

言語処理における識別モデルの発展

-- HMMからCRFまで --

坪井祐太, 鹿島久嗣 (IBM Research)
工藤 拓 (Google)

今日話した内容

- HMMの話
- 復号(Viterbi)
- 内積で書けるよ
- 内積で書ける別のモデル CRF, HM Perceptron, HM SVM

- 違いはパラメータ推定法
- Generative Model
 - HMM
 - HMM課題

- Discriminative Modelのメリット: 直接 $P(Y|X)$, 素性間の重なりを考慮できる

- Conditional Model
 - CRF
 -
- Discriminative Model
 - HM Perceptron
 - メリット:
- Generative Model vs Discriminative のまとめた表

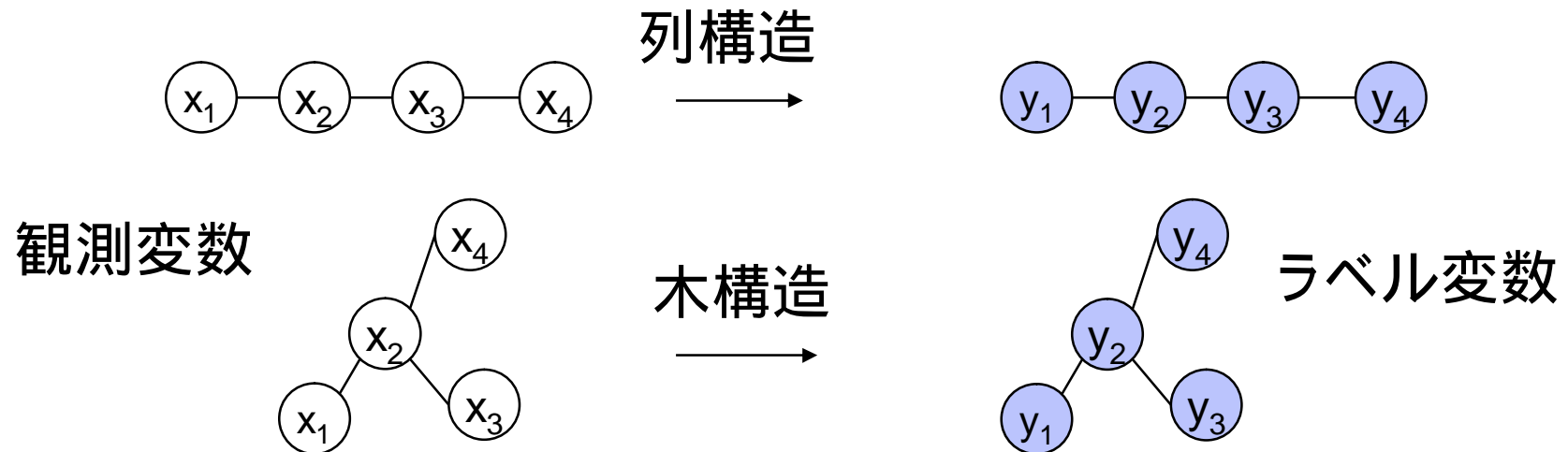
- 形態素解析の例

チュートリアルの流れ

- 構造のラベル付け問題の定義
- 2つのアプローチ: 生成モデル vs. 識別モデル
 - 生成モデル: 隠れマルコフモデル
 - 識別モデル: 条件付確率場
- そのほかの識別モデル
- 計算機実験による性能比較

構造ラベル付与問題とは？

- 観測されたデータ構造 x に対応するラベル構造 y への写像を学習する問題

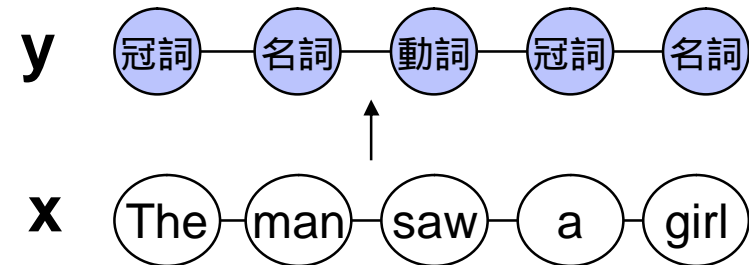


- 構造を持つ観測データの各ノードにラベルを付与する問題として捉えることができる応用が多く存在する。
 - 例：自然言語処理、バイオインフォマティクス

自然言語処理での構造ラベル付与問題の例(列構造)

- 品詞タグ付与タスク

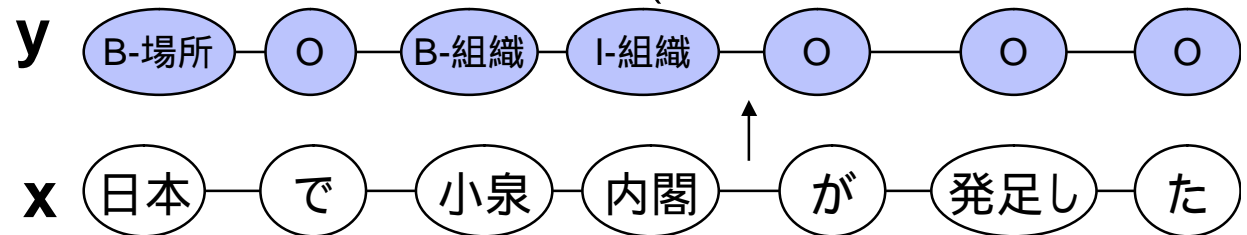
単語列に対して品詞ラベルを付与するタスク



- 固有表現抽出タスク

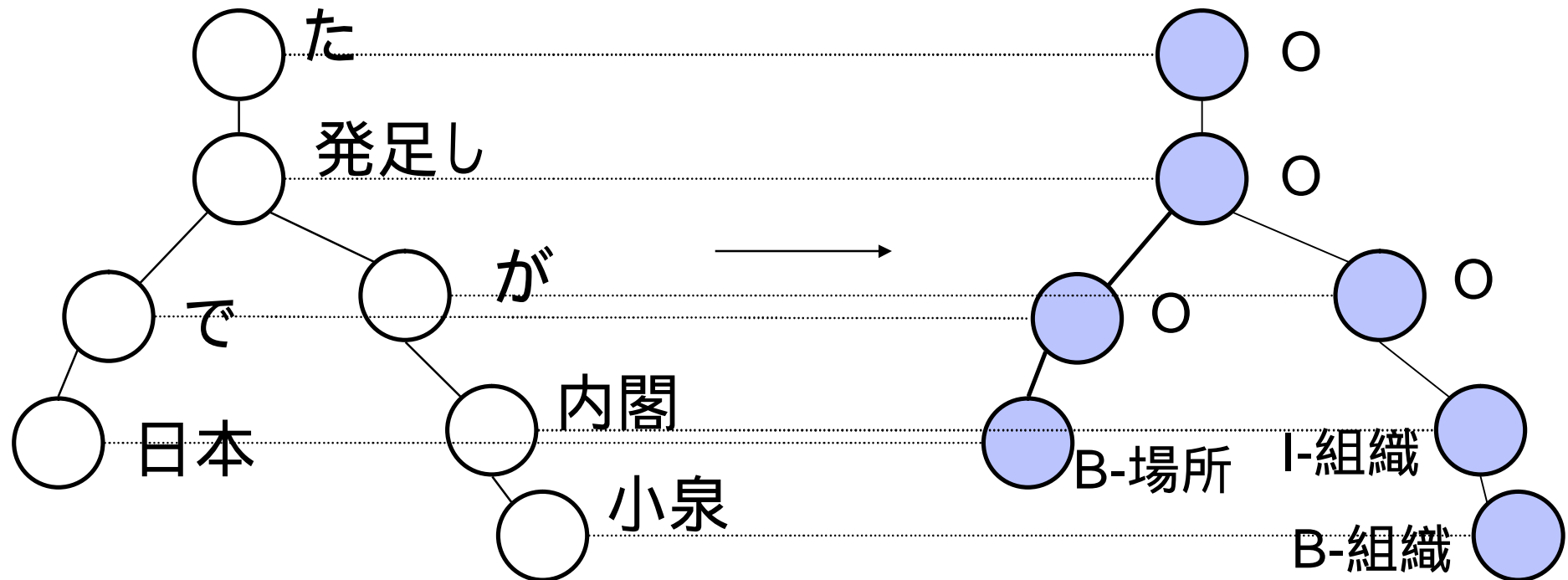
人名・組織名等の固有表現をテキスト中から抽出

単語列に対して固有表現の始まり(B-XXX)と続く固有表現(I-XXX)を示すラベルを付与するタスク(Oは固有表現以外)



