

Mixture Densities

$$p(\mathbf{x}) = \sum_{i=1}^k p(\mathbf{x} | \mathbf{G}_i) P(\mathbf{G}_i)$$

where \mathbf{G}_i the components / groups / clusters,
 $P(\mathbf{G}_i)$ mixture proportions (priors),
 $p(\mathbf{x} | \mathbf{G}_i)$ component densities

Gaussian mixture where $p(\mathbf{x} | \mathbf{G}_i) \sim \mathcal{N}(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$

parameters $\Phi = \{P(\mathbf{G}_i), \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i\}_{i=1}^k$

unlabeled sample $\mathbf{X} = \{\mathbf{x}^t\}_t$ (unsupervised learning)

Classes vs. Clusters

- Supervised: $\mathbf{x} = \{ \mathbf{x}^t, \mathbf{r}^t \}_t$
- Classes $C_i \ i=1, \dots, K$

$$p(\mathbf{x}) = \sum_{i=1}^K p(\mathbf{x} | C_i) P(C_i)$$

where $p(\mathbf{x} | C_i) \sim N(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$

- $\Phi = \{P(C_i), \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i\}_{i=1}^K$

$$\hat{P}(C_i) = \frac{\sum_t r_i^t}{N} \quad \mathbf{m}_i = \frac{\sum_t r_i^t \mathbf{x}^t}{\sum_t r_i^t}$$

$$\mathbf{S}_i = \frac{\sum_t r_i^t (\mathbf{x}^t - \mathbf{m}_i)(\mathbf{x}^t - \mathbf{m}_i)^T}{\sum_t r_i^t}$$

- Unsupervised: $\mathbf{x} = \{ \mathbf{x}^t \}_t$
- Clusters $G_i \ i=1, \dots, k$

$$p(\mathbf{x}) = \sum_{i=1}^k p(\mathbf{x} | G_i) P(G_i)$$

where $p(\mathbf{x} | G_i) \sim N(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$

- $\Phi = \{P(G_i), \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i\}_{i=1}^k$

Labels, \mathbf{r}^t ?

Expectation–Maximization (EM)

- Log likelihood with a mixture model

$$\begin{aligned} \mathcal{L}(\Phi | \mathbf{X}) &= \log \prod_t p(\mathbf{x}^t | \Phi) \\ &= \sum_t \log \sum_{i=1}^k p(\mathbf{x}^t | \mathbf{G}_i) P(\mathbf{G}_i) \end{aligned}$$

- Assume hidden variables \mathbf{z} , which when known, make optimization much simpler
- Complete likelihood, $\mathcal{L}_c(\Phi | \mathbf{X}, \mathbf{Z})$, in terms of \mathbf{x} and \mathbf{z}
- Incomplete likelihood, $\mathcal{L}(\Phi | \mathbf{X})$, in terms of \mathbf{x}

E- and M-steps

- Iterate the two steps

1. E-step: Estimate z given X and current Φ
2. M-step: Find new Φ' given z , X , and old Φ .

$$\text{E-step} : Q(\Phi | \Phi^l) = E[L_c(\Phi | X, Z) | X, \Phi^l]$$

$$\text{M-step} : \Phi^{l+1} = \arg \max_{\Phi} Q(\Phi | \Phi^l)$$

An increase in Q increases incomplete likelihood

$$L(\Phi^{l+1} | X) \geq L(\Phi^l | X)$$

EM in Gaussian Mixtures

- $z_i^t = 1$ if \mathbf{x}^t belongs to G_i , 0 otherwise (labels r^t_i of supervised learning); assume $p(\mathbf{x}|G_i) \sim \mathcal{N}(\mu_i, \Sigma_i)$

- E-step:
$$E[z_i^t | \mathcal{X}, \Phi^l] = \frac{p(\mathbf{x}^t | G_i, \Phi^l) P(G_i)}{\sum_j p(\mathbf{x}^t | G_j, \Phi^l) P(G_j)}$$
$$= P(G_i | \mathbf{x}^t, \Phi^l) \equiv h_i^t$$

