

TRUST-TECH BASED METHODS FOR
OPTIMIZATION AND LEARNING

A Dissertation

Presented to the Faculty of the Graduate School
of Cornell University

in Partial Fulfillment of the Requirements for the Degree of
Doctor of Philosophy

by

Chandankumar Reddy Karrem

May 2007

© 2007 Chandankumar Reddy Karrem

ALL RIGHTS RESERVED

TRUST-TECH BASED METHODS FOR
OPTIMIZATION AND LEARNING

Chandankumar Reddy Karrem, Ph.D.

Cornell University 2007

Many problems that arise in machine learning domain deal with nonlinearity and quite often demand users to obtain global optimal solutions rather than local optimal ones. Optimization problems are inherent in machine learning algorithms and hence many methods in machine learning were inherited from the optimization literature. Popularly known as the *initialization problem*, the ideal set of parameters required will significantly depend on the given initialization values. The recently developed TRUST-TECH (TRansformation Under STability-reTaining Equilibria CHaracterization) methodology systematically explores the subspace of the parameters to obtain a complete set of local optimal solutions. In this thesis work, we propose TRUST-TECH based methods for solving several optimization and machine learning problems. TRUST-TECH explores the dynamic and geometric characteristics of stability boundaries of a nonlinear dynamical system corresponding to the nonlinear function of interest. Basically, our method coalesces the advantages of the traditional local optimizers with that of the dynamic and geometric characteristics of the stability regions of the corresponding nonlinear dynamical system. Two stages namely, the local stage and the neighborhood-search stage, are repeated alternatively in the solution space to achieve improvements in the quality of the solutions. The local stage obtains the local maximum of the nonlinear function and the neighborhood-search stage helps to escape out of the local maximum by moving towards the neighboring stability regions. Our methods were tested on both synthetic and real datasets and the advantages of using this novel framework are

clearly manifested. This framework not only reduces the sensitivity to initialization, but also allows the flexibility for the practitioners to use various global and local methods that work well for a particular problem of interest. Other hierarchical stochastic algorithms like evolutionary algorithms and smoothing algorithms are also studied and frameworks for combining these methods with TRUST-TECH have been proposed and evaluated on several test systems.

BIOGRAPHICAL SKETCH

Chandan Reddy was born in Jammalamadugu, Andhra Pradesh, India on May 11, 1980. After obtaining Bachelor's degree from Pondicherry Engineering College in 2001, he moved to the United States to pursue his higher education. He obtained his master's degree from the Department of Computer Science and Engineering at Michigan State University in August 2003. Later, he joined the Department of Electrical and Computer Engineering at Cornell University to pursue his doctoral studies. His primary research interests are in the areas of Machine Learning, Data Mining, Optimization, Computational Statistics, Bioinformatics and Biomedical Imaging.

Dedicated to my parents

ACKNOWLEDGEMENTS

I would like to thank my thesis advisor, Prof. Hsiao-Dong Chiang, without whose guidance and involvement this work would not have been possible. His enthusiasm and inspiration were extremely helpful in successfully pursuing this research work. I would also like to thank my committee member Prof. John Hopcroft for his encouragement and useful discussions. I am grateful to Dr. Peter Doerschuk for serving on my thesis committee. Special thanks to Dr. Bala Rajaratnam for his help with a few statistical works related to my thesis.

I am deeply indebted to my father, who was my main inspiration for pursuing graduate studies and choosing to have an academic career. Needless to say I had plenty of great discussions and fun with many people in the department including Jeng-Heui, Warut, Choi, Seunghee, Siva, Bhavin, Wesley, Jay and Huachen. I was also fortunate to have a good company at Cornell who have kept my spirits up and my life interesting: Manpreet, Pallavi, Sudha, Raju, Kumar, Mahesh, Pankaj, Deepak, Pradeep, Ajay and Shobhit. A special thanks to the graduate secretaries for taking care of countless numbers of things without any hitch.

Last and the most important, my deepest gratitude goes to my dear mother for her love, and for making me the person I am.

TABLE OF CONTENTS

| | | |
|----------|--|-----------|
| 1 | Introduction | 1 |
| 1.1 | Machine Learning | 1 |
| 1.2 | Optimization Methods | 3 |
| 1.3 | Motivation for this Thesis | 6 |
| 1.4 | Contributions of this Thesis | 8 |
| 1.5 | Organization of this Thesis | 9 |
| 2 | Finding Saddle Points on Potential Energy Surfaces | 11 |
| 2.1 | Introduction | 11 |
| 2.2 | Relevant Background | 14 |
| 2.3 | Theoretical Background | 15 |
| 2.4 | A Stability Boundary based Method | 21 |
| 2.5 | Implementation Issues | 25 |
| 2.6 | Experimental Results | 30 |
| 2.6.1 | Test Case 1: Two-dimensional Potential Energy Surface | 30 |
| 2.6.2 | Test Case 2: Three-dimensional symmetric systems | 33 |
| 2.6.3 | Test Case 3: Higher Dimensional Systems | 36 |
| 2.6.4 | Special Cases : Eckhardt surface | 40 |
| 2.7 | Discussion | 43 |
| 3 | TRUST-TECH based Expectation Maximization for Learning Mixture Models | 47 |
| 3.1 | Relevant Background | 49 |
| 3.2 | Preliminaries | 51 |
| 3.2.1 | Mixture Models | 52 |

| | | |
|----------|---|-----------|
| 3.2.2 | Expectation Maximization | 53 |
| 3.2.3 | EM for GMMs | 55 |
| 3.2.4 | Nonlinear Transformation | 56 |
| 3.3 | TRUST-TECH based Expectation Maximization | 58 |
| 3.4 | Implementation Details | 60 |
| 3.5 | Results and Discussion | 62 |
| 3.5.1 | Synthetic Datasets | 64 |
| 3.5.2 | Real Datasets | 65 |
| 3.5.3 | Discussion | 66 |
| 4 | Motif Refinement using Neighborhood Profile Search | 71 |
| 4.1 | Relevant Background | 73 |
| 4.2 | Preliminaries | 75 |
| 4.2.1 | Problem Formulation | 75 |
| 4.2.2 | Dynamical Systems for the Scoring Function | 78 |
| 4.3 | Neighborhood Profile Search | 80 |
| 4.4 | Implementation Details | 85 |
| 4.5 | Experimental Results | 88 |
| 4.5.1 | Synthetic Datasets | 89 |
| 4.5.2 | Real Datasets | 93 |
| 5 | Component-wise Smoothing for Learning Mixture Models | 96 |
| 5.1 | Overview | 96 |
| 5.2 | Relevant Background | 97 |
| 5.3 | Preliminaries | 101 |
| 5.4 | Smoothing Log-Likelihood Surface | 103 |
| 5.4.1 | Kernel Parameters | 105 |

| | | |
|----------|---|------------|
| 5.4.2 | EM Updates | 106 |
| 5.5 | Algorithm and its implementation | 106 |
| 5.6 | Results and Discussion | 109 |
| 5.6.1 | Reduction in the number of local maxima | 112 |
| 5.6.2 | Smoothing for Initialization | 112 |
| 6 | TRUST-TECH based Neural Network Training | 121 |
| 6.1 | Overview | 121 |
| 6.2 | Relevant Background | 124 |
| 6.3 | Training Neural Networks | 128 |
| 6.4 | Problem Transformation | 131 |
| 6.5 | TRUST-TECH based Training | 132 |
| 6.6 | Implementation Details | 136 |
| 6.6.1 | Architecture and Local Methods | 136 |
| 6.6.2 | Initialization Schemes | 137 |
| 6.6.3 | TRUST-TECH | 137 |
| 6.7 | Experimental Results | 139 |
| 6.7.1 | Benchmark Datasets | 139 |
| 6.7.2 | Error Estimation | 141 |
| 6.7.3 | Classification Accuracy | 142 |
| 6.7.4 | Visualization | 144 |
| 7 | Evolutionary TRUST-TECH | 146 |
| 7.1 | Overview | 146 |
| 7.2 | Variants of Evolutionary Algorithms | 148 |
| 7.3 | Evolutionary TRUST-TECH Algorithm | 151 |
| 7.4 | Experimental Results | 154 |

| | | |
|----------|--|------------|
| 7.5 | Parallel Evolutionary TRUST-TECH | 155 |
| 8 | Conclusion and Future Work | 156 |
| 8.1 | Conclusion | 156 |
| 8.2 | Future Work | 159 |
| | Bibliography | 162 |

LIST OF FIGURES

| | | |
|-----|---|----|
| 1.1 | (a) Supervised learning with data points from two different classes separated by a hyperplane represented using a dashed line. (b) Un-supervised Learning or data clustering - two well separated clusters. | 3 |
| 1.2 | Different optimization methods. | 4 |
| 1.3 | Various stages of TRUST-TECH (a) The dark regions indicate the promising subspaces. (b) The dots indicate the local optimal solutions and the dotted arrows indicate the convergence of local optimization method from a given initial point. | 7 |
| 1.4 | Organization chart of this thesis. The main contributions are in the areas of optimization and learning. | 10 |
| 2.1 | The surface and contour plots of a two-dimensional energy function. A saddle point (x_d) is located between two local minima (x_s^1 and x_s^2). x_m^1 and x_m^2 are two local maxima located in the orthogonal direction. | 13 |
| 2.2 | Phase potrait of a gradient system. The solid lines with solid arrows represent the basin boundary. $\partial A_p(x_s^1) = \bigcup_{i=1}^3 \overline{W^s(x_d^i)}$. The local minima x_s^1 and x_s^2 correspond to the stable equilibrium points of the gradient system. The saddle point (x_d^1) corresponds to the dynamic decomposition point that connects the two stable equilibrium points. | 18 |

| | | |
|-----|---|----|
| 2.3 | Contour plot of a 2-D LEPS potential (described in appendix-A). Each line represents the values of a constant potential. A and B are the two local minima. DDP is the dynamic decomposition point to be computed. The search direction is the direction of the vector joining AB . The exit point (X_{ex}) is obtained by finding the peak of the function value along this vector. The dotted line indicates the stability boundary. The dashed line indicates the search direction. | 22 |
| 2.4 | Illustration of the step 3 of our algorithm. m_1 is integrated till m'_1 is reached and traced back along the vector m'_1B to get m_2 , and so on. m_n are the n exit points on the stability boundary. The n^{th} point is the Minimum Gradient Point where the magnitude of the gradient (G_M) is the minimum amongst all the computed gradient values along the stability boundary ($X_{ex,m_2,m_3..m_n}$). DDP is the dynamic decomposition point. | 24 |
| 2.5 | The plot of the function value along the vector AB . The curve monotonically increases starting from A and then decreases till it reaches B . X_{ex} is the exit point where the function attains its peak value. $a_1, a_2, \dots a_9$ are the interval points. The marked circles indicate that the function has been evaluated at these points. The empty circles indicate the points where the function is yet to be evaluated. | 26 |
| 2.6 | The interval point a_5 can be present either to the left or to the right of the peak. When a_6 is reached, the golden section search method is invoked with a_6 and a_4 as the interval. | 27 |

| | | |
|------|---|----|
| 2.7 | Two dimensional contour plot of the potential energy surface corresponding to the muller-brown function described in Eq. (2.12). A,B and C are stable equilibrium points and DDP1,DDP2 are two dynamic decomposition points. The dashed lines indicate the initial search direction. The dots indicate the results of the stability boundary following procedure. These points are the points on the stability boundary that reach the MGP. | 31 |
| 2.8 | The gradient curve corresponding to the various points obtained from the stability boundary following procedure. The graph shows that the magnitude of the gradient slowly increases in the initial phases and then starts to reduce. The highlighted point corresponds to the gradient at the MGP. | 33 |
| 2.9 | Characteristic curves of Lennard-Jones potential and the Morse potential with all parameters set to unity. Solid line represents Lennard-Jones potential (2.13) and the dashed line represents the Morse potential (2.15). | 35 |
| 2.10 | Top view of the surface and the seven atom island on the surface of an FCC crystal. The shading indicates the height of the atoms. . . . | 37 |
| 2.11 | Some sample configurations of the heptamer island on the surface of the FCC crystal. First row- saddle points configurations. Second row- other local minimum energy configurations corresponding to the above saddle point configurations. | 39 |
| 2.12 | The gradient curve obtained in one of the higher dimensional test cases. MGP indicates the minimum gradient point where the magnitude of the gradient reaches the minimum value. This point is used as a initial guess for a local optimization method to obtain DDP. . . | 39 |

| | | |
|------|--|----|
| 2.13 | Two dimensional contour plot of the potential energy surface of the Eckhardt energy function described on Eq. (2.18). A and B are stable equilibrium points and DDP1,DDP2 are two dynamic decomposition points. M is a source. The dots correspond to the points on the stability boundary during the stability boundary following procedure. | 41 |
| 2.14 | Gradient curve corresponding to the various points along the stability boundary on the Eckhardt surface. | 42 |
| 3.1 | Data consisting of three Gaussian components with different mean and variance values. Note that each data point doesn't have a hard membership that it belongs to only one component. Most of the points in the first component will have high probability with which they belong to it. In this case, the other components do not have much influence. Components 2 and 3 data points are not clearly separated. The problem of learning mixture models involves estimating the parameters of the Gaussian components and finding the probabilities with which each data sample belongs to the component. | 49 |
| 3.2 | Various stages of our algorithm in (a) Parameter space - the solid lines indicate the practical stability boundary. Points highlighted on the stability boundary are the decomposition points. The dotted arrows indicate the convergence of the EM algorithm. The dashed lines indicate the neighborhood-search stage. x_1 and x_2 are the exit points on the practical stability boundary (b) Different points in the function space and their corresponding log-likelihood function values. | 58 |

| | | |
|-----|--|----|
| 3.3 | Parameter estimates at various stages of our algorithm on the three component Gaussian mixture model (a) Poor random initial guess (b) Local maximum obtained after applying EM algorithm with the poor initial guess (c) Exit point obtained by our algorithm (d) The final solution obtained by applying the EM algorithm using the exit point as the initial guess. | 63 |
| 3.4 | Graph showing likelihood vs Evaluations. A corresponds to the original local maximum (L=-3235.0). B corresponds to the exit point (L=-3676.1). C corresponds to the new initial point (L=-3657.3) after moving out by ‘ ϵ ’. D corresponds to the new local maximum (L=-3078.7). | 64 |
| 4.1 | Synthetic DNA sequences containing some instance of the pattern ‘CCGATTACCGA’ with a maximum number of 2 mutations. The motifs in each sequence are highlighted in the box. We have a (11,2) motif where 11 is the length of the motif and 2 is the number of mutations allowed. | 72 |
| 4.2 | Diagram illustrates the TRUST-TECH method of escaping from the original solution (A) to the neighborhood local optimal solutions (a_{1i}) through the corresponding exit points (e_{1i}). The dotted lines indicate the local convergence of the EM algorithm. | 82 |
| 4.3 | A summary of escaping out of the local optimum to the neighborhood local optimum. Observe the corresponding trend of the alignment score at each step. | 83 |

| | | |
|-----|--|-----|
| 4.4 | 2-D illustration of first tier improvements in a $3l$ dimensional objective function. The original local maximum has a score of 163.375. The various Tier-1 solutions are plotted and the one with highest score (167.81) is chosen. | 88 |
| 4.5 | The average scores with the corresponding first tier and second tier improvements on synthetic data using the random starts with TRUST-TECH approach with different (l,d) motifs. | 92 |
| 4.6 | The average scores with the corresponding first tier and second tier improvements on synthetic data using the Random Projection with TRUST-TECH approach with different (l,d) motifs. | 92 |
| 5.1 | Block diagram of the traditional approach and the smoothing approach. | 99 |
| 5.2 | Smoothing a nonlinear surface at different levels. Tracing the global maxima (A, B and C) at different levels. | 100 |
| 5.3 | Different convolution kernels (a) Triangular function (b) Step function and (c) Gaussian function | 102 |
| 5.4 | The effects of smoothing a Gaussian density function with a Gaussian kernel. | 103 |
| 5.5 | Block Diagram of the smoothing approach. Smooth likelihood surface is obtained by convolving the original likelihood surface with a convolution kernel which is chosen to be a Gaussian kernel in our case. | 105 |
| 5.6 | Flowchart of the smoothing algorithm | 107 |
| 5.7 | True mixture of the three Gaussian components with 900 samples. . | 109 |
| 5.8 | True mixtures of the more complicated overlapping Gaussian case with 1000 samples. | 110 |

| | | |
|------|--|-----|
| 5.9 | Various stages during the smoothing process. (a) The original log-likelihood surface which is very rugged (b)-(c) Intermediate smoothed surfaces (d) Final smoothed surface with only two local maxima. | 111 |
| 5.10 | Reduction in the number of local maxima for various datasets. | 113 |
| 6.1 | The Architecture of MLP with single hidden layer having k hidden nodes and the output layer having a single output node. x_1, x_2, \dots, x_n is an n -dimensional input feature vector. w_{ij} are the weights and b_1, b_2, \dots, b_k are the biases for these k nodes. The activation function for the hidden nodes is sigmoidal and for the output node is linear. | 122 |
| 6.2 | Comparison between the two frameworks (a) Traditional approach and (b) TRUST-TECH based approach. The main difference is the inclusion of the stability region based neighborhood-search stage that can explore the neighborhood solutions. | 123 |
| 6.3 | Block Diagram of the TRUST-TECH based training method. | 127 |
| 6.4 | Flow chart of our method | 133 |
| 6.5 | Diagram Illustrating two tier exit point strategy. The 'A*' represents the initial guess. Dotted arrows represent the convergence of the local solver. Solid arrows represent the gradient ascent linear searches along eigenvector directions. 'X' indicates a new initial condition in the neighboring stability region. M represents the local minimum obtained by applying local methods from 'X'. A_{1i} indicates Tier-1 local minima. e_{1i} are the exit points between M and A_{1i} . Similarly, A_{2j} and e_{2j} are the second-tier local minima and their corresponding exit points respectively. | 135 |

| | | |
|-----|---|-----|
| 6.6 | Spider web diagrams showing the tier-1 and tier-2 improvements using TRUST-TECH method on various benchmark datasets. The basis two axes are chosen arbitrarily and the vertical axis represents the improvements in the classifier accuracy. The distances between each tier are normalized to unity. | 145 |
| 7.1 | Topology of the nonlinear surface discussed in the introduction chapter. The initial point obtained after the recombination is ‘s’. The original concept of mutations will randomly perturb the point in the parameter space and obtain different points (s_1, s_2 and s_3). Two different evolutionary models (a) Local refinement strategy where the local optimization method is applied with ‘s’ as initial guess to obtain ‘A’. (b) Evolutionary TRUST-TECH methodology where the surrounding solutions are explored systematically after the recombination. | 152 |

LIST OF TABLES

| | | |
|-----|---|----|
| 2.1 | Stable equilibrium points and DDPs for the Muller-Brown Surface described in Eq. (2.12). | 32 |
| 2.2 | Stable equilibrium points and dynamic decomposition points for the Muller-Brown Surface described in Eq. (2.12). | 33 |
| 2.3 | Stable equilibrium points and dynamic decomposition points for the three atom Lennard Jones Cluster described in Eq. (2.13). | 36 |
| 2.4 | Results of our algorithm on a heptamer island over the FCC crystal. The number of force evaluations made to compute the MGP is given in the last column. | 40 |
| 2.5 | Stable equilibrium points and dynamic decomposition points for the Eckhardt Surface. | 42 |
| 3.1 | Description of the Notations used | 52 |
| 3.2 | Performance of TRUST-TECH-EM algorithm on an average of 100 runs on various synthetic and real datasets compared with random start EM algorithm | 65 |
| 3.3 | Comparison of TRUST-TECH-EM with other methods | 66 |
| 3.4 | Number of iterations taken for the convergence of the best solution. | 67 |
| 4.1 | A count of nucleotides A, T, G, C at each position $k = 1..l$ in all the sequences of the data set. $k = 0$ denotes the background count. | 76 |
| 4.2 | A count of nucleotides $j \in \{A, T, G, C\}$ at each position $k = 1..l$ in all the sequences of the data set. C_k is the k^{th} nucleotide of the consensus pattern which represents the nucleotide with the highest value in that column. Let the consensus pattern be GACT...G and b_j be the background. | 79 |

| | | |
|-----|---|-----|
| 4.3 | The consensus patterns and their corresponding scores of the original local optimal solution obtained from multiple random starts on the synthetic data. The best first tier and second tier optimal patterns and their corresponding scores are also reported. | 91 |
| 4.4 | The results of performance coefficient with $m = 1$ on synthetically generated sequences. The IC scores are not normalized and the perfect score is 20 since there are 20 sequences. | 93 |
| 4.5 | Results of TRUST-TECH method on biological samples. The real motifs were obtained in all the six cases using the TRUST-TECH framework. | 94 |
| 5.1 | Comparison of smoothing algorithm with the random starts. Mean and standard deviations across 100 random starts are reported. . . . | 114 |
| 6.1 | Description of the notations used | 129 |
| 6.2 | Summary of Benchmark Datasets. | 141 |
| 6.3 | Percentage improvements in the classification accuracies over the training and test data using TRUST-TECH with multiple random restarts. | 143 |
| 6.4 | Percentage improvements in the classification accuracies over the training and test data using TRUST-TECH with MATLAB initialization. | 143 |
| 7.1 | Results of Evolutionary TRUST-TECH model | 154 |

Chapter 1

Introduction

The problem of finding a global optimal solution arise in many disciplines ranging from science to engineering. In real world applications, multi-dimensional objective functions usually contain a large number of local optimal solutions. Obtaining a global optimal solution is of primary importance in these applications and is a very challenging problem. Some examples of these applications are : molecular confirmation prediction [22], VLSI design in microelectronics [95], resource allocation problems [41], design of wireless networks [69], financial decision making [87], structural engineering [58] and parameter estimation problems [53]. In this thesis, the primary focus is on the parameter estimation problems that arise in the field of machine learning.

1.1 Machine Learning

Machine learning algorithms can be broadly classified into two categories [44]: (i) Supervised learning and (ii) Unsupervised learning. The primary goal in *supervised learning* is to learn a mapping from x to y given a training dataset which consists of pairs (x_i, y_i) , where $x_i \in \mathcal{X}$ are the data points and $y_i \in \mathcal{Y}$ are the labels (or targets). A standard assumption is that the pairs (x_i, y_i) are sampled i.i.d. from some distribution. If y takes values in a finite set (discrete values) then it is a classification problem and if it takes values in a continuous space, then it is a regression problem. Support vector machines [21], artificial neural networks [68] and boosting [61] are the most popular algorithms for supervised learning. All these algorithms will construct a classification (or regression) model based on certain training data available. Usually, the effectiveness of any algorithm is evaluated using testing data

which is separate from the training data. In this thesis, we will primarily focus on *artificial neural networks* and estimating the parameters of its model. Constructing a model using artificial neural network involves estimating the parameters of the model that can effectively exploit the potential of the model. These parameters are usually obtained by finding the global minimum on the error surface. More details on training neural networks will be presented in Chapter 6.

Unsupervised learning, on the other hand, will train models using only the data-points without the target values. In simple terms, only x values are available without y values. Problems like outlier detection, density estimation, data clustering and noise removal fall under this category. Data clustering is one of the widely studied unsupervised learning topics [82]. Density estimation is a more generalized notion of data clustering. It involves in estimating and understanding the underlying distribution of the data. Applications of density estimation include trend analysis and data compression. Typically, one would like to estimate the parameters of a model that consists of multiple components of varying densities. More details on these mixture models and the use of expectation maximization (EM) algorithm for parameter estimation of these models will be presented in Chapter 3.

Fig. 1.1 shows the supervised and unsupervised learning scenarios. In supervised learning, the main goal is to train a model such that a final target class can be estimated for a new (unseen) data point. In a simple binary classification problem, a hyperplane (indicated by a dashed line) separates both the classes. In clustering problems, the main goal is to form groupings of the data and obtain any interesting structure (or patterns) in the data (see Fig. 1.1(b)).

In both the models mentioned above (neural networks and mixture models), estimating the parameters correspond to obtaining a global optimal solution on a highly nonlinear surface. The surface can be generated based on a function that

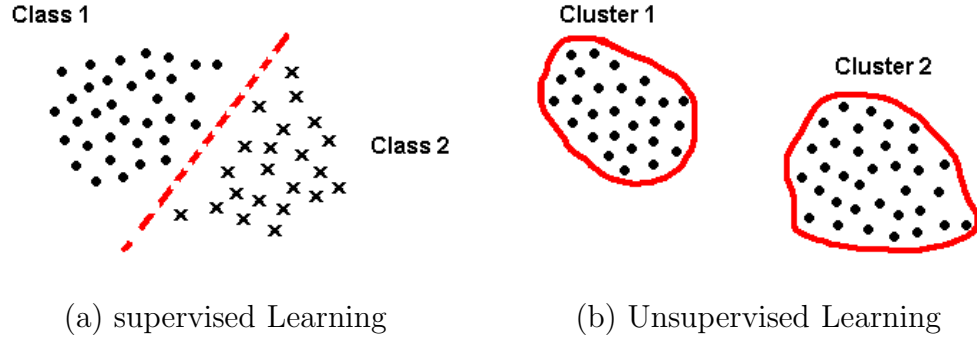


Figure 1.1: (a) Supervised learning with data points from two different classes separated by a hyperplane represented using a dashed line. (b) Unsupervised Learning or data clustering - two well separated clusters.

might represent the error in the training data or the likelihood of the data given the model. We will now introduce different categories of the nonlinear optimization methods studied in the literature.

1.2 Optimization Methods

Optimization methods can be broadly classified into two categories: *global methods* and *local methods*. Global methods explore the entire search space and obtain promising regions that have a higher probability of finding a global optimal solution. They can be either deterministic or stochastic in nature depending on the usage of some random component in the algorithm [78]. Deterministic global methods include branch and bound [79], homotopy based [59], interval analysis [66] and trajectory-based [39]. These methods are also known as exact methods. Stochastic global methods include techniques such as evolutionary algorithms [3], simulated annealing [88], tabu search [54] and ant colony optimization [42] which can give asymptotic guarantees of finding the global optimal solution. Many heuristic search algorithms can be incorporated into these stochastic methods to improve the performance of the algorithm in terms of quality of the solution and the speed of convergence.

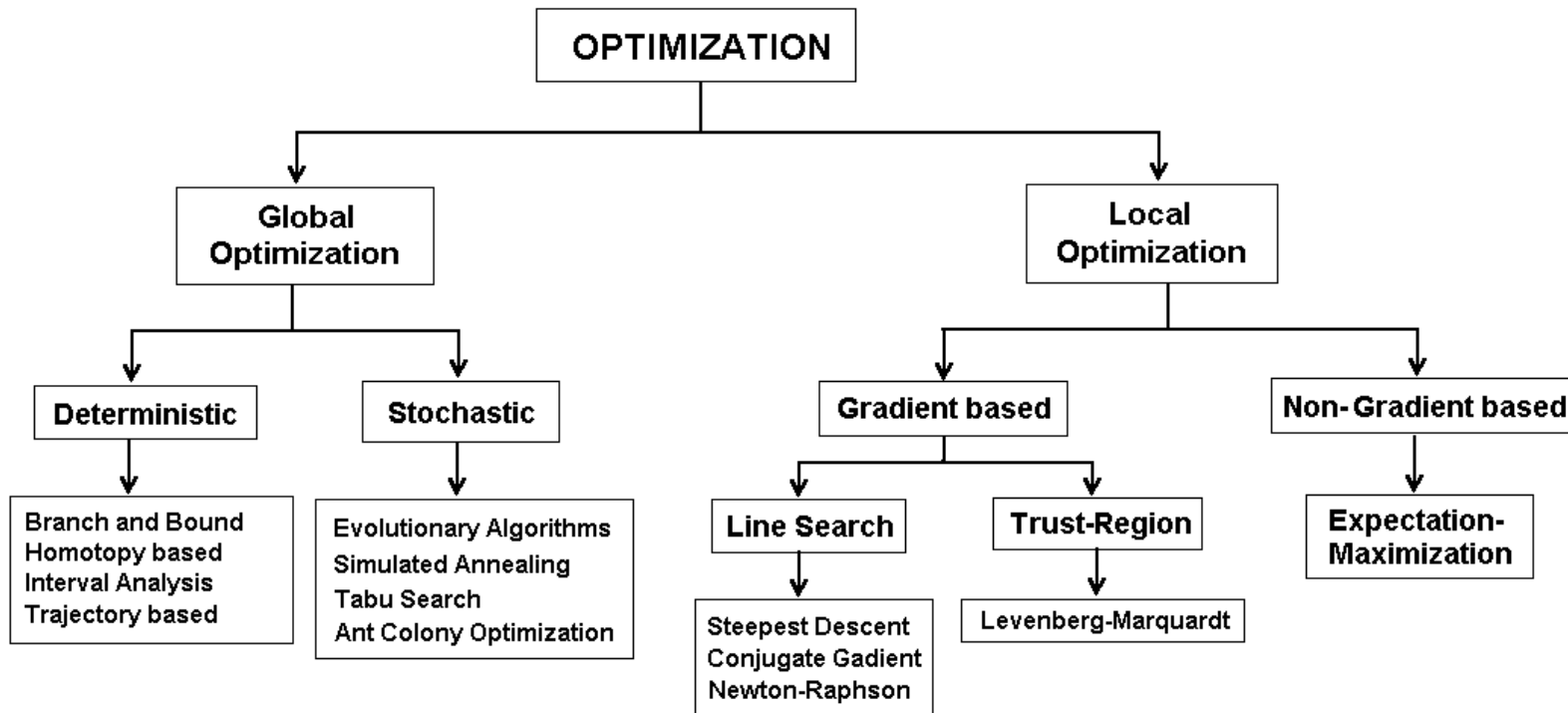


Figure 1.2: Different optimization methods.

Usually, local methods are deterministic in nature and can be grouped into two categories: *gradient based* and *non-gradient based methods*.

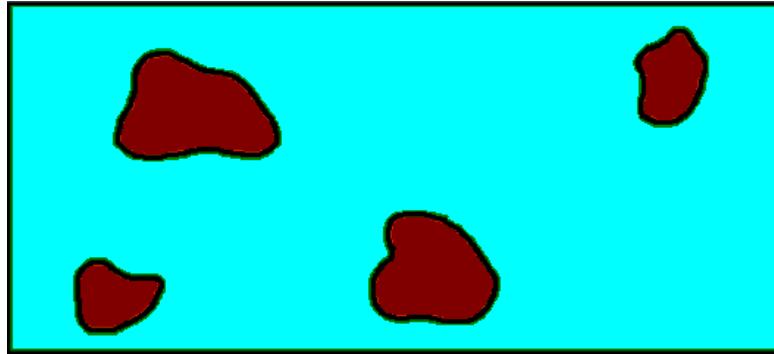
1. *Gradient based methods* can be further classified into (i) Line search methods and (ii) Trust-region methods. Line search algorithms usually select some descent direction (based on the gradient information) and minimize the function value along the chosen direction. This process is repeated until a local minimum is reached. The most popular line search algorithms include steepest descent [116], conjugate gradient [2] and Newton-Raphson method [110]. Trust region methods [24, 96], on the other hand, make an assumption about the nonlinear surface locally and do not use any form of line search. Typically, they assume the surface to be a simple quadratic model such that the minimum can be located directly if the model assumption is good which usually happens when the initial guess is close to the local minimum. If the model assumption is not accurate, then gradient information is used to guide the initial guess and after a certain period the model assumption is made again.
2. In contrast to the above mentioned gradient based methods, *non-gradient based methods* do not use any gradient information. Usually, these methods rely on a particular form of the nonlinear function and ensure that a chosen iterative parameter updating scheme results in the decrease of the function value [28]. For density estimation problems using maximum likelihood function, a popular class of iterative parameter estimation methods is the Expectation Maximization (EM) algorithm which converges to the maximum likelihood estimate of the mixture parameters locally [37, 123].

1.3 Motivation for this Thesis

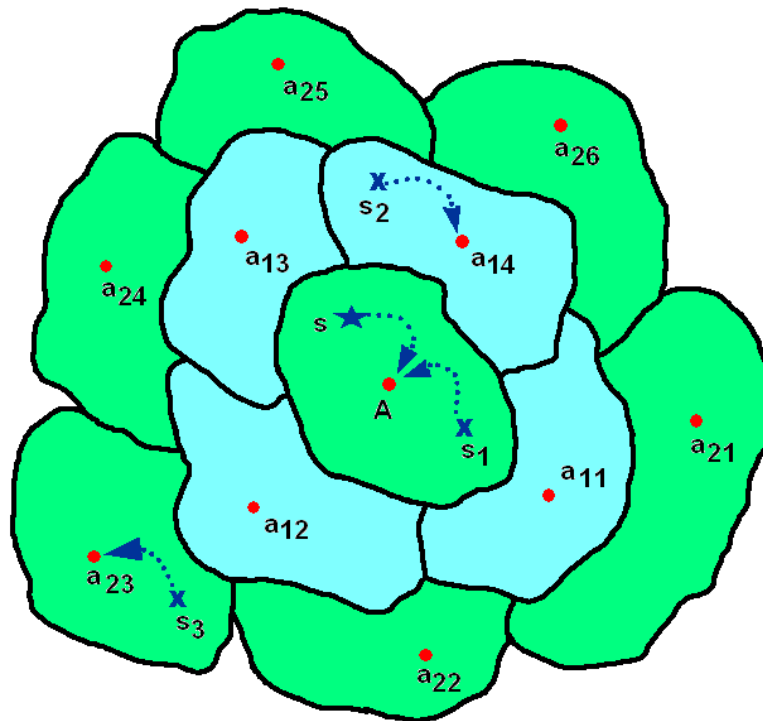
Finding the global optimal solution of a function is a lot more tedious and challenging for many problems. The task of finding such a solution is quite complex and increases rapidly with the dimensionality of the problem. Typically, the problem of finding the global optimal solution is solved in a hierarchical manner. Identifying some promising regions in a search space is relatively easier using certain global methods (such as genetic algorithms and simulated annealing) available in the literature. However, the fine tuning capability of these global methods is very poor and sometimes yield a comparatively less accurate solution even though the neighborhood region appears to be promising. Hence, there is an absolute necessity for exploring this surrounding to obtain a better solution.

Fig. 1.3 clearly shows the difficulties in dealing with nonlinear surfaces. Global methods can be used to obtain promising subspaces in the parameter space. These are indicated by dark shaded regions in Fig. 1.3(a). However, these promising regions are not convex in nature. i.e. they will have multiple local optimal solutions. Fig. 1.3(b) gives the top view of the nonlinear surface in the promising region. The dots indicate the local optimal solutions. ‘S’ is the initial point obtained from the global methods. Applying local method at ‘S’ will converge to ‘A’. There are other stochastic methods that can search the neighborhood regions e.g. mutations in genetic algorithms, low temperature annealing in simulated annealing.

We will now discuss the problems with these optimization methods mentioned above and motivate the necessity of TRUST-TECH (TRansformation Under STability-reTaining Equilibria CHaracterization) based methods. Finding a local optimum is relatively easier and straightforward using a local method. The stochastic methods randomly perturb a given point without much topological understanding of the



(a) Promising regions in the search space



(b) Promising Subspaces with multiple local optimal solutions

Figure 1.3: Various stages of TRUST-TECH (a) The dark regions indicate the promising subspaces. (b) The dots indicate the local optimal solutions and the dotted arrows indicate the convergence of local optimization method from a given initial point.

nonlinear surface. By transforming the nonlinear function into its corresponding dynamical system, TRUST-TECH can obtain neighboring local optimal solutions deterministically [29, 30]. TRUST-TECH not only guarantees that a new local maximum obtained is different from the original solution but also confirms that any solution in a particular direction will not be missed. As shown in Fig. 1.3(b), the given local optimal solution ‘A’ is randomly perturbed to obtain new initial points (s_1, s_2, s_3) . Applying local method using these initial points again, one can obtain the local optimal solutions A, a_{14} and a_{23} respectively. It can be observed that the solutions might appear again or might sometimes even miss some of the neighborhood solutions.

In this thesis, we develop TRUST-TECH based methods for systematically finding neighborhood solutions for problems that arise in the fields of optimization and machine learning. This method is more reliable and deterministic when compared to other stochastic approaches which merely use random moves to obtain new solutions. To begin with, the original nonlinear function is transformed into its corresponding dynamical system. There will be a one-to-one correspondence of all the critical points under this transformation. Also, this will allow us to define the concepts like stability boundaries which can be used to obtain the neighborhood solutions effectively.

1.4 Contributions of this Thesis

The main contributions of this thesis work are :

- Transform the problem of finding saddle points on potential energy surfaces to the problem of finding decomposition points on its corresponding nonlinear dynamical system. Apply a novel stability boundary based algorithm for

tracing the stability boundary and obtaining saddle points on potential energy surfaces that arise in the fields of computational chemistry and computational biology.

- Develop TRUST-TECH based Expectation Maximization (EM) algorithm for the learning finite mixture models.
- Apply the above mentioned TRUST-TECH based EM algorithm for the challenging motif finding problem in bioinformatics.
- Develop a component-wise kernel smoothing algorithm for learning Gaussian mixture models more efficiently. Demonstrate empirically that the number of unique local maxima on the likelihood surface can be reduced.
- Implement TRUST-TECH based training algorithm for artificial neural networks that can explore the topology of the error surface and obtain optimal set of parameters for the network model.
- Develop evolutionary TRUST-TECH methods that combines the advantages of the popular stochastic global optimization methods (like Genetic algorithms) with deterministic TRUST-TECH based search strategies. Demonstrate that the promising subspaces in the search space can be achieved at a much faster rate using evolutionary TRUST-TECH.

1.5 Organization of this Thesis

Fig. 1.4 shows the organization chart of this thesis. The main contributions are in the areas of nonlinear optimization relevant to machine learning. TRUST-TECH, smoothing methods and evolutionary algorithms are the optimization methods studied and developed in this thesis. In the context of machine learning, improvements

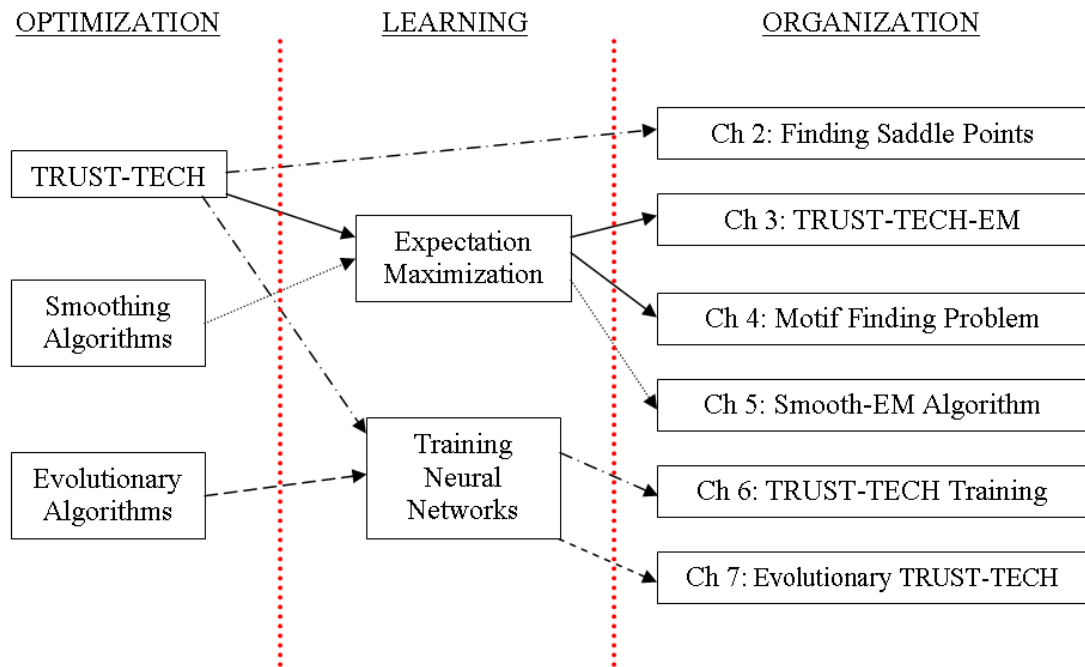


Figure 1.4: Organization chart of this thesis. The main contributions are in the areas of optimization and learning.

have been proposed for the two most widely studied algorithms, namely Expectation Maximization (in unsupervised learning) and training neural networks (in supervised learning). The remainder of this thesis is organized as follows: Chapter 2 describes a novel stability boundary based method to find saddle points on potential energy surfaces. A novel TRUST-TECH based Expectation Maximization algorithm for learning mixture models is proposed in Chapter 3 and this algorithm is applied to the motif finding problem in bioinformatics in Chapter 4. A Component-wise kernel smoothing algorithm for learning Gaussian mixture models is proposed in Chapter 5. Application of the TRUST-TECH method for efficient training of neural networks is discussed in Chapter 6. Chapter 7 proposes evolutionary TRUST-TECH model with some preliminary yet promising results. Finally, Chapter 8 concludes the discussion and proposes a few future research directions for the algorithms proposed in this thesis.

