

Кластеризация и сегментация без учителя – зачем и как?

Антон Колонин

akolonin@aigents.com

<https://github.com/aigents>

Telegram: akolonin



<https://www.nsu.ru/>



<https://agirussia.org>

Есть ли у нас план?

Что такое кластеризация и сегментация?

Зачем они нужны?

Повышение эффективности обучения с подкреплением

Прогнозная аналитика на временных рядах

Интерпретируемая обработка текстов

Кластеризация без учителя – это как?

Ставим опыты на участниках

Обсуждаем метрики и результаты эксперимента

Сегментация без учителя – а так можно?

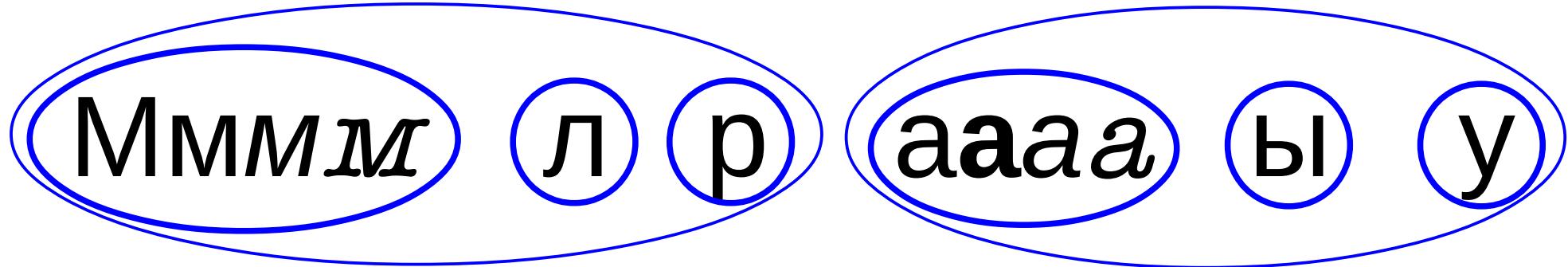
Загадываем загадку “по теме”

Обсуждаем метрики/параметры и результаты эксперимента

Подтверждаем отгадки

Делаем выводы!

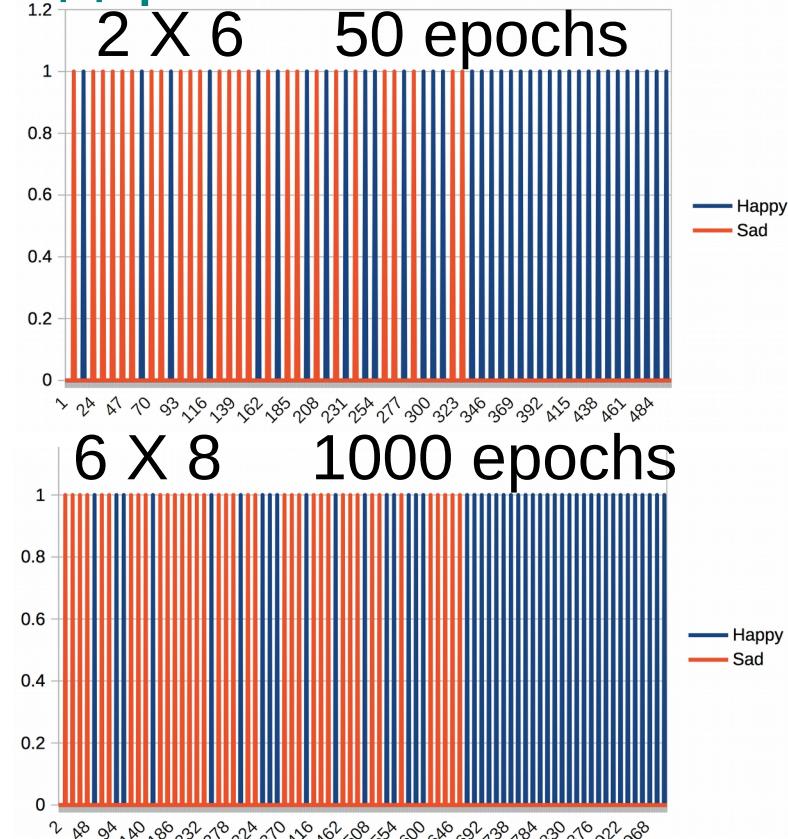
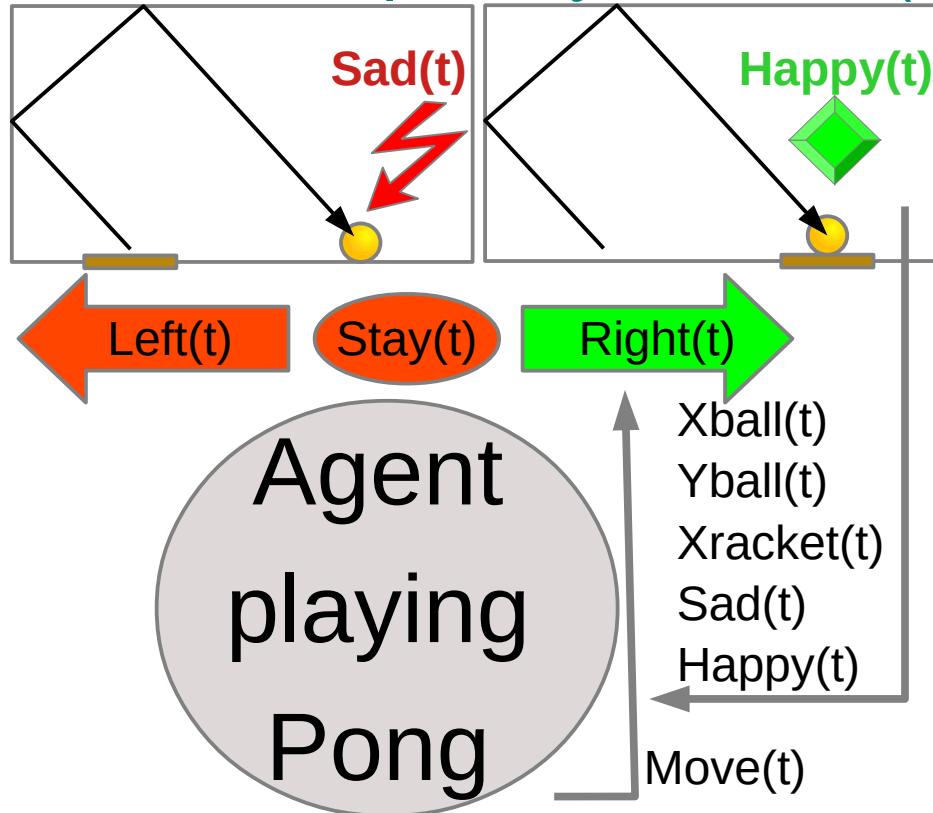
Кластеризация и сегментация



Мама мыла раму



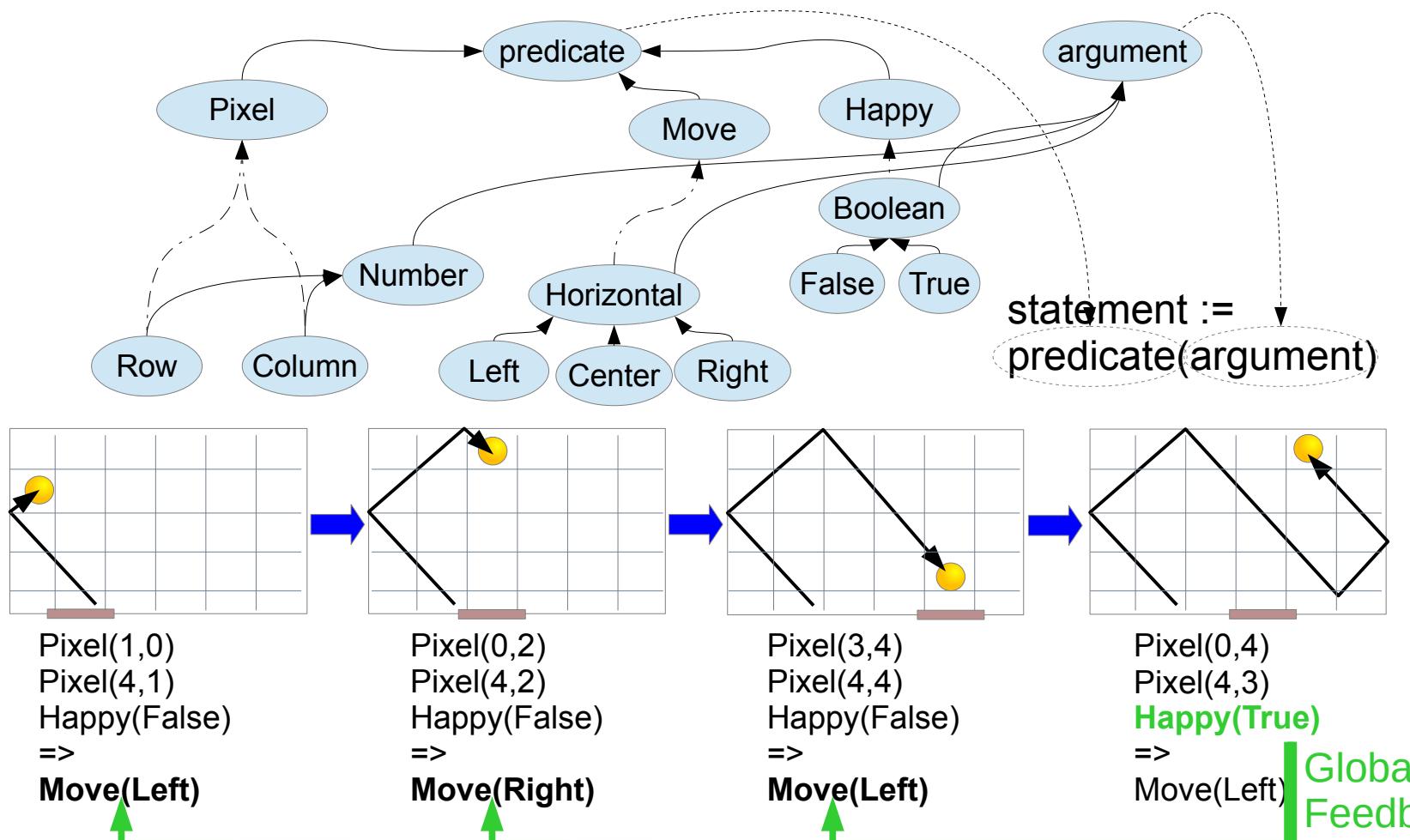
Выявление (не)успешных последовательностей действий при обучении с (само)подкреплением



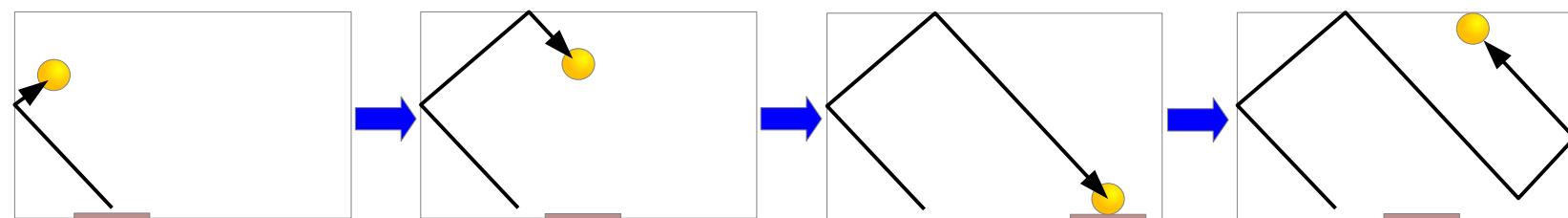
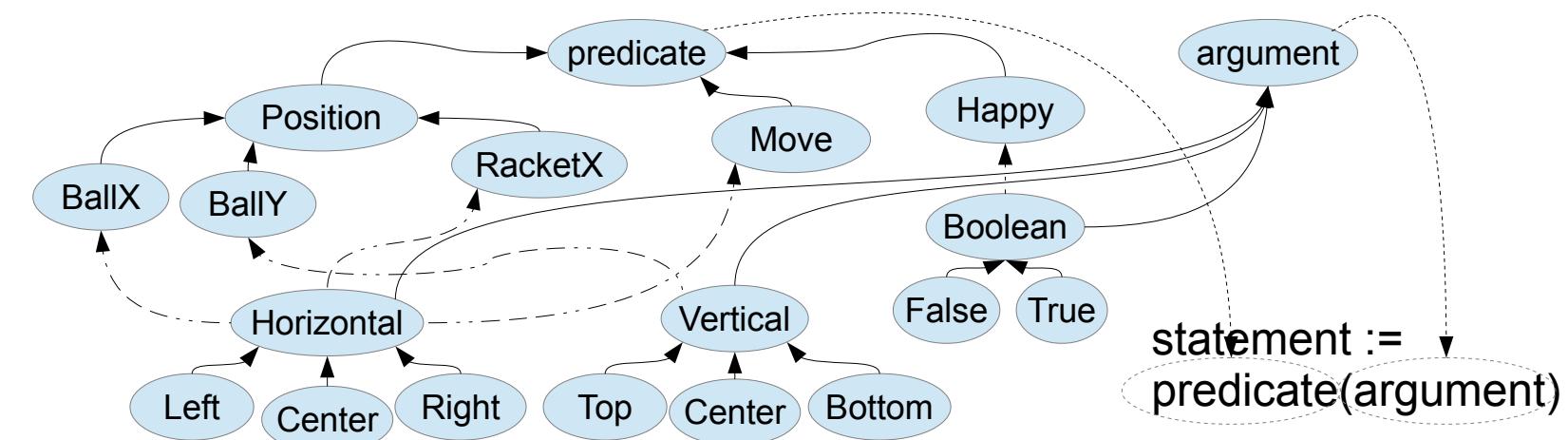
Unsupervised Learning of Temporal Abstractions with Slot-based Transformers
Anand Gopalakrishnan, Kazuki Irie, Jürgen Schmidhuber, Sjoerd van Steenkiste
<https://arxiv.org/abs/2203.13573>

