

# HIDDEN MARKOV MODELS

Source: Hastie et al. (2009), Daumé III. Thanks to D. Hsu.

Please do not distribute these slides publicly, beyond using them for this course.

# MARKOV MODELS

**Markov model:** a stochastic process  $\{Y_t\}_{t \in \mathbb{N}}$  where, for each  $t \in \mathbb{N}$ , the conditional distribution of the next state  $Y_{t+1}$  given all previous states  $\{Y_\tau : \tau \leq t\}$  only depends on the value of the current state  $Y_t$ .

# MARKOV MODELS

**Markov model:** a stochastic process  $\{Y_t\}_{t \in \mathbb{N}}$  where, for each  $t \in \mathbb{N}$ , the conditional distribution of the next state  $Y_{t+1}$  given all previous states  $\{Y_\tau : \tau \leq t\}$  only depends on the value of the current state  $Y_t$ .

Conditioned on present  $Y_t$ , past  $\{Y_\tau\}_{\tau < t}$  and future  $\{Y_\tau\}_{\tau > t}$  are independent.

$$\cdots \longrightarrow Y_{t-1} \longrightarrow Y_t \longrightarrow Y_{t+1} \longrightarrow \cdots$$

# MARKOV MODELS

**Markov model:** a stochastic process  $\{Y_t\}_{t \in \mathbb{N}}$  where, for each  $t \in \mathbb{N}$ , the conditional distribution of the next state  $Y_{t+1}$  given all previous states  $\{Y_\tau : \tau \leq t\}$  only depends on the value of the current state  $Y_t$ .

Conditioned on present  $Y_t$ , past  $\{Y_\tau\}_{\tau < t}$  and future  $\{Y_\tau\}_{\tau > t}$  are independent.

$$\cdots \longrightarrow Y_{t-1} \longrightarrow Y_t \longrightarrow Y_{t+1} \longrightarrow \cdots$$

Specifying a Markov chain (with discrete **state space**  $[K] = \{1, 2, \dots, K\}$ ):

- ▶ **Initial state distribution:**  $K$ -dimensional probability vector  $\pi$

$$\pi_i = \Pr(Y_1 = i).$$

- ▶ **Transition matrix:**  $K \times K$  matrix  $\mathbf{A}$

$$A_{i,j} = \Pr(Y_{t+1} = j \mid Y_t = i)$$

(rows of  $\mathbf{A}$  are probability vectors).

# HIDDEN MARKOV MODELS

**Hidden Markov model (HMM):** a Markov chain  $\{(X_t, Y_t)\}_{t \in \mathbb{N}}$ , where

# HIDDEN MARKOV MODELS

**Hidden Markov model (HMM):** a Markov chain  $\{(X_t, Y_t)\}_{t \in \mathbb{N}}$ , where

- ▶  $\{Y_t\}_{t \in \mathbb{N}}$  is also a Markov chain (with state space  $[K] = \{1, 2, \dots, K\}$ )  
(hidden state sequence);

# HIDDEN MARKOV MODELS

**Hidden Markov model (HMM):** a Markov chain  $\{(X_t, Y_t)\}_{t \in \mathbb{N}}$ , where

- ▶  $\{Y_t\}_{t \in \mathbb{N}}$  is also a Markov chain (with state space  $[K] = \{1, 2, \dots, K\}$ ) (hidden state sequence);
- ▶ conditioned on  $Y_t$ , corresponding  $X_t$  is *independent of all other variables*;

# HIDDEN MARKOV MODELS

**Hidden Markov model (HMM):** a Markov chain  $\{(X_t, Y_t)\}_{t \in \mathbb{N}}$ , where

- ▶  $\{Y_t\}_{t \in \mathbb{N}}$  is also a Markov chain (with state space  $[K] = \{1, 2, \dots, K\}$ ) (hidden state sequence);
- ▶ conditioned on  $Y_t$ , corresponding  $X_t$  is *independent of all other variables*;
- ▶ the  $Y_t$  are *hidden*, and the  $X_t$  are *observed*.