Chapter 2

Finding Saddle Points on Potential

Energy Surfaces

In this chapter, we will use the concepts of stability regions and stability boundaries to obtain saddle points on potential energy surfaces. The task of finding saddle points on potential energy surfaces plays a crucial role in understanding the dynamics of a micro-molecule as well as in studying the folding pathways of macromolecules like proteins. It is proposed that the problem of finding the saddle points on a high dimensional potential energy surface be transformed into the problem of finding dynamic decomposition points (DDP) of its corresponding nonlinear dynamical system. A novel stability boundary following procedure is used to trace the stability boundary to compute the DDP; hence the saddle points. The proposed method was successful in finding the saddle points on different potential energy surfaces of various dimensions. A simplified version of the algorithm has also been used to find the saddle points of symmetric systems with the help of some analytical knowledge. The main advantages and effectiveness of the method are clearly illustrated with some examples. Promising results of our method are shown on various problems with varied degrees of freedom.

2.1 Introduction

Recently, there has been a lot of interest across various disciplines to understand a wide variety of problems related to bioinformatics and computational biology. One of the most challenging problems in the field of computational biology is de-novo protein structure prediction where the structure of a protein is estimated from some

complex energy functions. Scientists have related the native structure of a protein structurally to the global minimum of the potential energy surface of its energy function [40]. If the global minimum could be found reliably from the primary amino acid sequence, it would provide us with new insights into the nature of protein folding. However, understanding the process of protein folding involves more than just predicting the folded structures of foldable sequences. The folding pathways in which the proteins attain their native structure can deliver some important information about the properties of the protein structure [101].

Proteins usually have multiple stable macrostates [50]. The conformations associated with one macrostate correspond to a certain biological function. Understanding the transition between these macrostates is important to comprehend the interactions of that protein with its environment and to understand the kinetics of the folding process, we need the structure of the transition state. Since, it is difficult to characterize these structures by manual experiments, simulations are an ideal tool for the characterization of the transition structures. Recently, biophysicists started exploring the computational methods that can be used to analyze conformational changes and identify possible reaction pathways [18]. In particular, the analysis of complex transitions in macromolecules has been widely studied [73].

From a computational viewpoint, transition state conformations correspond to saddle points. *Saddle points* are the points on a potential energy surface where the gradient is zero and where the Hessian of the potential energy function has only one negative eigenvalue [70]. Intuitively, this means that a saddle point is a maximum along one direction but a minimum along all other orthogonal directions. Fig. 1 shows a saddle point (x_d) located between two local minima (x_s^1 and x_s^2) and two local maxima (x_m^1 and x_m^2). As shown in the figure, the saddle point is a maximum

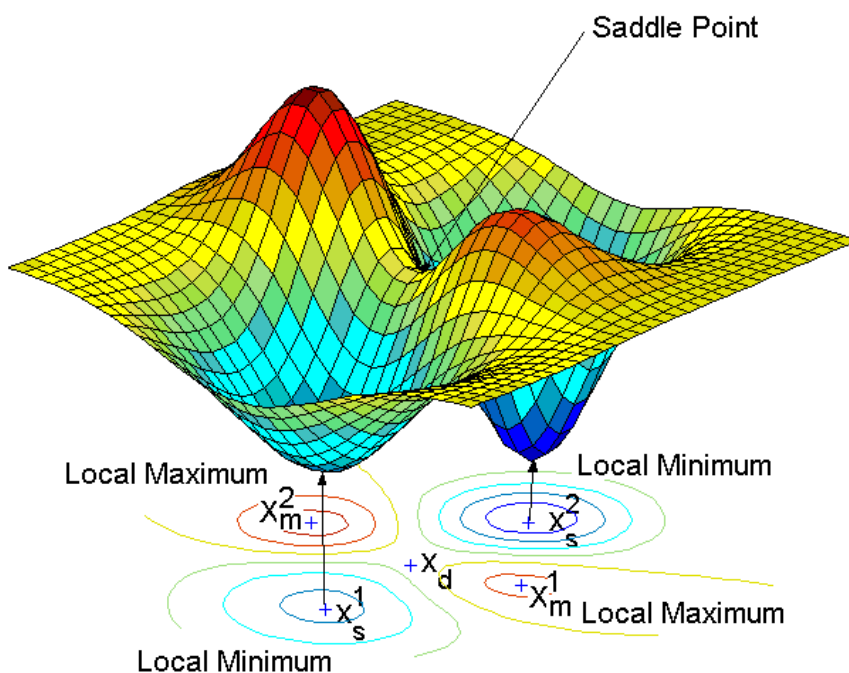


Figure 2.1: The surface and contour plots of a two-dimensional energy function. A saddle point (x_d) is located between two local minima (x_s^1 and x_s^2). x_m^1 and x_m^2 are two local maxima located in the orthogonal direction.

along the direction of the vector joining the two local minima and a minimum along its orthogonal direction (or the direction of the vector joining the two local maxima). The direction in which the saddle point is the maximum is usually unknown in most of the practical problems and is the direction of interest. This makes the problem of finding the saddle points more challenging than the problem of finding local minima on a potential energy surface. In terms of transition states, saddle points are local maxima with respect to the reaction coordinates for folding and local minima with respect to all other coordinates. The search for the optimal transition state becomes a search for the saddle points on the edge of the potential energy basin corresponding to the initial state. Finding these saddle points on potential energy surfaces can provide new insights about the folding mechanism of proteins. The primary focus of this chapter is to find the saddle points on different potential energy surfaces with varied degrees of freedom using TRUST-TECH based method.

2.2 Relevant Background

The task of finding saddle points has been a topic of active research in the field of computational chemistry for almost two decades. Recently, there has also been some interest in finding the saddle points of the Lennard-Jones clusters since it will give some idea about the dynamics of the system [43]. The properties of higher-index saddle points have been invoked in recent theories of the dynamics of supercooled liquids. Since the eigenvalues of the Hessian matrix can provide some information about the saddle points, several methods based on the idea of diagonalization of the Hessian matrix [86, 5, 71] were proposed in the literature. Some improved methods dealing with the updates of Hessian matrix have also been proposed [118]. Even though these methods appear to find saddle points accurately, they work mainly for low dimensional systems. These methods are not practical for higher dimensional problems because of the tremendous increase in the computational cost.

However, some methods that work without the necessity for computing the second derivatives have been developed. Because of the scalability issues, much more importance is given to algorithms that use only the first derivatives to compute the saddle points. A detailed description of the methods that work only based on first derivatives along with their advantages and disadvantages is given in a recent review paper [73]. The various methods that are used to find saddle points are drag method [73], dimer method [72], self penalty walk [35], activation relaxation technique [8], ridge method [81], conjugate peak refinement [56], DHS method [38], Nudged elastic band [83, 74], Step and slide [103]. Continuation methods for finding the saddle points are described in [89]. Almost all these methods except the dimer method are used to identify the saddle point between two given neighboring local minima. Though the dimer method successfully finds the saddle points in those cases where

only one minimum is given, it does not have a good control over which saddle point it intends to find.

All these methods start searching for saddle points from the local minimum itself and hence they need to compute the first derivative. However, our approach doesn't require the gradient information starting from the local minima. It will find the stability boundary in a given direction and then trace the stability boundary till the saddle point is reached [120]. This tracing of the stability boundary is more efficient than looking for saddle points in the entire search space. This work presents a completely novel *stability boundary* based approach to compute the saddle point between two given local minima. Our method is based on some of the fundamental results on stability regions of nonlinear dynamical systems [33, 32, 91].

2.3 Theoretical Background

Before presenting the details of the TRUST-TECH based methods, we review some fundamental concepts of nonlinear dynamical systems. The notations, definitions and theorems introduced in this section will hold for the rest of the thesis without any changes unless otherwise explicitly stated. Let us consider an unconstrained search problem on a nonlinear surface defined by the objective function

$$f(x) \tag{2.1}$$

where $f(x)$ is assumed to be in $C^2(\mathbb{R}^n, \mathbb{R})$.

Definition 1 \bar{x} is said to be a critical point of (2.1) if it satisfies the following condition

$$\nabla f(x) = 0 \tag{2.2}$$

A critical point is said to be *nondegenerate* if at the critical point $\bar{x} \in \mathfrak{R}^n$,

$$d^T \nabla_{xx}^2 f(\bar{x}) d \neq 0 \quad (\forall d \neq 0).$$

We construct the following *gradient system* in order to locate critical points of the objective function (2.1):

$$\frac{dx}{dt} = F(x) = -\nabla f(x) \quad (2.3)$$

where the state vector x belongs to the Euclidean space \mathfrak{R}^n , and the vector field $F : \mathfrak{R}^n \rightarrow \mathfrak{R}^n$ satisfies the sufficient condition for the existence and uniqueness of the solutions. The solution curve of Eq. (2.3) starting from x at time $t = 0$ is called a *trajectory* and it is denoted by $\Phi(x, \cdot) : \mathfrak{R} \rightarrow \mathfrak{R}^n$. A state vector x is called an *equilibrium point* of Eq. (2.3) if $F(x) = 0$.

Definition 2 *An equilibrium point is said to be hyperbolic if the Jacobian of F at point x has no eigenvalues with zero real part. A hyperbolic equilibrium point is called a (asymptotically) stable equilibrium point (SEP) if all the eigenvalues of its corresponding Jacobian have negative real part. Conversely, it is an unstable equilibrium point if some eigenvalues have a positive real part.*

An equilibrium point is called a *type- k equilibrium point* if its corresponding Jacobian has exact k eigenvalues with positive real part. When $k = 0$, the equilibrium point is (asymptotically) stable and it is called a *sink* (or *attractor*). If $k = n$, then the equilibrium point is called a *source* (or *repeller*).

A dynamical system is completely *stable* if every trajectory of the system leads to one of its stable equilibrium points. The *stable* ($W^s(\tilde{x})$) and *unstable* ($W^u(\tilde{x})$) manifolds of an equilibrium point, say \tilde{x} , is defined as:

$$W^s(\tilde{x}) = \{x \in \mathfrak{R}^n : \lim_{t \rightarrow \infty} \Phi(x, t) = \tilde{x}\} \quad (2.4)$$

$$W^u(\tilde{x}) = \{x \in \mathfrak{R}^n : \lim_{t \rightarrow -\infty} \Phi(x, t) = \tilde{x}\} \quad (2.5)$$

The *stability region* (also called *region of attraction*) of a stable equilibrium point x_s of a dynamical system (2.3) is denoted by $A(x_s)$ and is

$$A(x_s) = \{x \in \mathfrak{R}^n : \lim_{t \rightarrow \infty} \Phi(x, t) = x_s\} \quad (2.6)$$

The boundary of stability region is called the *stability boundary* of x_s and will be denoted by $\partial A(x_s)$. It has been shown that the stability region is an open, invariant and connected set [33]. From the topological viewpoint, the stability boundary is a $(n - 1)$ dimensional closed and invariant set. A new concept related to the stability regions namely the *quasi-stability region* (or *practical stability region*), was developed in [32].

The *practical stability region* of a stable equilibrium point x_s of a nonlinear dynamical system (2.3), denoted by $A_p(x_s)$ and is

$$A_p(x_s) = \text{int } \overline{A(x_s)} \quad (2.7)$$

where \bar{A} denotes the closure of A and $\text{int } \bar{A}$ denotes the interior of \bar{A} . $\text{int } \bar{A}(x_s)$ is an open set. The boundary of practical stability region is called the *practical stability boundary* of x_s and will be denoted by $\partial A_p(x_s)$.

It has been shown that the practical stability boundary $\partial A_p(x_s)$ is equal to $\partial \bar{A}(x_s)$ [33]. The practical stability boundary is a subset of its stability boundary.

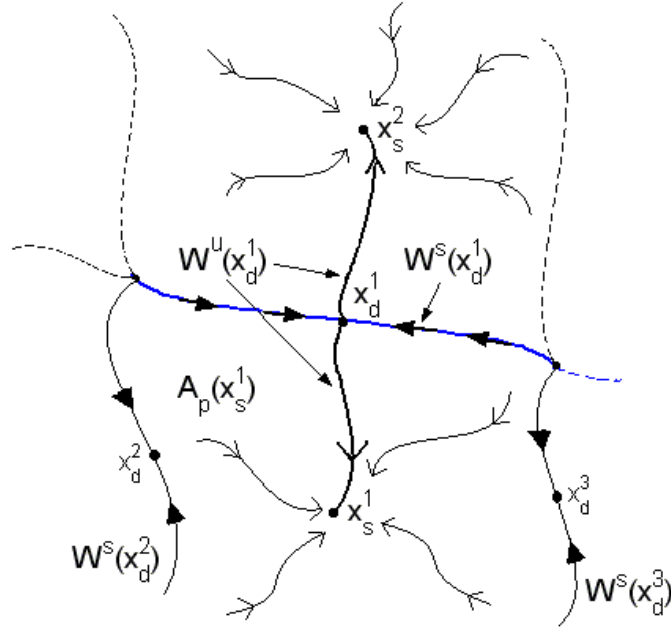


Figure 2.2: Phase potrait of a gradient system. The solid lines with solid arrows represent the basin boundary. $\partial A_p(x_s^1) = \bigcup_{i=1}^3 \overline{W^s(x_d^i)}$. The local minima x_s^1 and x_s^2 correspond to the stable equilibrium points of the gradient system. The saddle point (x_d^1) corresponds to the dynamic decomposition point that connects the two stable equilibrium points.

It eliminates the complex portion of the stability boundary which has no “contact” with the complement of the closure of the stability region. A complete characterization of the practical stability boundary for a large class of nonlinear dynamical systems can be found.

Definition 3 A type-1 equilibrium point x_d ($k=1$) on the practical stability boundary of a stable equilibrium point x_s is called a dynamic decomposition point.

To comprehend the transformation, we need to define *Lyapunov function*. A smooth function $V(\cdot) : \mathfrak{R}^n \rightarrow \mathfrak{R}$ satisfying $\dot{V}(\Phi(x, t)) < 0, \forall x \notin \{\text{set of equilibrium points (E)}\}$ and $t \in \mathfrak{R}^+$ is termed as Lyapunov function.

Theorem 2.3.1 [31]: $F(x)$ is a Lyapunov function for the negative quasi-gradient system (2.3).

Theorem 2.3.2 (*Characterization of stability boundary*)[32]: Consider a nonlinear dynamical system described by (2.3). Let $\sigma_i, i=1,2,\dots$ be the equilibrium points on the stability boundary $\partial A(x_s)$ of a stable equilibrium point, say x_s . Then

$$\partial A(x_s) \subseteq \bigcup_{\sigma_i \in \partial A} W^s(\sigma_i). \quad (2.8)$$

Theorem 2.3.2 completely characterizes the stability boundary for nonlinear dynamical systems by asserting that the stability boundary is the union of the stable manifolds of all critical elements on the stability boundary. This theorem gives an explicit description of the geometrical and dynamical structure of the stability boundary. This theorem can be extended to the characterization of the practical stability boundary in terms of the stable manifold of the dynamic decomposition point.

Theorem 2.3.3 (*Characterization of practical stability boundary*)[32]: Consider a nonlinear dynamical system described by (2.3). Let $\sigma_i, i=1,2,\dots$ be the dynamic decomposition points on the practical stability boundary $\partial A_p(x_s)$ of a stable equilibrium point, say x_s . Then

$$\partial A_p(x_s) \subseteq \bigcup_{\sigma_i \in \partial A_p} \overline{W^s(\sigma_i)}. \quad (2.9)$$

Theorem 2.3.3 asserts that the practical stability boundary is contained in the union of the closure of the stable manifolds of all the dynamic decomposition points on the practical stability boundary. Hence, if the dynamic decomposition points can be identified, then an explicit characterization of the practical stability boundary can be established using (2.9).

Theorem 2.3.4 (*Unstable manifold of type-1 equilibrium point*)[89]: Let x_s^1 be a stable equilibrium point of the gradient system (2.3) and x_d be a type-1 equilibrium point on the practical stability boundary $\partial A_p(x_s)$. Assume that there exist ϵ and δ such that $\|\nabla f(x)\| > \epsilon$ unless $x \in B_\delta(\hat{x}), \hat{x} \in \{x : \nabla f(x) = 0\}$. There exists another stable equilibrium point x_s^2 to which the one dimensional unstable manifold of x_d converges. Conversely, if $\overline{A_p(x_s^1)} \cap \overline{A_p(x_s^2)} \neq \emptyset$, then there exists a dynamic decomposition point x_d on $\partial A_p(x_s^1)$.

Theorem 2.3.4 is imperative to understand some of the underlying concepts behind the development of TRUST-TECH. It associates the notion of stable equilibrium points, practical stability regions ($A_p(x_s)$), practical stability boundaries ($\partial A_p(x_s)$) and type-1 equilibrium points. As shown in fig. 2.2, The unstable manifold (W^u) of the dynamic decomposition point x_d^1 converges to the two stable equilibrium points x_s^1 and x_s^2 . Also, it should be noted that x_d^1 is present on the stability boundary of x_s^1 and x_s^2 .

We also need to show that under the transformation from (2.1) to (2.3), the properties of the critical points remain unchanged. Theorem 2.3.5 illustrates the correspondence of the critical points of the original system.

Theorem 2.3.5 (*Critical Points and their correspondence*)[31]: An equilibrium point of (2.3) is hyperbolic if, and only if, the corresponding critical point is nondegenerate. Moreover, if \bar{x} is a hyperbolic equilibrium point of (2.3), then

1. \bar{x} is a stable equilibrium point of (2.3) if and only if \bar{x} is an isolated local minimum for (2.1)
2. \bar{x} is a source of (2.3) if and only if \bar{x} is an isolated local maximum for (2.1)
3. \bar{x} is a dynamic decomposition point of (2.3) if and only if \bar{x} is a saddle point for (2.1)

2.4 A Stability Boundary based Method

Our TRUST-TECH based boundary tracing method uses the theoretical concepts of dynamical systems presented in the previous section. The method described in this section finds the DDP when the two neighborhood local minima are given. Our method is illustrated on a two-dimensional LEPS potential energy surface [115]. The equations corresponding to the LEPS potential are given in the appendix-A. The two local minima are A and B and the dynamic decomposition point is DDP .

Given : Two neighborhood local minima (A , B)

Goal : To obtain the corresponding DDP

Algorithm :

Step1: *Initializing the Search direction* : Since the location of the neighborhood local minima is already given, the initial search direction becomes explicit. The vector that joins the two given local minima (A and B) is chosen to be the initial search direction.

Step 2: *Locating the exit point (X_{ex})* : (see fig. 2.3) Along the direction AB , starting from A , the function value is evaluated at different step intervals. Since the vector is between two given local minima, the function value will monotonically increase and then decrease till it reaches the other local minimum (B). The point where the energy value attains its peak is called the *exit point*.

Step 3: *Moving along the stability boundary to locate the Minimum Gradient Point* : We used a novel *stability boundary following procedure* to move along the practical stability boundary. Once the exit point is identified, the consecutive points on the stability boundary can be identified by this stability boundary following procedure. The exit point (X_{ex}) is integrated for a predefined number of times. Let

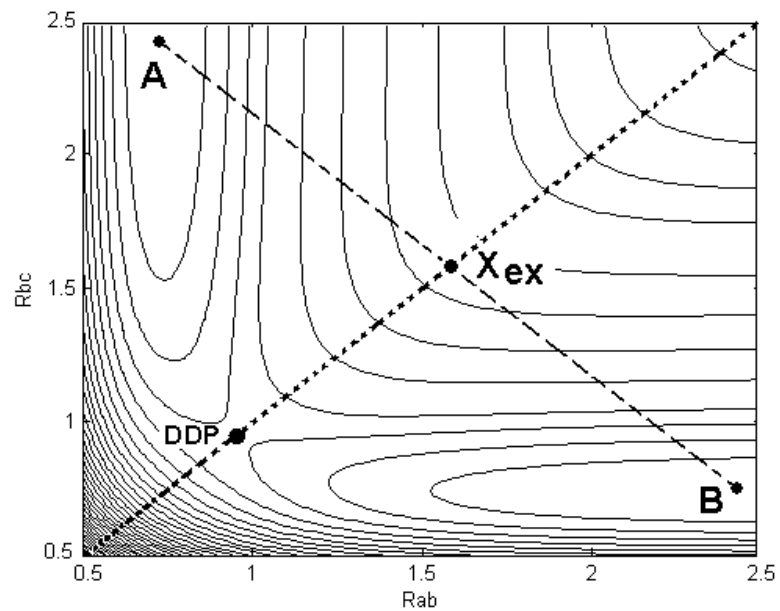


Figure 2.3: Contour plot of a 2-D LEPS potential (described in appendix-A). Each line represents the values of a constant potential. A and B are the two local minima. DDP is the dynamic decomposition point to be computed. The search direction is the direction of the vector joining AB . The exit point (X_{ex}) is obtained by finding the peak of the function value along this vector. The dotted line indicates the stability boundary. The dashed line indicates the search direction.

m'_1 be the new point obtained after integration. The function value between m'_1 and the local minimum is evaluated and the peak value is obtained. Let the new boundary point along the vector $m'_1 B$ starting from the point m'_1 and where the value attains the peak be m_2 . This process is repeated and several points on the stability boundary are obtained. During this traversal, the value of the gradient along the boundary points is noted and the process of moving along the boundary is terminated when the minimum gradient point (MGP) is obtained. In summary, the trajectory of integration is being modified so that it moves towards the MGP and will not converge to one of the local minima. This is an intelligent TRUST-TECH based scheme for following the stability boundary which is the *heart* of the proposed method. This step is named as the stability boundary following procedure.

Step 4: *Locating the Dynamic Decomposition point (DDP)* : The Minimum Gradient Point (m_n) obtained from the previous step will be located in the neighborhood of the dynamic decomposition point. A local minimizer to solve the system of nonlinear equations is applied with m_n as initial guess and this will yield the DDP. A detailed survey about different Local minimizations applied to a wide variety of areas is given in [130].

Remarks:

- The step size to be chosen during the step 2 of our algorithm is very critical for faster computation and accuracy of the exit points.
- The number of integrations to be performed from a point on the stability boundary (m_k) till the new point (m'_k) is reached is problem specific and it depends on the behaviour of the stability boundary near that point.
- The minimum gradient point is usually in the neighborhood of the DDP. Newton's method is a powerful local solver that can be used to obtain the DDP.

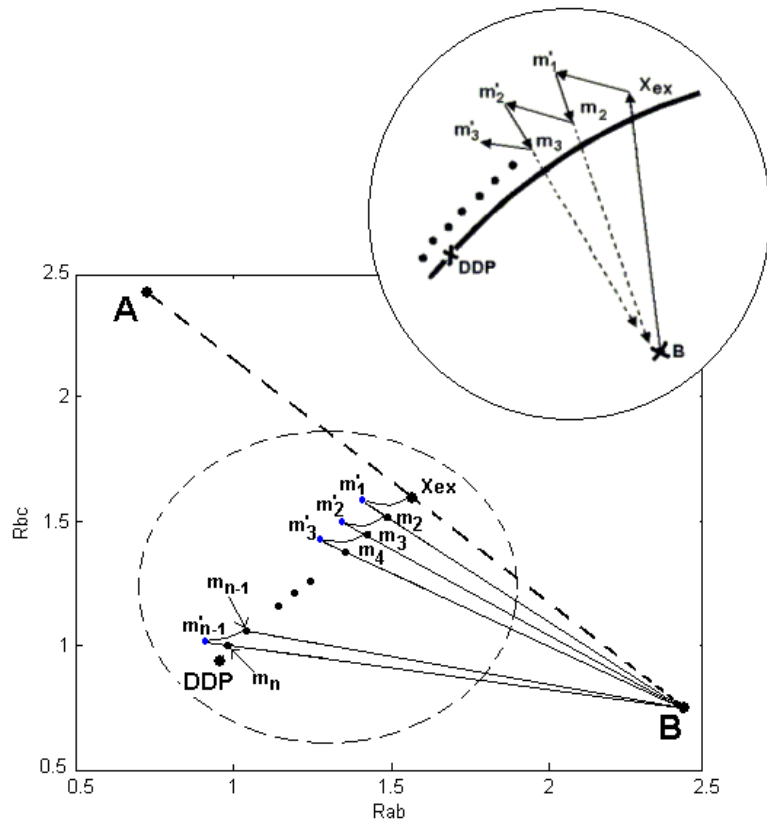


Figure 2.4: Illustration of the step 3 of our algorithm. m_1 is integrated till m'_1 is reached and traced back along the vector $m'_1 B$ to get m_2 , and so on. m_n are the n exit points on the stability boundary. The n^{th} point is the Minimum Gradient Point where the magnitude of the gradient (G_M) is the minimum amongst all the computed gradient values along the stability boundary ($X_{ex}, m_2, m_3, \dots, m_n$). DDP is the dynamic decomposition point.

2.5 Implementation Issues

For our illustration, we consider a N dimensional function F with variables X_i , where $i = 1 \dots N$. From the algorithmic viewpoint, our method consists of three stages: (i) Finding the exit point, (ii) Following the stability boundary and (iii) Computing the DDP. Let A and B are the given local minima, the pseudocode for finding the DDP is as follows:

```
point procedure locate_DDP( $A, B$ )
```

```
Initialize stepsize = 10 // initial evaluation step size
```

```
 $EP \leftarrow$  find_ExitPt( $A, B, stepsize$ ) // Exit point
```

```
 $MGP \leftarrow$  Boundary_Following( $EP$ ) // Trace the stability boundary
```

```
 $DDP \leftarrow$  local_minimizer( $MGP$ ) // Compute the DDP
```

```
return  $DDP$ 
```

The procedure *find_ExitPt* will find the point on the stability boundary between two given points A and B starting from A (*step 2* of our algorithm). If the function value is monotonically decreasing from the first step, then it indicates that there will not be any boundary in that direction, and the search is changed to a new direction. The new direction can be obtained by making $B = 2A - B$. Finding the exit point can be done more efficiently by first evaluating the function at comparatively large step intervals (See Fig. 2.5). The function evaluation is started from a_1 , a_2 and so on. Once a_6 is reached, the energy value starts to reduce indicating that the peak value has been reached. Golden section search algorithm (see appendix-B) is used to efficiently find the exit point within a small interval range where the peak value is present. The golden section search is applied to obtain the exit point X_{ex} within the intervals a_4 and a_6 . In fact, golden section search could have been used from the

two given local minima. We prefer to evaluate the function value at certain intervals because using this method the stability boundary can be identified without knowing the other local minimum as well.

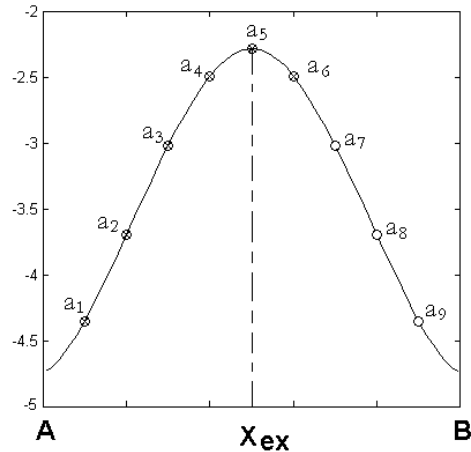


Figure 2.5: The plot of the function value along the vector AB . The curve monotonically increases starting from A and then decreases till it reaches B . X_{ex} is the exit point where the function attains its peak value. a_1, a_2, \dots, a_9 are the interval points. The marked circles indicate that the function has been evaluated at these points. The empty circles indicate the points where the function is yet to be evaluated.

point procedure $find_ExitPt(A, B, steps)$

- 1: Initialize eps // accuracy
- 2: $interval = (B - A)/steps$
- 3: $cur = eval(A)$
- 4: $tmp = A + interval$
- 5: **if** $eval(tmp) < cur$ **then**
- 6: $B = 2 * A - B$
- 7: **end if**
- 8: **for** $i = 1$ to $steps$ **do**
- 9: $tmp = A + i * interval$
- 10: $prev = cur$
- 11: $cur = eval(tmp)$

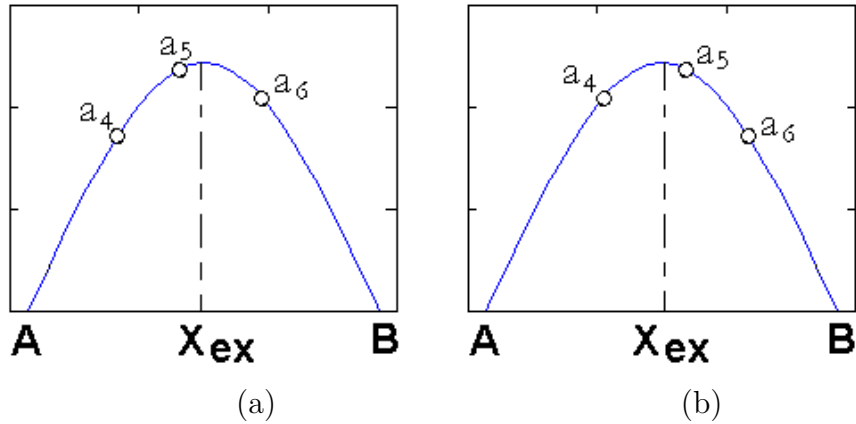


Figure 2.6: The interval point a_5 can be present either to the left or to the right of the peak. When a_6 is reached, the golden section search method is invoked with a_6 and a_4 as the interval.

```

12:  if  $prev > cur$  then
13:       $newPt = A + (i - 2) * interval$ 
14:       $BdPt = GoldenSectionSearch(tmp, newPt, eps)$ 
15:      return  $BdPt$ 
16:  end if
17: end for
18: return  $NULL$ 

```

Fig. 2.6 shows the other two possibilities of having the interval points. It must be noted that the golden section search procedure should not be invoked with two consecutive intervals between which the value starts to reduce. It should be applied to two intervals a_4 and a_6 because the point a_5 can be present on either side of the peak. From Fig. 2.6b, we can see that the value at a_5 is greater than a_4 and less than a_6 but still the peak is not present between a_5 and a_6 . This is the reason why the golden section search method should have a_4 and a_6 as its arguments.

In an ideal case, integration from the exit point will lead to the DDP. However, due to the numerical approximations made, the integration starting from the

exit point will eventually converge to its corresponding local minimum. Hence, we need to adjust the path of integration so that it can effectively trace the stability boundary.

```

point procedure Boundary_Following(ExPt)

1: Initialize dt // Integral step size
2: Initialize intsteps // No. of integration steps
3: Initialize smallstep //small step size
4: reduce_flag = OFF // To check if gradient reduced
5: Store BdPt = ExPt
6: GM_next = GM(BdPt)
7: while (1) do
8:   integrate BdPt for intsteps number of times and with dt step size to obtain
   NewPt
9:   prevPt = BdPt
10:  Obtain another exit point BdPt = find_ExitPt(NewPt, A, smallstep)
11:  GM_prev = GM_next
12:  Magnitude of the gradient for the new exit point GM_next = GM(BdPt)
13:  if GM_next < GM_prev then
14:    reduce_flag = ON // compare the magnitude of the gradient
15:  end if
16:  if GM_next > GM_prev & reduce_flag = ON then
17:    MGP has been reached MGPt = prevPt
18:  end if
19:  return MGPt
20: end while

```

The procedure *Boundary_Following* takes in the exit point as the argument

and returns the computed MGP obtained by tracing the stability boundary starting from the exit point (*step 3* of our algorithm). This function implements the stability boundary following procedure in our algorithm and is responsible for carefully tracing the practical stability boundary. The point on the stability boundary is integrated for a predefined number of times. From the exit point, it is integrated using the equation shown below.

$$X_{(i+1)} = X_{(i)} - \left(\frac{\partial F}{\partial X_{(i)}} \right) \cdot \Delta t \quad (2.10)$$

The exit point (X_{ex}) is integrated for a predefined number of times. Let m'_1 be the new point obtained after integration. The function value between m'_1 and the local minimum is evaluated and the peak value is obtained. Let the new boundary point along the vector m'_1B starting from the point m'_1 and where the value attains the peak be m_2 . This process is repeated and several points on the stability boundary are obtained. During this traversal, the value of the gradient along the boundary is noted and the process of moving along the boundary is terminated when the minimum gradient point (MGP) is obtained. In summary, the trajectory of integration is being modified so that it moves towards the MGP and will not converge to one of the local minima.

The procedure *integrate* computes a point (*newPt*) by integrating a certain number of integral steps (*intstep*) with a predefined integration step size (Δt) from the boundary point (*bdPt*). Once again a new exit point on the practical stability boundary is obtained from *newPt* using the procedure *find_ExitPt*. This process of integrating and tracing back is repeated for a certain number of steps. Another important issue is the stopping criterion for this tracing procedure. For this, we will have to compute the magnitude of the gradient at all points obtained on the stability

boundary. The magnitude of the gradient (G_M) is calculated using Eq. (2.11).

$$G_M = \sqrt{\sum_{i=1}^N \left(\frac{\partial F}{\partial X_i} \right)^2} \quad (2.11)$$

where Δt is the integral step size. The G_M value can either start increasing and then reduce or it might start reducing from the exit point. The *reduce_flag* indicates that the G_M value started to reduce before. GM_next and GM_prev are the two variables used to store the values of the current and previous G_M values respectively. The *MGP* is obtained when $GM_next > GM_prev$ and *reduce_flag* = *ON*.

2.6 Experimental Results

2.6.1 Test Case 1: Two-dimensional Potential Energy Surface

The *Muller-Brown surface* is a standard two-dimensional example of a potential energy function in theoretical chemistry [105]. This surface was designed for testing the algorithms that find saddle points. Eq. (2.12) gives the Muller-Brown energy function. Fig. 2.7 represents the two-dimensional contour plot of the potential energy surface of the muller-brown function.

$$C(x, y) = \sum_{i=1}^4 A_i \exp \left[a_i(x - x_i^o)^2 + b_i(x - x_i^o)(y - y_i^o) + c_i(y - y_i^o)^2 \right]. \quad (2.12)$$

where

$$\begin{aligned} A &= (-200.0, -100.0, -170.0, -15.0) & \mathbf{a} &= (-1.0, -1.0, -6.5, -0.7) \\ x^o &= (1.0, 0.0, -0.5, -1.0) & \mathbf{b} &= (0.0, 0.0, 11.0, 0.6) \\ y^o &= (0.0, 0.5, 1.5, 1.0) & \mathbf{c} &= (-10.0, -10.0, -6.5, 0.7) \end{aligned}$$

As shown in Fig. 2.7, there are three stable equilibrium points (A,B and C) and two dynamic decomposition points (DDP1,DDP2) on the muller-brown potential energy surface. DDP 1 is present between A and B and is more challenging to find, compared to DDP 2 which is present between B and C. Table 2.2 shows the exact locations and energy values of the local minima and the dynamic decomposition points.

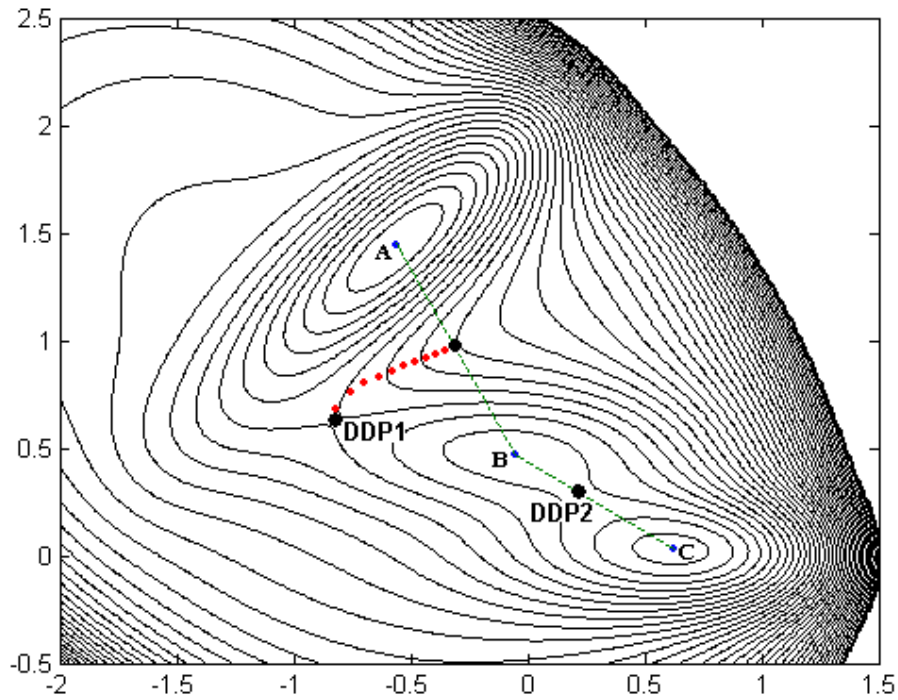


Figure 2.7: Two dimensional contour plot of the potential energy surface corresponding to the muller-brown function described in Eq. (2.12). A,B and C are stable equilibrium points and DDP1,DDP2 are two dynamic decomposition points. The dashed lines indicate the initial search direction. The dots indicate the results of the stability boundary following procedure. These points are the points on the stability boundary that reach the MGP.

Table 2.1: Stable equilibrium points and DDPs for the Muller-Brown Surface described in Eq. (2.12).

| Equilibrium Points | Location | Energy value $c(\cdot)$ |
|--------------------|----------------|-------------------------|
| SEP A | (-0.558,1.442) | -146.7 |
| DDP 1 | (-0.822,0.624) | -40.67 |
| SEP B | (-0.05,0.467) | -80.77 |
| DDP 2 | (0.212,0.293) | -72.25 |
| SEP C | (0.623,0.028) | -108.7 |

We construct the dynamical system corresponding to (2.12) as follows:

$$\begin{bmatrix} \dot{x}(t) \\ \dot{y}(t) \end{bmatrix} = - \begin{bmatrix} \frac{\partial C}{\partial x} \\ \frac{\partial C}{\partial y} \end{bmatrix}$$

$$\frac{\partial C}{\partial x} = \sum_{i=1}^4 A_i \cdot P \cdot [2a_i(x - x_i^o) + b_i(y - y_i^o)]$$

$$\frac{\partial C}{\partial y} = \sum_{i=1}^4 A_i \cdot P \cdot [b_i(x - x_i^o) + 2c_i(y - y_i^o)]$$

where

$$P = \exp [a_i(x - x_i^o)^2 + b_i(x - x_i^o)(y - y_i^o) + c_i(y - y_i^o)^2]$$

The exit point obtained between the local minima A and B is (-0.313, 0.971). Fig. 2.7 shows the results of our algorithm on Muller-Brown surface. The dashed lines indicate the initial search vector which is used to compute the exit point. The dots indicate the points along the stability boundary obtained during the stability boundary following procedure. These dots move from the initial exit point towards the MGP. The gradient curve corresponding to the points along this stability boundary is shown in Fig. 2.8. From the MGP, the local minimizer is applied to obtain the dynamic decomposition point (DDP 1). The exit point obtained between the local minima B and C is (0.218, 0.292). It converges to the dynamic decomposition point

(DDP 2) directly when the local minimizer is applied. Hence, given the three local minima, we are able to find the two saddle points present between them using our method.

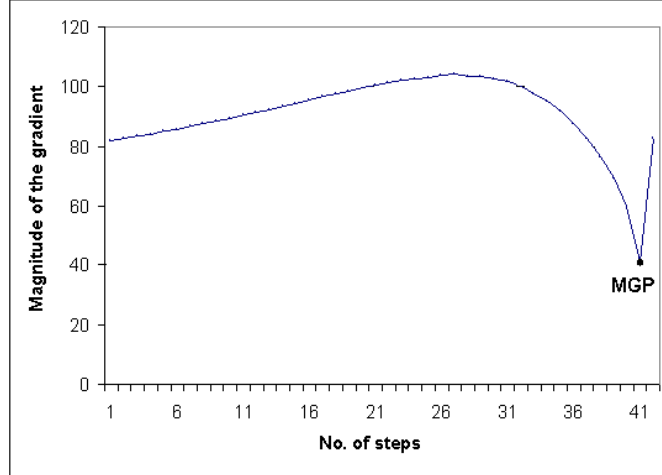


Figure 2.8: The gradient curve corresponding to the various points obtained from the stability boundary following procedure. The graph shows that the magnitude of the gradient slowly increases in the initial phases and then starts to reduce. The highlighted point corresponds to the gradient at the MGP.

