

Mining Heterogenous Relationships from Pubmed Abstracts Using Weak Supervision

This manuscript ([permalink](#)) was automatically generated from [greenelab/text_mined_hetnet_manuscript@d3e96c7](#) on July 8, 2019.

Authors

- **David N. Nicholson**

 [0000-0003-0002-5761](#) ·  [danich1](#)

Department of Systems Pharmacology and Translational Therapeutics, University of Pennsylvania · Funded by GBMF4552

- **Daniel S. Himmelstein**

 [0000-0002-3012-7446](#) ·  [dhimmel](#) ·  [dhimmel](#)

Department of Systems Pharmacology and Translational Therapeutics, University of Pennsylvania · Funded by GBMF4552

- **Casey S. Greene**

 [0000-0001-8713-9213](#) ·  [cgreene](#) ·  [GreeneScientist](#)

Department of Systems Pharmacology and Translational Therapeutics, University of Pennsylvania · Funded by GBMF4552 and R01 HG010067

Abstract

This is a **rough draft** of a manuscript on label function reuse for text mining heterogenous relationship from Pubmed Abstracts.

#Introduction Set introduction for paper here Talk about problem, goal, and significance of paper

Recent Work

Talk about what has been done in the field in regards to text mining and knowledge base integration

Materials and Methods

Hetionet

Hetionet [1] is a large heterogenous network that contains pharmacological and biological information. This network depicts information in the form of nodes and edges of different types: nodes that represent biological and pharmacological entities and edges which represent relationships between entities. Hetionet v1.0 contains 47,031 nodes with 11 different data types and 2,250,197 edges that represent 24 different relationship types (Figure 1). Edges in Hetionet were obtained from open databases, such as the GWAS Catalog [2] and DrugBank [2]. For this project, we analyzed performance over a subset of the Hetionet relationship types: disease associates with a gene (DaG), compound binds to a gene (CbG), gene interacts with gene (GiG) and compound treating a disease (CtD).