# HIDDEN MARKOV MODELS

**Hidden Markov model (HMM):** a Markov chain  $\{(X_t, Y_t)\}_{t \in \mathbb{N}}$ , where

- ▶  $\{Y_t\}_{t \in \mathbb{N}}$  is also a Markov chain (with state space  $[K] = \{1, 2, \dots, K\}$ ) (**hidden state sequence**);
- ▶ conditioned on  $Y_t$ , corresponding  $X_t$  is *independent of all other variables*;
- ▶ the  $Y_t$  are *hidden*, and the  $X_t$  are *observed*.

$$\begin{array}{ccccccc} \cdots & \longrightarrow & Y_{t-1} & \longrightarrow & Y_t & \longrightarrow & Y_{t+1} & \longrightarrow & \cdots \\ & & \downarrow & & \downarrow & & \downarrow & & \\ & & X_{t-1} & & X_t & & X_{t+1} & & \end{array}$$

# HIDDEN MARKOV MODELS

**Hidden Markov model (HMM):** a Markov chain  $\{(X_t, Y_t)\}_{t \in \mathbb{N}}$ , where

- ▶  $\{Y_t\}_{t \in \mathbb{N}}$  is also a Markov chain (with state space  $[K] = \{1, 2, \dots, K\}$ ) (**hidden state sequence**);
- ▶ conditioned on  $Y_t$ , corresponding  $X_t$  is *independent of all other variables*;
- ▶ the  $Y_t$  are *hidden*, and the  $X_t$  are *observed*.

$$\begin{array}{ccccccc} \cdots & \longrightarrow & Y_{t-1} & \longrightarrow & Y_t & \longrightarrow & Y_{t+1} & \longrightarrow & \cdots \\ & & \downarrow & & \downarrow & & \downarrow & & \\ & & X_{t-1} & & X_t & & X_{t+1} & & \end{array}$$

**Time-homogeneous HMM:** *conditional* distribution of  $X_t$  given  $Y_t$  does not depend on  $t$ . (We'll focus on these.)

# HIDDEN MARKOV MODELS

**Hidden Markov model (HMM):** a Markov chain  $\{(X_t, Y_t)\}_{t \in \mathbb{N}}$ , where

- ▶  $\{Y_t\}_{t \in \mathbb{N}}$  is also a Markov chain (with state space  $[K] = \{1, 2, \dots, K\}$ ) (**hidden state sequence**);
- ▶ conditioned on  $Y_t$ , corresponding  $X_t$  is *independent of all other variables*;
- ▶ the  $Y_t$  are *hidden*, and the  $X_t$  are *observed*.

$$\begin{array}{ccccccc} \cdots & \longrightarrow & Y_{t-1} & \longrightarrow & Y_t & \longrightarrow & Y_{t+1} \longrightarrow \cdots \\ & & \downarrow & & \downarrow & & \downarrow \\ & & X_{t-1} & & X_t & & X_{t+1} \end{array}$$

**Time-homogeneous HMM:** *conditional* distribution of  $X_t$  given  $Y_t$  does not depend on  $t$ . (We'll focus on these.)

**Useful subscript notation:**  $Y_{s:t} = (Y_s, Y_{s+1}, \dots, Y_t)$  for  $s \leq t$ .

# HMM PARAMETERS (DISCRETE OBSERVATIONS)

For time-homogeneous HMM where  $X_t$  takes values in  $[D] = \{1, 2, \dots, D\}$ :

# HMM PARAMETERS (DISCRETE OBSERVATIONS)

For time-homogeneous HMM where  $X_t$  takes values in  $[D] = \{1, 2, \dots, D\}$ :

- **Initial state distribution:**  $K$ -dimensional probability vector  $\pi$

$$\pi_i = \Pr(Y_1 = i).$$

# HMM PARAMETERS (DISCRETE OBSERVATIONS)

For time-homogeneous HMM where  $X_t$  takes values in  $[D] = \{1, 2, \dots, D\}$ :

- **Initial state distribution:**  $K$ -dimensional probability vector  $\pi$

$$\pi_i = \Pr(Y_1 = i).$$

- **Transition matrix:**  $K \times K$  matrix  $\mathbf{A}$

$$A_{i,j} = \Pr(Y_{t+1} = j \mid Y_t = i)$$

(rows of  $\mathbf{A}$  are probability vectors).

