

12. 系列データの識別

12.1 ラベル系列に対する識別

- ラベル系列に対する識別問題の分類
 - 入力の系列長と出力の系列長が等しい
 - 例) 形態素解析処理、固有表現抽出
 - 系列ラベリング問題 \Rightarrow CRF
 - 入力の系列長に関わらず出力の系列長が 1
 - 例) 動画像のラベル付け、単語音声認識
 - 系列識別問題 \Rightarrow HMM
 - 入力の系列長と出力の系列長に対応関係がない
 - 例) 連続音声認識
 - 系列識別と探索を組み合わせた複雑な処理

12.2 系列ラベリング問題— CRF—

- 系列ラベリング問題とは
 - 入力系列の個々の要素に対して、識別によるラベリングを行う問題
 - 入力系列の要素の出現確率は、前後の要素と独立ではないことが多いので、1 入力 1 出力の識別器を連続的に適用する方法では、性能が上がらない
 - ⇒ 入力や出力の系列としての特徴を使う
 - 可能な出力系列の組合せは膨大な数になるので、単純な事後確率最大法は使えない
 - ⇒ 探索によって最適解を求める

12.2 系列ラベリング問題— CRF—

- 系列ラベリング問題の事例

- 形態素解析

入力	系列	で	入力	さ	れる	各	要素
出力	名詞	助詞	名詞	動詞	接尾辞	接頭辞	名詞

- 固有表現抽出（人を指す表現の抽出例）

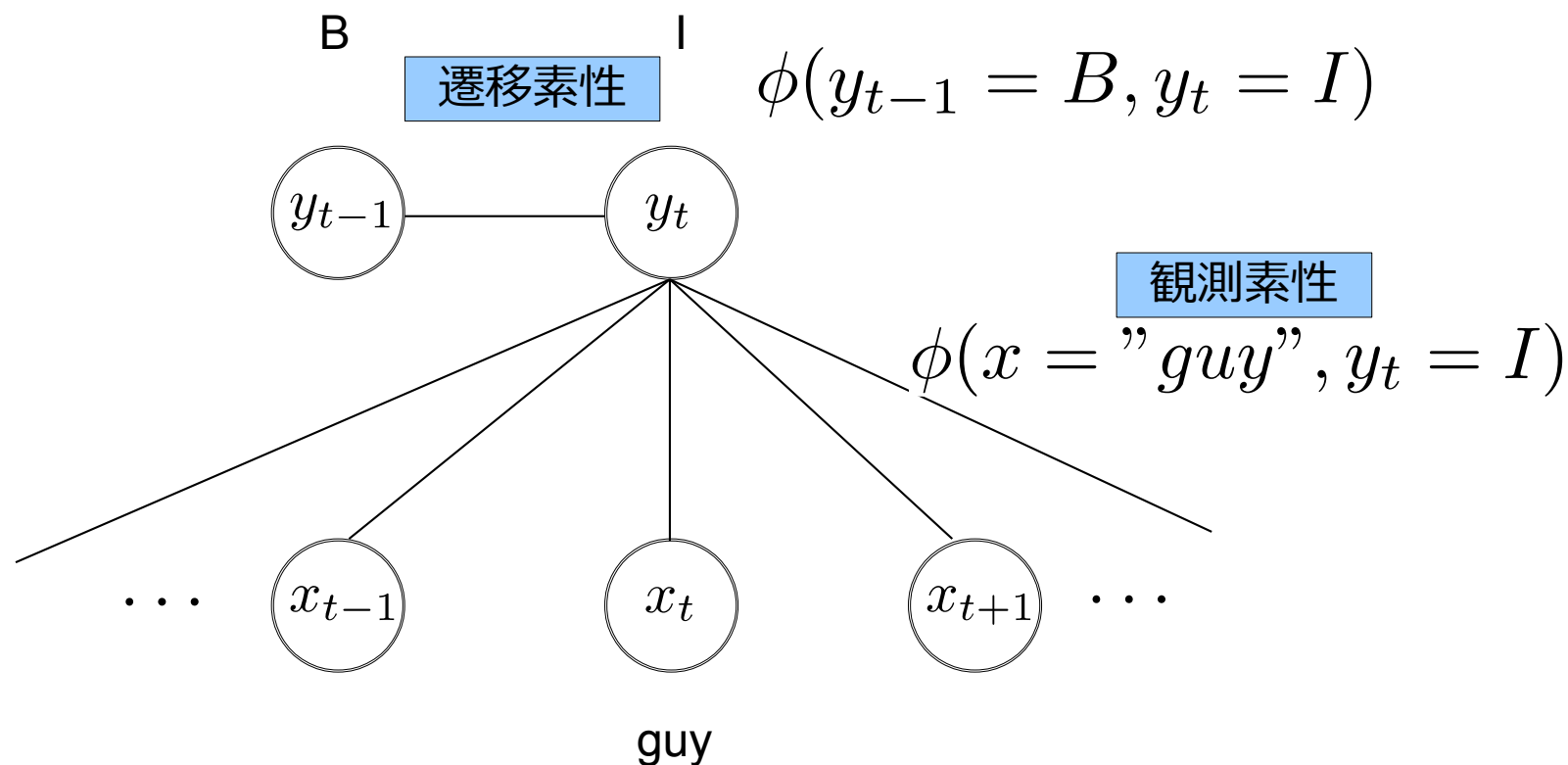
入力	Suddenly, the tall German guy talked to me							
出力	O	B	I	I	I	O	O	B

