

Machine Learning for Sensory Signals

Gaussian and Mixture Gaussian Models

02-03-2017

Gaussian Mixture Models

A Gaussian Mixture Model (GMM) is defined as

$$p(\mathbf{x}|\Theta) = \sum_{k=1}^K \alpha_k p(\mathbf{x}|\theta_k)$$

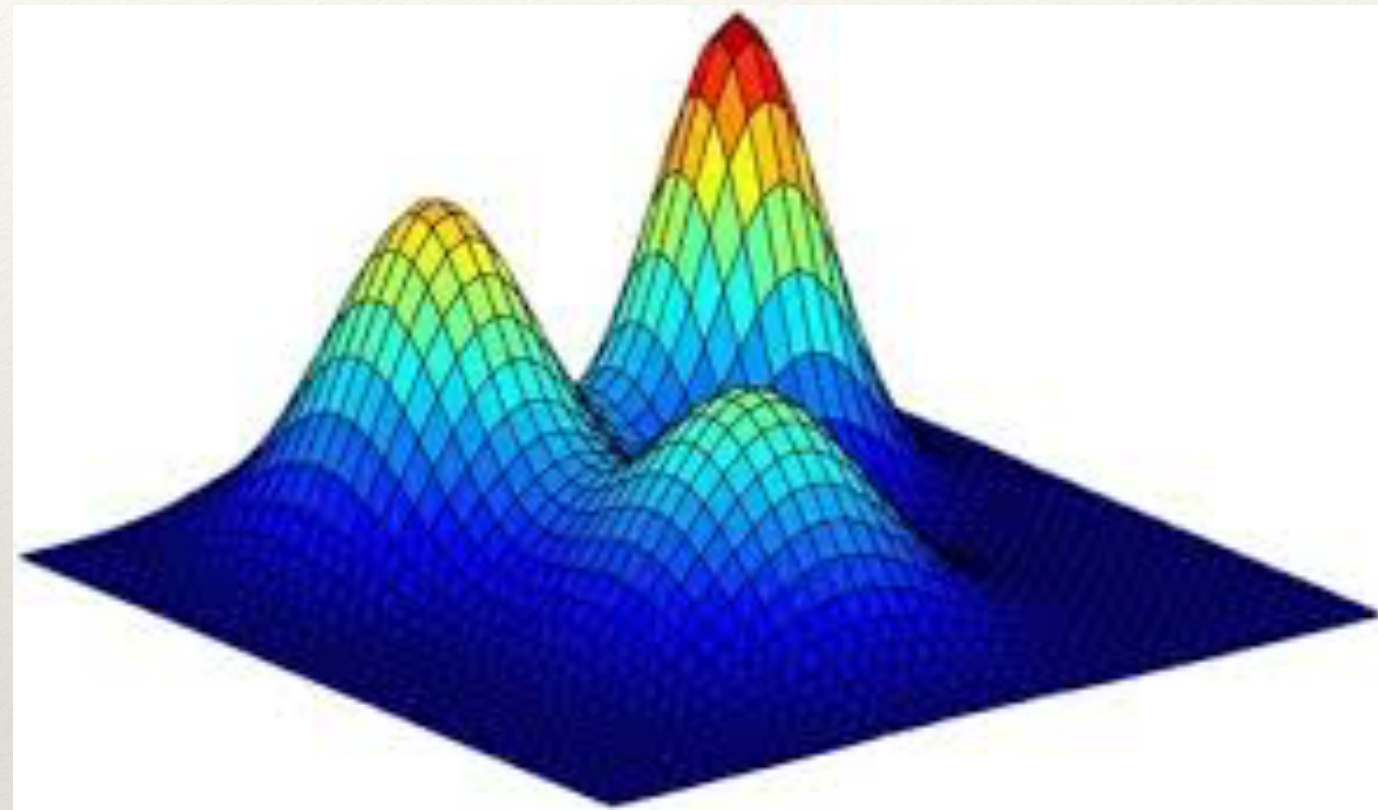
$$p(\mathbf{x}|\theta_k) = \frac{1}{\sqrt{(2\pi)^D |\Sigma_k|}} \exp\left\{ -\frac{1}{2}(\mathbf{x} - \mu_k)^* \Sigma_k^{-1} (\mathbf{x} - \mu_k) \right\}$$

The weighting coefficients have the property

$$\sum_{k=1}^K \alpha_k = 1$$

Gaussian Mixture Modeling

- Properties of GMM
 - Can model multi-modal data.
 - Identify data clusters.
 - Can model arbitrarily complex data distributions



The set of parameters for the model are

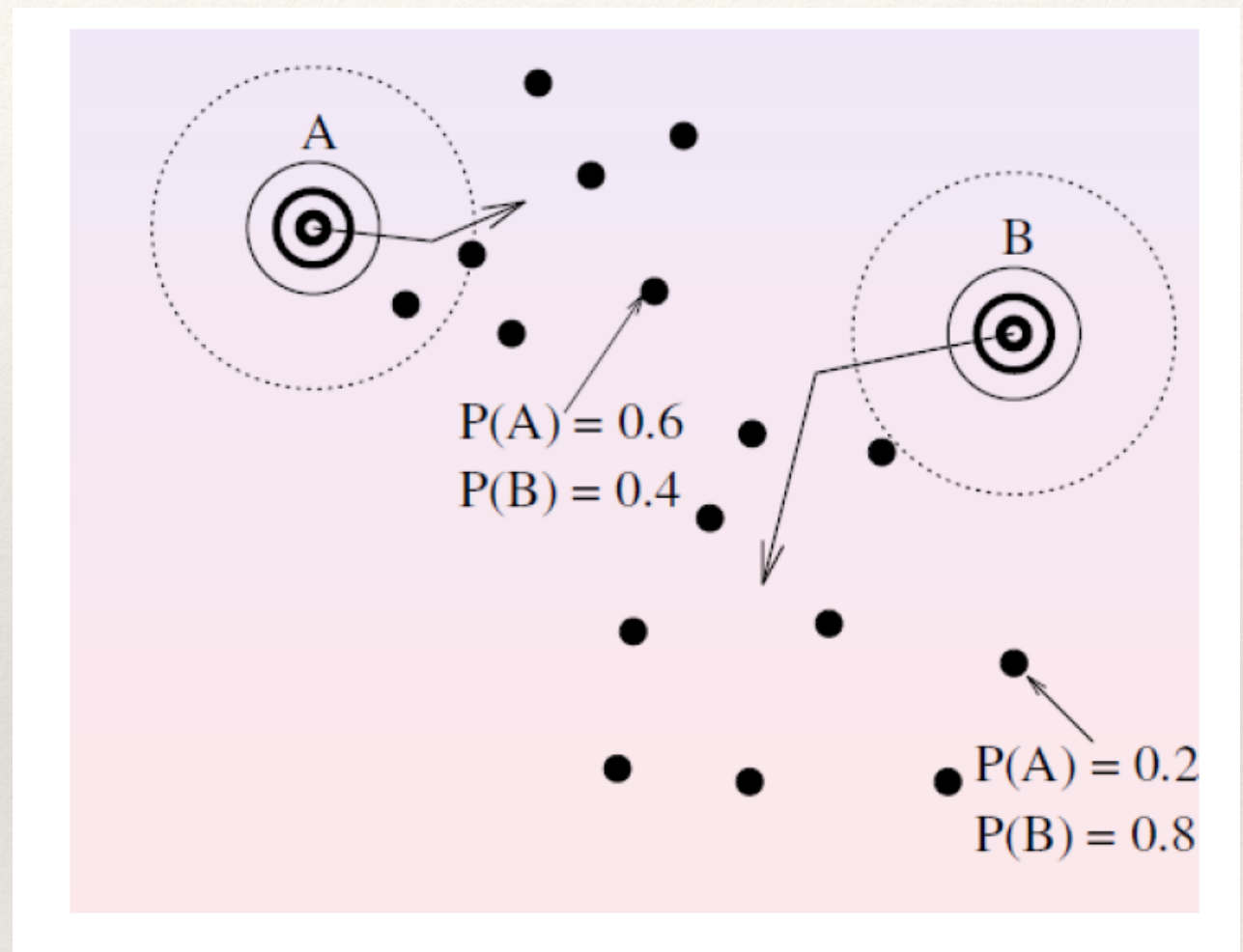
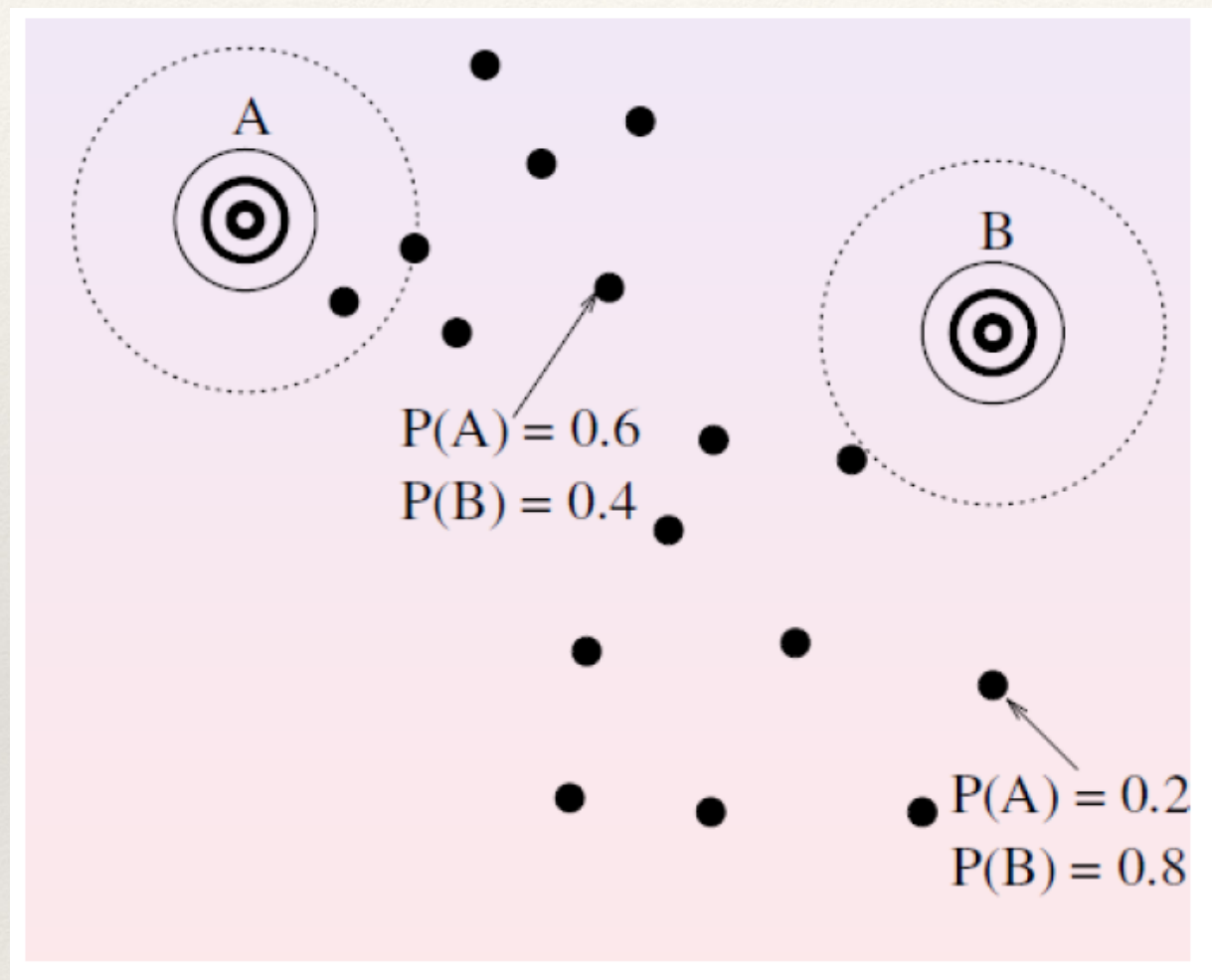
$$\Theta_k = \{\alpha_k, \theta_k\}_{k=1}^K \quad \theta_k = \{\mu_k, \Sigma_k\}$$