<https://www.youtube.com/watch?v=2LPLhJKh95g>
<https://www.springerprofessional.de/neuro-symbolic-architecture-for-experiential-learning-in-discret/20008336>
<https://github.com/agents/agents-java/tree/master/src/main/java/net/webstructor/agi>
Copyright © 2022 Anton Kolonin, Agents®

“Дискретное” операционное пространство Pong



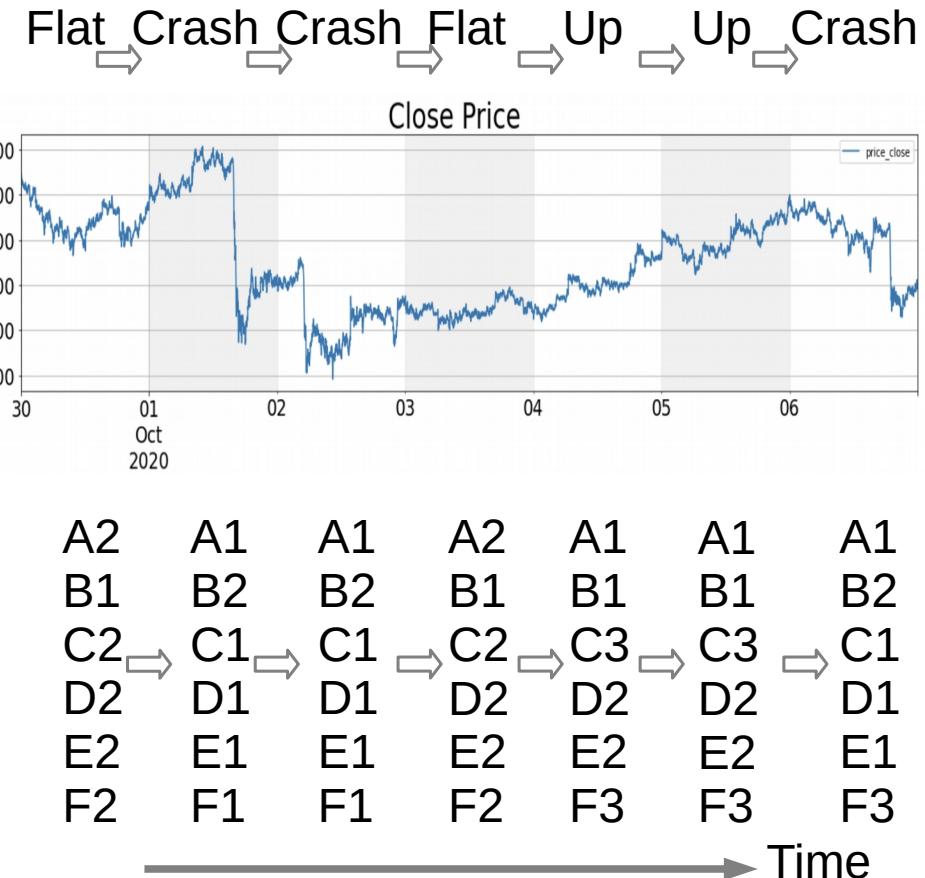
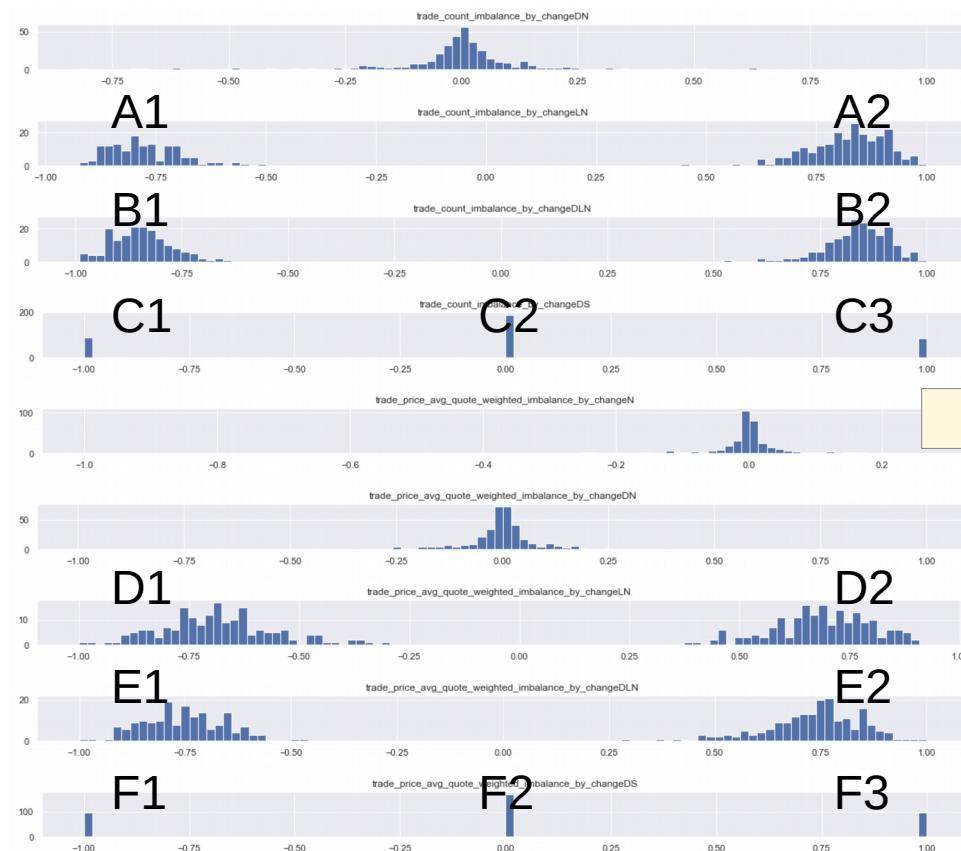
“Символьное” операционное пространство Pong



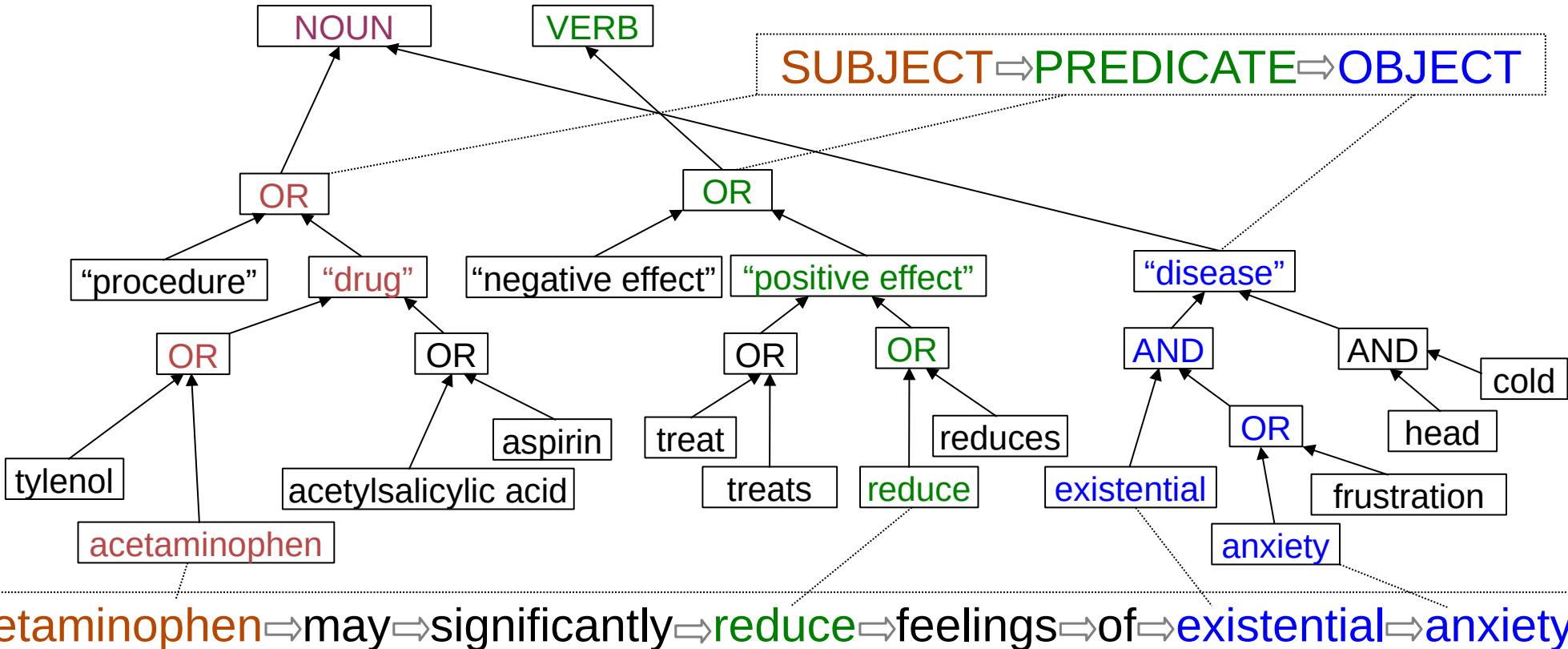
Global
Feedback

“Символизация” состояний рынка

“Квантификация” временных рядов



Построение гетерархий понятий, отношений, грамматических структур, слов и пунктуации при обучении языку без учителя



<https://ieeexplore.ieee.org/document/7361868>
<https://github.com/aigents/aigents-java>
<https://github.com/aigents/aigents-java-nlp>

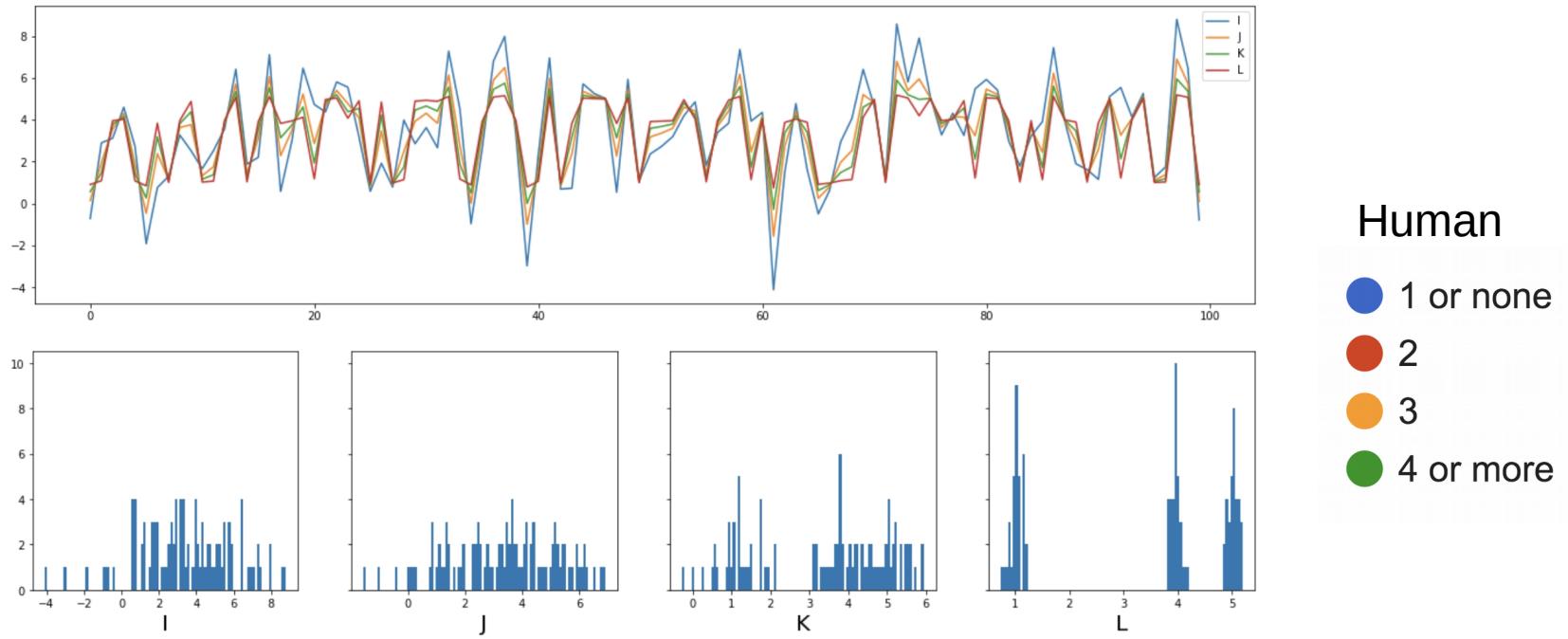
<https://www.springerprofessional.de/unsupervised-language-learning-in-opencog/15995030>
<https://www.springerprofessional.de/en/programmatic-link-grammar-induction-for-unsupervised-language-le/17020348>
<https://github.com/singnet/language-learning/>

