



Лингвистика для NLP

Ася Ройтберг, asya.roytberg@bigdatateam.org

ML Instructor at BigData Team, <http://bigdatateam.org/>
<https://www.facebook.com/bigdatateam/>

16.12.2020



MEGAFON

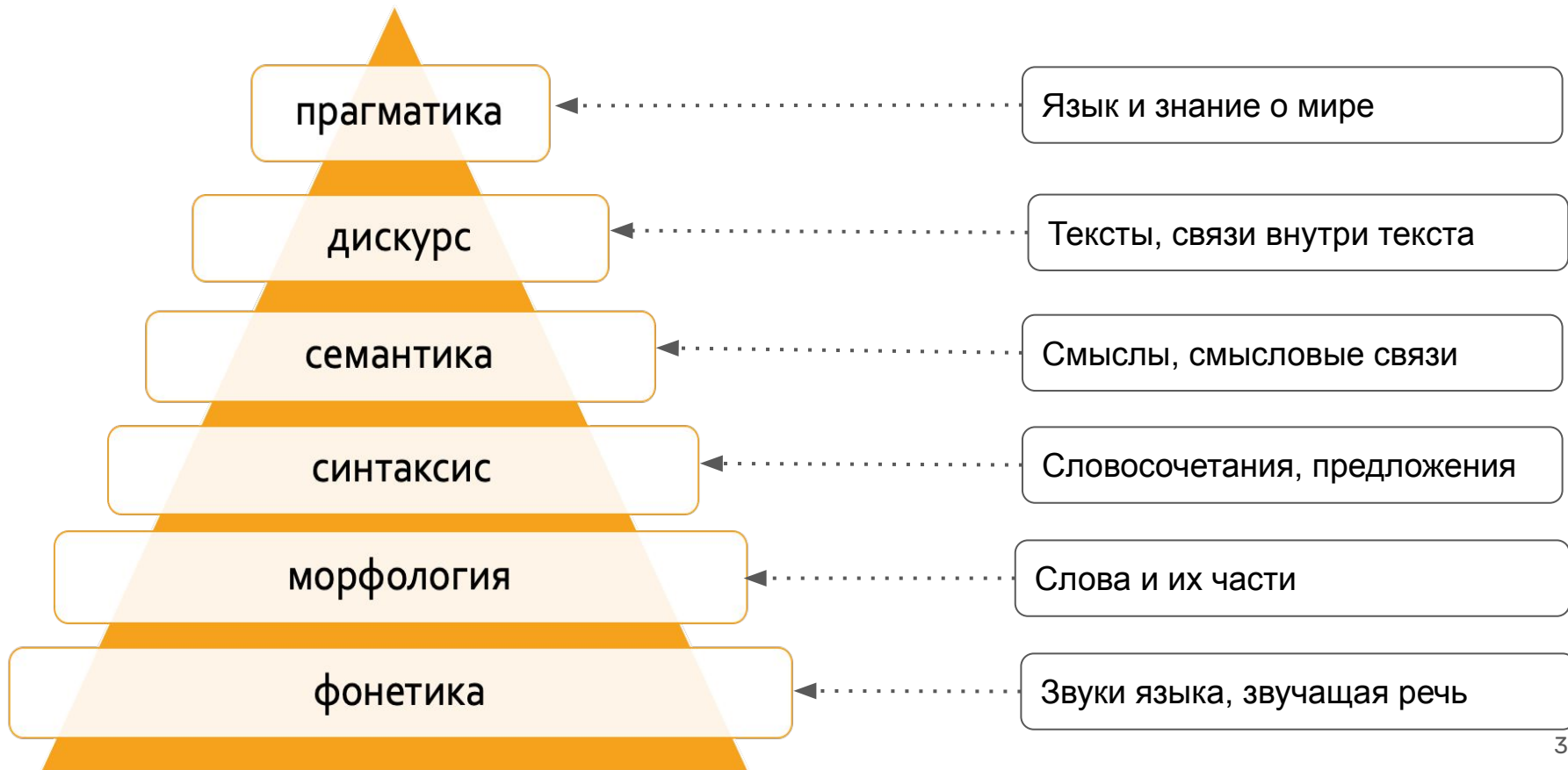


Зачем нам лингвистика?

- ▶ Читать SOTA статьи
- ▶ Улучшение моделей
- ▶ Извлечение “умных фичей”



Уровни языка





Что такое слово?

**Токенизация, стемминг,
лемматизация**



Одно слово или разные?

1

Написаны через пробел = разные слова?

Не присесть ли тебе на
диван-кровать?

2

Формы одного слова или разные слова

прыгал
прыгающий
прыгнувший
попрыгунчик
прыгнул



Разбиение на слова

- дефис

счет-фактура, сильно-сильно, English-speaking,
easy-to-remember, well-known vs specially designed

- апостроф

isn't - cannot, 's,

- слитные предлоги в иврите (, לערץ , ץ בער
מערץ)

- Word2Vec, TF\IDF, P.S. и т.д.



