



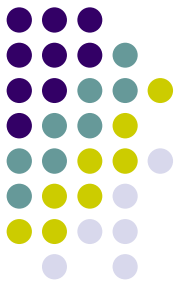
Adopting Semi-supervised Learning Algorithms for Mining Remote Sensing Imagery: Summary of Results and Open Research Problems

Ranga Raju Vatsavai^{1,3},
Shashi Shekhar¹, and Thomas E. Burk².

¹Spatial Database Research Group,
Department of Computer Science.

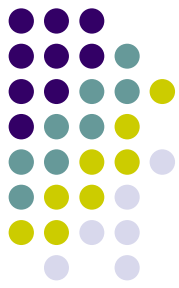
²Remote Sensing and Geospatial Analysis Lab,
Department of Forest Resources.
University of Minnesota.

³IBM-Research, India Research Laboratory
Indian Institute of Technology-Delhi, India.



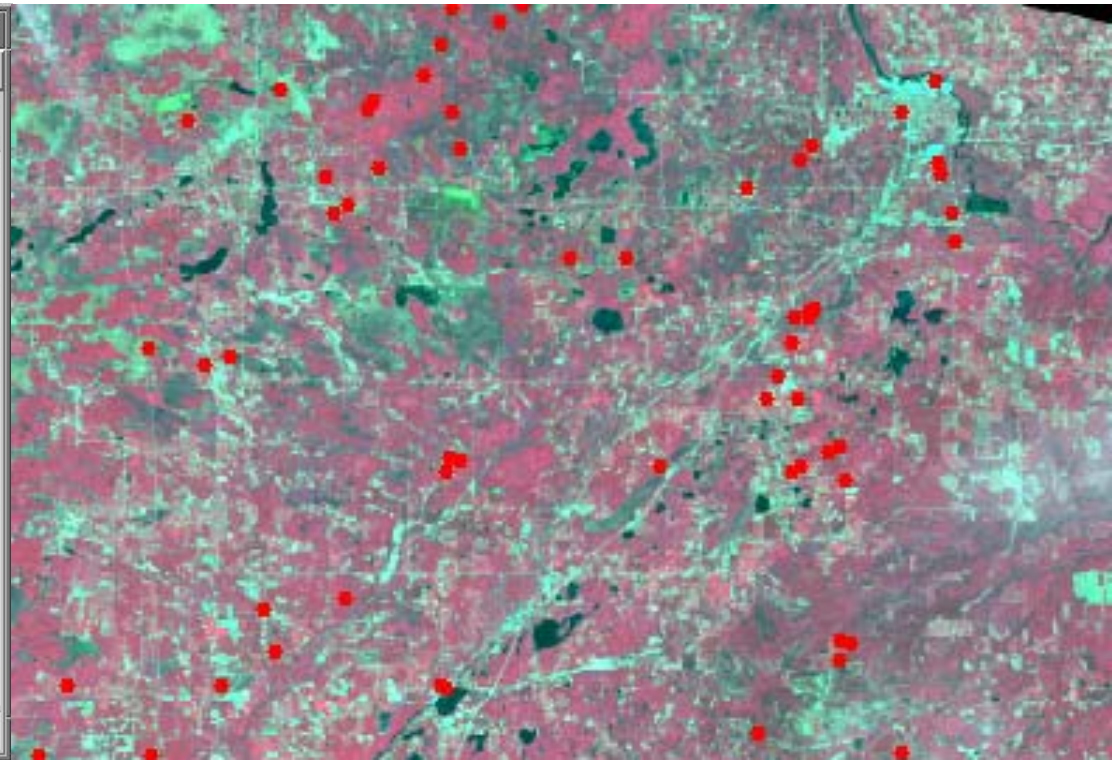
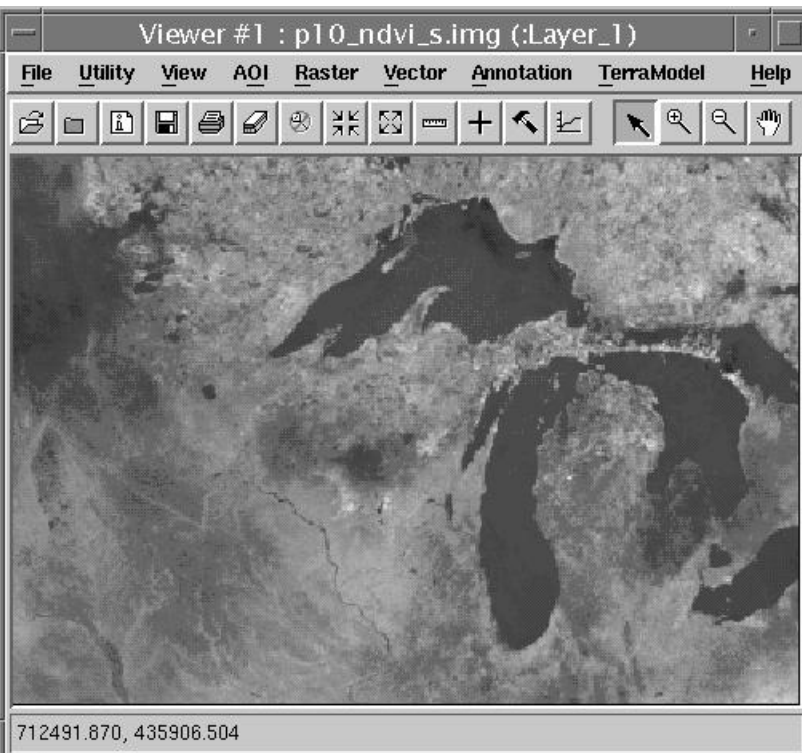
Outline

- Introduction
- Semi-supervised Classification
- Results
- Conclusions and Open Research Problems



Introduction

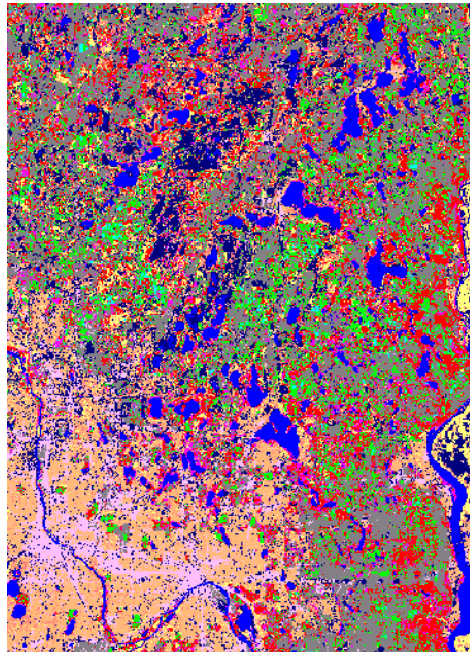
- Remote Sensing
 - Spectral Resolution – multi-spectral to hyper-spectral
 - Spatial Resolution – low (60m) to high (1m)








Introduction

- Classification*
 - Supervised – MLC, MAP, DT, NN, ...
 - Unsupervised – Clustering (kNN, kMeans, GMM...)