Кластеризация без учителя

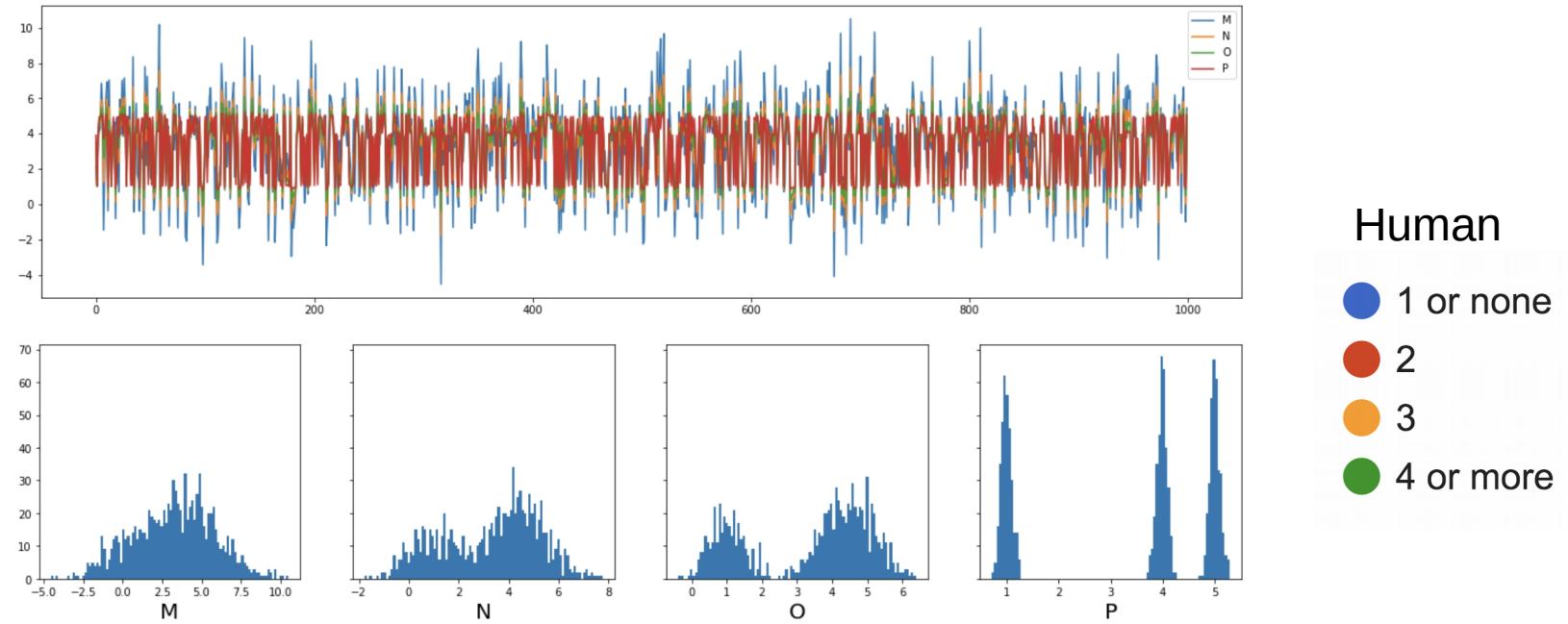
Не задаем магических значений
“K” (требуемого числа кластеров)
и порогов для агрегации!

https://github.com/aigents/pygents/blob/main/notebooks/cluster/distribution_modes.ipynb

Квантификация – Опыт 1



Квантификация – Опыт 2

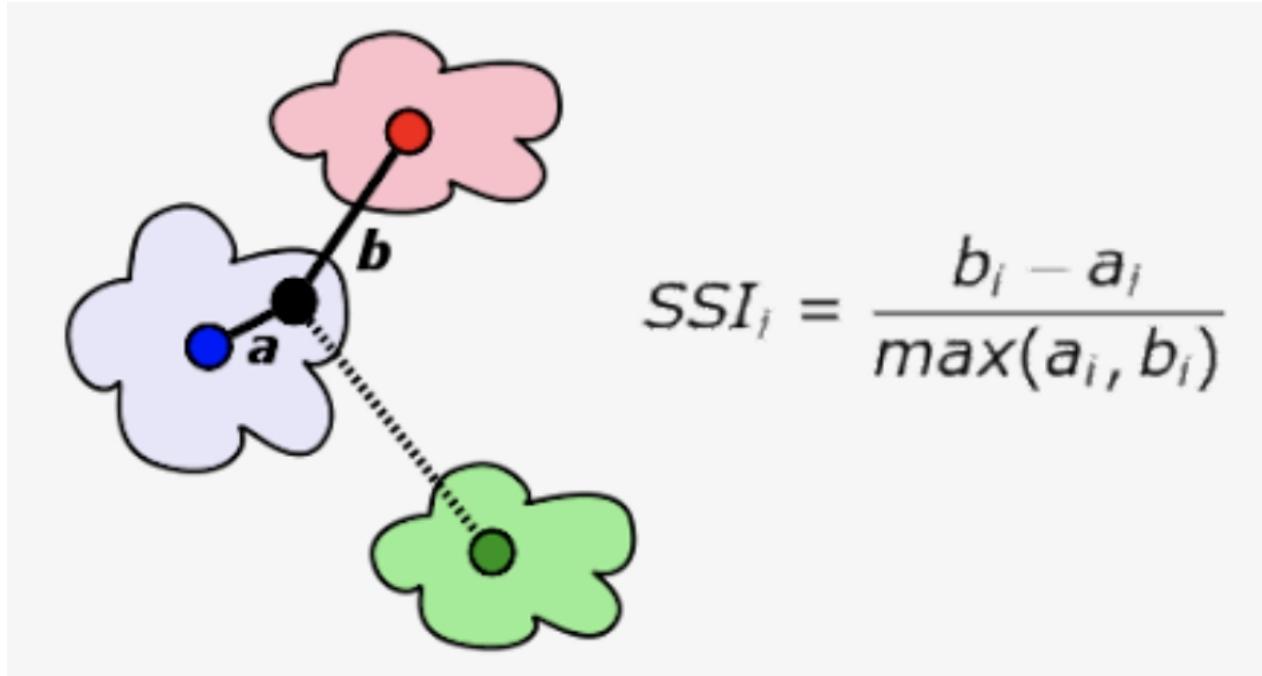


Метрики кластеризации/квантификации

1. Maximizing "Silhouette Coefficient" (SC), which appears more human-intuitive but does not work for K=1

[https://en.wikipedia.org/wiki/Silhouette_\(clustering\)](https://en.wikipedia.org/wiki/Silhouette_(clustering))

<https://towardsdatascience.com/silhouette-coefficient-validating-clustering-techniques-e976bb81d10c>



<http://sujitpal.blogspot.com/2018/03/an-implementation-of-silhouette-score.html>

Метрики кластеризации/квантификации

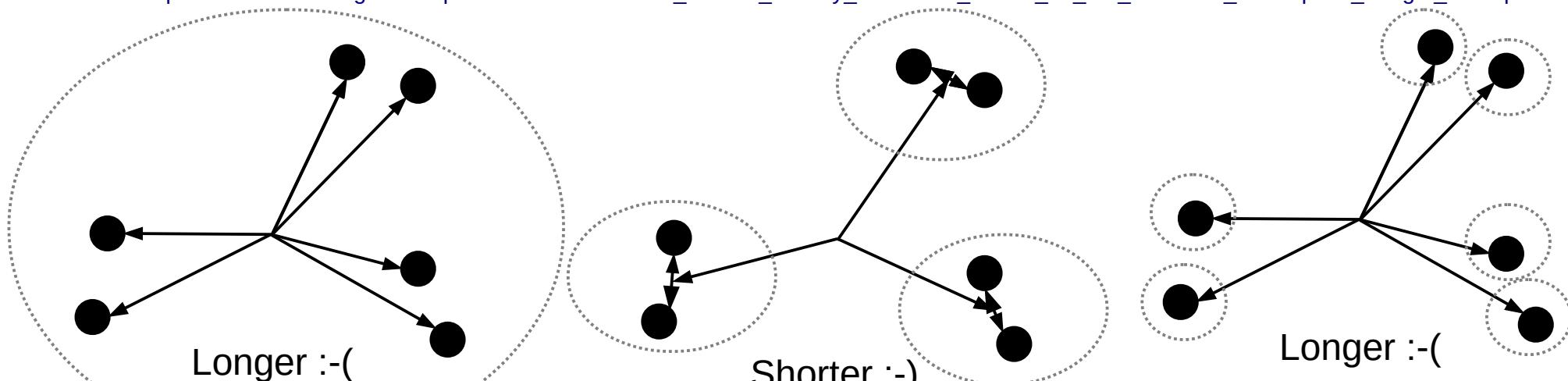
1. Maximizing "Silhouette Coefficient" (SC), which appears more human-intuitive but does not work for K=1

[https://en.wikipedia.org/wiki/Silhouette_\(clustering\)](https://en.wikipedia.org/wiki/Silhouette_(clustering))

<https://towardsdatascience.com/silhouette-coefficient-validating-clustering-techniques-e976bb81d10c>

2. Minimizing "Normalized Centroid Distance" (NCD) - based on "minimum description length" idea, works for K=1, does not align with human "reductionist" intuition for diverse distributions (tends to create more clusters than needed)

https://www.researchgate.net/publication/221020638_Cluster_Validity_Measures_Based_on_the_Minimum_Description_Length_Principle



Метрики кластеризации/квантификации

1. Maximizing "Silhouette Coefficient" (SC), which appears more human-intuitive but does not work for K=1

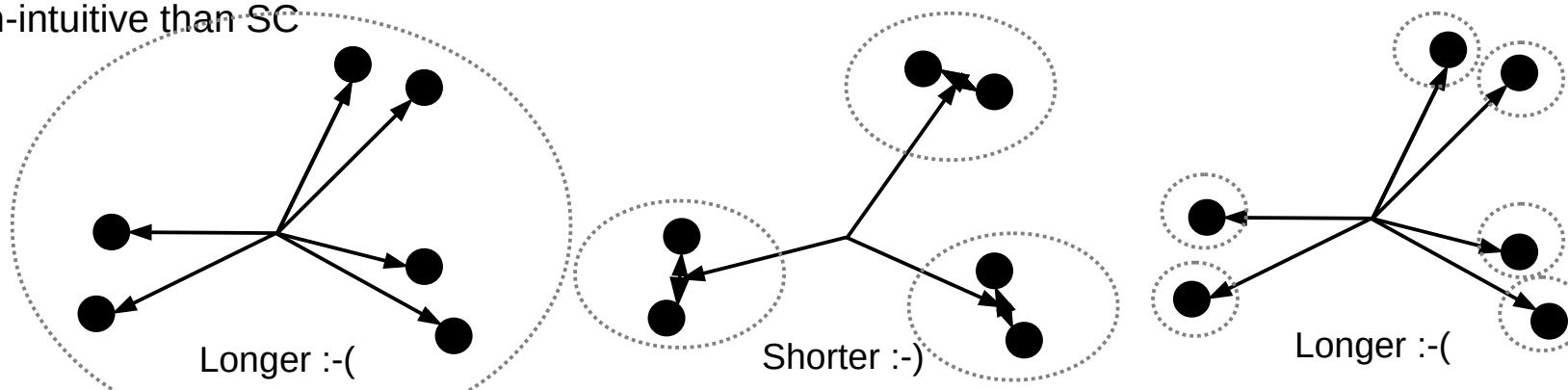
[https://en.wikipedia.org/wiki/Silhouette_\(clustering\)](https://en.wikipedia.org/wiki/Silhouette_(clustering))

<https://towardsdatascience.com/silhouette-coefficient-validating-clustering-techniques-e976bb81d10c>

