

カーネル法に基づく構造データのラベル付け学習アルゴリズム

鹿島 久嗣 坪井 祐太

† 日本アイ・ビー・エム株式会社 東京基礎研究所 〒 242-8502 神奈川県大和市下鶴間 1623-14

E-mail: †{hkashima,yutat}@jp.ibm.com

あらまし 本稿において、我々はパーセプトロンアルゴリズムに基づく、配列や木、グラフなどの構造を持ったデータのラベル付け学習アルゴリズムを提案する。また、ラベル付けに用いることのできるいくつかのカーネル関数とその効率的な計算法を与える。提案手法は、完全にカーネル化されており、かつ、点ごとのラベル予測を行うため、任意の数の非観測変数を含む、サイズの大きい素性を用いることができる。この点において、提案手法は最大エントロピーマルコフモデルや条件付確率場など少数の非観測変数をもつ素性しか扱うことのできない既存の手法と大きく異なる。
キーワード カーネル法, パーセプトロン, 周辺化カーネル, 固有表現抽出, 情報抽出

Kernel-based Discriminative Learning Algorithms for Labeling Structured Data

Hisashi KASHIMA and Yuta TSUBOI

† IBM Tokyo Research Laboratory Shimotsuruma 1623-14, Yamato-shi, Kanagawa, 242-8502 Japan

E-mail: †{hkashima,yutat}@jp.ibm.com

Abstract We introduce a new perceptron-based discriminative learning algorithm for labeling structural data such as sequences, trees and graphs. Since it is fully kernelized and employs the pointwise label prediction, large features including arbitrary number of hidden variables can be incorporated with polynomial time complexity. This is contrasted with existing labelers that can handle only features of a small number of hidden variables such as Maximum Entropy Markov Models and Conditional Random Fields. We also introduce several kernel functions for labeling sequences, trees and graphs and the efficient algorithms for them.

Key words Kernel Methods, Perceptron, Marginalized Kernel, Named Entity Recognition, Information Extraction

1. ま え が き

ラベル付け問題は、観測変数集合の値が与えられたときに、非観測変数集合の値を予測する問題であり、教師ありの多クラス分類問題を一般化した問題である。特に、配列ラベル付けは、自然言語処理、バイオインフォマティクス、Webデータの解析など多くの分野で見られる重要な問題である。配列ラベル付けの学習問題は長年活発に研究がなされてきたが、より一般的な木やグラフといった構造データのラベル付けについてはほとんど研究がなされていない。本稿では、これら一般的な構造データのラベル付け学習問題において、カーネル法によるアプローチを行う。

配列ラベル付け問題には長年、隠れマルコフモデル (HMM) が用いられ、一定の成功を収めてきた。しかしながら、HMMは観測変数と非観測変数の同時確率分布を推定するという、より難しい問題を解いており、多くの学習データを必要とする。また、観測変数は、非観測変数が与えられると互いに独立であ

ることが求められるため、単語と、その頭文字のように互いに独立でない素性を自然に扱うことができないという問題がある。一方、近年、最大エントロピーマルコフモデル (MEMM) [1] や条件付確率場 (CRF) [2] などの条件付モデルと呼ばれるモデルが、観測変数から非観測変数を予測するという目的に適しているため、注目を集めている。さらに、Collins [3] は隠れマルコフ (HM) パーセプトロンと呼ばれる、パーセプトロンアルゴリズムに基づく効率的なアルゴリズムを提案している。これを拡張する形で、サポートベクターマシン (SVM) に基づくアルゴリズム [4] やブースティングに基づくアルゴリズム [5] も提案されている。

いくつかの問題において、自然言語処理におけるイディオムや、バイオインフォマティクスにおけるモチーフなどのようにバイグラムなどの局所的な素性を組み合わせるだけでは広い範囲の文脈を捉えるのに充分でない場合がある。しかしながら、現在提案されているほとんどの手法では、少数の非観測変数を含む素性しか用いることができない。ラベル付けの段階は、通

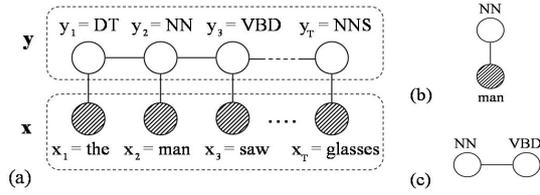


図 1 (a) 1 次のマルコフ仮定に基づく配列のグラフィカルモデル表現. 黒いノードは観測変数を, 白いノードは非観測変数を表す. 文 “the, man, saw, ..., glasses.” は, 観測変数のラベル列 \mathbf{x} , 品詞タグ列 “DT, NN, VBD, ..., NNS” は, 非観測変数のラベル列 \mathbf{y} を表す. (b) 観測変数と非観測変数のペアからなる素性 (c) 連続する 2 つの非観測変数からなる素性

常 Viterbi アルゴリズムのような動的計画法に基づく方法が用いられるが, ひとつの素性に含まれる非観測変数の数に対して計算量が指数的に増大する. HM パーセプトロンや HMSVM は双対型で書くこともできるが, やはりこの問題は解決されない.

