# COMS 4771 Lecture 22

1. Hidden Markov models

# HIDDEN MARKOV MODELS

# 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 $\{Y_\tau : \tau \le t\}$ 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 $[K] = \{1, 2, \ldots, K\}$):

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

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

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

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

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

**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 $[K] = \{1, 2, \ldots, K\}$) (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 $[K] = \{1, 2, \ldots, K\}$) (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 $[K] = \{1, 2, \ldots, K\}$) (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 $[K] = \{1, 2, \ldots, K\}$) (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*.

$$\cdots \longrightarrow Y_{t-1} \longrightarrow \ Y_t \ \longrightarrow Y_{t+1} \longrightarrow \cdots$$
$$\downarrow \qquad \quad \downarrow \qquad \quad \downarrow$$
$$X_{t-1} \qquad X_t \qquad X_{t+1}$$

# 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 $[K] = \{1, 2, \ldots, K\}$) (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*.

$$\cdots \longrightarrow Y_{t-1} \longrightarrow Y_t \longrightarrow Y_{t+1} \longrightarrow \cdots$$
$$\downarrow \qquad\quad \downarrow \qquad\quad \downarrow$$
$$X_{t-1} \qquad X_t \qquad X_{t+1}$$

**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 $[K] = \{1, 2, \ldots, K\}$) (hidden state sequence);

- ▶ conditioned on $Y_t$, corresponding $X_t$ is *independent of all other variables*;

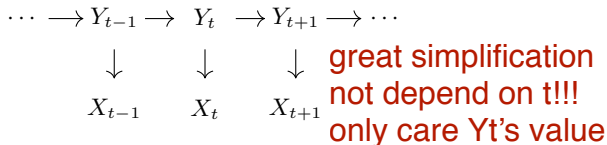- ▶ the $Y_t$ are *hidden*, and the $X_t$ are *observed*.

$$\cdots \longrightarrow Y_{t-1} \longrightarrow Y_t \longrightarrow Y_{t+1} \longrightarrow \cdots$$
$$\downarrow \qquad \downarrow \qquad \downarrow$$
$$X_{t-1} \qquad X_t \qquad X_{t+1}$$

great simplification
not depend on t!!!
only care Yt's value

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

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

# HMM parameters (discrete observations)

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

# HMM PARAMETERS (DISCRETE OBSERVATIONS)

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

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

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

# HMM parameters (discrete observations)

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

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

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

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

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

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

# HMM PARAMETERS (DISCRETE OBSERVATIONS)

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

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

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

▶ **Transition matrix**: $K \times K$ matrix $\boldsymbol{A}$     <span style="color:red">underlying thing</span>

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

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

▶ **Emission matrix**: $K \times D$ matrix $\boldsymbol{B}$     <span style="color:red">emission thing</span>

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

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

**Mixture model**

$$Y$$

$$\downarrow$$

$$\boldsymbol{X}$$

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

**Hidden Markov model**

$$Y_1 \longrightarrow Y_2 \longrightarrow \cdots \longrightarrow Y_\ell$$

$$\downarrow \qquad \downarrow \qquad\qquad\qquad \downarrow$$

$$X_1 \qquad X_2 \qquad\qquad\qquad X_\ell$$

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

| Mixture model | Hidden Markov model |
|---|---|
| $Y$ | $Y_1 \longrightarrow Y_2 \longrightarrow \cdots \longrightarrow Y_\ell$ |
| $\downarrow$ | $\downarrow \qquad \downarrow \qquad\qquad\qquad \downarrow$ |
| $\boldsymbol{X}$ | $X_1 \qquad X_2 \qquad\qquad\qquad X_\ell$ |

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

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

| Mixture model | Hidden Markov model |
|:---:|:---:|
| $Y$ | $Y_1 \longrightarrow Y_2 \longrightarrow \cdots \longrightarrow Y_\ell$ |
| $\downarrow$ | $\downarrow \quad\quad \downarrow \quad\quad\quad\quad\quad\quad \downarrow$ |
| $\boldsymbol{X}$ | $X_1 \quad\quad X_2 \quad\quad\quad\quad\quad X_\ell$ |
| ($Y$ is hidden, $\boldsymbol{X}$ is observed.) | ($Y_{1:\ell}$ is hidden, $X_{1:\ell}$ is observed.) |
| For $K$ component mixture model, $Y$ takes values in $[K]$. | For sequence of length $\ell$, $Y_{1:\ell}$ takes values in $[K]^\ell$. |

# CONNECTIONS TO MIXTURE MODELS

|  **Mixture model**  |  **Hidden Markov model**  |
|:---:|:---:|
| $Y$ | $Y_1 \rightarrow Y_2 \longrightarrow \cdots \longrightarrow Y_\ell$ |
| $\downarrow$ | $\downarrow \qquad \downarrow \qquad\qquad\qquad \downarrow$ |
| $\boldsymbol{X}$ | $X_1 \qquad X_2 \qquad\qquad\qquad X_\ell$ |

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

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

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

| **Mixture model** | **Hidden Markov model** |
|---|---|

$$Y \qquad\qquad Y_1 \;\rightarrow\; Y_2 \;\longrightarrow\; \cdots \;\longrightarrow\; Y_\ell$$

$$\downarrow \qquad\qquad \downarrow \qquad \downarrow \qquad\qquad\qquad \downarrow$$

$$\boldsymbol{X} \qquad\qquad X_1 \qquad X_2 \qquad\qquad X_\ell$$

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

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

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

▶ $Y_1 \to Y_2 \to X_2$

| **Mixture model** | **Hidden Markov model** |
|---|---|

$$Y \qquad\qquad Y_1 \;\rightarrow\; Y_2 \;\longrightarrow\; \cdots \;\longrightarrow\; Y_\ell$$

$$\downarrow \qquad\qquad \downarrow \qquad \downarrow \qquad\qquad\qquad \downarrow$$

$$\boldsymbol{X} \qquad\qquad X_1 \qquad X_2 \qquad\qquad\qquad X_\ell$$

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

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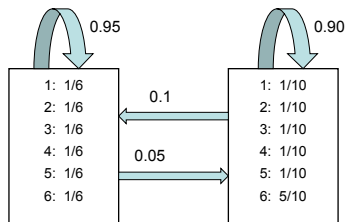      $Y$ takes values in $[K]$.         $Y_{1:\ell}$ takes values in $[K]^\ell$.

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**

<span style="color:red">sub-distribution</span>

$$Y$$
$$\downarrow$$
$$\boldsymbol{X}$$

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

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

<span style="color:red">hidden state(sub-distribution)</span>

**Hidden Markov model**

$$Y_1 \longrightarrow Y_2 \longrightarrow \cdots \longrightarrow Y_\ell$$
$$\downarrow \qquad \downarrow \qquad\qquad\qquad \downarrow$$
$$X_1 \qquad X_2 \qquad\qquad\qquad X_\ell$$

($Y_{1:\ell}$ 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!

