



CONSTRUCTING CO-OCCURRENCE NETWORK EMBEDDINGS TO ASSIST ASSOCIATION EXTRACTION FOR COVID-19 AND OTHER CORONAVIRUS INFECTIOUS DISEASES

Feichen Shen, Ph.D.

Assistant Professor of Biomedical Informatics

Division of Digital Health Sciences

David Oniani

R&D Intern

Kern Center for the Science of Health Care Delivery

BACKGROUND

-COVID-19

- An infectious disease caused by severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2)
- As of today, more than 8 million cases across 188 countries and territories, resulting in more than 435,000 deaths.

BACKGROUND

-CORD-19

- The White House and leading research groups have prepared the COVID-19 Open Research Dataset (CORD-19).
- Over 138,000 scholarly articles, including over 69,000 with full text, about COVID-19, SARS-CoV-2, and related coronaviruses.

BACKGROUND

-CORD-19-on-FHIR

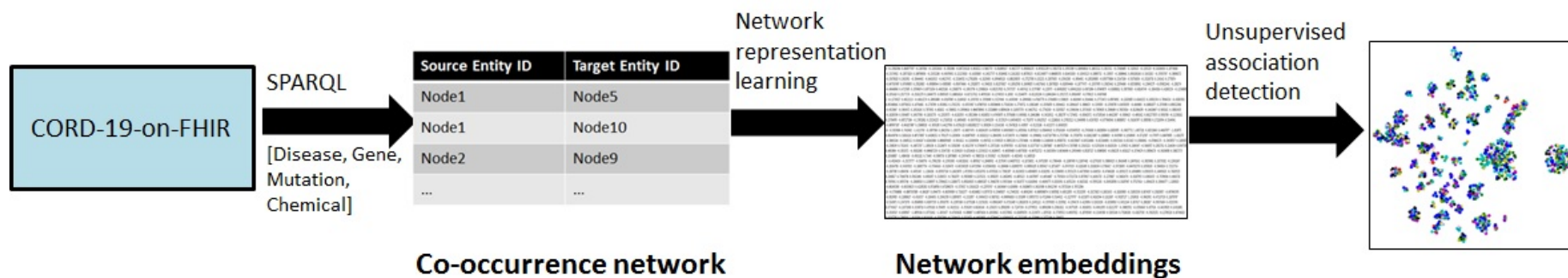
- A linked data version of CORD-19.
- Incorporate annotations from different sources, such as PubTator and LitCovid.
- Easy to retrieve co-occurrences among entities by using SPARQL.

OBJECTIVE

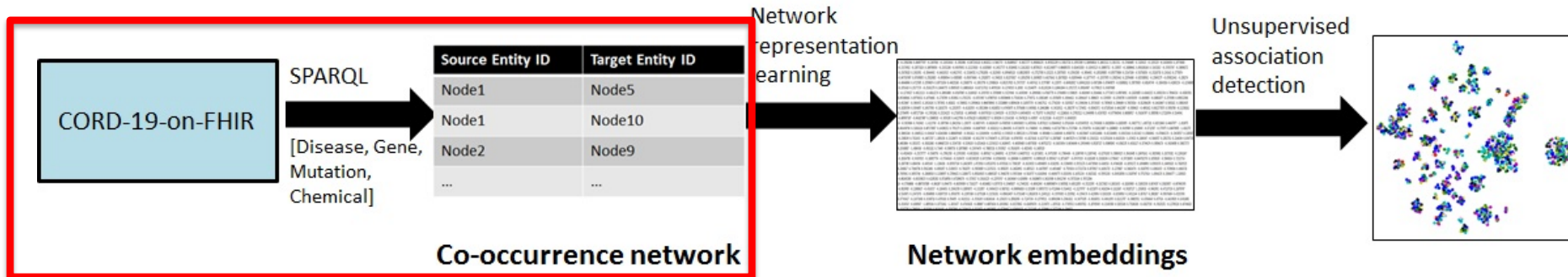
- Co-occurrence network
- Network embeddings
- Accelerate knowledge discovery for COVID-19

METHODS

-OVERVIEW



METHODS



METHODS

-CO-OCCURRENCE NETWORK

```
pmc:annotations [  
  pmc:id "1" ;  
    pmc:infons [ pmc:identifier "MESH:D003371" ; pmc:type "Disease" ] ;  
    pmc:locations [ pmc:length "5"^^xsd:int ; pmc:offset "20312"^^xsd:int ] ;  
    pmc:text "cough" ],  
  pmc:id "2" ;  
    pmc:infons [ pmc:identifier "MESH:C000657245" ; pmc:type "Disease" ] ;  
    pmc:locations [ pmc:length "19"^^xsd:int ; pmc:offset "14766"^^xsd:int ] ;  
    pmc:text "2019-nCoV infection" ],],
```

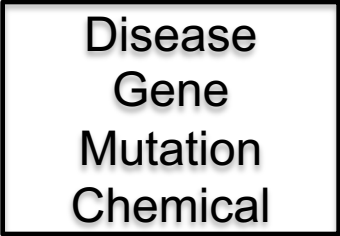
Disease
Gene
Mutation
Chemical

```
pmc:annotations [  
  pmc:id "5" ;  
    pmc:infons [ pmc:identifier "59272" ; pmc:ncbi_homologene "41448" ; pmc:type "Gene" ] ;  
    pmc:locations [ pmc:length "31"^^xsd:int ; pmc:offset "1986"^^xsd:int ] ;  
    pmc:text "angiotensin-converting enzyme 2" ],  
  pmc:id "7" ;  
    pmc:infons [ pmc:identifier "MESH:C000657245" ; pmc:type "Disease" ] ;  
    pmc:locations [ pmc:length "19"^^xsd:int ; pmc:offset "14766"^^xsd:int ] ;  
    pmc:text "2019-nCoV infection" ],],
```