2. Minimizing "Normalized Centroid Distance" (NCD) - based on "minimum description length" idea, works for K=1, does not align with human "reductionist" intuition for diverse distributions (tends to create more clusters than needed)

https://www.researchgate.net/publication/221020638_Cluster_Validity_Measures_Based_on_the_Minimum_Description_Length_Principle

3. Minimizing "Normalized Centroid Distance times Centroids" (NCDC) - extends NCD multiplying it by number of clusters to penalize creation of too many clusters, works for K=1, more human-intuitive than NCD but less human-intuitive than SC



Метрики кластеризации/квантификации

1. Maximizing "Silhouette Coefficient" (SC), which appears more human-intuitive but does not work for K=1

[https://en.wikipedia.org/wiki/Silhouette_\(clustering\)](https://en.wikipedia.org/wiki/Silhouette_(clustering))

<https://towardsdatascience.com/silhouette-coefficient-validating-clustering-techniques-e976bb81d10c>

2. Minimizing "Normalized Centroid Distance" (NCD) - based on "minimum description length" idea, works for K=1, does not align with human "reductionist" intuition for diverse distributions (tends to create more clusters than needed)

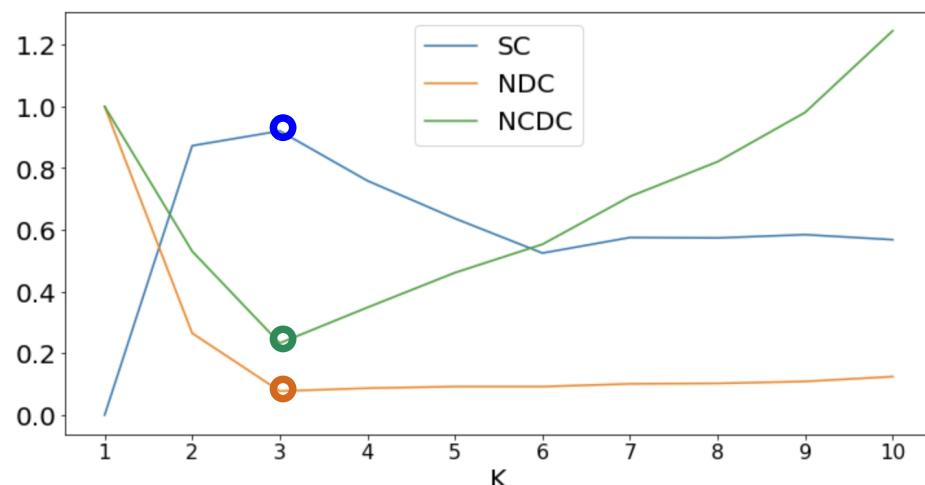
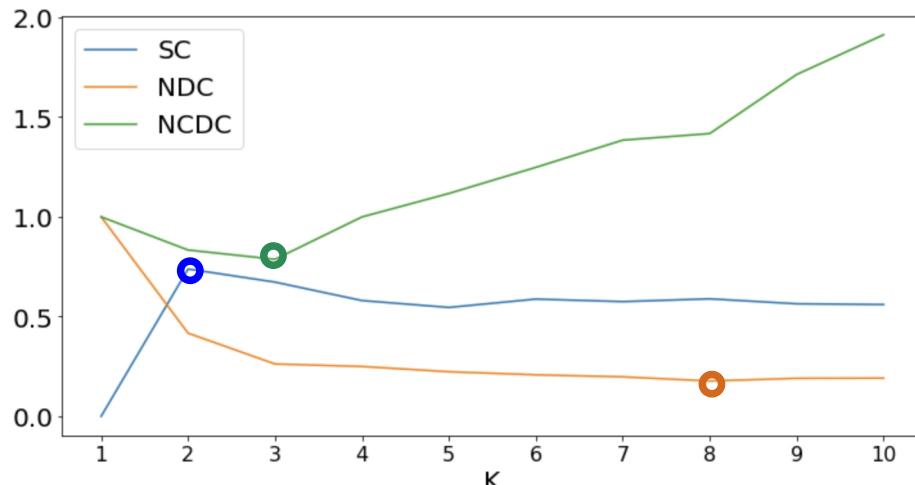
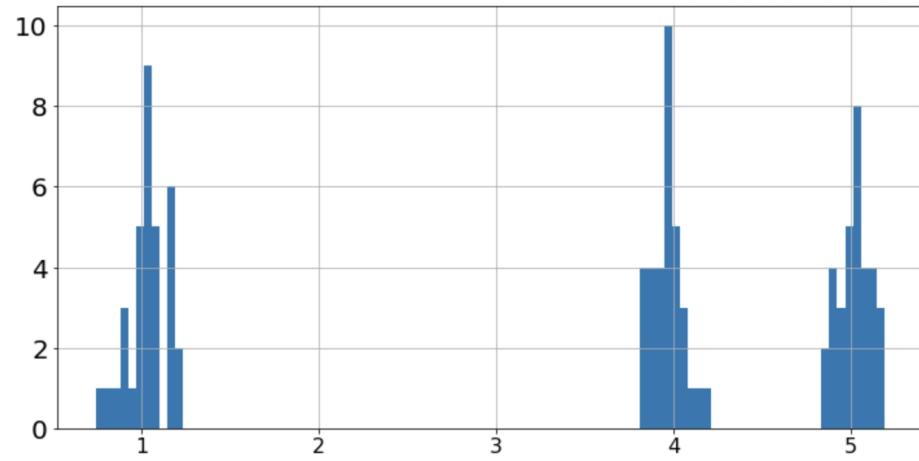
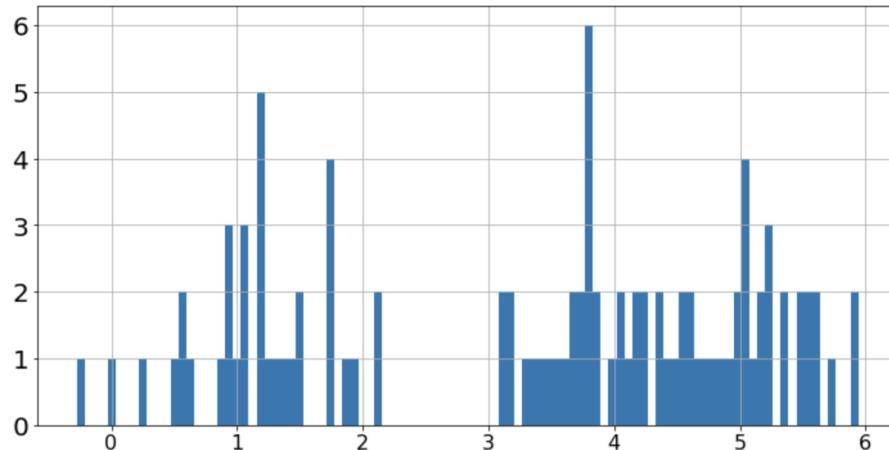
https://www.researchgate.net/publication/221020638_Cluster_Validity_Measures_Based_on_the_Minimum_Description_Length_Principle

3. Minimizing "Normalized Centroid Distance times Centroids" (NCDC) - extends NCD multiplying it by number of clusters to penalize creation of too many clusters, works for K=1, more human-intuitive than NCD but less human-intuitive than SC

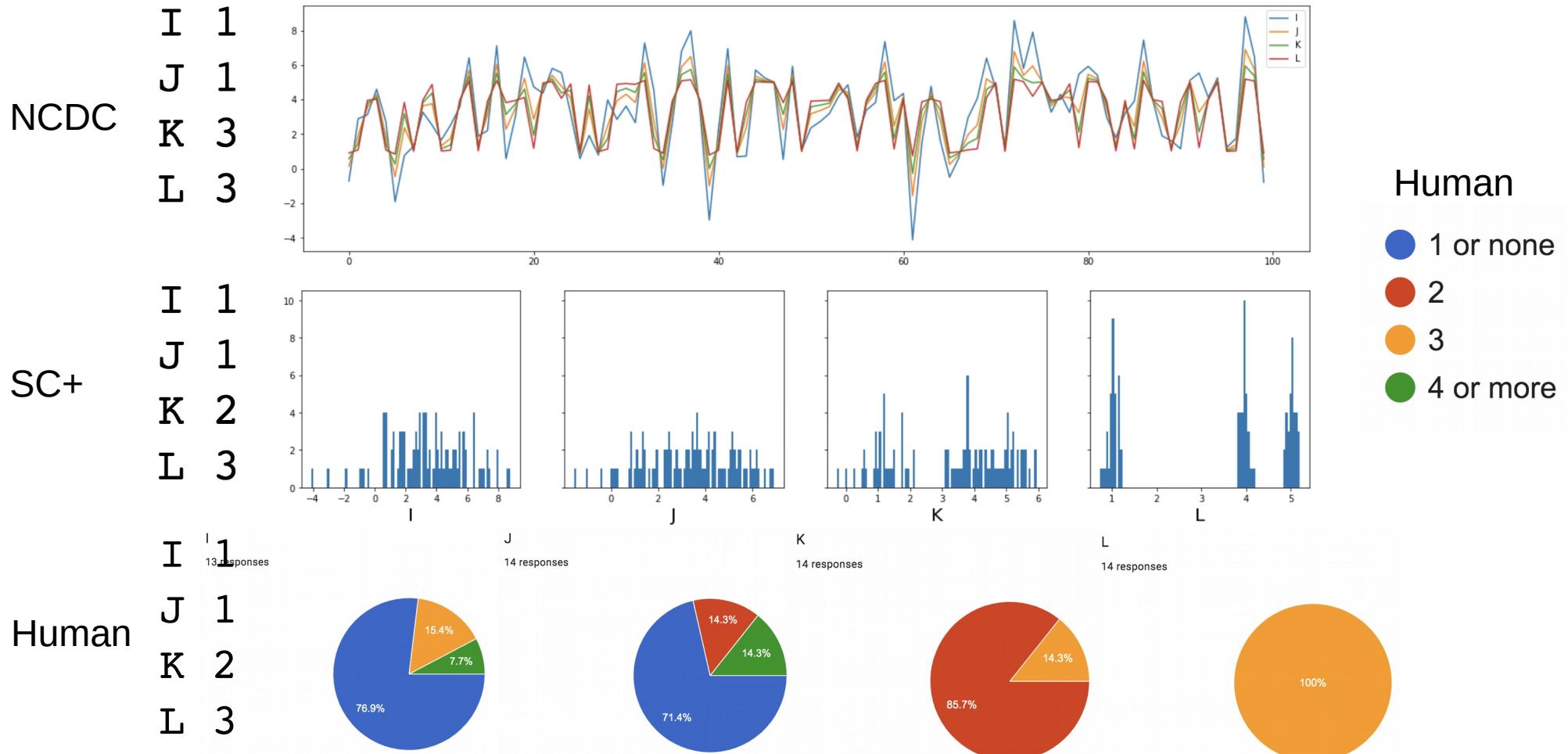
4. Using SC+ (maximize SC if it is above threshold 0.65 or minimize NCDC otherwise) - seems generally optimal from human intuition perspective

https://github.com/aigents/pygents/blob/main/notebooks/cluster/distribution_modes.ipynb

Метрики квантификации (1-D кластеризации)