本稿では, はじめにラベル付け候補を絞ることなく, 大域的な素性を用いてラベル付けを行うことのできる方法である **周辺化ラベル付けパーセプトロン** を提案する. これは, 周辺化カーネル [6] と, 点ごとのラベル予測 [7] という 2 つのアイデアに基づいている. 点ごとのラベル予測によって, 各非観測変数のラベルをそれぞれ予測できるため, 前述の問題が解決される. また, 周辺化カーネル [6] を用いたカーネル化によって, 任意の数の非観測変数を含むような大きい素性を扱うことができるようになる. さらに我々は, 配列や順序木といったいくつかの構造データにたいして, 実際の周辺化カーネル関数を提案する. 近年, カーネル法を用いた構造データの分類は精力的に研究されており, 配列カーネル [8], [9] や木カーネル [10], [11], グラフカーネル [12], [13] などが提案されている^(注1). これらは大抵, 広範囲にわたる文脈を取り入れるため, 部分列や部分グラフなどの部分構造が, 構造データ内に出現する回数などを用いて素性ベクトルを定義している. これらの素性ベクトルは通常, 高次元 (場合によっては無限次元) になるため, 接尾辞木や動的計画法, あるいは行列計算などを用いて, 素性を陽に数え上げることなく効率的にカーネル関数を計算する方法を提案されている. 我々はこれらのカーネルに基づき, 配列や順序木に対する周辺化カーネルと, その効率的な計算方法を与える.^(注2) また, 最後に実世界データを用いた, 配列と木に対するラベル付け問題の実験を行い, その有効性を検証する.

2. 隠れマルコフ (HM) パーセプトロン

ラベル付け問題は, 観測変数 $\mathbf{x} = (x_1, x_2, \dots, x_T)$ から非観測変数 $\mathbf{y} = (y_1, y_2, \dots, y_T)$ へのマッピングを求める問題として定義できる. ここで $x_t \in \Sigma_x$, $y_t \in \Sigma_y$ とする. 例えば, 品詞タグ付けなどのタスクにおいて, x_t は t 番目の単語を表わし, y_t は t 番目の単語の品詞を表す (図 1). 学習器は, 訓練データの集合 $E = (e^{(1)}, e^{(2)}, \dots, e^{(|E|)})$ を用いることができる. ここで, $e^{(i)} = (\mathbf{x}^{(i)}, \mathbf{y}^{(i)})$ は i 番目の訓練例で, $|\mathbf{x}^{(i)}| = |\mathbf{y}^{(i)}| = T^{(i)}$ とする.

Collins [3] は, CRF や MEMM の代替手法としてパーセプ

(注1): [14] に詳細なサーベイがある.

(注2): 非循環有向グラフに対する周辺化カーネルは [15] を参照されたい.

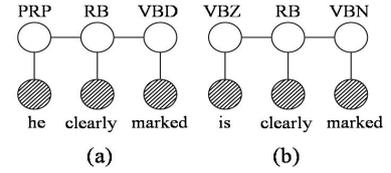


図 2 配列ラベリングのための素性 (長さ 3). バイグラムの素性では “marked” の品詞を決定するには充分でない.

トロンに基づくラベル付け学習アルゴリズムである HM パーセプトロンを提案した. HM パーセプトロンは訓練例をオンライン型で順番に処理できるため効率がよく, その性能は CRF に匹敵することが示されている [16]. HM パーセプトロンの重要な点は, 観測変数集合から非観測変数集合へのマッピングを, 観測変数集合の値と非観測変数集合の値の組の 2 値分類, すなわち $\Sigma_x^T \times \Sigma_y^T \rightarrow \{+1, -1\}$ として捉える点である. F を全ての素性の集合とし, $\phi_f(\mathbf{x}, \mathbf{y})$ を, 素性 $f \in F$ が (\mathbf{x}, \mathbf{y}) 内に出現する回数, $\Phi(\mathbf{x}, \mathbf{y})$ はそれをベクトルの形にあらわしたものとす. \mathbf{x} が入力として与えられると, HM パーセプトロンは現在の素性の重みに基づき, 予測 $\hat{\mathbf{y}}$ を出力する.

$$\hat{\mathbf{y}} = \operatorname{argmax}_{\mathbf{y} \in \Sigma_y^T} \sum_{f \in F} w_f \phi_f(\mathbf{x}, \mathbf{y}) = \operatorname{argmax}_{\mathbf{y} \in \Sigma_y^T} \langle \mathbf{w}, \Phi(\mathbf{x}, \mathbf{y}) \rangle \quad (1)$$

ここで w_f は素性 f の重みであり, \mathbf{w} はこれらをベクトルの形で表したものとす. 1 次のマルコフ仮定の下では, 観測変数と非観測変数のペアからなる素性 (図 1(b)) と, 連続する 2 つの非観測変数からなる素性 (図 1(c)) の, 2 つのタイプの素性が考えられる. 重みの学習は, 訓練例 $e^{(i)}$ に対する予測が誤り, すなわち $\hat{\mathbf{y}}^{(i)} \neq \mathbf{y}^{(i)}$ であったときに,

$$\mathbf{w}^{new} = \mathbf{w}^{old} + \Phi(\mathbf{x}^{(i)}, \mathbf{y}^{(i)}) - \Phi(\mathbf{x}^{(i)}, \hat{\mathbf{y}}^{(i)}) \quad (2)$$