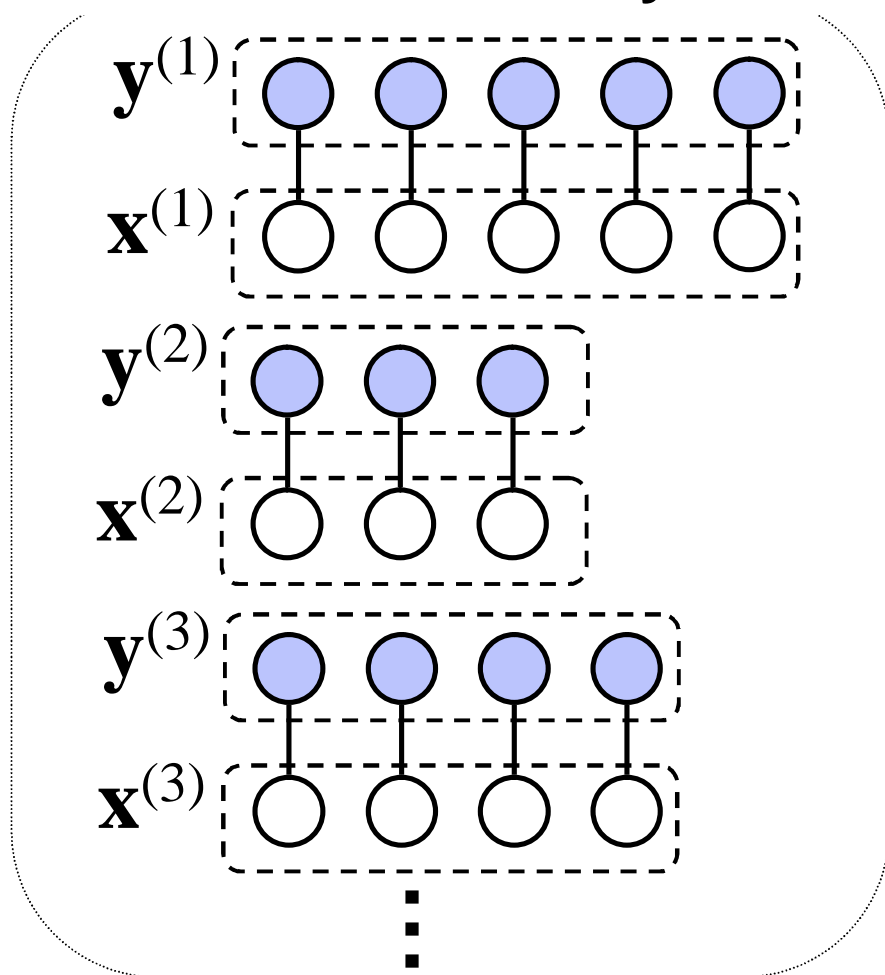
自然言語処理での構造ラベル付与問題の例(木構造)

- 係り受け木に対する固有表現抽出タスク
 - 係り受け解析によって生成された単語間の関係を表す係り受け木に対して、ラベルを付与。
 - 主語と述語の関係など言語構造を考慮したラベル付与

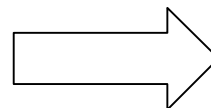


教師付き学習による構造ラベル付与問題

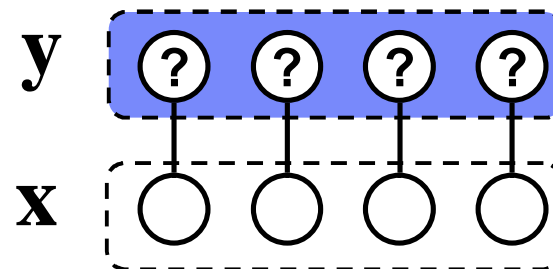
学習データ (正しい x, y ペア)



パラメータ推定
(学習)



復号
(予測)

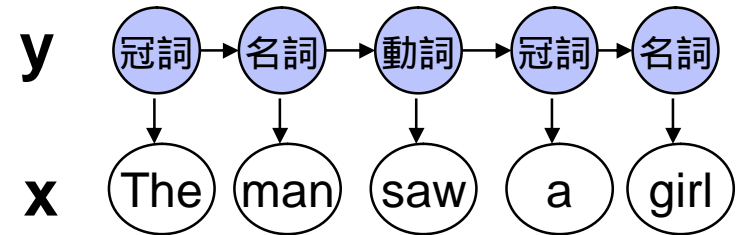


ラベル無しデータ

構造ラベル付与学習モデル

- 生成モデルに基づく手法
 - 隠れマルコフモデル(HMM)
- 識別モデルに基づく手法
 - 条件付確率場(CRF)

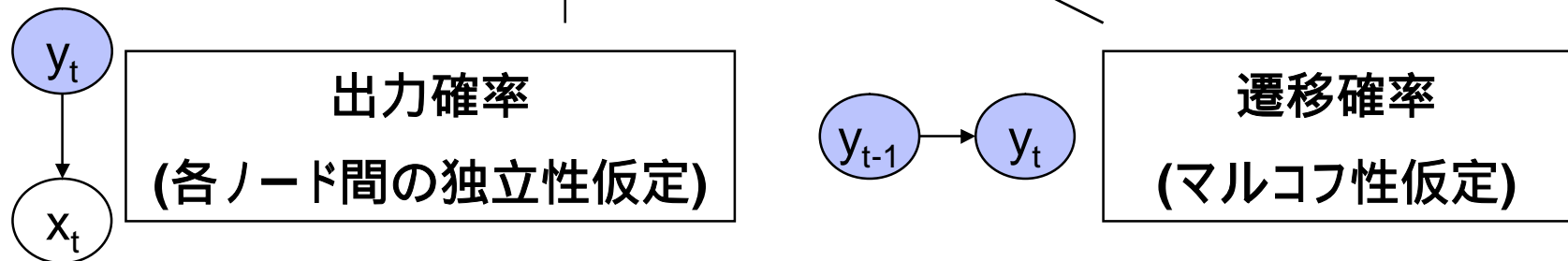
生成モデルによる構造ラベル付与学習
隠れマルコフモデル(HMM)



- **xとyの同時分布に基づくモデル**
- **生成確率を出力確率と遷移確率に分解してモデル化**

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

$$= \prod_{t=1}^T P(x_t / y_t) P(y_t | y_{t-1}) \quad (T \text{は構造} \mathbf{x}, \mathbf{y} \text{のサイズ})$$



隠れマルコフモデル(HMM):
ラベル列の予測 (復号問題)

predict

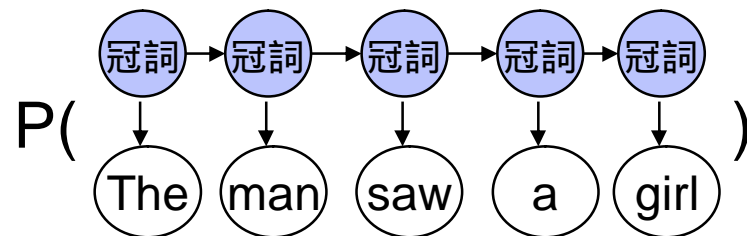
$$\mathbf{y} = \underset{\mathbf{y}}{\operatorname{argmax}} P(\mathbf{x} | \mathbf{y}) = \underset{\mathbf{y}}{\operatorname{argmax}} P(\mathbf{x}, \mathbf{y})$$

$$= \underset{\mathbf{y}}{\operatorname{argmax}} \prod_{t=1}^T P(x_t / y_t) P(y_t | y_{t-1})$$

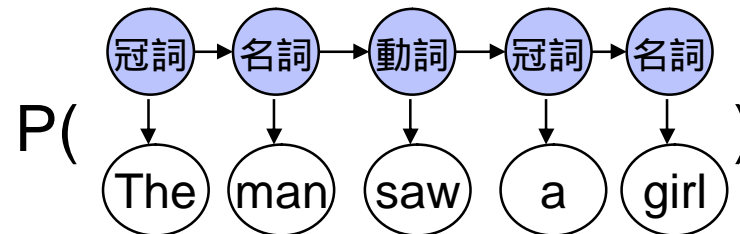
ある \mathbf{x} が与えられたときに、確率が最大になるラベル列 \mathbf{y} を見つきたい。

ありうるラベル列全て ($|Y|^T$ 個) の列挙は非効率 ($|Y|$ は目的ラベル集合のサイズ)

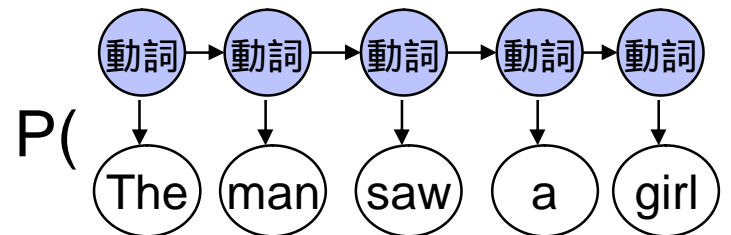
全部で $|Y|^T$ 個



...



