# CS 175: Project in Artificial Intelligence
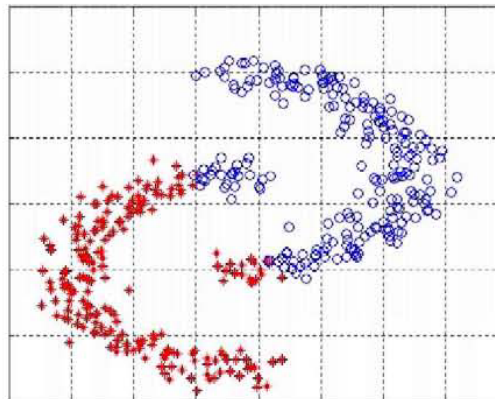
## Slides 5: Clustering

# **Topic 7: Clustering**

Some slides taken from Prof. Ihler

# Clustering

- So far: *supervised* learning
  - Given observed features and targets
    - Predict targent: class labels, stock prices, etc.
- Today: *Unsupervised* learning
  - Only attributes (features)
  - Want to discover **structure** in the data

- Ex: the data may be concentrated in **clusters**



*Is this a good clustering?*

# Why unsupervised learning?

- Often we've wanted to change data representations
  - Add more features (polynomials, etc.)
  - Select good features (boosting & decision stumps)

- Unsupervised learning is the same problem
  - Produce a new representation of same data
  - New representation should be more meaningful
  - Could be used in later steps (classification, etc)

- Smaller representation
  - Computationally less expensive
  - Low storage
    - Ex: store just cluster label, rather than attribute values
  - Might avoid overfitting

- Might simplify prediction problem as well
  - Netflix: predictions based on "most similar" users or movies

# K-Means Clustering

- A simple clustering algorithm

- Iterate between
  - Updating the assignment of data to clusters
  - Updating the cluster's summarization

- Suppose we have K clusters, c=1..K
  - Represent clusters by locations $\mu_c$
  - Example i has features $x_i$
  - Represent assignment of $i^{th}$ example $z_i \in 1..K$

- Iterate until convergence:
  - For each datum, find the closest cluster

  $$z_i = \arg\min_c \| x_i - \mu_c \|^2 \qquad \forall i$$

  - Set each cluster to the mean of all assigned data:

  $$\forall c, \qquad \mu_c = \frac{1}{N_c} \sum_{i \in S_c} x_i \qquad\qquad S_c = \{ i : z_i = c \}, \ N_c = |S_c|$$

# K-Means as optimization

- Optimizing the cost function

$$C(\underline{z}, \underline{\mu}) = \sum_i \|x_i - \mu_{z_i}\|^2$$

- Greedy descent technique
  - Steps
    - Choose closest cluster
    - Choose mean of assigned data
  - Each step only decreases the cost  (why?)

- As with any descent method, beware of local minima
  - Algorithm behavior depends significantly on initalization
  - Many heuristics
    - Random (not bad); Farthest (sensitive); some mix maybe?

# Choosing the number of clusters

- With cost function

$$C(\underline{z}, \underline{\mu}) = \sum_i \|x_i - \mu_{z_i}\|^2$$

  what is the optimal value of k?

- Can increasing k ever increase the cost?

- This is a model complexity issue
  - Much like choosing lots of features – they only (seem to) help
  - But we want our clustering to *generalize* to new data

- One solution is to penalize for complexity
  - Bayesian information criterion (BIC)
  - Add  (# parameters) * log(N) to the cost
  - Now more clusters can increase cost, if they don't help "enough"

# Hierarchical  Agglomerative Clustering
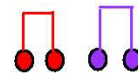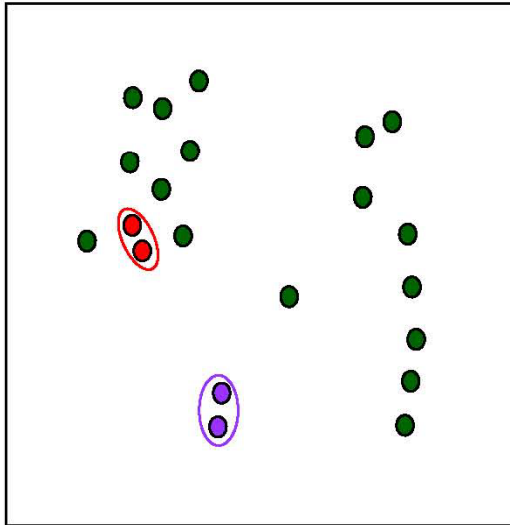
- Another simple clustering alg

Initially, every datum is a cluster

- Define a distance between clusters (return to this)

- Initialize: every example is a cluster

- Iterate:

  – Compute distances between all clusters
    (store for efficiency)

  – Merge two closest clusters

- Save both clustering and *sequence* of cluster ops
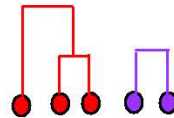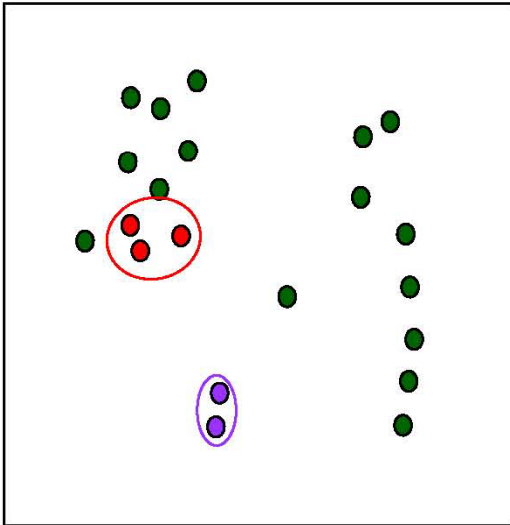
- "Dendrogram"

