

Expectation Maximization

Brandon Caie and Jonny Coutinho

Intro: Expectation Maximization Algorithm



- EM algorithm provides a general approach to learning in presence of unobserved variables.
- In many practical learning settings, only a subset of relevant features or variables might be observable.
 - Eg: Hidden Markov, Bayesian Belief Networks

Simple Example: Coin Flipping



- Suppose you have 2 coins, A and B, each with a certain bias of landing heads, θ_A, θ_B .
- Given data sets $X_A = \{x_{1,A}, \dots, x_{m_A,A}\}$ and $X_B = \{x_{1,B}, \dots, x_{m_B,B}\}$
Where $x_{i,j} = \begin{cases} 1 & ; \text{if heads} \\ 0 & ; \text{otherwise} \end{cases}$
- No hidden variables – easy solution. $\theta_j = \frac{1}{m_j} \sum_{i=1}^{m_j} x_{i,j}$; sample mean

Simplified MLE

5 sets of 10 tosses each

H T T T H H T H T H
H H H H T H H H H H
H T H H H H H T H H
H T H T T T H H T T
T H H H T H H H T H

5 sets, 10 tosses per set

Coin A	Coin B
	5 H, 5 T
9 H, 1 T	
8 H, 2 T	
	4 H, 6 T
7 H, 3 T	
24 H, 6 T	9 H, 11 T

Goal: determine coin parameters without knowing the identity of each data set's coin.

Solution: Expectation-maximization

$$\hat{\theta}_1 = \frac{24}{24 + 6} = 0.80$$
$$\hat{\theta}_2 = \frac{9}{9 + 11} = 0.45$$

Coin Flip With hidden variables

- What if you were given the same dataset of coin flip results, but no coin identities defining the datasets?

Here: $X = \{x_1, \dots, x_m\}$; the observed variable

$$Z = \begin{Bmatrix} z_{1,1} & \dots & z_{m,1} \\ \dots & z_{i,j} & \dots \\ z_{1,k} & \dots & z_{m,k} \end{Bmatrix} \quad \text{where } z_{i,j} = \begin{cases} 1 & ; \text{if } x_i \text{ is from } j^{\text{th}} \text{ coin} \\ 0 & ; \text{otherwise} \end{cases}$$

But Z is not known. (Ie: 'hidden' / 'latent' variable)

0) Initialize some arbitrary hypothesis of parameter values (θ):

$$\theta = \{ \theta_1, \dots, \theta_k \} \quad \text{coin flip example: } \theta = \{ \theta_A, \theta_B \} = \{0.6, 0.5\}$$

1) Expectation (E-step)

$$E[z_{i,j}] = \frac{p(x = x_i | \theta = \theta_j)}{\sum_{n=1}^k p(x = x_i | \theta = \theta_n)}$$

2) Maximization (M-step)

$$\theta_j = \frac{\sum_{i=1}^m E[z_{i,j}] x_i}{\sum_{i=1}^m E[z_{i,j}]}$$

If $z_{i,j}$ is known:

$$\theta_j = \frac{\sum_{i=1}^{m_j} x_i}{m_j}$$

EM- Coin Flip example

H T T T H H T H T H
H H H H T H H H H H
H T H H H H H T H H
H T H T T T H H T T
T H H H T H H H T H

5 sets, 10 tosses per set

- Initialize θ_A and θ_B to chosen value
 - Ex: $\theta_A=0.6$, $\theta_B=0.5$
- Compute a probability distribution of possible completions of the data using current parameters

EM- Coin Flip example

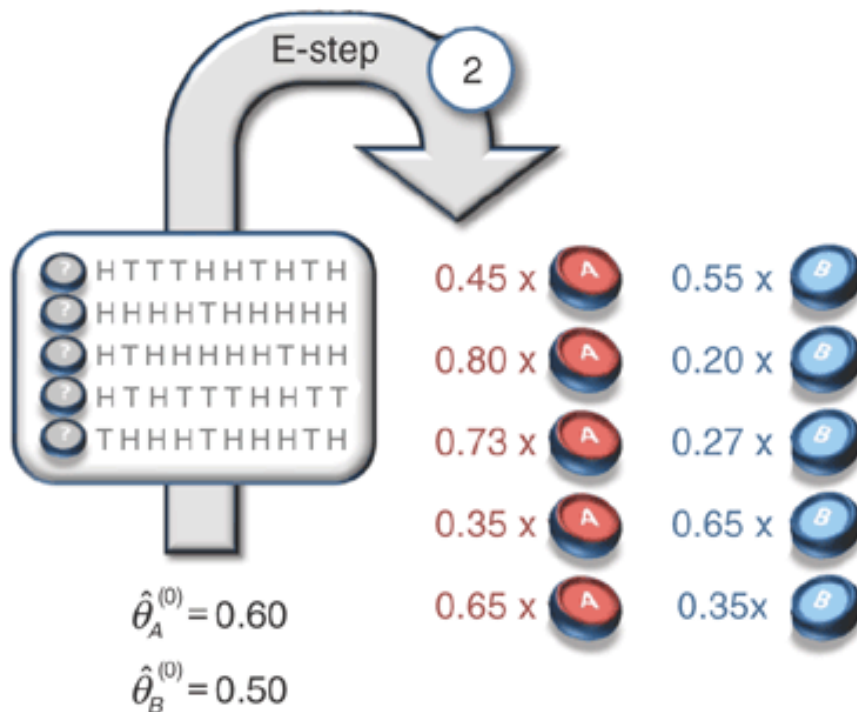
Set 1 H T T T H H T H T H

- What is the probability that I observe 5 heads and 5 tails in coin A and B *given* the initializing parameters $\theta_A=0.6$, $\theta_B=0.5$?
- Compute likelihood of set 1 coming from coin A or B using the binomial distribution with mean probability θ on n trials with k successes

$$p(k) = \binom{n}{k} \theta^k (1 - \theta)^{n-k}$$

- Likelihood of "A" = 0.00079
- Likelihood of "B" = 0.00097
- Normalize to get probabilities $\rightarrow A=0.45$, $B=0.55$