によって更新される. さらに (1) は訓練例の重み α を用いることで双対型で書くことができる.

$$\hat{\mathbf{y}} = \operatorname{argmax}_{\mathbf{y} \in \Sigma_y^T} \sum_{j=1}^{|E|} \sum_{\tilde{\mathbf{y}} \in \Sigma_y^T} \alpha_j(\tilde{\mathbf{y}}) \langle \Phi(\mathbf{x}^{(j)}, \tilde{\mathbf{y}}), \Phi(\mathbf{x}, \mathbf{y}) \rangle \quad (3)$$

初期値 $\alpha = \mathbf{0}$ から始めると, 重みの更新ルールは以下のように書き直せる.

- $\alpha_i^{new}(\mathbf{y}^{(i)}) = \alpha_i^{old}(\mathbf{y}^{(i)}) + 1$
- $\alpha_i^{new}(\hat{\mathbf{y}}^{(i)}) = \alpha_i^{old}(\hat{\mathbf{y}}^{(i)}) - 1$

$(\mathbf{x}^{(i)}, \hat{\mathbf{y}}^{(i)})$ は負例として働くため, 擬似負例と言われる.

さて, 問題によっては上で用いたサイズ 2 の素性 (図 1(b)(c)) では広範囲にわたる文脈をモデル化するには充分でない場合がある. そこで HM パーセプトロンにおいて, 長さ d までの素性 (図 2) を用いることを考える. $\mathbf{x}_{t_1}^{t_2}$ は \mathbf{x} の t_1 番目から t_2 番目までの観測変数集合を表すとし, $\mathbf{y}_{t_1}^{t_2}$ も同様に定義する. 全ての素性は高々 d 個の非観測変数に依存することを考慮すると, $\Phi(\mathbf{x}, \mathbf{y})$ は以下のように展開できる.

$$\begin{aligned} \Phi(\mathbf{x}, \mathbf{y}) &= \Phi(\mathbf{x}_1^d, \mathbf{y}_1^d) - \Phi(\mathbf{x}_2^d, \mathbf{y}_2^d) \\ &\quad + \Phi(\mathbf{x}_2^{d+1}, \mathbf{y}_2^{d+1}) - \Phi(\mathbf{x}_3^{d+1}, \mathbf{y}_3^{d+1}) \\ &\quad \dots \\ &\quad + \Phi(\mathbf{x}_{T-d}^{T-1}, \mathbf{y}_{T-d}^{T-1}) - \Phi(\mathbf{x}_{T-d+1}^{T-1}, \mathbf{y}_{T-d+1}^{T-1}) \\ &\quad + \Phi(\mathbf{x}_{T-d+1}^T, \mathbf{y}_{T-d+1}^T). \end{aligned} \quad (4)$$

(4) を (1) に代入することで、 \mathbf{y} の予測は動的計画法を用いて行うことができる。

$$\begin{aligned} \max_{\mathbf{y}} \langle \mathbf{w}, \Phi(\mathbf{x}, \mathbf{y}) \rangle &= \max_{y_1, \dots, y_d} s(y_1, \dots, y_d) \\ s(y_t, \dots, y_{t+d-1}) &= \max_{y_{t+d}} s(y_{t+1}, \dots, y_{t+d}) \\ &+ \langle \mathbf{w}, \Phi(\mathbf{x}_t^{t+d-1}, \mathbf{y}_t^{t+d-1}) - \Phi(\mathbf{x}_{t+1}^{t+d-1}, \mathbf{y}_{t+1}^{t+d-1}) \rangle \\ s(y_{T-d+1}, \dots, y_T) &= \langle \mathbf{w}, \Phi(\mathbf{x}_{T-d+1}^T, \mathbf{y}_{T-d+1}^T) \rangle \end{aligned}$$

しかしながら、各再帰のステップの計算時間は d について指数的に増加してしまう。たとえ双対型 (3) にしても、長さ d までの非観測変数の文脈は陽に考えなければならないため同様の問題は残る。これはカーネルを用いて任意のサイズの素性を導入する上で重大な障害となる。CRF も同様の分解に基づく動的計画法を用いるため同様である。さらには、素性の出現を数える際ギャップなどの曖昧さを許すようなカーネル [8] を用いるに到っては、長さ d の素性が、長さ d 以上の領域にまたがりうるため、分解 (4) でさえ用いることができない。Collins は、あらかじめ学習された、局所的な素性のみを用いる別のラベル付け器を用いて、適当な数のラベル付け候補を出力し、それをより広範囲にわたる素性を用いて再ランキングするという方法でこの問題を避けている [17] が、正しいラベル付けが、広範囲の素性に基づく始めの候補の中に含まれているという保証はない。