# Iteration 1

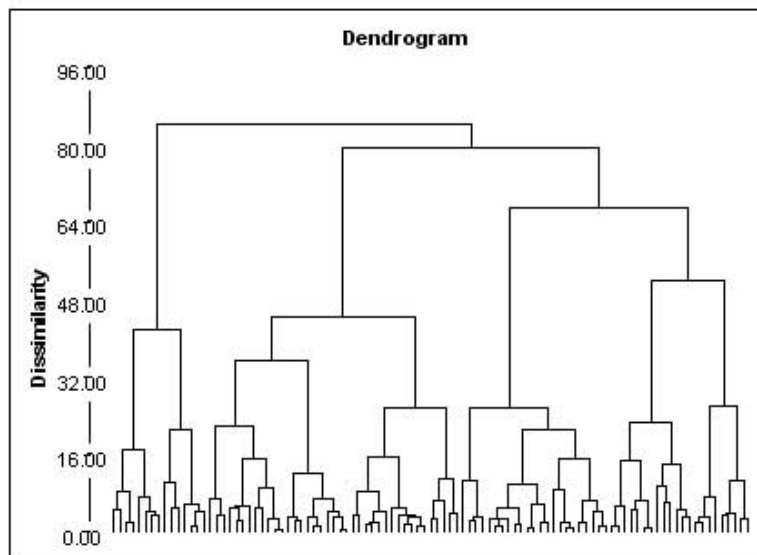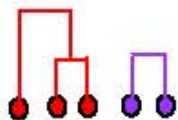# Iteration 2

# Iteration 3

- Builds up a sequence of clusters ("hierarchical")



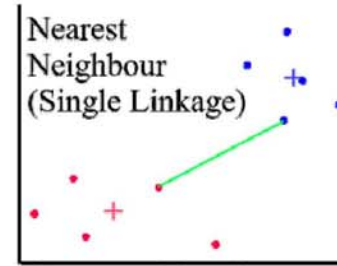- Algorithm complexity   $O(N^2)$ (Why?)

**In matlab: "linkage" function  (stats toolbox)**
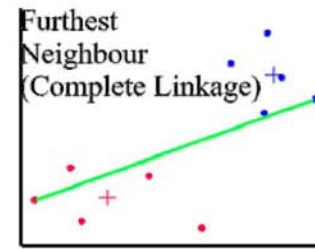
# Dendrogram

# Cluster Distances

$$D_{\min}(C_i, C_j) = \min_{x \in C_i, \ y \in C_j} \|x - y\|^2$$



Nearest Neighbour (Single Linkage)
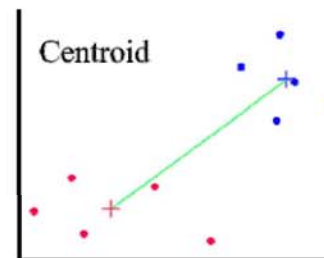
produces minimal spanning tree.

$$D_{\max}(C_i, C_j) = \max_{x \in C_i, \ y \in C_j} \|x - y\|^2$$



Furthest Neighbour (Complete Linkage)
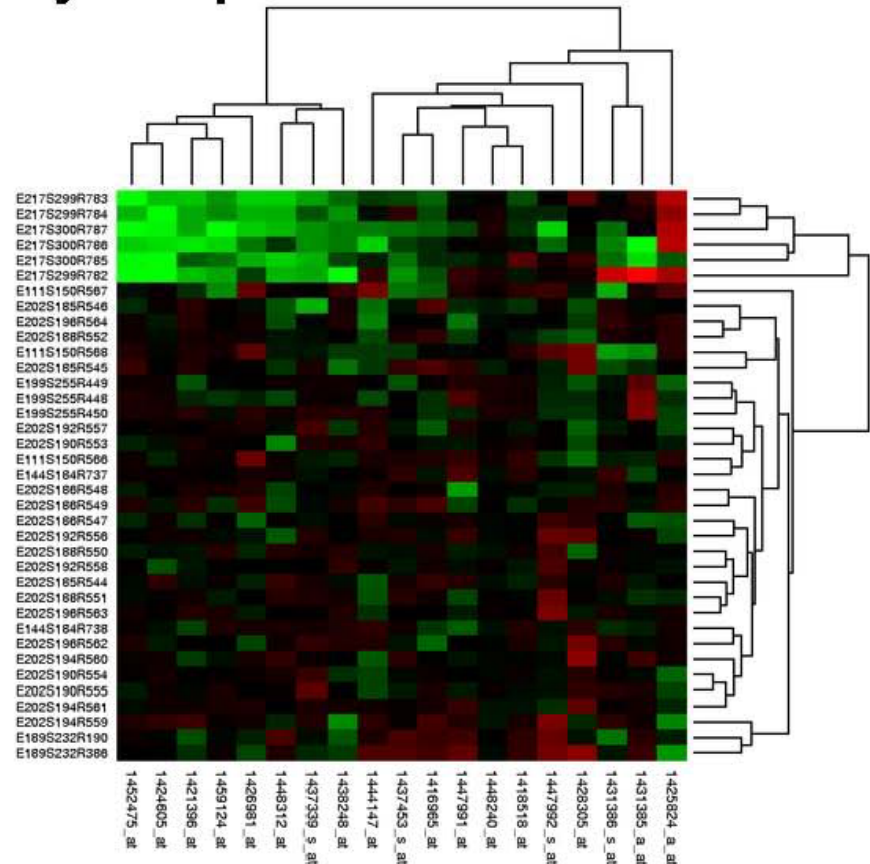
avoids elongated clusters.

$$D_{\mathrm{avg}}(C_i, C_j) = \frac{1}{|C_i||C_j|} \sum_{x \in C_i, \ y \in C_j} \|x - y\|^2$$

$$D_{\mathrm{means}}(C_i, C_j) = \|\mu_i - \mu_j\|^2$$
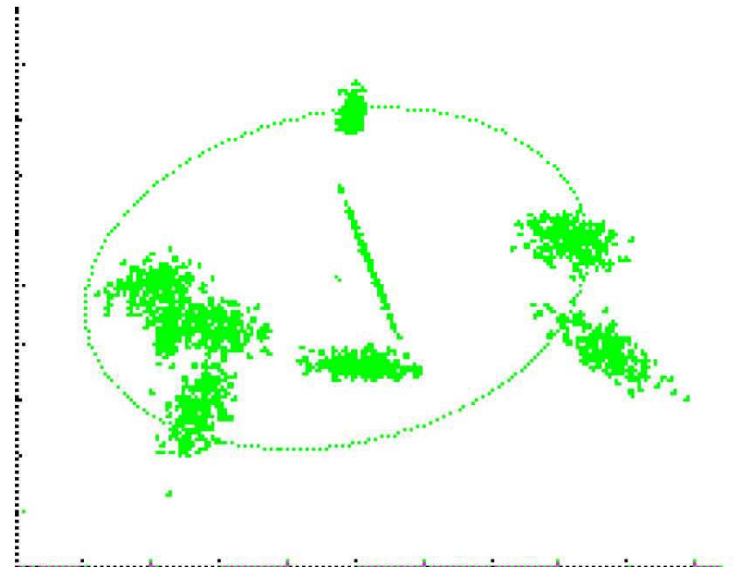


Centroid

# Example: microarray expression

- **Measure gene expression**

- **Various experimental conditions**
  - Cancer, normal
  - Time
  - Subjects

- **Explore similarities**
  - What genes change together?
  - What conditions are similar?

