

Introduction to Machine Learning

Ranga Raju Vatsavai, Ph.D.

Chancellors Faculty Excellence Associate Professor in Geospatial Analytics
Department of Computer Science, North Carolina State University (NCSU)

Feb. 25-27, 2019

Semi-supervised Learning

- Ground-truth is costly and time consuming
 - Empirical studies show that we need at least (10-30) x dimensions number of samples per class
- Unlabeled samples are plenty and cheap
- Can we combine small number of labeled samples and large number of unlabeled sample to improve learning model?
 - Yes: semi-supervised learning

2/26/19

© Raju Vatsavai

Self-training

- Key assumption
 - One's own high confidence predictions are correct
- Algorithm
 1. Train f from $D_l = (\mathbf{X}, Y)$
 2. Use f to predict on $D_{ul} = (\mathbf{x}, ?)$
 3. Add $(\mathbf{x}, f(\mathbf{x}))$ to labeled data D_l
 - Variations: All, most confident, weighted?
 4. Repeat

2/27/19

© Raju Vatsavai

Self-training

- Advantages
 - Simplest semi-supervised learning
 - Easy to implement (wrapper around existing algorithms)
- Disadvantages
 - Early mistakes may lead to degrade in performance

2/27/19

© Raju Vatsavai

Probabilistic Approach

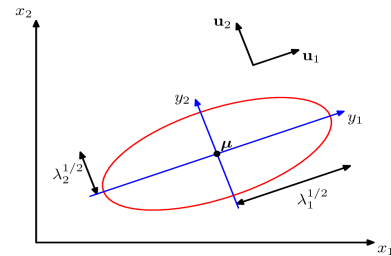
- Assume data is generated by a mixture model
 - E.g., Gaussian Mixture Model
- $D = D_I + D_{II}$
- How do you estimate parameters? MLE don't work with missing labels
- Solution
 - Expectation Maximization (EM) algorithm

2/27/19

© Raju Vatsavai

Multivariate Gaussian

- Bivariate covariance matrix representation



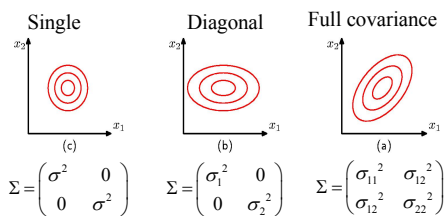
2/26/19

© Raju Vatsavai

6

Multivariate Gaussian

- Σ plays key role, but accurate estimation requires large number samples (10-30 x dim; per class)
- However, Σ can be simplified

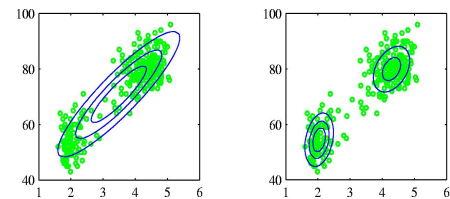


2/26/19

© Raju Vatsavai

7

What If The Data Is Multimodal

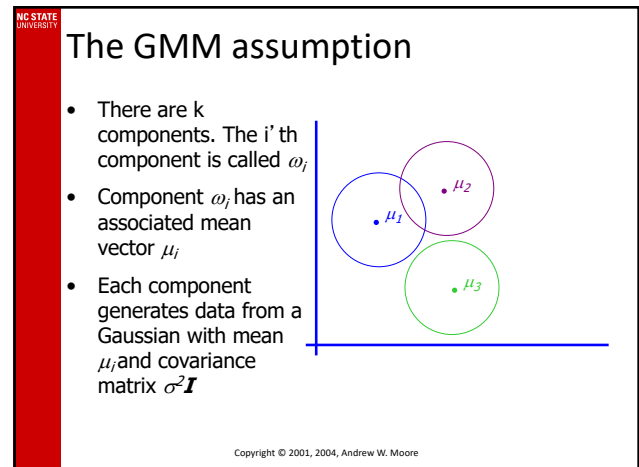
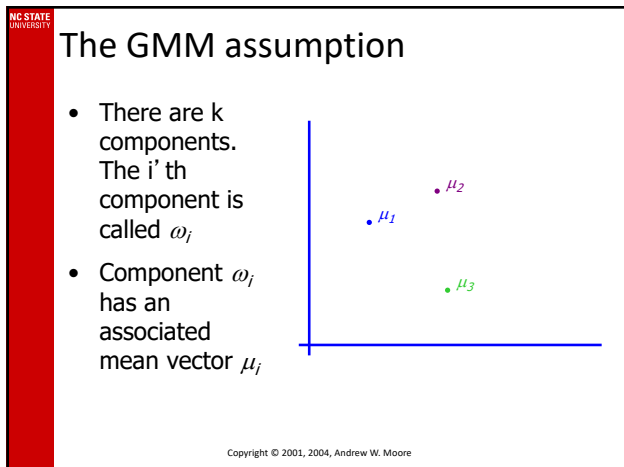
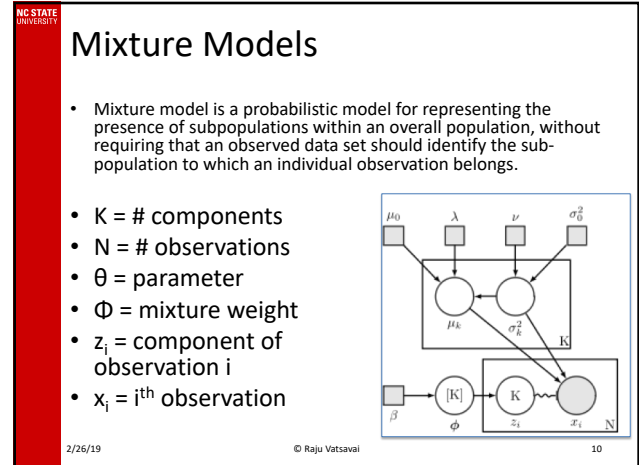
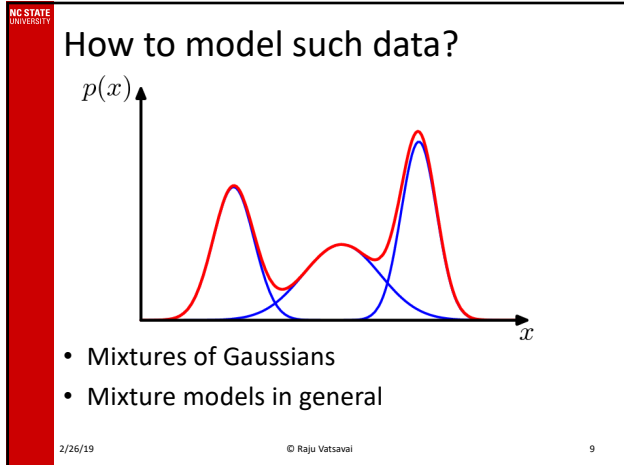