1

Не просто разбиение по пробелам

2

Сложные случаи – вопрос договоренностей



Слово, словоформа

прыгал
прыгающий
прыгнувший
попрыгунчик
прыгнул

вода
воды
водой
воде

water
waters
water

В каком языке у слов больше словоформ?
английский vs. русский

склонение и спряжение



Какие есть проблемы?



Стемминг: проблемы

1. Чередования букв: ~~кошк~~~~а~~ – ~~кош~~~~ек~~ ~~круг~~ – ~~круж~~~~ок~~
2. Нерегулярные формы: ~~ребен~~~~ок~~ – ~~дет~~~~и~~ ~~был~~ – ~~есть~~
3. Короткие слова: ~~е~~~~л~~, ~~был~~
4. Не различает, где аффикс, а где часть слова:

~~руки~~ vs. ~~при~~ ~~element~~ vs. ~~measurement~~



Стемминг: преимущества



СКОРОСТЬ



Лемматизация: к чему сводим?

- Для существительных – номинатив (Что это?)
- Для остальных частей речи – словарная форма

Словарная форма – вопрос традиции.

Вот так выглядит “словарная форма”
латинского глагола

facio, feci, factum, ere (делать)



Лемматизация: проблемы

1

Омонимия

зАмок - замОк

пятнистая рысь - бежал рысью

три – три | тереть

села - селО | сесть

2

Скорость





- ▶ Список частей речи
- ▶ Частеречная омонимия
- ▶ Знаменательные\служебные части речи
- ▶ Нас не интересуют сложные случаи



Части речи (POS)

UPOS – набор универсальных тегов.

XPOS – специфические наборы тегов для отдельных
ЯЗЫКОВ

- StanfordNLP
- Stanza

word: Barack	upos: PROPN	xpos: NNP
word: Obama	upos: PROPN	xpos: NNP
word: was	upos: AUX	xpos: VBD
word: born	upos: VERB	xpos: VBN
word: in	upos: ADP	xpos: IN
word: Hawaii	upos: PROPN	xpos: NNP
word: .	upos: PUNCT	xpos: .



Universal Dependencies POS

<https://universaldependencies.org>

Open class words	Closed class words	Other
<p><u>ADJ</u>: adjective (прилагательное)</p> <p><u>ADV</u>: adverb (наречие)</p> <p><u>INTJ</u>: interjection (междометие)</p> <p><u>NOUN</u>: noun (существительное)</p> <p><u>PROPN</u>: proper noun (имя собственное)</p> <p><u>VERB</u>: verb (глагол)</p>	<p><u>ADP</u>: adposition (пред\послелоги)</p> <p><u>AUX</u>: auxiliary (вспомогательный глагол)</p> <p><u>CCONJ</u>: coordinating conjunction (сочинительный союз)</p> <p><u>DET</u>: determiner (детерминатив)</p> <p><u>NUM</u>: numeral (числительные)</p> <p><u>PART</u>: particle (частица)</p> <p><u>SCONJ</u>: subordinating conjunction (подчинительный союз)</p>	<p><u>PUNCT</u>: punctuation (пунктуация)</p> <p><u>SYM</u>: symbol (символы)</p> <p><u>X</u>: other (другое)</p>



Universal Dependencies POS

UD POS	He UD POS
<u>SpaCy</u>	<u>nlTK</u>
<u>Flair</u>	
<u>UDPipe</u>	
<u>Stanza</u>	<u>Stanza</u>



Вспоминаем про НММ



Три задачи НММ

- ▶ Правдоподобие
- ▶ Декодирование
- ▶ Обучение



Три задачи НММ: декодирование

Дано:

модель $\lambda = (\{a_{ij}\}, \{b_j(k)\}, \{\pi_i\})$;

