

# Outline

- Bayesian concept learning: Discussion
- Probabilistic models for unsupervised and semi-supervised category learning

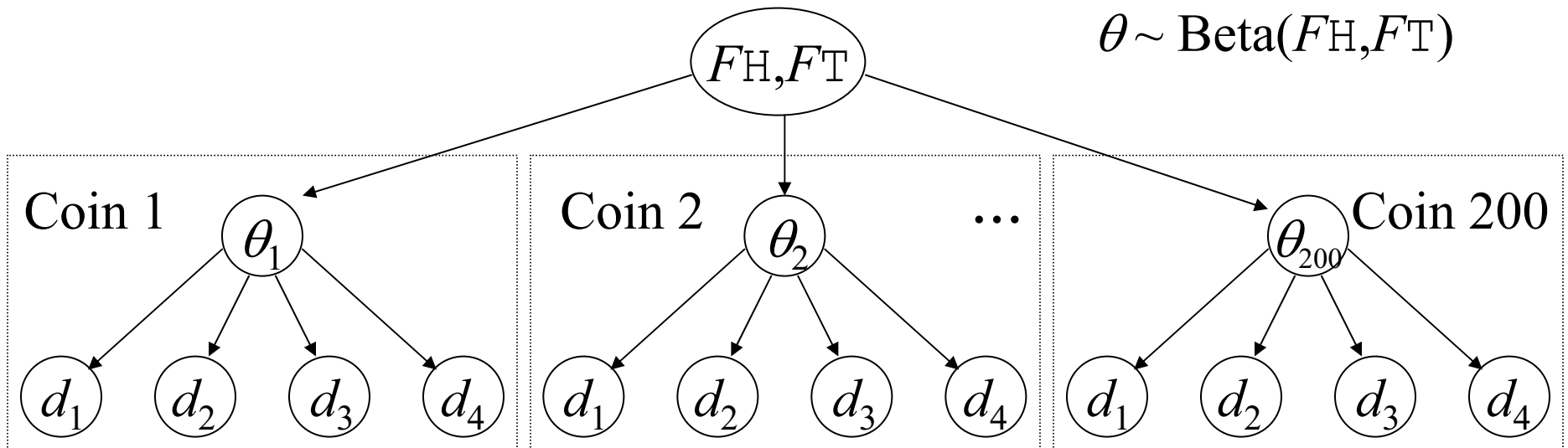
# Discussion points

- Relation to “Bayesian classification”?
- Relation to debate between rules / logic / symbols and similarity / connections / statistics?
- Where do the hypothesis space and prior probability distribution come from?

# Discussion points

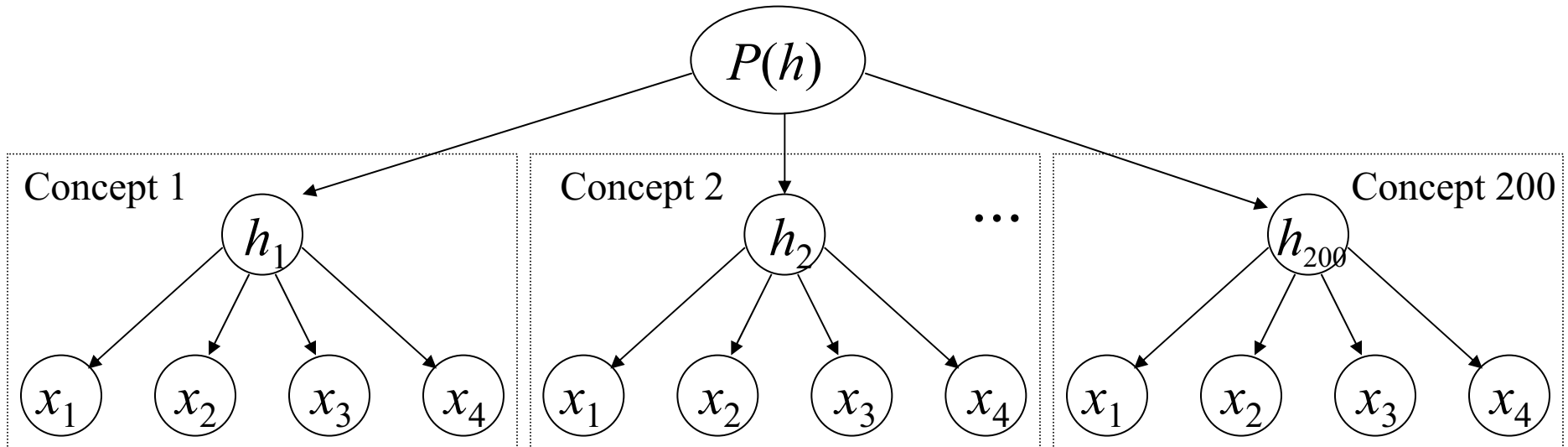
- Relation to “Bayesian classification”?
  - Causal attribution versus referential inference.
  - Which is more suited to natural concept learning?
- Relation to debate between rules / logic / symbols and similarity / connections / statistics?
- Where do the hypothesis space and prior probability distribution come from?

