
Diseases and Genes



Sky Martin and Kiersten Johnson



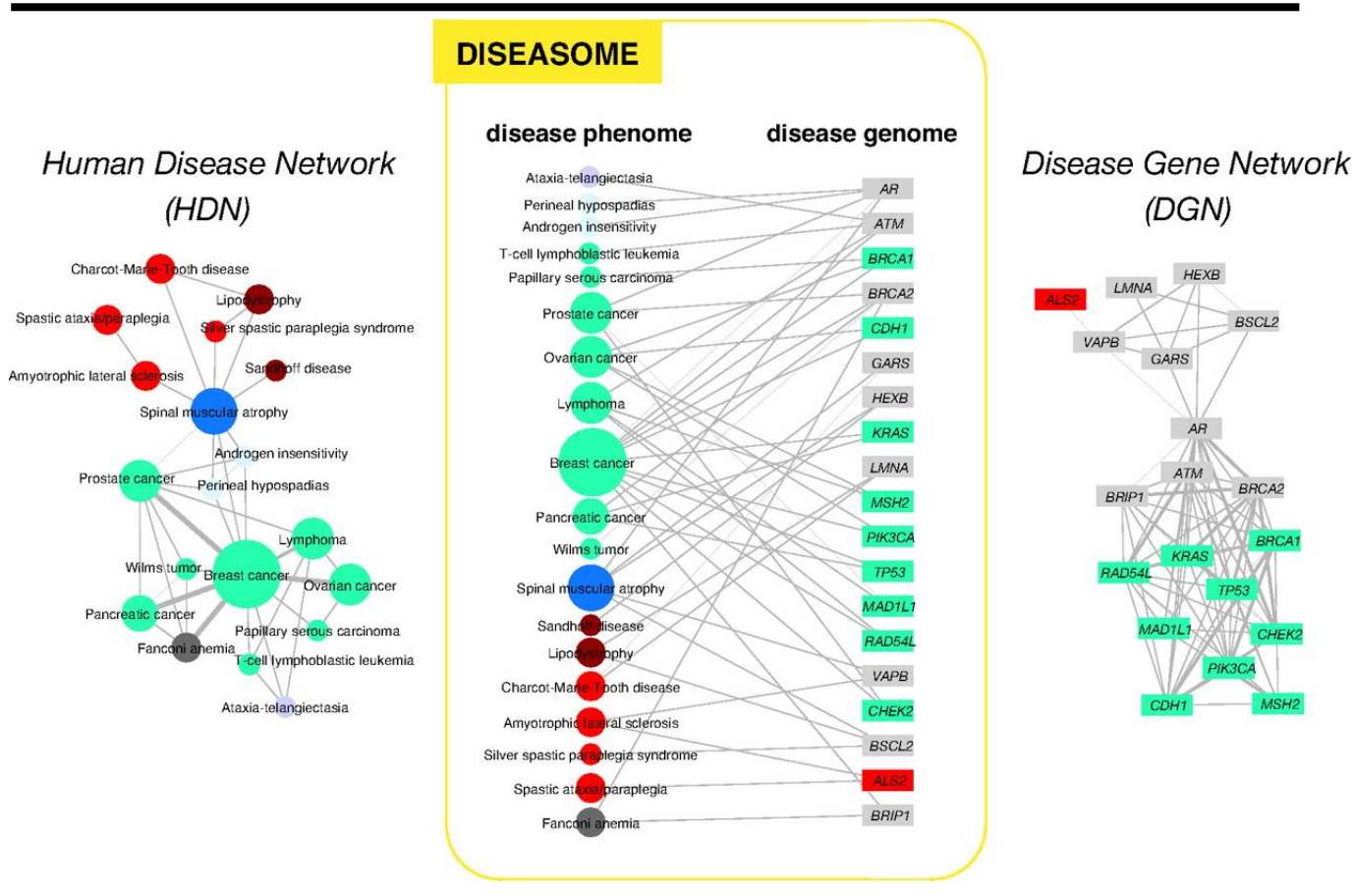
Online Mendelian Inheritance in Man

The Human Gene Disease Network

2007

Hand labeled subset of OMIM

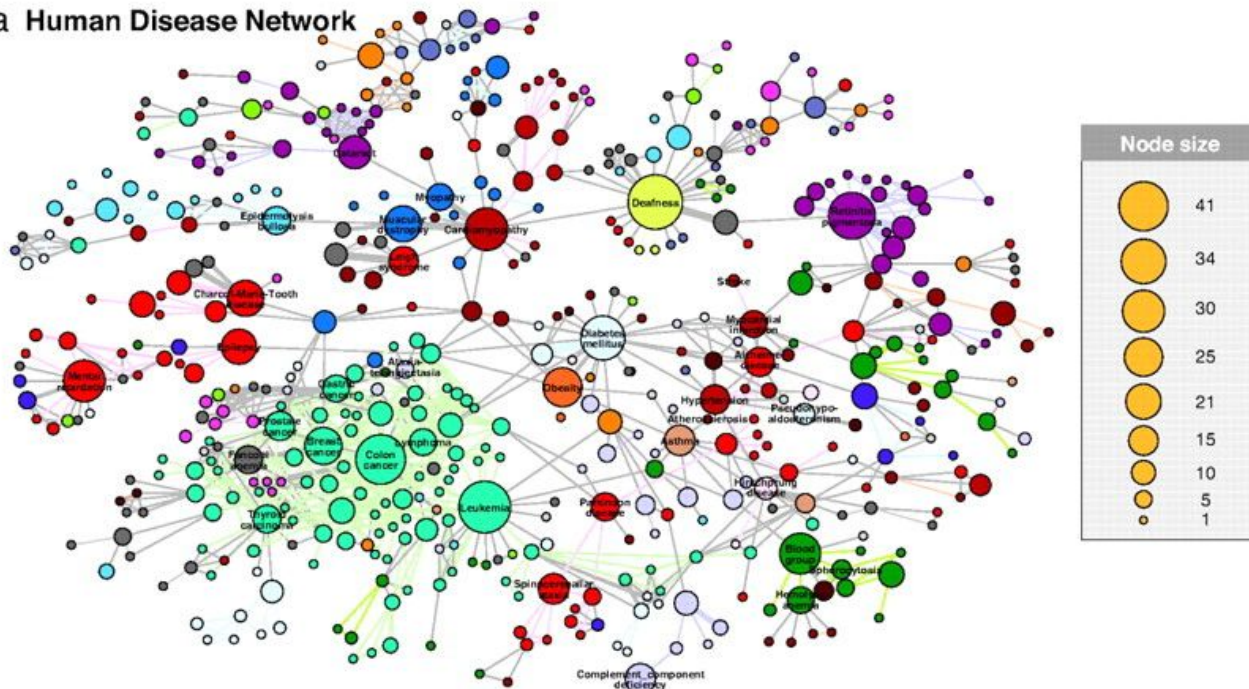
Diseases: 1,284
Genes: 1,777



Human Disease Network (HDN)



