

Available online at www.sciencedirect.com



Procedia Computer Science

Procedia Computer Science 00 (2010) 1-10

Generative Topographic Mapping by Deterministic Annealing

Jong Youl Choi^{a,b}, Judy Qiu^a, Marlon Pierce^a, Geoffrey Fox^{a,b}

^aPervasive Technology Institute, Indiana University, Bloomington IN, USA ^bSchool of Informatics and Computing, Indiana University, Bloomington IN, USA

Abstract

Generative Topographic Mapping (GTM) is an important technique for dimension reduction which has been successfully applied to many fields. However the usual Expectation-Maximization (EM) approach to GTM can easily get stuck in local minima and so we introduce a Deterministic Annealing (DA) approach to GTM which is more robust and less sensitive to initial conditions so we do not need to use many initial values to find good solutions. DA has been very successful in clustering, hidden Markov Models and Multidimensional Scaling but typically uses a fixed cooling schemes to control the temperature of the system. We propose a new cooling scheme which can adaptively adjust the choice of temperature in the middle of process to find better solutions. Our experimental measurements suggest that deterministic annealing improves the quality of GTM solutions.

Keywords: deterministic annealing, optimization, dimension reduction

1. Introduction

Visualization of high-dimensional data in a low-dimension space is the core of exploratory data analysis in which users seek the most meaningful information hidden under the intrinsic complexity of data mostly due to high dimensionality. Among many tools available thanks to the recent advances in machine learning techniques and statistical analysis, Generative Topographic Mapping (GTM) has been extensively studied and successfully used in many application areas: biology [1, 2], medicine [3, 4], feature selection [5], to name a few.

Generative Topographic Mapping (GTM) [6, 7] is an unsupervised learning method developed for modeling the probability density of data and finding a non-linear mapping of highdimensional data onto a low-dimensional space. Similar algorithms known for dimension reduction are Principle Component Analysis (PCA), Multidimensional Scaling (MDS) [8], and Self-Organizing Map (SOM) [9]. In contrast to the Self-Organizing Map (SOM) which does

Email addresses: jychoi@cs.indiana.edu (Jong Youl Choi), xqiu@indiana.edu (Judy Qiu), mpierce@cs.indiana.edu (Marlon Pierce), gcf@cs.indiana.edu (Geoffrey Fox)

not have any density model [7], GTM defines an explicit probability density model based on Gaussian distribution. For this reason, GTM is also known as a principled alternative to SOM.

The problem challenged by the GTM is to find the best set of parameters associated with Gaussian mixtures by using an optimization method, notably the Expectation-Maximization (EM) algorithm [6, 7]. Although the EM algorithm [10] has been widely used in many optimization problems, it has been shown a severe limitation, known as a local optima problem [11], in which the EM method can be easily trapped in local optima failing to return global optimum and so the outputs are very sensitive to initial conditions. To overcome such problem occurred in the original GTM [6, 7], we have applied a robust optimization method known as Deterministic Annealing (DA) [12] to seek global solutions, not local optima. Up to our best knowledge, this is the first work to use DA algorithm to solve the GTM problem in the literature.

The DA [12] has been successfully applied to solve many optimization problems in various machine learning algorithms and applied in many problems, such as clustering [13, 12, 14], visualization [15], protein alignment [16], and so on. The core of the DA algorithm is to seek an global optimal solution in a deterministic way [12], which contrasts to stochastic methods used in the simulated annealing [17], by controlling the level of randomness. This process, adapted from a physical process, is known as annealing in that an optimal solution is gradually revealing as lowering *temperature* which controls randomness. At each level of temperature , the DA algorithm chooses an optimal solution by using the principle of maximum entropy [18, 19, 20], a rational approach to choose the most unbiased and non-committal answers for given conditions.

Regarding the cooling process, which affects the speed of convergency, in general it should be slow enough to get the optimal solutions. Two types of fixed cooling schemes, exponential cooling schemes in which temperature T is reduced exponentially such that $T = \alpha T$ with a coefficient $\alpha < 1$ and linear scheme like $T = T - \delta$ for $\delta > 0$, are commonly used. Choosing the best cooling coefficient α or δ is very problem dependent and no common rule has not yet been reported. Beside of those conventional schemes, in this paper we also present a new cooling scheme, named *adaptive cooling schedule*, which can improve the performance of DA by dynamically computing the next temperatures during the annealing process with no use of fixed and predefined ones.