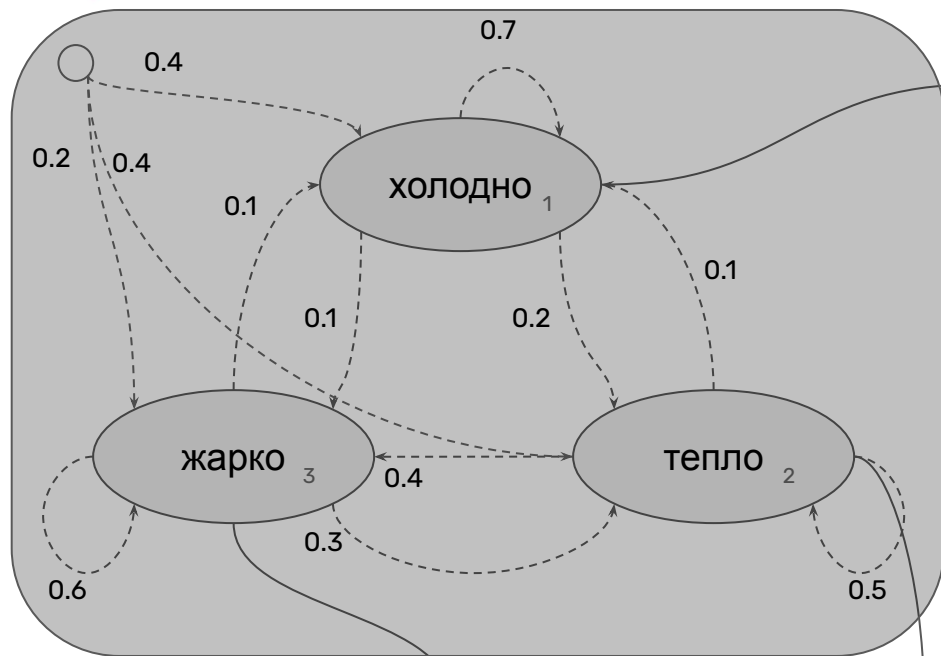
наблюдаемая последовательность $o_1 o_2 \dots o_T$

Необходимо найти наиболее вероятную последовательность
скрытых состояний

$$\operatorname{argmax} P[q_1 q_2 \dots, q_T \mid \lambda, o_1 o_2 \dots o_T]$$



Три задачи HMM: декодирование



$$\begin{aligned} b_1(0) &= P[0|\text{холодно}] = 0.1 \\ b_1(1) &= P[1|\text{холодно}] = 0.3 \\ b_1(2) &= P[2|\text{холодно}] = 0.4 \\ b_1(3) &= P[3|\text{холодно}] = 0.2 \end{aligned}$$

Какая последовательность погодных условий наиболее вероятна, при условии, что Петя съел:

1 июня - 2 мороженных
2 июня - 1 мороженое
3 июня - 0 мороженных?

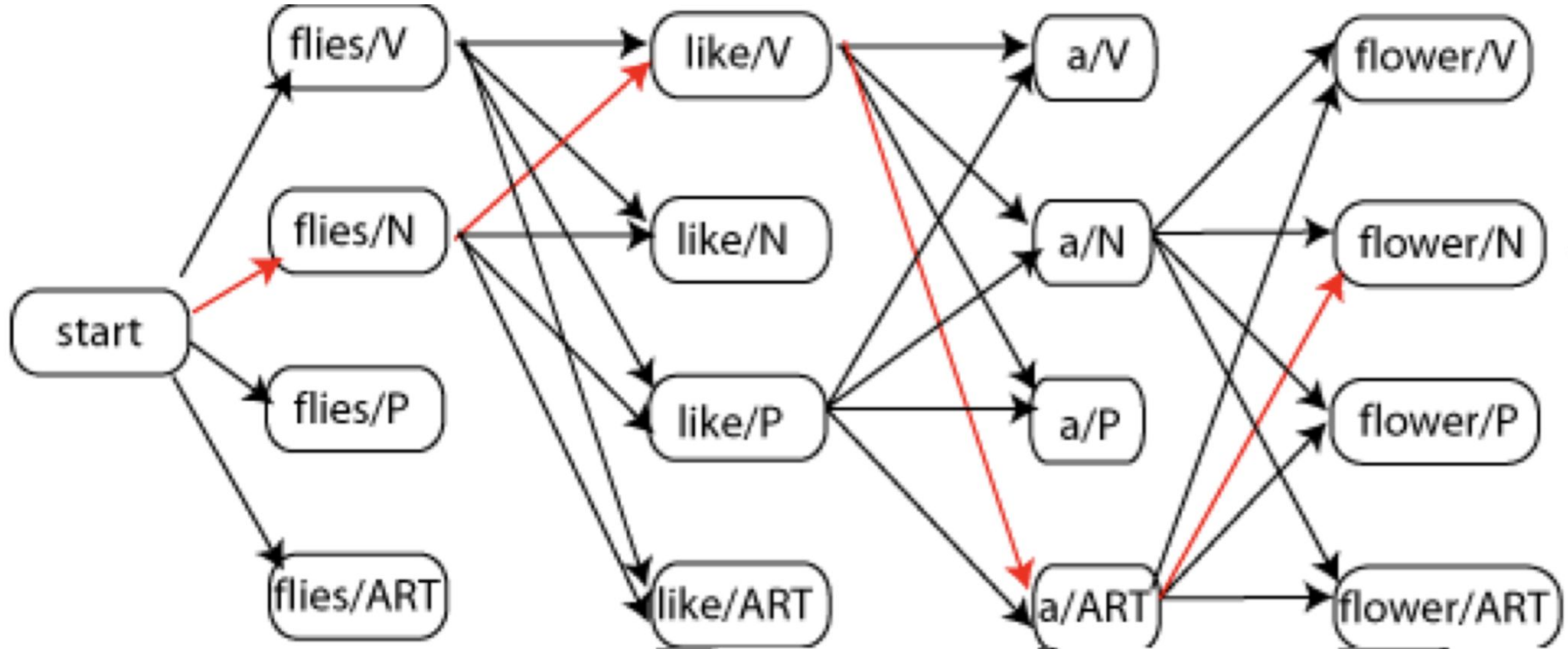
Перебирать все варианты вычислительно долго!