# Hierarchical priors

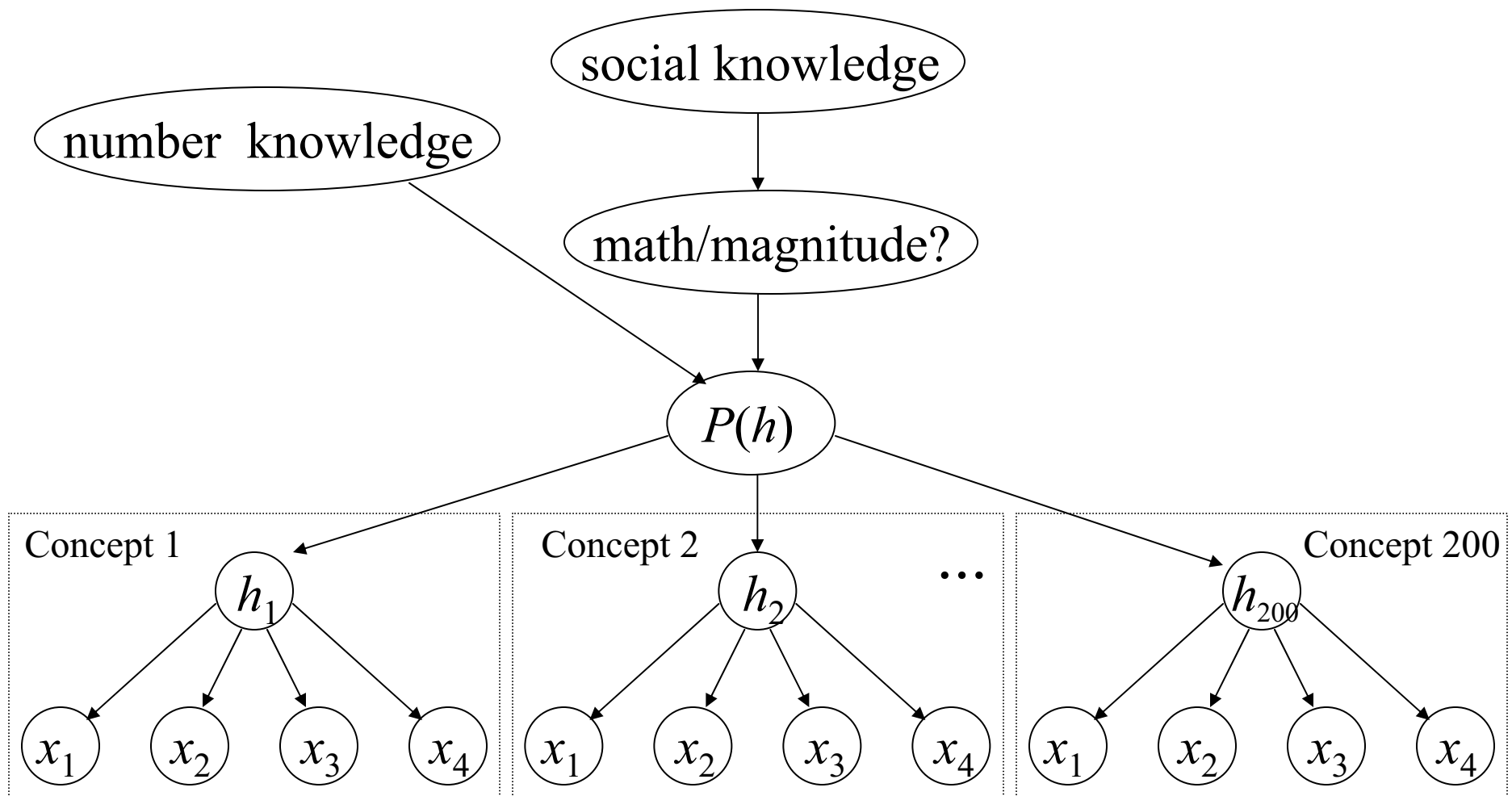


- Latent structure captures what is common to all coins, and also their individual variability

# Hierarchical priors



- Latent structure captures what is common to all concepts, and also their individual variability
- *Is this all we need?*



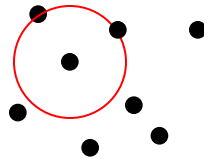
- Hypothesis space is not just an arbitrary collection of hypotheses, but a principled system.
- Far more structured than our experience with specific number concepts.

# Outline

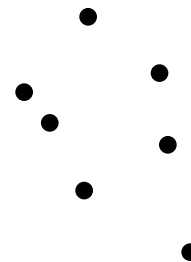
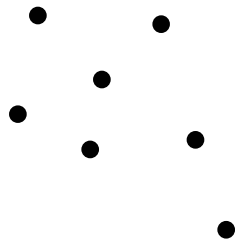
- Bayesian concept learning: Discussion
- Probabilistic models for unsupervised and semi-supervised category learning

# Simple model of concept learning

“This is a blicket.”




