

# Математические модели в морфологии

## Условные случайные поля. Система UDPipe.

Алексей Андреевич Сорокин

спецкурс, ОТИПЛ МГУ,  
осенний семестр 2017–2018 учебного года  
07 ноября 2017 г.

Определение наилучшей последовательности

## Определение наилучшей последовательности

- Морфологический разбор сводится к нахождению наилучшей последовательности состояний.
- Состояния нельзя находить полным перебором.

## Определение наилучшей последовательности

## Определение наилучшей последовательности

- Морфологический разбор сводится к нахождению наилучшей последовательности состояний.
- Состояния нельзя находить полным перебором.
- Логарифм вероятности выходной последовательности:

$$\log p(\mathbf{q}|\mathbf{w}) = -\log Z(\mathbf{w}) + \sum_{t=1}^n \sum_k \theta_k f_k(q_t, q_{t-1}, \mathbf{w}, t)$$

- Логарифм вероятности начала последовательности:

$$\log p(\mathbf{q}_{1,m}|\mathbf{w}) = -\log Z(\mathbf{w}) + \sum_{t=1}^m \sum_k \theta_k f_k(q_t, q_{t-1}, \mathbf{w}, t)$$

## Определение наилучшей последовательности

## Определение наилучшей последовательности

- Морфологический разбор сводится к нахождению наилучшей последовательности состояний.
- Состояния нельзя находить полным перебором.
- Логарифм вероятности выходной последовательности:

$$\log p(\mathbf{q}|\mathbf{w}) = -\log Z(\mathbf{w}) + \sum_{t=1}^n \sum_k \theta_k f_k(q_t, q_{t-1}, \mathbf{w}, t)$$

- Логарифм вероятности начала последовательности:

$$\log p(\mathbf{q}_{1,m}|\mathbf{w}) = -\log Z(\mathbf{w}) + \sum_{t=1}^m \sum_k \theta_k f_k(q_t, q_{t-1}, \mathbf{w}, t)$$

- Формула пересчёта вероятностей:

$$\log p(\mathbf{q}_{1,m}|\mathbf{w}) = \log p(\mathbf{q}_{1,m-1}|\mathbf{w}) + \sum_k \theta_k f_k(q_m, q_{m-1}, \mathbf{w}, m)$$

Определение наилучшей последовательности

## Пересчёт вероятностей

- Формула пересчёта вероятностей:

$$\log p(\mathbf{q}_{1:m} | \mathbf{w}) = \log p(\mathbf{q}_{1:m-1} | \mathbf{w}) + \sum_k \theta_k f_k(q_m, q_{m-1}, \mathbf{w}, m)$$

## Определение наилучшей последовательности

## Пересчёт вероятностей

- Формула пересчёта вероятностей:

$$\log p(\mathbf{q}_{1:m} | \mathbf{w}) = \log p(\mathbf{q}_{1:m-1} | \mathbf{w}) + \sum_k \theta_k f_k(q_m, q_{m-1}, \mathbf{w}, m)$$

- $L(q_{m-1}, q_m, m)$  — штраф за переход из  $q_{m-1}$  в  $q_m$  в момент  $m$ .
- Штраф не зависит от предыдущего пути (только от последнего состояния).

## Определение наилучшей последовательности

## Пересчёт вероятностей

- Формула пересчёта вероятностей:

$$\log p(\mathbf{q}_{1:m} | \mathbf{w}) = \log p(\mathbf{q}_{1:m-1} | \mathbf{w}) + \sum_k \theta_k f_k(q_m, q_{m-1}, \mathbf{w}, m)$$

- $L(q_{m-1}, q_m, m)$  — штраф за переход из  $q_{m-1}$  в  $q_m$  в момент  $m$ .
- Штраф не зависит от предыдущего пути (только от последнего состояния).
- Можно запомнить последнее состояние и применить алгоритм Витерби (почти так же, как в марковских моделях).

