

1. MODELOS OCULTOS DE MARKOV

Un Modelo Oculto de Markov (Hidden Markov Model HMM) es un proceso estocástico que consta de un proceso de Markov no observado (oculto) $\mathbf{q} = \{q_t\}_{t \in N}$ y un proceso observado $\mathbf{O} = \{o_t\}_{t \in N}$ cuyos estados son dependientes estocásticamente de los estados ocultos, es decir, es un proceso bivariado (\mathbf{q}, \mathbf{O}) . Los HMMs se pueden considerar también como sistemas generativos estocásticos, los cuales se emplean en la modelación de series de tiempo.

1.1 CADENAS DE MARKOV

Una cadena de Markov $\mathbf{q} = \{q_t\}_{t \in N}$ es un proceso estocástico de Markov discreto. Un proceso estocástico se llama de Markov si conocido el presente, el futuro no depende del pasado, esto quiere decir, que dada una variable estocástica q_{t-1} que denota el estado del proceso en el tiempo $t-1$, entonces la probabilidad de transición en el momento t se define como $P[q_t = \sigma_t \mid q_{t-1} = \sigma_{t-1}]$. Formalmente, una cadena de Markov se define como (Q, A) , donde $Q = \{1, 2, \dots, N\}$ son los posibles estados de la cadena y $A = (a_{ij})_{n \times n}$ es una matriz de transición de estados en el modelo. Si $A(t) = a_{ij}(t)_{n \times n}$ es independiente del tiempo entonces el proceso se llama homogéneo y las probabilidades de transición de estados son de la forma $a_{ij}(t) = P[q_t = j \mid q_{t-1} = i]$ con las siguientes propiedades:

- i) $0 \leq a_{ij} \leq 1, \quad 1 \leq i, j \leq N;$
- ii) $\sum_{j=1}^N a_{ij} = 1, \quad 1 \leq i \leq N.$

La condición fundamental de que sea una cadena de Markov establece que las probabilidades de transición y emisión dependen solamente del estado actual y no del pasado, esto es, $P[q_t = j \mid q_{t-1} = i, q_{t-2} = k, \dots] = P[q_t = j \mid q_{t-1} = i] = a_{ij}(t)$.

En este trabajo se considera el conjunto de estados finitos.

1.2 DEFINICION DE MODELOS OCULTOS DE MARKOV

Un modelo oculto de Markov es una cadena de q junto con un proceso estocástico que toma valores en un alfabeto Σ y el cual depende de q .

Estos sistemas evolucionan en el tiempo pasando aleatoriamente de estado a estado y emitiendo en cada momento al azar algún símbolo del alfabeto Σ . Cuando se encuentra en el estado $q_{t-1} = i$, tiene la probabilidad a_{ij} de moverse al estado $q_t = j$ en el siguiente instante y la probabilidad $b_j(k)$ de emitir el símbolo $o_t = v_k$ en el tiempo t .

Sólamente los símbolos emitidos por el proceso q son observables, pero no la ruta o secuencia de estados q , de ahí el calificativo de "oculto" de Markov, ya que el proceso de Markov q es no observado.

El siguiente ejemplo ilustra un proceso q independiente del tiempo. Supóngase que en un salón se encuentra un número N muy grande de urnas de vidrio. Dentro de cada urna se tiene una cantidad M de bolas de colores. Un mago está en el salón y de acuerdo con algún procedimiento aleatorio elige una urna inicial. De ésta saca al azar una bola y registra su color como una observación. La bola es retornada a la urna de la cual fué seleccionada. A continuación selecciona una nueva urna de acuerdo con un procedimiento aleatorio que depende de la urna actual y la elección de alguna bola es repetida. Este proceso completo se realiza en un tiempo T y genera una secuencia de observación finita de colores O de longitud T , la cual puede modelarse como la salida observable de un HMM. Se asume que las urnas son seleccionadas independientemente.

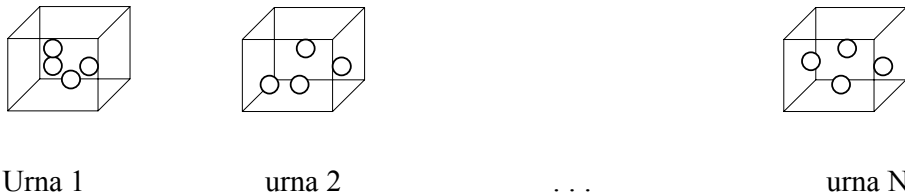


Figura. 1.1 Modelo de urnas y bolas de N estados que ilustra el caso general de un HMM con símbolos discretos.