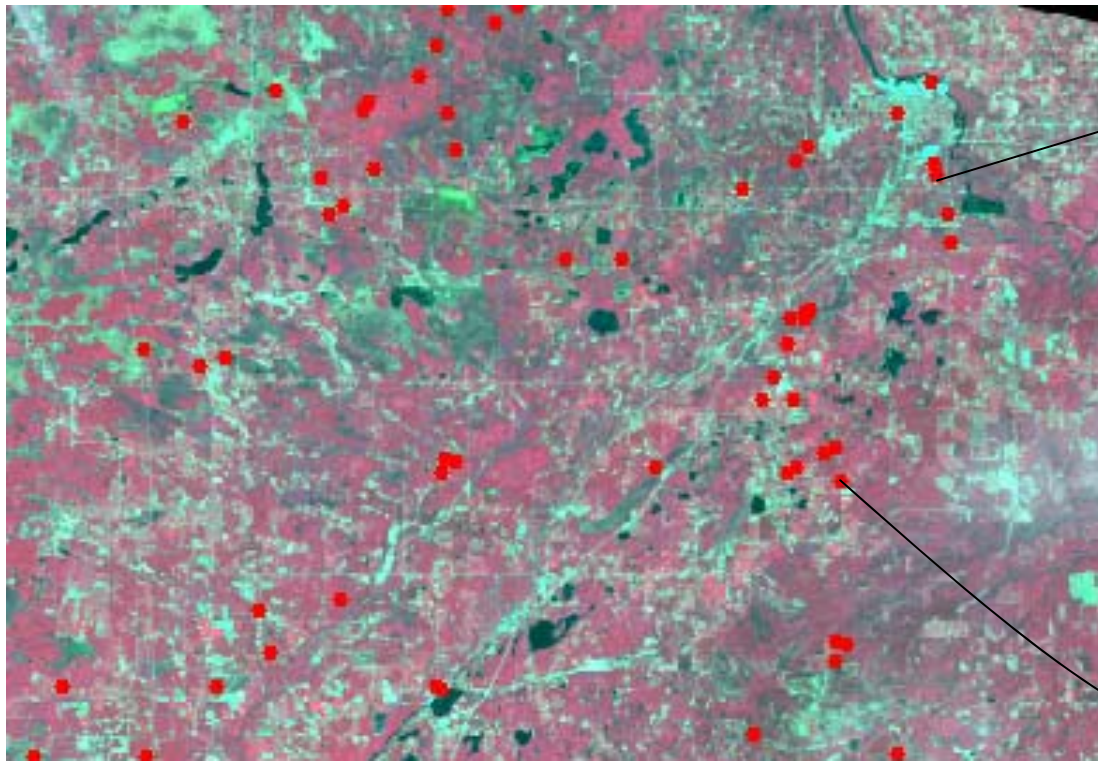


| | | | |
|---|-----------------|---|--------------------|
|  | Bare |  | Wetlands |
|  | Crop |  | Low Density Urban |
|  | Grass |  | High Density Urban |
|  | Upland Conifer |  | Lowland Conifer |
|  | Upland Hardwood |  | Lowland Hardwood |
|  | Water | | |

Supervised Classification Process



- Example Application - Thematic Mapping
 - 7Bands, 7000P×7000L, 30m Pixel size, 16day revisit, 170×183km.



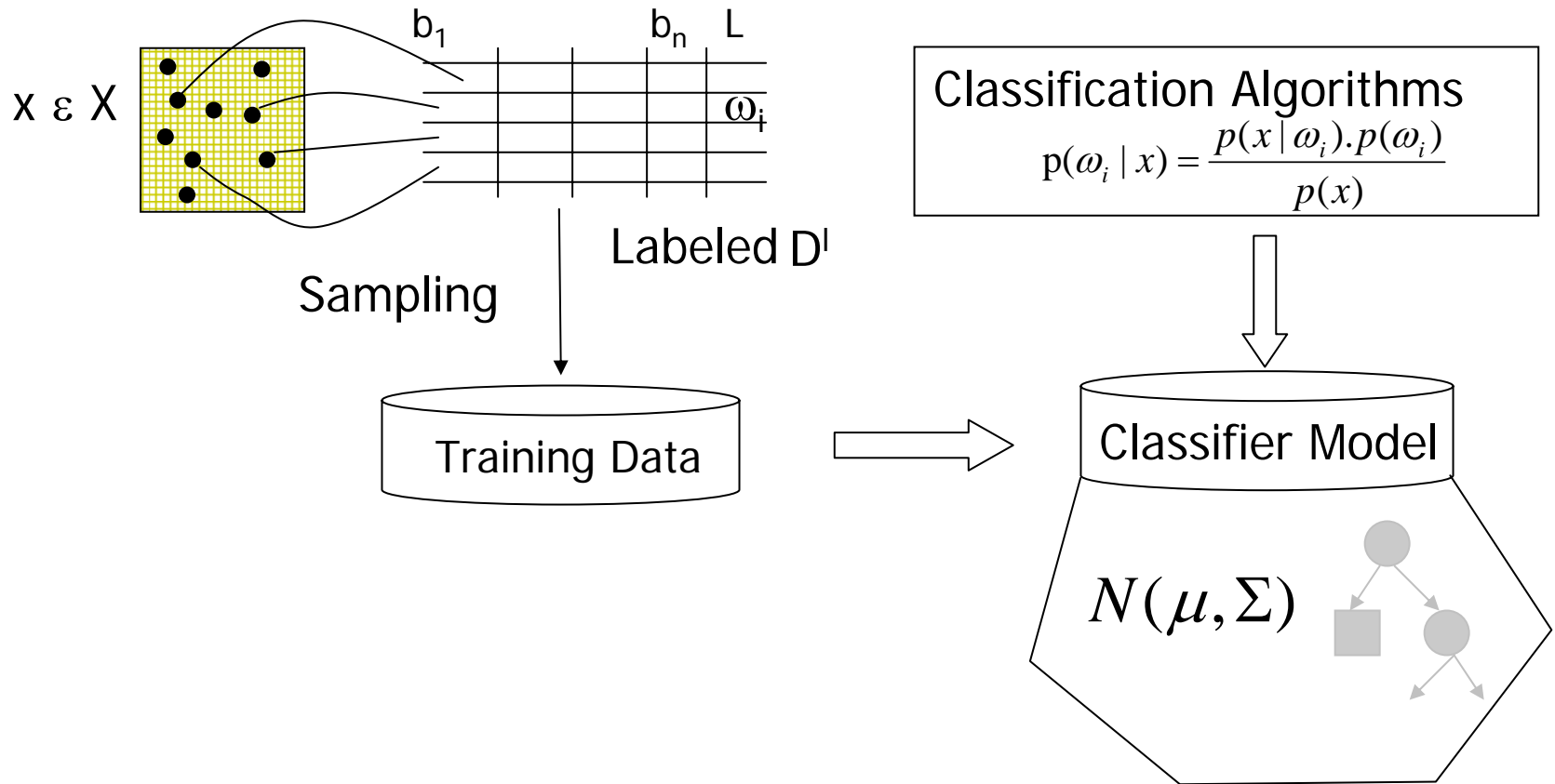
Sample plots
ID, Lat, Lon



Training Data
ID, b_1, \dots, b_n , Label



Supervised Classification Process





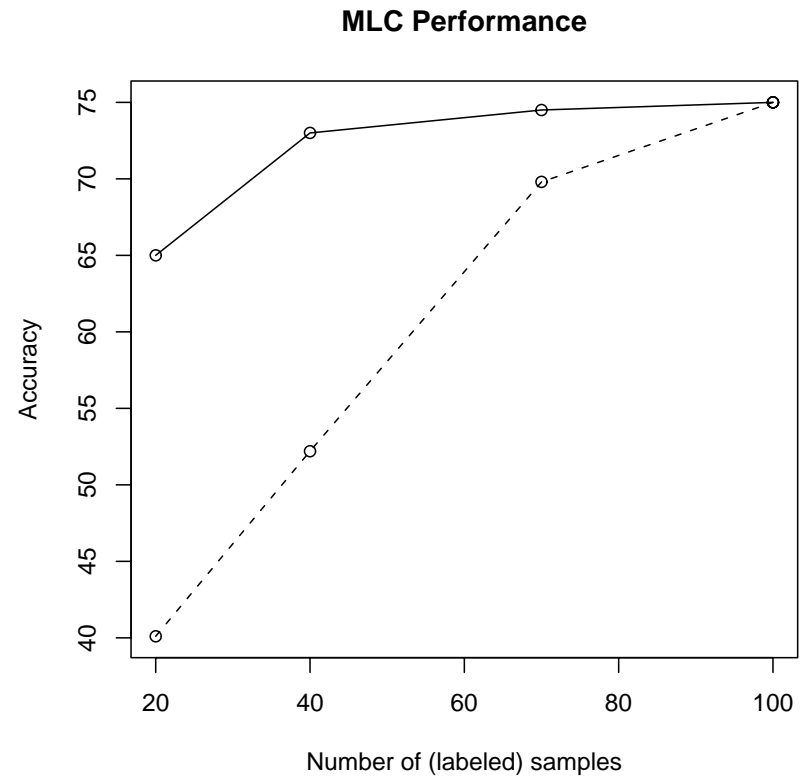
Supervised Classification Process

- Bayes' Theorem: $p(\omega_i | x) = \frac{p(x | \omega_i) \cdot p(\omega_i)}{p(x)}$
- Assuming $D_i = \{x_1, \dots, x_n\}$, and x 's i.i.d, the likelihood function is simply: $p(D_1 | \theta_1) = \prod_{k=1}^n p(x_k | \theta_1)$
- The maximum-likelihood estimation of θ_1 is the value that maximizes $p(D_1 | \theta_1)$. $\hat{\theta}_1 = \arg \max_{\theta_1} l(\theta_1)$
 $\hat{\mu}_1 = \frac{1}{n} \sum_{k=1}^n x_k$
- MLE gives familiar quantities: $\hat{\Sigma}_1 = \frac{1}{n} \sum_{k=1}^n (x_k - \hat{\mu}_1)(x_k - \hat{\mu}_1)'$
 - Asymptotically unbiased, consistent, efficient
 - All these properties are true if $n \rightarrow \infty$
 - $|D| \sim (10 \text{ to } 100) * (\text{number of dimensions})$