$$\begin{aligned} b_3(0) &= P[0|\text{жарко}] = 0 \\ b_3(1) &= P[1|\text{жарко}] = 0.1 \\ b_3(2) &= P[2|\text{жарко}] = 0.3 \\ b_3(3) &= P[3|\text{жарко}] = 0.6 \end{aligned}$$

$$\begin{aligned} b_2(0) &= P[0|\text{тепло}] = 0.1 \\ b_2(1) &= P[1|\text{тепло}] = 0.3 \\ b_2(2) &= P[2|\text{тепло}] = 0.4 \\ b_2(3) &= P[3|\text{тепло}] = 0.2 \end{aligned}$$



Что в случае задачи разметки частей речи скрытые состояния, а что наблюдаемые?





Морфология: summary

- ▶ Морфология занимается словами и их частями
- ▶ Не всегда ясно, где разные слова, а где одно слово
- ▶ У слов есть части речи
- ▶ Можно сделать POS таггер на основе HMM
- ▶ UD - популярная мультязычная POS разметка
- ▶ UD - не единственный способ размечать POS

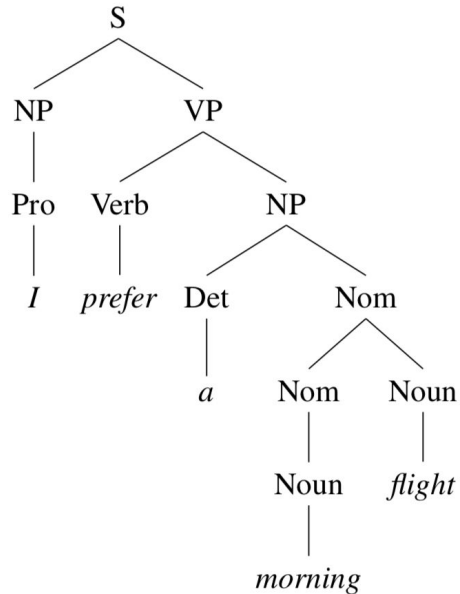


Составляющие и зависимости



Грамматика составляющих

Constituency Grammar



S- sentence

NP - Nominal Phrase

VP - Verbal Phrase

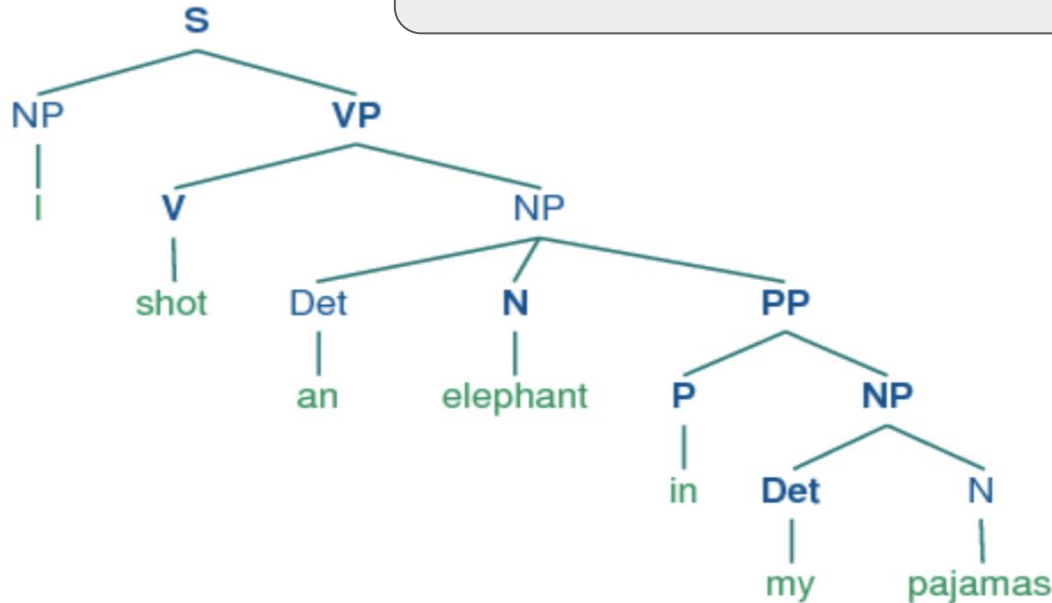
PP - Prepositional Phrase



Грамматика составляющих

Неоднозначность (Ambiguity)

I shot an elephant in my pajamas.

