

# Hidden Markov Models

## STATS 305B: Applied Statistics

Scott Linderman

February 12, 2025

# Recap

Last time...

- ▶ Mixture Models (Gaussian and general exp-fam)
- ▶ K-Means as MAP estimation
- ▶ EM as maximizing the marginal log likelihood
- ▶ Connecting EM and VI

# Outline

Today...

- ▶ Hidden Markov Models
- ▶ The forward-backward algorithm
- ▶ EM for HMMs
- ▶ Implementation Details

## Recap: Gaussian Mixture Models

Recall the basic Gaussian mixture model,

$$z_t \stackrel{\text{iid}}{\sim} \text{Cat}(\boldsymbol{\pi}) \quad (1)$$

$$x_t \mid z_t \sim \mathcal{N}(\boldsymbol{\mu}_{z_t}, \boldsymbol{\Sigma}_{z_t}) \quad (2)$$

where

- ▶  $z_t \in \{1, \dots, K\}$  is a **latent mixture assignment**
- ▶  $x_t \in \mathbb{R}^D$  is an **observed data point**
- ▶  $\boldsymbol{\pi} \in \Delta_K$ ,  $\boldsymbol{\mu}_k \in \mathbb{R}^D$ , and  $\boldsymbol{\Sigma}_k \in \mathbb{R}_{\geq 0}^{D \times D}$  are parameters

(Here we've switched to indexing data points by  $t$  rather than  $n$ .)

Let  $\Theta$  denote the set of parameters. We can be Bayesian and put a prior on  $\Theta$  and run Gibbs or VI, or we can point estimate  $\Theta$  with EM, etc.

## Recap: Gaussian Mixture Models II

Draw the graphical model.

## Recap: Gaussian Mixture Models III

Recall the EM algorithm for mixture models,

- **E step:** Compute the posterior distribution

$$q(\mathbf{z}_{1:T}) = p(\mathbf{z}_{1:T} \mid \mathbf{x}_{1:T}; \Theta) \quad (3)$$

$$= \prod_{t=1}^T p(z_t \mid \mathbf{x}_t; \Theta) \quad (4)$$

$$= \prod_{t=1}^T q_t(z_t) \quad (5)$$

- **M step:** Maximize the ELBO wrt  $\Theta$ ,

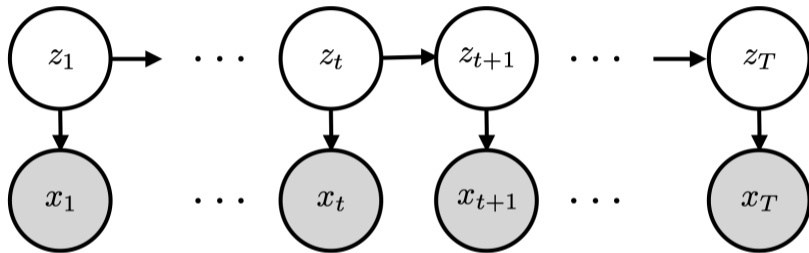
$$\mathcal{L}(\Theta) = \mathbb{E}_{q(\mathbf{z}_{1:T})} [\log p(\mathbf{x}_{1:T}, \mathbf{z}_{1:T}; \Theta) - \log q(\mathbf{z}_{1:T})] \quad (6)$$

$$= \mathbb{E}_{q(\mathbf{z}_{1:T})} [\log p(\mathbf{x}_{1:T}, \mathbf{z}_{1:T}; \Theta)] + c. \quad (7)$$

For exponential family mixture models, the M-step only requires expected sufficient statistics.

## Hidden Markov Models

Hidden Markov Models (HMMs) are like mixture models with temporal dependencies between the mixture assignments.



This graphical model says that the joint distribution factors as,

$$p(z_{1:T}, \mathbf{x}_{1:T}) = p(z_1) \prod_{t=2}^T p(z_t | z_{t-1}) \prod_{t=1}^T p(\mathbf{x}_t | z_t). \quad (8)$$

We call this an HMM because the *hidden* states follow a Markov chain,  $p(z_1) \prod_{t=2}^T p(z_t | z_{t-1})$ .