a Human Disease Network



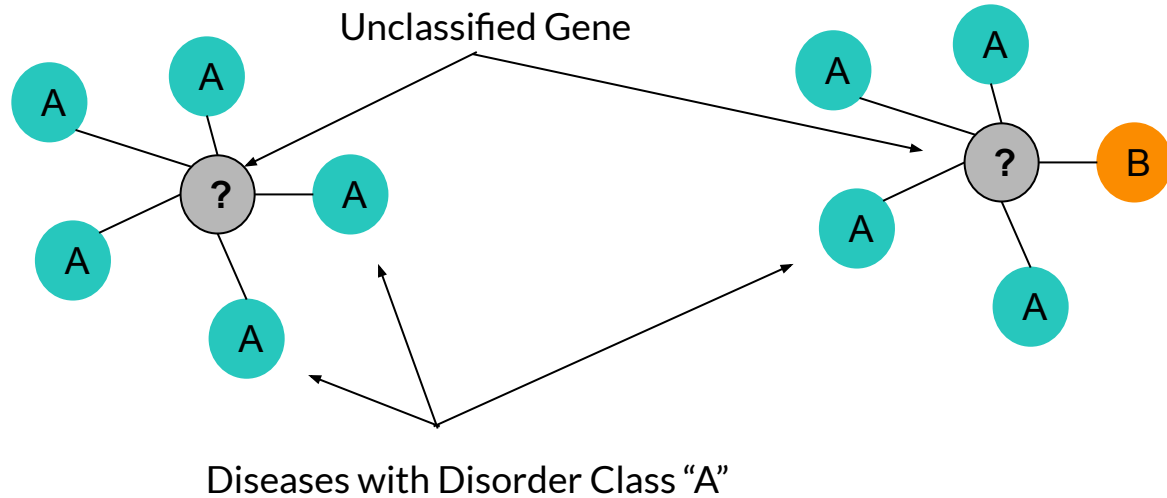
Node size



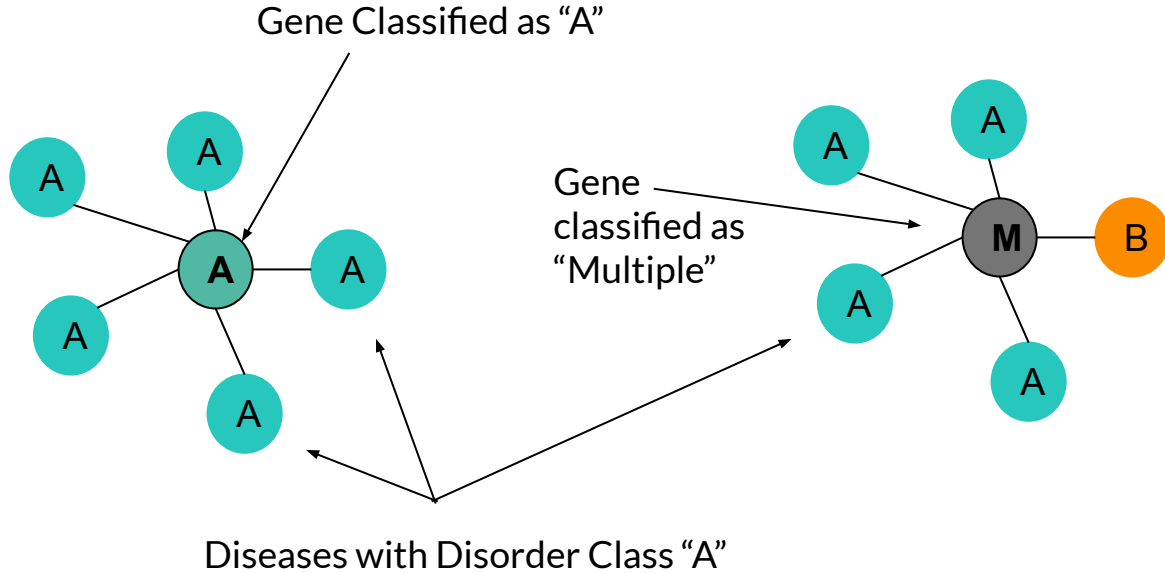
Node: Disease

Edge: Shared gene

Disorder Class Assignment to Genes

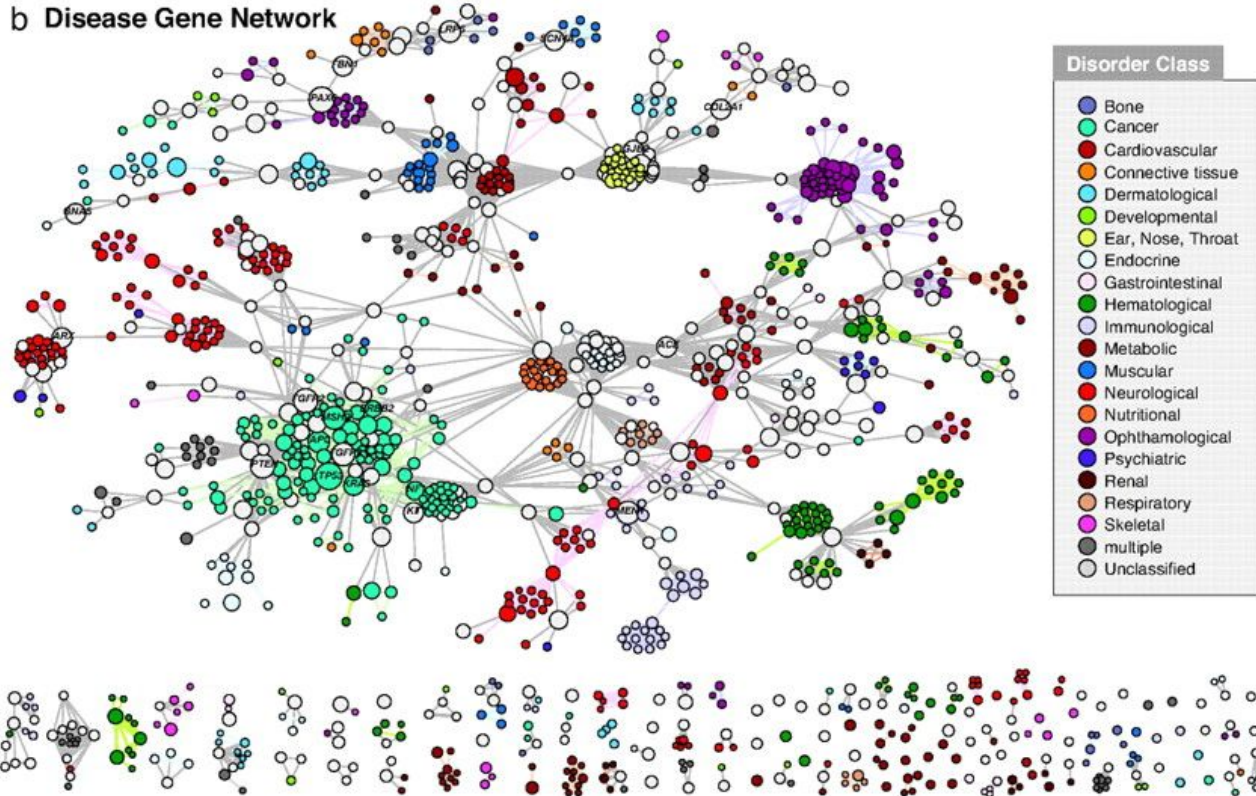


Disorder Class Assignment to Genes



Disease Gene Network (DGN)

