parameters \mathbf{C}' and \mathbf{C}'' we have to plug in the estimates derived above into our target function $Qg(\boldsymbol{\lambda}, \hat{\boldsymbol{\lambda}})$ which yields the expression

$$\sum_{k=1}^K \underbrace{\left(\sum_{t=1}^T z_t(k) \cdot \log \mathcal{N}(\mathbf{C}'\mathbf{x}_t \mid \mathbf{C}'\hat{\boldsymbol{\mu}}_k, \mathbf{C}'\hat{\boldsymbol{\Sigma}}_k\mathbf{C}'^\top) \right)}_{E_k} + \sum_{t=1}^T \log \mathcal{N}(\mathbf{C}''\mathbf{x}_t \mid \mathbf{C}''\hat{\boldsymbol{\mu}}, \mathbf{C}''\hat{\boldsymbol{\Sigma}}\mathbf{C}''^\top).$$

Since the weights $z_t(k)$ of the k th expectation term E_k are identical to the weights involved in the estimation of $\hat{\boldsymbol{\mu}}_k$ and $\hat{\boldsymbol{\Sigma}}_k$, the equation

$$E_k = -\frac{Z_k}{2} \cdot \left(\log |2\pi\mathbf{C}'\hat{\boldsymbol{\Sigma}}_k\mathbf{C}'^\top| - d \right)$$

can be shown to be valid. A quite similar argument holds for the right hand side of $Qg(\cdot, \cdot)$. After division by a factor $-T/2$ and elimination of additive terms irrelevant to the optimization problem at hand we are left with the expression

$$\ell(\mathbf{C}', \mathbf{C}'') = \sum_{k=1}^K \left(P_k \cdot \log |\mathbf{C}'\hat{\boldsymbol{\Sigma}}_k\mathbf{C}'^\top| \right) + \log |\mathbf{C}''\hat{\boldsymbol{\Sigma}}\mathbf{C}''^\top|$$

where the weights P_k denote the ratio Z_k/T . We conclude that in order to obtain the desired MLR matrix $\mathbf{C} = (\mathbf{C}'^\top, \mathbf{C}''^\top)^\top$ the log likelihood function $\ell(\mathbf{C}) = \ell(\mathbf{C}', \mathbf{C}'')$ has to be *minimized* subject to the orthonormality restrictions $\mathbf{C}\mathbf{C}^\top = \mathbf{I}_D$.

4. OPTIMIZATION OF THE MLR MATRIX

Unfortunately, the function ℓ presents a delicate structure, including log determinants of non-quadratic matrix products. Moreover, the situation is complicated by the implicit coupling of \mathbf{C}' and \mathbf{C}'' based on the orthonormality constraint of \mathbf{C} . This fact prevents us from applying a straightforward gradient descent algorithm, controlled by the partial derivatives

$$\mathbf{H}' = \nabla_{\mathbf{C}'} \ell(\mathbf{C}) = \sum_{k=1}^K \left[P_k \cdot \left(\mathbf{C}'\hat{\boldsymbol{\Sigma}}_k\mathbf{C}'^\top \right)^{-1} \mathbf{C}'\hat{\boldsymbol{\Sigma}}_k \right]$$

and

$$\mathbf{H}'' = \nabla_{\mathbf{C}''} \ell(\mathbf{C}) = \left(\mathbf{C}''\hat{\boldsymbol{\Sigma}}\mathbf{C}''^\top \right)^{-1} \mathbf{C}''\hat{\boldsymbol{\Sigma}}$$

(see [8, p.568]). In order to overcome the latter complication, our optimization task has to be transformed towards an unrestricted problem. Our crucial device for that enterprise is the observation that every rotation \mathbf{C} of \mathbb{R}^D may be decomposed into a product

$$\mathbf{C}(\phi) = \mathbf{C}(\phi_1, \dots, \phi_R) = \prod_{r=1}^R \mathbf{U}_{p_r q_r}(\phi_r)$$

of *elementary rotations* of the special form

$$\mathbf{U}_{pq}(\phi) = \mathbf{I}_D + \cos \phi (\mathbf{e}_p \mathbf{e}_p^\top + \mathbf{e}_q \mathbf{e}_q^\top) + \sin \phi (\mathbf{e}_p \mathbf{e}_q^\top - \mathbf{e}_q \mathbf{e}_p^\top)$$

where \mathbf{e}_p denotes the p th unit vector and the index pairs (p_r, q_r) range over all combinations $1 \leq p < q \leq D$, i.e., $R = (D-1) \cdot D/2$. From the geometrical point of view, $\mathbf{U}_{pq}(\phi)$ describes a rotation of \mathbb{R}^D in the (p, q) -plane by the angle ϕ . Henceforth, two basically different approaches are possible, according to whether the free parameters ϕ_r shall be improved independently from another or not.

In the former case we proceed quite similar to the Jacobian coordinate descent algorithm for eigenvector computation. Iteratively, a declining sequence $\mathbf{C}_0, \mathbf{C}_1, \dots$ of feature transformations

$$\mathbf{C}_0 = \mathbf{I}_D, \quad \mathbf{C}_s = \mathbf{C}_{s-1} \cdot \mathbf{U}_{p_s q_s}(\phi_s)$$

is created subject to the requirement $\ell(\mathbf{C}_s) \leq \ell(\mathbf{C}_{s-1})$. The task of generating an improved rotation in step s , then, reduces to a one-dimensional minimization problem in the variable ϕ_s . For the control of some standard descent algorithm computing the partial derivative of $\ell(\mathbf{C}_s)$ with respect ϕ_s is required. The total derivative results in

$$\frac{\partial \ell(\mathbf{C}_s)}{\partial \phi_s} = \sum_{i=1}^D \sum_{j=1}^D \frac{\partial \ell(\mathbf{C}_s)}{\partial (C_s)_{ij}} \cdot \frac{\partial (C_s)_{ij}}{\partial \phi_s}.$$

Due to the fact that the rotations depend on each other

the required number of iteratively applied elementary rotations will generally exceed R .

