



# **Advanced Parser for Biomedical Texts** Maxim Holmatov<sup>1)</sup>, Anton Karazeev<sup>2)</sup> <sup>1)</sup> Saint-Petersburg State Pediatric Medical University, <sup>2)</sup> Moscow Institute of Physics and Technology Laboratory of Functional Analysis of the Genome

## Introduction

Large amounts of biomedical data available to us today from various sources make it at least impractical and in many cases impossible to analyze by hand even if confined within a specific problem. On the other hand most of these data are stored in a natural language form which makes it hard to process automatically. Fortunately a vast experience gained in the field of natural language processing (NLP) can be utilized to automate this process. We developed an advanced parser for biomedical texts that should simplify both data retrieval and analysis.

We considered the following problems:

- 1. parsing of informative multiword phrases
- 2. parsing and detection of chemical names written in different notations - trivial notation and IUPAC and SMILES-like
- 3. assigning word embeddings for parsed words and phrases
- 4. analyzing complex syntactic dependencies between them

## Methods

To improve parsing quality we decided to learn to extract informative n-grams (e.g. instead of ['amino', 'acid', ...] we want to get ['amino\_acid', ...]) to account for existence of multiword biomedical terms.

To better identify informative n-grams and give a numerical estimate of their validity two main approaches were used.

First one relies on finding the most important edges in word collocation network for analyzed text. Word collocation networks are weighted directed graphs with each vertex corresponding to a word in the text and edge weights equal to the bigram frequency in the document. The most important edges are found by calculating centrality measures of network (degree, closeness, betweenness, etc.) or with the PageRank algorithm [Lahiri et al.]. This process can be applied to analyze documents separately or to generate a custom dictionary of ngrams from a large corpus of texts.

Second approach uses term frequency-inverse document frequency (TF-IDF) statistic. It rewards frequent terms inside a document but punishes words that are frequent in the whole corpus which helps to filter out the words that are just commonly used in a language.



 $D_{\mathrm{variatio}}$ 



Collocation graph based on the abstract of [Harris et al, Stimulation of bone formation in vivo by phosphate supplementation. Calcif Tissue Res. 1976 Nov 24;22(1): 85-98.]. Stop-words were removed. Arrows skipped for convenience even though the graph is directed. Size of the node is proportional to its PageRank score.

