

МФТИ

Юсупов Идрис (i.yusupov@phystech.edu)

РАСПОЗНАВАНИЕ ИМЕНОВАННЫХ СУЩНОСТЕЙ С ИСПОЛЬЗОВАНИЕМ
СИНТАКТИКО-СЕМАНТИЧЕСКИХ ПРИЗНАКОВ И НЕЙРОСЕТЕЙ

<h1>1. Цель работы</h1> <p>Разработать новостной агрегатор, который соответствует следующим требованиям:</p> <ul style="list-style-type: none">автоматическая классификация новости к городу,всё должно происходить без участия человека, но должна быть возможность вмешаться,пока для 5 городов.		<h1>2. Общая схема</h1>		
<h1>3. Парсер</h1> <p>Извлекает из текста:</p> <ul style="list-style-type: none">длиннейшие последовательности слов с большой буквы,однокоренные с городом слова(см. правила в словаре),слова в кавычках.	<h1>4. Вектор признаков</h1> <ul style="list-style-type: none">Мешок слов,доменное имя источника новости,фильтруется с помощью словаря.			
<h1>5. Словарь</h1> <p>В словаре хранятся токены и правила.</p> <p>Токен:</p> <ul style="list-style-type: none">нормализован(стемминг),относится к городам. <p>Правило:</p> <ul style="list-style-type: none">на основе регулярных выражений,может отображаться в токен,относится к городам. <p>Начальные данные извлекаются из OpenStreetMap(названия улиц, достопримечательности и т.п).</p>		<h1>6. Классификаторы</h1> <p>В качестве классификатора используется Multinomial Naive Bayes with ROSE-smoothing[1].</p> <p><i>Классификатор по городам</i></p> <ul style="list-style-type: none">"один против всех",количество классификаторов получается равным количеству городов. <p><i>Классификатор выбросов</i></p> <ul style="list-style-type: none">Необходимо определять новость, которая не относится к представленным городам.Для этого используем классификатор по 2м классам - "относится к городу" и "НЕ относится к городу".		<h1>7. Извлечение отношений</h1> <ul style="list-style-type: none">В классифицированном документе находятся пары вида(ГОРОД, СЛОВО_ИЗ_СЛОВАРЯ_ГОРОДА).По найденным вхождениям создается <i>вектор признаков</i>(левый и правый контексты, позиция слова и др.).На основе сформированных вектров признаков создаются паттерны.По паттернам находим новые слова. Похоже на DIPRE[2].
<h1>8. Результаты</h1>				
<h1>9. Заключение</h1> <div><div><p><i>Разработана расширяемая система:</i></p><ul style="list-style-type: none">с использованием машинного обучения,с начальным заполнением словаря геоданными из OpenStreetMap,оценка точности: > 85%,сайт: rbcitynews.ru.</div><div><p><i>В будущем планируется:</i></p><ul style="list-style-type: none">добавление новых городов,извлечение парсером имён собственных,кластеризация токенов в словаре.</div></div>				
<h1>Благодарности</h1> <ul style="list-style-type: none">Научному руководителю Макееву Г.А.Компании "Техинформ".		<h1>Ссылки</h1> <ul style="list-style-type: none">1. Alexander Y. Liu and Cheryl E. Martin. 2011. Smoothing multinomial naïve bayes in the presence of imbalance2. Sergey Brin. 1998. Extracting Patterns and Relations from the World Wide Web.3. http://www.openstreetmap.org/4. http://mit.spbau.ru/files/datamining1009.pdf		