- Real world data are rarely unimodal
- Many times we don't know the labels for all components

2/26/19

© Raju Vatsavai

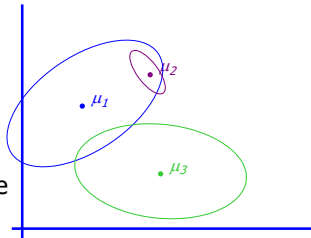
8



The General GMM assumption

Assume that each data point is generated according to the following recipe:

1. Pick a component at random. Choose component i with probability $P(\omega_i)$.
2. Data point $\sim N(\mu_i, \Sigma_i)$

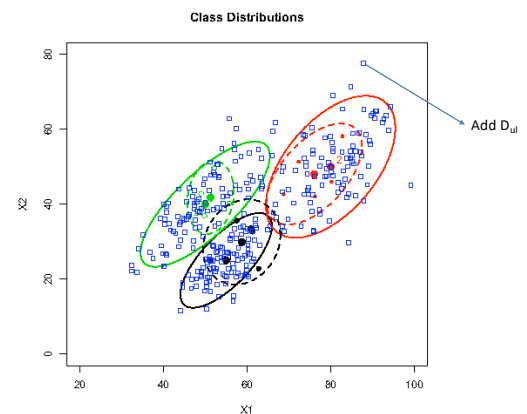
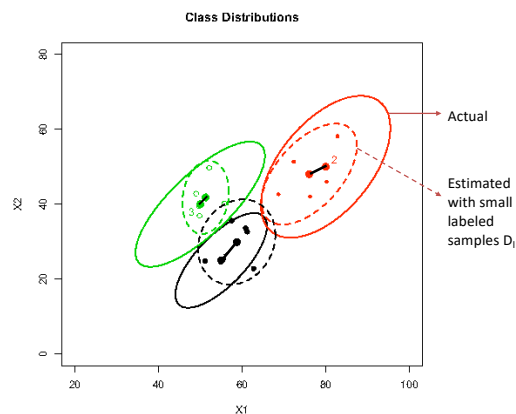