The E-step



Coin A	Coin B
$\approx 2.2 \text{ H}, 2.2 \text{ T}$	$\approx 2.8 \text{ H}, 2.8 \text{ T}$
$\approx 7.2 \text{ H}, 0.8 \text{ T}$	$\approx 1.8 \text{ H}, 0.2 \text{ T}$
$\approx 5.9 \text{ H}, 1.5 \text{ T}$	$\approx 2.1 \text{ H}, 0.5 \text{ T}$
$\approx 1.4 \text{ H}, 2.1 \text{ T}$	$\approx 2.6 \text{ H}, 3.9 \text{ T}$
$\approx 4.5 \text{ H}, 1.9 \text{ T}$	$\approx 2.5 \text{ H}, 1.1 \text{ T}$
$\approx 21.3 \text{ H}, 8.6 \text{ T}$	$\approx 11.7 \text{ H}, 8.4 \text{ T}$

The M-step

Coin A	Coin B
$\approx 2.2 \text{ H}, 2.2 \text{ T}$	$\approx 2.8 \text{ H}, 2.8 \text{ T}$
$\approx 7.2 \text{ H}, 0.8 \text{ T}$	$\approx 1.8 \text{ H}, 0.2 \text{ T}$
$\approx 5.9 \text{ H}, 1.5 \text{ T}$	$\approx 2.1 \text{ H}, 0.5 \text{ T}$
$\approx 1.4 \text{ H}, 2.1 \text{ T}$	$\approx 2.6 \text{ H}, 3.9 \text{ T}$
$\approx 4.5 \text{ H}, 1.9 \text{ T}$	$\approx 2.5 \text{ H}, 1.1 \text{ T}$
$\approx 21.3 \text{ H}, 8.6 \text{ T}$	$\approx 11.7 \text{ H}, 8.4 \text{ T}$

$$\hat{\theta}_A^{(1)} \approx \frac{21.3}{21.3 + 8.6} \approx 0.71$$

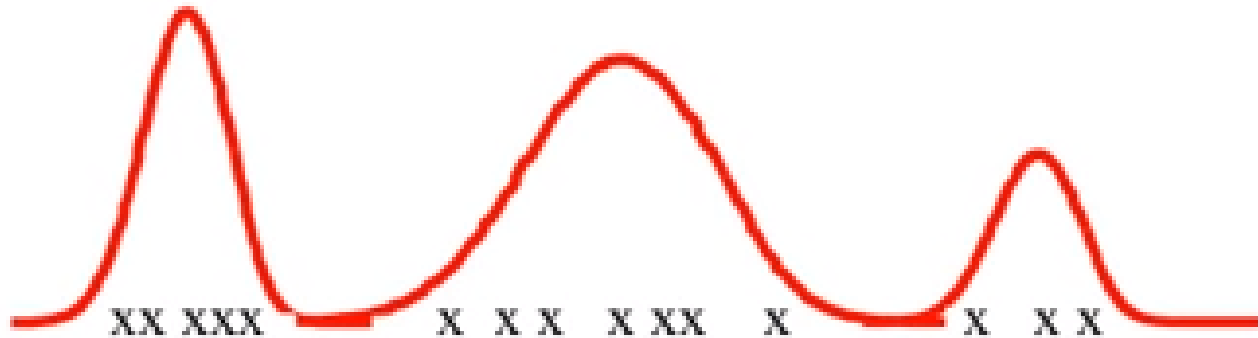
$$\hat{\theta}_B^{(1)} \approx \frac{11.7}{11.7 + 8.4} \approx 0.58$$

Summary



1. Choose starting parameters
2. Estimate probability using these parameters that each data set (x_i) came from j^{th} coin ($E[z_{i,j}]$)
3. Use these probability values ($E[z_{i,j}]$) as weights on each data point when computing a new θ_j to describe each distribution
4. Summate these expected values, use maximum likelihood estimation to derive new parameter values to repeat process

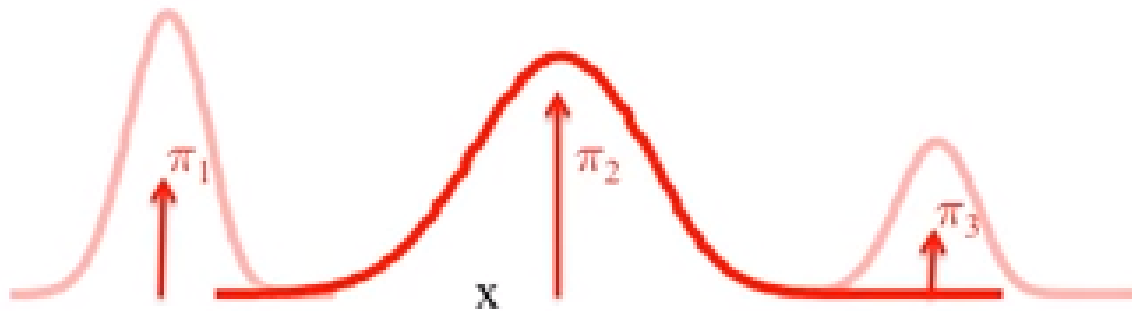
- When data is continuous, can be described by Normal Distributions