- $Y_1 \to Y_2 \to X_2$
- $X_2 \to Y_2 \to Y_3 \to X_3$
- $X_1 \to Y_1 \to Y_{2:\ell} \to X_{2:\ell}$
- $\cdots$

<span style="color:red">note the relationship between chain!</span>

**Casino die-rolling game**:

Randomly switch between two possible dice: one is fair, the other loaded.
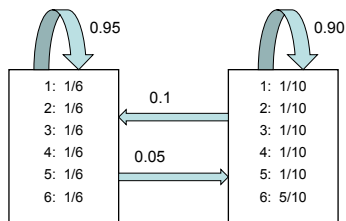
The dice are otherwise indistinguishable!

**Casino die-rolling game**:

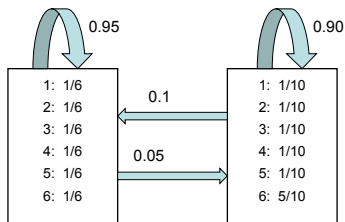Randomly switch between two possible dice: one is fair, the other loaded.

The dice are otherwise indistinguishable!

HMM parameters:

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

and $\boldsymbol{\pi} = (1, 0)$ 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:

emission matrix

$$\boldsymbol{A} = \begin{array}{c} \text{fair die} \\ \text{loaded die} \end{array} \overset{\begin{array}{cc} \text{fair die} & \text{loaded die} \end{array}}{\begin{pmatrix} 0.95 & 0.05 \\ 0.10 & 0.90 \end{pmatrix}}, \quad \boldsymbol{B} = \begin{array}{c} \text{fair die} \\ \text{loaded die} \end{array} \overset{\begin{array}{cccccc} 1 & 2 & 3 & 4 & 5 & 6 \end{array}}{\begin{pmatrix} \frac{1}{6} & \frac{1}{6} & \frac{1}{6} & \frac{1}{6} & \frac{1}{6} & \frac{1}{6} \\ \frac{1}{10} & \frac{1}{10} & \frac{1}{10} & \frac{1}{10} & \frac{1}{10} & \frac{1}{2} \end{pmatrix}},$$

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

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

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

▶ **Given**: parameters $\boldsymbol{\theta} = (\boldsymbol{\pi}, \boldsymbol{A}, \boldsymbol{B})$, observation sequence $x_{1:\ell} \in [D]^{\ell}$.

▶ **Goal**: conditional distribution of $Y_{s:t}$ given $X_{1:\ell} = x_{1:\ell}$ ($1 \leq s \leq t \leq \ell$):

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

# HMM INFERENCE/LEARNING PROBLEMS

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

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

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

### Most probable state sequence (decoding)

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

may be the current yt, thus we could predict next state

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

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

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

Most probable state sequence (decoding)

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

Parameter estimation

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

**Conditional probability**

Gray bars: Loaded dice used.

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

**Decoding**

Gray bars: Loaded dice used.

Blue: Most probable state $Z_t$.

▶ **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

# HMM PROBABILITY COMPUTATIONS

HMM can be used for sequences of **arbitrary length**.

HMM can be used for sequences of **arbitrary length**.

**Subtle difficulty**: The most straightforward formulae for

- observation sequence probabilities (e.g., $\Pr_{\boldsymbol{\theta}}(X_{1:\ell} = x_{1:\ell})$) and
- conditional state sequence probabilities (e.g., $\Pr_{\boldsymbol{\theta}}(Y_\ell = y_\ell \mid X_{1:\ell} = x_{1:\ell})$)

do not suggest efficient algorithms for computation.

# PROBABILITY COMPUTATIONS

HMM can be used for sequences of **arbitrary length**.

**Subtle difficulty**: The most straightforward formulae for

- observation sequence probabilities (e.g., $\Pr_{\boldsymbol{\theta}}(X_{1:\ell} = x_{1:\ell})$) and
- conditional state sequence probabilities (e.g., $\Pr_{\boldsymbol{\theta}}(Y_\ell = y_\ell \mid X_{1:\ell} = x_{1:\ell})$)

do not suggest efficient algorithms for computation.

**Need to exploit special structure of HMMs to get efficient algorithms.**

**Probability of observation sequence** $\Pr_{\boldsymbol{\theta}}(X_{1:\ell} = x_{1:\ell})$.

**Probability of observation sequence** $\Pr_{\boldsymbol{\theta}}(X_{1:\ell} = x_{1:\ell})$.

We'll use the fact that conditioned on $Y_{1:\ell}$, the $\{X_t\}_{t \in [\ell]}$ are independent.

**Probability of observation sequence** $\Pr_{\boldsymbol{\theta}}(X_{1:\ell} = x_{1:\ell})$.

We'll use the fact that conditioned on $Y_{1:\ell}$, the $\{X_t\}_{t \in [\ell]}$ are independent.

$$\Pr_{\boldsymbol{\theta}}(X_{1:\ell} = x_{1:\ell}) \quad = \quad \sum_{y_{1:\ell} \in [K]^{\ell}} \Pr_{\boldsymbol{\theta}}(X_{1:\ell} = x_{1:\ell} \wedge Y_{1:\ell} = y_{1:\ell})$$

**Probability of observation sequence** $\Pr_{\boldsymbol{\theta}}(X_{1:\ell} = x_{1:\ell})$.

We'll use the fact that conditioned on $Y_{1:\ell}$, the $\{X_t\}_{t \in [\ell]}$ are independent.

$$
\begin{aligned}
\Pr_{\boldsymbol{\theta}}(X_{1:\ell} = x_{1:\ell}) &= \sum_{y_{1:\ell} \in [K]^{\ell}} \Pr_{\boldsymbol{\theta}}(X_{1:\ell} = x_{1:\ell} \wedge Y_{1:\ell} = y_{1:\ell}) \\
&= \sum_{y_{1:\ell} \in [K]^{\ell}} \Pr_{\boldsymbol{\theta}}(Y_{1:\ell} = y_{1:\ell}) \cdot \Pr_{\boldsymbol{\theta}}\big(X_{1:\ell} = x_{1:\ell} \,|\, Y_{1:\ell} = y_{1:\ell}\big)
\end{aligned}
$$

**Probability of observation sequence** $\Pr_{\boldsymbol{\theta}}(X_{1:\ell} = x_{1:\ell})$.

We'll use the fact that conditioned on $Y_{1:\ell}$, the $\{X_t\}_{t \in [\ell]}$ are independent.

$$
\begin{aligned}
\Pr_{\boldsymbol{\theta}}(X_{1:\ell} = x_{1:\ell}) &= \sum_{y_{1:\ell} \in [K]^\ell} \Pr_{\boldsymbol{\theta}}(X_{1:\ell} = x_{1:\ell} \wedge Y_{1:\ell} = y_{1:\ell}) \\
&= \sum_{y_{1:\ell} \in [K]^\ell} \Pr_{\boldsymbol{\theta}}(Y_{1:\ell} = y_{1:\ell}) \cdot \Pr_{\boldsymbol{\theta}}\big(X_{1:\ell} = x_{1:\ell} \,|\, Y_{1:\ell} = y_{1:\ell}\big) \\
&= \sum_{y_{1:\ell} \in [K]^\ell} \Pr_{\boldsymbol{\theta}}(Y_{1:\ell} = y_{1:\ell}) \cdot \prod_{t=1}^{\ell} \Pr_{\boldsymbol{\theta}}\big(X_t = x_t \,|\, Y_t = y_t\big)
\end{aligned}
$$

Given Y1:l, {Xt} are independent.