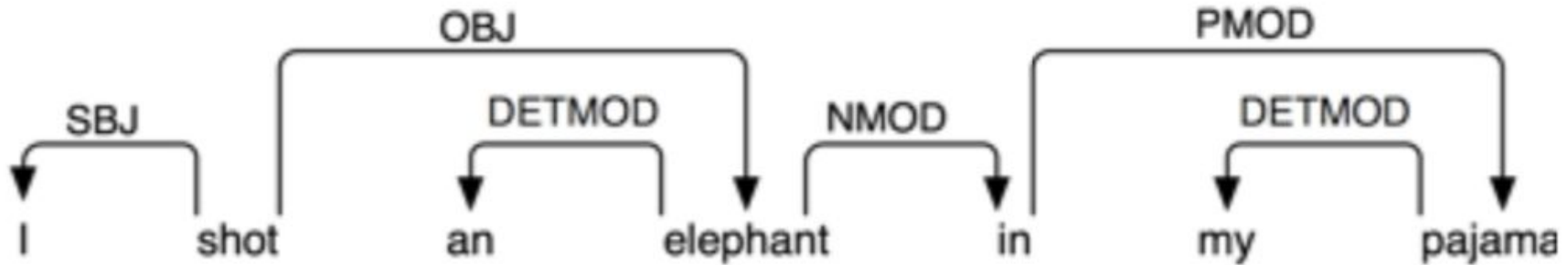


Who wears the
pajamas?