Table 2.2: Stable equilibrium points and dynamic decomposition points for the Muller-Brown Surface described in Eq. (2.12).

| Equilibrium Points | Location | Energy value $c(\cdot)$ |
|--------------------|----------------|-------------------------|
| SEP A | (-0.558,1.442) | -146.7 |
| DDP 1 | (-0.822,0.624) | -40.67 |
| SEP B | (-0.05,0.467) | -80.77 |
| DDP 2 | (0.212,0.293) | -72.25 |
| SEP C | (0.623,0.028) | -108.7 |

2.6.2 Test Case 2: Three-dimensional symmetric systems

Three atom Lennard Jones clusters : This system is mainly used to demonstrate a simplified version of our algorithm. Our method has some advantages when applied to energy surfaces that are symmetric in nature. To demonstrate this, we used

the Lennard Jones pair potential which is a simple and commonly used model for interaction between atoms. The Lennard Jones potential is given by the Eq. (2.13). For simplicity, we applied reduced units, i.e. the values of ϵ and r_0 are taken to be unity. Plot of Lennard-Jones Potential of interaction between two atoms generated using Eq. (2.13) is shown in Fig. 2.9. The original problem is to find the global minimum of the potential energy surface obtained from the interaction between N atoms with two-body central forces. In this example, we consider the potential energy surface corresponding to the three-atom cluster which exhibits symmetric behaviour along the x-axis.

$$V = \sum_{i=1}^{N-1} \sum_{j=i+1}^N v(r_{ij})$$

$$v(r_{ij}) = \epsilon \left[\left(\frac{r_0}{r_{ij}} \right)^{12} - 2 \left(\frac{r_0}{r_{ij}} \right)^6 \right] \quad (2.13)$$

where

$$r_{ij} = \sqrt{(x_i - x_j)^2 + (y_i - y_j)^2 + (z_i - z_j)^2}$$

The total potential energy (V) of the microcluster is the summation of all two-body interaction terms, $v(r_{ij})$ is the potential energy term corresponding to the interaction of atom i with atom j , and r_{ij} is the Euclidean distance between i and j . ϵ describes the strength of the interaction and r_0 is the distance at which the potential is zero.

For a three atom cluster, let the coordinates be (x_1, y_1, z_1) , (x_2, y_2, z_2) and (x_3, y_3, z_3) . Though, there are nine variables in this system, due to the translational and rotational variants, the effective dimension is reduced to three. This reduction can be done by setting the other six variables to zero. Hence, the effective variables are (x_2, x_3, y_3) .

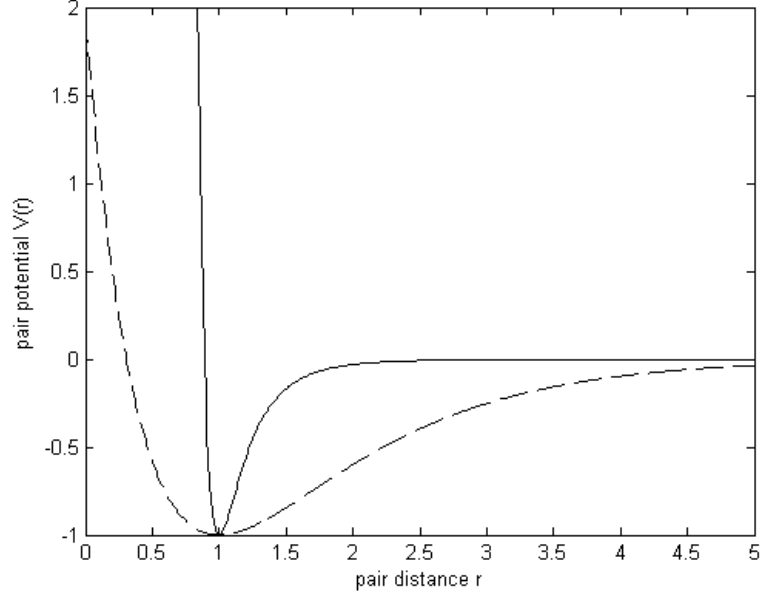


Figure 2.9: Characteristic curves of Lennard-Jones potential and the Morse potential with all parameters set to unity. Solid line represents Lennard-Jones potential (2.13) and the dashed line represents the Morse potential (2.15).

We construct the dynamical system corresponding to (2.13) as follows:

$$\begin{bmatrix} \dot{x}_2(t) \\ \dot{x}_3(t) \\ \dot{y}_3(t) \end{bmatrix} = - \begin{bmatrix} \frac{\partial V}{\partial x_2} \\ \frac{\partial V}{\partial x_3} \\ \frac{\partial V}{\partial y_3} \end{bmatrix} \quad (2.14)$$

where

$$\frac{\partial V}{\partial x_i} = \sum_{\substack{j=1 \\ j \neq i}}^3 \frac{12}{(r_{ij})^8} \left[1 - \left(\frac{1}{r_{ij}} \right)^6 \right] \cdot (x_i - x_j)$$

and

$$\frac{\partial V}{\partial y_i} = \sum_{\substack{j=1 \\ j \neq i}}^3 \frac{12}{(r_{ij})^8} \left[1 - \left(\frac{1}{r_{ij}} \right)^6 \right] \cdot (y_i - y_j)$$

Based on the two given local minima, one can compute the exit point analytically (not numerically) since the system is symmetric. Since the exit point is an exact (not an approximate) value, one can eventually reach the dynamic decomposition point by integrating the exit point. The exit point is $(0.0,0.0,0.0)$, $(2.0,0.0,0.0)$, $(1.0,0.0,0.0)$.

Table 2.3: Stable equilibrium points and dynamic decomposition points for the three atom Lennard Jones Cluster described in Eq. (2.13).

| Equilibrium Point | Location | Energy value $c(\cdot)$ |
|-------------------|--------------------|-------------------------|
| SEP A | $(1.0,0.5,0.866)$ | -3.000 |
| DDP 1 | $(2.0,1.0,0.0)$ | -2.031 |
| SEP B | $(1.0,0.5,-0.866)$ | -3.000 |

In this case, the stability boundary following procedure is not needed. Integrating from the exit point will eventually find the dynamic decomposition point. The last two steps of our method, stability boundary following procedure and the local minimizer are not required to obtain the dynamic decomposition point. It is clear from the example above that in all cases where we compute the exit point numerically, we get an approximate of the exit point which will eventually converge to one of the two local minima after integration. In such cases, the stability boundary following procedure will guide us to maintain the path along the stability boundary and prevent us from being trapped in one of the local minima. Hence, a simplified version of our algorithm is developed for finding saddle points on symmetric surfaces.

2.6.3 Test Case 3: Higher Dimensional Systems

Heptamer island on a crystal: To test our method on a higher dimensional system, we have chosen a heptamer island on the surface of an Face-Centered Cubic (FCC)

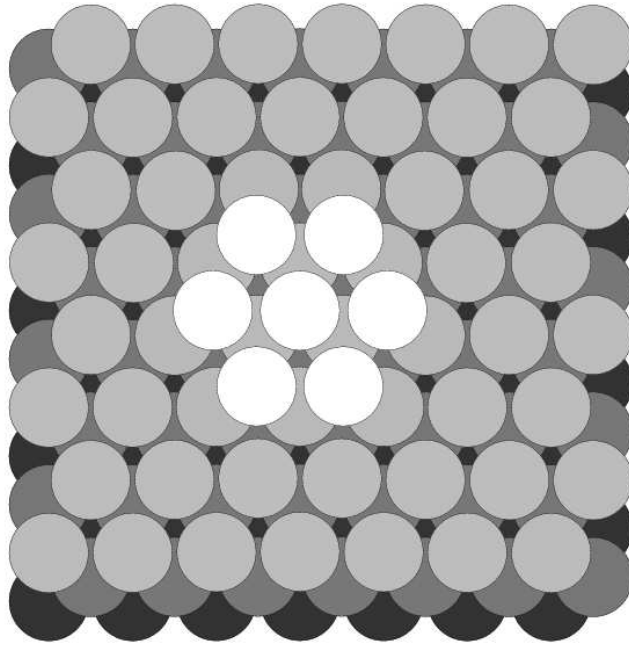


Figure 2.10: Top view of the surface and the seven atom island on the surface of an FCC crystal. The shading indicates the height of the atoms.

crystal. This system will not only illustrate the atomic scale mechanism of island diffusion on surfaces but also will help us to understand the kinetics of a process. The atoms interact via a pairwise additive Morse potential described by Eq. (2.15).

$$V(r) = A \left(e^{-2\alpha(r-r_0)} - 2e^{-\alpha(r-r_0)} \right) \quad (2.15)$$

where $A = 0.71$ eV, $\alpha = 1.61 \text{ \AA}^{-1}$, $r_0 = 2.9 \text{ \AA}$. These parameters were chosen in such a way that it will reproduce diffusion barriers on real surfaces. The potential was cut and shifted at 9.5 \AA . The surface is simulated with a 6 layer slab, each layer containing 56 atoms. The minimum energy lattice constant for the FCC solid is 2.74 \AA . The bottom three layers in the slab are held fixed. A total of $7 + 168 = 175$ atoms are allowed to move during the search for dynamic decomposition points. Hence, this is an example of 525 (175×3) dimensional search problem. This is the same system used in the review paper by Henkelman et al [73]. Fig. 2.10 shows the

top view of the initial configuration of the island with a compact heptamer sitting on top of the surface. The shading indicates the height of the atoms. The white ones are the heptamer island on the surface of the crystal and the black ones are at the bottom most part of the crystal. Fig. 2.11 shows some sample configurations corresponding to saddle points and local minima on the potential energy surface.

Let N be the number of atoms that can move (175) and N' be the total number of atoms (343).

$$V(r) = \sum_{i=1}^N \sum_{j=1}^{N'} A \left(e^{-2\alpha(r_{ij}-r_0)} - 2e^{-\alpha(r_{ij}-r_0)} \right) \quad (2.16)$$

where r_{ij} is the Euclidean distance between i and j . The *Dynamical System* is a $3N$ column matrix given by :

$$\begin{aligned} & [\dot{x}_1(t) \ \dot{x}_2(t) \ \dots \ \dot{x}_n(t) \ \dot{y}_1(t) \ \dot{y}_2(t) \ \dots \ \dot{y}_n(t) \ \dot{z}_1(t) \ \dot{z}_2(t) \ \dots \ \dot{z}_n(t)]^T \\ & = - \left[\begin{array}{cccccccccccc} \frac{\partial V}{\partial x_1} & \frac{\partial V}{\partial x_2} & \dots & \frac{\partial V}{\partial x_n} & \frac{\partial V}{\partial y_1} & \frac{\partial V}{\partial y_2} & \dots & \frac{\partial V}{\partial y_n} & \frac{\partial V}{\partial z_1} & \frac{\partial V}{\partial z_2} & \dots & \frac{\partial V}{\partial z_n} \end{array} \right]^T \end{aligned}$$

where

$$\frac{\partial V}{\partial x_i} = \sum_{\substack{j=1 \\ j \neq i}}^n 2\alpha A \left(e^{-\alpha(r-r_0)} - e^{-2\alpha(r-r_0)} \right) \cdot \frac{(x_i - x_j)}{r_{ij}} \quad (2.17)$$

and derivatives are computed with respect to y_i and z_i in a similar manner.

To illustrate the importance of the minimum gradient point and the effectiveness of the stability boundary tracing, we compared our results with other methods reported in [73]. Energy value at the given local minimum is -1775.7911. E_{min} is the energy value at the new local minimum. E_{saddle} is the energy value at the saddle point. E_{MGP} is the energy value at the Minimum Gradient Point. $RMSD$ is the root mean square distance of all the atoms at the MGP and the saddle point. Fig. 2.12

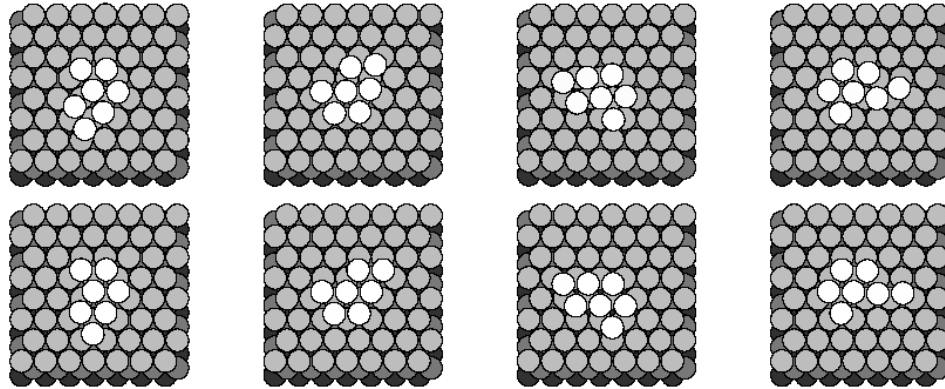


Figure 2.11: Some sample configurations of the heptamer island on the surface of the FCC crystal. First row- saddle points configurations. Second row- other local minimum energy configurations corresponding to the above saddle point configurations.

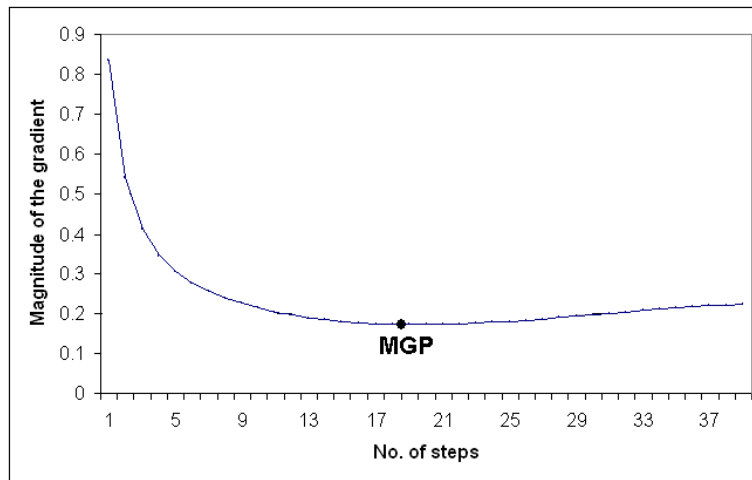


Figure 2.12: The gradient curve obtained in one of the higher dimensional test cases. MGP indicates the minimum gradient point where the magnitude of the gradient reaches the minimum value. This point is used as a initial guess for a local optimization method to obtain DDP.

Table 2.4: Results of our algorithm on a heptamer island over the FCC crystal. The number of force evaluations made to compute the MGP is given in the last column.

| No. | E_{min} | E_{saddle} | E_{MGP} | $RMSD$ | Δ Energy | Δ Force | fevals |
|-----|------------|--------------|------------|---------|-----------------|----------------|--------|
| 1 | -1775.7787 | -1775.19 | -1775.2139 | 0.00917 | 0.02376 | 0.01833 | 49 |
| 2 | -1775.7787 | -1775.1716 | -1775.1906 | 0.00802 | 0.01903 | 0.01671 | 43 |
| 3 | -1775.0079 | -1774.8055 | -1774.8213 | 0.0121 | 0.01578 | 0.01664 | 97 |
| 4 | -1775.006 | -1774.8041 | -1774.9228 | 0.03208 | 0.11868 | 0.02654 | 67 |
| 5 | -1775.0058 | -1774.8024 | -1774.8149 | 0.01229 | 0.01256 | 0.01704 | 91 |
| 6 | -1775.0942 | -1774.5956 | -1774.3819 | 0.04285 | -0.21362 | 0.04343 | 37 |
| 7 | -1775.0931 | -1774.5841 | -1774.3916 | 0.03296 | -0.19252 | 0.03905 | 43 |
| 8 | -1775.01 | -1774.3106 | -1775.0789 | 0.05287 | 0.76832 | 0.04031 | 97 |
| 9 | -1775.0097 | -1774.3082 | -1775.0848 | 0.05297 | 0.77662 | 0.04113 | 97 |
| 10 | -1774.3896 | -1774.2979 | -1774.9551 | 0.05623 | 0.65718 | 0.03103 | 79 |
| 11 | -1774.3928 | -1774.2997 | -1774.9541 | 0.05615 | 0.65439 | 0.03086 | 79 |
| 12 | -1774.3933 | -1774.2792 | -1774.2938 | 0.02262 | 0.01450 | 0.07092 | 19 |

shows the gradient curve for one of the higher dimensional test cases. Δ Energy is the energy difference between E_{MGP} and E_{saddle} . Δ Force is the magnitude of the gradient at the MGP. The results of our method is shown in Table 2.4. The last column indicates the number of gradient computations that were made to reach the minimum gradient point. As seen from the table, our method finds the saddle points with fewer number of gradient computations when compared to other methods. Typically, even the best available method takes at least 200-300 evaluations of the gradient. For detailed results about the performance of other methods refer to [73]. The MGP that was computed in our case varied from 0.01-0.05. The MGP can be treated as a saddle point for most of the practical applications. The RMSD (root mean square distance) value between the MGP and the saddle point is very low.

2.6.4 Special Cases : Eckhardt surface

The Eckhardt surface [47] is an exceptional case where we need to perturb the exit point in order to follow the stability boundary. Such cases almost never occur

in practice and hence dealing with such surfaces will not be given much importance. Eq. (2.18) gives the Eckhardt energy function. Fig. 2.13 represents the two-dimensional contour plot of the potential energy surface of the Eckhardt function.

$$C(x, y) = e^{-[x^2+(y+1)^2]} + e^{-[x^2+(y-1)^2]} + 4 e^{-[3(x^2+y^2)/2]} + y^2/2 \quad (2.18)$$

As shown in Fig. 2.13, there are two local minima (A,B), two dynamic decomposition points (1,2) on the Eckhardt potential energy surface. A local maximum is present exactly at the center of the vector joining the two local minima. The two dynamic decomposition points are on either side of the maximum. Table 2.5 shows the energy values at the local minima and the dynamic decomposition points.

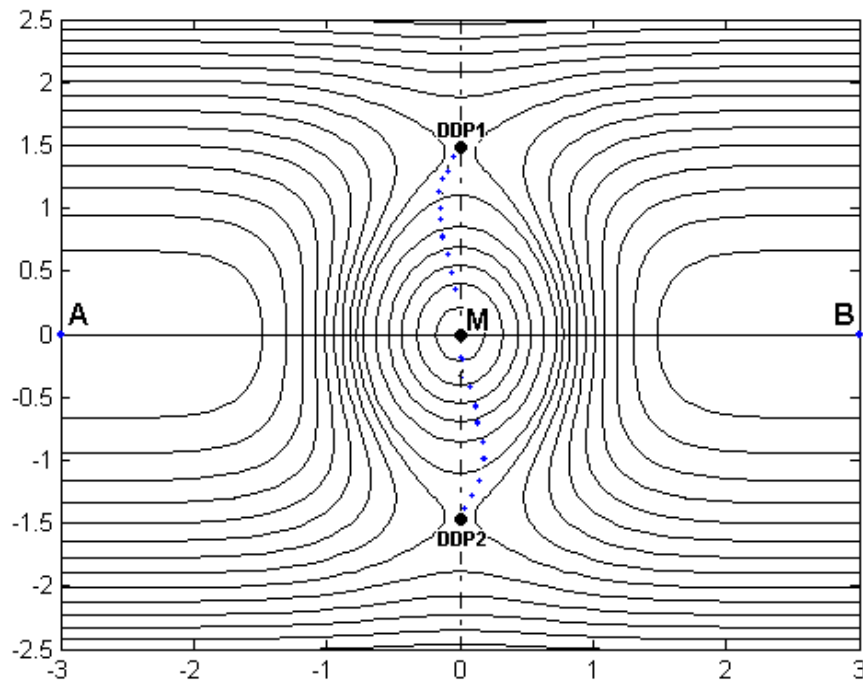


Figure 2.13: Two dimensional contour plot of the potential energy surface of the Eckhardt energy function described on Eq. (2.18). A and B are stable equilibrium points and DDP1,DDP2 are two dynamic decomposition points. M is a source. The dots correspond to the points on the stability boundary during the stability boundary following procedure.

Table 2.5: Stable equilibrium points and dynamic decomposition points for the Eckhardt Surface.

| Equilibrium Points | Location | Energy value $c(\cdot)$ |
|--------------------|---------------|-------------------------|
| SEP A | (-3.0,0.0) | 0.0 |
| DDP 1 | (0.0,1.4644) | 2.0409 |
| SEP B | (3.0,0.0) | 0.0 |
| DDP 2 | (0.0,-1.4644) | 2.0409 |
| SEP M | (0.0,0.0) | 4.7358 |

We construct the dynamical system corresponding to (2.18) as follows:

$$\begin{bmatrix} \dot{x}(t) \\ \dot{y}(t) \end{bmatrix} = - \begin{bmatrix} \frac{\partial C}{\partial x} \\ \frac{\partial C}{\partial y} \end{bmatrix}$$

$$\frac{\partial C}{\partial x} = -2x e^{-[x^2+(y+1)^2]} - 2x e^{-[x^2+(y-1)^2]} - 12x e^{-[3(x^2+y^2)/2]}$$

$$\frac{\partial C}{\partial y} = -2(y+1) e^{-[x^2+(y+1)^2]} - 2(y-1) e^{-[x^2+(y-1)^2]} - 12y e^{-[3(x^2+y^2)/2]} + y$$

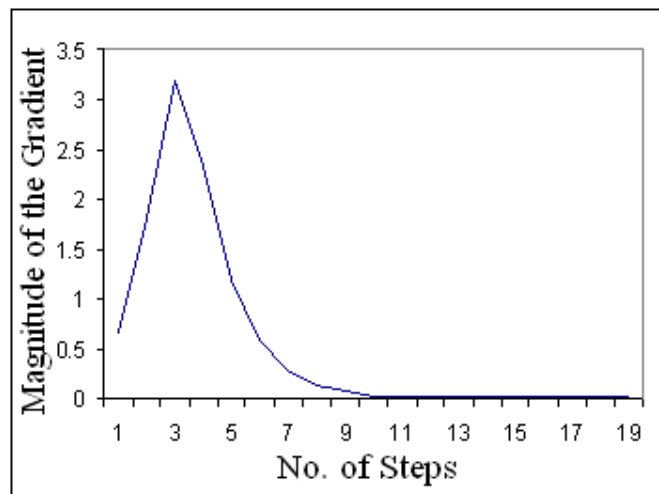


Figure 2.14: Gradient curve corresponding to the various points along the stability boundary on the Eckhardt surface.

This surface can be treated as a special case where there are different critical points lying on the vector joining the two local minima. As seen from the Fig. 2.13,

there is a local maximum at $(0,0)$, which is the exit point obtained. It should be noted that this system is also symmetric and hence one can obtain the exit point analytically. Since the exit point is also a critical point, the exit point is first perturbed and then the stability boundary following procedure is used to compute the MGP. The two dynamic decomposition points are at $(0.0,1.4644)$ and $(0.0,-1.4644)$. The gradient curve is shown in Fig. 2.14. From this MGP, the local minimizer is applied to obtain the DDP1. The other dynamic decomposition point (DDP2) is similarly obtained. Our method was also successful in finding saddle points on other two dimensional test surfaces like Minyaev Quapp [102], NFK (Neria-Fischer-Karplus) [107] etc.

2.7 Discussion

Saddle points play a vital role in realizing the folding pathways of a protein as well as in understanding the transition state structures during chemical reactions. This chapter primarily focuses on a new TRUST-TECH based method for finding saddle points on potential energy surfaces using stability boundaries. Our approach is based on some fundamental results of nonlinear dynamical systems. A novel *stability boundary following procedure* has been used to move along the stability boundary. Our method was able to find the saddle points on a wide range of surfaces with varying dimensions. The primary advantage of our method comes from the fact that the stability boundary following is computationally more efficient than directly searching for a saddle point from the given local minimum. Deterministic nature of the algorithm ensures that the same saddle point is obtained every run. Very few user-specific parameters makes it easy for a new user to implement. We have also explored the symmetric behaviour of some energy surfaces to obtain the exit

point analytically and developed a simplified version of our algorithm to compute the saddle points. The algorithm has also been tested successfully on a heptamer island over the surface of an FCC crystal.

Finding saddle points can be of importance for other problems related to global optimization. This can be done by using our method presented here as a tool for *escaping from a given local minimum to another local minimum in the neighborhood*. Though the current work assumes that the two local minima between which the saddle point is computed are given, it can be easily extended to find saddle points from a single local minimum.

APPENDIX-A: LEPS Potential

The model of LEPS potential simulates a reaction involving three atoms confined to motion along a line. Only one bond can be formed either between atoms A and B or between atoms B and C.

$$C(x, y) = \frac{Q_{AB}}{1+a} + \frac{Q_{BC}}{1+b} + \frac{Q_{AC}}{1+c} - \left[\frac{J_{AB}^2}{(1+a)^2} + \frac{J_{BC}^2}{(1+b)^2} + \frac{J_{AC}^2}{(1+c)^2} + \frac{J_{AB}J_{BC}}{(1+a)(1+b)} + \frac{J_{BC}J_{AC}}{(1+b)(1+c)} + \frac{J_{AB}J_{AC}}{(1+a)(1+c)} \right]^{\frac{1}{2}}$$

where the Q functions represent the Coulomb interactions between electron clouds and the nuclei and the J functions represent the quantum mechanical exchange interactions. The form of the Q and J functions is given below:

$$Q(r) = \frac{d}{2} \left(\frac{3}{2} e^{-2\alpha(r-r_0)} - e^{-\alpha(r-r_0)} \right)$$

$$J(r) = \frac{d}{4} \left(e^{-2\alpha(r-r_0)} - 6e^{-\alpha(r-r_0)} \right)$$

The parameters were chosen to be $a = b = c = 0.05$, $d_{AB} = d_{AC} = d_{BC} = 4.746$, $\alpha = 1.942$ and $r_0 = 0.742$. The details about the LEPS potential are given in [115].

APPENDIX-B: Golden section search

The following procedure describes the golden section search method. Let a and b be the two intervals between which the *exit point* is located. Golden section search method computes the exit point with an accuracy of $\pm\epsilon$. r is the *golden mean* ($\frac{3-\sqrt{5}}{2}$) [116]. $f(x)$ returns the function value at point x .

procedure *GoldenSectionSearch*(a, b, ϵ)

- 1: Initialize $r = 0.38197$ (golden mean)
- 2: $c = a + r(b - a)$
- 3: $d = b - r(b - a)$
- 4: **while** $|b - a| > \epsilon$ **do**
- 5: **if** $f(c) > f(d)$ **then**
- 6: $b = d, d = c, c = a + r(b - a)$
- 7: **else**
- 8: $a = c, c = d, d = b - r(b - a)$
- 9: **end if**
- 10: **end while**
- 11: **return** b

Chapter 3

TRUST-TECH based Expectation

Maximization for Learning Mixture

Models

In this chapter, we develop a TRUST-TECH based algorithm for solving the problem of mixture modeling. In the field of statistical pattern recognition, finite mixtures allow a probabilistic model-based approach to unsupervised learning [100]. One of the most popular methods used for fitting mixture models to the observed data is the *Expectation-Maximization* (EM) algorithm which converges to the maximum likelihood estimate of the mixture parameters locally [37, 123]. The usual steepest descent, conjugate gradient, or Newton-Raphson methods are too complicated for use in solving this problem [143]. EM has become a popular method since it takes advantage of problem specific properties. EM based methods have been successfully applied to solve a wide range of problems that arise in pattern recognition [10, 13], clustering [6], information retrieval [109], computer vision [23], data mining [134] etc.

Without loss of generality, we will consider the problem of learning parameters of Gaussian Mixture Models (GMM). Fig 3.1 shows data generated by three Gaussian components with different mean and variance. Note that every data point has a probabilistic (or soft) membership that gives the probability with which it belongs to each of the components. Points that belong to component 1 will have high probability of membership for component 1. On the other hand, data points belonging to components 2 and 3 are not well separated. The problem of learning mixture models involves estimating the parameters of these components and finding the probabilities with which each data point belongs to these components. Given

the number of components and an initial set of parameters, EM algorithm computes the optimal estimates of the parameters that maximize the likelihood of the data given the estimates of these components. However, the main problem with the EM algorithm is that it is a ‘*greedy*’ method which is very sensitive to the given initial set of parameters. To overcome this problem, a novel three-stage algorithm is proposed [121]. The main research concerns that motivated the new algorithm presented in this chapter are :

- EM algorithm converges to a local maximum of the likelihood function very quickly.
- There are several other promising local optimal solutions in the vicinity of the solutions obtained from the methods that provide good initial guesses of the solution.
- Model selection criteria usually assumes that the global optimal solution of the log-likelihood function can be obtained. However, achieving this is computationally intractable.
- Some regions in the search space do not contain any promising solutions. The promising and non-promising regions coexist and it becomes challenging to avoid wasting computational resources to search in non-promising regions.

Of all the concerns mentioned above, the fact that most of the local maxima are not distributed uniformly [138] makes it important to develop algorithms that can avoid searching in the low-likelihood regions and focus on exploring promising subspaces more thoroughly. This subspace search will also be useful for making the solution less sensitive to the initial set of parameters. Here, we propose a novel three-stage algorithm for estimating the parameters of mixture models. Using

TRUST-TECH method and EM algorithm simultaneously to exploit the problem specific features of the mixture models, the proposed three-stage algorithm obtains the optimal set of parameters by searching for the global maximum in a systematic manner.

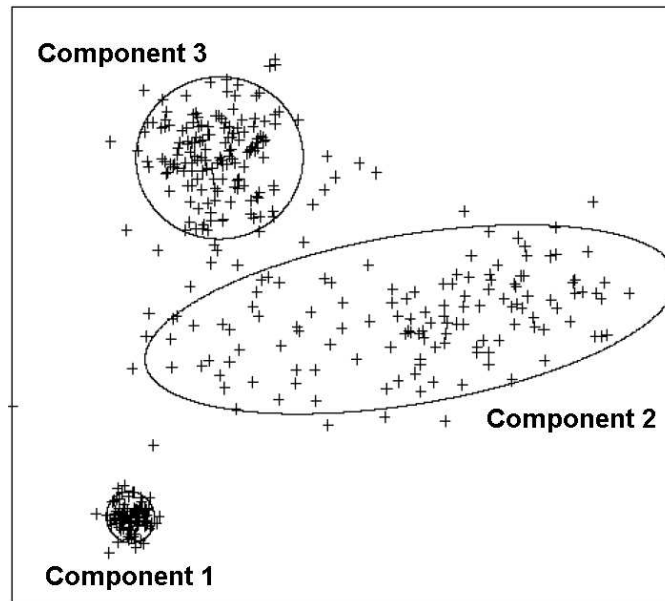


Figure 3.1: Data consisting of three Gaussian components with different mean and variance values. Note that each data point doesn't have a hard membership that it belongs to only one component. Most of the points in the first component will have high probability with which they belong to it. In this case, the other components do not have much influence. Components 2 and 3 data points are not clearly separated. The problem of learning mixture models involves estimating the parameters of the Gaussian components and finding the probabilities with which each data sample belongs to the component.

3.1 Relevant Background

Although EM and its variants have been extensively used for learning mixture models, several researchers have approached the problem by identifying new techniques that give good initial points. More generic techniques like deterministic annealing [127, 138], genetic algorithms [112, 97] have been applied to obtain a good set of parameters. Though, these techniques have asymptotic guarantees, they are very time

consuming and hence may not be used for most of the practical applications. Some problem specific algorithms like split and merge EM [139], component-wise EM [55], greedy learning [140], incremental version for sparse representations [106], parameter space grid [94] are also proposed in the literature. Some of these algorithms are either computationally very expensive or infeasible when learning mixtures in high dimensional spaces [94]. In spite of all the expense in these methods, very little effort has been taken to explore promising subspaces within the larger parameter space. Most of these algorithms eventually apply the EM algorithm to move to a locally maximal set of parameters on the likelihood surface. Simple approaches like running EM from several random initializations, and then choosing the final estimate that leads to the local maximum with higher value of the likelihood can be successful to certain extent [67, 126].

Though some of these methods apply other additional mechanisms (like perturbations [49]) to escape out of the local optimal solutions, systematic methods are yet to be developed for searching the subspace. The dynamical system of the log-likelihood function reveals more information about the topology of the nonlinear log-likelihood surface [31]. Hence, the difficulties of finding good solutions when the error surface is very rugged can be overcome by understanding the geometric and dynamic characteristics of the log-likelihood surface. Though this method might introduce some additional cost, one has to realize that existing approaches are much more expensive due to their stochastic nature. Specifically, for a problem in this context, where there is a non-uniform distribution of local maxima, it is difficult for most of the methods to search neighboring regions [144]. For this reason, it is more desirable to apply TRUST-TECH based Expectation Maximization (TT-EM) algorithm after obtaining some point in a promising region. The main advantages of the proposed algorithm are that it :

- Explores most of the neighborhood local optimal solutions unlike the traditional stochastic algorithms.
- Acts as a flexible interface between the EM algorithm and other global method. Sometimes, a global method will optimize an approximation of the original function. Hence, it is important to provide an interface between the EM algorithm and the global method.
- Allows the user to work with existing clusters obtained from the traditional approaches and improves the quality of the solutions based on the maximum likelihood criteria.
- Helps the expensive global methods to truncate early.
- Exploits the heuristics that the EM algorithm that it converges at a faster rate if the solutions are promising.

While trying to obtain multiple optimal solutions, TRUST-TECH can dynamically change the threshold for the number of iterations. For e.g. while computing Tier-1 solutions, if a promising solution has been obtained with a few iterations, then all the rest of the tier-1 solutions will use this value as their threshold.

3.2 Preliminaries

We will now introduce some necessary preliminaries on mixture models, EM algorithm and nonlinear transformation. Table 3.1 gives the notations used in this chapter :

Table 3.1: Description of the Notations used

| Notation | Description |
|---------------|---------------------------------------|
| d | number of features |
| n | number of data points |
| k | number of components |
| s | total number of parameters |
| Θ | parameter set |
| θ_i | parameters of i^{th} component |
| α_i | mixing weights for i^{th} component |
| \mathcal{X} | observed data |
| \mathcal{Z} | missing data |
| \mathcal{Y} | complete data |
| t | timestep for the estimates |

3.2.1 Mixture Models

Lets assume that there are k Gaussians in the mixture model. The form of the probability density function is as follows :

$$p(x|\Theta) = \sum_{i=1}^k \alpha_i p(x|\theta_i) \quad (3.1)$$

where $x = [x_1, x_2, \dots, x_d]^T$ is the feature vector of d dimensions. The α_k 's represent the *mixing weights*. Θ represents the parameter set $(\alpha_1, \alpha_2, \dots, \alpha_k, \theta_1, \theta_2, \dots, \theta_k)$ and p is a univariate Gaussian density parameterized by θ_i (i.e. μ_i and σ_i):

$$p(x|\theta_i) = \frac{1}{\sqrt{(2\pi)\sigma_i}} e^{-\frac{(x-\mu_i)^2}{2\sigma_i^2}} \quad (3.2)$$

Also, it should be noticed that being probabilities α_i must satisfy

$$0 \leq \alpha_i \leq 1, \forall i = 1, \dots, k, \text{ and } \sum_{i=1}^k \alpha_i = 1 \quad (3.3)$$

Given a set of n i.i.d samples $\mathcal{X} = \{x^{(1)}, x^{(2)}, \dots, x^{(n)}\}$, the log-likelihood corre-

sponding to a mixture is

$$\log p(\mathcal{X}|\Theta) = \log \prod_{j=1}^n p(x^{(j)}|\Theta) = \sum_{j=1}^n \log \sum_{i=1}^k \alpha_i p(x^{(j)}|\theta_i) \quad (3.4)$$

The goal of learning mixture models is to obtain the parameters $\hat{\Theta}$ from a set of n data points which are the samples of a distribution with density given by (3.1). The *Maximum Likelihood Estimate* (MLE) is given by :

$$\hat{\Theta}_{MLE} = \arg \max_{\Theta} \{ \log p(\mathcal{X}|\Theta) \} \quad (3.5)$$

where $\tilde{\Theta}$ indicates the entire parameter space. Since, this MLE cannot be found analytically for mixture models, one has to rely on iterative procedures that can find the global maximum of $\log p(\mathcal{X}|\Theta)$. The EM algorithm described in the next section has been used successfully to find the local maximum of such a function [98].

3.2.2 Expectation Maximization

The EM algorithm assumes \mathcal{X} to be *observed* data. The missing part, termed as *hidden* data, is a set of n labels $\mathcal{Z} = \{\mathbf{z}^{(1)}, \mathbf{z}^{(2)}, \dots, \mathbf{z}^{(n)}\}$ associated with n samples, indicating which component produced each sample [98]. Each label $\mathbf{z}^{(j)} = [z_1^{(j)}, z_2^{(j)}, \dots, z_k^{(j)}]$ is a binary vector where $z_i^{(j)} = 1$ and $z_m^{(j)} = 0 \forall m \neq i$, means the sample $x^{(j)}$ was produced by the i^{th} component. Now, the complete log-likelihood i.e. the one from which we would estimate Θ if the *complete data* $\mathcal{Y} = \{ \mathcal{X}, \mathcal{Z} \}$ is

$$\log p(\mathcal{X}, \mathcal{Z}|\Theta) = \sum_{j=1}^n \log \prod_{i=1}^k [\alpha_i p(x^{(j)}|\theta_i)]^{z_i^{(j)}}$$

$$\log p(\mathcal{Y}|\Theta) = \sum_{j=1}^n \sum_{i=1}^k z_i^{(j)} \log [\alpha_i p(x^{(j)}|\theta_i)] \quad (3.6)$$

The EM algorithm produces a sequence of estimates $\{\widehat{\Theta}(t), t = 0, 1, 2, \dots\}$ by alternately applying the following two steps until convergence :

- **E-Step** : Compute the conditional expectation of the hidden data, given \mathcal{X} and the current estimate $\widehat{\Theta}(t)$. Since $\log p(\mathcal{X}, \mathcal{Z}|\Theta)$ is linear with respect to the missing data \mathcal{Z} , we simply have to compute the conditional expectation $\mathcal{W} \equiv E[\mathcal{Z}|\mathcal{X}, \widehat{\Theta}(t)]$, and plug it into $\log p(\mathcal{X}, \mathcal{Z}|\Theta)$. This gives the Q -function as follows :

$$Q(\Theta|\widehat{\Theta}(t)) \equiv E_{\mathcal{Z}}[\log p(\mathcal{X}, \mathcal{Z}|\Theta)|\mathcal{X}, \widehat{\Theta}(t)] \quad (3.7)$$

Since \mathcal{Z} is a binary vector, its conditional expectation is given by :

$$w_i^{(j)} \equiv E [z_i^{(j)} | \mathcal{X}, \widehat{\Theta}(t)] = Pr [z_i^{(j)} = 1 | x^{(j)}, \widehat{\Theta}(t)] = \frac{\widehat{\alpha}_i(t)p(x^{(j)}|\widehat{\theta}_i(t))}{\sum_{i=1}^k \widehat{\alpha}_i(t)p(x^{(j)}|\widehat{\theta}_i(t))} \quad (3.8)$$

where the last equality is simply the Bayes law (α_i is the a priori probability that $z_i^{(j)} = 1$), while $w_i^{(j)}$ is the a posteriori probability that $z_i^{(j)} = 1$ given the observation $x^{(j)}$.

- **M-Step** : The estimates of the new parameters are updated using the following equation :

$$\widehat{\Theta}(t+1) = \underset{\Theta}{arg \max} \{Q(\Theta, \widehat{\Theta}(t))\} \quad (3.9)$$

3.2.3 EM for GMMs

Several variants of the EM algorithm have been extensively used to solve this problem. The convergence properties of the EM algorithm for Gaussian mixtures are thoroughly discussed in [143]. The Q – *function* for GMM is given by :

$$Q(\Theta|\widehat{\Theta}(t)) = \sum_{j=1}^n \sum_{i=1}^k w_i^{(j)} \left[\log \frac{1}{\sigma_i \sqrt{2\pi}} - \frac{(x^{(j)} - \mu_i)^2}{2\sigma_i^2} + \log \alpha_i \right] \quad (3.10)$$

where

$$w_i^{(j)} = \frac{\frac{\alpha_i(t)}{\sigma_i(t)} e^{-\frac{1}{2\sigma_i(t)^2}(x^{(j)} - \mu_i(t))^2}}{\sum_{i=1}^k \frac{\alpha_i(t)}{\sigma_i(t)} e^{-\frac{1}{2\sigma_i(t)^2}(x^{(j)} - \mu_i(t))^2}} \quad (3.11)$$

The maximization step is given by the following equation :

$$\frac{\partial}{\partial \Theta_k} Q(\Theta|\widehat{\Theta}(t)) = 0 \quad (3.12)$$

where Θ_k is the parameters for the k^{th} component. Because of the assumption made that each data point comes from a single component, solving the above equation becomes trivial. The updates for the maximization step in the case of GMMs are given as follows :

$$\mu_i(t+1) = \frac{\sum_{j=1}^n w_i^{(j)} x^{(j)}}{\sum_{j=1}^n w_i^{(j)}} \quad (3.13)$$

$$\sigma_i^2(t+1) = \frac{\sum_{j=1}^n w_i^{(j)} (x^{(j)} - \mu_i(t+1))^2}{\sum_{j=1}^n w_i^{(j)}} \quad (3.14)$$

$$\alpha_i(t+1) = \frac{1}{n} \sum_{j=1}^n w_i^{(j)} \quad (3.15)$$

3.2.4 Nonlinear Transformation

This section mainly deals with the transformation of the original log-likelihood function into its corresponding nonlinear dynamical system and introduces some terminology pertinent to comprehend our algorithm. This transformation gives the correspondence between all the critical points of the s -dimensional likelihood surface and that of its dynamical system. For the case of spherical Gaussian mixtures with k components, we have the number of unknown parameters $s = 3k - 1$. For convenience, the maximization problem is transformed into a minimization problem defined by the following objective function :

$$\max_{\Theta} \{ \log p(\mathcal{X}|\Theta) \} = \min_{\Theta} \{ -\log p(\mathcal{X}|\Theta) \} = \min_{\Theta} f(\Theta) \quad (3.16)$$

Lemma 1 $f(\Theta)$ is $C^2(\mathfrak{R}^s, \mathfrak{R})$.

