

Maximum Entropy Markov Models (MEMM) for Information Extraction and Segmentation

Roman Firstein

Vortrag

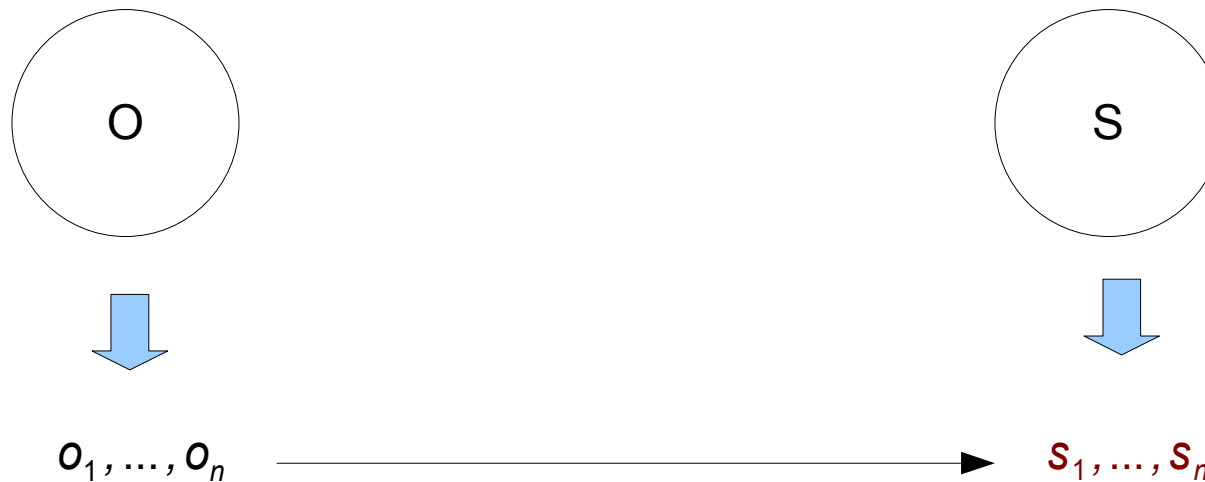
- Vorwort
- Modell
- (Maximale) Entropie am Beispiel der Urne
- Binäres Feature
- Bedingte Entropie
- Maximale Entropie für unseres Modell

Vorwort (1)

Worum geht es überhaupt?!

Gegeben:

- Menge O und S von Objekten (disjunkt)
- \mathcal{W} 'keitsraum als Teilmenge der Relationen $(s_1, \dots, s_j, o_1, \dots, o_j) \forall j$
- Stichprobe



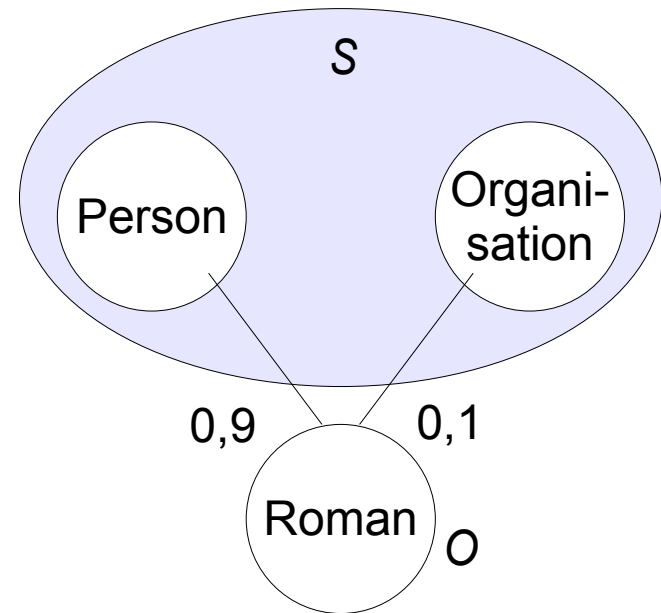
Gesucht:

- Wahrscheinlichkeitsfunktion $P(s_i, o_i)$, $1 \leq i \leq j$, so dass $(s_1, \dots, s_j, o_1, \dots, o_j)$ maximal ist.

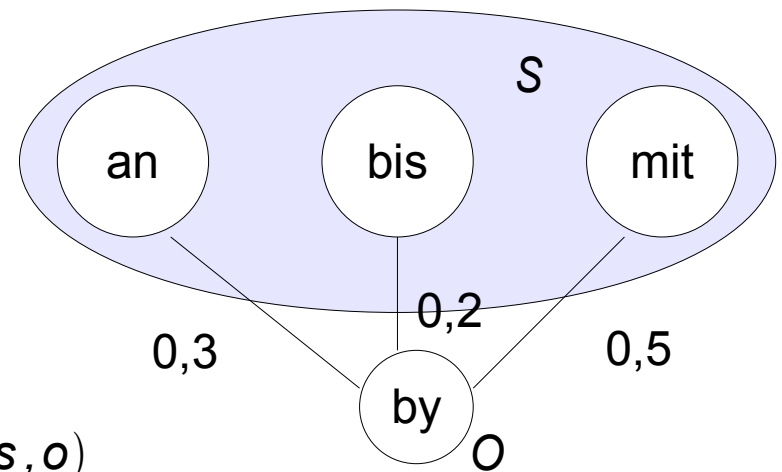
Vorwort (2)

Beispiele

1. Rolle Erkennung:

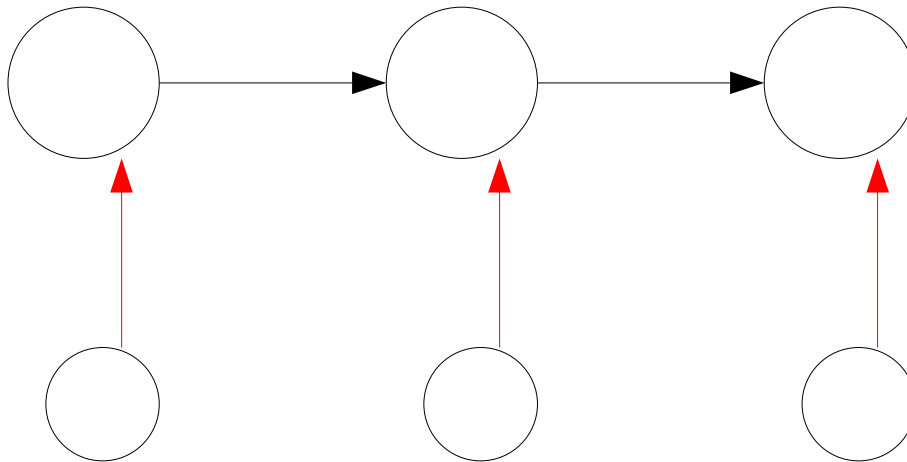


2. Übersetzung aus dem englischen:

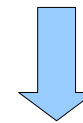


Zahlen sind die Wahrscheinlichkeiten $P(s, o)$

Das NEUE Modell (MEMM)

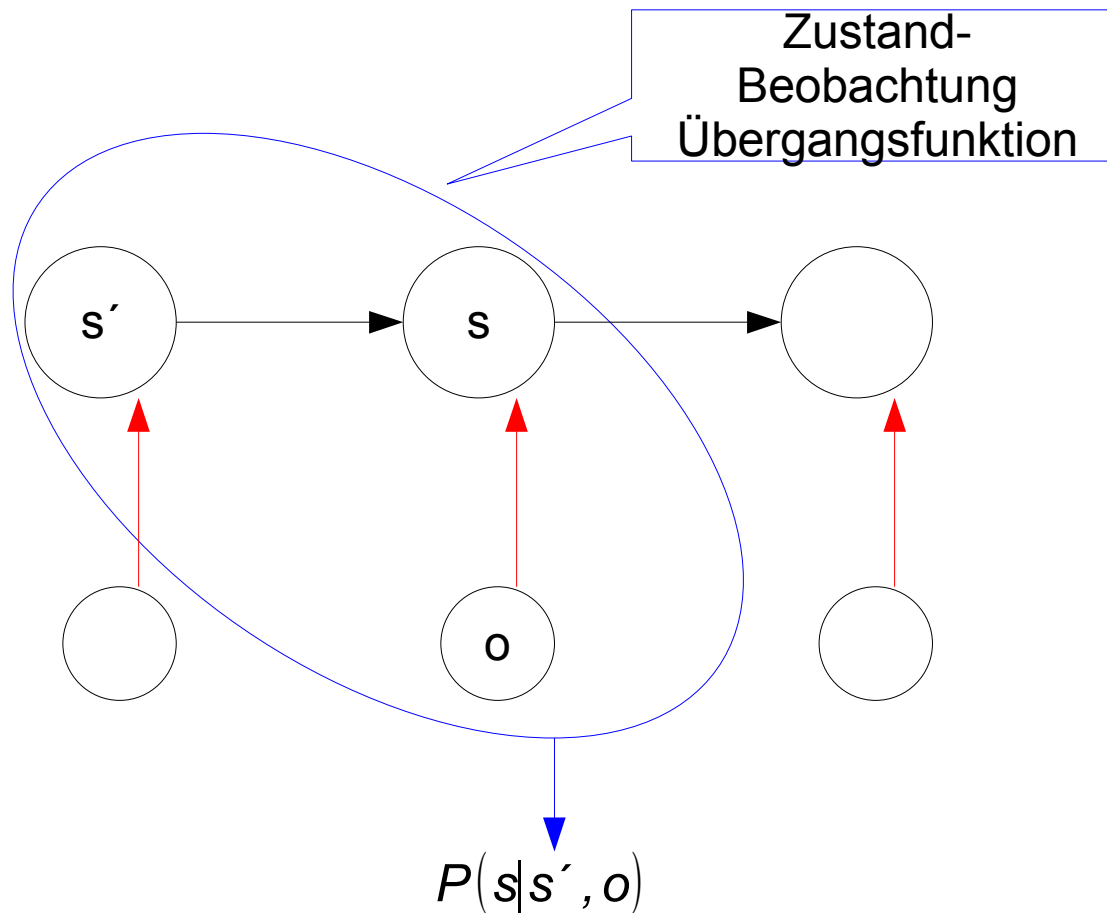


Zustand hängt jetzt
von der Beobachtung
und von dem
vorherigen Zustand ab



Vorteil:
bei der Übersetzung
kann man auch gucken
wie das vorherige Wort
übersetzt wurde.