Supervised Classification Process



- MLC performance as $|D|$ increases



Symbols



x = Feature Vector

n = #of features or
#dimensions

μ = Mean Vector

One μ per class

Σ = Covariance Matrix

One Σ per class

Subscripts are context
dependent

$$x = \begin{bmatrix} b_1 \\ \dots \\ b_n \end{bmatrix}_{n \times 1} = \begin{bmatrix} 63 \\ 53 \\ 35 \\ 121 \\ 76 \\ 31 \end{bmatrix}, \quad \mu_\omega = \begin{bmatrix} m_1 \\ \dots \\ m_n \end{bmatrix}_{n \times 1} = \begin{bmatrix} 67.33 \\ 60.44 \\ 41.44 \\ 89.77 \\ 79.66 \\ 40.44 \end{bmatrix}$$

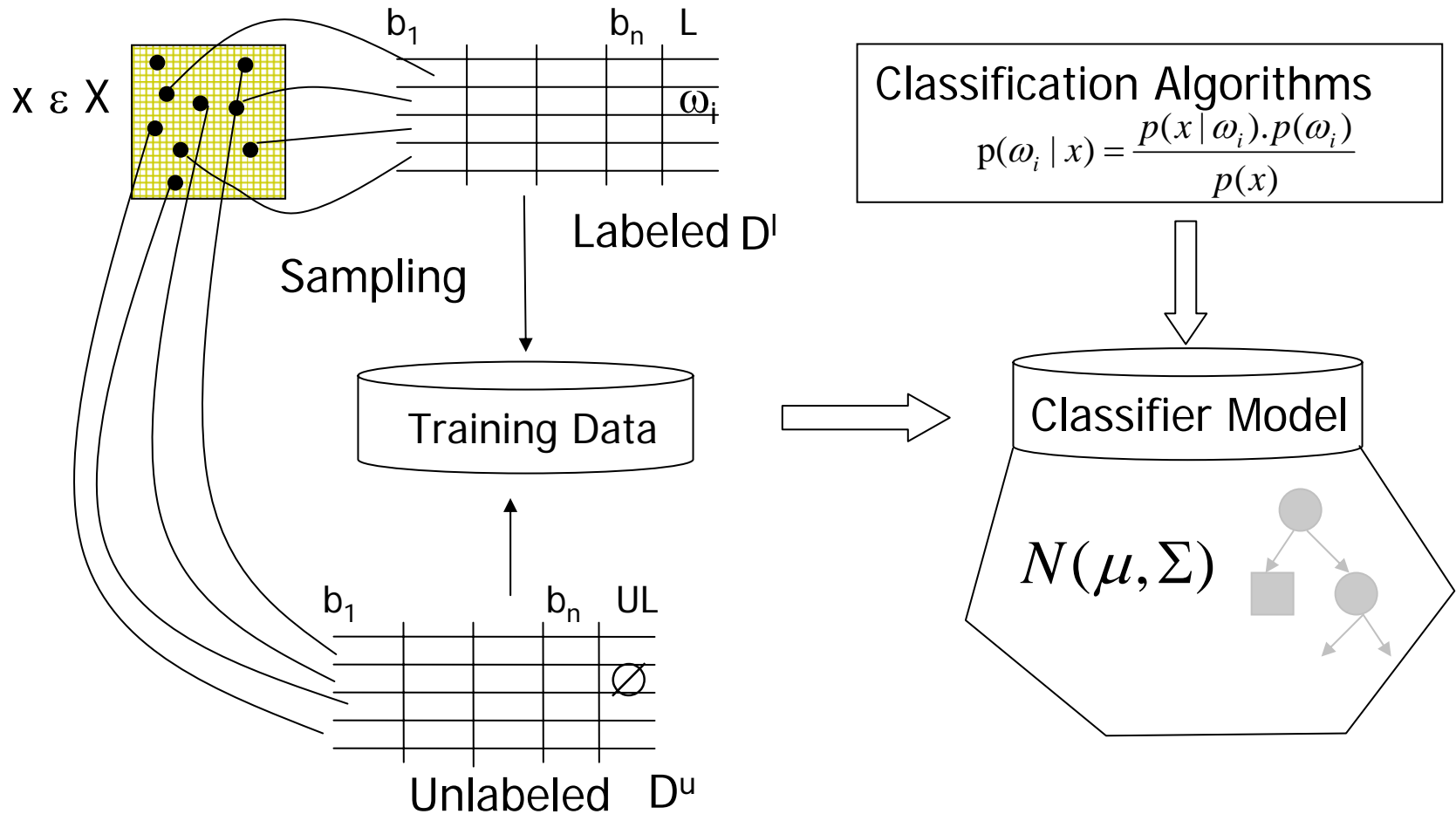
$$\Sigma_\omega = \begin{bmatrix} v_{1,1} & \dots & \dots \\ \dots & \dots & \dots \\ \dots & \dots & v_{n,n} \end{bmatrix}_{n \times n} = \begin{bmatrix} 2.0 & .2 & .33 & .58 & 1.13 & -0.7 \\ .41 & 1.7 & .6 & .2 & -.3 & -1.2 \\ & & & & & \\ & & & & & \\ -1.16 & & & & & \\ & & & & & 2.7 \end{bmatrix}$$

Semi-supervised Classification



- Acquiring Labels
 - Costly, Time consuming, Labor intensive
 - Error prone
 - May be impossible at times –
 - Emergency situations (fires, floods, cyclones, ...)
 - Accessibility, Privacy
- Samples are readily available, but not labels
- There are several classification problems that have to deal with insufficient learning samples
- Q? Can we still learn with partially labeled training data

Semi-supervised Classification Process





Learning with incomplete data $D = D^l \cup D^u$

- Assume current estimate of parameter is θ_i
- See what happens to L when new θ is estimated

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

- We would like to choose θ to maximize r.h.s
- EM*, introduce unobserved vars Z such that if Z is known, the θ is computed easily

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

*A. Dempster, N. M. 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.



Learning with incomplete data $\mathcal{D} = \mathcal{D}^l \cup \mathcal{D}^u$

- The EM algorithm first finds the expected value of the complete-data log-likelihood –

- E-step
$$Q(\theta; \theta^{(i)}) = E[\log p(D^l, D^u; \theta) | D^l; \theta^{(i)}]$$

- The second step (M-step) of the EM algorithm is to maximize the E[] of 1st step.

- M-step
$$\theta^{(i+1)} = \arg \max_{\theta} Q(\theta, \theta^{(i)})$$



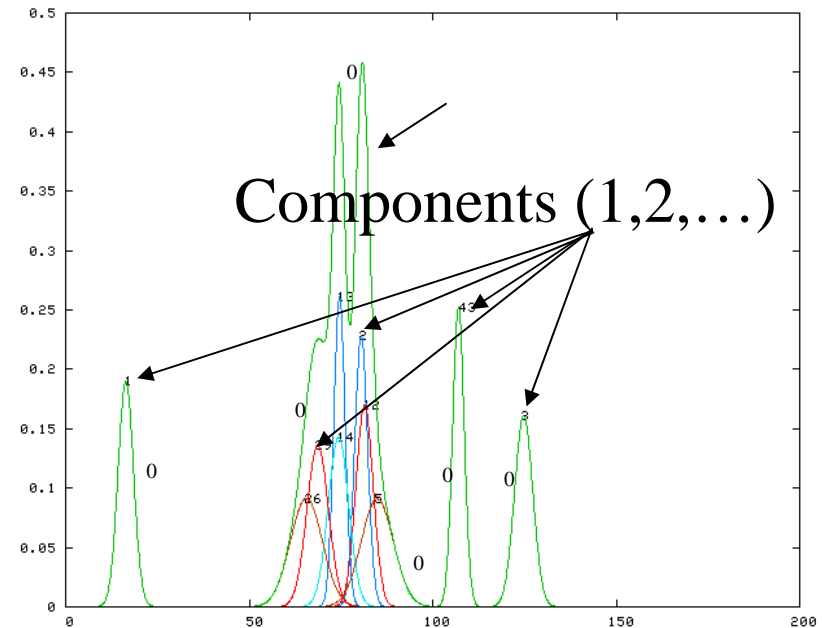
EM update equations for GMM*

- Assume that D is drawn from a Gaussian Mixture Model, given by

$$p(x | \theta) = \sum_{i=1}^M \alpha_i p_i(x | \theta_i)$$

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



*J. Bilmes. A gentle tutorial on the em algorithm and its application to parameter estimation for gaussian mixture and hidden markov models. Technical Report, University of Berkeley, ICSI-TR-97-021, 1997., 1997.