b Disease Gene Network



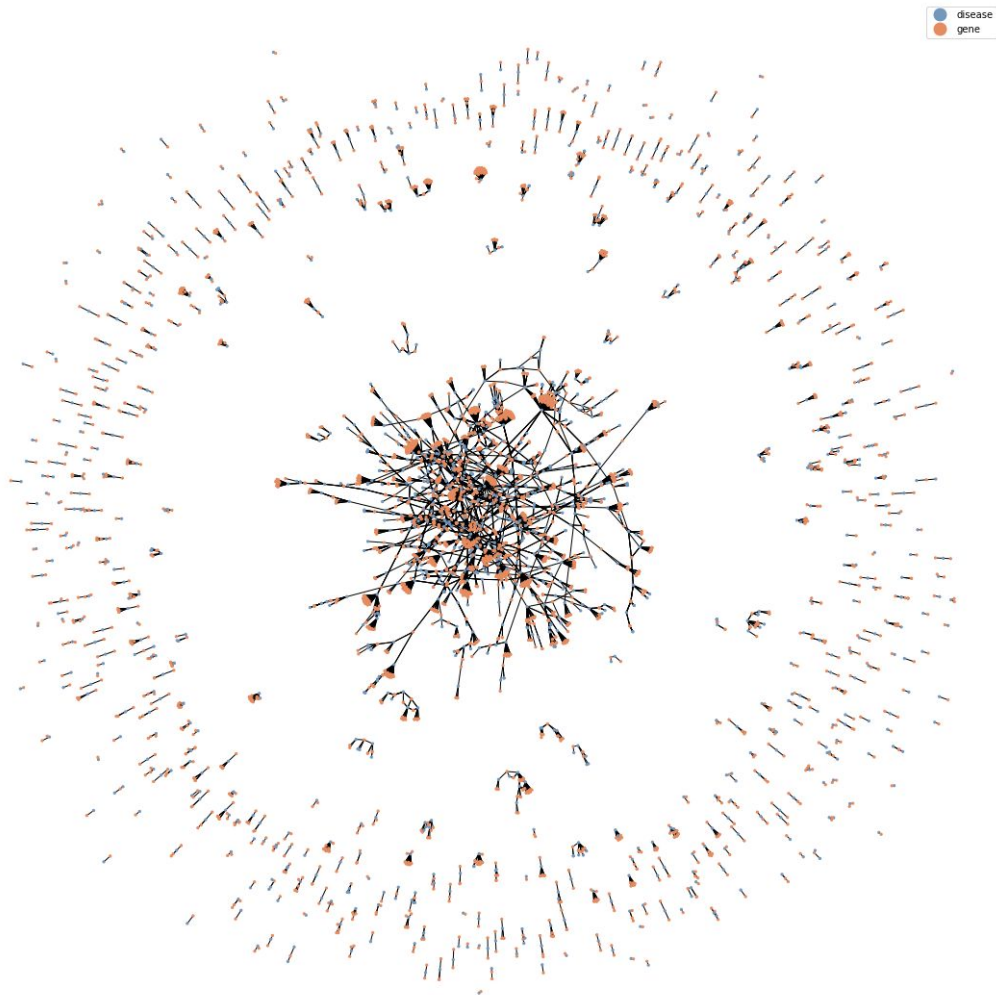
Node: Genes

Edge: Shared disease

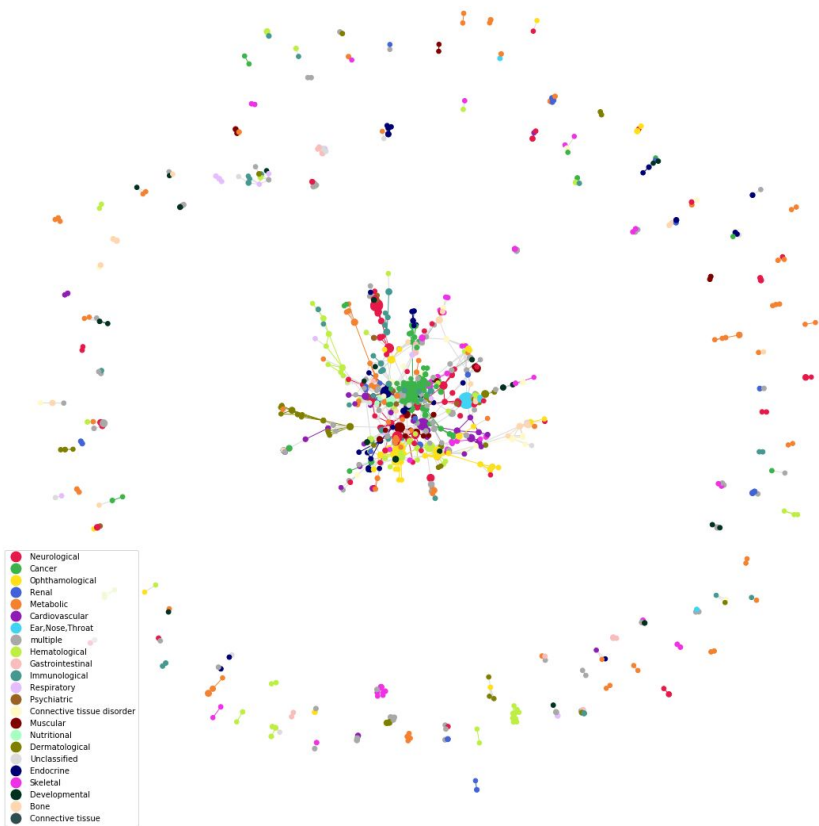
Our Visualizations and Analysis

Human Gene Disease Data

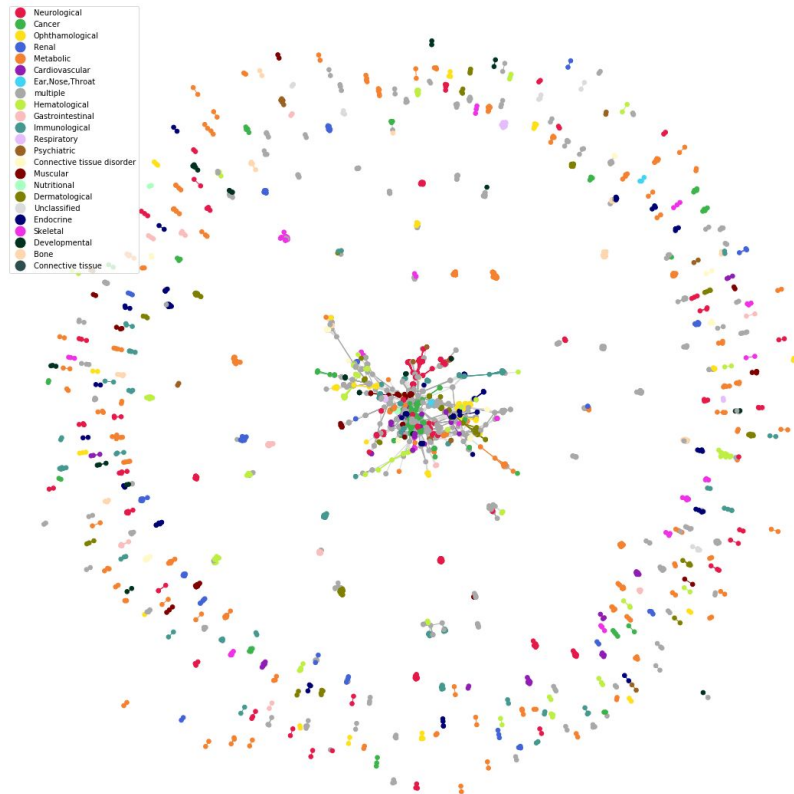