# HMM PARAMETERS (DISCRETE OBSERVATIONS)

For time-homogeneous HMM where  $X_t$  takes values in  $[D] = \{1, 2, \dots, D\}$ :

- ▶ **Initial state distribution:**  $K$ -dimensional probability vector  $\pi$

$$\pi_i = \Pr(Y_1 = i).$$

- ▶ **Transition matrix:**  $K \times K$  matrix  $\mathbf{A}$

$$A_{i,j} = \Pr(Y_{t+1} = j \mid Y_t = i)$$

(rows of  $\mathbf{A}$  are probability vectors).

- ▶ **Emission matrix:**  $K \times D$  matrix  $\mathbf{B}$

$$B_{i,j} = \Pr(X_t = j \mid Y_t = i)$$

(rows of  $\mathbf{B}$  are probability vectors).

Solution: Viterbi algorithm and several other approaches

# CONNECTIONS TO MIXTURE MODELS

## Mixture model

$Y$

$\downarrow$

$X$

( $Y$  is hidden,  $X$  is observed.)

## Hidden Markov model

$Y_1 \rightarrow Y_2 \rightarrow \cdots \rightarrow Y_\ell$

$\downarrow$

$\downarrow$

$\downarrow$

$X_1$

$X_2$

$X_\ell$

( $Y_{1:\ell}$  is hidden,  $X_{1:\ell}$  is observed.)



# CONNECTIONS TO MIXTURE MODELS

## Mixture model

 $Y$  $\downarrow$  $\mathbf{X}$ 

( $Y$  is hidden,  $\mathbf{X}$  is observed.)

For  $K$  component mixture model,  
 $Y$  takes values in  $[K]$ .

## Hidden Markov model

 $Y_1 \rightarrow Y_2 \rightarrow \cdots \rightarrow Y_\ell$  $\downarrow$  $\downarrow$  $\downarrow$  $X_1$  $X_2$  $X_\ell$ 

( $Y_{1:\ell}$  is hidden,  $X_{1:\ell}$  is observed.)

# CONNECTIONS TO MIXTURE MODELS

## Mixture model

 $Y$  $\downarrow$  $X$ 

( $Y$  is hidden,  $X$  is observed.)

For  $K$  component mixture model,  
 $Y$  takes values in  $[K]$ .

## Hidden Markov model

 $Y_1 \rightarrow Y_2 \rightarrow \cdots \rightarrow Y_\ell$  $\downarrow$  $\downarrow$  $\downarrow$  $X_1$  $X_2$  $X_\ell$ 

( $Y_{1:\ell}$  is hidden,  $X_{1:\ell}$  is observed.)

For sequence of length  $\ell$ ,  
 $Y_{1:\ell}$  takes values in  $[K]^\ell$ .

# CONNECTIONS TO MIXTURE MODELS

## Mixture model

$Y$

$\downarrow$

$\mathbf{X}$

( $Y$  is hidden,  $\mathbf{X}$  is observed.)

For  $K$  component mixture model,  
 $Y$  takes values in  $[K]$ .

## Hidden Markov model

$Y_1 \rightarrow Y_2 \rightarrow \cdots \rightarrow Y_\ell$

$\downarrow$

$\downarrow$

$\downarrow$

$X_1$

$X_2$

$X_\ell$

( $Y_{1:\ell}$  is hidden,  $X_{1:\ell}$  is observed.)

For sequence of length  $\ell$ ,  
 $Y_{1:\ell}$  takes values in  $[K]^\ell$ .

Graphical diagram for HMM correctly suggests that every path—even ignoring arrow directions—is a Markov chain!

# CONNECTIONS TO MIXTURE MODELS

## Mixture model

$Y$

$\downarrow$

$X$

( $Y$  is hidden,  $X$  is observed.)