METHODS

-CO-OCCURRENCE NETWORK

```
PREFIX rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#>
PREFIX fhir: <http://hl7.org/fhir/>
PREFIX pmc: <https://www.ncbi.nlm.nih.gov/pmc/articles#>
SELECT distinct ?pmc_id0 ?text0 ?pmc_id1 ?text1 (count(?text1) as ?count) WHERE {
  ?pmc pmc:annotations
  [ pmc:id ?id0 ; pmc:text ?text0 ; pmc:infos
    [ pmc:type ?type0 ; pmc:identifier ?pmc_id0 ] ] .
  FILTER ((?type0 = 'Biomedical_Entity')) .
  {select * where{
    ?pmc pmc:annotations
    [ pmc:id ?id1 ; pmc:text ?text1 ; pmc:infos
      [ pmc:type ?type1 ; pmc:identifier ?pmc_id1 ] ] .
    FILTER ((?type1='Disease') && (contains (lcase(str(?text1)), "coronavirus") || contains (lcase(str(?text1)), "sars") || contains (lcase(str(?text1)), "covid-19") || contains (lcase(str(?text1)), "pneumonia") || contains (lcase(str(?text1)), "fever") || contains (lcase(str(?text1)), "fibrosis") || contains (lcase(str(?text1)), "diarrhea") || contains (lcase(str(?text1)), "bronchitis") || contains (lcase(str(?text1)), "ebola") || contains (lcase(str(?text1)), "influenza") || contains (lcase(str(?text1)), "zika")))).
  }
}
}Group by ?pmc_id0 ?text0 ?pmc_id1 ?text1 Order by DESC(?count)
```



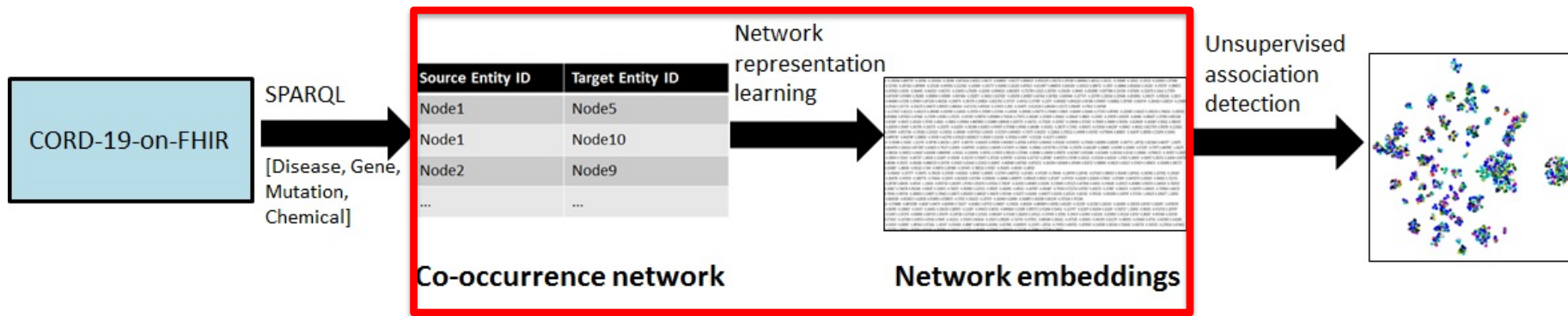
Disease
Gene
Mutation
Chemical

METHODS

-CO-OCCURRENCE NETWORK

| | text0 | text1 |
|----|---------------------------------|-----------------------|
| 1 | Hu B | pneumonia |
| 2 | albumin | fever |
| 3 | alanine aminotransferase | fever |
| 4 | aspartate aminotransferase | fever |
| 5 | SARS | pneumonia |
| 6 | SARS | 2019-nCoV pneumonia |
| 7 | Nam | 2019-nCoV pneumonia |
| 8 | Sri | 2019-nCoV pneumonia |
| 9 | bat | coronavirus infection |
| 10 | bat | SARS |
| 11 | bat | pneumonia |
| 12 | ACE2 | SARS |
| 13 | angiotensin-converting enzyme 2 | SARS |

METHODS



METHODS

-NETWORK EMBEDDINGS

- Similar to word embeddings
- Node \rightarrow Word, Neighborhood \rightarrow Context
- However, sliding window for text is not suitable for non-linear graph
- One solution is to use random walk to select “context” in graph

METHODS

-NETWORK EMBEDDINGS

Node2vec

- Scalable feature learning for networks.
- Grover A, Leskovec J.
- In Proceedings of the 22nd ACM SIGKDD international conference on Knowledge discovery and data mining 2016 Aug 13 (pp. 855-864). ACM.