## Определение наилучшей последовательности

## Пересчёт вероятностей

- Формула пересчёта вероятностей:

$$\log p(\mathbf{q}_{1,m} | \mathbf{w}) = \log p(\mathbf{q}_{1,m-1} | \mathbf{w}) + \sum_k \theta_k f_k(q_m, q_{m-1}, \mathbf{w}, m)$$

- $L(q_{m-1}, q_m, m)$  — штраф за переход из  $q_{m-1}$  в  $q_m$  в момент  $m$ .
- Штраф не зависит от предыдущего пути (только от последнего состояния).
- Можно запомнить последнее состояние и применить алгоритм Витерби (почти так же, как в марковских моделях).
- Введём частичные штрафы  $\alpha_{m,i}$ :

$$\alpha_{m,i} = \max_{\mathbf{q}_{1,m-1}, q_m=i} \log p(\mathbf{q}_{1,m} | \mathbf{w})$$

## Определение наилучшей последовательности

## Пересчёт вероятностей

- Формула пересчёта вероятностей:

$$\log p(\mathbf{q}_{1,m} | \mathbf{w}) = \log p(\mathbf{q}_{1,m-1} | \mathbf{w}) + \sum_k \theta_k f_k(q_m, q_{m-1}, \mathbf{w}, m)$$

- $L(q_{m-1}, q_m, m)$  — штраф за переход из  $q_{m-1}$  в  $q_m$  в момент  $m$ .
- Штраф не зависит от предыдущего пути (только от последнего состояния).
- Можно запомнить последнее состояние и применить алгоритм Витерби (почти так же, как в марковских моделях).
- Введём частичные штрафы  $\alpha_{m,i}$ :

$$\alpha_{m,i} = \max_{\mathbf{q}_{1,m-1}, q_m=i} \log p(\mathbf{q}_{1,m} | \mathbf{w})$$

- Формула пересчёта:

$$\begin{aligned}\alpha_{m+1,i} &= \max_j \alpha_{m,j} + L(j, i, m+1) \\ L(j, i, m+1) &= \sum_k \theta_k f_k(q_i, q_j, \mathbf{w}, m+1)\end{aligned}$$

Определение наилучшей последовательности

# Декодирование оптимальной последовательности

Нахождение оптимальной выходной последовательности:

- Вычислить частичные штрафы и обратные ссылки:

$$\alpha_{0,j} = \llbracket j = 0 \rrbracket,$$

$$\alpha_{m,i} = \max_j \alpha_{m,j} + \sum_k \theta_k f_k(j, i, \mathbf{w}, m), \quad m = 1, \dots, n,$$

$$\delta_{m,i} = \operatorname{argmax}_j \alpha_{m,j} + \sum_k \theta_k f_k(j, i, \mathbf{w}, m)$$

## Определение наилучшей последовательности

# Декодирование оптимальной последовательности

Нахождение оптимальной выходной последовательности:

- Вычислить частичные штрафы и обратные ссылки:

$$\alpha_{0,j} = \llbracket j = 0 \rrbracket,$$

$$\alpha_{m,i} = \max_j \alpha_{m,j} + \sum_k \theta_k f_k(j, i, \mathbf{w}, m), \quad m = 1, \dots, n,$$

$$\delta_{m,i} = \operatorname{argmax}_j \alpha_{m,j} + \sum_k \theta_k f_k(j, i, \mathbf{w}, m)$$

- Восстановить последовательность по обратным ссылкам:

$$q(n) = \operatorname{argmax}_i A_{n,i}$$

$$q(t-1) = \delta_{t,q(t)}$$

- По состояниям восстанавливается разбор.

## Определение наилучшей последовательности

## Декодирование оптимальной последовательности

- Сложность декодирования:  $L * M * R * K$ .

$L$  — длина последовательности,

$M$  — максимальное число “активных” состояний,

$R$  — количество способов продолжить состояние,

$K$  — число признаков.