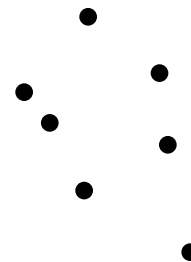
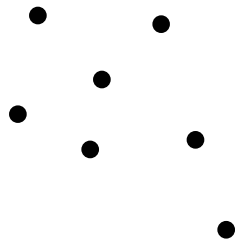
“Can you show me the other blickets?”





# Simple model of concept learning

 Other blickets.

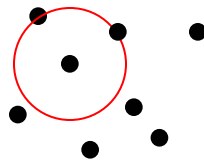


# The objects of planet Gazoob

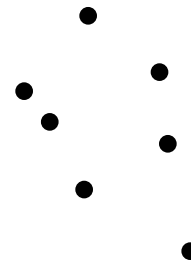
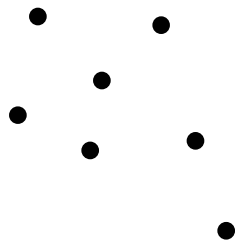
Image removed due to copyright considerations.

# Simple model of concept learning

“This is a blicket.”



“Can you show me the other blickets?”



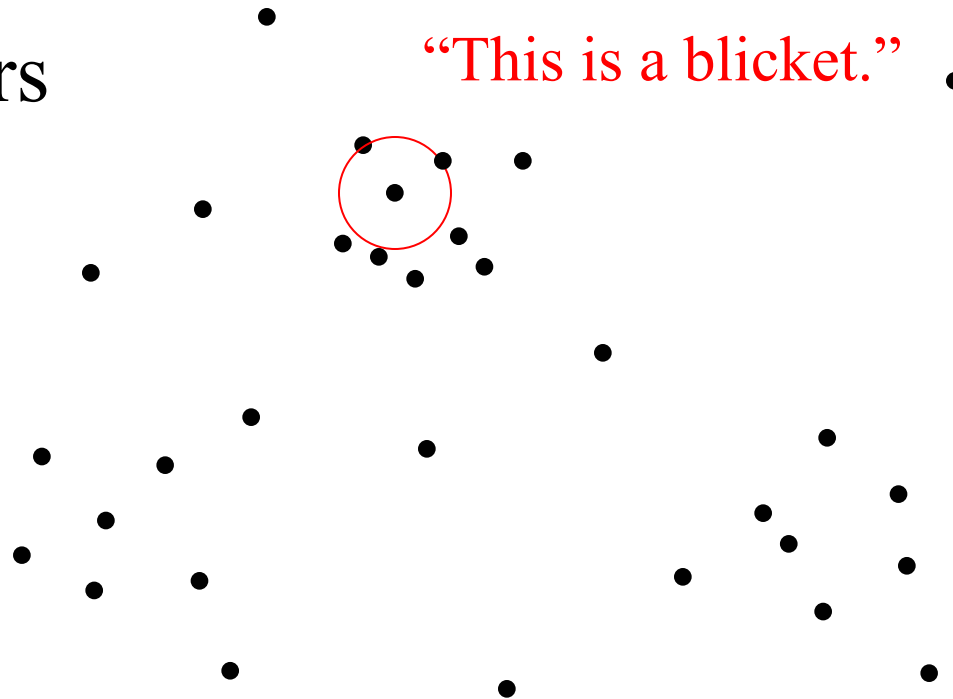
Learning from just one positive example is possible if:

- Assume concepts refer to clusters in the world.
- Observe enough unlabeled data to identify clear clusters.

