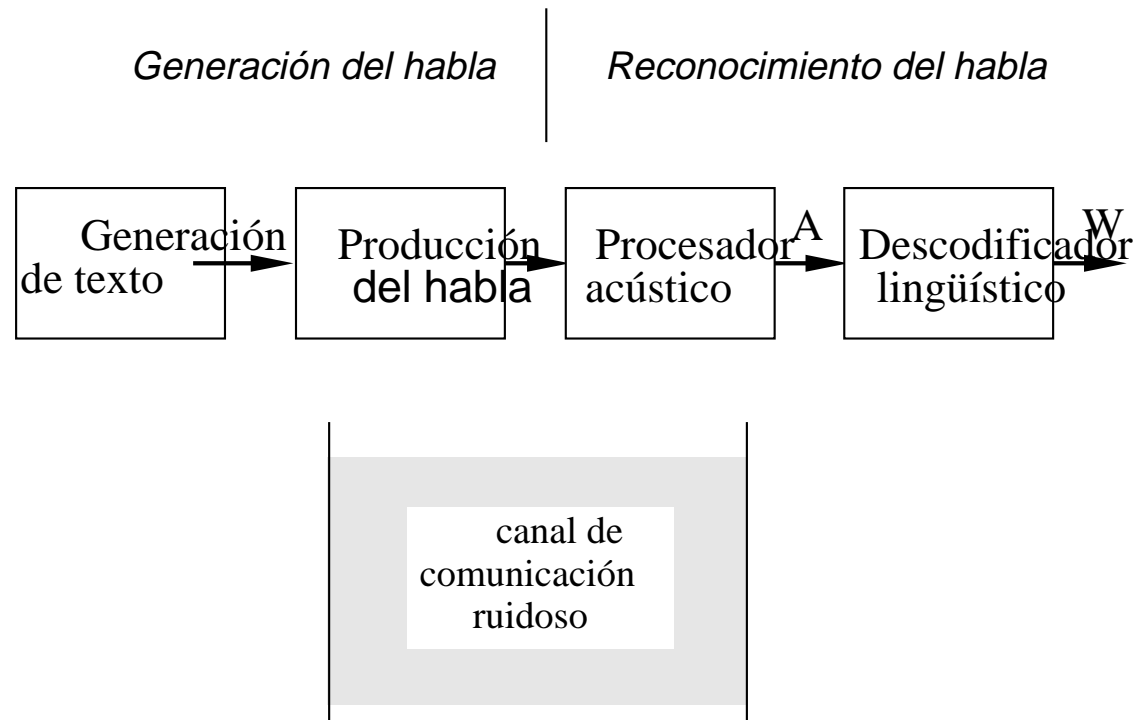


Modelos ocultos de Markov

- Introducción
- Formulación del problema
- Algoritmo de avance-retroceso
- Búsqueda de Viterbi
- Estimación del parámetro de Baum-Welch
- Otras consideraciones
 - Secuencias de observación múltiple
 - Modelos telefónicos para reconocimiento de voz continuo
 - Modelos ocultos de Markov (HMMs) de densidad continua
 - Cuestiones de implementación



El reconocimiento se consigue al maximizar la probabilidad de la cadena lingüística \mathbf{W} , dadas las pruebas acústicas \mathbf{A} , ej., elegir la secuencia lingüística $\hat{\mathbf{W}}$ tal que

$$P(\hat{\mathbf{W}}|\mathbf{A}) = \max_{\mathbf{W}} P(\mathbf{W}|\mathbf{A})$$

Enfoque teórico de información al ASR

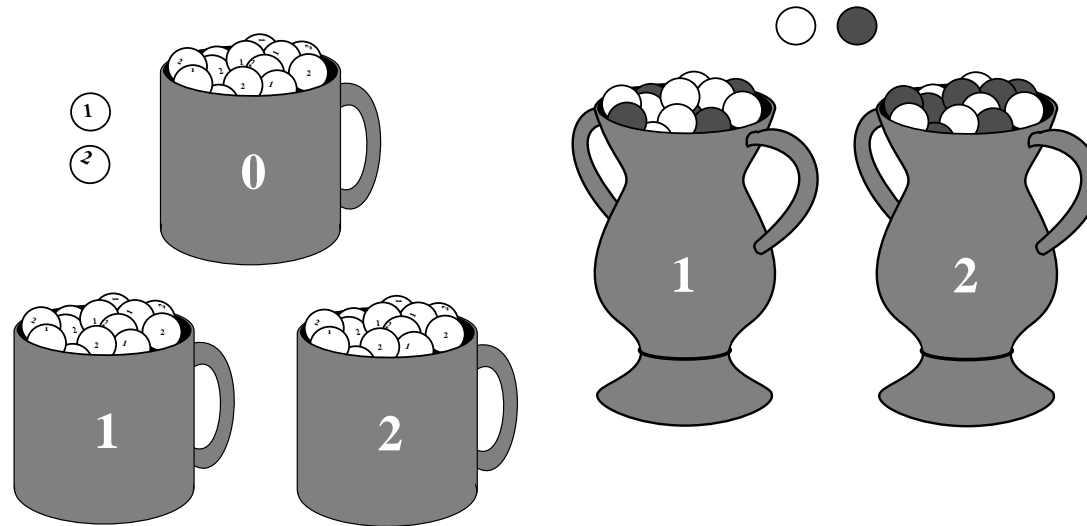
- A partir de la regla de Bayes:

$$P(\mathbf{W}|\mathbf{A}) = \frac{P(\mathbf{A}|\mathbf{W})P(\mathbf{W})}{P(\mathbf{A})}$$

- Los modelos ocultos de Markov (HMM) se ocupan de la cantidad $P(\mathbf{A}|\mathbf{W})$
- Cambio en la notación:

$$\begin{array}{lll} \mathbf{A} & \rightarrow & \mathbf{O} \\ \mathbf{W} & \rightarrow & \lambda \\ P(\mathbf{A}|\mathbf{W}) & \rightarrow & P(\mathbf{O}|\lambda) \end{array}$$

