## Ch I  Markov models on chains and acyclic graphs

### 1. Markov models on chains

#### LA  Definitions & basic properties

- Sequence $S = (S_1, \ldots, S_n)$ of $K$-valued random variables $s_i \in K$
- $K$ is a finite set, its elements are called states
- $p(s) = p(s_1, \ldots, s_n)$ is a joint probability distr. on $K^n$

W.l.o.g. we can write

$$p(s_1, \ldots, s_n) = p(s_n \mid s_1, \ldots, s_{n-1}) \, p(s_1, \ldots, s_{n-1})$$

$$= \ldots$$

$$= p(s_n \mid s_1, \ldots, s_{n-1}) \, p(s_{n-1} \mid s_1, \ldots, s_{n-2}) \cdots p(s_2 \mid s_1) \, p(s_1)$$

**Definition 1a**  A p.d. on $K^n$ is a Markov chain if

$$p(s) = p(s_1) \prod_{i=2}^{n} p(s_i \mid s_{i-1})$$

holds $\forall \, s \in K^n$

**Definition 1b**  A p.d. on $K^n$ is a Markov chain if

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

holds $\forall \, s \in K^n$, where $g_i : K^2 \to \mathbb{R}_+$ are some functions

Equivalence:
- a) $\to$ b)  trivial
- b) $\to$ a)  recursively apply the following step

$$p(s_{n-1}, s_n) = \left\{ \sum_{s_1, \ldots, s_{n-2}} \prod_{i=2}^{n-1} g_i(s_{i-1}, s_i) \right\} g_n(s_{n-1}, s_n)$$

$$\quad \hookrightarrow \quad g_n(s_{n-1}, s_n) = p(s_n \mid s_{n-1}) \, b_{n-1}(s_{n-1}) \quad \text{with same } b_{n-1}$$

Therefore, we have

$$p(s_1,...,s_n) = \underbrace{\left[ \prod_{i=2}^{n-1} g_i(s_{i-1}, s_i) \right] b_{n-1}(s_{n-1})}_{p(s_1,...,s_{n-1})} \cdot p(s_n | s_{n-1})$$

Another useful formula

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

Example 1 (Ehrenfest model)

Consider $N$ particles in two containers. At each discrete time $t = 1, 2, ...$, independently from the past, a particle is selected at random and moved to the other container. Let $S_t$ denote the number of particles in the first container at time $t$. Then we have

$$p(S_t = k \mid S_{t-1} = \ell) = \begin{cases} \frac{N-\ell}{N} & \text{if } k = \ell+1 \\ \frac{\ell}{N} & \text{if } k = \ell-1 \\ 0 & \text{otherwise} \end{cases}$$

Q: How does $p(S_t = k)$, $k = 1, ..., N$ behave for $t \to \infty$?

Example 2 (Random walk on a graph)

Consider a random walk on an undirected graph $V, E$

- $K = V$ states, $S_t \in V$ position of the walker at time $t$
- $p(S_1)$ some p.d. for the start vertex
- $p(S_t = i \mid S_{t-1} = j) = \begin{cases} W_{ij} & \text{if } \{i,j\} \in E \\ 0 & \text{otherwise} \end{cases}$

where the $W_{ij}$ fulfill $\sum_{i \in N(j)} W_{ij} = 1 \quad \forall j \in V$

## B. Homogeneous Markov chains, stationary distributions

**Definition 2** A Markov chain is *homogeneous* if its conditional prob's $P(s_i | s_{i-1})$ do not depend on the position $i$, i.e.

$$P(s_i = k | s_{i-1} = k') = g(k, k') \quad \forall \; i = 2, .., n.$$

We know that

$$P(s_i = k) = \sum_{k' \in K} P(s_i = k | s_{i-1} = k') \, p(s_{i-1} = k').$$

Consider $p(s_i = k)$, $k \in K$ as components of a vector $\overline{\pi}_i \in \mathbb{R}_+^K$ and $p(s_i = k | s_{i-1} = k')$, $k, k' \in K$ as elements of a $K \times K$ matrix $P$. Then the previous eq. reads

$$\overline{\pi}_i = P \, \overline{\pi}_{i-1}$$

and more general, we have $\overline{\pi}_i = P^{i-1} \overline{\pi}_1$.

It may happen that there $\exists$ a p.d. $\overline{\pi}^*$ on $K$ s.t. $P\overline{\pi}^* = \overline{\pi}^*$. We call it a *stationary* p.d. of $P$.