Los siguientes son ejemplos de posibles secuencias de observación del modelo de las urnas y las bolas:

\mathcal{O}^1 = (amarillo, verde, azul, verde, rojo, amarillo, naranja, rojo, verde, azul, amarillo),

\mathcal{O}^2 = (amarillo, rojo, verde, rojo, azul, naranja, verde, rojo, azul, amarillo, rojo, verde),

\mathcal{O}^3 = (rojo, azul, amarillo, rojo, azul, verde, rojo, amarillo, naranja, naranja, verde, rojo),

\mathcal{O}^4 = (rojo, verde, naranja, rojo, rojo, azul, verde, amarillo, azul, rojo, verde, rojo).

El alfabeto es:

Σ = verde, azul, rojo, amarillo, naranja

Los estados ocultos son:

$\mathcal{Q} = \{1, 2, \dots, N\}$

Las probabilidades de obtener un color en cada urna son:

urna 1	urna 2	...	urna N
P(rojo) = $b_1(1)$	P(rojo) = $b_2(1)$...	P(rojo) = $b_N(1)$
P(azul) = $b_1(2)$	P(azul) = $b_2(2)$...	P(azul) = $b_N(2)$
P(verde) = $b_1(3)$	P(verde) = $b_2(3)$...	P(verde) = $b_N(3)$
P(amarillo) = $b_1(4)$	P(amarillo) = $b_2(4)$...	P(amarillo) = $b_N(4)$
...
P(naranja) = $b_1(M)$	P(naranja) = $b_2(M)$...	P(naranja) = $b_N(M)$

Las probabilidades de pasar de una urna a otra son:

P(1,1) = a_{11}	P(2,1) = a_{21}	P(3,1) = a_{31}	...	P(N,1) = a_{N1}
P(1,2) = a_{12}	P(2,2) = a_{22}	P(3,2) = a_{32}	...	P(N,2) = a_{N2}
...
P(1,N) = a_{1N}	P(2,N) = a_{2N}	P(3,N) = a_{3N}	...	P(N,N) = a_{NN}

El primer problema consiste en decidir cual proceso es representado por los estados y después decidir cuantos estados pueden estar en el modelo.

Como se ilustró antes, el HMM más simple que corresponda al comportamiento de este proceso es aquel en el cual cada estado representa una urna específica y cada color representa un posible símbolo de observación. Por cada estado se define una probabilidad de extraer una bola (color) y una probabilidad de pasar a la siguiente urna. Los colores de las bolas dentro de cada urna pueden o no ser los mismos y pueden existir números diferentes de bolas de cada color en cada urna. Por lo tanto, una observación aislada de un color en particular no dice inmediatamente de cuál urna procede.

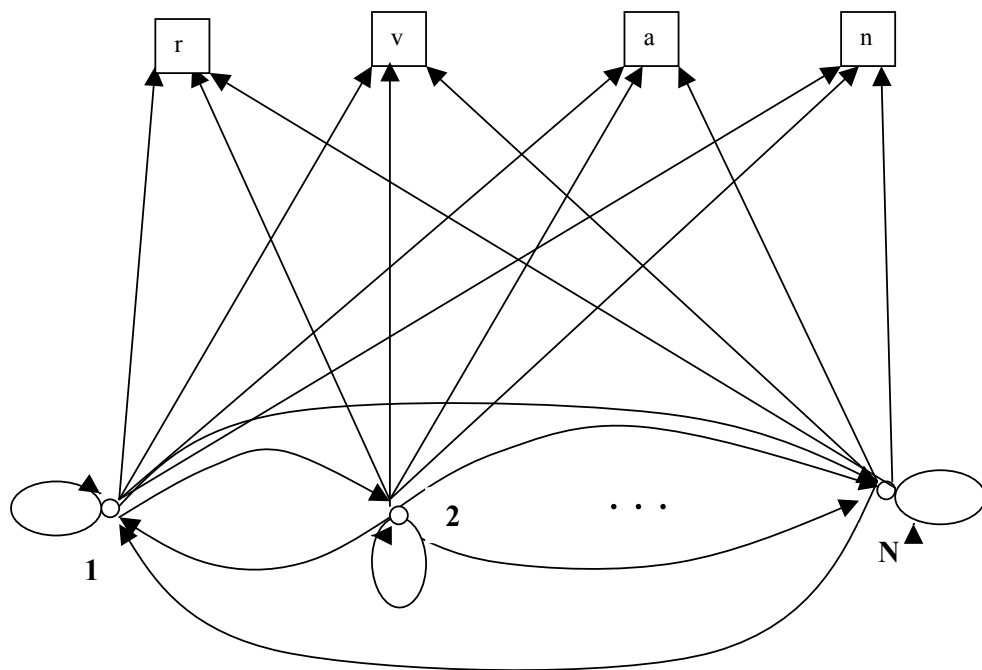


Figura 1.2 Arquitectura del grafo del modelo de urnas y bolas

1.3 ARQUITECTURAS DE HMMS

Un HMM puede ser representado como un grafo dirigido de transiciones/emisiones como se ilustra en la figura 1.2. La arquitectura específica que permita modelar de la mejor forma posible las propiedades observadas depende en gran medida de las características del problema. Las arquitecturas mas usadas son:

- 1) *Ergódicas o completamente conectadas* en las cuales cada estado del modelo puede ser alcanzado desde cualquier otro estado en un número finito de pasos (figura 1.2).
- 2) *Izquierda-derecha, hacia adelante o Bakis* las cuales tienen la propiedad de que en la medida que el tiempo crece se avanza en la secuencia de observación asociada O , y en esa misma medida el índice que señala el estado del modelo permanece o crece, es decir, los estados del sistema van de izquierda a derecha (figura 1.3). En secuencias biológicas y en reconocimiento de la voz estas arquitecturas modelan bien los aspectos lineales de las secuencias.
- 3) *Izquierda-derecha paralelas*, son dos arquitecturas *izquierda-derecha* conectadas entre sí (figura 1.4).

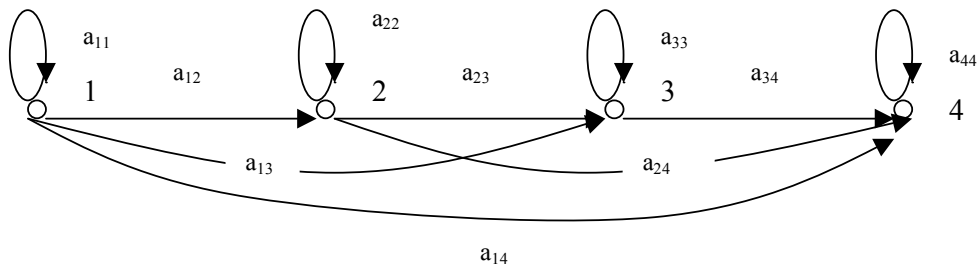


Figura 1.3 Modelo izquierda-derecha con 4 estados

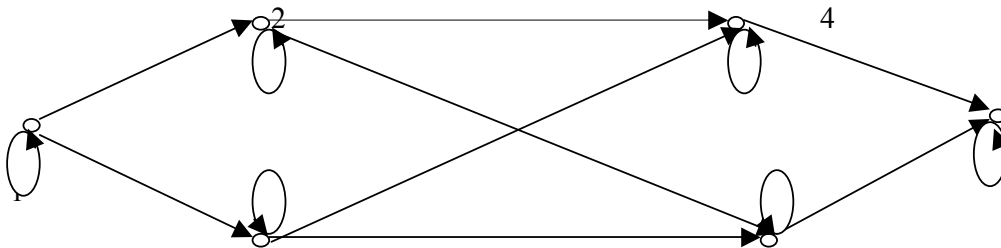


Figura 1.4 Modelo izquierda derecha paralelo con 6 estados

1.4 DEFINICION DE LOS ELEMENTOS DE UN HMM

Formalmente, un HMM discreto de primer orden se define como una cinco-tupla

$$\lambda = (\Sigma, Q, A, B, \pi)$$

donde

- i) $\Sigma = \{v_1, v_2, \dots, v_M\}$ es un alfabeto o conjunto discreto finito de M símbolos.
- ii) $Q = \{1, 2, \dots, N\}$ es un conjunto finito de N estados.
- iii) $A = (a_{ij})_{N \times N}$ es una matriz de probabilidades de transición donde a_{ij} es la probabilidad de transición desde el estado i al estado j , para todo $i, j \in N$.
- iv) $B = (b_j(o_t))_{N \times M}$ es un vector de probabilidades de emisión de símbolos, uno por cada estado, donde $b_j = (b_{j1}, b_{j2}, \dots, b_{jM})$ es la probabilidad de emisión del símbolo v_k del alfabeto en el estado j .
- v) $\pi = (\pi_1, \pi_2, \dots, \pi_N)$ es un vector de probabilidades del estado inicial q_o en Q .

Las probabilidades de iniciación, transición y emisión son los parámetros del modelo.

Un HMM puede ser usado como un generador de secuencias de observaciones $O = (o_1, o_2, \dots, o_T)$ donde:

- o_t es uno de los símbolos de Σ para $t \in \{1, 2, \dots, T\}$
- T es la longitud de la secuencia de observación O , es decir, el número de observaciones en la secuencia
- $\lambda = (A, B, \pi)$ son los parámetros del modelo.

Un HMM define una medida de probabilidad μ , sobre el espacio de secuencias Σ^* .