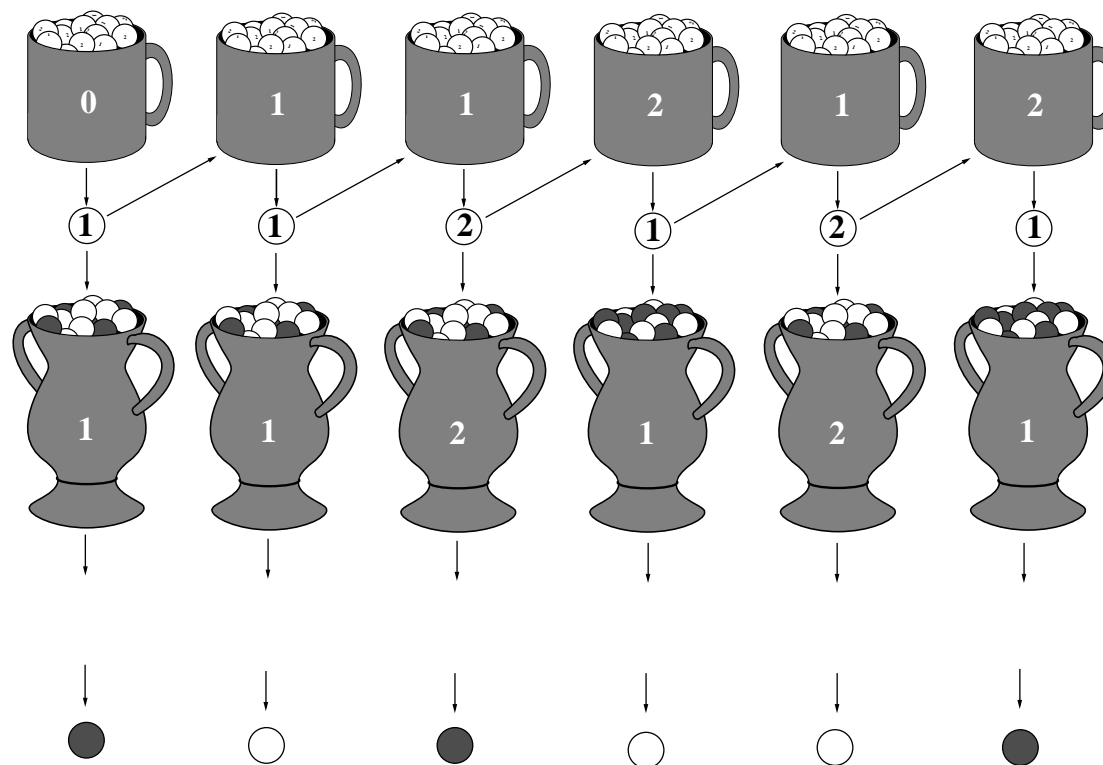
HMM: Un ejemplo



- Considere 3 tazas, cada una conteniendo mezclas de piedras con *estados* 1 y 2
- Las fracciones para la taza *ith* son a_{i1} y a_{i2} , y $a_{i1} + a_{i2} = 1$
- Considere 2 urnas, cada una con mezclas de bolas blancas y negras.
- Las fracciones para la urna *ith* son $b_i(B)$ y $b_i(W)$; $b_i(B) + b_i(W) = 1$
- El vector de parámetro para este modelo es:

$$\lambda = \{a_{01}, a_{02}, a_{11}, a_{12}, a_{21}, a_{22}, b_1(B), b_1(W), b_2(B), b_2(W)\}$$

HMM: Un ejemplo (continuación)



Secuencia de observación: $\mathbf{O} = \{B, W, B, W, W, B\}$

Secuencia de estado: $\mathbf{Q} = \{1, 1, 2, 1, 2, 1\}$

Objetivo: Dado el modelo λ y la secuencia de observación \mathbf{O} , ¿cómo se puede determinar la secuencia de estado subyacente \mathbf{Q} ?

Elementos de un modelo oculto de Markov discreto

- N : número de estados del modelo
 - estados, $\mathbf{S} = \{s_1, s_2, \dots, s_N\}$
 - estado en tiempo t , $q_t \in \mathbf{S}$
- M : número de símbolos de observación (ej., observaciones discretas)
 - símbolos de observación, $\mathbf{V} = \{v_1, v_2, \dots, v_M\}$
 - observación en tiempo t , $o_t \in \mathbf{V}$
- $\mathbf{A} = \{a_{ij}\}$: distribución de la probabilidad de la transición del estado
 - $a_{ij} = P(q_{t+1} = s_j | q_t = s_i), 1 \leq i, j \leq N$
- $\mathbf{B} = \{b_j(k)\}$: distribución de la probabilidad del símbolo de observación del estado j
 - $b_j(k) = P(v_k \text{ at } t | q_t = s_j), 1 \leq j \leq N, 1 \leq k \leq M$
- $\pi = \{\pi_i\}$: distribución del estado inicial
 - $\pi_i = P(q_1 = s_i), 1 \leq i \leq N$

Desde una perspectiva notacional, un HMM se escribe típicamente como: $\lambda = \{\mathbf{A}, \mathbf{B}, \pi\}$

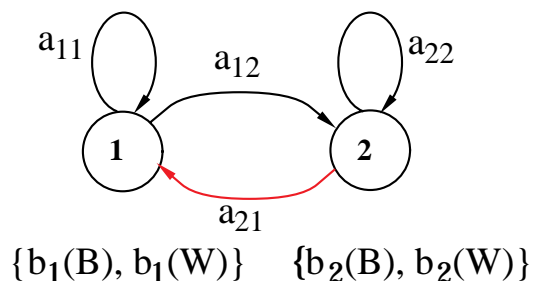
HMM: Un ejemplo (continuación)

Para nuestro ejemplo simple:

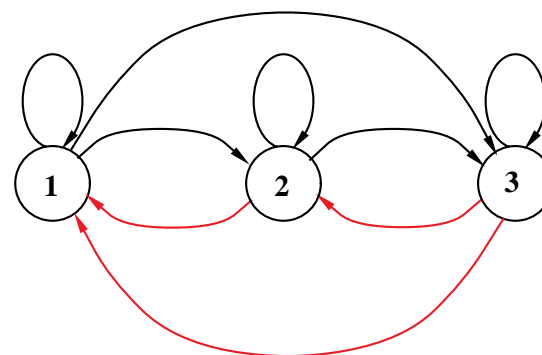
$$\pi = \{a_{01}, a_{02}\}, \quad \mathbf{A} = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix}, \text{ and } \mathbf{B} = \begin{bmatrix} b_1(B) & b_1(W) \\ b_2(B) & b_2(W) \end{bmatrix}$$