# Complications

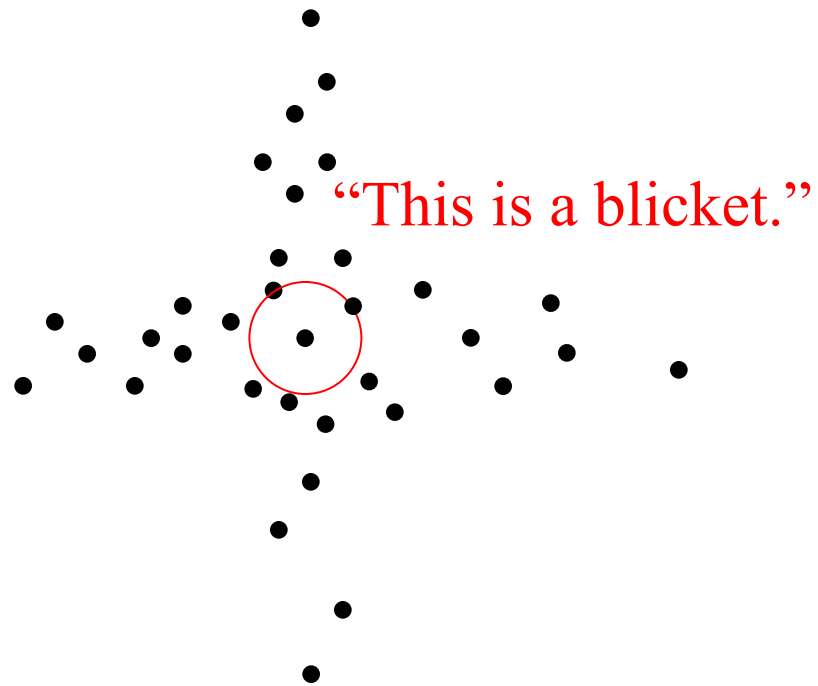
# Complications

- Outliers



# Complications

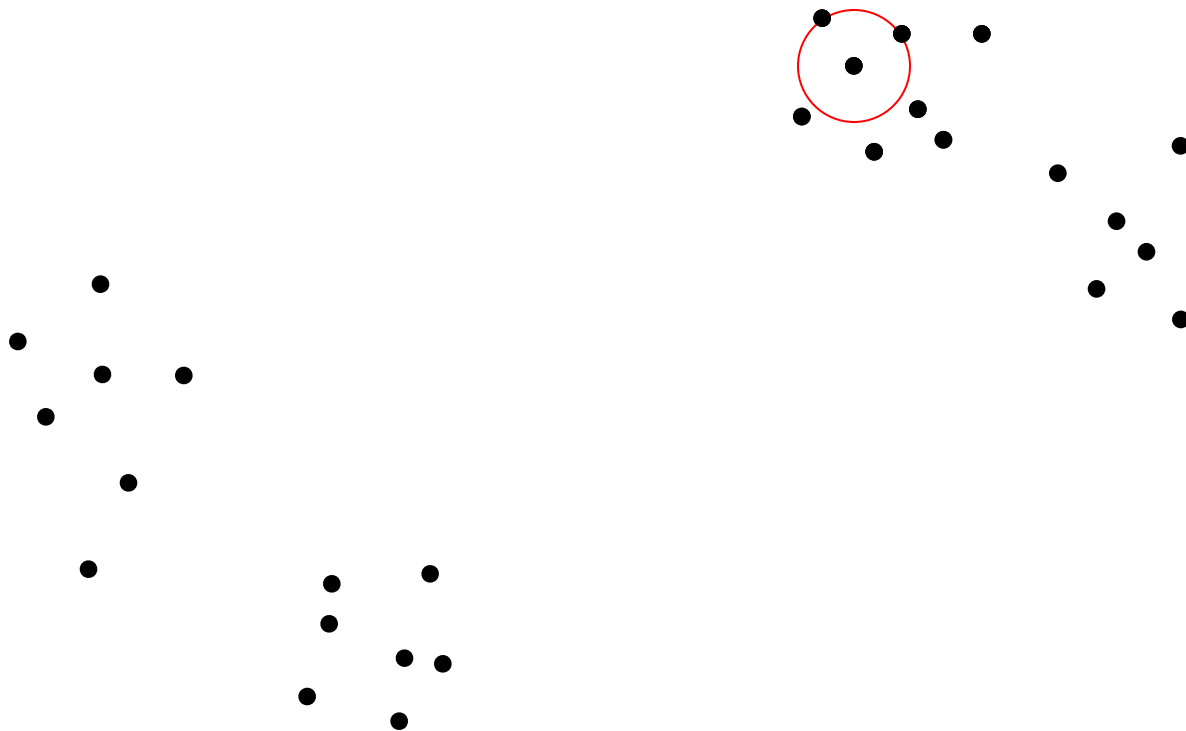
- Overlapping clusters



# Complications

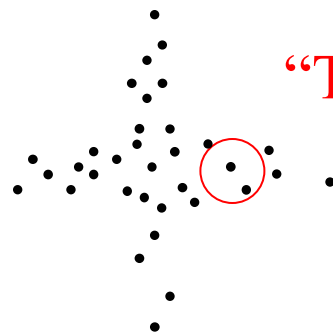
- How many clusters?

“This is a blicket.”

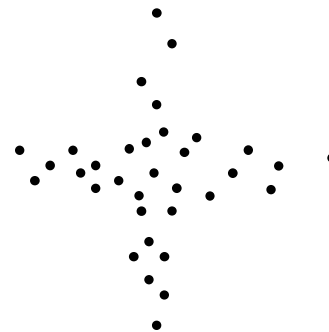


# Complications

- Clusters that are not simple blobs



“This is a blicket.”

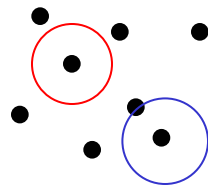




# Complications

- Concept labels inconsistent with clusters

“This is a blicket.”



“This is a gazzer.”