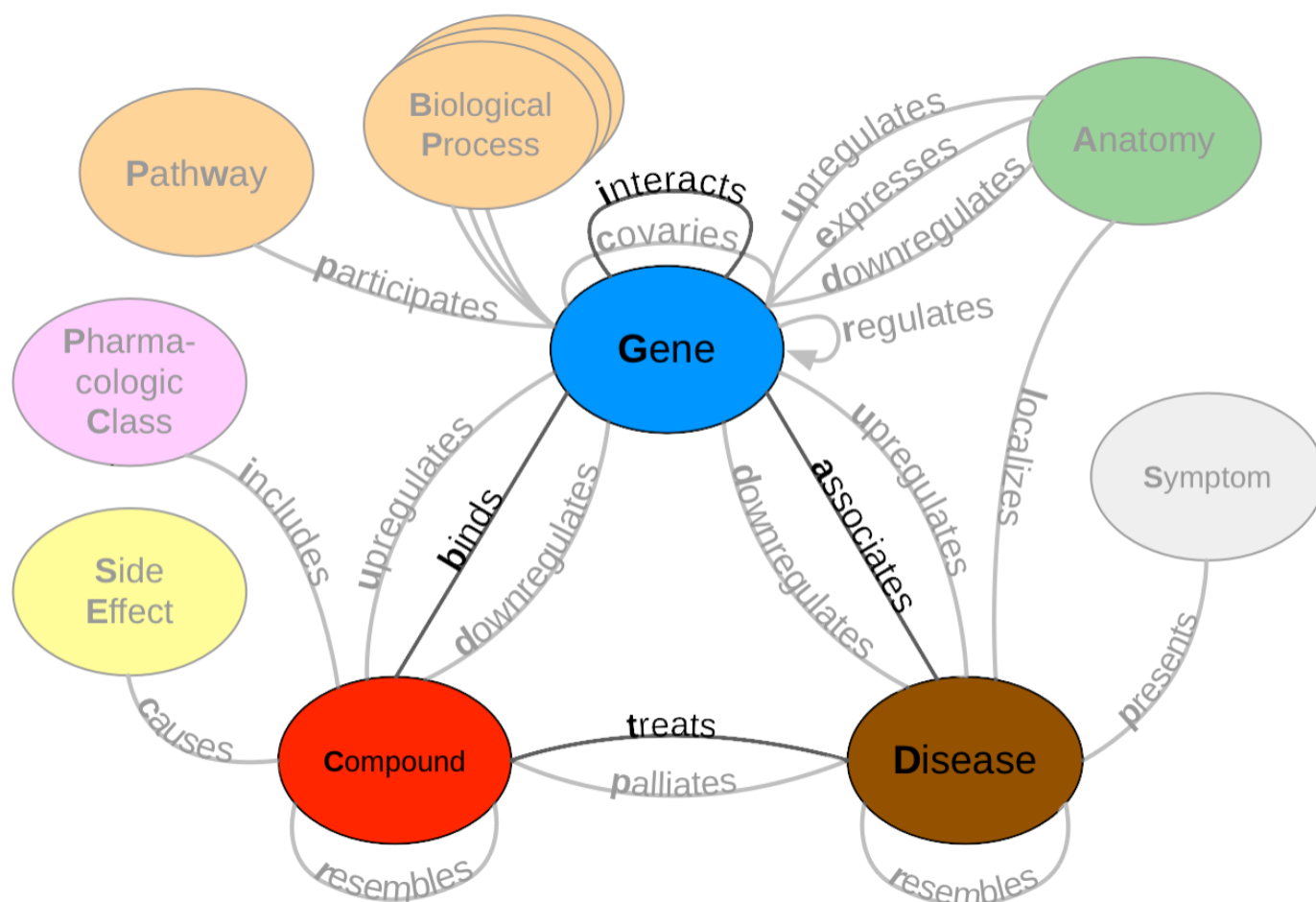


Figure 1: A metagraph (schema) of Hetionet where pharmacological, biological and disease entities are represented as nodes and the relationships between them are represented as edges. This project only focuses on the information shown in bold; however, we can extend this work to incorporate the faded out information as well.

Dataset

We used PubTator [3] as input to our analysis. PubTator provides MEDLINE abstracts that have been annotated with well-established entity recognition tools including DNorm [4] for disease mentions, GeneTUKit [5] for gene mentions, Gnorm [6] for gene normalizations and a dictionary based look system for compound mentions [7]. We downloaded PubTator on June 30, 2017, at which point it contained 10,775,748 abstracts. Then we filtered out mention tags that were not contained in hetionet. We used the Stanford CoreNLP parser [8] to tag parts of speech and generate dependency trees. We extracted sentences with two or more mentions, termed candidate sentences. Each candidate sentence was stratified by co-mention pair to produce a training set, tuning set and a testing set (shown in Table 1). Each unique co-mention pair is sorted into four categories: (1) in hetionet and has sentences, (2) in hetionet and doesn't have sentences, (3) not in hetionet and does have sentences and (4) not in hetionet and doesn't have sentences. Within these four categories each pair receives their own individual partition rank (continuous number between 0 and 1). Any rank lower than 0.7 is sorted into training set, while any rank greater than 0.7 and lower than 0.9 is assigned to tuning set. The rest of the pairs with a rank greater than or equal to 0.9 is assigned to the test set. Sentences that contain more than one co-mention pair are treated as multiple individual candidates. We hand labeled five hundred to a thousand candidate sentences of each relationship to obtain a ground truth set (Table 1, [dataset](#)).

Table 1: Statistics of Candidate Sentences. We sorted each candidate sentence into a training, tuning and testing set. Numbers in parentheses show the number of positives and negatives that resulted from the hand-labeling process.

| Relationship | Train | Tune | Test |
|-------------------------|--------|--------------------|-------------------|
| Disease Associates Gene | 2.35 M | 31K (397+, 603-) | 313K (351+, 649-) |
| Compound Binds Gene | 1.7M | 468K (37+, 463-) | 227k (31+, 469-) |
| Compound Treats Disease | 1.013M | 96K (96+, 404-) | 32K (112+, 388-) |
| Gene Interacts Gene | 12.6M | 1.056M (60+, 440-) | 257K (76+, 424-) |

Label Functions

describe what a label function is and how many we created for each relation

Training Models

Generative Model

talk about generative model and how it works ### Word Embeddings mention facebook's fastText model and how we used it to train word vectors ### Discriminator Model talk about the discriminator model and how it works ### Discriminator Model Calibration talk about calibrating deep learning models with temperature smoothing

Experimental Design

talk about sampling experiment

Results

Random Sampling of Generative Model

place the grid aurocs here for generative model

Discriminator Model Builds Off Generative Model

place the grid of aurocs here for discriminator model

Random Noise Generative Model

place the results of random label function experiment

Reconstructing Hetionet

place figure of number of new edges that can be added to hetionet as well as edges we can reconstruct using this method

Discussion

Here mention why performnace increases in the beginning for the generative model then decreases

Discuss discriminator model performance given generative model

Mention Take home messages

1. have a centralized set of negative label functions and focus more on contstructing positive label functions

Conclusion and Future Direction

Recap the original problem - takes a long time to create useful label function

Proposed solution - reuse label functions

Mention incorporating more relationships Mention creating a centralized multitask text extractor using this method.