- Cluster data as Gaussians, with parameters: $(\mu_j, \sigma_j^2, \pi_j)$

$$p(z = j) = \pi_j$$

$$p(x|z = j) = N(x; \mu_j, \sigma_j^2)$$



EM algorithm in Gaussian Mixtures



Step 0) Initialize $\theta = \begin{Bmatrix} \mu_1, \dots, \mu_k \\ \sigma_1^2, \dots, \sigma_k^2 \\ \pi_1, \dots, \pi_k \end{Bmatrix}$ (assuming k clusters)

Step 1) Expectation: compute $r_{i,j}$ for each x_i

$$r_{i,j} = \frac{\pi_{i,j} p(x|z = j)}{\sum_{n=1}^k \pi_{i,n} p(x|z = n)}$$

EM algorithm for Gaussian Mixture

Step 2) Maximization:

$$m_j = \sum_i r_{i,j}$$

$$\pi_j = \frac{m_j}{m}$$

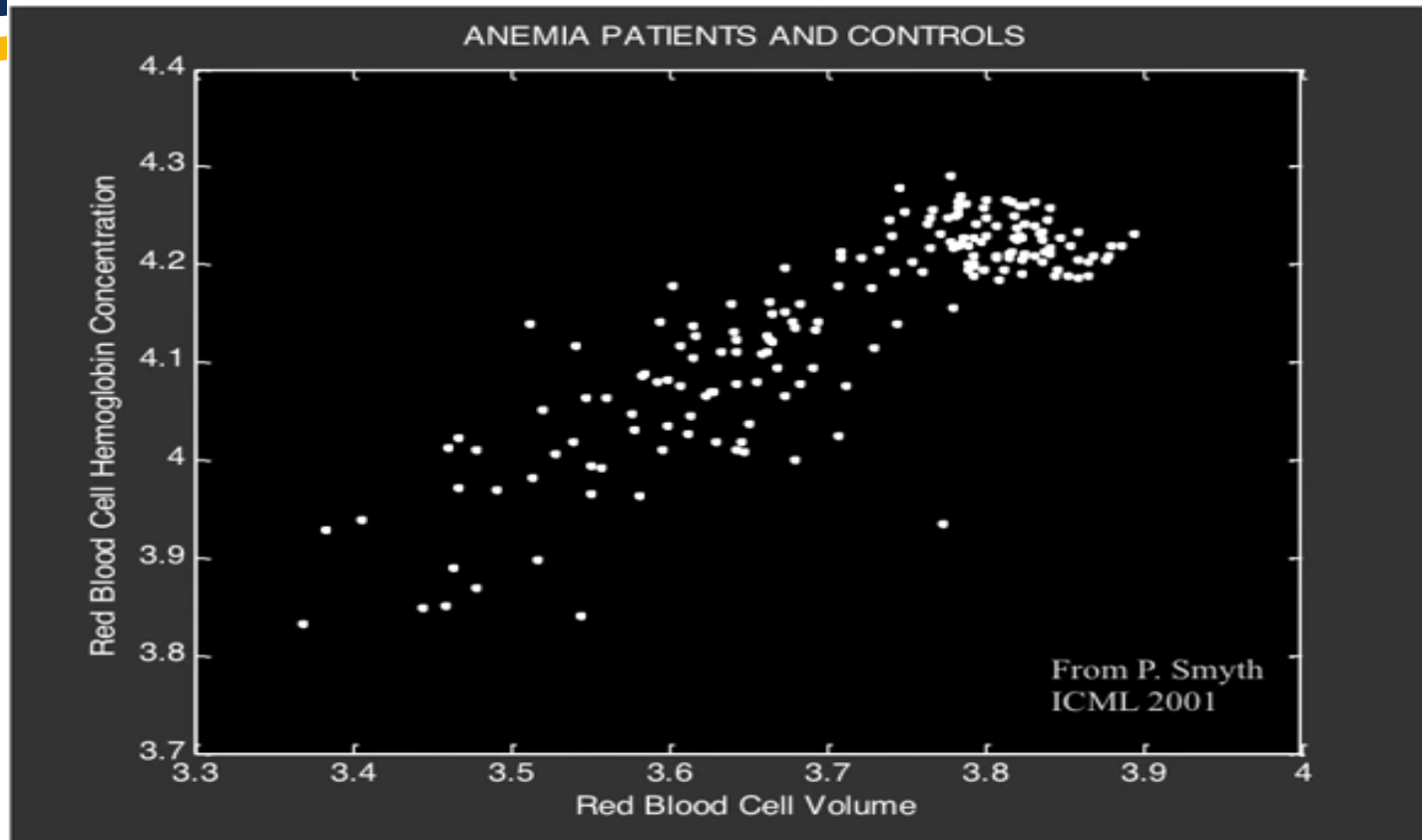
$$\mu_j = \frac{1}{m_j} \sum_i r_{i,j} x_i$$