# Hidden Markov Models II

An HMM consists of three components:

- 1. Initial distribution:**  $z_1 \sim \text{Cat}(\boldsymbol{\pi}_0)$
- 2. Transition matrix:**  $z_t \sim \text{Cat}(\mathbf{P}_{z_{t-1}})$  where  $\mathbf{P} \in [0, 1]^{K \times K}$  is a *row-stochastic* transition matrix with rows  $\mathbf{P}_k$ .
- 3. Emission distribution:**  $\mathbf{x}_t \sim p(\cdot \mid \boldsymbol{\theta}_{z_t})$



## Example: The *occasionally* dishonest casino

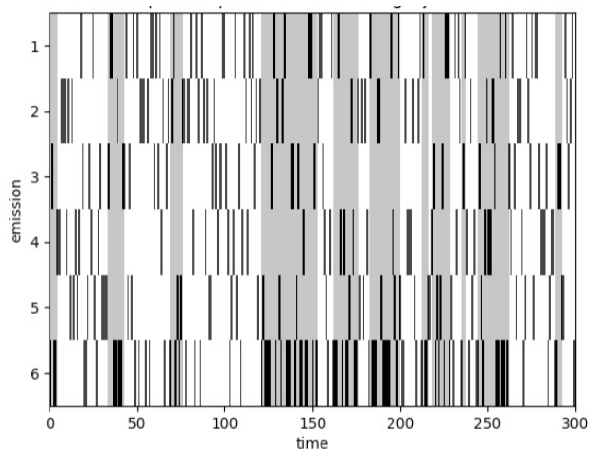
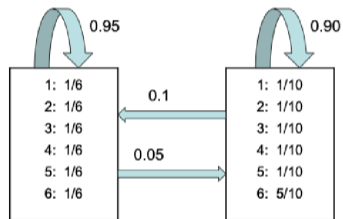


Figure: An *occasionally* dishonest casino that sometimes throws loaded dice.

From [https://probml.github.io/dynamax/notebooks/hmm/casino\\_hmm\\_inference.html](https://probml.github.io/dynamax/notebooks/hmm/casino_hmm_inference.html)

## Example: HMM for splice site recognition

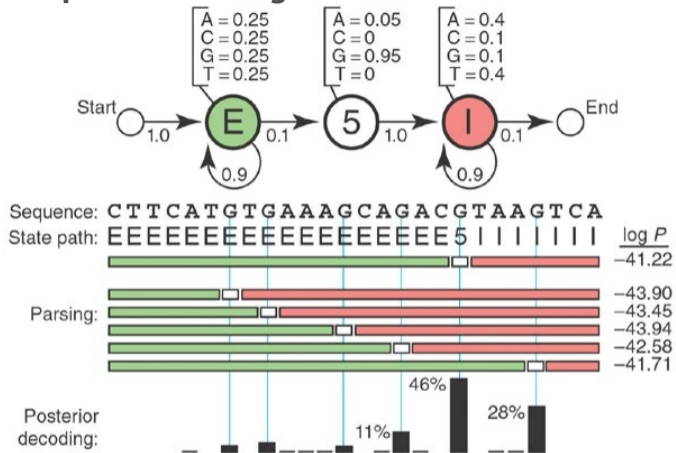


Figure: A toy model for parsing a genome to find 5' splice sites. From ?.