Diagrama de estado

estado-2

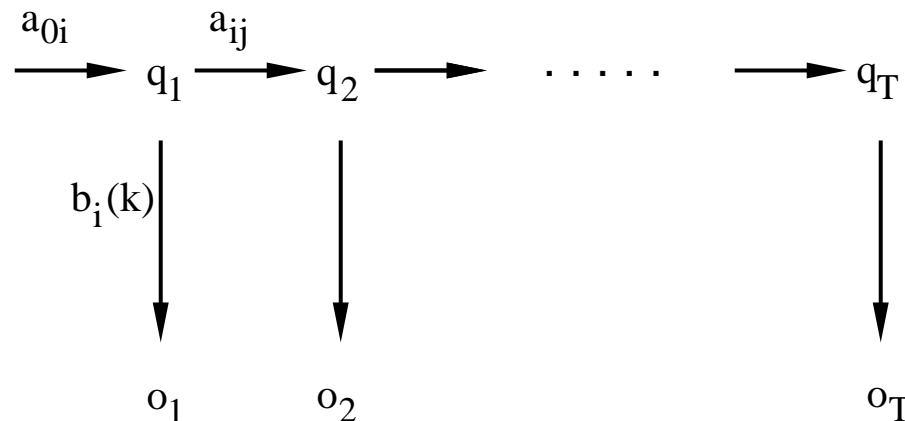


estado-3

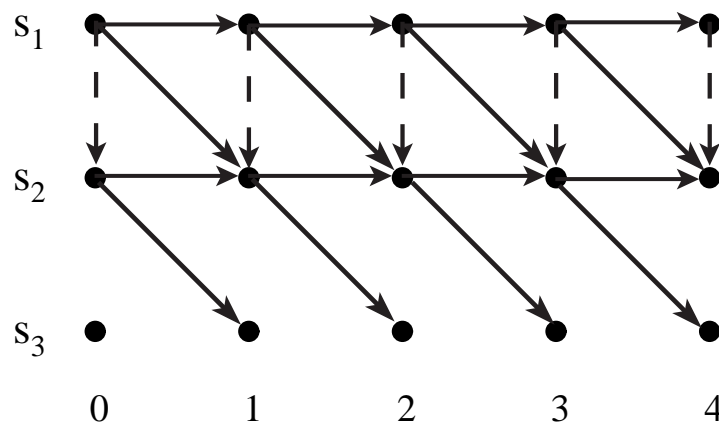
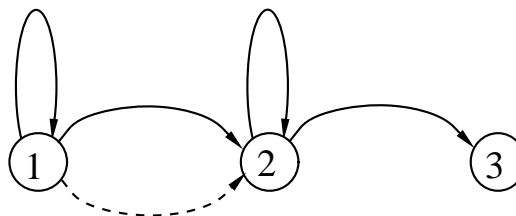


Generación de observaciones de HMM

1. Elegir un estado inicial $q_1 = s_i$, basado en la distribución del estado inicial, π
2. Para $t = 1$ a T :
 - Elegir $o_t = v_k$ en función de la distribución de probabilidad del símbolo en el estado s_i , $b_i(k)$
 - Transición a un nuevo estado $q_{t+1} = s_j$ según la distribución de probabilidad de la transición de estado para el estado s_i , a_{ij}
3. Incrementar en 1, volver al paso 2 si $t \leq T$; de lo contrario, terminar



Representación del diagrama de estado de Trellis



La línea con guiones representa una transición *nula*, en la que no se genera ningún símbolo de observación.

Tres problemas básicos de HMM

1. **Puntuación:** Dada una secuencia de observación $\mathbf{O} = \{o_1, o_2, \dots, o_T\}$ y un modelo $\lambda = \{\mathbf{A}, \mathbf{B}, \pi\}$, ¿cómo calculamos $P(\mathbf{O} | \lambda)$, la probabilidad de la secuencia de observación?
==> Algoritmo de avance-retroceso
2. **Ajuste:** Dada una secuencia de observación $\mathbf{O} = \{o_1, o_2, \dots, o_T\}$, ¿cómo elegimos una secuencia de estado $\mathbf{Q} = \{q_1, q_2, \dots, q_T\}$ que de algún modo sea optimal?
==> Algoritmo de Viterbi
3. **Entrenamiento:** ¿Cómo ajustamos los parámetros del modelo $\lambda = \{\mathbf{A}, \mathbf{B}, \pi\}$ para maximizar $P(\mathbf{O} | \lambda)$?
==> Procedimientos de reestimación de Baum-Welch

Cálculo de $P(\mathbf{O}|\lambda)$

$$P(\mathbf{O}|\lambda) = \sum_{all \mathbf{Q}} P(\mathbf{O}, \mathbf{Q}|\lambda)$$

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

- Considere la secuencia del estado *fijo*: $\mathbf{Q} = q_1 q_2 \dots q_T$

$$P(\mathbf{O}|\mathbf{Q}, \lambda) = b_{q_1}(o_1)b_{q_2}(o_2)\dots b_{q_T}(o_T)$$

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

Por tanto:

$$P(\mathbf{O}|\lambda) = \sum_{q_1, q_2, \dots, q_T} \pi_{q_1} b_{q_1}(o_1) a_{q_1 q_2} b_{q_2}(o_2) \dots a_{q_{T-1} q_T} b_{q_T}(o_T)$$

- Cálculo requerido $\approx 2T \cdot NT$ (existen NT secuencias)
Para $N = 5, T = 100 \Rightarrow 2 \cdot 100 \cdot 5100 \approx 1072$ cálculos

El algoritmo de avance

- Definamos la variable de avance $\alpha_t(i)$, como la probabilidad de la secuencia de observación parcial hasta el tiempo t y estado s_i en el tiempo t , dado el modelo, ej.:

$$\alpha_t(i) = P(o_1 o_2 \dots o_t, q_t = s_i | \lambda)$$

- Se puede demostrar fácilmente que:

$$\alpha_1(i) = \pi_i b_i(o_1), \quad 1 \leq i \leq N$$

$$P(\mathbf{O} | \lambda) = \sum_{i=1}^N \alpha_T(i)$$

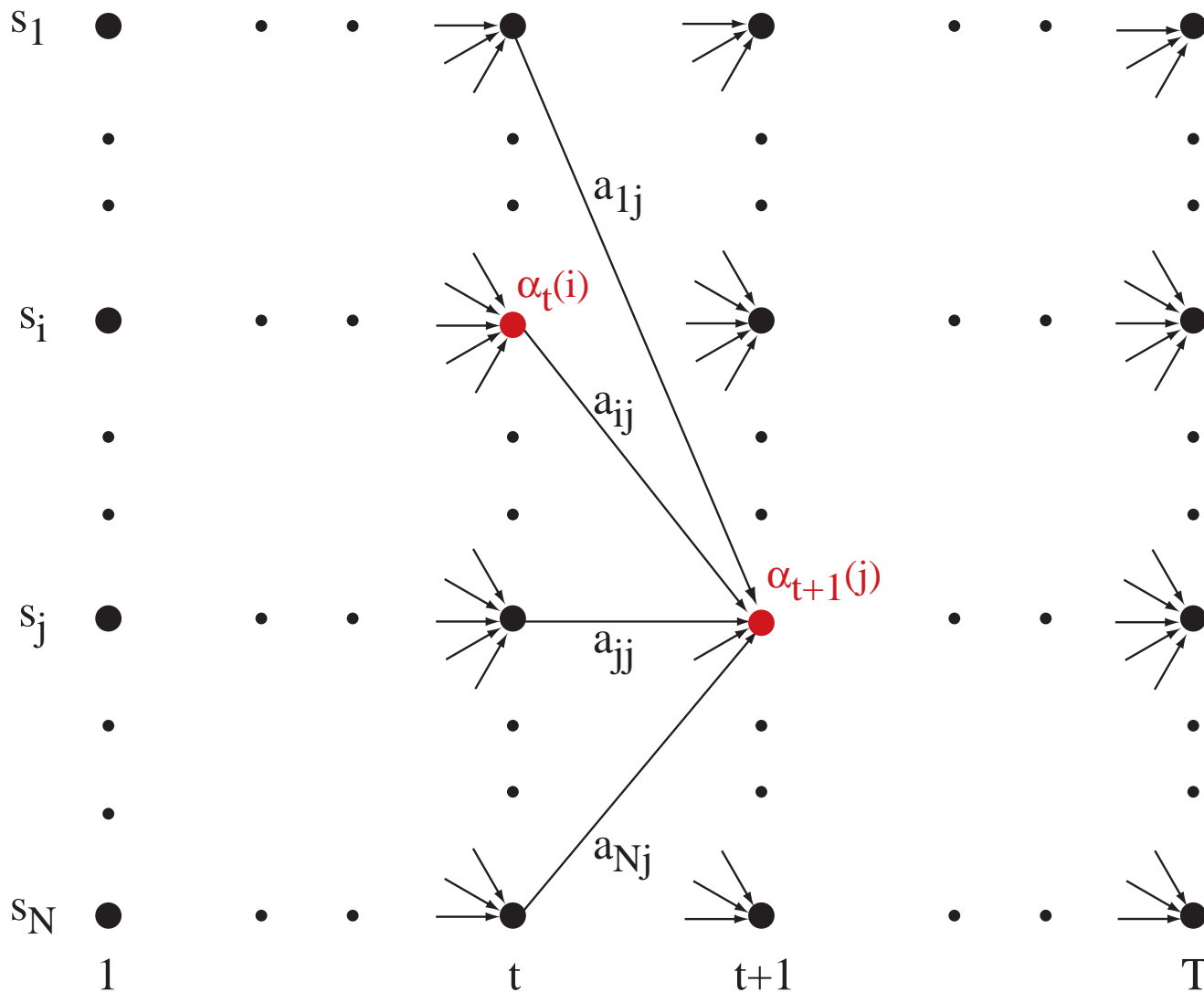
- Por inducción :

