

is also a directed graphical model.

- Settings:
- * Q_t : hidden units random variable, q_t : observed hidden unit
 - O_t : observation random variable, o_t : observed observation.
 - * N : # of possible states, denoted by S_1, S_2, \dots, S_N
 - * M : # of possible observations denoted by V_1, V_2, \dots, V_M
 - * $\pi = (\pi_1, \pi_2, \dots, \pi_N)$ where $\pi_i = P(Q_1 = S_i)$
initial state.
 - * A : transition matrix, where $A_{ij} = P(Q_t = S_j | Q_{t-1} = S_i)$.
 - * B : observation probability matrix,
where $B_{jk} \equiv b_j(k) = P(O_t = V_k | Q_t = S_j)$.

we put $\lambda = (\pi, A, B)$ as parameter of HMM.

- * Merkov Property:

$$P(Q_t | Q_{1:t-1}, O_{1:t-1}) = P(Q_t | Q_{t-1}) \quad \forall t$$

- * Joint Probability:

$$\begin{aligned} P(q_{1:T}, o_{1:T}) &= P(q_1) \cdot P(o_1 | q_1) \cdot \prod_{t=2}^T P(q_t | q_{t-1}) \cdot P(o_t | q_t) \\ &= \pi_1 \cdot B_{11} \cdot \prod_{t=1}^T A_{q_{t-1}, q_t} \cdot B_{q_T} \end{aligned}$$

We can sample a sequence from $\lambda = (\pi, A, B)$.

Three Basic Problems in HMM Learning.

* Note M, N are fixed in all three basic problems

1. Evaluation Problem.

Given fixed parameter $\lambda = (\pi, A, B)$, evaluate the probability of an observed sequence $O^{1:T}$

i.e. $P(O^{1:T} | \lambda)$.

Both λ and $O^{1:T}$ are fixed

Example * given λ , $O^{1:T}$ and $O'^{1:T}$, which sequence is more likely to happen under λ ?

* given $O^{1:T}$ and λ, λ' , which parameter is a better fit to $O^{1:T}$ Solve by Dynamic Programming

2. Decoding problem:

Given fixed λ and fixed $O^{1:T}$, what are the best hidden states sequence ($q_{1:T}$) corresponding to these observations?

Solve by Dynamic Programming.

3. Learning Parameters:

Given N, M and $O^{1:T}$, what is best parameters $\lambda = (\pi, A, B)$?

Find a λ that maximises the likelihood function $P(O^{1:T} | \lambda)$.

(if T is large enough, one training sample is enough to learn λ)

Solve by EM

1: Evaluation Problem:

Given $N, M, \lambda, o_{1:T}$, Compute $P(o_{1:T}|\lambda)$

Let Σ be the set of all possible sequences of hidden states.

From the law of total probability:

$$\begin{aligned} P(o_{1:T}|\lambda) &= \sum_{Q_{1:T} \in \Sigma} P(o_{1:T}, q_{1:T}|\lambda) \\ &= \sum_{Q_{1:T} \in \Sigma} P(o_{1:T}|q_{1:T}, \lambda) \cdot P(q_{1:T}|\lambda). \quad \dots \text{--- } \textcircled{\times} \end{aligned}$$

$$\text{where } P(o_{1:T}|q_{1:T}, \lambda) = \prod_{t=1}^T P(o_t|q_t, \lambda) = \prod_{t=1}^T B_{tq_t}$$

$$P(q_{1:T}|\lambda) = \pi_1 \cdot \prod_{t=2}^T A_{q_{t-1}, q_t}$$

Note we need to generate $|\Sigma|$ sequences of hidden states and observations. Each state can take values in N possible states and there are T time steps, $|\Sigma| = N^T$

This makes $\textcircled{\times}$ too complex to solve.

Solution: Dynamic Programming

Target: find $P(o_{1:T}|\lambda)$.

Note O_t only depends on Q_t .

$$\Rightarrow P(o_{1:T}|\lambda) \xrightarrow{\quad} P(o_{1:T-1}|\lambda) \quad \xrightarrow{\quad} P(o_T|\lambda).$$

Q_t only depends on Q_{t-1}

Then combine them by enumerating all possible states

of Q_{T-1} and Q_T complexity: N^2

$$\text{Similarly, } P(o_{1:T-1}|\lambda) \xrightarrow{\quad} P(o_{1:T-2}|\lambda) \quad \xrightarrow{\quad} P(o_{T-1}|\lambda).$$

This then becomes a Dynamic Programming Setting.