The main contributions of our paper are as follow:

- i. Developing the DA-GTM algorithm which uses the DA algorithm to solve GTM problem. Our DA-GTM can give more robust answers than the original EM-based GTM, not suffering from the local optima problem (Section 3).
- ii. Deriving closed-form equations to predict phase transitions which is a characteristic behavior of DA [21]. Our phase transition formula can give a guide to set the initial temperature of the DA-GTM algorithm. (Section 4)
- iii. Developing an adaptive cooling schedule scheme, in which the DA-GTM algorithm can automatically adjust the cooling schedule of DA in an on-line manner. With this scheme, users are not required to set parameters related with cooling schedule in DA. (Section 5)

Our experiment results showing the performance of DA-GTM in real-life application will be shown in Section 6 followed by our conclusion (Section 7)

2. GTM Reviews

We start by briefly reviewing the original GTM [6, 7]. The GTM algorithm is to find a nonlinear manifold embedding of K latent variables z_k in low L-dimension space or latent space,



such that $z_k \in \mathbb{R}^L(k = 1, ..., K)$, which can optimally represent the given N data points $x_n \in \mathbb{R}^D(n = 1, ..., N)$ in the higher D-dimension space or data space (usually $L \ll D$) (Figure 1). This is achieved by two steps: First, mapping the latent variables in the latent space to the data space with respect to a non-linear mapping $f : \mathbb{R}^L \mapsto \mathbb{R}^D$. Let us denote the mapped points in the data space as y_k . Secondly, estimating probability density between the mapped points y_k and the data points x_n by using the Gaussian noise model in which the distribution is defined as an isotropic Gaussian centered on y_k having variance β^{-1} (β is known as precision). More specifically, the probability density $p(x_n|y_k)$ is defined by the following Gaussian distribution:

$$\mathcal{N}(\boldsymbol{x}_n|\boldsymbol{y}_k,\boldsymbol{\beta}) = \left(\frac{\beta}{2\pi}\right)^{D/2} \exp\left\{-\frac{\beta}{2}\|\boldsymbol{x}_n-\boldsymbol{y}_k\|^2\right\}.$$
 (1)

The choice of the non-linear mapping $f : \mathbb{R}^L \mapsto \mathbb{R}^D$ can be made from any parametric, nonlinear model. In the original GTM [6, 7], for example, a generalized linear regression model has been used, in which the map is a linear combination of a set of fixed *M* basis functions, such that,

$$\mathbf{y}_k = \boldsymbol{\phi}_k^{\text{Ir}} \mathbf{W},\tag{2}$$

where a column vector $\boldsymbol{\phi}_k = (\phi_{k1}, ..., \phi_{kM})$ is a mapping of z_k by the *M* basis function $\phi_m : \mathbb{R}^L \mapsto \mathbb{R}$ for m = 1, ..., M, such that $\phi_{km} = \phi_m(z_k)$ and **W** is a $M \times D$ matrix containing weight parameters. A^{Tr} represents a transpose of A.

With this model setting, the GTM corresponds to a Gaussian mixture problem to find an optimal set of y_k 's which maximizes the following log-likelihood:

$$\mathcal{L}(\boldsymbol{W},\boldsymbol{\beta}) = \operatorname{argmax}_{\boldsymbol{W}\boldsymbol{\beta}} \sum_{n=1}^{N} \ln \left\{ \frac{1}{K} \sum_{k=1}^{K} p(\boldsymbol{x}_{n} | \boldsymbol{y}_{k}) \right\}$$
(3)

In other words, it is the problem of finding *K* centers for *N* data points, known as *K*-clustering problem. Since *K*-clustering problem is NP-hard [22], the GTM uses the EM method which starts with initial random weight matrix W and iteratively refines it to maximize Eq. (3), which can be easily trapped in local optima. More details can be found in the original GTM paper [6, 7].