$$\alpha_{t+1}(j) = \left[\sum_{i=1}^N \alpha_t(i) a_{ij} \right] b_j(o_{t+1}), \quad \begin{matrix} 1 \leq t \leq T-1 \\ 1 \leq j \leq N \end{matrix}$$

- El cálculo está en el orden de $N^2 T$.

Para $N = 5, T = 100 \Rightarrow 100 \cdot 5^2$ cálculos, en vez de 10^7

Ilustración del algoritmo de avance



El algoritmo de retroceso

- Del mismo modo, definamos la variable de retroceso $\beta_t(i)$, como la probabilidad de la secuencia de observación parcial desde el tiempo $t + 1$ hasta el final, dado el estado s_i en tiempo t y el modelo, ej.,

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

- Puede demostrarse fácilmente que:

$$\beta_T(i) = 1, \quad 1 \leq i \leq N$$

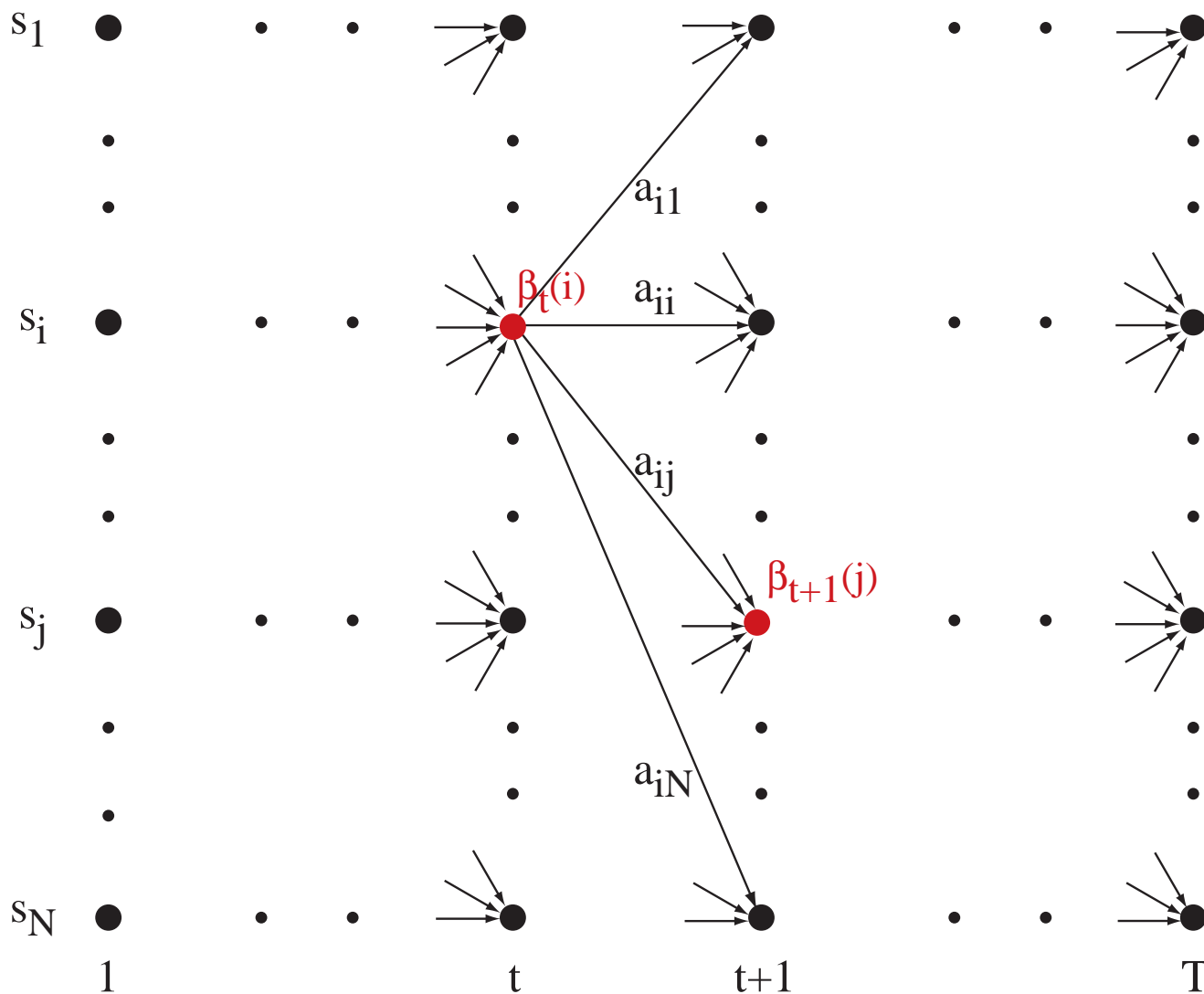
y:

$$P(\mathbf{O}|\lambda) = \sum_{i=1}^N \pi_i b_i(o_1) \beta_1(i)$$

- Por inducción:

$$\beta_t(i) = \sum_{j=1}^N a_{ij} b_j(o_{t+1}) \beta_{t+1}(j), \quad \begin{array}{l} t = T-1, T-2, \dots, 1 \\ 1 \leq i \leq N \end{array}$$