**Probability of observation sequence** $\Pr_{\boldsymbol{\theta}}(X_{1:\ell} = x_{1:\ell})$.

We'll use the fact that conditioned on $Y_{1:\ell}$, the $\{X_t\}_{t \in [\ell]}$ are independent.

$$
\begin{aligned}
\Pr_{\boldsymbol{\theta}}(X_{1:\ell} = x_{1:\ell}) \quad &= \sum_{y_{1:\ell} \in [K]^{\ell}} \Pr_{\boldsymbol{\theta}}(X_{1:\ell} = x_{1:\ell} \wedge Y_{1:\ell} = y_{1:\ell}) \\
&= \sum_{y_{1:\ell} \in [K]^{\ell}} \Pr_{\boldsymbol{\theta}}(Y_{1:\ell} = y_{1:\ell}) \cdot \Pr_{\boldsymbol{\theta}}(X_{1:\ell} = x_{1:\ell} \mid Y_{1:\ell} = y_{1:\ell}) \\
&\qquad\qquad \text{this step, assume yt is know for current sequence} \\
&= \sum_{y_{1:\ell} \in [K]^{\ell}} \Pr_{\boldsymbol{\theta}}(Y_{1:\ell} = y_{1:\ell}) \cdot \prod_{t=1} \Pr_{\boldsymbol{\theta}}(X_t = x_t \mid Y_t = y_t)
\end{aligned}
$$

**But summation is over $K^{\ell}$ terms—seems intractable for large $\ell$.**

**Probability of observation sequence** $\Pr_{\boldsymbol{\theta}}(X_{1:\ell} = x_{1:\ell})$.

We'll use the fact that conditioned on $Y_{1:\ell}$, the $\{X_t\}_{t \in [\ell]}$ are independent.

$$
\begin{aligned}
\Pr_{\boldsymbol{\theta}}(X_{1:\ell} = x_{1:\ell}) \quad &= \sum_{y_{1:\ell} \in [K]^\ell} \Pr_{\boldsymbol{\theta}}(X_{1:\ell} = x_{1:\ell} \wedge Y_{1:\ell} = y_{1:\ell}) \\
&= \sum_{y_{1:\ell} \in [K]^\ell} \Pr_{\boldsymbol{\theta}}(Y_{1:\ell} = y_{1:\ell}) \cdot \Pr_{\boldsymbol{\theta}}(X_{1:\ell} = x_{1:\ell} \mid Y_{1:\ell} = y_{1:\ell}) \\
&= \sum_{y_{1:\ell} \in [K]^\ell} \Pr_{\boldsymbol{\theta}}(Y_{1:\ell} = y_{1:\ell}) \cdot \prod_{t=1}^{\ell} \Pr_{\boldsymbol{\theta}}(X_t = x_t \mid Y_t = y_t)
\end{aligned}
$$

marigina

**But summation is over $K^\ell$ terms—seems intractable for large $\ell$.**

**Fortunately, the summation can be computed iteratively in time linear in $\ell$.**

HMM parameters:
$$\pi_y = \Pr(Y_1 = y); \quad A_{y,z} = \Pr(Y_{t+1} = z \mid Y_t = y); \quad B_{y,x} = \Pr(X_t = x \mid Y_t = y).$$

**Example: probability of observation triplet** $x_{1:3} \in [D]^3$

$$\Pr_{\boldsymbol{\theta}}(X_{1:3} = x_{1:3})$$

# SEQUENCE PROBABILITY COMPUTATION

HMM parameters:
$\pi y = \Pr(Y_1 = y); \ Ay, z = \Pr(Y_{t+1} = z \mid Y_t = y); \ By, x = \Pr(X_t = x \mid Y_t = y).$

**Example: probability of observation triplet** $x_{1:3} \in [D]^3$

$$\Pr_{\boldsymbol{\theta}}(X_{1:3} = x_{1:3})$$

$$= \sum_{y_{1:3} \in [K]^3} \Pr_{\boldsymbol{\theta}}(Y_{1:3} = y_{1:3}) \cdot \prod_{t=1}^{3} \Pr_{\boldsymbol{\theta}}(X_t = x_t \mid Y_t = y_t)$$

# Sequence probability computation

HMM parameters:
$$\pi_y = \Pr(Y_1 = y); \quad A_{y,z} = \Pr(Y_{t+1} = z \mid Y_t = y); \quad B_{y,x} = \Pr(X_t = x \mid Y_t = y).$$

**Example: probability of observation triplet** $x_{1:3} \in [D]^3$

$$
\begin{aligned}
&\Pr_{\boldsymbol{\theta}}(X_{1:3} = x_{1:3}) \\
&= \sum_{y_{1:3} \in [K]^3} \Pr_{\boldsymbol{\theta}}(Y_{1:3} = y_{1:3}) \cdot \prod_{t=1}^{3} \Pr_{\boldsymbol{\theta}}\left(X_t = x_t \mid Y_t = y_t\right) \\
&= \sum_{y_{1:3} \in [K]^3} \underbrace{\left(\pi_{y_1} \cdot A_{y_1, y_2} \cdot A_{y_2, y_3}\right)}_{\text{Markov chain probabilities}} \cdot \underbrace{\left(B_{y_1, x_1} \cdot B_{y_2, x_2} \cdot B_{y_3, x_3}\right)}_{\text{emission probabilities}}
\end{aligned}
$$

HMM parameters:
$$\pi_y = \Pr(Y_1 = y); \quad A_{y,z} = \Pr(Y_{t+1} = z \mid Y_t = y); \quad B_{y,x} = \Pr(X_t = x \mid Y_t = y).$$

**Example: probability of observation triplet** $x_{1:3} \in [D]^3$

$$\Pr_{\boldsymbol{\theta}}(X_{1:3} = x_{1:3})$$

$$= \sum_{y_{1:3} \in [K]^3} \Pr_{\boldsymbol{\theta}}(Y_{1:3} = y_{1:3}) \cdot \prod_{t=1}^{3} \Pr_{\boldsymbol{\theta}}(X_t = x_t \mid Y_t = y_t)$$

$$= \sum_{y_{1:3} \in [K]^3} \underbrace{(\pi_{y_1} \cdot A_{y_1, y_2} \cdot A_{y_2, y_3})}_{\text{Markov chain probabilities}} \cdot \underbrace{(B_{y_1, x_1} \cdot B_{y_2, x_2} \cdot B_{y_3, x_3})}_{\text{emission probabilities}}$$

$$= \underbrace{\sum_{y_1 \in [K]} \pi_{y_1} \cdot B_{y_1, x_1}}_{O(K) \text{ time}} \underbrace{\sum_{y_2 \in [K]} A_{y_1, y_2} \cdot B_{y_2, x_2}}_{O(K) \text{ time for each } y_1 \in [K]} \underbrace{\sum_{y_3 \in [K]} A_{y_2, y_3} \cdot B_{y_3, x_3}}_{O(K) \text{ time for each } y_2 \in [K]}$$