-
-
- The log-likelihood function over the entire data in this case will have a **logarithm of a summation**

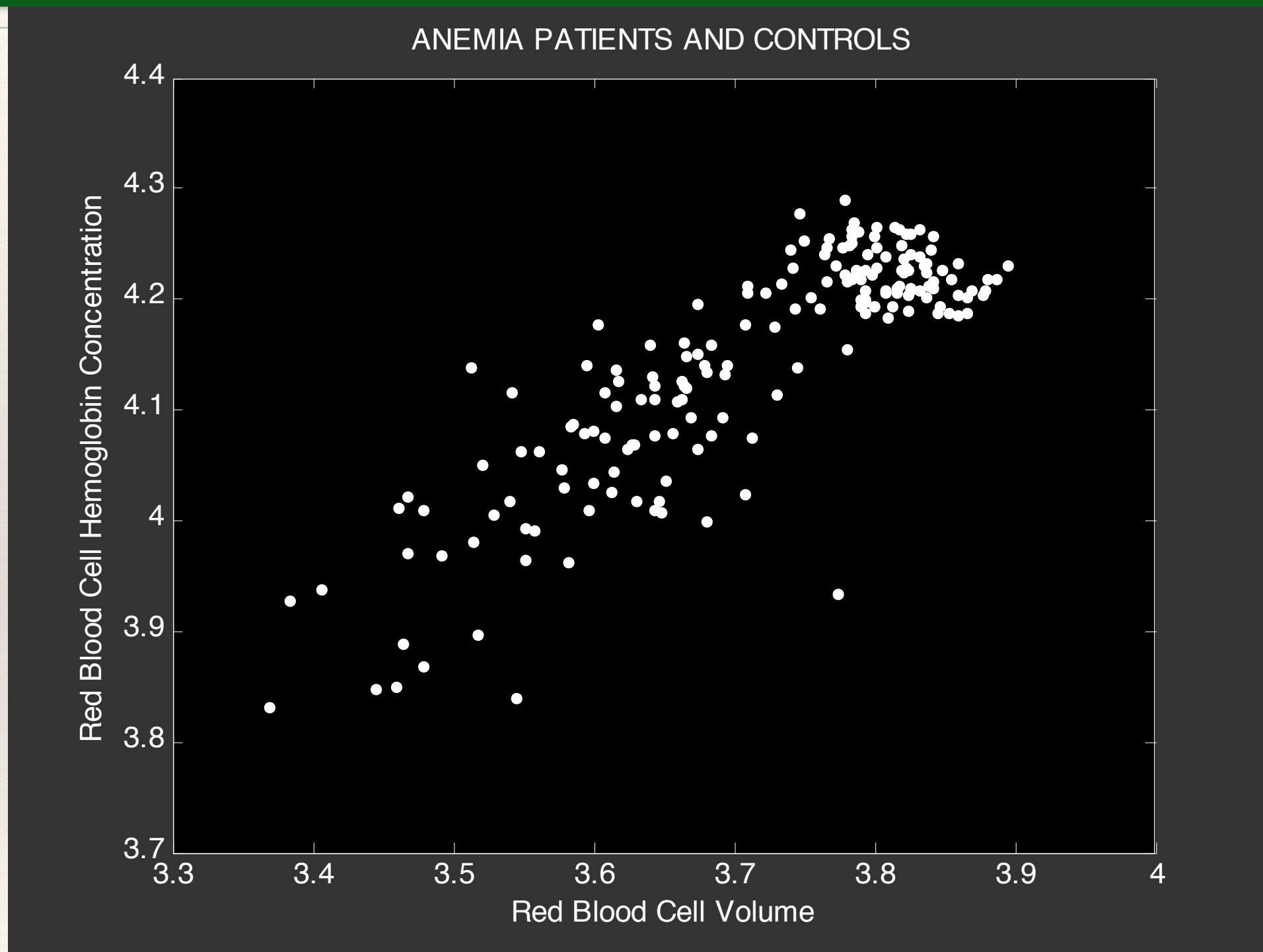
$$\log L(\Theta) = \sum_{i=1}^N \log \left(\sum_{k=1}^K \alpha_k p(\mathbf{x}_i | \theta_k) \right)$$

- Solving for the optimal parameters using MLE for GMM is not straight forward.
- Resort to the **Expectation Maximization (EM)** algorithm