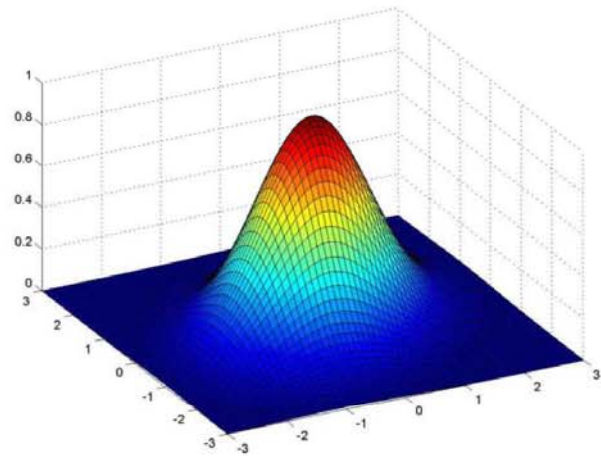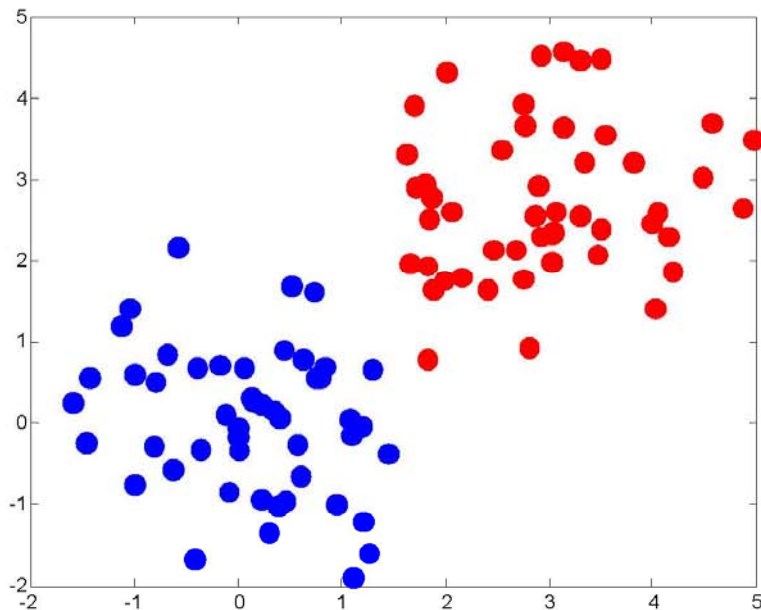- **Cluster on both genes and conditions**

# Mixtures of Gaussians

- ## K-means algorithm
  - Assigned each example to exactly one cluster
  - What if clusters are overlapping?
    - Hard to tell which cluster is right
    - Maybe we should try to remain uncertain

- ## Gaussian mixture models
  - Clusters modeled as Gaussians
    - Not just by their mean
  - EM algorithm: assign data to cluster with some *probability*

# Multivariate Gaussian models

$$\mathcal{N}(\underline{x} \; ; \; \underline{\mu}, \Sigma) = \frac{1}{(2\pi)^{d/2}} |\Sigma|^{-1/2} \exp\left\{ -\frac{1}{2}(\underline{x} - \underline{\mu})^T \Sigma^{-1} (\underline{x} - \underline{\mu}) \right\}$$





We'll model each cluster using one of these Gaussian "bells"…

# EM Algorithm: E-step

- Start with parameters describing each cluster
- Mean $\mu_c$, Covariance $\Sigma_c$, "size" $\pi_c$

- E-step ("Expectation")
  - For each datum (example) $x_i$,
  - Compute "$r_{ic}$", the probability that it belongs to cluster c
    - Compute its probability under model c
    - Normalize to sum to one (over clusters c)

$$r_{ic} = \frac{\pi_c \, \mathcal{N}(x_i \; ; \; \mu_c, \Sigma_c)}{\sum_{c'} \pi_{c'} \, \mathcal{N}(x_i \; ; \; \mu_{c'}, \Sigma_{c'})}$$

  - If $x_i$ is very likely under the $c^{th}$ Gaussian, it gets high weight
  - Denominator just makes r's sum to one

# EM Algorithm: M-step

- Start with assignment probabilities $r_{ic}$
- Update parameters: mean $\mu_c$, Covariance $\Sigma_c$, "size" $\pi_c$

- M-step ("Maximization")
  - For each cluster (Gaussian) x_c,
  - Update its parameters using the (weighted) data points

$$N_c = \sum_i r_{ic}$$    **Total responsibility allocated to cluster c**

$$\pi_c = \frac{N_c}{N}$$    **Fraction of total assigned to cluster c**

$$\mu_c = \frac{1}{N_c} \sum_i r_{ic} x_i$$      $$\Sigma_c = \frac{1}{N_c} \sum_i r_{ic}(x_i - \mu_c)^T (x_i - \mu_c)$$