NODE2VEC

2 steps

- Sampling strategy

Apply random walk on graph to prepare input data

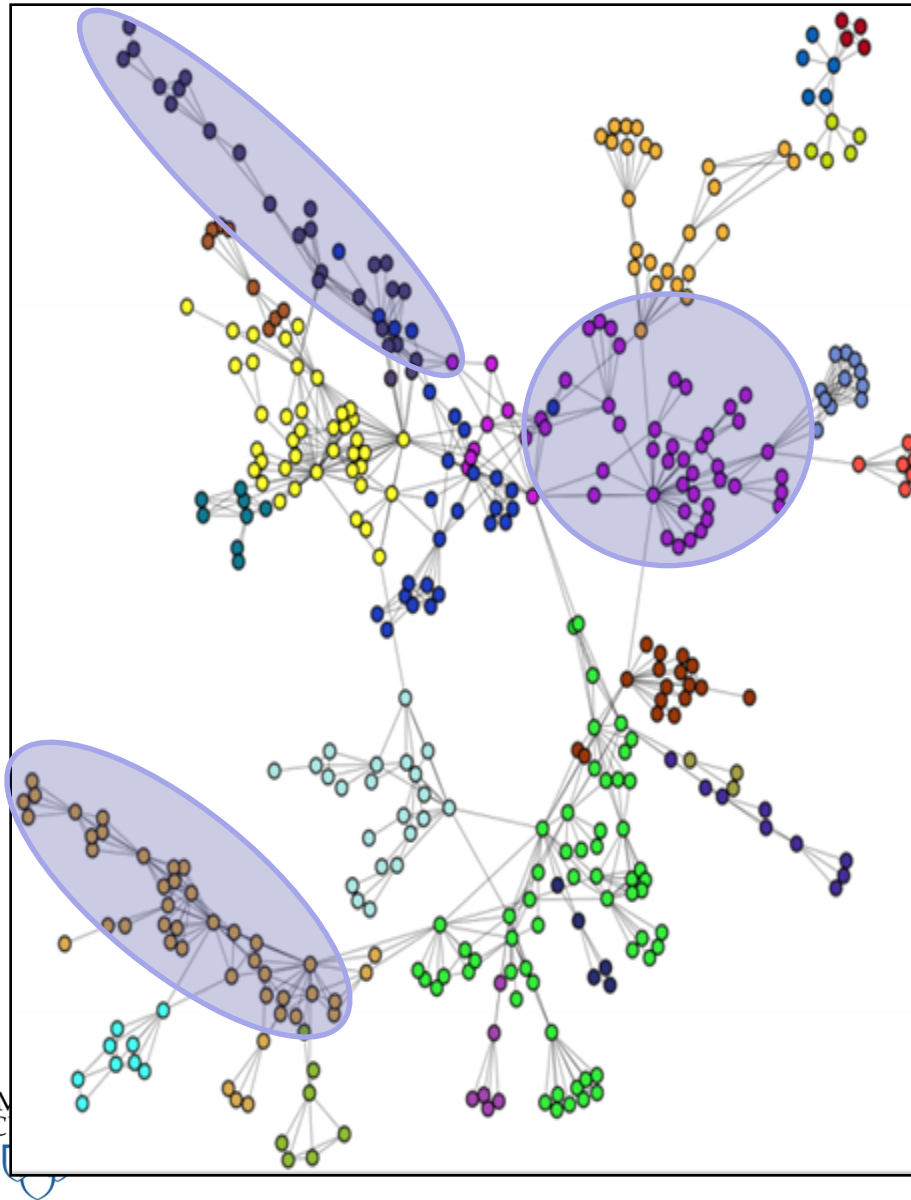
- Node embeddings

Apply word2vec on prepared input data to generate embeddings for node

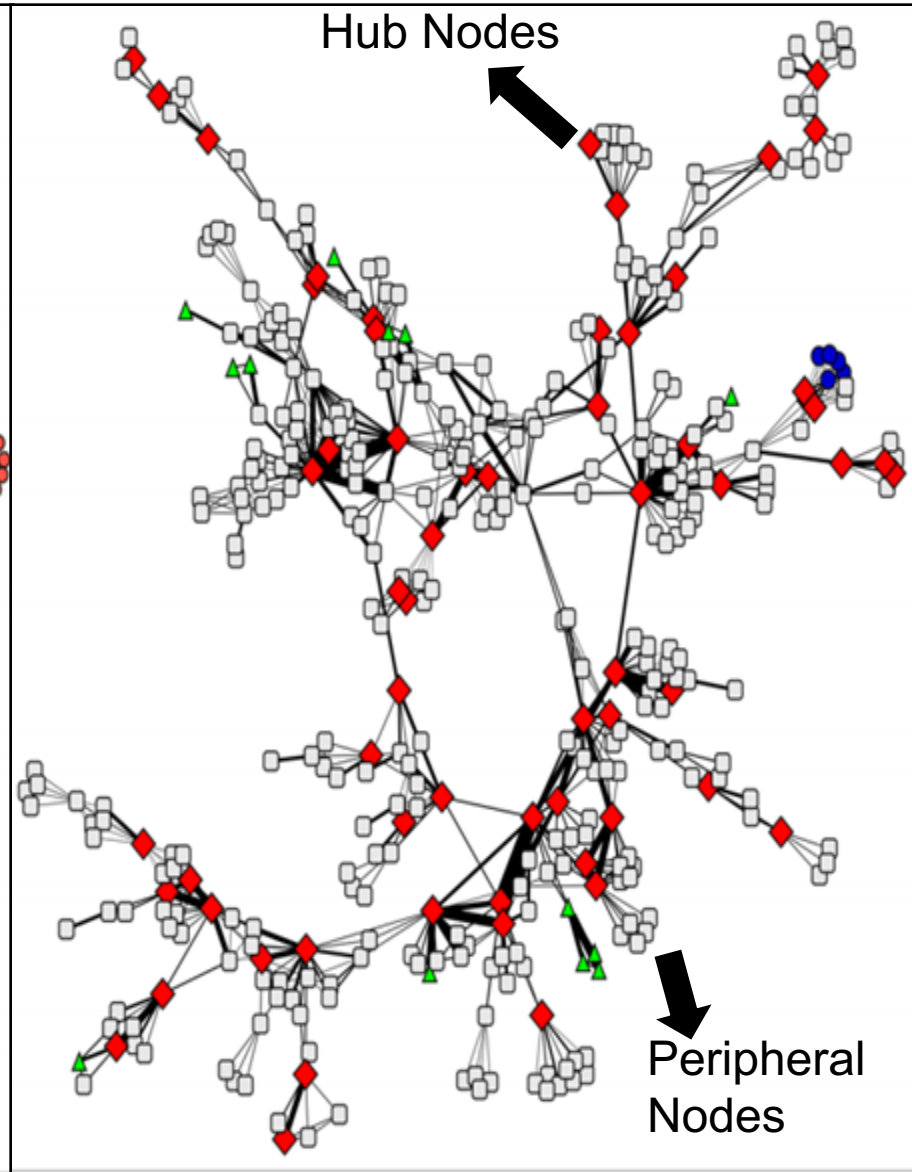
1. SAMPLING STRATEGY

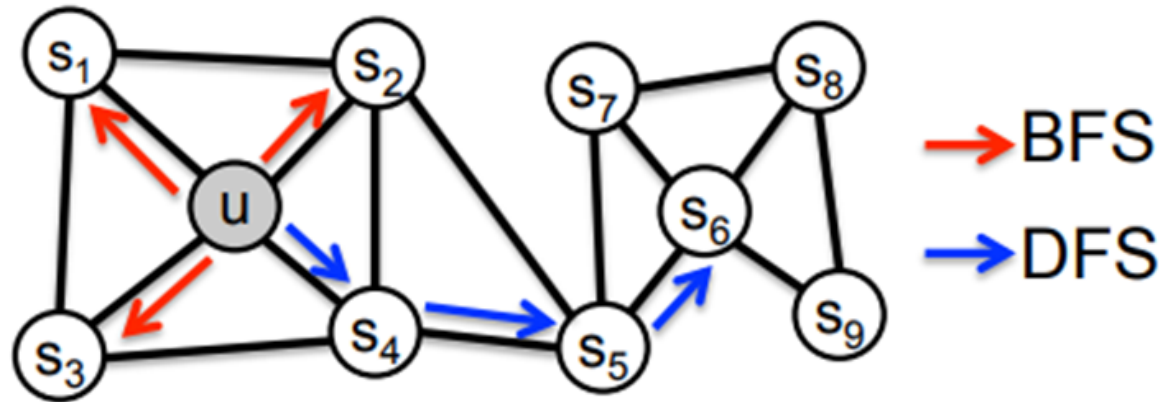
- 1. **Homophily equivalence:** embed nodes from the same network community closely together