Updated genes from OMIM



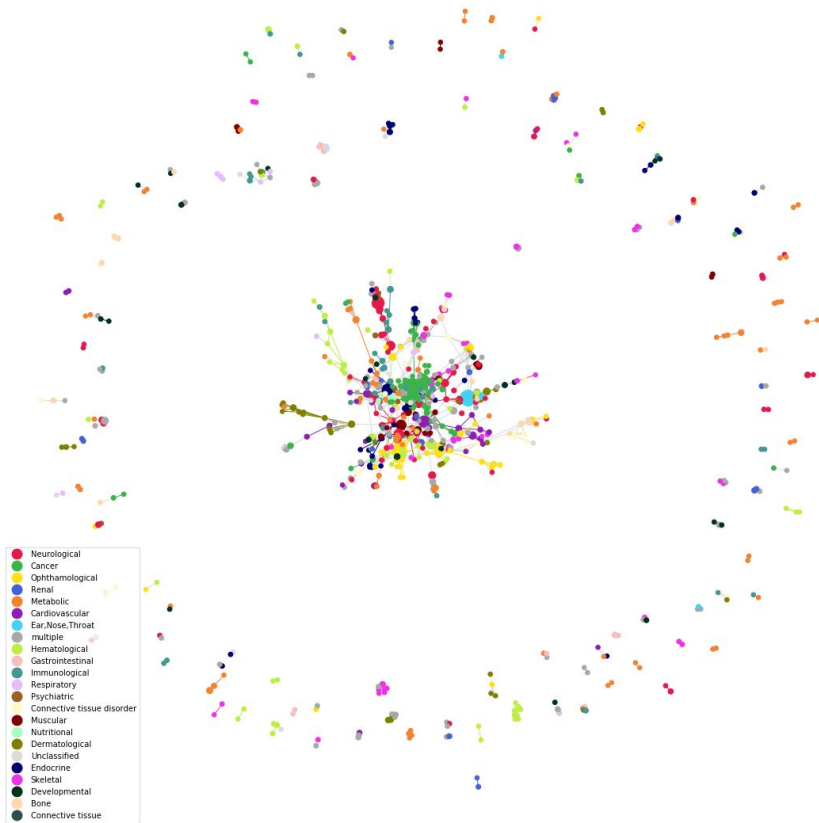
Disease Network



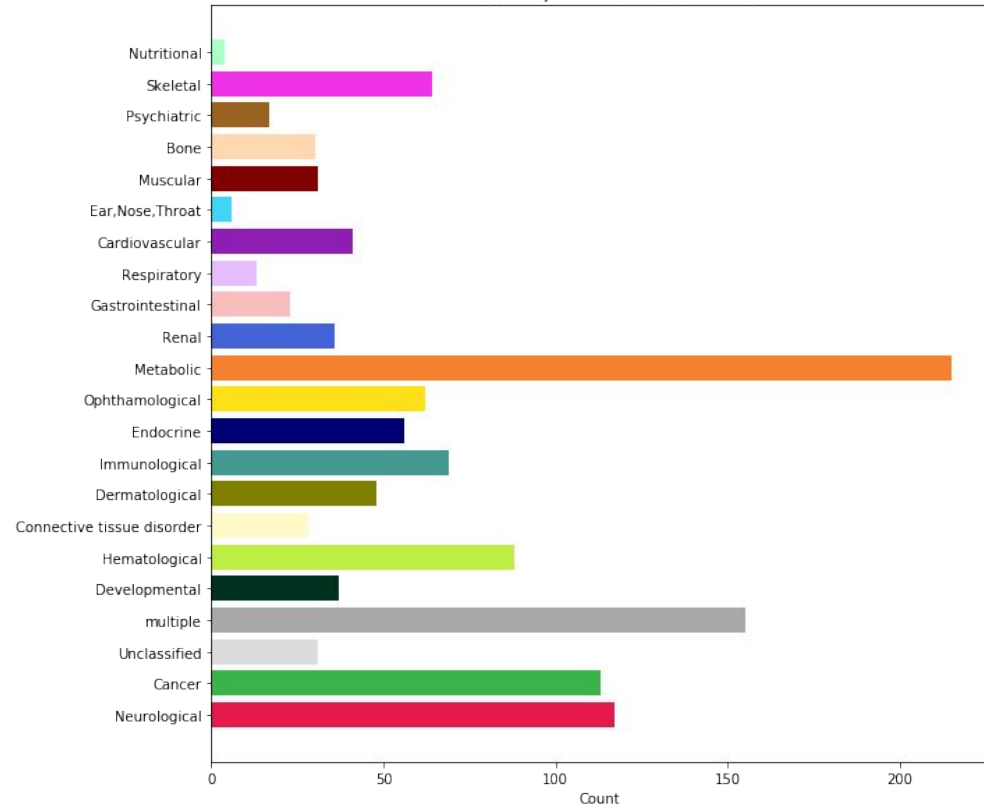
Gene Network



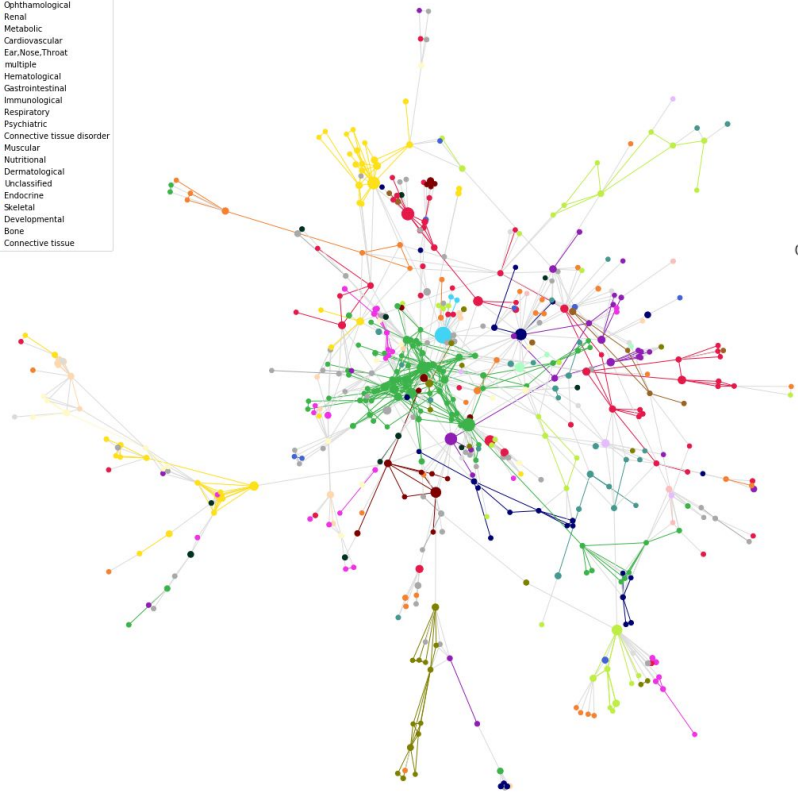
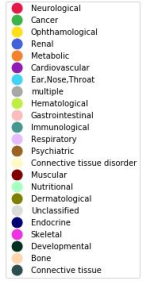
Disease Network