...



隠れマルコフモデル(HMM):

Viterbi 復号法による最適ラベル列の求め方

- 位置tまでのラベル列の確率が最大になる値を記憶するテーブルを使い再計算を避けることで、最適なラベル列を効率的に計算

$$\delta_1(y) = P(x_1 | y)$$

$$\delta_t(y) = \max_{\tilde{y} \in Y} \delta_{t-1}(\tilde{y}) P(x_t | y) P(y | \tilde{y})$$

英語品詞タグ付けタスクでの

$y \backslash x$	the	man	saw	...
冠詞	$P(\text{the} \text{冠詞})$	$\max_{\tilde{y} \in Y} \delta_1(\tilde{y}) \times P(\text{man} \text{冠詞}) P(\text{冠詞} \tilde{y})$	$\max_{\tilde{y} \in Y} \delta_2(\tilde{y}) \times P(\text{saw} \text{冠詞}) P(\text{冠詞} \tilde{y})$...
名詞	$P(\text{the} \text{名詞})$	$\max_{\tilde{y} \in Y} \delta_1(\tilde{y}) \times P(\text{man} \text{名詞}) P(\text{名詞} \tilde{y})$	$\max_{\tilde{y} \in Y} \delta_2(\tilde{y}) \times P(\text{saw} \text{名詞}) P(\text{名詞} \tilde{y})$...
動詞	$P(\text{the} \text{動詞})$	$\max_{\tilde{y} \in Y} \delta_1(\tilde{y}) \times P(\text{man} \text{動詞}) P(\text{動詞} \tilde{y})$	$\max_{\tilde{y} \in Y} \delta_2(\tilde{y}) \times P(\text{saw} \text{動詞}) P(\text{動詞} \tilde{y})$...

隠れマルコフモデル(HMM):
Viterbi 復号法による y_t の計算例

- y_t を最大にする y_{t-1} から y_t への遷移 (矢印) が決まった時の例 (英語品詞タグ付け)

$x \backslash y$	the	man	saw	...
冠詞	$P(\text{the} \text{冠詞})$	$P(\text{the} \text{動詞})$ $\times P(\text{man} \text{冠詞}) P(\text{冠詞} \text{動詞})$	$P(\text{the} \text{名詞})$ $\times P(\text{man} \text{動詞}) P(\text{動詞} \text{名詞})$ $\times P(\text{saw} \text{冠詞}) P(\text{冠詞} \text{動詞})$...
名詞	$P(\text{the} \text{名詞})$	$P(\text{the} \text{冠詞})$ $\times P(\text{man} \text{名詞}) P(\text{名詞} \text{冠詞})$	$P(\text{the} \text{名詞})$ $\times P(\text{man} \text{動詞}) P(\text{動詞} \text{名詞})$ $\times P(\text{saw} \text{名詞}) P(\text{名詞} \text{動詞})$...
動詞	$P(\text{the} \text{動詞})$	$P(\text{the} \text{名詞})$ $\times P(\text{man} \text{動詞}) P(\text{動詞} \text{名詞})$	$P(\text{the} \text{冠詞})$ $\times P(\text{man} \text{名詞}) P(\text{名詞} \text{冠詞})$ $\times P(\text{saw} \text{動詞}) P(\text{動詞} \text{名詞})$...

隠れマルコフモデル(HMM):

テーブルを用いた最適ラベル列の求め方 (Viterbi 復号法)

- 位置tまでのラベル列の確率が最大になる位置t-1のラベルを記憶するテーブルも同時に計算
- 末端(T)において確率($\pi_T(y)$)が最大になるyの $\pi_T(y)$ からバックトラックすることで、確率が最大になるラベル列を得ることができる。

$$\pi_t(y) = \operatorname{argmax}_{\tilde{y} \in \Sigma_y} \delta_{t-1}(\tilde{y}) P(x_t | y) P(y | \tilde{y})$$

$y \backslash x$	the	man	saw	...
冠詞		$_2(\text{冠詞}) = \text{動詞}$	$_3(\text{冠詞}) = \text{動詞}$...
名詞		$_2(\text{名詞}) = \text{冠詞}$	$_3(\text{名詞}) = \text{動詞}$...
動詞		$_2(\text{動詞}) = \text{名詞}$	$_3(\text{動詞}) = \text{名詞}$...

隠れマルコフモデルのパラメタ推定(学習)

- パラメータ

- 出力確率 $P(x_t|y_t)$

- 遷移確率 $P(y_t|y_{t-1})$

- 最尤推定

$$\hat{\theta} = \operatorname{argmax}_{\theta} \prod_i P_{\theta}(\mathbf{x}^{(i)}, \mathbf{y}^{(i)}) \quad (i \text{は学習データの索引})$$

$$= \operatorname{argmax}_{\theta} \prod_i \prod_t^{T^{(i)}} P_{\theta}(x_t^{(i)} | y_t^{(i)}) P_{\theta}(y_t^{(i)} | y_{t-1}^{(i)})$$

- 最尤なパラメータは共起頻度のカウンタで計算可能

生成モデルの課題点

- 同時分布の推定をすることで、間接的に分類問題を解いている。(P(x)を余分に推定している)

predict

$$\mathbf{y} = \underset{\mathbf{y}}{\operatorname{argmax}} P(\mathbf{y} \mid \mathbf{x}) = \underset{\mathbf{y}}{\operatorname{argmax}} P(\mathbf{x}, \mathbf{y})$$

- 同時分布を推定する十分なデータを得るのは難しいため、変数間の独立性を仮定することが多い。

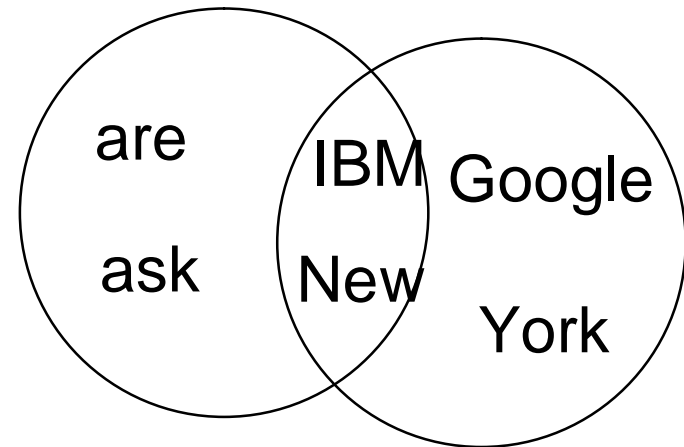
predict

$$\mathbf{y} = \underset{\mathbf{y}}{\operatorname{argmax}} \sum_{t=1}^T P(x_t/y_t)P(y_t \mid y_{t-1}) \quad (\text{HMM})$$

相互作用のある観測変数(素性)をうまく扱えない

相互作用のある素性の例

- 自然言語処理では単語それ自身以外に、単語の部分文字列等が素性に使われることが多い
- 品詞タグ付け
 - P(beautiful|形容詞), P(fulで終わる単語|形容詞)
 - P(beで始まる単語|形容詞)
 - P(immediately|副詞), P(lyで終わる単語|副詞)
- 固有表現抽出
 - P(3文字の単語 | B-組織), P(最初が大文字の単語 | B-組織)
 - P(小泉|B-人名), P(漢字だけからなる単語| B-人名)