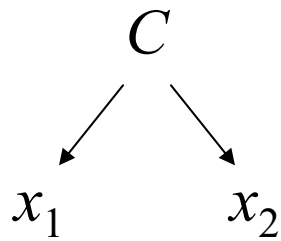


# Simple model of concept learning

- Can infer a concept from just one positive example if:
  - Assume concepts refer to clusters in the world.
  - Observe lots of unlabeled data, in order to identify clusters.
- How do we identify the clusters?
  - With no labeled data (“unsupervised learning”)
  - With sparsely labeled data (“semi-supervised learning”)

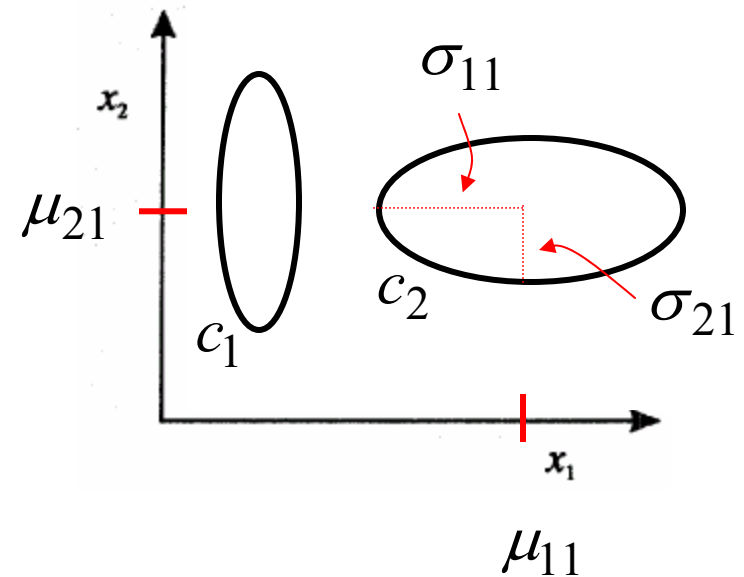
# Unsupervised clustering with probabilistic models

- Assume a simple parametric probabilistic model for clusters, e.g., Gaussian.



$$p(x | c_j) = p(x_1 | c_j) \times p(x_2 | c_j)$$

$$p(x_i | c_j) \propto e^{-\frac{(x_i - \mu_{ij})^2}{2\sigma_{ij}^2}}$$



# Unsupervised clustering with probabilistic models

- Assume a simple parametric probabilistic model for clusters, e.g., Gaussian.
- Two ways to characterize  $j$ th cluster:
  - Parameters:  $\mu_{ij}, \sigma_{ij}$
  - Assignments:  $z_j^{(k)} = 1$  if  $k$ th point belongs to cluster  $j$ , else 0.

# Unsupervised clustering with probabilistic models

- Chicken-and-egg problem:
  - Given assignments we could solve for maximum likelihood parameters:

$$\mu_{ij} = \frac{\sum_k z_j^{(k)} x_i^{(k)}}{\sum_k z_j^{(k)}} \quad \sigma^2_{ij} = \frac{\sum_k z_j^{(k)} \left( x_i^{(k)} - \mu_{ij} \right)^2}{\sum_k z_j^{(k)}}$$

# Unsupervised clustering with probabilistic models

- Chicken-and-egg problem:
  - Given parameters we could solve for assignments  $z_j^{(k)}$ :

$$z_j^{(k)} = \begin{cases} 1, & j = \arg \max_{j'} p(c_{j'} | \mathbf{x}^{(k)}) \\ 0, & \text{else} \end{cases}$$

$$p(c_j | \mathbf{x}^{(k)}) \propto p(\mathbf{x}^{(k)} | c_j) p(c_j)$$

$$\prod_i \frac{1}{\sqrt{2\pi}\sigma_{ij}} e^{-\frac{(x_i^{(k)} - \mu_{ij})^2}{2\sigma_{ij}^2}} p(c_j)$$

Solve for “base rate” parameters:

$$p(c_j) = \sum_k z_j^{(k)}$$

# Alternating optimization algorithm

0. Guess initial parameter values.
1. Given parameter estimates, solve for maximum a posteriori assignments  $z_j^{(k)}$ :

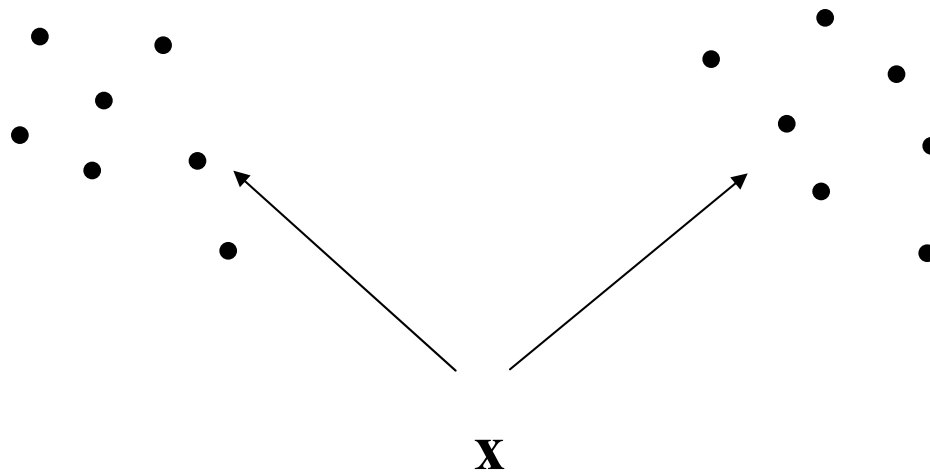
$$p(c_j | \mathbf{x}^{(k)}) \propto \prod_i \frac{1}{\sqrt{2\pi}\sigma_{ij}} e^{-(x_i^{(k)} - \mu_{ij})^2 / (2\sigma_{ij}^2)} p(c_j) \quad z_j^{(k)} = \begin{cases} 1, & j = \arg \max_{j'} p(c_{j'} | \mathbf{x}^{(k)}) \\ 0, & \text{else} \end{cases}$$