GMM – Update Equations

- E-Step

$$e_{ij} = \frac{|\hat{\Sigma}_j^k|^{-1/2} \exp\left\{-\frac{1}{2}(x_i - \hat{\mu}_j^k)^T \hat{\Sigma}_j^{-1,k} (x_i - \hat{\mu}_j^k)\right\}}{\sum_{l=1}^M |\hat{\Sigma}_l^k|^{-1/2} \exp\left\{-\frac{1}{2}(x_i - \hat{\mu}_l^k)^T \hat{\Sigma}_l^{-1,k} (x_i - \hat{\mu}_l^k)\right\}}$$

- M-Step

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

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

i^{th} data vector, j^{th} class



GMM - New Update Equations

$$\alpha_j = \frac{\left(\lambda_l m_j + \lambda_u \sum_{i=1}^n e_{ij}\right)}{\lambda_l m + \lambda_u n},$$

$$\hat{\mu}_j^{k+1} = \frac{\left(\lambda_l \sum_{i=1}^{m_j} y_{ij} + \lambda_u \sum_{i=1}^n e_{ij} x_i\right)}{\lambda_l m_j + \lambda_u \sum_{i=1}^n e_{ij}},$$

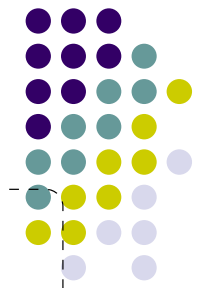
$$\hat{\Sigma}_j^{k+1} = \frac{\lambda_l \sum_{i=1}^{m_j} (y_{ij} - \hat{\mu}_j^{k+1})(y_{ij} - \hat{\mu}_j^{k+1})^T + \lambda_u \sum_{i=1}^n e_{ij} (x_i - \hat{\mu}_j^{k+1})(x_i - \hat{\mu}_j^{k+1})^T}{\lambda_l m_j + \lambda_u \sum_{i=1}^n e_{ij}}$$

i^{th} data vector, j^{th} class

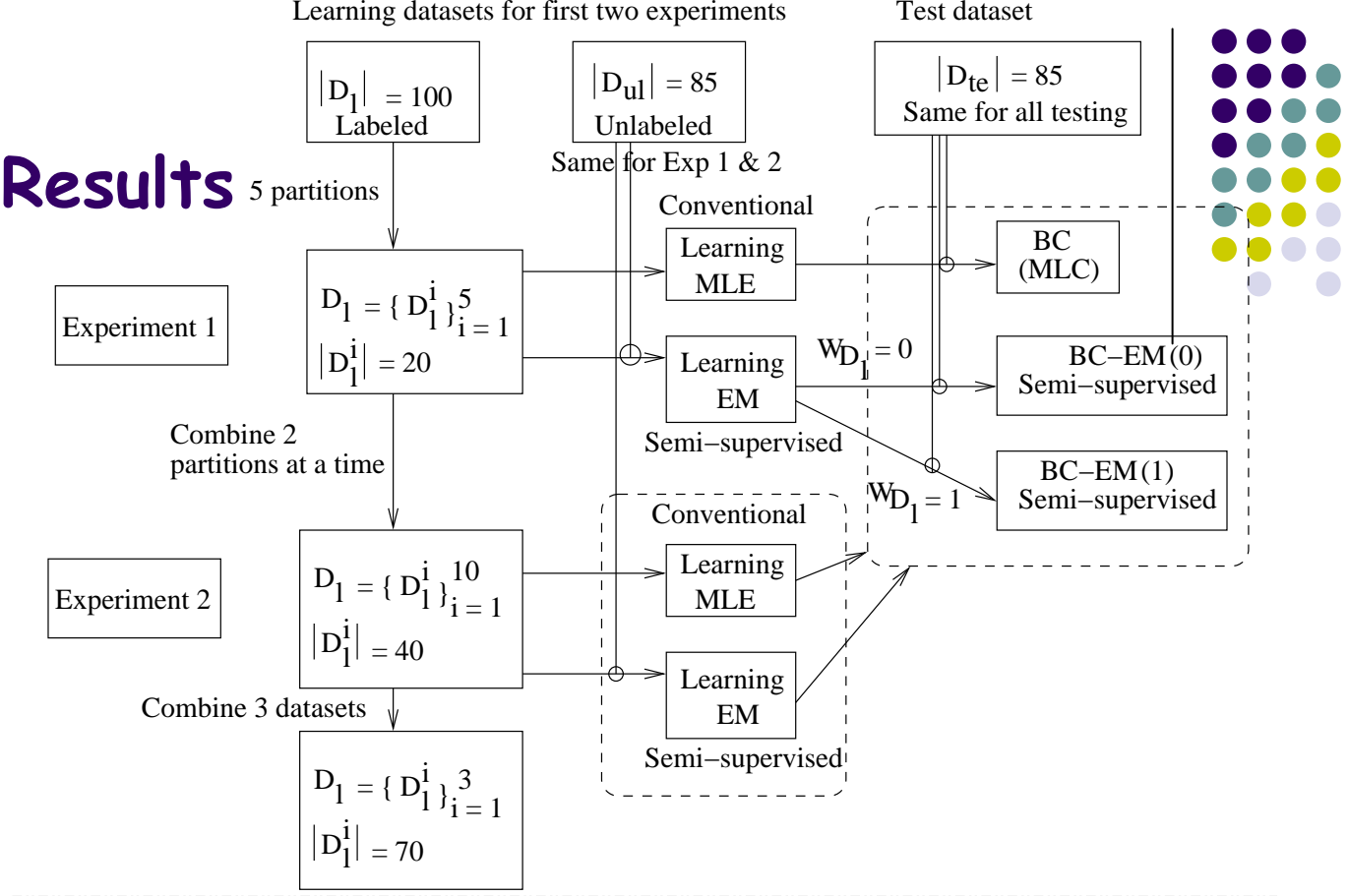


Semi-supervised learning Algorithm - Outline

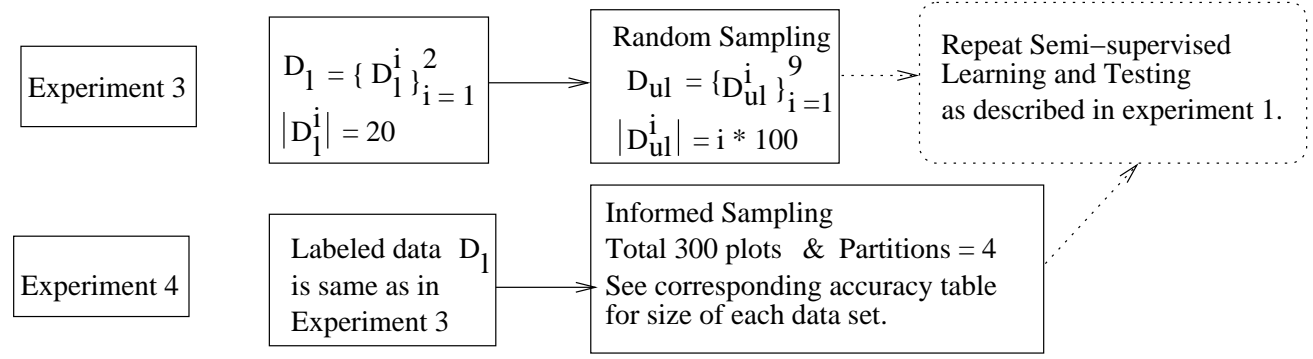
- Inputs: Collection D consisting of D^l and D^u
- Build initial classifier, estimate $\hat{\theta} = \arg \max_{\theta} l(\theta)$
- Loop while parameter estimate improves
 - (E-Step): Use current classifier estimate, $\hat{\theta}$, to estimate component membership of each data vector, i.e., the prob. that each mixture component (class) generated each sample
 - (M-Step): Re-estimate the classifier parameter, given the estimated component membership of each sample
- Output: A classifier, with improved estimates, that takes any new data sample and predicts a class label.