## Определение наилучшей последовательности

## Декодирование оптимальной последовательности

- Сложность декодирования:  $L * M * R * K$ .  
 $L$  — длина последовательности,  
 $M$  — максимальное число “активных” состояний,  
 $R$  — количество способов продолжить состояние,  
 $K$  — число признаков.
- Если состояние — энграммма порядка  $m$ , то  $M \sim R^m$ , где  $R$  — максимальная степень неоднозначности слова.

## Определение наилучшей последовательности

## Декодирование оптимальной последовательности

- Сложность декодирования:  $L * M * R * K$ .  
 $L$  — длина последовательности,  
 $M$  — максимальное число “активных” состояний,  
 $R$  — количество способов продолжить состояние,  
 $K$  — число признаков.
- Если состояние — энграммма порядка  $m$ , то  $M \sim R^m$ , где  $R$  — максимальная степень неоднозначности слова.
- Число признаков тоже растёт экспоненциально по  $m$ .
- Следствие:  $m > 2$  нереализуемо на практике, с  $m = 2$  проблемы для языков с развитой морфологией.

## Определение наилучшей последовательности

## Декодирование оптимальной последовательности

- Сложность декодирования:  $L * M * R * K$ .  
 $L$  — длина последовательности,  
 $M$  — максимальное число “активных” состояний,  
 $R$  — количество способов продолжить состояние,  
 $K$  — число признаков.
- Если состояние — энграммма порядка  $m$ , то  $M \sim R^m$ , где  $R$  — максимальная степень неоднозначности слова.
- Число признаков тоже растёт экспоненциально по  $m$ .
- Сложность обучения:  $T * N * C_0$ , где  $C_0$  — сложность декодирования,  $T$  — число эпох обучения,  $N$  — размер обучающей выборки.
- Следствие:  $m > 2$  нереализуемо на практике, с  $m = 2$  проблемы для языков с развитой морфологией.
- Часто CRF реализуют иерархически: сначала грубая классификация (части речи), потом более точная.

## Недостатки CRF

- Преимущества CRF:

- Большая гибкость по сравнению с марковскими моделями.
- Локальные признаки произвольной природы (в том числе лексические).

## Недостатки CRF

- Преимущества CRF:
  - Большая гибкость по сравнению с марковскими моделями.
  - Локальные признаки произвольной природы (в том числе лексические).
- Недостатки CRF:
  - Большие затраты (время и память) на обучение.
  - Большое количество признаков (в том числе избыточных).

## Недостатки CRF

- Преимущества CRF:
  - Большая гибкость по сравнению с марковскими моделями.
  - Локальные признаки произвольной природы (в том числе лексические).
- Недостатки CRF:
  - Большие затраты (время и память) на обучение.
  - Большое количество признаков (в том числе избыточных).
  - Невозможность учитывать удалённый контекст.

## Недостатки CRF

- Преимущества CRF:
  - Большая гибкость по сравнению с марковскими моделями.
  - Локальные признаки произвольной природы (в том числе лексические).
- Недостатки CRF:
  - Большие затраты (время и память) на обучение.
  - Большое количество признаков (в том числе избыточных).
  - Невозможность учитывать удалённый контекст.
- CRF применимы к произвольной задаче разметке последовательно: распознавание именованных сущностей, разбиение на морфемы и т. д.

## Недостатки CRF

- Преимущества CRF:
  - Большая гибкость по сравнению с марковскими моделями.
  - Локальные признаки произвольной природы (в том числе лексические).
- Недостатки CRF:
  - Большие затраты (время и память) на обучение.
  - Большое количество признаков (в том числе избыточных).
  - Невозможность учитывать удалённый контекст.
- CRF применимы к произвольной задаче разметке последовательно: распознавание именованных сущностей, разбиение на морфемы и т. д.
- В последние годы уступают нейронным моделям.

## Система UDPipe

- Наиболее эффективная из доступных в открытом доступе систем — система проекта UDPipe (Карлов университет, Прага, UFAL).

## Система UDPipe

- Наиболее эффективная из доступных в открытом доступе систем — система проекта UDPipe (Карлов университет, Прага, UFAL).
- Основана на условных случайных полях.
- Обучается с помощью усреднённого персептрона.