**Weighted mean of assigned data**      **Weighted covariance of assigned data (use new weighted means here)**
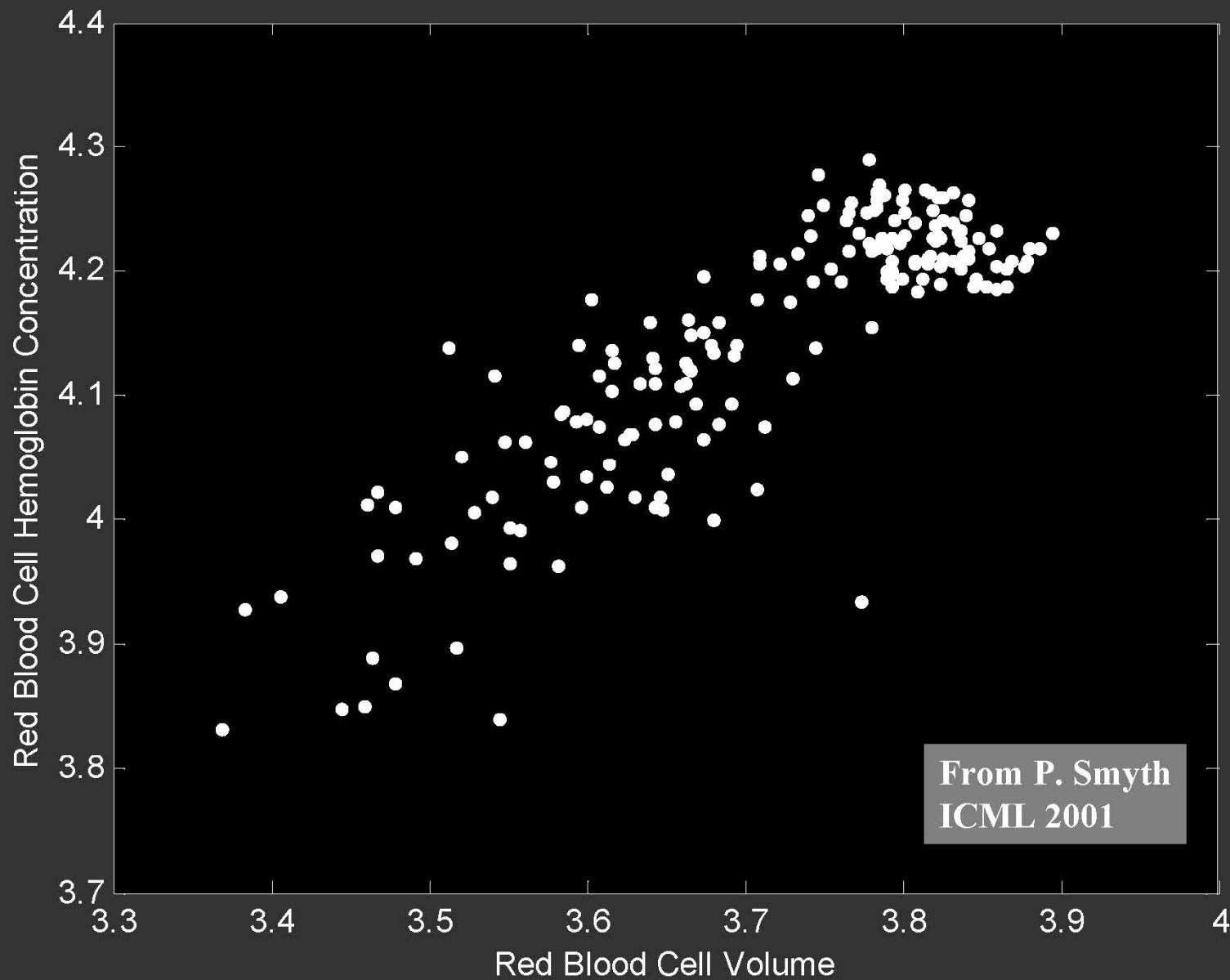
# Expectation-Maximization

- Each step increases the log-likelihood of our model

$$\log p(\underline{X}) = \sum_i \log \left[ \sum_c \pi_c \, \mathcal{N}(x_i \; ; \; \mu_c, \Sigma_c) \right]$$

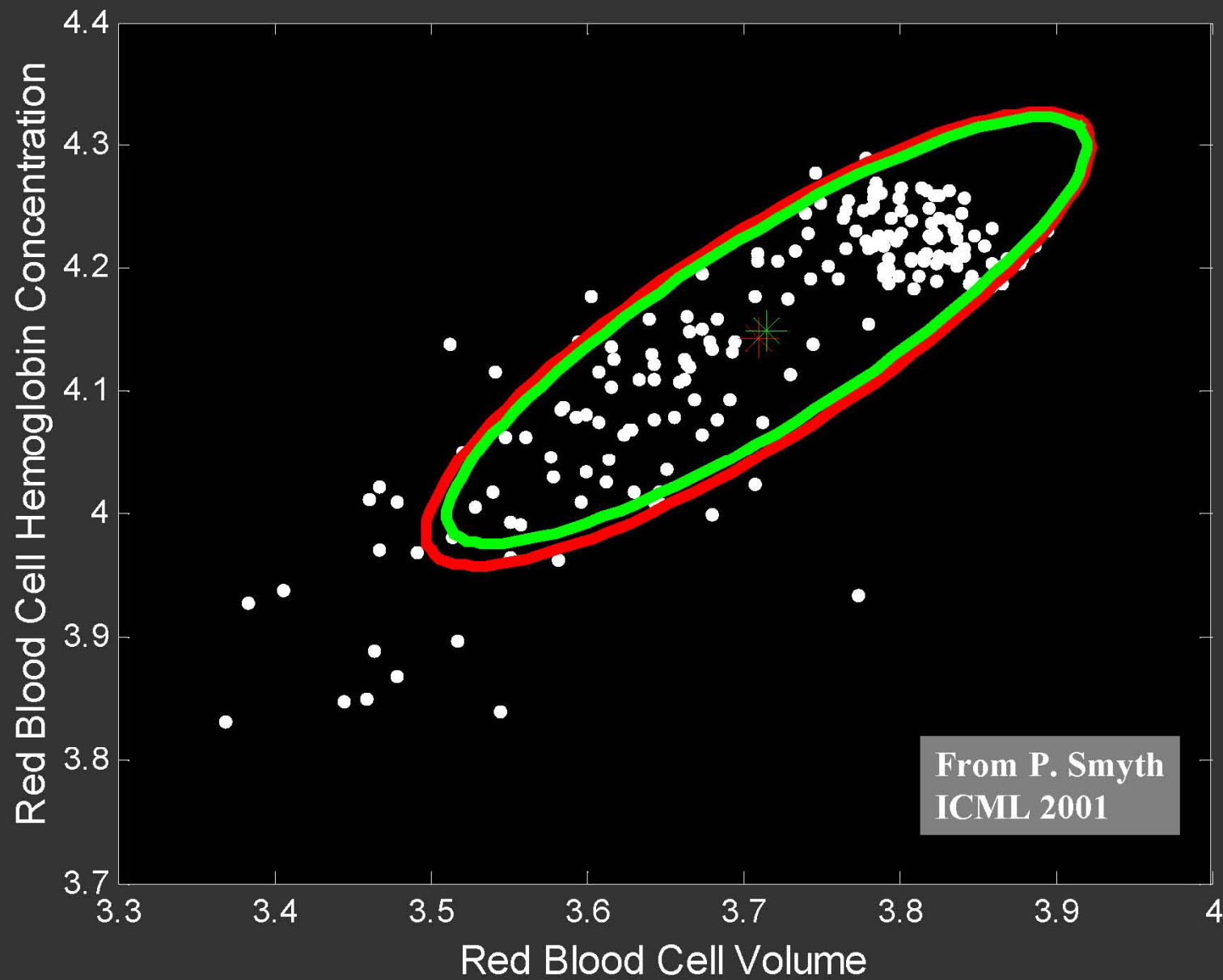(we won't derive this, though)