Ilustración del procedimiento de retroceso



Hallando secuencias optimales de estado

- Un criterio selecciona estados q_t , que son *individualmente* los más probables
 - Esto maximiza el número esperado de estados correctos
- Definamos $\gamma_t(i)$ como la probabilidad de estar en el estado s_i en el tiempo t , dada la secuencia de observación y el modelo, ej..

$$\gamma_t(i) = P(q_t = s_i | \mathbf{O}, \lambda) \quad \sum_{i=1}^N \gamma_t(i) = 1, \quad \forall t$$

- Luego el estado individualmente más probable q_t , en tiempo t es:

$$q_t = \operatorname{argmax}_{1 \leq i \leq N} \gamma_t(i) \quad 1 \leq t \leq T$$

- Observe que se puede demostrar que:

$$\gamma_t(i) = \frac{\alpha_t(i)\beta_t(i)}{P(\mathbf{O}|\lambda)}$$

Hallando secuencias optimales de estado

- El criterio de optimalidad individual presenta el problema de que la secuencia del estado optimal puede no obedecer a las restricciones de transición de estado.
- Otro criterio de optimalidad consiste en elegir la secuencia de estado que maximice $P(\mathbf{Q}, \mathbf{O}|\lambda)$; Esto se puede hallar mediante el algoritmo de *Viterbi*.
- Definamos $\delta_t(i)$ como la probabilidad más alta a lo largo de una trayectoria simple en tiempo t , que da cuenta de las primeras observaciones t , ej..

$$\delta_t(i) = \max_{q_1, q_2, \dots, q_{t-1}} P(q_1 q_2 \dots q_{t-1}, q_t = s_i, o_1 o_2 \dots o_t | \lambda)$$

- Por inducción: $\delta_{t+1}(j) = [\max_i \delta_t(i) a_{ij}] b_j(o_{t+1})$
- Para recuperar la secuencia de estado, debemos seguir la pista de la secuencia de estado que proporcionó la mejor trayectoria, en tiempo t , al estado s_i
 - Hacemos esto en una matriz aparte $\psi_t(i)$

El algoritmo de Viterbi

1. Inicialización:

$$\begin{aligned}\delta_1(i) &= \pi_i b_i(o_1), & 1 \leq i \leq N \\ \psi_1(i) &= 0\end{aligned}$$

2. Recursión:

$$\begin{aligned}\delta_t(j) &= \max_{1 \leq i \leq N} [\delta_{t-1}(i) a_{ij}] b_j(o_t), & 2 \leq t \leq T & \quad 1 \leq j \leq N \\ \psi_t(j) &= \operatorname{argmax}_{1 \leq i \leq N} [\delta_{t-1}(i) a_{ij}], & 2 \leq t \leq T & \quad 1 \leq j \leq N\end{aligned}$$

3. Terminación:

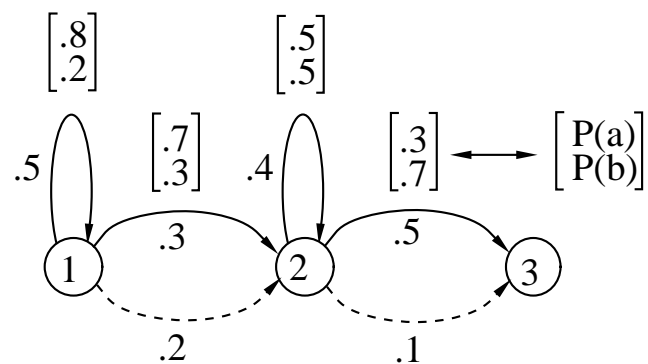
$$\begin{aligned}P^* &= \max_{1 \leq i \leq N} [\delta_T(i)] \\ q_T^* &= \operatorname{argmax}_{1 \leq i \leq N} [\delta_T(i)]\end{aligned}$$

4. Trayectoria inversa (secuencia de estado):

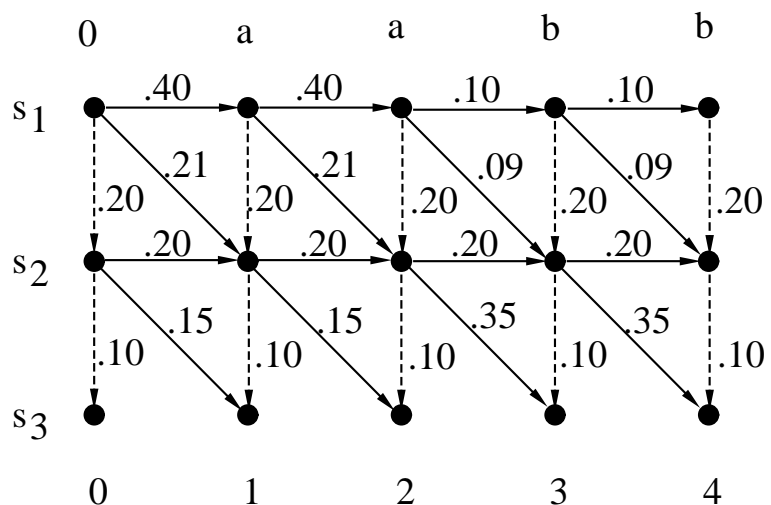
$$q_t^* = \psi_{t+1}(q_{t+1}^*), \quad t = T-1, T-2, \dots, 1$$

