

12.1 HIDDEN MARKOV MODELS

Si parte dall'esigenza di voler descrivere il comportamento di sistemi a stati finiti in cui le probabilità di transizione tra stati non dipendono dal particolare valore assunto dal tempo, ma soltanto dalla coppia di stati interessati nella transizione.

In una catena di Markov ogni stato può essere messo in corrispondenza con un evento fisico osservabile, per cui una situazione reale (ovvero una successione di eventi) può essere modellizzata attraverso una sequenza ordinata di stati assunti dalla catena in determinati istanti temporali.

Si può estendere il modello di Markov in modo da contemplare la situazione in cui l'osservazione sia legata allo stato attraverso una funzione probabilistica: in altre parole, il modello risultante è un processo stocastico che emette simboli osservabili appoggiandosi su un altro processo stocastico non osservabile (nascosto) che determina il cambio di stato. Le probabilità di emissione dei simboli cambiano da uno stato all'altro.

Un modello (catena) di Markov nascosto può essere descritto attraverso cinque elementi:

- l'insieme dei possibili stati che può assumere la catena di Markov

$$x \in V = \{v_1, \dots, v_N\}$$

- il vettore delle probabilità che la catena si trovi in un particolare stato all'istante iniziale

$$\pi_0 \in \mathbb{R}^{N \times 1}, \quad \pi_0(i) \geq 0 \quad \forall i \quad \text{e} \quad \sum_{i=1}^N \pi_0(i) = 1$$

$$\pi_0(i) = p(x_0 = i) = p(x_0 = v_i)$$

- la matrice di transizione

$$A \in \mathbb{R}^{N \times N}$$

il cui elemento di posizione $[A]_{ij} = a_{ij}$ indica la probabilità che la catena si trovi nello stato v_j condizionatamente al fatto che all'istante precedente si trovasse nello stato v_i , ovvero

$$p(x_{t+1} = j \mid x_t = i) = a_{ij} \geq 0;$$

Per tale matrice vale la relazione $\sum_{j=1}^N a_{ij} = 1$, ovvero la catena di Markov dallo stato i transita (all'istante successivo) con probabilità 1 in uno degli N stati (può rimanere anche in i). Ponendo $\mathbb{1}^T = [1 \ \cdots \ 1]$ la condizione precedente si può riscrivere come $A \mathbb{1} = \mathbb{1}$.

Le matrici con elementi non negativi le cui righe sommano a 1 si definiscono matrici stocastiche per righe.

- l'insieme delle uscite osservabili

$$y \in W = \{w_1, \dots, w_M\}$$

- la matrice $C \in \mathbb{R}^{N \times M}$, cosiddetta matrice delle emissioni, i cui elementi $[C]_{ij} = c_{ij}$ rappresentano la probabilità che l'uscita osservata sia w_j condizionatamente al fatto che la catena si trova nello stato v_i , ovvero

$$p(y_t = j \mid x_t = i) = c_{ij} \geq 0$$

Anche per la matrice C vale la proprietà $\sum_{j=1}^M c_{ij} = 1$ che può essere espressa come $C \mathbb{1}_M = \mathbb{1}_N$.

Nella descrizione dell'evoluzione nel tempo di una catena di Markov nascosta si incontrano i concetti di cammino e osservazione: il primo è un insieme ordinato di stati in cui la catena si viene a trovare sequenzialmente mentre l'osservazione è un insieme ordinato di simboli (uscite) che vengono osservate in sequenza.

Indicando con x_t lo stato in cui si trova la catena all'istante t , mentre con $x^t = (x_0, \dots, x_t)$ la sequenza di stati in cui la catena si viene a trovare in $t + 1$ istanti successivi e con $Y^h = (y_0, \dots, y_h)$ la successione di uscite osservate in $h + 1$ istanti successivi; in queste ipotesi si possono formulare diversi problemi sull'evoluzione della catena di Markov:

- definendo $p(x_t \mid Y^h) = \hat{p}_{t|h} \in \mathbb{R}^N$ come la probabilità che la catena si trovi all'istante t nello stato x_t sulla base delle osservazioni Y^h , si può cercare lo stato che massimizza tale probabilità (calcolo della probabilità 'forward'), ovvero

$$\hat{x}_{t|h} = \arg \max_i \hat{p}_{t|h}(i).$$

- si può cercare la sequenza di stati che massimizzano ad ogni istante la probabilità che la catena si trovi in essi sulla base delle osservazioni (calcolo della probabilità 'backward'), ovvero

$$\widehat{x}_{t|h}^t = (\widehat{x}_{0|h}, \dots, \widehat{x}_{t|h}) \in \underbrace{V \times \dots \times V}_{t+1 \text{ volte}}$$