- Iterate until convergence
  - Convergence guaranteed – another ascent method

- What should we do
  - If we want to choose a single cluster for an "answer"?
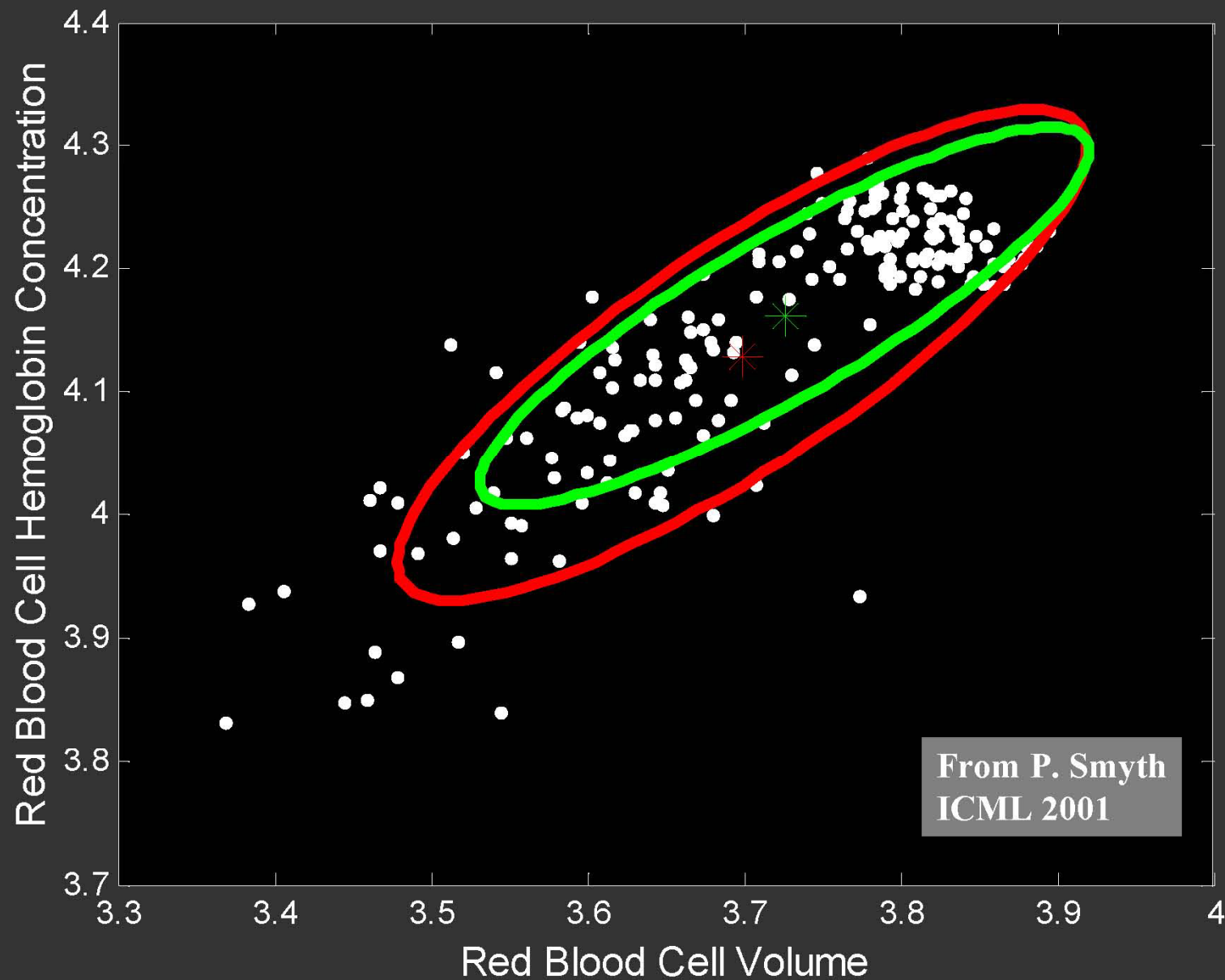  - With new data we didn't see during training?