Грамматика зависимостей

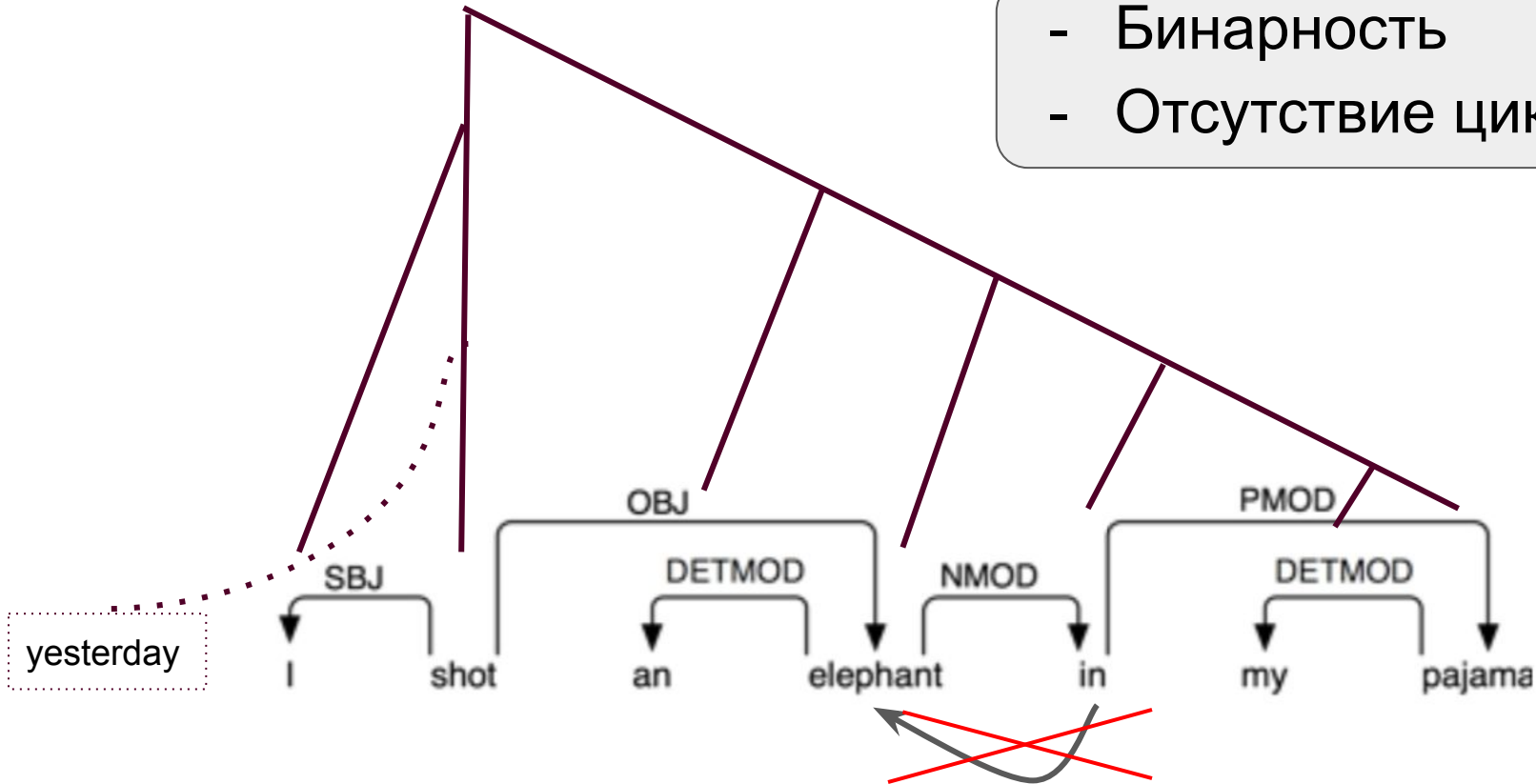
Dependency Grammar





Грамматика зависимостей

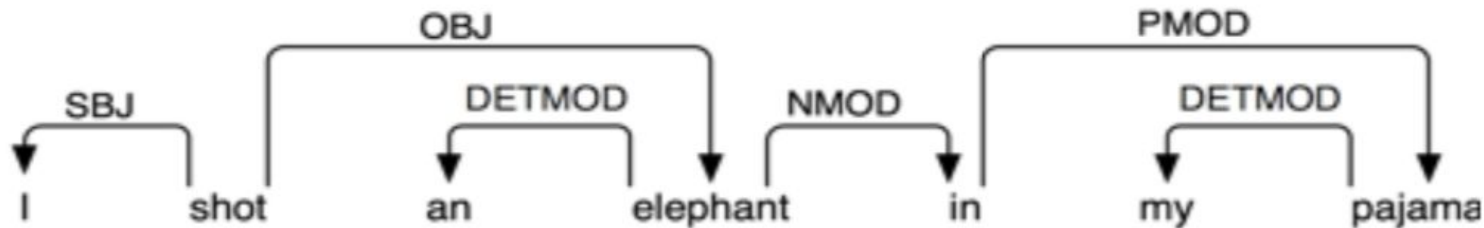
- Бинарность
- Отсутствие циклов





List of UD relationships

Subject – Predicate – Object – центр большинства предложений



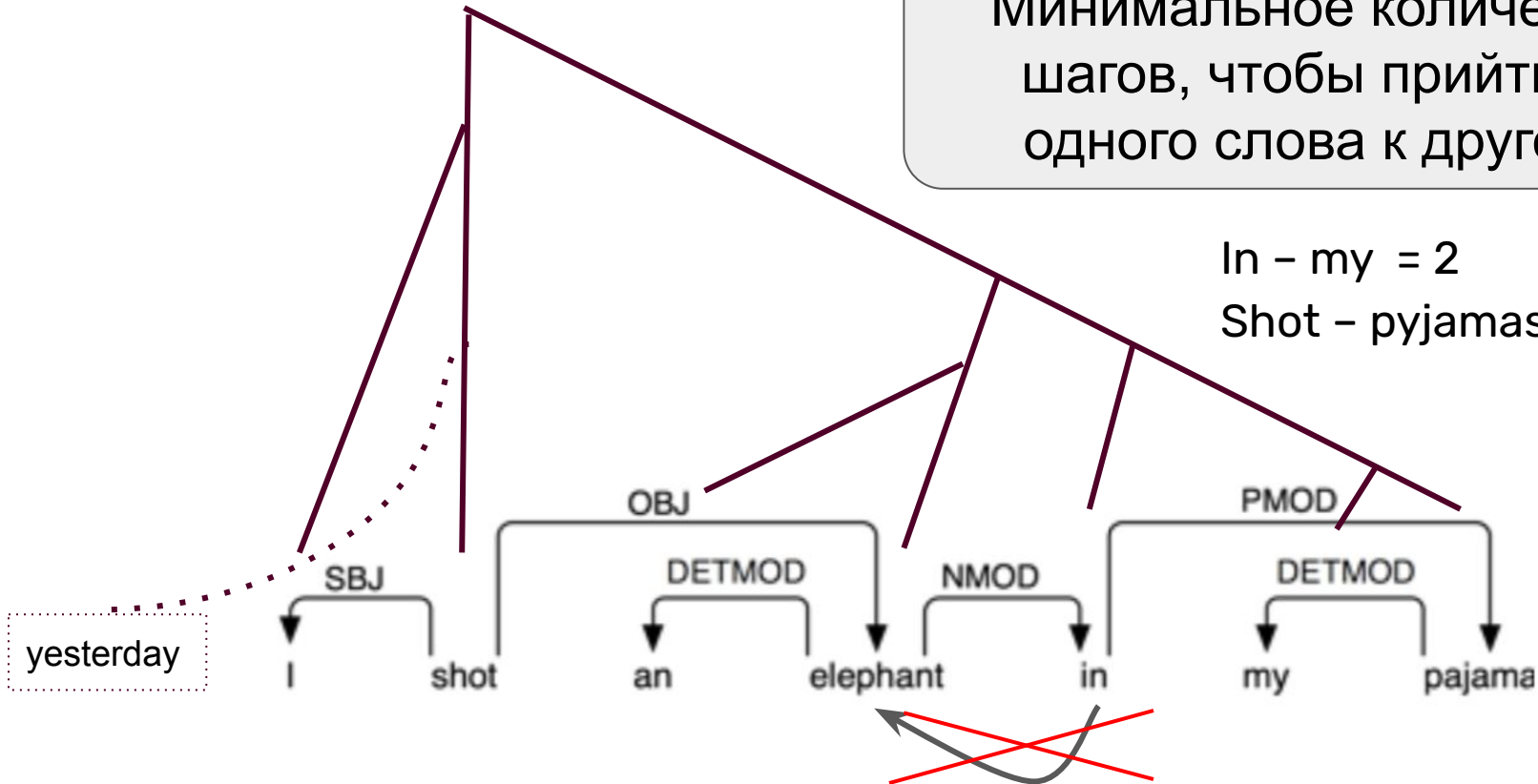


Синтаксическое расстояние

Минимальное количество шагов, чтобы прийти от одного слова к другому

In - my = 2

Shot - pyjamas = ?





Наименьший общий предок (LSA)

Расстояние между вершинами a и b =
расстояние от a до LCA + расстояние от b до LCA

Процедура $LCA(u, v)$:

$h1 := \text{depth}(u)$ $// \text{depth}(x) = \text{глубина вершины } x$

$h2 := \text{depth}(v)$

while $h1 \neq h2$:

if $h1 > h2$:

$u := \text{parent}(u)$

$h1 := h1 - 1$

else:

$v := \text{parent}(v)$

$h2 := h2 - 1$

while $u \neq v$:

$u := \text{parent}(u)$

$// \text{parent}(x) = \text{непосредственный предок вершины } x$

$v := \text{parent}(v)$

return u



- У предложений есть очень четкая структура
- Структура предложения различна в разных языках
- Есть 2 основных формализма :
 - a.** Грамматика составляющих
 - b.** Грамматика зависимостей
- Понятие синтаксического расстояния
- Синтаксическое расстояние не прямо зависит от линейного расстояния



Слова, смыслы и связи



Дистрибутивная семантика

- ▶ Лексическая единица описывается вектором
- ▶ Измерения (или компоненты) – другие слова лексикона
- ▶ Значения этих компонентов – частота совместной встречаемости “нашей” единицы с другими словами в данном корпусе