For  $K$  component mixture model,  
 $Y$  takes values in  $[K]$ .

## Hidden Markov model

$Y_1 \rightarrow Y_2 \rightarrow \cdots \rightarrow Y_\ell$

$\downarrow$

$\downarrow$

$\downarrow$

$X_1$

$X_2$

$X_\ell$

( $Y_{1:\ell}$  is hidden,  $X_{1:\ell}$  is observed.)

For sequence of length  $\ell$ ,  
 $Y_{1:\ell}$  takes values in  $[K]^\ell$ .

Graphical diagram for HMM correctly suggests that every path—even ignoring arrow directions—is a Markov chain!

►  $Y_1 \rightarrow Y_2 \rightarrow X_2$

# CONNECTIONS TO MIXTURE MODELS

## Mixture model

$Y$

$\downarrow$

$X$

( $Y$  is hidden,  $X$  is observed.)

For  $K$  component mixture model,  
 $Y$  takes values in  $[K]$ .

## Hidden Markov model

$Y_1 \rightarrow Y_2 \rightarrow \cdots \rightarrow Y_\ell$

$\downarrow$

$\downarrow$

$\downarrow$

$X_1$

$X_2$

$X_\ell$

( $Y_{1:\ell}$  is hidden,  $X_{1:\ell}$  is observed.)

For sequence of length  $\ell$ ,  
 $Y_{1:\ell}$  takes values in  $[K]^\ell$ .

Graphical diagram for HMM correctly suggests that every path—even ignoring arrow directions—is a Markov chain!

- ▶  $Y_1 \rightarrow Y_2 \rightarrow X_2$
- ▶  $X_2 \rightarrow Y_2 \rightarrow Y_3 \rightarrow X_3$

# CONNECTIONS TO MIXTURE MODELS

## Mixture model

 $Y$  $\downarrow$  $X$ 

( $Y$  is hidden,  $X$  is observed.)

For  $K$  component mixture model,  
 $Y$  takes values in  $[K]$ .

## Hidden Markov model

 $Y_1 \rightarrow Y_2 \rightarrow \cdots \rightarrow Y_\ell$  $\downarrow$  $\downarrow$  $\downarrow$  $X_1$  $X_2$  $X_\ell$ 

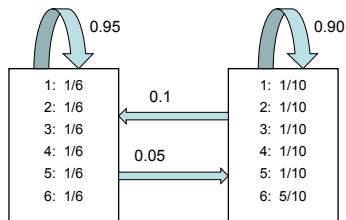
( $Y_{1:\ell}$  is hidden,  $X_{1:\ell}$  is observed.)

For sequence of length  $\ell$ ,  
 $Y_{1:\ell}$  takes values in  $[K]^\ell$ .

Graphical diagram for HMM correctly suggests that every path—even ignoring arrow directions—is a Markov chain!

- ▶  $Y_1 \rightarrow Y_2 \rightarrow X_2$
- ▶  $X_2 \rightarrow Y_2 \rightarrow Y_3 \rightarrow X_3$
- ▶  $X_1 \rightarrow Y_1 \rightarrow Y_{2:\ell} \rightarrow X_{2:\ell}$
- ▶ ...

# EXAMPLE: DISHONEST CASINO

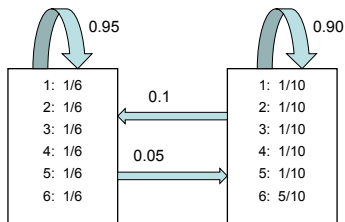


## Casino die-rolling game:

Randomly switch between two possible dice:  
one is fair, the other loaded.

The dice are otherwise indistinguishable!

# EXAMPLE: DISHONEST CASINO



## Casino die-rolling game:

Randomly switch between two possible dice:  
one is fair, the other loaded.

The dice are otherwise indistinguishable!