Das NEUE Modell (MEMM)



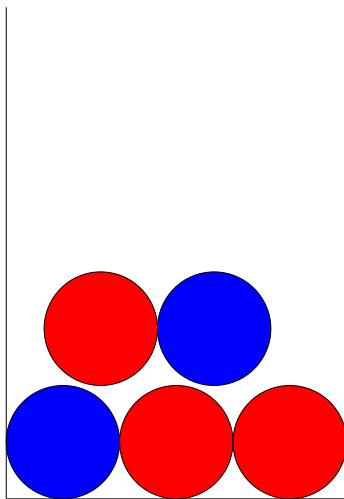
$P(s|s', o)$ unterteilen wir in $|S|$ Funktionen $P_{s'}(s|o)$.

Entropie

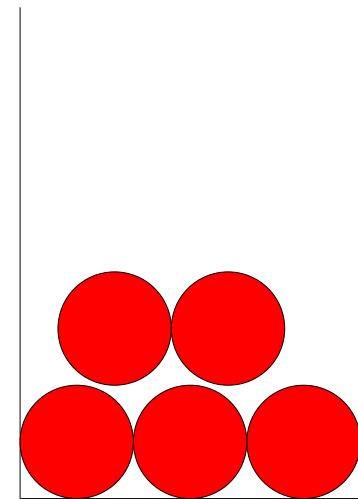
Frage: Wenn man die Verteilungen in der Urne nicht kennt, bei welcher der beiden Urnen wird die Überraschung größer?

Entropie H einer Wahrscheinlichkeitsverteilung X : Maß für die Unvorhersagbarkeit eines Ereignisses.

$$H(X) = - \sum_{i=1}^{|\{rot, blau\}|} p_i \cdot \log_2(p_i)$$



$$H(X) = -(0,6 \log_2(0,6) + 0,4 \log_2(0,4)) \approx 0,97$$



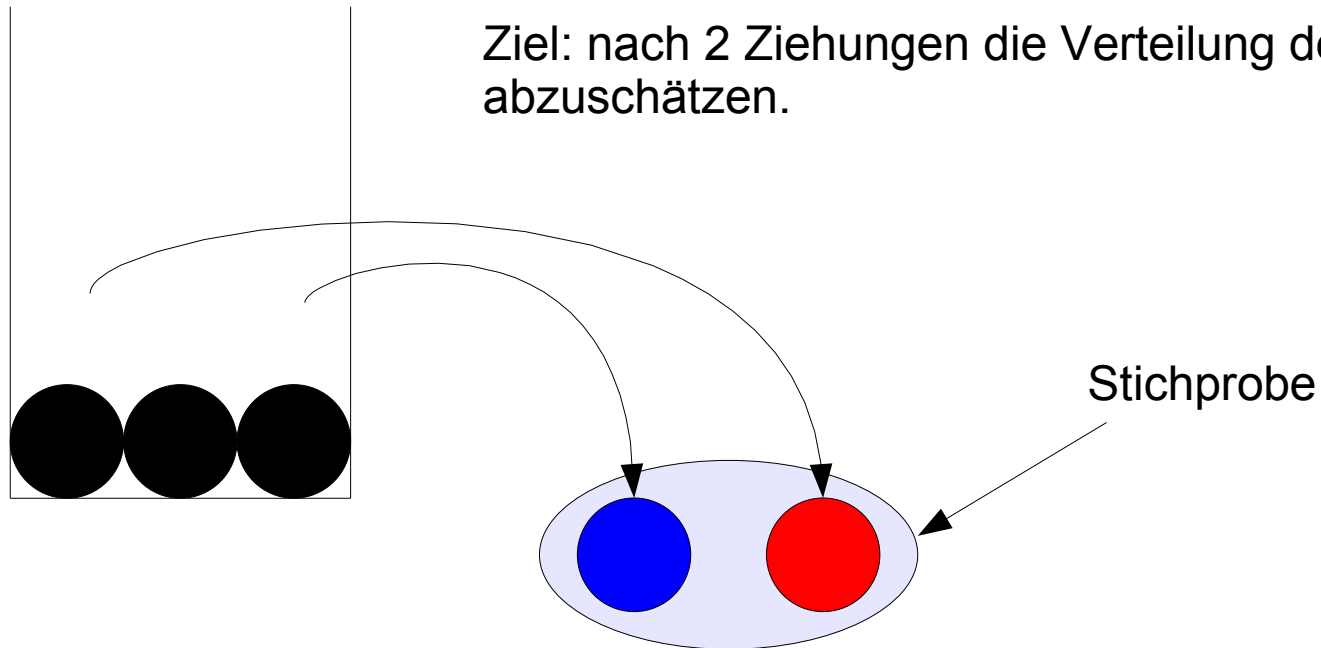
$$H(X) = -(\log_2(1)) = 0$$

Bei der (möglichst) Gleichverteilung wird die maximale Entropie erreicht !!!

Maximale Entropie

Gegeben: Die Urne enthält 5 Kugeln (rote und blaue). Die Anzahl der Kugeln sei bekannt, aber der jeweiligen Farben nicht.

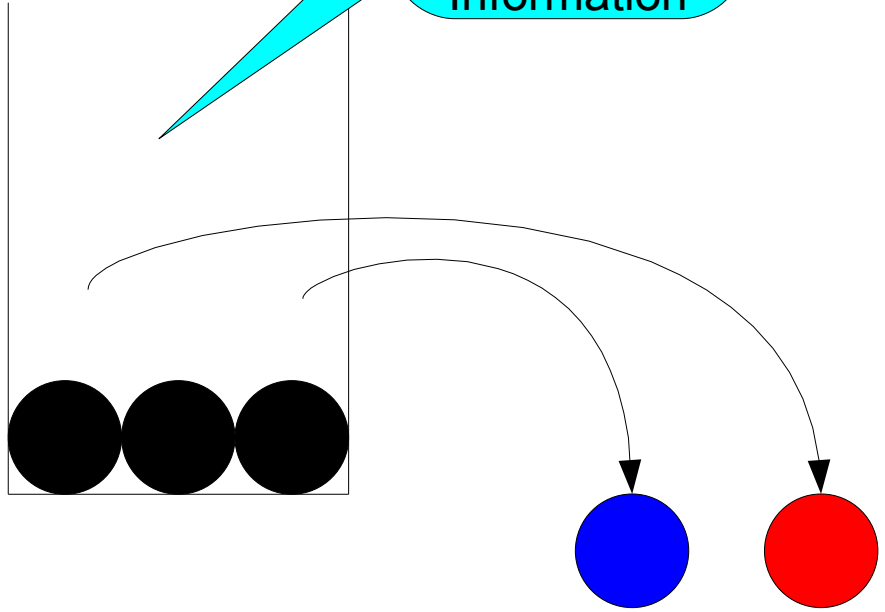
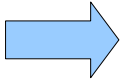
Ziel: nach 2 Ziehungen die Verteilung der Kugeln abzuschätzen.