3. GTM with Deterministic Annealing (DA-GTM)

Since the use of EM, the original GTM suffers from the local optima problem in which the GTM map can vary depending on the initial parameters. Instead of using the EM, we have applied the DA approach to find a global optimum solution. With the help of DA, we can have more robust GTM maps against the random initial value problem.

In fact, the GTM's objective function (3) to be maximized is the exactly the same problem, so called deterministic annealing clustering, discussed by K. Rose and G. Fox in [12, 23, 21, 24]. The problem is to find optimal clusters for a given distance or distortion measure by using the deterministic annealing approach . In Rose's paper, squared Euclidean distance has been used as distortion measurement. However, in GTM, distances are measured by the Gaussian probability as defined in (1). Thus, by plugging the Gaussian probability into Rose's DA method, we can solve the GTM problem with DA. By using Rose's equations, we can drive the following objective function, known as *free energy*, for the DA-GTM:

$$F(\boldsymbol{W},\boldsymbol{\beta},T) = \operatorname{argmin}_{\boldsymbol{W},\boldsymbol{\beta},T} - T \sum_{n=1}^{N} \ln\left\{ \left(\frac{1}{K}\right)^{\frac{1}{T}} \sum_{k=1}^{K} p(\boldsymbol{x}_{n}|\boldsymbol{y}_{k})^{\frac{1}{T}} \right\}$$
(4)

which we want to minimize as lowering temperature such that $T \rightarrow 1$.

Notice that the GTM's objective function (3) differs only the use of temperature T with our free energy function $F(W, \sigma, T)$. Especially, at T = 1, $\mathcal{L}(W, \beta) = -F(W, \beta, T)$ and so the original GTM's target function can be considered as a special case of the DA-GTM's.

To minimize (4), we need to find parameters which make the following derivatives be zeo.

$$\frac{\partial F}{\partial \mathbf{y}_k} = \beta \sum_{n=1}^N \rho_{kn}(\mathbf{x}_n - \mathbf{y}_k)$$
(5)

$$\frac{\partial F}{\partial \beta} = \sum_{n=1}^{N} \sum_{k=1}^{K} \rho_{kn} \left(\frac{D}{2\beta} - \frac{1}{2} || \mathbf{x}_n - \mathbf{y}_k ||^2 \right)$$
(6)

where ρ_{kn} is a property, known as *responsibility*, such that, $\rho_{kn} = p(\mathbf{x}_n | \mathbf{y}_k)^{\frac{1}{T}} / \sum_{k'=1}^{K} p(\mathbf{x}_n | \mathbf{y}_{k'})^{\frac{1}{T}}$.

By using the same matrix notations used in the original GTM paper [6, 7], the DA-GTM can be written as a process to seek an optimal weights W and precision β at each temperature T.

$$W = (\Phi^{\mathrm{Tr}} G' \Phi)^{-1} \Phi^{\mathrm{Tr}} R' X$$
(7)

$$\beta = \frac{1}{ND} \sum_{n}^{N} \sum_{k}^{K} \rho_{kn} \|\mathbf{x}_{n} - \mathbf{y}_{k}\|^{2}$$
(8)

where *X* is a *N*×*D* data matrix, Φ is a *K*×*M* basis matrix, *G'* is a *K*×*K* diagonal matrix with elements $g'_{kk} = \sum_{n}^{N} (\rho_{kn})^{\frac{1}{T}}$.

4. Phase Transitions of DA-GTM

As a characteristic behavior of the DA algorithm explained by Rose in [12], the DA undergoes phase transitions as lowering the temperatures. At some temperature in the DA, we can not obtain all solutions but, instead, we can only obtain effective number of solutions. All solutions will gradually pop out while the annealing process proceeds as with lowering the temperature.

In the DA-GTM, we can observe the same behavior. As an extreme example, at very high temperature, the DA-GTM gives only one effective latent point in which all y_k 's are converged into the same point \bar{x} which is the center of data points, such that $\bar{x} = \sum_{n=1}^{N} x_n/N$. As lowering the temperature under a certain point, y_k 's settled in the same position start to "explode". We call this temperature as the first critical temperature, denoted by T_c . As we further lowering

temperature, we can observe subsequent phase transitions and so existence of multiple critical temperatures. Computing the first phase transition is an important task since we should start our annealing process with the initial temperature bigger than T_c .