1.5 PROBLEMAS BASICOS DE LAS HMMs

Existen tres problemas básicos relacionados con los HMMs:

1. Calcular eficientemente $P(O|\lambda)$ la probabilidad de la secuencia de observación O dado el modelo $\lambda = (A, B, \pi)$ y la secuencia de observación $O = (o_1 o_2 \dots o_T)$.
2. Encontrar la trayectoria mas probable $q = (q_1 q_2 \dots q_T)$ dado el modelo λ y la secuencia de observación $O = (o_1 o_2 \dots o_T)$, es decir, $q = \arg \max_{r \in Q^*} P(r)$.
3. Ajustar los parámetros A, B, π para maximizar $P(O|\lambda)$

1.5.1 Estimativos iniciales de los parámetros del HMM

Existen diferentes formas de incorporar la información inicial para el diseño del HMM y estimar sus parámetros. La experiencia ha mostrado [9] que los estimativos iniciales uniformes o aleatorios de los parámetros π y A son adecuados en la mayoría de los casos. Y que los buenos estimativos iniciales para el parámetro B son útiles en el caso de secuencias con símbolos discretos y esenciales en el caso de distribuciones continuas. Por ejemplo, en biología computacional [1] se utiliza la distribución de Dirichlet para las probabilidades de transición y emisión y las distribuciones uniformes o aleatorias para las probabilidades del estado inicial.

1.5.2 Evaluación de la probabilidad de la secuencia

El problema 1 de los HMMs consiste en calcular la probabilidad de la secuencia de observación $\mathbf{O} = (o_1 o_2 \dots o_T)$, dado el modelo λ , es decir, $P(\mathbf{O}|\lambda)$. La forma más simple de resolver el problema 1 consiste en enumerar todas las posibles secuencias de estado de longitud T . Una de dichas secuencias es de la forma

$$q = (q_1 q_2 \dots q_T)$$

donde q_1 es el estado inicial.

La probabilidad de la secuencia de observación \mathbf{O} dada la anterior secuencia q es:

$$P(\mathbf{O}|q, \lambda) = \prod_{t=1}^T P(o_t | q_t, \lambda) = b_{q_1}(o_1) b_{q_2}(o_2) \dots b_{q_T}(o_T)$$

La probabilidad de la secuencia de estados q es:

$$P(q | \lambda) = \pi_{q_1} a_{q_1 q_2} a_{q_2 q_3} \dots a_{q_{T-1} q_T}$$

La probabilidad de que \mathbf{O} y q ocurran simultáneamente está dada por:

$$P(\mathbf{O}q | \lambda) = P(\mathbf{O}|q, \lambda)P(q | \lambda)$$

La probabilidad de \mathbf{O} sobre todas las posibles secuencias de estados Q , es el cálculo de:

$$P(O | \lambda) = P(O | q, \lambda)P(q | \lambda) = \pi_{q_1} b_{q_1}(o_1) a_{q_1 q_2} \dots a_{q_{T-1} q_T} b_{q_T}(o_T)$$

Esta expresión no permite un cómputo eficiente de la probabilidad debido a que el número de rutas en una arquitectura es exponencial. En cada tiempo $t = 1, 2, \dots, T$ se tienen N posibles estados alcanzables, por lo tanto N^T operaciones (figura 1.5). Existe un procedimiento mas eficiente para calcular dicha probabilidad, denominado *algoritmo de avance (forward algorithm)*.