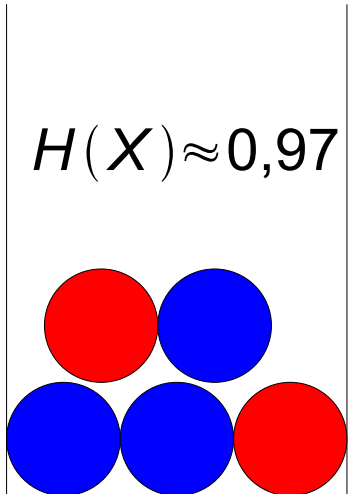
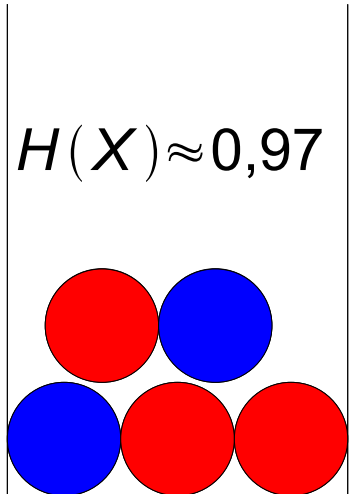
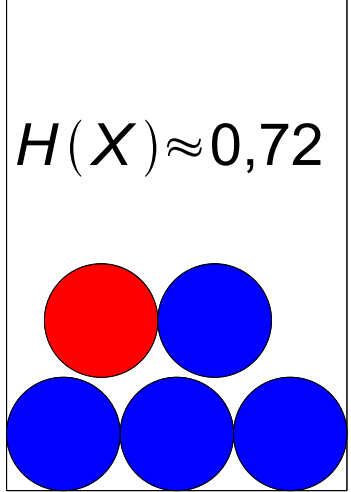
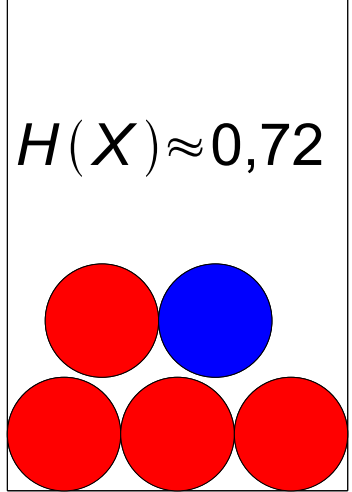
Beobachtung: die Urne enthält mindestens eine rote und eine blaue Kugel.

Maximale Entropie

Über diese drei Kandidaten keine Information



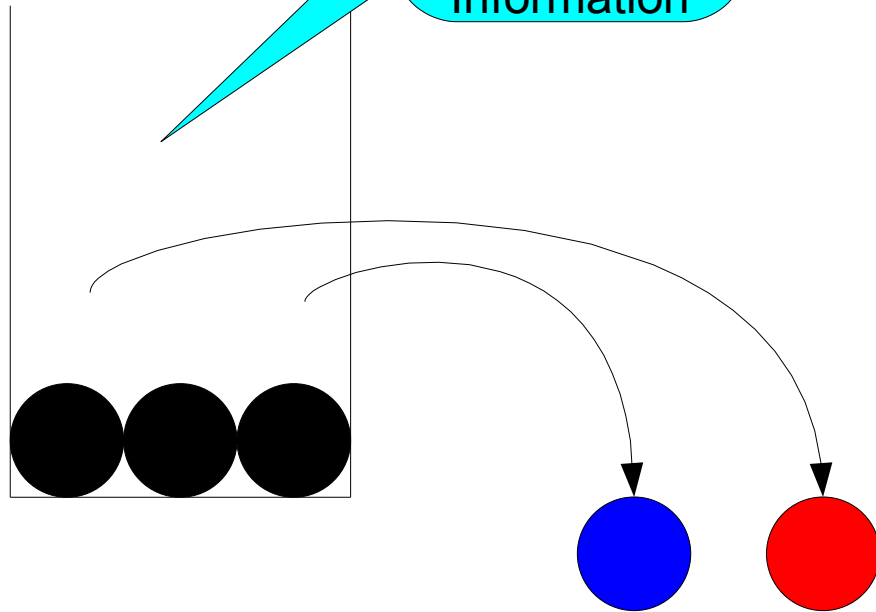
Verteilungen, die mit der Stichprobe vereinbart sind



$$H(X) = - \sum_{i=1}^{|\{rot, blau\}|} p_i \cdot \log_2(p_i)$$

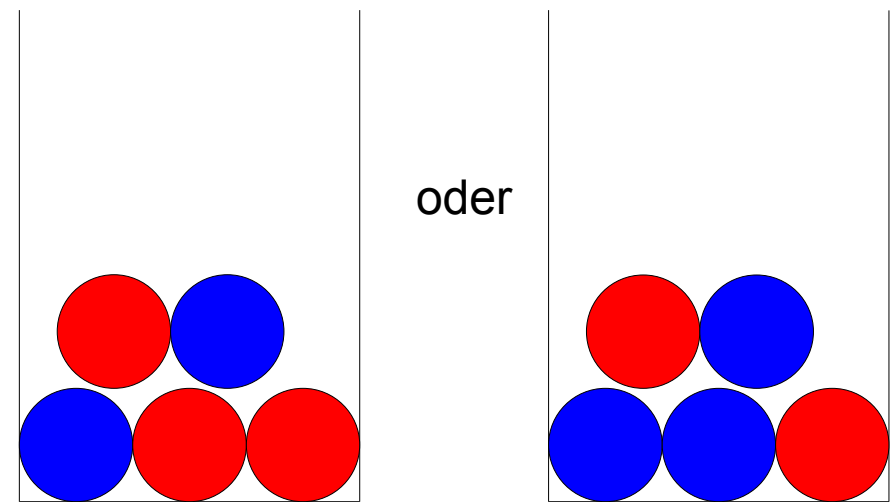
Maximale Entropie

Über diese drei Kandidaten keine Information



Schätzen wir die Verteilung in der Urne so ab, dass die Entropie H maximal wird.

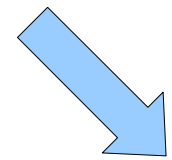
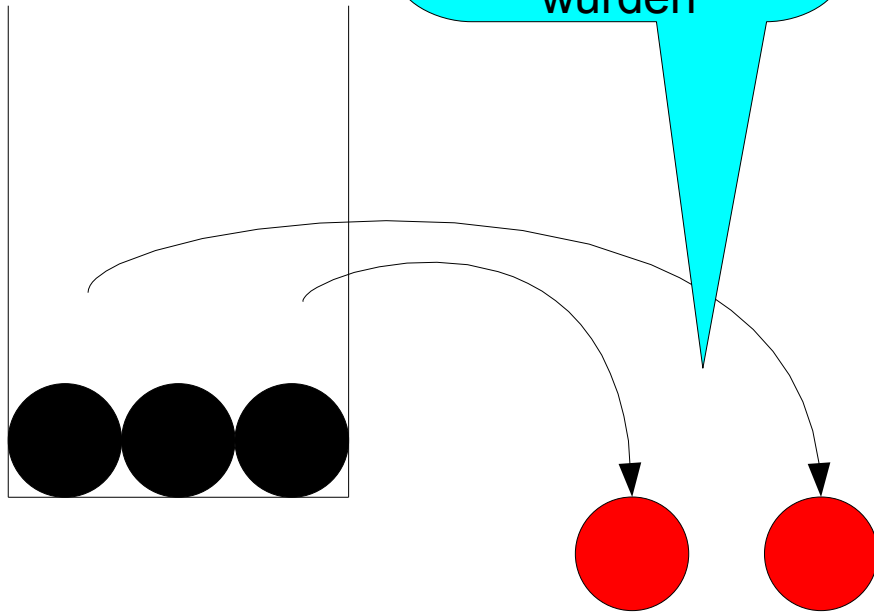
(möglichst) Gleichverteilung



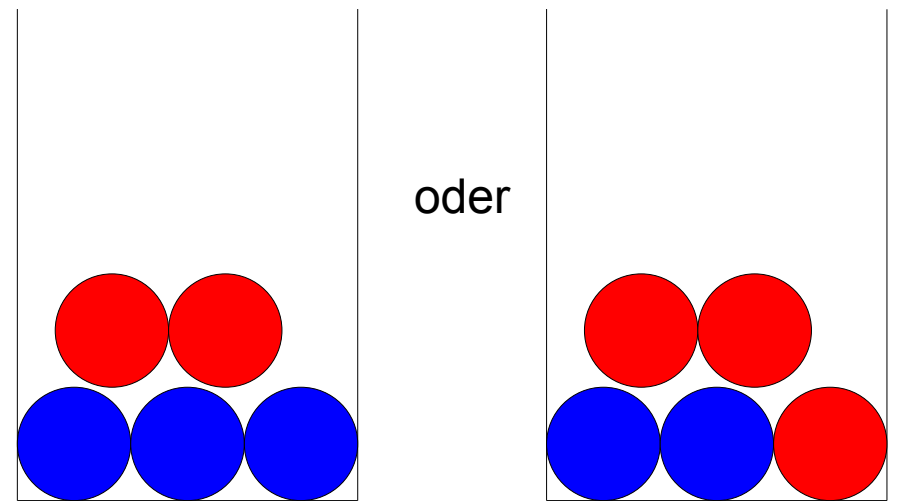
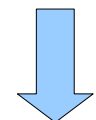
Verteilungen mit der maximalen Entropie

