

CS 340 Lec. 21: Hidden Markov Models

AD

April 2011

Modelling Dependent Data

- For the time being, we have always assumed that available data $\{\mathbf{x}_t\}_{t=1}^T$ are independent.
- In numerous applications, we only have access to data which are statistically dependent; i.e.

$$p\left(\{\mathbf{x}_t\}_{t=1}^T\right) \neq \prod_{t=1}^T p(\mathbf{x}_t).$$

- Typical applications include: speech processing, tracking, stock prices.
- Most popular model for time dependent data is Hidden Markov Models = Mixture Models + Markov chain on the “cluster labels”.

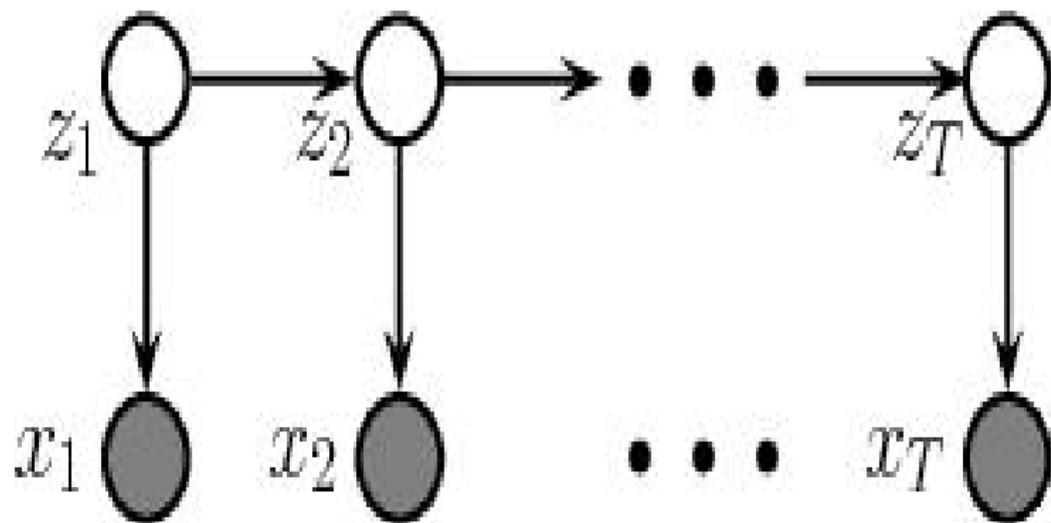
- In a standard mixture models, we have independent cluster labels $\{z_t\}_{t=1}^T$ and data $\{\mathbf{x}_t\}_{t=1}^T$ so

$$\begin{aligned} p\left(\{z_t\}_{t=1}^T, \{\mathbf{x}_t\}_{t=1}^T\right) &= \prod_{t=1}^T p(z_t, \mathbf{x}_t) = \prod_{t=1}^T p(z_t) p(\mathbf{x}_t | z_t) \\ &= \prod_{t=1}^T p(z_t) \prod_{t=1}^T p(\mathbf{x}_t | z_t) \end{aligned}$$

- In an HMM model, the cluster labels $\{z_t\}_{t=1}^T$ follow a Markov chain

$$p\left(\{z_t\}_{t=1}^T, \{\mathbf{x}_t\}_{t=1}^T\right) = p(z_1) \prod_{t=2}^T p(z_t | z_{t-1}) \prod_{t=1}^T p(\mathbf{x}_t | z_t)$$

Graphical Representation



HMM as a directed graphical model

- Assume $z_t \in \{1, \dots, K\}$ then the Markov chain is defined by its initial distribution

$$p(z_1 = k) = \pi_k$$

and the transition probabilities

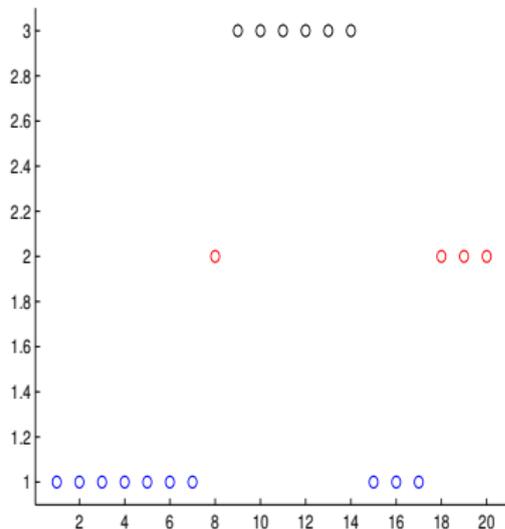
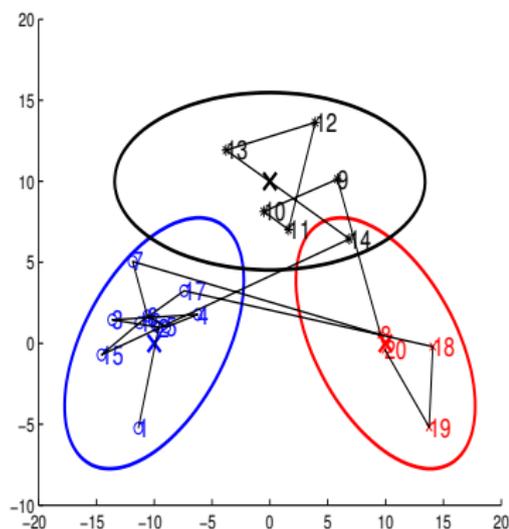
$$p(z_t = l | z_{t-1} = k) = P_{k,l}$$