- 2. **Structural equivalence:** nodes share similar roles have similar embeddings

Homophily equivalence

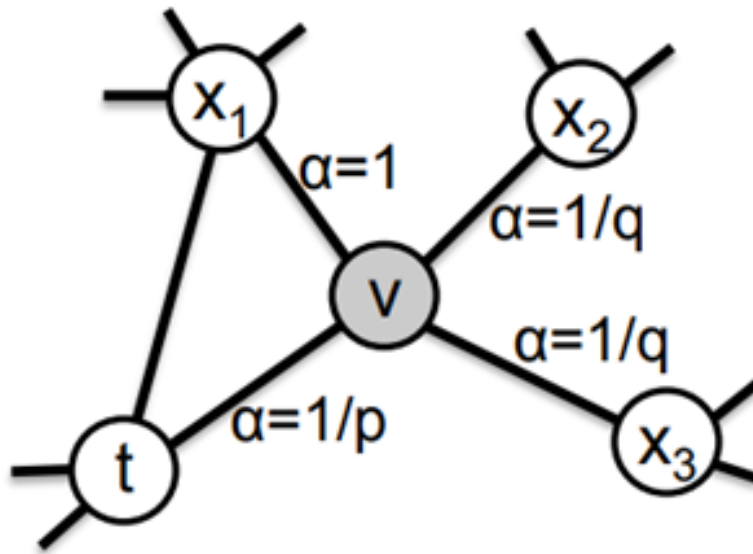


Structural equivalence





- **Breadth-first search (BFS):** structural analysis, microscopic view of neighborhood of every node
- **Depth-first search (DFS):** homophily analysis, macroscopic view of the neighborhood among different communities
- Real-world network has a lot of such mixture of BFS and DFS