1.1 Forward variable and algorithm

$$\begin{aligned} P(O_{1:T} | \lambda) &= \sum_{i=1}^N P(O_{1:T}, Q_T = s_i | \lambda) \\ &= \sum_{i=1}^N P(O_{1:T-1}, Q_T = s_i | \lambda) B_{iT} \end{aligned}$$

then:

$$\begin{aligned} &P(O_{1:T-1}, Q_T = s_i | \lambda) \\ &= \sum_{j=1}^N P(O_{1:T-1}, Q_T = s_i, Q_{T-1} = s_j | \lambda). \\ &= \sum_{j=1}^N P(O_{1:T-1}, Q_{T-1} = s_j | \lambda) \cdot P(Q_T = s_i | O_{1:T-1}, Q_{T-1} = s_j | \lambda) \\ &= \sum_{j=1}^N P(O_{1:T-1}, Q_{T-1} = s_j | \lambda) \cdot P(Q_T = s_i | Q_{T-1} = s_j, \lambda) \\ &= \sum_{j=1}^N P(O_{1:T-1}, Q_{T-1} = s_j | \lambda) \cdot A_{ji} \end{aligned}$$

Note that:

$$P(O_{1:T-1} | \lambda) = \sum_{j=1}^N P(O_{1:T-1}, Q_{T-1} = j | \lambda)$$

if we can compute $P(O_{1:T-1}, Q_{T-1} = j | \lambda) \quad \forall j \in \{1, \dots, N\}$.

we can evaluate $P(O_{1:T-1} | \lambda)$ hence $P(O_{1:T} | \lambda)$ for time T.

The FORWARD ALGORITHM:

$$\alpha_t(i) = P(O_{1:t}, Q_t = s_i | \lambda)$$

the recursion:

$$\alpha_{t+1}(i) = \left(\sum_{j=1}^N \alpha_t(j) \cdot A_{ji} \right) \cdot B_{it+1}$$

Pseudo-Code:

1. Initialise: $\alpha_1(i) = \pi_i B_{1i} \quad \forall i \in \{1, \dots, N\}$.

2. Forward recursion: $\forall t \in \{1, 2, \dots, T-1\}, \forall i \in \{1, \dots, N\}$.

$$\alpha_{t+1}(i) = \left(\sum_{j=1}^N \alpha_t(j) \cdot A_{ji} \right) \cdot B_{it+1}$$

3. Output: $P(O_{1:T} | \lambda) = \sum_{i=1}^N \alpha_T(i)$

Complexity here is $O(TN^2)$, much better than N^T

1.2 Backward variable and algorithm

When at time t , look back:

$\alpha_t(i)$ is the probability that $Q_t = S_i$ and $O_{1:t}$ are observed.

What about the future? Still have $O_{t+1:T}$.

Define: $f_t(i) = P(O_{t+1:T} | Q_t = S_i, \lambda)$

$f_t(i)$ is the probability of observing future output sequence $O_{t+1:T}$ if

$Q_t = S_i$ at time t

Recursive relationship:

$$f_t(i) = \sum_{j=1}^N A_{ij} B_{j,t+1} f_{t+1}(j).$$

This must be initialised at $t=T$ and move backward.

and $f_T(i)$ is the probability of observing NOTHING after time T

give $Q_T = S_i$. Here $f_T(i) = 1$.

$$\text{Finally } P(O_{1:T} | \lambda) = \sum_{i=1}^N \pi_i B_{i1} f_1(i).$$

Forward - Code:

1. Initialisation: $f_T(i) = 1, \forall i \in \{1, \dots, N\}$.

2. Backward Recursion: $\forall t \in \{T-1, T-2, \dots, 2, 1\}$ and $\forall i \in \{1, \dots, N\}$

$$f_t(i) = \sum_{j=1}^N A_{ij} B_{j,t+1} f_{t+1}(j).$$

3. Output

$$P(O_{1:T} | \lambda) = \sum_{i=1}^N \pi_i B_{i1} f_1(i).$$

Note: Forward and Backward algorithms are calculating the same thing: $P(O_{1:T} | \lambda)$. so they must match.

