

## 8. Supervised learning of HMMs: Empirical risk minimisation

Given: i.i.d. training data  $\tilde{T} = \{(x_i, s_i) \mid x_i \in \mathcal{X}^n, s_i \in K^n, i=1, \dots, m\}$   
and the loss function  $\ell(s, s') = \llbracket s \neq s' \rrbracket$

Recall that the optimal predictor  $h_u: \mathcal{X}^n \rightarrow K^n$  for 0/1 loss is

$$h_u(x) \in \operatorname{argmax}_{s \in K^n} p_{h_u}(x, s) = \operatorname{argmax}_{s \in K^n} \langle \Phi(x, s), u \rangle$$

Learning the model parameters  $u$  by empirical risk minimisation

$$\frac{1}{m} \sum_{(x, s) \in \tilde{T}} \llbracket s \neq h_u(x) \rrbracket \rightarrow \min_u$$

is not tractable because the objective function is piece-wise const.

Special case: Suppose  $\exists u^*$  s.t. the empirical risk on  $\tilde{T}$  is zero.  
Conditions for such  $u^*$  are

$$s \in \operatorname{argmax}_{s' \in K^n} p_{u^*}(x, s') \quad \forall (x, s) \in \tilde{T}$$

or, equivalently

$$\langle \Phi(x, s), u^* \rangle > \langle \Phi(x, s'), u^* \rangle \quad \forall s' \neq s, \forall (x, s) \in \tilde{T}$$

This is a system of linear inequalities  $\Rightarrow$  perceptron alg.  
Start with  $u=0$  and iterate

$$(1) \text{ find } \tilde{s} = \operatorname{argmax}_{s' \in K^n} \langle \Phi(x, s'), u \rangle \quad \forall (x, s) \in \tilde{T}.$$

This can be done by dynamic programming (Sec. 4)

$$(2) \text{ if for some } (x, s) \in \tilde{T}: \tilde{s} \neq s, \text{ update } u \text{ by}$$

$$u \rightarrow u + \Phi(x, s) - \Phi(x, \tilde{s})$$

The algorithm converges to a solution  $u^*$  in a finite number of steps (provided it exists).

General case: overcome the intractability by replacing the loss (as a function of  $u$ ) by a convex upper bound.  
E.g. „margin rescaling“ loss

$$\begin{aligned}\mathbb{E}[S \neq h_u(x)] &= \mathbb{E}[S \neq \operatorname{argmax}_{S' \in K^n} \langle \Phi(x, S'), u \rangle] \leq \\ &\leq \max_{S' \in K^n} \left\{ \mathbb{E}[S \neq S'] + \langle \Phi(x, S') - \Phi(x, S), u \rangle \right\}\end{aligned}$$

The empirical risk minimisation for this loss reads

$$\frac{1}{m} \sum_i^m \max_{(x, s) \in T} \max_{S' \in K^n} \left\{ \mathbb{E}[S \neq S'] + \langle \Phi(x, S') - \Phi(x, S), u \rangle \right\} \rightarrow \min_u$$

Solve it by (stochastic) subgradient descent or cutting plane algorithm or ... . The inner optimisation tasks  $\max_{S' \in K^n} \{\dots\}$  are solved by an algorithm like the one in Sec. 4

Remark 1 This approach is designated as „structured output SVM“ and can be generalised for more complex losses as e.g. the Hamming distance.

## 9. Unsupervised learning: EM algorithm for HMMs

Given: i.i.d. training data  $T = \{x^j \in \mathcal{X}^n \mid j=1, \dots, m\}$

The ML estimator reads  $u^* \in \arg \max_u \frac{1}{|T|} \sum_{x \in T} \log \sum_{s \in K^n} p_u(x, s)$

Recall the EM algorithm

$$L(u) = \frac{1}{m} \sum_{x \in T} \log \sum_{s \in K^n} \frac{\alpha(s|x)}{\alpha(s|x)} p_u(x, s)$$

for any  $\alpha(s|x) \geq 0$  s.t.  $\sum_{s \in K^n} \alpha(s|x) = 1 \quad \forall x \in T$ .

Using concavity of  $\log$ , we get the lower bound

$$\underline{L_B}(u, \alpha) = \frac{1}{m} \sum_{x \in T} \sum_{s \in K^n} \alpha(s|x) \log \frac{p_u(x, s)}{\alpha(s|x)}$$

or equivalently

$$L_B(u, \alpha) = \mathbb{E}_T [\log p_u(x) - D_{KL}(\alpha(s|x) \parallel p_u(s|x))]$$

The EM algorithm maximizes  $L_B(u, \alpha)$  by block-coordinate ascent w.r.t.  $\alpha$  and  $u$ . Start with some  $u^{(0)}$ .

E-step maximise  $L_B(u^{(t)}, \alpha)$  w.r.t.  $\alpha \Rightarrow$

$$\alpha^{(t+1)}(s|x) = p_{u^{(t)}}(s|x) \quad \forall s \in K^n, \forall x \in T$$

M-step maximise  $L_B(u, \alpha^{(t)})$  w.r.t.  $u \Rightarrow$

$$u^{(t+1)} \in \arg \max_u \frac{1}{m} \sum_{x \in T} \sum_{s \in K^n} \alpha^{(t)}(s|x) \log p_u(x, s)$$

Let us analyse the M-step for HMMs. The objective is

$$\frac{1}{m} \sum_{x \in T} \sum_{s \in K^n} \alpha^{(t)}(s|x) \langle \Phi(x, s), u \rangle - \log Z(u) \rightarrow \max_u$$

By denoting

$$\Psi = \frac{1}{m} \sum_{x \in T} \sum_{s \in K^n} \alpha^{(t)}(s|x) \Phi(x, s)$$

we get

$$\langle \Psi, u \rangle - \log Z(u) \rightarrow \max_u.$$

This is equivalent to the supervised learning task in Sec. 7. We know how to solve it, provided we can compute the components of  $\Psi$ .

### Computing $\Psi$

For each  $x \in T$  we want to compute

$$\Psi(x) = \sum_{s \in K^n} \alpha^{(t)}(s|x) \Phi(x, s) = \mathbb{E}_{p_{u(t)}(s|x)} [\Phi(x, s)]$$

i.e. we have to compute the posterior pairwise marginals  $p(s_{i-1}, s_i | x)$   $\forall i=2, \dots, n$  and  $s_{i-1}, s_i \in K$ .

This can be done by an algorithm similar to the one discussed in Sec. 5. The components of  $\Psi$  are then obtained by averaging the components of  $\Psi(x)$  over all  $x \in T$ , i.e.  $\Psi = \mathbb{E}_p [\Psi(x)]$ .

### Theorem 1 (w/o proof)

The sequence  $L(u^{(t)})$  is monotonously increasing and the sequence  $\alpha^{(t)}$  is convergent.

Remark 1 The EM algorithm for HMMs is known as Baum-Welch algorithm.