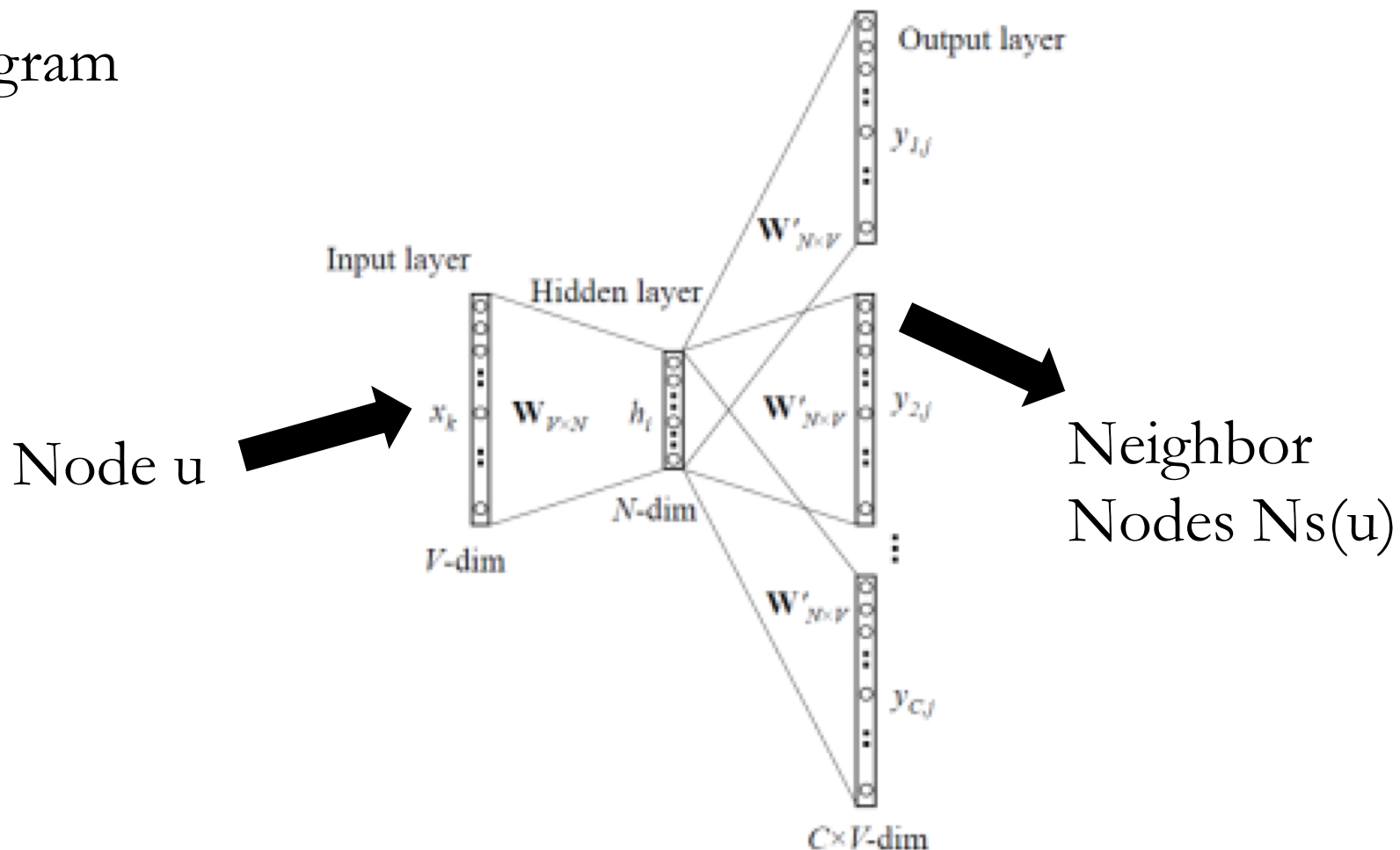
$$\alpha_{pq}(t, x) = \begin{cases} \frac{1}{p} & \text{if } d_{tx} = 0 \\ 1 & \text{if } d_{tx} = 1 \\ \frac{1}{q} & \text{if } d_{tx} = 2 \end{cases} \quad \text{Bias Term}$$

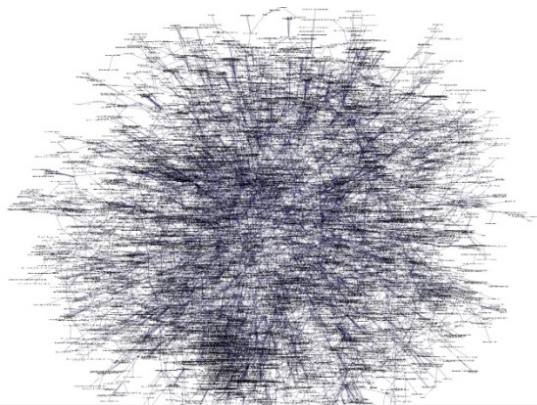
$$\pi_{vx} = \alpha_{pq}(t, x) \cdot w_{vx} \quad \text{Transition Probability}$$

- p and q are controller to balance between BFS and DFS
- p controls the return step, and q controls the step of walk to outside world