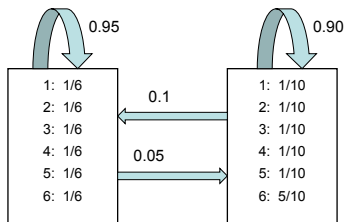
HMM parameters:

$$\mathbf{A} = \begin{matrix} & \begin{matrix} \text{fair die} & \text{loaded die} \end{matrix} \\ \begin{matrix} \text{fair die} \\ \text{loaded die} \end{matrix} & \begin{pmatrix} 0.95 & 0.05 \\ 0.10 & 0.90 \end{pmatrix} \end{matrix}, \quad \mathbf{B} = \begin{matrix} & \begin{matrix} 1 & 2 & 3 & 4 & 5 & 6 \end{matrix} \\ \begin{matrix} \text{fair die} \\ \text{loaded die} \end{matrix} & \begin{pmatrix} \frac{1}{6} & \frac{1}{6} & \frac{1}{6} & \frac{1}{6} & \frac{1}{6} & \frac{1}{6} \\ \frac{1}{10} & \frac{1}{10} & \frac{1}{10} & \frac{1}{10} & \frac{1}{10} & \frac{1}{2} \end{pmatrix} \end{matrix},$$

and  $\pi = (1, 0)$  if the casino starts out with the fair die.



# EXAMPLE: DISHONEST CASINO



## Casino die-rolling game:

Randomly switch between two possible dice:  
one is fair, the other loaded.

The dice are otherwise indistinguishable!

## HMM parameters:

### Transition matrix

$$A = \begin{matrix} & \begin{matrix} \text{fair die} & \text{loaded die} \end{matrix} \\ \begin{matrix} \text{fair die} \\ \text{loaded die} \end{matrix} & \begin{pmatrix} 0.95 & 0.05 \\ 0.10 & 0.90 \end{pmatrix} \end{matrix},$$

### Emission matrix

$$B = \begin{matrix} & \begin{matrix} 1 & 2 & 3 & 4 & 5 & 6 \end{matrix} \\ \begin{matrix} \text{fair die} \\ \text{loaded die} \end{matrix} & \begin{pmatrix} \frac{1}{6} & \frac{1}{6} & \frac{1}{6} & \frac{1}{6} & \frac{1}{6} & \frac{1}{6} \\ \frac{1}{10} & \frac{1}{10} & \frac{1}{10} & \frac{1}{10} & \frac{1}{10} & \frac{1}{2} \end{pmatrix} \end{matrix},$$

and  $\pi = (1, 0)$  if the casino starts out with the fair die.

**Problem:** Based on a sequence of rolls, guess which die was used at each time.

# HMM INFERENCE/LEARNING PROBLEMS

## Conditional probabilities (e.g., filtering/smoothing)

- ▶ **Given:** parameters  $\theta = (\pi, \mathbf{A}, \mathbf{B})$ , observation sequence  $x_{1:\ell} \in [D]^\ell$ .
- ▶ **Goal:** conditional distribution of  $Y_{s:t}$  given  $X_{1:\ell} = x_{1:\ell}$  ( $1 \leq s \leq t \leq \ell$ ):

$$\Pr_{\theta}(Y_{s:t} = y_{s:t} \mid X_{1:\ell} = x_{1:\ell}), \quad \text{for each } y_{s:t} \in [K]^{t-s+1}.$$

# HMM INFERENCE/LEARNING PROBLEMS

## Conditional probabilities (e.g., filtering/smoothing)

- ▶ **Given:** parameters  $\theta = (\pi, \mathbf{A}, \mathbf{B})$ , observation sequence  $x_{1:\ell} \in [D]^\ell$ .
- ▶ **Goal:** conditional distribution of  $Y_{s:t}$  given  $X_{1:\ell} = x_{1:\ell}$  ( $1 \leq s \leq t \leq \ell$ ):

$$\Pr_{\theta}(Y_{s:t} = y_{s:t} \mid X_{1:\ell} = x_{1:\ell}), \quad \text{for each } y_{s:t} \in [K]^{t-s+1}.$$

## Most probable state sequence (decoding)

- ▶ **Given:** parameters  $\theta = (\pi, \mathbf{A}, \mathbf{B})$ , observation sequence  $x_{1:\ell} \in [D]^\ell$ .
- ▶ **Goal:**  $\arg \max_{y_{1:\ell} \in [K]^\ell} \Pr_{\theta}(Y_{1:\ell} = y_{1:\ell} \mid X_{1:\ell} = x_{1:\ell})$ .