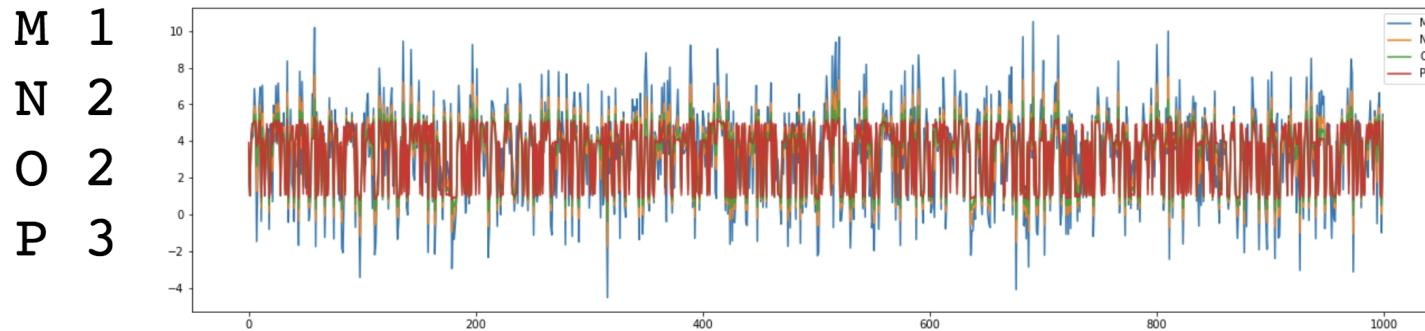


Квантификация – Опыт 1

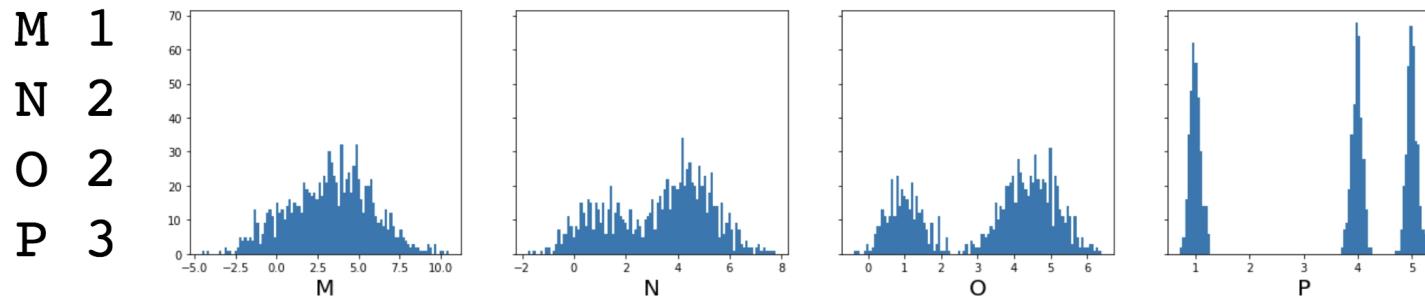


Квантификация – Опыт 2

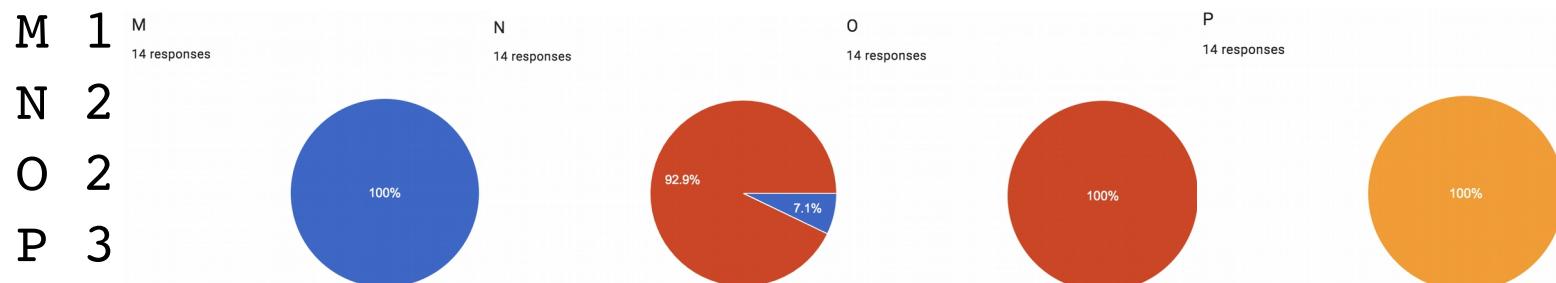
NCDC



SC+



Human



Human

- 1 or none
- 2
- 3
- 4 or more

Квантификация: Люди против Машины

Evaluate with "Fleiss' kappa" (FK) and "Krippendorff's alpha" (KA)

https://en.wikipedia.org/wiki/Fleiss%27_kappa

<https://stackoverflow.com/questions/51919897/is-fleiss-kappa-a-reliable-measure-for-interannotator-agreement-the-following-r>

<https://bmcmedresmethodol.biomedcentral.com/articles/10.1186/s12874-016-0200-9>

https://www.statsmodels.org/dev/generated/statsmodels.stats.inter_rater.fleiss_kappa.html

NCDC vs. SC+: 0.55, 0.56 (Moderate agreement)

Humans vs. humans: 0.59, 0.59 (Moderate agreement)

NCDC vs. humans: 0.47, 0.48 (Moderate agreement)

SC+ vs. humans: **0.92, 0.92** (Almost perfect agreement)

Сегментация (текста) без учителя

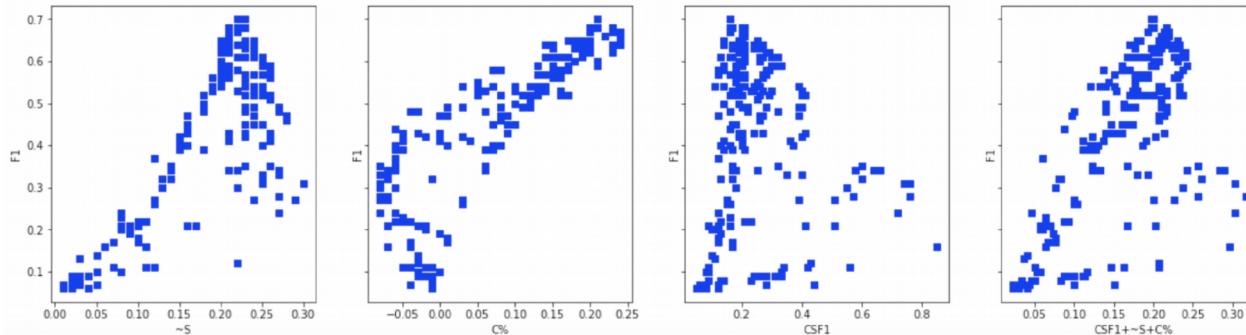
Заранее не знаем ничего о старт/стоп токенах,
разделителях, лексиконах и правилах пунктуации!!!

<https://github.com/aigents/pygents/tree/main/notebooks/nlp>

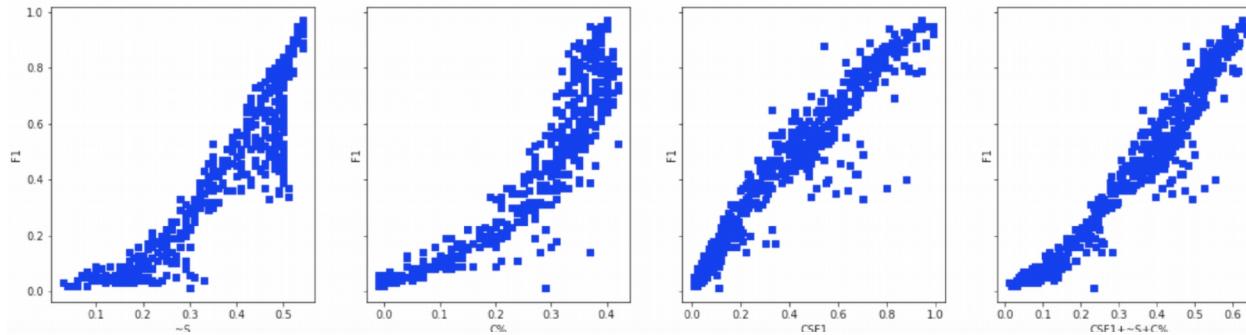
<https://arxiv.org/abs/2205.11443>

Угадай язык – где Русский, Китайский и Английский?

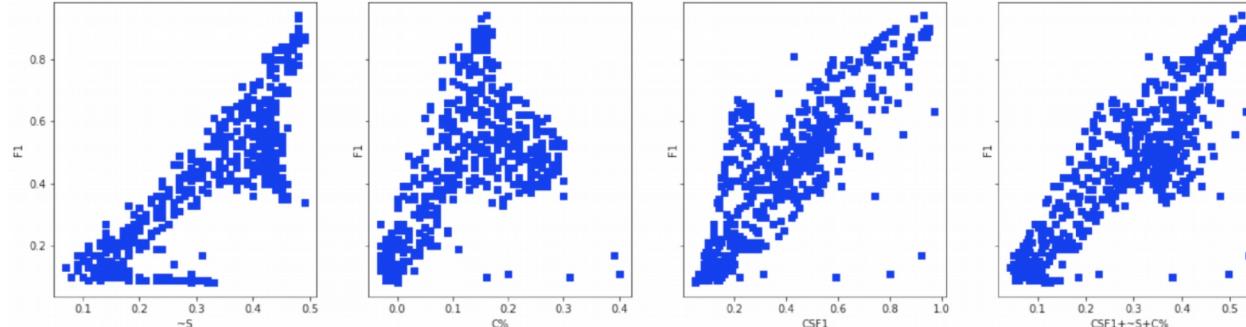
Язык 1



Язык 2



Язык 3



Токенизация или сегментация текста без учителя – подходы предшественников

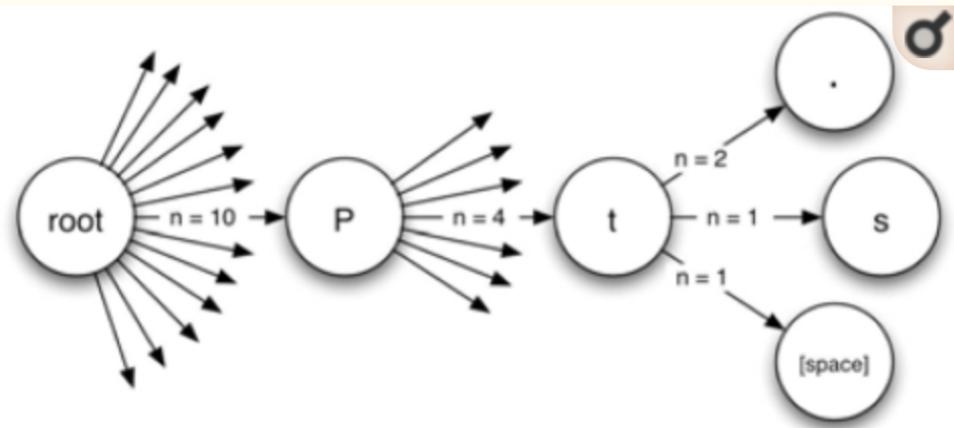


Figure 1

Trie data structure. The probability of observing an ‘s’ given the preceding string “Pt” is $\frac{1}{4}$, or 25%. The freedom following “pt” is 3.

Metrics/Indicators:

Mutual Information¹
Conditional Probability^{1,2}
Transition Freedom^{2,3}

¹ <https://scholarsarchive.byu.edu/cgi/viewcontent.cgi?article=6983&context=etd>

² <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC2655800/>

³ Karl Friston. The free-energy principle: a unified brain theory? <https://www.nature.com/articles/nrn2787>

Эксперимент – корпуса и методология

Train corpora

Chinese

CLUE News 2016 Validation – 270M

CLUE News 2016 Train – 8,500M

English

Brown – 6M

Gutenberg Children – 29M

Gutenberg Adult – 140M

Social Media – 68M

All above – combined

Russian

RusAge Test – 141M

RusAge Previews – 825M

Test corpus

Parallel Chinese/English/Russian

– 100 multi-sentence statements on finance

Metrics/Indicators:

Ngram (Character)

Probability or Conditional Transition Probability (p-/p+)

Deviation (dvp-/dvp+) from mean

Derivative (dp-/dp+) and “Peak”

Transition Freedom (f-/f+)

Deviation (dvf-/dvf+) from mean

Derivative (df-/df+) and “Peak”

Hyper-parameters:

Combination of Ngram ranks N ([1],[2],[3],[1,2],[1,2,3],...)

Threshold for model compression

Threshold for segmentation

Evaluations:

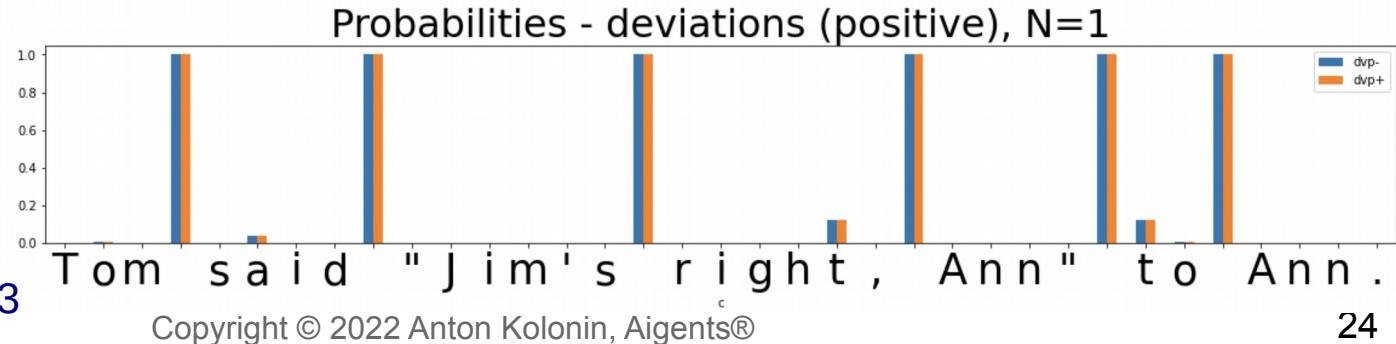
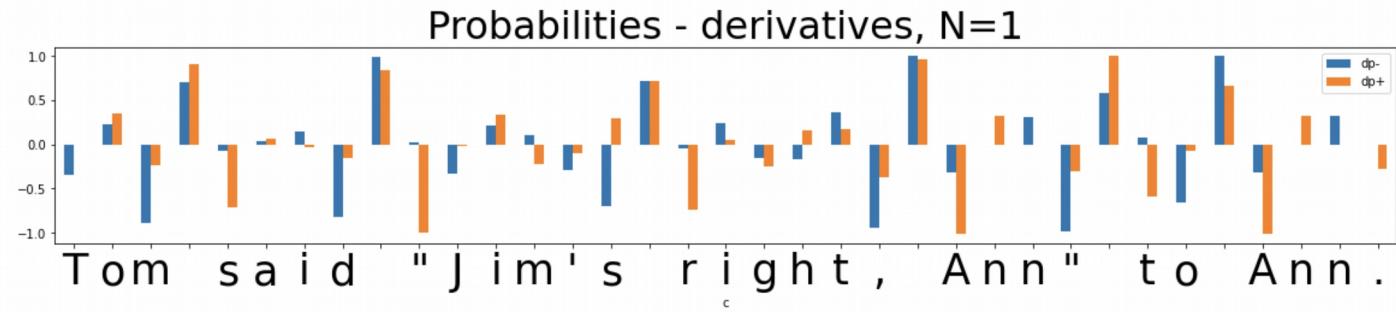
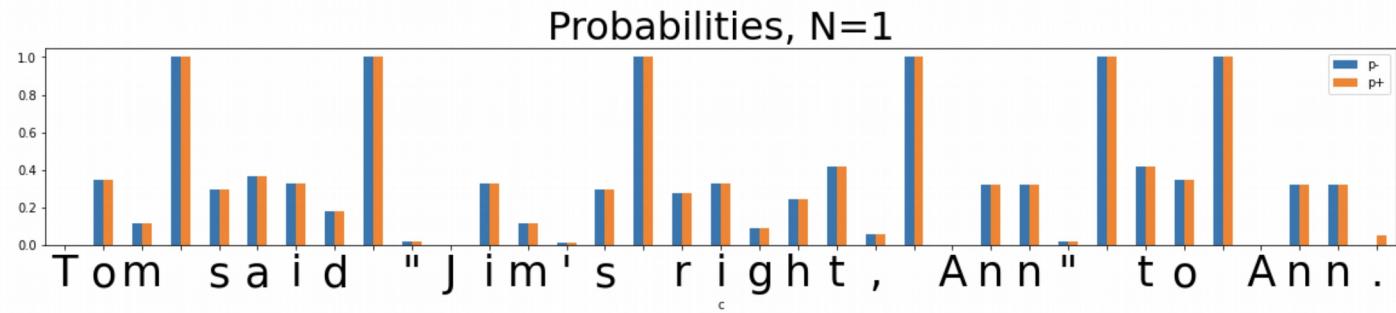
Tokenization F1, on set of tokens found comparing to delimiter-based (English/Russian) or Jieba (Chinese)

Precision on set of tokens found comparing to reference lexicons

<https://arxiv.org/abs/2205.11443>

Токенизация по вероятностям токенов/N-gram

Metrics/Indicators:
Ngram (Character)
Probability



Токенизация по условным вероятностям N-gram

Metrics/Indicators:

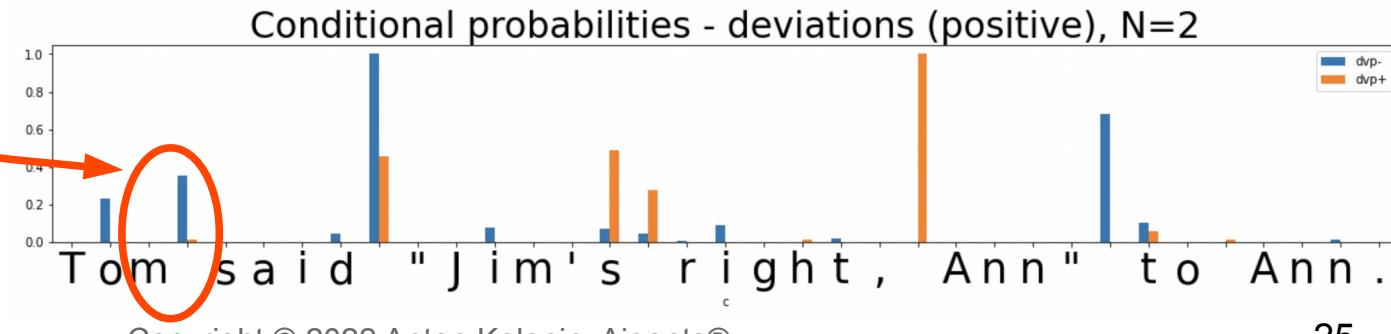
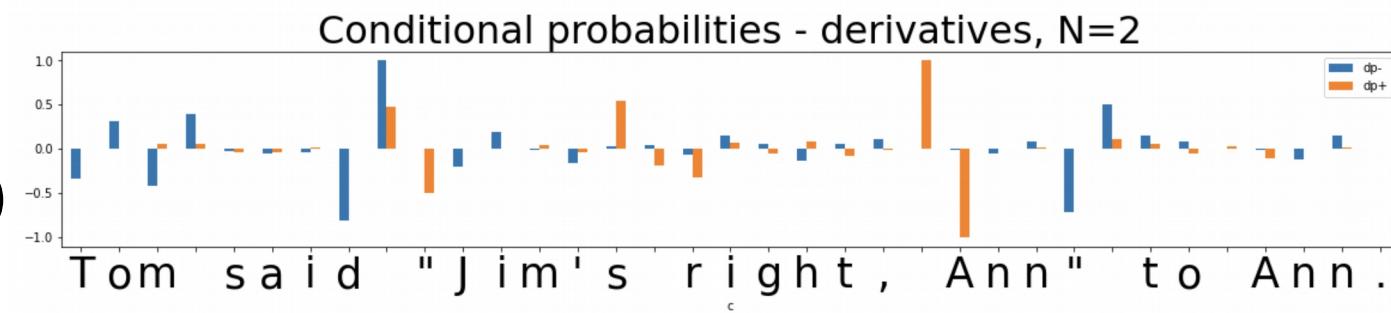
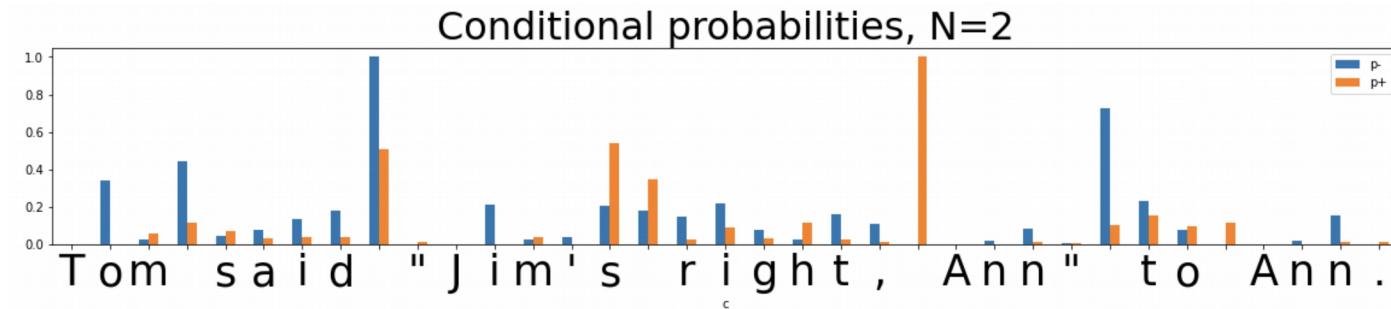
Ngram (Character)

Conditional
Probability

(of Transition)

$P(\text{Ngram}_{n+1})/P(\text{Ngram}_n)$

$P("m")/P(m")$



Токенизация по “свободам перехода” (производным)

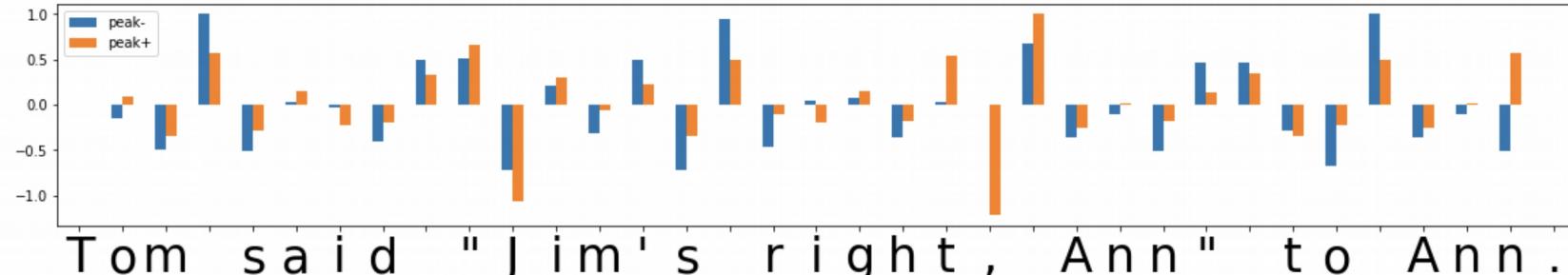
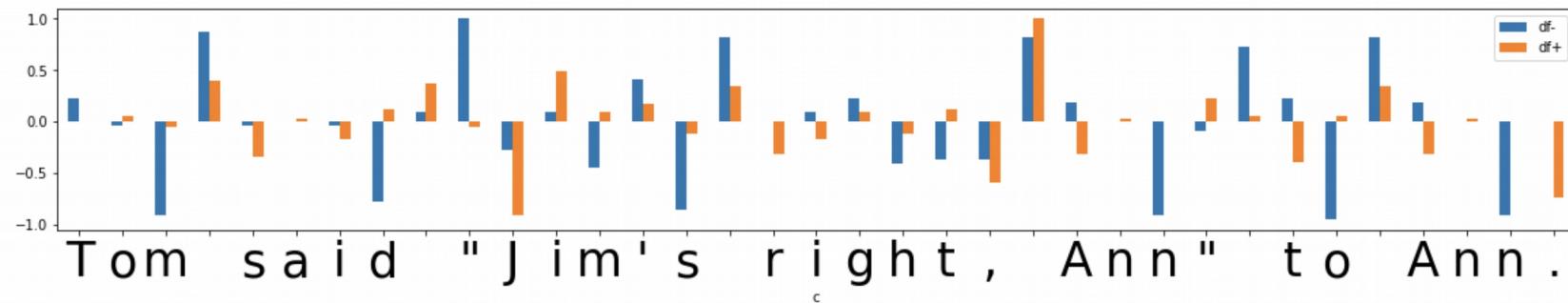
Metrics/ Indicators:

Transition
Freedom
(Freedom of
Transition)
Derivative

Transition
Freedom
“Peak”)

Threshold 0.25
Tom said "Jim's right, Ann" to Ann.
['Tom', ' ', 'said', ' ', "'", "Jim's", ' ', 'right', ' ', ' ', 'Ann', ' ", " ', 'to', ' ', 'Ann', '.']
['Tom', ' ', 'said', ' ', "'", 'Ji', 'm', "'s", ' ', 'right', ' ', ' ', 'Ann', ' ", " ', 'to', ' ', 'Ann', '.']
0.89

Threshold 0.35
Tom said "Jim's right, Ann" to Ann.
['Tom', ' ', 'said', ' ', "'", "Jim's", ' ', 'right', ' ', ' ', 'Ann', ' ", " ', 'to', ' ', 'Ann', '.']
['Tom', ' ', 'said', ' ', "'", 'Jim', "'s", ' ', 'right', ' ', ' ', 'Ann', ' ", " ', 'to', ' ', 'Ann', '.']
0.82



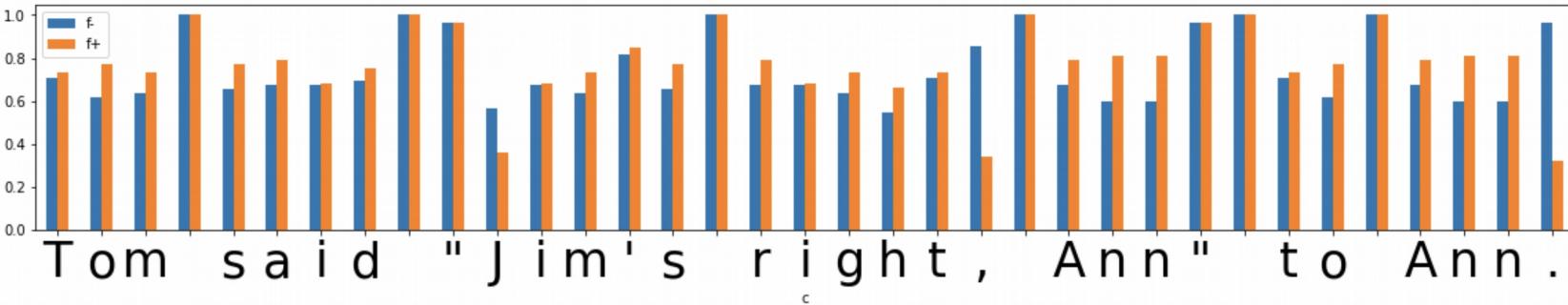
Токенизация по “свободам перехода” (отклонениям)

```
Threshold 0.25
Tom said "Jim's right, Ann" to Ann.
['Tom', ' ', 'said', ' ', ' ', "Jim's", ' ', 'right', ' ', ' ', 'Ann', ' ', ' ', 'to', ' ', ' ', 'Ann', '.']
['Tom', ' ', 'said', ' ', ' ', "Jim", " ", "s", ' ', 'right', ' ', ' ', 'Ann', ' ', ' ', 'to', ' ', ' ', 'Ann', '.']
0.89
```