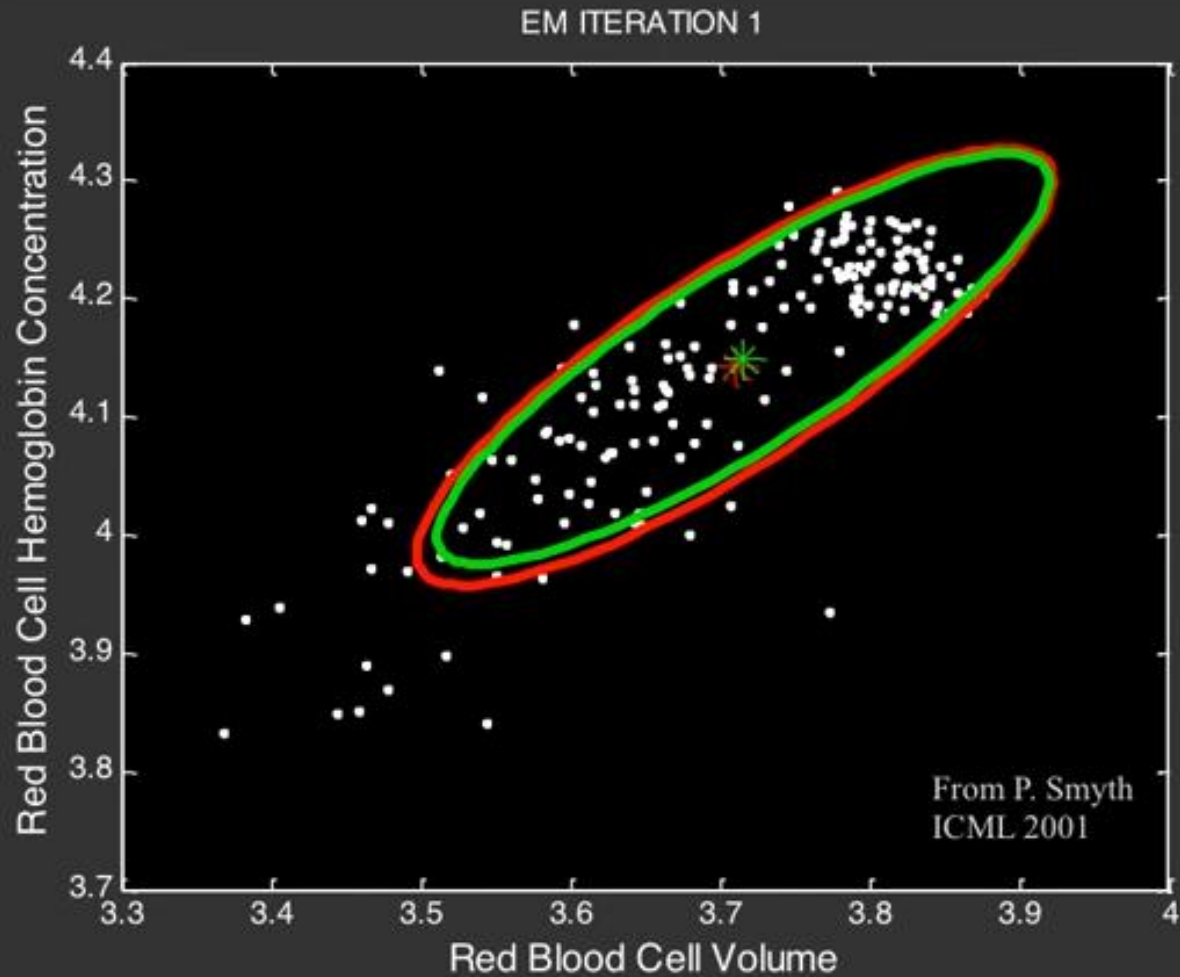
$$\sigma_j^2 = \frac{1}{m_j} \sum_i r_{i,j} (x_i - \mu_j)^2$$