Cómputo $\approx N^2 T$

El algoritmo de Viterbi: Un ejemplo

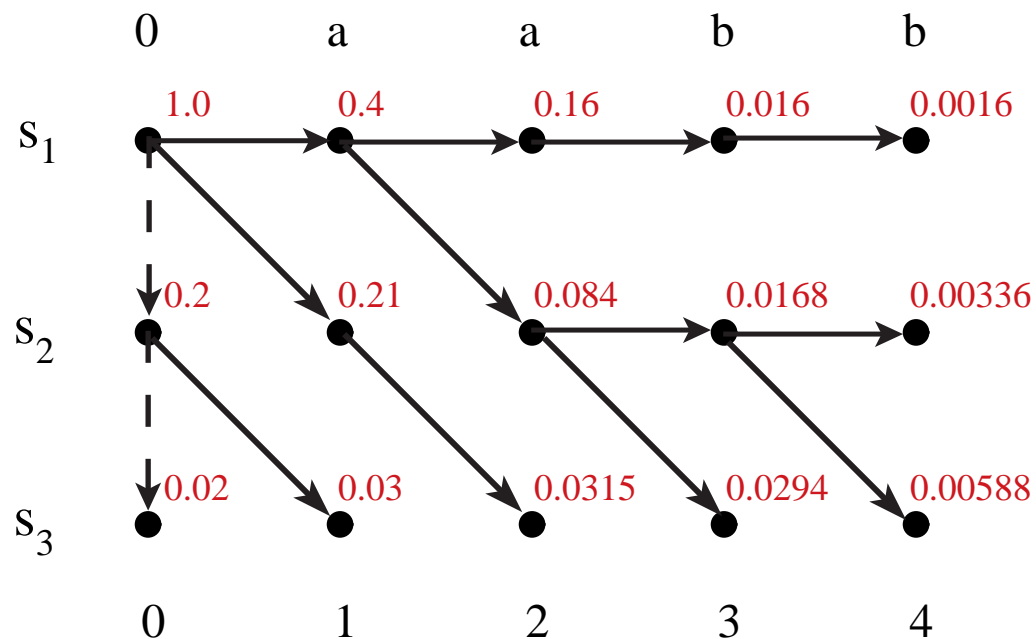


$O=\{a\ a\ b\ b\}$



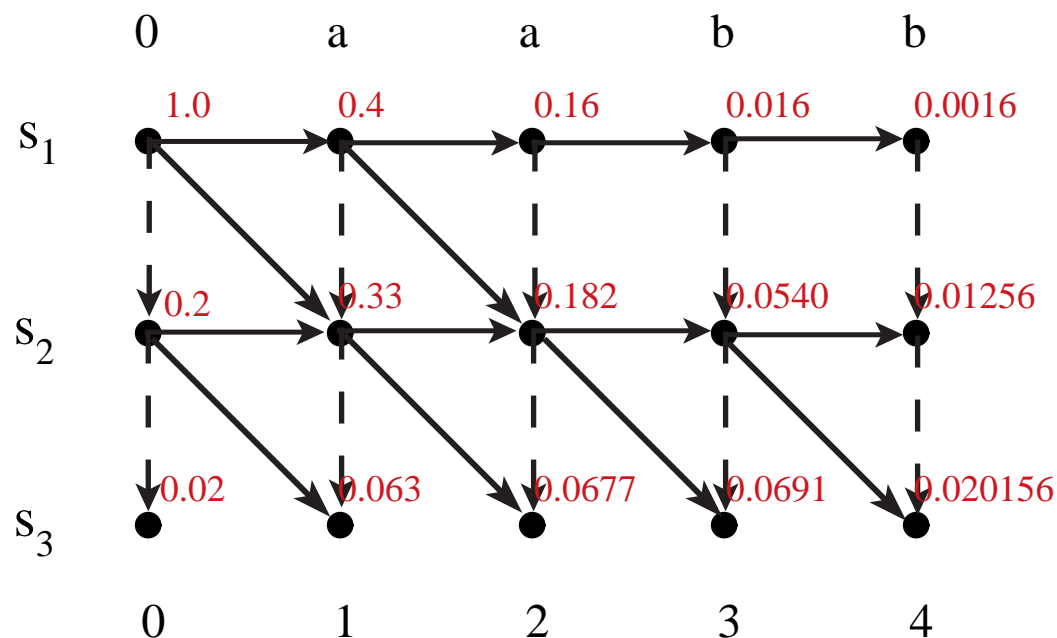
El algoritmo de Viterbi: Un ejemplo (continuación)

	0	a	aa	aab	$aabb$
s_1	1.0	s_1, a .4	s_1, a .16	s_1, b .016	s_1, b .0016
s_2	$s_1, 0$.2	$s_1, 0$.08	$s_1, 0$.032	$s_1, 0$.0032	$s_1, 0$.00032
		s_1, a .21	s_1, a .084	s_1, b .0144	s_1, b .00144
s_3	$s_2, 0$.02	s_2, a .04	s_2, a .042	s_2, b .0168	s_2, b .00336
		$s_2, 0$.021	$s_2, 0$.0084	$s_2, 0$.00168	$s_2, 0$.000336
s_3	$s_2, 0$.02	s_2, a .03	s_2, a .0315	s_2, b .0294	s_2, b .00588
		$s_2, 0$.021	$s_2, 0$.0084	$s_2, 0$.00168	$s_2, 0$.000336



Correspondencia mediante el algoritmo de avance-retroceso

	0	<i>a</i>	<i>aa</i>	<i>aab</i>	<i>aabb</i>
s_1	1.0	s_1, a .4	s_1, a .16	s_1, b .016	s_1, b .0016
s_2	$s_1, 0$.2	$s_1, 0$.08	$s_1, 0$.032	$s_1, 0$.0032	$s_1, 0$.00032
		s_1, a .21	s_1, a .084	s_1, b .0144	s_1, b .00144
		s_2, a .04	s_2, a .066	s_2, b .0364	s_2, b .0108
s_3	$s_2, 0$.02	$s_2, 0$.033	$s_2, 0$.0182	$s_2, 0$.0054	$s_2, 0$.001256
		s_2, a .03	s_2, a .0495	s_2, b .0637	s_2, b .0189



Reestimación de Baum-Welch

- La reestimación de Baum-Welch emplea el algoritmo de EM para determinar los parámetros de ML
- Definir $\xi_t(i, j)$ como la probabilidad de estar en el estado s_i en tiempo t y en el estado s_j en tiempo $t + 1$, dado el modelo y la secuencia de observación

$$\xi_t(i, j) = P(q_t = s_i, q_{t+1} = s_j | \mathbf{O}, \lambda)$$

- Luego:

$$\xi_t(i, j) = \frac{\alpha_t(i) a_{ij} b_j(o_{t+1}) \beta_{t+1}(j)}{P(\mathbf{O} | \lambda)}$$