3. 周辺化ラベル付けパーセプトロン

前節で述べた問題点を解決するために、カーネル法に基づくラベル付け器を提案する。通常、MEMM や CRF は、モデル $P_s(\mathbf{y}|\mathbf{x})$ は対数尤度の和 $\sum_i \log P_s(\mathbf{y}^{(i)}|\mathbf{x}^{(i)})$ を最大化することによって学習する。この目的関数は、ひとつの配列上の全部の非観測変数を正しくラベル付することを目指しているといえる。一方、Kakade ら [7] は多くのタスクにおいて、個々の非観測関数のラベルを正しくラベル付けできればよいという点に注目し、点ごとに周辺化された尤度に基づく目的関数 $\sum_i \sum_t \log P_p(y_t = y_t^{(i)}|\mathbf{x}^{(i)})$ を提案した。ここで、

$$P_p(y_t = y_t^{(i)}|\mathbf{x}^{(i)}) = \sum_{\mathbf{y}: y_t = y_t^{(i)}} P_s(\mathbf{y}^{(i)}|\mathbf{x}^{(i)})$$

は、 t 番目の非観測変数のラベルを $\hat{y}_t \in \Sigma_y$ に固定し、それ以外の非観測変数について周辺化した $P_s(\mathbf{y}|\mathbf{x})$ である。この目的関数は、ももとの配列ごとの尤度に基づく目的関数に匹敵する性能をもつことが示されている [16]。ここで重要な点は、 P_p は y_t にのみ依存するため、ラベル予測も点ごとに行われるということである。我々はこのアイデアと、パーセプトロンに基づくラベル付け器とを組み合わせることで、以下に定義される周辺化ラベル付けパーセプトロンを提案する。

$$\hat{y}_t = \operatorname{argmax}_{\hat{y}_t \in \Sigma_y} \sum_{\mathbf{y}: y_t = \hat{y}_t} P(\mathbf{y}|\mathbf{x}) \langle \mathbf{w}, \Phi(\mathbf{x}, \mathbf{y}) \rangle \quad (5)$$

argmax 操作は \hat{y}_t にのみ依存するので、点ごとのラベル予測が実現されている。基本的な考え方としては以下のとおりである。 y_t を予測する際、ももとのラベル付けパーセプトロンの出力 $\langle \mathbf{w}, \Phi(\mathbf{x}, \mathbf{y}) \rangle$ は、 t 番目の非観測変数のラベルを \hat{y}_t に固定し、そのほかの非観測変数について周辺化される。 $P(\mathbf{y}|\mathbf{x})$ は、非観測変数についての事前分布であり、これはあらかじめ訓練さ

れた MEMM や CRF などでもよい。また、周辺化された素性ベクトル $\sum_{\mathbf{y}: y_t = \hat{y}_t} P(\mathbf{y}|\mathbf{x})\Phi(\mathbf{x}, \mathbf{y})$ は、新しい素性ベクトルと考えることができるので、 $e^{(i)}$ の t 番目の非観測変数についての予測が誤りだったとき、すなわち、 $\hat{y}_t^{(i)} \neq y_t^{(i)}$ のとき、重みベクトル \mathbf{w} は、以下の更新則によって更新される。

$$\mathbf{w}^{new} = \mathbf{w}^{old} + \sum_{\mathbf{y}: y_t = y_t^{(i)}} P(\mathbf{y}|\mathbf{x})\Phi(\mathbf{x}, \mathbf{y}) - \sum_{\mathbf{y}: y_t = \hat{y}_t^{(i)}} P(\mathbf{y}|\mathbf{x})\Phi(\mathbf{x}, \mathbf{y}) \quad (6)$$

次に、周辺化ラベル付けパーセプトロンの双対型を導くことで、これをカーネル化する。 $\phi_f(\mathbf{x}, \mathbf{y}; t)$ を、素性 f が (\mathbf{x}, \mathbf{y}) 内に、 (\mathbf{x}, \mathbf{y}) の t 番目の非観測変数を含む形で出現する回数とし、また $\Phi_f(\mathbf{x}, \mathbf{y}; t)$ をそのベクトル表現とする。まず、ももとの形 (5) を \hat{y}_t に依存しない項を足すことで以下のように変形する。

$$\begin{aligned} \hat{y}_t &= \operatorname{argmax}_{\hat{y}_t \in \Sigma_y} \sum_{\mathbf{y}: y_t = \hat{y}_t} P(\mathbf{y}|\mathbf{x}) \langle \mathbf{w}, \Phi(\mathbf{x}, \mathbf{y}) \rangle \\ &+ \sum_{\hat{y}_t \in \Sigma_y} \sum_{\mathbf{y}: y_t = \hat{y}_t} P(\mathbf{y}|\mathbf{x}) \langle \mathbf{w}, \Phi(\mathbf{x}, \mathbf{y}; t) \rangle \\ &- \sum_{\mathbf{y} \in \Sigma_y^T} P(\mathbf{y}|\mathbf{x}) \langle \mathbf{w}, \Phi(\mathbf{x}, \mathbf{y}) \rangle \\ &= \operatorname{argmax}_{\hat{y}_t \in \Sigma_y} \sum_{\mathbf{y}: y_t = \hat{y}_t} P(\mathbf{y}|\mathbf{x}) \langle \mathbf{w}, \Phi(\mathbf{x}, \mathbf{y}; t) \rangle \quad (7) \end{aligned}$$