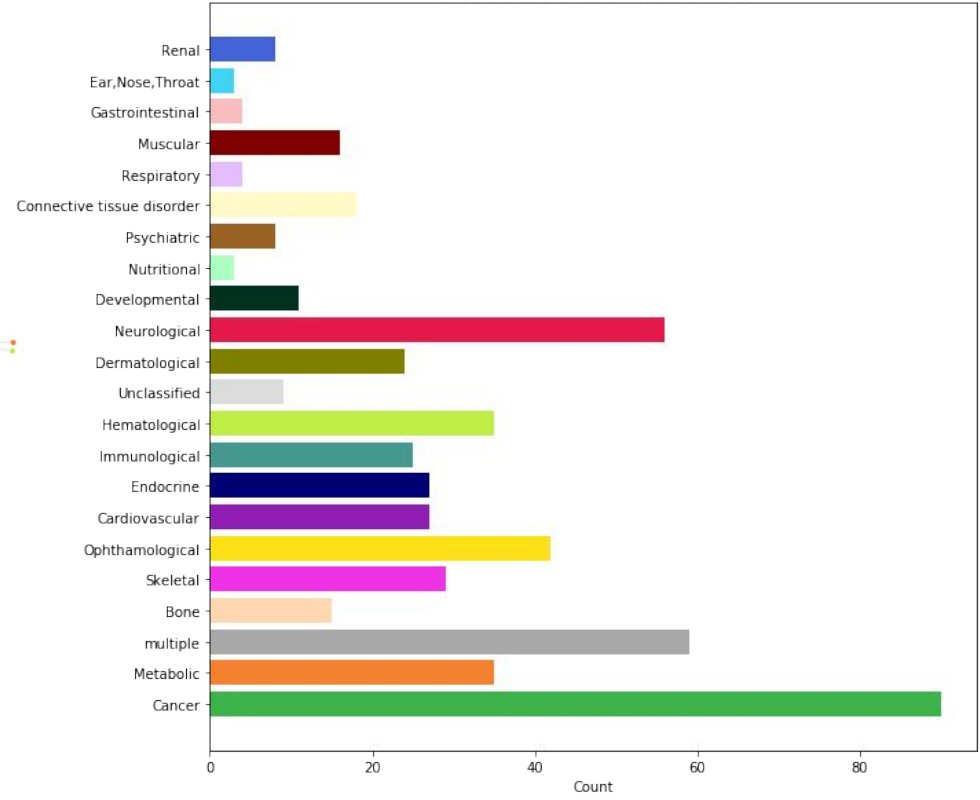
Disease Projection Class Distribution



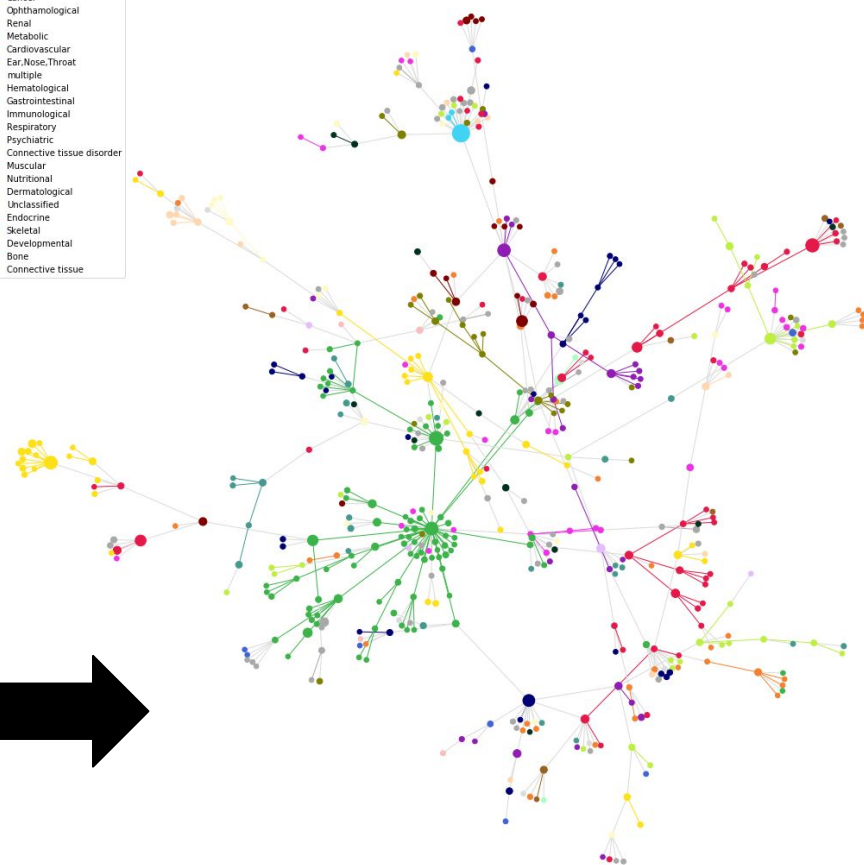
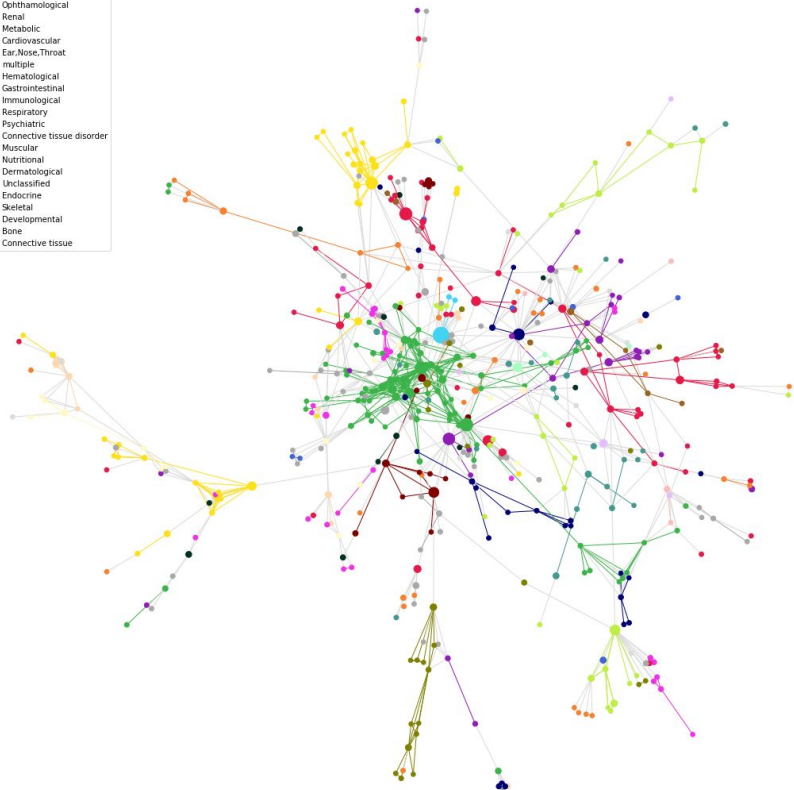
Connected Component: Disease Network



