

6. Representing an HMM as an exponential family

According to Def. 16 in sec. 1, the joint p.d. for a Markov model on a chain can be written as

$$p(s) = p(s_1, \dots, s_n) = \prod_{i=2}^n g_i(s_{i-1}, s_i)$$

To allow arbitrary non-negative g -s, we introduce a normalising factor Z . If, in addition, all g -s are strictly positive, we may write

$$p(s) = \frac{1}{Z(u)} \exp \left[\sum_{i=2}^n u_i(s_{i-1}, s_i) \right]$$

Remark 1 Notice that the normalising factor Z depends on the model parameters u , i.e. Z is a function of u . It is defined by

$$Z(u) = \sum_{s \in K^n} \exp \left[\sum_{i=2}^n u_i(s_{i-1}, s_i) \right]$$

and can be computed by the algorithm described in sec. 3

Remark 2 The formula

$$p(s) = p(s_1, \dots, s_n) = \frac{p(s_1, s_2) \cdot p(s_2, s_3) \cdot \dots \cdot p(s_{n-1}, s_n)}{p(s_2) \cdot p(s_3) \cdot \dots \cdot p(s_{n-1})}$$

for a Markov chain (see sec. 1) provides an easy way to compute the u -s given the pairwise marginal prob's.

Remark 3 The factors g_i resp. the potentials u_i define a Markov model uniquely but are themselves not unique

$$u_2(s_1, s_2) + \dots + \underbrace{u_i(s_{i-1}, s_i) + V_i(s_i)}_{\tilde{u}_i} + \underbrace{u_{i+1}(s_i, s_{i+1}) - V_i(s_i)}_{\tilde{u}_{i+1}} + \dots + u_n(s_{n-1}, s_n)$$

To emphasize the structure in which the model depends on its parameters u , we may write

$$p(s) = \frac{1}{Z(u)} \exp \left[\sum_{i=2}^n u_i (s_{i-1}, s_i) \right] = \frac{1}{Z(\vec{u})} \exp \langle \vec{\Psi}(s), \vec{u} \rangle$$

The vectors $\vec{\Psi}(s), \vec{u}$ are elements of $\mathbb{R}^{(n-1)K^2}$

$\vec{\Psi}(s)$ is a binary valued indicator vector with components defined as follows

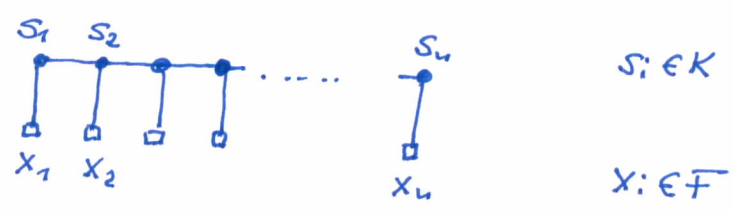
$$\Psi_{ikk'}(s) = \Psi_{ikk'}(s_{i-1}, s_i) = \delta(s_{i-1}, k) \cdot \delta(s_i, k')$$

$i=2, \dots, n, \quad k, k' \in K$

where δ denotes the Kronecker delta.

\vec{u} is the parameter vector with components $u_{ikk'} = u_i(k, k')$

A Hidden Markov model (on a chain) is defined on the graph



Its joint p.d. can be written as

$$\begin{aligned} \varphi(x, s) &= \prod_{i=2}^n g_i(s_{i-1}, s_i) \prod_{i=1}^n \tilde{g}_i(x_i, s_i) \\ &= \frac{1}{Z(u, \tilde{u})} \exp \left[\sum_{i=2}^n u_i (s_{i-1}, s_i) + \sum_{i=1}^n \tilde{u}_i (x_i, s_i) \right] \\ &= \frac{1}{Z(\vec{u})} \exp \langle \vec{\Psi}(x, s), \vec{u} \rangle \end{aligned}$$