2. Given assignments  $z_j^{(k)}$ , solve for maximum likelihood parameter estimates:

$$\mu_{ij} = \frac{\sum_k z_j^{(k)} x_i^{(k)}}{\sum_k z_j^{(k)}} \quad \sigma_{ij}^2 = \frac{\sum_k z_j^{(k)} (x_i^{(k)} - \mu_{ij})^2}{\sum_k z_j^{(k)}} \quad p(c_j) = \frac{\sum_k z_j^{(k)}}{k}$$

3. Go to step 1.

# Alternating optimization algorithm

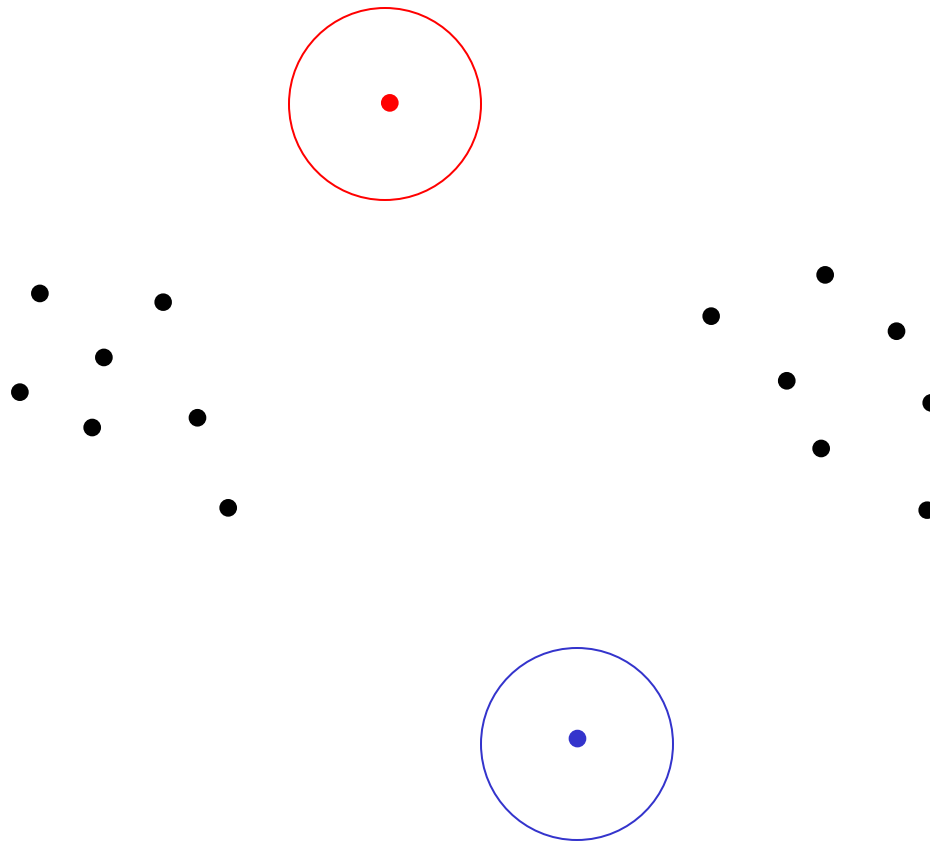


$\mathbf{z}$ : assignments to cluster  
 $\mu, \sigma, p(c_j)$ : cluster parameters

[For simplicity, assume  $\sigma, p(c_j)$  fixed.]