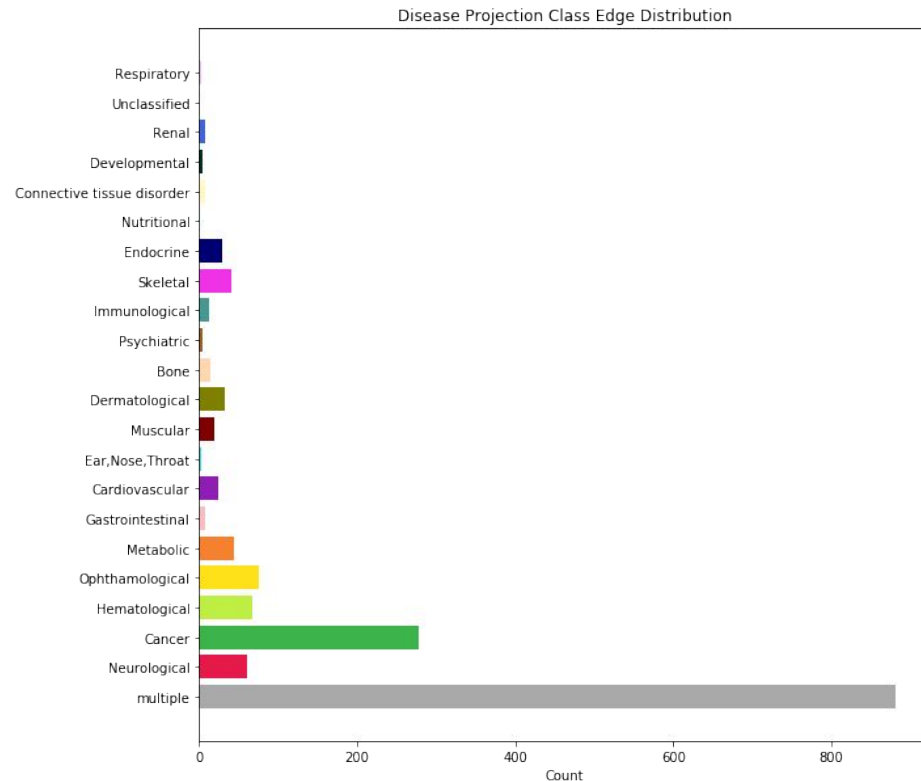
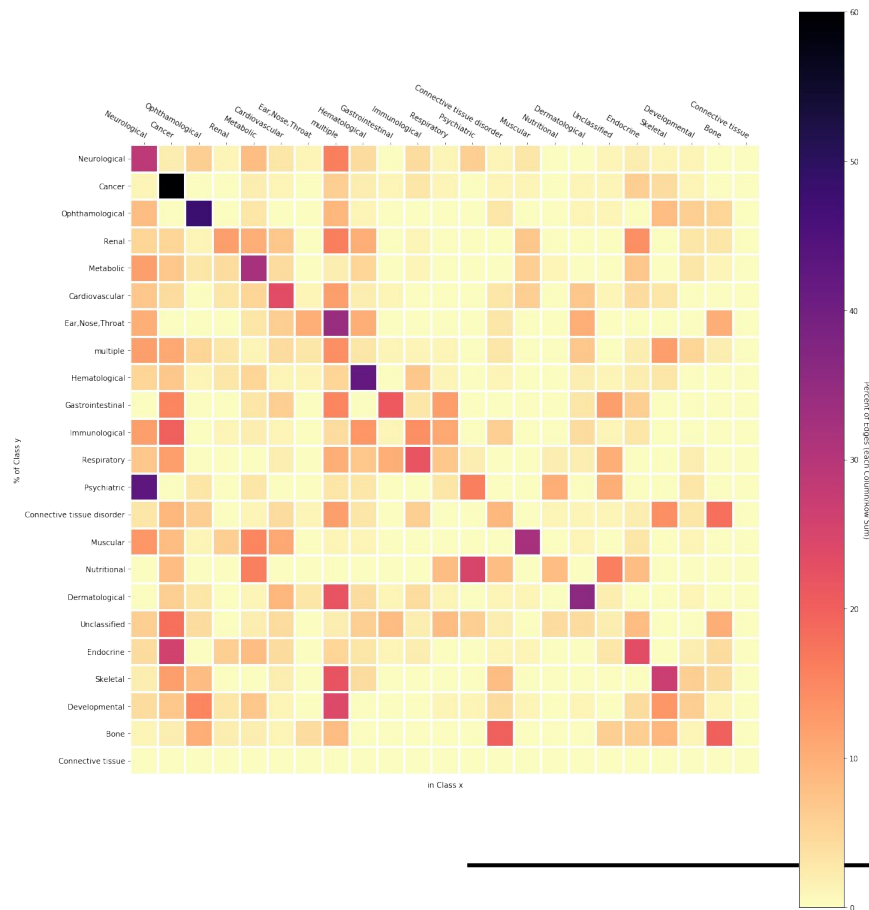
Disease Projection Class Distribution



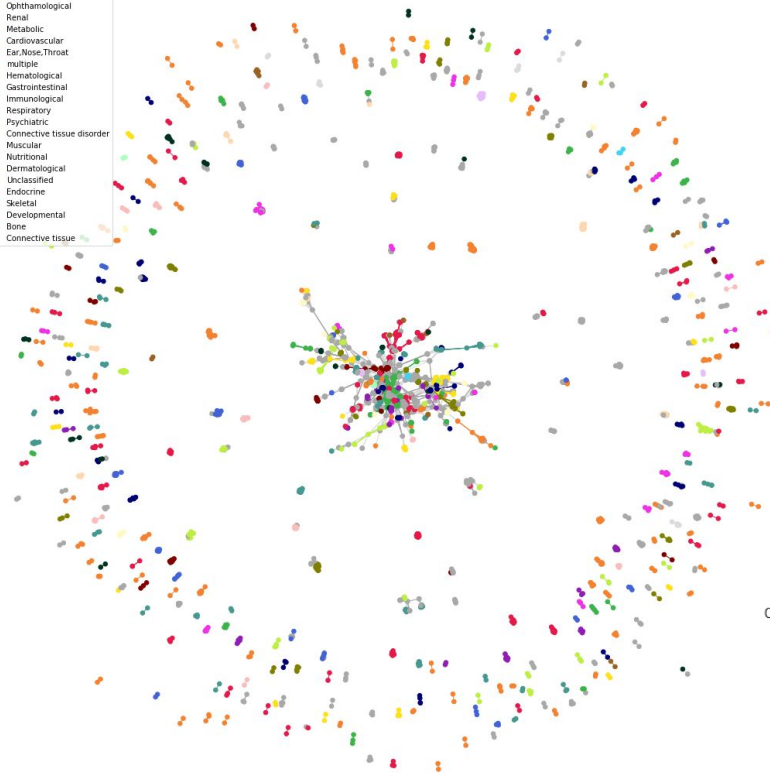
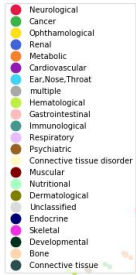
Breadth First Tree Traversal: Disease Network