2. Decoding Problem: Find the Best Hidden States.

Idea: find the maximum likelihood state for each time step independently.

i.e. $\forall t \in \{1, \dots, T\}$, find $P(Q_t = s_i | o_{1:T}, \lambda)$.

Define: $\tau_t(i) = P(Q_t = s_i | o_{1:T}, \lambda)$

Set: $g_t = \underset{i \in \{1, \dots, N\}}{\operatorname{argmax}} \tau_t(i)$

Note: $\tau_t(i)$ is the probability of Q_t being s_i when we observed the complete sequence $o_{1:T}$

$$\begin{aligned} \tau_t(i) &= P(Q_t = s_i | o_{1:T}, \lambda) = \frac{P(Q_t = s_i, o_{1:T} | \lambda)}{P(o_{1:T} | \lambda)} \xrightarrow{\text{By product:}} \\ &= \frac{P(o_{1:t}, o_{t+1:T} | Q_t = s_i, \lambda) \cdot P(Q_t = s_i | \lambda)}{\sum_{j=1}^N P(Q_t = s_j, o_{1:T} | \lambda)} \quad P(o_{1:T} | \lambda) = \sum_{i=1}^N \alpha_t(i) \cdot \beta_t(i). \\ &\quad \forall t. \\ &= \frac{P(o_{1:t} | Q_t = s_i, \lambda) \cdot P(o_{t+1:T} | Q_t = s_i, \lambda) \cdot P(Q_t = s_i | \lambda)}{\sum_{j=1}^N P(Q_t = s_j, o_{1:T} | \lambda)} \\ &\quad \underbrace{P_{t+1}(i)}_{\substack{\sum_{j=1}^N P(Q_t = s_j, o_{1:t} | \lambda)}} \quad \underbrace{\alpha_t(i)}_{\substack{\alpha_{t+1}(i)}} \\ &= \frac{P(o_{t+1:T} | Q_t = s_i, \lambda) \cdot P(Q_t = s_i, o_{1:t} | \lambda)}{\sum_{j=1}^N P(Q_t = s_j, o_{1:T} | \lambda)} \\ &= \frac{\alpha_t(i) \cdot \beta_t(i)}{\sum_{j=1}^N \alpha_t(j) \cdot \beta_t(j)}. \end{aligned}$$

To compute τ , we first set $\tau_t(i) = \alpha_t(i) \cdot \beta_t(i)$. Then we normalise τ $\forall t$. Note we can find g_t by finding the maximum element in unnormalised τ values. Since normalisation constant $P(o_{1:T} | \lambda)$ will not change the index of largest τ value at time t .

Practical Problem: This finds the best state $\forall t \in \{1, \dots, T\}$ independently, which may lead to the state sequences that should not appear. e.g. we set $S_1, S_3 = 0$. but if we use τ to decode, we will find solutions for $S_1' \neq 0$ and $S_3' \neq 0$. Hence this has serious problems

2.1 Viterbi Decoding

To eliminate the previous problem, we use jointly decode the entire optimal state sequence:

$$q_{1:t} = \underset{Q_{1:T}}{\operatorname{argmax}} P(Q_{1:T} | O_{1:T}, \lambda) = \underset{Q_{1:T}}{\operatorname{argmax}} P(Q_{1:T}, O_{1:T} | \lambda).$$

Dynamic Programming Solution:

The recursive relationship:

Given $O_{1:t}$, it is useful to know the optimal paths for N subproblems:

$$\max_{Q_{1:t-1}} P(Q_{1:t-1}, O_{1:t}, Q_t = s_i | \lambda), \quad i \in \{1, \dots, N\}.$$

We can solve at least the following two problems by using the answer to these N problems.

1. The optimal hidden state sequence till time t .

if the i^* subproblem has the largest probability (likelihood).

then the optimal parameter of this subproblem $q_{1:t-1}$ plus
 $q_t = i^*$ is the optimal hidden state sequence for $O_{1:t}$.

2. The optimal $q_{1:t+1}$ for $O_{1:t+1}$ can be divided into:

$q_{1:t-1}, q_t, q_{t+1}$. If we know $q_t = i^*$, the $q_{1:t-1}$ part can be found by the i^* -th subproblem. By Markov Assumption, we only need to consider N possible transitions $S_{q_t} \rightarrow S_{q_{t+1}}$ to decide which state is optimal for q_{t+1} .