## Система UDPipe

- Наиболее эффективная из доступных в открытом доступе систем — система проекта UDPipe (Карлов университет, Прага, UFAL).
- Основана на условных случайных полях.
- Обучается с помощью усреднённого персептрана.
- Использует развёрнутые шаблоны признаков (изначально разрабатывалась для чешского).

# Система UDPipe

- Наиболее эффективная из доступных в открытом доступе систем — система проекта UDPipe (Карлов университет, Прага, UFAL).
- Основана на условных случайных полях.
- Обучается с помощью усреднённого персептрона.
- Использует развёрнутые шаблоны признаков (изначально разрабатывалась для чешского).
- Предсказываются несколько меток: UPOSTAG (часть речи), XPOSTAG (детализованная часть речи: наклонение глагола, одушевлённость сущ-го, ...), FEATS (признаковое описание).

# Система UDPipe

- Наиболее эффективная из доступных в открытом доступе систем — система проекта UDPipe (Карлов университет, Прага, UFAL).
- Основана на условных случайных полях.
- Обучается с помощью усреднённого персептрона.
- Использует развёрнутые шаблоны признаков (изначально разрабатывалась для чешского).
- Предсказываются несколько меток: UPOSTAG (часть речи), XPOSTAG (детализованная часть речи: наклонение глагола, одушевлённость сущ-го, ...), FEATS (признаковое описание).
- Лемматизатор использует ту же математическую модель, но предсказывает шаблон словоизменения.

# Признаковое описание для английского языка

Context predicting whole tag	
Tags	Previous tag Previous two tags First letter of previous tag
Word forms	Current word form Previous word form Previous two word forms Following word form Following two word forms Last but one word form
Current word affixes	Prefixes of length 1-9 Suffixes of length 1-9
Current word features	Contains number Contains dash Contains upper case letter

Всего ~ 2,0 млн. признаков.

## Признаковое описание для чешского языка

Признак	Описание
$t_0$	текущая метка
$t_{-1}t_0$	биграмма меток
$t_{-2}; t_{-2}t_{-1}t_0$	триграмма меток
$w_{-2}, w_{-1}, w_0, w_1, w_2$	слово и его контекст
$I(V_L), t(V_L)$	позиция слова
$I(V_L), t(V_L)$	метка и лемма ближайшего глагола слева
$I(V_R), t(V_R)$	метка и лемма ближайшего глагола справа
	капитализация

Всего  $\sim 8,4$  млн. признаков.

## Условные случайные поля

- Условные случайные поля могут решать любую задачу, сводящуюся к разметке последовательностей.

## Условные случайные поля

- Условные случайные поля могут решать любую задачу, сводящуюся к разметке последовательностей.
- Пример: выделение именованных сущностей.
- Границы групп задаются метками *B*(начало), *M*(середина), *E*(конец), *S*(однословная группа), *O*(вне групп).

Андрей	Николаевич	Колмогоров	работал	в	Московском	Университете
B-PERS	M-PERS	E-PERS	O	O	B-ORG	E-ORG

## Условные случайные поля

- Условные случайные поля могут решать любую задачу, сводящуюся к разметке последовательностей.
- Пример: выделение именованных сущностей.
- Границы групп задаются метками *B*(начало), *M*(середина), *E*(конец), *S*(однословная группа), *O*(вне групп).

Андрей	Николаевич	Колмогоров	работал	в	Московском	Университете
B-PERS	M-PERS	E-PERS	O	O	B-ORG	E-ORG

- Основные признаки: капитализация, частеречные метки, сами слова (в окне ширины 5).

## Лемматизация: постановка задачи

- Пока мы что решали только задачу нахождения морфологической метки.
- Однако чаще всего нужна ещё и начальная форма (извлечение информации).