Experimental Results



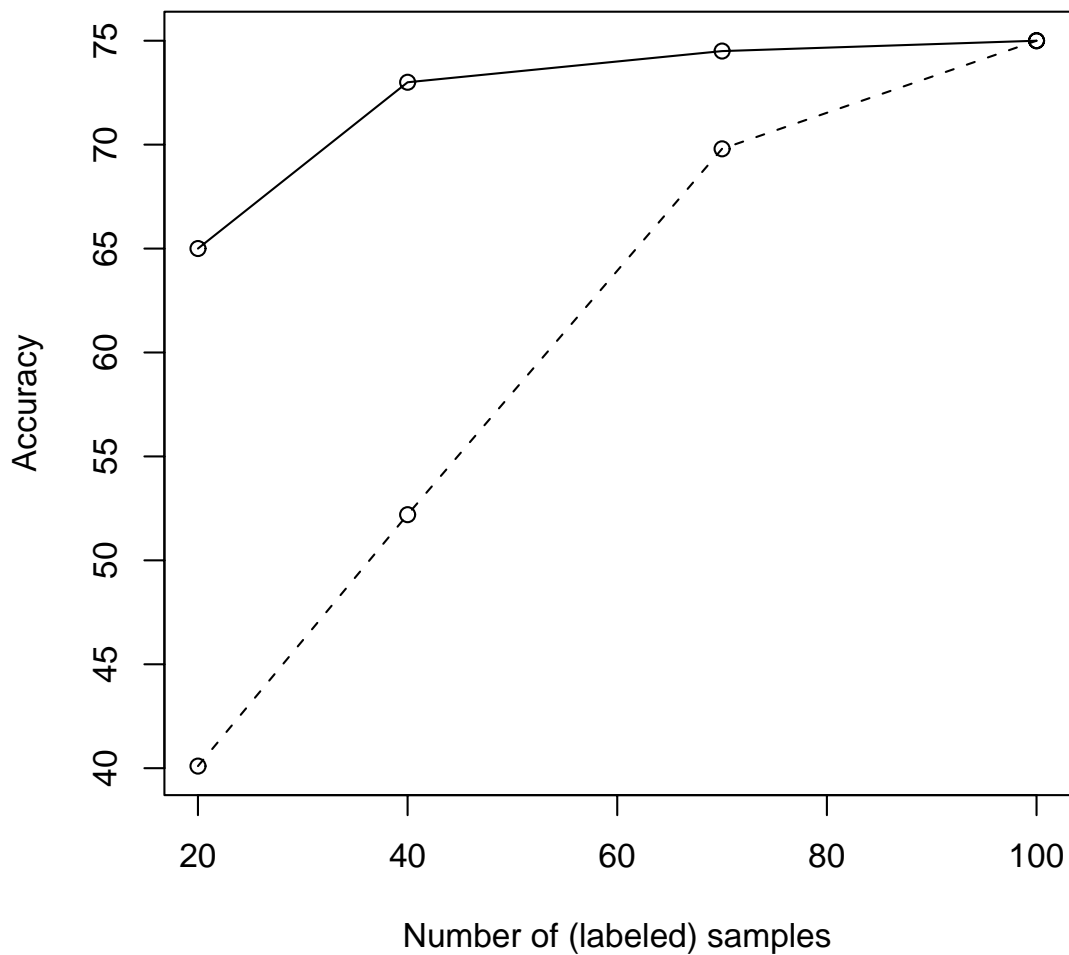
Learning Datasets for experiments 2 and 3 Test dataset is same as experiments 1 and 2
 Two learning datasets from Experiments 1 are chosen based on best (B20) and worst (W20) accuracies