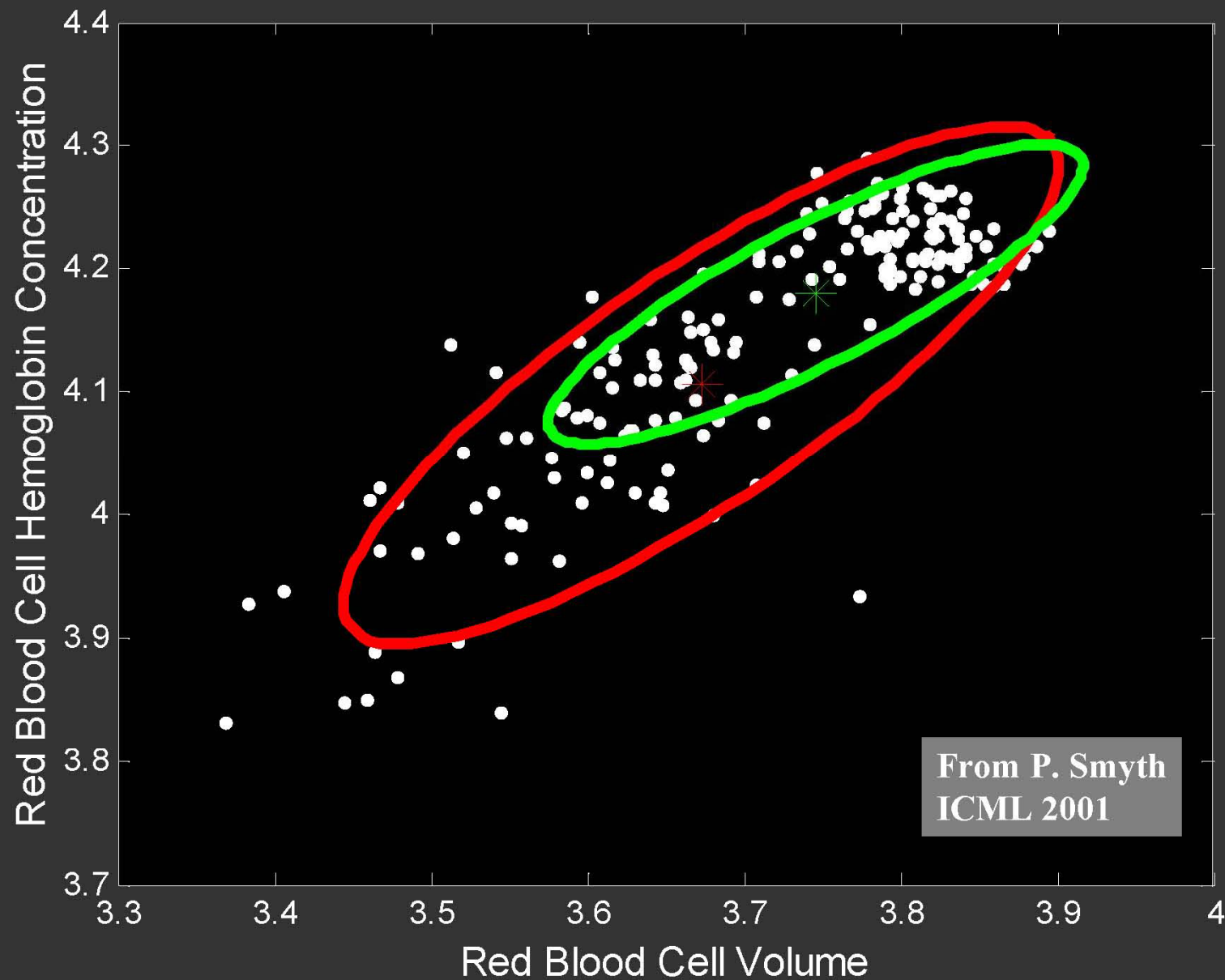ANEMIA PATIENTS AND CONTROLS

From P. Smyth
ICML 2001
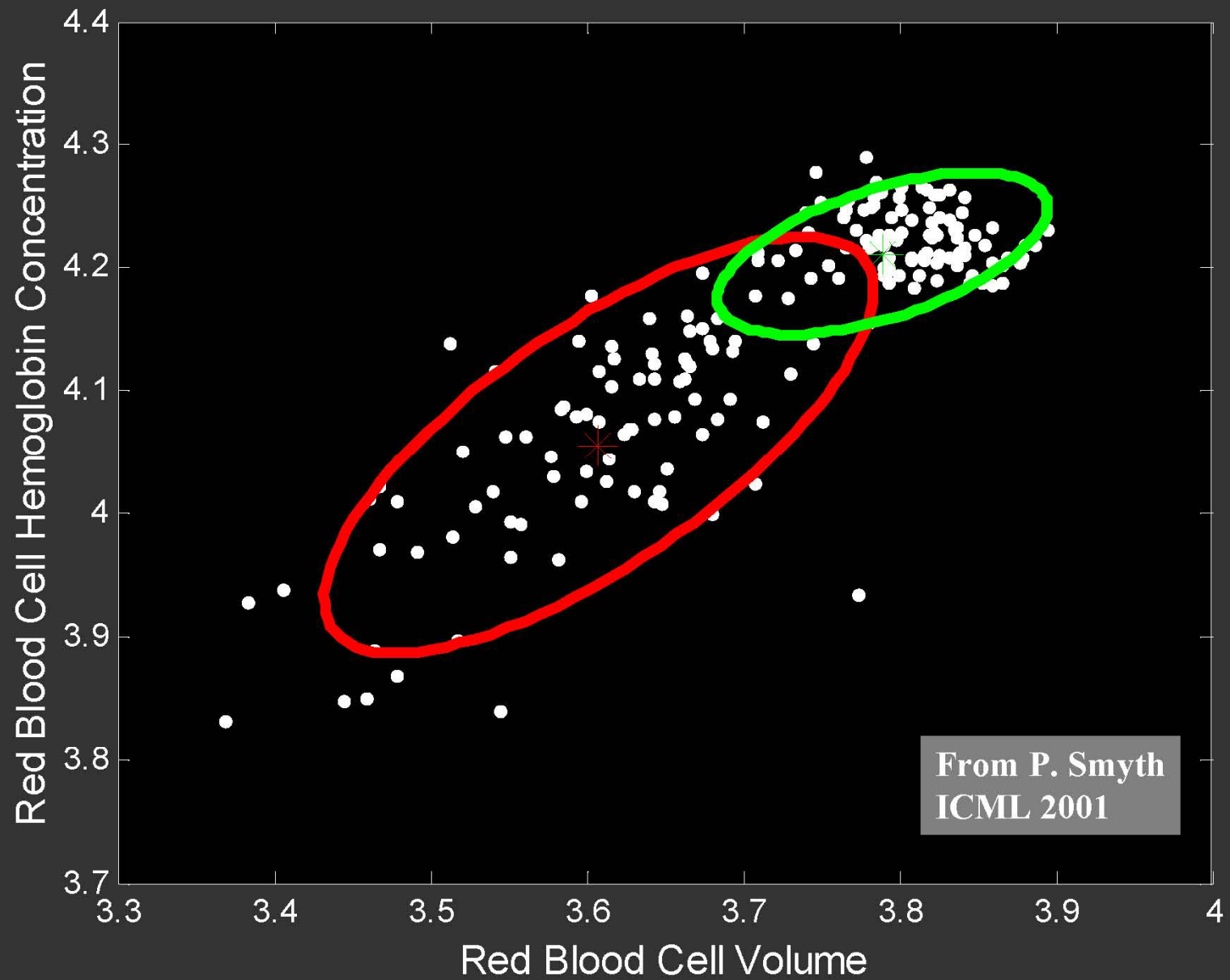
EM ITERATION 3

From P. Smyth
ICML 2001

EM ITERATION 5

Red Blood Cell Hemoglobin Concentration