2. WORD2VEC

Skip-gram



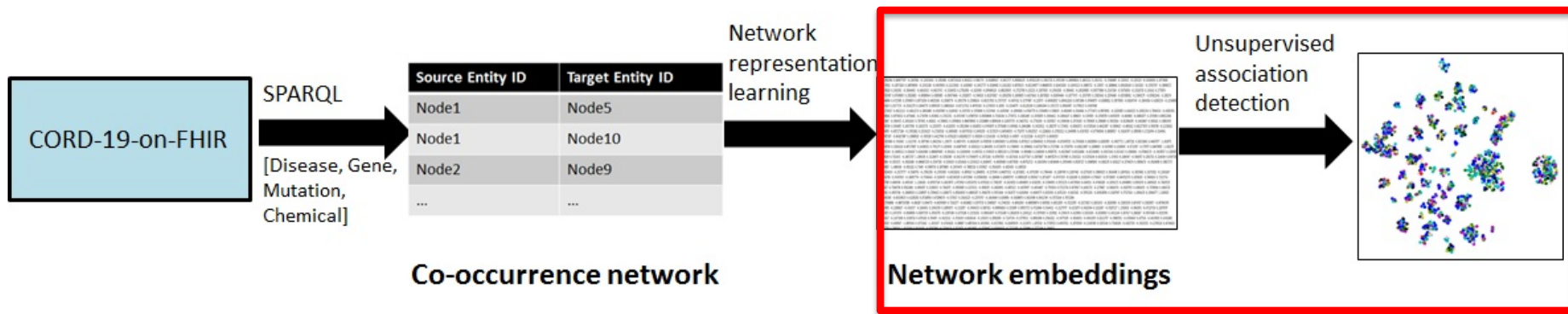


0 -0.240338 0.850522 0.153959 -0.002026 -0.064240 0.235587 -0.108960 -0.338055 -0.119517 0.175563 0.391702 -0.056000 -0.124363 -0.063387 0.
3890 -0.272510 0.076228 0.161161 0.821780 -0.870825 0.270818 -0.195049 -0.394763 -0.122105 0.208492 0.445066 -0.051002 -0.165230 -0.087852
1389 -0.260807 0.857100 0.186200 0.006649 -0.870756 0.245432 -0.184220 -0.374630 -0.125654 0.183740 0.413510 -0.069082 -0.153400 -0.066026
1833 -0.259970 0.851797 0.163549 -0.002428 -0.877161 0.241342 -0.191172 -0.370209 -0.126249 0.182140 0.415996 -0.066236 -0.157282 -0.070804
007 -0.259210 0.855444 0.166569 0.014121 -0.870808 0.240032 -0.191375 -0.365271 -0.127956 0.181081 0.414304 -0.065014 -0.157044 -0.058468
2001 -0.246176 0.847520 0.161245 0.009027 -0.868087 0.238785 -0.185408 -0.356083 -0.124393 0.178436 0.380874 -0.061958 -0.149504 -0.05083
838 -0.224584 0.849370 0.152085 0.004368 -0.863852 0.225471 -0.169435 -0.322289 -0.111722 0.172882 0.371441 -0.056873 -0.120084 -0.056818 0
1872 -0.223858 0.820879 0.124206 0.009716 -0.851364 0.234236 -0.187603 -0.329637 -0.122821 0.177479 0.365900 -0.040611 -0.120084 -0.062653
327 -0.228755 0.854087 0.141838 0.012875 -0.858397 0.180548 -0.155800 -0.317608 -0.105438 0.162329 0.344005 -0.057205 -0.122585 -0.063517 0
92 -0.225669 0.850669 0.132185 0.000005 -0.852382 0.215648 -0.157785 -0.317897 -0.104145 0.161382 0.355888 -0.054407 -0.118013 -0.052913 0
490 -0.215478 0.842010 0.120533 0.004403 -0.850679 0.205168 -0.150854 -0.306908 -0.100467 0.159951 0.343334 -0.057942 -0.112162 -0.047310 0
1759 -0.215291 0.850100 0.125453 0.005003 -0.863868 0.205241 -0.160292 -0.308093 -0.111927 0.162504 0.351751 -0.065556 -0.115052 -0.050709
2158 -0.224805 0.841183 0.145884 0.005946 -0.857137 0.219395 -0.160601 -0.338106 -0.102587 0.167780 0.370981 -0.050820 -0.129200 -0.059591
1401 -0.222384 0.853975 0.122621 0.014570 -0.857000 0.212194 -0.163410 -0.312887 -0.101367 0.159752 0.354470 -0.061913 -0.115980 -0.056708
1826 -0.234396 0.845661 0.148524 0.007663 -0.852708 0.224044 -0.171780 -0.334390 -0.112482 0.178457 0.378480 -0.064150 -0.127882 -0.050900
3870 -0.217018 0.850422 0.134899 0.012042 -0.856442 0.205347 -0.155844 -0.318195 -0.110612 0.159405 0.345328 -0.054990 -0.113820 -0.047748
1115 -0.223408 0.854321 0.140383 0.002554 -0.863774 0.200668 -0.162022 -0.321131 -0.111940 0.155479 0.353809 -0.051009 -0.119995 -0.055056