Experimental Results



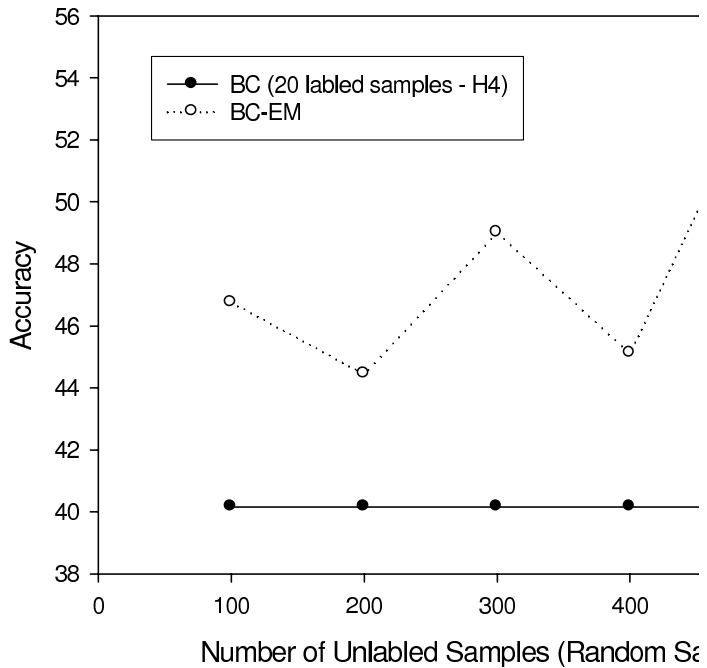
MLC Performance



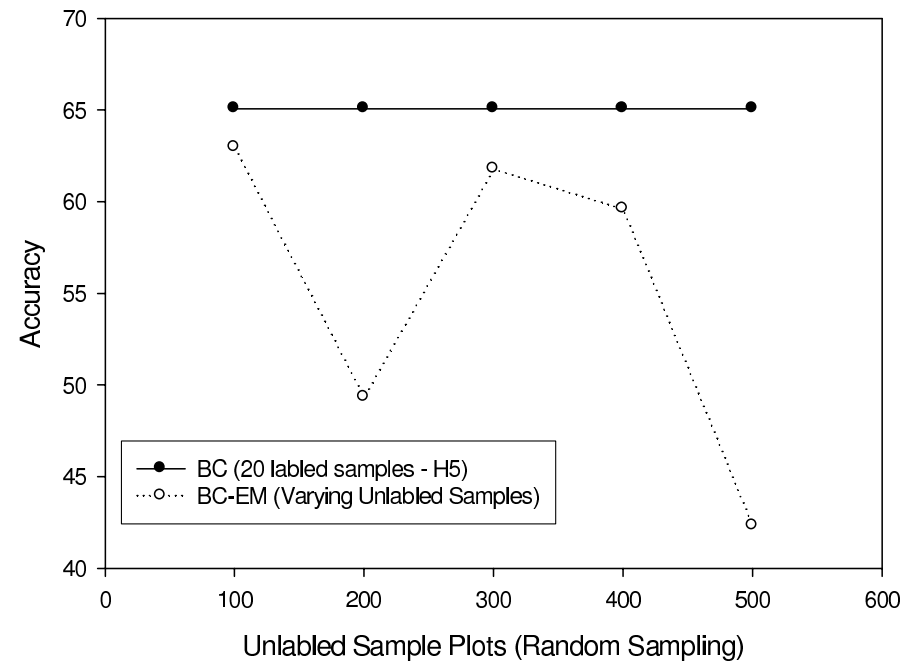
Experimental Results



Varying Number of Unlabeled (Random) Samples

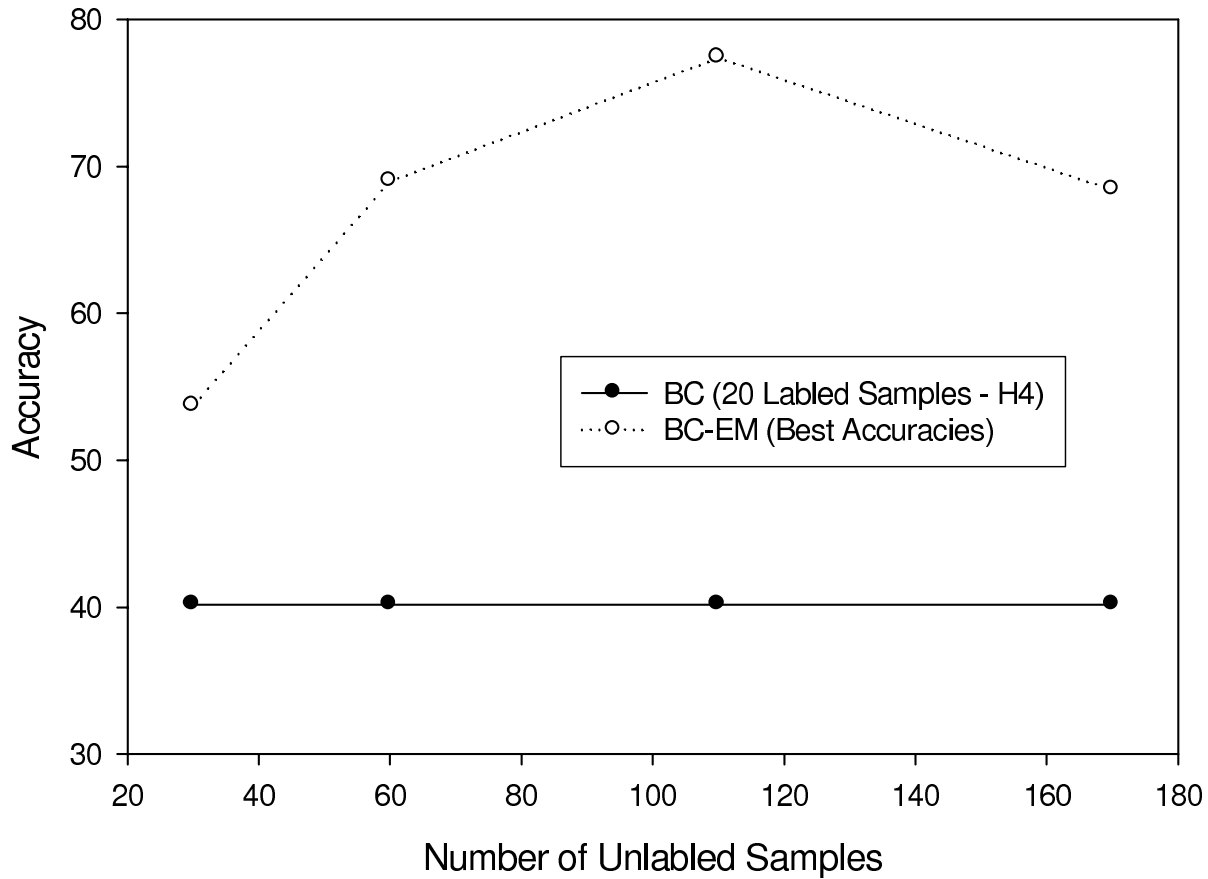


Varying Number of Unlabeled (Random) Sample Plots





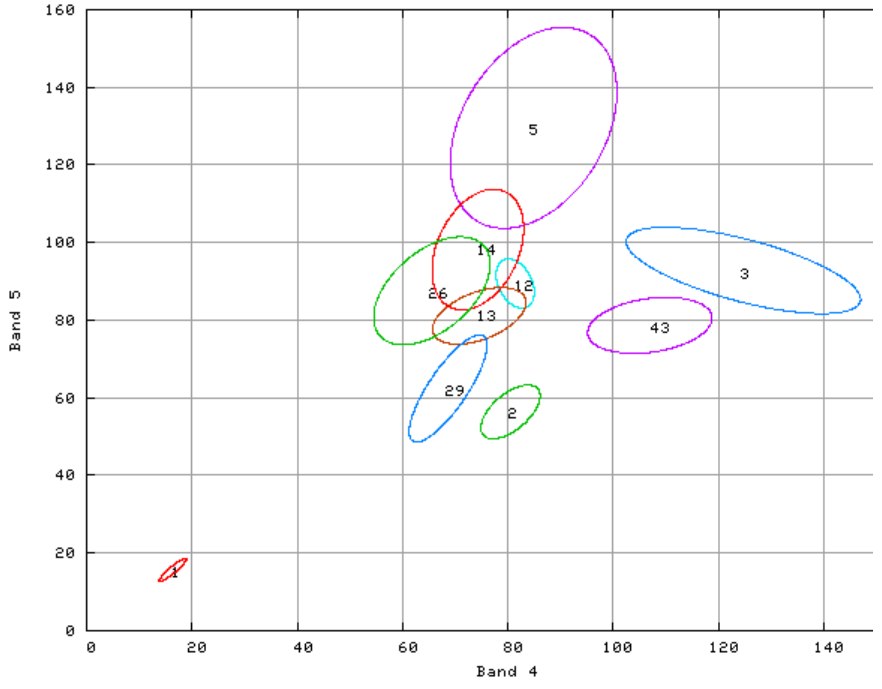
Varying Unlabeled Sample Plots (Informed Sampling)





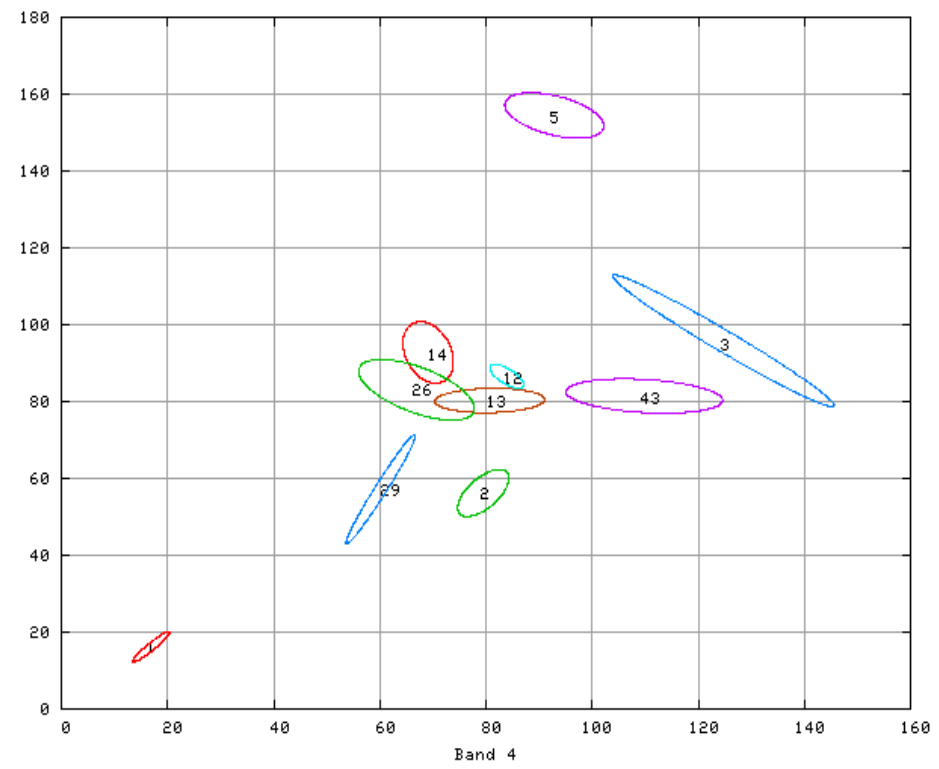
Experimental Results

Bivariate Normal Density Plot



MLE (100 labeled)