HMM parameters:
$$\pi_y = \Pr(Y_1 = y); \quad A_{y,z} = \Pr(Y_{t+1} = z \,|\, Y_t = y); \quad B_{y,x} = \Pr(X_t = x \,|\, Y_t = y).$$

**Example: probability of observation triplet** $x_{1:3} \in [D]^3$

the beautiful part
there is a
transformation from
any i to any j
thus we could reduce
through following form

$$\Pr_{\boldsymbol{\theta}}(X_{1:3} = x_{1:3})$$

$$= \sum_{y_{1:3} \in [K]^3} \Pr_{\boldsymbol{\theta}}(Y_{1:3} = y_{1:3}) \cdot \prod_{t=1}^{3} \Pr_{\boldsymbol{\theta}}(X_t = x_t \,|\, Y_t = y_t)$$

note the x
is fixed!!!

$$= \sum_{y_{1:3} \in [K]^3} \underbrace{(\pi_{y_1} \cdot A_{y_1, y_2} \cdot A_{y_2, y_3})}_{\text{Markov chain probabilities}} \cdot \underbrace{(B_{y_1, x_1} \cdot B_{y_2, x_2} \cdot B_{y_3, x_3})}_{\text{emission probabilities}}$$

$$= \underbrace{\sum_{y_1 \in [K]} \pi_{y_1} \cdot B_{y_1, x_1}}_{O(K) \text{ time}} \underbrace{\sum_{y_2 \in [K]} A_{y_1, y_2} \cdot B_{y_2, x_2}}_{O(K) \text{ time for each } y_1 \in [K]} \underbrace{\sum_{y_3 \in [K]} A_{y_2, y_3} \cdot B_{y_3, x_3}}_{O(K) \text{ time for each } y_2 \in [K]}$$

**Computing sums from right-to-left**: total time is $O(K^2 \ell)$ for length $\ell$.

note the complexity: constrain to its own component now!!!

# Conditional probabilities

**A simple case**: $\Pr_{\boldsymbol{\theta}}\big(Y_\ell = y_\ell \,|\, X_{1:\ell} = x_{1:\ell}\big) \; = \; \dfrac{\Pr_{\boldsymbol{\theta}}(X_{1:\ell} = x_{1:\ell} \,\wedge\, Y_\ell = y_\ell)}{\Pr_{\boldsymbol{\theta}}(X_{1:\ell} = x_{1:\ell})}.$

**A simple case**: $\Pr_{\boldsymbol{\theta}}\big(Y_\ell = y_\ell \,|\, X_{1:\ell} = x_{1:\ell}\big) \;=\; \dfrac{\Pr_{\boldsymbol{\theta}}(X_{1:\ell} = x_{1:\ell} \,\wedge\, Y_\ell = y_\ell)}{\Pr_{\boldsymbol{\theta}}(X_{1:\ell} = x_{1:\ell})}.$

$$
\begin{aligned}
&\Pr_{\boldsymbol{\theta}}(X_{1:\ell} = x_{1:\ell}) \\
&= \sum_{y_1 \in [K]} \pi_{y_1} \cdot B_{y_1, x_1} \sum_{y_2 \in [K]} A_{y_1, y_2} \cdot B_{y_2, x_2} \cdots \\
&\qquad \cdots \sum_{y_{\ell-1} \in [K]} A_{y_{\ell-2}, y_{\ell-1}} \cdot B_{y_{\ell-1}, x_{\ell-1}} \sum_{y_\ell \in [K]} A_{y_{\ell-1}, y_\ell} \cdot B_{y_\ell, x_\ell}.
\end{aligned}
$$

**A simple case**: $\Pr_{\boldsymbol{\theta}}\big(Y_\ell = y_\ell \,|\, X_{1:\ell} = x_{1:\ell}\big) \;=\; \dfrac{\Pr_{\boldsymbol{\theta}}(X_{1:\ell} = x_{1:\ell} \,\wedge\, Y_\ell = y_\ell)}{\Pr_{\boldsymbol{\theta}}(X_{1:\ell} = x_{1:\ell})}.$

$$\Pr_{\boldsymbol{\theta}}(X_{1:\ell} = x_{1:\ell} \,\wedge\, Y_\ell = y_\ell)$$

$$= \sum_{y_1 \in [K]} \pi_{y_1} \cdot B_{y_1, x_1} \sum_{y_2 \in [K]} A_{y_1, y_2} \cdot B_{y_2, x_2} \cdots$$

$$\cdots \sum_{y_{\ell-1} \in [K]} A_{y_{\ell-2}, y_{\ell-1}} \cdot B_{y_{\ell-1}, x_{\ell-1}} \qquad A_{y_{\ell-1}, y_\ell} \cdot B_{y_\ell, x_\ell}.$$

**A simple case**: $\Pr_{\boldsymbol{\theta}}\big(Y_\ell = y_\ell \mid X_{1:\ell} = x_{1:\ell}\big) \;=\; \dfrac{\Pr_{\boldsymbol{\theta}}(X_{1:\ell} = x_{1:\ell} \,\wedge\, Y_\ell = y_\ell)}{\Pr_{\boldsymbol{\theta}}(X_{1:\ell} = x_{1:\ell})}.$

$$
\begin{aligned}
&\Pr_{\boldsymbol{\theta}}(X_{1:\ell} = x_{1:\ell} \,\wedge\, Y_\ell = y_\ell) \\
&= \sum_{y_1 \in [K]} \pi_{y_1} \cdot B_{y_1}, x_1 \sum_{y_2 \in [K]} A_{y_1}, y_2 \cdot B_{y_2}, x_2 \cdots \\
&\quad \cdots \sum_{y_{\ell-1} \in [K]} A_{y_{\ell-2}}, y_{\ell-1} \cdot B_{y_{\ell-1}}, x_{\ell-1} \qquad A_{y_{\ell-1}}, y_\ell \cdot B_{y_\ell}, x_\ell.
\end{aligned}
$$

**Forward inductive computation**:
Keep track of $\alpha_t(y_t) := \Pr_{\boldsymbol{\theta}}(X_{1:t} = x_{1:t} \,\wedge\, Y_t = y_t)$ for each $y_t \in [K]$.

# CONDITIONAL PROBABILITIES

**A simple case**: $\Pr_{\boldsymbol{\theta}}\big(Y_\ell = y_\ell \mid X_{1:\ell} = x_{1:\ell}\big) = \dfrac{\Pr_{\boldsymbol{\theta}}(X_{1:\ell} = x_{1:\ell} \wedge Y_\ell = y_\ell)}{\Pr_{\boldsymbol{\theta}}(X_{1:\ell} = x_{1:\ell})}.$

$$\Pr_{\boldsymbol{\theta}}(X_{1:\ell} = x_{1:\ell} \wedge Y_\ell = y_\ell)$$

$$= \sum_{y_1 \in [K]} \pi_{y_1} \cdot B_{y_1, x_1} \sum_{y_2 \in [K]} A_{y_1, y_2} \cdot B_{y_2, x_2} \cdots$$

note here, yl is fixed! to get xl

$$\cdots \sum_{y_{\ell-1} \in [K]} A_{y_{\ell-2}, y_{\ell-1}} \cdot B_{y_{\ell-1}, x_{\ell-1}} \qquad \boxed{A_{y_{\ell-1}, y_\ell} \cdot B_{y_\ell, x_\ell} \cdot}$$

**Forward inductive computation**:

Keep track of $\alpha_t(y_t) := \Pr_{\boldsymbol{\theta}}(X_{1:t} = x_{1:t} \wedge Y_t = y_t)$ for each $y_t \in [K]$.

