

4-16 系列ラベリングの多層化

東 藍

松本 裕治

奈良先端科学技術大学院大学

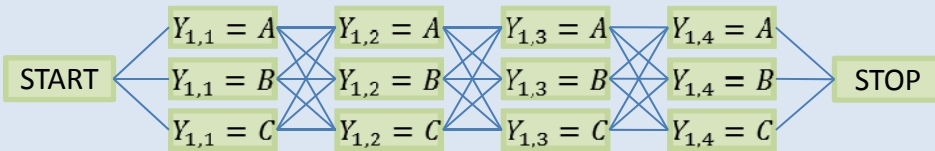


研究の背景

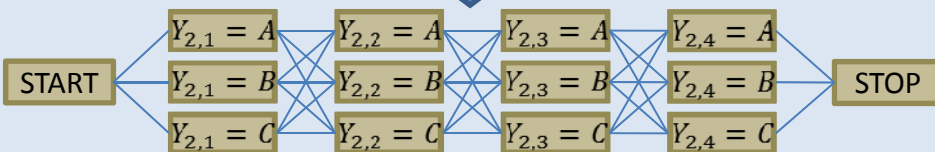
系列ラベリングが様々なタスクで顕著な成果を挙げている
→ 自然言語処理 → 音声認識 → 遺伝子解析

複数の系列ラベリングの段階適用の必要性
(e.g. 入力自然言語文に対する品詞付与と基本句同定)

系列 Y_1 の推定



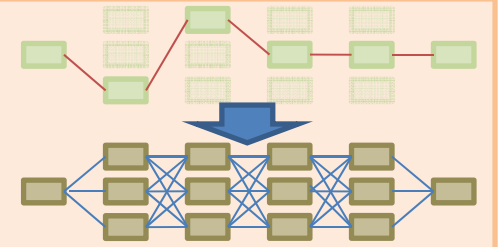
系列 Y_2 の推定



既存手法と問題点

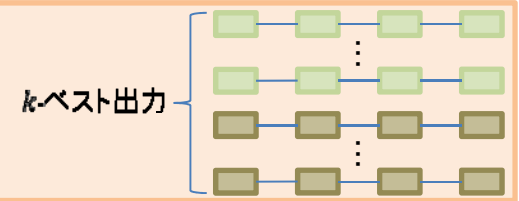
1-ベストパイプライン

実装コストが低い
解析誤りが後続で伝搬・拡大



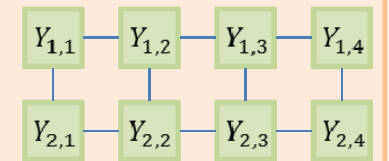
k-ベスト, リランキング

解析誤りの伝搬を減じるが
同時学習は考慮されない



グラフィカルモデルによる同時学習

同時学習可能だが
exact inference が不可能
近似手法が必要

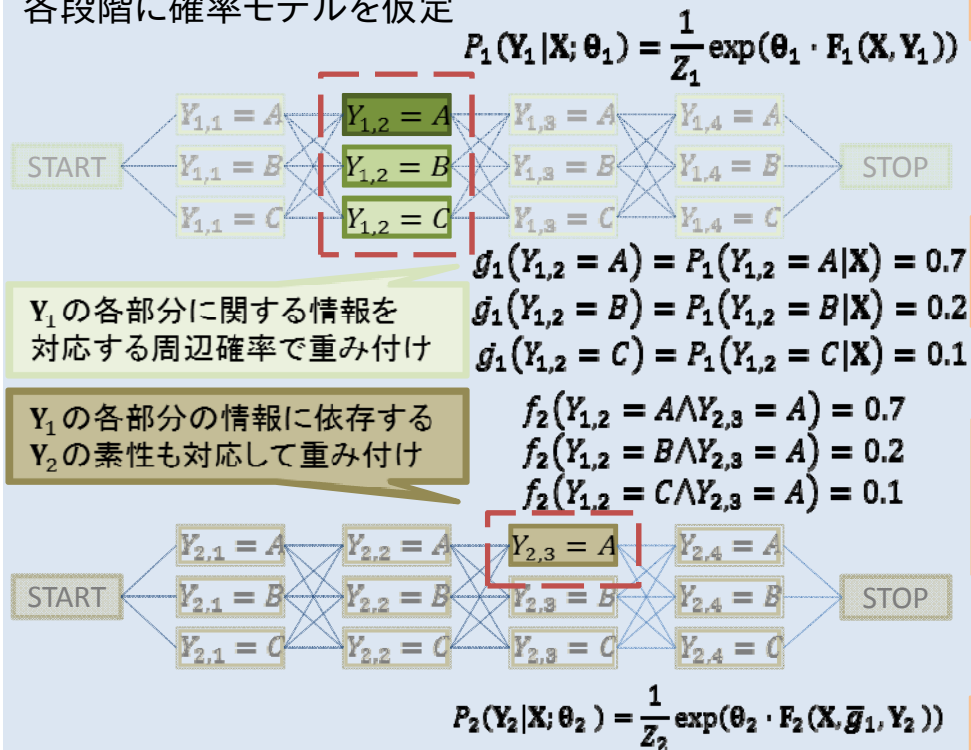


研究の目的

解析誤りの伝搬を防ぎ、かつ動的計画法による exact inference の利点を維持した同時学習手法の提案

採用する既存手法(素性の周辺化)

[Bunescu, 2008] のアイデアを採用
各段階に確率モデルを仮定



実験

英文に対する品詞タグ付と基本句同定の2つの系列ラベリングタスクに提案手法を適用し提案手法の有効性を検証

	1-ベストパイプライン	周辺化素性を利用した1-ベストパイプライン ([Bunescu, 2008] 相当)	提案手法 (周辺化素性+同時最適化)
品詞タグ付	0.956	(0.956)	0.958
基本句同定	0.921	0.927	0.931

実験により提案手法の有効性を確認. 今後は既存手法との精緻な比較や他タスクへの適用が課題.

本研究の手法

パラメタ θ_1 により周辺確率が変化

$\Rightarrow \bar{g}_1$ および f_2 の値も変化 $\Rightarrow \theta_1$ に依存した関数とみなす

$$\begin{cases} g_1(Y_{1,2} = A; \theta_1) \\ g_1(Y_{1,2} = B; \theta_1) \\ g_1(Y_{1,2} = C; \theta_1) \\ \vdots \end{cases} \quad \begin{cases} f_2(Y_{1,2} = A \wedge Y_{2,3} = A; \theta_1) \\ f_2(Y_{1,2} = B \wedge Y_{2,3} = A; \theta_1) \\ f_2(Y_{1,2} = C \wedge Y_{2,3} = A; \theta_1) \\ \vdots \end{cases}$$

後続の段階のモデルも θ_1 に間接的に依存しているとみなして最適化

$$P_2(Y_2|X; \theta_2, \theta_1) = \frac{1}{Z_2} \exp(\theta_2 \cdot F_2(X, \bar{g}_1(\theta_1), Y_2))$$

適当な目標関数 (e.g. P_1, P_2 の対数尤度関数の和) で同時最適化
最適化の際にはパラメタ勾配を
合成関数の微分法で分解して計算

$$\frac{\partial P_2}{\partial \theta_1} = \frac{\partial P_2}{\partial f_2} \cdot \frac{\partial f_2}{\partial \bar{g}_1} \cdot \frac{\partial \bar{g}_1}{\partial \theta_1} \quad \Rightarrow \quad \text{誤差逆伝搬法に類似}$$

パラメタ勾配の各因子は動的計画法で exact に計算可能