## 6. Representing an HMM as an exponential family

Def. 1b, Sec. 1 $\Rightarrow$ the joint p.d. for a Markov chain model with strictly positive prob's can be written as

$$P(s) = p(s_1, \ldots, s_n) = \frac{1}{Z} \prod_{i=2}^{n} g_i(s_{i-1}, s_i) = \frac{1}{Z} \exp\left[\sum_{i=2}^{n} u_i(s_{i-1}, s_i)\right]$$

**Remark 1**  The factors $g_i$, resp. the potentials $u_i$ define the model uniquely. The inverse is not true.

**Remark 2**  The normalising factor $Z$ 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 an algorithm similar to the one described in Sec. 3

Denote:

- $\vec{\varphi_i}(s_i) \in \mathbb{R}^K$ a binary valued indicator vector, which denotes the state $s_i$ ("one out of $K$" coding), i.e.

$$\vec{\varphi_i}(s_i = k) = (0, \ldots, \underset{\underset{pos. k}{\uparrow}}{1}, \ldots 0)$$

- $U_i$ denotes the $K \times K$ matrix with values $u_i(s_{i-1}, s_i)$

Then, the joint p.d. can be written as

$$P(s) = \frac{1}{Z(u)} \exp \sum_{i=2}^{n} \langle \vec{\varphi}_{i-1}(s_{i-1}), U_i \cdot \vec{\varphi_i}(s_i)\rangle$$

$$= \frac{1}{Z(u)} \exp \sum_{i=2}^{n} \langle \varphi_i(s_{i-1}, s_i), U_i\rangle$$

where

$\Phi_i(s_{i-1}, s_i) = \vec{\Phi}_{i-1}(s_{i-1}) \otimes \vec{\Phi}_i(s_i)$ is a $K \times K$ binary valued indicator matrix and

$\langle \Phi, u \rangle = tr(\Phi^T u)$ denotes the Frobenius inner product.

Finally, denote $\Phi = (\Phi_2, \Phi_3, .., \Phi_u)$ and $u = (u_1, u_2, .., u_u)$.

The joinet p.d. of a Markov chain model can be written as

$$P(s) = \frac{1}{Z(u)} exp \langle \Phi(s), u \rangle$$

The joint p.d. of an HMM can be written as

$$P(s) = \frac{1}{Z(u)} exp \langle \Phi(x,s), u \rangle$$

by using similar notations.

## 7. Supervised learning, ML-estimator

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

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

estimate the model parameters of the HMM by the maximum likelihood estimator

$$u^* \in \underset{u}{\arg\max} \prod_{(x,s) \in \mathcal{T}} p_u(x,s)$$

$$= \underset{u}{\arg\max} \frac{1}{|\mathcal{T}|} \sum_{(x,s) \in \mathcal{T}}^{\prime} \log p_u(x,s)$$

i.e. find optimal $u_i^*(s_{i-1}, s_i)$, $\tilde{u}_i^*(x_i, s_i)$, or, equivalently, $p(s_{i-1}, s_i)$, $p(x_i, s_i)$

Intuitive answer $u^*$ is given by

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

where $\beta$-s denote frequencies of corresponding events in $\mathcal{T}$.
Let us prove correctness. Log-likelihood of $\mathcal{T}$ is

$$L(u) = \frac{1}{|\mathcal{T}|} \sum_{(x,s) \in \mathcal{T}}^{\prime} \left[ \langle \varphi(x,s), u \rangle - \log Z(u) \right]$$

$$= \langle \psi, u \rangle - \log Z(u),$$

where

$$\psi = \mathbb{E}_{\mathcal{T}}(\varphi) = \frac{1}{|\mathcal{T}|} \sum_{(x,s) \in \mathcal{T}}^{\prime} \varphi(x,s)$$

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

Proof

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

The components of $\mathbb{E}_u(\Phi)$ are the pairwise marginal prob's on the edges of the model.

$$\nabla^2 \log Z(u) = \mathbb{E}_u(\Phi \otimes \Phi) - \mathbb{E}_u(\Phi) \otimes \mathbb{E}_u(\Phi)$$

$$= \mathbb{E}_u\left[ (\Phi - \mathbb{E}_u(\Phi)) \otimes (\Phi - \mathbb{E}_u(\Phi)) \right]$$

The expectation of a positive semidefinite matrix is p.s.d. $\Rightarrow$ $\log Z(u)$ is convex.

The log likelihood $L(u)$ is concave as a consequence and has global maxima only, which are given by

$$\nabla L(u^*) = \frac{1}{|T|} \sum_{(x,s) \in T}' \Phi(x,s) - \mathbb{E}_{u^*}(\Phi)$$

$$= \mathbb{E}_T(\Phi) - \mathbb{E}_{u^*}(\Phi) = 0$$

Recall that the components of $\mathbb{E}_u(\Phi)$ are the pairwise marginal prob's of the model $p_u(x,s)$. Hence, the optimiser $u^*$ defines the model that has precisely the same pairwise marginal prob's as the empirical marginal frequencies of $T$.

The concavity of $L(u)$ also ensures the __consistency__ of the estimator

__Theorem 1__ (w/o proof)

The maximum likelihood estimator for HMMs is consistent, i.e.

$$\mathbb{P}_u \left( \| u^*(\gamma) - u \| > \varepsilon \right) \xrightarrow{|\gamma| \to \infty} 0$$

holds for every $\varepsilon > 0$.