足して1にならないといけないのになってない

識別モデル

- x から y を直接推定するモデル
- 識別モデルの利点
 - 直接分類問題を解くことが出来る
 - 素性間の重なりを考慮して重みを学習
- 多クラスのロジスティック回帰モデル(最大エントロピーモデル)
 - 確率分布の形をした識別モデル
 - 条件付分布を直接モデル化

内積だよ

$$P(y | \mathbf{x}) = \frac{\exp(\langle \boldsymbol{\theta}, \boldsymbol{\Phi}(\mathbf{x}, y) \rangle)}{\sum_{\tilde{y} \in \mathbf{Y}} \exp(\langle \boldsymbol{\theta}, \boldsymbol{\Phi}(\mathbf{x}, \tilde{y}) \rangle)}$$

($\boldsymbol{\theta}$ は x, y の素性、 $\boldsymbol{\Phi}$ は素性に対する重み)

復号は一緒

識別モデルのパラメータ推定(学習)

- 最尤推定

$$\begin{aligned}\hat{\boldsymbol{\theta}} &= \arg \max_{\boldsymbol{\theta}} \prod_{i \in \text{trainingdata}} P_{\boldsymbol{\theta}}(y^{(i)} | \mathbf{x}^{(i)}) \\ &= \arg \max_{\boldsymbol{\theta}} \prod_{i \in \text{trainingdata}} \frac{\exp(\langle \boldsymbol{\theta}, \Phi(\mathbf{x}^{(i)}, y^{(i)}) \rangle)}{\sum_{\tilde{y} \in \mathbf{Y}} \exp(\langle \boldsymbol{\theta}, \Phi(\mathbf{x}^{(i)}, \tilde{y}) \rangle)}\end{aligned}$$

- パラメータの計算

- 素性間の相互作用を考慮してパラメータ推定をする必要があり、計算手順が複雑
- 損失関数の偏微分を用いて損失最小化問題を解く。

識別モデルによる構造ラベル付与学習: 条件付確率場 (Conditional Random Fields: CRF)

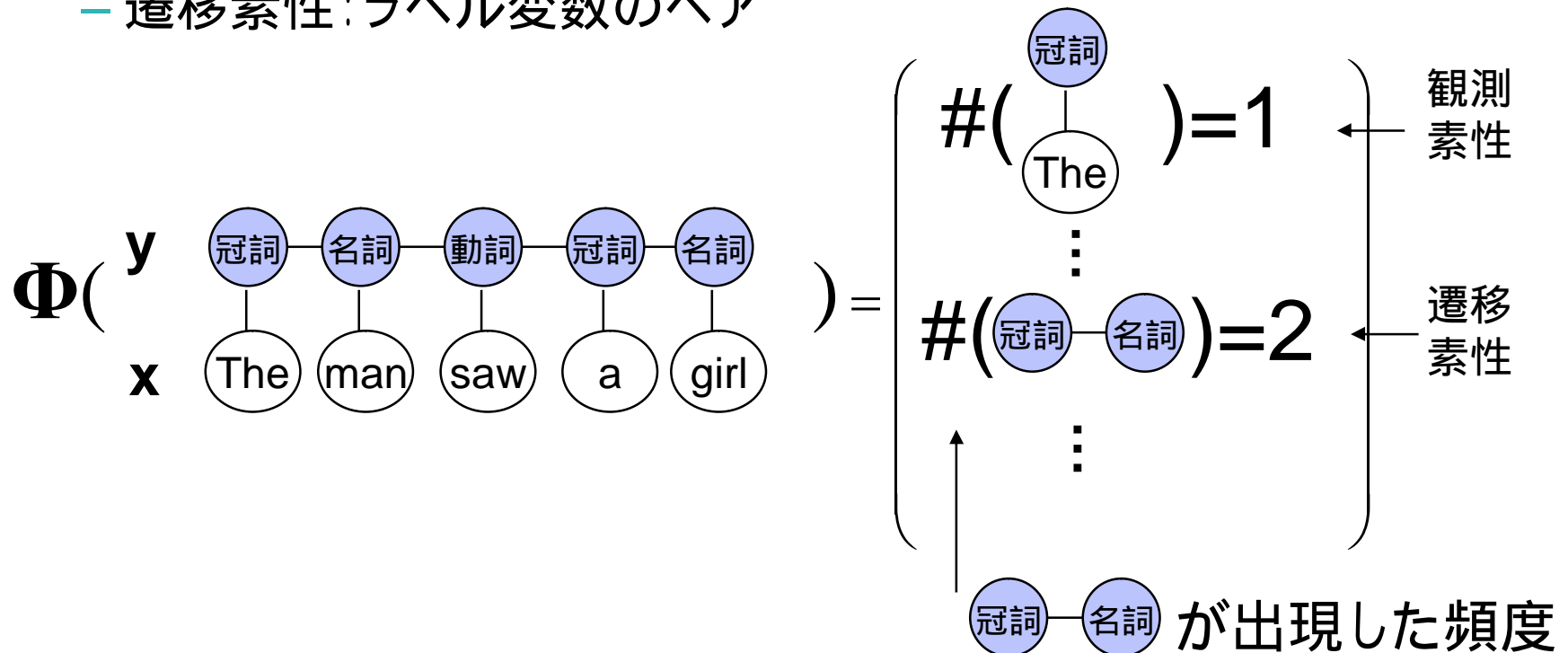
- ロジスティック回帰モデルを基に、ローカルな変数間の関係を素性(遷移素性)で表現したモデル

$$P(\mathbf{y} | \mathbf{x}) = \frac{\exp(\langle \boldsymbol{\theta}, \phi(\mathbf{x}, \mathbf{y}) \rangle)}{\sum_{\tilde{\mathbf{y}}} \exp(\langle \boldsymbol{\theta}, \phi(\mathbf{x}, \tilde{\mathbf{y}}) \rangle)}$$
$$= \frac{\exp\left(\sum_{\tau=1}^T \langle \boldsymbol{\theta}, \phi(\mathbf{x}, \mathbf{y}_{\tau}^{\tau+1}) \rangle\right)}{\sum_{\tilde{\mathbf{y}}} \exp\left(\sum_{\tau=1}^T \langle \boldsymbol{\theta}, \phi(\mathbf{x}, \tilde{\mathbf{y}}_{\tau}^{\tau+1}) \rangle\right)} \cdot \mathbf{y}_{\tau}^{\tau+1} = (y_{\tau}, y_{\tau+1})$$

(ϕ は \mathbf{x}, \mathbf{y} の素性、 $\boldsymbol{\theta}$ は素性に対する重み)

条件付確率場の素性

- 素性ベクトルの各要素は、素性が構造中に出現した頻度
 - 観測素性: 観測変数とラベル変数のペア
 - 遷移素性: ラベル変数のペア



その他の識別モデル：分類手法に基づくアプローチ

- **CRFのパラメタ推定の別アプローチ**
 - 分類手法(パーセプトロン, SVM, ...)にもとづいたアプローチ
- **何がうれしいか？ 通常のアプローチでは解けない問題に適用できる**

そのほかの識別モデル：隠れマルコフパーセプトロン

- CRFの推定は、訓練データの $y^{(i)}$ が出力される確率を「最大にする」ように学習される

$$\hat{\theta} = \arg \max_{\theta} \prod_{i \in \text{training data}} P_{\theta}(y^{(i)} | \mathbf{x}^{(i)})$$

- 別の考え方：訓練データの $y^{(i)}$ が出力される確率が、「ほかの $\tilde{y}^{(i)} \neq y^{(i)}$ の出力確率よりも大きければよい」

$$\log P_{\theta}(y^{(i)} | \mathbf{x}^{(i)}) > \log P_{\theta}(\tilde{y}^{(i)} | \mathbf{x}^{(i)})$$

$$\Leftrightarrow \langle \theta, \Phi(\mathbf{x}^{(i)}, y^{(i)}) \rangle - \langle \theta, \Phi(\mathbf{x}^{(i)}, \tilde{y}^{(i)}) \rangle > 0$$

