

# Hidden Markov Models

11-711: Algorithms for NLP

Fall 2017

- 1 Notations
- 2 Hidden Markov Model
- 3 Computing the Likelihood: Forward-Pass Algorithm
- 4 Finding the Hidden Sequence: Viterbi Algorithm
- 5 Estimating Parameters: Baum-Welch Algorithm

- 1 Notations
- 2 Hidden Markov Model
- 3 Computing the Likelihood: Forward-Pass Algorithm
- 4 Finding the Hidden Sequence: Viterbi Algorithm
- 5 Estimating Parameters: Baum-Welch Algorithm

- $\mathbb{R}$ : set of real numbers
- Cartesian products:
  - $\mathbb{R}^{D_1 \times D_2}$ : set of matrices of size  $D_1 \times D_2$ , with real entries

- $\mathbb{R}$ : set of real numbers
- Cartesian products:
  - $\mathbb{R}^{D_1 \times D_2}$ : set of matrices of size  $D_1 \times D_2$ , with real entries
- Vectors:
  - Lower case:  $a, b, c, \dots$
  - Row major:  $a \in \mathbb{R}^D$  means  $a \in \mathbb{R}^{1 \times D}$

- $\mathbb{R}$ : set of real numbers
- Cartesian products:
  - $\mathbb{R}^{D_1 \times D_2}$ : set of matrices of size  $D_1 \times D_2$ , with real entries
- Vectors:
  - Lower case:  $a, b, c, \dots$
  - Row major:  $a \in \mathbb{R}^D$  means  $a \in \mathbb{R}^{1 \times D}$
- Matrices:
  - Upper case:  $A, B, C, \dots$
  - Row major

- 1 Notations
- 2 Hidden Markov Model
- 3 Computing the Likelihood: Forward-Pass Algorithm
- 4 Finding the Hidden Sequence: Viterbi Algorithm
- 5 Estimating Parameters: Baum-Welch Algorithm





---

Name	Notation	Meaning/Property

---

---

Name	Notation	Meaning/Property
State space	$Q = \{q_1, q_2, \dots, q_N\}$	Set of $N$ states
Observation space	$V = \{w_1, w_2, \dots, w_V\}$	Set of $V$ states

---

---

---

Name	Notation	Meaning/Property
State space	$Q = \{q_1, q_2, \dots, q_N\}$	Set of $N$ states
Observation space	$V = \{w_1, w_2, \dots, w_V\}$	Set of $V$ states
State sequence	$S = \{s_1, s_2, \dots, s_T\}$	Sequence of $T$ steps. $s_i \in Q$
Observation sequence	$O = \{o_1, o_2, \dots, o_T\}$	Sequence of $T$ steps. $o_i \in V$

---

---

Name	Notation	Meaning/Property
State space	$Q = \{q_1, q_2, \dots, q_N\}$	Set of $N$ states
Observation space	$V = \{w_1, w_2, \dots, w_V\}$	Set of $V$ states
State sequence	$S = \{s_1, s_2, \dots, s_T\}$	Sequence of $T$ steps. $s_i \in Q$
Observation sequence	$O = \{o_1, o_2, \dots, o_T\}$	Sequence of $T$ steps. $o_i \in V$
Transition probs	$A \in \mathbb{R}^{N \times N}$	Each row is a valid distribution

---

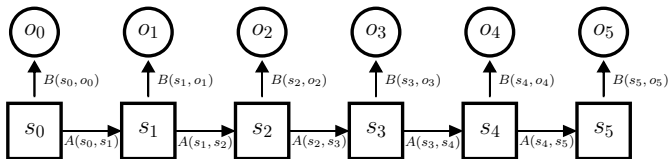
Valid (probability) distribution:

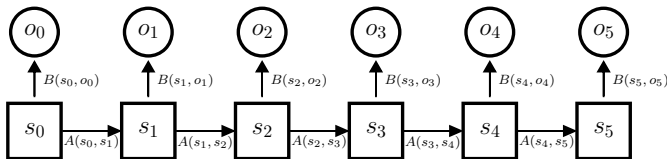
- $A(i, j) \geq 0$
- $\sum_{j=1}^N A(i, j) = 1$