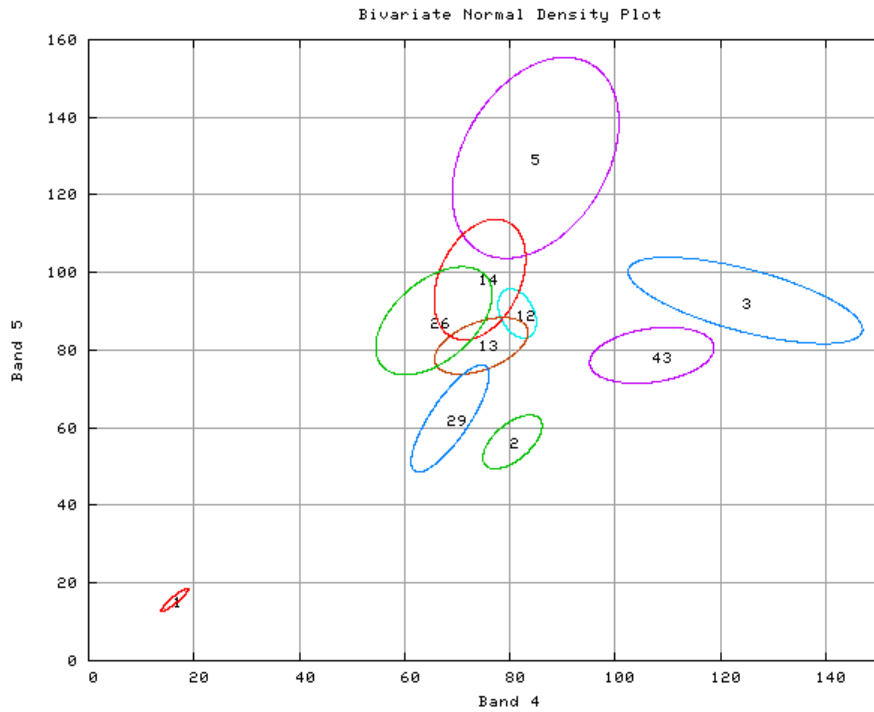
Bivariate Normal Density Plot



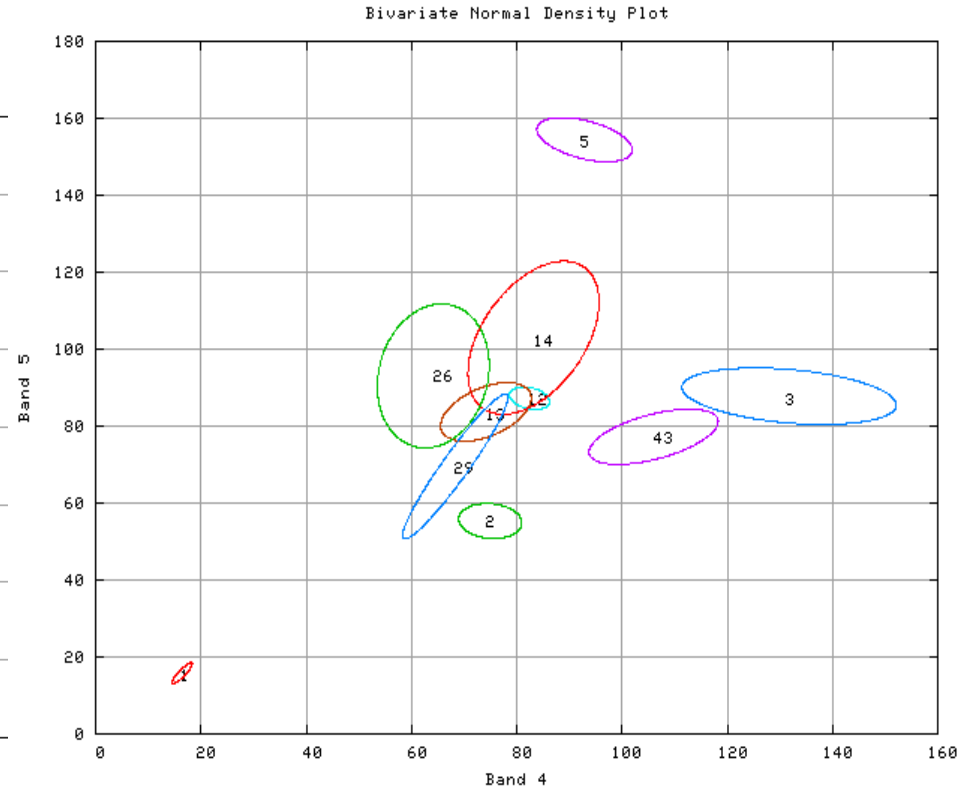
MLE (20 lab)



Experimental Results



MLE (100 labeled)



EM (20 lab + 85 ul)



Spatial Semi-supervised

- iid assumptions are not valid
- MAP/MRF
 - $p(c|x) = p(x|c)P(c)/p(x)$

$$\begin{aligned}P(c) &= P(c(i, j) | c(k, l); \{k, l\} \neq \{i, j\}). \\ &= P(c(i, j) | c(k, l); \{k, l\} \in s). \\ &= \frac{1}{Z} e^{\frac{-U(\omega)}{T}}\end{aligned}$$



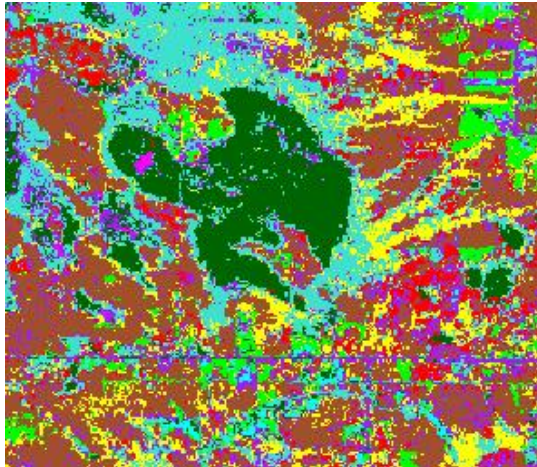
- Z is not computable
 - Assume
 - 10 classes,
 - 256x256 image
 - Total configurations
 - $10^{(256 \times 256)}$

$$U_s(C(i, j)) = \sum_{\{k, l\} \in s_{ij}} \beta I(C(i, j), C(k, l))$$

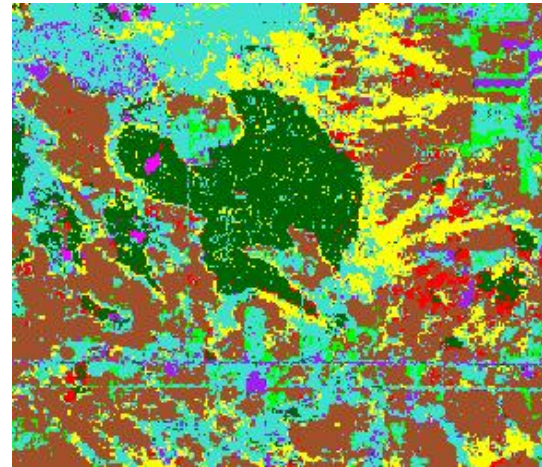
Where

$$I(C(i, j), C(k, l)) = \begin{cases} 1 & \text{if } C(i, j) = C(k, l) \\ 0 & \text{if } C(i, j) \neq C(k, l) \end{cases}$$

Spatial Semi-supervised

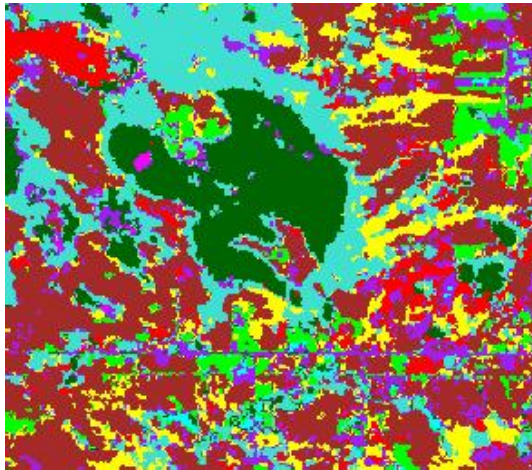


BC (60%)

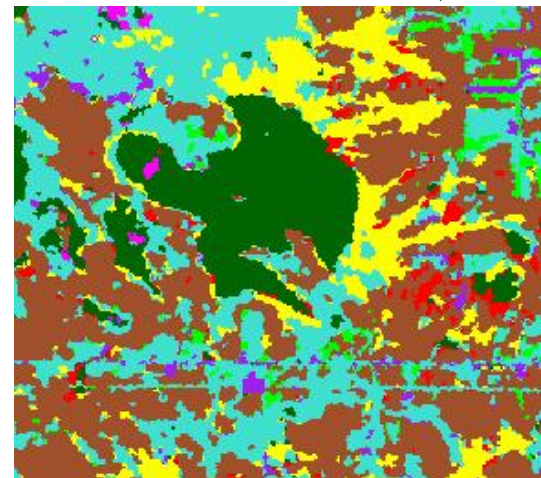


BC-EM (68%)