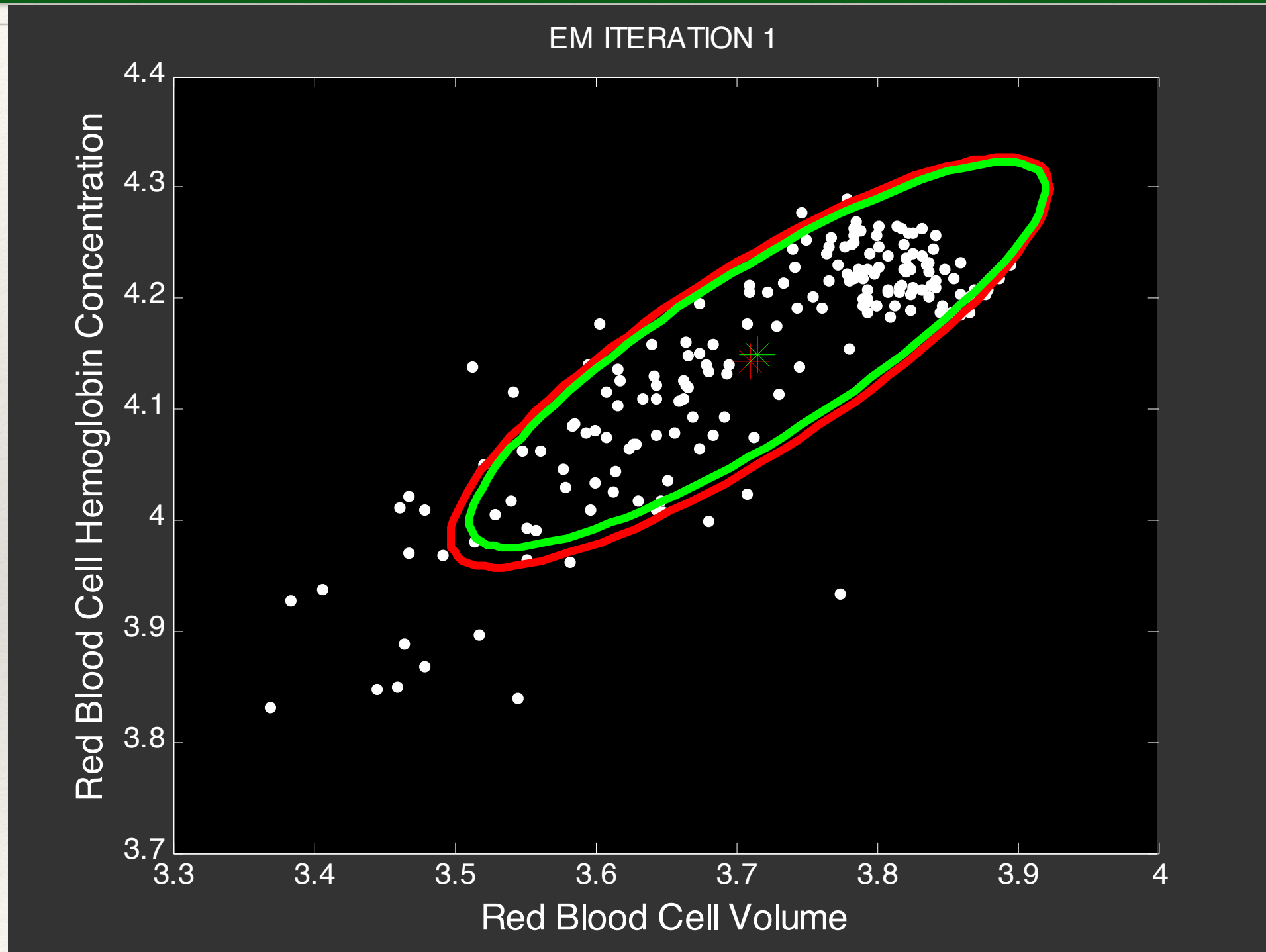
MLE for GMM



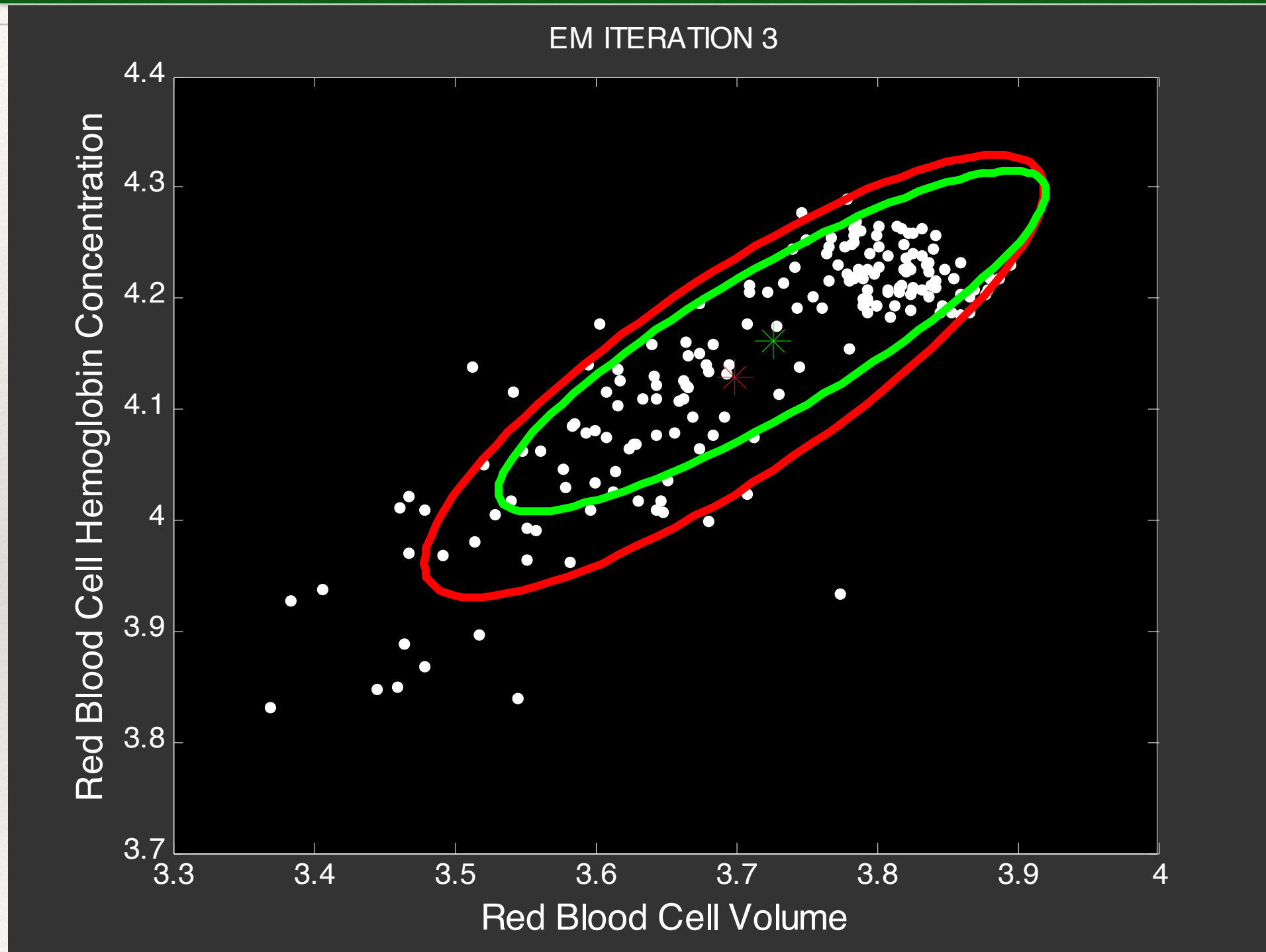
EM Algorithm for GMM



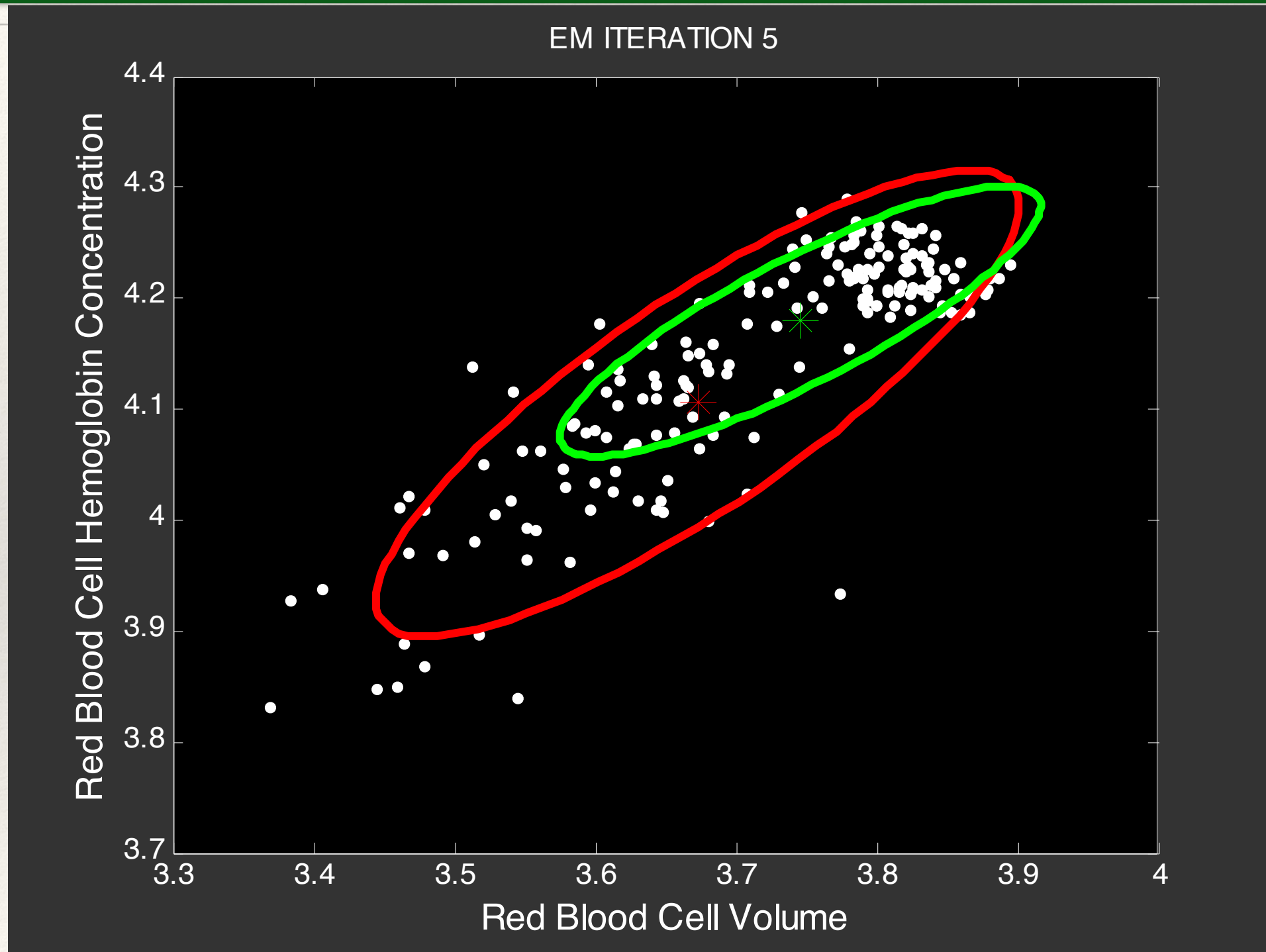
EM Algorithm for GMM



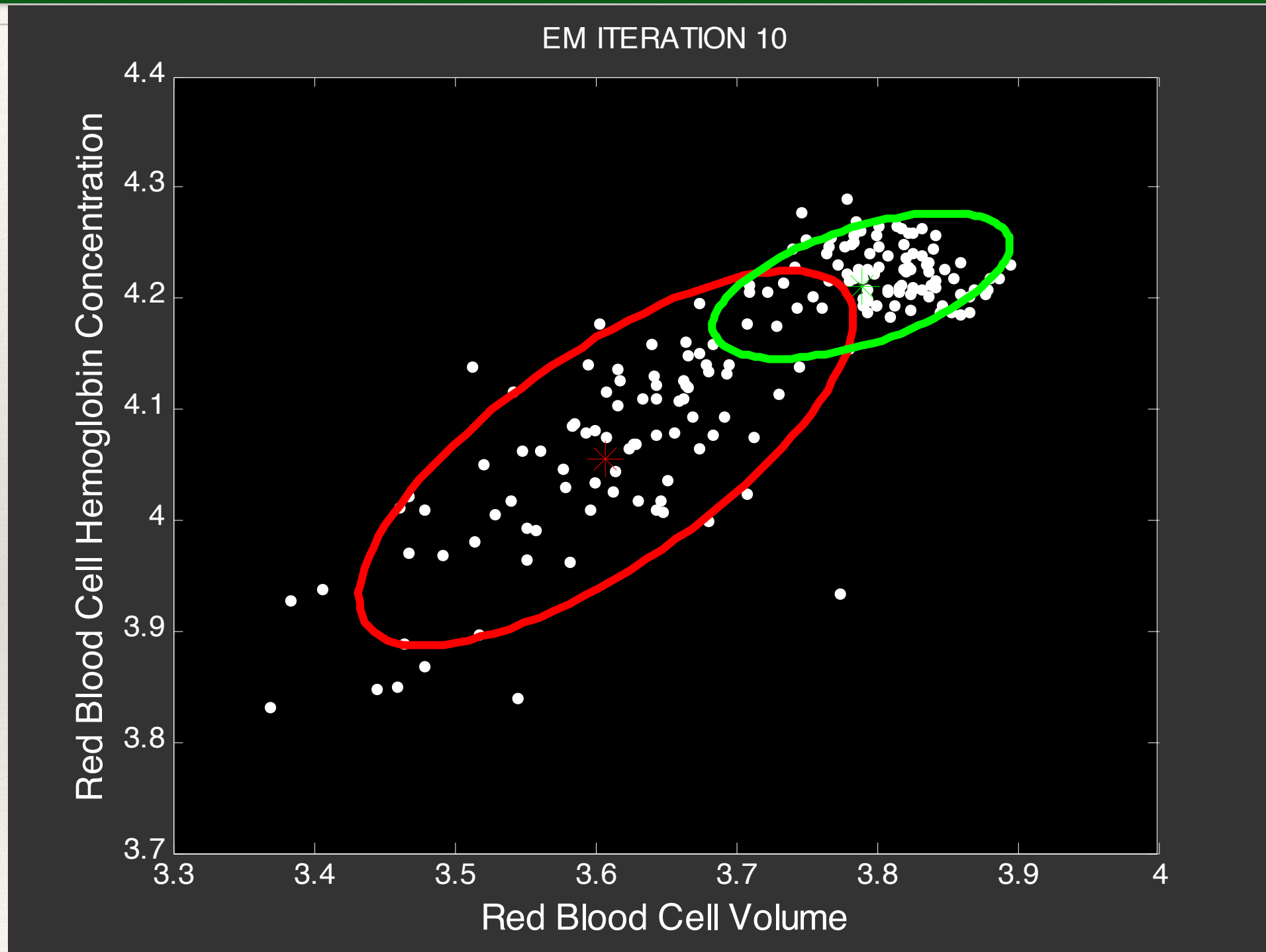
EM Algorithm for GMM



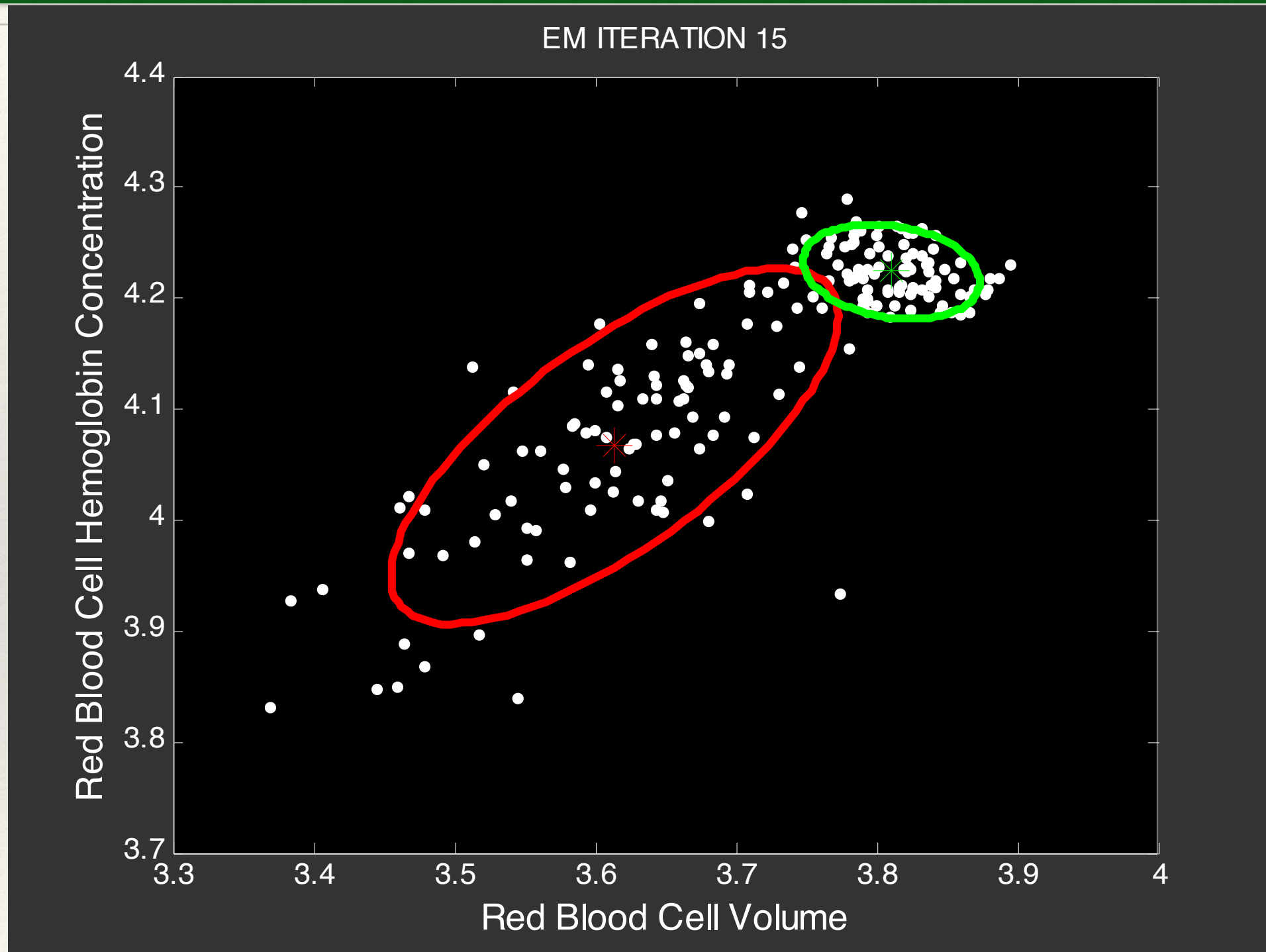
EM Algorithm for GMM



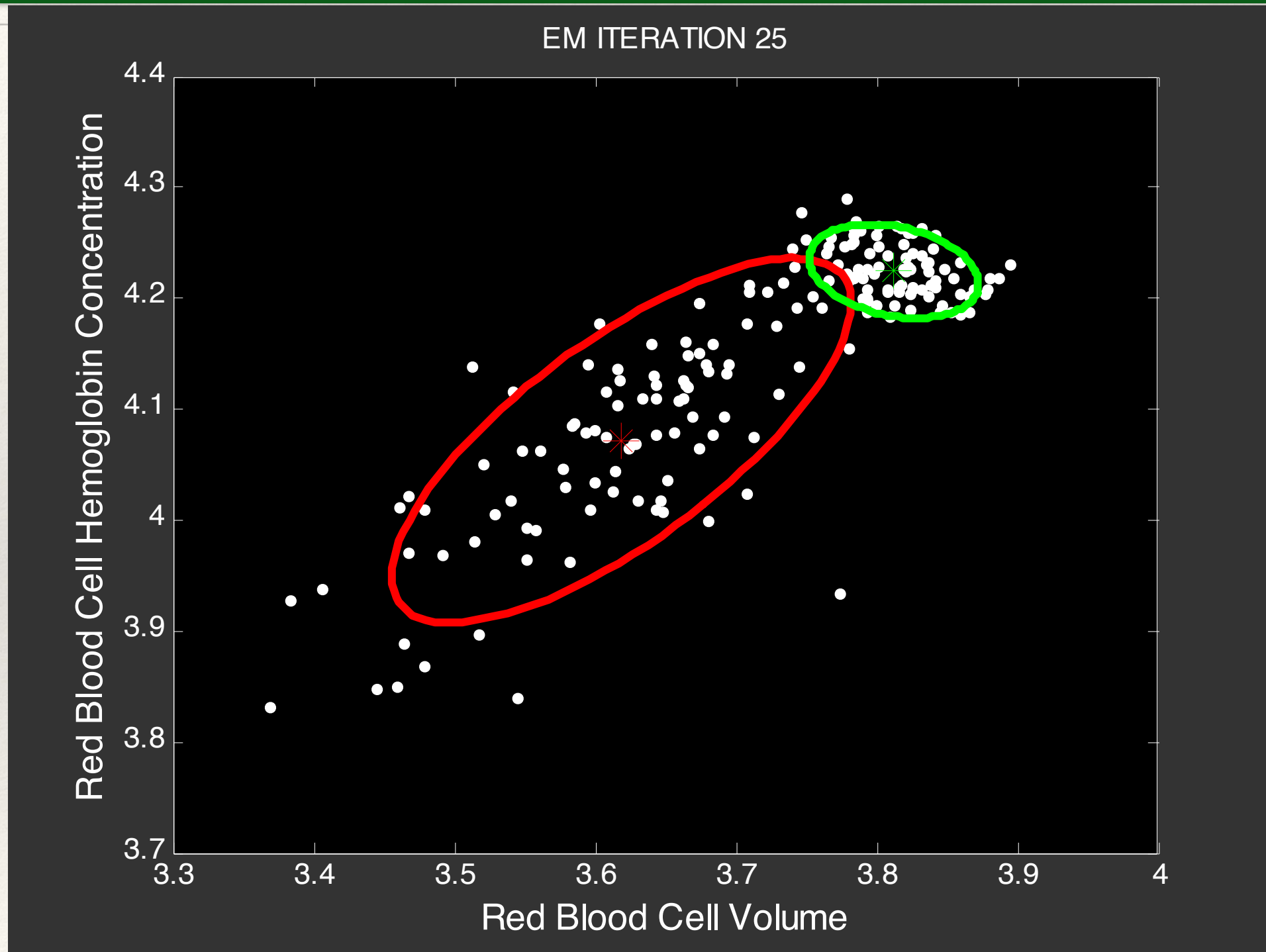
EM Algorithm for GMM



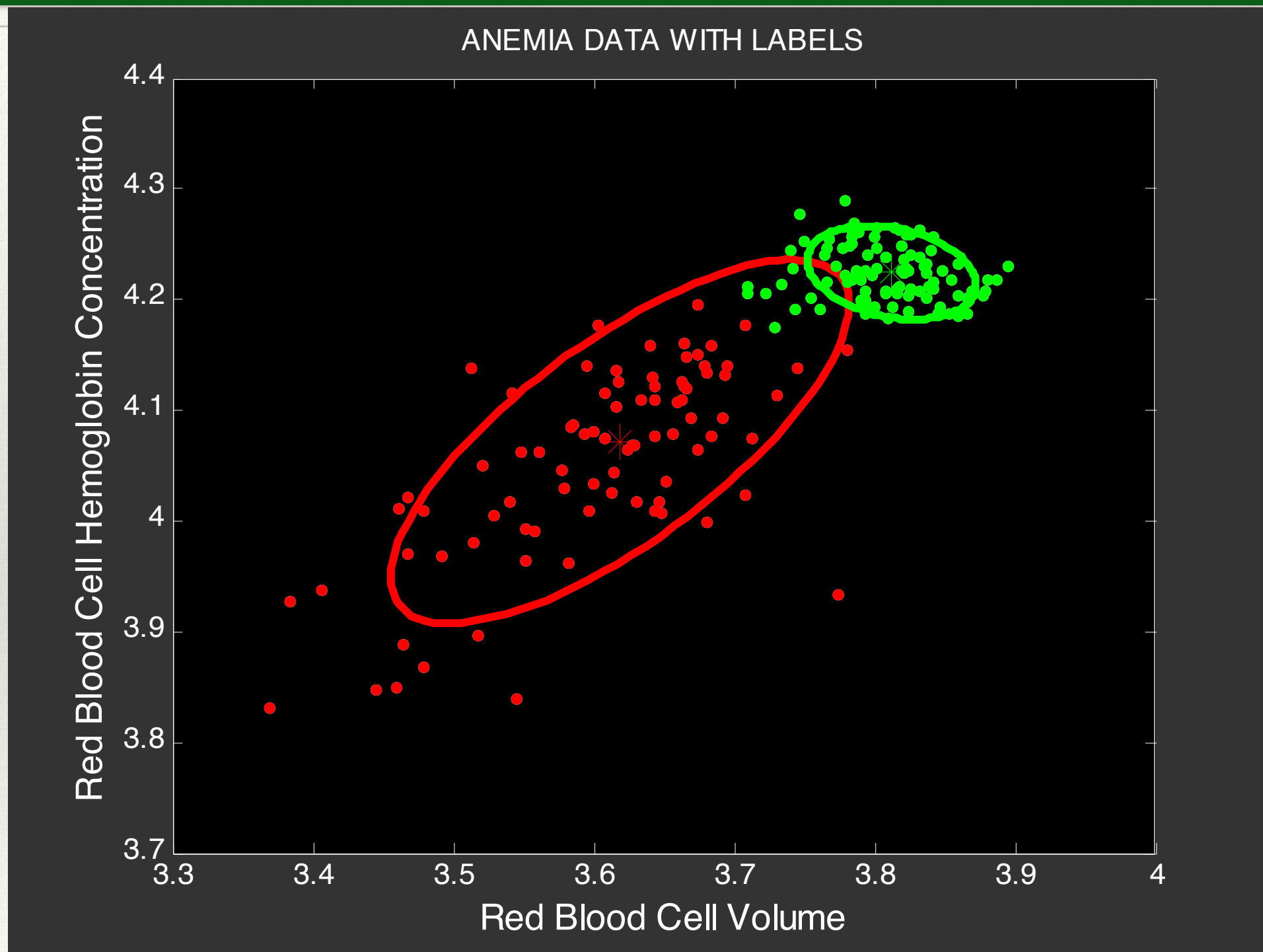