Maximale Entropie

Wie sieht die Verteilung aus, wenn 2 roten Kugeln gezogen wurden



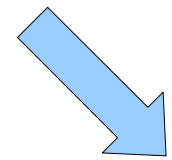
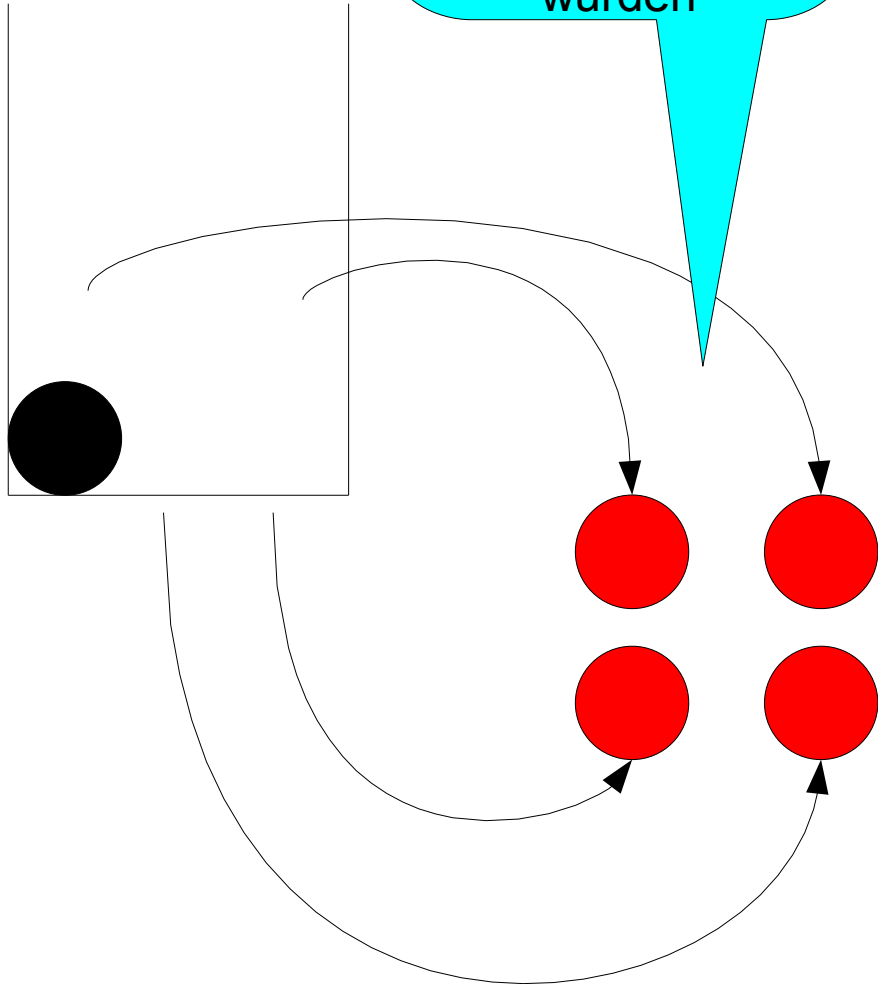
Wieder (möglichst) Gleichverteilung



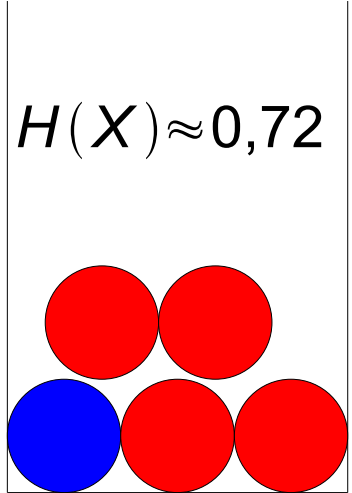
Verteilungen mit der maximalen Entropie

Maximale Entropie

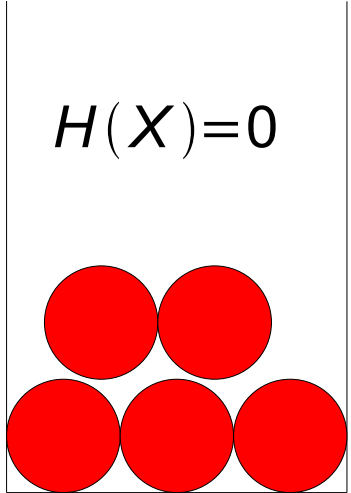
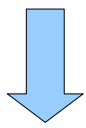
Wie sieht die Verteilung aus, wenn 4 roten Kugeln gezogen wurden



Wieder (möglichst) Gleichverteilung



maximale Entropie



keine maximale Entropie

Analogie mit der Urne und zurück zum MEMM

Die Urne enthält Relationen
der Form $(s^j, o^j) \forall j$

STICHPROBE



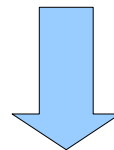
Zuordnung
per Hand
ausgeführt



$(s_1, \dots, s_n, o_1, \dots, o_n)$

$(s_1, \dots, s_m, o_1, \dots, o_m)$

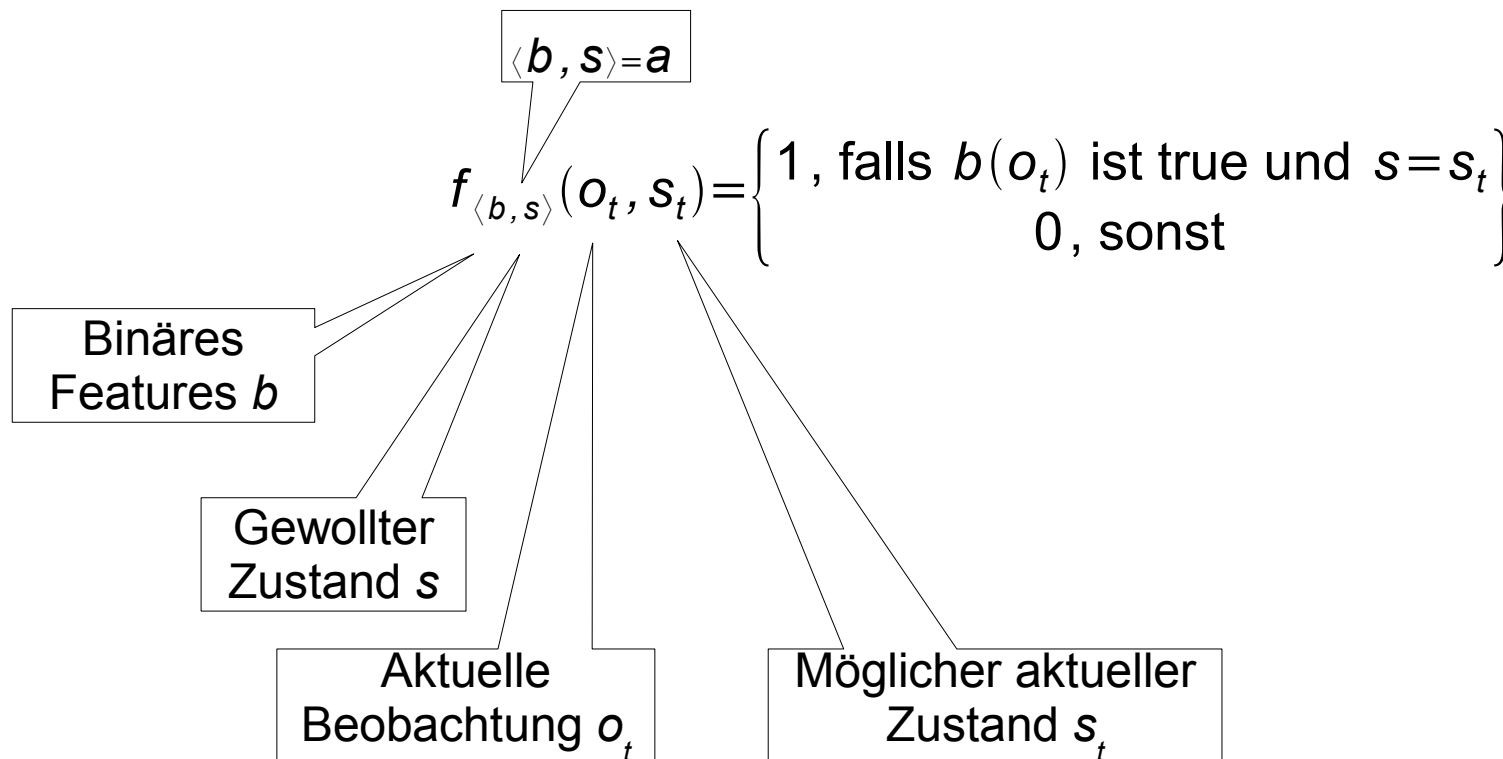
$(s_1, \dots, s_k, o_1, \dots, o_k)$



da man hier eine s-Folge und eine o-Folge hat, geht es hier um eine bedingte Entropie

Das binäre „Feature“

Stichprobe durch Charakteristische Funktion f kodieren, die einer Beobachtung „1“ oder „0“ zuordnet:



Beispiel

$O = \{she, is, by, me\}$ und $S = \{sie, ist, bei, mit, mir\}$

Stichprobe:

...she is by me... → ...sie ist *mit* mir...

...she is by me... → ...sie ist *bei* mir...

...she is by me... → ...sie ist *bei* mir...

...she is by me... → ...sie ist *mit* mir...