Example of EM in Gaussian Mixtures

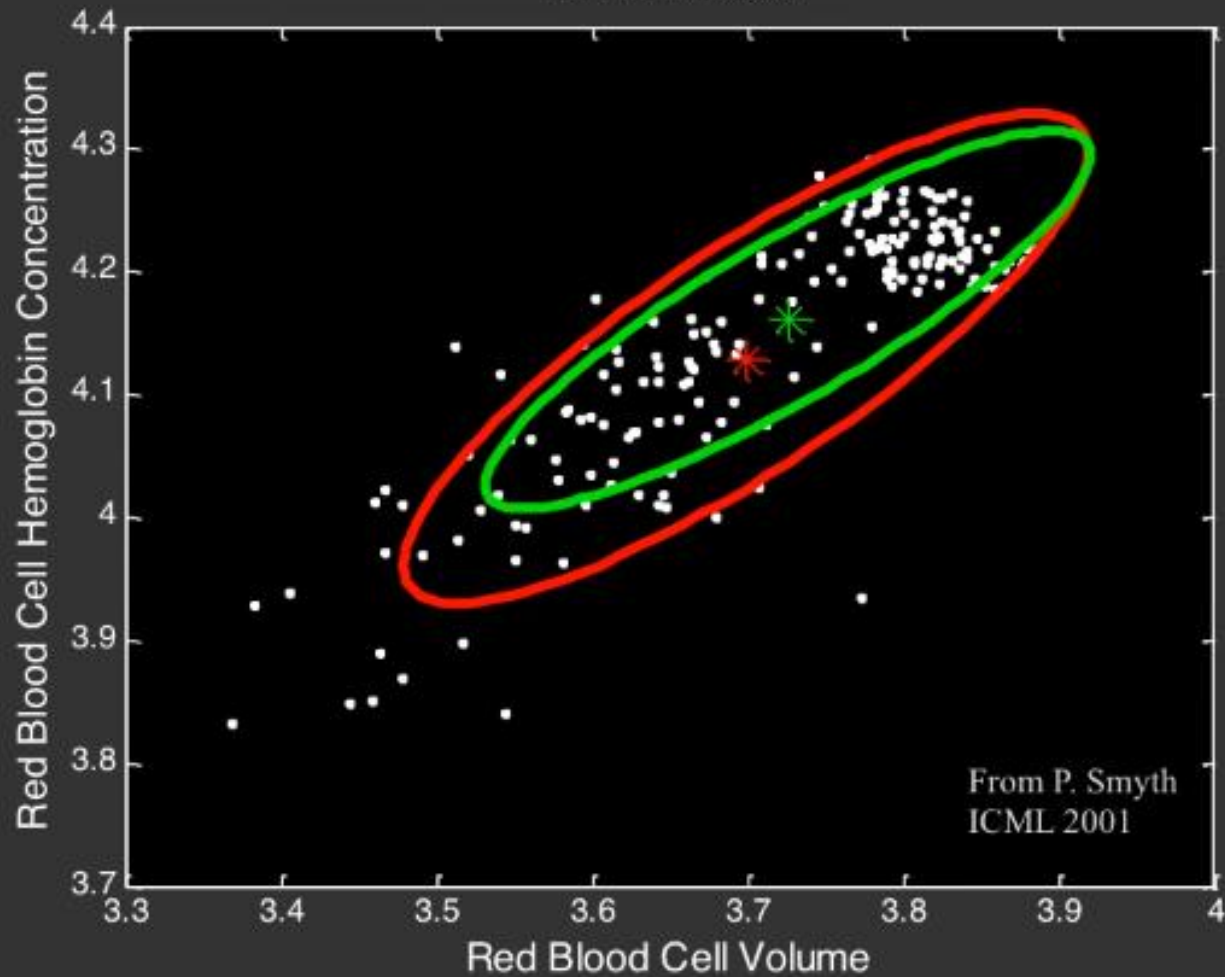


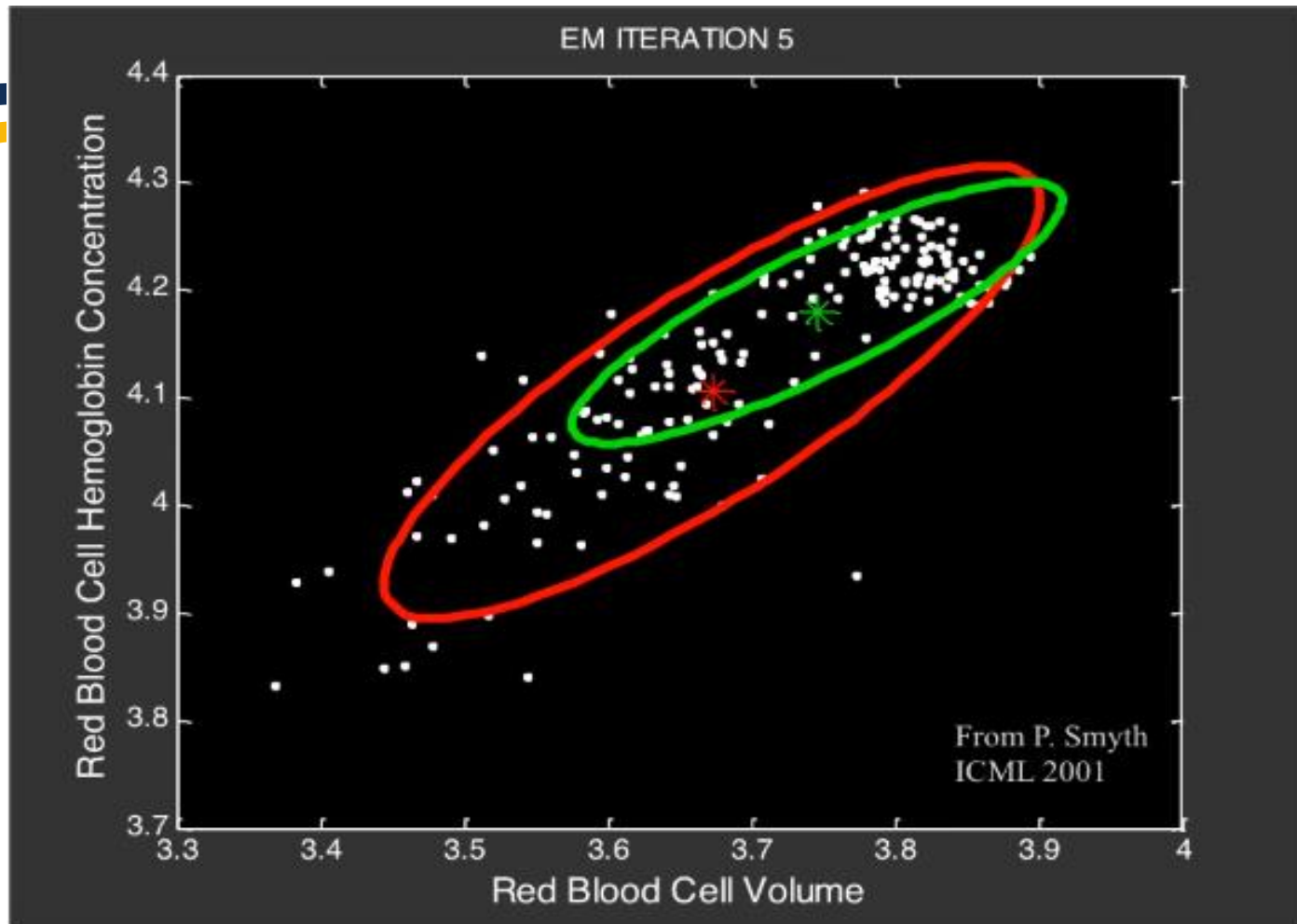
UNIVERSITY OF SYDNEY

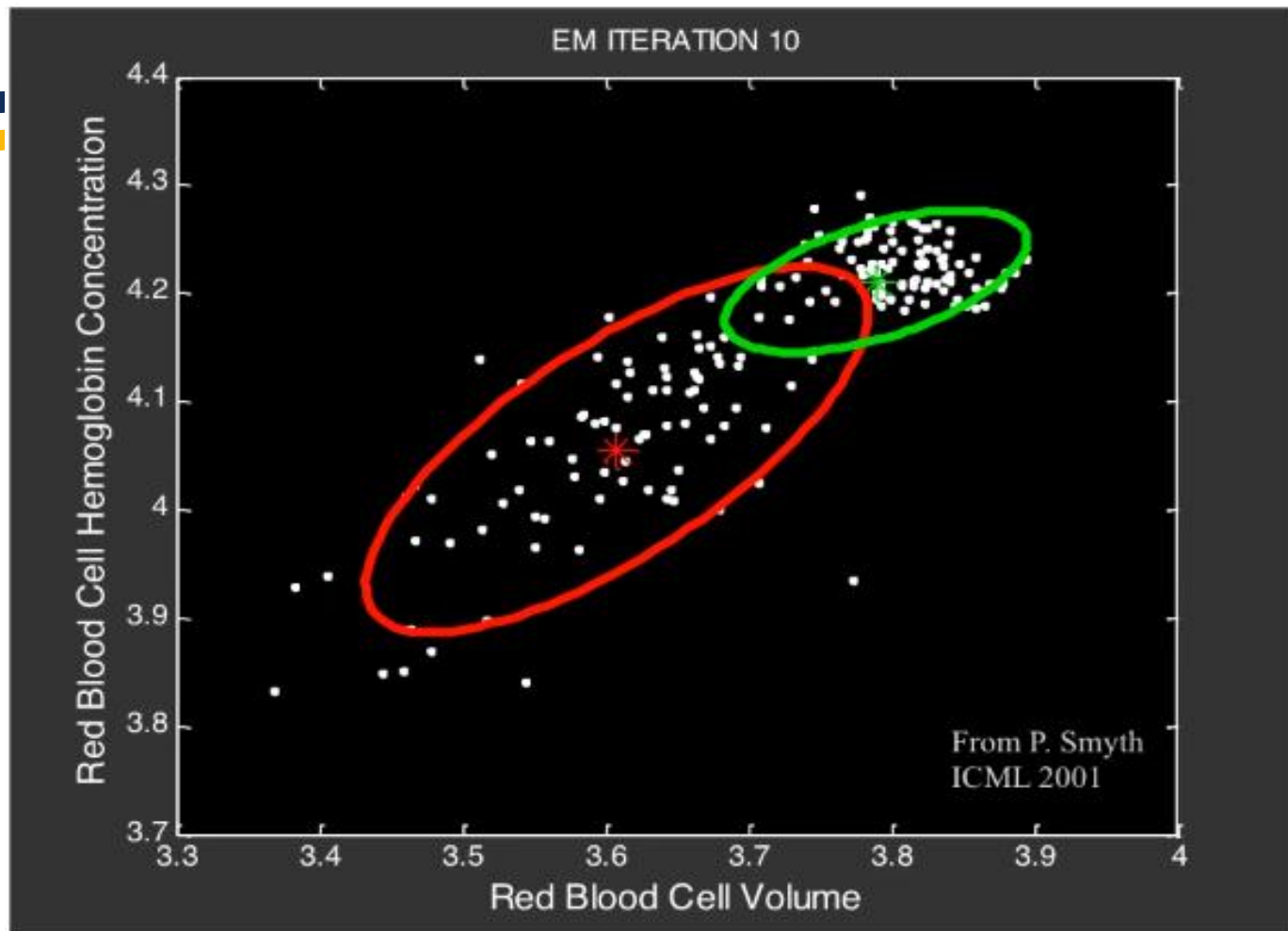


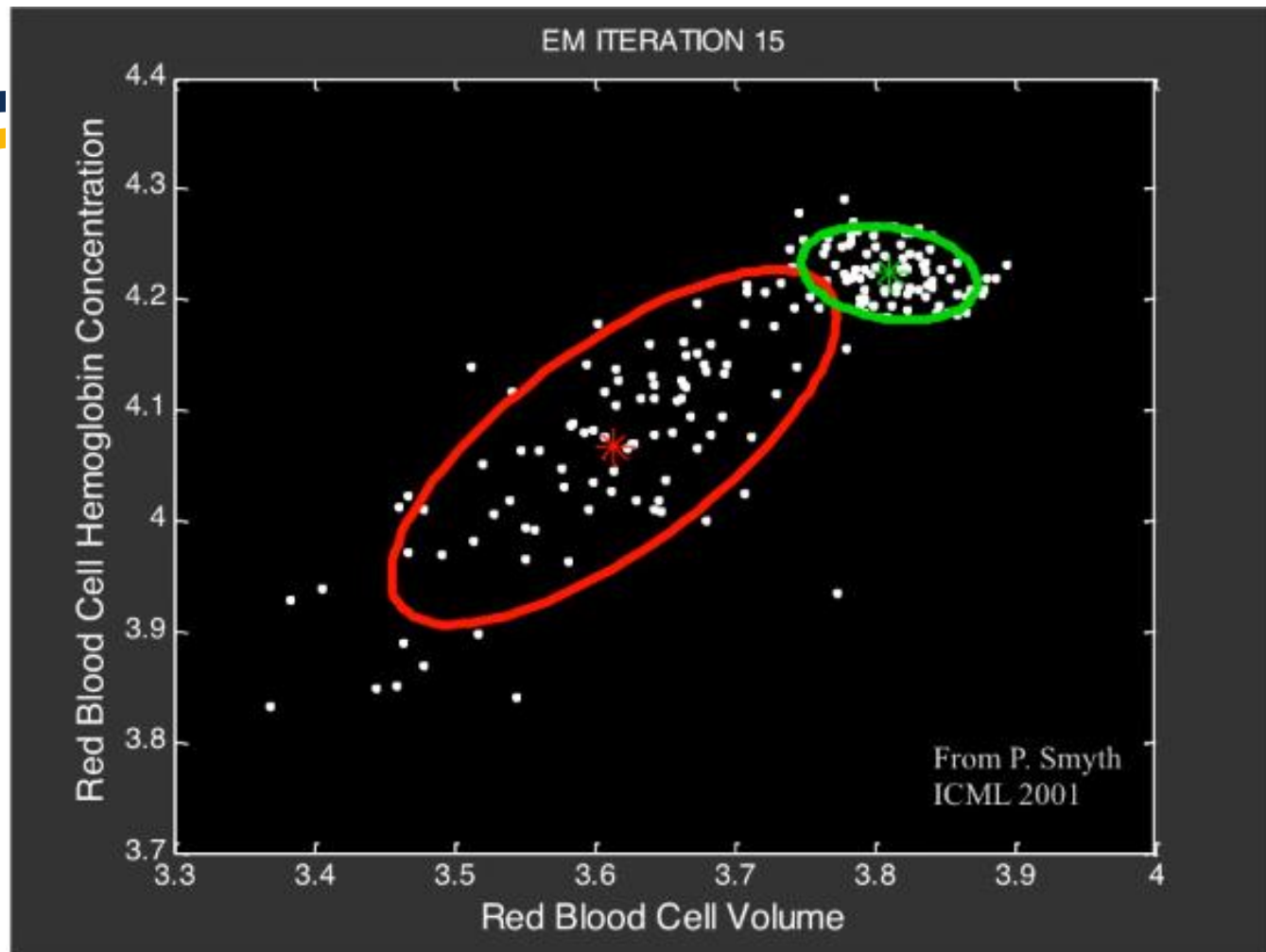


