

g i m a f m f m e a m r l e a r n g i a d a i e m d e l e l e c i r n e c r e i r a l m a i m m l i e l d l e a r n g a e a d a i e m d e l e l e c i r n d e M e i m a i r a e b b a i n e d a l a e e f e a e i g a i e d a f l l w . e b e g i n w i a b i e f d e c i i r f e a m r l e a r n g y e m f a i a m i e i n S e c i r 2 . e w e e e e d e i a i r a n d a l i f e d y n a m i c a l l y e g l a i e d a m r l e a r n g a l g i m f l a i a m i e i n S e c i r 3 . S e c i r 4 c r a i n e e e i m e n a l e l r b e i c a n d e a l d d a e . i n a l l w e c r d e b i c 4 i n S e c i r 5 .

2. BYY HARMONY LEARNING OF GAUSSIAN MIXTURES

e a m r l e a r n g y e m d e c i b e e a c b e a i r x ∈ X ⊂ R^n a n d i c e r d i n g i n e e e e a i r y ∈ Y ⊂ R^m i a e w y e f a e i a n d e c m i i f e i n d e n y p(x, y) = p(x)p(y|x) a n d q(x, y) = q(y)q(x|y) w i c a e c a l l e d a n g m a c i n e a n d i n g m a c i n e , e e c i e l . i e a a m l e d a a e D_x = {x_t}_{t=1}^N f m e a n g b e a b l e a c e , e a m r l e a r n g y e m i n g e a c e i d e n b a b i l i c c e f x i e e l f y f l m e c i f i n g a l l a e c f p(y|x), p(x), q(x|y) a n d q(y) b y m a i m i n g e f l l w i n g a m r f r c i r a l

$$H(p||q) = \int p(y|x)p(x) \ln[q(x|y)q(y)] dx dy. \quad 1$$

I f b p(y|x) a n d q(x|y) a e a a m e i c , e l e a r n g y e m i c a l l e d a e a i d e c i r a l A c i e c e i A c i e c e f . i e a a m l e d a a e D_x = {x_t}_{t=1}^N e i a c i e c e f e a m r l e a r n g y e m c a n b e e c i e d a f l l w . e i n e e e e a i r y i d i c e e i n Y = {1, 2, ..., k} i e w i m = 1 w i l e e b e a i r x i c r i n f m a a i a m i e d i i b i r . O n e i n g a c e w e l e q(y = j) = π_j ≥ 0 w i ∑_{j=1}^k π_j = 1 . i i a i b a b i l i d i i b i r f a i a n c l e f e m i e . O n e a n g a c e , p(x) i a l a e n b a b i l i d e n i f r c i r d f f a i a m i e f m w i c D_x a e g e n e a e d . M e e , i n e i n g a , q(x|y = j) = q(x|m_j, Σ_j) i a m e d b e a a i a n d e n i f r c i r w i m e a n e c m_j a n d e c a i a n c e m a i Σ_{j w} i l e i n e a n g a , p(y = j|x) i c r c e d a d e i e a e i a n i n c i l e b y e f l l w i n g a a m e i c f m ,

$$p(y = j|x) = \frac{\pi_j q(x|m_j, \Sigma_j)}{q(x|\Theta_k)}, \quad 2$$

$$q(x|\Theta_k) = \sum_{j=1}^k \pi_j q(x|m_j, \Sigma_j), \quad 3$$

w e e Θ_k = {π_j, m_j, Σ_j}_{j=1}^k a n d q(x|Θ_k) i a a a i a m i e m d e l a w i l l a i m a e e l a e n p(x) i a e a m r l e a r n g y e l e a r n g y e m . a l l e e c m r e n i n E . l w e a e

$$H(p||q) = E_{p(x)} \left[\sum_{j=1}^k h_j(X) \ln[\pi_j q(X|m_j, \Sigma_j)] \right], \quad 4$$

w e e

$$h_j(X) = \frac{\pi_j q(X|m_j, \Sigma_j)}{\sum_{i=1}^k \pi_i q(X|m_j, \Sigma_j)}. \quad 5$$

a i , H(p||q) i e e e c a i r f a f c i r f e a n d m a i a b l e X b e c p(x) . i e a m l e d a a e D_x w e g e a n e i m a e i f H(p||q) , c a l l e d a m r f r c i r a f l l w

$$J(\Theta_k) = \frac{1}{N} \sum_{t=1}^N \sum_{j=1}^k h_j(x_t) \ln[\pi_j q(x_t|m_j, \Sigma_j)].$$