7. Supervised learning, ML-estimator

Given: i.i.d. sample of pairs of sequences

$$\mathcal{T} = \{(x^j, s^j) \mid x^j \in F^n, s^j \in K^n, j = 1, \dots, l\}$$

Maximum likelihood estimator:

$$\begin{aligned} \vec{u}^* &= \operatorname{argmax}_{\vec{u}} \prod_{(x,s) \in \mathcal{T}} p_{\vec{u}}(x,s) \\ &= \operatorname{argmax}_{\vec{u}} \frac{1}{|\mathcal{T}|} \sum_{(x,s) \in \mathcal{T}} \log p_{\vec{u}}(x,s) \end{aligned}$$

i.e. find optimal

$$\vec{u}: \quad u_i(s_{i-1}, s_i), \quad \tilde{u}_i(x_i, s_i) \quad \text{or, equivalently}$$

$$p(s_{i-1}, s_i), \quad p(x_i, s_i)$$

Intuitive answer: \vec{u}^* is given by

$$p_{\vec{u}^*}(s_{i-1}, s_i) = \beta(s_{i-1}, s_i)$$

$$p_{\vec{u}^*}(x_i, s_i) = \beta(x_i, s_i)$$

where β -s denote frequencies of corresponding events in \mathcal{T} .

Let us prove that this is correct. The objective function for the learning task is

$$\begin{aligned} L(\vec{u}) &= \frac{1}{|\mathcal{T}|} \sum_{(x,s) \in \mathcal{T}} [\langle \Psi(x,s), \vec{u} \rangle - \log Z(\vec{u})] \\ &= \langle \vec{\Psi}, \vec{u} \rangle - \log Z(\vec{u}) \end{aligned}$$

where $\vec{\Psi}$ denotes the empirical mean of the random vector $\vec{\Phi}$:

$$\vec{\Psi} = \frac{1}{|\mathcal{T}|} \sum_{(x,s) \in \mathcal{T}} \vec{\Phi}(x,s)$$

The first term in $L(\vec{u})$ is linear and thus also concave. Let us prove that $\log Z(\vec{u})$ is convex.

$$\log Z(\vec{u}) = \log \sum_{x,s} \exp \langle \vec{\Phi}(x,s), \vec{u} \rangle$$

$$\nabla \log Z(\vec{u}) = \frac{1}{Z(\vec{u})} \sum_{x,s} \exp \langle \vec{\Phi}(x,s), \vec{u} \rangle \vec{\Phi}(x,s) \stackrel{!}{=} \mathbb{E}_{\vec{u}}(\vec{\Phi})$$

i.e. $\nabla \log Z(\vec{u}) \hat{=}$ expectation of the random vector $\vec{\Phi}$
components: pairwise marginal prob's
on edges

$$\begin{aligned} \nabla^2 \log Z(\vec{u}) &= \mathbb{E}_{\vec{u}}(\vec{\Phi} \otimes \vec{\Phi}) - \mathbb{E}_{\vec{u}}(\vec{\Phi}) \otimes \mathbb{E}_{\vec{u}}(\vec{\Phi}) \\ &= \mathbb{E}_{\vec{u}}[(\vec{\Phi} - \mathbb{E}_{\vec{u}}(\vec{\Phi})) \otimes (\vec{\Phi} - \mathbb{E}_{\vec{u}}(\vec{\Phi}))] \end{aligned}$$

i.e. $\nabla^2 \log Z(\vec{u}) \hat{=}$ covariance matrix of the random vector $\vec{\Phi} \Rightarrow$ symmetric and positive semidefinite

Lemma 1 The partition function $\log Z(\vec{u})$ of an HMM (with strictly positive p.d.) is convex in \vec{u} . \square

Hence, the objective function $L(\vec{u})$ is concave and has global maxima only. They are given by

$$\nabla L(\vec{u}) = \frac{1}{|\mathcal{T}|} \sum_{(x,s) \in \mathcal{T}} \vec{\Phi}(x,s) - \mathbb{E}_{\vec{u}^*}(\vec{\Phi}) = 0$$

But, the components of $E_{\vec{u}}(\vec{\Phi})$ are the pairwise marginals of the model $p_{\vec{u}}(x, s)$. The optimiser \vec{u}^* defines the model which has precisely the same pairwise marginals as the empirical marginal frequencies of \mathcal{T} .

Concavity of log-Likelihood $L(\vec{u})$ also ensures consistency of the estimator.

Theorem 1 (w/o proof) The maximum likelihood estimator for HMM-s is consistent, i.e.

$$P_{\vec{u}}(\|\vec{u}_*(\mathcal{T}_\epsilon) - \vec{u}\| > \epsilon) \xrightarrow{l \rightarrow \infty} 0$$

for every $\epsilon > 0$