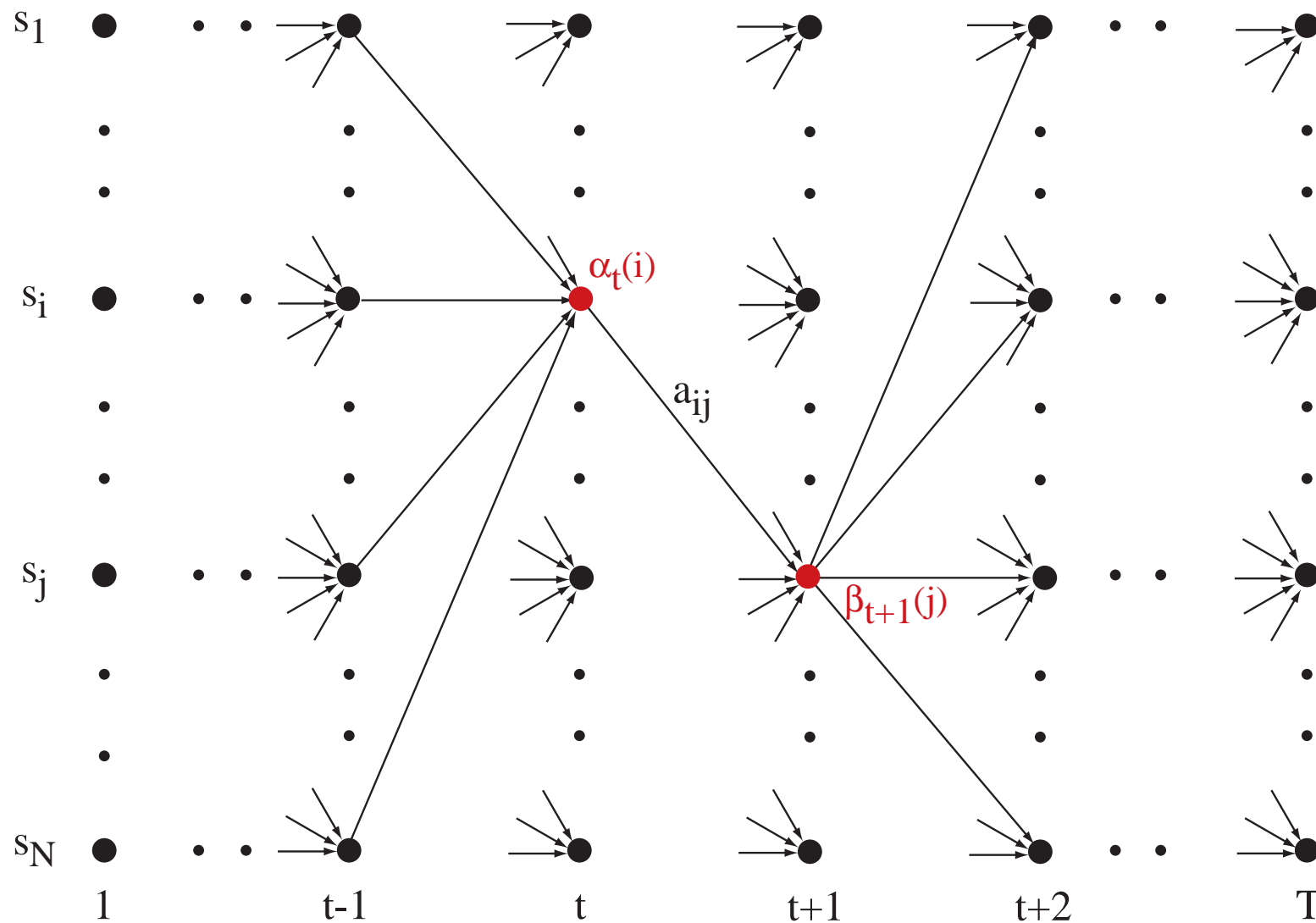
$$\gamma_t(i) = \sum_{j=1}^N \xi_t(i, j)$$

- Sumando $\gamma_t(i)$ y $\xi_t(i, j)$, obtenemos:

$$\sum_{t=1}^{T-1} \gamma_t(i) = \text{número esperado de transiciones desde } s_i$$

$$\sum_{t=1}^{T-1} \xi_t(i, j) = \text{número esperado de transiciones desde } s_i \text{ a } s_j$$

Procedimientos de reestimación de Baum-Welch



Fórmulas de reestimación de Baum-Welch

$$\begin{aligned}\bar{\pi} &= \text{número esperado de tiempos en el estado } s_i \text{ en } t = 1 \\ &= \gamma_1(i)\end{aligned}$$

$$\begin{aligned}\bar{a}_{ij} &= \frac{\text{número esperado de transiciones desde el estado } s_i \text{ a } s_j}{\text{número esperado de transiciones desde el estado } s_i} \\ &= \frac{\sum_{t=1}^{T-1} \xi_t(i, j)}{\sum_{t=1}^{T-1} \gamma_t(i)}\end{aligned}$$

$$\begin{aligned}\bar{b}_j(k) &= \frac{\text{número esperado de tiempos en el estado } s_j \text{ con símbolo } v_k}{\text{número esperado de tiempos en el estado } s_j} \\ &= \frac{\sum_{\substack{t=1 \\ o_t=v_k}}^T \gamma_t(j)}{\sum_{t=1}^T \gamma_t(j)}\end{aligned}$$

Fórmulas de reestimación de Baum-Welch

- Si $\lambda = (\mathbf{A}, \mathbf{B}, \pi)$ es el modelo inicial, y $\bar{\lambda} = (\bar{\mathbf{A}}, \bar{\mathbf{B}}, \bar{\pi})$ es el modelo reestimado. Se puede demostrar entonces que:
 1. El modelo inicial λ , define un punto crítico de la función de probabilidad, en cuyo caso $\lambda = \bar{\lambda}$, o
 2. El modelo $\bar{\lambda}$ es más probable que λ en el sentido de que $P(\mathbf{O}|\bar{\lambda}) > P(\mathbf{O}|\lambda)$, ej., hemos encontrado un nuevo modelo $\bar{\lambda}$, a partir del cual es más probable que la secuencia de observación se haya producido.
- Por tanto, podemos mejorar la probabilidad de que \mathbf{O} sea observado a partir del modelo, si utilizamos iterativamente $\bar{\lambda}$ en lugar de λ y repetimos la reestimación hasta que se alcance algún punto restrictivo. El modelo resultante es conocido como el HMM con máxima probabilidad.

Secuencias de observación múltiple

- En reconocimiento de voz se utilizan normalmente los HMM de izquierda a derecha. Estos HMM no pueden entrenarse mediante una secuencia de observación simple, porque únicamente se encuentran disponibles un pequeño número de observaciones para entrenar a cada estado. Para obtener estimaciones fiables de parámetros del modelo, se deben emplear las secuencias de observación múltiples. En este caso, el procedimiento de reestimación debe ser modificado.

- Denotemos el conjunto de las secuencias de observación K como

$$\mathbf{O} = \{\mathbf{O}^{(1)}, \mathbf{O}^{(2)}, \dots, \mathbf{O}^{(K)}\}$$

donde $\mathbf{O}^{(k)} = \{O^{(k)}_1, O^{(k)}_2, \dots, O^{(k)}_{T_k}\}$ es la secuencia de observación k -th.

- Suponga que las secuencias de observación son mutuamente independientes, y que queremos calcular los parámetros con el fin de maximizar

$$P(\mathbf{O} \mid \lambda) = \prod_{k=1}^K P(\mathbf{O}^{(k)} \mid \lambda) = \prod_{k=1}^K P_k$$

Secuencias de observación múltiple (continuación)

- Dado que las fórmulas de reestimación están basadas en la frecuencia de aparición de varios eventos, podemos modificarlos mediante la adición de frecuencias individuales de aparición para cada secuencia.