Proof.

Note from Eq.(3.4), we have

$$f(\Theta) = -\log p(\mathcal{X}|\Theta) = -\sum_{j=1}^n \log \sum_{i=1}^k \alpha_i p(\mathbf{x}^{(j)}|\boldsymbol{\theta}_i) \quad (3.17)$$

Each of the simple functions which appear in Eq. (3.17) are twice differentiable and continuous in the interior of the domain over which $f(\Theta)$ is defined. The function $f(\Theta)$ is composed of arithmetic operations of these simple functions and from basic results in analysis, we can conclude that $f(\Theta)$ is twice continuously differentiable. \triangleleft

Lemma 1 and the preceding arguments guarantee the existence of the gradient system associated with $f(\Theta)$ for the log-likelihood function in the case of spherical

Gaussians and allows us to construct the following negative gradient system :

$$\begin{aligned} & [\dot{\mu}_1(t) \dots \dot{\mu}_k(t) \dot{\sigma}_1(t) \dots \dot{\sigma}_k(t) \dot{\alpha}_1(t) \dots \dot{\alpha}_{k-1}(t)]^T \\ = & - \left[\frac{\partial f}{\partial \mu_1} \dots \frac{\partial f}{\partial \mu_k} \frac{\partial f}{\partial \sigma_1} \dots \frac{\partial f}{\partial \sigma_k} \frac{\partial f}{\partial \alpha_1} \dots \frac{\partial f}{\partial \alpha_{k-1}} \right]^T \end{aligned} \quad (3.18)$$

Theorem 3.2.1 (*Stability*): *The gradient system 3.18 is completely stable.*

Proof: See Appendix-A.

Developing a gradient system is one of the simplest transformation possible. One can think of a more complicated nonlinear transformations as well. We will now describe three main guidelines that must be satisfied by the transformation :

- The original log-likelihood function must be a Lyapunov function for the dynamical system.
- The location of the critical points must be preserved under this transformation.
- The system must be completely stable. In other words, every trajectory $\Phi(x, t)$ must be bounded.

From the implementation point of view, it is not required to construct this gradient system. However, to understand the details of our method, it is necessary to obtain this gradient system. For simplicity, we show the construction of the gradient system for the case of spherical Gaussians. It can be easily extended to the full covariance Gaussian mixture case. It should be noted that only $(k-1)$ α values are considered in the gradient system because of the unity constraint. The dependent variable α_k is written as follows :

$$\alpha_k = 1 - \sum_{j=1}^{k-1} \alpha_j \quad (3.19)$$

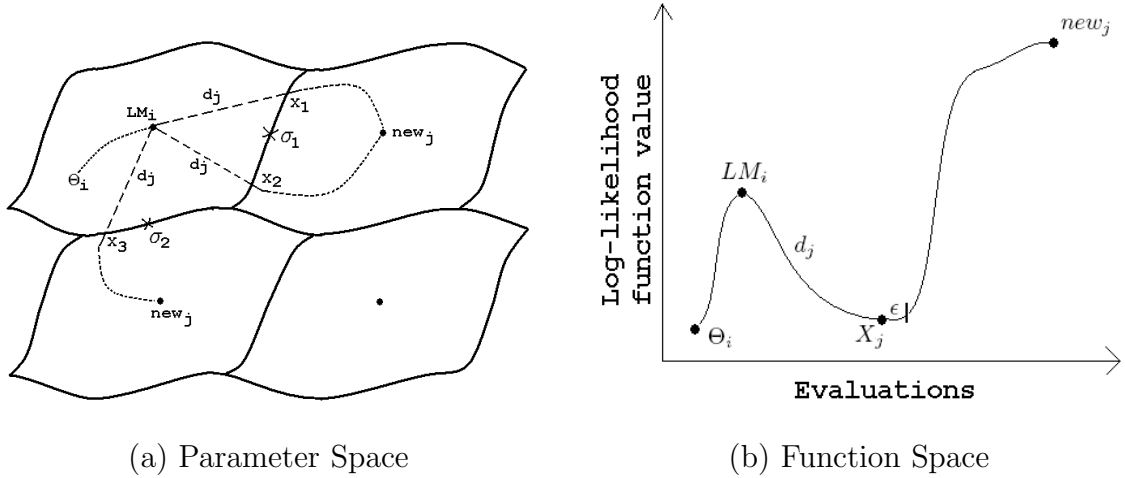


Figure 3.2: Various stages of our algorithm in (a) Parameter space - the solid lines indicate the practical stability boundary. Points highlighted on the stability boundary are the decomposition points. The dotted arrows indicate the convergence of the EM algorithm. The dashed lines indicate the neighborhood-search stage. x_1 and x_2 are the exit points on the practical stability boundary (b) Different points in the function space and their corresponding log-likelihood function values.

This gradient system and the decomposition points on the practical stability boundary of the stable equilibrium points will enable us to define *Tier-1 stable equilibrium point*.

Definition 4 For a given stable equilibrium point (x_s), a *Tier-1 stable equilibrium point* is defined as a stable equilibrium point whose stability boundary intersects with the stability boundary of x_s .

3.3 TRUST-TECH based Expectation Maximization

Our framework consists three stages namely: (i) global stage, (ii) local stage and (iii) neighborhood-search stage. The last two stages are repeated in the solution space to obtain promising solutions. Global method obtains promising subspaces of the solution space. The next stage is the local stage (or the EM stage) where the results

from the global methods are refined to the corresponding locally optimal parameter set. Then, during the neighborhood search stage, the exit points are computed and the neighborhood solutions are systematically explored through these exit points. Fig. 3.2 shows the different steps of our algorithm both in (a) the parameter space and (b) the function space.

It is beneficial to use the TRUST-TECH based algorithm at the promising subspaces. In this sense, the neighborhood-search stage can act as an interface between global methods for initialization and the EM algorithm which gives the local maxima. This approach differs from traditional local methods by computing multiple local maxima in the neighborhood region. This also enhances user flexibility in choosing between different sets of good clusterings. Though global methods can identify promising subsets, it is important to explore this more thoroughly especially in problems like parameter estimation.

Algorithm 1 TRUST-TECH based EM Algorithm

Input: Parameters Θ , Data \mathcal{X} , tolerance τ , Step S_p

Output: $\hat{\Theta}_{MLE}$

Algorithm:

Apply global method and store the q promising solutions $\Theta_{init} = \{\Theta_1, \Theta_2, \dots, \Theta_q\}$

Initialize $E = \phi$

while $\Theta_{init} \neq \phi$ **do**

Choose $\Theta_i \in \Theta_{init}$, set $\Theta_{init} = \Theta_{init} \setminus \{\Theta_i\}$

$LM_i = EM(\Theta_i, \mathcal{X}, \tau)$ $E = E \cup \{LM_i\}$

Generate promising direction vectors d_j from LM_i

for each d_j **do**

Compute Exit Point (X_j) along d_j starting from LM_i by evaluating the log-likelihood function given by (3.4)

$New_j = EM(X_j + \epsilon \cdot d_j, \mathcal{X}, \tau)$

if $new_j \notin E$ **then**

$E = E \cup New_j$

end if

end for

end while

$\hat{\Theta}_{MLE} = \max\{val(E_i)\}$

In order to escape out of a found local maximum, our method needs to compute certain promising directions based on the local behaviour of the function. One can realize that generating these promising directions is one of the important aspects of our algorithm. Surprisingly, choosing random directions to move out of the local maximum works well for this problem. One might also use other directions like eigenvectors of the Hessian or incorporate some domain-specific knowledge (like information about priors, approximate location of cluster means, user preferences on the final clusters) depending on the application that they are working on and the level of computational expense that they can afford. We used random directions in our work because they are very cheap to compute. Once the promising directions are generated, exit points are computed along these directions. *Exit points* are points of intersection between any given direction and the practical stability boundary of that local maximum along that particular direction. If the stability boundary is not encountered along a given direction, then there is a guarantee that one will not be able to find any new local maximum in that direction. With a new initial guess in the vicinity of the exit points, EM algorithm is applied again to obtain a new local maximum. Sometimes, this new point $(X_j + \epsilon \cdot d_j)$ might have convergence problems. In such cases, TRUST-TECH can help the convergence by integrating the dynamical system and obtaining another point that is much closer to the local optimal solution. However, this is not done here because of the fact that the computation of gradient for log-likelihood function is expensive.

3.4 Implementation Details

Our program is implemented in MATLAB and runs on Pentium IV 2.8 GHz machine. The main procedure implemented is *TT_EM* described in Algorithm 2. The

Algorithm 2 Params[] *TT_EM*(*Pset*, *Data*, *Tol*, *Step*)

Val = *eval*(*Pset*)
Dir[] = *Gen_Dir*(*Pset*)
Eval_MAX = 500
for *k* = 1 to *size*(*Dir*) **do**
 Params[*k*] = *Pset* *ExtPt* = *OFF*
 Prev_Val = *Val* *Cnt* = 0
 while (! *ExtPt*) && (*Cnt* < *Eval_MAX*) **do**
 Params[*k*] = *update*(*Params*[*k*], *Dir*[*k*], *Step*)
 Cnt = *Cnt* + 1
 Next_Val = *eval*(*Params*[*k*])
 if (*Next_Val* > *Prev_Val*) **then**
 ExtPt = *ON*
 end if
 Prev_Val = *Next_Val*
 end while
 if *count* < *Eval_MAX* **then**
 Params[*k*] = *update*(*Params*[*k*], *Dir*[*k*], *ASC*)
 Params[*k*] = *EM*(*Params*[*k*], *Data*, *Tol*)
 else
 Params[*k*] = *NULL*
 end if
end for
Return *max*(*eval*(*Params*[]))

algorithm takes the mixture data and the initial set of parameters as input along with step size for moving out and tolerance for convergence in the EM algorithm. It returns the set of parameters that correspond to the Tier-1 neighboring local optimal solutions. The procedure *eval* returns the log-likelihood score given by Eq. (3.4). The *Gen_Dir* procedure generates promising directions from the local maximum. Exit points are obtained along these generated directions. The procedure *update* moves the current parameter to the next parameter set along a given k^{th} direction $Dir[k]$. Some of the directions might have one of the following two problems: (i) exit points might not be obtained in these directions. (ii) even if the exit point is obtained it might converge to a less promising solution. If the exit points are not found along these directions, search will be terminated after *Eval_MAX* number of evaluations. For all exit points that are successfully found, *EM* procedure is applied and all the corresponding neighborhood set of parameters are stored in the *Params[]*. To ensure that the new initial points are in a new convergence region of the EM algorithm, one should move (along that particular direction) ‘ ϵ ’ away from the exit points. Since, different parameters will be of different ranges, care must be taken while multiplying with the step sizes. It is important to use the current estimates to get an approximation of the step size with which one should move out along each parameter in the search space. Finally, the solution with the highest likelihood score amongst the original set of parameters and the Tier-1 solutions is returned.

3.5 Results and Discussion

Our algorithm has been tested on both synthetic and real datasets. The initial values for the centers and the covariances were chosen uniformly random. Uniform

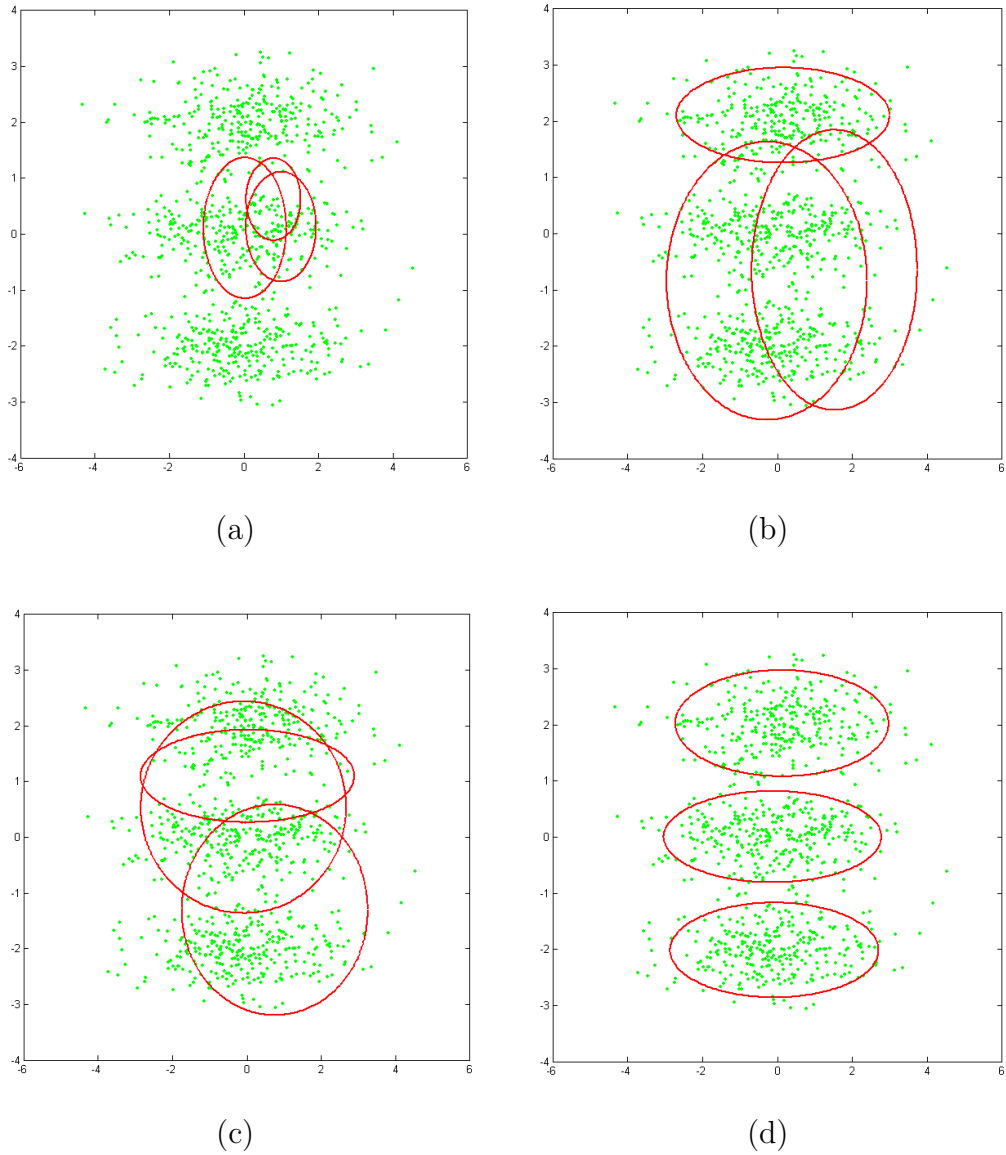


Figure 3.3: Parameter estimates at various stages of our algorithm on the three component Gaussian mixture model (a) Poor random initial guess (b) Local maximum obtained after applying EM algorithm with the poor initial guess (c) Exit point obtained by our algorithm (d) The final solution obtained by applying the EM algorithm using the exit point as the initial guess.

priors were chosen for initializing the components. For real datasets, the centers were chosen randomly from the sample points.

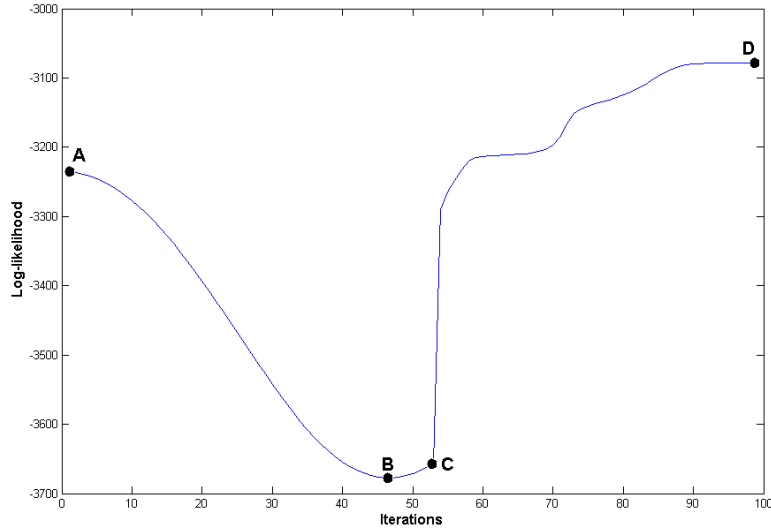


Figure 3.4: Graph showing likelihood vs Evaluations. A corresponds to the original local maximum (L=-3235.0). B corresponds to the exit point (L=-3676.1). C corresponds to the new initial point (L=-3657.3) after moving out by ‘ ϵ ’. D corresponds to the new local maximum (L=-3078.7).

3.5.1 Synthetic Datasets

A simple synthetic data with 40 samples and 5 spherical Gaussian components was generated and tested with our algorithm. Priors were uniform and the standard deviation was 0.01. The centers for the five components are given as follows: $\mu_1 = [0.3 \ 0.3]^T$, $\mu_2 = [0.5 \ 0.5]^T$, $\mu_3 = [0.7 \ 0.7]^T$, $\mu_4 = [0.3 \ 0.7]^T$ and $\mu_5 = [0.7 \ 0.3]^T$.

The second dataset was that of a diagonal covariance case containing $n = 900$ data points. The data generated from a two-dimensional, three-component Gaussian mixture distribution with mean vectors at $[0 \ -2]^T$, $[0 \ 0]^T$, $[0 \ 2]^T$ and same diagonal covariance matrix with values 2 and 0.2 along the diagonal [138]. All the three mixtures have uniform priors. Fig. 5.9 shows various stages of our algorithm and demonstrates how the clusters obtained from existing algorithms are improved

using our algorithm. The initial clusters obtained are of low quality because of the poor initial set of parameters. Our algorithm takes these clusters and applies the neighborhood-search stage and the EM stage simultaneously to obtain the final result. Fig. 3.4 shows the value of the log-likelihood during the neighborhood-search stage and the EM iterations.

In the third synthetic dataset, a more complicated overlapping Gaussian mixtures are considered [55]. The parameters are as follows: $\mu_1 = \mu_2 = [-4 \ -4]^T$, $\mu_3 = [2 \ 2]^T$ and $\mu_4 = [-1 \ -6]^T$. $\alpha_1 = \alpha_2 = \alpha_3 = 0.3$ and $\alpha_4 = 0.1$.

$$\Sigma_1 = \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1 \end{bmatrix} \quad \Sigma_2 = \begin{bmatrix} 6 & -2 \\ -2 & 6 \end{bmatrix}$$

$$\Sigma_3 = \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix} \quad \Sigma_4 = \begin{bmatrix} 0.125 & 0 \\ 0 & 0.125 \end{bmatrix}$$

Table 3.2: Performance of TRUST-TECH-EM algorithm on an average of 100 runs on various synthetic and real datasets compared with random start EM algorithm

| Dataset | Samples | Clusters | Features | EM(mean \pm std) | TT-EM(mean \pm std) |
|------------|---------|----------|----------|----------------------|-----------------------|
| Spherical | 40 | 5 | 2 | 38.07 \pm 2.12 | 43.55 \pm 0.6 |
| Elliptical | 900 | 3 | 2 | -3235 \pm 0.34 | -3078.7 \pm 0.03 |
| FC1 | 500 | 4 | 2 | -2345.5 \pm 175.13 | -2121.9 \pm 21.16 |
| FC2 | 2000 | 4 | 2 | -9309.9 \pm 694.74 | -8609.7 \pm 37.02 |
| Iris | 150 | 3 | 4 | -198.13 \pm 27.25 | -173.63 \pm 11.72 |
| Wine | 178 | 3 | 13 | -1652.7 \pm 1342.1 | -1618.3 \pm 1349.9 |

3.5.2 Real Datasets

Two real datasets obtained from the UCI Machine Learning repository [16] were also used for testing the performance of our algorithm. Most widely used Iris data with 150 samples, 3 classes and 4 features was used. Wine data set with 178 samples

was also used for testing. Wine data had 3 classes and 13 features. For these real data sets, the class labels were deleted thus treating it as an unsupervised learning problem. Table 3.2 summarizes our results over 100 runs. The mean and the standard deviations of the log-likelihood values are reported. The traditional EM algorithm with random starts is compared against our algorithm on both synthetic and real data sets. Our algorithm not only obtains higher likelihood value but also produces it with high confidence. The low standard deviation of our results indicates the robustness of obtaining the global maximum. In the case of the wine data, the improvements with our algorithm are not much significant compared to the other datasets. This might be due to the fact that the dataset might not have Gaussian components. Our method assumes that the underlying distribution of the data is mixture of Gaussians. Table 3.3 gives the results of TRUST-TECH-EM compared with other methods like split and merge EM and k-means+EM proposed in the literature.

Table 3.3: Comparison of TRUST-TECH-EM with other methods

| Method | Elliptical | Iris |
|---------------|------------------|-----------------|
| RS+EM | -3235 ± 14.2 | -198 ± 27 |
| K-Means+EM | -3195 ± 54 | -186 ± 10 |
| SMEM | -3123 ± 54 | -178.5 ± 6 |
| TRUST-TECH-EM | -3079 ± 0.03 | -173.6 ± 11 |

3.5.3 Discussion

It will be effective to use TRUST-TECH-EM for those solutions that appear to be promising. Due to the nature of the problem, it is very likely that the nearby solutions surrounding the existing solution will be more promising. One of the primary advantages of our method is that it can be used along with other popular methods

available and improve the quality of the existing solutions. In clustering problems, it is an added advantage to perform refinement of the final clusters obtained. Most of the focus in the literature was on new methods for initialization or new clustering techniques which often do not take advantage of the existing results and completely start the clustering procedure “*from scratch*”. Though shown only for the case of multivariate Gaussian mixtures, our technique can be effectively applied to any parametric finite mixture model.

Table 3.4: Number of iterations taken for the convergence of the best solution.

| Dataset | Avg. no. of iterations | No. of iterations for the best solution |
|-----------------|------------------------|---|
| Spherical | 126 | 73 |
| Elliptical | 174 | 86 |
| Full covariance | 292 | 173 |

Table 3.4 summarizes the average number of iterations taken by the EM algorithm for the convergence to the local optimal solution. We can see that the most promising solution produced by our TRUST-TECH methodology converges much faster. In other words, our method can effectively take advantage of the fact that the convergence of the EM algorithm is much faster for high quality solutions. We exploit this inherent property of the EM algorithm to improve the efficiency of our algorithm. Hence, for obtaining the Tier-1 solutions using our algorithm, the threshold for the number of iterations can be significantly lowered.

APPENDIX-A: Proof of Theorem 3.2.1

Proof. First, we will show that every bounded trajectory will converge to one of the equilibrium points. Second, we will show that every trajectory is bounded [31].

1. Let $\Phi(x, t)$ denote the bounded trajectory starting at x . Computing the time derivative along the trajectory, we get

$$\frac{d}{dt}f(\Phi(x, t)) = -(\nabla f(\Phi(x, t)))^T(\nabla f(\Phi(x, t))) \leq 0$$

Also, we know that $\frac{d}{dt}f(\Phi(x, t)) = 0$ if, and only if, $x \in E$. Hence, $f(x)$ is a Lyapunov function of the gradient system (3.18) and the ω -limit point of any bounded trajectory consists of equilibrium points only, i.e. any bounded trajectory will approach one of the equilibrium point.

2. Following the proof of proposition 1 presented in [31], we can show that every trajectory $\Phi(x, t)$ is bounded. However, we will have to show that the magnitude of the gradient of the log-likelihood function for the Gaussian mixture model is bounded on the entire domain of the parameter space.

$$\log p(\mathcal{Y}|\Theta) = - \sum_{j=1}^n \log \sum_{i=1}^k \alpha_i p(y^{(j)}|\theta_i) \quad (3.20)$$

Now, the domain of the parameter space is given as follows:

$-\infty < \mu_i < \infty$, Σ_i is positive definite and $0 \leq \alpha_i \leq 1$ where $\sum_{i=1}^k \alpha_i = 1$.

First, let us focus on α because it is a constrained variable.

Derivative with α :

$$\frac{\partial f}{\partial \alpha_r} = \sum_{j=1}^n \left[\frac{p(y^{(j)}|\theta_r)}{\sum_{i=1}^k \alpha_i p(y^{(j)}|\theta_i)} \right] \quad (3.21)$$

As $\alpha \rightarrow 1$, we have

$$\frac{\partial f}{\partial \alpha_r} = \sum_{j=1}^n \left[\frac{p(y^{(j)}|\theta_r)}{1 \cdot p(y^{(j)}|\theta_i)} \right] = n < \infty$$

As $\alpha \rightarrow 0$, we have

$$\frac{\partial f}{\partial \alpha_r} = \sum_{j=1}^n \left[\frac{p(y^{(j)}|\theta_r)}{\sum_{i=1, i \neq r}^k \alpha_i p(y^{(j)}|\theta_i)} \right] < \infty$$

Hence, the derivatives with respect to α are bounded.

Derivative with μ :

$$\frac{\partial f}{\partial \mu_r} = \sum_{j=1}^n \left[\frac{\alpha_r \frac{1}{\sqrt{(2\pi)\sigma_r}} e^{-\frac{(x^{(j)}-\mu_r)^2}{2\sigma_r^2}} \cdot \frac{1}{\sigma_r^2} (x^{(j)} - \mu_r)}{\sum_{i=1}^k \alpha_i p(y^{(j)}|\theta_i)} \right] \quad (3.22)$$

This is obviously bounded for and $\mu \in \mathfrak{R}$.

Derivative with σ :

$$\frac{\partial f}{\partial \sigma_r} = \sum_{j=1}^n \left[\frac{\frac{1}{\sigma_r} e^{-\frac{(x^{(j)}-\mu_r)^2}{2\sigma_r^2}} \cdot \frac{(x^{(j)}-\mu_r)^2}{\sigma_r^3} - \frac{1}{\sigma_r^2} e^{-\frac{(x^{(j)}-\mu_r)^2}{2\sigma_r^2}}}{\sum_{i=1}^k \alpha_i p(y^{(j)}|\theta_i)} \right] \quad (3.23)$$

As $\sigma_r \rightarrow 0$ the exponential factor goes to zero faster than $\frac{1}{\sigma_r}$ goes to infinity.

Hence, it is bounded. So, the gradient of the log-likelihood function is bounded in the entire domain of the parameter space.

Chapter 4

Motif Refinement using Neighborhood

Profile Search

As a case study of the algorithm developed in the previous chapter, we describe its application to the motif finding problem in bioinformatics. The main goal of the motif finding problem is to detect novel, over-represented unknown signals in a set of sequences (e.g. transcription factor binding sites in a genome). The most widely used algorithms for finding motifs obtain a generative probabilistic representation of these over-represented signals and try to discover profiles that maximize the information content score. Although these profiles form a very powerful representation of the signals, the major difficulty arises from the fact that the best motif corresponds to the global maximum of a non-convex continuous function. Popular algorithms like Expectation Maximization (EM) and Gibbs sampling tend to be very sensitive to the initial guesses and are known to converge to the nearest local maximum very quickly. In order to improve the quality of the results, EM is used with multiple random starts or any other powerful stochastic global methods that might yield promising initial guesses (like projection algorithms). Global methods usually do not give initial guesses in the convergence region of the best local maximum but rather give some point that is in a promising region. In this chapter, we apply the TRUST-TECH based Expectation Maximization (TT-EM) algorithm proposed in the previous chapter to this motif finding problem. It has the capability to search for alignments corresponding to Tier-1 local optimal solutions in the profile space. This search is performed in a systematic manner to explore the multiple local optimal solutions. This effective search is achieved by transforming the original optimization

problem into a dynamical system with certain properties and obtaining more useful information about the nonlinear likelihood surface via the dynamical and topological properties of the dynamic system.

```

GAATTCATACCAGATCACCGGATTCCCGACTCCAAATGTGTCCCCCTCACAC
TCCC CGGATTACCGTCTTCTGCTCTTAGACCACTCTACCCTATTCCCCACACT
CACCGGAGCCAAAGCCGCGGCCCTTCCGTTCCGATTACCGAAAAGACCCCA
CCCGTAGGTGGCAAGCTAGCTTAAGTAACGCCACTTCGATTAACCGAGGAAA
AATACATAACTGACCTATTATCGAGTTCAGATCAAGGTCAGGAACAAAGAA
ACA CCGATTACCGTAACCGTAAGATARTGGTATCGATACGTAGACAGTTTA

```

Figure 4.1: Synthetic DNA sequences containing some instance of the pattern ‘CCGATTACCGA’ with a maximum number of 2 mutations. The motifs in each sequence are highlighted in the box. We have a (11,2) motif where 11 is the length of the motif and 2 is the number of mutations allowed.

Recent developments in DNA sequencing have allowed biologists to obtain complete genomes for several species. However, knowledge of the sequence does not imply the understanding of how genes interact and regulate one another within the genome. Many transcription factor binding sites are highly conserved throughout the sequences and the discovery of the location of such binding sites plays an important role in understanding gene interaction and gene regulation. Although there are several variations of the motif finding algorithms, the problem studied in this chapter is defined as follows: without any previous knowledge of the consensus pattern, discover all the occurrences of the motifs and then recover a pattern for which all of these instances are within a given number of mutations (or substitutions). Despite the significant amount of literature available on the motif finding problem, many do not exploit the probabilistic models used for motif refinement [90, 4]. In this chapter, we consider a precise version of the motif discovery problem in computational biology as discussed in [20, 114]. The planted (l,d) motif problem [114] considered here is described as follows: Suppose there is a fixed but unknown nucleotide sequence

M (the *motif*) of length l . The problem is to determine M , given t sequences with t_i being the length of the i^{th} sequence and each one containing a planted variant of M . More precisely, each such planted variant is a substring that is M with exactly d point substitutions (see Fig. 4.1). More details about the complexity of the motif finding problem is given in [113]. A detailed assessment of different motif finding algorithms was published recently in [137].

4.1 Relevant Background

Existing approaches used to solve the motif finding problem can be classified into two main categories [51]. The first group of algorithms utilizes a generative probabilistic representation of the nucleotide positions to discover a consensus DNA pattern that maximizes the information content score. In this approach, the original problem of finding the best consensus pattern is formulated as finding the global maximum of a continuous non-convex function. The main advantage of this approach is that the generated profiles are highly representative of the signals being determined [46]. The disadvantage, however, is that the determination of the “best” motif cannot be guaranteed and is often a very difficult problem since finding global maximum of any continuous non-convex function is a challenging task. Current algorithms converge to the nearest local optimum instead of the global solution. Gibbs sampling [90], MEME [4], greedy CONSENSUS algorithm [75] and HMM based methods [48] belong to this category.

The second group uses patterns with ‘mismatch representation’ which defines a signal to be a consensus pattern and allows up to a certain number of mismatches to occur in each instance of the pattern. The goal of these algorithms is to recover the consensus pattern with the highest number of instances. These methods view

the representation of the signals as discrete and the main advantage of these algorithms is that they can guarantee that the highest scoring pattern will be the global optimum for any scoring function. The disadvantage, however, is that consensus patterns are not as expressive of the DNA signal as profile representations. Recent approaches within this framework include Projection methods [20, 119], string based methods [114], Pattern-Branching [117], MULTIPROFILER [85] and other branch and bound approaches [52, 51].

A hybrid approach could potentially combine the expressiveness of the profile representation with convergence guarantees of the consensus pattern. An example of a hybrid approach is the Random Projection [20] algorithm followed by EM algorithm [4]. It uses a global solver to obtain promising alignments in the discrete pattern space followed by further local solver refinements in continuous space [7, 131]. Currently, only few algorithms take advantage of a combined discrete and continuous space search [20, 51, 119]. We consider the profile representation of the motif and a new hybrid algorithm is developed to escape out of the local maxima of the likelihood surface. Some motivations to develop the new hybrid algorithm are :

- A motif refinement stage is vital and popularly used by many pattern based algorithms (like PROJECTION, MITRA etc) which try to find optimal motifs.
- The traditional EM algorithm used in the context of motif finding converges very quickly to the nearest local optimal solution (within 5-8 iterations) [17].
- There are many other promising local optimal solutions in the close vicinity of the profiles obtained from the global methods.

In spite of the importance placed on obtaining a global optimal solution in the context of motif finding, little work has been done in the direction of finding such

solutions [142]. There are several proposed methods to escape out of the local optimal solution to find better solutions in machine learning [49] and optimization [25] related problems. Most of them are stochastic in nature and usually rely on perturbing either the data or the hypothesis. These stochastic perturbation algorithms are inefficient because they will sometimes miss a neighborhood solution or obtain an already existing solution. To avoid these problems, we will apply TRUST-TECH-EM algorithm that can systematically explore multiple local maxima in a tier-by-tier manner. Our method is primarily based on transforming the log-likelihood function into its corresponding dynamical system and obtaining multiple local optimal solutions. The underlying theoretical details of TRUST-TECH is described in [31, 91].

4.2 Preliminaries

We will first describe our problem formulation and the details of the EM algorithm in the context of motif finding problem. We will then describe some details of the dynamical system of the log-likelihood function which enables us to search for the nearby local optimal solutions.

4.2.1 Problem Formulation

In this section, we transform the the problem of finding the best possible motif into a problem of finding the global maximum of a highly nonlinear log-likelihood scoring function obtained from its profile representation. The log-likelihood surface is made of $3l$ variables which are treated as the unknown parameters that are to be estimated. We will now describe these parameters $(Q_{k,j})$ and then represent the scoring function in terms of these parameters.

Let t be the total number of sequences and n be the average length of the

sequences. Let $S = \{S_1, S_2 \dots S_t\}$ be the set of t sequences. Let P be a single alignment containing the set of segments $\{P_1, P_2, \dots, P_l\}$. l is the length of the consensus pattern. For further discussion, we use the following variables

- $i = 1 \dots t$ -- for t sequences
- $k = 1 \dots l$ -- for positions within an l -mer
- $j \in \{A, T, G, C\}$ -- for each nucleotide

Table 4.1: A count of nucleotides A, T, G, C at each position $k = 1..l$ in all the sequences of the data set. $k = 0$ denotes the background count.

| j | $k = 0$ | $k = 1$ | $k = 2$ | $k = 3$ | $k = 4$ | ... | $k = l$ |
|-----|-----------|-----------|-----------|-----------|-----------|-----|-----------|
| A | $C_{0,1}$ | $C_{1,1}$ | $C_{2,1}$ | $C_{3,1}$ | $C_{4,1}$ | ... | $C_{l,1}$ |
| T | $C_{0,2}$ | $C_{1,2}$ | $C_{2,2}$ | $C_{3,2}$ | $C_{4,2}$ | ... | $C_{l,2}$ |
| G | $C_{0,3}$ | $C_{1,3}$ | $C_{2,3}$ | $C_{3,3}$ | $C_{4,3}$ | ... | $C_{l,3}$ |
| C | $C_{0,4}$ | $C_{1,4}$ | $C_{2,4}$ | $C_{3,4}$ | $C_{4,4}$ | ... | $C_{l,4}$ |

The count matrix can be constructed from the given alignments as shown in Table 4.1. We define $C_{0,j}$ to be the non-position specific background count of each nucleotide in all of the sequences where $j \in \{A, T, C, G\}$ is the running total of nucleotides occurring in each of the l positions. Similarly, $C_{k,j}$ is the count of each nucleotide in the k^{th} position (of the $l - mer$) in all the segments in P .

$$Q_{0,j} = \frac{C_{0,j}}{\sum_{J \in \{A, T, G, C\}} C_{0,J}} \tag{4.1}$$

$$Q_{k,j} = \frac{C_{k,j} + b_j}{t + \sum_{J \in \{A, T, G, C\}} b_J} \tag{4.2}$$

Eq. (4.1) shows the background frequency of each nucleotide where b_j is known as the Laplacian or Bayesian correction and is equal to $d * Q_{0,j}$ where d is some

constant usually set to unity. Eq. (4.2) gives the weight assigned to the type of nucleotide at the k^{th} position of the motif.

A Position Specific Scoring Matrix (PSSM) can be constructed from one set of instances in a given set of t sequences. From Eqs. (4.1) and (4.2), it is obvious that the following relationship holds:

$$\sum_{j \in \{A, T, G, C\}} Q_{k,j} = 1 \quad \forall k = 0, 1, 2, \dots, l \quad (4.3)$$

For a given k value in Eq. (4.3), each Q can be represented in terms of the other three variables. Since the length of the motif is l , the final objective function (i.e. the information content score) would contain $3l$ independent variables¹.

To obtain the score, every possible $l - mer$ in each of the t sequences must be examined. This is done so by multiplying the respective $Q_{i,j}/Q_{0,j}$ dictated by the nucleotides and their respective positions within the $l - mer$. Only the highest scoring $l - mer$ in each sequence is noted and kept as part of the alignment. The total score is the sum of all the best scores in each sequence.

$$A(Q) = \sum_{i=1}^t \log(A)_i = \sum_{i=1}^t \log \left(\prod_{k=1}^l \frac{Q_{k,j}}{Q_b} \right)_i = \sum_{i=1}^t \sum_{k=1}^l \log(Q'_{k,j})_i \quad (4.4)$$

$Q'_{k,j}$ is the ratio of the nucleotide probability to the corresponding background probability, i.e. $Q_{k,j}/Q_b$. $\log(A)_i$ is the score at each individual i^{th} sequence where t is the total number of sequences. In Eq. (4.4), we see that A is composed of the product of the weights for each individual position k . We consider this to be the Information Content (IC) score which we would like to maximize. $A(Q)$ is the non-convex $3l$ dimensional continuous function for which the global maximum

¹Although, there are $4l$ variables in total, because of the constraints obtained from Eq. (4.3), the parameter space will contain only $3l$ independent variables. Thus, the constraints help in reducing the dimensionality of the search problem.

corresponds to the best possible motif in the dataset. In summary, we transform the problem of finding the optimal motif into a problem of finding the global maximum of a non-convex continuous $3l$ dimensional function.

4.2.2 Dynamical Systems for the Scoring Function

In order to present our algorithm, we define the dynamical system corresponding to the log-likelihood function and the PSSM. The key contribution here is the development of this nonlinear dynamical system which will enable us to realize the geometric and dynamic nature of the likelihood surface. We construct the following *gradient system* in order to locate critical points of the objective function (4.4):

$$\dot{Q}(t) = -\nabla A(Q) \quad (4.5)$$

One can realize that this transformation preserves all of the critical points [31]. Now, we will describe the construction of the gradient system and the Hessian in detail. In order to reduce the dominance of one variable over the other, the values of each of the nucleotides that belong to the consensus pattern at the position k will be represented in terms of the other three nucleotides in that particular column. Let P_{ik} denote the k^{th} position in the segment P_i . This will also minimize the dominance of the eigenvector directions when the Hessian is obtained. The variables in the scoring function are transformed into new variables described in Table 4.2.

$$A(Q) = \sum_{i=1}^t \sum_{k=1}^l \log f_{ik}(w_{3k-2}, w_{3k-1}, w_{3k})_i \quad (4.6)$$

where f_{ik} can take the values $\{w_{3k-2}, w_{3k-1}, w_{3k}, 1 - (w_{3k-2} + w_{3k-1} + w_{3k})\}$ depending on the P_{ik} value.