Name	Notation	Meaning/Property
State space	$Q = \{q_1, q_2, \dots, q_N\}$	Set of $N$ states
Observation space	$V = \{w_1, w_2, \dots, w_V\}$	Set of $V$ states
State sequence	$S = \{s_1, s_2, \dots, s_T\}$	Sequence of $T$ steps. $s_i \in Q$
Observation sequence	$O = \{o_1, o_2, \dots, o_T\}$	Sequence of $T$ steps. $o_i \in V$
Transition probs	$A \in \mathbb{R}^{N \times N}$	Each row is a valid distribution
Emission probs	$B \in \mathbb{R}^{N \times V}$	Each row is a valid distribution

Valid (probability) distribution:

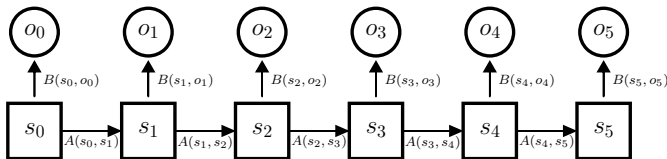
- $A(i, j) \geq 0$
- $\sum_{j=1}^N A(i, j) = 1$





- Markov assumption

$$P(s_0, s_1, s_2, \dots, s_T) = \prod_{t=1}^T P(s_t | s_{<t}) = \prod_{t=1}^T P(s_t | s_{t-1}) \stackrel{\text{def}}{=} \prod_{t=1}^T A(s_{t-1}, s_t)$$



- Markov assumption

$$P(s_0, s_1, s_2, \dots, s_T) = \prod_{t=1}^T P(s_t | s_{<t}) = \prod_{t=1}^T P(s_t | s_{t-1}) \stackrel{\text{def}}{=} \prod_{t=1}^T A(s_{t-1}, s_t)$$

- Independent assumption

$$P(o_t | o_{<t}, s_{\leq t}) = P(o_t | s_t) \stackrel{\text{def}}{=} B(s_t, o_t)$$



- Prior:  $P(A, B)$ 
  - Without seeing anything, how do you believe  $A, B$  should look like?

- Prior:  $P(A, B)$ 
  - Without seeing anything, how do you believe  $A, B$  should look like?
- Likelihood:  $P(O|A, B)$ 
  - Suppose you know  $A, B$ , how likely do you see  $O$ ?

- Prior:  $P(A, B)$ 
  - Without seeing anything, how do you believe  $A, B$  should look like?
- Likelihood:  $P(O|A, B)$ 
  - Suppose you know  $A, B$ , how likely do you see  $O$ ?
- Posterior:  $P(A, B|O)$ 
  - Suppose you see  $O$ , how do you believe  $A, B$  should look like?

- Prior:  $P(A, B)$ 
  - Without seeing anything, how do you believe  $A, B$  should look like?
- Likelihood:  $P(O|A, B)$ 
  - Suppose you know  $A, B$ , how likely do you see  $O$ ?
- Posterior:  $P(A, B|O)$ 
  - Suppose you see  $O$ , how do you believe  $A, B$  should look like?
- Bayes rule

$$P(A, B|O) = \frac{P(O|A, B) \cdot P(A, B)}{P(O)}$$

$$\text{posterior} = \frac{\text{likelihood} \cdot \text{prior}}{\text{data\_distribution}}$$

- Settings:
  - Prior  $P(A, B)$  are uniform.
  - $P(O)$  is unknown, but  $O$  will be observed.

- Settings:
  - Prior  $P(A, B)$  are uniform.
  - $P(O)$  is unknown, but  $O$  will be observed.
- **Question 1:** Compute the **likelihood**:  $P(O|A, B)$