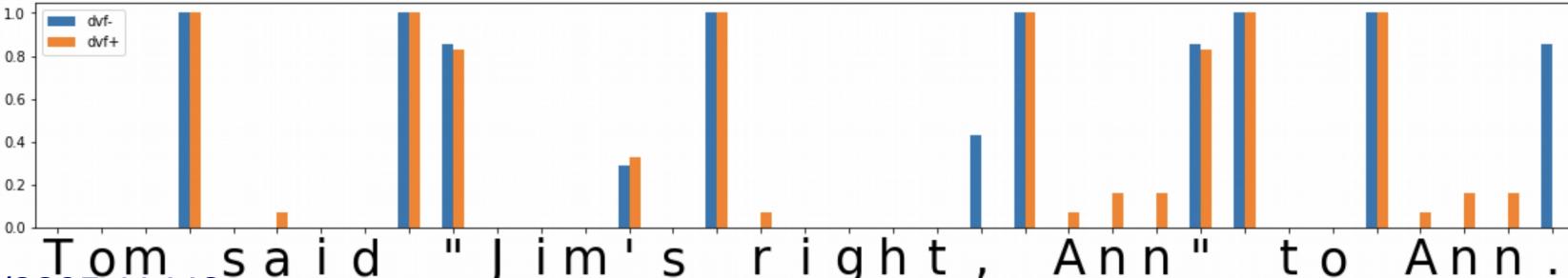
```
Threshold 0.35
Tom said "Jim's right, Ann" to Ann.
['Tom', ' ', 'said', ' ', ' ', "Jim's", ' ', 'right', ' ', ' ', 'Ann', ' ', ' ', 'to', ' ', ' ', 'Ann', '.']
['Tom', ' ', 'said', ' ', ' ', "Jim's", ' ', 'right', ' ', ' ', 'Ann', ' ', ' ', 'to', ' ', ' ', 'Ann', '.']
1.0
```

Metrics/ Indicators:

Transition
Freedom
(Freedom of
Transition)



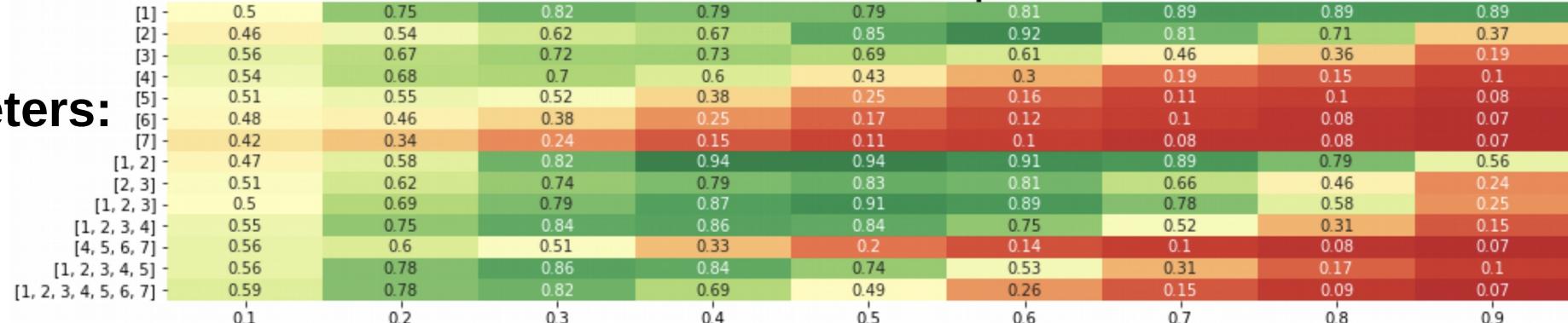
Transition Freedom Deviation



Токенизация – поиск гиперпараметров

English

F1 - Brown ddf- & ddf+ filter=0 parameters=10967135



Hyper-Parameters:

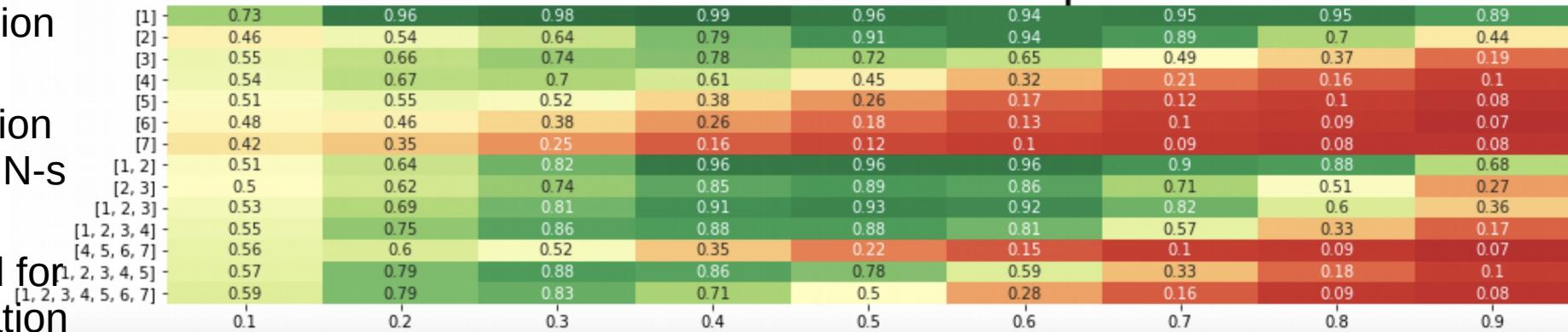
Metric:
Transition
Freedom

Threshold
for model
compression

Combination
of Ngram N-s

Threshold for
segmentation

F1 - Brown ddf- & ddf+ filter=0.0001 parameters=8643703

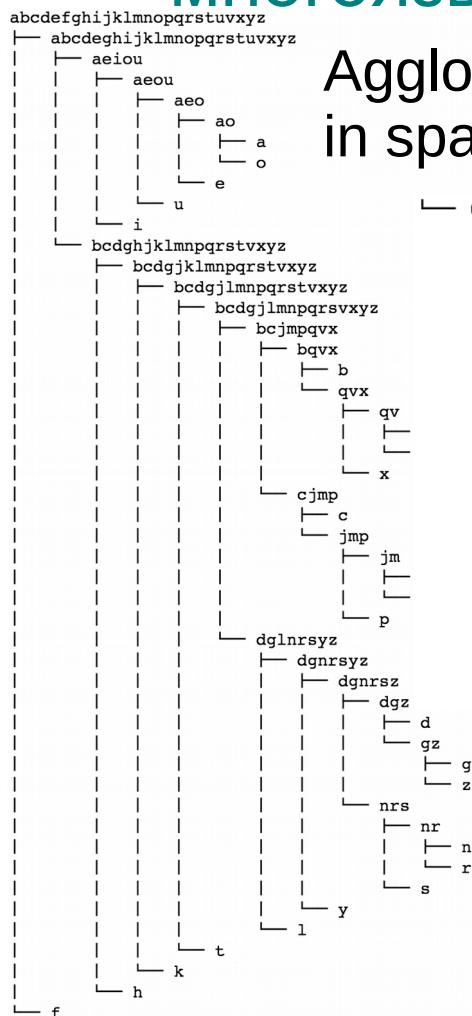


Оценки точности токенизации и обнаружения лексикона

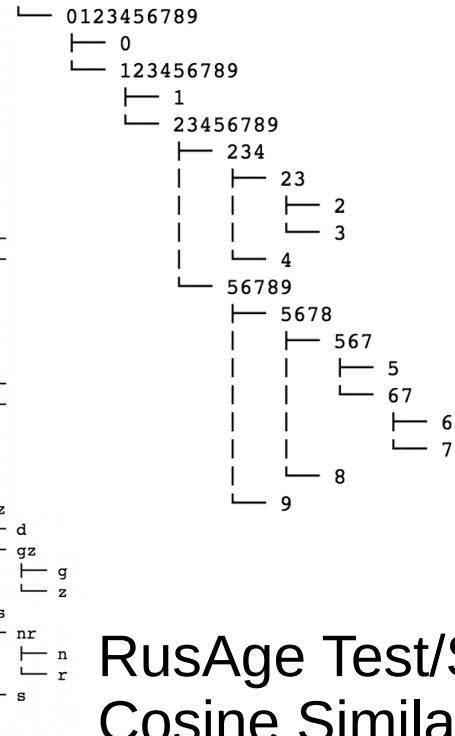
Language	Tokenizer	Tokenization F1	Lexicon Discovery Precision
English	Freedom-based	0.99	0.99 (vs 1.0)
English	Lexicon-based	0.99	-
English no spaces	Freedom-based	0.42	-
English no spaces	Lexicon-based	0.79	-
Russian	Freedom-based	1.0	1.0 (vs 1.0)
Russian	Lexicon-based	0.94	-
Russian no spaces	Freedom-based	0.26	-
Russian no spaces	Lexicon-based	0.72	-
Chinese	Freedom-based	0.71	0.92 (vs 0.94)
Chinese	Lexicon-based	0.83	-

Lexicon-based Tokenization - greedy/beam search on word length (optimal) or frequency

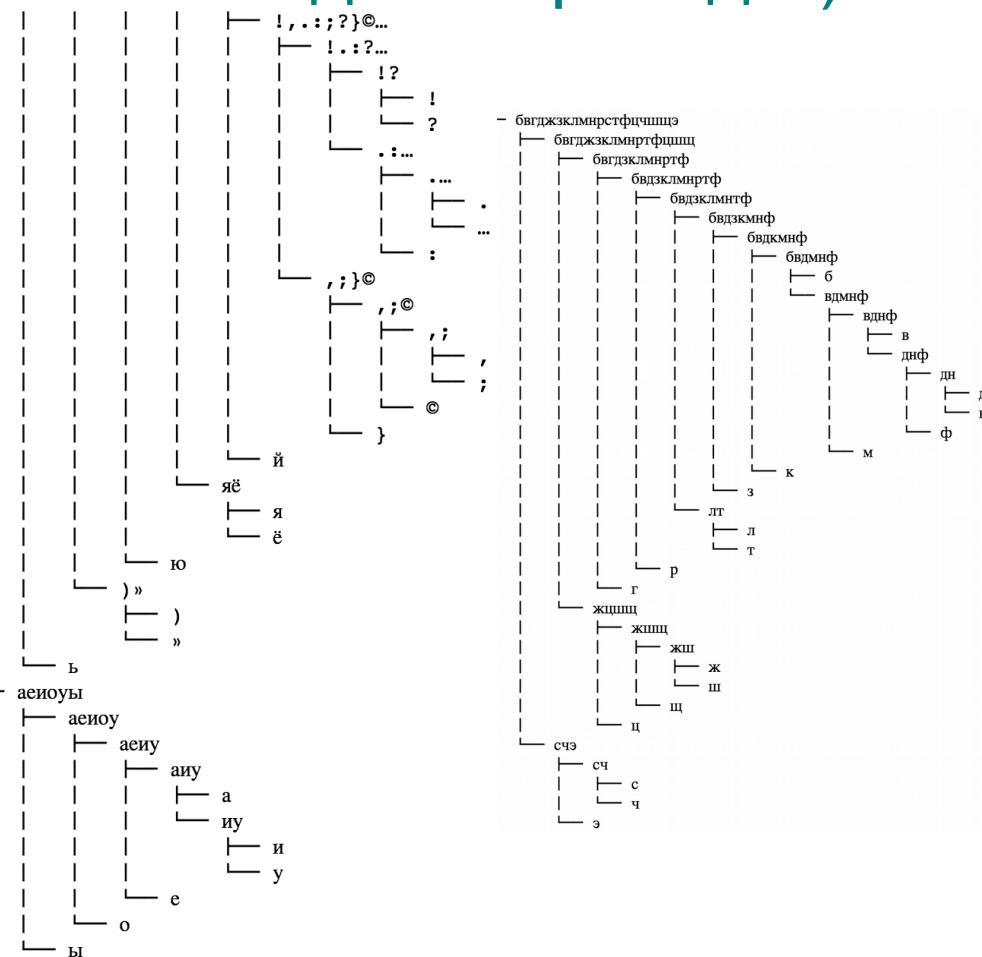
Классификация знаков пунктуации, гласных и согласных в многоязычных текстах (по “свободам перехода”)



