Hidden Markov Model

Statistical Machine Learning (ENGG*6600*02)

School of Engineering, University of Guelph, ON, Canada

Course Instructor: Benyamin Ghojogh Summer 2023

Probabilistic Graphical Models

Markov and Bayesian Networks

r(a 2)

- A <u>Probabilistic Graphical Model (PGM)</u> is a grpah-based representation of a complex distribution in the possibly high dimensional space [1].
- In other words, PGM is a combination of graph theory and probability theory.
- In a PGM, the <u>random variables</u> are represented by <u>nodes or vertices</u>. There exist <u>edges</u> between two variables which have interaction with one another in terms of <u>probability</u>. Different conditional probabilities can be represented by a PGM.
- There exist two types of PGM which are <u>Markov network</u> (also called <u>Markov random</u> <u>field</u>) and <u>Bayesian network</u> [1]. In the Markov network and Bayesian network, the edges of graph are <u>undirected</u> and <u>directed</u>, respectively.

The Markov Property

Consider a times series of random variables X₁, X₂,..., X_n. In general, the joint probability of these random variables can be written as:

$$\mathbb{P}(X_1, X_2, \dots, X_n) = \mathbb{P}(X_1) \mathbb{P}(X_2 \mid X_1) \mathbb{P}(X_3 \mid X_2, X_1) \dots \mathbb{P}(X_n \mid X_{n-1}, \dots, X_2, X_1), \quad (1)$$

according to chain (or multiplication) rule in probability.

 (The first order) Markov property is an assumption which states that in a time series of random variables X₁, X₂,..., X_n, every random variable is merely dependent on the latest previous random variable and not the others. In other words:

$$\mathbb{P}(\widehat{X_{i}}|\widehat{X_{i-1},X_{i-2},\ldots,X_{2},X_{1}}) = \mathbb{P}(\widehat{X_{i}}|\widehat{X_{i-1}}).$$
(2)

Hence, with Markov property, the chain rule is simplied to:

$$\mathbb{P}(X_1, X_2, \dots, X_n) = \mathbb{P}(X_1) \mathbb{P}(X_2 \mid X_1) \mathbb{P}(X_3 \mid X_2) \dots \mathbb{P}(X_n \mid X_{n-1}).$$
(3)

- The Markov property can be of any order. For example, in a <u>second order Markov</u> property, a random variable is dependent on the latest and one-to-latest variables. Usually, the default Markov property is of order one.
- A stochastic process which has the Markov process is called a Markovian process (or Markov process).

Discrete Time Markov Chain

- A Markov chain is a PGM which has Markov property.
- The Markov chain can be either directed or undirected.
- Usually, Markov chain is a **Bayesian network** where the edges are **directed**.
- It is important not to confuse Markov chain with Markov network.
- There are two types of Markov Chain which are <u>Discrete Time Markov Chain (DTMC)</u>
 [2] and <u>Continuous Time Markov Chain (CTMC)</u>
 [3]. As it is obvious from their names, in DTMC and CTMC, the time of transitions from a random variable to another one is and is not partitioned into discrete slots, respectively.
- If the variables in a <u>DTMC</u> are considered as <u>states</u>, the DTMC can be viewed as a **Finite-State Machine (FSM)** or a **Finite-State Automaton (FSA)** [4].

• HMM is a <u>DTMC</u> which contains a sequence of hidden variables (named states) in addition to a sequence of emitted observation symbols (outputs).



• We have an observation sequence of length τ which is the number of clock times, $t \in \{1, \dots, \tau\}$. Let <u>n</u> and <u>m</u> denote the number of states and observation symbols, respectively.



- We show the sets of states and possible observation symbols by $\mathcal{S} = \{s_1, \dots, s_n\}$ and $\mathcal{O} = \{o_1, \dots, o_m\}$, respectively.
- We show being in state s_i and in observation symbol o_i at time t by $\underline{s_i(t)}$ and $\underline{o_i(t)}$, respectively.
- Let $\mathbb{R}^{n \times n} \ni \mathbf{A} = [a_{i,j}]$ be the state **Transition Probability Matrix (TPM)**, where:

$$\underbrace{a_{i,j}}_{i=1} = \mathbb{P}(\widehat{s_j(t+1)} | \widehat{s_i(t)}).$$
(4)

We have:

$$\sum_{j=1}^{n} \mathbf{a}_{i,j} = 1.$$
(5)

• The Emission Probability Matrix (EPM) is denoted by as $\mathbb{R}^{n \times m} \ni \mathbf{B} = [b_{i,j}]$ where: $b_{i,j} := \mathbb{P}(o_j(t) | s_i(t)),$ (6)

which is the probability of emission of the observation symbols from the states.

We have:

$$\left[\sum_{j=1}^{n} b_{i,j} = 1. \right]$$

$$(7)$$

• Let the initial state distribution be denoted by the vector $\mathbb{R}^n \ni \pi = [\pi_1, \dots, \pi_n]$ where:

$$\pi_i := \mathbb{P}(s_i(1)), \tag{8}$$

and:

$$\sum_{i=1}^{n} \pi_i = 1, \tag{9}$$

to satisfy the probability properties.

- An HMM model is denoted by the tuple $\lambda = (\pi, A, B)$.
- Assume that a <u>sequence of states</u> is generated by the HMM according to the <u>TPM</u>. We denote this generated sequence of states by S^g := s^g(1),..., s^g(τ) where s^g(t) ∈ S, ∀t.
- Likewise, a sequence of outputs (observations) is generated by the HMM according to EPM. We denote this generated sequence of output symbols by $\mathcal{O}^g := o^g(1), \ldots, o^g(\tau)$ where $o^g(t) \in \mathcal{O}, \forall t$.

• We denote the probability of transition from state $s^{g}(t)$ to $s^{g}(t+1)$ by $a_{s^{g}(t),s^{g}(t+1)}$. So:

$$\overbrace{q_{s^{\mathcal{G}}(t)}^{s^{\mathcal{G}}(t+1)} := \mathbb{P}(s^{\mathcal{G}}(t+1) \mid s^{\mathcal{G}}(t)).$$

$$(10)$$

• Note that
$$s^g(t) \in S$$
 and $s^g(t+1) \in S$.

We also denote:

$$\pi_{s^{g}(i)} := \mathbb{P}(s^{g}(i)).$$
(11)

• Likewise, we denote the probability of state $\underline{s^g(t)}$ emitting the observation $\underline{o^g(t)}$ by $b_{\overline{s^g(t)}, \overline{o^g(t)}}$. So:

$$b_{\mathcal{F}}(t) \circ \mathcal{F}(t) := \mathbb{P}(o^{\mathcal{G}}(t) \mid s^{\mathcal{G}}(t)).$$
(12)

• Note that $s^g(t) \in S$ and $o^g(t) \in O$.