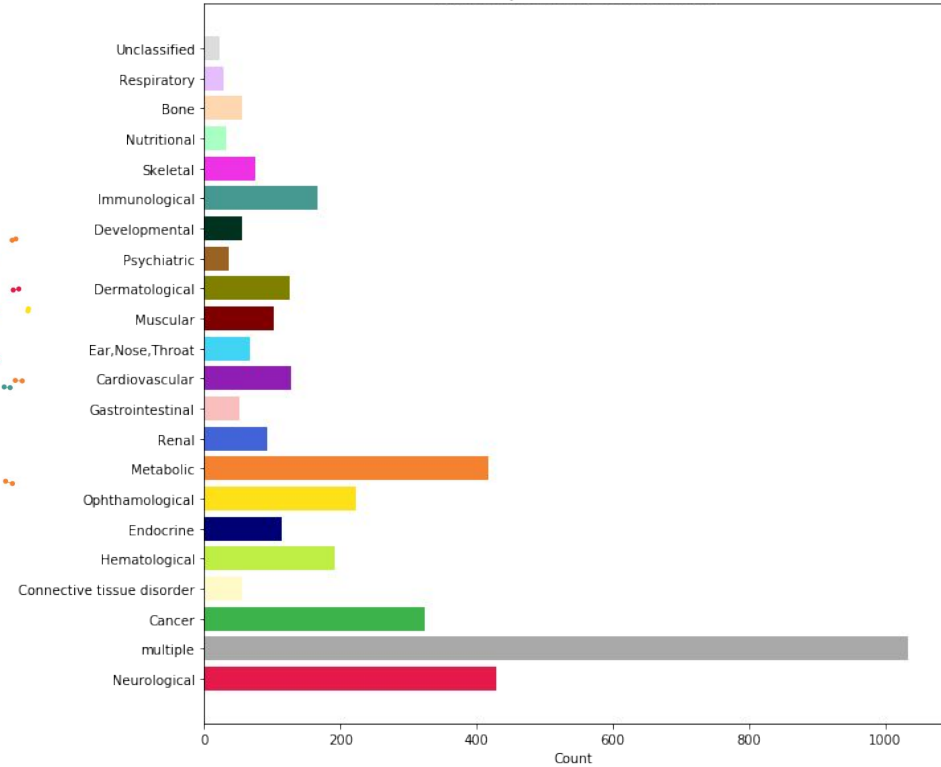
In v. Out Edges: Disease Network



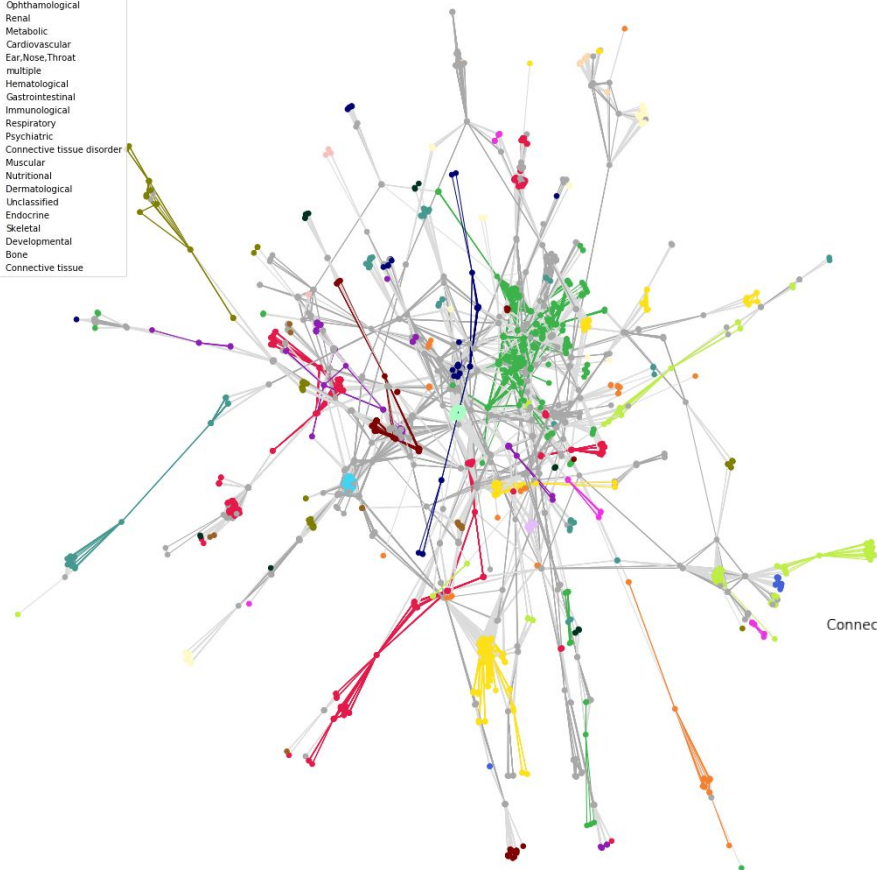
Gene Network



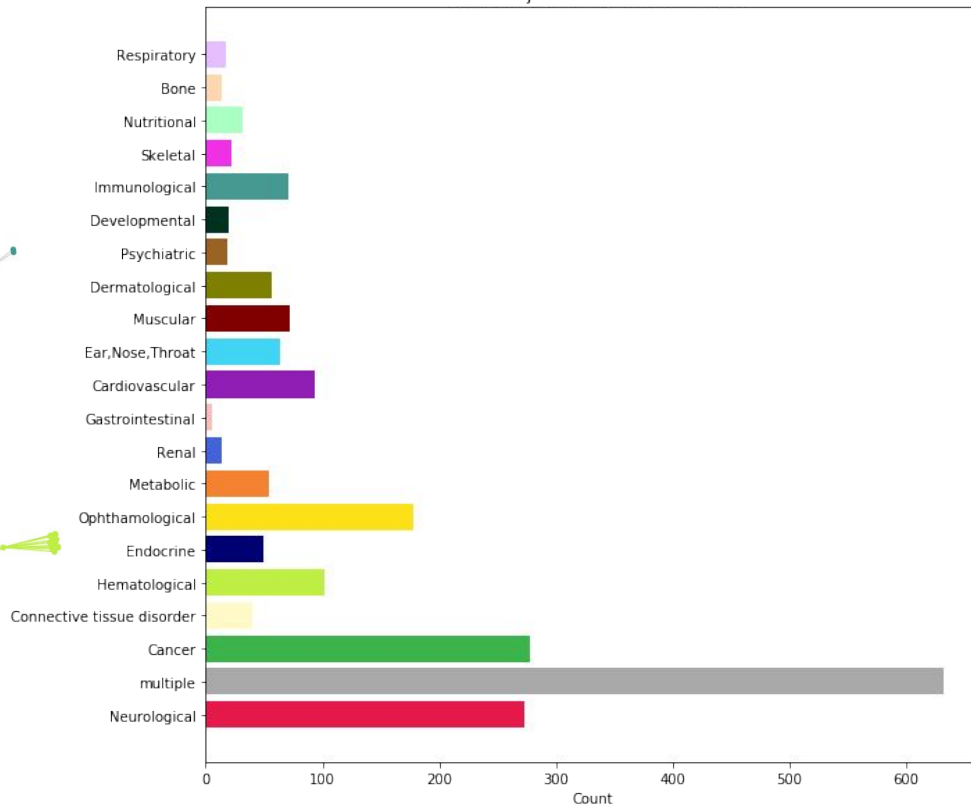
Gene Projection Class Distribution



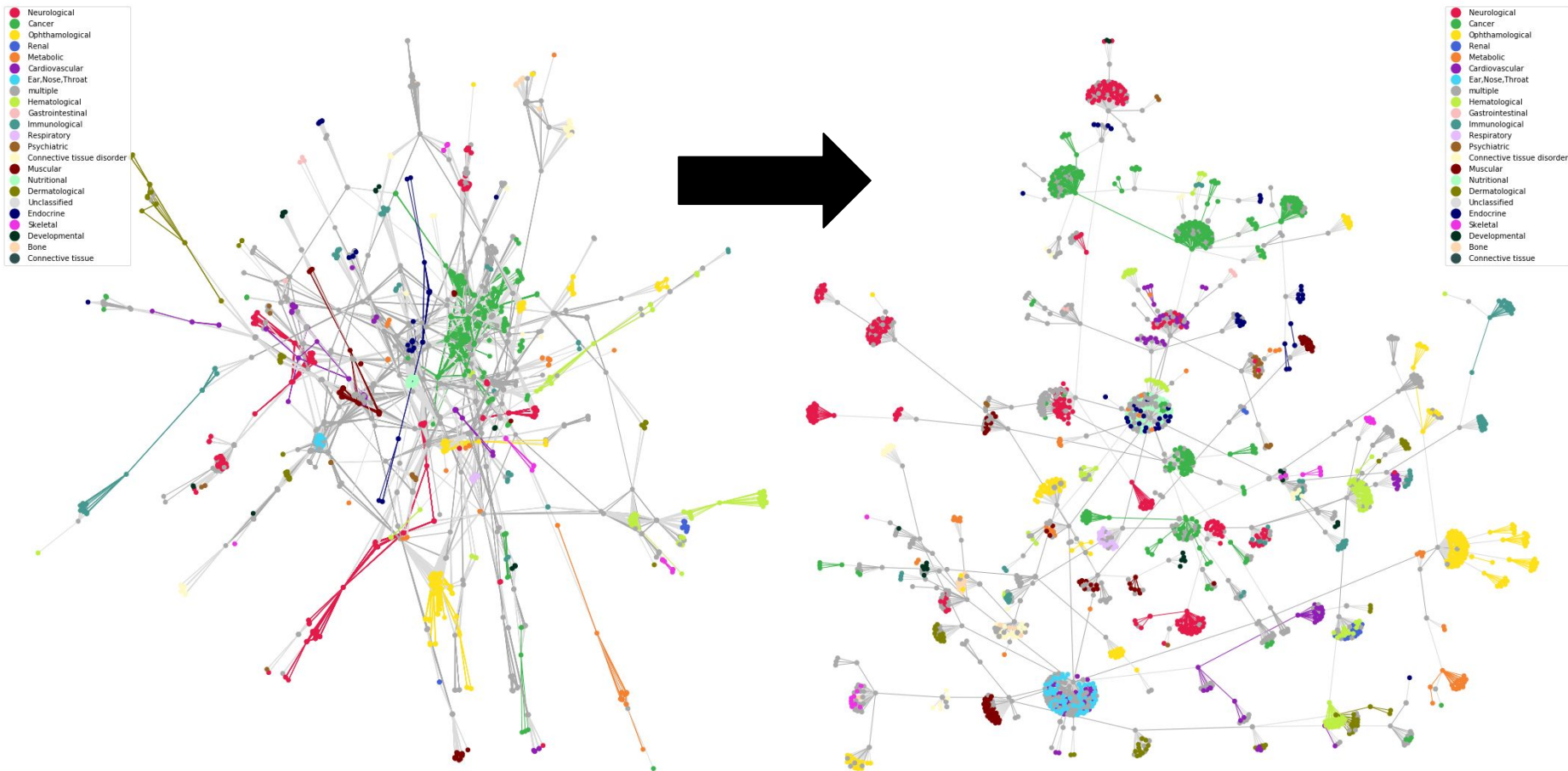
Connected Component: Gene Network



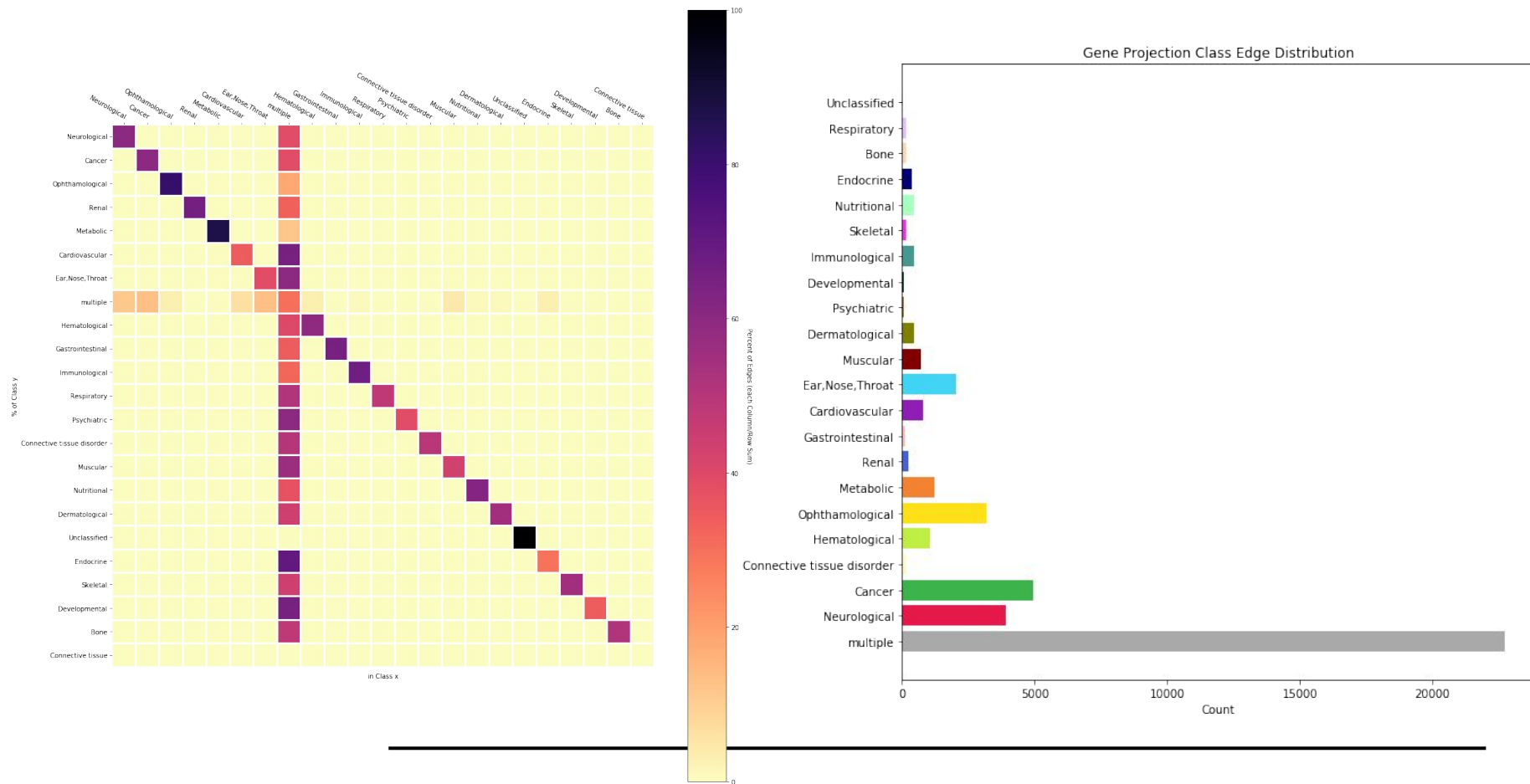
Gene Projection Class Distribution