隠れマルコフパーセプトロンのアルゴリズム

1. i 番目の訓練データに対して予測してみる

$$\underset{\underset{\text{predict}}{y}}{y^{(i)}} = \operatorname{argmax} P(y | \mathbf{x}^{(i)})$$

2. 当たっていれば ($\underset{\text{predict}}{y^{(i)}} \neq y^{(i)}$) 何もしない

3. 外れていたら、パラメータを修正

$$\boldsymbol{\theta}^{new} \leftarrow \boldsymbol{\theta}^{old} + \eta \left(\Phi(\mathbf{x}^{(i)}, y^{(i)}) - \Phi(\mathbf{x}^{(i)}, \tilde{y}^{(i)}) \right)$$

4. 1.に戻って繰り返す

最尤推定と隠れマルコフパーセプトロンの違いは、argmax操作のみで実現できるところ

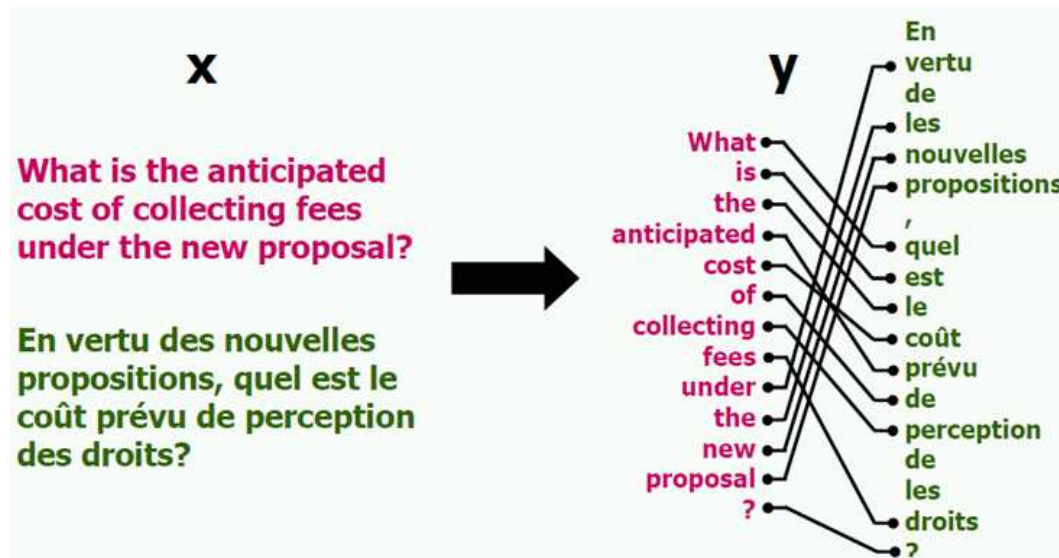
- 隠れマルコフパーセプトロンは訓練と予測が両方argmax操作

	CRF	隠れマルコフパーセプトロン
訓練		argmax
予測	argmax	argmax

- 通常、argmaxも動的計画法で実現できるのであまり変わりはない...しかし...

隠れマルコフパーセプトロンは、最尤推定では解けない問題が解ける？

- 問題によっては、argmax操作は、動的計画法以外の多項式時間アルゴリズムで実現できる
- 動的計画法で多項式時間で解けない問題は、最尤推定では解けない
 - たとえば、異なる言語の文章間での単語のマッチング
 - argmax操作が線形計画法(多項式時間)で解ける



図引用:

Dan Klein and Ben Tasker, "Max-Margin Methods for NLP: Estimation, Structure, and Applications", ACL 2005 Tutorial より

End of the Presentation

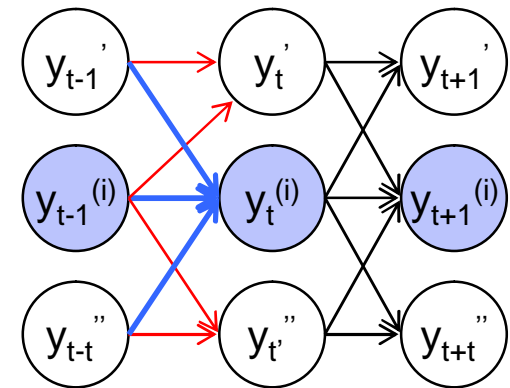
Reference

- **Dan Klein and Ben Tasker. Max-Margin Methods for NLP: Estimation, Structure, and Applications. ACL-05 Tutorial.**
- **Dan Klein and Chris Manning. Maxent Models, Conditional Estimation, and Optimization, without the Magic. ACL-03 Tutorial.**

Gradient of Point-wise Log-Loss

$$\frac{L}{\partial \boldsymbol{\theta}} = - \sum_{i \in \text{training data}} \left(\sum_t \sum_{\hat{\mathbf{y}} \in \hat{y}_t = y_t^{(i)}} P(\hat{\mathbf{y}} | \mathbf{x}^{(i)}) \Phi(\mathbf{x}^{(i)}, \hat{\mathbf{y}}) - \sum_{\check{\mathbf{y}} \in \mathbf{Y}} P(\check{\mathbf{y}} | \mathbf{x}^{(i)}) \Phi(\mathbf{x}^{(i)}, \check{\mathbf{y}}) \right)$$

$$\sum_{\hat{\mathbf{y}} \in \hat{y}_t = y_t^{(i)}} P(\hat{\mathbf{y}} | \mathbf{x}^{(i)}) \Phi(\mathbf{x}^{(i)}, \hat{\mathbf{y}}) = \frac{\sum_{\hat{\mathbf{y}} \in \hat{y}_t = y_t^{(i)}} P(\hat{\mathbf{y}} | \mathbf{x}^{(i)})}{\sum_{\check{\mathbf{y}} \in \mathbf{Y}} P(\check{\mathbf{y}} | \mathbf{x}^{(i)})} \Phi(\mathbf{x}^{(i)}, \hat{\mathbf{y}})$$



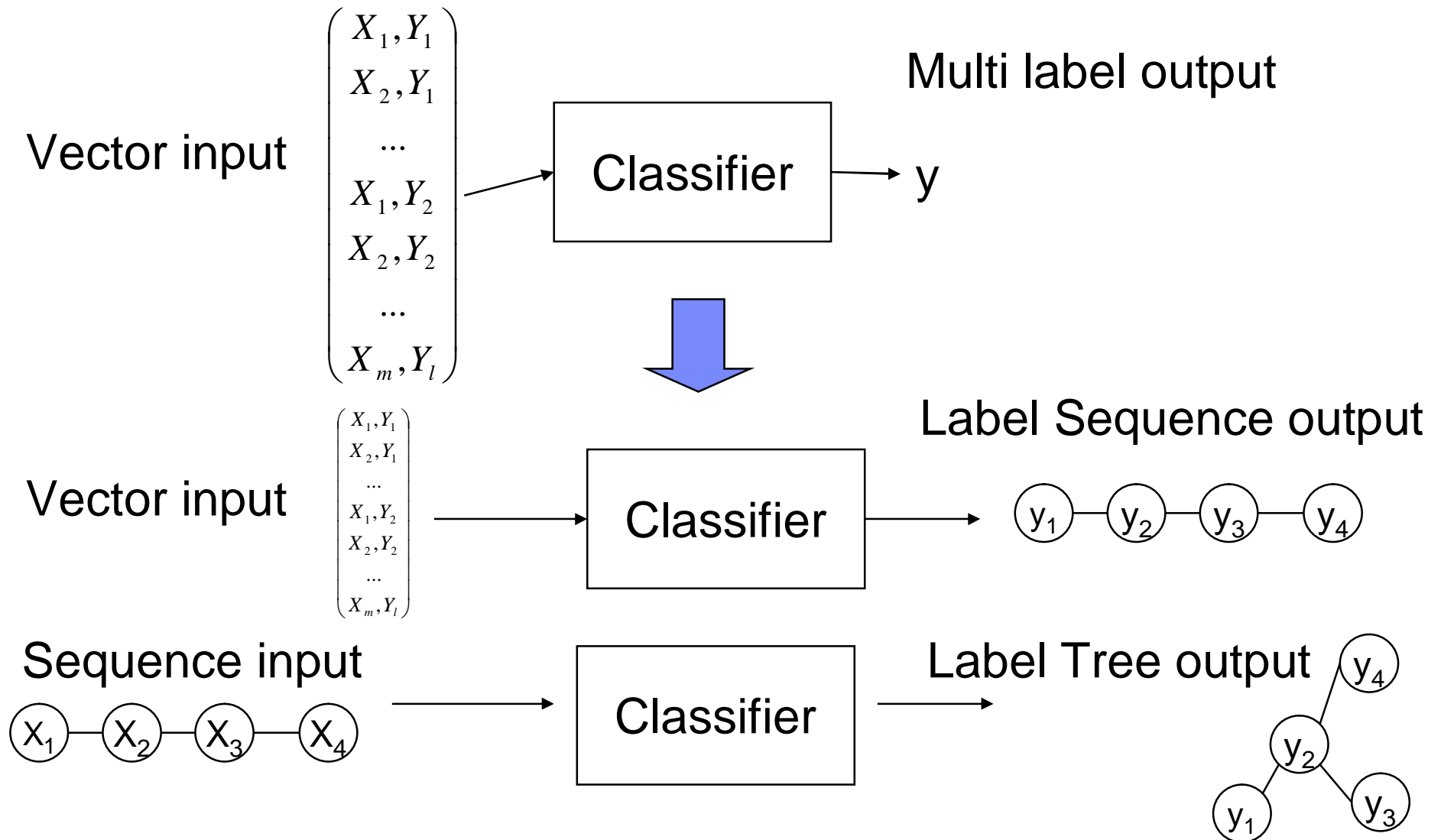
Classification of Structured Output (e.g. HMMs, CRFs)