更新則は以下のように書き換えられる。

$$\mathbf{w}^{new} = \mathbf{w}^{old} + \sum_{\mathbf{y}: y_t = y_t^{(i)}} P(\mathbf{y}|\mathbf{x})\Phi(\mathbf{x}, \mathbf{y}; t) - \sum_{\mathbf{y}: y_t = \hat{y}_t^{(i)}} P(\mathbf{y}|\mathbf{x})\Phi(\mathbf{x}, \mathbf{y}; t)$$

例の重み α を導入して、 \mathbf{w} を線形和の形で書くと、

$$\mathbf{w} = \sum_{j=1}^{|E|} \sum_{\tau=1}^{T^{(j)}} \sum_{\hat{y}_\tau \in \Sigma_y} \alpha_{j\tau}(\hat{y}_\tau) \sum_{\mathbf{y}: y_\tau = \hat{y}_\tau} P(\mathbf{y}|\mathbf{x})\Phi(\mathbf{x}^{(j)}, \mathbf{y}; \tau) \quad (8)$$

となるので、これを (7) に代入することで、周辺化カーネルラベル付けパーセプトロン

$$\hat{y}_t = \operatorname{argmax}_{\hat{y}_t \in \Sigma_y} \sum_{j=1}^{|E|} \sum_{\tau=1}^{T^{(j)}} \sum_{\hat{y}_\tau \in \Sigma_y} \alpha_{j\tau}(\hat{y}_\tau) K(\mathbf{x}^{(j)}, \mathbf{x}, \tau, t, \hat{y}_\tau, \hat{y}_t) \quad (9)$$

を得る。ここで、

$$\begin{aligned} K(\mathbf{x}^{(j)}, \mathbf{x}, \tau, t, \hat{y}_\tau, \hat{y}_t) &= \sum_{\mathbf{y}: y_\tau = \hat{y}_\tau} \sum_{\mathbf{y}': y'_t = \hat{y}_t} P(\mathbf{y}|\mathbf{x}^{(j)}) P(\mathbf{y}'|\mathbf{x}) \langle \Phi(\mathbf{x}^{(j)}, \mathbf{y}; \tau), \Phi(\mathbf{x}, \mathbf{y}'; t) \rangle \quad (10) \end{aligned}$$

は 2 つの例の τ 番目と t 番目の非観測変数のラベルをそれぞれ固定し周辺化されたカーネル関数であるため周辺化カーネル [6] と呼ばれる。アルゴリズムを図 3 に示す。

4. 周辺化カーネル関数の計算

本節では配列及び順序木のラベル付けのための周辺化カーネル (10) を提案する。近年、さまざまな構造データに対するカーネル関数は精力的に研究されている [14] が、これをラベル付けに使えるよう拡張する。ここでは簡単のため、事前分布 $P(\mathbf{y}|\mathbf{x})$ は、 $P(\mathbf{y}|\mathbf{x}) = \prod_{t=1}^T P(y_t|x_t)$ のように分解できるとするが、本稿の結果は CRF などのより一般的なモデルを事前分布として用いる場合にも拡張可能である。また、表記上の簡単

```

// Initialization
 $\alpha := \mathbf{0}$ 
// Training
for round = 1, ..., max_round,
  for i = 1, ..., |E|,
    // kernel computation
    for j = 1, ..., |E|,
      for t = 1, ..., T(i), and  $\tau = 1, ..., T(j)$ ,
        for  $\check{y}_\tau \in \Sigma_y$ , and  $\check{y}_t \in \Sigma_y$ ,
          compute  $K(\mathbf{x}^{(j)}, \mathbf{x}^{(i)}, \tau, t, \check{y}_\tau, \check{y}_t)$ 
    // prediction & update
    for t = 1, ..., T(i),
      predict  $\hat{y}_t^{(i)}$  by using Equation (9)
      if  $\hat{y}_t^{(i)} \neq y_t^{(i)}$ ,
         $\alpha_{it}(y_t^{(i)}) := \alpha_{it}(y_t^{(i)}) + 1$ 
         $\alpha_{it}(\hat{y}_t^{(i)}) := \alpha_{it}(\hat{y}_t^{(i)}) - 1$ 