Likelihood and Expectation Maximization in HMM

Likelihood



• According to the figure, the likelihood of occurrence of the state sequence S^{g} and the observation sequence \mathcal{O}^{g} is [5]: $\underbrace{\mathcal{L}} = \underbrace{\mathbb{P}(S^{g}, \mathcal{O}^{g})}_{\tau-1} = \underbrace{\mathbb{P}(s^{g}(1))}_{\tau-1} \underbrace{\mathbb{P}(s^{g}(t+1) \mid s^{g}(t))}_{t=1} \underbrace{\prod_{t=1}^{\tau} \mathbb{P}(o^{g}(t) \mid s^{g}(t))}_{t=1} = \underbrace{\mathbb{P}(s^{g}(1))}_{\tau-1} \underbrace{\prod_{t=1}^{\tau-1} \mathbb{P}(s^{g}(t+1) \mid s^{g}(t))}_{\tau-1} \underbrace{\prod_{t=1}^{\tau} \mathbb{P}(o^{g}(t) \mid s^{g}(t))}_{t=1} = \pi_{s^{g}(1)} \underbrace{\prod_{t=1}^{\tau-1} \mathbb{P}(s^{g}(t) \mid s^{g}(t))}_{\tau} \underbrace{\prod_{t=1}^{\tau} \mathbb{P}(o^{g}(t) \mid s^{g}(t))}_{\tau}.$ (13) Likelihood

We found:

$$\checkmark \quad L = \pi_{s^{g}(1)} \prod_{t=1}^{\tau} \widehat{a_{s^{g}(t), s^{g}(t+1)}} \prod_{t=1}^{\tau} b_{s^{g}(t), o^{g}(t)}.$$

- The log-likelihood is: $\mathbf{4} \quad \ell = \log(L) = \log \pi_{s^{g}(1)} + \sum_{t=1}^{\tau-1} \log(a_{s^{g}(t), s^{g}(t+1)}) + \sum_{t=1}^{\tau} \log(b_{s^{g}(t), o^{g}(t)}).$ (14)
- Let 1_i be a vector with entry one at index i, i.e., 1_i := [0,0,...,0,1,0...,0]^T. Also, 1_{s^g(1)} means the vector with entry one at the index of the first state in the sequence S^g. For example, if there are three possible states and a sequence of length three, s^g(1) = (2)s^g(2) = 1, s^g(3) = 3, we have 1_{s^g(1)} = [0,1]0]^T.

 The terms in this log-likelihood are:

$$a_{s^{g}(t-1),s^{g}(t)} = \prod_{i=1}^{n} \prod_{j=1}^{n} (a_{i,j})^{\mathbf{1}_{i}[j]} (a$$

Likelihood

Also:

where $\mathbf{1}_i[i] = \mathbf{1}_j[j] = 1$.

• Hence, we can write the log-likelihood:

$$\bigstar \quad \ell = \log \pi_{s^{g}(1)} + \sum_{t=1}^{\tau-1} \log(a_{s^{g}(t), s^{g}(t+1)}) + \sum_{t=1}^{\tau} \log(b_{s^{g}(t), o^{g}(t)}),$$

as:

$$\ell = \mathbf{1}_{s^{\mathcal{G}}(1)}^{\top} \log \pi + \sum_{t=1}^{\tau-1} \mathbf{1}_{s^{\mathcal{G}}(t-1)}^{\top} (\log \mathbf{A}) \mathbf{1}_{s^{\mathcal{G}}(t)} + \sum_{t=1}^{\tau} \mathbf{1}_{s^{\mathcal{G}}(t)}^{\top} (\log \mathbf{B}) \mathbf{1}_{o^{\mathcal{G}}(t)}.$$
(18)

We found:

$$\bigstar \quad \ell = \mathbf{1}_{s^{\mathcal{G}}(1)}^{\top} \log \pi + \sum_{t=1}^{\tau-1} \mathbf{1}_{s^{\mathcal{G}}(t-1)}^{\top} (\log \mathbf{A}) \mathbf{1}_{s^{\mathcal{G}}(t)} + \sum_{t=1}^{\tau} \mathbf{1}_{s^{\mathcal{G}}(t)}^{\top} (\log \mathbf{B}) \mathbf{1}_{o^{\mathcal{G}}(t)}.$$

• The missing variables in the log-likelihood are $\underline{1}_{s^g(t)}$, $\underline{1}_{s^g(t-1)}$, $\underline{1}_{s^g(t)}$, and $\underline{1}_{o^g(t)}$. The expectation of the log-likelihood with respect to the missing variables is:

$$\begin{array}{c}
\left(Q(\pi, \mathbf{A}, \mathbf{B})\right) = \mathbb{E}(\ell) \\
= \mathbb{E}(\mathbf{1}_{s^{\mathcal{E}}(1)}^{\top} \log \pi) + \sum_{t=1}^{\tau-1} \mathbb{E}(\mathbf{1}_{s^{\mathcal{E}}(t)}^{\top} (\log \mathbf{A}) \mathbf{1}_{s^{\mathcal{E}}(t+1)}) + \sum_{t=1}^{\tau} \mathbb{E}(\mathbf{1}_{s^{\mathcal{E}}(t)}^{\top} (\log \mathbf{B}) \mathbf{1}_{o^{\mathcal{E}}(t)}). \quad (19)
\end{array}$$

• We maximize the $Q(\pi, \mathbf{A}, \mathbf{B})$ with respect to the parameters $\pi_i(a_{i,i})$ and b_i .

$$\begin{aligned}
\widetilde{\mathcal{M}_{i}} ? \land i_{2} ; ? ; \\
\widetilde{\mathcal{M}_{i}} ? \land j_{2} ; ? ; \\
\widetilde{\mathcal{M}_{i}} ? : \\
\widetilde{\mathcal{M}_{i}} ? ; \\ \\
\widetilde{\mathcal{M}_{i}} ? ; \\
\widetilde{\mathcal{M}_{i}} ? ; \\
\widetilde{\mathcal{M}_{i}} ? ; \\
\widetilde{\mathcal{M}_{i}} ? ; \\ \\ \\
\widetilde{\mathcal{M}_{i}} ? ; \\ \\
\widetilde{\mathcal{M}_{i}} ? ; \\ \\ \\
\widetilde{\mathcal{M}_{i}} ?$$

where the constraints ensure that the probabilities in the initial states, the transition matrix, and the emission matrix add to one.

• The Lagrangian [6] for this optimization problem is:

$$\mathcal{L} = \underbrace{\left(\mathbb{E}(\mathbf{1}_{s^{\mathcal{G}}(1)}^{\top} \log \pi) + \sum_{t=1}^{\tau-1} \mathbb{E}(\mathbf{1}_{s^{\mathcal{G}}(t)}^{\top} (\log \mathbf{A}) \mathbf{1}_{s^{\mathcal{G}}(t+1)})\right) + \sum_{t=1}^{\tau} \mathbb{E}(\mathbf{1}_{s^{\mathcal{G}}(t)}^{\top} (\log \mathbf{B}) \mathbf{1}_{o^{\mathcal{G}}(t)}) \underbrace{\left(\sum_{i=1}^{n} \pi_{i} - 1\right) - (\eta_{i})}_{i=1} \sum_{j=1}^{n} a_{i,j} - 1) - (\eta_{i}) \underbrace{\left(\sum_{j=1}^{n} b_{i,j} - 1\right)}_{i=1}, \quad (21)$$

where η_1 , η_2 , and η_3 are the Lagrange multipliers.