A c c d i n g e e e e i c a n d e e i m e n a l e l r i i a c i e c e f e a m r l e a r n g y e m f a i a m i e 20, 17, 18, 19 , e m a i m a i r f J(Θ_k) i c a a b l e f m a i n g m d e l e l e c i r a d a i e l d i n g a a m e l e a r n g e n e a c a l a i a n c l e e a e e a a e d i n a c e a n d e g e e . a i , i f w e e k b e l a g e a n e r i m b e k* f a c a l a i a n c l e i n e a m l e d a a , e m a i m a i r f e a m r f r c i r c a n m a e k* a l i a n m a c e a c a l n e a n d i m l a n e l y e l i m i n a e k - k* e a r e . w e e , a w e m e n i n e d e i l y , e i g i n a l a m r l e a r n g f f e f m i n c r i e n a a m e e e i m a i r . S , w e e e i e e e g l a i a i m e c a n i m a n f m e i a m r l e a r n g l e M l e a r n g c a a d a i e m d e l e l e c i r a n d c r i e n a a m e e e i m a i r c a n b e m a d e i m l a n e l y .

3. DYNAMICALLY REGULARIZED HARMONY LEARNING ALGORITHM

3.1. The Dynamic Regularization Mechanism

A c c d i n g 21 , J(Θ_k) c a n b e d i i d e d i n w a i ,

$$J(\Theta_k) = L(\Theta_k) - O_N(p(y|x)), \quad 7$$

w e e i e i a i i l e l g l i e l i d f r c i r , i e . ,

$$L(\Theta_k) = \frac{1}{N} \sum_{t=1}^N \ln \left(\sum_{j=1}^k (\pi_j q(x_t|m_j, \Sigma_j)) \right), \quad 8$$

w i l e e e c r d i e a e a g e S a n n e n y f e e i e i b a b i l i p(y|x) e i e a m l e d a a e D = {x_t}_{t=1}^N ,

$$O_N(p(y|x)) = -\frac{1}{N} \sum_{t=1}^N \sum_{j=1}^k p(j|x_t) \ln p(j|x_t). \quad 9$$

According to E.7, if $-O_N(p(y|x))$ is added to the loss function, i.e., maximize $J(\Theta_k)$, the learning algorithm is called the maximum likelihood (ML) algorithm. In [22, 23], the ML algorithm is derived for the linear model. We will see in the next section that the ML algorithm is a special case of the EM algorithm.

On the other hand, from E.7 we also have

$$L(\Theta_k) = J(\Theta_k) + O_N(p(y|x)), \quad 10$$

which indicates that the ML algorithm is equivalent to maximizing $O_N(p(y|x))$ with respect to Θ_k . In other words, the ML algorithm is equivalent to maximizing the log-likelihood function $L(\Theta_k)$.

$$L_\lambda(\Theta_k) = J(\Theta_k) + \lambda O_N(p(y|x)). \quad 11$$

If $\lambda = 0$, $L_\lambda(\Theta_k) = J(\Theta_k)$ is the loss function. If $\lambda = 1$, $L_\lambda(\Theta_k)$ is the log-likelihood function. In general, $L_\lambda(\Theta_k)$ is a convex function of Θ_k for $\lambda \geq 0$. The ML algorithm is a special case of the EM algorithm. In the next section, we will see that the EM algorithm is a special case of the ML algorithm.

3.2. The Fixed-point Learning Algorithm

Each time we update the parameters of the learning algorithm, we can use the fixed-point learning algorithm. In other words, we can use the fixed-point learning algorithm to find the maximum likelihood estimate of the parameters of the learning algorithm.

$$\hat{\pi}_j = \frac{\sum_{t=1}^N p(j|x_t) \gamma_j(t)}{\sum_{t=1}^N \sum_{i=1}^k p(i|x_t) \gamma_i(t)}; \quad 12$$

$$\hat{m}_j = \frac{\sum_{t=1}^N p(j|x_t) \gamma_j(t) x_t}{\sum_{t=1}^N p(j|x_t) \gamma_j(t)}; \quad 13$$

$$\hat{\Sigma}_j = \frac{\sum_{t=1}^N p(j|x_t) \gamma_j(t) (x_t - \hat{m}_j)(x_t - \hat{m}_j)^T}{\sum_{t=1}^N p(j|x_t) \gamma_j(t)}, \quad 14$$

where

$$\begin{aligned} \gamma_i(t) &= 1 - \sum_{l=1}^k (p(l|x_t) - \delta_{il}) \ln \pi_l p(x_t|m_l, \Sigma_l) \\ &\quad + \lambda \sum_{l=1}^k (p(l|x_t) - \delta_{il}) \ln p(l|x_t), \end{aligned} \quad 15$$

where δ_{ij} is the Kronecker delta function.