**Question:** Suppose the splice site always had a GT sequence. How would you change the model to detect such sites?

## Example: Autoregressive HMM for video segmentation

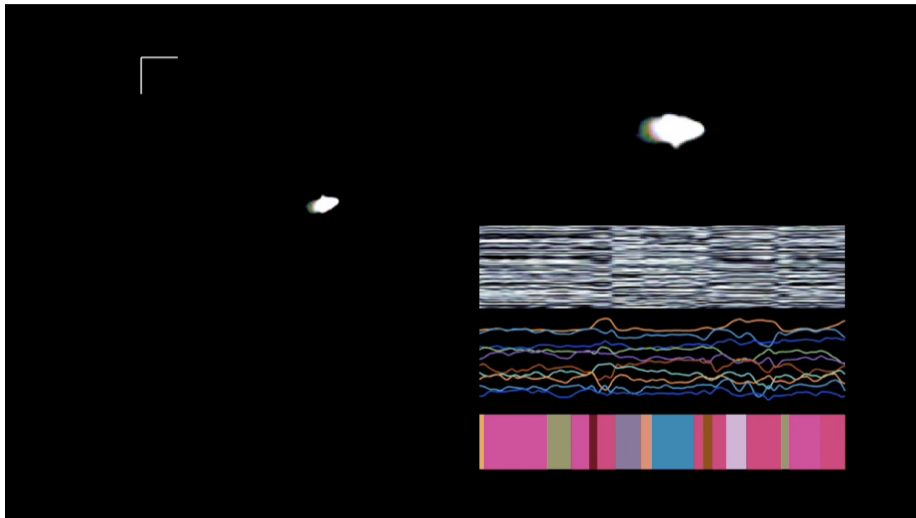


Figure: Segmenting videos of freely moving mice [?]. (Show video.)

## Hidden Markov Models III

We are interested in questions like:

- ▶ What are the *predictive distributions* of  $p(z_{t+1} \mid \mathbf{x}_{1:t})$ ?
- ▶ What is the *posterior marginal* distribution  $p(z_t \mid \mathbf{x}_{1:T})$ ?
- ▶ What is the *posterior pairwise marginal* distribution  $p(z_t, z_{t+1} \mid \mathbf{x}_{1:T})$ ?
- ▶ What is the *posterior mode*  $\mathbf{z}_{1:T}^* = \arg \max p(\mathbf{z}_{1:T} \mid \mathbf{x}_{1:T})$ ?
- ▶ How can we *sample the posterior*  $p(\mathbf{z}_{1:T} \mid \mathbf{x}_{1:T})$  of an HMM?
- ▶ What is the *marginal likelihood*  $p(\mathbf{x}_{1:T})$ ?
- ▶ How can we *learn the parameters* of an HMM?

**Question:** Why might these sound like hard problems?

## Computing the predictive distributions

The predictive distributions give the probability of the latent state  $z_{t+1}$  given observations *up to but not including* time  $t + 1$ . Let,

$$\alpha_{t+1}(z_{t+1}) \triangleq p(z_{t+1}, \mathbf{x}_{1:t}) \quad (9)$$

$$= \sum_{z_1=1}^K \cdots \sum_{z_t=1}^K p(z_1) \prod_{s=1}^t p(\mathbf{x}_s | z_s) p(z_{s+1} | z_s) \quad (10)$$

$$= \sum_{z_t=1}^K \left[ \left( \sum_{z_1=1}^K \cdots \sum_{z_{t-1}=1}^K p(z_1) \prod_{s=1}^{t-1} p(\mathbf{x}_s | z_s) p(z_{s+1} | z_s) \right) p(\mathbf{x}_t | z_t) p(z_{t+1} | z_t) \right] \quad (11)$$

$$= \sum_{z_t=1}^K \alpha_t(z_t) p(\mathbf{x}_t | z_t) p(z_{t+1} | z_t). \quad (12)$$

We call  $\alpha_t(z_t)$  the *forward messages*. We can compute them recursively! The base case is  $p(z_1 | \emptyset) \triangleq p(z_1)$ .