# HMM INFERENCE/LEARNING PROBLEMS

## Conditional probabilities (e.g., filtering/smoothing)

- ▶ **Given:** parameters  $\theta = (\pi, \mathbf{A}, \mathbf{B})$ , observation sequence  $x_{1:\ell} \in [D]^\ell$ .
- ▶ **Goal:** conditional distribution of  $Y_{s:t}$  given  $X_{1:\ell} = x_{1:\ell}$  ( $1 \leq s \leq t \leq \ell$ ):

$$\Pr_{\theta}(Y_{s:t} = y_{s:t} \mid X_{1:\ell} = x_{1:\ell}), \quad \text{for each } y_{s:t} \in [K]^{t-s+1}.$$

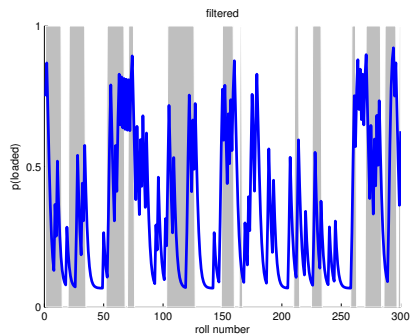
## Most probable state sequence (decoding)

- ▶ **Given:** parameters  $\theta = (\pi, \mathbf{A}, \mathbf{B})$ , observation sequence  $x_{1:\ell} \in [D]^\ell$ .
- ▶ **Goal:**  $\arg \max_{y_{1:\ell} \in [K]^\ell} \Pr_{\theta}(Y_{1:\ell} = y_{1:\ell} \mid X_{1:\ell} = x_{1:\ell})$ .

## Parameter estimation

- ▶ **Given:**  $n$  observation sequences  $x_{1:\ell}^{(s)}$  for  $s \in [n]$ .
- ▶ **Goal:** parameter estimates  $\hat{\theta} = (\hat{\pi}, \hat{\mathbf{A}}, \hat{\mathbf{B}})$ .

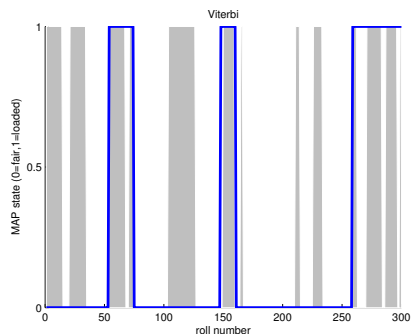
# EXAMPLE: DISHONEST CASINO



## Conditional probability

Gray bars: Loaded dice used.

Blue:  $\Pr_{\theta}(Y_t = \text{loaded} | X_{1:\ell} = x_{1:\ell})$



## Decoding

Gray bars: Loaded dice used.

Blue: Most probable state  $Z_t$ .

# SOME APPLICATIONS

- ▶ **Bioinformatics**

*Observations:* amino acids in a protein

*Hidden states:* indicators of evolutionary conservation

# SOME APPLICATIONS

- ▶ **Bioinformatics**

*Observations:* amino acids in a protein

*Hidden states:* indicators of evolutionary conservation

- ▶ **Natural language processing**

*Observations:* words in a sentence

*Hidden states:* words' part-of-speech or other word-type semantics

# SOME APPLICATIONS

- ▶ **Bioinformatics**

*Observations:* amino acids in a protein

*Hidden states:* indicators of evolutionary conservation

- ▶ **Natural language processing**

*Observations:* words in a sentence

*Hidden states:* words' part-of-speech or other word-type semantics

- ▶ **Speech recognition**

*Observations:* recorded speech at various (discrete) times

*Hidden states:* phonemes that the speaker intended to vocalize

- ▶ **Financial market cycles forecast:**

*Observations:* index of stock market

*Hidden states:* bull or bear market