References

1. **Systematic integration of biomedical knowledge prioritizes drugs for repurposing**

Daniel Scott Himmelstein, Antoine Lizee, Christine Hessler, Leo Brueggeman, Sabrina L Chen, Dexter Hadley, Ari Green, Pouya Khankhanian, Sergio E Baranzini

eLife (2017-09-22) <https://doi.org/cdfk>

DOI: [10.7554/elife.26726](https://doi.org/10.7554/elife.26726) · PMID: [28936969](https://pubmed.ncbi.nlm.nih.gov/28936969/) · PMCID: [PMC5640425](https://pubmed.ncbi.nlm.nih.gov/PMC5640425/)

2. **The new NHGRI-EBI Catalog of published genome-wide association studies (GWAS Catalog)**

Jacqueline MacArthur, Emily Bowler, Maria Cerezo, Laurent Gil, Peggy Hall, Emma Hastings, Heather Junkins, Aoife McMahon, Annalisa Milano, Joannella Morales, ... Helen Parkinson

Nucleic Acids Research (2016-11-29) <https://doi.org/f9v7cp>

DOI: [10.1093/nar/gkw1133](https://doi.org/10.1093/nar/gkw1133) · PMID: [27899670](https://pubmed.ncbi.nlm.nih.gov/27899670/) · PMCID: [PMC5210590](https://pubmed.ncbi.nlm.nih.gov/PMC5210590/)

3. **PubTator: a web-based text mining tool for assisting biocuration**

Chih-Hsuan Wei, Hung-Yu Kao, Zhiyong Lu

Nucleic Acids Research (2013-05-22) <https://doi.org/f475th>

DOI: [10.1093/nar/gkt441](https://doi.org/10.1093/nar/gkt441) · PMID: [23703206](https://pubmed.ncbi.nlm.nih.gov/23703206/) · PMCID: [PMC3692066](https://pubmed.ncbi.nlm.nih.gov/PMC3692066/)

4. **DNorm: disease name normalization with pairwise learning to rank**

R. Leaman, R. Islamaj Dogan, Z. Lu

Bioinformatics (2013-08-21) <https://doi.org/f5gj9n>

DOI: [10.1093/bioinformatics/btt474](https://doi.org/10.1093/bioinformatics/btt474) · PMID: [23969135](https://pubmed.ncbi.nlm.nih.gov/23969135/) · PMCID: [PMC3810844](https://pubmed.ncbi.nlm.nih.gov/PMC3810844/)

5. **GeneTUKit: a software for document-level gene normalization**

M. Huang, J. Liu, X. Zhu

Bioinformatics (2011-02-08) <https://doi.org/dng2cb>

DOI: [10.1093/bioinformatics/btr042](https://doi.org/10.1093/bioinformatics/btr042) · PMID: [21303863](https://pubmed.ncbi.nlm.nih.gov/21303863/) · PMCID: [PMC3065680](https://pubmed.ncbi.nlm.nih.gov/PMC3065680/)

6. **Cross-species gene normalization by species inference**

Chih-Hsuan Wei, Hung-Yu Kao

BMC Bioinformatics (2011-10-03) <https://doi.org/dnmvds>

DOI: [10.1186/1471-2105-12-s8-s5](https://doi.org/10.1186/1471-2105-12-s8-s5) · PMID: [22151999](https://pubmed.ncbi.nlm.nih.gov/22151999/) · PMCID: [PMC3269940](https://pubmed.ncbi.nlm.nih.gov/PMC3269940/)

7. **Collaborative biocuration–text-mining development task for document prioritization for curation**

T. C. Wiegers, A. P. Davis, C. J. Mattingly

Database (2012-11-22) <https://doi.org/gbb3zw>

DOI: [10.1093/database/bas037](https://doi.org/10.1093/database/bas037) · PMID: [23180769](https://pubmed.ncbi.nlm.nih.gov/23180769/) · PMCID: [PMC3504477](https://pubmed.ncbi.nlm.nih.gov/PMC3504477/)

8. **The Stanford CoreNLP Natural Language Processing Toolkit**

Christopher Manning, Mihai Surdeanu, John Bauer, Jenny Finkel, Steven Bethard, David McClosky
Proceedings of 52nd Annual Meeting of the Association for Computational Linguistics: System Demonstrations (2014)

<https://doi.org/gf3xhp>

DOI: [10.3115/v1/p14-5010](https://doi.org/10.3115/v1/p14-5010)