In DA, we can define such phase transitions as a moment of loosing stability of the objective function, the free energy F, and turning to be unstable. Mathematically, that moment corresponds to the point in which the Hessian of the object function looses its positive definiteness.

For our DA-GTM, we can write the following Hessian matrix as a block matrix:

$$H = \begin{bmatrix} H_{11} & \cdots & H_{1K} \\ \vdots & & \vdots \\ H_{K1} & \cdots & H_{KK} \end{bmatrix},$$
(9)

where a sub matrix H_{ij} is a second derivative of the DA-GTM's free energy F Eq. (4). More specifically, H_{ij} can be written as follows:

$$H_{kk} = \frac{\partial^2 F}{\partial \mathbf{y}_k \partial \mathbf{y}_k^{\mathsf{Tr}}} = -\sum_n^N \left\{ \frac{\beta^2}{T} \rho_{kn} (1 - \rho_{kn}) (\mathbf{x}_n - \mathbf{y}_k)^{\mathsf{Tr}} (\mathbf{x}_n - \mathbf{y}_k) - \beta \rho_{kn} \mathbf{I}_D \right\}, \text{ or } (10)$$

$$H_{kk'} = \frac{\partial^2 F}{\partial \mathbf{y}_k \partial \mathbf{y}_{k'}^{\text{Tr}}} = \sum_n^N \left\{ \frac{\beta^2}{T} \rho_{kn} \rho_{k'n} (\mathbf{x}_n - \mathbf{y}_k)^{\text{Tr}} (\mathbf{x}_n - \mathbf{y}_{k'}) \right\} (k \neq k'), \tag{11}$$

where k, k' = 1, ..., K, and I_D is an identity matrix of size D. Note that H_{kk} and $H_{kk'}$ are $D \times D$ matrices and thus, $H \in \mathbb{R}^{KD \times KD}$.

With the Hessian matrix above, we can compute the first phase transition point occurred at T_c . Assuming that the system hasn't undergone the first phase transition and the current temperature is high enough, then we will have all y_k 's settled in the center of the data point, denoted by $y_0 = \bar{x} = \sum_{n=1}^{N} x_n/N$, and equal responsibilities as follows:

$$\mathbf{y}_k = \mathbf{y}_0 \text{ and } \rho_{kn} = 1/K \tag{12}$$

for all k = 1, ..., K. Then, the second derivatives can be rewritten by

$$H_{kk} = -\frac{\beta^2 N}{TK^2} \left\{ (K-1) \mathbf{S}_{\mathbf{x}|\mathbf{y}_0} - \frac{TK}{\beta} \mathbf{I}_D \right\} \text{ and } H_{kk'} = \frac{\beta^2 N}{TK^2} \mathbf{S}_{\mathbf{x}|\mathbf{y}_0}$$
(13)

where $S_{x|y_0}$ represents a covariance matrix of centered data set such that,

$$S_{\mathbf{x}|\mathbf{y}_{0}} = \frac{1}{N} \sum_{n=1}^{N} (\mathbf{x}_{n} - \mathbf{y}_{0})^{\text{Tr}} (\mathbf{x}_{n} - \mathbf{y}_{0})$$
(14)

and the Hessian matrix also can be rewritten by

$$H = \frac{-\beta^2 N}{TK^2} \left\{ \begin{bmatrix} (K-1)\mathbf{S}_{\mathbf{x}|\mathbf{y}_0} & \cdots & -\mathbf{S}_{\mathbf{x}|\mathbf{y}_0} \\ \vdots & \ddots & \vdots \\ -\mathbf{S}_{\mathbf{x}|\mathbf{y}_0} & \cdots & (K-1)\mathbf{S}_{\mathbf{x}|\mathbf{y}_0} \end{bmatrix} - \frac{TK}{\beta} \mathbf{I}_{KD} \right\}$$
(15)