Compute $\alpha_{t+1}$ using $\alpha_t$ in $\boxed{O(K^2) \text{ time:}}$    something amazing is happening!

$$\alpha_{t+1}(y_{t+1}) = \left( \sum_{y_t \in [K]} \alpha_t(y_t) \cdot A_{y_t, y_{t+1}} \right) \boxed{\cdot B_{y_{t+1}, x_{t+1}}} \text{ for each } y_{t+1} \in [K].$$

actually, need to compute for each y(t+1), thus the cost is O(k^2)

powerful! each pair from yt to yt+1

For any $1 \le t < \ell$,

$$\Pr_{\boldsymbol{\theta}}\big(Y_t = y_t \mid X_{1:\ell} = x_{1:\ell}\big)$$
$$= \frac{\Pr_{\boldsymbol{\theta}}(X_{1:t} = x_{1:t} \wedge Y_t = y_t) \cdot \Pr_{\boldsymbol{\theta}}\big(X_{t+1:\ell} = x_{t+1:\ell} \mid Y_t = y_t\big)}{\Pr_{\boldsymbol{\theta}}(X_{1:\ell} = x_{1:\ell})}$$

(since $X_{t+1:\ell}$ is conditionally independent of $X_{1:t}$ given $Y_t$)

For any $1 \le t < \ell$,

$$\Pr_{\boldsymbol{\theta}}\big(Y_t = y_t \,|\, X_{1:\ell} = x_{1:\ell}\big)$$

$$= \frac{\Pr_{\boldsymbol{\theta}}(X_{1:t} = x_{1:t} \,\wedge\, Y_t = y_t) \cdot \Pr_{\boldsymbol{\theta}}\big(X_{t+1:\ell} = x_{t+1:\ell} \,|\, Y_t = y_t\big)}{\Pr_{\boldsymbol{\theta}}(X_{1:\ell} = x_{1:\ell})}$$

(since $X_{t+1:\ell}$ is conditionally independent of $X_{1:t}$ given $Y_t$)

$$= \frac{\alpha_t(y_t) \cdot \beta_t(y_t)}{\text{normalization term}}$$

where

$$\alpha_t(y_t) := \Pr_{\boldsymbol{\theta}}(X_{1:t} = x_{1:t} \,\wedge\, Y_t = y_t),$$
$$\beta_t(y_t) := \Pr_{\boldsymbol{\theta}}\big(X_{t+1:\ell} = x_{t+1:\ell} \,|\, Y_t = y_t\big).$$

For any $1 \le t < \ell$,

$$\Pr_{\boldsymbol{\theta}}\big(Y_t = y_t \mid X_{1:\ell} = x_{1:\ell}\big)$$

$$= \frac{\Pr_{\boldsymbol{\theta}}(X_{1:t} = x_{1:t} \wedge Y_t = y_t) \cdot \Pr_{\boldsymbol{\theta}}\big(X_{t+1:\ell} = x_{t+1:\ell} \mid Y_t = y_t\big)}{\Pr_{\boldsymbol{\theta}}(X_{1:\ell} = x_{1:\ell})}$$

(since $X_{t+1:\ell}$ is conditionally independent of $X_{1:t}$ given $Y_t$)

$$= \frac{\alpha_t(y_t) \cdot \beta_t(y_t)}{\text{normalization term}}$$

where

$$\alpha_t(y_t) := \Pr_{\boldsymbol{\theta}}(X_{1:t} = x_{1:t} \wedge Y_t = y_t),$$
$$\beta_t(y_t) := \Pr_{\boldsymbol{\theta}}\big(X_{t+1:\ell} = x_{t+1:\ell} \mid Y_t = y_t\big).$$

We already saw how to compute $\alpha_t(y_t)$ for each $y_t \in [K]$.

$$\beta_t(y_t) = \Pr_{\boldsymbol{\theta}}\big(X_{t+1:\ell} = x_{t+1:\ell} \,|\, Y_t = y_t\big)$$

$$
\begin{aligned}
\beta_t(y_t) &= \mathrm{Pr}_{\boldsymbol{\theta}}\big(X_{t+1:\ell} = x_{t+1:\ell} \,|\, Y_t = y_t\big) \\
&= \sum_{y_{t+1} \in [K]} \mathrm{Pr}_{\boldsymbol{\theta}}\big(X_{t+1:\ell} = x_{t+1:\ell} \,\wedge\, Y_{t+1} = y_{t+1} \,|\, Y_t = y_t\big)
\end{aligned}
$$

$$
\begin{aligned}
\beta_t(y_t) &= \Pr_{\boldsymbol{\theta}}\big(X_{t+1:\ell} = x_{t+1:\ell} \,|\, Y_t = y_t\big) \\
&= \sum_{y_{t+1} \in [K]} \Pr_{\boldsymbol{\theta}}\big(X_{t+1:\ell} = x_{t+1:\ell} \,\wedge\, Y_{t+1} = y_{t+1} \,|\, Y_t = y_t\big) \\
&= \sum_{y_{t+1} \in [K]} \Pr_{\boldsymbol{\theta}}\big(Y_{t+1} = y_{t+1} \,|\, Y_t = y_t\big) \\
&\qquad\qquad \cdot \Pr_{\boldsymbol{\theta}}\big(X_{t+1:\ell} = x_{t+1:\ell} \,|\, Y_{t+1} = y_{t+1}\big)
\end{aligned}
$$

$$\beta_t(y_t) = \Pr_{\boldsymbol{\theta}}\big(X_{t+1:\ell} = x_{t+1:\ell} \,|\, Y_t = y_t\big)$$

$$= \sum_{y_{t+1}\in[K]} \Pr_{\boldsymbol{\theta}}\big(X_{t+1:\ell} = x_{t+1:\ell} \,\wedge\, Y_{t+1} = y_{t+1} \,|\, Y_t = y_t\big)$$

$$= \sum_{y_{t+1}\in[K]} \Pr_{\boldsymbol{\theta}}\big(Y_{t+1} = y_{t+1} \,|\, Y_t = y_t\big)$$

$$\cdot \Pr_{\boldsymbol{\theta}}\big(X_{t+1:\ell} = x_{t+1:\ell} \,|\, Y_{t+1} = y_{t+1}\big)$$

$$= \sum_{y_{t+1}\in[K]} \Pr_{\boldsymbol{\theta}}\big(Y_{t+1} = y_{t+1} \,|\, Y_t = y_t\big)$$