**Definition 3** A homogeneous Markov chain is *irreducible* if for each pair of states $k, k'$ there is an $m > 0$ s.t. $P_{k k'}^m > 0$. I.e. there is a non-zero probability to reach state $k$ starting from ~~from~~ state $k'$ (after $m$ transitions)

A somewhat stronger condition ensures the existence & uniqueness of a stationary distribution and convergence to it.

**Theorem 1** (w/o proof) If for some $m > 0$ all elements of the matrix $P^m$ are strictly positive, then the Markov chain has a *unique* stationary distribution $\overline{\pi}^*$, which is a fixpoint

$$P^n \overline{\pi} \xrightarrow{n \to \infty} \overline{\pi}^* \quad \forall \; \overline{\pi}$$

Moreover

$$P^n = \pi^* \otimes e + E(n),$$

where $e = (1, \ldots, 1)$ and $E_{kk'}(n) = \mathcal{O}(h^n)$ with some $0 < h < 1$. ∎

Definition 4   A Markov chain satisfies the <u>detailed balance</u> condition if it has a stationary distribution $\pi \in \mathbb{R}^k_+$ s.t.

$$p(s_i | s_{i-1}) \pi(s_{i-1}) = p(s_{i-1} | s_i) \pi(s_i).$$

This means that the reverse Markov chain has the same transition probability matrix as the forward chain.
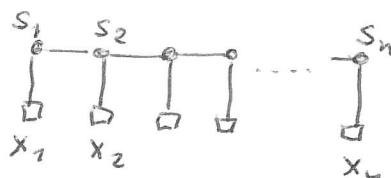
## C.   Hidden Markov models on chains

Common models in pattern recognition

$X = (X_1, \ldots, X_n)$ sequence of features (observable)

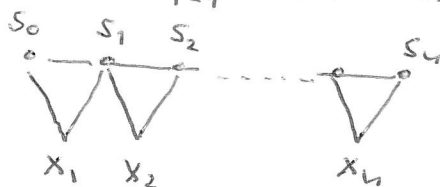$S = (S_1, \ldots, S_n)$ sequence of states (hidden)

<u>Hidden Markov model</u> (HMM):   a p.d. on pairs $(x, s)$ s.t.

a)   $$p(x,s) = \underbrace{\prod_{i=1}^{n} p(x_i | s_i)}_{p(x|s)} \cdot \underbrace{p(s_1) \cdot \prod_{i=2}^{n} p(s_i | s_{i-1})}_{p(s) \,-\, \text{Markov model}}$$



b) or slightly more general

$$P(x,s) = p(s_o) \prod_{i=1}^{n} p(x_i, s_i | s_{i-1})$$



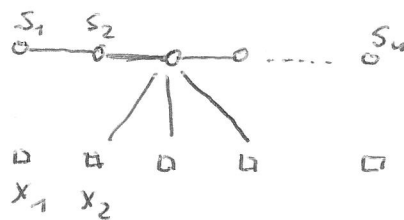Remark 1   This describes a stochastic regular language.

## Conditional HMM

As before, $X = (x_1, .., x_n)$ – sequence of features and
$S = (s_1, .., s_n)$ – sequence of hidden states

Discriminative model $\rightarrow$ we model only $p(s/x)$

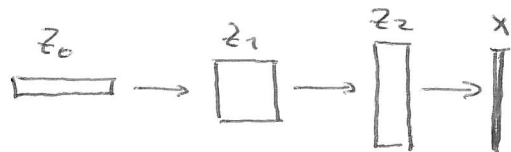$$P(S/x) = \frac{1}{Z(x)} \prod_{i=2}^{n} g_i (s_{i-1}, s_i, X),$$

where $Z(x)$ is a normalisation constant



Such models allow to model a direct dependence of $s_i$ on a larger context window of features

## Hierarchical variational autoencoders & diffusion models

Generative latent variable models (deep learning)



$z_0, .., z_n$ – latent variables (vectors, tensors)     $X$ – image

The model is specified by

$p(z_0)$ – simple distribution (uniform, standard Gaussian, etc.)

$p_\theta(z_i/z_{i-1})$ – parametrised conditional distributions

$p_\theta(x/z_n)$ – conditional distribution on images

If $z_k \in \mathbb{B}^{n_k}$, i.e. $z_k$ is a binary valued vector $\Rightarrow$

$$\log p(z_k/z_{k-1}) = \langle z_k, f(z_{k-1}, \theta) \rangle - C(z_{k-1}),$$

where $f(z_{k-1}, \theta)$ is modelled by a (deep) network.
$C(z_{k-1})$ is the log-partition function (normalising constant)