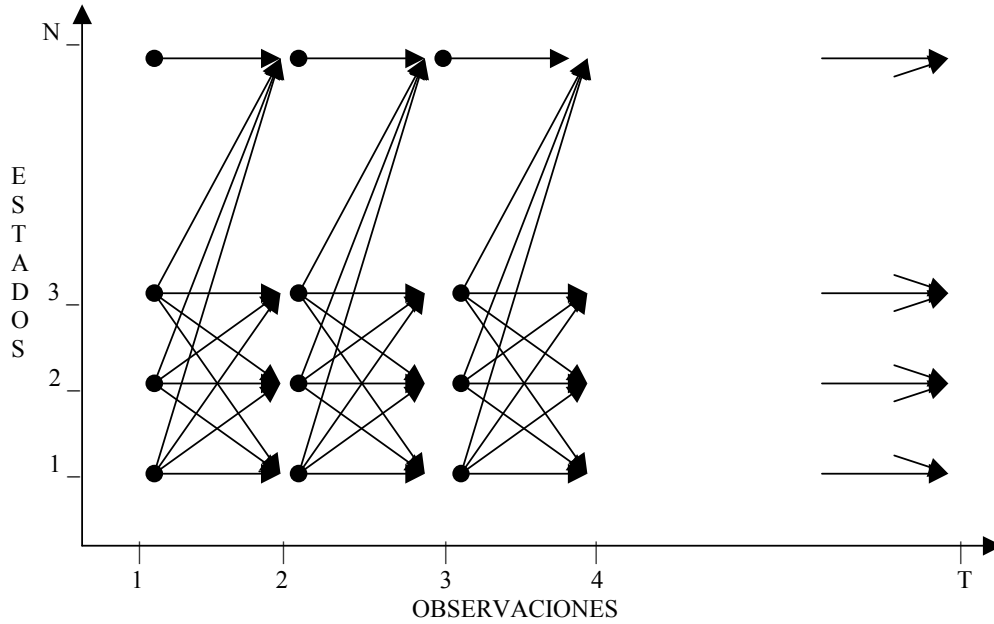


Figura 1.5 Rutas de un modelo ergódico con N estados y T observaciones

Algoritmo de avance (forward algorithm)

$$\text{Sea } \alpha(i) = P(o_1 o_2 \dots o_t, q_t = i | \lambda)$$

la probabilidad de la secuencia de observación parcial $o_1 o_2 \dots o_t$ en el estado i hasta el tiempo t , dado el modelo λ . Puede calcularse $\alpha_t(i)$, así:

$$1. \text{ Inicializa } : \alpha_1(i) = \pi_i b_i(o_1) \quad 1 \leq i \leq N$$

$$2. \text{Inducción : } \alpha_{t+1}(j) = \left[\sum_{i=1}^N \alpha_t(i) a_{ij} \right] b_j(o_{t+1}) \quad 1 \leq t \leq T-1, 1 \leq j \leq N$$

$$3. \text{Terminaci3n : } P(\mathbf{O} | \lambda) = \sum_{i=1}^N \alpha_T(i)$$

El n3mero de operaciones requeridas para calcular $\alpha_t(j)$, $1 \leq t \leq T$, $1 \leq j \leq N$, ilustradas en la figura 1.6, es exactamente de $N(N+1)(T-1)+N$ multiplicaciones y $N(N-1)(T-1)$ sumas, entonces es del orden de N^2T .

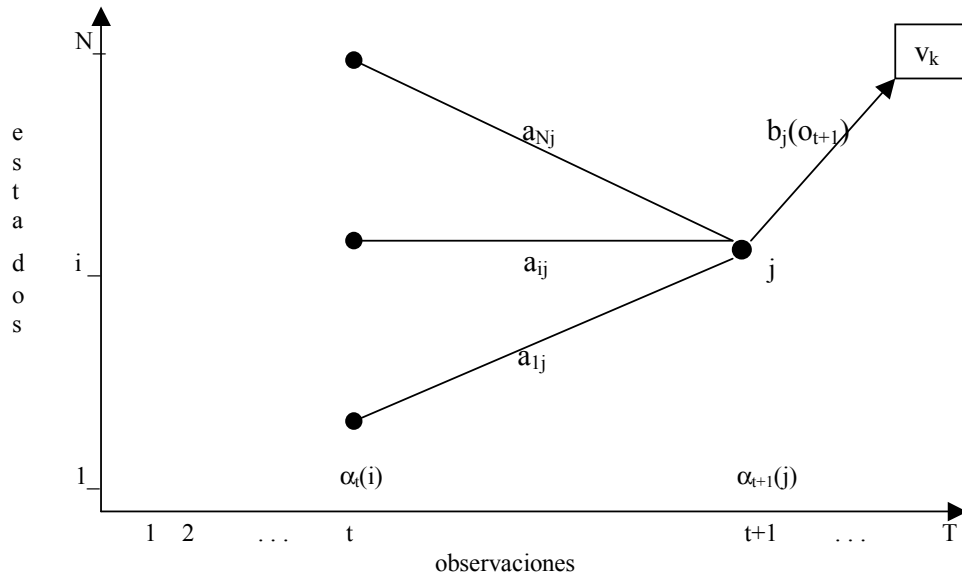


Figura 1.6 Secuencia de operaciones necesarias para el c3mputo de la variable $\alpha_{t+1}(j)$

Para el aprendizaje del HMM se deben propagar las probabilidades en forma inversa. El *algoritmo de retroceso (backward algorithm)* es la versi3n inversa del *algoritmo de forward* se define como sigue:

Algoritmo de retroceso (backward algorithm)

Sea $\beta_t(i) = P(o_{t+1}o_{t+2}\dots o_T \mid q_t = i, \lambda)$

La probabilidad de la secuencia de observación parcial desde $t+1$ hasta el final, dado el estado i en el tiempo t y el modelo λ . La variable $\beta_t(i)$ puede resolverse como sigue:

1. Inicializa $\beta_T(i) = 1, \quad 1 \leq i \leq N$
2. Inducción: $\beta_t(i) = \sum_{j=1}^N a_{ij}b_j(o_{t+1})\beta_{t+1}(j), \quad t=T-1, T-2, \dots, 1, \quad 1 \leq i \leq N$

Los cálculos necesarios de $\beta_t(i)$, $1 \leq t \leq T$, $1 \leq i \leq N$, es del orden de N^2T operaciones, según se ilustra en la figura 1.7.

