

HMM, MEMM, CRF

PRELIMINARIES

• MARKOV CHAINS

Probabilistic models that represent a sequence of random variables $\{q_t\}_{t=1, \dots, T}$ satisfying Markov property.

→ row-stochastic transition matrix $\underline{A} \in [\varnothing; 1]^{N \times N}$

→ set of states $S = \{s_1, \dots, s_N\}$

→ initial probability distribution $\underline{\pi} = \begin{bmatrix} \pi_1 \\ \vdots \\ \pi_N \end{bmatrix}$

• MIXTURE MODELS

Distribution models that arise from convex combinations of simple distributions.

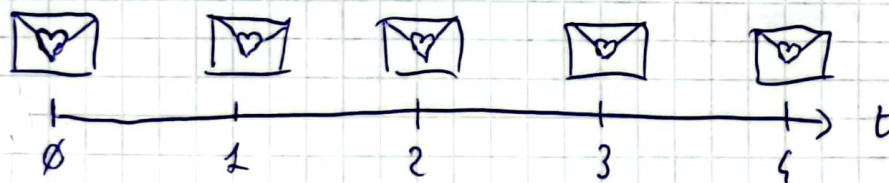
→ probability distributions $b_i(x)$

→ observations x

TOTAL PROBABILITY THEOREM $\Rightarrow IP(x) = \sum_i IP(x | s_i) IP(s_i) = \sum_i b_i(x) p_i$

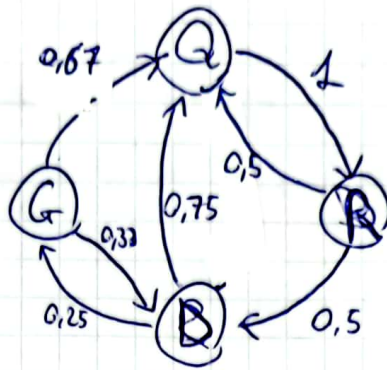
HIDDEN MARKOV MODELS

Imagine a sequence of love letters (or messages):



Each message can contain a question, a general consideration, or a love declaration.

$$S = \{Q, G, D, R\}$$



Question \rightarrow surely Reply

Reply \rightarrow maybe another Question, maybe a random Declaration

Declaration \rightarrow likely a "do you really love me" Question, or a General statement

General \rightarrow Question more likely than random Declaration

Problem: we can't observe the state itself, but we can measure the characteristics of the messages.

$$\underline{X}_t = \begin{bmatrix} \text{length} \\ \text{count("love")} \\ \text{count("you") - count("I")} \\ \text{count("?")} \\ \text{TFIDF(some words)} \end{bmatrix}$$

FINAL GOAL:

WE WANT TO ESTIMATE THE HIDDEN STATE (Q, R, D or G) OF EACH MESSAGE GIVEN THE FEATURES \underline{X}_t OF THE MESSAGE.

Lawrence Rabiner (1989) systematized tasks on such problems:

- ① LIKELIHOOD: $IP(\{\underline{X}_t\}_t | \lambda) \rightarrow$ parameters $\underline{A}, \underline{\pi}, b$
- ③ DECODING: $\{q_t^*\}_t = \operatorname{argmax} IP(\{q_t\}_t | \{\underline{X}_t\}_t, \lambda)$
- ② LEARNING: $\lambda^* = \operatorname{argmax} IP(\lambda | \{\underline{X}_t\}_t)$

① LIKELIHOOD (model \rightarrow prob. of observed sequence)
 $A, b(\cdot), \pi, \gamma(x_i), S$ $IP(\{x_i\} | \lambda)$

Recall:

- X_t depends only on q_t
- q_t depends only on q_{t-1} (Markov property)

If we had $\{q_t\}_{t=1, \dots, T}$, then:

$$IP(\{X_t\} | \{q_t\}, \lambda) = \prod_{t=1}^T IP(X_t | q_t, \lambda)$$

but we don't know them. So:

$$\begin{aligned} IP(\{X_t\}, \{q_t\} | \lambda) &= IP(\{X_t\} | \{q_t\}, \lambda) \cdot IP(\{q_t\} | \lambda) = \\ &= \prod_{t=1}^T IP(X_t | q_t, \lambda) \prod_{t=1}^T IP(q_t | q_{t-1}, \lambda) \end{aligned}$$

for each possible hidden state sequence $\{q_t\}$. $\# = N^T$
 \Downarrow
INTRACTABLE

DYNAMIC PROGRAMMING FRAMEWORK

\rightarrow Break down a complex problem in ^{overlapping} subproblems that can be solved recursively.