190 -0.237900 0.855337 0.145518 0.005572 -0.854828 0.230677 -0.169945 -0.331308 -0.111711 0.161080 0.370970 -0.063387 -0.125004 -0.052272 0
677 -0.214353 0.852189 0.134712 0.018339 -0.852341 0.208794 -0.155401 -0.303248 -0.100087 0.160090 0.344840 -0.055018 -0.114605 -0.057170 0
907 -0.197369 0.850839 0.125470 0.005000 -0.840182 0.184252 -0.136348 -0.277059 -0.097910 0.140903 0.304756 -0.043073 -0.105947 -0.043055 0
327 -0.215443 0.851710 0.145057 0.010977 -0.859160 0.210862 -0.160730 -0.315210 -0.103005 0.151150 0.346830 -0.046008 -0.122740 -0.051400 0
250 -0.212490 0.849053 0.139140 0.001940 -0.861963 0.199641 -0.150844 -0.313782 -0.103023 0.149482 0.346885 -0.049574 -0.123820 -0.050471 0
008 -0.214639 0.840417 0.130425 0.005720 -0.858724 0.203819 -0.150882 -0.312844 -0.104439 0.156508 0.340755 -0.062049 -0.115408 -0.045725 0
4307 -0.213807 0.845414 0.125482 0.010044 -0.859401 0.206374 -0.154959 -0.313508 -0.104100 0.154504 0.345739 -0.061927 -0.113131 -0.052572
1008 -0.207112 0.850631 0.113010 0.002981 -0.850075 0.156813 -0.140940 -0.285359 -0.104100 0.141817 0.330695 -0.046988 -0.112400 -0.040817
996 -0.215350 0.841208 0.125453 0.012029 -0.850889 0.203927 -0.155883 -0.305718 -0.100120 0.155509 0.339007 -0.054132 -0.110824 -0.040225 0
1224 -0.217405 0.840022 0.140100 0.011862 -0.856477 0.200958 -0.160410 -0.310002 -0.104738 0.159651 0.355519 -0.061271 -0.115044 -0.050507
2002 -0.216393 0.852970 0.141161 0.004912 -0.856215 0.200668 -0.162001 -0.318708 -0.102070 0.147595 0.344340 -0.050995 -0.109240 -0.057845
332 -0.209703 0.830971 0.131613 0.005140 -0.840731 0.200631 -0.162045 -0.304422 -0.100730 0.150964 0.345540 -0.050013 -0.110919 -0.040142 0
2901 -0.205702 0.845253 0.129652 0.010000 -0.840508 0.190658 -0.152954 -0.287070 -0.090575 0.152010 0.324420 -0.040287 -0.110820 -0.040171
004 -0.203708 0.839433 0.129529 0.012260 -0.851414 0.180383 -0.144608 -0.287998 -0.096271 0.140931 0.320889 -0.040625 -0.108652 -0.047109 0
1828 -0.205504 0.851870 0.130605 0.000029 -0.850687 0.190051 -0.151474 -0.300000 -0.100494 0.142002 0.319341 -0.052939 -0.100017 -0.047203
1301 -0.199030 0.830256 0.125451 0.010022 -0.850147 0.200504 -0.157500 -0.295043 -0.100004 0.156233 0.318045 -0.050007 -0.100020 -0.050013
1705 -0.207183 0.841820 0.132277 0.009582 -0.856525 0.196352 -0.156907 -0.291991 -0.095514 0.152518 0.336234 -0.051432 -0.111460 -0.053085
1625 -0.208008 0.841435 0.137097 0.005170 -0.850694 0.199951 -0.160806 -0.297397 -0.090223 0.157134 0.337223 -0.051472 -0.116200 -0.050408