• We had: $\mathcal{L} = \mathbb{E}(\mathbf{1}_{s^{\mathcal{G}}(1)}^{\top} \log \pi) + \sum_{t=1}^{t} \mathbb{E}(\mathbf{1}_{s^{\mathcal{G}}(t)}^{\top} (\log \mathbf{A}) \mathbf{1}_{s^{\mathcal{G}}(t+1)}) + \sum_{t=1}^{\tau} \mathbb{E}(\mathbf{1}_{s^{\mathcal{G}}(t)}^{\top} (\log \mathbf{B}) \mathbf{1}_{o^{\mathcal{G}}(t)}) - \eta_1(\sum_{i=1}^{n} \pi_i - 1) - \eta_2(\sum_{j=1}^{n} a_{i,j} - 1) - \eta_3(\sum_{j=1}^{n} b_{i,j} - 1).$ • Similarly, we have:



• We had: $\mathcal{L} = \mathbb{E}(\mathbf{1}_{s\mathfrak{f}(1)}^{\top} \log \pi) + \sum_{t=1}^{\tau-1} \mathbb{E}(\mathbf{1}_{s\mathfrak{f}(t)}^{\top} (\log \mathbf{A}) \mathbf{1}_{s\mathfrak{f}(t+1)}) + \sum_{t=1}^{\tau} \mathbb{E}(\mathbf{1}_{s\mathfrak{f}(t)}^{\top} (\log \mathbf{B}) \mathbf{1}_{o\mathfrak{f}(t)}) - \eta_1(\sum_{i=1}^n \pi_i - 1) - \eta_2(\sum_{j=1}^n a_{i,j} - 1) - \eta_3(\sum_{j=1}^n b_{i,j} - 1).$ • Likewise, we have:



Evaluation in HMM

Evaluation in HMM

• Evaluation in HMM means the following [7, 8]: Given the observation sequence $\underline{\mathcal{O}^g = \mathcal{O}^g(1), \ldots, \mathcal{O}^g(\tau)}$ and the HMM model $\lambda = (\pi, \mathbf{A}, \mathbf{B})$, we want to compute $\mathbb{P}(\mathcal{O}^g \mid \lambda)$, i.e., the probability of the generated observation sequence.

In summary:

$$\mathcal{O}^{g}, \text{given:} \left(\lambda \right) \Longrightarrow \left(\mathbb{P}(\mathcal{O}^{g} \mid \lambda) = ? \right)$$
 (31)

• Note that $\mathbb{P}(\mathcal{O}^g | \lambda)$ can also be denoted by $\mathbb{P}(\mathcal{O}^g; \lambda)$. The $\mathbb{P}(\mathcal{O}^g | \lambda)$ is sometimes referred to as the **likelihood**.

Direct Calculation

• Assume that the state sequence $S^{g} = s^{g}(1), \dots, s^{g}(\tau)$ has caused the observation sequence $\mathcal{O}^{g} = o^{g}(1), \dots, o^{g}(\tau)$. Hence, we have: $\begin{array}{c} \mathbb{P}(\mathcal{O}^{g} \mid S^{g}, \lambda) = t_{S^{g}(1)} (s(1)) (s^{g}(2)) (s^{g}(2)) \cdots (s^{g}(\tau)), s^{g}(\tau) \\ \mathbb{P}(\mathcal{O}^{g} \mid S^{g}, \lambda) = t_{S^{g}(1)} (s^{g}(2)) (s^{g}(2)) \cdots (s^{g}(\tau)), s^{g}(\tau) \\ \mathbb{P}(S^{g} \mid \lambda) = \pi_{s^{g}(1)} (s^{g}(2)) (s^{g}(2)) (s^{g}(2)) (s^{g}(2)) (s^{g}(2)) (s^{g}(\tau)) \\ \mathbb{P}(S^{g} \mid \lambda) = \pi_{s^{g}(1)} (s^{g}(2)) (s^{g}(2)) (s^{g}(2)) (s^{g}(2)) (s^{g}(2)) (s^{g}(2)) \\ \mathbb{P}(\mathcal{O}^{g}, S^{g} \mid \lambda) = \mathbb{P}(\mathcal{O}^{g} \mid S^{g}, \lambda) (s^{g}(2)) (s^{g}(2)) (s^{g}(2)) (s^{g}(2)) (s^{g}(2)) (s^{g}(2)) (s^{g}(2)) \\ \mathbb{P}(S^{g}(1)) (s^{g}(1)) (s^{g}(2)) (s^{g}(2)) (s^{g}(2)) (s^{g}(2)) (s^{g}(2)) (s^{g}(2)) (s^{g}(2)) (s^{g}(2)) \\ \mathbb{P}(S^{g}(1)) (s^{g}(1)) (s^{g}(2)) (s^{g}(2)) (s^{g}(2)) (s^{g}(2)) (s^{g}(2)) (s^{g}(2)) \\ \mathbb{P}(S^{g}(1)) (s^{g}(2)) (s^{g}(2)) (s^{g}(2)) (s^{g}(2)) (s^{g}(2)) (s^{g}(2)) \\ \mathbb{P}(S^{g}(1)) (s^{g}(2)) (s^{g}(2)) (s^{g}(2)) (s^{g}(2)) \\ \mathbb{P}(S^{g}(1)) (s^{g}(2)) (s^{g}(2)) (s^{g}(2)) (s^{g}(2)) \\ \mathbb{P}(S^{g}(1)) (s^{g}(2)) (s^{g}(2)) \\ \mathbb{P}(S^{g}(1)) (s^{g}(2)) (s^{g}(2)) \\ \mathbb{P}(S^{g}(1)) \\ \mathbb{P}(S^{g}(1)) \\ \mathbb{P}(S^{g}(1)) (s^{g}(2)) \\ \mathbb{P}(S^{g}(1)) \\ \mathbb{P$

which is the probability of occurrence of both the observation sequence \mathcal{O}^g and state sequence \mathcal{S}^g .

 Any state sequence may have caused the observation sequence O^g. Therefore, according to the law of total probability, we have:

$$\mathbb{P}(\mathcal{O}^{g}|\lambda) = \sum_{\forall S^{g}} \mathbb{P}(\mathcal{O}^{g}, S^{g}|\lambda) \stackrel{(34)}{=} \sum_{\forall S^{g}} \mathbb{P}(\mathcal{O}^{g}|S^{g}, \lambda) \mathbb{P}(S^{g}|\lambda)$$

$$\stackrel{(35)}{=} \sum_{\forall S^{g}(1)} \sum_{\forall S^{g}(2)} \cdots \sum_{\forall S^{g}(\tau)} \overline{\pi_{S^{g}(1)} b_{S^{g}(1), o^{g}(1)} a_{S^{g}}(\mathbf{s}) g^{g}(2) b_{S^{g}(\tau), o^{g}(\tau)} a_{S^{g}(\tau)} g^{g}(\tau)}, g^{g}(2) b_{S^{g}(\tau), o^{g}(\tau)} a_{S^{g}(\tau)} g^{g}(\tau), g^{g}(\tau)}, g^{g}(2) b_{S^{g}(\tau), o^{g}(\tau)} a_{S^{g}(\tau)} g^{g}(\tau)}, g^{g}(2) b_{S^{g}(\tau)} g^{g}(\tau) g^{g}(\tau)}, g^{g}(2) b_{S^{g}(\tau)} g^{g}(\tau), g^{g}(\tau)}, g^{g}(2) b_{S^{g}(\tau)} g^{g}(\tau)}, g^{g}(2) b_{S^{g}(\tau)} g^{g}(\tau), g^{g}(\tau)}, g^{g}(2) b_{S^{g}(\tau)} g^{g}(\tau)}, g^{g}(\tau)}, g^{g}(2) b_{S^{g}(\tau)} g^{g}(\tau)}, g^{g}(2) b_{S^{g}(\tau)} g^{g}(\tau)}, g$$

which means that we start with the first state, then output the first observation, and then go to the next state. This procedure is repeated until the last state. The summations are over all possible states in the state sequence.