Table 4.2: A count of nucleotides $j \in \{A, T, G, C\}$ at each position $k = 1..l$ in all the sequences of the data set. C_k is the k^{th} nucleotide of the consensus pattern which represents the nucleotide with the highest value in that column. Let the consensus pattern be GACT...G and b_j be the background.

| j | $k = b$ | $k = 1$ | $k = 2$ | $K = 3$ | $k = 4$ | ... | $k = l$ |
|-----|---------|---------|---------|---------|----------|-----|------------|
| A | b_A | w_1 | C_2 | w_7 | w_{10} | ... | w_{3l-2} |
| T | b_T | w_2 | w_4 | w_8 | C_4 | ... | w_{3l-1} |
| G | b_G | C_1 | w_5 | w_9 | w_{11} | ... | C_l |
| C | b_C | w_3 | w_6 | C_3 | w_{12} | ... | w_{3l} |

The first derivative of the scoring function is a one dimensional vector with $3l$ elements.

$$\nabla A = \left[\frac{\partial A}{\partial w_1} \quad \frac{\partial A}{\partial w_2} \quad \frac{\partial A}{\partial w_3} \quad \cdots \quad \frac{\partial A}{\partial w_{3l}} \right]^T \quad (4.7)$$

and each partial derivative is given by

$$\frac{\partial A}{\partial w_p} = \sum_{i=1}^t \frac{\frac{\partial f_{ip}}{\partial w_p}}{f_{ik}(w_{3k-2}, w_{3k-1}, w_{3k})} \quad (4.8)$$

$$\forall p = 1, 2 \dots 3l \text{ and } k = \text{round}(p/3) + 1$$

The Hessian $\nabla^2 A$ is a block diagonal matrix of block size 3×3 . For a given sequence, the entries of the 3×3 block will be the same if that nucleotide belongs to the consensus pattern (C_k). This nonlinear transformation will preserve all the critical points on the likelihood surface. The theoretical details of the proposed method and their advantages are published in [31]. If we can identify all the saddle points on the stability boundary of a given local maximum, then we will be able to find all the tier-1 local maxima. However, finding all of the saddle points is computationally intractable and hence we have adopted a heuristic by generating the eigenvector directions of the PSSM at the local maximum.

4.3 Neighborhood Profile Search

We will now describe our novel neighborhood search framework which is applied in the profile space. Our framework consists of the following three stages:

- *Global stage* in which the promising solutions in the entire search space are obtained.
- *Refinement stage* (or *local stage*) where a local method is applied to the solutions obtained in the previous stage in order to refine the profiles.
- *Neighborhood-search stage* where the exit points are computed and the Tier-1 and Tier-2 solutions are explored systematically.

In the global stage, a branch and bound search is performed on the entire dataset. All of the profiles that do not meet a certain threshold (in terms of a given scoring function) are eliminated in this stage. Some promising initial alignments are obtained by applying these methods (like projection methods) on the entire dataset. Most promising set of alignments are considered for further processing. The promising patterns obtained are transformed into profiles and local improvements are made to these profiles in the refinement stage. The consensus pattern is obtained from each nucleotide that corresponds to the largest value in each column of the PSSM. The $3l$ variables chosen are the nucleotides that correspond to those that are not present in the consensus pattern. Because of the probability constraints discussed in the previous section, the largest weight can be represented in terms of the other three variables.

To solve Eq. (4.4), current algorithms begin at random initial alignment positions and attempt to converge to an alignment of $l - mers$ in all of the sequences that maximize the objective function. In other words, the $l - mer$ whose $\log(A)_i$ is

the highest (with a given PSSM) is noted in every sequence as part of the current alignment. During the maximization of $A(Q)$ function, the probability weight matrix and hence the corresponding alignments of $l - mers$ are updated. This occurs iteratively until the PSSM converges to the local optimal solution. The consensus pattern is obtained from the nucleotide with the largest weight in each position (column) of the PSSM. This converged PSSM and the set of alignments correspond to a local optimal solution. The neighborhood-search stage where the neighborhood of the original solution is explored in a systematic manner is shown below:

Input: Local Maximum (A).

Output: Best Local Maximum in the neighborhood region.

Algorithm:

Step 1: Construct the PSSM for the alignments corresponding to the local maximum (A) using Eqs.(4.1) and (4.2).

Step 2: Calculate the eigenvectors of the Hessian matrix for this PSSM.

Step 3: Find exit points (e_{1i}) on the practical stability boundary along each eigenvector direction.

Step 4: For each exit point, the corresponding Tier-1 local maxima (a_{1i}) are obtained by applying the EM algorithm after the ascent step.

Step 5: Repeat this process for promising Tier-1 solutions to obtain Tier-2 (a_{2j}) local maxima.

Step 6: Return the solution that gives the maximum information content score among $\{A, a_{1i}, a_{2j}\}$.

To escape out of this local optimal solution, our approach requires the computation of a Hessian matrix (i.e. the matrix of second derivatives) of dimension $(3l)^2$ and the $3l$ eigenvectors of the Hessian. these directions were chosen as a general

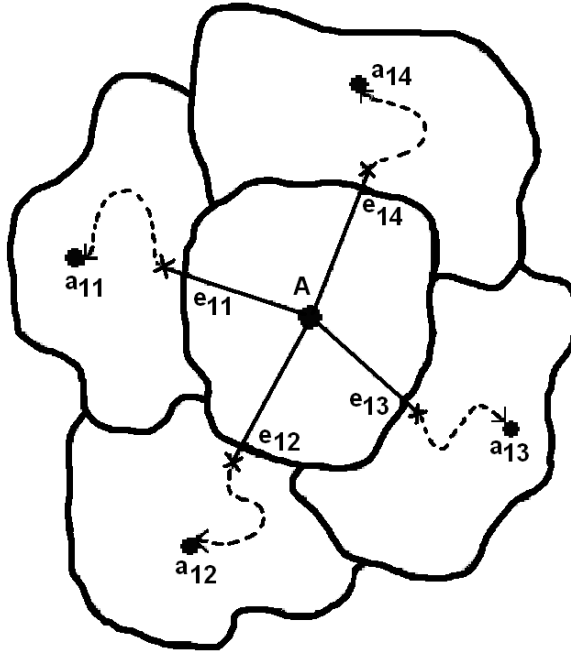


Figure 4.2: Diagram illustrates the TRUST-TECH method of escaping from the original solution (A) to the neighborhood local optimal solutions (a_{1i}) through the corresponding exit points (e_{1i}). The dotted lines indicate the local convergence of the EM algorithm.

heuristic and are not problem dependent. Depending on the dataset that is being worked on, one can obtain even more promising directions. The main reasons for choosing the eigenvectors of the Hessian as search directions are:

- Computing the eigenvectors of the Hessian is related to finding the directions with extreme values of the second derivatives, i.e., directions of extreme normal-to-isosurface change.
- The eigenvectors of the Hessian will form the basis vectors for the search directions. Any other search direction can be obtained by a linear combination of these directions.
- This will make our algorithm deterministic since the eigenvector directions are always unique.

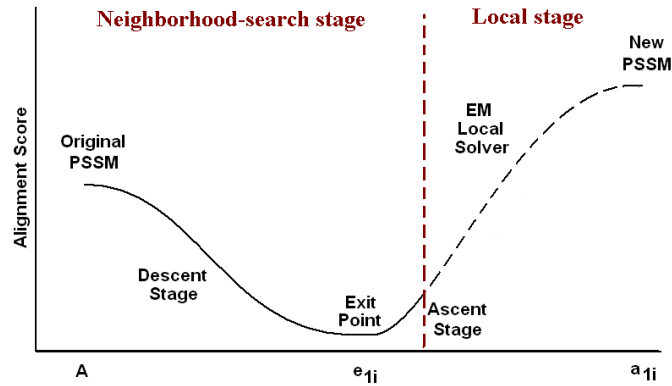


Figure 4.3: A summary of escaping out of the local optimum to the neighborhood local optimum. Observe the corresponding trend of the alignment score at each step.

The value of the objective function is evaluated along these eigenvector directions with some small step size increments. Since the starting position is a local optimal solution, one will see a steady decline in the function value during the initial steps; we call this the *descent stage*. Since the Hessian is obtained only once during the entire procedure, it is more efficient compared to Newton's method where an approximate Hessian is obtained for every iteration. After a certain number of evaluations, there may be an increase in the value indicating that the stability boundary is reached. The point along this direction intersecting the stability boundary is called the *exit point*. Once the exit point has been reached, few more evaluations are made in the direction of the same eigenvector to improve the chances of reaching a new convergence region. This procedure is clearly shown in Fig 4.3. Applying the local method directly from the exit point may give the original local maximum. The ascent stage is used to ensure that the new guess is in a different convergence zone. Hence, given the best local maximum obtained using any current local methods, this framework allows us to systematically escape out of the local maximum to explore surrounding local maxima. The complete algorithm is shown below :

Input: The DNA sequences, length of the motif(l), Maximum No. of Mutations(d)

Output: Motif (s)

Algorithm:

Step 1: Given the sequences, apply Random Projection algorithm to obtain different set of alignments.

Step 2: Choose the promising buckets and apply EM algorithm to refine these alignments.

Step 3: Apply the TRUST-TECH method to obtain nearby promising local optimal solutions.

Step 4: Report the consensus pattern that corresponds to the best alignments and their corresponding PSSM.

The new framework can be treated as a hybrid approach between global and local methods. It differs from traditional local methods by the ability to explore multiple local solutions in the neighborhood region in a systematic and effective manner. It differs from global methods by working completely in the profile space and searching a subspace efficiently in a deterministic manner. However, the main difference of this work compared to the algorithm presented in the previous chapter is that the global method is performed in discrete space and the local method is performed in the continuous space. In other words, the both the global and local method do not optimize the same function. In such cases, it is even more important to search for neighborhood local optimal solutions in the continuous space.

4.4 Implementation Details

Our program was implemented on Red Hat Linux version 9 and runs on a Pentium IV 2.8 GHz machine. The core algorithm that we have implemented is *TT_EM* described in Algorithm 3. *TT_EM* obtains the initial alignments and the original data sequences along with the length of the motif. This procedure constructs the PSSM, performs EM refinement, and then computes the Tier-1 and Tier-2 solutions by calling the procedure *Next_Tier*. The eigenvectors of the Hessian were computed using the source code obtained from [116]. *Next_Tier* takes a PSSM as an input and computes an array of PSSMs corresponding to the next tier local maxima using the TRUST-TECH methodology.

Algorithm 3 Motif *TT_EM*(*init_aligns*, *seqs*, *l*)

```
PSSM = Construct_PSSM(init_aligns)
New_PSSM = Apply_EM(PSSM, seqs)
TIER1 = Next_Tier(seqs, New_PSSM, l)
for i = 1 to 3l do
  if TIER1[i] <> zeros(4l) then
    TIER2[i][ ] = Next_Tier(seqs, TIER1[i], l)
  end if
end for
Return best(PSSM, TIER1, TIER2)
```

Given a set of initial alignments, Algorithm 3 will find the best possible motif in the profile space in a tier-by-tier manner. For implementation considerations, we have shown only for two tiers. Initially, a PSSM is computed using *construct_PSSM* from the given alignments. The procedure *Apply_EM* will return a new PSSM that corresponds to the alignments obtained after the EM algorithm has been applied to the initial PSSM. The details of the procedure *Next_Tier* are given in Algorithm 4. From a given local solution (or PSSM), *Next_Tier* will compute all the $3l$ new PSSMs corresponding to the tier-1 local optimal solutions. The second tier patterns

are obtained by calling the *Next_Tier* from the first tier solutions ². Finally, the pattern with the highest score amongst all the PSSMs is returned.

Algorithm 4 $PSSMs[] \text{ Next_Tier}(seqs, PSSM, l)$

```

Score = eval(PSSM)
Hess = Construct_Hessian(PSSM)
Eig[ ] = Compute_EigVec(Hess)
MAX_Iter = 100
for k = 1 to 3l do
    PSSMs[k] = PSSM    Count = 0
    Old_Score = Score    ep_reached = FALSE
    while (! ep_reached) && (Count < MAX_Iter) do
        PSSMs[k] = update(PSSMs[k], Eig[k], step)
        Count = Count + 1
        New_Score = eval(PSSMs[k])
        if (New_Score > Old_Score) then
            ep_reached = TRUE
        end if
        Old_Score = New_Score
    end while
    if count < MAX_Iter then
        PSSMs[k] = update(PSSMs[k], Eig[k], ASC)
        PSSMs[k] = Apply_EM(PSSMs[k], Seqs)
    else
        PSSMs[k] = zeros(4l)
    end if
end for
Return PSSMs[ ]

```

The procedure *Next_Tier* takes a PSSM, applies the TRUST-TECH method and computes an array of PSSMs that corresponds to the next tier local optimal solutions. The procedure *eval* evaluates the scoring function for the PSSM using (4.4). The procedures *Construct_Hessian* and *Compute_EigVec* compute the Hessian matrix and the eigenvectors respectively. *MAX_iter* indicates the maximum number of uphill evaluations that are required along each of the eigenvector direc-

²New PSSMs might not be obtained for certain search directions. In those cases, a zero vector of length $4l$ is returned. Only those new PSSMs which do not have this value will be used for any further processing.

tions. The neighborhood PSSMs will be stored in an array variable *PSSMs* (this is a vector). The original PSSM is updated with a small step until an exit point is reached or the number of iterations exceeds the *MAX_Iter* value. If the exit point is reached along a particular direction, few more function evaluations are made to ensure that the PSSM has exited the original convergence region and has entered a new one. The EM algorithm is then used during this ascent stage to obtain a new PSSM ³.

The initial alignments are converted into the profile space and a PSSM is constructed. The PSSM is updated (using the EM algorithm) until the alignments converge to a local optimal solution. The TRUST-TECH methodology is then employed to escape out of this local optimal solution to compute nearby first tier local optimal solutions. This process is then repeated on promising first tier solutions to obtain second tier solutions. As shown in Fig. 4.2, from the original local optimal solution, various exit points and their corresponding new local optimal solutions are computed along each eigenvector direction. Sometimes, two directions may yield the same local optimal solution. This can be avoided by computing the saddle point corresponding to the exit point on the stability boundary [122]. There can be many exit points, but there will only be a unique saddle point corresponding to the new local minimum. For computational efficiency, the TRUST-TECH approach is only applied to promising initial alignments (i.e. random starts with higher Information Content score). Therefore, a threshold $A(Q)$ score is determined by the average of the three best first tier scores after 10-15 random starts; any current and first tier solution with scores greater than the threshold is considered for further exploration.

³For completeness, the entire algorithm has been shown in this section. However, during the implementation, several heuristics have been applied to reduce the running time of the algorithm. For example, if the first tier solution is not very promising, it will not be considered for obtaining the corresponding second tier solutions.

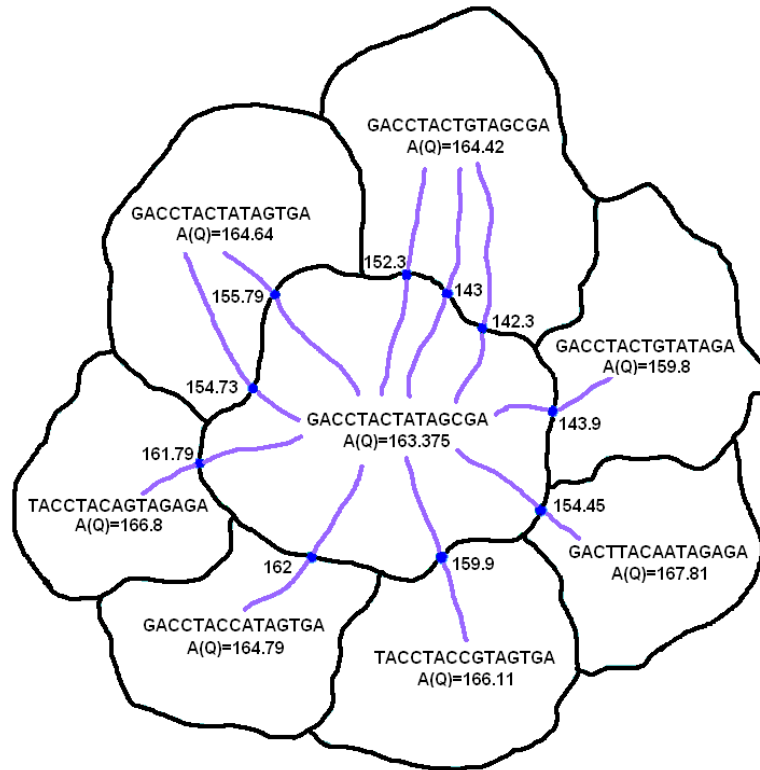


Figure 4.4: 2-D illustration of first tier improvements in a $3l$ dimensional objective function. The original local maximum has a score of 163.375. The various Tier-1 solutions are plotted and the one with highest score (167.81) is chosen.

Additional random starts are carried out in order to aggregate at least ten first tier solutions. The TRUST-TECH method is repeated on all first tier solutions above a certain threshold to obtain second-tier local optimal solutions.

4.5 Experimental Results

Experiments were performed on both synthetic data and real data. Two different methods were used in the global stage: random start and random projection. The main purpose of our work is not to demonstrate that our algorithm can outperform the existing motif finding algorithms. Rather, the main work here focuses on improving the results that are obtained from other efficient algorithms. We have chosen to demonstrate the performance of our algorithm on the results obtained from

the random projection method which is a powerful global method that has outperformed other traditional motif finding approaches like MEME, Gibbs sampling, WINNOWER, SP-STAR, etc. [20]. Since the comparison was already published, we mainly focus on the performance improvements of our algorithm as compared to the random projection algorithm. For the random start experiment, a total of N random numbers between 1 and $(t - l + 1)$ corresponding to initial set of alignments are generated. We then proceeded to evaluate our TRUST-TECH methodology from these alignments.

Fig. 4.4 shows the Tier-1 solutions obtained from a given consensus pattern. Since the exit points are being used instead of saddle points, our method might sometimes find the same local optimal solution obtained before. As seen from the figure, the Tier-1 solutions can differ from the original pattern by more than just one nucleotide position. Also, the function value at the exit points is much higher than the original value.

4.5.1 Synthetic Datasets

The synthetic datasets were generated by implanting some motif instances into $t = 20$ sequences each of length 600. Let m correspond to one full random projection + EM cycle. We have set $m = 1$ to demonstrate the efficiency of our approach. We compared the performance coefficient (PC) which gives a measure of the average performance of our implementation compared to that of Random Projection. The PC is given by :

$$PC = \frac{|K \cap P|}{|K \cup P|} \quad (4.9)$$

where K is the set of the residue positions of the planted motif instances, and P

is the corresponding set of positions predicted by the algorithm. Table 4.4 gives an overview of the performance of our method compared to the random projection algorithm on the (l, d) motif problem for different l and d values.

Our results show that by branching out and discovering multiple local optimal solutions, higher m values are not needed. A higher m value corresponds to more computational time because projecting the l -mers into k -sized buckets is a time consuming task. Using our approach, we can replace the need for randomly projecting l -mers repeatedly in an effort to converge to a global optimum by deterministically and systematically searching the solution space modeled by our dynamical system and improving the quality of the existing solutions. The improvements of our algorithm are clearly shown in Table 4.3. We can see that, for higher length motifs, the improvements are more significant.

As opposed to stochastic processes like mutations in genetic algorithms, our approach eliminates the stochastic nature and obtains the nearby local optimal solutions systematically. Fig. 4.5 shows the performance of the TRUST-TECH approach on synthetic data for different (l, d) motifs. The average scores of the ten best solutions obtained from random starts and their corresponding improvements in Tier-1 and Tier-2 are reported. One can see that the improvements become more prominent as the length of the motif is increased. Table 4.3 shows the best and worst of these top ten random starts along with the consensus pattern and the alignment scores.

With a few modifications, more experiments were conducted using the Random Projection method. The Random Projection method will eliminate non-promising regions in the search space and gives a number of promising sets of initial patterns. EM refinement is applied to only the promising initial patterns. Due to the robustness of the results, the TRUST-TECH method is employed only on the top five local

Table 4.3: The consensus patterns and their corresponding scores of the original local optimal solution obtained from multiple random starts on the synthetic data. The best first tier and second tier optimal patterns and their corresponding scores are also reported.

| (l,d) | Initial Pattern | Score | First Tier Pattern | Score | Second Tier Pattern | Score |
|--------|----------------------|-------|----------------------|-------|----------------------|-------|
| (11,2) | AACGGTCGCAG | 125.1 | CCCGGTCGCTG | 147.1 | CCCGGGAGCTG | 153.3 |
| (11,2) | ATACCAGTTAC | 145.7 | ATACCAGTTTC | 151.3 | ATACCAGGGTC | 153.6 |
| (13,3) | CTACGGTCGTCTT | 142.6 | CCACGGTTGTCTC | 157.8 | CCTCGGGTTTGTC | 158.7 |
| (13,3) | GACGCTAGGGGGT | 158.3 | GAGGCTGGGCAGT | 161.7 | GACCTTGGGTATT | 165.8 |
| (15,4) | CCGAAAAGAGTCCGA | 147.5 | CCGCAATGACTGGGT | 169.1 | CCGAAAGGACTGCGT | 176.2 |
| (15,4) | TGGGTGATGCCTATG | 164.6 | TGGGTGATGCCTATG | 166.7 | TGAGAGATGCCTATG | 170.4 |
| (17,5) | TTGTAGCAAAGGCTAAA | 143.3 | CAGTAGCAAAGACTACC | 173.3 | CAGTAGCAAAGACTTCC | 175.8 |
| (17,5) | ATCGCGAAAGGTTGTGG | 174.1 | ATCGCGAAAGGATGTGG | 176.7 | ATTGCGAAAGAATGTGG | 178.3 |
| (20,6) | CTGGTGATTGAGATCATCAT | 165.9 | CAGATGGTTGAGATCACCTT | 186.9 | CATTTAGCTGAGTTCACCTT | 194.9 |
| (20,6) | GGTCACTTAGTGGCGCCATG | 216.3 | GGTCACTTAGTGGCGCCATG | 218.8 | CGTCACTTAGTCGCGCCATG | 219.7 |

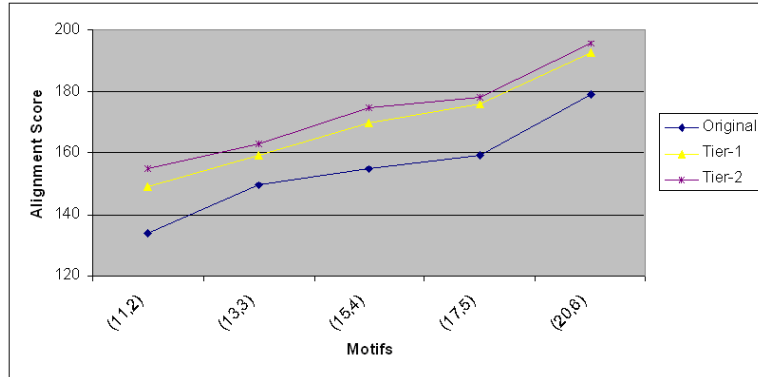


Figure 4.5: The average scores with the corresponding first tier and second tier improvements on synthetic data using the random starts with TRUST-TECH approach with different (l, d) motifs.

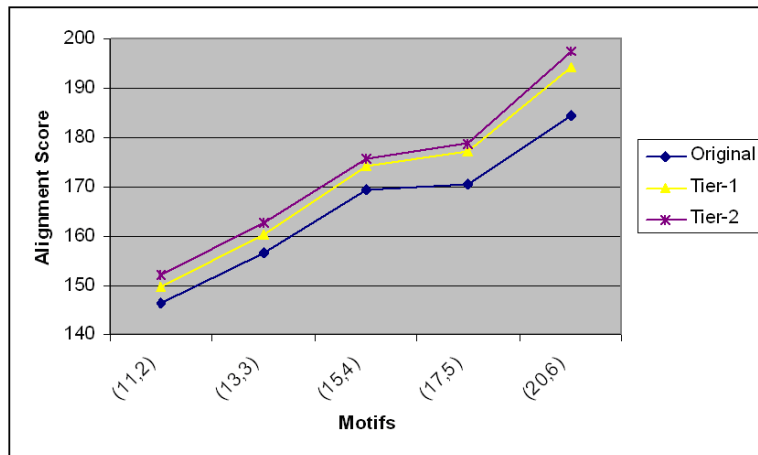


Figure 4.6: The average scores with the corresponding first tier and second tier improvements on synthetic data using the Random Projection with TRUST-TECH approach with different (l, d) motifs.

optima. The TRUST-TECH method is again repeated on the top scoring first tier solutions to arrive at the second tier solutions. Fig. 4.6 shows the average alignment scores of the best random projection alignments and their corresponding improvements in tier-1 and tier-2 are reported. In general, the improvement in the first tier solution is more significant than the improvements in the second tier solutions.

Table 4.4: The results of performance coefficient with $m = 1$ on synthetically generated sequences. The IC scores are not normalized and the perfect score is 20 since there are 20 sequences.

| Motif (1,d) | PC obtained using Random Projection | PC obtained using TRUST-TECH method |
|----------------|--|--|
| (11,2) | 20 | 20 |
| (15,4) | 14.875 | 17 |
| (20,6) | 12.667 | 18 |

4.5.2 Real Datasets

Table 4.5 shows the results of the TRUST-TECH methodology on real biological sequences. We have chosen $l = 20$ and $d = 2$. ‘ t ’ indicates the number of sequences in the real data. For the biological samples taken from [20, 117], the value m once again is the average number of random projection + EM cycles required to discover the motif. All other parameter values (like projection size $k=7$ and threshold $s=4$) are chosen to be the same as those used in the Random projection paper [20]. All of the motifs were recovered with $m = 1$ using the TRUST-TECH strategy. Without the TRUST-TECH strategy, the Random Projection algorithm needed *multiple cycles* ($m=8$ in some cases and $m=15$ in others) in order to retrieve the correct motif. This elucidates the fact that global methods can only be used to a certain extent and should be combined with refined local heuristics in order to obtain better efficiency. Since the random projection algorithm has outperformed other

prominent motif finding algorithms like SP-STAR, WINNOWER, Gibbs sampling etc., we did not repeat the same experiments that were conducted in [20]. Running one cycle of random projection + EM is much more expensive computationally. The main advantage of our strategy comes from the deterministic nature of our algorithm in refining motifs.

Table 4.5: Results of TRUST-TECH method on biological samples. The real motifs were obtained in all the six cases using the TRUST-TECH framework.

| Sequence | Size | t | Best (20,2) Motif | Reference Motif |
|-----------------|------|----|-----------------------------|------------------|
| E. coli CRP | 1890 | 18 | <u>TGTGAAATAGATCACATTTT</u> | TGTGANNNGNTCACA |
| preproinsulin | 7689 | 4 | <u>GGAAATTGCAGCCTCAGCCC</u> | CCTCAGCCC |
| DHFR | 800 | 4 | <u>CTGCAATTCGCGCCAAACT</u> | ATTCNNGCCA |
| metallothionein | 6823 | 4 | <u>CCCTCTGCGCCCGGACCGGT</u> | TGCRYCGG |
| c-fos | 3695 | 5 | <u>CCATATTAGGACATCTGCGT</u> | CCATATTAGAGACTCT |
| yeast ECB | 5000 | 5 | <u>GTATTTCCCGTTTAGGAAAA</u> | TTCCCNNTNAGGAAA |

The TRUST-TECH framework broadens the search region in order to obtain an improved solution which may potentially correspond to a better motif. In most of the profile based algorithms, EM is used to obtain the nearest local optimum from a given starting point. In our approach, we obtain promising results by computing multiple local optimal solutions in a tier-by-tier manner. We have shown on both real and synthetic data sets that beginning from the EM converged solution, the TRUST-TECH approach is capable of searching in the neighborhood regions for another solution with an improved information content score. This will often translate into finding a pattern with less hamming distance from the resulting alignments in each sequence. Our approach has demonstrated an improvement in the score on all datasets that it was tested on. One of the primary advantages of the TRUST-TECH methodology is that it can be used with different global and local methods. The main contribution of our work is to demonstrate the capability of this hybrid EM algorithm in the context of the motif finding problem. The TRUST-TECH approach can potentially use any global method and improve its results efficiently.

From these simulation results, we observe that motif refinement stage plays a vital role and can yield accurate results deterministically.

In the future, we would like to continue our work by combining other global methods available in the literature with existing local solvers like EM or GibbsDNA that work in continuous space. By following the example of [137], we may improve the chances of finding more promising patterns by combining our algorithm with different global and local methods.

Chapter 5

Component-wise Smoothing for Learning Mixture Models

5.1 Overview

The task of obtaining an optimal set of parameters to fit a mixture model to a given data can be formulated as a problem of finding the global maximum of its highly nonlinear likelihood surface. This chapter introduces a new smoothing algorithm for learning a mixture model from multivariate data. Our algorithm is based on the conventional Expectation Maximization (EM) approach applied to a smoothed likelihood surface. A family or hierarchy of smooth log-likelihood surfaces is constructed using convolution based approaches. In simple terms, our method first smoothens the likelihood function and then applies the EM algorithm to obtain a promising solution on the smooth surface. This solution is used as an initial guess for the EM algorithm applied to the next smooth surface in the hierarchy. This process is repeated until the original likelihood surface is solved. The smoothing process reduces the overall gradient of the surface and the number of local maxima. This effective optimization procedure eliminates extensive search in the non-promising regions of the parameter space. Results on benchmark datasets demonstrate significant improvements of the proposed algorithm compared to other approaches. Reduction in the number of local maxima has also been demonstrated empirically on several datasets.

As mentioned in Chapter 3, one of the key problems with the EM algorithm is that it is a ‘greedy’ method which is sensitive to initialization. The log-likelihood surface on which the EM algorithm is applied is very rugged with many local max-

ima. Because of its greedy nature, EM algorithm tends to get stuck at a local maximum that corresponds to erroneous set of parameters for the mixture components. Obtaining an improved likelihood function value not only provides better parameter estimates but also enhances the generalization capability of the given mixture models [99]. The fact that the local maxima are not uniformly distributed makes it important for us to develop algorithms that help in avoiding search in non-promising regions. More focus needs to be given for searching the promising subspace by obtaining promising initial estimates. This can be achieved by smoothing the surface and obtaining promising regions and then gradually trace back these solutions onto the original surface. In this work, we develop a hierarchical smoothing algorithm for the mixture modeling problem using convolution-based approach. We propose the following desired properties for smoothing algorithms :

- Preserve the presence of the global optimal solution.
- Enlarge the convergent region (with respect to the EM algorithm) of the global optimal solutions and other promising sub-optimal solutions.
- Reduce the number of local maximum or minimum.
- Smooth different regions of the search space differently.
- Avoid over smoothing which might make the surface too flat and cause convergence problems for the local solvers.

5.2 Relevant Background

Since EM is very sensitive to the initial set of parameters that it begins with, several methods are proposed in the literature to identify good initial points. All these

methods are mentioned in Chapter 3. In this section, we primarily focus on the literature relevant to smoothing algorithms.

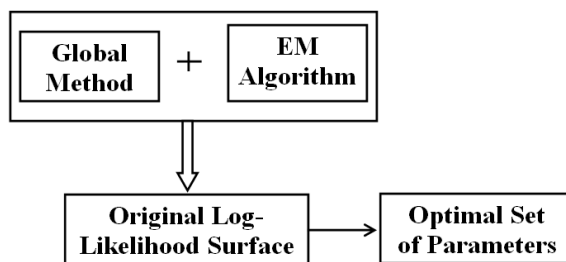
Different smoothing strategies have been successfully used in various applications for solving a diverse set of problems. Smoothing techniques are used to reduce irregularities or random fluctuations in time series data [134, 11]. In the field of natural language processing, smoothing techniques are also used for adjusting maximum likelihood estimate to produce more accurate probabilities for language models [27]. Convolution based smoothing approaches are predominantly used in the field of digital image processing for image enhancement by noise removal [15, 34]. Other variants of smoothing techniques include continuation methods [125, 45] which are used successfully in various applications. Different multi-level procedures other than smoothing and its variants are clearly illustrated in [136].

For optimization problems, smoothing procedure helps in reducing the ruggedness of the surface and helps the local methods to prevent the local minima problem. It was used for the structure prediction of molecular clusters [133]. This smoothing procedure obtains a hierarchy of smooth surfaces with a fewer and fewer local maxima. Promising initial points can be obtained by tracing back promising solutions at each level. This is an initialization procedure which has the capability to avoid searching non-promising regions. Obtaining the number of components using some model selection criterion [99] is not the primary focus of this work. In other words, our algorithm assumes that the number of components are known before hand. In summary, the main contributions are :

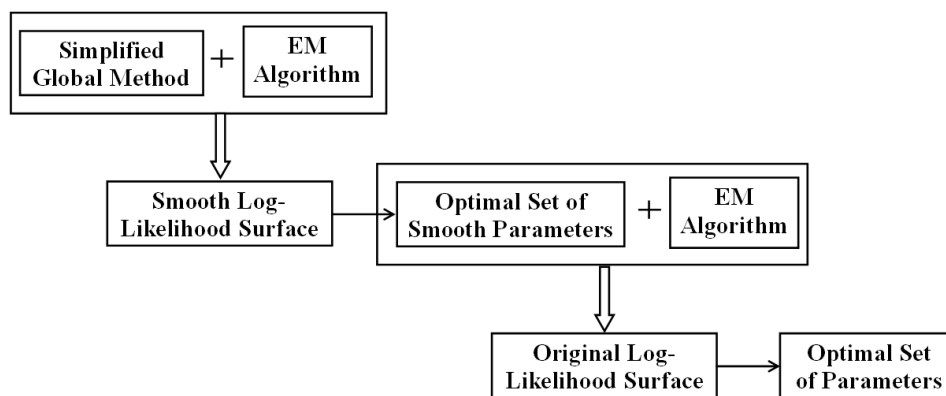
- Develop convolution-based smoothing algorithms for obtaining optimal set of parameters.
- Demonstrate that the density-based convolution on the entire dataset will

result smoothing the likelihood surface with respect to the parameters.

- Empirically show that the number of local maxima on the log-likelihood surface is reduced.
- Show that smoothing is effective in obtaining promising initial set of parameters



(a) Conventional method



(b) Smoothing approach

Figure 5.1: Block diagram of the traditional approach and the smoothing approach.

Fig. 5.1 compares the conventional approach with the smoothing approach. In the traditional approach, a global method in combination with the EM algorithm is used to find the optimal set of parameters on the log-likelihood surface. In the smoothing approach, a simplified version of the global method is applied in combination with the EM algorithm to obtain an optimal set of parameters on the smooth

surface which are again used in combination with the EM algorithm to obtain optimal set of parameters on the original log-likelihood surface. Since the smoothed log-likelihood surface is easy to traverse (has fewer local maxima), one can gain significant computational benefits by applying a simplified global method compared to that of the conventional global method on the original log-likelihood surface which are usually very expensive.

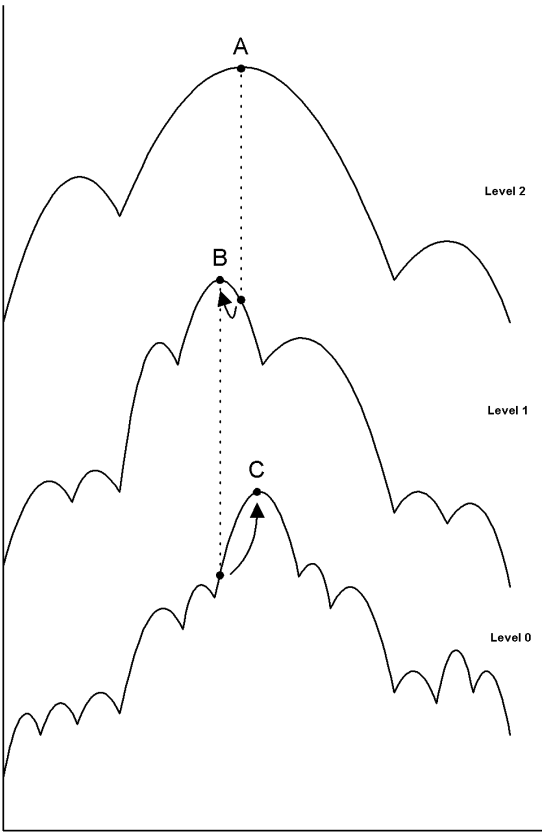


Figure 5.2: Smoothing a nonlinear surface at different levels. Tracing the global maxima (A, B and C) at different levels.

Fig. 5.2 shows a nonlinear surface and its smoothed versions in one dimension. During the smoothing process, there is no guarantee that the global maximum will retain its location. However, if the smoothing factor is not changed significantly, one can carefully traceback the global maximum by applying the EM algorithm at each level. For example, ‘C’ is the global maximum of the original nonlinear sur-

face which contains 11 local maxima and is indicated by level 0. Two smoothed versions at level 1 (with 7 local maxima) and level 2 (with 3 local maxima) can be constructed using kernel smoothing. ‘B’ indicates the global maximum in level 1 and ‘A’ indicates the global maximum in level 2. Obtaining the global maximum in level 2 is relatively easier compared to the global maximum in level 0 because of the presence of fewer local maxima. Once the point ‘A’ is obtained, it is used as an initial guess for the EM algorithm at level 1 to obtain ‘B’. Similarly ‘C’ in level 0 is obtained by applying EM algorithm initialized with ‘B’.

5.3 Preliminaries

We will now introduce some preliminaries on convolution kernels. For smoothing the mixture model, any kernel can be used for convolution if it can yield a closed form solution in each E and M step. Three widely used kernels are shown in Fig. 5.3. We choose to use Gaussian kernel for smoothing the original log-likelihood function for the following reasons :

- When the underlying distribution is assumed to be generated from Gaussian components, Gaussian kernels give the optimal performance.
- The analytic form of the likelihood surface obtained after smoothing is very similar to the original likelihood surface.
- Since the parameters of the original components and the kernels will be of the same scale, changing the parameters correspondingly to scale will be much easier.

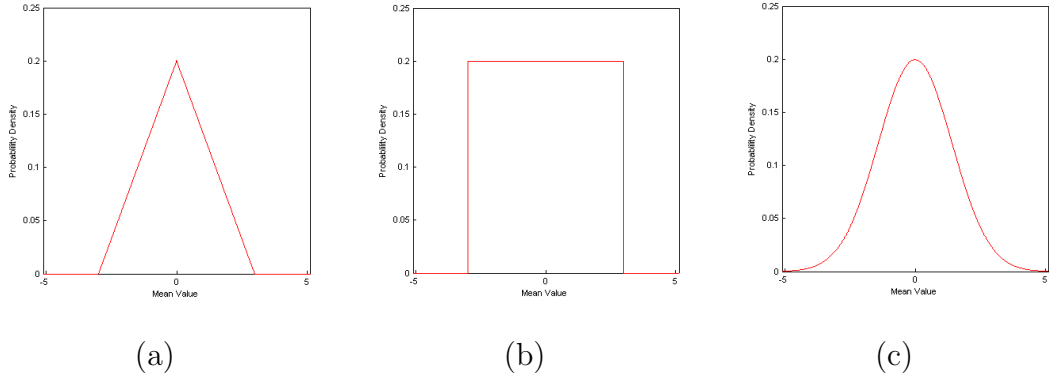


Figure 5.3: Different convolution kernels (a) Triangular function (b) Step function and (c) Gaussian function

Instead of convolving the entire likelihood surface directly, we convolve individual components of the mixture model separately. Lets consider a Gaussian density function parameterized by θ_i (i.e. μ_i and σ_i) under the assumption that all the components have the same functional form (d-variate Gaussians):

$$p(x|\theta_i) = \frac{1}{\sigma_i\sqrt{2\pi}} e^{-\frac{(x-\mu_i)^2}{2\sigma_i^2}} \quad (5.1)$$

Lets consider the following Gaussian kernel :

$$g(x) = \frac{1}{\sigma_0\sqrt{2\pi}} e^{-\frac{(x-\mu_0)^2}{2\sigma_0^2}} \quad (5.2)$$

$$\begin{aligned} p'(x|\theta_i) &= p(x|\theta_i) \otimes g(x) = \frac{1}{\sigma_i\sqrt{2\pi}} e^{-\frac{(x-\mu_i)^2}{2\sigma_i^2}} \otimes \frac{1}{\sigma_0\sqrt{2\pi}} e^{-\frac{(x-\mu_0)^2}{2\sigma_0^2}} \\ &= \frac{1}{\sqrt{2\pi}(\sigma_i^2 + \sigma_0^2)} e^{-\frac{(x-(\mu_i+\mu_0))^2}{2(\sigma_i^2 + \sigma_0^2)}} \end{aligned} \quad (5.3)$$

Convolution of Gaussians: When a Gaussian density function with parameters μ_1 and σ_1 is convolved with a Gaussian kernel with parameters μ_0 and σ_0 , then

the resultant density function is also Gaussian with mean $(\mu_1 + \mu_0)$ and variance $(\sigma_1^2 + \sigma_0^2)$. The proof for this claim is given in Appendix-A.