```

図3 周辺化カーネルラベル付けパーセプトロン

のため本節の残りの部分では、 $K(\tau, t, \check{y}_\tau, \check{y}_t)$ を i 番目の例と j 番目の例が与えられたときのカーネル関数とする。我々のカーネル関数を計算する上で重要な点は、すべての素性が、より小さい素性を組み合わせることで表現できる点である。たとえば、図4で素性 (b) が、その一番右端を t 番目の非観測変数に、また素性 (c) が、その一番左端を t 番目の非観測変数に置いて出現しているとする、素性 (b) と (c) を組み合わせた素性 (a) は、その2番目の変数を t 番目の非観測変数に置いて出現していることになる。従って、本節で提案する周辺化カーネルはすべて上流カーネル $K_U(\tau, t)$ 、下流カーネル $K_D(\tau, t)$ 、点カーネル $K_P(\tau, t, \check{y}_\tau, \check{y}_t)$ の3つに分解できるとする。

$$K(\tau, t, \check{y}_\tau, \check{y}_t) = K_U(\tau, t) \cdot K_P(\tau, t, \check{y}_\tau, \check{y}_t) \cdot K_D(\tau, t)$$

$$K_U(\tau, t) = \sum_{\mathbf{y}_U(\tau)} \sum_{\mathbf{y}'_U(t)} P(\mathbf{y}_U(\tau) | \mathbf{x}_U^{(j)}(\tau)) P(\mathbf{y}'_U(t) | \mathbf{x}_U^{(i)}(t)) \cdot \langle \Phi(\mathbf{x}_U^{(j)}(\tau), \mathbf{y}_U(\tau); \tau), \Phi(\mathbf{x}_U^{(i)}(t), \mathbf{y}'_U(t); t) \rangle \quad (11)$$

$$K_D(\tau, t) = \sum_{\mathbf{y}_D(\tau)} \sum_{\mathbf{y}'_D(t)} P(\mathbf{y}_D(\tau) | \mathbf{x}_D^{(j)}(\tau)) P(\mathbf{y}'_D(t) | \mathbf{x}_D^{(i)}(t)) \cdot \langle \Phi(\mathbf{x}_D^{(j)}(\tau), \mathbf{y}_D(\tau); \tau), \Phi(\mathbf{x}_D^{(i)}(t), \mathbf{y}'_D(t); t) \rangle \quad (12)$$

$$K_P(\tau, t, \check{y}_\tau, \check{y}_t) = \frac{P(\check{y}_\tau | \mathbf{x}_\tau^{(j)}) P(\check{y}_t | \mathbf{x}_t^{(i)}) \langle \Phi(\mathbf{x}_\tau^{(j)}, \check{y}_\tau; \tau), \Phi(\mathbf{x}_t^{(i)}, \check{y}_t; t) \rangle}{\left(\sum_{\check{y}_\tau} \sum_{\check{y}_t} P(\check{y}_\tau | \mathbf{x}_\tau^{(j)}) P(\check{y}_t | \mathbf{x}_t^{(i)}) \langle \Phi(\mathbf{x}_\tau^{(j)}, \check{y}_\tau; \tau), \Phi(\mathbf{x}_t^{(i)}, \check{y}_t; t) \rangle \right)^2} \quad (13)$$

ここで、 $\mathbf{y}_U(\tau)$ と $\mathbf{y}_D(\tau)$ は、それぞれ τ 番目から上流あるいは下流に位置する非観測変数の集合とする ($\mathbf{y}_U(\tau) \cap \mathbf{y}_D(\tau) = \{y_\tau\}$)。また $K_P(\tau, t, \check{y}_\tau, \check{y}_t)$ の分母は $K_U(\tau, t)$ と $K_D(\tau, t)$ の重なり部分を打ち消すための項である。 K_U や K_D における周辺化を陽に計算することは現実的ではないが、多くの場合、動的計画法などをつかうことで効率的に計算可能である。この分解によって、 i 番目と j 番目の例が与えられたときに、全ての $(\tau, t, \check{y}_\tau, \check{y}_t)$ のとりうる値に対し周辺化カーネルを1回の計算で計算できる。

4.1 配列のラベル付けのための周辺化カーネル

ここでは配列における広範囲にわたる依存関係を考慮するた

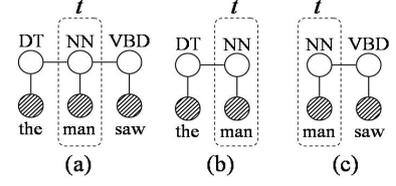


図4 上流の素性と下流の素性を組み合わせることで、大きい素性が構成される例

め、素性は図2に示すような、任意の長さ連結した非観測変数と観測変数のペアとする。素性のサイズに制限を設けないため、様々なサイズの素性が入り混じることになるが、そこで、混合定数 $c > 0$ を用いて、 ϕ_f を属性 f のサイズ d に応じて重み付ける。例えば f が (\mathbf{x}, \mathbf{y}) に N 回出現したとすると、 $\phi_f(\mathbf{x}, \mathbf{y})$ は $N \cdot c^d$ と定義する。 $K_U(\tau, t)$ や $K_D(\tau, t)$ において、ある点 τ より上流の変数 $\mathbf{y}_U(\tau)$ とは、その点とそれより左にある変数のことである。同様に下流の変数 $\mathbf{y}_D(\tau)$ とは、その点とそれより右にある変数のことであるとする。アルゴリズムは、丁度 HMM の前向き・後ろ向きアルゴリズムのように、それぞれ K_U と K_D を計算するための2つの動的計画法のループからなる。

$$K_U(\tau, t) = c^2 k(x_\tau^{(j)}, x_t^{(i)}) (K_U(\tau - 1, t - 1) + 1) \quad (14)$$

$$K_D(\tau, t) = c^2 k(x_\tau^{(j)}, x_t^{(i)}) (K_D(\tau + 1, t + 1) + 1) \quad (15)$$