Direct Calculation

We found:

$$\bigstar \quad \mathbb{P}(\mathcal{O}^{g} \mid \lambda) = \sum_{\forall s^{g}(1)} \sum_{\forall s^{g}(2)} \cdots \sum_{\forall s^{g}(\tau)} \pi_{s^{g}(1)} b_{s^{g}(1), o^{g}(1)} a_{s^{g}}(\mathbf{p}) g^{g}(2) b_{s^{g}(\tau), o^{g}(\tau)} a_{s^{g}(\tau), o^{g}(\tau)} a_{s^{g}(\tau)} a_{s$$

- The time complexity of this direct calculation of P(O^g | λ) is in the order of O(2τn^τ) because at every time clock t ∈ {1,...,τ}, there are n possible states to go through [7].
- Because of n^{τ} , this is very inefficient especially for long sequences (large τ).

- A more efficient algorithm for evaluation in HMM is forward-backward procedure [7].
- The forward-backward procedure includes two stages, i.e., forward and backward belief propagation stages.
- The Forward Belief Propagation:

• Similar to the belief propagation procedure, we define the forward message until time t as:

$$\alpha_i(t) := \mathbb{P}\left(o^g(1), o^g(2), \dots, o^g(t)\right) s^g(t) = s |\lambda\rangle, \quad \bigstar$$
(37)

which is the probability of partial observation sequence until time t and being in state s_i at time t.

• Algorithm 2 shows the forward belief propagation from state one to the state τ .

$$\begin{array}{c|c} \mathbf{I} \quad \mathbf{Input:} \ \underline{\lambda} = (\pi, \mathbf{A}, \mathbf{B}) \\ \mathbf{2} \quad \underline{\alpha_i(1)} = \overline{\pi_i} \ \underline{b_{i,o^g(1)}}, \quad \forall i \in \{1, \dots, n\} \\ \mathbf{3} \quad \mathbf{for} \ state \ j \ from 1 \ to \ n \ \mathbf{do} \\ \mathbf{4} \\ \mathbf{5} \\ \hline \\ \mathbf{6} \quad \underbrace{\mathbb{P}(\mathcal{O}^g \mid \lambda)}_{\mathbf{7}} = \underbrace{\sum_{i=1}^n \alpha_i(\tau)}_{\mathbf{7} \quad \mathbf{Return}} \underbrace{\mathbb{P}(\mathcal{O}^g \mid \lambda), \forall i, \forall t: \alpha_i(t)}_{\mathbf{7} \quad \mathbf{N} \quad \mathbf{c}} \\ \mathbf{7} \quad \mathbf{Return} \\ \end{array}$$

Algorithm 2: The forward belief propagation in the forward-backward procedure

$$\begin{array}{l} \mathbf{1} \ \mathbf{Input:} \ \lambda = (\boldsymbol{\pi}, \boldsymbol{A}, \boldsymbol{B}) \\ \mathbf{2} \ \alpha_i(1) = \pi_i \ b_{i,o^g(1)}, \quad \forall i \in \{1, \dots, n\} \\ \mathbf{3} \ \mathbf{for} \ state \ j \ from 1 \ to \ n \ \mathbf{do} \\ \mathbf{4} \ \left[\begin{array}{c} \mathbf{for} \ time \ t \ from 1 \ to \ (\tau - 1) \ \mathbf{do} \\ \mathbf{5} \ \left[\begin{array}{c} \alpha_j(t+1) = \left[\sum_{i=1}^n \alpha_i(t) \ a_{i,j} \right] \ b_{j,o^g(t+1)} \\ \mathbf{6} \ \mathbb{P}(\mathcal{O}^g \mid \lambda) = \sum_{i=1}^n \alpha_i(\tau) \\ \mathbf{7} \ \mathbf{Return} \ \mathbb{P}(\mathcal{O}^g \mid \lambda), \forall i, \forall t : \alpha_i(t) \end{array} \right]$$

Algorithm 2: The forward belief propagation in the forward-backward procedure

• In this algorithm, $\alpha_i(t)$ is solved inductively. The initial forward message is:

$$\alpha_i(1) = \pi_i (i, o^{\mathcal{E}}(1)) \qquad \forall i \in \{1, \dots, n\},$$
(38)

which is the probability of occurrence of the initial state s_i and the observation symbol $o^g(1)$.

• The next forward messages are calculated as:

$$\alpha_j(t+1) = \left[\sum_{i=1}^n \alpha_i(t) a_{i,j}\right] b_{j,o^g(t+1)},\tag{39}$$

which is the probability of occurrence of observation sequence o^g(1),..., o^g(t), being in state s_i at time t, going to state j at time t + 1, and the observation symbol o^g(t + 1).
The summation is because, at time t, the state s^g(t) can be any state so we should use the law of total probability.

$$\begin{array}{l} \mathbf{1} \ \mathbf{Input:} \ \lambda = (\pi, \mathbf{A}, \mathbf{B}) \\ \mathbf{2} \ \alpha_i(1) = \pi_i \ b_{i,o^g(1)}, \quad \forall i \in \{1, \dots, n\} \\ \mathbf{3} \ \mathbf{for} \ state \ j \ from 1 \ to \ n \ \mathbf{do} \\ \mathbf{4} \ \left[\begin{array}{c} \mathbf{for} \ time \ t \ from 1 \ to \ (\tau - 1) \ \mathbf{do} \\ \mathbf{5} \ \left[\begin{array}{c} \alpha_j(t+1) = \left[\sum_{i=1}^n \alpha_i(t) \ a_{i,j} \right] \ b_{j,o^g(t+1)} \\ \mathbf{6} \ \overline{\mathbb{P}(\mathcal{O}^g \mid \lambda) = \sum_{i=1}^n \alpha_i(\tau)} \\ \mathbf{7} \ \mathbf{Return} \ \mathbb{P}(\mathcal{O}^g \mid \lambda), \forall i, \forall t : \alpha_i(t) \end{array} \right]$$

Algorithm 2: The forward belief propagation in the forward-backward procedure

• Finally, using the law of total probability, we have:

$$\underbrace{\mathbb{P}(\mathcal{O}^{g} \mid \lambda) = \sum_{i=1}^{n} \mathbb{P}(\mathcal{O}^{g}, s^{g}(\tau) = s_{i} \mid \lambda) = \sum_{i=1}^{n} \alpha_{i}(\tau),}_{(40)}$$

which is the desired probability in the evaluation for HMM. Hence, the forward belief propagation suffices for evaluation.

• The Backward Belief Propagation:

 Again similar to the belief propagation procedure, we define the backward message since time τ to t + 1 as:

$$\beta_{i}(t) := \mathbb{P}(o^{g}(t+1), o^{g}(t+2), \dots, o^{g}(\tau) | s^{g}(t) = s_{i}, \lambda),$$

$$(41)$$

which is the probability of partial observation sequence from t + 1 until the end time τ given being in state s_i at time t.