B: Begin
I: inside
O: outside

12.2 系列ラベリング問題— CRF—

- 対数線型モデルによる系列ラベリング

- 素性関数の導入



12.2 系列ラベリング問題— CRF—

- 対数線型モデル

$$P(\boldsymbol{y}|\boldsymbol{x}) = \frac{1}{Z_{\boldsymbol{x},\boldsymbol{w}}} \exp(\boldsymbol{w} \cdot \boldsymbol{\phi}(\boldsymbol{x}, \boldsymbol{y}))$$

- 出力の決定

$$\begin{aligned} \boldsymbol{y}^* &= \arg \max_{\boldsymbol{y}} P(\boldsymbol{y}|\boldsymbol{x}) \\ &= \arg \max_{\boldsymbol{y}} \frac{1}{Z_{\boldsymbol{x},\boldsymbol{w}}} \exp(\boldsymbol{w} \cdot \boldsymbol{\phi}(\boldsymbol{x}, \boldsymbol{y})) \\ &= \arg \max_{\boldsymbol{y}} \boldsymbol{w} \cdot \boldsymbol{\phi}(\boldsymbol{x}, \boldsymbol{y}) \end{aligned}$$

12.2 系列ラベリング問題— CRF—

- 素性関数の制限

$$\mathbf{y}^* = \arg \max_{\mathbf{y}} \sum_t \mathbf{w} \cdot \phi(\mathbf{x}, y_t, y_{t-1})$$

- ビタビアルゴリズムによって探索が可能

Algorithm 12.1 ビタビアルゴリズム

for $t = 2$ to $|x|$ do

 for all y_t do

$$\alpha(t, y_t) = \max_{y_{t-1}} \{ \mathbf{w} \cdot \phi(\mathbf{x}, y_t, y_{t-1}) + \alpha(t-1, y_{t-1}) \}$$

$$B(t, y_t) = \arg \max_{y_{t-1}} \{ \mathbf{w} \cdot \phi(\mathbf{x}, y_t, y_{t-1}) + \alpha(t-1, y_{t-1}) \}$$