The first phase transition occurs when the system is getting unstable so that the above Hessian matrix is loosing its positive definiteness. I.e., the first phase transition is the moment when the

Hessian matrix becomes singular and so its determinant equals 0(zero), such that $\det(H) = 0$ at $T = T_c$, which holds the following:

$$\operatorname{eig}\left(\begin{bmatrix} (K-1)S_{\boldsymbol{x}|\boldsymbol{y}_0} & \cdots & -S_{\boldsymbol{x}|\boldsymbol{y}_0} \\ \vdots & \ddots & \vdots \\ -S_{\boldsymbol{x}|\boldsymbol{y}_0} & \cdots & (K-1)S_{\boldsymbol{x}|\boldsymbol{y}_0} \end{bmatrix}\right) = \frac{T_c K}{\beta}$$
(16)

where eig(A) is an eigenvalue of A.

We can further simplify the above equation by using the Kronecker product:

$$\operatorname{eig}\left(\left[\begin{array}{ccc} K-1 & \cdots & -1\\ \vdots & \ddots & \vdots\\ -1 & \cdots & K-1 \end{array}\right] \otimes S_{\boldsymbol{x}|\boldsymbol{y}_0}\right) = \operatorname{eig}\left(\left[\begin{array}{ccc} K-1 & \cdots & -1\\ \vdots & \ddots & \vdots\\ -1 & \cdots & K-1 \end{array}\right]\right) \otimes \operatorname{eig}\left(S_{\boldsymbol{x}|\boldsymbol{y}_0}\right) (17)$$

Since the first critical temperature is the most largest one, we can use only the maximum eigenvalue among multiple of them. Thus, the first critical temperature can be obtained by the following equation:

$$T_c = \beta \lambda_{\max} \tag{18}$$

where λ_{\max} is the largest value of $eig(S_{x|y_0})$. Computing the subsequent critical temperature will be discuss in the next.

5. Adaptive cooling schedule

The DA has been applied in many applications and proved its success to find global optimal solutions by avoiding local minima. However, up to our knowledge, no literature has been found to research on the cooling schedule of DA. Commonly used cooling schedule is exponential, such as $T = \alpha T$, or linear, such as $T = T - \delta$ for some fixed coefficient α and δ . Those scheduling schemes are fixed in that cooling temperatures are pre-defined and the coefficient α or δ will not be changed during the process, regardless of the complexity of a given problem.

However, as we discussed previously, the DA algorithm undergoes the phase transitions in which the solution space can change dramatically. One may try to use very small δ near 0 or *alpha* near 1 to make smooth transitions avoiding such drastic changes. However, the procedure can go too long to be used in practice but still there is no guarantee to produce global optimal solutions (You will see an example in our experiment result in the next).

To overcome such problem, we propose an adaptive cooling schedule in which next cooling temperatures are determined dynamically during the annealing process. More specifically, at every iteration of DA algorithm, we predict the next phase transition temperature and move to the point as quickly as possible. Figure 2 shows an example, comparing fixed cooling schedules ((a) and (b)) versus an adaptive cooling schedule.

General approach to find next critical temperatures T_c at any given temperature T will be same with the previous one in that we need to find T_c to make det(H) = 0. However, in contrast to the method used in computing the first critical points, for this problem we can't have any simplified closed-form equations in general. Instead, we can try to divide the problem into pieces to derive closed-form equations and choose the best solution from multiple candidate answers. This method could not provide an exact solution, compared to the method to solve the problem



Figure 2: Various cooling schedule schemes. While exponential (a) and linear (b) is fixed and predefined, our new cooling scheme (c) is adaptive that the next temperature is determined in the on-line manner.

in a whole set, but rather give us an approximation. However, our experiment has showed that such approximation can be a good solution. Also, another advantage we can expect in using our adaptive cooling schedule is that users has no need to set any coefficient for cooling schedule. Cooling process is automatically adjust to the problem.

At a given state in GTM, we have found K soft clusters in which each group is represented by y_k . Instead of finding a global critical temperature from K clusters at once, we can find each critical temperature at each group and choose one among K candidates which is the most biggest yet lower than current temperature so that it can be used as a next temperature.