- Settings:
  - Prior  $P(A, B)$  are uniform.
  - $P(O)$  is unknown, but  $O$  will be observed.
- **Question 1:** Compute the **likelihood**:  $P(O|A, B)$
- **Question 2:** (only for HMMs) Find the **hidden sequence**  $S$

$$S^* = \operatorname{argmax}_S P(O|S, A, B)$$

- Settings:
  - Prior  $P(A, B)$  are uniform.
  - $P(O)$  is unknown, but  $O$  will be observed.
- **Question 1:** Compute the **likelihood**:  $P(O|A, B)$
- **Question 2:** (only for HMMs) Find the **hidden sequence**  $S$

$$S^* = \operatorname{argmax}_S P(O|S, A, B)$$

- **Question 3:** Find the **posterior**

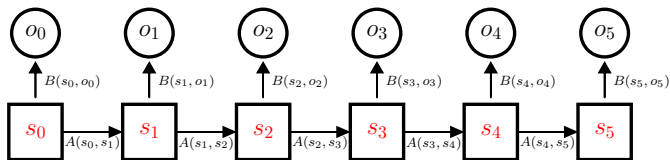
$$A^*, B^* = \operatorname{argmax}_{A, B} P(A, B|O)$$



- 1 Notations
- 2 Hidden Markov Model
- 3 Computing the Likelihood: Forward-Pass Algorithm**
- 4 Finding the Hidden Sequence: Viterbi Algorithm
- 5 Estimating Parameters: Baum-Welch Algorithm

Compute the **likelihood**:  $P(o_1, o_2, \dots, o_T | A, B)$

$O$  depends on the hidden sequence  $S$



# Compute $P(O|A, B)$ for a particular $S$

---

$$S = (s_1, s_2, \dots, s_T)$$

$$\begin{aligned} P(O|S, A, B) &= P(o_1|s_1)P(o_2|s_2) \cdots P(o_T|s_T) \\ &= \prod_{i=1}^T P(o_i|s_i) \end{aligned}$$

$$P(O|A, B) = \sum_S P(O|S, A, B) \cdot P(S)$$

$$\begin{aligned} P(O|A, B) &= \sum_S P(O|S, A, B) \cdot P(S) \\ &= \sum_{s_1, \dots, s_T} \left( \prod_{i=1}^T P(o_i|s_i) \cdot \prod_{i=1}^T P(s_i|s_{i-1}) \right) \end{aligned}$$

$$\begin{aligned} P(O|A, B) &= \sum_S P(O|S, A, B) \cdot P(S) \\ &= \sum_{s_1, \dots, s_T} \left( \prod_{i=1}^T P(o_i | s_i) \cdot \prod_{i=1}^T P(s_i | s_{i-1}) \right) \\ &= \sum_{s_1, \dots, s_T} \left( \prod_{i=1}^T B(s_i, o_i) \cdot \prod_{i=1}^T A(s_{i-1}, s_i) \right) \end{aligned}$$

$$\begin{aligned} P(O|A, B) &= \sum_S P(O|S, A, B) \cdot P(S) \\ &= \sum_{s_1, \dots, s_T} \left( \prod_{i=1}^T P(o_i | s_i) \cdot \prod_{i=1}^T P(s_i | s_{i-1}) \right) \\ &= \sum_{s_1, \dots, s_T} \left( \prod_{i=1}^T B(s_i, o_i) \cdot \prod_{i=1}^T A(s_{i-1}, s_i) \right) \end{aligned}$$

But there are  $\Theta(N^T)$  possible sequences  $S$



- Compute harder instances based on easier instances
- Cache easier instances

- $T = 1$

$$P(o_1|A, B) = \sum_{s_1} B(s_1, o_1) \cdot A(s_0, s_1)$$

- $T = 1$

$$P(o_1|A, B) = \sum_{s_1} B(s_1, o_1) \cdot A(s_0, s_1)$$

- $T = 2$

$$P(o_1, o_2|A, B) = \sum_{s_2} \sum_{s_1} B(s_1, o_1)B(s_2, o_2) \cdot A(s_0, s_1)A(s_1, s_2)$$

- $T = 1$

$$P(o_1|A, B) = \sum_{s_1} B(s_1, o_1) \cdot A(s_0, s_1)$$

- $T = 2$

$$\begin{aligned} P(o_1, o_2|A, B) &= \sum_{s_2} \sum_{s_1} B(s_1, o_1) B(s_2, o_2) \cdot A(s_0, s_1) A(s_1, s_2) \\ &= \sum_{s_2} B(s_2, o_2) \left( \sum_{s_1} B(s_1, o_1) \cdot A(s_0, s_1) A(s_1, s_2) \right) \end{aligned}$$