Graph space

Embeddings space

METHODS



METHODS

-NODE CLUSTERING

- t-distributed stochastic neighbor embedding (t-SNE)
- Density-based spatial clustering of applications with noise (DBSCAN)

RESULTS

- 3,626 diseases
- 5,741 genes
- 524 mutations
- 6,878 chemicals
- 16,769 nodes and 49,696 edges

RESULTS

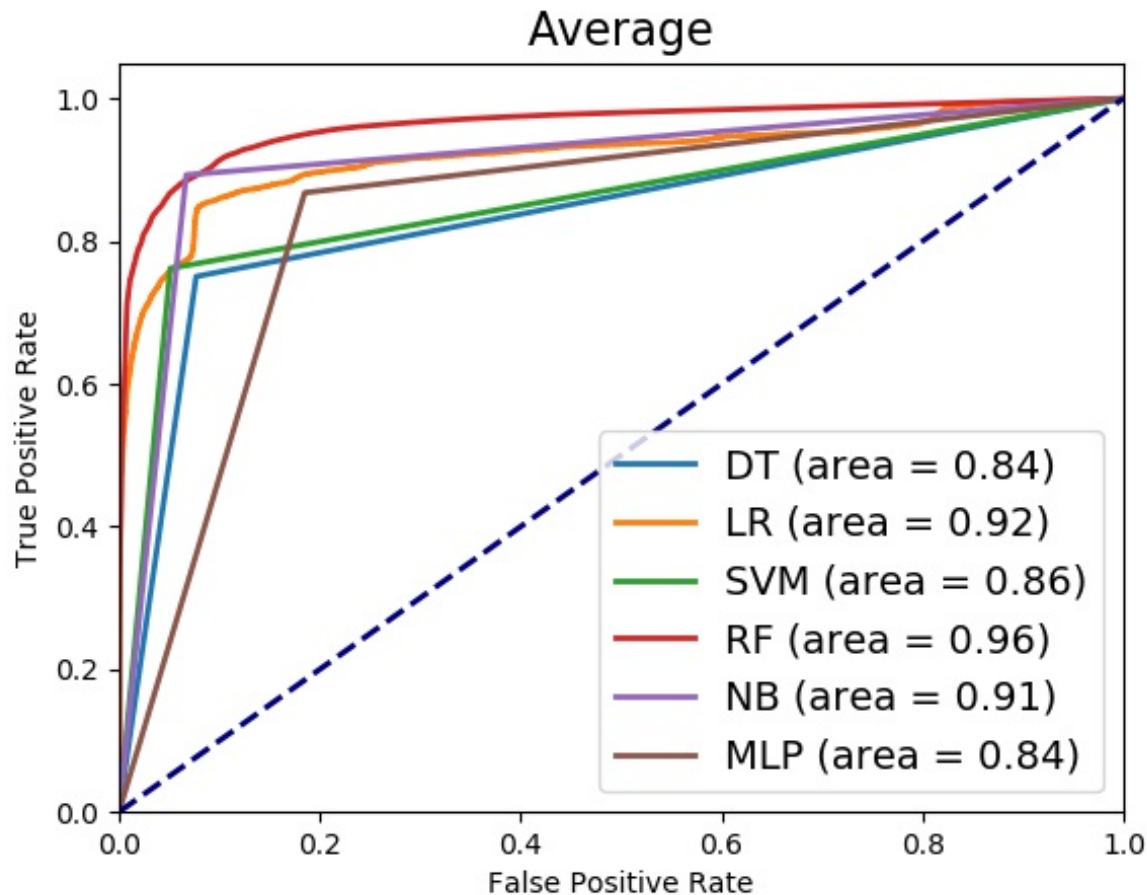
Link prediction

- Generate edge embedding based on node embedding
- Given nodes u and v , $f(u)$ and $f(v)$ denotes their feature representations

| Operator | Definition |
|-------------|-------------------------|
| Average | $\frac{f(u) + f(v)}{2}$ |
| Hadamard | $f(u) * f(v)$ |
| Weighted-L1 | $ f(u) - f(v) $ |
| Weighted-L2 | $ f(u) - f(v) ^2$ |

RESULTS

Decision Tree (DT), Logistic Regression (LR), Support Vector Machine (SVM), Random Forest (RF), Naive Bayes (NB), and Multi-layer Perceptron (MLP)



RESULTS

| Algorithms | Precision | Recall | F1 score |
|------------|-------------|-------------|-------------|
| DT | 0.85 | 0.84 | 0.84 |
| LR | 0.87 | 0.85 | 0.85 |
| SVM | 0.87 | 0.86 | 0.85 |
| RF | 0.91 | 0.91 | 0.90 |
| NB | 0.91 | 0.91 | 0.91 |
| MLP | 0.84 | 0.84 | 0.84 |

RESULTS

COVID-19 (cluster #6)

Top 10 most relevant entities

VP35 (Gene)

HD11 (Gene)

Coronavirus infection process (Disease)

Fibroblast growth factor (FGF)-2 (Gene)

Acute respiratory infection illness (Disease)

PIGS (Gene)

TGF alpha (Gene)

SFPQ (Gene)

Tumour necrosis factor (TNF) (Gene)

Praziquantel (Chemical)

RESULTS

Pulmonary coronavirus infection (Cluster #1)

Top 10 most relevant entities

PTP (Gene)

SARS-CoV-infected human airway epithelia cell cultures (Disease)

"5'-tgg gat tca aca" (Chemical)

Tracheanasal respiratory epithelial cells nd llamas (lama glama) in (Disease)

Suppressor of cytokine signaling 3 (Gene)

KAT (Gene)

CD32 (Gene)

Maternal SARS infection (Disease)

Respiratory syndrome coronavirus (MERS-CoV) infections (Disease)

S27 (Gene)

RESULTS

Sars-cov infection damages lung (Cluster #2)

Top 10 most relevant entities

IL-1-alpha (Gene)

Sucralfate prn (Chemical)

Acute respiratory syndrome-cov infection (Disease)

IL-5- and IL-13-producing ilc-iis (Gene)

HAP1 (Gene)

FSK (Chemical)

Low fever (Disease)

HIV and Ebola virus infection (Disease)

YKL-40 (Gene)

ETF (Gene)

RESULTS

Coronavirus upper respiratory infection (Cluster #23)

Top 10 most relevant entities

Viruses actinobacillus pleuropneumoniae (Disease)

Plasmin (Gene)

JAM-1 (Gene)

TNF receptor-associated factor 6 (Gene)

GPC3 (Gene)

Renin (Gene)

ZO-1 (Gene)

Cathepsin G (Gene)

rs5743313 (Mutation)

Alpha1 antitrypsin (Gene)

RESULTS

Coronavirus-infected pneumonia (Cluster #10)

Top 10 most relevant entities

Respiratory syncytial viral infection (Disease)

Pegylated interferon-alpha (Chemical)

IFITM6 (Gene)

Feline b (Chemical)

E119V (Mutation)

Epac2 (Gene)

GFTP2 (Gene)

Hepatitis coronavirus infection (Disease)

Ouabain (Chemical)

LY6G (Gene)

CONCLUSION

- The construction of co-occurrence network embeddings for COVID-19 and related coronavirus infectious diseases.
- Published in the Journal of the American Medical Informatics Association. 2020 May 27.
- <https://academic.oup.com/jamia/advance-article/doi/10.1093/jamia/ocaa117/5847598>
- <https://github.com/shenfc/COVID-19-network-embeddings>
- <https://www.davidoniani.com/covid-19-network>

FUTURE WORK

- Adding more keywords to SPARQL query
- Including weights over edge for training purpose
- Manual evaluation by clinical investigators

ACKNOWLEDGEMENT

- CORD-19, CORD-19-on-FHIR
- David Oniani, Guoqian Jiang, Hongfang Liu
- National Institute of Health (NIH) U01TR0062-1.