Agglomerative Clustering in space of Transitions

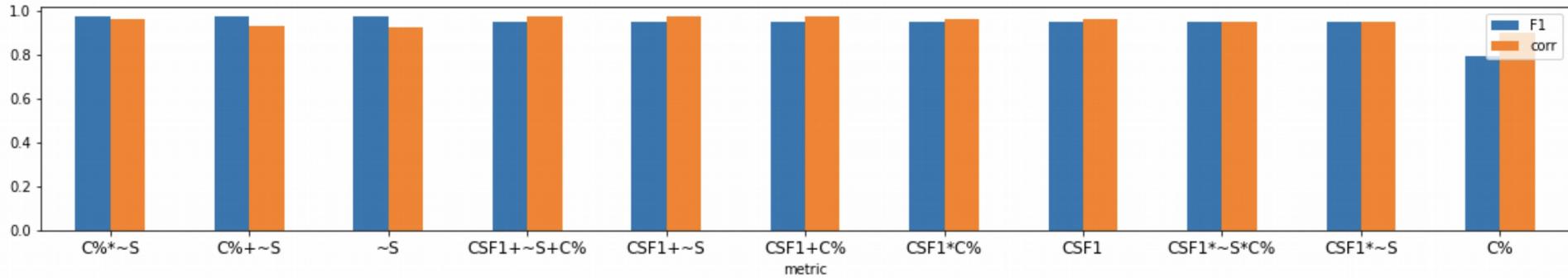


RusAge Test/Small, Cosine Similarity

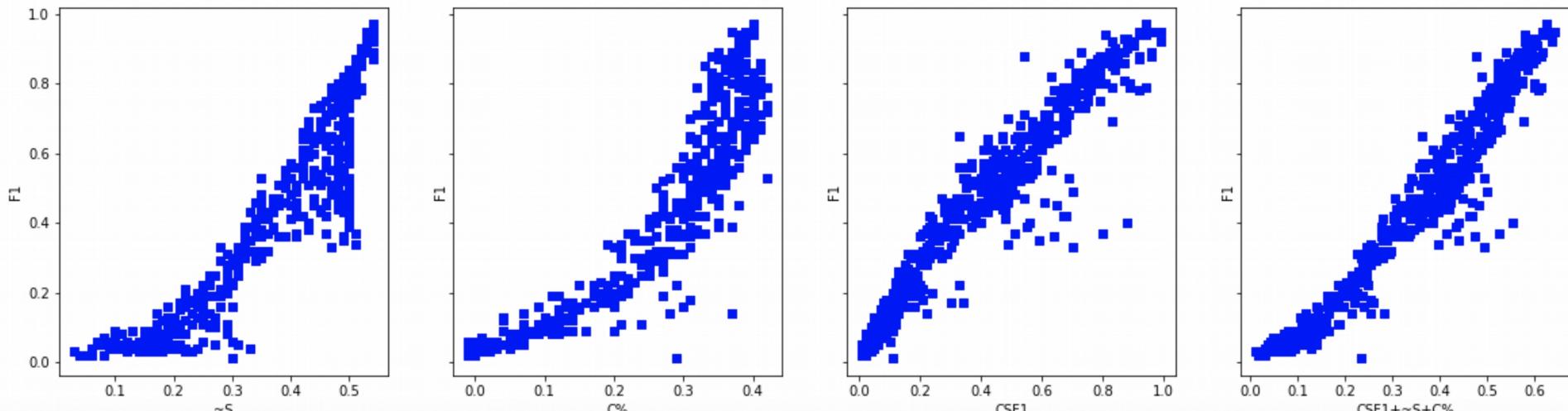


Самонастройка гиперпараметров – Английский (отклонение TF)

Test 1000

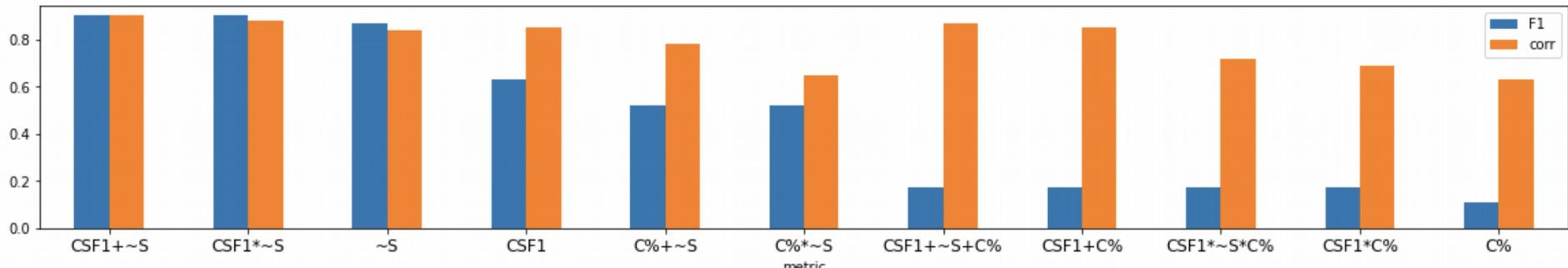


F1 as function of $\sim S$, C% and CSF1 used for hyper-parameter selection

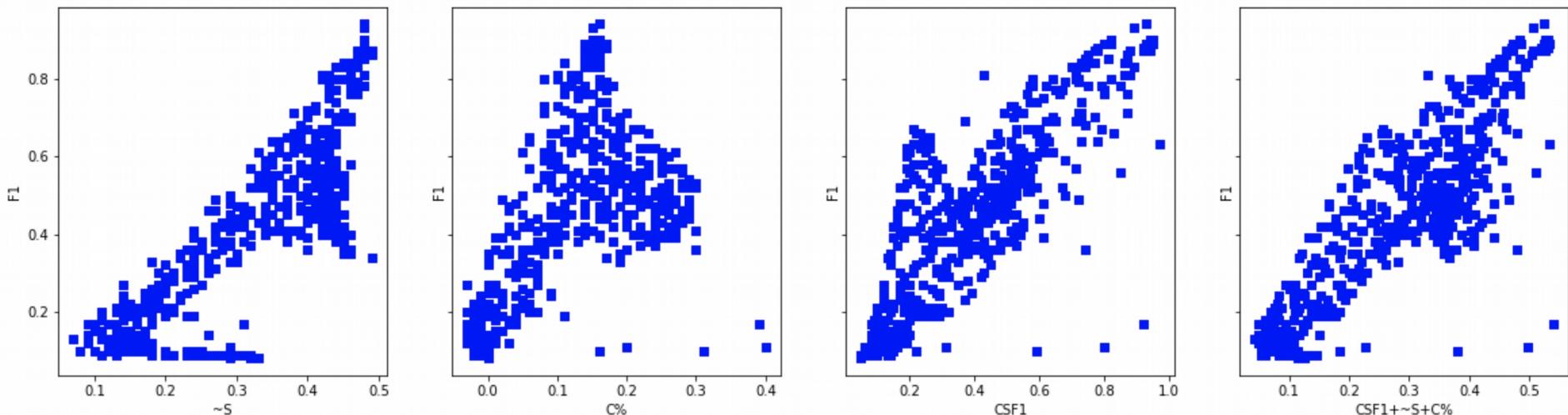


Самонастройка гиперпараметров – Русский (отклонение TF)

Test 1000

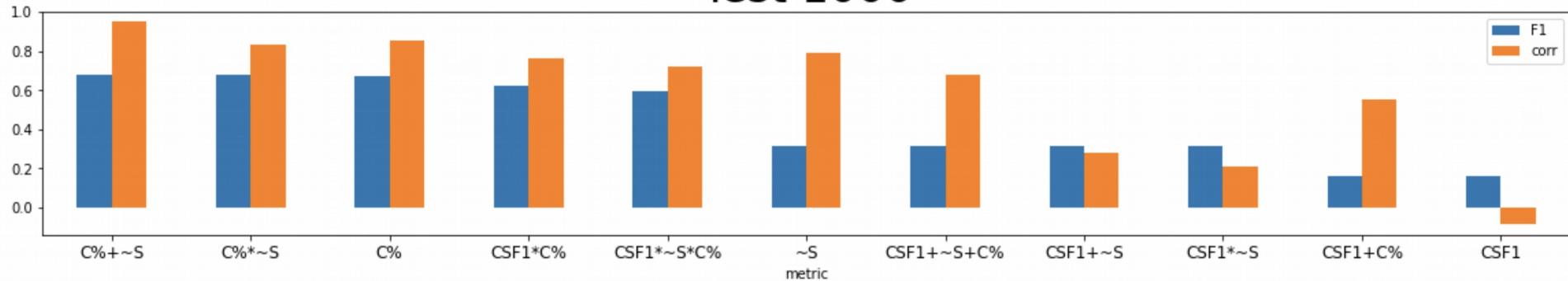


F1 as function of $\sim S$, C% and CSF1 used for hyper-parameter selection

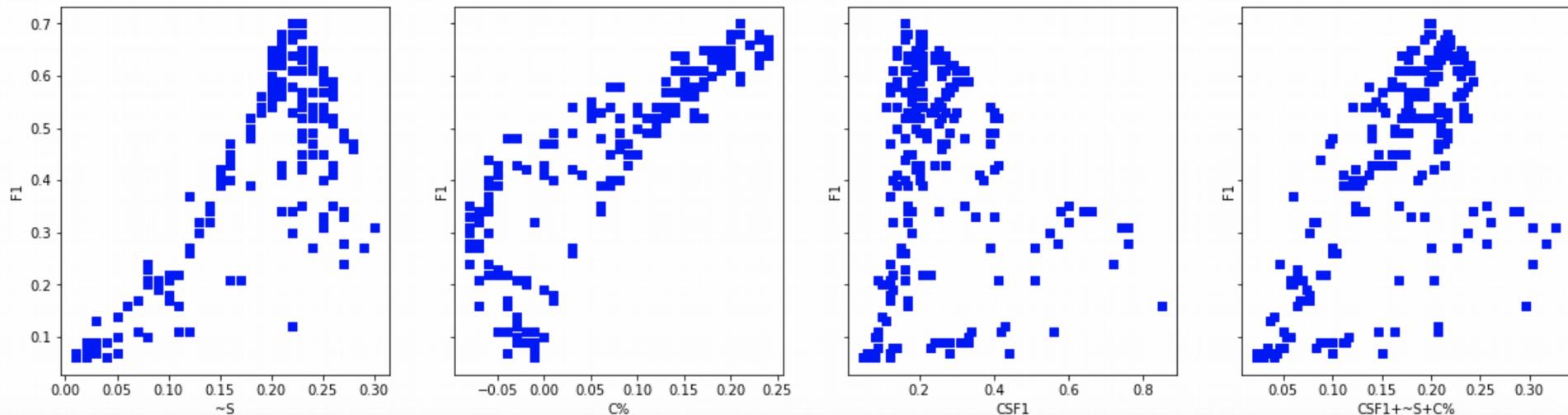


Самонастройка гиперпараметров – Китайский (“пики” TF)

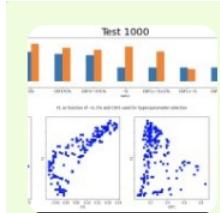
Test 1000



F1 as function of $\sim S$, C% and CSF1 used for hyper-parameter selection



Вы все правильно угадали!



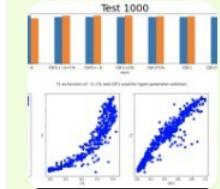
Screen Shot 2022-06-16 at 11.08.54.png

247.8 KB

OPEN WITH

Язык 1

11:22 ✓



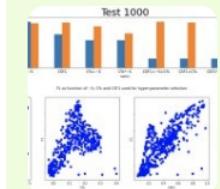
Screen Shot 2022-06-16 at 11.09.45.png

256.8 KB

OPEN WITH

Язык 2

11:23 ✓



Screen Shot 2022-06-16 at 11.09.59.png

276.4 KB

OPEN WITH

Язык 3

11:23 ✓

Poll 

11% 1 - русский

23% 2 - русский

70% 3 - русский

5% 1 - английский

70% 2 - английский

17% 3 - английский

82% 1 - китайский

0% 2 - китайский

5% 3 - китайский

VIEW RESULTS

↪ 3 ⚡ 11:25 ✓

<https://t.me/agibots/2695>

Выводы?

“Интуитивная” (человекоподобная) кластеризация – возможна!

По крайней мере, для временных рядов с мультимодальными распределениями.

Высокоточная сегментация и выявление лексиконов и знаков пунктуации без учителя для языков типа русского и английского – возможна!

С языками типа китайского – надо разбираться.

Посмотрим, как это все будет работать для обучения с подкреплением, прогнозирования временных рядов и построения интерпретируемых полноценных языковых моделей...

<https://github.com/aigents/pygents>

Спасибо! Вопросы?



Антон Колонин
akolonin@aigents.com
<https://github.com/aigents>

Telegram: [akolonin](#)



<https://agirussia.org>