- $T = 1$

$$P(o_1|A, B) = \sum_{s_1} B(s_1, o_1) \cdot A(s_0, s_1)$$

- $T = 2$

$$\begin{aligned} P(o_1, o_2|A, B) &= \sum_{s_2} \sum_{s_1} B(s_1, o_1) B(s_2, o_2) \cdot A(s_0, s_1) A(s_1, s_2) \\ &= \sum_{s_2} B(s_2, o_2) \left( \sum_{s_1} B(s_1, o_1) \cdot A(s_0, s_1) A(s_1, s_2) \right) \end{aligned}$$

- $T = 2$

$$P(o_1, o_2 | A, B) = \sum_{s_2} \sum_{s_1} B(s_1, o_1) B(s_2, o_2) \cdot A(s_0, s_1) A(s_1, s_2)$$

- $T = 2$

$$P(o_1, o_2 | A, B) = \sum_{s_2} \sum_{s_1} B(s_1, o_1) B(s_2, o_2) \cdot A(s_0, s_1) A(s_1, s_2)$$

- $T = 3$

$$P(o_1, o_2, o_3 | A, B) = \sum_{s_3} \sum_{s_2} \sum_{s_1} B(s_1, o_1) B(s_2, o_2) B(s_3, o_3) \cdot A(s_0, s_1) A(s_1, s_2) A(s_2, s_3)$$

- $T = 2$

$$P(o_1, o_2 | A, B) = \sum_{s_2} \sum_{s_1} B(s_1, o_1) B(s_2, o_2) \cdot A(s_0, s_1) A(s_1, s_2)$$

- $T = 3$

$$\begin{aligned} P(o_1, o_2, o_3 | A, B) &= \sum_{s_3} \sum_{s_2} \sum_{s_1} B(s_1, o_1) B(s_2, o_2) B(s_3, o_3) \cdot A(s_0, s_1) A(s_1, s_2) A(s_2, s_3) \\ &= \sum_{s_3} B(s_3, o_3) \sum_{s_2} \sum_{s_1} B(s_1, o_1) B(s_2, o_2) \cdot A(s_0, s_1) A(s_1, s_2) A(s_2, s_3) \end{aligned}$$



- $T = 2$

$$P(o_1, o_2 | A, B) = \sum_{s_2} \sum_{s_1} B(s_1, o_1) B(s_2, o_2) \cdot A(s_0, s_1) A(s_1, s_2)$$

- $T = 3$

$$\begin{aligned} P(o_1, o_2, o_3 | A, B) &= \sum_{s_3} \sum_{s_2} \sum_{s_1} B(s_1, o_1) B(s_2, o_2) B(s_3, o_3) \cdot A(s_0, s_1) A(s_1, s_2) A(s_2, s_3) \\ &= \sum_{s_3} B(s_3, o_3) \sum_{s_2} \sum_{s_1} B(s_1, o_1) B(s_2, o_2) \cdot A(s_0, s_1) A(s_1, s_2) A(s_2, s_3) \end{aligned}$$

- $T = 2$

$$P(o_1, o_2 | A, B) = \sum_{s_2} \sum_{s_1} B(s_1, o_1) B(s_2, o_2) \cdot A(s_0, s_1) A(s_1, s_2)$$

- $T = 3$

$$\begin{aligned} P(o_1, o_2, o_3 | A, B) &= \sum_{s_3} \sum_{s_2} \sum_{s_1} B(s_1, o_1) B(s_2, o_2) B(s_3, o_3) \cdot A(s_0, s_1) A(s_1, s_2) A(s_2, s_3) \\ &= \sum_{s_3} B(s_3, o_3) \sum_{s_2} \sum_{s_1} B(s_1, o_1) B(s_2, o_2) \cdot A(s_0, s_1) A(s_1, s_2) A(s_2, s_3) \end{aligned}$$

- We can catch the blue quantities

$$f[t, s] \stackrel{\text{def}}{=} P(o_1, o_2, \dots, o_t, s_t = s | A, B)$$

- Caching values:

$$f[t, s] \stackrel{\text{def}}{=} P(o_1, o_2, \dots, o_t, s_t = s | A, B)$$

- Caching values:

$$f[t, s] \stackrel{\text{def}}{=} P(o_1, o_2, \dots, o_t, s_t = s | A, B)$$

- How to compute  $f[t, s]$ ?