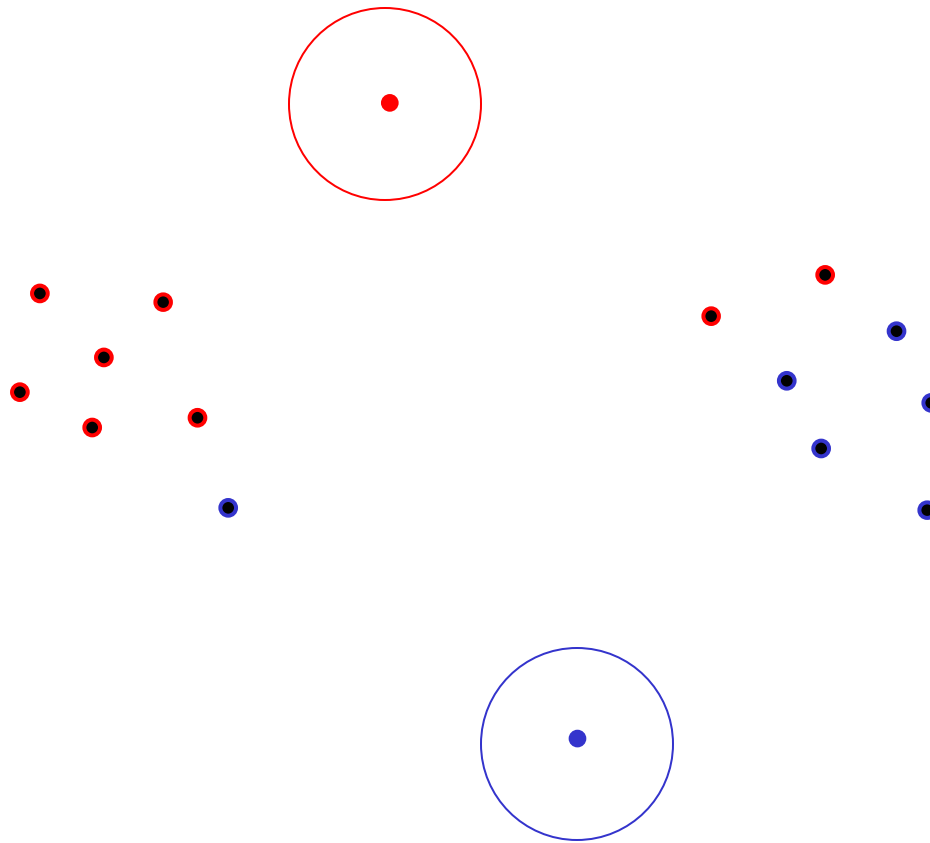


# Alternating optimization algorithm



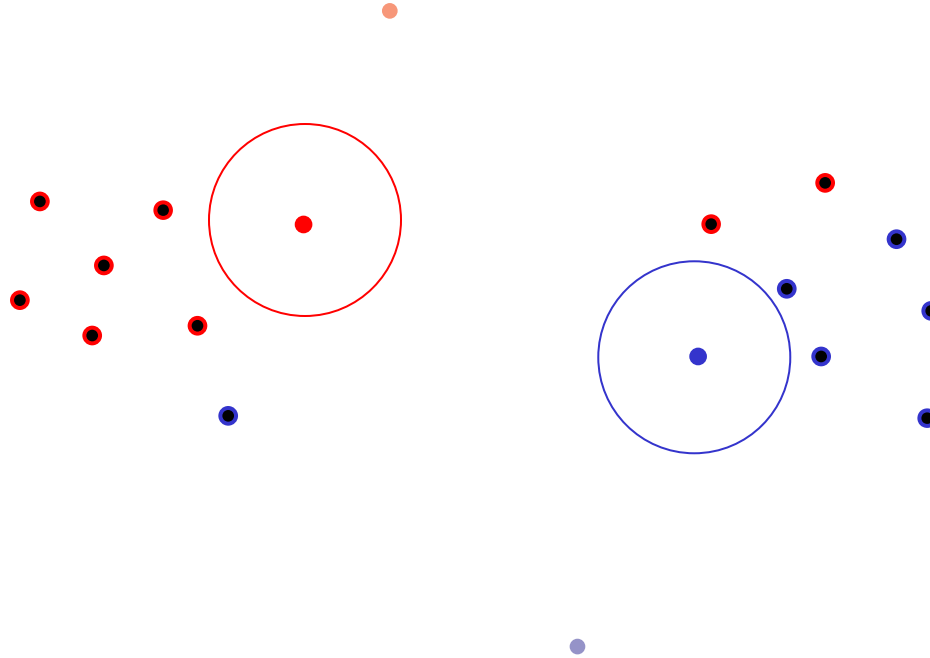
Step 0: initial parameter values

# Alternating optimization algorithm



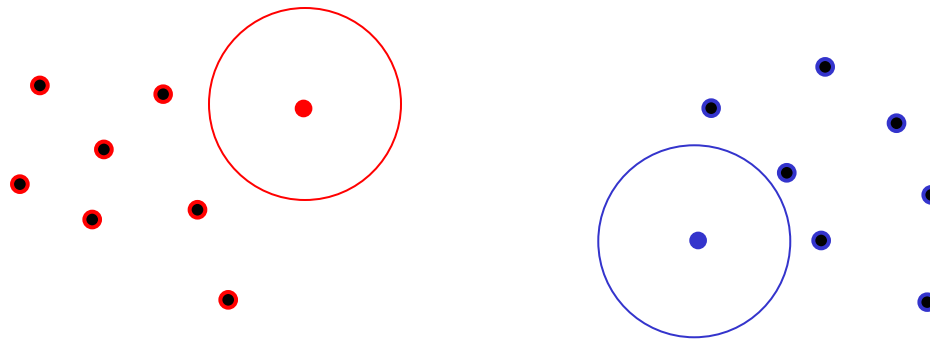
Step 1: update assignments

# Alternating optimization algorithm



Step 2: update parameters

# Alternating optimization algorithm



Step 1: update assignments

# Alternating optimization algorithm



Step 2: update parameters

# Alternating optimization algorithm

0. Guess initial parameter values.
1. Given parameter estimates, solve for maximum a posteriori assignments  $z_j^{(k)}$ :

$$p(c_j | \mathbf{x}^{(k)}) \propto \prod_i \frac{1}{\sqrt{2\pi}\sigma_{ij}} e^{-(x_i^{(k)} - \mu_{ij})^2 / (2\sigma_{ij}^2)} p(c_j) \quad z_j^{(k)} = \begin{cases} 1, & j = \arg \max_{j'} p(c_{j'} | \mathbf{x}^{(k)}) \\ 0, & \text{else} \end{cases}$$