- We have the conditional densities/distribution

$$p(\mathbf{x}_t | z_t = k) = p_k(\mathbf{x}_t)$$

where we could have for example $p_k(\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_t; \mu_k, \Sigma_k)$ or $p_k(\mathbf{x}_t = 1) = \alpha_k$.

Data Sampled from an HMM



(left) Some data sampled from a 3 state HMM. Each state emits from a 2d Gaussian. (right) The hidden state sequence $\{z_t\}$.

- *Automatic speech recognition*: \mathbf{x}_t is the speech signal, z_t represents the word that is being spoken.
- *Activity recognition*: \mathbf{x}_t represents a video frame, z_t is the class of activity the person is engaged in (e.g., running, walking, sitting, etc.)
- *Part of speech tagging*: \mathbf{x}_t represents a word, z_t represents its part of speech (noun, verb, adjective, etc.)
- *Gene finding*: \mathbf{x}_t represents the DNA nucleotides (A,C,G,T), z_t represents whether we are inside a gene-coding region or not.

- Assume for the time being that the parameters of the models are known, we want to estimate z_t given observations $\{\mathbf{x}_t\}$.
- **Filtering**: compute $p(z_t = k | \mathbf{x}_{1:t})$
- **Prediction**: compute $p(z_{t+L} = k | \mathbf{x}_{1:t})$ for $L > 0$
- **Smoothing**: compute $p(z_t = k | \mathbf{x}_{1:T})$
- In the independent case,
 $p(z_t = k | \mathbf{x}_{1:t}) = p(z_t = k | \mathbf{x}_{1:T}) = p(z_t = k | \mathbf{x}_t)$ and
 $p(z_{t+L} = k | \mathbf{x}_{1:t}) = p(z_{t+L} = k)$.

Inference in HMM: Filtering

- Given the filter $p(z_{t-1} | \mathbf{x}_{1:t-1})$ at time $t - 1$, we compute $p(z_t | \mathbf{x}_{1:t})$ as follows.
- Prediction:*

$$\begin{aligned} p(z_t = k | \mathbf{x}_{1:t-1}) &= \sum_{l=1}^K p(z_{t-1} = l, z_t = k | \mathbf{x}_{1:t-1}) \\ &= \sum_{l=1}^K p(z_t = k | \mathbf{x}_{1:t-1}, z_{t-1} = l) p(z_{t-1} = l | \mathbf{x}_{1:t-1}) \\ &= \sum_{l=1}^K p(z_t = k | z_{t-1} = l) p(z_{t-1} = l | \mathbf{x}_{1:t-1}) \\ &= \sum_{l=1}^K P_{l,k} p(z_{t-1} = l | \mathbf{x}_{1:t-1}) \end{aligned}$$

- Update:*

$$\begin{aligned} p(z_t = k | \mathbf{x}_{1:t}) &= \frac{p(\mathbf{x}_t | z_t = k) p(z_t = k | \mathbf{x}_{1:t-1})}{\sum_{l=1}^K p(\mathbf{x}_t | z_t = l) p(z_t = l | \mathbf{x}_{1:t-1})} \\ &= \frac{p_k(\mathbf{x}_t) p(z_t = k | \mathbf{x}_{1:t-1})}{\sum_{l=1}^K p_l(\mathbf{x}_t) p(z_t = l | \mathbf{x}_{1:t-1})} \end{aligned}$$

- This has computational complexity $O(K^2 T)$.

Inference in HMM: Prediction

- We want to compute $p(z_{t+L} = k | \mathbf{x}_{1:t})$ for $L > 0$.
- We have

$$p(z_{t+1} = k | \mathbf{x}_{1:t}) = \sum_{l=1}^K P_{l,k} p(z_t = l | \mathbf{x}_{1:t-1})$$

and similarly

$$p(z_{t+m} = k | \mathbf{x}_{1:t}) = \sum_{l=1}^K P_{l,k} p(z_{t+m-1} = l | \mathbf{x}_{1:t-1})$$

Inference in HMM: Smoothing

- We have for $1 \leq t < T$

$$p(z_t = k | \mathbf{x}_{1:T}) = \frac{p(z_t = k | \mathbf{x}_{1:t-1}) p(\mathbf{x}_{t:T} | z_t = k)}{\sum_{l=1}^K p(z_t = l | \mathbf{x}_{1:t-1}) p(\mathbf{x}_{t:T} | z_t = l)}$$