FORWARD ALGORITHM

$$\begin{aligned} \text{Let } \alpha_t(j) &= IP(q_t = s_j, x_1, \dots, x_t | \lambda) = \\ &= \sum_i \alpha_{t-1}(i) A_{ij} b_j(x_t) = \\ &= \sum_i IP(q_{t-1} = s_i, x_1, \dots, x_{t-1} | \lambda) IP(q_t = s_j | q_{t-1} = s_i) IP(x_t | q_t) \end{aligned}$$

$$\begin{aligned}
 & \text{since } IP(q_t = s_j, x_t, x_{t-1}, \dots, x_1) = \\
 &= \sum_i IP(q_t = s_j, q_{t-1} = s_i, x_t, x_{t-1}, \dots, x_1) = \\
 &= \sum_i IP(q_t = s_j, x_t \mid q_{t-1} = s_i, x_{t-1}, \dots, x_1) \cdot \underbrace{IP(q_{t-1} = s_i, x_{t-1}, \dots, x_1)}_{\substack{\text{Markov property} \\ \text{doesn't depend on } q_t}} \\
 &= \sum_i IP(q_t = s_j, x_t \mid q_{t-1} = s_i) \alpha_{t-1}(i)
 \end{aligned}$$

At the first step $\alpha_1(j) = \pi_j b_j(x_1)$.

In the end:

$$IP(\{x_t\}_{t=1 \dots T} \mid \lambda) = \sum_{j=1}^N \alpha_T(j)$$

② DECODING (model, $\underbrace{\text{observed sequence}}_{\substack{\text{state set} \\ S}} \xrightarrow{\quad} \text{hidden state sequence}$
 $\underline{A}, \underline{\pi}, \underline{b}(\cdot)$ $\{x_t\}$ $\{q_t\}$)

Naive solution: for each of the N^T possible state sequences $\{q_t\}_{t=1 \dots T}$ run the forward algorithm and take the sequence with highest likelihood.

→ INTRACTABLE

VITERBI ALGORITHM

$$\begin{aligned}
 \text{Let } \delta_t(j) &= \max_{q_1, \dots, q_{t-1}} IP(q_1, \dots, q_{t-1}, q_t = s_j, x_1, \dots, x_t \mid \lambda) = \\
 &= \max_{q_1, \dots, q_{t-1}} \underbrace{IP(x_t, q_t = s_j \mid q_{t-1} = q_t)}_{\text{doesn't depend on } q_1, \dots, q_{t-2}} \cdot \underbrace{IP(q_1, \dots, q_{t-1}, x_1, \dots, x_{t-1})}_{\text{for Markov property and since } \underline{A} \text{ and } \underline{b}(\cdot) \text{ are given}}
 \end{aligned}$$

$$\Rightarrow \delta_t(j) = \max_i \underbrace{(A_{ij} \delta_{t-1}(i))}_{P(q_t=s | q_{t-1}=s, i)} \underbrace{b_j(x_t)}_{P(x_t | q_t, s_j)}$$

At the first step: $\delta_1(j) = \pi_j b_j(x_1)$

By iteratively computing $\delta_t(j)$ and the state s_j^* maximizing it, it is possible to reconstruct the optimal state sequence (not in the forward run, but in BACKTRACKING).

$$\begin{cases} q_T^* = \operatorname{argmax}_j \delta_T(j) \\ q_t^* = \psi_{t+1}(q_{t+1}^*) \end{cases} \text{ for } t < T$$

where $\psi_{t+1}(j) = \operatorname{argmax}_i (\delta_t(i) A_{ij})$

↳ notice it's outgoing, while in the definition of δ it was incoming

② LEARNING (observed $\{x_t\}$ sequence, $\xrightarrow[S]{\text{state set}}$ model parameters $A, \pi, b(\cdot)$)

Unsupervised learning (even supervised learning is possible).

→ learning A , $b(\cdot)$ and π

Main solution: Knowing the state sequence, we would be able to estimate parameters immediately.

→ but we don't know $\{q_t\}$

Let's define:

BACKWARD PROBABILITY

$$\beta_t(j) = P(X_T, \dots, X_{t+2} \mid q_t = j, \lambda) \quad \text{with } \beta_T(j) = 1 \quad \forall j$$

Note: forward probability $\alpha_t(j)$ is the probability of a certain state and observations up to it; backward prob. is the probability of observations after a certain state given it.