$$\cdot \Pr_{\boldsymbol{\theta}}\big(X_{t+1} = x_{t+1} \,|\, Y_{t+1} = y_{t+1}\big) \cdot \Pr_{\boldsymbol{\theta}}\big(X_{t+2:\ell} = x_{t+2:\ell} \,|\, Y_{t+1} = y_{t+1}\big)$$

$$
\begin{aligned}
\beta_t(y_t) &= \Pr_{\boldsymbol{\theta}}\big(X_{t+1:\ell} = x_{t+1:\ell} \,|\, Y_t = y_t\big) \\
&= \sum_{y_{t+1} \in [K]} \Pr_{\boldsymbol{\theta}}\big(X_{t+1:\ell} = x_{t+1:\ell} \,\wedge\, Y_{t+1} = y_{t+1} \,|\, Y_t = y_t\big) \\
&= \sum_{y_{t+1} \in [K]} \Pr_{\boldsymbol{\theta}}\big(Y_{t+1} = y_{t+1} \,|\, Y_t = y_t\big) \\
&\qquad\qquad \cdot \Pr_{\boldsymbol{\theta}}\big(X_{t+1:\ell} = x_{t+1:\ell} \,|\, Y_{t+1} = y_{t+1}\big) \\
&= \sum_{y_{t+1} \in [K]} \Pr_{\boldsymbol{\theta}}\big(Y_{t+1} = y_{t+1} \,|\, Y_t = y_t\big) \\
&\qquad \cdot \Pr_{\boldsymbol{\theta}}\big(X_{t+1} = x_{t+1} \,|\, Y_{t+1} = y_{t+1}\big) \cdot \Pr_{\boldsymbol{\theta}}\big(X_{t+2:\ell} = x_{t+2:\ell} \,|\, Y_{t+1} = y_{t+1}\big) \\
&= \sum_{y_{t+1} \in [K]} A_{y_t, y_{t+1}} \cdot B_{y_{t+1}, x_{t+1}} \cdot \beta_{t+1}(y_{t+1}).
\end{aligned}
$$

# CONDITIONAL PROBABILITIES FOR $t < \ell$

$$
\begin{aligned}
\beta_t(y_t) &= \Pr_{\boldsymbol{\theta}}\big(X_{t+1:\ell} = x_{t+1:\ell} \,\big|\, Y_t = y_t\big) \\
&= \sum_{y_{t+1} \in [K]} \Pr_{\boldsymbol{\theta}}\big(X_{t+1:\ell} = x_{t+1:\ell} \,\wedge\, Y_{t+1} = y_{t+1} \,\big|\, Y_t = y_t\big) \\
&= \sum_{y_{t+1} \in [K]} \Pr_{\boldsymbol{\theta}}\big(Y_{t+1} = y_{t+1} \,\big|\, Y_t = y_t\big) \\
&\qquad\qquad \cdot \Pr_{\boldsymbol{\theta}}\big(X_{t+1:\ell} = x_{t+1:\ell} \,\big|\, Y_{t+1} = y_{t+1}\big) \\
&= \sum_{y_{t+1} \in [K]} \Pr_{\boldsymbol{\theta}}\big(Y_{t+1} = y_{t+1} \,\big|\, Y_t = y_t\big) \\
&\qquad \cdot \Pr_{\boldsymbol{\theta}}\big(X_{t+1} = x_{t+1} \,\big|\, Y_{t+1} = y_{t+1}\big) \cdot \Pr_{\boldsymbol{\theta}}\big(X_{t+2:\ell} = x_{t+2:\ell} \,\big|\, Y_{t+1} = y_{t+1}\big) \\
&= \sum_{y_{t+1} \in [K]} A_{y_t, y_{t+1}} \cdot B_{y_{t+1}, x_{t+1}} \cdot \beta_{t+1}(y_{t+1}).
\end{aligned}
$$

Keep in mind, xt+1 is known! it could be emissioned by any yt+1

**Backward inductive computation**: Compute $\beta_t$ using $\beta_{t+1}$ in $O(K^2)$ time.

Given parameters $\boldsymbol{\theta} = (\boldsymbol{\pi}, \boldsymbol{A}, \boldsymbol{B})$ and sequence $x_{1:\ell} \in [D]^{\ell}$:

Given parameters $\theta = (\pi, A, B)$ and sequence $x_{1:\ell} \in [D]^\ell$:

▶ (**Forward pass**)

Starting with $\alpha_1(y_1) = \pi y_1 \cdot B y_1, x_1$ for each $y_1 \in [K]$,

$$\alpha_{t+1}(y_{t+1}) = \left( \sum_{y_t \in [K]} \alpha_t(y_t) \cdot A_{y_t, y_{t+1}} \right) \cdot B y_{t+1}, x_{t+1} \quad \text{for each } y_{t+1} \in [K].$$

Given parameters $\boldsymbol{\theta} = (\boldsymbol{\pi}, \boldsymbol{A}, \boldsymbol{B})$ and sequence $x_{1:\ell} \in [D]^{\ell}$:

▶ (**Forward pass**)
Starting with $\alpha_1(y_1) = \pi_{y_1} \cdot B_{y_1, x_1}$ for each $y_1 \in [K]$,

$$\alpha_{t+1}(y_{t+1}) = \left( \sum_{y_t \in [K]} \alpha_t(y_t) \cdot A_{y_t, y_{t+1}} \right) \cdot B_{y_{t+1}, x_{t+1}} \quad \text{for each } y_{t+1} \in [K].$$

▶ (**Backward pass**)
Starting with $\beta_\ell(y_\ell) = 1$ for each $y_\ell \in [K]$,

$$\beta_{t-1}(y_{t-1}) = \sum_{y_t \in [K]} A_{y_{t-1}, y_t} \cdot B_{y_t, x_t} \cdot \beta_t(y_t) \quad \text{for each } y_{t-1} \in [K].$$

# FORWARD-BACKWARD ALGORITHM

Given parameters $\boldsymbol{\theta} = (\boldsymbol{\pi}, \boldsymbol{A}, \boldsymbol{B})$ and sequence $x_{1:\ell} \in [D]^\ell$:

▶ (**Forward pass**)
  Starting with $\alpha_1(y_1) = \pi_{y_1} \cdot B_{y_1, x_1}$ for each $y_1 \in [K]$,

$$\alpha_{t+1}(y_{t+1}) = \left( \sum_{y_t \in [K]} \alpha_t(y_t) \cdot A_{y_t, y_{t+1}} \right) \cdot B_{y_{t+1}, x_{t+1}} \quad \text{for each } y_{t+1} \in [K].$$

▶ (**Backward pass**)
  Starting with $\boxed{\beta_\ell(y_\ell) = 1 \text{ for each } y_\ell \in [K]},$

$$\beta_{t-1}(y_{t-1}) = \sum_{y_t \in [K]} A_{y_{t-1}, y_t} \cdot B_{y_t, x_t} \cdot \beta_t(y_t) \quad \text{for each } y_{t-1} \in [K].$$

▶ (**Also in backward pass**)
  Compute conditional probabilities:

$$\Pr_{\boldsymbol{\theta}}\big(Y_t = y_t \mid X_{1:\ell} = x_{1:\ell}\big) = \frac{\alpha_t(y_t) \cdot \beta_t(y_t)}{\text{normalization term}} \quad \text{for each } y_t \in [K].$$

Can also compute

$$\Pr_{\boldsymbol{\theta}}\big(Y_{t:t+1} = y_{t:t+1} \mid X_{1:\ell} = x_{1:\ell}\big) \quad \text{for each } y_{t:t+1} \in [K]^2$$

using forward-backward.