The sketch of the algorithm is as follows: i) For each k cluster (k = 1, ..., K), find a candidate of next critical temperature $T_{c,k}$ which should satisfy det $(H_{kk}) = 0$ as defined in (10). ii) Choose the most biggest yet lower than the current T among $\{T_{c,k}\}$.

To find $det(H_{kk}) = 0$, we need to rewrite Eq. (10) as follows:

$$H_{kk} = \frac{-\beta^2}{T} \left(U_{\mathbf{x}|\mathbf{y}_k} - V_{\mathbf{x}|\mathbf{y}_k} - \frac{Tg'_{kk}}{\beta} I_D \right)$$
(19)

with the following definition:

$$\boldsymbol{U}_{\boldsymbol{x}|\boldsymbol{y}_{k}} = \sum_{n=1}^{N} \rho_{kn} (\boldsymbol{x}_{n} - \boldsymbol{y}_{k})^{\mathrm{Tr}} (\boldsymbol{x}_{n} - \boldsymbol{y}_{k})$$
(20)

$$\boldsymbol{V}_{\boldsymbol{x}|\boldsymbol{y}_{k}} = \sum_{n=1}^{N} (\rho_{kn})^{2} (\boldsymbol{x}_{n} - \boldsymbol{y}_{k})^{\mathrm{Tr}} (\boldsymbol{x}_{n} - \boldsymbol{y}_{k})$$
(21)

and $g'_{kk} = \sum_{n=1}^{N} \rho_{kn}$. Then, the condition det $(H_{kk}) = 0$ at $T = T_{c,k}$ will hold the following:

$$\operatorname{eig}\left(\boldsymbol{U}_{\boldsymbol{x}|\boldsymbol{y}_{k}}-\boldsymbol{V}_{\boldsymbol{x}|\boldsymbol{y}_{k}}\right) = \frac{g_{kk}'}{\beta}T_{c,k}$$
(22)

Then, the next critical temperature $T_{c,k}$ for the cluster k can computed by

$$T_{c,k} = \frac{\beta}{g'_{kk}} \lambda_{\max,k}$$
(23)

where $\lambda_{\max,k}$ is the largest but less than Tg'_{kk}/β eigenvalue of the matrix $(U_{x|y_k} - V_{x|y_k})$. The overall pseudo code is shown in Algorithm 1 and Algorithm 2.

6. Experiment Results

To compare the performances of our DA-GTM with the original EM-based GTM (EM-GTM hereafter for short), we have performed a set of experiments by using two dataset: i) the oil

Algorithm 1 GTM with Deterministic Annealing	Algorithm 2 Find the next critical temperature
DA-GTM	NEXTCRITICALTEMP
1: Set $T > T_c$ by using Eq. (18)	1: for $k = 1$ to K do
2: Choose randomly <i>M</i> basis function $\phi_m(m = 1,, M)$	2: $\Lambda_k \leftarrow \{\emptyset\}$
3: Compute Φ whose element $\phi_{km} = \phi_m(z_k)$	3: for each $\lambda \in \operatorname{eig} \left(U_{\mathbf{x} \mathbf{y}_k} - V_{\mathbf{x} \mathbf{y}_k} \right)$ do
4: Initialize randomly W	4: if $\lambda < Tg'_{\mu\nu}/\beta$ then
5: Compute β by Eq. (8)	5: $\Lambda_k \leftarrow \stackrel{\sim}{\Lambda}_k \cup \lambda$
6: while $T \ge 1$ do	6: end if
7: Update W by Eq. (7)	7: end for
8: Update β by Eq. (8)	8: $\lambda_{\max,k} \leftarrow \max(\Lambda_k)$
9: $T \leftarrow \text{NextCriticalTemp}$	9: $T_{c,k} \leftarrow \beta \lambda_{\max,k} / g'_{kk}$
10: end while	10: end for
11: return Φ, W, β	11: return $T_c \leftarrow \max(\{T_{c,k}\})$

flow data used in the original GTM paper [6, 7] obtained from the GTM website¹, which has 1,000 points having 12 dimensions for 3-phase clusters and ii) chemical compound dataset obtained from PubChem database², which is a NIH-funded repository for over 26 million chemical molecules and provides their chemical structure fingerprints and biological activities, for the purpose of chemical information mining and exploration. Among 26 Million PubChem dataset, in this paper we have randomly selected 1,000 subset having 166 dimensions.