## Лемматизация: постановка задачи

- Пока мы что решали только задачу нахождения морфологической метки.
- Однако чаще всего нужна ещё и начальная форма (извлечение информации).
- Если словоформа словарная и известна морфологическая метка, то достаточно просто посмотреть в словаре.
- Исключений очень мало:

*лечу+V+Pres+Sg+1 ↪ лечить / лететь.*

## Лемматизация: постановка задачи

- Пока мы что решали только задачу нахождения морфологической метки.
- Однако чаще всего нужна ещё и начальная форма (извлечение информации).
- Если словоформа словарная и известна морфологическая метка, то достаточно просто посмотреть в словаре.
- Исключений очень мало:

*лечу+V+Pres+Sg+1 ↪ лечить / лететь.*

- Что делать с несловарными словами?

## Лемматизация: постановка задачи

- Пока мы что решали только задачу нахождения морфологической метки.
- Однако чаще всего нужна ещё и начальная форма (извлечение информации).
- Если словоформа словарная и известна морфологическая метка, то достаточно просто посмотреть в словаре.
- Исключений очень мало:

$\text{лечу} + V + \text{Pres} + \text{Sg} + 1 \mapsto \text{лечить} / \text{лететь}.$

- Что делать с несловарными словами?
- Для них априори неизвестен даже список возможных меток и оценки их вероятностей.

# Морфологический анализ как задача классификации

- Задача: дано несловарное слово  $w$ , найти наиболее вероятную морфологическую метку  $t(w)$ .
- Для марковских моделей и CRF нужно чуть больше: распределение вероятностей  $p(t|w)$  для возможных меток  $t$ .

# Морфологический анализ как задача классификации

- Задача: дано несловарное слово  $w$ , найти наиболее вероятную морфологическую метку  $t(w)$ .
- Для марковских моделей и CRF нужно чуть больше: распределение вероятностей  $p(t|w)$  для возможных меток  $t$ .
- Можно решать как задачу классификации:
  - Описать каждое слово вектором признаков. Морфологическая метка — предсказываемый класс.
  - Словарные/корпусные слова — обучающая выборка.

# Морфологический анализ как задача классификации

- Стандартный набор признаков: суффиксы слова до  $m$  символов ( $m \sim 5 - 8$ ), префиксы (если нужно), капитализация.
- Каждый признак кодируется двоичной схемой:

	-'	-а	-ть	-ль	-ка	-ать	...	Класс
играть	1	0	1	0	0	1	...	VERB, Mood=Inf,...
бегать	1	0	1	0	0	1	...	VERB, Mood=Inf,...
плыть	1	0	1	0	0	0	...	VERB, Mood=Inf,...
кровать	1	0	1	0	0	1	...	NOUN, case=Nom,...
тень	1	0	0	0	0	0	...	NOUN, case=Nom,...
рука	0	1	0	0	1	0	...	NOUN, case=Nom,...
броска	0	1	0	0	1	0	...	NOUN, case=Gen,...

# Морфологический анализ как задача классификации

- Стандартный набор признаков: суффиксы слова до  $m$  символов ( $m \sim 5 - 8$ ), префиксы (если нужно), капитализация.
- Каждый признак кодируется двоичной схемой:

	-'	-а	-ть	-ль	-ка	-ать	...	Класс
играть	1	0	1	0	0	1	...	VERB, Mood=Inf,...
бегать	1	0	1	0	0	1	...	VERB, Mood=Inf,...
плыть	1	0	1	0	0	0	...	VERB, Mood=Inf,...
кровать	1	0	1	0	0	1	...	NOUN, case=Nom,...
тень	1	0	0	0	0	0	...	NOUN, case=Nom,...
рука	0	1	0	0	1	0	...	NOUN, case=Nom,...
броска	0	1	0	0	1	0	...	NOUN, case=Gen,...

- После этого применим любой алгоритм машинного обучения (метод опорных векторов, персепtron, логистическая регрессия, деревья решений...)

## Лемматизация как задача классификации

- Данный подход также позволяет осуществлять лемматизацию.
- Для этого надо предсказывать тип преобразования от словоформы к лексеме.