 end for

end for

$\mathbf{y}^* = \alpha$ の最大値に対応する B を逆に辿る

12.3 系列識別問題— HMM—

- 例題

- PC 操作系列による熟練度の判定

- k: キーボード、 g: マウス、 e: エラー
- 初心者の入力系列例

k e k g k e k g g k g k k e g e e k e e e g e

- 熟練者の入力系列例

k k e k g k k k e k g k g g g e g k g

- 判定したい入力系列

k g e k g k k g e k g e k e e k e g e k

12.3 系列識別問題— HMM—

- 生成モデルによるアプローチ
 - 系列識別問題ではクラスの事前確率を得られることが多い

$$y^* = \arg \max_y P(y|\mathbf{x})$$

$$= \arg \max_y \frac{P(\mathbf{x}, y)}{P(\mathbf{x})}$$

$$= \arg \max_y \frac{P(\mathbf{x}|y)P(y)}{P(\mathbf{x})}$$

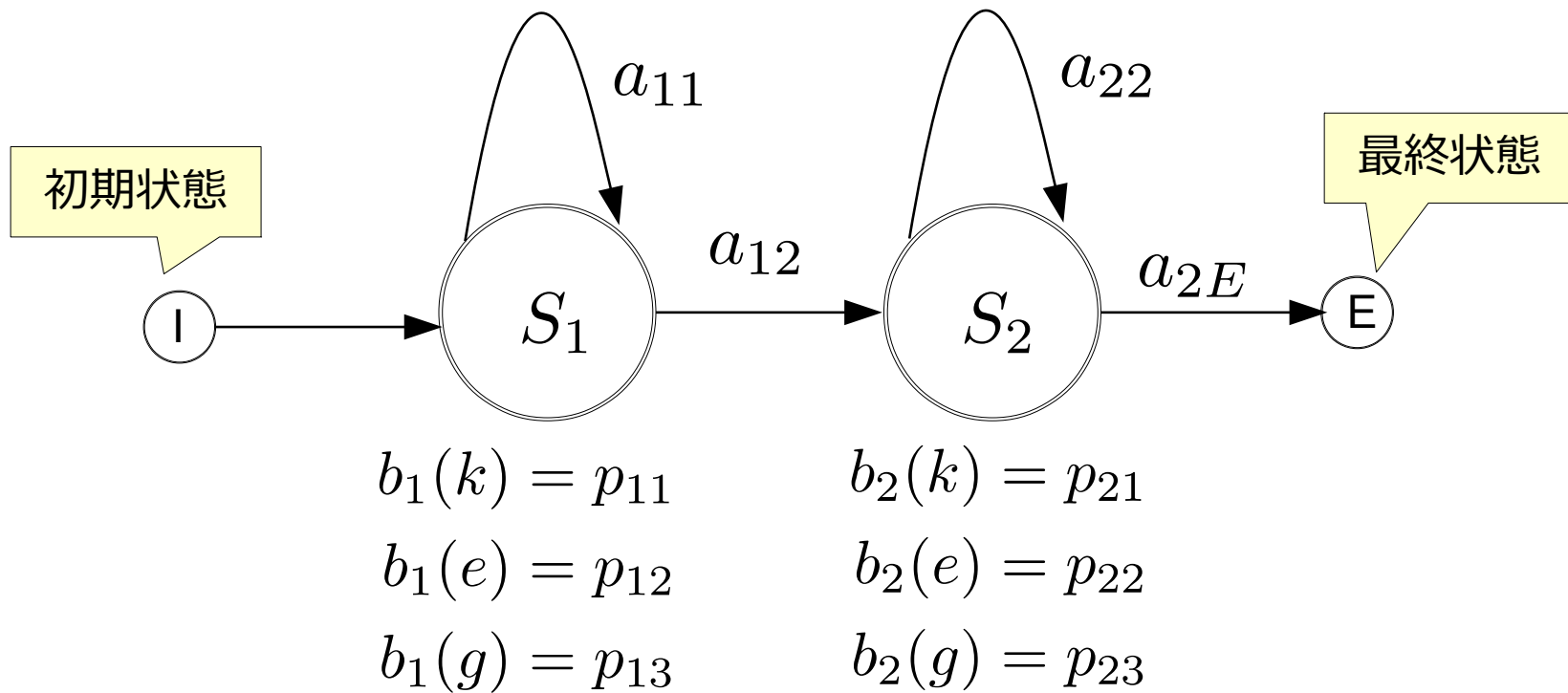
$$= \arg \max_y P(\mathbf{x}|y)P(y)$$

尤度

事前確率

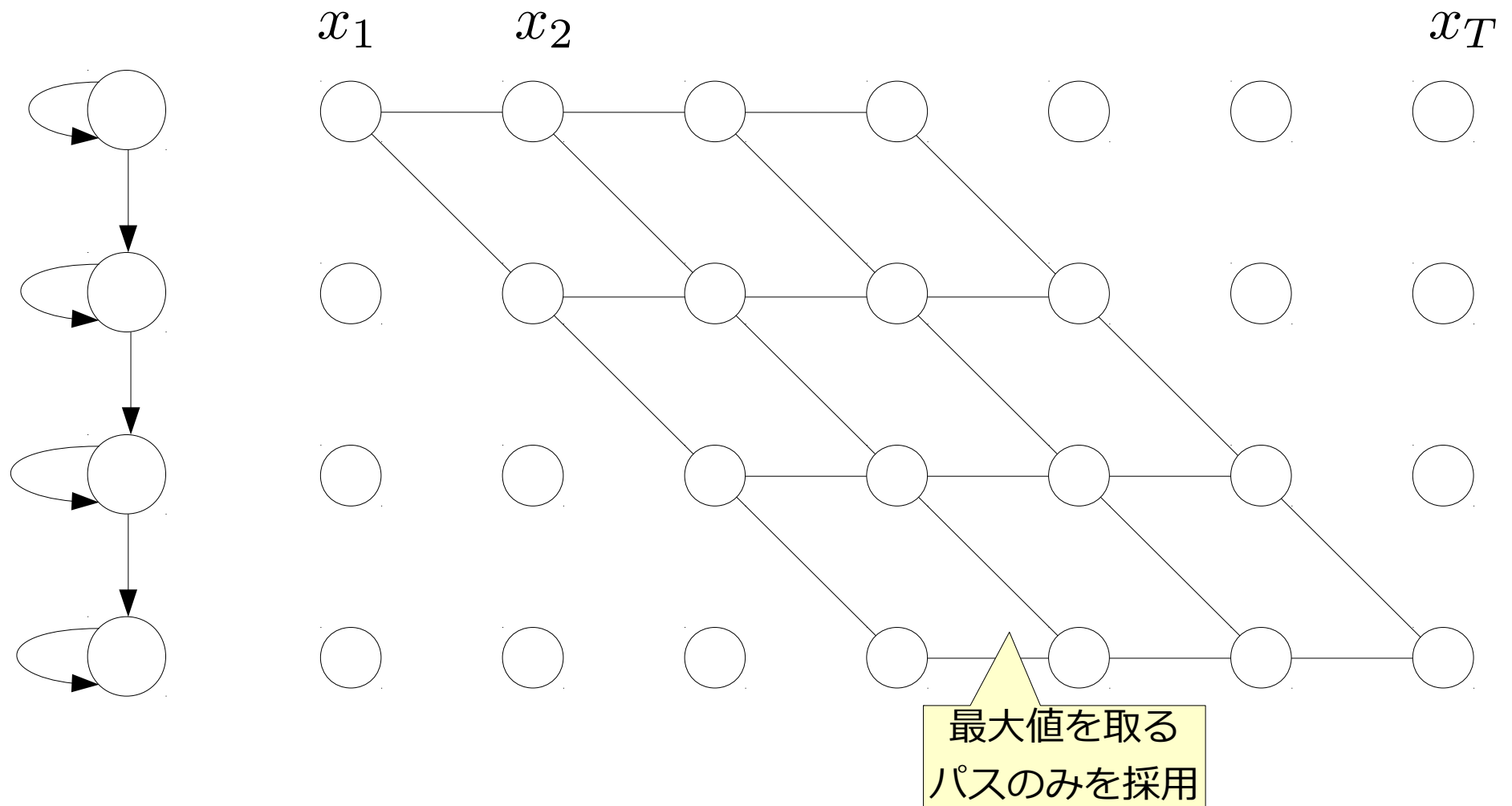
12.3 系列識別問題—HMM—

- 不定長入力に対する尤度計算法
 - 自己遷移を持つ確率オートマトンを用いる



12.3 系列識別問題—HMM—

- ビタビアルゴリズムを用いた探索



12.3 系列識別問題—HMM—

- HMM の学習 : EM アルゴリズム