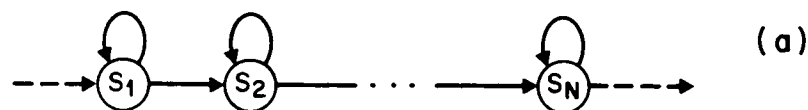
$$\bar{a}_{ij} = \frac{\sum_{k=1}^K \sum_{t=1}^{T_k-1} \xi_t^k(i, j)}{\sum_{k=1}^K \sum_{t=1}^{T_k-1} \gamma_t^k(i)} = \frac{\sum_{k=1}^K \frac{1}{P_k} \sum_{t=1}^{T_k-1} \alpha_t^k(i) a_{ij} b_j(o_{t+1}^{(k)}) \beta_{t+1}^k(j)}{\sum_{k=1}^K \frac{1}{P_k} \sum_{t=1}^{T_k-1} \alpha_t^k(i) \beta_t^k(i)}$$

$$\bar{b}_j(\ell) = \frac{\sum_{k=1}^K \sum_{\substack{t=1 \\ o_t^{(k)} = v_\ell}}^{T_k} \gamma_t^k(j)}{\sum_{k=1}^K \sum_{t=1}^{T_k} \gamma_t^k(j)} = \frac{\sum_{k=1}^K \frac{1}{P_k} \sum_{\substack{t=1 \\ o_t^{(k)} = v_\ell}}^{T_k} \alpha_t^k(i) \beta_t^k(i)}{\sum_{k=1}^K \frac{1}{P_k} \sum_{t=1}^{T_k} \alpha_t^k(i) \beta_t^k(i)}$$

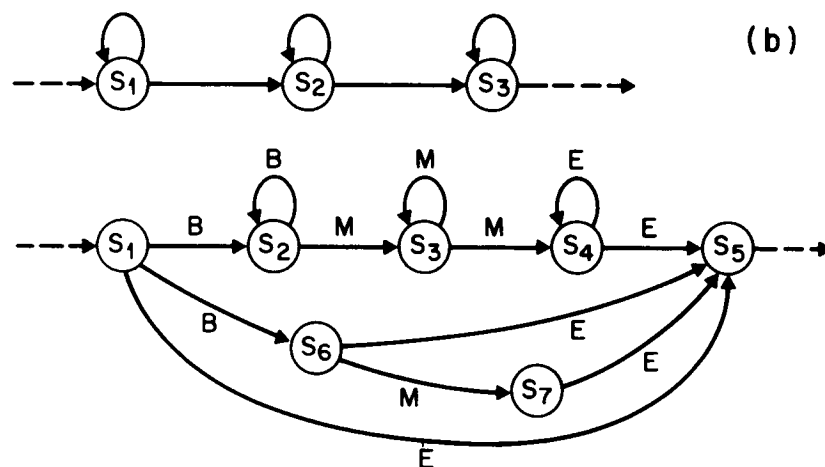
Los HMM en entornos telefónicos

- Los HMM basados en palabras son adecuados para el reconocimiento de pequeños vocabularios. Para el reconocimiento automático (ASR) de extensos vocabularios, los modelos basados en subpalabras (ej., entornos telefónicos) son más apropiados.

MODELO DE PALABRA

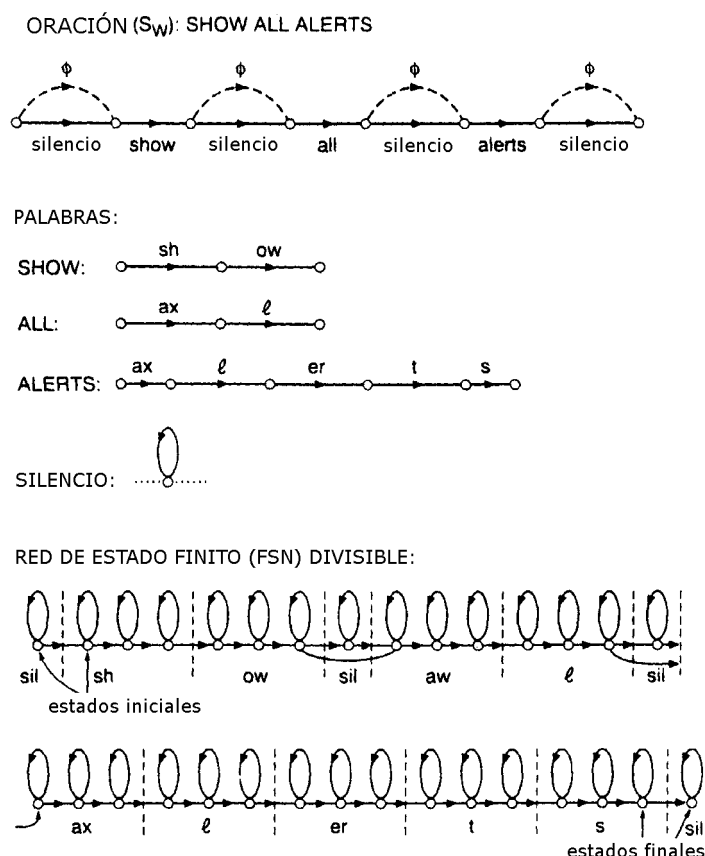


UNIDAD DE SUBPALABRA



HMMs en entornos telefónicos (continuación)

- Los modelos telefónicos pueden presentar muchos estados, y las palabras se forman a partir de una concatenación de modelo telefónicos.



Modelos ocultos de Markov de densidad continua

- Un HMM de *densidad continua* sustituye las probabilidades de observación discreta $b_j(k)$, mediante un PDF continuo $b_j(\mathbf{x})$
- Una práctica común es representar $b_j(\mathbf{x})$ como una mezcla de gaussianas:

$$b_j(\mathbf{x}) = \sum_{k=1}^M c_{jk} N[\mathbf{x}, \mu_{jk}, \Sigma_{jk}] \quad 1 \leq j \leq N$$

donde c_{jk} es el peso de mezcla

$$c_{jk} \geq 0 \quad (1 \leq j \leq N, 1 \leq k \leq M, \text{ y } \sum_{k=1}^M c_{jk} = 1, 1 \leq j \leq N),$$

N es la densidad normal, y

μ_{jk} y Σ_{jk} son el vector de medias y la matriz de covarianza asociados con el estado j y la mezcla k .

Variaciones en el modelado acústico

- Los HMM *semicontinuos* computan primero un libro de código de VQ de tamaño M
 - El libro de código de VQ se modela entonces como una familia de PDF gaussianos.
 - Cada palabra e código está representada por un PDF gaussiano, y se puede utilizar junto con otros para modelar los vectores acústicos.
 - Desde la perspectiva de los CD-HMM (HMM de densidad continua), esto equivale a utilizar el mismo conjunto de mezclas M para modelar todos los estados.
 - Por tanto, se le conoce normalmente como un HMM de *mezcla enlazada*.
- Los tres métodos se han utilizado en muchas tareas de reconocimiento de voz, con diversos resultados.
- Para extensos vocabularios, en el reconocimiento de voz continuo con una cantidad suficiente (ej., decenas de horas) de entrenamiento de datos, los sistemas CD-HMM actualmente producen el mejor rendimiento, pero con un aumento considerable el el cómputo.

Cuestiones de implementación

- Escalamiento: para prevenir la generación de un valor inferior al mínimo aceptable.
- *Entrenamiento de K -medias segmentales*: para entrenar probabilidades de observación, ejecutando en primer lugar el alineamiento de Viterbi.
- Estimaciones iniciales de λ : para facilitar modelos robustos.
- Recorte: para reducir el cómputo de búsqueda.

- X. Huang, A. Acero y H. Hon, *Spoken Language Processing*, Prentice-Hall, 2001.
- F. Jelinek, *Statistical Methods for Speech Recognition*. MIT Press, 1997.
- L. Rabiner y B. Juang, *Fundamentals of Speech Recognition*, Prentice-Hall, 1993.