- **Viterbi Decoding**
- **Finding the most probable label sequence.**
- **Avoids re-calculation**

Parameter Estimation of Structured Output (CRFs)

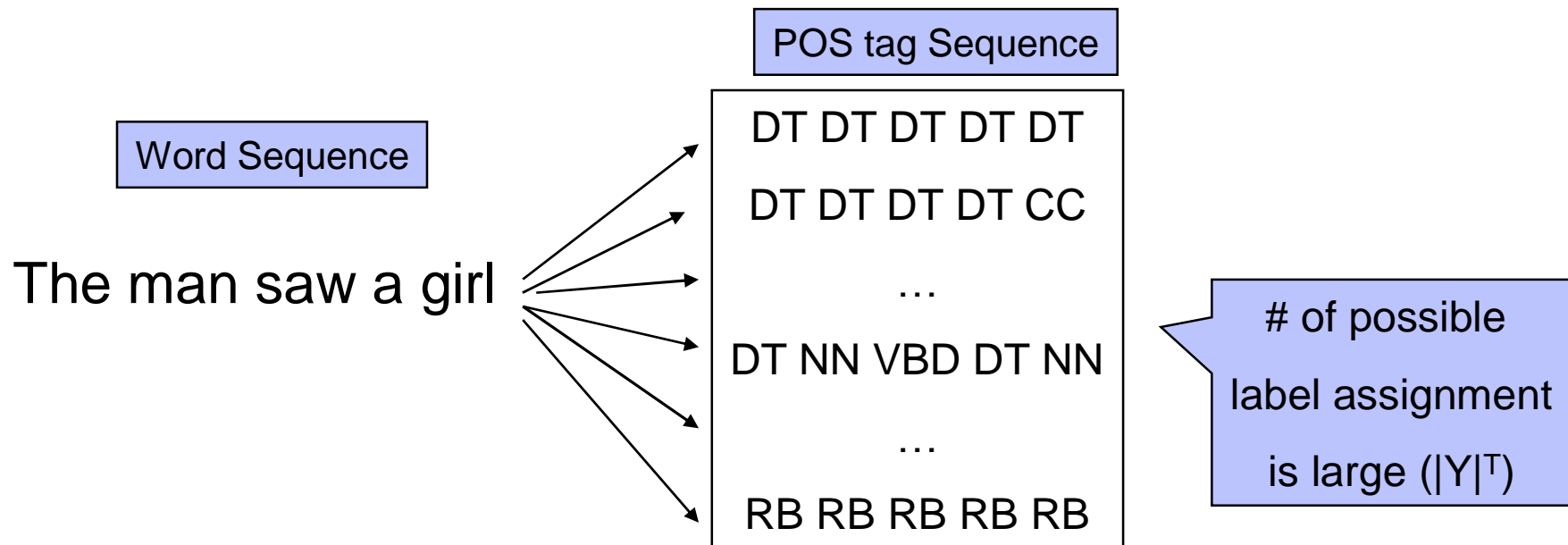
$$\begin{aligned} -\sum_i \log P(Y | \mathbf{X}) &= -\sum_i \log \frac{\exp(\langle \boldsymbol{\theta}, \boldsymbol{\varphi}(\mathbf{X}, Y) \rangle)}{\sum_{\tilde{Y} \in \mathbf{Y}} \exp(\langle \boldsymbol{\theta}, \boldsymbol{\varphi}(\mathbf{X}, \tilde{Y}) \rangle)} \\ &= -\sum_i \left(\langle \boldsymbol{\theta}, \boldsymbol{\varphi}(\mathbf{X}, Y) \rangle - \log \sum_{\tilde{Y} \in \mathbf{Y}} \exp(\langle \boldsymbol{\theta}, \boldsymbol{\varphi}(\mathbf{X}, \tilde{Y}) \rangle) \right) \end{aligned}$$

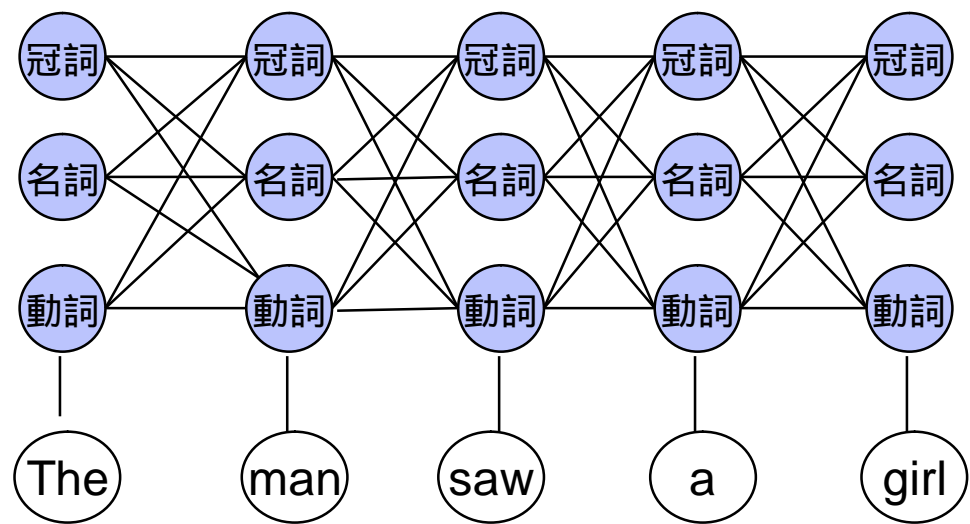
Structured Output Learning (e.g. POS tagging, NER, Parsing)



Structured Output Learning

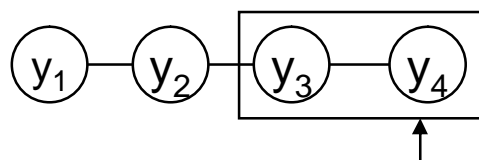
- In general, structured output problem is one of multi label problems.
- In practice, all the labels are not enumerable and sparse in training data.





条件付確率場(CRF): 識別モデルによる構造ラベル付与学習

- **Dynamic Programming (Belief Propagation) algorithm**
 - solves graphical inference problems via a series of local message-passing operations.
 - a label depends only its neighborhood labels (Markov assumption)
 - reuses previous calculation to avoid re-calculation



When we decide y_4 , we only see the decision of y_3 .