Figure 3: Comparison of EM-GTM (a) and DA-GTM with different cooling scheme: exponential (b) and adaptive (c) for oil-flow data with 3-phases (A=Homogeneous, B=Annular, and C=Stratified configuration). Average of 10 randominitialized execution results shows DA-GTM with adaptive (c) is the best (the biggest log-likelihood) followed by DA-GTM with conventional exponential-scheme (b). EM-GTM (a) shows worst performance.

First we have compared various GTM maps produced by EM-GTM, DA-GTM with different cooling schemes: exponential and adaptive (Figure 3). We have executed each algorithm 10 times with random setup and selected median results among them. As a result shown in Figure 3, DA-GTM with adaptive cooling schedule has shown the best performance.

In the next, we have compared the performance of EM-GTM and DA-GTM with 3 cooling schedule schemes: i) Adaptive, which we have prosed in this paper, ii) Exponential with cooling coefficients $\alpha = 0.95$ (denoted Exp-A hereafter), and iii) Exponential with cooling coefficients $\alpha = 0.99$ (denoted Exp-B hereafter). For each DA-GTM setting, we have also applied 3 different starting temperature 5, 6, and 7, which are all bigger than the 1st critical temperature which is

¹GTM homepage, http://www.ncrg.aston.ac.uk/GTM/

²PubChem project, http://pubchem.ncbi.nlm.nih.gov/





Figure 4: Comparison of EM-GTM with DA-GTM in various settings. Average of 50 random initialized runs are measured for EM-GTM, DA-GTM with 3 cooling schemes (adaptive, exponential with α = 0.95 (Exp-A) and α = 0.99 (Exp-B).

Figure 5: A progress of DA-GTM with adaptive cooling schedule. This example show how DA-GTM with adaptive cooling schedule progresses through iterations

about 4.64 computed by Eq. (18). Figure 4 shows the summary of our experiment results in which numbers are estimated by the average of 50 executions with random initialization.

As a result shown in Figure 4(a), DA-GTM shows strong robustness against local minima since the mean of log-likelihood from DA-GTM's outperforms EM-GTM's by about 11.15 % even with smaller deviations. Interestingly, at the starting temperature 7 and 9, the performance of Exp-B, which used more finer-grain cooling schedule ($\alpha = 0.99$) than Exp-A ($\alpha = 0.95$), is lower than Exp-A and even under EM-GTM. This shows that fine-grained cooling schedule can not guarantee the best solutions than corse-grained ones in using DA. However, our adaptive cooling scheme mostly outperforms among other cooling schemes. Figure 5 shows an example of execution of DA-GTM with adaptive cooling schedule.



Figure 6: GTM map of 1,000 PubChem dataset by using EM-GTM (a) and DA-DTM (b). Chosen as a median map of 10 random initialized execution of EM-GTM and DA-GTM execution, DA-GTM results show better than EM-GTM with bigger log-likelihood.

We have also compared EM-GTM and DA-GTM for the 1,000 PubChem dataset. As shown in Figure 6, DA-GTM outperforms EM-GTM since DA-GTM's average log-likelihood (-36,454) is bigger than EM-GTM's (-36,584).

7. Conclusion

We have solved the GTM problem, originally using EM, by using the deterministic annealing (DA) algorithm which is more resilient against the local optima problem and less sensitive to initial conditions, from which the original EM method was suffered. We have also developed a new cooling scheme, called *adaptive cooling schedule*. In contrast to the conventional cooling schemes such as linear and exponential ones, all of them are pre-defined and fixed, our adaptive cooling scheme can adjust granularity of cooling speed in an on-line matter. In our experiment, our cooling scheme can outperform other conventional methods.