ここで、

$$k(x_\tau^{(j)}, x_t^{(i)}) = \sum_{\check{y}_\tau \in \Sigma_y} \sum_{\check{y}_t \in \Sigma_y} P(\check{y}_\tau | x_\tau^{(j)}) P(\check{y}_t | x_t^{(i)}) \cdot k_y(\check{y}_\tau, \check{y}_t) k_x(x_\tau^{(j)}, x_t^{(i)})$$

とする。境界条件は $K_U(\tau, 0) = K_U(0, t) = K_D(\tau, |T^{(i)}| + 1) = K_D(|T^{(j)}| + 1, t) = 0$ である。 K_P は以下によって計算できる。

$$K_P(\tau, t, \check{y}_\tau, \check{y}_t) = \frac{c^2 P(\check{y}_\tau | x_\tau^{(j)}) P(\check{y}_t | x_t^{(i)}) k_y(\check{y}_\tau, \check{y}_t) k_x(x_\tau^{(j)}, x_t^{(i)})}{(c^2 k(x_\tau^{(j)}, x_t^{(i)}))^2} \quad (16)$$

ここで、 k_x と k_y は2つの引数が同一であるときに1を、そうでないときに0を返す関数とするが、ここを2つの変数の間のカーネル関数で置き換えることも可能である。明らかに、このカーネル関数の計算量は $O(T^{(i)} T^{(j)})$ である。^(注3)

4.2 順序木のラベル付けのための周辺化カーネル

順序木 (図5) のラベル付けのためのカーネルを、順序木カーネル [11] を拡張することによって設計する。素性は図5に示すような任意の順序木の形をしているとする。配列ラベル付けのときと同様、混合定数 c を用いて $\phi_f(\mathbf{x}, \mathbf{y})$ を定義する。 $K_U(\tau, t)$ や $K_D(\tau, t)$ において、ある点 τ より上流の変数 $\mathbf{y}_U(\tau)$ とは、その点とそれより外側にある変数のことであり、下流の変数 $\mathbf{y}_D(\tau)$ とは、その点とそれより内側にある変数のことであるとする。アルゴリズムは、丁度、確率文脈自由文法の内側・外側アルゴリズムのように、それぞれ K_U と K_D を計算するための2つの動的計画法のループからなる。 $ch(\tau, v)$ を τ 番目のノー

(注3) : ここでの結果は [8] にみられるようなギャップを許した曖昧なマッチングを許すように拡張可能である。詳しくは [15] を参照されたい。

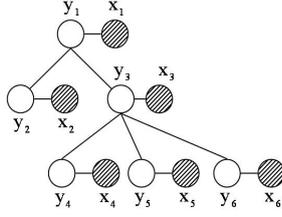


図5 順序木の例：黒いノードは観測変数を，白いノードは非観測変数を表す。

ドの v 番目の子を表すインデクスとし， $pa(\tau)$ を τ 番目のノードの親ノードを表すインデクスとする．また， $\#ch(\tau)$ を τ 番目のノードの子供の数とし， $chID(\tau)$ を， τ 番目のノードがその親ノードの何番目の子ノードであるかを表すインデクスとする．

順序木カーネル [11] の再帰的計算と同様にして， K_D は後行順で動的計画法により，

$$K_D(\tau, t) = c^2 k(x_\tau^{(j)}, x_t^{(i)}) S_F(\tau, t, \#ch(\tau), \#ch(t))$$

によって計算される．ここで S_F も以下によって再帰的に定義される．

$$\begin{aligned} S_F(\tau, t, v, u) &= K_D(ch(\tau, v), ch(t, u)) S_F(\tau, t, v-1, u-1) \\ &\quad + S_F(\tau, t, v-1, u) + S_F(\tau, t, v, u-1) \\ &\quad - S_F(\tau, t, v-1, u-1) \end{aligned}$$

境界条件は $S_F(\tau, t, 0, u) = S_F(\tau, t, v, 0) = 1$ である．直感的には $S_F(\tau, t, v, u)$ は $ch(\tau, 1)$ 番目から $ch(\tau, v)$ 番目の子集合と， $ch(t, 1)$ 番目から $ch(t, u)$ 番目まで子集合の全ての，順序を保存したマッチングについての和をとっているといえる．

同時に，あとで K_U を計算するときのために， $ch(\tau, v)$ 番目から $ch(\tau, \#ch(\tau))$ 番目の子集合と， $ch(t, u)$ 番目から $ch(t, \#ch(t))$ 番目まで子集合の全てのマッチングについての和をとったものを表す $S_B(\tau, t, v, u)$ を

$$\begin{aligned} S_B(\tau, t, v, u) &= K_D(ch(\tau, v), ch(t, u)) S_B(\tau, t, v+1, u+1) \\ &\quad + S_B(\tau, t, v+1, u) + S_B(\tau, t, v, u+1) \\ &\quad - S_B(\tau, t, v+1, u+1) \end{aligned}$$

と定義する．境界条件は $S_B(\tau, t, \#ch(\tau) + 1, u) = S_B(\tau, t, v, \#ch(t) + 1) = 1$ である．

K_U は先行順の動的計画法で計算される． $K_U(\tau, t)$ は，親ノード同士の上流カーネルと，兄弟同士の下流カーネルを組み合わせることで計算できる (図 6)．

$$\begin{aligned} K_U(\tau, t) &= c^2 k(x_\tau^{(j)}, x_t^{(i)}) \left(1 + K_U(pa(\tau), pa(t)) \right. \\ &\quad \cdot S_F(pa(\tau), pa(t), chID(\tau) - 1, chID(t) - 1) \\ &\quad \left. \cdot S_B(pa(\tau), pa(t), chID(\tau) + 1, chID(t) + 1) \right) \end{aligned}$$

K_P は (16) と同様である．[11] と同様の解析によってこのカーネルの計算量は $O(T^{(i)} T^{(j)})$ であることが容易に分かる．

5. 実験

本節では自然言語処理の課題 (固有表現抽出と製品使用情報

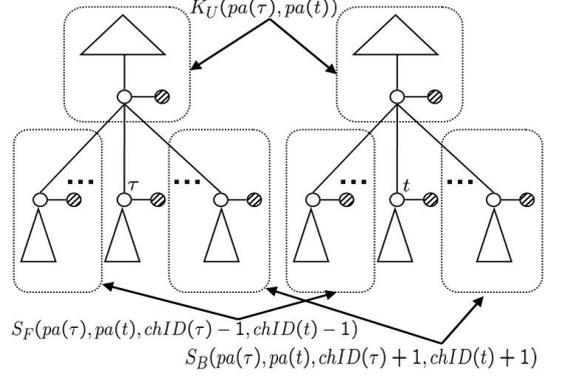


図6 順序木における $K_U(\tau, t)$ の計算のイメージ。

表1 固有表現抽出結果 (括弧内は標準偏差)

	精度	適合率	再現率	F1
配列カーネル	88.7% (3.4)	49.0% (6.0)	23.1% (8.1)	30.5 (6.7)
HM パーセプトロン	80.2% (11.5)	23.8% (14.6)	17.9% (3.0)	18.6 (5.2)

表2 製品使用情報抽出結果 (括弧内は標準偏差)

	精度	適合率	再現率	F1
SEQUENCE KERNEL	89.7% (2.0)	52.2% (9.5)	29.6% (4.0)	37.5(4.1)
TREE KERNEL	89.9% (2.7)	51.4% (10.9)	32.5% (14.4)	38.9 (12.1)
HM-PERCEPTRON	89.7% (1.8)	51.5% (8.5)	24.0% (21.4)	28.9(20.3)

抽出) での実験結果を示し，提案手法により効率的に扱うことが可能になった大きな構造素性の有用性を検証する．本実験では周辺化ラベル付けパーセプトロンと HM パーセプトロンの性能比較を行った．

5.1 固有表現抽出

固有表現抽出は情報抽出課題の一つで，文書内の人名・地名・組織名などを識別する問題である．実験では，CoNLL2002 固有表現抽出課題で提供されたスペイン語のデータの最初の 300 文 (8,541 単語) を使用した．この課題は，各単語に固有表現の種類を示すラベルを付与する課題である．ラベルの数は，各固有表現の種類と表現の境界種別の組み合わせ 8 ラベルと非固有表現を示すダミーラベルを合わせた 9 種類 ($|\Sigma_y| = 9$) である．観測変数のカーネルは単語とそのスペリング (詳細は [16] の S2 素性を参照のこと) を素性とした 2 次の多項式カーネルを使用した．ただし，HM パーセプトロンは [4] に倣い前後 1 単語の素性も含めた．周辺化ラベル付けパーセプトロンでは配列ラベル付けカーネル (配列カーネル) を用いた．混合定数は事前実験により求めた $c = 1$ とした．また，非観測変数に対する事前分布は一様分布 ($P(y_t|x_t) = 1/|\Sigma_y|$) として設計した．性能評価にはダミーラベルを含めたラベル付けの精度，固有表現ラベルの適合率・再現率，適合率と再現率の調和平均値である F1 尺度を用いた．

表 1 は 5 分割交差検定の結果である．提案手法は HM パーセプトロンを大きく上回る性能を示した．配列ラベル付けカーネルを用いた提案手法が良い性能を示した理由は，固有表現抽出の性質によるものと考えられる．固有表現抽出課題は自然言語処理の問題の中で広い文脈情報を考慮する必要がある問題の一つである．例えば，「ホワイトハウス」は文脈により場所または組織を示すが異なる．この実験結果は，人手による素性の選択なしに幅広い文脈情報を効率的に扱うことが出来るカーネル化されたラベル付け学習アルゴリズムの優位性を示すものと言える．