...she is by me... → ...sie ist *bei* mir...

$$f_{\langle b=\text{das Wort ist by}, s=\text{mit} \rangle}(o_t, s_t) = \begin{cases} 1, & \text{falls } o_t = \text{by} \text{ und die Übersetzung von } s_t \text{ ist } \textit{mit} \\ 0, & \text{sonst} \end{cases}$$

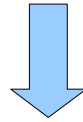
$$f_{\langle b=\text{das Wort ist by}, s=\text{bei} \rangle}(o_t, s_t) = \begin{cases} 1, & \text{falls } o_t = \text{by} \text{ und die Übersetzung von } s_t \text{ ist } \textit{bei} \\ 0, & \text{sonst} \end{cases}$$

$$f_{\langle b=\text{das Wort ist she}, s=\text{sie} \rangle}(o_t, s_t) = \begin{cases} 1, & \text{falls } o_t = \text{she} \text{ und die Übersetzung von } s_t \text{ ist } \textit{sie} \\ 0, & \text{sonst} \end{cases}$$

...

Bedingte Entropie

$$H_{s'}(X) = - \sum_s P_{s'}(s|o) \cdot \log_2(P_{s'}(s|o)) \quad \text{bei gegebenem } o$$



Ziel: Maximierung dieser Funktion unter den Nebenbedingungen:

$$1 = \sum_s P_{s'}(s|o) \quad \forall s', \text{ bei gegebenem } o$$

$$f_a(o, s_k) = \sum_s f_a(o, s) \cdot P_{s'}(s|o) \quad \forall s', s_k, a \text{ bei gegebenem } o$$

Vollständig
bzgl. P

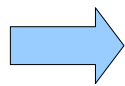
Vollständig
bzgl. f_a



LAGRANGE

Lagrangesche Multiplikatorenmethode zur Maximierung einer Funktion unter Nebenbedingungen

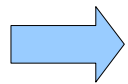
$$L_{s'}(\lambda, P_{s'}) = \underbrace{-\sum_s P_{s'}(s|o) \cdot \log_2(P_{s'}(s|o))}_{H_{s'}} + \sum_a \lambda_a \underbrace{\left(\sum_s f_a(o, s) \cdot P_{s'}(s|o) - f_a(o, s_k) \right)}_0$$



Erste Ableitung nach P bilden und gleich Null setzen:

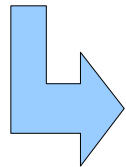
$$\frac{\partial L_{s'}(\lambda, P_{s'})}{\partial P_{s'}(s|o)} = \left(-1 - \log_2 P_{s'}(s|o) + \sum_s \lambda_a \cdot f_a(o, s) \cdot P_{s'}(s|o) \right) = 0$$

Ohne
Zwischenschritte



$$P_{s'}(s|o) = \frac{1}{\sum_s \exp\left(\sum_a \lambda_a f_a(o, s)\right)} \exp\left(\sum_a \lambda_a f_a(o, s)\right)$$

Verteilung mit der maximalen Entropie (Ohne Beweis):



$$P_{s'}(s|o) = \frac{1}{Z(o, s')} \exp\left(\sum_a \lambda_a f_a(o, s)\right)$$

Normalisierung, damit die Summe gleich 1 wird

Lagrange - Multiplikator, der mit dem GIS Verfahren berechnet wird

Exponentielle Darstellung der W'keit

$$Z(o, s') = \sum_s \exp\left(\sum_a \lambda_a f_a(o, s)\right)$$

Die Abschätzung der Zustandsfolge aus der Beobachtungsfolge(1)

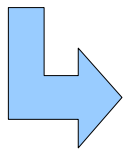
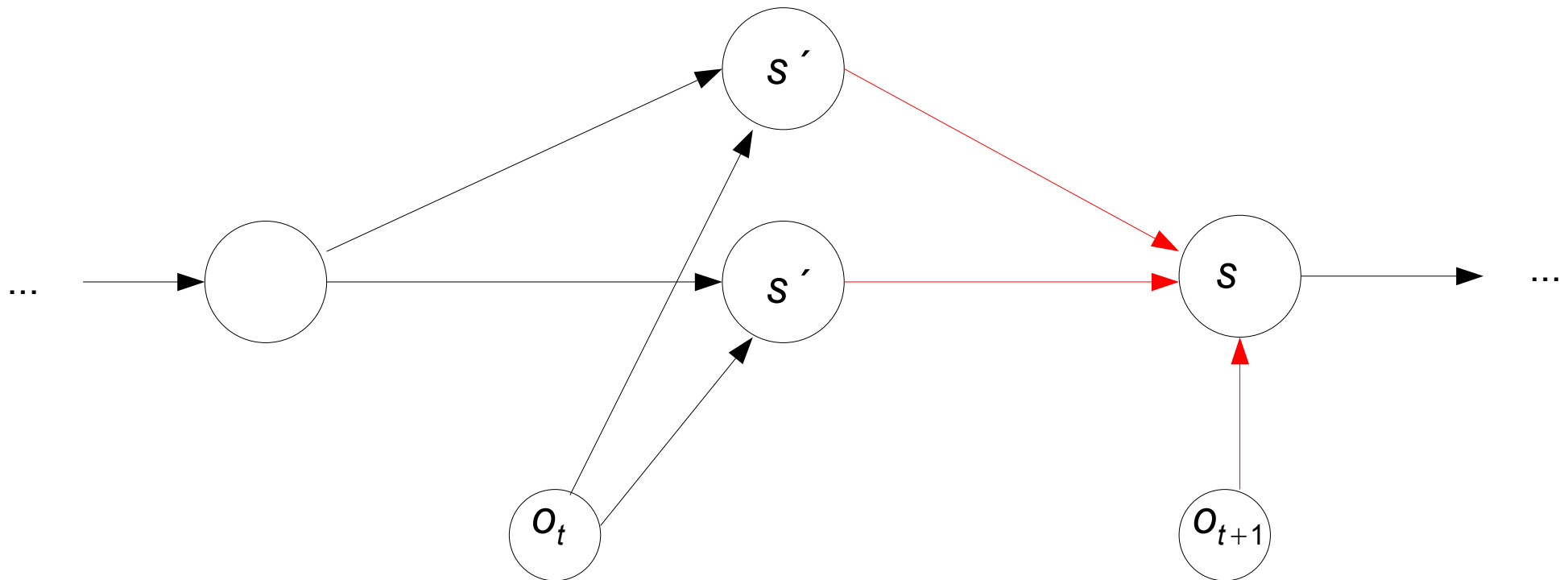
Läuft analog wie bei einem HMM:

$$\delta_{t+1}(\mathbf{s}) = \max_{\mathbf{s}' \in \mathcal{S}} \left\{ \delta_t(\mathbf{s}') \cdot P_{\mathbf{s}'}(\mathbf{s} | o_{t+1}) \right\} \quad \text{wobei}$$

$\delta_t(\mathbf{s}')$: Wahrscheinlichste Zustandsfolge bis zum o_t endet im \mathbf{s}'

Die Abschätzung der Zustandsfolge aus der Beobachtungsfolge(2)

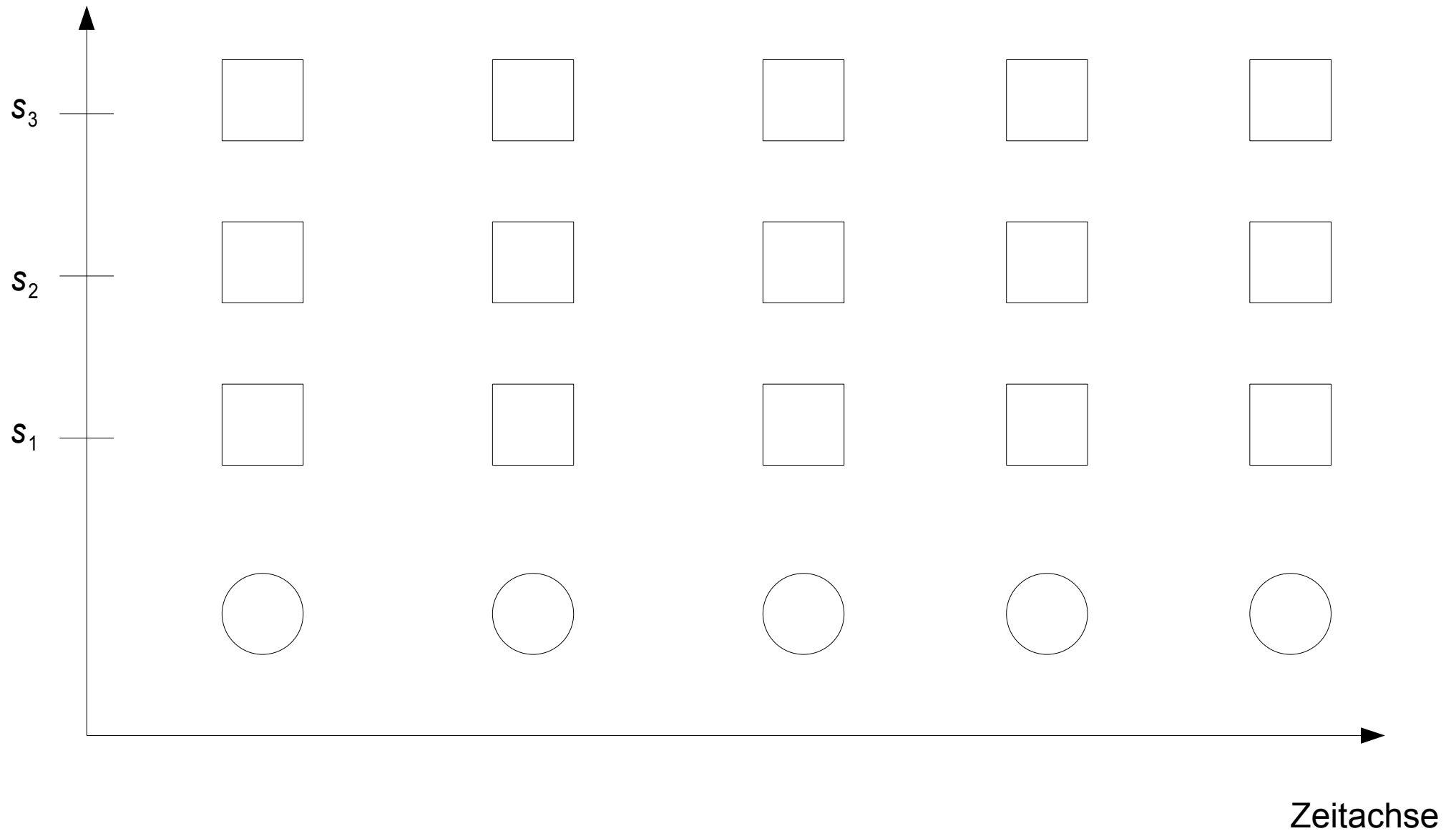
$$\delta_{t+1}(s) = \max_{s' \in S} \{ \delta_t(s') \cdot P_{s'}(s | o_{t+1}) \}$$