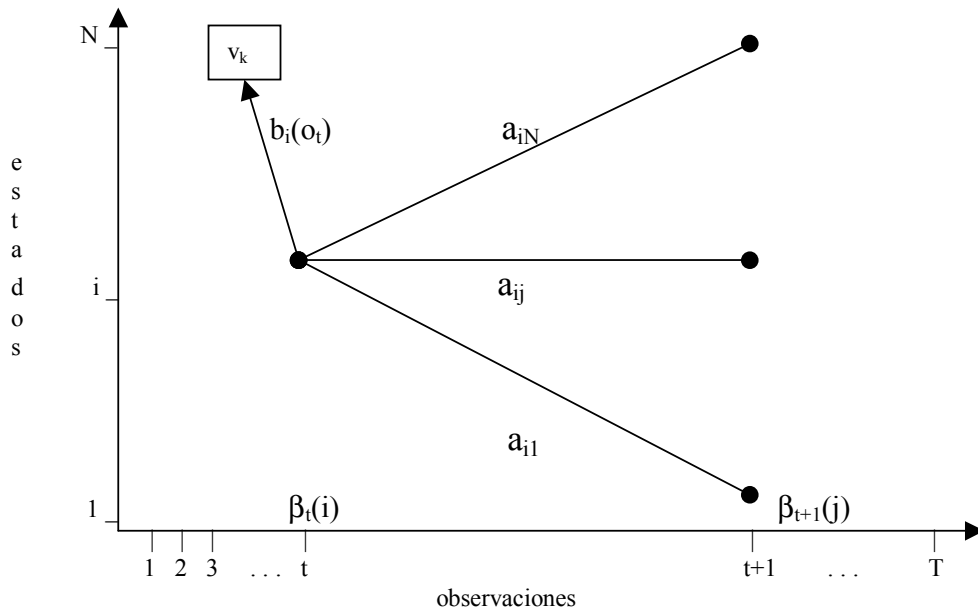


Figura 1.7 Secuencia de operaciones requeridas para el cálculo de la variable de retroceso $\beta_t(i)$

1.5.3 Secuencia óptima de estados

La ruta mas probable en un HMM es útil para el aprendizaje y para el alineamiento de secuencias con el modelo. La ruta mas probable $q = (q_1 q_2 \dots q_T)$ para la secuencia de observación $O = (o_1 o_2 \dots o_T)$ puede ser calculada utilizando el *algoritmo de Viterbi*.

$$\text{Sea } \delta_t(i) = \max_{q^1, q^2, \dots, q^{t-1}} P[q^1 q^2 \dots q^{t-1}, q_t = i, o_1 o_2 \dots o_t | \lambda]$$

la probabilidad asociada con la ruta mas probable a lo largo de un camino simple que toma en cuenta las primeras t observaciones y finaliza en el estado i . Por inducción se tiene que

$$\delta_{t+1}(j) = \left[\max_i \delta_t(i) a_{ij} \right] b_j(o_{t+1})$$

Para conservar la secuencia de estados se define la variable $\psi_t(j)$. El procedimiento para encontrar la mejor secuencia de estados se establece así:

1. Inicialización : $\delta_1(i) = \pi b_i(o_1) \quad 1 \leq i \leq N$

$$\Psi_1(i) = 0, \quad 1 \leq i \leq N$$

2. Inducción : $\delta_t(j) = \max_{1 \leq i \leq N} [\delta_{t-1}(i) a_{ij}] b_j(o_t), \quad 2 \leq t \leq T, \quad 1 \leq j \leq N$

$$\Psi_t(j) = \arg \max_{1 \leq i \leq N} [\delta_{t-1}(i) a_{ij}], \quad 2 \leq t \leq T, \quad 1 \leq j \leq N$$

3. Terminación : $P^* = \max_{1 \leq i \leq N} [\delta_T(i)]$

$$Q_T^* = \arg \max_{1 \leq i \leq N} [\delta_T(i)]$$

4. Ruta inversa : $q_t^* = \Psi_{t+1}(q_{t+1}^*), \quad t = T-1, T-2, \dots, 1.$

El algoritmo de Viterbi está basado en el método de la programación dinámica. La estructura de retículo implementa eficientemente los cálculos del algoritmo (figura 1.8).