$$f[t + 1, s] = P(o_1, o_2, \dots, o_t, o_{t+1}, s_{t+1} = s | A, B)$$

- Caching values:

$$f[t, s] \stackrel{\text{def}}{=} P(o_1, o_2, \dots, o_t, s_t = s | A, B)$$

- How to compute  $f[t, s]$ ?

$$\begin{aligned} f[t + 1, s] &= P(o_1, o_2, \dots, o_t, o_{t+1}, s_{t+1} = s | A, B) \\ &= \sum_{s'} P(o_1, o_2, \dots, o_t, s_t = s' | A, B) \cdot P(s_{t+1} = s | s_t = s') P(o_{t+1} | s) \end{aligned}$$

- Caching values:

$$f[t, s] \stackrel{\text{def}}{=} P(o_1, o_2, \dots, o_t, s_t = s | A, B)$$

- How to compute  $f[t, s]$ ?

$$\begin{aligned} f[t + 1, s] &= P(o_1, o_2, \dots, o_t, o_{t+1}, s_{t+1} = s | A, B) \\ &= \sum_{s'} P(o_1, o_2, \dots, o_t, s_t = s' | A, B) \cdot P(s_{t+1} = s | s_t = s') P(o_{t+1} | s) \\ &= \sum_{s'} f[t, s'] \cdot A(s', s) B(s, o_{t+1}) \end{aligned}$$

- How to compute  $f[t, s]$ ?

$$f[t + 1, s] = \sum_{s'} f[t, s'] \cdot A(s', s)B(s, o_{t+1})$$

- How to compute  $f[t, s]$ ?

$$f[t + 1, s] = \sum_{s'} f[t, s'] \cdot A(s', s)B(s, o_{t+1})$$

- What can we do with  $f[t, s]$ ?

$$P(o_1, o_2, \dots, o_t | A, B) = \sum_s P(o_1, o_2, \dots, o_t, s_t = s | A, B) \cdot P(o_t | s_t = s)$$



- How to compute  $f[t, s]$ ?

$$f[t + 1, s] = \sum_{s'} f[t, s'] \cdot A(s', s) B(s, o_{t+1})$$

- What can we do with  $f[t, s]$ ?

$$\begin{aligned} P(o_1, o_2, \dots, o_t | A, B) &= \sum_s P(o_1, o_2, \dots, o_t, s_t = s | A, B) \cdot P(o_t | s_t = s) \\ &= \sum_s f[t, s] \cdot B(s, o_t) \end{aligned}$$

(1) Initialize:

- For each hidden state  $s$ :

$$f[1, s] = P(o_1, s_1 = s | A, B) \leftarrow B(s, o_1) \cdot A(s_0, s)$$

(2) For  $t = 2$  to  $T$ :

- For each hidden state  $s$ :

$$f[t, s] \leftarrow \sum_{s'} f[t-1, s'] \cdot A(s', s) B(s, o_t)$$

(3) Finally:

$$P(o_1, o_2, \dots, o_T | A, B) \leftarrow \sum_s f[T, s] B(s, o_T)$$

(1) Initialize:

- For each hidden state  $s$ :

$$f[1, s] = P(o_1, s_1 = s | A, B) \leftarrow B(s, o_1) \cdot A(s_0, s)$$

(2) For  $t = 2$  to  $T$ :

- For each hidden state  $s$ :

$$f[t, s] \leftarrow \sum_{s'} f[t-1, s'] \cdot A(s', s) B(s, o_t)$$

(3) Finally:

$$P(o_1, o_2, \dots, o_T | A, B) \leftarrow \sum_s f[T, s] B(s, o_T)$$

Complexity:  $O(T \cdot N^2)$

- 1 Notations
- 2 Hidden Markov Model
- 3 Computing the Likelihood: Forward-Pass Algorithm
- 4 Finding the Hidden Sequence: Viterbi Algorithm**
- 5 Estimating Parameters: Baum-Welch Algorithm