EM Algorithm for GMM



EM Algorithm for GMM

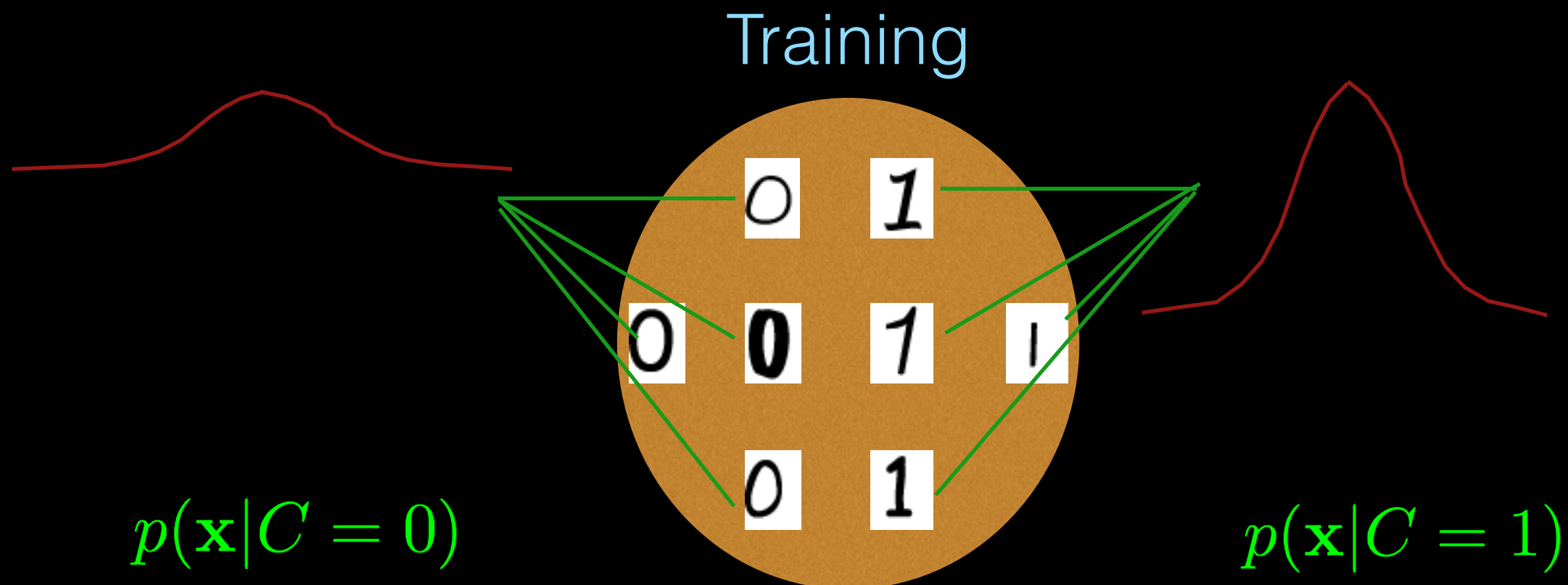


EM Algorithm for GMM



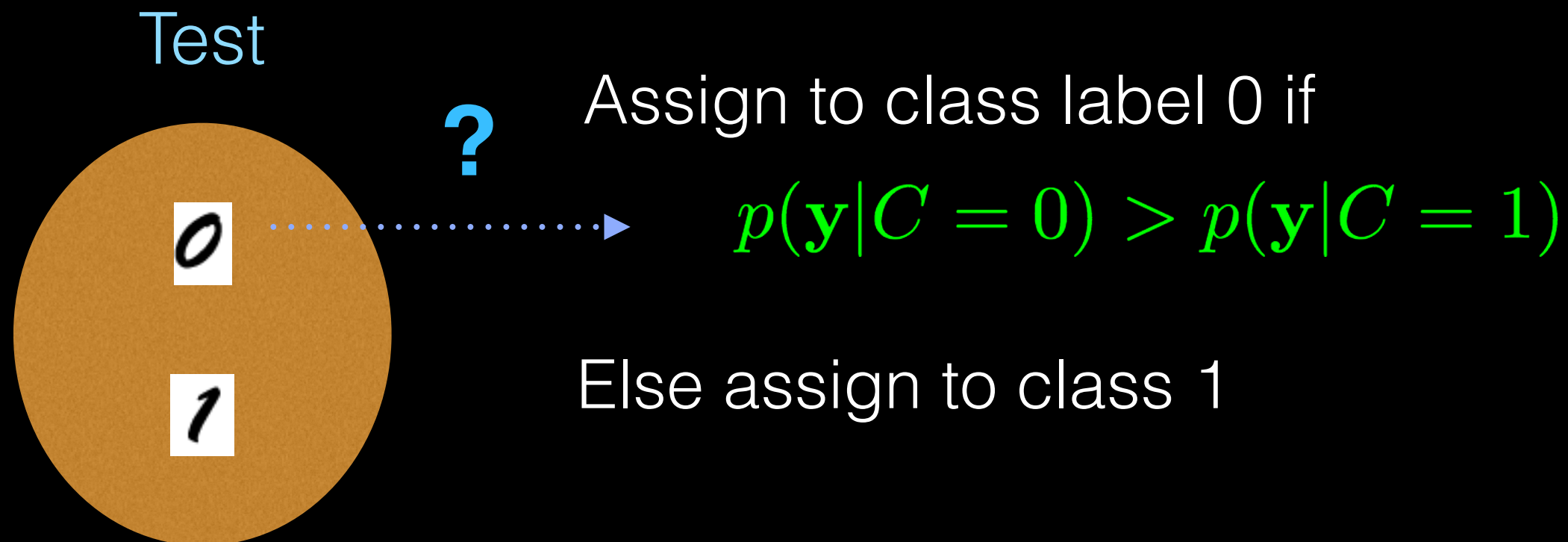
Generative classifier

- Model the two classes separately using probability distributions -
 - Make each sample \mathbf{X}_i (28x28) as a vector \mathbf{x}_i of size 784.
 - Build class dependent probability $p(\mathbf{x}|C = 0)$ & $p(\mathbf{x}|C = 1)$