• Algorithm 3 shows the backward belief propagation.

$$\begin{array}{c|c} \mathbf{Input:} & \underline{\lambda} = (\pi, \underline{A}, \underline{B}) \\ \mathbf{2} & \beta_i(\tau) = \underline{1}, \quad \forall i \in \{1, \dots, n\} \\ \mathbf{3} & \text{for state } i from 1 to n \, \mathbf{do} \\ \mathbf{4} & \\ \mathbf{5} & \begin{bmatrix} for \underbrace{time t from (\tau - 1) to 1 \, \mathbf{do}} \\ \beta_i(t) = \sum_{j=1}^n a_{i,j} \, b_{j,o^g(t+1)} \\ \beta_i(t) : \beta_i(t) \end{bmatrix} \\ \mathbf{6} & \text{Return } \forall i, \forall t : \beta_i(t) \\ \end{array}$$

Algorithm 3: The backward belief propagation in the forward-backward procedure

1 Input:
$$\lambda = (\pi, A, B)$$

2 $\beta_i(\tau) = 1, \quad \forall i \in \{1, \dots, n\}$
3 for state i from 1 to n do
4 $\left[\begin{array}{c} \text{for time t from } (\tau - 1) \text{ to 1 do} \\ \beta_i(t) = \sum_{j=1}^n a_{i,j} b_{j,o^s(t+1)} \\ 6 \text{ Return } \forall_i, \forall t : \beta_i(t) \end{array} \right]$

Algorithm 3: The backward belief propagation in the forward-backward procedure

• In this algorithm, the initial backward message is:

$$\beta_i(\tau) = 1, \quad \forall i \in \{1, \dots, n\}.$$
(42)

• The next backward messages are calculated as:

$$\beta_i(t) = \sum_{j=1}^{n} (a_{i,j}) b_{j,o^g(t+1)},$$
(43)
which is the probability of being in state s_i at time t , going to state j at time $t+1$ and
the observation symbol $o^g(t+1)$
The summation is because, at time $t+1$, the state $s^g(t+1)$ can be any state so we
should use the law of total probability.

1 Input: $\lambda = (\pi, A, B)$ 2 $\beta_i(\tau) = 1, \quad \forall i \in \{1, \dots, n\}$ 3 for state *i* from 1 to *n* do 4 5 $\begin{bmatrix} for time t from (\tau - 1) to 1 do \\ \beta_i(t) = \sum_{j=1}^n a_{i,j} b_{j,o^g(t+1)} \end{bmatrix}$ 6 Return $\forall i, \forall t : \beta_i(t)$

Algorithm 3: The backward belief propagation in the forward-backward procedure

 It is noteworthy that for very long sequences, the α_i(t) and β(t) become extremely small, recursively (see Algorithms 2 and 3). Hence, some people normalize them at every iteration of algorithm [5]:

$$\begin{array}{c} \alpha_{i}(t) \leftarrow \alpha_{i}(t) \\ \hline \begin{array}{c} \alpha_{i}(t) \leftarrow \overline{\sum_{j=1}^{\tau} \alpha_{j}(t)}, \\ \beta_{i}(t) \leftarrow \overline{\sum_{j=1}^{\tau} \beta_{j}(t)}, \\ \end{array} \end{array}$$
(44)
(45)

in order to sum to one. However, note that if this normalization is done, we will have $\mathbb{P}(\mathcal{O}^g \mid \lambda) = 1$.

Estimation in HMM

Estimation in HMM

$$\mathcal{S}^{g}, \text{given: } \mathcal{O}^{g}, \lambda \implies \mathbb{P}(\mathcal{S}^{g} \mid \mathcal{O}^{g}, \lambda) = ?$$

$$(46)$$

Greedy Approach

• Let the probability of being in state s_i at time t given the observation sequence \mathcal{O}^g and the HMM model λ be denoted by:

$$\underline{\gamma_i(t)} := \mathbb{P}(s^g(t) = s_i \mid \mathcal{O}^g, \lambda). \quad \bigstar$$
(47)

• We can say:

$$\gamma_{i}(t) \underbrace{\mathbb{P}(\mathcal{O}^{g} \mid \lambda)}_{(a)} = \underbrace{\mathbb{P}(s^{g}(t) = s_{i} \mid \mathcal{O}^{g}, \lambda)}_{(a)} \underbrace{\mathbb{P}(\mathcal{O}^{g} \mid \lambda)}_{(a)} \underbrace{\mathbb{P}(s^{g}(t) = s_{i}, \mathcal{O}^{g}, \lambda)}_{(b)} \underbrace{\mathbb{P}(s^{g}(t) \mid \beta_{i}(t))}_{(c)} \underbrace{\mathbb{P}(s^{g}(t) \mid \beta_{i}(t))}_{($$

• where (a) is because of chain rule in probability and (b) is because:

t

$$\approx \alpha_i(t) \beta_i(t) \underbrace{\mathbb{P}(o^g(1), o^g(2), \dots, o^g(t), s^g(t) = s_i \mid \lambda)}_{\times \mathbb{P}(o^g(t+1), o^g(t+2), \dots, o^g(\tau) \mid s^g(t) = s_i, \lambda)} = \underbrace{\mathbb{P}(o^g(1), o^g(2), \dots, o^g(\tau), s^g(t) = s_i, \lambda)}_{\mathbb{P}(o^g(1), o^g(2), \dots, o^g(\tau), s^g(t) = s_i, \lambda)} = \mathbb{P}(\mathcal{O}^g, s^g(t) = s_i, \lambda).$$

• The reason of (b) can also be interpreted in this way: it is because of the forward-backward procedure which states the **belief over a variable as product of the** forward and backward messages.

Greedy Approach

In the greedy approach, at every time t, we select a state with maximum probability of occurrence without considering the other states in the sequence. Therefore, we have:

$$s^{g}(t) = \arg\max_{1 \le i \le n} \gamma_{i}(t), \quad \forall t \in \{1, \dots, \tau\},$$
(50)

- The greedy approach does not optimize over the whole path but greedily chooses the best state at every time step.
- Another approach is to find the best state sequence which has the highest probability of occurrence, i.e., maximizing P(O^g, S^g | λ) [7]. The Viterbi algorithm (1967) [9, 10] can be used to find this path of states [11].
- Different works, such as [12], have worked on using Viterbi algorithm for HMM.
- Algorithm 4 shows the Viterbi algorithm for estimation in HMM [7].

$$\begin{array}{|c|c|c|c|c|c|} & \mathbf{I} \mbox{ Introduct} & \lambda = (\pi, A, B) \\ & 2 & // \mbox{ Introduction:} \\ & 3 & \delta_i(1) = \pi_i b_{i,o^g(1)}, & \forall i \in \{1, \dots, n\} \\ & \overline{\psi}_i(1) \equiv 0, & \forall i \in \{1, \dots, n\} \\ & 5 & \sqrt{7} \mbox{ Recursion:} \\ & 6 \mbox{ for state } j \mbox{ from 1 to n } \mathbf{do} \\ & 7 & \mathbf{Tor time tfrom 2 to \tau } \mathbf{do} \\ & 7 & \mathbf{Tor time tfrom 2 to \tau } \mathbf{do} \\ & 8 & \mathbf{f}_j(t) = \max_{1 \leq i \leq n} (\delta_i(t-1) a_{i,j}) \ b_{j,o^g(t)} \\ & 9 & \mathbf{f}_i(t) = \max_{1 \leq i \leq n} \delta_i(\tau) \\ & 11 & p^* = \max_{1 \leq i \leq n} \delta_i(\tau) \\ & 12 & s^*(\tau) = \arg\max_{1 \leq i \leq n} \delta_i(\tau) \\ & 13 & // \mbox{ Backtracking:} \\ & 14 \ \mathbf{for time tfrom \tau - 1 to 1 } \mathbf{do} \\ & 15 & \lfloor s^*(t) = \psi_{s^*(t+1)}(t+1) \\ & 16 & \mathbb{P}(\mathcal{O}^g, S^g \mid \lambda), S^g = s^*(1), \dots, s^*(\tau) \\ \end{array}$$

Algorithm 4: The Viterbi algorithm for estimation in HMM



except that $\alpha_j(t)$ in the forward belief propagation uses sum-product algorithm [13] while $\delta_j(t)$ in the Viterbi algorithm uses max-product algorithm [14, 15].

