

Обучение КОНЕЧНЫХ АВТОМАТОВ

Автор: Степанов Всеволод
Руководитель: Кураленок Игорь Евгеньевич

5 семестр

Цель

Разработать алгоритмы обучения конечных автоматов для задач классификации и регрессии

Конечные автоматы

- Простая структура
- Легкая интерпретируемость

- У обычных автоматов состояния либо терминальные либо нет
- Сопоставим каждому состоянию число
- Значение автомата на строке:
 - DFA: значение в конечной вершине
 - PFA: матожидание значения

Общая схема

1. Расширение словаря, снижение зависимости между символами
 2. Обучение автомата
 3. Gradient boosting
- Алгоритм не требует каких-либо знаний предметной области или специализированного подхода

Детерминированный автомат

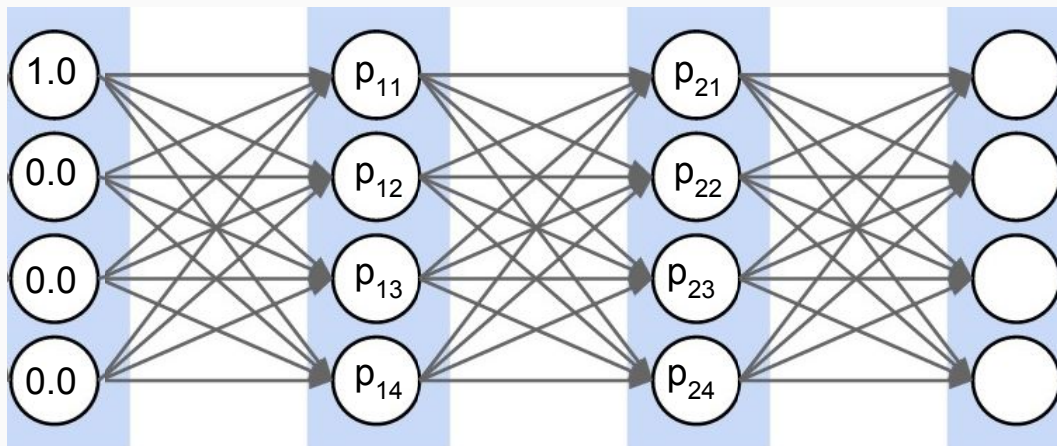
- Жадное итеративное построение
- Модификации: добавить/удалить ребро, разделить состояние
- Минимизируем дисперсию, как в деревьях

Детерминированный автомат

- На маленьких датасетах хорошо работает
- На больших слишком долго

Вероятностный автомат

- Для каждой буквы целая матрица переходов
- Можно представить как полносвязную нейронную сеть
- Обучаем градиентным спуском



Результаты

1. Splice-junction dataset

- a. DFA: error rate ~ 0.055 .
- b. PFA: error rate ~ 0.2
- c. Деревья: $\sim 0.05-0.4$.

2. IMDB sentiment analysis

- a. DFA: слишком долго обучать
- b. PFA: error rate ~ 0.3 (улучшается)
- c. Bag of words: ~ 0.11

Ссылка на репозиторий:

<https://github.com/spbsu-ml-community/jml/tree/seva>