Can also compute

$$\Pr_{\boldsymbol{\theta}}\big(Y_{t:t+1} = y_{t:t+1} \mid X_{1:\ell} = x_{1:\ell}\big) \quad \text{for each } y_{t:t+1} \in [K]^2$$

using forward-backward.

Using Markov property, can string together these probabilities to get

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

Can also compute

$$\Pr_{\boldsymbol{\theta}}\big(Y_{t:t+1} = y_{t:t+1} \,|\, X_{1:\ell} = x_{1:\ell}\big) \quad \text{for each } y_{t:t+1} \in [K]^2$$

using forward-backward.

Using Markov property, can string together these probabilities to get

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

Similar procedure for computing most likely state sequence:

$$\arg\max_{y_{1:\ell} \in [K]^\ell} \Pr_{\boldsymbol{\theta}}\big(Y_{1:\ell} = y_{1:\ell} \,|\, X_{1:\ell} = x_{1:\ell}\big)$$

(Viterbi algorithm).

Can also compute

$$\Pr_{\boldsymbol{\theta}}\big(Y_{t:t+1} = y_{t:t+1} \mid X_{1:\ell} = x_{1:\ell}\big) \quad \text{for each } y_{t:t+1} \in [K]^2$$

using forward-backward.

Using Markov property, can string together these probabilities to get

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

Similar procedure for computing most likely state sequence:

$$\underset{y_{1:\ell} \in [K]^{\ell}}{\arg\max} \Pr_{\boldsymbol{\theta}}\big(Y_{1:\ell} = y_{1:\ell} \mid X_{1:\ell} = x_{1:\ell}\big)$$

(Viterbi algorithm).

**See Rabiner's tutorial for details.**

# HMM PARAMETER ESTIMATION

# PARAMETER ESTIMATION

**Parameter estimation problem**:

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

**Parameter estimation problem**:

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

As is the case for mixture models, MLE for HMMs is generally intractable.

# PARAMETER ESTIMATION

**Parameter estimation problem**:

▶ **Given**: $n$ observation sequences $x_{1:\ell}^{(s)}$ for $s \in [n]$.

▶ **Goal**: parameter estimates $\hat{\boldsymbol{\theta}} = (\hat{\boldsymbol{\pi}}, \widehat{\boldsymbol{A}}, \widehat{\boldsymbol{B}})$.

As is the case for mixture models, MLE for HMMs is generally intractable.

Nevertheless, we can use **Expectation-Maximization** to find a local maximizer of the likelihood function. (Called the **Baum-Welch** algorithm in this context.)

Suppose we have current guess for parameters $\hat{\boldsymbol{\theta}} = (\hat{\boldsymbol{\pi}}, \widehat{\boldsymbol{A}}, \widehat{\boldsymbol{B}})$.

# EXPECTATION-MAXIMIZATION FOR HMMS

Suppose we have current guess for parameters $\hat{\boldsymbol{\theta}} = (\hat{\boldsymbol{\pi}}, \widehat{\boldsymbol{A}}, \widehat{\boldsymbol{B}})$.

Compute, for each training sequence $x_{1:\ell}^{(s)}$,

$$
\gamma_t^{(s)}(y_t) := \mathrm{Pr}_{\hat{\boldsymbol{\theta}}}\Big(Y_t = y_t \,|\, X_{1:\ell} = x_{1:\ell}^{(s)}\Big) \quad \text{for all } y_t \in [K]
$$

$$
\xi_t^{(s)}(y_{t-1}, y_t) := \mathrm{Pr}_{\hat{\boldsymbol{\theta}}}\Big(Y_{t-1:t} = y_{t-1:t} \,|\, X_{1:\ell} = x_{1:\ell}^{(s)}\Big) \quad \text{for all } y_{t-1:t} \in [K]^2
$$

using Forward-Backward (see Rabiner tutorial for full details).

# EXPECTATION-MAXIMIZATION FOR HMMs

Suppose we have current guess for parameters $\hat{\boldsymbol{\theta}} = (\hat{\boldsymbol{\pi}}, \widehat{\boldsymbol{A}}, \widehat{\boldsymbol{B}})$.

Compute, for each training sequence $x_{1:\ell}^{(s)}$,

$$
\begin{aligned}
\gamma_t^{(s)}(y_t) &:= \operatorname{Pr}_{\hat{\boldsymbol{\theta}}}\Big(Y_t = y_t \,|\, X_{1:\ell} = x_{1:\ell}^{(s)}\Big) \quad \text{for all } y_t \in [K] \\
\xi_t^{(s)}(y_{t-1}, y_t) &:= \operatorname{Pr}_{\hat{\boldsymbol{\theta}}}\Big(Y_{t-1:t} = y_{t-1:t} \,|\, X_{1:\ell} = x_{1:\ell}^{(s)}\Big) \quad \text{for all } y_{t-1:t} \in [K]^2
\end{aligned}
$$

using Forward-Backward (see Rabiner tutorial for full details).

---

Expected complete log likelihood of $\boldsymbol{\theta} = (\boldsymbol{\pi}, \boldsymbol{A}, \boldsymbol{B})$:

$$
\sum_{s=1}^{n} \Bigg\{ \sum_{y_1 \in [K]} \gamma_1^{(s)}(y_1) \ln \pi_{y_1} + \sum_{t=2}^{\ell} \sum_{y_{t-1} \in [K]} \sum_{y_t \in [K]} \xi_t^{(s)}(y_{t-1}, y_t) \ln A_{y_{t-1}, y_t} \\
+ \sum_{t=1}^{\ell} \sum_{y_t \in [K]} \gamma_t^{(s)}(y_t) \sum_{j=1}^{D} \mathbb{1}\{x_t^{(s)} = j\} \ln B_{y_t, j} \Bigg\}.
$$

# EXPECTATION-MAXIMIZATION FOR HMMS

Suppose we have current guess for parameters $\hat{\boldsymbol{\theta}} = (\hat{\boldsymbol{\pi}}, \widehat{\boldsymbol{A}}, \widehat{\boldsymbol{B}})$.

Compute, for each training sequence $x_{1:\ell}^{(s)}$,

$$\gamma_t^{(s)}(y_t) \quad := \quad \mathrm{Pr}_{\hat{\boldsymbol{\theta}}}\Big(Y_t = y_t \,|\, X_{1:\ell} = x_{1:\ell}^{(s)}\Big) \quad \text{for all } y_t \in [K]$$

$$\xi_t^{(s)}(y_{t-1}, y_t) \quad := \quad \mathrm{Pr}_{\hat{\boldsymbol{\theta}}}\Big(Y_{t-1:t} = y_{t-1:t} \,|\, X_{1:\ell} = x_{1:\ell}^{(s)}\Big) \quad \text{for all } y_{t-1:t} \in [K]^2$$

using Forward-Backward (see Rabiner tutorial for full details).