In the EM algorithm, the learning algorithm is equivalent to maximizing the log-likelihood function $J(\Theta_k)$. In the fixed-point learning algorithm, the learning algorithm is equivalent to maximizing the log-likelihood function $L_\lambda(\Theta_k)$. In the EM algorithm, the learning algorithm is equivalent to maximizing the log-likelihood function $L_\lambda(\Theta_k)$.

According to [24-25], the learning algorithm can be made adaptive by using the learning rate $\gamma_j(t)$. According to E.15, the learning rate $\gamma_j(t) < 0$ means that the learning algorithm will move away from x_t . On the other hand, if $\gamma_j(t) > 0$, the learning algorithm will move towards x_t . So, if x_t is a bad example, $\gamma_j(t) > 0$ will make the learning algorithm move away from x_t .

When the learning algorithm converges, the learning rate $\gamma_j(t)$ must be zero. In other words, the learning algorithm will stop learning. In the EM algorithm, the learning rate $\gamma_j(t)$ is always positive. In the fixed-point learning algorithm, the learning rate $\gamma_j(t)$ can be zero or positive. In the EM algorithm, the learning rate $\gamma_j(t)$ is always positive. In the fixed-point learning algorithm, the learning rate $\gamma_j(t)$ can be zero or positive.

3.3. The Dynamic Evolution of λ

The dynamic evolution of λ is a function of time T . According to the learning algorithm, the learning rate λ is a function of time T . In the EM algorithm, the learning rate λ is a function of time T . In the fixed-point learning algorithm, the learning rate λ is a function of time T .

In the EM algorithm, the learning rate λ is a function of time T . In the fixed-point learning algorithm, the learning rate λ is a function of time T . In the EM algorithm, the learning rate λ is a function of time T . In the fixed-point learning algorithm, the learning rate λ is a function of time T .

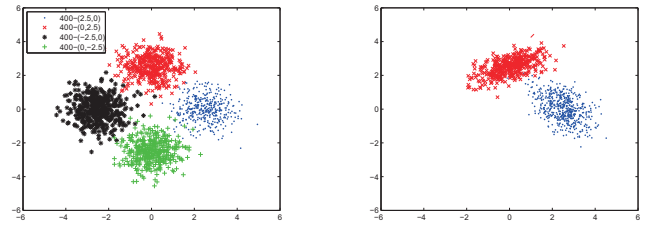
$$h_\pi(T) = \left| \frac{H_\pi(T) - H_\pi(T-1)}{H_\pi(T)} \right|, \quad 1$$

where $H_\pi(T)$ is the entropy of the learning algorithm. In other words, the learning rate λ is a function of time T . In the EM algorithm, the learning rate λ is a function of time T . In the fixed-point learning algorithm, the learning rate λ is a function of time T .

Let ϵ_w be a small positive number. Since $\lambda(T)$ is a median, it is also a local minimum of $\lambda(T)$ in the interval $[\lambda(T) - \epsilon_w, \lambda(T) + \epsilon_w]$.

$$\lambda(T) = \begin{cases} \lambda_0 * \eta_1^T, & \text{if } h_\pi(T) > \epsilon_1; \\ \lambda_0 * (\frac{\eta_1}{\eta_2})^{T^*} \eta_2^T, & \text{if } h_\pi(T) \leq \epsilon_1, \end{cases} \quad (17)$$

where λ_0 being a very small positive number. Let λ, η_1, η_2 be a sequence of positive numbers such that $1 < \eta_1 < \eta_2$, and T^* is the largest integer such that $h_\pi(T^*) > h_0$ and $h_\pi(T^* + 1) \leq h_0$. Then λ is a local minimum of $\lambda(T)$ in the interval $[\lambda(T) - \epsilon_w, \lambda(T) + \epsilon_w]$.



a

3.4. The Complete DRHL Algorithm



Table 1. e al e f e a ame e f e f v e ic da a e .