Similar to Eq. (38):

the initial
$$\delta_{j}(t)$$
 is:

$$\delta_{i}(1) = \pi_{i} b_{i,o^{\mathcal{E}}(1)}, \quad \forall i \in \{1, \dots, n\},$$

$$\delta_{i}(1) = \pi_{i} b_{i,o^{\mathcal{E}}(1)}, \quad \forall i \in \{1, \dots, n\}.$$
(52)

We define the index maximizing in Eq. (51):

$$\delta_j(t) = \overbrace{\substack{1 \le i \le n}}^{\max} \left[\delta_i(t-1) \, \mathbf{a}_{i,j} \right] b_{j,o^{\mathcal{B}}(t)},$$

as:

state
$$\psi_j(t) = \arg \max_{1 \le i \le n} \left(\delta_i(t-1) a_{i,j} \right).$$
 (53)

• Then, a <u>backward analysis</u> is done starting from the end of state sequence:

$$\mathbf{Y} = \max_{1 \le i \le n} \delta_i(\tau),$$

$$s^*(\tau) = \arg \max_{1 \le i \le n} \delta_i(\tau),$$

$$(54)$$

$$(55)$$

and the other states in the sequence are <u>backtracked</u> as:

$$s^{*}(t) = \psi_{s^{*}(t+1)}(t+1).$$
(56)

The states S^g = s*(1), s*(2), ..., s*(τ) are the desired state sequence in the estimation.
 Therefore, the states in the state sequence are maximizing the forward belief propagation in a max-product setting.

• The probability of this path of states with maximum probability of occurrence is:

$$\mathbb{P}(\mathcal{O}^{g}, \mathcal{S}^{g} \mid \lambda) = p^{*}.$$
(57)

Note that the Viterbi algorithm can be visualized using a trellis structure (see Appendix A in [16]).

Training the HMM

Training the HMM

- Training HMM means the following [7, 8]: Given the observation sequence $\overline{\mathcal{O}^g} = o_1^g, \dots, o_T^g$, we want to adjust the HMM model parameters $\lambda = (\pi, A, B)$ in order to maximize $\mathbb{P}(\mathcal{O}^g \mid \lambda)$.
- In summary:

given:
$$\mathcal{O}^{g}, \mathcal{O}, \mathcal{S} \implies \bigwedge^{\lambda} = \arg \max_{\lambda} \mathbb{P}(\mathcal{O}^{g} \mid \lambda).$$

(58)

• We can solve for Eq. (58):

using maximum likelihood estimation using Expectation Maximization (EM).

- The <u>Baum-Welch algorithm</u> (1970) [17] is the <u>most well-known method</u> for training <u>HMM</u>. It makes use of the <u>EM</u> results.
- We define the probability of occurrence of a path being in states s_i and s_j, respectively, at times t and t+1 by:

• We can say:

$$\begin{aligned}
\xi_{i,j}(t) &:= \mathbb{P}(s^{g}(t) = s_{i})s^{g}(t+1) = s_{j} | \mathcal{O}^{g}, \lambda \rangle \quad (59) \\
\end{aligned}$$

$$\begin{aligned}
\underbrace{\xi_{i,j}(t) \mathbb{P}(\mathcal{O}^{g} | \lambda) = \mathbb{P}(s^{g}(t) = s_{i}, s^{g}(t+1) = s_{j} | \mathcal{O}^{g}, \lambda) \mathbb{P}(\mathcal{O}^{g} | \lambda) \\
\overset{(a)}{=} \mathbb{P}(s^{g}(t) = s_{i}, s^{g}(t+1) = s_{j}, \mathcal{O}^{g}, \lambda) \stackrel{(b)}{=} \alpha_{i}(t) a_{i,j} b_{j,os(t+1)} \beta_{j}(t+1) \\
&= \underbrace{\xi_{i,j}(t) = \frac{\alpha_{i}(t) a_{i,j} b_{j,os(t+1)} \beta_{j}(t+1)}{\mathbb{P}(\mathcal{O}^{g} | \lambda)}}_{\underset{r=1}{\sim} \underbrace{\xi_{i,j}(t) = \frac{\alpha_{i}(t) a_{i,j} b_{j,os(t+1)} \beta_{j}(t+1)}{\mathbb{P}(\mathcal{O}^{g} | \lambda)}}_{(61)} \\
\end{aligned}$$

where (a) is because of chain rule in probability and the reason of (b) is in the next slide.

We found:

$$\xi_{i,j}(t) \mathbb{P}(\mathcal{O}^{g} \mid \lambda) = \mathbb{P}(s^{g}(t) = s_{i}, s^{g}(t+1) = s_{j} \mid \mathcal{O}^{g}, \lambda) \mathbb{P}(\mathcal{O}^{g} \mid \lambda)$$

$$\stackrel{(a)}{=} \mathbb{P}(s^{g}(t) = s_{i}, s^{g}(t+1) = s_{j}, \mathcal{O}^{g}, \lambda) \stackrel{(b)}{=} \alpha_{i}(t) a_{i,j} b_{j,o^{g}(t+1)} \beta_{j}(t+1)$$

• (b) is because of the following: According to Eqs. (37), (4), (6), and (41), we have:

$$\begin{array}{c} \bigstar & \alpha_{i}(t) = \mathbb{P}(o^{g}(1), \dots, o^{g}(t), s^{g}(t) = s_{i} \mid \lambda), \\ \bigstar & a_{i,j} = \mathbb{P}(s_{j}(t+1) \mid s_{i}(t)), \\ \bigstar & b_{j,o^{g}(t+1)} = \mathbb{P}(o^{g}(t+1) \mid s_{j}(t+1)), \\ \bigstar & \beta_{j}(t+1) = \mathbb{P}(o^{g}(t+2), \dots, o^{g}(\tau) \mid s^{g}(t+1) = s_{j}, \lambda). \end{array}$$

Therefore, we have:

$$\alpha_i(t) a_{i,j}(t) b_{j,o^g(t+1)} \beta_j(t+1) = \mathbb{P}\big(s^g(t) = s_i, s^g(t+1) = s_j, \mathcal{O}^g, \lambda\big).$$

• Note that we have $\sum_{i=1}^{n} \sum_{j=1}^{n} \xi_{i,j}(t) = 1$. Also, note that $\mathbb{P}(\mathcal{O}^{g} | \lambda)$ in Eq. (60) can be obtained from either the denominator of Eq. (61) or line 6 in Algorithm 2 (i.e., Eq. (40)):

$$\mathbb{P}(\mathcal{O}^{g} \mid \lambda) = \sum_{i=1}^{n} \mathbb{P}(\mathcal{O}^{g}, s^{g}(\tau) = s_{i} \mid \lambda) = \sum_{i=1}^{n} \alpha_{i}(\tau).$$

In Eq. (60):

 $\xi_{i,j}(t) = \frac{\alpha_i(t) a_{i,j} b_{j,o^g(t+1)} \beta_j(t+1)}{\mathbb{P}(\mathcal{O}^g \mid \lambda)},$ the terms $\alpha_i(t) a_{i,j}(t) b_{j,o^g(t+1)}$ and $\beta_j(t+1)$ stand for the probability of the first t observations ending in state s_i at time t, the probability of transitioning from state s_i (at time t) to state s_j (at time t+1), the probability of observing $o^g(t+1)$ from state s_j at time t+1, and the probability of the remainder of the observation sequence, respectively.