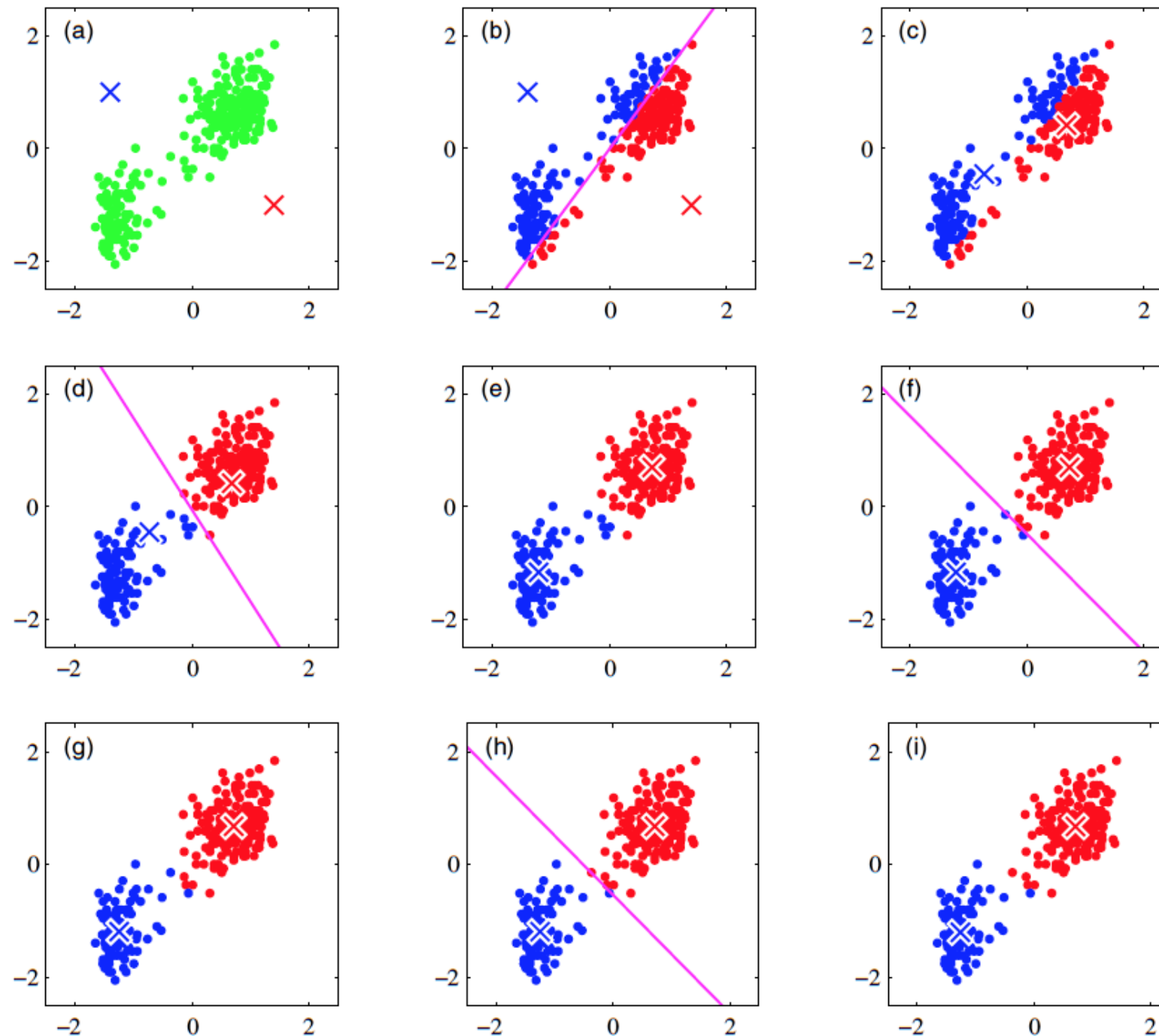


Generative classifier

- For the test sample
 - Make each sample \mathbf{Y} (28x28) as a vector \mathbf{y} of size 784.
 - Compute the probability of generating sample \mathbf{y} for each class.



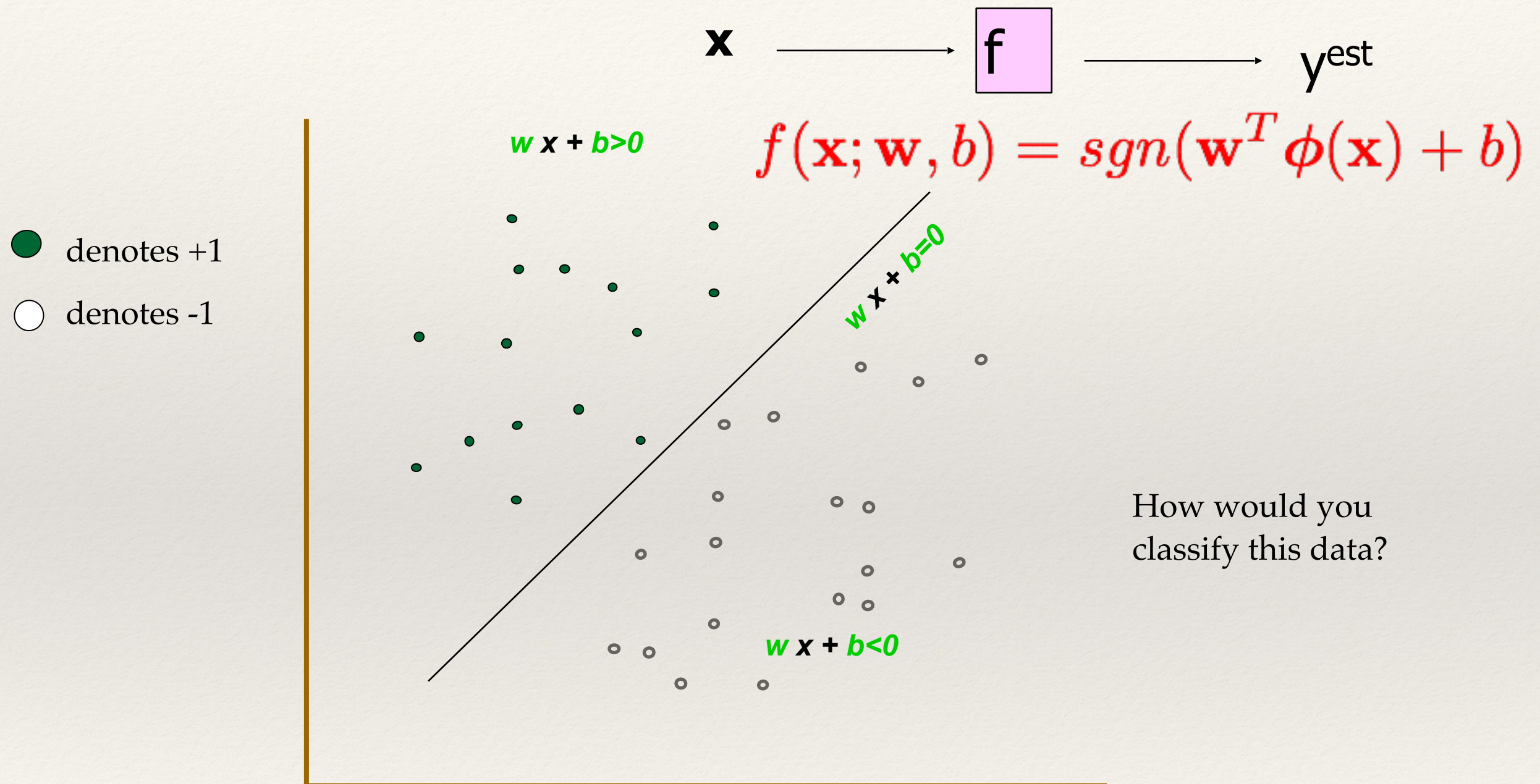
K-means Algorithm for Initialization



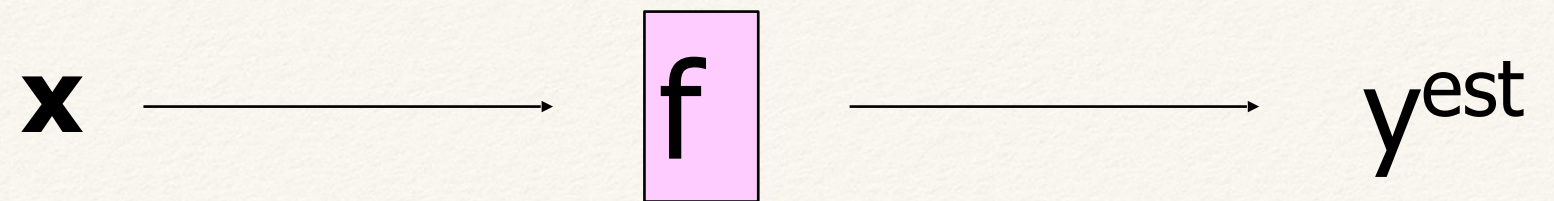
Discriminative models and Discriminant Functions