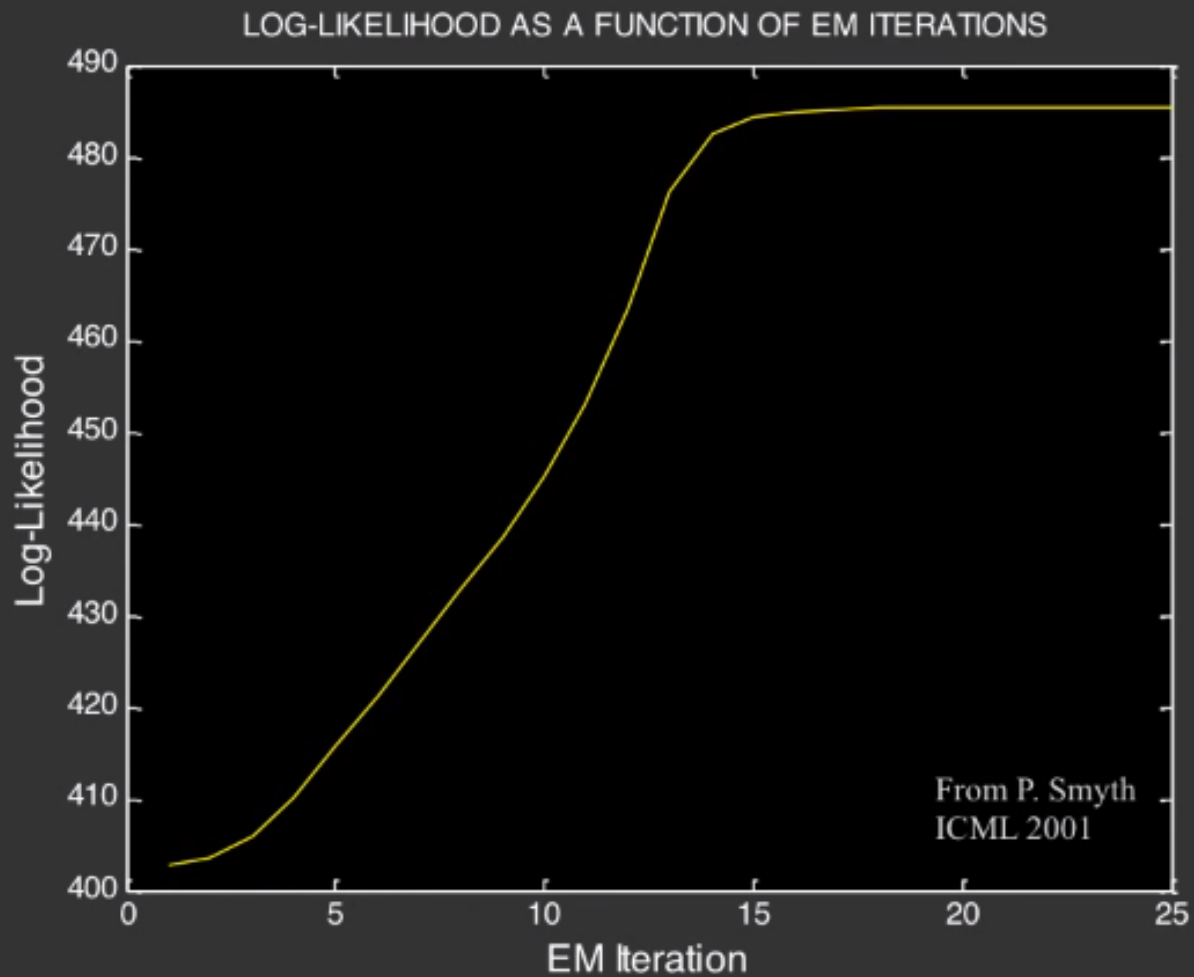
EM ITERATION 3



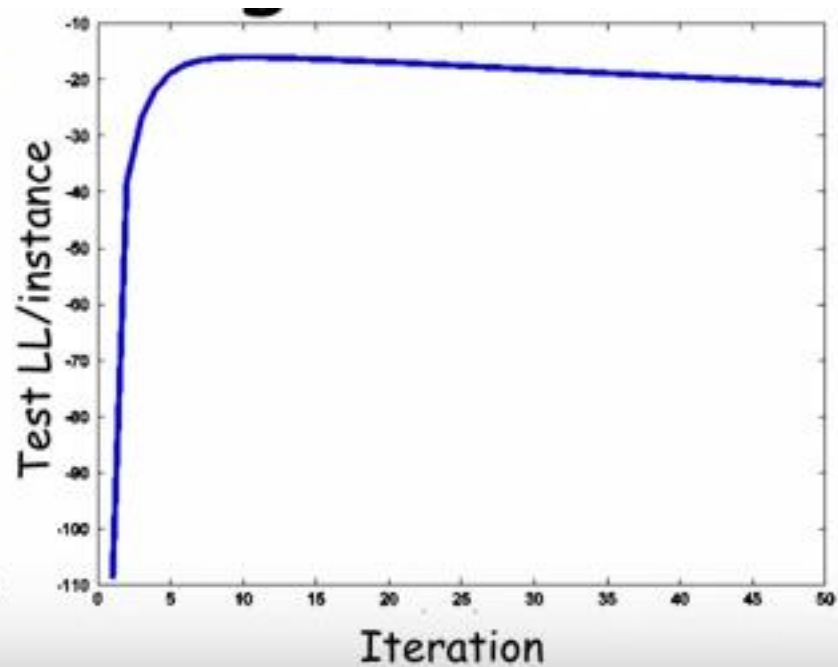
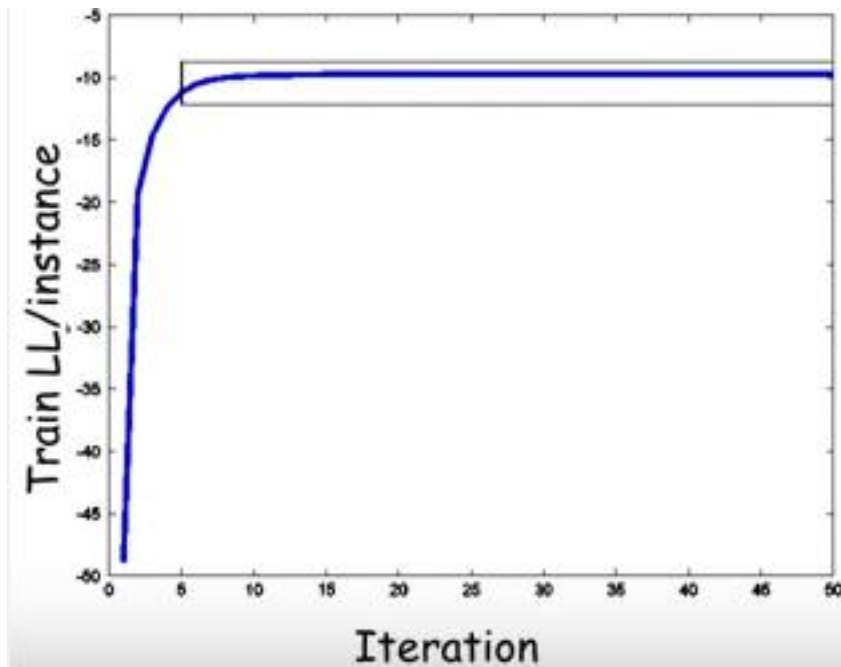








Overfitting through convergence



- Hidden variables and incomplete data lead to more complex likelihood functions w/ many local optima
- Since EM only solves for a single local optima, choosing a good initial parameter estimation is critical
- Strategies to improve initialization
 - Multiple random restarts
 - Use prior knowledge
 - Output of a simpler, though less robust algorithm

- [Matlab EM Algorithm](#)
- Tom Mitchell- Machine Learning: Chapter 6 (on lab wiki)
- [EM Algorithm Derivation, Convergence, Hidden Markov and GMM Applications](#)
- [Nature Review Article](#)