Recursively:

$$\begin{aligned} \beta_t(j) &= P(X_T, \dots, X_{t+2} \mid q_t = j, \lambda) = \\ &= \sum_{i=1}^N P(X_T, \dots, X_{t+2} \mid q_{t+2} = i, \lambda) P(q_{t+2} = i \mid q_t = j, \lambda) \\ &= \sum_{i=1}^N P(X_T, \dots, X_{t+2} \mid q_{t+2} = i, \lambda) P(X_{t+2} \mid q_{t+2} = i, \lambda) A_{ji} \\ &= \sum_{i=1}^N \beta_{t+2}(i) b_i(X_{t+2}) A_{ji} \end{aligned}$$

Annotation: In the second line, $q_t = j, \lambda$ is crossed out with a line and an arrow pointing to it with the text "doesn't count anymore".

Maximum likelihood estimation in an iterative way:

EXPECTATION-MAXIMIZATION ALGORITHM

E-STEP (EXPECTATION):

we current estimates \hat{A} , \hat{b} , $\hat{\pi}$ to compute
 $IP(q_t = s_i, q_{t+1} = s_j | \{x_t\}) =: \xi_t(i, j)$

$$\xi_t(i, j) = \frac{IP(q_t = s_i, q_{t+1} = s_j, x_1, \dots, x_T)}{IP(\{x_t\}_t)}$$

where $IP(q_t = s_i, q_{t+1} = s_j, x_1, \dots, x_T) =$

given q_t and q_{t+1} , we can split the sequence in independent segments

$$\begin{aligned}
 &= IP(x_1, \dots, x_T | q_{t+1} = s_j, q_t = s_i) IP(q_{t+1} = s_j, q_t = s_i) \\
 &= IP(x_1, \dots, x_t | q_t = s_i) IP(x_{t+1} | q_{t+1} = s_j, q_t = s_i) \\
 &\quad \cdot IP(x_{t+2}, \dots, x_T | q_{t+1} = s_j, q_t = s_i) IP(q_{t+1} = s_j, q_t = s_i) \\
 &\quad \quad \quad \swarrow \text{irrelevant} \\
 &= P(x_1, \dots, x_t | q_t = s_i) b_j(x_{t+1}) \beta_{t+1}(j) \\
 &\quad \cdot IP(q_{t+1} = s_j | q_t = s_i) IP(q_t = s_i) \\
 &= \alpha_t(i) b_j(x_{t+1}) \beta_{t+1}(j) A_{ij}
 \end{aligned}$$

↳ with $\xi_t(i, j)$ we can estimate \hat{A}_{ij}

Then, we define:

$$\begin{aligned}
 \gamma_t(j) &= IP(q_t = s_j | \{x_t\}) = \\
 &= \frac{IP(q_t = s_j, x_1, \dots, x_T)}{IP(x_1, \dots, x_T)} = \frac{\alpha_t(j) \beta_t(j)}{IP(x_1, \dots, x_T)}
 \end{aligned}$$

↳ with $\gamma_t(j)$ we can estimate \hat{b}_j

M-STEP

Update parameters:

$$\hat{A}_{ij} = \frac{\mathbb{E} \left[\sum_{t=1}^{T-1} \mathbb{1}_{\{q_{t+1}=s_j, q_t=s_i\}} | \{x_t\} \right]}{\mathbb{E} \left[\sum_{t=1}^{T-1} \mathbb{1}_{\{q_t=s_i\}} | \{x_t\} \right]} = \frac{\sum_{t=1}^{T-1} \mathbb{P}(q_{t+1}=s_j, q_t=s_i | \{x_t\})}{\sum_{t=1}^{T-1} \sum_{j=1}^N \mathbb{P}(q_{t+1}=s_j, q_t=s_i | \{x_t\})}$$

probability of transitioning from s_i to s_j

$$= \frac{\sum_{t=1}^{T-1} \xi_t(i,j)}{\sum_{t=1}^{T-1} \sum_{j=1}^N \xi_t(i,j)}$$

$$\hat{b}_j(x) = \frac{\mathbb{E} \left[\sum_{t=1}^T \mathbb{1}_{\{x_t=x; q_t=s_j\}} | \{x_t\} \right]}{\mathbb{E} \left[\sum_{t=1}^T \mathbb{1}_{\{q_t=s_j\}} | \{x_t\} \right]} = \frac{\sum_{t=1}^T \mathbb{P}(x_t=x; q_t=s_j | \{x_t\})}{\sum_{t=1}^T \mathbb{P}(q_t=s_j | \{x_t\})}$$

some symbol

$$= \frac{\sum_{t=1}^T \gamma_t(j) \mathbb{P}(x_t=x | \{x_t\})}{\sum_{t=1}^T \gamma_t(j)}$$

either 0 or 1

Thorough discussion about EM is complex, but in general it is about maximizing expected likelihood (expected since we don't know hidden states) iteratively:

