

2중기능조건부우연마당을 리용한 고유실체인식정확도 개선의 한가지 방법

우 준 민

고유실체인식에 많이 리용되는 통계모형들[3]로는 숨은마르코브모형(HMM), 최대엔트로피마르코브모형(MEMM), 조건부우연마당(CRF) 등이 있다.

선행연구[1]에서는 조선어고유실체인식정확도를 개선할수 있게 실체지시단어모임을 리용한 특징들을 추가하고 CRF를 리용하여 인식을 진행하였다.

선행연구[2]에서는 클래스 n 그램언어모형을 구축하고 그것에 의한 어휘분리방법으로 형태부해석과 고유실체인식의 일체화를 실현하는 방법으로 조선어고유실체인식을 진행하였다. 이 방법은 형태부해석과 고유실체인식을 동시에 진행한다는 우점을 가지고있지만 고유실체인식의 고유한 언어적정보들을 모형으로 부호화하기가 힘들고 훈련본문에서 매우 적은 빈도수를 가지고 나타나거나 전혀 나타나지 않는 고유실체의 특성을 원만히 고려하지 못하는 결함을 가지고있다.

이러한 문제를 해결하기 위하여 다음과 같은 문제를 설정한다.

첫째로, 조선어고유실체의 문맥특성을 고려한 2중기능CRF모형을 제기한다.

둘째로, 2중기능CRF모형의 학습과 추론방법을 제기하고 모형을 리용한 고유실체인식방법을 제기한다.

셋째로, 실험을 통하여 제안한 방법의 효과성을 검증한다.

1. 2중기능조건부우연마당을 리용한 조선어고유실체인식방법

관측렬 $x = ([x_t^T]_{t=1}^T)^T$ 이 주어졌을 때 전체 렬에 대한 클래스값 $h \in H$ 와 요구되는 표식렬 $y = [y_t]_{t=1}^T$ 들을 동시에 얻는 문제를 생각하자.

이 문제를 해결하기 위하여 다음과 같은 형식의 CRF에 기초한 모형을 설정한다.(그림)

$$p(h, y | x) = \frac{1}{Z(x)} \exp \left[\sum_{t=2}^T \phi_t(h, y_t, y_{t-1}, x_t) + \phi_1(h, y_1, x_1) \right]$$

여기서

$$Z(x) = \sum_h \sum_y \exp \left[\sum_{t=2}^T \phi_t(h, y_t, y_{t-1}, x_t) + \phi_1(h, y_1, x_1) \right]$$

이다.

론문에서는 다음과 같은 형식의 클래스조건부포텐살함수들의 모임을 도입한다.

$$\begin{aligned} \phi_t(h, y_t, y_{t-1}, x_t) &= \phi_t^1(h, y_t, x_t) + \phi_t^2(h, y_t, y_{t-1}) \\ \phi_1(h, y_1, x_1) &= \phi_1^1(h, y_1, x_1) + \phi_1^2(h, y_1) \end{aligned}$$

여기서 $\phi_t^2(h, y_t, y_{t-1})$ 과 $\phi_t^1(h, y_t, x_t)$ 는 클라스조건부적인 전이포텐셜함수와 한번수포텐셜함수이다.

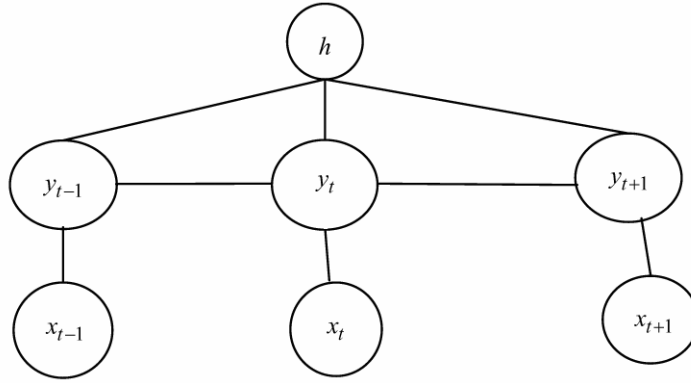


그림. 2중기능조건부우연마당모형

선형제한을 주면 다음과 같이 쓸수 있다.

$$\begin{aligned}\phi_t^1(h, y_t, x_t) &= \sum_{i=1}^K \delta(y_t - i) \omega_i^h \cdot x_t \\ \phi_t^2(h, y_t, y_{t-1}) &= \sum_{i=1}^K \sum_{j=1}^K \delta(y_t - j) \delta(y_{t-1} - i) \xi_{ij}^h \\ \phi_1^2(h, y_1) &= \sum_{i=1}^K \delta(y_1 - j) \xi_i^h \\ \phi_1(h, y_1, x_1) &= \phi_1^1(h, y_1, x_1) + \phi_1^2(h, y_1)\end{aligned}$$

모형학습은 모형포텐셜함수들의 파라메터들에 대한 추정으로 이루어진다.

클라스조건포텐셜함수를 고찰할 때 이것은 다음과 같은 파라메터들의 추정으로 된다.

$$\{\omega_i^h\}_{h \in H, i \in y}, \{\xi_{ij}^h\}_{h \in H, i, j \in y}, \{\xi_i^h\}_{h \in H, i \in y}$$

모형학습을 위하여 다음과 같은 로그우도최량화문제를 풀어보자.

$$\log p(h, y | x) = \sum_{t=2}^T \phi_t(h, y_t, y_{t-1}, x_t) + \phi_1(h, y_1, x_1) - \log Z(x)$$

여기서 모형의 로그우도를 계산하는것은 y 와 h 의 전체 값에 대하여 합인 $Z(x)$ 에 대한 계산을 요구한다. 이 계산을 쉽게 하기 위하여 다음과 같은 형식으로 $Z(x)$ 를 표현한다.

$$Z(x) = \sum_{h \in H} Z(x | h)$$

여기서

$$Z(x | h) = \sum_y \exp \left[\sum_{t=2}^T \phi_t(h, y_t, y_{t-1}, x_t) + \phi_1(h, y_1, x_1) \right]$$

이다.

이것은 앞방향—뒤방향알고리즘에 의하여 쉽게 계산할수 있다.

$$Z(\mathbf{x} | h) = \sum_{j=1}^K \alpha_t(h, j)$$

웃식에서 클라스조건부앞방향확률 $\alpha_t(h, j)$ 는 다음과 같이 계산된다.

$$\alpha_t(h, j) = \sum_{j=1}^K \alpha_{t-1}(h, j) \exp\{\phi_t(h, y_t = j, y_{t-1} = i, \mathbf{x}_t)\}, t \geq 2$$

초기값은

$$\alpha_1(h, j) = \exp\{\phi_1(h, y_1 = j, \mathbf{x}_1)\}$$

이다.

N 개의 훈련렬에 대하여 모형의 로그우도를 최대화하는 파라미터들은 L-BFGS와 같은 방법으로 구할수 있다.

렬에 대한 분류 다시말하여 관측렬 \mathbf{x} 가 주어졌을 때 최량클라스표식 h 는 다음과 같은 최량화문제를 풀어 구한다.

$$\hat{h} = \arg \max_{h \in H} p(h | \mathbf{x})$$

여기서

$$p(h | \mathbf{x}) = \sum_y p(h, \mathbf{y} | \mathbf{x}) = \frac{\sum_y \exp \left[\sum_{t=2}^T \phi_t(h, y_t, y_{t-1}, \mathbf{x}_t) + \phi_1(h, y_1, \mathbf{x}_1) \right]}{Z(\mathbf{x})} = \frac{Z(\mathbf{x} | h)}{Z(\mathbf{x})}$$

이다.

따라서 이 문제는 다음과 같이 쓸수 있다.

$$\hat{h} = \arg \max_{h \in H} Z(\mathbf{x} | h)$$

관측렬 \mathbf{x} 가 주어졌을 때 렬토막화는 다음과 같은 최량화문제로 구성되는 렬복호화문제로 볼수 있다.

$$\hat{y} = \arg \max_y \log p(\mathbf{y} | \mathbf{x})$$

이렇게 구성한 2중기능조건부우연마당을 리용하면 고유실체와 형태부들이 결합하여 생기는 고유실체들을 인식할수 있다.

이때 x 는 단어렬이고 y 는 형태부렬에 대한 표기이며 h 는 그 형태부렬에 대한 인식된 고유실체인식표식이다.

다음으로 고유실체인식을 위하여 조건부우연마당에서 리용하는 징표형타선택문제를 생각하자.

우선 관측렬에 대한 징표선택에서 가장 중요한것은 이전의 표식이다.

실례로 《황해북도 사리원시 ××협동농장》이라고 하면 《××협동농장》이 조직실체인것으로 해서 전체적으로 조직실체지시단어가 된다.

또한 중요한 징표의 하나는 의미정보이다.

실례로 《황해남도 옹진군 xx 수산사업소》라고 하면 《xx 수산사업소》의 의미범주가 조직에 속하기때문에 전체적으로 조직실체지시단어가 된다.

조선어고유실체인식에 리용하는 징표형태는 표 1과 같다.

표 1. 조선어고유실체인식에 리용하는 징표형태

징표형태	징표의 의미
CUR_ORG	현재 단어가 조직실체이고 렬의 마지막단어인 경우
CUR_MEAN_ORG	현재의 단어의 의미범주가 조직의 하위개념이고 렬의 마지막단어인 경우
PREV_ORG	이전단어의 표식과 현재단어의 표식
PREV_MEAN_ORG	이전단어의 의미와 현재단어의 의미
PREV_TAG_MEAN_ORG	이전단어의 표식과 현재단어의 의미

이와 같은 형태의 징표를 장소와 사람범주에 대하여 선택하고 2중기능CRF를 리용하여 단어렬에 대하여 고유실체인식을 진행하면 우에서 제기한 문제가 해결될 수 있다.

2. 실험 및 성능평가

숨은마르코프모형과 규칙에 기초한 고유실체인식기와의 인식성능비교를 진행한 결과는 표 2와 같다. 실험에서는 확장된 고유실체지시단어모임을 리용하여 고유실체인식을 진행하였다.

표 2. 인식성능비교

구분	적중률	완전률	F-척도
HMM	88.78	82.35	85.44
선행방법[2]	98.9	69	83.86
제안방법	93.75	98.14	91.39

표 2에서 보는바와 같이 추출규칙에 기초한 방법은 적중률은 높지만 완전률이 낮다는 것을 알 수 있다. 그것은 추출규칙에 반영하지 못한 고유실체들이 존재하기때문이다. 한편 최대엔트로피모형에 기초한 방법은 숨은마르코프모형에 기초한 방법보다 적중률과 완전률을 둘 다 개선하였다. 논문에서 제기한 방법은 또한 추출규칙에 기초한 방법에 비하여 적중률은 약간 떨어지지만 완전률을 훨씬 높였으며 전체적인 체계의 성능평가지표인 F-척도를 개선하였다.

맺 는 말

조선어문맥특성을 반영하여 고유실체인식정확도를 개선할 수 있는 2중기능조건부우연마당을 설계하고 그것을 리용한 고유실체인식방법을 제기하고 효과성을 검증하였다.

참 고 문 헌

- [1] 김일성종합대학창립70돛기념 전국부문별 과학기술토론회논문집(정보, 자동화), 31, 주체105(2016).
- [2] 홍충성; 정보과학, 4, 45, 주체104(2015).
- [3] D. Kaur, V. Gupta; International Journal of Computer Science Issues, 7, 239, 2010.

주체107(2018)년 11월 5일 원고접수

A Method of Improving the Accuracy of the Named Entity Recognition Using Dual-functional CRF

U Jun Min

In this paper we designed the dual-functional CRF by taking into account Korean named entity characteristics and proposed a method of improving the performance of NER systems using it.

Key words: named entity recognition, conditional random field