Given  $O$ . Find the hidden sequence  $S$

$$S^* = \operatorname{argmax}_S P(O|S, A, B)$$

$$S = (s_1, s_2, \dots, s_T)$$

$$\begin{aligned}\Rightarrow P(O|S, A, B) &= P(o_1|s_1)P(o_2|s_2) \cdots P(o_T|s_T) \\ &= \prod_{i=1}^T P(o_i|s_i)\end{aligned}$$

$$S = (s_1, s_2, \dots, s_T)$$

$$\begin{aligned}\Rightarrow P(O|S, A, B) &= P(o_1|s_1)P(o_2|s_2) \cdots P(o_T|s_T) \\ &= \prod_{i=1}^T P(o_i|s_i)\end{aligned}$$

But there are  $\Theta(N^T)$  possible sequences  $S$

Replace argmax with max:

- $T = 1$

$$\begin{aligned}\max P(o_1|A, B) &= \max_{s_1} P(o_1|s_1) \cdot P(s_1|s_0) \\ &= \max_{s_1} B(s_1, o_1) \cdot A(s_0, s_1)\end{aligned}$$



Replace argmax with max:

- $T = 1$

$$\begin{aligned}\max P(o_1|A, B) &= \max_{s_1} P(o_1|s_1) \cdot P(s_1|s_0) \\ &= \max_{s_1} B(s_1, o_1) \cdot A(s_0, s_1)\end{aligned}$$

- $T = 2$

$$\max P(o_1, o_2|A, B) = \max_{s_2} \max_{s_1} B(s_1, o_1)B(s_2, o_2) \cdot A(s_0, s_1)A(s_1, s_2)$$

Replace argmax with max:

- $T = 1$

$$\begin{aligned}\max P(o_1|A, B) &= \max_{s_1} P(o_1|s_1) \cdot P(s_1|s_0) \\ &= \max_{s_1} B(s_1, o_1) \cdot A(s_0, s_1)\end{aligned}$$

- $T = 2$

$$\begin{aligned}\max P(o_1, o_2|A, B) &= \max_{s_2} \max_{s_1} B(s_1, o_1) B(s_2, o_2) \cdot A(s_0, s_1) A(s_1, s_2) \\ &= \max_{s_2} B(s_2, o_2) \left( \max_{s_1} B(s_1, o_1) \cdot A(s_0, s_1) A(s_1, s_2) \right)\end{aligned}$$

Replace argmax with max:

- $T = 1$

$$\begin{aligned}\max P(o_1|A, B) &= \max_{s_1} P(o_1|s_1) \cdot P(s_1|s_0) \\ &= \max_{s_1} B(s_1, o_1) \cdot A(s_0, s_1)\end{aligned}$$

- $T = 2$

$$\begin{aligned}\max P(o_1, o_2|A, B) &= \max_{s_2} \max_{s_1} B(s_1, o_1)B(s_2, o_2) \cdot A(s_0, s_1)A(s_1, s_2) \\ &= \max_{s_2} B(s_2, o_2) \left( \max_{s_1} B(s_1, o_1) \cdot A(s_0, s_1)A(s_1, s_2) \right)\end{aligned}$$

- $T = 2$

$$\max_{s_1, s_2} P(o_1, o_2 | A, B) = \max_{s_2} \max_{s_1} B(s_1, o_1) B(s_2, o_2) \cdot A(s_0, s_1) A(s_1, s_2)$$

- $T = 2$

$$\max_{s_1, s_2} P(o_1, o_2 | A, B) = \max_{s_2} \max_{s_1} B(s_1, o_1) B(s_2, o_2) \cdot A(s_0, s_1) A(s_1, s_2)$$

- $T = 3$

$$\max_{s_1, s_2, s_3} P(o_1, o_2, o_3 | A, B) = \max_{s_3} \max_{s_2} \max_{s_1} B(s_1, o_1) B(s_2, o_2) B(s_3, o_3) \cdot A(s_0, s_1) A(s_1, s_2) A(s_2, s_3)$$

- $T = 2$

$$\max_{s_1, s_2} P(o_1, o_2 | A, B) = \max_{s_2} \max_{s_1} B(s_1, o_1) B(s_2, o_2) \cdot A(s_0, s_1) A(s_1, s_2)$$