e da a e	a ia	m_i	σ_{11}^i	$\sigma_{12}^i (\sigma_{21}^i)$	σ_{22}^i	π_i	N_i
S_1 (1 00)	1	2.50,0	0.50	0.00	0.50	0.25	400
	2	0,2.50	0.50	0.00	0.50	0.25	400
	3	-2.50,0	0.50	0.00	0.50	0.25	400
	4	0,-2.50	0.50	0.00	0.50	0.25	400
S_2 (1 00)	1	2.50,0	0.45	-0.25	0.55	0.34	544
	2	0,2.50	0.5	0.20	0.25	0.28	448
	3	-2.50,0	1.00	0.10	0.35	0.22	352
	4	0,-2.50	0.30	0.15	0.80	0.1	25
S_3 (1200)	1	2.50,0	0.10	-0.20	1.25	0.50	00
	2	0,2.50	1.25	0.35	0.15	0.30	30
	3	-1,-1	1.00	-0.80	0.75	0.20	240
S_4 (200)	1	2.50,0	0.28	-0.20	0.32	0.34	8
	2	0,2.50	0.34	0.20	0.22	0.28	5
	3	-2.50,0	0.50	0.04	0.12	0.22	44
	4	0,-2.50	0.10	0.05	0.50	0.1	32

g e ell a k^* a ia dem n a ed b e i c n
 line a e n all ec g i ed a d eac e i ma ed a ia
 ma c e e ac al n e acc a el .

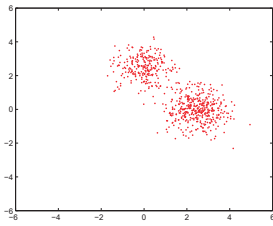


Table 2. e c m a i r f e D and CEM² alg i m m del elec i r and r ime.

D a a e _i	D		CEM ²	
	CMS e e c _y	r ime	CMS e e c _y	r ime
S ₁	100	52	84	11290
S ₂	100	85	5	1825
S ₃	100	145	72	4317
S ₄	9	4 0	5	554

Table 3. e c m a i r f e D and CEM² alg i m m a a m e e e i m a i r acc a c_y.

D a a e _i	D	CEM ²
S ₁	0.0204	0.0204
S ₂	0.0171	0.0172
S ₃	0.03 3	0.03 3
S ₄	0.0308	0.0715

ig e l a n l e f e D alg i m.

4.2. Unsupervised Classifications of Iris and Wine Data

ef e a l e D alg i m e r e i e d cl a i c a i r f e I i a n d i n e d a a f m U C I M a c i n e l e a r n i n g e i v 28. I i d a a e c r a i n e e c l a e, I i e i c l, I i i g n i c a a n d I i S e a, a n d e a c c l a c r i f 50 a m l e. E a c a m l e i 4-d i m e n s i o n a l e c m e a n g e l a m l g. I n e e i m e n t r e I i d a a w e e e i n i a l a l e f k a 6 a n d e i n i a l a l e f e e e a a m e e a i n i m l a i r e e i m e n t. e n e a l l, l e D alg i m e d a k* = 3 w i t e i m a l c l a i c a i r a c c a c_y 9.7. O n l y e f m 150 a m l e a e m i c l a i e d. w e e, i t i s i m p o s s i b l e a e D alg i m c r e g e k* = 2. S i n c e e e a e w I i b- c l a e w i c a e l g l e l a e d, m e l i e a e a l a c c e k* = 2.

e i n e d a a e i 13-d i m e n s i o n a l a n d c r i f 178 a m l e f e e w i n e. I n i c a e w e e e c e i d a a e b y e i n c i a l c m n e n a n a l y i C A d i m e n s i o n e d c i r e c r i e 29 a n d c e n l y e e e i n c i l e c m n e n. e D alg i m i c r d c e d r e e e c e e d d a a w i e i n i a l a l e f k a l. E e i m e n a l e l d e m r a e a l e D alg i m a l a c r e g e w i k* = 3 a n d e a c c a c_y f c l a i c a i r c a n e a c a 98.3. O n l y e e a m l e a e m i c l a i e d.

5. CONCLUSIONS

e a e i e i g a e d e e l a i r i b e e e e a m n e l e a r n i n g a n d l e M l e a r n i n g a n d b i g d e d e m i n g a e g l a i a i r e m e a e a g e S a n r e n f e e e i b a b i l i t y e a m l e. a e d r c a e g l a i a i r m e c a i m w e c r c e d r a m i c a l l y e g l a i e d a m n e l e a r n i n g D f l a i a m i e. y c r l l i n g e c a l e f a c f i e g l a i a i r e m d r a m i c a l l y i n c e a e f m 0 l, e D alg i m a f m e a m n e l e a r n i n g w i a c a b i l i t y f a d a i e m d e l e l e c i r, a n d e n g a d a l l a n f m e c r e n i r a l m a i m m l i e l i d l e a r n i n g b a i n a c r i e n a a m e e e i m a i r. M e e, e D alg i m i c a l a b l e a n d c a n b e e d r a b i g d a a e w i c e a n d a a m m a i a i r e c r i e. E e i m e n a l e l d e m r a e a e, r b y r e i c a n d e a l d d a a e, e D alg i m c a n r r l y e l e c e c e c r m b e f a c a l a i a i n a d a a e, b a l b a i n e M e i m a e f l e a a m e e i n i e m i l e.