BC-MRF (65%)



BC-EM-MRF (72%)





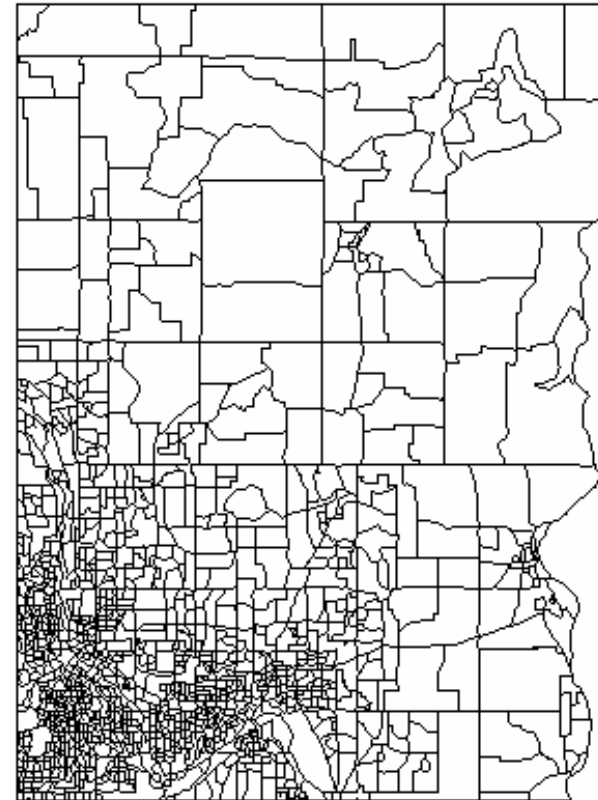
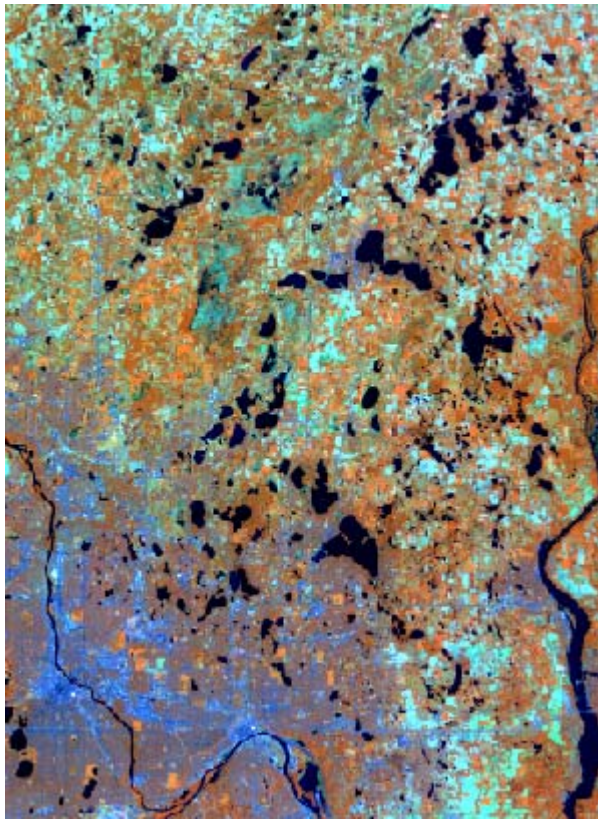
Open Research Problems

- Care should be exercised when collecting unlabeled samples
 - No samples from small classes
 - Mixed plots
 - Samples from unknown components
- Extend the model for multi-source data

Open Research Problems



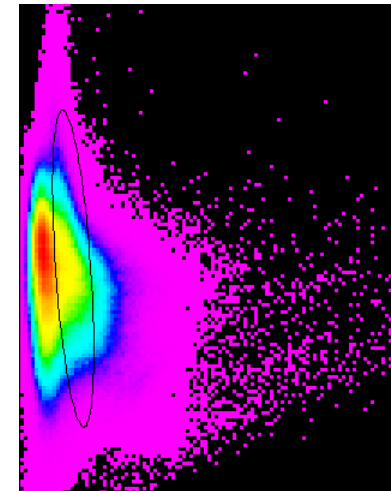
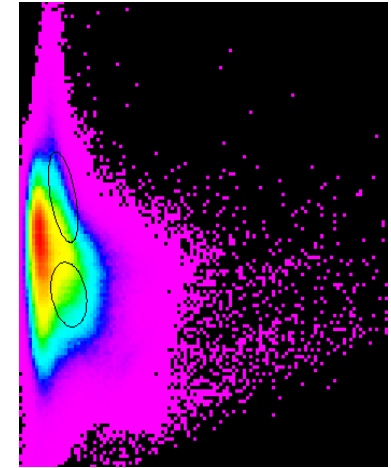
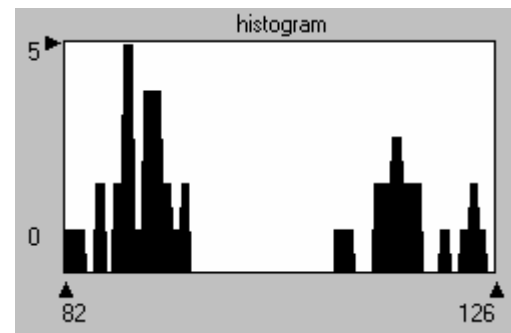
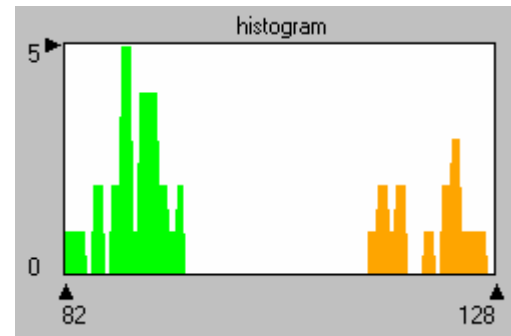
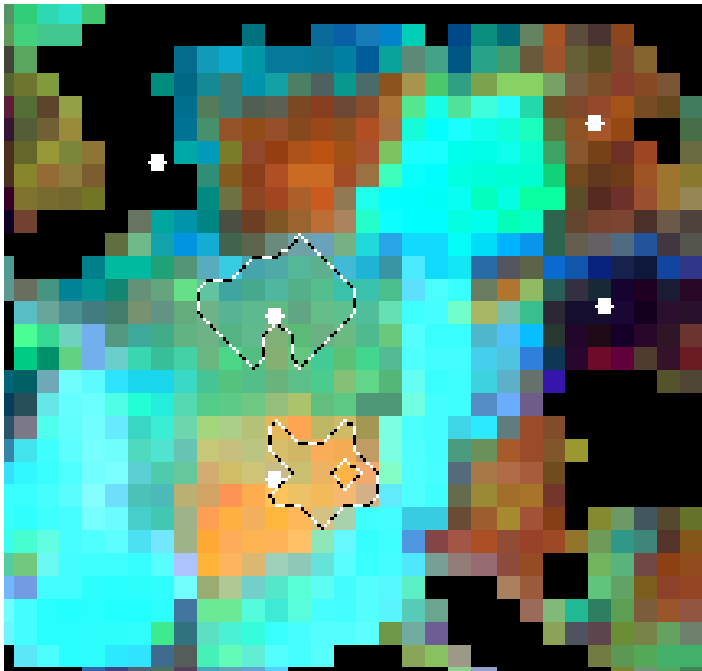
- Extend the model for multi-source data





Open Research Problems

- Discovering sub-components (e.g., Anderson Level 2 and 3 classes)





Open Research Problems

- Efficient Algorithms
- Global Optimization (Stochastic versions)
- Semi-supervised MCC
- Semi-supervised + Active Learning

- Not a research problem, but
 - Can we create bench mark dataset(s) similar to CMU or UCI for Remote Sensing Data Mining.

- Conclusions
 - Results were promising



Acknowledgements

- Prof. Joydeep Ghosh, UT Austin
- Team @ Remote Sensing & Geospatial Analysis Lab
- Team @ Spatial Databases Research Group, C.Sci. Dept. U of Minnesota.