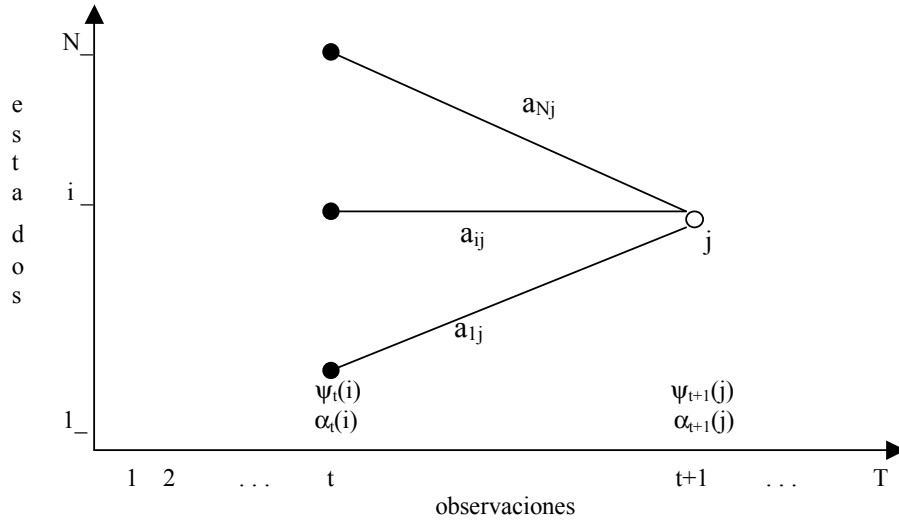


Figura 1.8 Secuencia de operaciones necesarias para calcular las variables $\alpha_t(j)$ y $\psi_t(j)$.

1.5.4 Aprendizaje del modelo

El problema mas difícil de los HMMs es determinar un método para ajustar los parámetros (A, B, π) del modelo para satisfacer los criterios de optimización. No se conoce una forma analítica para fijar los parámetros que maximice la probabilidad de la secuencia de observación. Varios algoritmos están disponibles para el entrenamiento de un HMM, entre ellos, Baum-Welch o EM (expectation maximization), GEM (EM generalizado) y diferentes formas de descenso por gradiente. El procedimiento de entrenamiento con EM [9] para reestimar los parámetros A, B y π , utiliza la variable $\xi_t(i, j)$, que es la probabilidad de encontrarse en el estado i en el tiempo t y en el estado j en el tiempo $t+1$, dado el modelo λ y la secuencia de observación O , esto es

$$\xi_t(i, j) = P[q_t = i, q_{t+1} = j | O, \lambda] = \frac{P[q_t = i, q_{t+1} = j, O | \lambda]}{P(O | \lambda)}$$

La variable $\xi_t(i, j)$ puede reescribirse a partir de las variables $\alpha_t(j)$ y $\beta_t(i)$, así:

$$\xi_t(i, j) = \frac{\alpha_t(i) a_{ij} b_j(o_{t+1}) \beta_{t+1}(j)}{\sum_{i=1}^N \sum_{j=1}^N \alpha_t(i) a_{ij} b_j(o_{t+1}) \beta_{t+1}(j)}$$

La sumatoria de $\xi_t(i, j)$ sobre t puede interpretarse como el número esperado de transiciones desde el estado i hasta el estado j , es decir,

$$\sum_{t=1}^{T-1} \sum_{j=1}^N \xi_t(i, j) = \text{número esperado de transiciones desde el estado } i \text{ en } \mathbf{O}$$

$$\sum_{t=1}^{T-1} \xi_t(i, j) = \text{número esperado de transiciones desde el estado } i \text{ hasta el estado } j \text{ en } \mathbf{O}$$

Con las anteriores expresiones se pueden enunciar las siguientes fórmulas de reestimación para A, B y π .

π'_i = frecuencia esperada en el estado i en el tiempo $t = 1$, entonces es igual a

$$\pi'_i = \sum_{j=1}^N \xi_1(i, j)$$

$$a'_{ij} = \frac{\text{número esperado de transiciones desde el estado } i \text{ hasta el estado } j}{\text{número esperado de transiciones desde el estado } i}$$

$$a_{ij} = \frac{\sum_{t=1}^{T-1} \alpha(i) a_{ij} b_j(o_{t+1}) \beta_{t+1}(j)}{\sum_{t=1}^{T-1} \sum_{j=1}^N \alpha(i) a_{ij} b_j(o_{t+1}) \beta_{t+1}(j)}$$

$$b_j(k)' = \frac{\text{número esperado de veces en el estado } j \text{ y observando el símbolo } v_k}{\text{número esperado de veces en el estado } j}$$

$$= \frac{\sum_{t=1}^T \sum_{i=1}^N \sum_{j=1}^N \alpha_t(j) a_{ij} b_j(o_t + 1) \beta_{t+1}(j)}{\sum_{t=1}^T \sum_{i=1}^N \sum_{j=1}^N \alpha_t(j) a_{ij} b_j(o_t + 1) \beta_{t+1}(j)}$$

Las trayectorias que satisfacen las condiciones requeridas por la ecuación se ilustran en la figura 1.9.