Дистрибутивная семантика

Проблема: многозначность



Для эмбедингов тоже
проблема

Гуляли по стрелка Васильевского острова.

Повесили указатели в виде стрелок.

Соедини главное и зависимое слова стрелкой.



“Атомы смысла”

Концепция

См. напр.

- Модель Смысл ↔ Текст
- Теория семантических примитивов

Есть некоторый ограниченный набор смыслов и отношений между этими смыслами, из которых можно собрать значение любого слова.

Очень упрощенно:

Блюдце - тарелка с размером меньше среднего

Бежать - идти быстро



Смысл суммы \neq сумме смыслов

► Устойчивые выражения:

красна девица, rain cats and dogs

► Эвфемизмы: белая полярная лиса

► Коллокации:

крепкий чай – *сильный чай

принять душ – *получить душ

powerful engine – *strong engine

$$P(AB) > P(A) * P(B)$$

и многое
другое



Семантические классы слов

- Группы слов одной части речи общим компонентом значения

стол
стул
трюмо
этажерка

идти
ползти
бежать
ковылять

карандаш
веревка
подушка
провод
ковёр



Семантическое поле

Слова разных частей речи, объединенные общим компонентом значения

стрелять
лук пушка
ружье
стрелок курок
пуля

палочка
колдовать
ведьма магия
волшебник

роутер пинг
маршрутизация
фильтровать
мониторинг
витая пара



Именованные сущности

Традиционное понимание: ~ имена собственные

- Персоны
- Названия организация
- Названия географических объектов
- Названия белков

.....



Именованные сущности Custom edition

Современное понимание:

некоторый класс существительных + даты

- Названия политических должностей
- Термины родства
- Тип организации
- и т.д.