Now, the new smooth density can be obtained by convolving with the Gaussian kernel given by Eq. (5.2). Convolving two Gaussians to obtain another Gaussian is shown graphically in Fig. 5.4. It can also be observed that if the mean of one of the Gaussians is zero, then the mean of the resultant Gaussian is not shifted. The only change is in the variance parameter. Since, shifting mean is not a good choice for optimization problems and we are more interested in reducing the peaks, we chose to increase the variance parameter without shifting the mean.

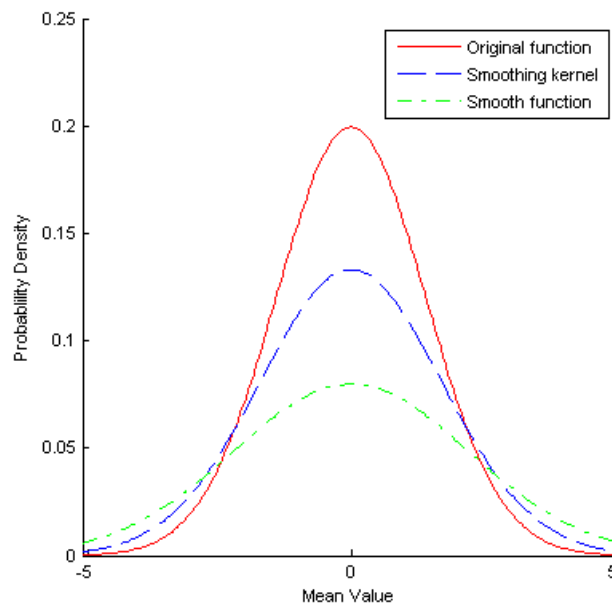


Figure 5.4: The effects of smoothing a Gaussian density function with a Gaussian kernel.

5.4 Smoothing Log-Likelihood Surface

The overall log-likelihood surface can be convolved using a Gaussian kernel directly.

This is not a feasible approach because of the following reasons :

- It results in an analytic expression that is not easy to work on and computing the EM updates will become cumbersome.
- It is Computationally very expensive.
- Different regions of search space must be smoothed differently. Choosing parameters to do this task is hard.

To avoid the first problem, we exploit the structure of the problem. Since the log-likelihood surface is obtained from individual densities, smoothing each component's individual density function will smoothen the overall log-likelihood surface. This will also give the flexibility to chose the kernel parameters which is discussed in following subsection.

After computing the new density (p'), we can define the

$$p'(x|\Theta) = \sum_{i=1}^k \alpha_i p'(x|\theta_i) \quad (5.4)$$

Now, the smooth log-likelihood function is given by:

$$f'(\mathcal{X}, \Theta) = \sum_{j=1}^n \log \sum_{i=1}^k \alpha_i p'(x^{(j)}|\theta_i) \quad (5.5)$$

Theorem 5.4.1 (*Density Smoothing*): *Convolution of a Gaussian density function with parameters μ_1 and σ_1 with a Gaussian kernel with parameters $\mu_0 = 0$ and σ_0 is equivalent to convolving the function with respect to μ_1 .*

Proof: See Appendix-B.

Fig. 5.5 shows the block diagram of the smoothing procedure. The original likelihood surface is obtained from the initial set of parameters and the given dataset.

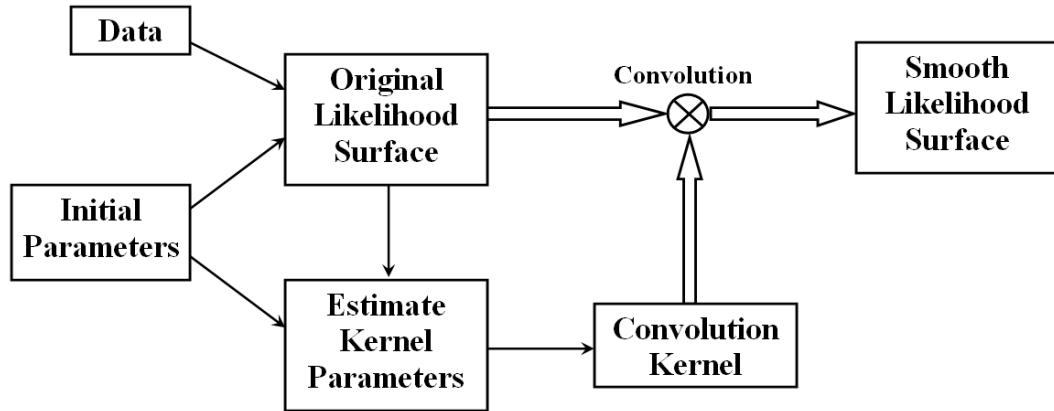


Figure 5.5: Block Diagram of the smoothing approach. Smooth likelihood surface is obtained by convolving the original likelihood surface with a convolution kernel which is chosen to be a Gaussian kernel in our case.

The kernel parameters are chosen from the initial set of parameters and the original log-likelihood surface. The kernel is then convolved with the original log-likelihood surface to obtain smooth log-likelihood surface.

5.4.1 Kernel Parameters

The parameters of the smoothing kernel can be chosen to be fixed so that they need not depend on the parameters of individual components. Fixed kernels will be effective when the underlying distribution comes from similar components. The main disadvantage of choosing a fixed kernel is that some of the components might not be smoothed while others might be over smoothed. Since, the Gaussian kernel has the property that the convolution sums up the parameters, this can also be treated as *Additive smoothing*. To avoid the problems of fixed kernel smoothing, we introduce the concept of variable kernel smoothing. Each component will be treated differently and smoothed according to the existing parameter values. This smoothing strategy is much more flexible and works effectively in practice. Since, the kernel parameters are effectively multiplied, this smoothing can be considered as *Multiplicative Smoothing*. In other words, σ_0 must be chosen individually for

different components and it must be a function of σ_i . Both these approaches don't allow for smoothing the mixing weight parameters (α 's).

5.4.2 EM Updates

For both of the above mentioned smoothing kernels, the following equations are valid. The complete derivations of these EM equations for the case of fixed kernel smoothing is given in Appendix-C. The Q – *function* of the EM algorithm applied to the smoothed log-likelihood surface is given by:

$$Q(\Theta|\hat{\Theta}(t)) = \sum_{j=1}^n \sum_{i=1}^k w_i^{(j)} \left[\log \frac{1}{\sqrt{2\pi(\tilde{\sigma}_i^2)}} - \frac{(x^{(j)} - \tilde{\mu}_i)^2}{2\tilde{\sigma}_i^2} + \log \alpha_i \right] \quad (5.6)$$

where

$$w_i^{(j)} = \frac{\frac{\alpha_i(t)}{\tilde{\sigma}_i} e^{-\frac{1}{2\tilde{\sigma}_i^2}(x^{(j)} - \tilde{\mu}_i(t))^2}}{\sum_{i=1}^k \frac{\alpha_i(t)}{\tilde{\sigma}_i} e^{-\frac{1}{2\tilde{\sigma}_i^2}(x^{(j)} - \tilde{\mu}_i(t))^2}} \quad (5.7)$$

$\tilde{\Theta}$ represents the smoothed parameters. The updates for the maximization step in the case of GMMs are given as follows :

$$\tilde{\mu}_i(t+1) = \frac{\sum_{j=1}^n w_i^{(j)} x^{(j)}}{\sum_{j=1}^n w_i^{(j)}} \quad (5.8)$$

$$\tilde{\sigma}_i^2(t+1) = \frac{\sum_{j=1}^n w_i^{(j)} (x^{(j)} - (\tilde{\mu}_i(t+1)))^2}{\sum_{j=1}^n w_i^{(j)}} \quad (5.9)$$

$$\tilde{\alpha}_i(t+1) = \frac{1}{n} \sum_{j=1}^n w_i^{(j)} \quad (5.10)$$

5.5 Algorithm and its implementation

This section describes the smoothing algorithm in detail and explains the implementation details. The basic advantage of the smoothing approach is that a simplified

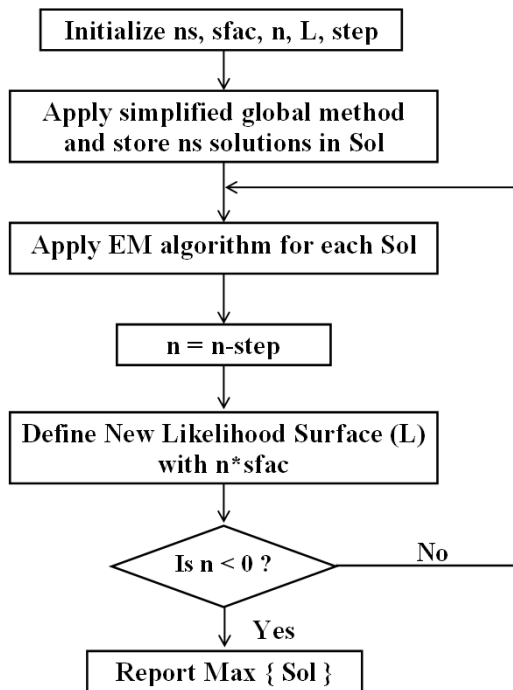


Figure 5.6: Flowchart of the smoothing algorithm

version of the global method can be used to explore fewer promising local maxima on the smoothed surface. These solutions are used as initial guesses for the EM algorithm which is again applied to the next level of smoothing. Smoothing will help to avoid search in non-promising areas of the parameter space. Fig. 5.6 gives the flow chart of our smoothing algorithm.

Before describing the algorithm, we first introduce certain variables that are used. The likelihood surface (defined by L) depends on the parameters and the available data. The smoothing factor ($sfac$) determines the extent to which the likelihood surface needs to be smoothed (which is usually chosen by trial-and-error). ns denotes the number of solutions that will be traced. n determines number of levels in the smoothing hierarchy. It is clear that there is a trade-off between the number of levels and the accuracy of this method. Having many levels might increase the accuracy of the solutions, but it is computationally expensive. On the other hand, having few levels is computationally very cheap, but we might have to forgo the

quality of the final solution. Deciding these parameters is not only user-specific but also depends significantly on the data that is being modeled. Algorithm 5 describes the smoothing approach.

Algorithm 5 Smooth-EM Algorithm

Input: Parameters Θ , Data \mathcal{X} , Tolerance τ , Smooth factor $Sfac$, number of levels nl , number of solutions ns

Output: $\hat{\Theta}_{MLE}$

Algorithm:

step=1/nl Sfac=Sfac/nl

L=Smooth(\mathcal{X} , Θ ,nl*Sfac)

Sol=Global(\mathcal{X} , Θ , L,ns)

while $n \geq 0$ **do**

 nl=nl-step

 L=Smooth(\mathcal{X} , Θ ,nl*Sfac)

for $i=1:ns$ **do**

 Sol(i)=EM(Sol(i), \mathcal{X} ,L, τ)

end for

end while

$\hat{\Theta}_{MLE} = \max\{\text{Sol}\}$

The algorithm takes smoothing factor, number of levels, number of solutions, parameters set and the data as input and computes the global maximum on the log-likelihood surface. Smooth function returns the likelihood surface corresponding to smoothing factor at each level. Initially, a simple global method is used to identify promising solutions (ns) on the smooth likelihood surface which are stored in Sol . With these solutions as initial estimates, we then apply EM algorithm on the likelihood surface corresponding to the next level smooth surface. The EM algorithm also returns ns number of solutions corresponding to the ns number of initial estimates. At every iteration, new likelihood surface is constructed with a reduced smoothing factor. This process is repeated until the smoothing factor becomes zero which corresponds to the original likelihood surface. Though, it appears to be a daunting task, it can be easily implemented in practice. The main idea is to construct a family or hierarchy of surfaces and carefully trace the promising solutions

from the top most surface to the bottom most one. In terms of tracing back the solutions to uncoarsened models, our method resembles other multi-level methods proposed in [84, 36]. The main difference is that the dimensionality of the parameter space is not changed during the smoothing (or coarsening) process.

5.6 Results and Discussion

Our algorithm has been tested on three different datasets. The initial values for the centers were chosen from the available data points randomly. The covariances were chosen randomly and uniform prior is assumed for initializing the components.

A simple synthetic data with 40 samples and 5 spherical Gaussian components was generated and tested with our algorithm. Priors were uniform and the standard deviation was 0.01. The centers for the five components are given as follows : $\mu_1 = [0.3 \ 0.3]^T$, $\mu_2 = [0.5 \ 0.5]^T$, $\mu_3 = [0.7 \ 0.7]^T$, $\mu_4 = [0.3 \ 0.7]^T$ and $\mu_5 = [0.7 \ 0.3]^T$.

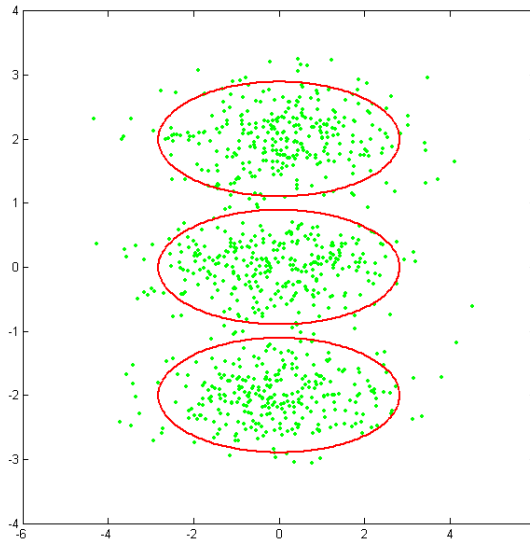


Figure 5.7: True mixture of the three Gaussian components with 900 samples.

The second dataset was that of a diagonal covariance case. The data generated from a two-dimensional, three-component Gaussian mixture distribution [138] with

mean vectors at $[0 \ -2]^T$, $[0 \ 0]^T$, $[0 \ 2]^T$ and same diagonal covariance matrix with values 2 and 0.2 along the diagonal. All the three mixtures have uniform priors. The true mixtures with data generated from these three components are shown in Fig. 5.7. In the third synthetic dataset, a more complicated overlapping Gaussian mixtures are considered [55]. It has four components with 1000 data samples (see Fig. 5.8). The parameters are as follows : $\mu_1 = \mu_2 = [-4 \ -4]^T$, $\mu_3 = [2 \ 2]^T$ and $\mu_4 = [-1 \ -6]^T$. $\alpha_1 = \alpha_2 = \alpha_3 = 0.3$ and $\alpha_4 = 0.1$.

$$C_1 = \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1 \end{bmatrix} \quad C_2 = \begin{bmatrix} 6 & -2 \\ -2 & 6 \end{bmatrix}$$

$$C_3 = \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix} \quad C_4 = \begin{bmatrix} 0.125 & 0 \\ 0 & 0.125 \end{bmatrix}$$

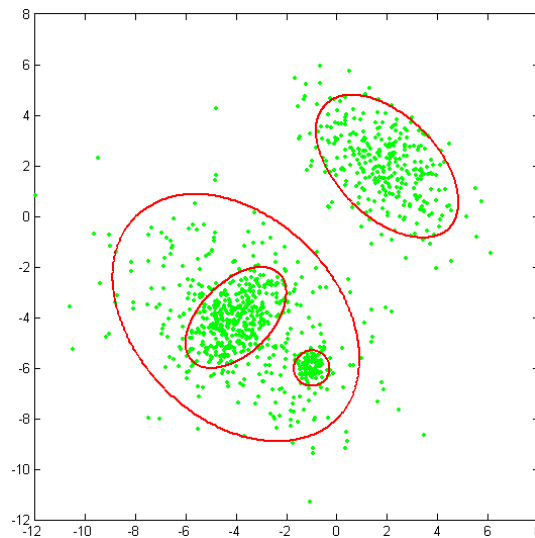
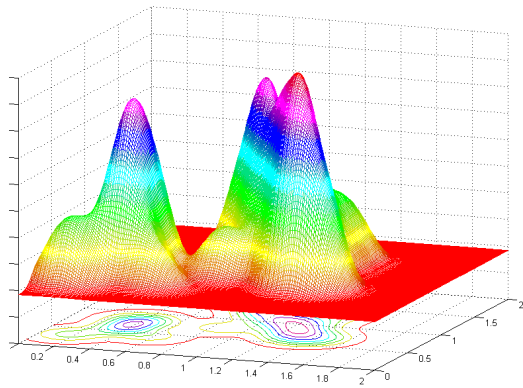
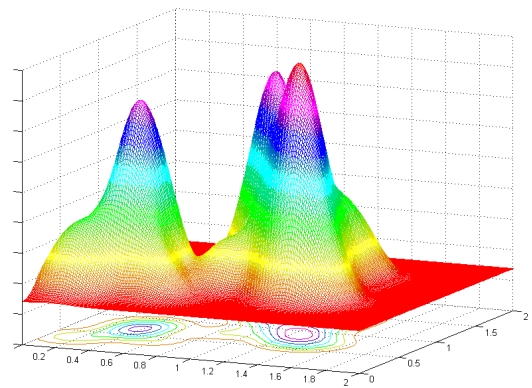


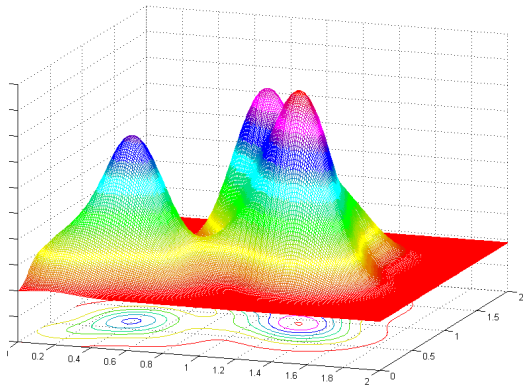
Figure 5.8: True mixtures of the more complicated overlapping Gaussian case with 1000 samples.



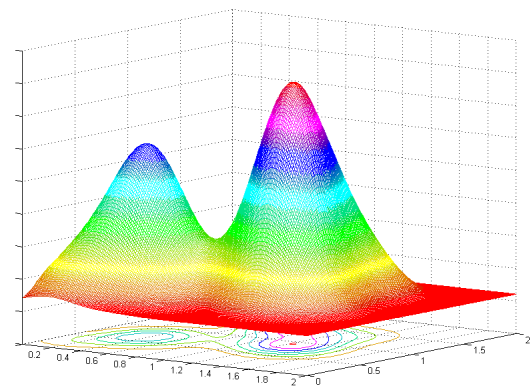
(a)



(b)



(c)



(d)

Figure 5.9: Various stages during the smoothing process. (a) The original log-likelihood surface which is very rugged (b)-(c) Intermediate smoothed surfaces (d) Final smoothed surface with only two local maxima.

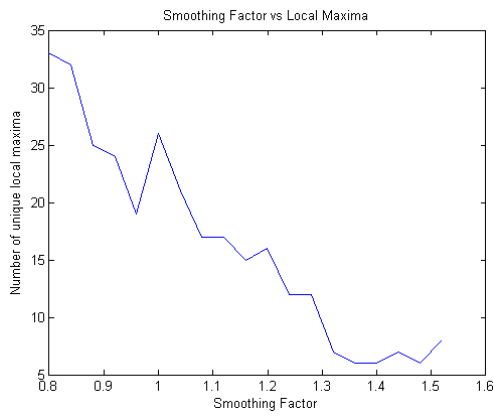
5.6.1 Reduction in the number of local maxima

One of the main advantages of the proposed smoothing algorithm is to ensure that the number of local maxima on the likelihood surface has been reduced. To the best of our knowledge, there is no theoretical way of estimating the amount of reduction in the number of unique local maximum on the likelihood surface. We hence use empirical simulations to justify the fact that the procedure indeed reduces the number of local maxima. Fig. 5.9 demonstrates the capability of our algorithm to reduce the number of local maxima. In this simple case, there were six local maxima originally, which were reduced to two local maxima after smoothing. Other stages during the transformation are also shown.

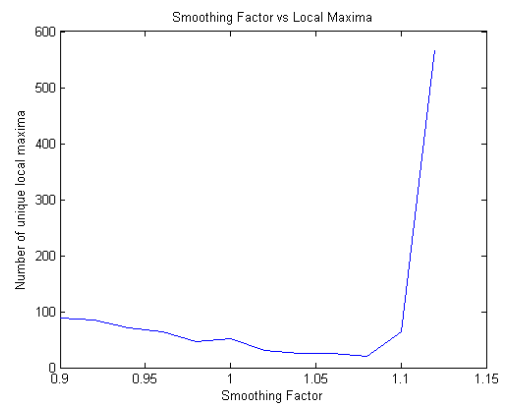
Fig. 5.10 shows the variation of the number of local maxima with respect to the smoothing factor for different datasets. One can see that if the smoothing factor is increased beyond a certain threshold value (σ_{opt}), the number of local maxima increases rapidly. This might be due to the fact that over-smoothing the surface will make the surface flat, thus making it difficult for the EM to converge. Experiments were conducted using 1000 random starts and the number of unique local maxima were stored.

5.6.2 Smoothing for Initialization

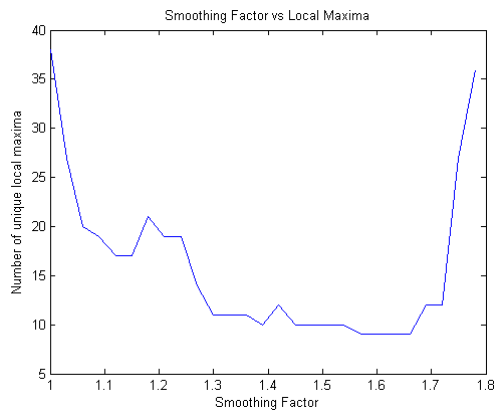
Smoothing the likelihood surface also helps in the optimization procedure. Experiments were conducted using 100 random starts. The average across all the starts is reported. The surface is then smoothed and some promising solutions are used to trace the local optimal solutions and the average across all these starts are reported. Table 5.1 summarizes the results obtained directly with the original likelihood and the smoothed likelihood. We have used only two levels and tracked three solutions



(a) Spherical Dataset



(b) Elliptical Dataset



(c) Iris Dataset

Figure 5.10: Reduction in the number of local maxima for various datasets.

for each level. The two main claims (reduction in the number of local maximum and better initial estimates) about the contributions have been justified.

Table 5.1: Comparison of smoothing algorithm with the random starts. Mean and standard deviations across 100 random starts are reported.

| Dataset | RS+EM | Smooth+EM |
|-----------------|---------------------|---------------------|
| Spherical | 36.3 ± 2.33 | 41.22 ± 0.79 |
| Elliptical | -3219 ± 0.7 | -3106 ± 12 |
| Full covariance | -2391.3 ± 35.3 | -2164.3 ± 18.56 |
| Iris | -196.34 ± 15.43 | -183.51 ± 2.12 |

More sophisticated global methods like Genetic algorithms, simulated annealing, adaptive partitioning [135] etc. and their simplified versions can also be used in combination with our approach. Since the main focus of our work is to demonstrate the smoothing capability, we used multiple random restarts as our global method.

Our algorithm is based on the conventional Expectation Maximization (EM) approach applied to a smoothed likelihood surface. A hierarchy of smooth surfaces is constructed and optimal set of parameters are obtained by tracing back the promising solutions at each level. The basic idea here is to obtain promising solutions in the smoothed surface and carefully trace the corresponding solutions in the intermediate surfaces by applying EM algorithm at every level. This smoothing process not only reduces the overall gradient of the surface but also reduces the number of local maxima. This is an effective optimization procedure that eliminates extensive search in the non-promising areas of the parameter space. Benchmark results demonstrate a significant improvement of the proposed algorithm compared to other existing methods.

One can apply TRUST-TECH method for tracing the optimal solutions across different smooth levels. Applying an EM algorithm might not guarantee to trace the optimal solution at each level. Especially, if the smoothing procedure distorts

the surface to a considerable amount, it will be difficult to trace the optimal solution by merely applying the EM algorithm. Using TRUST-TECH, one can avoid this problem. Even if there is a considerable distortion in the surface, TRUST-TECH can help in searching the neighborhood regions effectively. In summary, we can use TRUST-TECH (instead of EM) to trace the optimal solutions.

APPENDIX-A: Convolution of two Gaussians

Proof. Lets consider two gaussian density functions with parameters θ_1 and θ_0 .

$$P(x|\theta_1) = \frac{1}{\sqrt{2\pi}\sigma_1} e^{-\frac{(x-\mu_1)^2}{2\sigma_1^2}} \quad (5.11)$$

$$P(x|\theta_0) = \frac{1}{\sqrt{2\pi}\sigma_0} e^{-\frac{(x-\mu_0)^2}{2\sigma_0^2}} \quad (5.12)$$

$$(5.13)$$

By definition of convolution, we have

$$g(t) \otimes h(t) = \int_{-\infty}^{\infty} g(\tau)h(t - \tau)d\tau \quad (5.14)$$

$$c(x|\theta_1, \theta_0) = p(x|\theta_1) \otimes p(x|\theta_0) \quad (5.15)$$

$$= \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}\sigma_1} e^{-\frac{((x-\tau)-\mu_1)^2}{2\sigma_1^2}} \frac{1}{\sqrt{2\pi}\sigma_0} e^{-\frac{(\tau-\mu_0)^2}{2\sigma_0^2}} d\tau \quad (5.16)$$

$$(5.17)$$

$$= \frac{1}{\sqrt{2\pi}\sigma_1} \frac{1}{\sqrt{2\pi}\sigma_0} \int_{-\infty}^{\infty} e^{-\frac{\sigma_0^2(\tau-x+\mu_1)^2 + \sigma_1^2(\tau-\mu_0)^2}{2\sigma_1^2\sigma_0^2}} d\tau \quad (5.18)$$

After rearranging the terms that are independent of τ and further simplification, we get

$$c(x|\theta_1, \theta_0) = \frac{e^{-\frac{(x-(\mu_1+\mu_0))^2}{2(\sigma_1^2+\sigma_0^2)}}}{\sqrt{2\pi}\sigma_1 \sqrt{2\pi}\sigma_0} \int_{-\infty}^{\infty} e^{-\frac{(\tau-\mu_\tau)^2}{2\sigma_\tau^2}} d\tau \quad (5.19)$$

where

$$\mu_\tau = \frac{\mu_0\sigma_1^2 + (x - \mu_1)\sigma_0^2}{(\sigma_1^2 + \sigma_0^2)} \sigma_\tau = \frac{\sigma_1^2\sigma_0^2}{(\sigma_1^2 + \sigma_0^2)} \quad (5.20)$$

Hence, we have

$$\begin{aligned} c(x|\theta_1, \theta_0) &= \frac{e^{-\frac{(x-(\mu_1+\mu_0))^2}{2(\sigma_1^2+\sigma_0^2)}}}{\sqrt{2\pi(\sigma_1^2+\sigma_0^2)}} \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi} \sigma_\tau} e^{-\frac{(\tau-\mu_\tau)^2}{2\sigma_\tau^2}} d\tau \\ &= \frac{1}{\sqrt{2\pi(\sigma_1^2+\sigma_0^2)}} e^{-\frac{(x-(\mu_1+\mu_0))^2}{2(\sigma_1^2+\sigma_0^2)}} \end{aligned} \tag{5.21}$$

(because the quantity inside the integral is 1 for a Gaussian density function.) \triangleleft

APPENDIX-B: Proof of Theorem 5.4.1

Proof. Convolution of Gaussian density with respect to the mean is shown below:

$$\tilde{c}(x, \theta_0) = \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}\sigma_1} e^{-\frac{(x-(\mu_1-\tau))^2}{2\sigma_1^2}} \frac{1}{\sqrt{2\pi}\sigma_0} e^{-\frac{\tau^2}{2\sigma_0^2}} d\tau \quad (5.22)$$

Now, consider Eq. (5.15). Substituting $\mu_0 = 0$ and replacing τ with $-\tau$, we get

$$c(x, \theta_0) = \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}\sigma_1} e^{-\frac{(x+\tau-\mu_1)^2}{2\sigma_1^2}} \frac{1}{\sqrt{2\pi}\sigma_0} e^{-\frac{\tau^2}{2\sigma_0^2}} d\tau \quad (5.23)$$

From Eqs. (5.22) and (5.23), we can see that convolution of a Gaussian density function with a Gaussian density with zero mean is equivalent to convolving the function with respect to mean.◁

APPENDIX-C: Derivations for EM updates

For simplicity, we show the derivations for EM updates in the fixed kernel case. Lets consider the case where a fixed Gaussian kernel with parameters μ_0 and σ_0 which will be used to convolve each component of the GMM. We know that

$$\log p(\mathcal{X}, \mathcal{Z}|\Theta) = \log \prod_{j=1}^n p(x^{(j)}|z^{(j)}, \Theta) \cdot p(z^{(j)}) \quad (5.24)$$

For the j^{th} data point, we have

$$p(x^{(j)}|z^{(j)}, \Theta) \cdot p(z^{(j)}) = \prod_{i=1}^k \left[\frac{1}{\sqrt{2\pi(\sigma_i^2 + \sigma_0^2)}} e^{-\frac{(x - (\mu_i + \mu_0))^2}{2(\sigma_i^2 + \sigma_0^2)}} p(z_i^{(j)} = 1) \right]^{z_i^{(j)}} \quad (5.25)$$

Hence,

$$\begin{aligned} \log p(\mathcal{X}, \mathcal{Z}|\Theta) &= \sum_{j=1}^n \log p(x^{(j)}|z^{(j)}, \Theta) \cdot p(z^{(j)}) \\ &= \sum_{j=1}^n \sum_{i=1}^k \log \left[\frac{1}{\sqrt{2\pi(\sigma_i^2 + \sigma_0^2)}} e^{-\frac{(x - (\mu_i + \mu_0))^2}{2(\sigma_i^2 + \sigma_0^2)}} p(z_i^{(j)} = 1) \right]^{z_i^{(j)}} \\ &= \sum_{j=1}^n \sum_{i=1}^k z_i^{(j)} \left[-\log(\sqrt{2\pi(\sigma_i^2 + \sigma_0^2)}) - \frac{(x - (\mu_i + \mu_0))^2}{2(\sigma_i^2 + \sigma_0^2)} + \log \alpha_i \right] \end{aligned} \quad (5.26)$$

Expectation Step : For this step, we need to compute the Q -function which is the expected value of Eq. (5.26) with respect to the hidden variables.

$$\begin{aligned} Q(\Theta|\Theta^{(t)}) &= E_z [\log p(\mathcal{X}, \mathcal{Z}|\Theta)|\mathcal{X}, \Theta^{(t)}] \\ &= \sum_{j=1}^n \sum_{i=1}^k E_z[z_i^{(j)}] \left[-\log(\sqrt{2\pi(\sigma_i^2 + \sigma_0^2)}) \right. \\ &\quad \left. - \frac{(x - (\mu_i + \mu_0))^2}{2(\sigma_i^2 + \sigma_0^2)} + \log \alpha_i \right] \end{aligned} \quad (5.27)$$

To compute the Expected value of the hidden variables ($w_i^{(j)}$),

$$\begin{aligned}
w_i^{(j)} &= E_z[z_i^{(j)}] = \sum_{c=0}^1 c * p(z_i^{(j)} = c | \Theta^{(t)}, x^{(j)}) \\
&= \frac{p(x^{(j)} | \Theta^{(t)}, z_i^{(j)} = 1) p(z_i^{(j)} = 1 | \Theta^{(t)})}{p(x^{(j)} | \Theta^{(t)})} \\
&= \frac{\frac{1}{\sqrt{(\sigma_i^2 + \sigma_0^2)}} e^{-\frac{(x - (\mu_i + \mu_0))^2}{2(\sigma_i^2 + \sigma_0^2)}} \alpha_i^{(t)}}{\sum_{m=1}^k \frac{1}{\sqrt{(\sigma_m^2 + \sigma_0^2)}} e^{-\frac{(x - (\mu_m + \mu_0))^2}{2(\sigma_m^2 + \sigma_0^2)}} \alpha_m^{(t)}}
\end{aligned} \tag{5.28}$$

Maximization Step : The maximization step is given by the following equation :

$$\frac{\partial}{\partial \Theta_i} Q(\Theta | \hat{\Theta}(t)) = 0 \tag{5.29}$$

where Θ_i are the parameters of the i^{th} component. Due to the assumption made that each data point comes from a single component, solving the above equation becomes trivial. The updates for the maximization step in the case of GMMs are given as follows :

$$(\mu_i + \mu_0) = \frac{\sum_{j=1}^n w_i^{(j)} x^{(j)}}{\sum_{j=1}^n w_i^{(j)}} \tag{5.30}$$

$$(\sigma_i^2 + \sigma_0^2) = \frac{\sum_{j=1}^n w_i^{(j)} (x^{(j)} - (\mu_i + \mu_0))^2}{\sum_{j=1}^n w_i^{(j)}} \tag{5.31}$$

$$\alpha_i = \frac{1}{n} \sum_{j=1}^n w_i^{(j)} \tag{5.32}$$

Chapter 6

TRUST-TECH based Neural Network

Training

Supervised learning using artificial neural networks has numerous applications in various domains of science and engineering. Efficient training mechanisms in a neural network play a vital role in deciding the network architecture and the accuracy of the classifier. Most popular training algorithms tend to be greedy and hence get stuck at the nearest local minimum of the error surface. To overcome this problem, some global methods (like multiple restarts, genetic algorithms, simulated annealing etc.) for efficient training make use of stochastic approaches in combination with local methods to obtain an effective set of training parameters. Due to the stochastic nature and lack of effective *fine tuning* capability, these algorithms often fail to obtain an optimal set of training parameters. In this chapter, a new method to improve the local search capability of training algorithms is proposed. This new method takes advantage of TRUST-TECH to compute neighborhood local minimum of the error surface. The proposed approach obtains multiple local optimal solutions surrounding the current local optimal solution in a systematic manner. Empirical results on different machine learning datasets indicate that the proposed algorithm outperforms current algorithms available in the literature.

6.1 Overview

Artificial neural networks (ANN) were developed analogous to the human brain for the purpose of improving conventional learning capabilities. They are used for a wide variety of applications in diverse areas such as function approximation, time-series

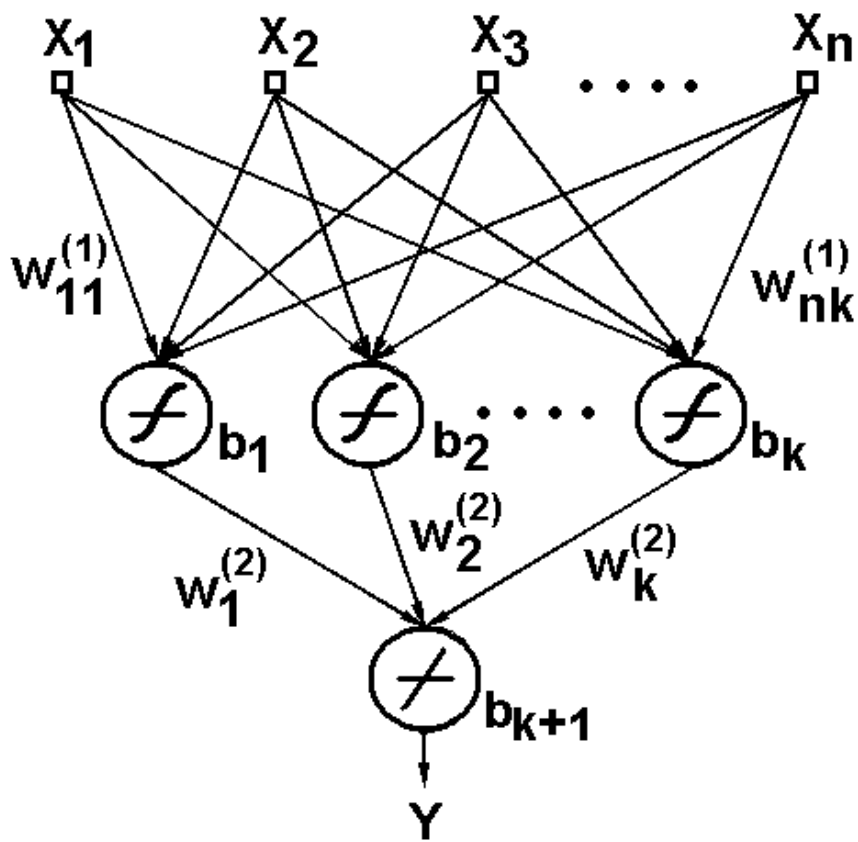


Figure 6.1: The Architecture of MLP with single hidden layer having k hidden nodes and the output layer having a single output node. x_1, x_2, \dots, x_n is an n -dimensional input feature vector. w_{ij} are the weights and b_1, b_2, \dots, b_k are the biases for these k nodes. The activation function for the hidden nodes is sigmoidal and for the output node is linear.

prediction, medical diagnosis, character recognition, load forecasting, speaker identification and risk management. These networks serve as excellent approximators of nonlinear continuous functions [93]. However, using an artificial neural network to model a system usually involves dealing with certain difficulties in achieving the best representation of the classification problem.

The two challenging tasks in the process of learning using ANNs are network architecture selection and optimal training. In deciding the architecture for the feedforward neural network (also known as Multi-Layer Perceptron, MLP), a larger network will always provide better prediction accuracy for the data available. However, such a large network that is too complicated and customized to some given problem will lose its generalization capability for the unseen data [19]. Also, every additional neuron translates to increased hardware cost. Hence, it is vital to develop algorithms that can exploit the potential of a given architecture which can be achieved by obtaining the global minimum of the error on the training data. Hence, the goal of optimal training of the network is to find a set of weights that achieves the global minimum of the mean square error (MSE) [68]. Fig. 6.1 shows the architecture of a single hidden layer neural network with n input nodes, k hidden nodes and 1 output node. The network is trained to deliver the output value (Y_i) for the i^{th} sample at the output node which will be compared to the actual target value (t_i).

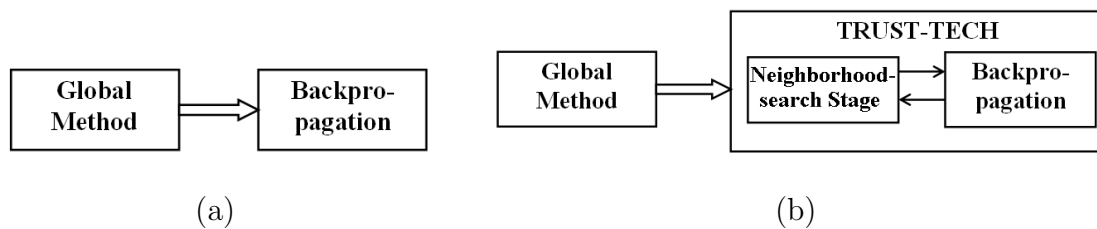


Figure 6.2: Comparison between the two frameworks (a) Traditional approach and (b) TRUST-TECH based approach. The main difference is the inclusion of the stability region based neighborhood-search stage that can explore the neighborhood solutions.

The main focus of this work is to develop a robust training algorithm for obtaining the optimal set of weights of an artificial neural network. Several training algorithms have been extensively studied in the literature [68]. Backpropagation (BP) algorithm is a very robust deterministic local method that have received significant attention. Though BP is comparatively cheaper in terms of time and easy to implement, it can only obtain local optimal solutions. On the contrary, some global methods like multiple random starts, genetic algorithms and simulated annealing can identify promising regions of the weight space, but are essentially stochastic in nature and computationally expensive. Expecting such stochastic algorithms to fine-tune the training weights will be even more time consuming. Thus, there is a necessity to efficiently search for good solutions in promising regions of the solution space, which can be accomplished by the newly proposed TRUST-TECH based algorithm. In this chapter, we introduce a novel algorithm that will search the weight space in a systematic tier-by-tier manner. Fig. 6.2 compares the traditional approach with our proposed approach. The main difference between the two approaches is the inclusion of the neighborhood-search stage, where an improved set of training weights are obtained by systematically exploring the neighborhood local optimal solutions in a promising subspace.

6.2 Relevant Background

The performance of a feedforward neural network is usually gauged by measuring the MSE of its outputs from the expected target values. The goal of optimal training is to find a set of parameters that achieves the global minimum of the MSE [14, 132, 93]. For a n -dimensional dataset, the MSE over Q samples in the training set is given as

follows :

$$C(W) = \frac{1}{Q} \sum_{i=1}^Q [t(i) - y(X, W)]^2 \quad (6.1)$$

where $t(i)$ is the target output for the i^{th} sample, X is the input vector and W is the weight vector. The MSE as a function of the weight parameters is highly nonlinear containing several local minima. The network's weights and thresholds must be set so as to minimize the prediction error made by the network. Since it is not possible to analytically determine the global minimum of the error surface, the neural network training is essentially an exploration of the error surface for an optimal set of parameters that attains this globally optimal solution.