## Computing the predictive distributions II

We can also write these recursions in a vectorized form. Let

$$\boldsymbol{\alpha}_t = \begin{bmatrix} \alpha_t(z_t = 1) \\ \vdots \\ \alpha_t(z_t = K) \end{bmatrix} = \begin{bmatrix} p(z_t = 1, \mathbf{x}_{1:t-1}) \\ \vdots \\ p(z_t = K, \mathbf{x}_{1:t-1}) \end{bmatrix} \quad \text{and} \quad \mathbf{l}_t = \begin{bmatrix} p(\mathbf{x}_t | z_t = 1) \\ \vdots \\ p(\mathbf{x}_t | z_t = K) \end{bmatrix} \quad (13)$$

both be vectors in  $\mathbb{R}_+^K$ . Then,

$$\boldsymbol{\alpha}_{t+1} = \mathbf{P}^\top (\boldsymbol{\alpha}_t \odot \mathbf{l}_t) \quad (14)$$

where  $\odot$  denotes the Hadamard (elementwise) product and  $\mathbf{P}$  is the transition matrix.

## Computing the predictive distributions III

Finally, to get the predictive distributions we just have to normalize,

$$p(z_{t+1} | \mathbf{x}_{1:t}) \propto p(z_{t+1}, \mathbf{x}_{1:t}) = \alpha_{t+1}(z_{t+1}). \quad (15)$$

**Question:** What does the normalizing constant tell us?

## Computing the posterior marginal distributions

The posterior marginal distributions give the probability of the latent state  $z_t$  given *all the observations* up to time  $T$ .

$$p(z_t | \mathbf{x}_{1:T}) = \sum_{z_1=1}^K \cdots \sum_{z_{t-1}=1}^K \sum_{z_{t+1}=1}^K \cdots \sum_{z_T=1}^K p(\mathbf{z}_{1:T}, \mathbf{x}_{1:T}) \quad (16)$$

$$\begin{aligned} &= \left[ \sum_{z_t=1}^K \cdots \sum_{z_{t-1}=1}^K p(z_1) \prod_{s=1}^{t-1} p(\mathbf{x}_s | z_s) p(z_{s+1} | z_s) \right] \times p(\mathbf{x}_t | z_t) \\ &\quad \times \left[ \sum_{z_{t+1}=1}^K \cdots \sum_{z_T=1}^K \prod_{u=t+1}^T p(z_u | z_{u-1}) p(\mathbf{x}_u | z_u) \right] \end{aligned} \quad (17)$$

$$= \alpha_t(z_t) \times p(\mathbf{x}_t | z_t) \times \beta_t(z_t) \quad (18)$$

where we have introduced the *backward messages*  $\beta_t(z_t)$ .



## Computing the backward messages

The backward messages can be computed recursively too,

$$\beta_t(z_t) \triangleq \sum_{z_{t+1}=1}^K \cdots \sum_{z_T=1}^K \prod_{u=t+1}^T p(z_u | z_{u-1}) p(\mathbf{x}_u | z_u) \quad (19)$$

$$= \sum_{z_{t+1}=1}^K p(z_{t+1} | z_t) p(\mathbf{x}_{t_1} | z_{t+1}) \left( \sum_{z_{t+2}=1}^K \cdots \sum_{z_T=1}^K \prod_{u=t+2}^T p(z_u | z_{u-1}) p(\mathbf{x}_u | z_u) \right) \quad (20)$$

$$= \sum_{z_{t+1}=1}^K p(z_{t+1} | z_t) p(\mathbf{x}_{t_1} | z_{t+1}) \beta_{t+1}(z_{t+1}). \quad (21)$$

For the base case, let  $\beta_T(z_T) = 1$ .

## Computing the backward messages (vectorized)

Let

$$\boldsymbol{\beta}_t = \begin{bmatrix} \beta_t(z_t = 1) \\ \vdots \\ \beta_t(z_t = K) \end{bmatrix} \quad (22)$$

be a vector in  $\mathbb{R}_+^K$ . Then,

$$\boldsymbol{\beta}_t = \mathbf{P}(\boldsymbol{\beta}_{t+1} \odot \mathbf{l}_{t+1}). \quad (23)$$

Let  $\boldsymbol{\beta}_T = \mathbf{1}_K$ .