2. Given assignments  $z_j^{(k)}$ , solve for maximum likelihood parameter estimates:

Why hard assignments?

$$\mu_{ij} = \frac{\sum_k z_j^{(k)} x_i^{(k)}}{\sum_k z_j^{(k)}} \quad \sigma_{ij}^2 = \frac{\sum_k z_j^{(k)} (x_i^{(k)} - \mu_{ij})^2}{\sum_k z_j^{(k)}} \quad p(c_j) = \frac{\sum_k z_j^{(k)}}{k}$$

3. Go to step 1.

# EM algorithm

0. Guess initial parameter values  $\theta = \{\mu, \sigma, p(c_j)\}$ .
1. **“Expectation” step:** Given parameter estimates, compute expected values of assignments  $z_j^{(k)}$

$$h_j^{(k)} \quad p(c_j | \mathbf{x}^{(k)}; \theta) \propto \prod_i \frac{1}{\sqrt{2\pi}\sigma_{ij}} e^{-(x_i^{(k)} - \mu_{ij})^2 / (2\sigma_{ij}^2)} p(c_j)$$

2. **“Maximization” step:** Given expected assignments, solve for maximum likelihood parameter estimates:

$$\mu_{ij} = \frac{\sum_k h_j^{(k)} x_i^{(k)}}{\sum_k h_j^{(k)}} \quad \sigma_{ij}^2 = \frac{\sum_k h_j^{(k)} (x_i^{(k)} - \mu_{ij})^2}{\sum_k h_j^{(k)}} \quad p(c_j) = \frac{\sum_k h_j^{(k)}}{k}$$

# What EM is really about

- Define a single probabilistic model for the whole data set:

$$p(\mathbf{X} | \theta) = \prod_k p(\mathbf{x}^{(k)} | \theta)$$

$$\prod_k \sum_j p(\mathbf{x}^{(k)} | c_j; \theta) p(c_j; \theta)$$

“mixture  
model”

$$\prod_k \sum_j \prod_i \frac{1}{\sqrt{2\pi\sigma_{ij}}} e^{-\frac{(x_i^{(k)} - \mu_{ij})^2}{2\sigma_{ij}^2}} p(c_j)$$



# What EM is really about

- Define a single probabilistic model for the whole data set:

$$p(\mathbf{X} | \theta) = \prod_k p(\mathbf{x}^{(k)} | \theta)$$

$$\prod_k \sum_j p(\mathbf{x}^{(k)} | c_j; \theta) p(c_j; \theta)$$

“mixture  
model”

$$\prod_k \sum_j \prod_i \frac{1}{\sqrt{2\pi\sigma_{ij}}} e^{-\frac{(x_i^{(k)} - \mu_{ij})^2}{2\sigma_{ij}^2}} p(c_j)$$

- How do we maximize w.r.t.  $\theta$ ?

# What EM is really about

- Maximization would be simpler if we introduced new labeling variables  $\mathbf{Z} = \{z_j^{(k)}\}$ :

$$p(\mathbf{X}, \mathbf{Z} | \theta) = \prod_k \prod_j \left( p(\mathbf{x}^{(k)} | c_j; \theta) p(c_j; \theta) \right)^{z_j^{(k)}}$$

$$\begin{aligned} \log p(\mathbf{X}, \mathbf{Z} | \theta) &= \sum_k \sum_j z_j^{(k)} \sum_i \log \left( p(x_i^{(k)} | c_j; \theta) p(c_j; \theta) \right) \\ &= - \sum_k \sum_j z_j^{(k)} \sum_i (x_i^{(k)} - \mu_{ij})^2 / (2\sigma_{ij}^2) + \log p(c_j) \end{aligned}$$

- Problem: we don't know  $\mathbf{Z} = \{z_j^{(k)}\}$ !

# What EM is really about

- Maximization expected value of the “complete data” loglikelihood,  $\log p(\mathbf{X}, \mathbf{Z} | \theta)$ :
  - **E-step:** Compute expectation

$$Q(\theta | \theta^{(t)}) = \sum_{\mathbf{Z}} p(\mathbf{Z} | \mathbf{X}, \theta^{(t)}) \log p(\mathbf{X}, \mathbf{Z} | \theta)$$

- **M-step:** Maximize

$$\theta^{(t+1)} = \arg \max_{\theta} Q(\theta | \theta^{(t)})$$