Training algorithms can be broadly classified into '*local*' and '*global*' methods. Local methods begin at some initial points and deterministically move towards a local minimum. From an initial random configuration of weights and thresholds, these local training methods incrementally (greedily) seek for improved solution until they reach a local minimum. Typically, some form of the gradient information at the current point on the error surface is calculated and used to make a downhill move. Eventually, the algorithm stops at a low point, which usually is a local minimum. In the context of training neural networks, this local minima problem is a well-studied research topic [63]. The most commonly used training method in MLP is the backpropagation algorithm [129] which has been tested successfully for different kinds of problems. Despite having many variants, BP faces the problem of stopping at local minimum instead of proceeding towards the global minimum [76, 132]. Modifications [92] to the basic BP model have been suggested to help the algorithm escape from being trapped in a local minimum. However, while these improved methods reduce the tendency to sink into local minimum by providing some form of

perturbations to the search direction, it does not train the network to converge to a global minimum within a reasonable number of iterations [80, 141]. Based on the movement towards improved solutions, local methods can be subdivided into two categories:

1. **Line search methods:** These algorithms select some descent direction (based on the gradient information) and minimize the error function value along this particular direction. This process is repeated until a local minimum is reached. Most popular choices for the descent directions are Newton's direction or conjugate direction. In the context of neural networks, apart from the obvious steepest descent methods, other widely used line search algorithms are Newton's method [9], the BFGS method [111] and conjugate gradient methods [26, 104].
2. **Trust region methods:** Trust region methods are by far the fastest convergent methods compared to the above mentioned line-search methods. The surface is assumed to be a simple model (like a parabola) such that the minimum can be located directly if the model assumption is good which usually happens when the initial guess is close to the local minimum. They require more storage space compared to conjugate gradient methods [65] and hence are not effective for large-scale applications.

All these local methods discussed so far assume that they already have an initial guess to begin with. Usually the quality of the final solution depends significantly on the initial set of parameters available. Hence, in practice, none of these local methods are used by themselves. They are usually combined with stochastic global methods which yield a promising set of parameters in the weight space. These global methods explore the entire error surface and thus the chance of attaining a near-

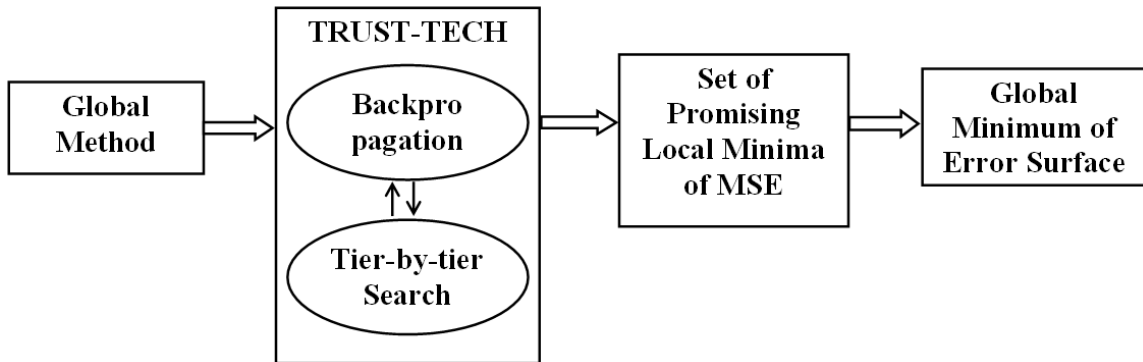


Figure 6.3: Block Diagram of the TRUST-TECH based training method.

global optimal solution is high. More advanced techniques like Genetic algorithms [19] and simulated annealing [1] are applied in combination with standard BP in order to obtain more promising solutions and avoid being stuck at local minimum [124]. The use of various global optimization algorithms for finding an effective set of training parameters is comprehensively given in [132]. Although these methods (asymptotically) guarantee convergence to the global minimum, it usually exhibits very slow convergence even for simple learning tasks. Though methods can explore the entire solution space effectively and obtain promising local optimal solutions, it lacks fine-tuning capabilities to obtain a precise final solution and requires local methods like BP to be employed. Other traditional methods like Monte Carlo method, Tabu search, ant colony optimization and particle swarm optimization are also stochastic in nature and suffer from the same problems described above.

From the above discussion, one can realize that there is a clear gap between global and local methods. Typically, most of the successful practical algorithms are a combination of these global and local methods. In other words, these two approaches do not communicate well between each other. Approaches that might resemble our methodology are TRUST [25] and dynamic tunneling [128]. These methods attempt to move out of the local minimum in a stochastic manner. The training algorithm

proposed in this chapter differs from these two methods by deterministically escaping out of the local minimum and systematically exploring multiple local minima on the error surface in a tier-by-tier manner in order to advance towards the global minimum. This approach is based on the fundamental concepts of stability regions that were established in [31, 91]. Fig. 6.3 shows the block diagram of the TRUST-TECH methodology. Basically, a global method yields points in certain promising regions of the search space. These points are used as initial guesses to search the neighborhood subspace in a systematic manner. TRUST-TECH relies on a robust, fast local method to obtain a local optimal solution. It explores the parameter subspace in a tier-by-tier manner by transforming the function into its corresponding dynamical system and exploring the neighboring stability regions. Thus, it gives a set of promising local optimal solutions from which a global minimum is selected. In this manner, TRUST-TECH can be treated as an effective interface between the global and local methods, which enables the communication between these two methods. It also allows the flexibility of choosing different global and local methods depending on their availability and performance for certain specific classification tasks.

6.3 Training Neural Networks

Without loss of generality, we consider a feedforward neural network with one input layer, one hidden layer and one output layer. Specifically, the output layer contains only one node that will yield all the possible target values depending on its activation function. Table 6.1 gives the notations used in the rest of this chapter.

Let k be the number of hidden nodes in the hidden layer and the input vector is

Table 6.1: Description of the notations used

| Notation | Description |
|----------|---|
| Q | Number of training samples |
| X | Input vector |
| W | Weight vector |
| n | Number of features |
| k | Number of hidden nodes |
| w_{0j} | weight between the output node and the j^{th} hidden node |
| w_{ij} | weight between the i^{th} input node and the j^{th} hidden node |
| b_0 | bias of the output node |
| b_j | bias of the j^{th} hidden node |
| ϕ_1 | Activation function of the hidden nodes |
| ϕ_2 | Activation function of the output node |
| t_i | target value of the i^{th} input sample |
| y | output of the network |
| e_i | Error for the i^{th} input sample |

n -dimensional. Then the final nonlinear mapping of our model is given by :

$$y(W, X) = \phi_2 \left(\sum_{j=1}^k w_{0j} \phi_1 \left(\sum_{i=1}^n w_{ij} x_i + b_j \right) + b_0 \right) \quad (6.2)$$

where ϕ_1 and ϕ_2 are the activation functions of the hidden nodes and the output nodes respectively. ϕ_1 and ϕ_2 can be same functions or can be different functions. We have chosen to use ϕ_1 to be sigmoidal and ϕ_2 to be linear. Results in the literature [64], suggest that this set of activation functions yield the best results for feedforward neural networks. As shown in Fig. 6.1, w_{0j} indicate the weights between the hidden layer and the output layer and w_{ij} indicate the weights between the input layer and the hidden layer. b_j are the biases of the k hidden nodes and b_0 is the bias of the output node. x_i is the n -dimensional input feature vector and X_i indicates the i^{th} training sample. The task of the network is to learn associations between the input-output pairs $(X_1, t_1), (X_2, t_2), \dots, (X_Q, t_Q)$. The weight vector to be optimized is constructed as follows:

$$w = (w_{01}, w_{02}, \dots, w_{0k}, \dots, w_{n1}, w_{n2}, \dots, w_{nk}, b_0, b_1, b_2, \dots, b_k)^T$$

which includes all the weights and biases that are to be computed. Hence, the problem of training neural networks is s -dimensional unconstrained minimization problem where $s = (n + 2)k + 1$.

$$\min_w C(w) \tag{6.3}$$

The mean squared error which is to be minimized can be written as

$$C(w) = \frac{1}{Q} \sum_{i=1}^Q e_i^2(w) \tag{6.4}$$

where the error

$$e_i(w) = t_i - y(w, x_i) \tag{6.5}$$

The error cost function $C(\cdot)$ averaged over all training data is a highly nonlinear function of the synaptic vector w . Ignoring the constant for simplicity, it can be shown that

$$\nabla C(w) = J^T(w)e(w) \tag{6.6}$$

$$\nabla^2 C(w) = J^T(w)J(w) + S(w) \tag{6.7}$$

where $J(w)$ is the Jacobian matrix

$$J(w) = \begin{bmatrix} \frac{\partial e_1}{\partial W_1} & \frac{\partial e_1}{\partial W_2} & \cdot & \cdot & \frac{\partial e_1}{\partial W_N} \\ \frac{\partial e_2}{\partial W_1} & \frac{\partial e_2}{\partial W_2} & \cdot & \cdot & \frac{\partial e_2}{\partial W_N} \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ \frac{\partial e_Q}{\partial W_1} & \frac{\partial e_Q}{\partial W_2} & \cdot & \cdot & \frac{\partial e_Q}{\partial W_N} \end{bmatrix}$$

and

$$S(w) = \sum_{i=1}^Q e_i(w) \nabla^2 e_i(w) \quad (6.8)$$

Generally, if we would like to minimize $J(w)$ with respect to the parameter vector w , any variation of Newton's method can be written as

$$\begin{aligned} \Delta w &= - [\nabla^2 C(w)]^{-1} \nabla C(w) \\ &= - [J^T(w)J(w) + S(w)]^{-1} J^T(w)e(w) \end{aligned} \quad (6.9)$$

6.4 Problem Transformation

We explore the geometrical structure of the error surface to explore multiple local optimal solutions in a systematic manner. Firstly, we describe the transformation of the original minimization problem into its corresponding nonlinear dynamical system and then propose a new TRUST-TECH based training algorithm for finding multiple local optimal solutions.

This section mainly deals with the transformation of the original error function into its corresponding nonlinear dynamical system and introduces some terminology pertinent to comprehend our algorithm. This transformation gives the correspondence between all the critical points of the error surface and that of its corresponding

gradient system. To analyze the geometric structure of the error surface, we build a *generalized gradient system* described by

$$\frac{dw}{dt} = -A(w)\nabla C(w) \quad (6.10)$$

where the error function C is assumed to be twice differentiable to guarantee unique solution for each initial condition $w(0)$ and $A(w)$ is a positive definite symmetric matrix for all $w \in \mathfrak{R}^n$. It is interesting to note the relationship between Eqs. (6.10) and (6.9) and obtain different local solving methods used to find the nearest local optimal solution with guaranteed convergence. For example, if $A(w) = I$, then it is a naive error back-propagation algorithm. If $A(w) = [J(w)^T J(w)]$ then it is the Gauss-Newton method and if $A(w) = [J(w)^T J(w) + \mu I]$ then it is the Levenberg-Marquardt method. This transformation of the original error function to its corresponding dynamical system will enable us to transform the problem of finding multiple local minima on the error surface into the problem of finding multiple stable equilibrium points of its corresponding dynamical system. This will enable us to apply TRUST-TECH method for training neural networks and obtain promising solutions.

6.5 TRUST-TECH based Training

The proposed TRUST-TECH based algorithm for training neural networks, uses a promising starting point (A^*) as input and outputs the best local minimum of the neighborhood in the weight space. Figure 6.4 shows the flowchart of our approach.

Input: Initial guess(A^*), Tolerance (τ), Step size (s)

Output: Best local minimum (A_{ij}) in the neighborhood

Algorithm:

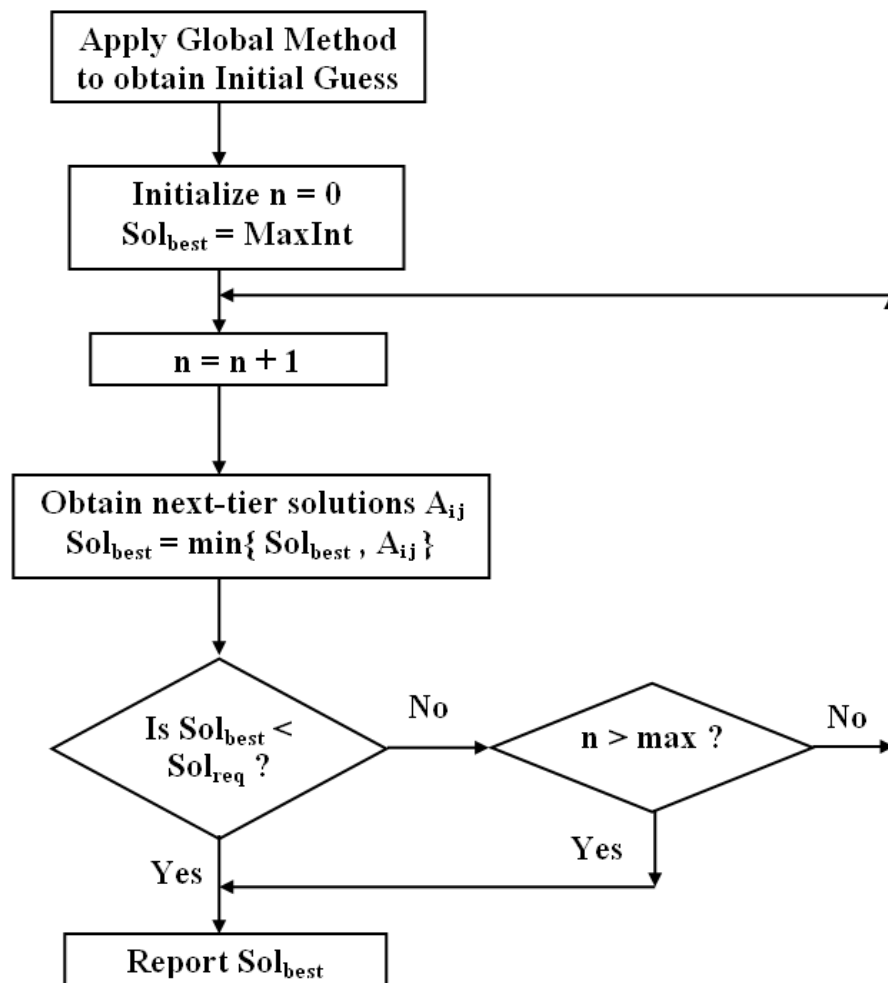


Figure 6.4: Flow chart of our method

Step 1: Obtaining good initial guess (A^):* The initial guess for the algorithm can be obtained from global search methods or from a purely random start. Some domain knowledge about the specific dataset that the network is being trained on, might help in eliminating non-promising set of initial weights.

Step 2: Moving to the local minimum (M): Using an appropriate local solver (such as conjugate-gradient, quasi-Newton or Levenberg-Marquardt), the local optimum M is obtained using A^* as the initial guess.

Step 3: Determining the search direction (d_j): The eigenvectors d_j of the Jacobian are computed at m_i . These eigenvector directions might lead to promising regions of the subspace. Other search directions can also be chosen based on the specific problem that is being dealt.

Step 4: Escaping from the local minimum: Taking small step sizes away from m_i along the d_j directions increases the objective function value till it hits the stability boundary. However, the objective function value then decreases after the search trajectory moves away from the exit point. This new point is used as initial guess and local solver is applied again (go to Step 2).

Step 5: Finding Tier-1 local minima (A_{1i}): Exploring the neighborhood of the local optimal solution corresponding to the initial guess leads to tier-1 local minima. Exploring from tier- k local minima leads to tier- $k + 1$ local minima.

Step 6: Exploring Tier- k local minima (A_{kj}): Explore all other tiers in the similar manner described above (see Fig. 6.5). From all these solutions, the best one is chosen to be the desired global optimum.

Step 7: Termination Criteria: The procedure can be terminated when the best solution obtained so far is satisfactory (lesser than sol_{req}) or a predefined maximum number of tiers is explored.

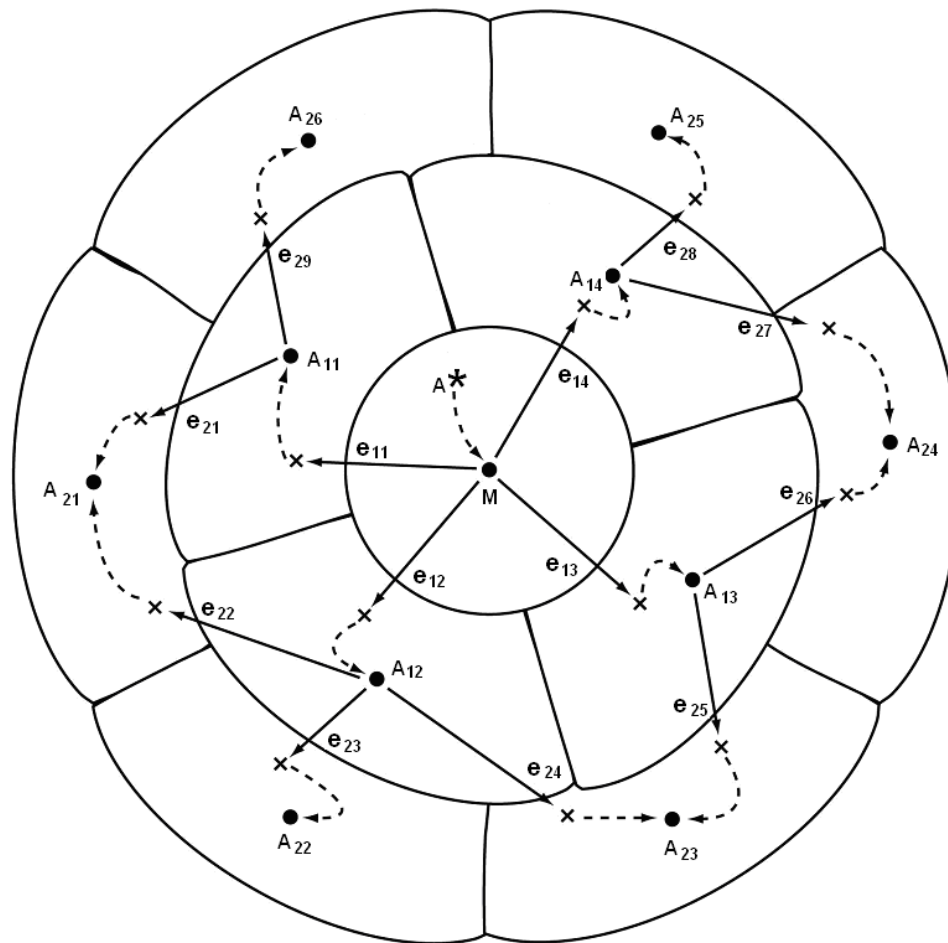


Figure 6.5: Diagram Illustrating two tier exit point strategy. The ' A^* ' represents the initial guess. Dotted arrows represent the convergence of the local solver. Solid arrows represent the gradient ascent linear searches along eigenvector directions. ' X ' indicates a new initial condition in the neighboring stability region. M represents the local minimum obtained by applying local methods from ' X '. A_{1i} indicates Tier-1 local minima. e_{1i} are the exit points between M and A_{1i} . Similarly, A_{2j} and e_{2j} are the second-tier local minima and their corresponding exit points respectively.

6.6 Implementation Details

All programs were implemented in MATLAB v6.5 and run on Pentium IV 2.8 GHz machines. This section describes the various implementation details used in our simulations. The following issues are discussed in detail : (i) Architecture and local methods, (ii) Initialization Schemes and (iii) TRUST-TECH.

6.6.1 Architecture and Local Methods

As described in the introduction section, we have chosen to demonstrate the capability of our newly proposed TRUST-TECH algorithm on a network with single hidden layer and an output layer containing only one output node. This architecture is not complicated and has the capability to precisely demonstrate the problems with the existing approaches. A network described here contains n (number of attributes) input nodes which is equal to the number of features available in the dataset, one hidden layer with k nodes and one output node. Thus, each network has nk weights and k biases to the hidden layer, and k weights and one bias to the output node. Hence, training a neural network is necessarily a search problem of dimensionality $(n + 2)k + 1$. Each hidden node has a tangent-sigmoid transfer function and the output node has a pure linear transfer function. The number of nodes in the hidden layer is determined by incrementally adding hidden nodes, and selecting the architecture that achieves a compromise between minimal error value and minimal number of nodes. The trust region based Levenberg-Marquardt algorithm is chosen because of efficiency in terms of time and space consumption. It utilizes the approximation of the Jacobian in its iterative gradient descent, which will be used for generating promising directions in TRUST-TECH.

6.6.2 Initialization Schemes

Two different initialization schemes were implemented. The most basic global method which is multiple random starts with initial set of parameters between -1 and 1. More effective global method namely Nguyen-Widrow (NW) algorithm [108] has also been used to test the performance of our algorithm. The NW algorithm is implemented as the standard initialization procedure in MATLAB. In both cases, the best initial set of parameters in terms of training error is chosen and improved with our TRUST-TECH algorithm.

Algorithm 6 *New_Wts TRUST_TECH(NE T, Wts, s, τ)*

```
Wts = Train(NE T, Wts,  $\tau$ )
Error = Estimate(NE T, Wts)
Thresh = c * Error
Wts1[ ] = Neighbors(NE T, Wts, s,  $\tau$ )
for k = 1 to size(Wts1) do
    if Estimate(NE T, Wts1[k]) < Thresh then
        Wts2[k][ ] = Neighbors(NE T, Wts1, s,  $\tau$ )
    end if
end for
Return best(Wts, Wts1, Wts2)
```

6.6.3 TRUST-TECH

It is effective to apply the TRUST-TECH methodology to those promising solutions obtained from stochastic global methods. Algorithm 6 describes the two-tier TRUST-TECH algorithm. *NE T* assumes to have a fixed architecture with a single output node. *s* is the step size used for evaluating the objective function value till it obtains an exit point. τ is the tolerance of error used for the convergence of the local method. *Weights* give the initial set of weight parameter values. *Train* function implements the Levenberg-Marquardt method that obtains the local optimal solution from the initial condition. The procedure *Estimate* computes the MSE value

of the network model. A threshold value (*Thresh*) is set based on this MSE value. The procedure *Neighbors* returns all the next tier local optimal solutions from a given solution. After obtaining all the tier-1 solutions, the procedure *Neighbors* is again invoked (only for promising solutions) to obtain the second-tier solutions. The algorithm finally compares the initial solution, tier-1 and tier-2 solutions and returns the network corresponding to the lowest error amongst all these solutions.

Algorithm 7 *Wts[] Neighbors (NET, Wts, s, τ)*

```

[Wts, Hess] = Train(NET, Wts,  $\tau$ )
evec = Eig_Vec(Hess)
Wts[ ] = NULL
for  $k = 1$  to size(evec) do
    Old_Wts = Wts
    ext_Pt = Find_Ext(NET, Old_Wts, s, evec[k])
    if (ext_Pt) then
        New_Wts = Move(NET, Old_Wts, evec[k])
        New_Wts = Train(NET, New_Wts,  $\tau$ )
        Errors = Estimate(NET, New_Wts)
        Wts[ ] = Append(Wts[ ], New_Wts, Errors)
    end if
end for
Return Wts[ ]

```

The approximate Hessian matrix obtained during the updation in the Levenberg-Marquardt method used for computing the search direction. Since there is no optimal way of obtaining the search directions, the Eigen vectors of this Hessian matrix are used as search directions. Along each search direction, the exit point is obtained by evaluating the function value along that particular direction. The step size for evaluation is chosen to be the average step size taken during the convergence of the local procedure. The function value increases initially and then starts to reduce indicating the presence of exit point on the stability boundary. *Move* function ensures that a new point (obtained from the exit point) is located in a different (neighboring) stability region. From this new initial guess, the local method is applied again

to obtain a new local optimal solution. For certain directions, an exit point might not be encountered. For these directions, the search for exit points will be stopped after evaluating the function for certain number of steps. This avoids inefficient use of resources required to search in non-promising directions.

6.7 Experimental Results

6.7.1 Benchmark Datasets

The newly proposed training method is evaluated using seven benchmark datasets taken from the UCI machine learning repository available at [16]. Since the main focus of our work is the development of TRUST-TECH based training algorithm, only simple experiments were conducted for choosing the architecture of the neural network. The hidden nodes in the hidden layer are added incrementally and the training error is computed. The final architecture is chosen with a fixed number of hidden nodes after which there is no significant improvement in the training error even when a hidden node is added. The number of nodes where the improvement in the training error is not significant is chosen as the final architecture. Table 6.2 summarizes the datasets. It gives the number of samples, input features, output classes along with the number of hidden nodes of the optimal architecture. These datasets have varying degrees of complexity in terms of sample size, output classes and the class overlaps. Here is the description of the datasets:

1. *Cancer* : This dataset contains data from cancer patients. It has 683 samples out of which 444 are benign cases and 239 are malignant cases. 9 attributes describing the tumor were used for classification.

2. *Diabetes* : This dataset gives information about patients who have some signs of diabetes according to World Health Organization criteria. Each sample has 8 real valued attributes. A total of 768 samples with 500 negative cases and 268 positive cases are available.
3. *Image* : This dataset contains images which were drawn randomly from a database of 7 outdoor images. The images were hand-segmented to create a classification for every pixel. There are 19 attributes that describe each instance (which is a 3x3 region) of a given image. The dataset contains a total of 2310 samples.
4. *Ionosphere* : This radar data was collected by a system consisting a phased array of 16 high-frequency antennas with total transmitted power on the order of 6.4 kilowatts. The targets were free electrons in the ionosphere. The dataset consists of 351 samples with 34 attributes. The classification task here is to separate good radar signals from that of the bad ones.
5. *Iris* : This dataset contains 3 classes of 50 samples each, where each class refers to a type of iris plant. It is relatively simple dataset where one class is linearly separable from the other two, but the other two have significant overlap and are not linearly separable from each other. The four attributes considered for classification are sepal length, sepal width, petal length and petal width. All attributes are measured in centimeters.
6. *Sonar* : This dataset is used for the classification of sonar signals. The task is to discriminate sonar signals bounced off a metal cylinder from those bounced off a roughly cylindrical rock. The dataset contains a total of 208 samples (111 for mines and 97 for rocks). The data set contains signals that were obtained from a variety of different aspect angles, spanning 90 degrees for the

cylinder and 180 degrees for the rock. Each pattern is a set of 60 numbers in the range 0.0 to 1.0 that represents the energy within a particular frequency band, integrated over a certain period of time.

7. *Wine* : This dataset was obtained from the results of a chemical analysis of wines derived from three different cultivars. The analysis determined the quantities of 13 constituents found in each of the three types of wines. A total of 178 samples with the following distribution (59,71,48).

Table 6.2: Summary of Benchmark Datasets.

| Dataset (\mathcal{D}) | Sample Size (Q) | Input Features (n) | Output Classes (p) | Hidden Nodes (H) | Search Variables ($(n+2)k+1$) |
|------------------------------|---------------------------|------------------------------|------------------------------|----------------------------|---------------------------------------|
| Cancer | 683 | 9 | 2 | 5 | 56 |
| Diabetes | 178 | 8 | 3 | 4 | 61 |
| Image | 2310 | 19 | 7 | 8 | 169 |
| Ionosphere | 351 | 34 | 2 | 9 | 325 |
| Iris | 150 | 4 | 3 | 3 | 19 |
| Sonar | 208 | 60 | 2 | 8 | 497 |
| Wine | 178 | 13 | 3 | 4 | 61 |

6.7.2 Error Estimation

To demonstrate the generalization capability (and hence the robustness) of the training algorithm, ten-fold cross validation is performed on each dataset. This practice of cross validation effectively removes any bias in the dataset segmentation. The use of the validation dataset allows early stopping of the local method and prevents over-fitting to a particular dataset. Essentially, each dataset is partitioned into ten folds of approximately equal size. Let these folds are denoted by T_1, T_2, \dots, T_{10} . Each time, the validation set will be T_i in which the the target labels will be deleted. The

test set is T_j for $j = (i + 1) \bmod 10$. The training set comprises of the rest of the dataset and is given by :

$$\sum_{\substack{k=1 \\ k \neq i, k \neq j}}^{10} T_k \quad (6.11)$$

The final MSE is the average of all the errors obtained across each of the ten folds. Usually, training error is much lower than the test error because the network is modeled using the training data and this data will be more accurately classified compared to the unseen test data. All the network parameters including the architecture and the set of weights are obtained using the training data. Once the final model is fixed, the accuracy on the test data will provide an estimate of the generalization capability of the network model and the training algorithm.

6.7.3 Classification Accuracy

The criteria of evaluation is given by the classification accuracy of the network model. The classification accuracy is given by the following formula :

$$\% \text{ accuracy} = \frac{\text{diff}(t(i), y(W, X))}{Q} * 100 \quad (6.12)$$

where *diff* gives the number of misclassified samples. Tables 6.3 and 6.4 shows the improvements in the train error and the test error using TRUST-TECH methodology. For effective implementation, only the best five tier-1 and corresponding tier-2 solutions were obtained using the TRUST-TECH strategy. For some of the datasets, there had been considerable improvement in the classifier performance.

Table 6.3: Percentage improvements in the classification accuracies over the training and test data using TRUST-TECH with multiple random restarts.

| Dataset | Train Error | | | Test Error | | |
|------------|-------------|------------|--------|------------|------------|--------|
| | MRS+BP | TRUST-TECH | Gain | MRS+BP | TRUST-TECH | Gain |
| Cancer | 2.21 | 1.74 | 27.01 | 3.95 | 2.63 | 50.19 |
| Image | 9.37 | 8.04 | 16.54 | 11.08 | 9.74 | 13.76 |
| Ionosphere | 2.35 | 0.57 | 312.28 | 10.25 | 7.96 | 28.77 |
| Iris | 1.25 | 1.00 | 25.00 | 3.33 | 2.67 | 24.72 |
| Diabetes | 22.04 | 20.69 | 6.52 | 23.83 | 20.58 | 15.79 |
| Sonar | 1.56 | 0.72 | 116.67 | 19.17 | 12.98 | 47.69 |
| Wine | 4.56 | 3.58 | 27.37 | 14.94 | 6.73 | 121.99 |

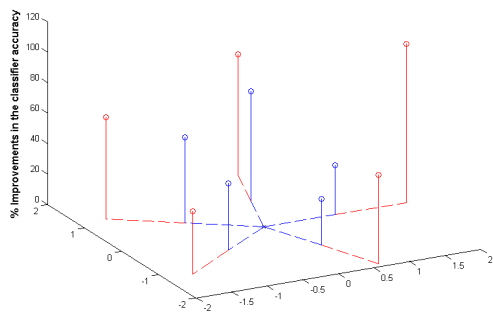
Table 6.4: Percentage improvements in the classification accuracies over the training and test data using TRUST-TECH with MATLAB initialization.

| Dataset | Train Error | | | Test Error | | |
|------------|-------------|------------|--------|------------|------------|--------|
| | NW+BP | TRUST-TECH | Gain | NW+BP | TRUST-TECH | Gain |
| Cancer | 2.25 | 1.57 | 42.99 | 3.65 | 3.06 | 19.06 |
| Image | 7.48 | 5.17 | 44.82 | 9.39 | 7.40 | 26.90 |
| Ionosphere | 1.56 | 0.92 | 69.57 | 8.67 | 6.54 | 32.57 |
| Iris | 1.33 | 0.67 | 100.00 | 3.33 | 2.67 | 25.00 |
| Diabetes | 21.41 | 19.55 | 9.53 | 23.70 | 21.09 | 12.37 |
| Sonar | 2.35 | 0.42 | 456.96 | 17.26 | 14.38 | 20.03 |
| Wine | 7.60 | 1.62 | 370.06 | 14.54 | 4.48 | 224.82 |

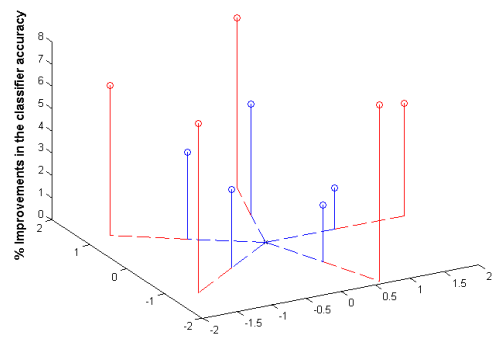
6.7.4 Visualization

The improvements of the TRUST-TECH method are demonstrated using spider web diagrams. Spiderweb diagram (shown in Fig. 6.6) is a pretentious way to demonstrate the accuracy improvements in a tier-by-tier manner. The circle in the middle of the plot represents the starting local optimal solution. The basic two dimensions are chosen arbitrarily for effective visualization and the vertical axis is the percentage improvement in the classification accuracy. Unit distances are used between the tiers and the improvements are averaged out for 10 folds. The five vertical lines surrounding the center circle are the best five local minima obtained from a tier-1 search across all folds. The tier-2 improvements are also plotted. It should be noted that the best tier-1 solution need not give the best second tier solution.

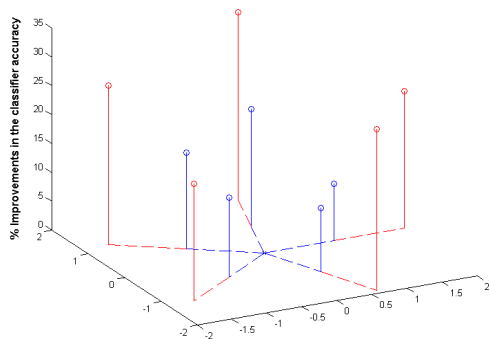
Most successful algorithms for training artificial neural networks make use of some stochastic approaches in combination with backpropagation to obtain an effective set of training parameters. Due to the limited fine-tuning capability of these algorithms, even the best solutions that they can provide are locally optimal. In this chapter, a new TRUST-TECH based method for improving the local search capability of these training algorithms is proposed. This method improves the neural network model thus allowing improved classification accuracies by providing a better set of training parameters. Because of the non-probabilistic in nature, multiple runs of our method from a given initial guess will provide exactly the same results. Different global and local methods work effectively on different datasets. The proposed TRUST-TECH based training algorithm allows the user to have the flexibility of choosing different global and local techniques for training.



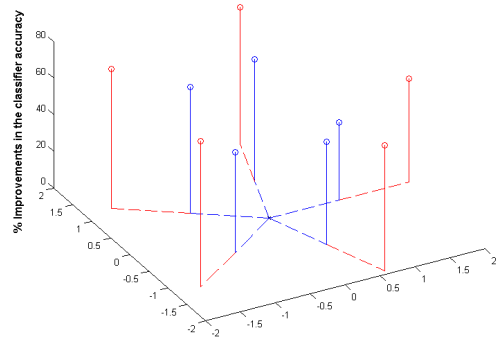
(a) Wine Dataset



(b) Diabetes Dataset



(c) Cancer Dataset



(d) Image Dataset

Figure 6.6: Spider web diagrams showing the tier-1 and tier-2 improvements using TRUST-TECH method on various benchmark datasets. The basis two axes are chosen arbitrarily and the vertical axis represents the improvements in the classifier accuracy. The distances between each tier are normalized to unity.

Chapter 7

Evolutionary TRUST-TECH

This chapter discusses the advantages of using the proposed TRUST-TECH based schemes in combination with the widely used evolutionary algorithms. The main focus of this work is to demonstrate the use of having a deterministic search compared to a stochastic one for exploring the neighborhood during the global stage. Most of the popular stochastic global optimization methods use some probabilistic neighborhood search algorithms for exploring the promising subspaces. Due to this probabilistic nature, one might obtain solutions that were already found or sometimes, even miss some promising solutions. Adding a deterministic search schemes like TRUST-TECH will help this neighborhood search to find promising solutions more effectively. First, we provide an overview of evolutionary computation, and describe the evolutionary algorithms in detail.

7.1 Overview

Research efforts on developing computational models have been rapidly growing in recent times. Amongst the several ways of developing effective computational models, evolutionary computation have become very popular. The evolutionary computational models [60] use the well-studied computational models of evolutionary processes as key elements. There are a variety of evolutionary computational models that have been proposed and studied which we will refer to as evolutionary algorithms. In simple terms, they simulate the process of evolution of the individual components via processes of selection and reproduction. These processes depend on the fitness of the individuals as defined by an environment. Evolutionary algorithms maintain a population of individuals that evolve according to rules of selection and

other genetic operators, such as recombination and mutation. Each individual in the population receives a measure of its fitness in the environment. Of all these operations mentioned above, selection is the main one that can exploit the available fitness information and mainly considers those individuals with high fitness value. Recombination and mutation perturb those individuals, providing general search strategies and heuristics for exploring the solution space. These algorithms are sufficiently complex to provide robust and powerful adaptive search mechanisms especially when the number of optimal solutions grow exponentially with the problem complexity.

Algorithm 8 Evolutionary Algorithm

Input: Initial Population, no. of iterations

Output: promising solutions

Algorithm:

$t = 0$

Initialize population $P(t)$

Evaluate $P(t)$

while not done **do**

$t = t + 1$

 parent selection $P(t)$

 recombine $P(t)$

 mutate $P(t)$

 evaluate $P(t)$

 survive $P(t)$

end while

Algorithm 8 outlines a simple evolutionary algorithm (EA). A population of individual structures is initialized and then evolved from generation to generation by repeated applications of evaluation, selection, recombination, and mutation. The initial and final population size N is generally constant in an evolutionary algorithm. In other words, the initial fixed number of solutions are chosen and they evolve into more and more promising solutions as the time progresses. In our case, the time will be indicated in terms of the number of iterations taken to evolve.

An evolutionary algorithm typically initializes its population randomly, although some apriori information and domain specific knowledge can also be used to obtain promising starting points. If promising starting values are chosen then good solutions can be obtained in a fewer iterations compared to the number of iterations taken for an algorithm that was started with random points. Evaluation measures the fitness of each individual according to its worth in the environment. The complexity of the evaluation process is highly problem dependent. It may be as simple as computing a fitness function or as complex as running an elaborate simulation. Selection is usually performed in two steps, parent selection and survival. Parent selection decides who becomes parents and how many children the parents can have. Children are created via recombination, which exchanges information between parents. This recombination mechanism is the vital component of an evolutionary algorithm because it will provide new individuals in the environment. Mutation is an important step in the algorithm which perturbs the children to obtain minor changes in the solution. The children are then evaluated using the fitness function and the survive step decides who survives in the population.

7.2 Variants of Evolutionary Algorithms

The algorithm described above is the basic backbone of a simple evolutionary algorithm. Several variants and improvements for this basic architecture have been proposed in the literature and is currently an active topic of research. We will now discuss the three most popular and well-studied variations of evolutionary algorithms. The three methodologies are : 1. Evolutionary programming [57], 2. Evolution strategies [12] and 3. Genetic algorithms [77, 62]. At a higher level, all these methods implement an evolutionary algorithm but the details of their imple-

mentation are completely different. They differ in the choice of problem representation, types of selection mechanism, forms of genetic operators, and performance measures.

Evolutionary programming (EP), developed by Fogel et al. [57] traditionally used problem representations that are specific to the application domain. For example, in real-valued optimization problems, the individuals within the population are real-valued vectors. Similarly, ordered lists are used for traveling salesman problems, and graphs for applications with finite state machines. The basic scheme is similar to an evolutionary algorithm except that the recombination is generally not performed since the forms of mutation used are adaptive. These mutations are quite flexible that can produce perturbations that are similar to recombination, if desired. One of the heavily studied aspects of this approach is the extent to which an evolutionary algorithm is affected by its choice of the perturbation rates used to produce variability and the novelty in evolving populations.

Genetic algorithms (GAs), developed by Holland [77, 62] are arguably the most well known form of evolutionary algorithms. They have been traditionally used in a more domain independent setting, namely, bit-strings. Those individuals with higher relative fitness are more likely to be selected as parents. N children are created via recombination from the N parents. The N children are then mutated and the best survivors will replace the N parents in the population. It should be noted that there is a strong emphasis on mutation and crossover.

The third category is an evolution strategy (ES). It follows the basic EA architecture and the number of children created is usually greater than N . After initialization and evaluation, individuals are selected uniformly random. Survival is deterministic and allows the N best children to survive. Like EP, considerable effort is made on adapting mutation as the algorithm runs by allowing each variable within an

individual to have an adaptive mutation rate that is normally distributed with a zero expectation. Unlike EP, however, recombination does play an important role in evolution strategies, especially in adapting mutation.

These three approaches (EP, ESs, and GAs) have inspired an increasing amount of research and development of new forms of evolutionary algorithms for use in specific problem solving contexts. In this chapter, we focus on the aspect of mutations and improve their performance. The mutations usually allow us to obtain individuals with minor changes. In terms of the parameter space, mutations will merely perturb the solutions to obtain new solutions that might potentially be more promising than the original solutions. Most of the stochastic optimization methods which try to obtain global optimal solutions perform some kind of neighborhood search. In the popularly used simulated annealing technique [88], the neighborhood search is performed by local moves.