- $T = 3$

$$\begin{aligned} \max_{s_1, s_2, s_3} P(o_1, o_2, o_3 | A, B) &= \max_{s_3} \max_{s_2} \max_{s_1} B(s_1, o_1) B(s_2, o_2) B(s_3, o_3) \cdot A(s_0, s_1) A(s_1, s_2) A(s_2, s_3) \\ &= \max_{s_3} B(s_3, o_3) \max_{s_2} \max_{s_1} B(s_1, o_1) B(s_2, o_2) \cdot A(s_0, s_1) A(s_1, s_2) A(s_2, s_3) \end{aligned}$$

- $T = 2$

$$\max_{s_1, s_2} P(o_1, o_2 | A, B) = \max_{s_2} \max_{s_1} B(s_1, o_1) B(s_2, o_2) \cdot A(s_0, s_1) A(s_1, s_2)$$

- $T = 3$

$$\begin{aligned} \max_{s_1, s_2, s_3} P(o_1, o_2, o_3 | A, B) &= \max_{s_3} \max_{s_2} \max_{s_1} B(s_1, o_1) B(s_2, o_2) B(s_3, o_3) \cdot A(s_0, s_1) A(s_1, s_2) A(s_2, s_3) \\ &= \max_{s_3} B(s_3, o_3) \max_{s_2} \max_{s_1} B(s_1, o_1) B(s_2, o_2) \cdot A(s_0, s_1) A(s_1, s_2) A(s_2, s_3) \end{aligned}$$

- $T = 2$

$$\max_{s_1, s_2} P(o_1, o_2 | A, B) = \max_{s_2} \max_{s_1} B(s_1, o_1) B(s_2, o_2) \cdot A(s_0, s_1) A(s_1, s_2)$$

- $T = 3$

$$\begin{aligned} \max_{s_1, s_2, s_3} P(o_1, o_2, o_3 | A, B) &= \max_{s_3} \max_{s_2} \max_{s_1} B(s_1, o_1) B(s_2, o_2) B(s_3, o_3) \cdot A(s_0, s_1) A(s_1, s_2) A(s_2, s_3) \\ &= \max_{s_3} B(s_3, o_3) \max_{s_2} \max_{s_1} B(s_1, o_1) B(s_2, o_2) \cdot A(s_0, s_1) A(s_1, s_2) A(s_2, s_3) \end{aligned}$$

- We can catch the blue quantities

$$g[t, s] \stackrel{\text{def}}{=} \max_{s_1, s_2, \dots, s_{t-1}} P(o_1, o_2, \dots, o_t, s_t = s | A, B)$$



- Caching values:

$$g[t, s] \stackrel{\text{def}}{=} \max_{s_1 \dots s_{t-1}} P(o_1, o_2, \dots, o_t, s_t = s | A, B)$$

- Caching values:

$$g[t, s] \stackrel{\text{def}}{=} \max_{s_1 \dots s_{t-1}} P(o_1, o_2, \dots, o_t, s_t = s | A, B)$$

- How to compute  $g[t, s]$ ? Same!

$$g[t + 1, s] = \max_{s_1, \dots, s_t} P(o_1, o_2, \dots, o_t, o_{t+1}, s_{t+1} = s | A, B)$$

- Caching values:

$$g[t, s] \stackrel{\text{def}}{=} \max_{s_1 \dots s_{t-1}} P(o_1, o_2, \dots, o_t, s_t = s | A, B)$$

- How to compute  $g[t, s]$ ? Same!

$$\begin{aligned} g[t+1, s] &= \max_{s_1, \dots, s_t} P(o_1, o_2, \dots, o_t, o_{t+1}, s_{t+1} = s | A, B) \\ &= \max_{s'} P(o_1, o_2, \dots, o_t, s_t = s' | A, B) \cdot P(s_{t+1} = s | s_t = s') P(o_{t+1} | s) \end{aligned}$$

- Caching values:

$$g[t, s] \stackrel{\text{def}}{=} \max_{s_1 \dots s_{t-1}} P(o_1, o_2, \dots, o_t, s_t = s | A, B)$$

- How to compute  $g[t, s]$ ? Same!

$$\begin{aligned} g[t+1, s] &= \max_{s_1, \dots, s_t} P(o_1, o_2, \dots, o_t, o_{t+1}, s_{t+1} = s | A, B) \\ &= \max_{s'} P(o_1, o_2, \dots, o_t, s_t = s' | A, B) \cdot P(s_{t+1} = s | s_t = s') P(o_{t+1} | s) \\ &= \max_{s'} g[t, s'] \cdot A(s', s) B(s, o_{t+1}) \end{aligned}$$