Viterbi Algorithmus für die Abschätzung der wahrscheinlichsten Zustandsfolge aus der Beobachtungsfolge.

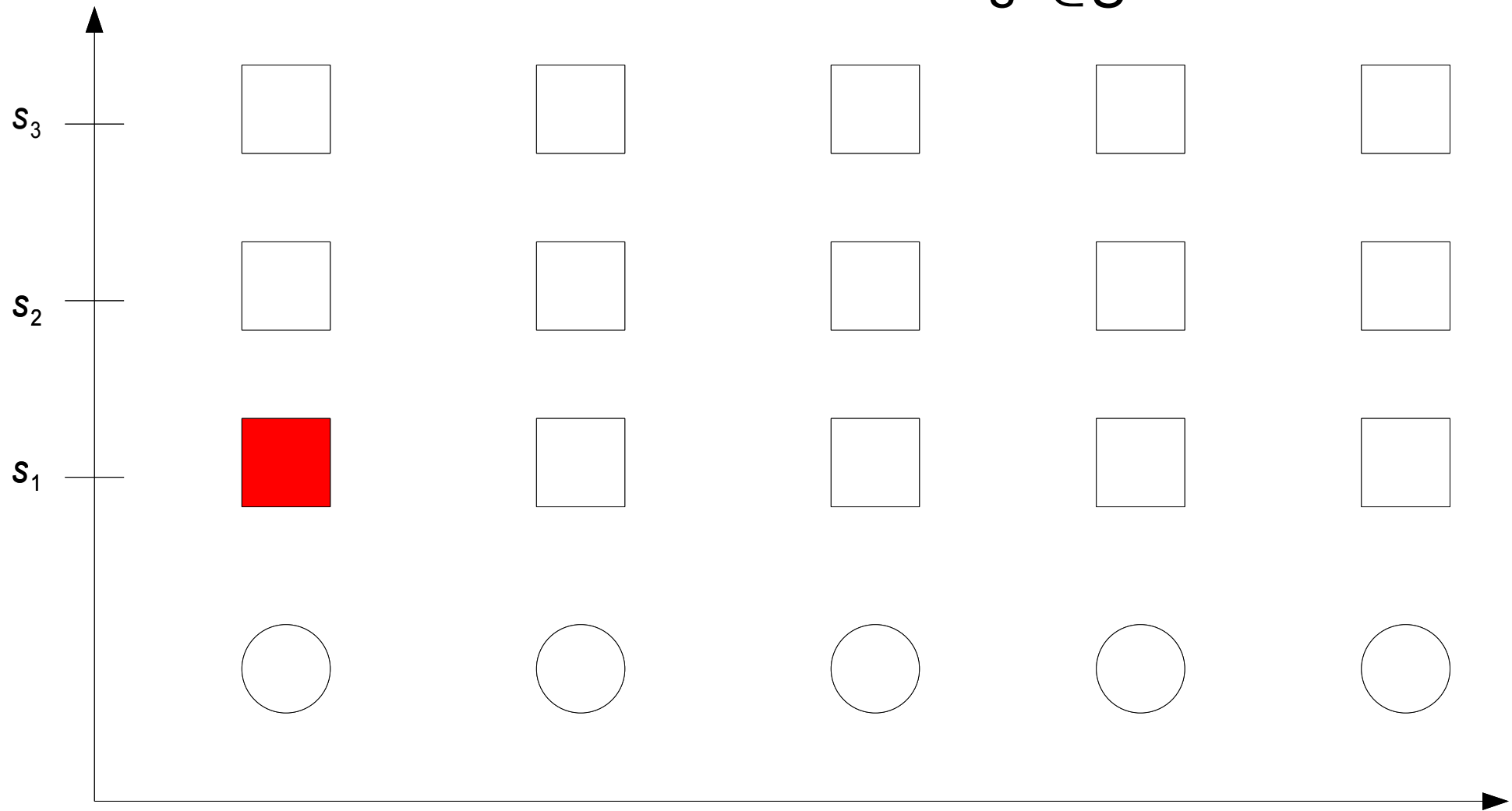
Viterbi Algorithmus

$$\delta_{t+1}(s) = \max_{s' \in S} \{ \delta_t(s') \cdot P_{s'}(s | o_{t+1}) \}$$



Viterbi Algorithmus

$$\delta_{t+1}(s) = \max_{s' \in S} \{ \delta_t(s') \cdot P_{s'}(s | o_{t+1}) \}$$

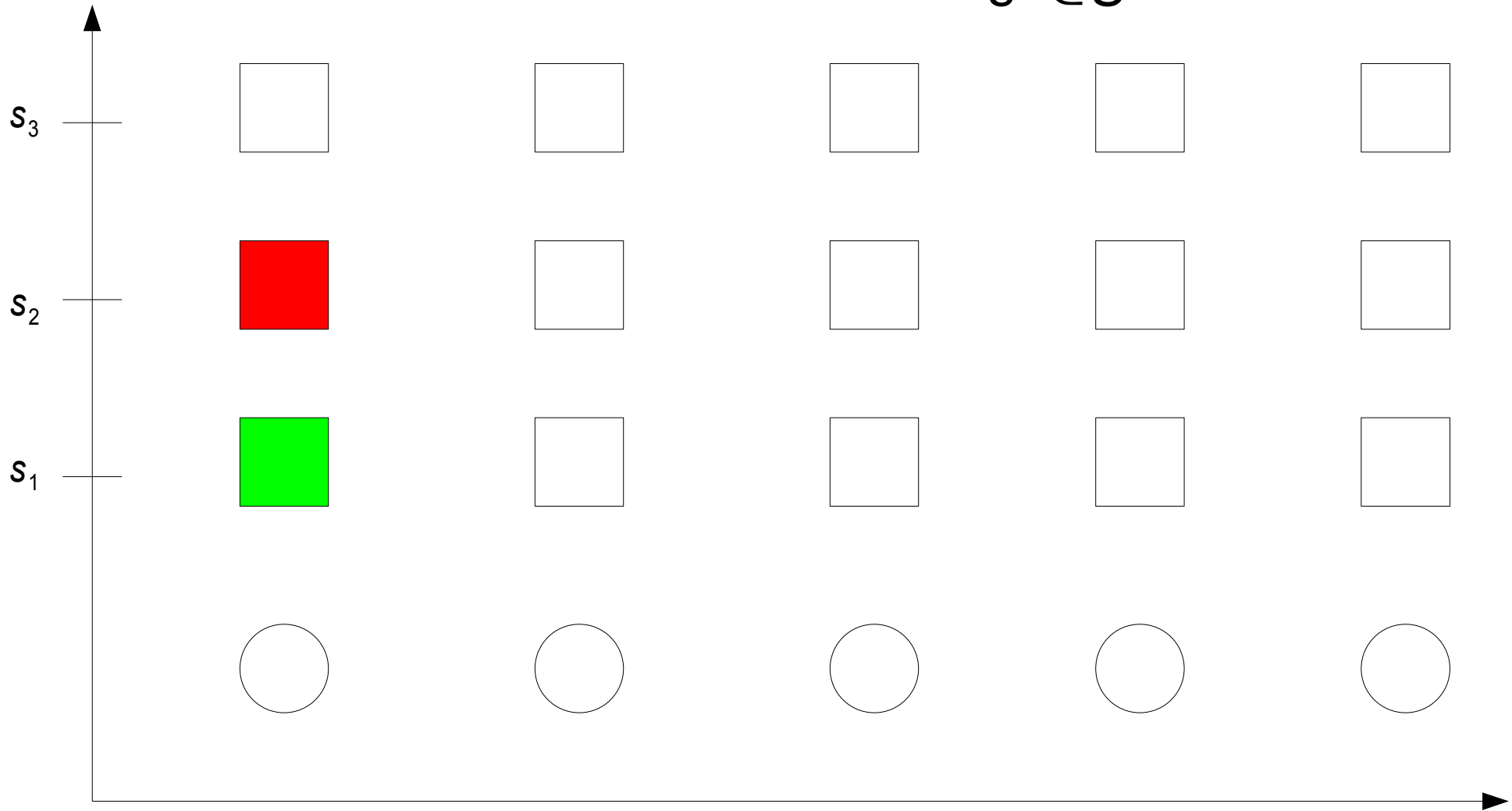


$$\delta_1(s_1) = \max_{s' \in S} \{ \delta_0(s') \cdot P_{s'}(s_1 | o_1) \}$$

Zeitachse

Viterbi Algorithmus

$$\delta_{t+1}(s) = \max_{s' \in S} \{ \delta_t(s') \cdot P_{s'}(s | o_{t+1}) \}$$