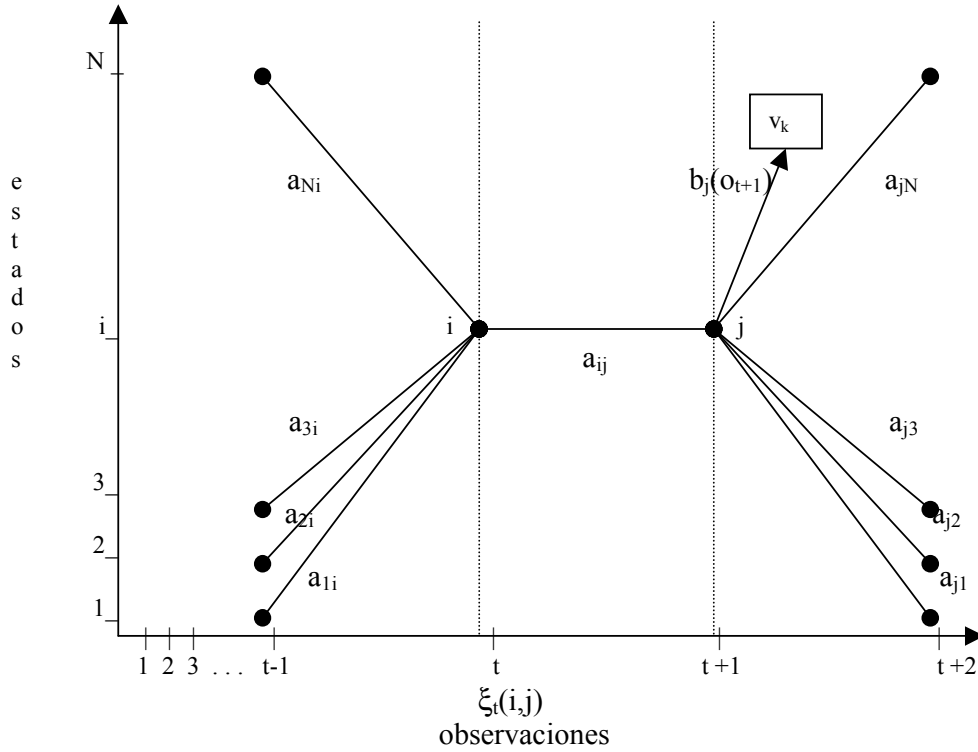


Figura 1.9 Secuencia de operaciones requeridas para el cálculo de $\xi_t(i,j)$.

Si designamos al modelo reestimado como $\lambda' = (A', B', \pi')$ pueden presentarse dos posibilidades (1) el modelo inicial $\lambda = (A, B, \pi)$ define un punto crítico de la función de probabilidad, en cuyo caso, $\lambda' = \lambda$; o (2) modelo λ' es mas probable que el modelo λ , es decir, $P(O | \lambda') > P(O | \lambda)$ y se ha conseguido un nuevo modelo λ' a partir del cual la secuencia de observación es mas probable.

1.6 SECUENCIAS DE OBSERVACIÓN CONTINUAS Y SECUENCIAS DE OBSERVACIÓN MÚLTIPLES

Hasta ahora se han considerado secuencias de observación caracterizadas por símbolos discretos que pertenecen a un alfabeto finito y que usan probabilidades discretas en cada estado del modelo, no obstante, en algunos problemas las secuencias de observación son señales continuas y por lo tanto es conveniente usar HMMs con densidades de observación continuas y funciones de densidad de probabilidades que aseguren la reestimación consistente de los parámetros del modelo. Por otra parte, se ha revisado el entrenamiento de los parámetros del modelo con una sola secuencia de observación, pero en la práctica existen muchas aplicaciones, tales como reconocimiento de la voz y alineamiento de secuencias biológicas en las que se debe trabajar con múltiples secuencias de observación para hacer mas confiable el modelo, esto es,

$$\mathbf{O} = [\mathbf{O}^{(1)}, \mathbf{O}^{(2)}, \dots, \mathbf{O}^{(k)}]$$

donde

$$\mathbf{O}^{(k)} = (O_1^{(k)} O_2^{(k)} \dots O_T^{(k)}) \text{ es la } k\text{-ésima secuencia de observación}$$