- How to compute  $g[t, s]$ ?

$$g[t + 1, s] = \max_{s'} g[t, s'] \cdot A(s', s)B(s, o_{t+1})$$

- How to compute  $g[t, s]$ ?

$$g[t + 1, s] = \max_{s'} g[t, s'] \cdot A(s', s)B(s, o_{t+1})$$

- What can we do with  $g[t, s]$ ? Same!

$$\max_{s_1, \dots, s_t} P(o_1, o_2, \dots, o_t | A, B) = \max_s P(o_1, o_2, \dots, o_t, s_t = s | A, B) \cdot P(o_t | s_t = s)$$

- How to compute  $g[t, s]$ ?

$$g[t + 1, s] = \max_{s'} g[t, s'] \cdot A(s', s)B(s, o_{t+1})$$

- What can we do with  $g[t, s]$ ? Same!

$$\begin{aligned} \max_{s_1, \dots, s_t} P(o_1, o_2, \dots, o_t | A, B) &= \max_s P(o_1, o_2, \dots, o_t, s_t = s | A, B) \cdot P(o_t | s_t = s) \\ &= \max_s g[t, s] \cdot B(s, o_t) \end{aligned}$$

(1) Initialize:

- For each hidden state  $s$ :

$$g[1, s] \leftarrow B(s, o_1) \cdot A(s_0, s)$$

(2) For  $t = 2$  to  $T$ :

- For each hidden state  $s$ :

$$g[t, s] \leftarrow \max_{s'} g[t-1, s'] \cdot A(s', s) B(s, o_t)$$

(3) Finally:

$$\max_{s_1 \dots s_T} P(o_1, o_2, \dots, o_T | A, B) \leftarrow \max_s g[T, s] B(s, o_T)$$



(1) Initialize:

- For each hidden state  $s$ :

$$g[1, s] \leftarrow B(s, o_1) \cdot A(s_0, s)$$

(2) For  $t = 2$  to  $T$ :

- For each hidden state  $s$ :

$$g[t, s] \leftarrow \max_{s'} g[t-1, s'] \cdot A(s', s) B(s, o_t)$$

(3) Finally:

$$\max_{s_1 \dots s_T} P(o_1, o_2, \dots, o_T | A, B) \leftarrow \max_s g[T, s] B(s, o_T)$$

We want **argmax**, not **max** !!!

(1) Initialize:

- For each hidden state  $s$ :

$$g[1, s] \leftarrow B(s, o_1) \cdot A(s_0, s)$$

(2) For  $t = 2$  to  $T$ :

- For each hidden state  $s$ :

$$g[t, s] \leftarrow \max_{s'} g[t-1, s'] \cdot A(s', s) B(s, o_t)$$

$$h[t, s] \leftarrow \operatorname{argmax}_{s'} g[t-1, s'] \cdot A(s', s) B(s, o_t)$$

(3) Follow  $h[t, s]$  to find  $s_T^*, s_{T-1}^*, \dots, s_1^*$ .

# + traceback: Viterbi Algorithm

---

(1) Initialize:

- For each hidden state  $s$ :

$$g[1, s] \leftarrow B(s, o_1) \cdot A(s_0, s)$$

(2) For  $t = 2$  to  $T$ :

- For each hidden state  $s$ :

$$g[t, s] \leftarrow \max_{s'} g[t-1, s'] \cdot A(s', s) B(s, o_t)$$

$$h[t, s] \leftarrow \operatorname{argmax}_{s'} g[t-1, s'] \cdot A(s', s) B(s, o_t)$$

(3) Follow  $h[t, s]$  to find  $s_T^*, s_{T-1}^*, \dots, s_1^*$ .

Complexity:  $O(T \cdot N^2)$

- 1 Notations
- 2 Hidden Markov Model
- 3 Computing the Likelihood: Forward-Pass Algorithm
- 4 Finding the Hidden Sequence: Viterbi Algorithm
- 5 Estimating Parameters: Baum-Welch Algorithm

Find the **posterior**

$$A^*, B^* = \operatorname{argmax}_{A, B} P(A, B | O)$$