- **Transition features**
 - represent structural label consistency

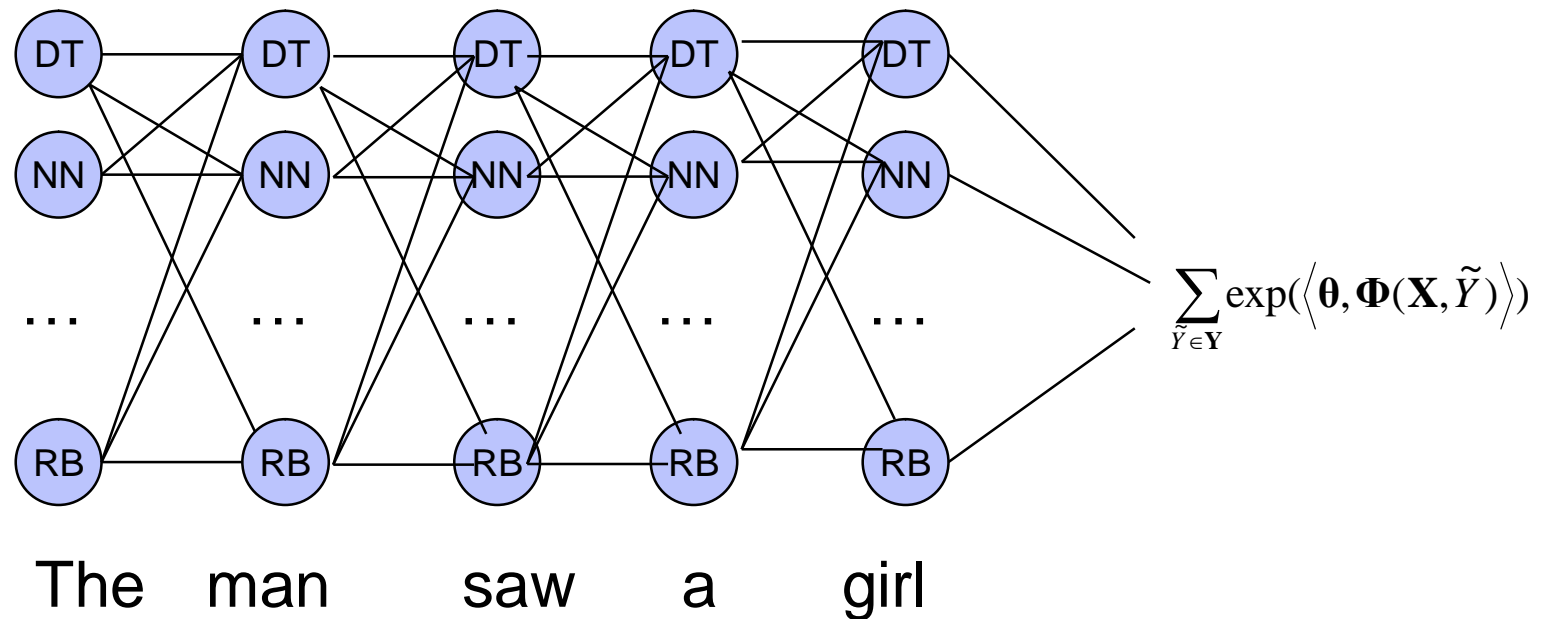
Machine Learning on Structured data

$$\begin{aligned} P(\mathbf{y} \mid \mathbf{x}) \\ &= \frac{\exp(\langle \boldsymbol{\theta}, \Phi(\mathbf{x}, \mathbf{y}) \rangle)}{\sum_{\tilde{\mathbf{y}} \in \mathbf{Y}} \exp(\langle \boldsymbol{\theta}, \Phi(\mathbf{x}, \tilde{\mathbf{y}}) \rangle)}, \end{aligned}$$

Dynamic Programming

$$P(Y | \mathbf{X}) = \frac{\exp(\langle \boldsymbol{\theta}, \Phi(\mathbf{X}, Y) \rangle)}{\sum_{\tilde{Y} \in \mathcal{Y}} \exp(\langle \boldsymbol{\theta}, \Phi(\mathbf{X}, \tilde{Y}) \rangle)}$$

- reuses previous calculation to avoid re-calculation
- 計算したい値
 - 存在しうる全てのLabel Pathのスコアの和
 - あるポジションtにある素性が出てくるLabel Pathのスコアの和



分類問題で使われる生成モデル (例: Naïve Bayes)

ベイズの法則

- ベイズ分類器
predict

$$y = \underset{y}{\operatorname{argmax}} P(y | \mathbf{x}) = \underset{y}{\operatorname{argmax}} P(\mathbf{x} | y) P(y)$$

(\mathbf{x} は観測変数ベクトル、 y はラベル変数)

- 結合分布

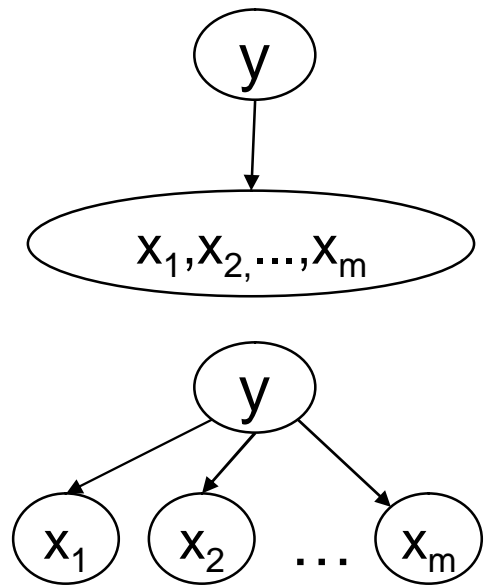
$$P(\mathbf{x} | y) = P(x_1, x_2, \dots, x_m | y)$$

- ナイーブ分布

(結合条件付分布を緩め、観測変数間の独立性を仮定)

$$P(\mathbf{x} | y) = \prod_j^{|\mathbf{x}|} P(x_j | y)$$

($|\mathbf{x}|$ はベクトルのサイズ)



テーブルを用いた最適ラベル列の求め方 (Viterbi 復号法)

- 位置tまでのラベル列の確率が最大になる位置t-1のラベルを記憶するテーブルも同時に計算
- 末端(T)において確率($\pi_T(y)$)が最大になるyの $\pi_T(y)$ からバックトラックすることで、確率が最大になるラベル列を得ることができる。

$$\pi_t(y) = \operatorname{argmax}_{\tilde{y} \in \Sigma_y} \delta_{t-1}(\tilde{y}) P(x_t | y) P(y | \tilde{y})$$

$y \backslash x$	the	man	saw	...
冠詞		$\delta_2(\text{冠詞}) = \text{動詞}$	$\delta_3(\text{冠詞}) = \text{動詞}$...
名詞		$\delta_2(\text{名詞}) = \text{冠詞}$	$\delta_3(\text{名詞}) = \text{動詞}$...
動詞		$\delta_2(\text{動詞}) = \text{名詞}$	$\delta_3(\text{動詞}) = \text{名詞}$...

Another interpretation of HMM

- Log likelihood can be written as the inner product of **weight vector w** & **feature vector**

$$\log \Pr \left(\begin{array}{cccc} \text{DT} & \text{NN} & \text{VBD} & \text{NN} \\ \circ & \circ & \circ & \circ \\ | & | & | & | \\ \text{the} & \text{man} & \text{saw} & \text{glasses} \end{array} \right) = \langle \mathbf{w} \cdot \Phi(\mathbf{x}, \mathbf{y}) \rangle$$

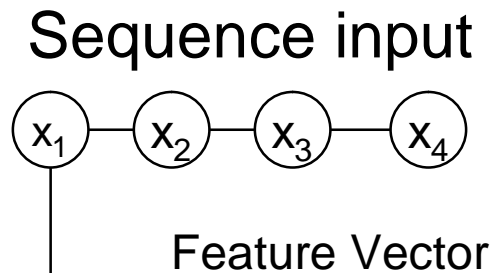
$$\mathbf{W} = \begin{pmatrix} \log \Pr \left(\begin{array}{c} \text{DT} \\ \circ \\ | \\ \text{the} \\ \text{NN} \end{array} \right) \\ \log \Pr \left(\begin{array}{c} \circ \\ | \\ \text{man} \\ \text{VBD NN} \end{array} \right) \\ \vdots \end{pmatrix} \quad \Phi(\mathbf{x}, \mathbf{y}) = \begin{pmatrix} \# \left(\begin{array}{c} \text{DT} \\ \circ \\ | \\ \text{the} \\ \text{NN} \end{array} \right) \\ \# \left(\begin{array}{c} \circ \\ | \\ \text{man} \\ \text{VBD NN} \end{array} \right) \\ \vdots \end{pmatrix}$$

weight vector (model parameters)

feature vector (labeled sequence)