Now we have everything we need to compute the posterior marginal,

$$p(z_t = k \mid \mathbf{x}_{1:T}) = \frac{\alpha_{t,k} l_{t,k} \beta_{t,k}}{\sum_{j=1}^K \alpha_{t,j} l_{t,j} \beta_{t,j}}. \quad (24)$$

We just derived the **forward-backward algorithm** for HMMs [?].

## What do the backward messages represent?

**Question:** If the forward messages represent the predictive probabilities  $\alpha_{t+1}(z_{t+1}) = p(z_{t+1}, \mathbf{x}_{1:t})$ , what do the backward messages represent?

## Computing the posterior pairwise marginals

**Exercise:** Use the forward and backward messages to compute the posterior pairwise marginals

$$p(z_t, z_{t+1} \mid \mathbf{x}_{1:T}).$$

## Normalizing the messages for numerical stability

If you're working with long time series, especially if you're working with 32-bit floating point, you need to be careful.

The messages involve products of probabilities, which can quickly overflow.

There's a simple fix though: after each step, re-normalize the messages so that they sum to one. I.e. replace

$$\alpha_{t+1} = \mathbf{P}^\top (\alpha_t \odot l_t) \quad (25)$$

with

$$\tilde{\alpha}_{t+1} = \frac{1}{A_t} \mathbf{P}^\top (\tilde{\alpha}_t \odot l_t) \quad (26)$$

$$A_t = \sum_{k=1}^K \sum_{j=1}^K P_{jk} \tilde{\alpha}_{t,j} l_{t,j} \equiv \sum_{j=1}^K \tilde{\alpha}_{t,j} l_{t,j} \quad (\text{since } \mathbf{P} \text{ is row-stochastic}). \quad (27)$$

This leads to a nice interpretation: The normalized messages are predictive likelihoods  $\tilde{\alpha}_{t+1,k} = p(z_{t+1} = k \mid \mathbf{x}_{1:t})$ , and the normalizing constants are  $A_t = p(\mathbf{x}_t \mid \mathbf{x}_{1:t-1})$ .

## EM for Hidden Markov Models

Now we can put it all together. To perform EM in an HMM,

- **E step:** Compute the posterior distribution

$$q(\mathbf{z}_{1:T}) = p(\mathbf{z}_{1:T} \mid \mathbf{x}_{1:T}; \Theta). \quad (28)$$

(Really, run the **forward-backward algorithm** to get posterior marginals and pairwise marginals.)

- **M step:** Maximize the ELBO wrt  $\Theta$ ,

$$\mathcal{L}(\Theta) = \mathbb{E}_{q(\mathbf{z}_{1:T})} [\log p(\mathbf{x}_{1:T}, \mathbf{z}_{1:T}; \Theta)] + c \quad (29)$$

$$\begin{aligned} &= \mathbb{E}_{q(\mathbf{z}_{1:T})} \left[ \sum_{k=1}^K \mathbb{I}[z_1 = k] \log \pi_{0,k} \right] + \mathbb{E}_{q(\mathbf{z}_{1:T})} \left[ \sum_{t=1}^{T-1} \sum_{i=1}^K \sum_{j=1}^K \mathbb{I}[z_t = i, z_{t+1} = j] \log P_{i,j} \right] \\ &\quad + \mathbb{E}_{q(\mathbf{z}_{1:T})} \left[ \sum_{t=1}^T \sum_{k=1}^K \mathbb{I}[z_t = k] \log p(\mathbf{x}_t; \theta_k) \right] \end{aligned} \quad (30)$$

For exponential family observations, the M-step only requires expected sufficient statistics.

# What else?

- ▶ How can we sample the posterior?
- ▶ How can we find the posterior mode?
- ▶ How can we choose the number of states?
- ▶ What if my transition matrix is sparse?

# References