Define: $\delta_t(i) = \max_{Q_{1:t-1}} P(Q_{1:t-1}, O_{1:t}, Q_t = s_i | \lambda)$.

Recursive: $\delta_{t+1}(i) = \max_{j \in \{1, \dots, N\}} (\delta_t(j) A_{ji} B_{it+1})$. --- \otimes

where $\delta_t(j)$ is the probability of the j -th subproblem.

A_{ji} tranmits from S_j to S_i at time t and $t+1$ respectively.

$B_{i,t+1}$ is the probability of observing O_{t+1} when state is S_i .

Then $\delta_t(j) A_{ji} B_{i,t+1}$ is the probability of the optimal state sequence when we know $Q_t = S_j$ for an observation sequence $O_{1:t+1}$.

Forward Recursion: Start with $t=1$

After S values are calculated, we find q_T by:

$$q_T = \underset{i \in \{1, \dots, N\}}{\operatorname{argmax}} \delta_T(i)$$

From eqn \otimes , if we know the optimal state at time $t+1$ is q_{t+1} ,

we just need to find which j leads to the largest $\delta_t(j) A_{ji} B_{i,t+1}$.

Then S_j is the optimal state for Q_t . That is, we need to record the optimal transitions from t to $t+1$.

We use $\psi_{t+1}(i)$ to denote the optimal state at time t if the optimal state is i at time $t+1$, here.

$$\psi_{t+1}(i) = \underset{j \in \{1, \dots, N\}}{\operatorname{argmax}} \left(\delta_t(j) A_{ji} B_{i,t+1} \right) = \underset{j \in \{1, \dots, N\}}{\operatorname{argmax}} (\delta_t(j) A_{ji})$$

Dyno-Code: (Viterbi decoding)

1. Initialise: $\delta_1(i) = \pi_i B_{i1}$, $\psi_1(i) = 0$, $\forall i \in \{1, \dots, N\}$.

2. Recur: For $t=2, 3, \dots, T-1$, and $\forall i \in \{1, \dots, N\}$

$$\delta_{t+1}(i) = \max_{j \in \{1, \dots, N\}} (\delta_t(j) A_{ji} B_{i,t+1}).$$

$$\psi_{t+1}(i) = \max_{j \in \{1, \dots, N\}} (\delta_t(j) A_{ji}).$$

3. Output: the optimal state q_T : $q_T = \underset{i \in \{1, \dots, N\}}{\operatorname{argmax}} \delta_T(i)$.

and the rest of optimal path is determined by: for $t=T-1, T-2, \dots, 2, 1$,

$$q_t = \psi_{t+1}(q_{t+1}).$$

Complexity $\mathcal{O}(TN^2)$.

3. Learning Parameters in HMM.

Given $N, M, O_{1:T}$, learn $\lambda = (\pi, A, B)$

Classical Algo: Baum-Welch. (Maximum likelihood technique). (EM)

$$\lambda^* = \arg\max_{\lambda} P(O_{1:T} | \lambda).$$

Iterative algorithm:

1. Random initialise $\lambda^{(0)}$, compute $\ell^{(0)} = \log P(O_{1:T} | \lambda)$.

2. Update parameters to get $\lambda^{(1)}$.

⋮

Define new variable $\xi_t(i, j) = P(Q_t = S_i, Q_{t+1} = S_j | O_{1:T}, \lambda)$.

$\xi_t(i, j)$ is the expected proportion of transition from S_i (time t) to S_j (time $t+1$)

Then it is natural to use $\xi_t(i, j)$ to update the value for A_{ij} in $\lambda^{(t+1)}$

\uparrow
computed based on $\lambda^{(t)}$

$$\xi_t(i, j) \cdot P(O_{1:T} | \lambda) = P(Q_t = S_i, Q_{t+1} = S_j, O_{1:T} | \lambda)$$

$$\text{Then we have } \xi_t(i, j) = \frac{\alpha_t(i) A_{ij} B_{j,t+1} \beta(j)}{P(O_{1:T} | \lambda)}$$

$$\text{Since } \xi_t(i, j) \text{ is a probability, we have } \sum_{i=1}^N \sum_{j=1}^M \xi_t(i, j) = 1$$

$$\text{Hence } \sum_{i=1}^N \sum_{j=1}^M \frac{\alpha_t(i) A_{ij} B_{j,t+1} \beta(j)}{P(O_{1:T} | \lambda)} = 1$$