- si può cercare la traiettoria dello stato della catena con più alta probabilità sulla base delle osservazioni Y^h (problema della decodifica), ovvero

$$\widehat{x}_{t|h}^t = \arg \max_{x_0, \dots, x_t} p(x^t | Y^h) \in \underbrace{V \times \dots \times V}_{t+1 \text{ volte}}$$

Da notare che in generale $\widehat{x}_{t|h}^t \neq \widehat{x}_{t|h}$.

Sulla base dell'analogia con il caso continuo per cui vale

$$p(x_t | Y^t) \propto p(y_t | x_t) \int_{x_{t-1}} p(x_t | x_{t-1}) p(x_{t-1} | Y^{t-1}) dx_{t-1}$$

si può scrivere la seguente relazione che definisce esplicitamente la probabilità che la catena si trovi nello stato $x_t = i$ sulla base delle osservazioni Y^t :

$$\underbrace{p(x_t = i | Y^t)}_{\widehat{p}_{i|t}(i)} \propto (y_t = \bar{y}_t | x_t = i) \left(\sum_{j=1}^N p(x_t = i | x_{t-1} = j) \underbrace{p(x_{t-1} = j | Y^{t-1})}_{\widehat{p}_{t-1|t-1}(j)} \right)$$

da cui

$$\widehat{p}_{i|t}(i) \propto \underbrace{c_i(\bar{y}_t)}_{\text{aggiornamento}} \underbrace{\sum_{j=1}^N a_{ij} \widehat{p}_{t-1|t-1}(j)}_{\text{predizione}}$$

dove

$$c_i(\bar{y}_t) = p(y_t = \bar{y}_t | x_t = i)$$

e

$$C(\bar{y}_t) = [c_1(\bar{y}_t) \quad \dots \quad c_N(\bar{y}_t)]^T$$

è la colonna della matrice C relativa all'uscita \bar{y}_t , $\bar{y}_t \in W$. Scrivendo in forma compatta otteniamo:

- la predizione

$$\hat{p}_{t|t-1}^T = \hat{p}_{t-1|t-1}^T A$$

- la predizione a più passi

$$\hat{p}_{t|h}^T = \hat{p}_{h|h}^T A^{t-h}$$

- e il filtraggio

$$\hat{p}_{t|t} = \frac{\hat{p}_{t|t-1} * c(\bar{y}_t)}{\|\hat{p}_{t|t-1} * c(\bar{y}_t)\|_1} \equiv p_{t|t}^T = \frac{\text{diag}[c(\bar{y}_t)]p_{t|t-1}}{\|\text{diag}[c(\bar{y}_t)]p_{t|t-1}\|_1}$$

dove $*$ indica il prodotto elemento per elemento di due vettori colonna (simbologia Matlab).

Consideriamo ora il caso più corrispondente allo smoothing, nel quale vengono considerate tutte le misure (anche future) per la predizione dello stato, cioè $p(x_t|Y^T)$. In particolare l'applicazione delle usali regole sull'indipendenza condizionata e di Bayes si ottiene:

$$\begin{aligned} p(x_t = i | Y^T) &= \frac{p(Y^t, x_t = i)p(\tilde{Y}^{t+1} | x_t = i)}{p(Y^T)} \\ &= \frac{\alpha_{t|t}(i)\beta_{t|t+1}(i)}{p(Y^T)} \\ &= \frac{\alpha_{t|t}(i)\beta_{t|t+1}(i)}{\sum_{i=1}^N \alpha_{t|t}(i)\beta_{t|t+1}(i)} \end{aligned}$$

Il termine relativo ad $\alpha_{t|t}$ si ottiene facilmente in maniera simile a $\hat{p}_{t|t}$ visto in precedenza e che si può riassumere nelle seguenti equazioni:

- $\alpha_{t+1|t}^T = \alpha_{t|t}^T A$
- $\alpha_{t+1|t+1} = \alpha_{t+1|t} * c(\bar{y}_{t+1})$
- $\alpha_{0|-1} = \pi_0$