Support Vector Machines

Linear Classifiers

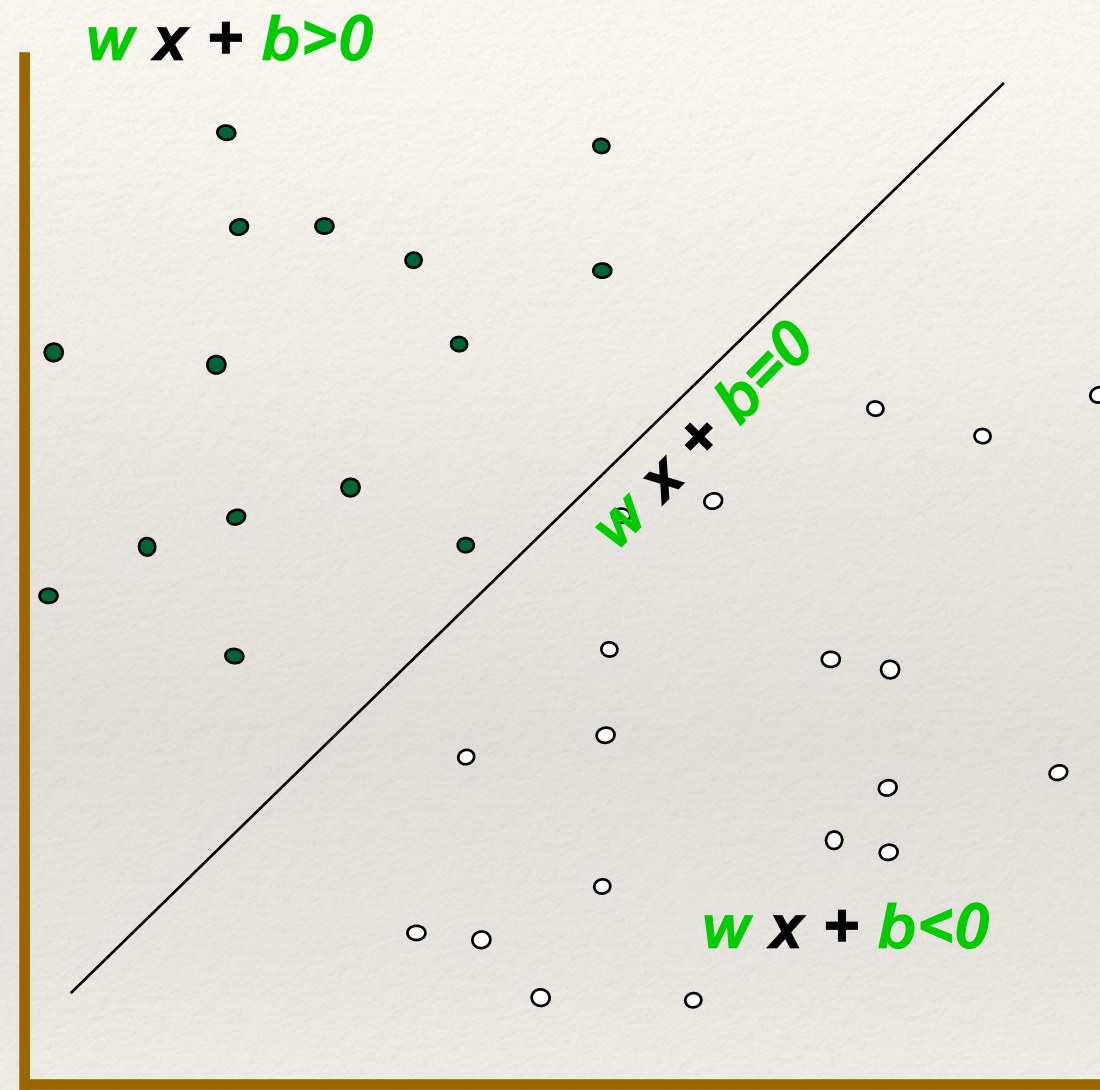


Linear Classifiers



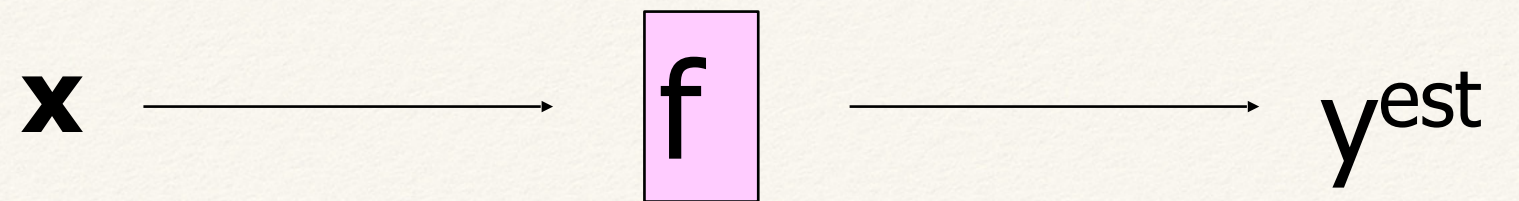
$$f(\mathbf{x}; \mathbf{w}, b) = \text{sgn}(\mathbf{w}^T \phi(\mathbf{x}) + b)$$

- denotes +1
- denotes -1



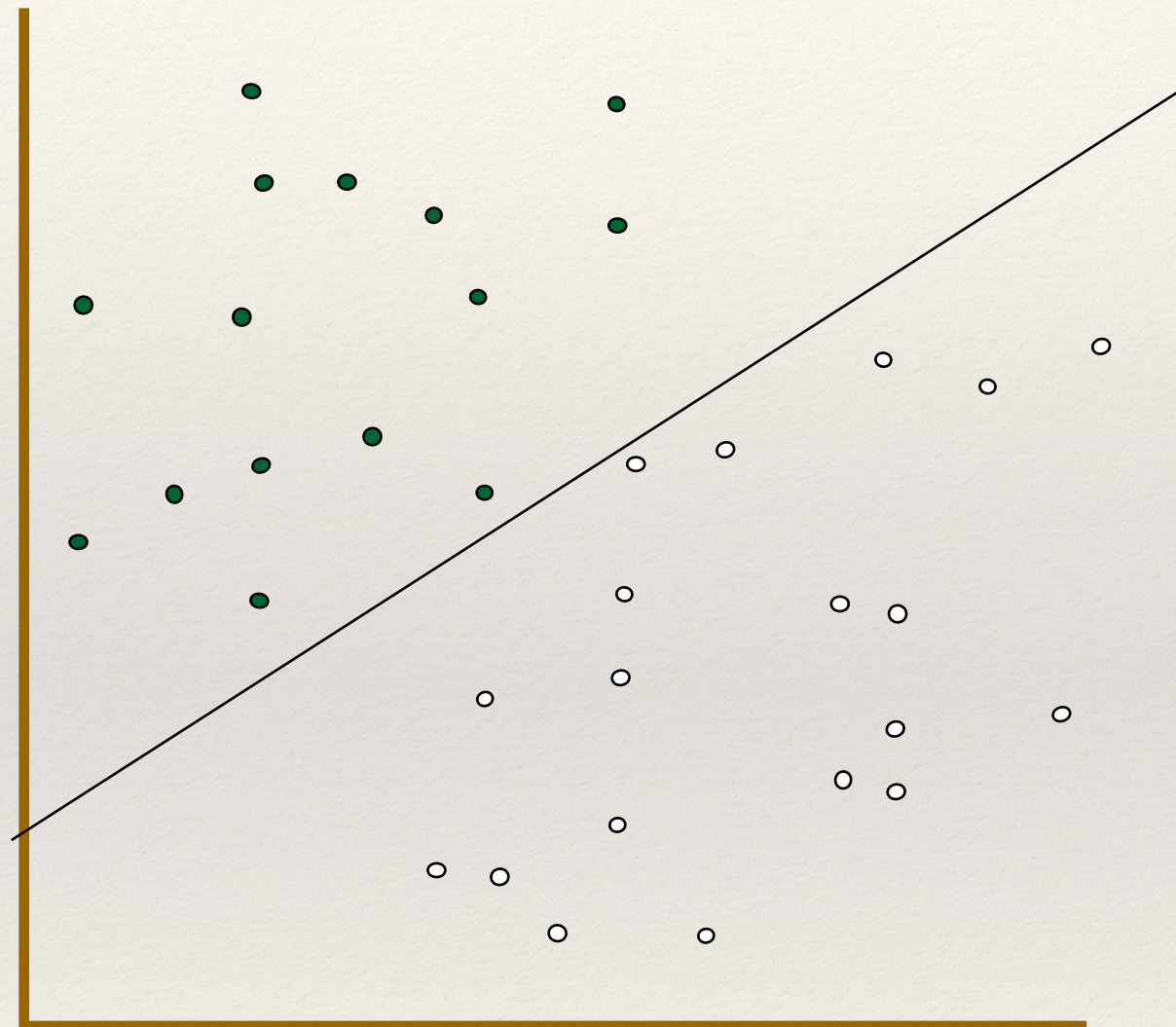
How would you classify this data?

Linear Classifiers



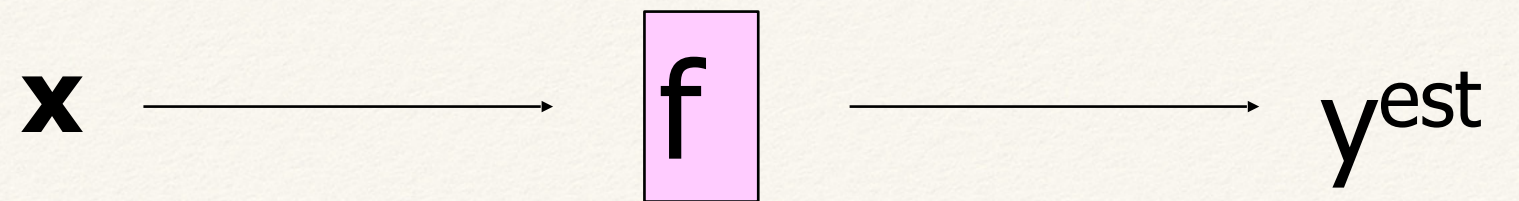
$$f(\mathbf{x}; \mathbf{w}, b) = \text{sgn}(\mathbf{w}^T \phi(\mathbf{x}) + b)$$

- denotes +1
- denotes -1

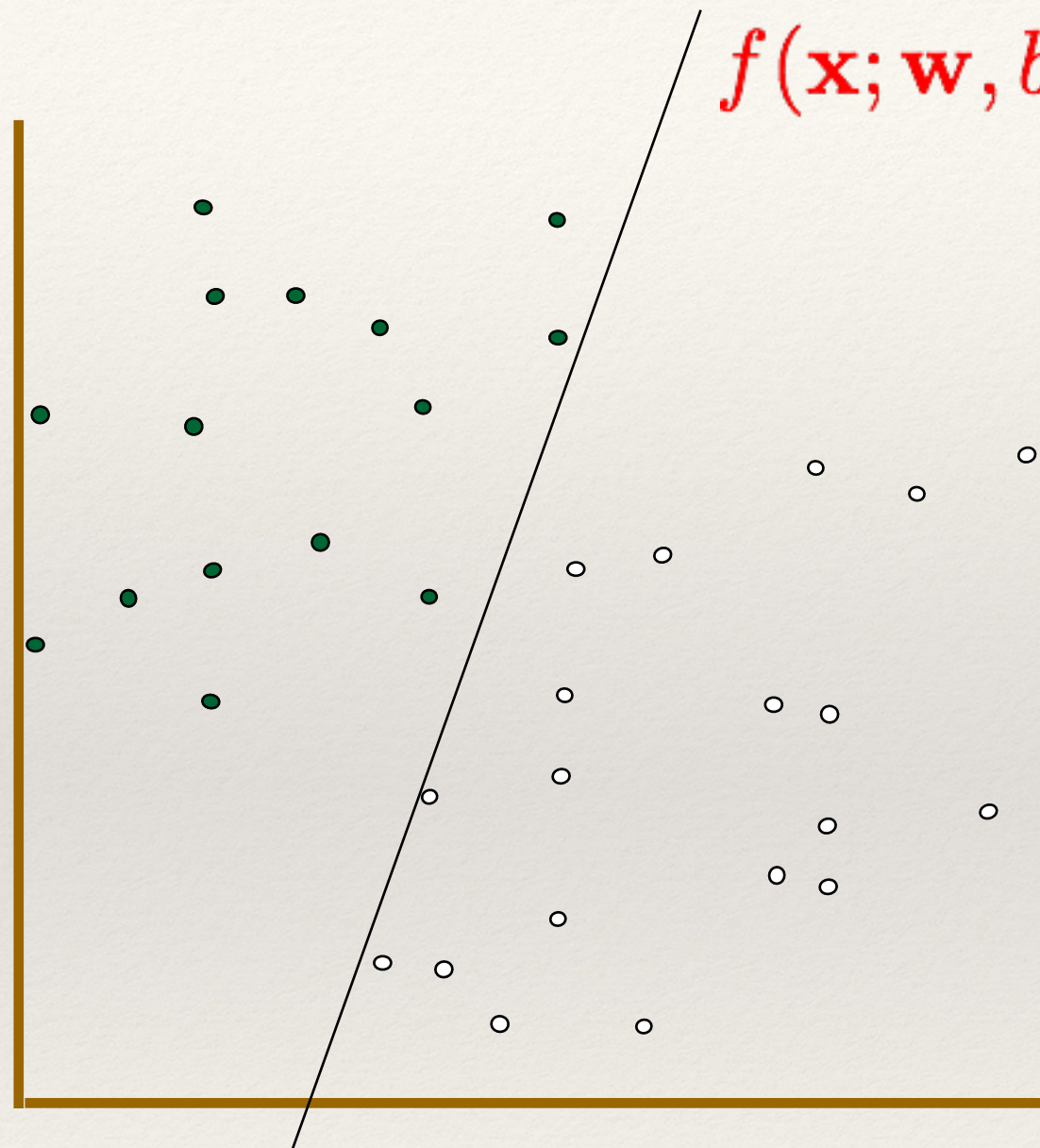


How would you classify this data?