## Лемматизация как задача классификации

- Данный подход также позволяет осуществлять лемматизацию.
- Для этого надо предсказывать тип преобразования от словоформы к лексеме.
- Простейший тип преобразования — изменения суффикса:

коленями  $\mapsto$  колено

X+ями  $\mapsto$  X+о

## Лемматизация как задача классификации

- Данный подход также позволяет осуществлять лемматизацию.
- Для этого надо предсказывать тип преобразования от словоформы к лексеме.
- Простейший тип преобразования — изменения суффикса:

коленями  $\mapsto$  колено

X+ями  $\mapsto$  X+о

- Можно добавить изменения префиксов:

наикраснейший  $\mapsto$  красный

наи+X+ейший  $\mapsto$  X+ый

## Лемматизация как задача классификации

- Данный подход также позволяет осуществлять лемматизацию.
- Для этого надо предсказывать тип преобразования от словоформы к лексеме.
- Простейший тип преобразования — изменения суффикса:

коленями  $\mapsto$  колено

X+ями  $\mapsto$  X+о

- Можно добавить изменения префиксов:

наикраснейший  $\mapsto$  красный

наи+X+ейший  $\mapsto$  X+ый

- Иногда изменения не описываются “непрерывным” шаблоном:

vuelvo  $\mapsto$  volver

броска  $\mapsto$  бросок

yiktubun  $\mapsto$  kataba

## Абстрактные парадигмы

- Можно ввести “разрывные” шаблоны словоизменения:

vuelvo	↔	volver
1+ue+2+o	↔	1+o+2+er
броска	↔	бросок
1+2+a	↔	1+o+2
yiktubun	↔	kataba
yi+1+2+u+3+un	↔	1+a+2+a+3+a

- Переменная часть шаблона (основа) — наибольшая общая подпоследовательность (НОП) начальной и производной словоформы.

## Абстрактные парадигмы

- Можно ввести “разрывные” шаблоны словоизменения:

vuelvo	↔	volver
1+ue+2+o	↔	1+o+2+er
броска	↔	бросок
1+2+a	↔	1+o+2
yiktubun	↔	kataba
yi+1+2+u+3+un	↔	1+a+2+a+3+a

- Переменная часть шаблона (основа) — наибольшая общая подпоследовательность (НОП) начальной и производной словоформы.
- Лингвистическое обоснование:

“Основу любой словоформы можно представить как состоящую из двух частей: неизменяемой и изменяемой, то есть такой, которой хотя бы в одной из прочих словоформ парадигмы соответствует некоторая иная цепочка букв”.

(А. А. Зализняк, “Русское именное словоизменение”)

## Абстрактные парадигмы

- Алгоритм извлечения абстрактной парадигмы словоизменения:

- Найти НОП: *yiktbun, kataba*  $\mapsto$  k-t-b
- Выделить в словах элементы НОП:

*yik**tub**un, kata**ba***

# Абстрактные парадигмы

- Алгоритм извлечения абстрактной парадигмы словоизменения:

- Найти НОП: *yiktbun, kataba*  $\mapsto$  k-t-b
- Выделить в словах элементы НОП:

*yikt**ub**un, kata**ba***

- Заменить *i*-ый элемент НОП на переменную  $x_i$ :

k  $\mapsto$   $x_1$ , t  $\mapsto$   $x_2$ , b  $\mapsto$   $x_3$ ,

у i  $x_1x_2$  и  $x_3$  un, а  $x_2$  а  $x_3$  а

# Абстрактные парадигмы

- Алгоритм извлечения абстрактной парадигмы словоизменения:
  - Найти НОП: *yiktbun, kataba*  $\mapsto$  k-t-b
  - Выделить в словах элементы НОП:  
*yik**tub**un, kata**ba***
  - Заменить *i*-ый элемент НОП на переменную  $x_i$ :  
 $k \mapsto x_1, t \mapsto x_2, b \mapsto x_3,$   
 $y \mapsto x_1 x_2 \text{ и } x_3 \text{ un, а } x_2 \text{ а } x_3 \text{ a}$
  - Вернуть получившийся шаблон
- Проблема: как находить НОП.