References

- A. Staiano, L. De Vinco, A. Ciaramella, G. Raiconi, R. Tagliaferri, G. Longo, G. Miele, R. Amato, C. Del Mondo, C. Donalek, et al., Probabilistic principal surfaces for yeast gene microarray data-mining, in: Data Mining (ICDM 2004). Proceedings of Fourth IEEE International Conference, 2004, pp. 202–208.
- [2] D. D'Alimonte, D. Lowe, I. Nabney, V. Mersinias, C. Smith, MILVA: An interactive tool for the exploration of multidimensional microarray data (2005).
- [3] A. Vellido, P. Lisboa, Handling outliers in brain tumour MRS data analysis through robust topographic mapping, Computers in Biology and Medicine 36 (10) (2006) 1049–1063.
- [4] D. Maniyar, I. Nabney, Visual data mining using principled projection algorithms and information visualization techniques, in: Proceedings of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining, ACM, 2006, p. 648.
- [5] A. Vellido, Assessment of an Unsupervised Feature Selection Method for Generative Topographic Mapping, Lecture Notes in Computer Science 4132 (2006) 361.
- [6] C. Bishop, M. Svensén, C. Williams, GTM: The generative topographic mapping, Neural computation 10 (1) (1998) 215–234.
- [7] C. Bishop, M. Svensén, C. Williams, GTM: A principled alternative to the self-organizing map, Advances in neural information processing systems (1997) 354–360.
- [8] J. Kruskal, Multidimensional scaling by optimizing goodness of fit to a nonmetric hypothesis, Psychometrika 29 (1) (1964) 1–27.
- [9] T. Kohonen, The self-organizing map, Neurocomputing 21 (1-3) (1998) 1-6.
- [10] A. Dempster, N. Laird, D. Rubin, Maximum likelihood from incomplete data via the EM algorithm, Journal of the Royal Statistical Society. Series B (Methodological) 39 (1) (1977) 1–38.
- [11] N. Ueda, R. Nakano, Deterministic annealing EM algorithm, Neural Networks 11 (2) (1998) 271–282.
- [12] K. Rose, Deterministic annealing for clustering, compression, classification, regression, an; related optimization problems, Proceedings of the IEEE 86 (11) (1998) 2210–2239.
- [13] T. Hofmann, J. Buhmann, Pairwise data clustering by deterministic annealing, IEEE Transactions on Pattern Analysis and Machine Intelligence 19 (1) (1997) 1–14.
- [14] X. Yang, Q. Song, Y. Wu, A robust deterministic annealing algorithm for data clustering, Data & Knowledge Engineering 62 (1) (2007) 84–100.
- [15] H. Klock, J. Buhmann, Multidimensional scaling by deterministic annealing, Lecture Notes in Computer Science 1223 (1997) 245–260.
- [16] L. Chen, T. Zhou, Y. Tang, Protein structure alignment by deterministic annealing, Bioinformatics 21 (1) (2005) 51–62.
- [17] S. Kirkpatric, C. Gelatt, M. Vecchi, Optimization by simulated annealing, Science 220 (4598) (1983) 671-680.
- [18] E. Jaynes, Information theory and statistical methods I, Physics Review 106 (1957) (1957) 620-630.
- [19] E. Jaynes, Information theory and statistical mechanics. II, Physical review 108 (2) (1957) 171–190.
- [20] E. Jaynes, On the rationale of maximum-entropy methods, Proceedings of the IEEE 70 (9) (1982) 939–952.
- [21] K. Rose, E. Gurewitz, G. Fox, Statistical mechanics and phase transitions in clustering, Physical Review Letters 65 (8) (1990) 945–948.
- [22] D. Aloise, A. Deshpande, P. Hansen, P. Popat, NP-hardness of Euclidean sum-of-squares clustering, Machine Learning 75 (2) (2009) 245–248.
- [23] K. Rose, E. Gurewitz, G. Fox, A deterministic annealing approach to clustering., Pattern Recognition Letters 11 (9) (1990) 589–594.
- [24] K. Rose, E. Gurewitz, G. Fox, Vector quantization by deterministic annealing, IEEE Transactions on Information Theory 38 (4) (1992) 1249–1257.