Da un punto di vista teorico tali espressioni sono soddisfacenti, però da un punto di vista computazionale tale calcolo potrebbe risultare troppo oneroso. Non è necessario normalizzare ad ogni passo, meglio se lo si fa ogni tanto: in letteratura esistono tecniche che operano proprio in questo senso. Si veda per esempio L. Rabiner, A Tutorial on Hidden Markov Models and Selected Applications in Speech Recognition, Proceedings of IEEE, vol 77(2), pp. 257-286

Analizzando invece il secondo termine relativo a $\beta_{t|t+1}$ otteniamo:

$$\begin{aligned}
\beta_{t|t+1}(i) &= p(\tilde{Y}^{t+1} | x_t = i) = \sum_{j=1}^N p(\tilde{Y}^{t+1}, x_{t+1} = j | x_t = i) \\
&= \sum_{j=1}^N p(\tilde{Y}^{t+1} | x_{t+1} = j, x_t = i) p(x_{t+1} = j | x_t = i) \\
&= \sum_{j=1}^N p(\tilde{Y}^{t+1} | x_{t+1} = j) p(x_{t+1} = j | x_t = i) \\
&= \beta_{t+1|t+1}(j) a_{ij}
\end{aligned}$$

nella quale (secondo passaggio), prima utilizziamo la formula di Bayes, e poi (penultimo passaggio) sfruttiamo il fatto che le misure future condizionando rispetto a x_{t+1} e x_t sono indipendenti da x_t . Alla fine si ottiene una espressione moltiplicativa tra due termini, dove il primo per definizione è proprio $\beta_{t+1|t+1}(j)$, mentre il secondo a_{ij} . Questo significa che l'elemento j -esimo del vettore $\beta_{t|t+1}$ si ottiene prendendo la riga i di a moltiplicato per il vettore $\beta_{t+1|t+1}$.

Scrivendo in forma compatta otteniamo

$$\beta_{t|t+1} = A\beta_{t+1|t+1}$$

In pratica per mantenera la stessa analogia del vettore $\alpha_{t|t}$ è opportuno prendere il trasposto di tutta l'espressione:

$$\beta_{t|t+1}^T = \beta_{t+1|t+1}^T A^T$$

in cui si nota che nelle catene di Markov l'inversione della dinamica comporta a prendere la trasposta della matrice A .

Invece, per quel che riguarda l'aggiornamento

$$\beta_{t+1|t+1}(i) = p(\tilde{Y}^{t+1} | x_{t+1} = i) = \underbrace{p(y_{t+1} | x_{t+1} = i)}_{c_i(\bar{y}_{t+1})} \underbrace{p(\tilde{Y}^{t+2} | x_{t+1} = i)}_{\beta_{t+1|t+2}(i)}$$

dove si nota che la misura t -esima e le successive sono indipendenti tra di loro una volta che ho condizionato rispetto a x_{t+1} . Trascrivendo sempre in forma compatta otteniamo l'espressione

$$\beta_{t|t} = c(\bar{y}_t) .* \beta_{t|t+1}$$

mentre l'inizializzazione (che in questo caso invece è la condizione finale) è data da

$$\beta_{T|T} = \underbrace{p(y_T | x_T = i)}_{c_i(\bar{y}_T)}$$

Quindi, riassumendo l'algoritmo di forward-backward calcola lo stato con maggiore probabilità in un certo istante date tutte le misure fino all'istante T , cioè la sequenza $\{\hat{x}_{0|T}, \dots, \hat{x}_{T|T}\}$ per ogni $t = 0, \dots, T$. Tale sequenza è generalmente diversa dalla seguente, data da:

$$\{\hat{x}_{0|T}, \dots, \hat{x}_{T|T}\} = \arg \max_{x_0, \dots, x_T} p(x_0, x_1, \dots, x_T | Y^T)$$

Quest'ultima infatti è la probabilità di una particolare sequenza: tra tutte le possibili prendo quella con probabilità maggiore. Si verifica velocemente che sono vere le due seguenti relazioni:

$$p(\hat{x}_{0|T}, \dots, \hat{x}_{T|T} | Y^T) \geq p(\hat{x}_{0|T}, \dots, \hat{x}_{T|T} | Y^T)$$

$$p(\hat{x}_{t|T} | Y^T) \leq p(\hat{x}_{t|T} | Y^T)$$

Per calcolare la sequenza $\{\hat{x}_{0|T}, \dots, \hat{x}_{T|T}\}$ si utilizza un algoritmo particolare della programmazione dinamica: l'algoritmo di Viterbi, che verrà analizzato nella lezione seguente.