## Results

	1	1	1	
PageRank	Gaussian KL(bigram, token)	Gaussian KL(token, bigram)	Variational KL(bigram, mixture)	Variational KL(mixture, bigram)
breast_cancer	ang_iii	citron_kinase	coli_isolates	early_disease
cancer_cells	citron_kinase	biliary_complications	liver_cancer	hpv_dna
gene_expression	biliary_complications	vte_prophylaxis	hpv_dna	liver_cancer
cell_lines	vte_prophylaxis	serum_calcium	model_group	coli_isolates
tumor_cells	new_drugs	dsrna_binding	molecular_target	viral_rna
stem_cells	status_epilepticus	acute_ethanol	cardiac_fibroblasts	reported_cases
prostate_cancer	tuberculosis_isolates	hand_hygiene	early_disease	molecular_target
gastric_cancer	mrsa_strains	status_epilepticus	reported_cases	model_group
cell_cycle	serum_calcium	ang_iii	genetic_studies	meningococcal_disease
patients_treated	acute_ethanol	synthesized_compounds	meningococcal_disease	molecular_data
TF-IDF	Gaussian KL(bigram, token)	Gaussian KL(token, bigram)	Variational KL(bigram, mixture)	Variational KL(mixture, bigram)
TF-IDF gene_expression	Gaussian KL(bigram, token) beta_sheet	Gaussian KL(token, bigram) beneficial_effects	Variational KL(bigram, mixture) related_protein	Variational KL(mixture, bigram) related_protein
TF-IDF gene_expression wild_type	Gaussian KL(bigram, token) beta_sheet disease_ad	Gaussian KL(token, bigram) beneficial_effects results_mean	Variational KL(bigram, mixture) related_protein combination_therapy	Variational KL(mixture, bigram) related_protein viral_rna
TF-IDF gene_expression wild_type present_study	Gaussian KL(bigram, token) beta_sheet disease_ad beneficial_effects	Gaussian KL(token, bigram) beneficial_effects results_mean self_renewal	Variational KL(bigram, mixture) related_protein combination_therapy significant_reduction	Variational KL(mixture, bigram) related_protein viral_rna hiv_positive
TF-IDF gene_expression wild_type present_study cell_lines	Gaussian KL(bigram, token) beta_sheet disease_ad beneficial_effects self_renewal	Gaussian KL(token, bigram) beneficial_effects results_mean self_renewal remains_unclear	Variational KL(bigram, mixture) related_protein combination_therapy significant_reduction viral_rna	Variational KL(mixture, bigram) related_protein viral_rna hiv_positive significant_reduction
TF-IDF gene_expression wild_type present_study cell_lines amino_acid	Gaussian KL(bigram, token) beta_sheet disease_ad beneficial_effects self_renewal old_woman	Gaussian KL(token, bigram) beneficial_effects results_mean self_renewal remains_unclear efficacy_safety	Variational KL(bigram, mixture) related_protein combination_therapy significant_reduction viral_rna study_performed	Variational KL(mixture, bigram) related_protein viral_rna hiv_positive significant_reduction combination_therapy
TF-IDF gene_expression wild_type present_study cell_lines amino_acid results_suggest	Gaussian KL(bigram, token) beta_sheet disease_ad beneficial_effects self_renewal old_woman insulin_sensitivity	Gaussian KL(token, bigram) beneficial_effects results_mean self_renewal remains_unclear efficacy_safety studies_performed	Variational KL(bigram, mixture) related_protein combination_therapy significant_reduction viral_rna study_performed rat_model	Variational KL(mixture, bigram) related_protein viral_rna hiv_positive significant_reduction combination_therapy tissue_specific
TF-IDF gene_expression wild_type present_study cell_lines amino_acid results_suggest breast_cancer	Gaussian KL(bigram, token) beta_sheet disease_ad beneficial_effects self_renewal old_woman insulin_sensitivity et_al	Gaussian KL(token, bigram) beneficial_effects results_mean self_renewal remains_unclear efficacy_safety studies_performed beta_sheet	Variational KL(bigram, mixture) related_protein combination_therapy significant_reduction viral_rna study_performed rat_model tissue_specific	Variational KL(mixture, bigram) related_protein viral_rna hiv_positive significant_reduction combination_therapy tissue_specific study_performed
TF-IDF gene_expression wild_type present_study cell_lines amino_acid results_suggest breast_cancer long_term	Gaussian KL(bigram, token) beta_sheet disease_ad beneficial_effects self_renewal old_woman insulin_sensitivity et_al false_positive	Gaussian KL(token, bigram) beneficial_effects results_mean self_renewal remains_unclear efficacy_safety studies_performed beta_sheet negative_bacteria	Variational KL(bigram, mixture) related_protein combination_therapy significant_reduction viral_rna study_performed rat_model tissue_specific terminal_region	Variational KL(mixture, bigram) related_protein viral_rna hiv_positive significant_reduction combination_therapy tissue_specific study_performed methods_total
TF-IDF gene_expression wild_type present_study cell_lines amino_acid results_suggest breast_cancer long_term mg_kg	Gaussian KL(bigram, token) beta_sheet disease_ad beneficial_effects self_renewal old_woman insulin_sensitivity et_al false_positive therapeutic_targets	Gaussian KL(token, bigram) beneficial_effects results_mean self_renewal remains_unclear efficacy_safety studies_performed beta_sheet negative_bacteria old_woman	Variational KL(bigram, mixture) related_protein combination_therapy significant_reduction viral_rna study_performed rat_model tissue_specific terminal_region hiv_positive	Variational KL(mixture, bigram) related_protein viral_rna hiv_positive significant_reduction combination_therapy tissue_specific study_performed methods_total using_different

#### **Kullback-Leibler Divergence**

In the context of machine learning,  $D_{KL}(P||Q)$  is often called the information gain achieved if P is used instead of Q.

$$D_{ ext{KL}}(P\|Q) = \sum_i P(i) \, \log rac{P(i)}{Q(i)}.$$

$$\sum_{a} (f \| g) = \sum_{a} \pi_{a} \log \frac{\sum_{a'} \pi_{a'} e^{-D(f_{a} \| f_{a'})}}{\sum_{b} \omega_{b} e^{-D(f_{a} \| g_{b})}}$$

KL-divergence method allows us to determine which sets of words are better to replace with an ngram as we can calculate the informativeness of ngram









## References

- Kullback Leibler Divergence Between Gaussian Mixture Models, In IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), 2007. Gaussian Embedding, 2014. seomoz/word2gauss.
- 1. John R. Hershey and Peder A. Olsen, *Approximating the* 2. Luke Vilnis and Andrew McCallum, Word Representations via 3. Moz, Gaussian Word Embeddings, <u>https://github.com/</u>