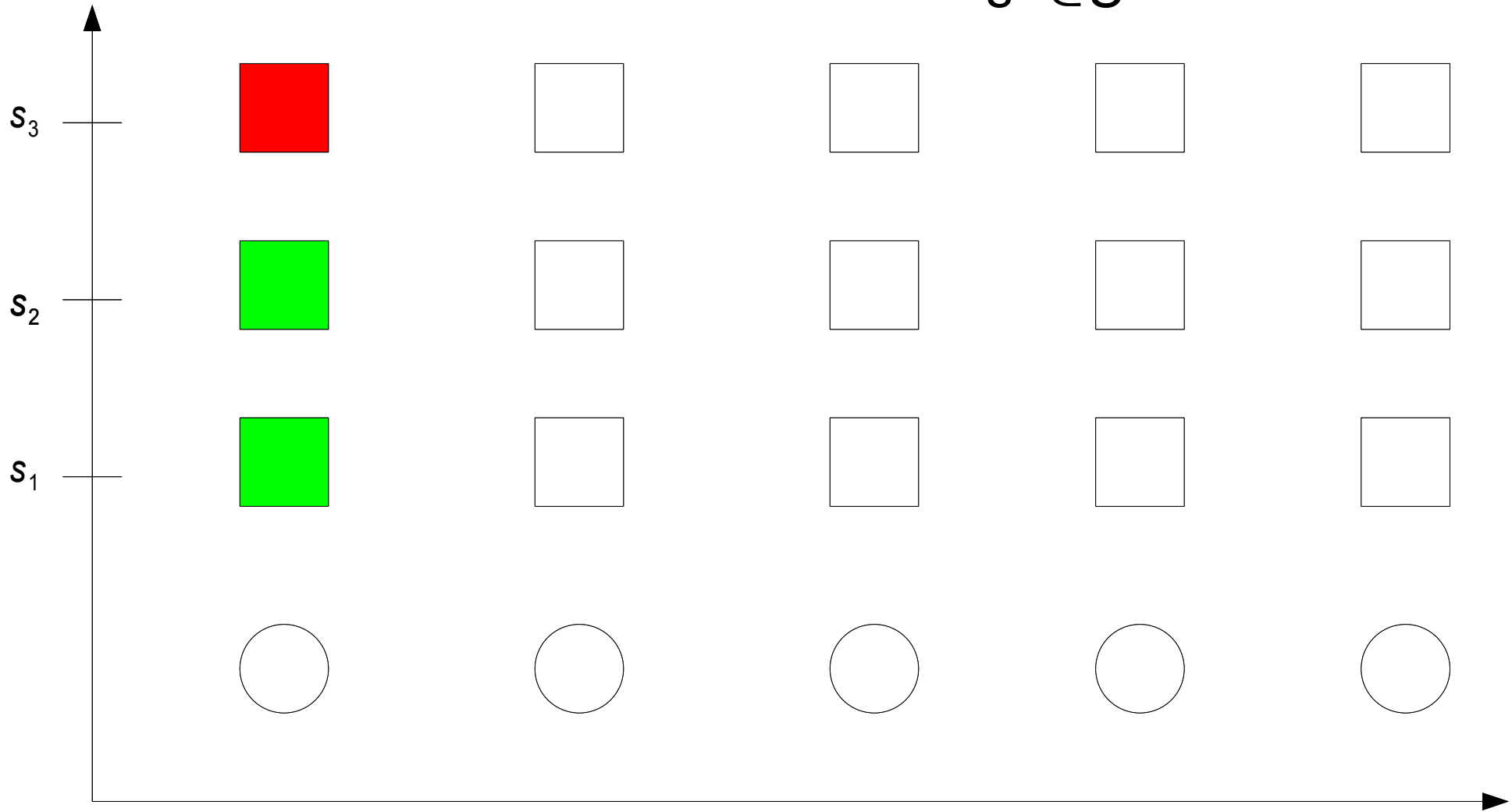


$$\delta_1(s_2) = \max_{s' \in S} \{ \delta_0(s') \cdot P_{s'}(s_2 | o_1) \}$$

Zeitachse

Viterbi Algorithmus

$$\delta_{t+1}(s) = \max_{s' \in S} \{ \delta_t(s') \cdot P_{s'}(s | o_{t+1}) \}$$

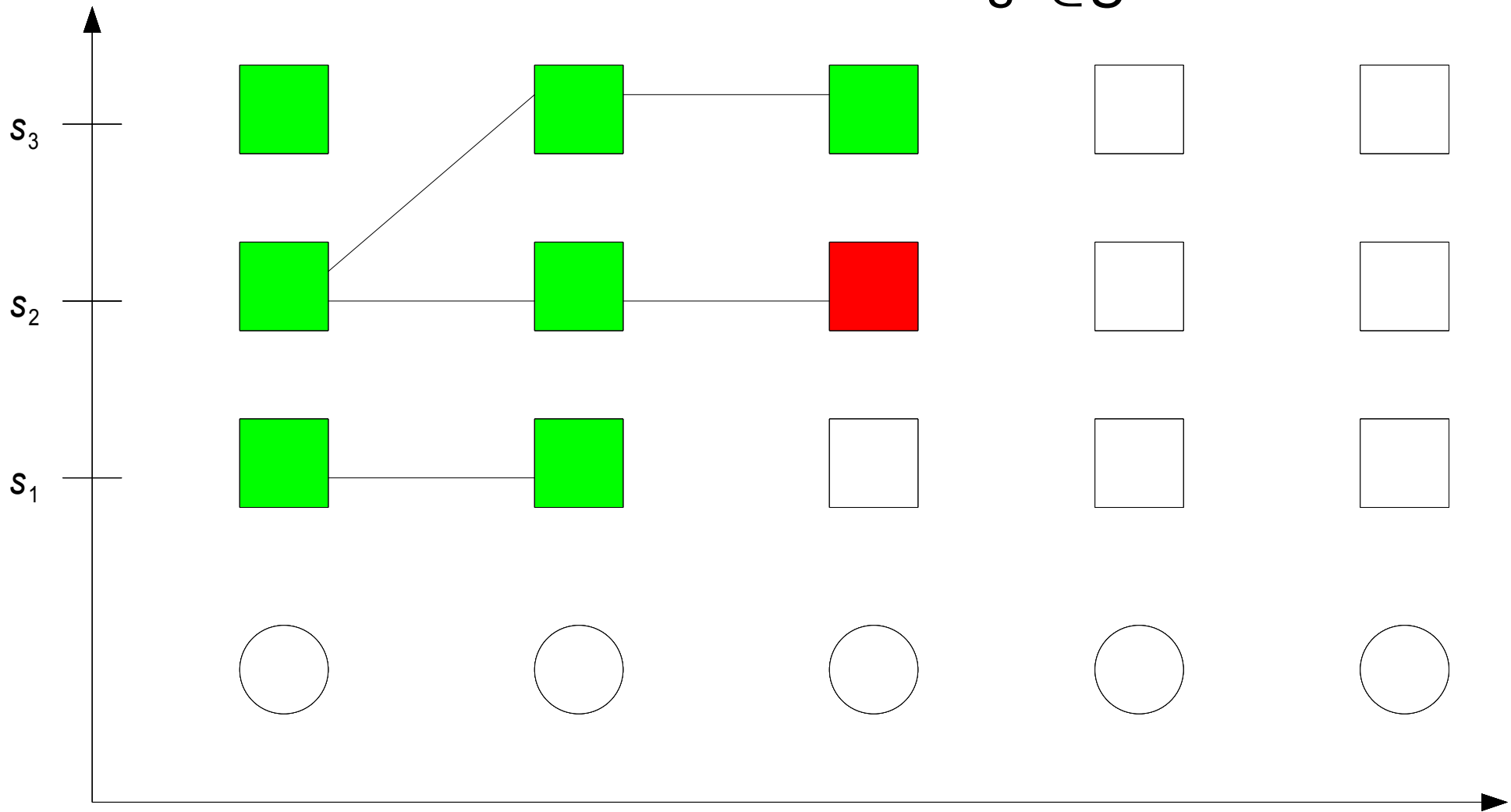


$$\delta_1(s_3) = \max_{s' \in S} \{ \delta_0(s') \cdot P_{s'}(s_3 | o_1) \}$$

Zeitachse

Viterbi Algorithmus

$$\delta_{t+1}(s) = \max_{s' \in S} \{ \delta_t(s') \cdot P_{s'}(s | o_{t+1}) \}$$