Основные отношения семантического поля

- **Синонимия** (милый – симпатичный)
- **Антонимия** (входить - выходить)
- **Часть-целое** (ветка – дерево, окно – дом)
- **Гипонимия** (~ включение в класс: человек - живые существа)
 - *человек* – гипоним, *живые существа* – гипероним
- **Несовместимость** (мать – отец)



Семантические отношения

Любые отношения, которые можно формализовать

- Person **родиться в** Location
- Person **быть сотрудником** Organization
- Organization **быть акционером** Organization
- Location **находиться в** Location
- Person **быть родителем** Person



Онтологии и базы знаний

Понятия	<u>Airbus</u>	<u>Bill Gates</u>
Атрибуты	<u>aerospace manufacturer</u>	human
Отношения	member of → Linux Foundation	Spouse → Melinda Gates



Извлечение отношений

Дистанционная разметка и [Wikidata](#)

place of birth (P19)

Джон Леннон родился в портовом английском городе Ливерпуль

Q1203

Q24826

Employer (P108)

Стив Джобс – один из основателей корпорации Apple и киностудии Pixar

Q19837

Q312

Q127552



Проблемы. False positive

Джон Леннон уехал из Ливерпуля в 1960-х.

Уехать в город из аэропорта Ливерпуля «Джон Леннон» (LPL).

Разругавшись с Apple, Джобс создал фирму NeXT.



- Компоненты значения (“атомы смыслов”)
- Многозначность и эмбединги
- Сумма смыслов \neq смыслу суммы
- Семантические классы
- Семантические поля
- Именованные сущности
- Семантические отношения
- Использование баз знаний для извлечения отношений



Тексты и отношения в них



Референция и кореференция

Референт - объект реального мира, называемый выражением

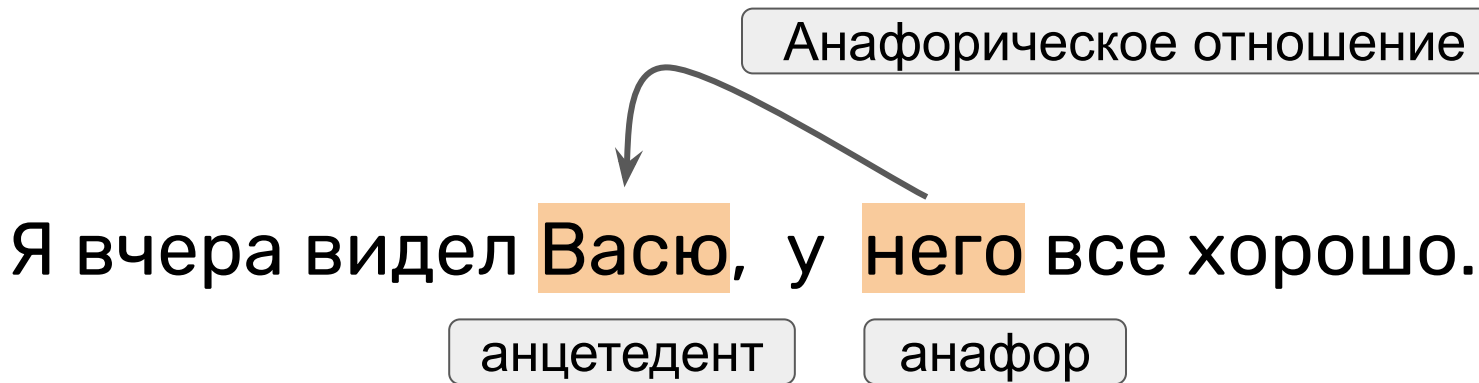
Кореференция - ситуация, когда несколько выражений в тексте, отсылают к одному объекту

Дом на Бейкер стрит – легендарное жилище Шерлока Холмса. Именно в нем были найдены разгадки к самым запутанным делам.



Анафора (+ Катафора)

Для понимания значения элемента необходима отсылка к другому элементу из текста.



Синглтон – объект, который только один раз упоминается в тексте



► Не только местоимения

IPO QIWI состоялось сегодня на бирже Nasdaq. Торги расписками компании стартовали в 17:00.

► Не только кореференция

В Уругвае прошли выборы. Новый президент – Луис Альберто Локалле





Связность текста

Coherence

- ▶ **Coherence relations:** отношения причины, контраста ...
- ▶ **Entity-based coherence:** есть ядерные (salient) слова вокруг которых строится текст
- ▶ **Topic-based coherence:** предложения об одной теме обычно оказываются рядом
- ▶ **Global coherence:** у каждого типа текста есть ожидаемая структура



Чем тексты могут отличаться?

- ▶ Измеряемые параметры: длина текста, длины предложений и т. п.
- ▶ Сбалансированность лексики
- ▶ Структура текста
- ▶ Стилль текста



Как выбрать правильный корпус?

Главный вопрос всего:

С какими текстами предстоит работать в “боевых условиях” ?



- Текст состоит из единиц, связанных отношениями
- Понятие связности текста
- Понятие структуры текста
- Нужно внимательно выбирать тексты, на которых вы учитесь



Thank you! Any questions?

Feedback: <https://forms.gle/ZvjrXYt9UKD8BkK48>

Ася Ройтберг, asya.roytberg@bigdatateam.org [LinkedIn](#)

ML Instructor at BigData Team, <http://bigdatateam.org/>

<https://www.facebook.com/bigdatateam/>