Linear Classifiers



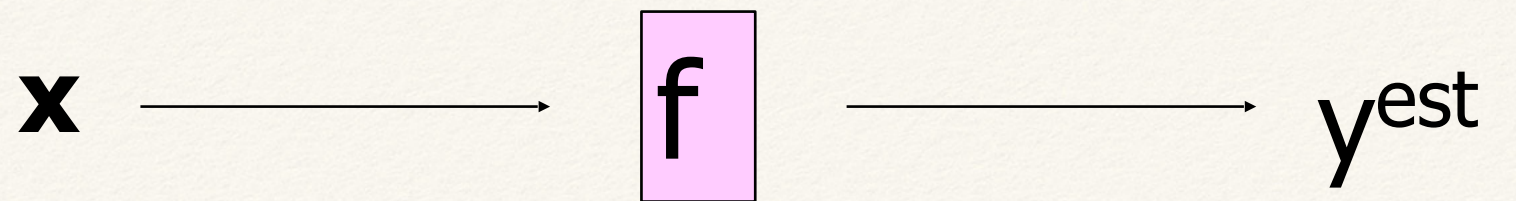
- denotes +1
- denotes -1



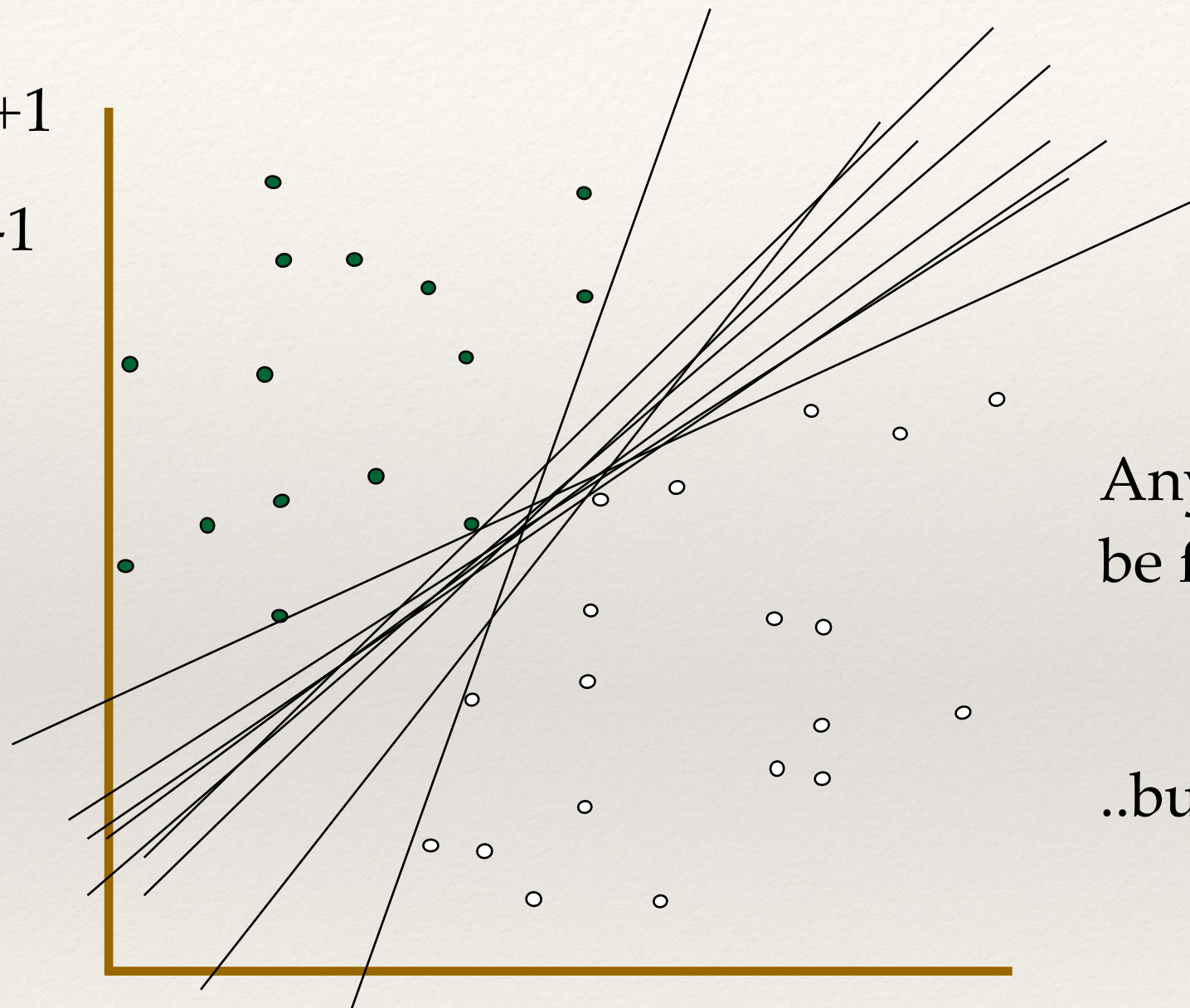
$$f(\mathbf{x}; \mathbf{w}, b) = \text{sgn}(\mathbf{w}^T \phi(\mathbf{x}) + b)$$

How would you
classify this data?

Linear Classifiers



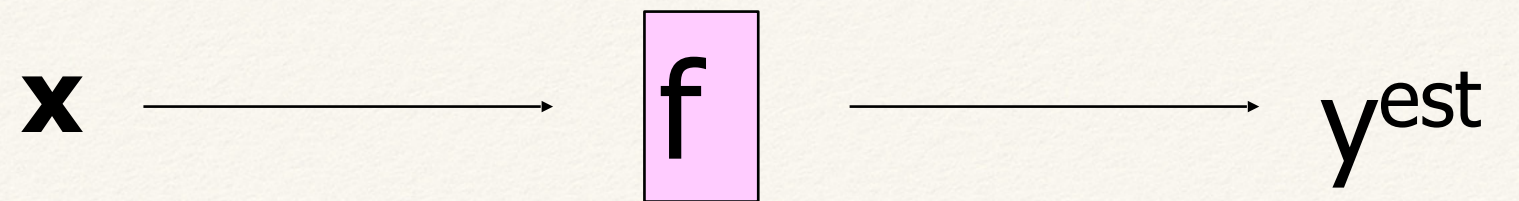
- denotes +1
- denotes -1



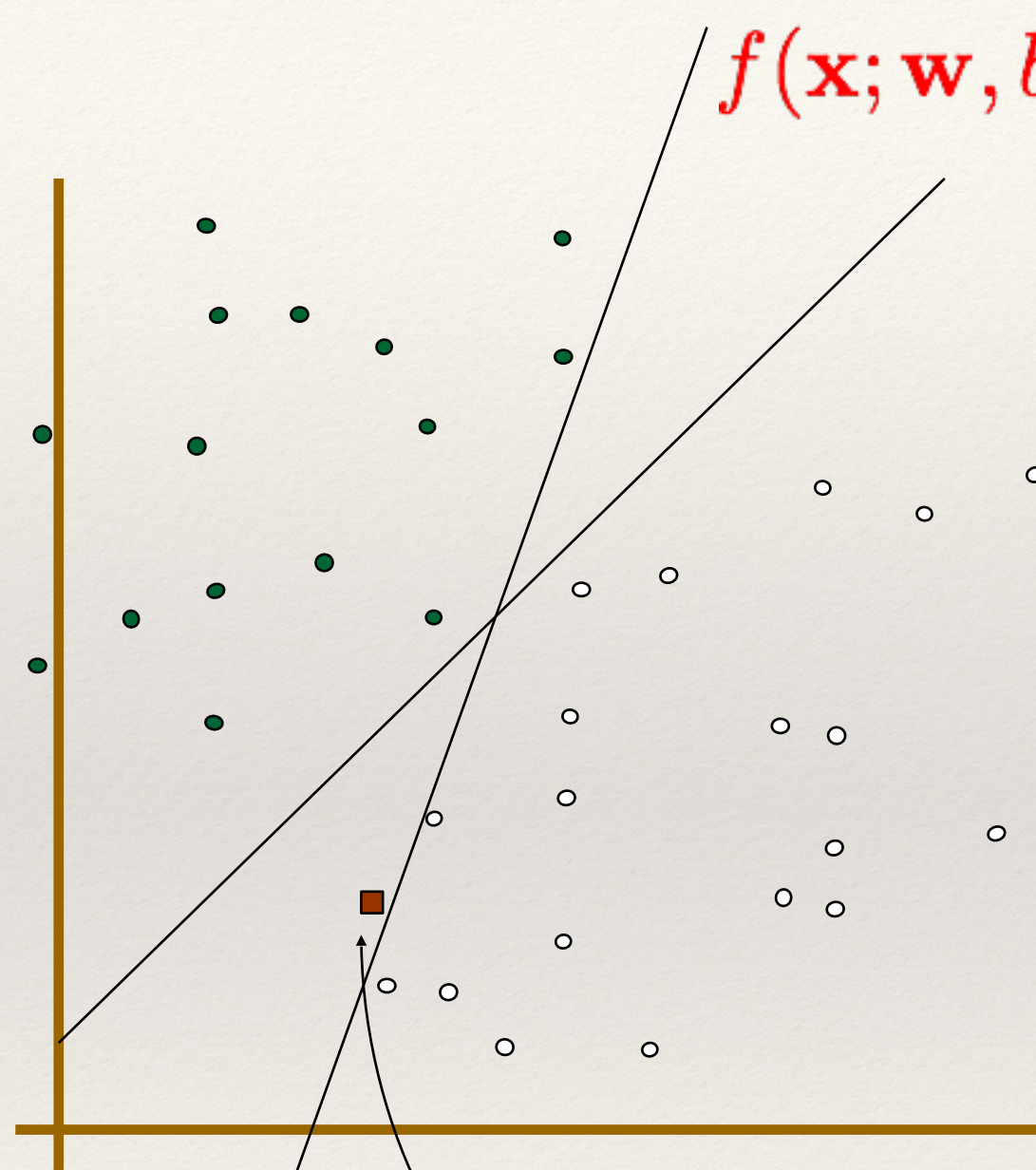
Any of these would be fine..

..but which is best?

Linear Classifiers



- denotes +1
- denotes -1

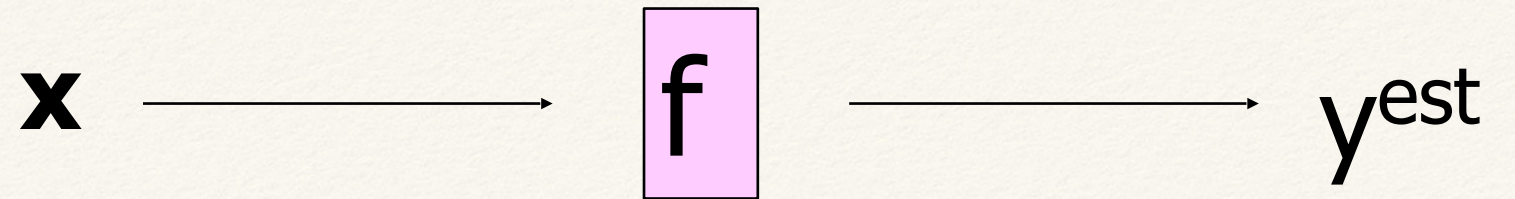


$$f(\mathbf{x}; \mathbf{w}, b) = \text{sgn}(\mathbf{w}^T \phi(\mathbf{x}) + b)$$

How would you classify this data?