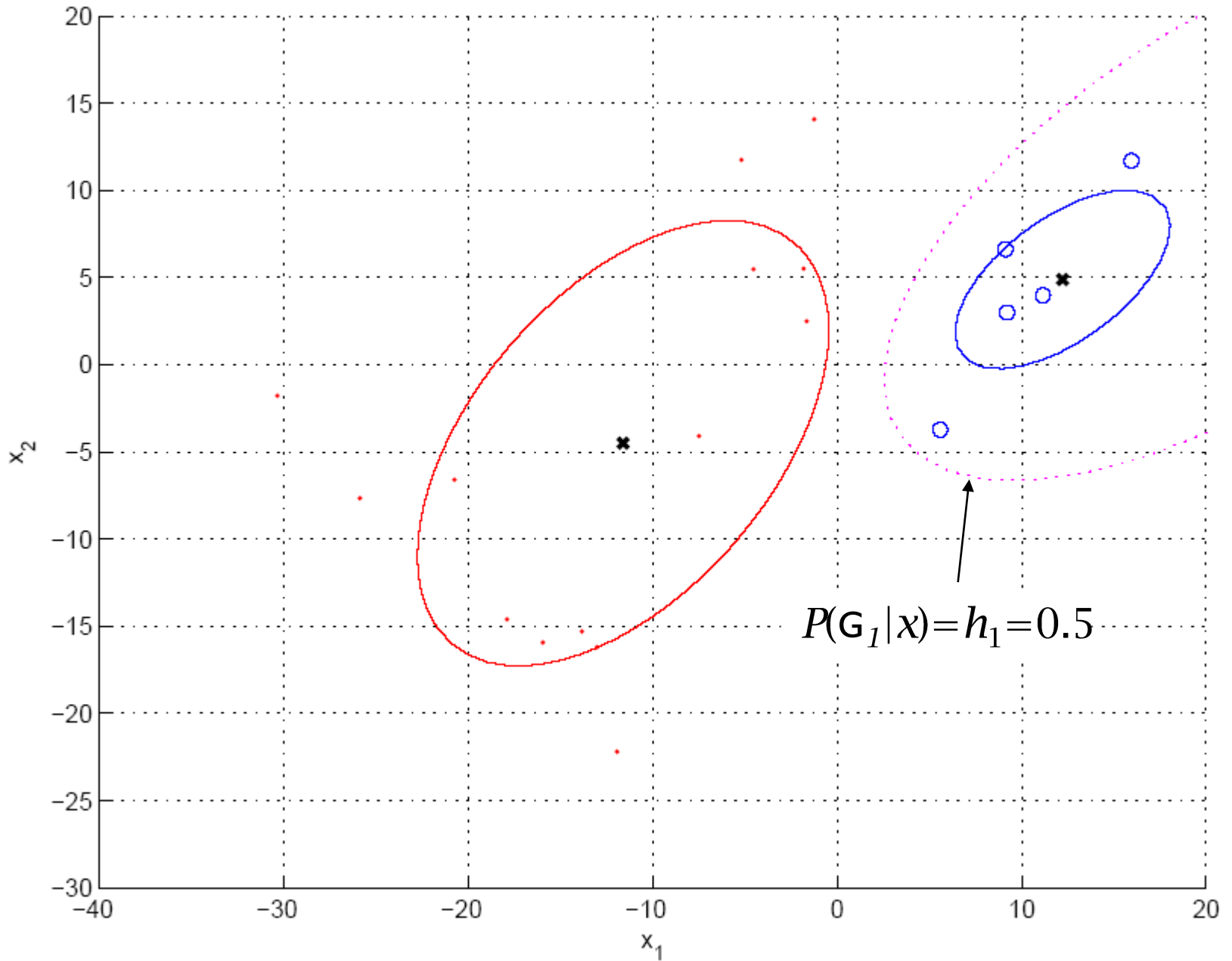
- M-step:

$$P(G_i) = \frac{\sum_t h_i^t}{N} \quad \mathbf{m}_i^{l+1} = \frac{\sum_t h_i^t \mathbf{x}^t}{\sum_t h_i^t}$$

*Use estimated labels
in place of
unknown labels*

$$\mathbf{S}_i^{l+1} = \frac{\sum_t h_i^t (\mathbf{x}^t - \mathbf{m}_i^{l+1})(\mathbf{x}^t - \mathbf{m}_i^{l+1})^T}{\sum_t h_i^t}$$

EM solution



Mixtures of Latent Variable Models

- Regularize clusters
 1. Assume shared/diagonal covariance matrices
 2. Use PCA/FA to decrease dimensionality: Mixtures of PCA/FA

$$p(\mathbf{x}_t | \mathbf{G}_i) = \mathbf{N}(\mathbf{m}_i, \mathbf{V}_i \mathbf{V}_i^T + \psi_i)$$

Can use EM to learn \mathbf{V}_i (Ghahramani and Hinton, 1997; Tipping and Bishop, 1999)



After Clustering

- Dimensionality reduction methods find correlations between features and group features
- Clustering methods find similarities between instances and group instances
- Allows knowledge extraction through
 - number of clusters,
 - prior probabilities,
 - cluster parameters, i.e., center, range of features.Example: CRM, customer segmentation

Clustering as Preprocessing

- Estimated group labels h_j (soft) or b_j (hard) may be seen as the dimensions of a new k dimensional space, where we can then learn our discriminant or regressor.
- **Local representation** (only one b_j is 1, all others are 0; only few h_j are nonzero) vs **Distributed representation** (After PCA; all z_j are nonzero)

Mixture of Mixtures

- In classification, the input comes from a mixture of classes (supervised).
- If each class is also a mixture, e.g., of Gaussians, (unsupervised), we have a mixture of mixtures:

$$p(\mathbf{x} | \mathbf{C}_i) = \sum_{j=1}^{k_i} p(\mathbf{x} | \mathbf{G}_{ij})P(\mathbf{G}_{ij})$$

$$p(\mathbf{x}) = \sum_{i=1}^K p(\mathbf{x} | \mathbf{C}_i)P(\mathbf{C}_i)$$