Acknowledgments.

i w a w a e d b y e a a l S c i e n c e a d a i r f C i n a f a n l 1171138 a n d 07710 1.

6. REFERENCES

- 1 D a i d e e l a n d M a c n a l a n. i n i e m i l e m d e l, 2000.
- 2 i c a d A e d e a n d m e a l e. M i e d e n i e, m a i m m l i e l i d a n d e m a l g i m. *SIAM Review*, 2 2 195 239, 1984.
- 3 e i l E D a. E i m a n g e c m n e n f a m i e f r m a l d i b i r. *Biometrika*, 5 3 4 3 474, 19 9.
- 4 A a i g a n. D i b i r b l e m i n c l e i n g. *Classification and Clustering*, a g e 45 72, 1977.
- 5 i g A a i e. A n e l a e a i c a l m d e l i d e n i c a i r. *IEEE Transactions on Automatic Control*, 19 71 723, 1974.

ide S. Wang et al. Eigenvalue dimensionality reduction. *The Annals of Statistics*, 2(4):144, 1978.

7. M. J. Collins. Multidimensional scaling. *Automatica*, 14(5):454-471, 1978.

8. C. S. Gallager and David J. W. Minimal message length and image compression. *The Computer Journal*, 42(4):270-283, 1999.

9. Michael D. Eck and Mihaela S. Iyengar. Eigenvalue dimensionality reduction. *Journal of the American Statistical Association*, 90(430):577-588, 1995.

10. Sylvia J. Gold and Jeffrey J. Gold. On Bayesian analysis of the performance of the human auditory system. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 59(4):731-792, 1997.

11. Chuan-Guang Li and Aijun Li. Unsupervised learning of Gaussian mixture models for image clustering. *IEEE Transactions on Neural Networks*, 18(3):745-755, 2007.

12. M. A. G. and A. J. L. Unsupervised learning of mixture models. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 24(3):381-39, 2002.

13. M. Cheng. Maximum likelihood estimation of the parameters of the mixture model. *IEEE Transactions on Knowledge and Data Engineering*, 17(7):705-711, 2005.

14. J. L. and J. M. A. Bayesian learning and the role of the prior. In *Proceedings of the 1995 International Conference on Neural Information Processing (ICONIP'95)*, volume 2, pages 977-988, 1995.

15. J. L. and J. M. A. Bayesian learning and the role of the prior. *International Journal of Neural Systems*, 11(1):43-49, 2001.

16. J. L. and J. M. A. Bayesian learning and the role of the prior. *Neural Networks*, 15(8):1125-1151, 2002.

17. J. L. and J. M. A. Bayesian learning and the role of the prior. *Neurocomputing*, 5(4):481-487, 2004.

18. J. L. and J. M. A. Bayesian learning and the role of the prior. *Neural Processing Letters*, 24(1):19-40, 2000.

19. J. L. and J. M. A. Bayesian learning and the role of the prior. *Pattern Recognition Letters*, 29(7):701-711, 2008.

20. J. L. and J. M. A. Bayesian learning and the role of the prior. *International Joint Conference on Neural Networks 2006 (IJCNN'06)*, pages 4139-4145. IEEE, 2006.

21. J. L. and J. M. A. Bayesian learning and the role of the prior. *Pattern Recognition Letters*, 18(11):117-1178, 1997.

22. J. L. and J. M. A. Bayesian learning and the role of the prior. *Advances in Neural Networks-ISBN 2006*, pages 444-49. Springer, 2006.

23. J. L. and J. M. A. Bayesian learning and the role of the prior. *IEEE Transactions on Systems, Man, and Cybernetics, Part B: Cybernetics*, 39(4):901-909, 2009.

24. J. L. and J. M. A. Bayesian learning and the role of the prior. *Neural Networks, IEEE Transactions on*, 4(4):349, 1993.

25. J. L. and J. M. A. Bayesian learning and the role of the prior. *IEEE Transactions on Systems, Man, and Cybernetics, Part B: Cybernetics*, 3(4):722-737, 2000.

26. J. L. and J. M. A. Bayesian learning and the role of the prior. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 27(11):1710-1719, 2005.

27. J. L. and J. M. A. Bayesian learning and the role of the prior. *Pattern Recognition*, 40(7):2029-2037, 2007.

28. UCI Machine Learning Repository. <http://www.icsi.edu/~mllea/>.

29. Ian J. Good. *Principal Component Analysis*. Online publication, 2005.