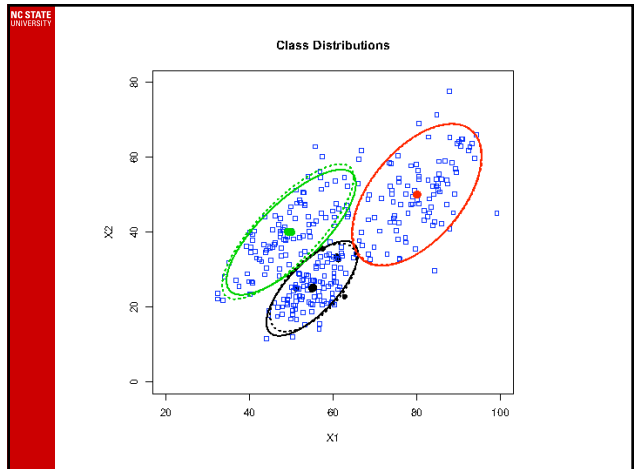
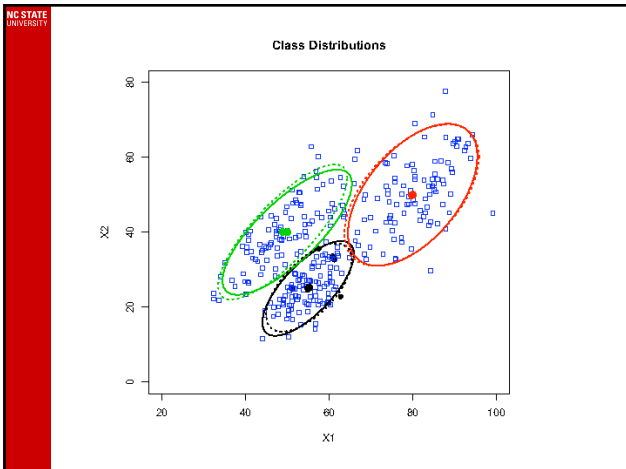
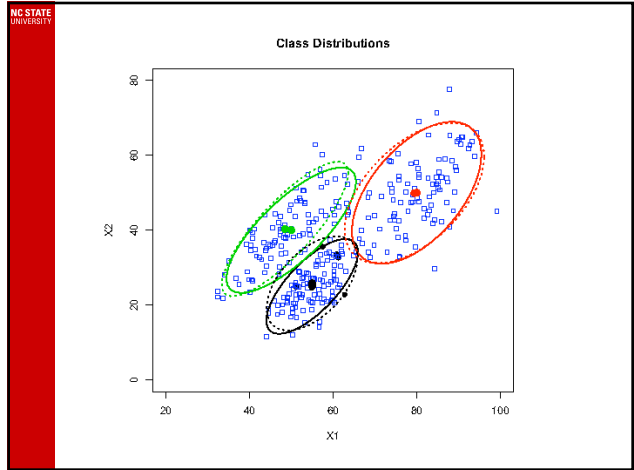
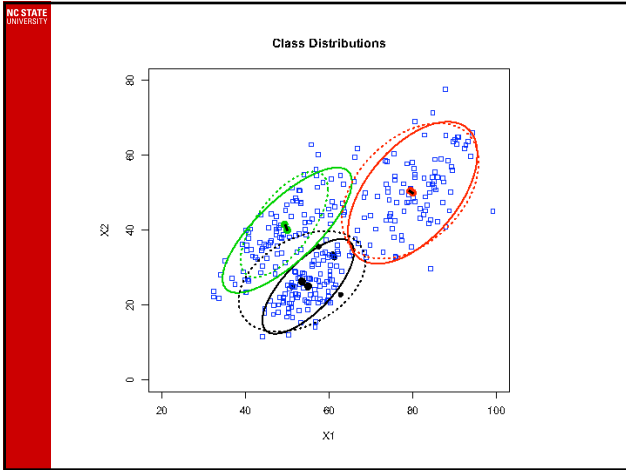
Copyright © 2001, 2004, Andrew W. Moore

GMM

- GMM ($K=M$) $p(x|\theta) = \sum_{i=1}^M \alpha_i p_i(x|\theta_i)$
 where $\theta = (\alpha_1, \dots, \alpha_M; \theta_1, \dots, \theta_M)$
 such that $\sum_{i=1}^M \alpha_i = 1$, $0 < \alpha_i < 1$ and
 p_i pdf parameterized by θ_i
- Maximize

$$L(\theta) - L(\theta_i) = \ln \frac{\sum_z p(x|z, \theta) p(z|\theta)}{p(x|\theta_i)}$$





NC STATE UNIVERSITY

Semi-supervised Learning

EM to estimate GMM parameters

- E-Step**

$$e_{ij} = \frac{|\Sigma_j|^{-1/2} \exp\left\{-\frac{1}{2}(x_i - \hat{\mu}_j)^T \hat{\Sigma}_j^{-1} (x_i - \hat{\mu}_j)\right\}}{\sum_{i=1}^N |\Sigma_i|^{-1/2} \exp\left\{-\frac{1}{2}(x_i - \hat{\mu}_i)^T \hat{\Sigma}_i^{-1} (x_i - \hat{\mu}_i)\right\}}$$
- M-Step**

$$\alpha_j = \frac{\sum_{i=1}^N e_{ij}}{N}, \quad \hat{\mu}_j^{k+1} = \frac{\sum_{i=1}^N e_{ij} x_i}{\sum_{i=1}^N e_{ij}}$$

and $\hat{\Sigma}_j^{k+1} = \frac{\sum_{i=1}^N e_{ij} (x_i - \hat{\mu}_j^{k+1})(x_i - \hat{\mu}_j^{k+1})^T}{\sum_{i=1}^N e_{ij}}$

ith data vector, jth class

NC STATE UNIVERSITY

Semi-supervised Learning

10 Classes, 100 Training Samples (10-30) x No of dimensions / class

Small Subset of 20 Training Samples

20 labeled + 80 unlabeled samples

Supervised (IC) vs. Semi-supervised (IC-EM)

Accuracy

Fixed Unlabeled (85) and Varying (Increasing) Labeled

Ranga Raju Vatsavai, Shashi Shekhar, Thomas E. Burk: A Semi-Supervised Learning Method for Remote Sensing Data Mining, ICITAI 2005: 207-211