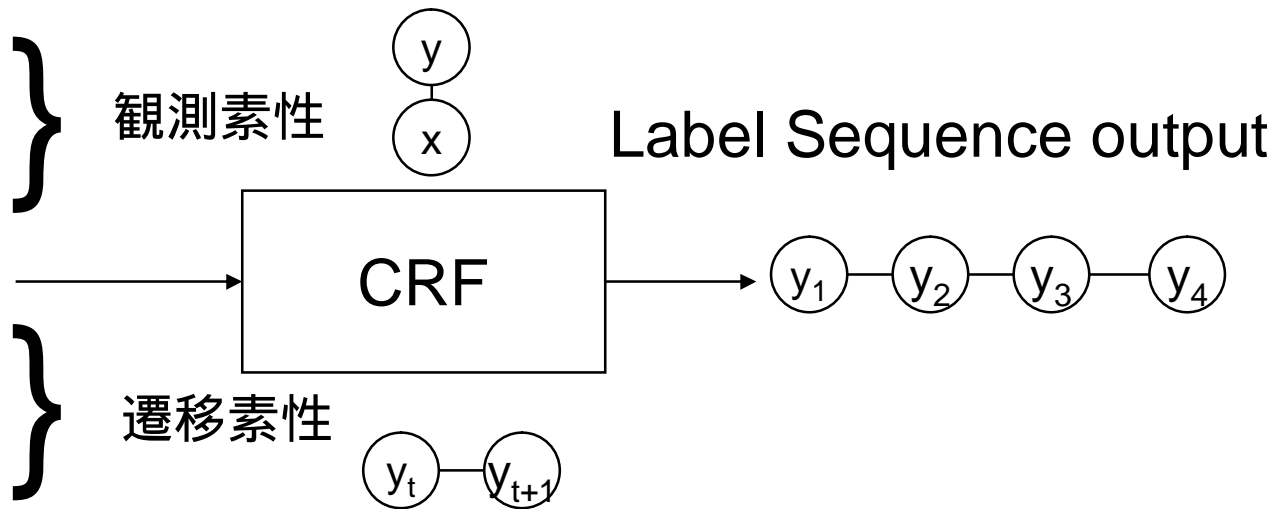
Number of times $\begin{array}{c} \text{VBD} \text{ NN} \\ \circ - \circ \end{array}$ appearing

識別モデルによる構造ラベル付与学習: 条件付確率場 (Conditional Random Fields: CRF)



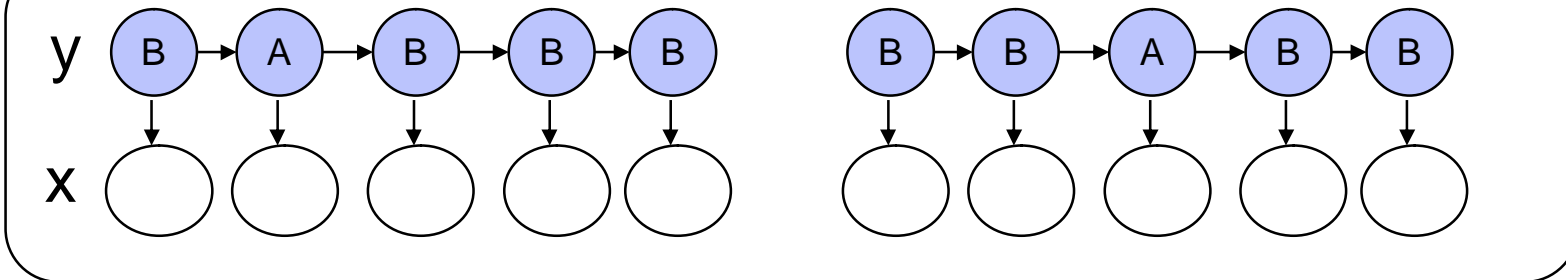
最大エントロピーモデルで素性はこんな
ローカルな変数間の関係を素性で表現

$$\left(\begin{array}{l} \sum_t x_t = a \wedge y_t = y_a \\ \sum_t x_t = x_a \wedge y_t = y_1 \\ \dots \\ \sum_t X = x_m \wedge Y = y_l \\ \sum_t Y_t = y_a \wedge Y_{t+1} = Y_a \\ \dots \\ \sum_t Y_t = y_l \wedge Y_{t+1} = Y_l \end{array} \right)$$



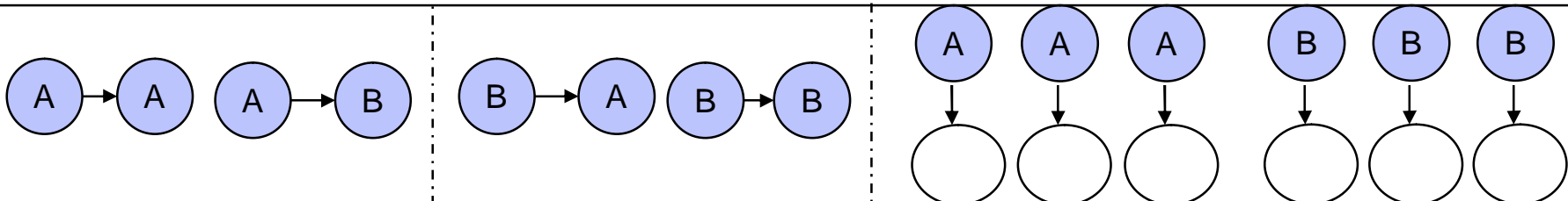
生成モデルと識別モデルの素性の重みの違い

学習データ



パラメータ推定

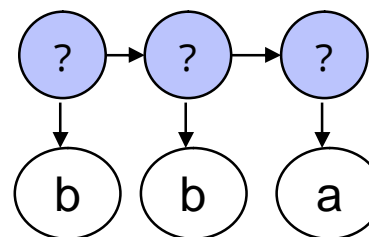
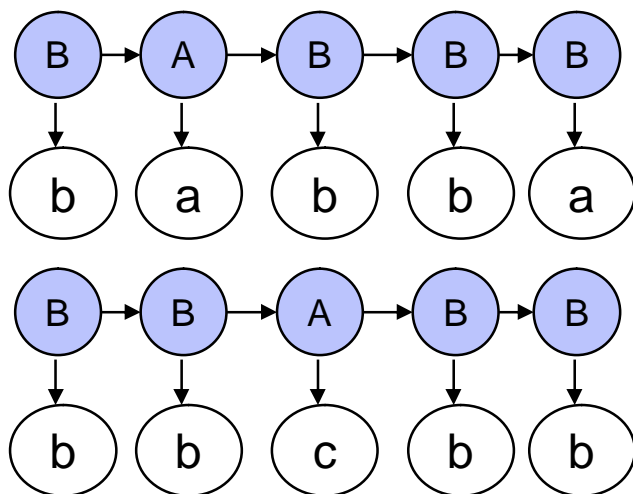
HMMの重み(1/2以上を青、1/2未満を赤)



CRFの重み(プラスを青、マイナスを赤)

識別に影響が少ない B → B の重みが低い

生成モデルと識別モデル



$$9/10 * 6/8 * 9/10 * 6/8 * 1/10$$

$$= 2916 / 64000$$

$$= 0.4556$$

$$9/10 * 6/8 * 9/10 * 2/8 * 1/2$$

$$= 972 / 12800$$

$$= 0.0759$$

■ HMM

	0/2				
	2/2				
	6/8				
	2/8	1/2	1/2	1/10	9/10

- **# weights for observation**
- **A-a:5.96 A-b:-13.41 A-c:5.96**
- **B-a:-5.96 B-b:13.41 B-c:-5.96**
- **# weights for transition**
- **A-A:-2.14 A-B:2.28**
- **B-A:1.91 B-B:-2.05**

- **Y-X**

- A-a:6.09 A-b:-20.67

- B-a:-6.09 B-b:20.67

- **Y-Y**

- A-A:-3.69 A-B:10.08

- B-A:-2.45 B-B:-3.93