$$\Rightarrow P(O_{1:T} | \lambda) = \sum_{i=1}^N \sum_{j=1}^M \alpha_t(i) A_{ij} B_{j,t+1} \beta(j). \quad \forall t \in \{1, 2, \dots, T-1\}.$$

$$\text{We also get that: } \gamma_t(i) = \sum_{j=1}^M \xi_t(i, j).$$

$$\text{Recursive Relationship: } \xi_t(i, j) = \frac{\alpha_t(i) A_{ij} B_{j,t+1} \beta_{t+1}(j)}{\sum_{i=1}^N \sum_{j=1}^M \alpha_t(i) A_{ij} B_{j,t+1} \beta_{t+1}(j)}. \quad \begin{matrix} \text{--- Responsibility} \\ \text{in EM.} \end{matrix}$$

The parameters $\lambda = (\Pi, \Delta, \Xi)$ can be updated using τ and ξ .

* Since $\tau_t(i)$ is the expected proportion of $Q_t = S_i$, we update τ_{ti} using $\tau_t(i)$.

$$* A_{ij} = \frac{\sum_{t=1}^{T-1} \xi_t(i, j)}{\sum_{t=1}^{T-1} \tau_t(i)}$$

Interpretation:

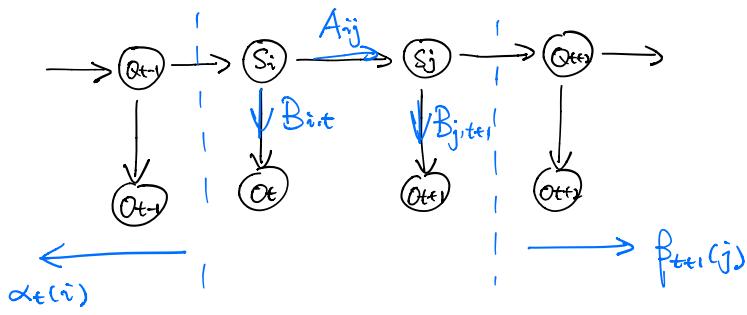
- $\xi_t(i, j)$ is the expected probability of transition from S_i to S_j at t .

Hence the expected number of transitions from S_i to S_j in the training sequence is $\sum_{t=1}^{T-1} \xi_t(i, j)$.

- $\sum_{t=1}^{T-1} \tau_t(i)$ is the expected number of times that any hidden states being S_i

Hence A_{ij} , the transition probability from S_i to S_j , is the proportion of transitions $S_i \rightarrow S_j$ in all transitions starting from S_i .

$$* B_{jk} = \frac{\sum_{t=1}^T \tau_t(j)}{\sum_{t=1}^T [\alpha_t = k] \tau_t(j)}, \text{ where } [\cdot] \text{ is the indicator function.}$$



Summary :

	Definition	Recurrsion
α	$\alpha_t(i) = P(O_{1:T}, Q_t = s_i \lambda)$	$\alpha_{t+1}(i) = \left(\sum_{j=1}^n \alpha_t(j) A_{ji} \right) B_{ji, t+1}$
β	$\beta_t(i) = P(O_{t+1:T} Q_t = s_i, \lambda)$	$\beta_t(i) = \sum_{j=1}^n A_{ij} \beta_{j, t+1} f_{t+1}(j)$
γ	$\gamma_t(i) = P(Q_t = s_i O_{1:T}, \lambda)$	$\gamma_t(i) = \frac{\alpha_t(i) \beta_t(i)}{\sum_{j=1}^n \alpha_t(j) \beta_t(j)}$
δ	$\delta_t(i) = \max_{Q_{1:t-1}} P(Q_{1:t-1}, O_{1:t}, Q_t = s_i \lambda)$	$\delta_{t+1}(i) = \max_{j \in \{1, \dots, n\}} (\delta_t(j) A_{ji} B_{ji, t+1})$
ξ	$\xi_t(i, j) = P(Q_t = s_i, Q_{t+1} = s_j O_{1:T}, \lambda)$	$\xi_t(i, j) = \frac{\alpha_t(i) A_{ij} B_{ji, t+1} f_{t+1}(j)}{\sum_{i=1}^n \sum_{j=1}^n \alpha_t(i) A_{ij} B_{ji, t+1} f_{t+1}(j)}$