However, performing the neighborhood search in a stochastic manner will have many problems like:

- One might not know the extent to which the mutation has to be performed on an individual.
- It is difficult to understand the locations of the phenotype where the mutation should occur.

For both the problems mentioned above, one can realize that performing a more systematic neighborhood search for obtaining better solutions will help in performing this mutation step in a better manner. In this chapter, we replace this concept of mutation by other local search techniques. Our first model will use a local refinement strategy instead of mutations. Our second model will use the TRUST-TECH based neighborhood search for searching the nearby local optimal solutions. Of course,

it might incur some additional computational cost to perform these sophisticated neighborhood search strategies, but they will almost surely have the advantage of performing a nearly perfect local search which can have the potential of reducing the total number of iterative steps taken for convergence. In other words, we alter the mutation aspect of the algorithm so that it can result in the reduction of the total number of generations (computational time) required by the complete algorithm.

7.3 Evolutionary TRUST-TECH Algorithm

Algorithm 9 Evolutionary Algorithm with Local Refinement

Input: Initial Population, no. of iterations

Output: promising solutions

Algorithm:

$t = 0$

Initialize population $P(t)$

Evaluate $P(t)$

while not done **do**

$t = t + 1$

 parent selection $P(t)$

 recombine $P(t)$

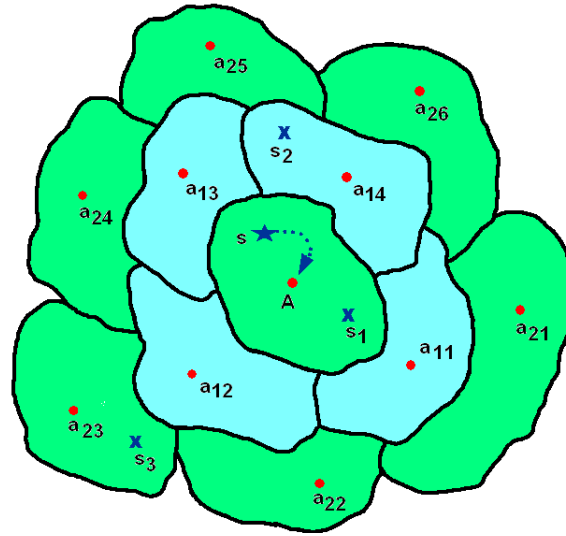
 Local Refinement $P(t)$

 evaluate $P(t)$

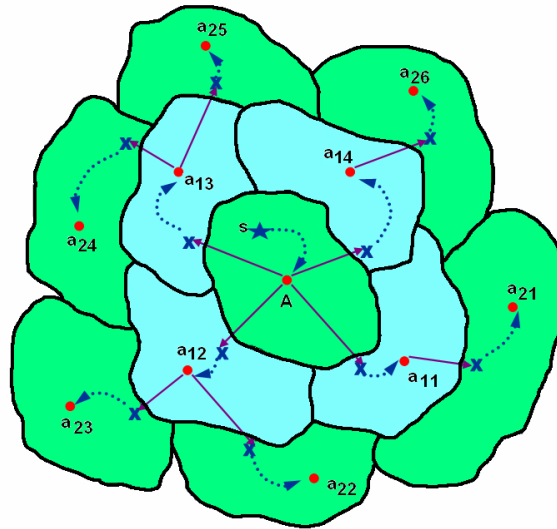
 survive $P(t)$

end while

First, we will describe the evolutionary algorithm with local refinement strategy. In this model, the mutation is replaced with a local refinement strategy. This can be either in discrete space or a continuous space. The important point here is to ensure that a solution which is better than the existing solution is always obtained. The local refinement strategy will obtain the corresponding local optimal solution. In the case of mutations, there is no guarantee that the new point has a higher fitness function value than the original point. To ensure a better fitness value, one can apply this local refinement which is a greedy procedure to obtain a better solution.



(a) Evolutionary Local Refinement



(b) Evolutionary TRUST-TECH

Figure 7.1: Topology of the nonlinear surface discussed in the introduction chapter. The initial point obtained after the recombination is ‘s’. The original concept of mutations will randomly perturb the point in the parameter space and obtain different points (s_1, s_2 and s_3). Two different evolutionary models (a) Local refinement strategy where the local optimization method is applied with ‘s’ as initial guess to obtain ‘A’. (b) Evolutionary TRUST-TECH methodology where the surrounding solutions are explored systematically after the recombination.

All other aspects of the evolutionary algorithm remain unchanged. Algorithm 9 describes the evolutionary algorithm with local refinement strategy.

Algorithm 10 Evolutionary TRUST-TECH

Input: Initial Population, no. of iterations

Output: promising solutions

Algorithm:

$t = 0$

Initialize population $P(t)$

Evaluate $P(t)$

while not done **do**

$t = t + 1$

 parent selection $P(t)$

 recombine $P(t)$

 TRUST-TECH $P(t)$

 evaluate $P(t)$

 survive $P(t)$

end while

Fig 7.1 clearly demonstrates both the proposed models. Let ‘s’ denote the point obtained after the recombination process. Mutation will randomly perturb this point to obtain another point. It can be either s_1, s_2 or s_3 . In all these cases, there is no guarantee that the new point has a higher fitness function value. However, applying a local refinement strategy, ‘s’ will converge to ‘A’ which has either equal or higher fitness value. The mutations are completely replaced using this local refinement method. There might be other promising solutions with higher fitness function values. Hence, in the evolutionary TRUST-TECH mechanism, this local refinement strategy is replaced by the TRUST-TECH methodology. The TRUST-TECH used in this step is identical to the algorithm presented in the previous chapter. The only difference being that the starting point is obtained from the recombination operator. Using the TRUST-TECH strategy, neighborhood solutions are obtained and the best one is retained for the next iteration in the evolutionary process.

7.4 Experimental Results

For our implementation, we started with 10 initial points and refined them through the evolutionary process. During each selection stage, two parents are chosen randomly and recombined. Amongst the four individuals (2 parents and 2 children), only two individuals with higher fitness function value are retained and the other two are discarded. The population is refined for 100 generations. Each time only 10 recombinations are performed. For the local refinement strategy, we used the Levenberg-Marquardt method. For the TRUST-TECH implementation, only the first-tier solutions were obtained for computational efficiency. Amongst all the tier one solutions, the best solution is chosen and all the rest are discarded.

Table 7.1: Results of Evolutionary TRUST-TECH model

| Dataset | EA | EA+LR | EA+TRUST-TECH |
|---------|--------|---------|---------------|
| Pima | 0.1642 | 0.1432 | 0.1367 |
| Cancer | 0.0193 | 0.01654 | 0.01423 |
| Wine | 0.0775 | 0.0651 | 0.0472 |

Table 7.1 reports the results of our approach. For the Local refinement and TRUST-TECH strategies, we used only 50 generations. The MSE over the training has been reported. One of the main observations from our experiments is that it is not necessary that every iteration during the local refinement will yield a better score than the original model. For example, if the local refinement strategy yields a better score at a particular generation, then its improvement during the next few generations might not be significant. This applies to even TRUST-TECH method as well. However, as it is evident from the results, there can be significant improvements in terms of MSE and this can be achieved at fewer number of iterations of the evolutionary algorithm.

7.5 Parallel Evolutionary TRUST-TECH

From the results shown in the previous section, one can see that TRUST-TECH can help to provide faster convergence of the traditional evolutionary algorithms. The proposed framework can be easily extended to work on parallel machines. This will work in a similar manner to any other parallel evolutionary algorithm. The local refinement and neighborhood strategies proposed here can work independent for each of the selection and recombination individuals chosen. Two individuals are selected for recombination and then the local refinement can take place on a different local machine. This way, all the local refinements can be performed in different machines and the final results can be evaluated and the next generation parents are chosen in a centralized machine. This centralized machine (or the server) can obtain the results from all the local machines which perform the neighborhood search.

Chapter 8

Conclusion and Future Work

This chapter concludes our discussion and highlights the most important contributions of this thesis. It also discusses the future research directions that one might want to pursue using the models presented in this thesis.

8.1 Conclusion

In this thesis work, we develop TRUST-TECH based methods for various problems related to areas of heuristic search, optimization and learning. We demonstrate the applicability and effectiveness of these methods for practical and high-dimensional nonlinear optimization problems. One of the main ideas of this framework is to transform the original optimization problem into a dynamical system with certain properties and obtain more useful information about the nonlinear surface via the dynamical and topological properties of the dynamical system.

In Chapter 2, we apply a stability boundary following procedure for obtaining saddle points on various potential energy surfaces that arise in the field of computational biology and computational chemistry. To find a saddle point, following the stability boundary is computationally more efficient than directly searching for a saddle point from a given local minimum. This algorithm works deterministically and requires very few user-specific parameters. The procedure has been successfully tested on a 525-dimensional problem. For energy surfaces that show symmetric behaviour, we propose a simplified version of this procedure.

Nonlinear optimization problems arise in several domains in science and engineering. Several algorithms had been proposed in the optimization literature for solving these problems efficiently. Typically, optimization methods can be classified

into two categories: (1) Global methods and (2) Local methods. Global methods are powerful stochastic methods that search the entire parameter space and obtain promising regions. Local methods, on the other hand, are deterministic methods and usually converge to a locally optimal solution that is nearest to a given initial point. There is a clear gap between these two methods and there is a need for a method that can search in the neighborhood regions.

The proposed TRUST-TECH based methods systematically search for neighborhood local optimal solutions by exploring the dynamic and geometric characteristics of stability boundaries of a nonlinear dynamical system corresponding to the nonlinear function of interest. This framework consists of three stages namely: (i) Global stage, (ii) Local stage and (iii) Neighborhood-search stage. These methods have been successfully used for various machine learning problems and are demonstrated in the context of both supervised (training artificial neural networks in Chapter 6) and unsupervised (expectation maximization in Chapter 3) machine learning problems. In both these scenarios, obtaining a global optimal solution in the parameter space corresponds to exploiting the complete potential of the given model. More complicated models might achieve the same function value but, they tend to lose the generalization capability. Our methods are tested on both synthetic and real datasets and the advantages of using this TRUST-TECH based framework are clearly manifested. The improvements in the performance of the expectation maximization algorithm are demonstrated in the context of mixture modeling and general likelihood problems such as the motif finding problem.

This framework not only reduces the sensitivity to initialization, but also allows the flexibility for the practitioners to use various global and local methods that work well for a particular problem of interest. Thus, it can act as a flexible interface between the local method (EM) and other global methods. This interface plays

a vital role in problems where the functions optimized by the global method and the local method are not the same. For example, in the motif finding problem, a global method is used in the discrete space and a local method is used in the continuous space. The points obtained as a result of the global method need not be in the convergence region (of the local method) of the most promising solutions. In such cases, applying tier-by-tier search can significantly improve the quality of the solutions as demonstrated by the results of finding optimal motifs in Chapter 4. Also, this framework has the potential to work with any global and local method by treating them as a black-box (without knowing their detailed algorithms).

In Chapter 5, we propose a novel smoothing framework in the context of Gaussian mixture models. The proposed component-wise kernel smoothing approach can reduce the number of local maxima on the log-likelihood surface and can potentially obtain promising initial points for the EM algorithm. Performing component-wise smoothing will maintain the structure of the log-likelihood function and hence the EM can be applied directly with a few modifications.

Frameworks for combining TRUST-TECH with other hierarchical stochastic algorithms such as smoothing algorithms and evolutionary algorithms are also proposed and tested on various datasets. In Chapter 7, it has been shown that the neighborhood-search stage can be incorporated into the global stage in order to improve the performance of the global stage in terms of the quality of the solutions and the speed of convergence.

8.2 Future Work

The algorithm proposed for finding saddle points can be applied to the problem of finding pseudo-native like structures and their corresponding transition states during the protein folding process. The TRUST-TECH based Expectation-Maximization algorithm can be extended to other widely used EM related problems for the family of probabilistic graphical models such as k-means clustering, training Hidden Markov Models, Mixture of Factor Analyzers, Probabilistic Principal Component Analysis, Bayesian Networks etc. Extension of these techniques to Markov Chain Monte Carlo methods (like Gibbs sampling) is also feasible. Several real-world applications such as image segmentation, gene finding, speech processing and text classification can benefit significantly from these methods. Extensions to constrained optimization problems appears to be a promising direction as well. Different global methods and local solvers can be used along with TRUST-TECH framework to study the flexibility of this framework. For machine learning problems, automatically choosing a model is an important and difficult problem to tackle. TRUST-TECH based methods are generic enough to incorporate any model selection criterion (such as Akaike information criterion (AIC) and Bayesian information criterion (BIC)) into the objective function. Basically, this term is added in the objective function in order to penalize complex models.

As a continuation of the smoothing work, the effects of convolving Gaussian components with other kernels must be investigated. Efficient algorithms for choosing the smoothing parameter automatically based on the available data can be developed. Though applied for Gaussian mixture models in this paper, convolution based smoothing strategies can be treated as powerful optimization tools that can enhance the search capability significantly. The novel neural network training algorithm can

be extended to the problem of simultaneously deciding the architecture and the training parameters. This problem is characterized by a combination of a set of discrete (for architecture) and continuous (for training) variables. Its performance on large scale applications like character recognition, load forecasting etc. must be tested. The power of evolutionary TRUST-TECH has been demonstrated only in the case of training neural networks problem. The algorithm is not specific to neural network and can be demonstrated for general nonlinear programming problems in the future.

In the context of optimization, TRUST-TECH can also help the local solvers to converge to desired optimal solution. Especially, when a local method is applied to a point near the stability boundary, it might diverge or might converge to a different local optimal solution at a very slow rate. TRUST-TECH methodology can be used to integrate the system and obtain new points that are much closer to the desired local optimal solution. Using this new point as initial condition, there are higher chances that the local method will converge to the desired local optimal solution. The number of integration steps required and the step size for integration are research topics that can be pursued in the future.

Most of the problems that were discussed in this thesis primarily focus on finding critical points. Usually, the trajectory of convergence is not of much importance in optimization and machine learning problems. Transformation of the nonlinear objective function to its corresponding gradient system is proposed in this thesis. One can extend this work by scaling the dynamical system so that it preserves the location of all the critical points. Scaling can modify the trajectories significantly and can potentially improve the speed of convergence. The characterization of the scaling factor required for various systems is a potential research topic.

Though proposed for a few global methods like smoothing and evolutionary algo-

rithms, TRUST-TECH can be used to improve the performance of any hierarchical stochastic global optimization method. For example, multi-level optimization methods are popular tools that are powerful in solving problems in various domains of engineering. The basic idea is to reduce the dimensionality of the problem and obtain a solution and carefully trace back this solution to the original problem. During this trace back procedure, local methods might be inefficient in tracing the solution accurately and the use of TRUST-TECH based tier-by-tier search can help in exploring the neighborhood and obtaining the global optimal solution accurately.

BIBLIOGRAPHY

- [1] S. Amato, B. Apolloni, P. Caporali, U. Madiesani, and A. Zanaboni. Simulated annealing approach in backpropagation. *Neurocomputing*, 3(5):207–220, 1991.
- [2] M. Avriel. *Nonlinear Programming: Analysis and Methods*. Dover Publishing, 2003.
- [3] T. Back. *Evolutionary Algorithms in Theory and Practice - Evolution Strategies, Evolutionary Programming, Genetic Algorithms*. Oxford University Press, 1996.
- [4] T. Bailey and C. Elkan. Fitting a mixture model by expectation maximization to discover motifs in biopolymers. *The First International Conference on Intelligent Systems for Molecular Biology*, pages 28–36, 1994.
- [5] J. Baker. An algorithm for the location of transition states. *Journal of Computational Chemistry*, 7(4):385 – 395, 1986.
- [6] J. D. Banfield and A. E. Raftery. Model-based gaussian and non-gaussian clustering. *Biometrics*, 49(3):803– 821, 1993.
- [7] Y. Barash, G. Bejerano, and N. Friedman. A simple hyper-geometric approach for discovering putative transcription factor binding sites. *Proc. of First International Workshop on Algorithms in Bioinformatics*, 2001.
- [8] G.T. Barkema and N. Mousseau. Identification of relaxation and diffusion mechanisms in amorphous silicon. *Physics Review Letters*, 77:43–58, 1996.
- [9] T. Battiti. First and second-order methods for learning: Between steepest descent and Newton’s method. *Neural Computation*, 4(2):141–166, 1992.
- [10] L. Baum, T. Petrie, G. Soules, and N. Weiss. A maximization technique occurring in the statistical analysis of probabilistic functions of markov chains. *Annals of Mathematical Statistics*, 41(1):164–171, 1970.
- [11] J. Beran and G. Mazzola. Visualizing the relationship between two time series by hierarchical smoothing. *Journal of Computational and Graphical Statistics*, 8(2):213–238, 1999.
- [12] H.-G. Beyer and H.-P. Schwefel. Evolution strategies: A comprehensive introduction. *Journal of Natural Computing*, 1(1):3–52, 2002.
- [13] J. A. Bilmes. A gentle tutorial of the EM algorithm and its application to parameter estimation for gaussian mixture and hidden markov models. Technical report, U. C. Berkeley, April 1998.
- [14] C.M. Bishop. *Neural Networks for Pattern Recognition*. Oxford University Press, 1995.

- [15] A. Blake and A. Zisserman. *Visual reconstruction*. MIT Press, Cambridge, MA, 1987.
- [16] C.L. Blake and C.J. Merz. UCI repository of machine learning databases. <http://www.ics.uci.edu/mllearn/MLRepository.html>, University of California, Irvine, Dept. of Information and Computer Sciences, 1998.
- [17] K. Blekas, D. Fotiadis, and A. Likas. Greedy mixture learning for multiple motif discovery in biological sequences. *Bioinformatics*, 19(5):607–617, 2003.
- [18] G. Bokinsky, D. Rueda, V. K. Misra, A. Gordus, M. M. Rhodes, H. P. Babcock, N. G. Walter, and X. Zhuang. Single-molecule transition-state analysis of rna folding. *Proceedings of the National Academy of Sciences*, 100:9302–9307, 2003.
- [19] J. Branke. Evolutionary algorithms for neural network design and training. *Proceedings of the 1st Nordic Workshop on Genetic Algorithms and its Applications*, 3:145–163, 1995.
- [20] J. Buhler and M. Tompa. Finding motifs using random projections. *Proceedings of the fifth annual international conference on Research in computational molecular biology*, pages 69–76, 2001.
- [21] C. J. C. Burges. A tutorial on support vector machines for pattern recognition. *Data Mining and Knowledge Discovery*, 2:121–167, 1998.
- [22] R. H. Byrd, E. Eskow, A. van der Hoek, R. B. Schnabel, C.S. Shao, and Z. H. Zou. Global optimization methods for protein folding problems. In *Global minimization of nonconvex energy functions: molecular conformation and protein folding*, pages 29–39. Amer. Math. Soc., 1996.
- [23] C. Carson, S. Belongie, H. Greenspan, and J. Malik. Blobworld: Image segmentation using expectation-maximization and its application to image querying. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 24(8):1026–1038, 2002.
- [24] M. Celis, J. E. Dennis, and R. A. Tapia. *Numerical Optimization*, chapter A trust region strategy for nonlinear equality constrained optimization, pages 71–82. Philadelphia: SIAM, 1985.
- [25] B.C. Cetin, J. Barhen, and J.W. Burdick. Terminal repeller unconstrained subenergy tunnelling (TRUST) for fast global optimization. *Journal Optimization Theory and Applications*, 77(1):97–126, 1993.
- [26] C. Charalambous. Conjugate gradient algorithm for efficient training of artificial neural networks. *IEEE Proceedings*, 139(3):301–310, 1992.
- [27] S. F. Chen and J. T. Goodman. An empirical study of smoothing techniques for language modeling. In *In Proceedings of the 34th Annual Meeting of the ACL*, pages 310–318, June 1996.

- [28] V. Cherkassky and F. Mulier. *Learning from Data: Concepts, Theory, and Methods*. John Wiley and Sons, 1998.
- [29] H. D. Chiang. *Dynamical method for obtaining global optimal solution of general nonlinear programming problems*. United States Patent, 2002.
- [30] H. D. Chiang and J. Lee. *Heuristic techniques and its applications to Electric Power Systems*, chapter TRUST-TECH based Methodology for nonlinear optimization. IEEE Press, 2007.
- [31] H.D. Chiang and C.C. Chu. A systematic search method for obtaining multiple local optimal solutions of nonlinear programming problems. *IEEE Transactions on Circuits and Systems : I Fundamental Theory and Applications*, 43(2):99–109, 1996.
- [32] H.D. Chiang and A. Fekih-Ahmed. Quasi-stability regions of nonlinear dynamical systems:Theory. *IEEE Transactions on Circuits and Systems*, 43:627–635, 1996.
- [33] H.D. Chiang, M. Hirsch, and F. F. Wu. Stability regions of nonlinear autonomous dynamical systems. *IEEE Transactions on Automatic Control*, 33:16–27, 1988.
- [34] C.K. Chu, I. Glad, F. Godtlielsen, and J.S. Marron. Edge-preserving smoothers for image processing. *Journal of the American Statistical Association*, 93(442):526–541, 1998.
- [35] R. Czerminski and R. Elber. Reaction path study of conformational transitions in flexible systems: Applications to peptides. *Journal of Chemical Physics*, 92:5580–5601, 1990.
- [36] S. Dasgupta and L. J. Schulman. A two-round variant of em for gaussian mixtures. *Proceedings of the 16th Conference in Uncertainty in Artificial Intelligence*, pages 152–159, 2000.
- [37] A. P. Dempster, N. A. Laird, and D. B. Rubin. Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society Series B*, 39(1):1–38, 1977.
- [38] M.J.S. Dewar, E.F. Healy, and J.J.P. Stewart. Location of transition states in reaction mechanisms. *Journal of Chemical Society Faraday Transactions*, 80:227–233, 1984.
- [39] I. Diener. *Handbook of Global Optimization: Nonconvex Optimization and Its Applications*, chapter Trajectory Methods in Global Optimization, pages 649–668. Dordrecht, Netherlands: Kluwer, 1995.

- [40] K.A. Dill, A.T. Phillips, and J. B. Rosen. Protein structure prediction and potential energy landscape analysis using continuous global minimization. *Proceedings of the first annual international conference on Computational molecular biology*, pages 109 – 117, 1997.
- [41] C. Dong and I. Michel. A global optimization approach to resource allocation problems with consideration of economy of scale. *Journal of Information and Decision Technologies*, 19(5):257–266, 1994.
- [42] M. Dorigo and T. Sttzle. *Ant Colony Optimization*. MIT Press, 2004.
- [43] J.P.K. Doye and D.J. Wales. Saddle points and dynamics of lennard-jones clusters, solids and supercooled liquids. *Journal of Chemical Physics*, 116:3777–3788, 2002.
- [44] R. O. Duda, P. E. Hart, and D. G. Stork. *Pattern Classification*. Wiley, New York, 2001.
- [45] D. M. Dunlavy, D. P. O’leary, D. Klimov, and D. Thirumalai. Hope: A homotopy optimization method for protein structure prediction. *Journal of Computational Biology*, 12(10):1275–1288, 2005.
- [46] R. Durbin, S.R. Eddy, A. Krogh, and G. Mitchison. *Biological Sequence Analysis : Probabilistic Models of Proteins and Nucleic Acids*. Cambridge University Press, 1999.
- [47] B. Eckhardt. Irregular scattering. *Physica*, D33:89–98, 1988.
- [48] S. R. Eddy. Profile hidden markov models. *Bioinformatics*, 14(9):755–763, 1998.
- [49] G. Elidan, M. Ninio, N. Friedman, and D. Schuurmans. Data perturbation for escaping local maxima in learning. *In Proceedings of the Eighteenth National Conference on Artificial Intelligence*, pages 132 – 139, 2002.
- [50] B. Erman, I. Bahar, and R. L. Jernigan. Equilibrium states of rigid bodies with multiple interaction sites. application to protein helices. *Journal of Chemical Physics*, 107:2046–2059, 1997.
- [51] E. Eskin. From profiles to patterns and back again: A branch and bound algorithm for finding near optimal motif profiles. *Proceedings of the eighth annual international conference on Research in computational molecular biology*, pages 115–124, 2004.
- [52] E. Eskin and P. Pevzner. Finding composite regulatory patterns in DNA sequences. *Tenth International Conference on Intelligent Systems for Molecular Biology*, pages 354–363, 2002.

- [53] W. R. Esposito and C. A. Floudas. Global optimization for the parameter estimation of differential-algebraic systems. *Industrial and engineering chemistry research*, 39(5):1291–1310, 2000.
- [54] F. Glover and M. Laguna. *Modern Heuristic Techniques for Combinatorial Problems*, chapter Tabu Search. John Wiley and Sons Inc., 1993.
- [55] M. Figueiredo and A.K. Jain. Unsupervised learning of finite mixture models. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 24(3):381–396, 2002.
- [56] S. Fischer and M. Karplus. Conjugate peak refinement : an algorithm for finding reaction paths and accurate transition states in systems with many degrees of freedom. *Chemical Physics Letters*, 192:252–261, 1992.
- [57] L. J. Fogel, A. J. Owens, and M. J. Walsh. *Artificial Intelligence through Simulated Evolution*. John Wiley, 1966.
- [58] C.M. Foley and D. Schinler. Automated design of steel frames using advanced analysis and object-oriented evolutionary computation. *Journal of Structural Engineering, American Society of Civil Engineers*, 129(5):648–656, 2003.
- [59] W. Forster. *Handbook of Global Optimization: Nonconvex Optimization and Its Applications*, chapter Homotopy Methods, pages 669–750. Dordrecht, Netherlands: Kluwer, 1995.
- [60] A. S. Fraser. Simulation of genetic systems by automatic digital computers. i. introduction. *Australian Journal of Biological Sciences*, 10:484–491, 1957.
- [61] J. H. Friedman, T. Hastie, and R. Tibshirani. Additive logistic regression: A statistical view of boosting. *Annals of Statistics*, 28(2):337–407, 2000.
- [62] D.E. Goldberg. *Genetic Algorithms in Search, Optimization and Machine Learning*. Addison-Wesley Longman Publishing Co. Inc., 1989.
- [63] M. Gori and A. Tesi. On the problem of local minima in backpropagation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 14(1):76–86, 1992.
- [64] D. Gorse, A.J. Shepherd, and J.G. Taylor. The new era in supervised learning. *Neural Networks*, 10(2):343–352, 1997.
- [65] M. T. Hagan and M. Menhaj. Training feedforward networks with the marquart algorithm. *IEEE Transactions on Neural Networks*, 5(6):989–993, 1994.
- [66] E. Hansen. *Global Optimization Using Interval Analysis*. Dekker, New York, 1992.
- [67] T. Hastie and R. Tibshirani. Discriminant analysis by Gaussian mixtures. *Journal of the Royal Statistical Society series B*, 58:158–176, 1996.

- [68] S. Haykin. *Neural Networks: A Comprehensive Foundation*. Prentice Hall, 1999.
- [69] J. He, A. Verstak, L. T. Watson, T. S. Rappaport, C. R. Anderson, N. Ramakrishnan, C. A. Shaffer, W. H. Tranter, K. Bae, and J. Jiang. Global optimization of transmitter placement in wireless communication systems. *Proceedings of the 2002 High Performance Computing Symposium (HPC'02)*, pages 328–333, 2002.
- [70] D. Heidrich, W. Kliesch, and W. Quapp. *Properties of Chemical Interesting Potential Energy Surfaces*. Springer, Lecture Notes in Chemistry 56, Berlin, 1991.
- [71] T. Helgaker. Transition state optimizations by trust-region image minimization. *Chemical Physics Letters*, 182(5):503–510, 1991.
- [72] G. Henkelman, G. Johansson, and H. Jonsson. A dimer method for finding saddle points on high dimensional potential surfaces using only first derivatives. *Journal of Chemical Physics*, 111:7010–7022, 1999.
- [73] G. Henkelman, G. Johansson, and H. Jonsson. Methods for finding saddle points and minimum energy paths. *Progress on Theoretical Chemistry and Physics*, 111:269–300, 2000.
- [74] G. Henkelman, B.P. Uberuaga, and H. Jonsson. A climbing image nudged elastic band method for finding saddle points and minimum energy paths. *Journal of Chemical Physics*, 113:9901–9904, 2000.
- [75] G. Hertz and G. Stormo. Spelling approximate or repeated motifs using a suffix tree. *Lecture Notes in Computer Science*, 1380:111–127, 1999.
- [76] G. E. Hinton. How neural networks learn from experience. *Scientific American*, 267(3):144–151, 1992.
- [77] John Holland. *Adaptation in Natural and Artificial Systems*. University of Michigan Press, 1975.
- [78] R. Horst and P.M. Pardalos (eds.). *Handbook of Global Optimization*. Kluwer, Dordrecht, 1995.
- [79] R. Horst and H. Tuy. *Global Optimization: Deterministic Approaches*. Berlin: Springer-Verlag, 1996.
- [80] D. Ingman and Y. Merlis. Local minimization escape using thermodynamic properties of neural networks. *Neural Networks*, 4(3):395–404, 1991.
- [81] I.V. Ionova and E.A. Carter. Ridge method for finding saddle points on potential energy surfaces. *Journal of Chemical Physics*, 98:6377–6386, 1993.

- [82] A. K. Jain and R. C. Dubes. *Algorithms for Clustering Data*. Prentice Hall, Englewood Cliffs, NJ, 1988.
- [83] H. Jonsson, G. Mills, and K. W. Jacobsen. *Classical and Quantum Dynamics in Condensed Phase Simulations* ed. B. J. Berne, G. Ciccotti and D. F. Coker, chapter Nudged Elastic Band Method for Finding Minimum Energy Paths of Transitions. World Scientific, 1998.
- [84] G. Karypis and V. Kumar. A fast and high quality multilevel scheme for partitioning irregular graphs. *SIAM Journal on Scientific Computing*, 20(1):359–392, 1999.
- [85] U. Keich and P. Pevzner. Finding motifs in the twilight zone. *Bioinformatics*, 18:1374–1381, 2002.
- [86] Y.G. Khait, A.I. Panin, and A.S. Averyanov. Search for stationary points of arbitrary index by augmented hessian method. *International Journal of Quantum Chemistry*, 54:329–336, 1995.
- [87] K.J. Kim. Toward global optimization of case-based reasoning systems for financial forecasting. *Applied Intelligence*, 21(3):239–249, 2004.
- [88] S. Kirkpatrick, C.D. Gelatt, and M.P. Vecchi. Optimization by simulated annealing. *Science*, 220(4598):671–680, 1983.
- [89] L. Lastras. Continuation methods and saddle points : a new framework. Master’s thesis, Cornell University, 1998.
- [90] C. Lawrence, S. Altschul, M. Boguski, J. Liu, A. Neuwald, and J. Wootton. Detecting subtle sequence signals: a gibbs sampling strategy for multiple alignment. *Science*, 262:208–214, 1993.
- [91] J. Lee and H.D. Chiang. A dynamical trajectory-based methodology for systematically computing multiple optimal solutions of general nonlinear programming problems. *IEEE Transactions on Automatic Control*, 49(6):888 – 899, 2004.
- [92] A. V. Lehmen, E. G. Paek, P. F. Liao, A. Marrakchi, and J. S. Patel. Factors influencing learning by backpropagation. *Proceedings of IEEE International Conference on Neural Networks*, pages 335–341, 1988.
- [93] F.H.F. Leung, H.K. Lam, S.H. Ling, and P.K.S. Tam. Tuning of the structure and parameters of a neural network using an improved genetic algorithm. *IEEE Transactions on Neural Networks*, 14(1):79–88, 2003.
- [94] J. Q. Li. *Estimation of Mixture Models*. PhD thesis, Department of Statistics, Yale University, 1999.
- [95] V. Litovski and M. Zwolinski. *VLSI Circuit Simulation and Optimization*. Chapman and Hall, 1997.

- [96] D. Marquardt. An algorithm for least-squares estimation of nonlinear parameters. *SIAM Journal on Applied Mathematics*, 11:431–441, 1963.
- [97] A. M. Martnez and J. Vitri. Learning mixture models using a genetic version of the EM algorithm. *Pattern Recognition Letters*, 21(8):759–769, 2000.
- [98] G. McLachlan and T. Krishnan. *The EM Algorithm and Extensions*. John Wiley and Sons, New York, 1997.
- [99] G. McLachlan and D. Peel. *Finite Mixture Models*. John Wiley and Sons, 2000.
- [100] G. J. McLachlan and K. E. Basford. *Mixture models: Inference and applications to clustering*. Marcel Dekker, New York, 1988.
- [101] C. Merlo, Ken A. Dill, and T. R. Weikl. Phi values in protein folding kinetics have structural and energetic components. *Proceedings of the National Academy of Sciences*, 102(29):10171–10175, 2005.
- [102] R.M. Minyaev, W. Quapp, G. Subramanian, P.V.R. Schleyer, and Y. Ho. Internal conrotation and disrotation in H₂BCH₂BH₂ and diborylmethane 1,3 H exchange. *Journal of Computational Chemistry*, 18(14):1792–1803, 1997.
- [103] R.A. Miron and K.A. Fichthorn. The step and slide method for finding saddle points on multidimensional potential surfaces. *Journal of Chemical Physics*, 115(19):8742–8747, 2001.
- [104] A. F. Moller. A scaled conjugate gradient algorithm for fast supervised learning. *Neural Networks*, 6(4):525–533, 1993.
- [105] K. Muller and L. D. Brown. Location of saddle points and minimum energy paths by a constrained simplex optimization procedure. *Theoret. Chim. Acta*, 53:75–93, 1979.
- [106] R. M. Neal and G. E. Hinton. A new view of the EM algorithm that justifies incremental, sparse and other variants. In M. I. Jordan, editor, *Learning in Graphical Models*, pages 355–368. Kluwer Academic Publishers, 1998.
- [107] E. Neria, S. Fischer, and M. Karplus. Simulation of activation free energies in molecular systems. *Journal of Chemical Physics*, 105(5):1902–1921, 1996.
- [108] D. Nguyen and B. Widrow. Improving the learning speed of 2-layer neural networks by choosing initial values of the adaptive weights. *Proceedings of the International Joint Conference on Neural Networks*, 1990.
- [109] K. Nigam, A. McCallum, S. Thrun, and T. Mitchell. Text classification from labeled and unlabeled documents using EM. *Machine Learning*, 39(2-3):103 – 134, 2000.
- [110] J. Nocedal and S. Wright. *Numerical Optimization*. Springer, New York, 1999.

- [111] S. Osowski, P. Bojarczak, and M. Stodolski. Fast second order learning algorithm for feedforward multilayer neural networks and its application. *Neural Networks*, 9(9):1583–1596, 1996.
- [112] F. Pernkopf and D. Bouchaffra. Genetic-based EM algorithm for learning gaussian mixture models. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 27(8):1344–1348, 2005.
- [113] P. Pevzner. *Computational Molecular Biology - an algorithmic approach*, chapter Finding Signals in DNA, pages 133–152. MIT Press, 2000.
- [114] P. Pevzner and S-H. Sze. Combinatorial approaches to finding subtle signals in DNA sequences. *The Eighth International Conference on Intelligent Systems for Molecular Biology*, pages 269–278, 2000.
- [115] J.C. Polanyi and W.H. Wong. Location of energy barriers. I. effect on the dynamics of reactions $A + BC$. *The Journal of Chemical Physics*, 51(4):1439–1450, 1969.
- [116] W.H. Press, S.A. Teukolsky, W.T. Vetterling, and B.P. Flannery. *Numerical Recipes in C: The Art of Scientific Computing*. Cambridge University Press, 1992.
- [117] A. Price, S. Ramabhadran, and P.A. Pevzner. Finding subtle motifs by branching from sample strings. *BIOINFORMATICS*, 1(1):1–7, 2003.
- [118] W. Quapp, M. Hirsch, O. Imig, and D. Heidrich. Searching for saddle points of potential energy surfaces by following a reduced gradient. *Journal of Computational Chemistry*, 19:1087–1100, 1998.
- [119] B. Raphael, L.T. Liu, and G. Varghese. A uniform projection method for motif discovery in DNA sequences. *IEEE Transactions on Computational biology and Bioinformatics*, 1(2):91–94, 2004.
- [120] C. K. Reddy and H. D. Chiang. A stability boundary based method for finding saddle points on potential energy surfaces. *Journal of Computational Biology*, 13(3):745–766, 2006.
- [121] C. K. Reddy, H. D. Chiang, and B. Rajaratnam. TRUST-TECH based expectation maximization for learning finite mixture models. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, under revision, 2007.
- [122] C. K. Reddy, Y. C. Weng, and H. D. Chiang. Refining motifs by improving information content scores using neighborhood profile search. *BMC Algorithms for Molecular Biology*, 1(23):1–14, 2006.
- [123] R. A. Redner and H. F. Walker. Mixture densities, maximum likelihood and the EM algorithm. *SIAM Review*, 26:195–239, 1984.

- [124] C.R. Reeves. *Modern Heuristic Techniques for Combinatorial Problems*. John Wiley and Sons, New York, NY, 1993.
- [125] S.L. Richter and R.A. DeCarlo. Continuation methods: Theory and applications. *IEEE Transactions on Circuits and Systems*, 30(6):347–352, 1983.
- [126] S. J. Roberts, D. Husmeier, I. Rezek, and W. Penny. Bayesian approaches to gaussian mixture modeling. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 20(11):1133 – 1142, 1998.
- [127] K. Rose. Deterministic annealing for clustering, compression, classification, regression, and related optimization problems. *Proceedings of the IEEE*, 86(11):2210–2239, 1998.
- [128] P. RoyChowdhury, Y. P. Singh, and R. A. Chansarkar. Dynamic tunneling technique for efficient training of multi-layer perceptrons. *IEEE transactions on Neural Networks*, 10(1):48–55, 1999.
- [129] D. E. Rumelhart, G. E. Hinton, and Ronald J. Williams. Learning representations by back-propagating errors. *Nature*, 323(9):533–536, 1986.
- [130] T. Schlick. *Molecular Modeling and Simulation: An Interdisciplinary Guide*, chapter Multivariate Minimization in Computational Chemistry. Springer Verlag, New York, 2002.
- [131] E. Segal, Y. Barash, I. Simon, N. Friedman, and D. Koller. From promoter sequence to expression: a probabilistic framework. In *Proceedings of the sixth annual international conference on Computational biology*, pages 263 – 272, Washington, DC, USA, 2002.
- [132] Y. Shang and B. W. Wah. Global optimization for neural network training. *IEEE Computer*, 29(3):45–54, 1996.
- [133] C.S. Shao, R. Byrd, E. Eskow, and R.B. Schnabel. Global optimization for molecular clusters using a new smoothing approach. *Journal of Global Optimization*, 16(2):167–196, 2000.
- [134] R.H. Shumway and D.S. Stoffer. An approach to time series smoothing and forecasting using the EM algorithm. *Journal of Time Series Analysis*, 3(4):253–264, 1982.
- [135] Z. B. Tang. Adaptive partitioned random search to global optimization. *IEEE Transactions on Automatic Control*, 39(11):2235–2244, 1994.
- [136] S. H. Teng. Coarsening, sampling, and smoothing: Elements of the multilevel method. *Algorithms for Parallel Processing, IMA Volumes in Mathematics and its Applications*, 105, Springer Verlag:247–276, 1999.

- [137] M. Tompa and N. Li et al. Assessing computational tools for the discovery of transcription factor binding sites. *Nature Biotechnology*, 23(1):137 – 144, 2005.
- [138] N. Ueda and R. Nakano. Deterministic annealing EM algorithm. *Neural Networks*, 11(2):271–282, 1998.
- [139] N. Ueda, R. Nakano, Z. Ghahramani, and G.E. Hinton. SMEM algorithm for mixture models. *Neural Computation*, 12(9):2109–2128, 2000.
- [140] J. J. Verbeek, N. Vlassis, and B. Krose. Efficient greedy learning of gaussian mixture models. *Neural Computation*, 15(2):469–485, 2003.
- [141] C. Wang and J. C. Principe. Training neural networks with additive noise in the desired signal. *IEEE Transactions on Neural Networks*, 10(6):1511–1517, 1999.
- [142] E. Xing, W. Wu, M.I. Jordan, and R. Karp. LOGOS: A modular bayesian model for de novo motif detection. *Journal of Bioinformatics and Computational Biology*, 2(1):127–154, 2004.
- [143] L. Xu and M. I. Jordan. On convergence properties of the EM algorithm for gaussian mixtures. *Neural Computation*, 8(1):129–151, 1996.
- [144] B. Zhang, C. Zhang, and X. Yi. Competitive EM algorithm for finite mixture models. *Pattern Recognition*, 37(1):131–144, 2004.