- We can compute $p(\mathbf{x}_{t:T} | z_t = k)$ using the following backward recursion initialized at $p(\mathbf{x}_T | z_T = k) = p_k(\mathbf{x}_T)$

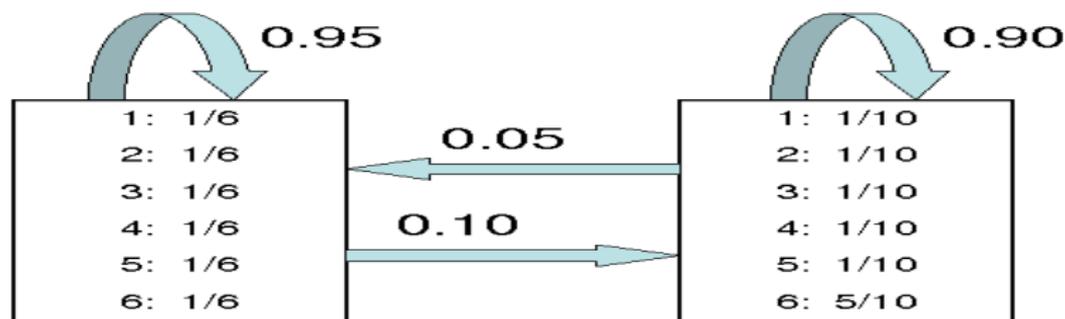
$$\begin{aligned} p(\mathbf{x}_{t+1:T} | z_t = k) &= \sum_{l=1}^K p(\mathbf{x}_{t+1:T}, z_{t+1} = l | z_t = k) \\ &= \sum_{l=1}^K p(\mathbf{x}_{t+1:T} | z_t = k, z_{t+1} = l) p(z_{t+1} = l | z_t = k) \\ &= \sum_{l=1}^K p(\mathbf{x}_{t+1:T} | z_{t+1} = l) P_{k,l} \end{aligned}$$

and

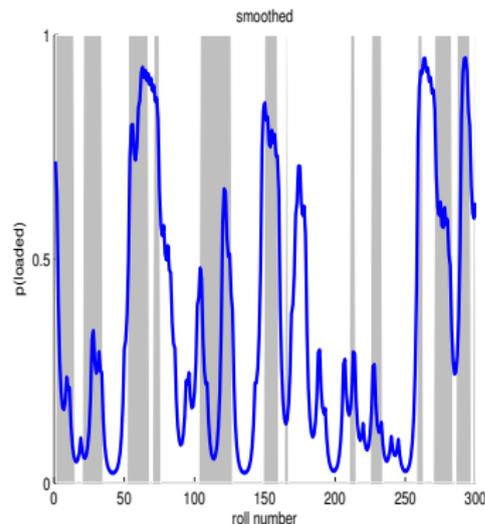
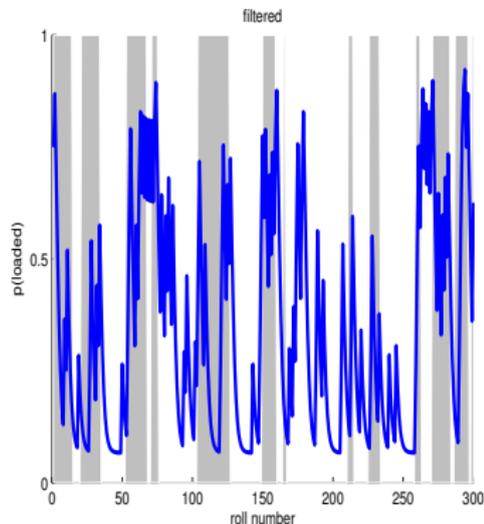
$$\begin{aligned} p(\mathbf{x}_{t:T} | z_t = k) &= p(\mathbf{x}_t, \mathbf{x}_{t+1:T} | z_t = k) \\ &= p_k(\mathbf{x}_t) p(\mathbf{x}_{t+1:T} | z_t = k) \end{aligned}$$

Example: Casino

- In this model, $x_t \in \{1, 2, \dots, 6\}$ represents which dice face shows up, and z_t represents the identity of the dice that is being used. Most of the time the casino uses a fair dice, $z = 1$, but occasionally it switches to a loaded dice, $z = 2$, for a short period.
- If $z = 1$ the observation distribution is a uniform distribution over $\{1, 2, \dots, 6\}$. If $z = 2$, the observation distribution is skewed towards face 6.



Example: Casino



Inference in the dishonest casino. Vertical gray bars denote the samples that we generated using a loaded die. (left) Filtered estimate of probability of using a loaded dice. We hope to see a spike up whenever there is a gray bar. (right) Smoothed estimates.