Red Blood Cell Volume

From P. Smyth
ICML 2001

EM ITERATION 10

From P. Smyth
ICML 2001

EM ITERATION 15

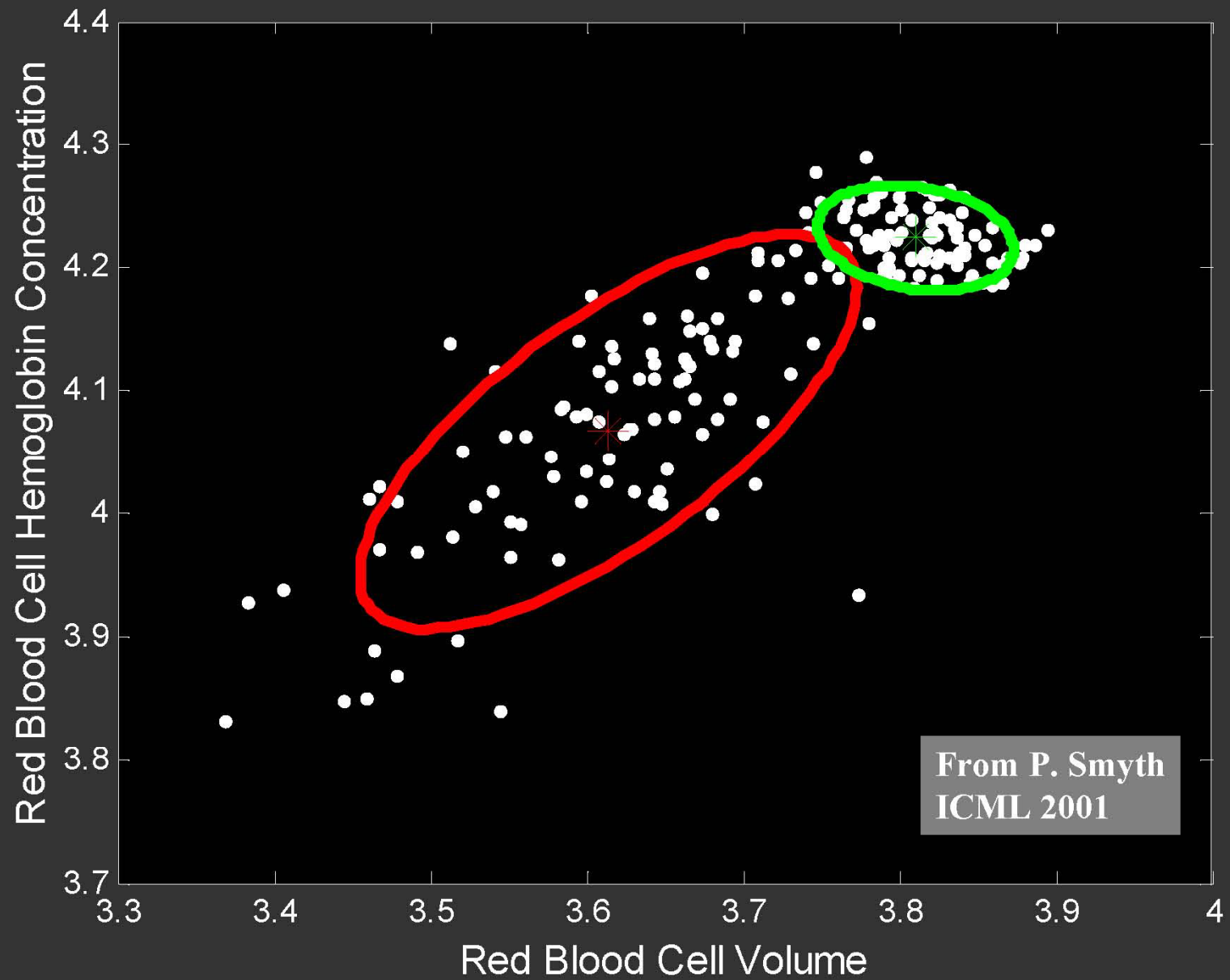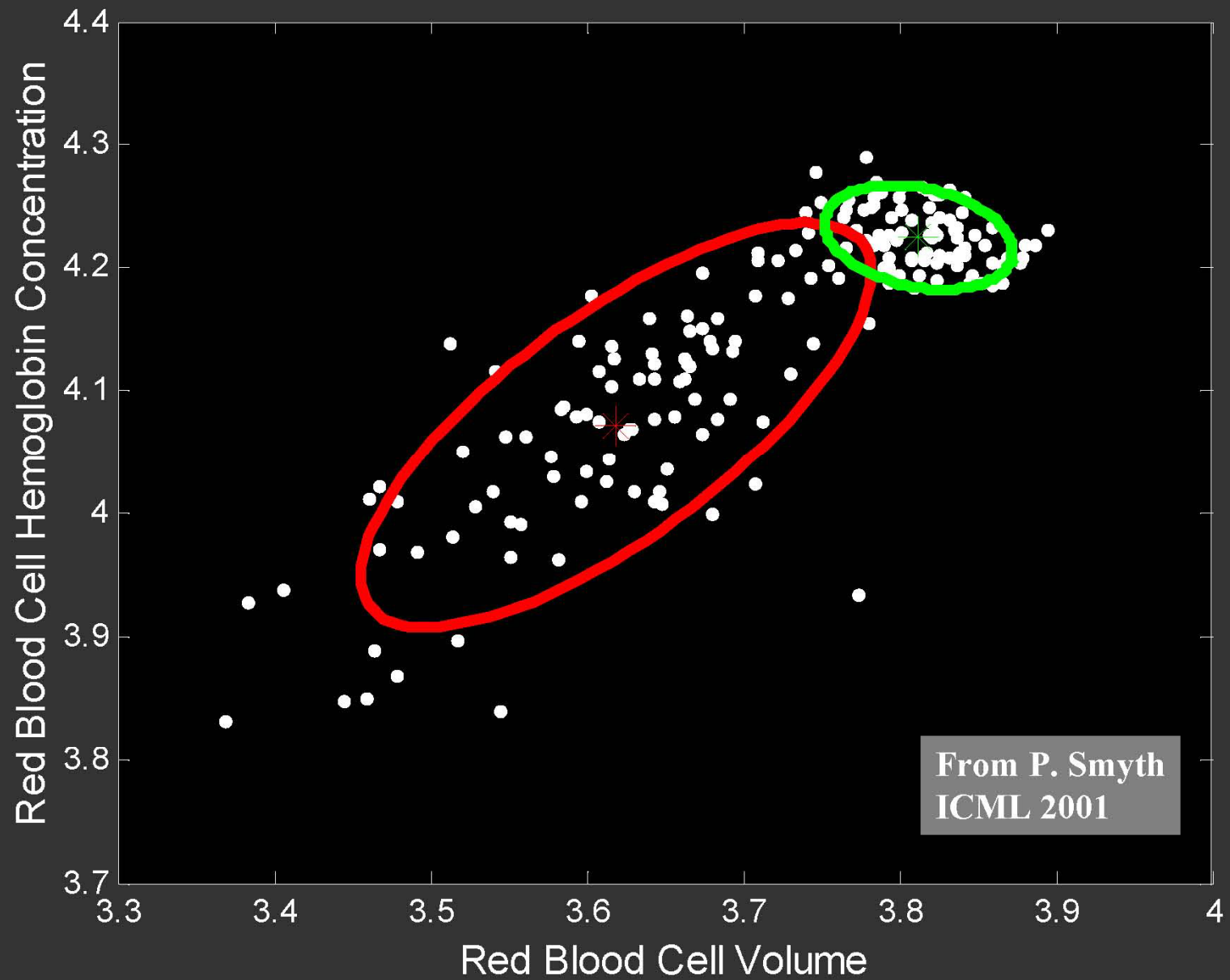Red Blood Cell Hemoglobin Concentration