• Now, recall the Eqs. (24), (27), and (30) from EM algorithm for HMM:

$$\bigstar_{i,j} = \mathbb{E}(\mathbf{1}_{s^{g}(1)}[i]), \quad a_{i,j} = \frac{\sum_{t=1}^{\tau-1} \mathbb{E}(\mathbf{1}_{s^{g}(t)}[i] \, \mathbf{1}_{s^{g}(t+1)}[j])}{\sum_{t=1}^{\tau-1} \mathbb{E}(\mathbf{1}_{s^{g}(t)}[i])}, \quad b_{i,j} = \frac{\sum_{t=1}^{\tau} \mathbb{E}(\mathbf{1}_{s^{g}(t)}[i] \, \mathbf{1}_{o^{g}(t)}[j])}{\sum_{t=1}^{\tau} \mathbb{E}(\mathbf{1}_{s^{g}(t)}[i])}.$$

On the other hand, recall Eqs. (47) and (59):

 \mathbb{E}

$$\gamma_i(t) = \mathbb{P}(s^g(t) = s_i | \mathcal{O}^g, \lambda), \quad \xi_{i,j}(t) = \mathbb{P}(s^g(t) = s_i, s^g(t+1) = s_j | \mathcal{O}^g, \lambda).$$

$$(\mathbf{1}_{s^{\mathcal{G}}(1)}[i]) = \gamma_i(1), \mathcal{A}$$
(62)

$$\mathbb{E}(\mathbf{1}_{s^{\mathcal{G}}(t)}[i]\,\mathbf{1}_{s^{\mathcal{G}}(t+1)}[j]) = \xi_{i,j}(t), \not\triangleleft$$
(64)

$$\mathbb{E}(\mathbf{1}_{s^{g}(t)}[i] \mathbf{1}_{o^{g}(t)}[j]) = \gamma_{i}(t), \text{ where } o^{g}(t) = j \text{ in } \gamma_{i}(t).$$
(65)

Hence:

$$\left(\underbrace{\pi_{i}=\gamma_{i}(1)}_{\tau=1}\overset{(47)}{=} \mathbb{P}\left(s^{g}(1)=s_{i} \mid \mathcal{O}^{g}, \lambda\right), \quad \forall i \in \{1, \dots, n\}, \qquad (66)$$

$$\underbrace{ b_{i,j} = \frac{\sum_{t=1, o^g(t)=j}^{\tau} \gamma_i(t)}{\sum_{t=1}^{\tau} \gamma_i(t)}, \quad \forall i \in \{1, \dots, n\}, \forall j \in \{1, \dots, m\}.$$

• We found:

$$\bigstar \quad b_{i,j} = \frac{\sum_{t=1, o^{g}(t)=j}^{\tau} \gamma_i(t)}{\sum_{t=1}^{\tau} \gamma_i(t)}, \quad \forall i \in \{1, \dots, n\}, \forall j \in \{1, \dots, m\}.$$

• With change of variable, we have:

$$\bigstar \quad b_{j,k} = \frac{\sum_{t=1, o^g(t)=k}^{\tau} \gamma_j(t)}{\sum_{t=1}^{\tau} \gamma_j(t)}, \quad \forall j \in \{1, \dots, n\}, \forall k \in \{1, \dots, m\}.$$
(68)

• The algorithm is shown in Algorithm 5 [7]. In this algorithm the initial probabilities of being in state s_i at time t = 1 is according to Eq. (66). Then the a_{i,j} is then calculated using Eq. (67). According to counting in probability, it can also be interpreted as the ratio of the expected number of transitions from state s_i to s_j over the expected number of transitions out of state s_i. Finally, the b_{j,k} is calculated using Eq. (68). Likewise, using counting in probability, the b_{j,k} can be interpreted as the ratio of the expected number of times being in state s_j and seeing observation o_k over the expected number of times being in state s_i.

Algorithm 5: The Baum-Welch algorithm for training the HMM

• For calculating Eq. (68):

we use Eq. (49):

$$\bigstar \quad \gamma_i(t) = \frac{\alpha_i(t)\,\beta_i(t)}{\sum_{j=1}^n \alpha_j(t)\,\beta_j(t)},$$

where:

$$\alpha_{j}(t+1) = \left[\sum_{i=1}^{n} \alpha_{i}(t) \mathbf{a}_{i,j}\right] \mathbf{b}_{j,k},$$

$$\beta_{i}(t) = \sum_{j=1}^{n} \mathbf{a}_{i,j} \mathbf{b}_{j,k},$$
(69)
(70)

in line 5 in Algorithm 2 and in line 5 in Algorithm 3, respectively. We use the obtained $\alpha_i, \forall i, \forall t$ and $\beta_i(t), \forall i, \forall t$ for calculating Eq. (49).

Applications of HMM

Application in Speech Recognition

- Assume we have a dictionary of words consisting of |W| words.
- For every word indexed by $w \in \{1, ..., |\mathcal{W}|\}$, we have $|\mathcal{Q}_w|$ training instances spoken by one or several people. The training instances for a word are indexed by q where $\overline{q} \in \{1, ..., |\mathcal{Q}_w|\}$.
- Every training instance is a sequence of observation symbols obtained from formants [18].
- We consider an HMM model for every word in the dictionary.
- Training the HMMs are as [7, 8]:
 For every word w_i, consider the training sequences, O^g_w, indexed by <u>q</u>, i.e., O^g₁, ..., O^g_{|Q_w|};
 Train the HMM for the <u>w-th word</u> using Algorithm 6 to obtain λ_w.
 For an <u>unknown test word</u> with sequence O^g_t = o^g₁, ..., o^g_{|Q_t|}, we recognize the word as [7, 8]:
 - **(**) Calculate $\mathbb{P}(\mathcal{O}_t^{\mathcal{G}} | \lambda_w)$ for all $w \in \{1, \dots, |\mathcal{W}|\}$ using the forward belief propagation, i.e., Algorithm 2 (see Eq. (40)).
 - 2 The test word is recognized as:

$$w^* = \arg\max_{w} \mathbb{P}(\mathcal{O}_t^g \mid \lambda_w).$$
(71)

So, the test word is recognized as the w-th word in the dictionary.

Application in Speech Recognition

- In test phase for speech recognition, usually, the <u>Viterbi algorithm is used</u> [8]. Hence, another approach for recognition of the test word O^g_t is:
 - Calculate $\mathbb{P}(\mathcal{O}_t^{\mathcal{S}}, \mathcal{S}_t^{\mathcal{G}} | \lambda_w)$ for all $w \in \{1, \dots, |W|\}$ using the Viterbi algorithm, i.e., Algorithm 4 (see Eq. (57)).
 - 2 The test word is recognized as:

$$w^* = \arg\max_{w} \mathbb{P}(\mathcal{O}_t^g, \mathcal{S}_t^g \mid \lambda_w).$$
(72)

So, the test word is recognized as the *w*-th word in the dictionary.

Note that the words are pronounced with <u>different lengths</u> (fast or slowly) by different people. As <u>HMM is robust to different repetitions of states</u>, the recognition of words with different pacing is possible.

Application in Action Recognition

- In action recognition, every action can be seen as a sequence of poses where every pose may be repeated for several frames [19]. Hence, HMM can be used for action recognition [20].
- Assume we have a set of actions denoted by \mathcal{W} where the actions are indexed by $w \in \{1, \dots, |\mathcal{W}|\}$.
- We have |Q_w| training instances for every action where the training instances are indexed by q ∈ {1,..., |Q_w|}.
- Every training instance is a **sequence of observation symbols** where the symbols are the **poses**, e.g., sitting, standing, etc.
- An HMM is trained for every action [19, 21].
- Training and testing HMMs for action recognition is the same as training and test phases explained for speech recognition.
- The actions are performed with **different sequence lengths (fast or slowly)** by different people. As **HMM is robust to different repetitions of states**, the recognition of actions with different pacing is possible.