If, however, the entries of the parameter vector ϕ representing the transformation matrix are to be jointly optimized, i.e., a successor ϕ^+ of simultaneously improved angles is sought for in each iteration step, the problem stays multidimensional. Similar to the Jacobian method the gradient vector $\nabla_{\phi} \ell(\mathbf{C}(\phi))$ for this multivariate function is given by

$$\left(\frac{\partial \ell(\mathbf{C}(\phi))}{\partial \phi_r} \right)_{1 \leq r \leq R} = \left(\sum_{i=1}^D \sum_{j=1}^D \frac{\partial \ell(\mathbf{C}(\phi))}{\partial C_{ij}} \cdot \frac{\partial C_{ij}}{\partial \phi_r} \right)_{1 \leq r \leq R}.$$

The left hand side derivatives inside the summation can be drawn from the matrices \mathbf{H}' , \mathbf{H}'' as defined above. The rightmost expression involves the calculation of $\partial \mathbf{C} / \partial \phi_r$ for each r . It is easily shown that $\mathbf{C}_r(\phi)$ shares the product form with the original matrix \mathbf{C} except for the r th factor which amounts to $\partial \mathbf{U}_{p_r q_r}(\phi_r) / \partial \phi_r$ and is available through elementary calculus.

5. FIRST EXPERIMENTAL RESULTS

The MLR approach was evaluated running a simple frame-by-frame phone recognition task. For that purpose, a multivariate Gaussian classifier was designed using ten hours of speech for training and one hour for testing. The recognition rates achieved when 44 different phonetic labels had to be distinguished are shown in Table 1. Three different orthonormal mappings were employed in order to reduce the 36-dimensional space (12 cepstral coefficients taken from 3 neighboring speech frames) to the \mathbb{R}^{12} , decreasing the number of distribution parameters as well as computational complexity by roughly one order of magnitude.

rotation matrix	KL	LD	MLR
recognition rate	52.7%	53.2%	54.1%

Table 1: Labeling accuracy and feature rotations

As indicated in Table 1, the MLR mapping allows a slight improvement of recognition accuracy over the Karhunen-Loève (KL) or linear discriminant (LD) transforms. Whilst the MLR approach certainly presents a new target function which is superior to the variance criteria of KL and LD, the ML rotation matrix is undoubtedly much harder to optimize. The gradient descent methods tested so far proved quite impractical due to the enormous expense in computing the partial derivatives. Moreover, local optimization is a questionable strategy when confronted with an objective function as rugged as ℓ . Thus we moved to global minimizers without reference to derivatives such as the simplex algorithm [9] or combinatorial optimization procedures; the above result relates to the great deluge algorithm [10].

6. CONCLUSION

The FTHMM outlined above is an extension of the SCHMM formalism which incorporates a feature rotation matrix \mathbf{C} as a Maximum-Likelihood (ML) trainable component of the probabilistic model. The ML parameter estimation can be done straightforwardly by applying the EM algorithm for the computation of the free transform parameters, too. The MLR generalizes work on linear feature selection, for instance the KL and LD transformations.

Nevertheless, in contrast to classical linear transformation methods the resulting parameter estimation problem forces the use of numerical or combinatorial procedures for global optimization. Thus, the actual slight increase of recognition accuracy correlates with the solution of a more complex optimization within the off-line training stage.

7. REFERENCES

- [1] S.E. Levinson, L.R. Rabiner, and M.M. Sondhi. An Introduction to the Application of the Theory of Probabilistic Functions of a Markov Process to Automatic Speech Recognition. *Bell Systems Technical Journal*, 62(4):1035–1074, 1983.
- [2] X.D. Huang, Y. Ariki, and M.A. Jack. *Hidden Markov Models for Speech Recognition*. Number 7 in Information Technology Series. Edinburgh University Press, Edinburgh, 1990.
- [3] S. Furui. Speaker-Independent Isolated Word Recognition Using Dynamic Features of Speech Spectrum. *IEEE Trans. on Acoustics, Speech, and Signal Processing*, 34(1):52–59, 1986.
- [4] S.B. Davis and P. Mermelstein. Comparison of Parametric Representation for Monosyllabic Word Recognition in Continuously Spoken Sentences. *IEEE Trans. on Acoustics, Speech, and Signal Processing*, 28(4):357–366, 1980.
- [5] J.T. Tou and R.P. Heydorn. Some Approaches to Optimum Feature Extraction. In J.T. Tou, editor, *Computer and Information Sciences*, volume 2, pages 57–89. Academic Press, New York, 1967.
- [6] H.P. Friedman and J. Rudin. On Some Invariant Criteria for Grouping Data. *American Statistical Association Journal*, 12:1159–1178, 1967.
- [7] A.P. Dempster, N.M. Laird, and D.B. Rubin. Maximum Likelihood from Incomplete Data via the EM Algorithm. *J. Royal Statist. Soc. Ser. B*, 39(1):1–22, 1977.
- [8] K. Fukunaga. *Introduction to Statistical Pattern Recognition*. Academic Press, Boston, 1990.
- [9] H. A. Eiselt, G. Pederzoli, and C.-L. Sandblom. *Continuous Optimization Models*. Walter de Gruyter, Berlin, 1987.
- [10] G. Dueck, T. Scheuer, and H. Wallmeier. Tolerance Threshold and Great Deluge: New Ideas for Optimization (in German). *Spektrum der Wissenschaft*, pages 42–51, März 1993.