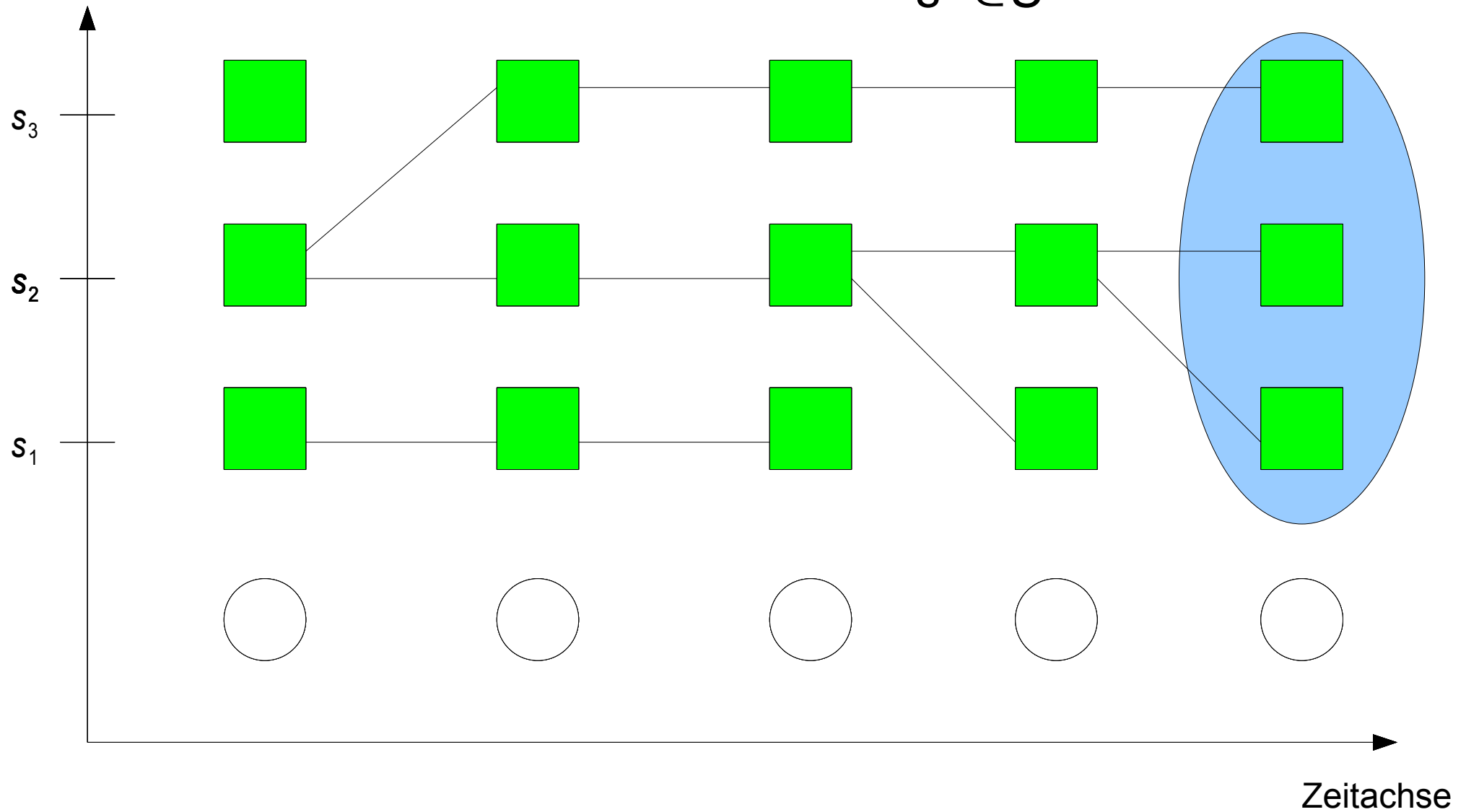


$$\delta_3(s_2) = \max_{s' \in S} \{ \delta_2(s') \cdot P_{s'}(s_2 | o_3) \}$$

Zeitachse

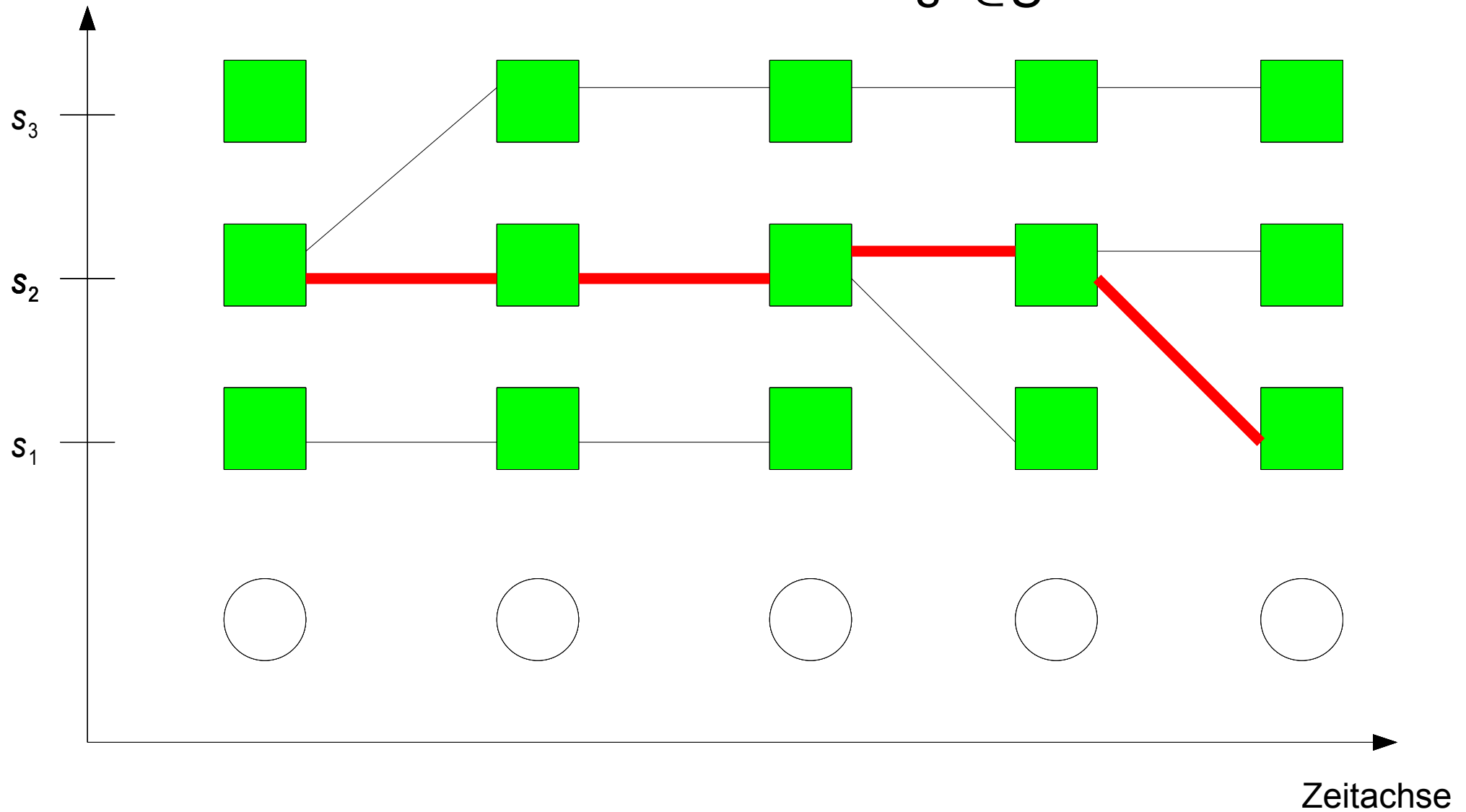
Viterbi Algorithmus

$$\delta_{t+1}(s) = \max_{s' \in S} \{ \delta_t(s') \cdot P_{s'}(s | o_{t+1}) \}$$



Viterbi Algorithmus

$$\delta_{t+1}(s) = \max_{s' \in S} \{ \delta_t(s') \cdot P_{s'}(s | o_{t+1}) \}$$

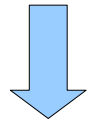


Wie gut ist MEMM?

Ist schon viel besser als das HMM.

Mögliche Verbesserung vom MEMM:

- jedes Element in der Zustandsfolge sollte von der kompletten Beobachtungsfolge abhängen.



Conditional Random Fields (nächster Vortrag)