Application in Action Recognition

 In action recognition, we have a dataset of actions consisting of several defined poses [19, 21]. For example, if the dataset includes three actions sit, stand, and turn, the format of actions is as follows:



where the actions are modeled as sequences of some poses, i.e., stand, sit, and tilt.

- The actions can have different lengths or pacing.
- An example training dataset with its instances is shown in Table 1. In some sequences of dataset, there are some noisy poses in the middle of sequences of correct poses for making a difficult instance.
- An example test dataset is also shown in Table 2. The three test sequences are different from the training sequences to check the generalizability of the HMM models.
- Three HMM models can be trained for the three actions in this dataset and then, the test
 action sequences can be fed to the HMM models to be recognized.
- See our paper "Fisherposes for human action recognition using Kinect sensor data" (2017) [19].

Application in Action Recognition

Action	Sequence	t = 1	t = 2	t = 3	t = 4	t = 5	t = 6	t = 7	t = 8	t = 9	t = 10	t = 11	t = 12	_
	1	stand	stand	stand	sit	sit	sit	×	×	×	×	×	×	-
	2	stand	stand	sit	sit	sit	sit	sit	sit	×	×	×	×	
Sit	3	stand	stand	stand	stand	stand	sit	sit	×	×	×	×	×	
	4	stand	stand	stand	stand	stand	sit	sit	sit	sit	sit	×	×	
	5	stand	stand	stand	stand	stand	sit	sit	sit	sit	stand	sit	sit	_
Stand	1	sit	sit	sit	stand	stand	stand	×	×	×	×	×	×	
	2	sit	sit	sit	sit	sit	sit	stand	stand	×	×	×	×	
	3	sit	sit	stand	stand	stand	stand	stand	×	×	×	×	×	
	4	sit	sit	sit	sit	stand	stand	stand	stand	stand	×	×	×	
	5	sit	sit	sit	sit	stand	stand	stand	sit	stand	stand	stand	×	
	1	stand	stand	stand	tilt	tilt	tilt	×	×	×	×	×	×	
	2	stand	stand	tilt	tilt	tilt	tilt	tilt	tilt	×	×	×	×	z
Turn	3	stand	stand	stand	stand	stand	tilt	tilt	×	×	×	×	×	1
	4	stand	stand	stand	stand	tilt	tilt	tilt	tilt	tilt	tilt	×	×	
	5	stand	stand	stand	stand	tilt	tilt	tilt	stand	tilt	tilt	tilt	×	

Table 1. An example training dataset for action recognition

Action	t = 1	t = 2	t = 3	t = 4	t = 5	t = 6	t = 7	t = 8	t = 9	t = 10	t = 11
Sit	stand	stand	stand	sit	sit	sit	sit	×	×	×	×
Stand	sit	sit	sit	sit	stand	stand	stand	×	×	×	×
Turn	stand	stand	stand	stand	tilt	tilt	stand	stand	tilt	tilt	tilt

Table 2. An example test dataset for action recognition

Acknowledgment

- Some slides are based on our tutorial paper: "Hidden Markov Model: Tutorial" [22]
- For more information on HMM, refer to our tutorial paper [22].
- Some slides of this slide deck are inspired by teachings of <u>Prof. Mehdi Molkaraie</u>, Prof. <u>Kevin Granville</u>, and <u>Prof. Mu Zhu</u> at University of Waterloo, Department of Statistics.
- HMM in sklearn: https://scikit-learn.sourceforge.net/stable/modules/hmm.html

References

- D. Koller and N. Friedman, Probabilistic graphical models: principles and techniques. MIT press, 2009.
- [2] S. M. Ross, *Introduction to probability models*. Academic press, 2014.
- [3] G. F. Lawler, *Introduction to stochastic processes*. Chapman and Hall/CRC, 2018.
- [4] T. L. Booth, Sequential machines and automata theory. New York, NY: Wiley, 1967.
- [5] Z. Ghahramani, "An introduction to hidden Markov models and Bayesian networks," in Hidden Markov models: applications in computer vision, pp. 9–41, World Scientific, 2001.
- [6] S. Boyd and L. Vandenberghe, *Convex optimization*. Cambridge university press, 2004.
- [7] L. R. Rabiner and B.-H. Juang, "An introduction to hidden Markov models," IEEE ASSP Magazine, vol. 3, no. 1, pp. 4–16, 1986.
- [8] L. R. Rabiner, "A tutorial on hidden markov models and selected applications in speech recognition," *Proceedings of the IEEE*, vol. 77, no. 2, pp. 257–286, 1989.
- [9] A. Viterbi, "Error bounds for convolutional codes and an asymptotically optimum decoding algorithm," *IEEE transactions on Information Theory*, vol. 13, no. 2, pp. 260–269, 1967.

References (cont.)

- [10] G. D. Forney, "The Viterbi algorithm," *Proceedings of the IEEE*, vol. 61, no. 3, pp. 268–278, 1973.
- [11] P. Blunsom, "Hidden Markov models," tech. rep., 2004.
- [12] Y. He, "Extended Viterbi algorithm for second order hidden markov process," in [1988 Proceedings] 9th International Conference on Pattern Recognition, pp. 718–720, IEEE, 1988.
- [13] F. R. Kschischang, B. J. Frey, and H.-A. Loeliger, "Factor graphs and the sum-product algorithm," *IEEE Transactions on information theory*, vol. 47, no. 2, pp. 498–519, 2001.
- [14] Y. Weiss and W. T. Freeman, "On the optimality of solutions of the max-product belief-propagation algorithm in arbitrary graphs," *IEEE Transactions on Information Theory*, vol. 47, no. 2, pp. 736–744, 2001.
- [15] J. Pearl, Probabilistic reasoning in intelligent systems: networks of plausible inference. Elsevier, 2014.
- [16] D. Jurafsky and J. H. Martin, Speech and Language Processing: An introduction to speech recognition, computational linguistics and natural language processing. Pearson Prentice Hall, 2019.
- [17] L. E. Baum, T. Petrie, G. Soules, and N. Weiss, "A maximization technique occurring in the statistical analysis of probabilistic functions of Markov chains," *The annals of mathematical statistics*, vol. 41, no. 1, pp. 164–171, 1970.

References (cont.)

- [18] I. R. Titze and D. W. Martin, *Principles of voice production*. Acoustical Society of America, 1998.
- [19] B. Ghojogh, H. Mohammadzade, and M. Mokari, "Fisherposes for human action recognition using kinect sensor data," *IEEE Sensors Journal*, vol. 18, no. 4, pp. 1612–1627, 2017.
- [20] J. Yamato, J. Ohya, and K. Ishii, "Recognizing human action in time-sequential images using hidden Markov model," in *Proceedings 1992 IEEE Computer Society conference on computer vision and pattern recognition*, pp. 379–385, IEEE, 1992.
- [21] M. Mokari, H. Mohammadzade, and B. Ghojogh, "Recognizing involuntary actions from 3d skeleton data using body states," *Scientia Iranica*, 2018.
- [22] B. Ghojogh, F. Karray, and M. Crowley, "Hidden Markov model: Tutorial," 2019.