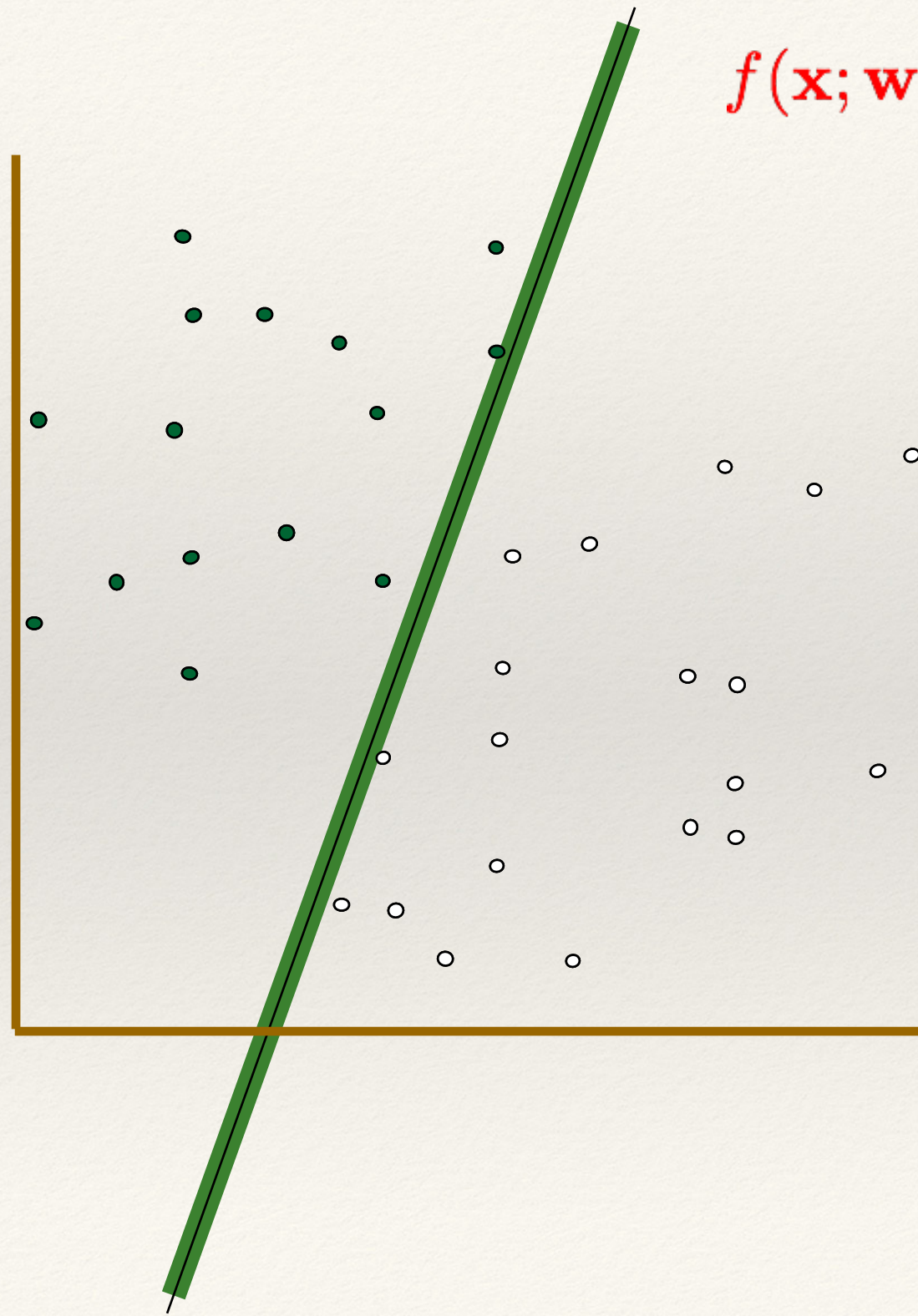
Misclassified
to +1 class

Linear Classifiers



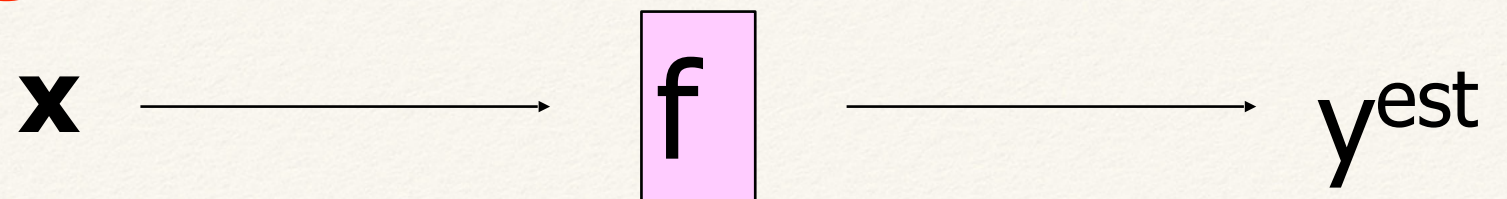
$$f(\mathbf{x}; \mathbf{w}, b) = \text{sgn}(\mathbf{w}^T \phi(\mathbf{x}) + b)$$

- denotes +1
- denotes -1



Define the **margin** of a linear classifier as the width that the boundary could be increased by before hitting a datapoint.

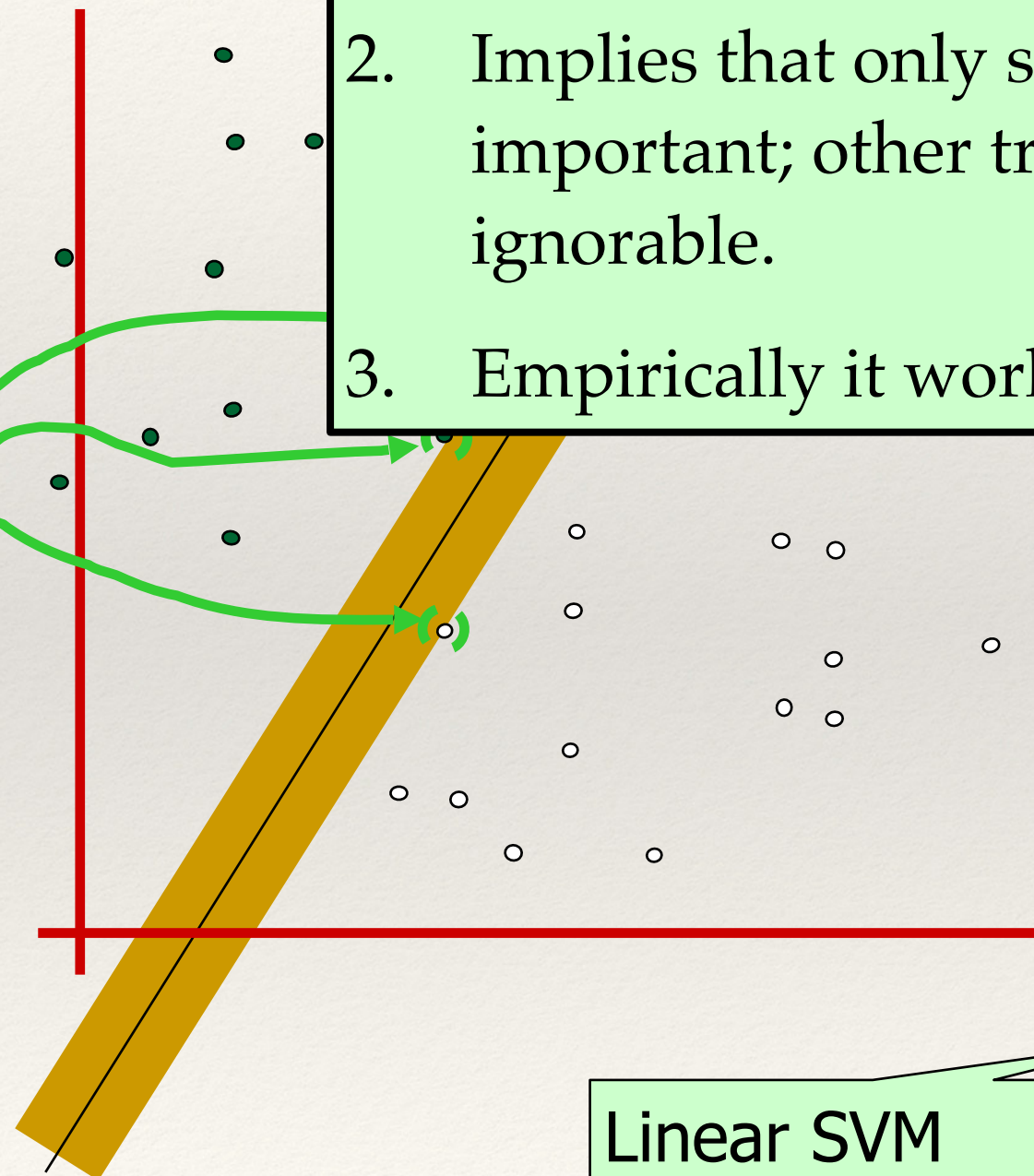
Maximum Margin



- denotes +1
- denotes -1

1. Maximizing the margin is good according to intuition
2. Implies that only support vectors are important; other training examples are ignorable.
3. Empirically it works very very well.

Support Vectors
are those data points that the margin pushes up against



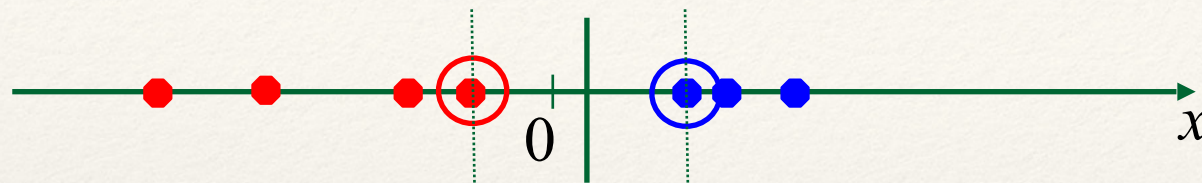
linear classifier
with the, um,
maximum margin.

This is the simplest
kind of SVM
(Called an LSVM)

Linear SVM

Non-linear SVMs

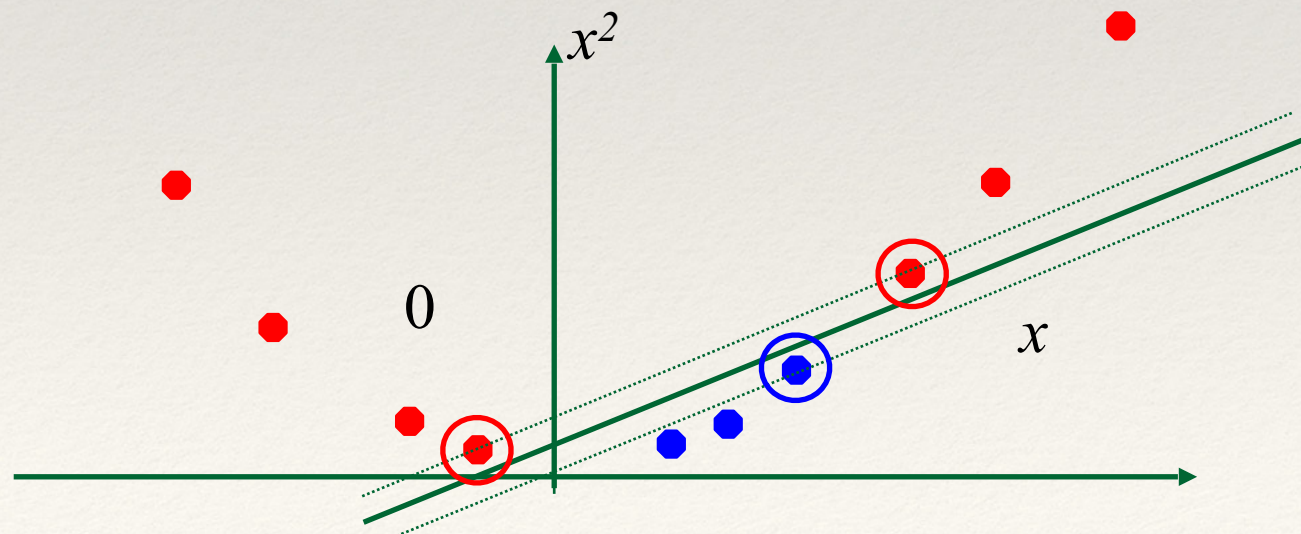
- Datasets that are linearly separable with some noise work out great:



- But what are we going to do if the dataset is just too hard?



- How about... mapping data to a higher-dimensional space:



Non-linear SVMs: Feature spaces

- General idea: the original input space can always be mapped to some higher-dimensional feature space where the training set is separable:

