

LTAT.02.004 MACHINE LEARNING II

**Missing roadmap to  
Expectation-maximisation algorithm**

Sven Laur  
University of Tartu

## Hard clustering versus soft clustering

	Hard clustering	Soft clustering
Maximisation goal	$p[\mathbf{x}_1, \dots, \mathbf{x}_n   \mathbf{z}, \Theta]$	$p[\Theta] \cdot \sum_{\mathbf{z}} p[\mathbf{x}_1, \dots, \mathbf{x}_n, \mathbf{z}   \Theta]$
Optimisation method	Two-step maximisation algorithm	
Tactical objective	$F(\mathbf{z}, \Theta)$	$F(\mathbf{q}, \Theta)$
Mixture proportions	Ignored by design	Core of the model
Cluster labels	Search goal	Integrated out

## Desired properties of the tactical objective

**Property I.** Let  $q_{\Theta}(\cdot)$  be the optimal probability distribution for label vectors  $\mathbf{z}$  for fixed model parameters  $\Theta$ . Then the tactical objective coincides with the actual objective:

$$F(q_{\Theta}, \Theta) = \log p[\Theta | \mathbf{x}_1, \dots, \mathbf{x}_n] \ .$$

**Property II.** For fixed model parameters  $\Theta$  the optimal probability distribution can be found as the posterior probability of label vectors:

$$q_{\Theta}(\mathbf{z}) = p[\mathbf{z} | \mathbf{x}_1, \dots, \mathbf{x}_n, \Theta] \ .$$

### Rationale

- ▷ The first property is essential for obtaining a local maxima.
- ▷ The second property is needed to justify the practical algorithm.

## The derivation of the tactical objective

Let  $q(\mathbf{z})$  be an arbitrary probability distribution over label vectors  $\mathbf{z}$ . Then tautology together with Jensen's inequality assures

$$\begin{aligned}\log p[\Theta | \mathbf{x}_1, \dots, \mathbf{x}_n] &= \log \left( \sum_{\mathbf{z}} q(\mathbf{z}) \cdot \frac{p[\Theta, \mathbf{z} | \mathbf{x}_1, \dots, \mathbf{x}_n]}{q(\mathbf{z})} \right) \\ &\geq \sum_{\mathbf{z}} q(\mathbf{z}) \cdot \log \left( \frac{p[\Theta, \mathbf{z} | \mathbf{x}_1, \dots, \mathbf{x}_n]}{q(\mathbf{z})} \right) = F(q, \Theta)\end{aligned}$$

For the probability assignment  $q_{\Theta}(\mathbf{z}) = p[\mathbf{z} | \mathbf{x}_1, \dots, \mathbf{x}_n, \Theta]$  we get

$$\begin{aligned}F(q_{\Theta}, \Theta) &= \sum_{\mathbf{z}} q_{\Theta}(\mathbf{z}) \cdot \log \left( \frac{p[\Theta, \mathbf{z} | \mathbf{x}_1, \dots, \mathbf{x}_n]}{p[\mathbf{z} | \mathbf{x}_1, \dots, \mathbf{x}_n, \Theta]} \right) \\ &= \sum_{\mathbf{z}} q_{\Theta}(\mathbf{z}) \cdot \log (p[\Theta | \mathbf{x}_1, \dots, \mathbf{x}_n]) = \log p[\Theta | \mathbf{x}_1, \dots, \mathbf{x}_n] .\end{aligned}$$

## Tactical objective as a linearisation

The expectation-maximisation algorithm can be viewed as follows:

- ▷ Guess model parameters  $\Theta^{(i)}$ .
- ▷ Compute probability assignments  $q_{\Theta^{(i)}}(z) = p[z|x_1, \dots, x_n, \Theta^{(i)}]$ .
- ▷ Approximate  $p[z|x_1, \dots, x_n, \Theta]$  with a linearisation  $F(q_{\Theta^{(i)}}, \Theta)$ .
- ▷ Fix a new guess  $\Theta^{(i+1)}$  that maximises  $F(q_{\Theta^{(i)}}, \Theta)$ .

As the actual value and linearisation can be expressed as

$$\log p[\Theta|x_1, \dots, x_n] = \sum_z q_{\Theta^{(i)}}(z) \cdot \log \left( \frac{p[\Theta, z|x_1, \dots, x_n]}{p[z|x_1, \dots, x_n, \Theta]} \right)$$
$$F(q_{\Theta^{(i)}}, \Theta) = \sum_z q_{\Theta^{(i)}}(z) \cdot \log \left( \frac{p[\Theta, z|x_1, \dots, x_n]}{p[z|x_1, \dots, x_n, \Theta^{(i)}]} \right)$$

## Tactical objective as a linearisation

- ▷ Guess model parameters  $\Theta^{(i)}$ .
- ▷ Compute probability assignments  $q_{\Theta^{(i)}}(z) = p[z|\mathbf{x}_1, \dots, \mathbf{x}_n, \Theta^{(i)}]$ .
- ▷ Approximate  $p[z|\mathbf{x}_1, \dots, \mathbf{x}_n, \Theta]$  with the linear function  $F(q_{\Theta^{(i)}}, \Theta)$ .
- ▷ Fix a new guess  $\Theta^{(i+1)}$  that maximises  $F(q_{\Theta^{(i)}}, \Theta)$

Kullback-Leibler divergence between probability assignments for label vectors

$$D(q_{\Theta^{(i)}} || q_{\Theta}) = \sum_z q_{\Theta^{(i)}}(z) \cdot \log \left( \frac{p[z|\mathbf{x}_1, \dots, \mathbf{x}_n, \Theta^{(i)}]}{p[z|\mathbf{x}_1, \dots, \mathbf{x}_n, \Theta]} \right)$$

measures the linearisation error  $p[z|\mathbf{x}_1, \dots, \mathbf{x}_n, \Theta^{(i)}] - F(q_{\Theta^{(i)}}, \Theta)$ .

## Simplification of the lower bound

**Observation I.** The distribution  $q_{\Theta}$  decomposes into a product of posteriors:

$$q_{\Theta}(\mathbf{z}) = \prod_{i=1}^n p[z_i | \mathbf{x}_i, \Theta] \ .$$

**Observation II.** Let  $\mathbf{W}$  be matrix of weights  $w_{ij} = p[z_i = j | \mathbf{x}_i, \Theta^{(*)}]$ . The lower bound can be expressed only in terms of  $\mathbf{W}$  and  $\Theta$ :

$$\begin{aligned} F(q_{\Theta^{(*)}}, \Theta) &= \log p[\Theta] - \sum_{i=1}^n \sum_{j=1}^k w_{ij} \log w_{ij} \\ &\quad + \sum_{i=1}^n \sum_{j=1}^k w_{ij} \cdot \log \lambda_j + \sum_{i=1}^n \sum_{j=1}^k w_{ij} \cdot \log (p[\mathbf{x}_i | \Theta_j]) \ . \end{aligned}$$

## Parameter optimisation

Hard clustering finds model parameters of the  $j$ th cluster by solving

$$\sum_{i=1}^n [z_i = j] \cdot \log(p[\mathbf{x}_i | \Theta_j]) \rightarrow \max \quad .$$

Soft clustering finds model parameters of the  $j$ th cluster by solving

$$\sum_{i=1}^n \sum_{j=1}^k w_{ij} \cdot \log(p[\mathbf{x}_i | \Theta_j]) \rightarrow \max \quad .$$

and additionally updates mixture proportions  $\lambda_1, \dots, \lambda_k$ .