---

Expected complete log likelihood of $\boldsymbol{\theta} = (\boldsymbol{\pi}, \boldsymbol{A}, \boldsymbol{B})$:

$$\sum_{s=1}^{n} \Bigg\{ \sum_{y_1 \in [K]} \gamma_1^{(s)}(y_1) \ln \pi_{y_1} + \sum_{t=2}^{\ell} \sum_{y_{t-1} \in [K]} \sum_{y_t \in [K]} \xi_t^{(s)}(y_{t-1}, y_t) \ln A_{y_{t-1}, y_t}$$
$$+ \sum_{t=1}^{\ell} \sum_{y_t \in [K]} \gamma_t^{(s)}(y_t) \sum_{j=1}^{D} \mathbb{1}\{x_t^{(s)} = j\} \ln B_{y_t, j} \Bigg\}.$$

Can easily find maximizing parameters $\boldsymbol{\theta}$ (subject to constraints that $\boldsymbol{\pi}$ and rows of $\boldsymbol{A}$ and $\boldsymbol{B}$ are probability distributions).

**Input**: $n$ observation sequences $x_{1:\ell}^{(s)}$ for $s \in [n]$.

Initialize $\boldsymbol{\theta} = (\boldsymbol{\pi}, \boldsymbol{A}, \boldsymbol{B})$ somehow.

# EXPECTATION-MAXIMIZATION FOR HMMS

**Input**: $n$ observation sequences $x_{1:\ell}^{(s)}$ for $s \in [n]$.

Initialize $\boldsymbol{\theta} = (\boldsymbol{\pi}, \boldsymbol{A}, \boldsymbol{B})$ somehow. Then repeat:

- ▶ **E step**: For each $s \in [n]$, compute

$$
\begin{aligned}
\gamma_t^{(s)}(y_t) &:= \Pr_{\boldsymbol{\theta}}\Big(Y_t = y_t \mid X_{1:\ell} = x_{1:\ell}^{(s)}\Big) \\
\xi_t^{(s)}(y_{t-1}, y_t) &:= \Pr_{\boldsymbol{\theta}}\Big(Y_{t-1:t} = y_{t-1:t} \mid X_{1:\ell} = x_{1:\ell}^{(s)}\Big).
\end{aligned}
$$

using Forward-Backward.

# EXPECTATION-MAXIMIZATION FOR HMMS

**Input**: $n$ observation sequences $x_{1:\ell}^{(s)}$ for $s \in [n]$.

Initialize $\boldsymbol{\theta} = (\boldsymbol{\pi}, \boldsymbol{A}, \boldsymbol{B})$ somehow. Then repeat:

▶ **E step**: For each $s \in [n]$, compute

$$
\begin{aligned}
\gamma_t^{(s)}(y_t) &:= \Pr_{\boldsymbol{\theta}}\Big(Y_t = y_t \mid X_{1:\ell} = x_{1:\ell}^{(s)}\Big) \\
\xi_t^{(s)}(y_{t-1}, y_t) &:= \Pr_{\boldsymbol{\theta}}\Big(Y_{t-1:t} = y_{t-1:t} \mid X_{1:\ell} = x_{1:\ell}^{(s)}\Big).
\end{aligned}
$$

using Forward-Backward.

▶ **M step**: Update parameters

$$
\begin{aligned}
\pi_i &:= \frac{\sum_{s=1}^n \gamma_1^{(s)}(i)}{\sum_{s=1}^n \sum_{j \in [K]} \gamma_1^{(s)}(j)} \\
A_{i,j} &:= \frac{\sum_{s=1}^n \sum_{t=2}^\ell \xi_t^{(s)}(i,j)}{\sum_{s=1}^n \sum_{t=2}^\ell \sum_{k \in [K]} \xi_t^{(s)}(i,k)} \\
B_{i,j} &:= \frac{\sum_{s=1}^n \sum_{t=1}^\ell \gamma_t^{(s)}(i) \cdot \mathbb{1}\{x_t^{(s)} = j\}}{\sum_{s=1}^n \sum_{t=1}^\ell \gamma_t^{(s)}(i)}.
\end{aligned}
$$

Can have HMMs with continuous observations $X_t$, say, taking values in $\mathbb{R}^d$.

# Generalization to continuous observations

Can have HMMs with continuous observations $\boldsymbol{X}_t$, say, taking values in $\mathbb{R}^d$.

Specify conditional densities $p_i$ of $\boldsymbol{X}_t$ given $Y_t = i$ for each $i \in [K]$
(e.g., Gaussians $\mathrm{N}(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$).

Can have HMMs with continuous observations $\boldsymbol{X}_t$, say, taking values in $\mathbb{R}^d$.

Specify conditional densities $p_i$ of $\boldsymbol{X}_t$ given $Y_t = i$ for each $i \in [K]$
(e.g., Gaussians $\mathrm{N}(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$).

- Forward-Backward remains the same, except with $B_{y_t, x_t}$ replaced by density value $p_{y_t}(\boldsymbol{x}_t)$.

Can have HMMs with continuous observations $\boldsymbol{X}_t$, say, taking values in $\mathbb{R}^d$.

Specify conditional densities $p_i$ of $\boldsymbol{X}_t$ given $Y_t = i$ for each $i \in [K]$
(e.g., Gaussians $\mathrm{N}(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$).

- Forward-Backward remains the same, except with $B_{y_t, x_t}$ replaced by density value $p_{y_t}(\boldsymbol{x}_t)$.
- "M step" in E-M maximizes expected complete log likelihood of conditional density parameters (e.g., $\boldsymbol{\mu}_i$ and $\boldsymbol{\Sigma}_i$ for Gaussian densities).

$$\sum_{s=1}^{n} \left\{ \sum_{y_1 \in [K]} \gamma_1^{(s)}(y_1) \ln \pi_{y_1} + \sum_{t=2}^{\ell} \sum_{y_{t-1} \in [K]} \sum_{y_t \in [K]} \xi_t^{(s)}(y_{t-1}, y_t) \ln A_{y_{t-1}, y_t} \right.$$
$$\left. + \sum_{t=1}^{\ell} \sum_{y_t \in [K]} \gamma_t^{(s)}(y_t) \ln p_{y_t}(\boldsymbol{x}_t^{(s)}) \right\}.$$

# RECAP

- HMM = Markov chain $\{(X_t, Y_t)\}_{t \in \mathbb{N}}$ where hidden state sequence $\{Y_t\}_{t \in \mathbb{N}}$ is a discrete Markov chain; and conditioned on $Y_t$, observation $X_t$ is independent of everything else.

- Computing sequence probabilities and hidden state conditional probabilities avoids exponential computation due to Markov chain structure.

- Key algorithms: Forward-Backward algorithm (computing conditional probabilities), Viterbi (for most probably hidden state sequence), Baum-Welch (same as E-M for HMMs).

- Many applications: heavily used in speech recognition, bioinformatics, natural language processing, etc.