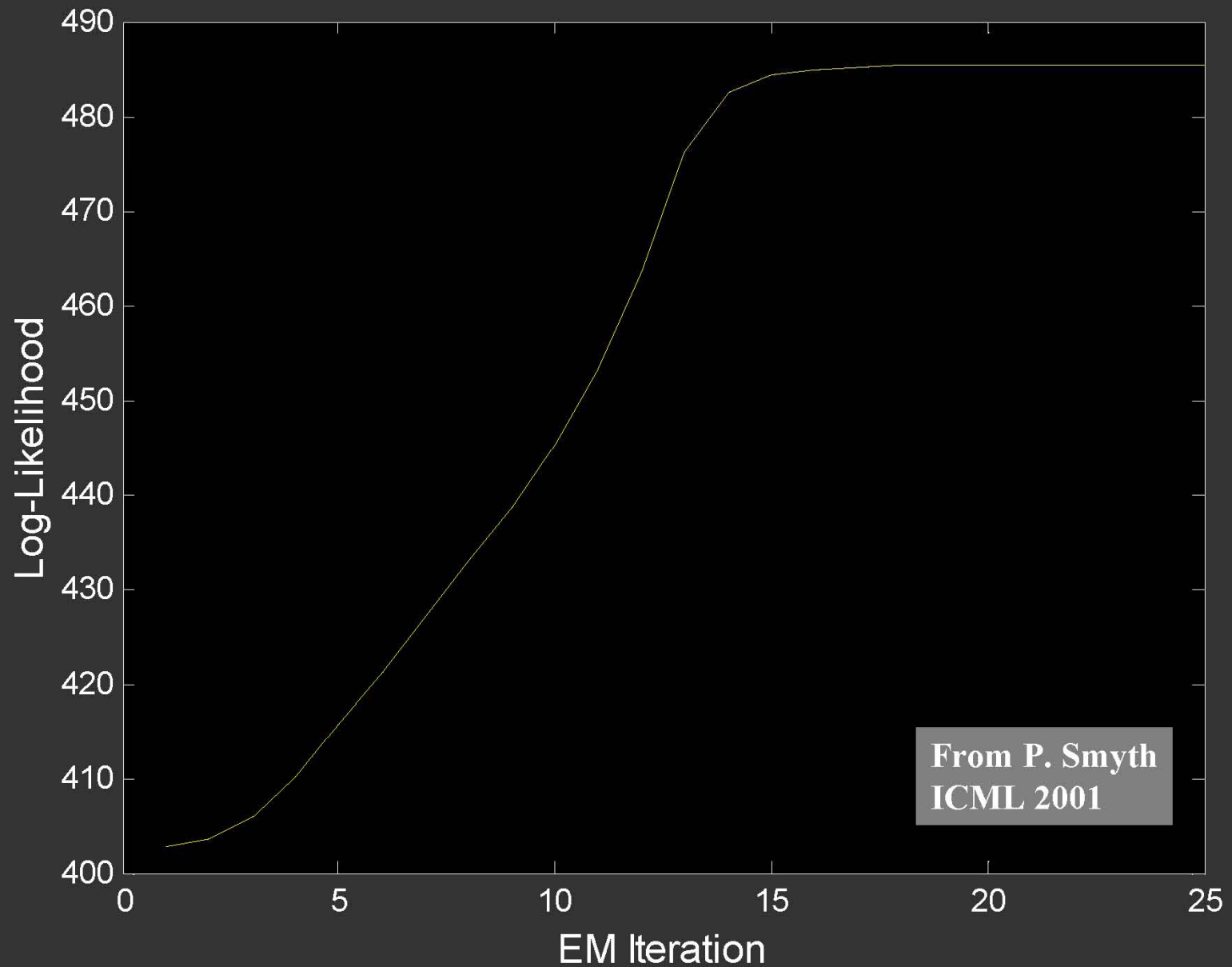Red Blood Cell Volume

From P. Smyth
ICML 2001

EM ITERATION 25

From P. Smyth
ICML 2001

LOG-LIKELIHOOD AS A FUNCTION OF EM ITERATIONS

From P. Smyth
ICML 2001

# EM and missing data

- EM is a general framework for partially observed data
  - "Complete data" $x_i$, $z_i$ – features and assignments
  - Assignments $z_i$ are missing (unobserved)

- EM corresponds to
  - Computing the distribution over all $z_i$ given the parameters
  - Maximizing the "expected complete" log likelihood
  - GMMs = plug in "soft assignments", but not always so easy

- Alternative: Stochastic EM
  - Instead of expectations, just sample the $z_i$ (often easier)
  - Called "imputing" the values of $z_i$
  - Behaves similarly, but with extra randomness
    - Not obvious when it has converged

# Gibbs sampling for clustering

- Another technique for inferring uncertain cluster assignments
  - K-means: take the best assignment
  - EM: assign "partially"
  - Stochastic EM: sample assignment
  - All: choose best cluster descriptions given assignments

- Gibbs sampling ("Markov chain Monte Carlo")
  - Assign randomly, probability equal to EM's weight
  - *Sample* a cluster description given assignment
  - Requires a probability model over cluster parameters

- This doesn't really find the "best" clustering
  - It eventually samples almost all "good" clusterings
  - Converges "in probability", randomness helps us explore configurations
  - Also tells us about uncertainty of clustering
  - Disadvantage: not obvious when "done"