Breadth First Tree Traversal: Gene Network



In v. Out Edges: Gene Network



Attribute - Disorder Class Prediction

why

Assign disorder class to full OMIM data set.

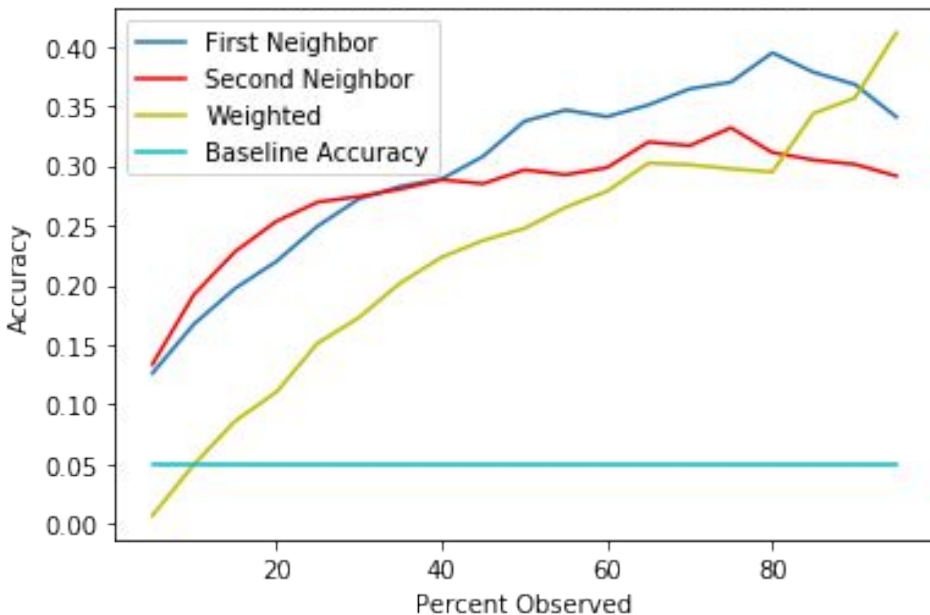
Prediction Algorithms

1. First Neighbor
2. Second Neighbor
3. Weighted
4. Baseline (random)

Predicting Classes of Disorders

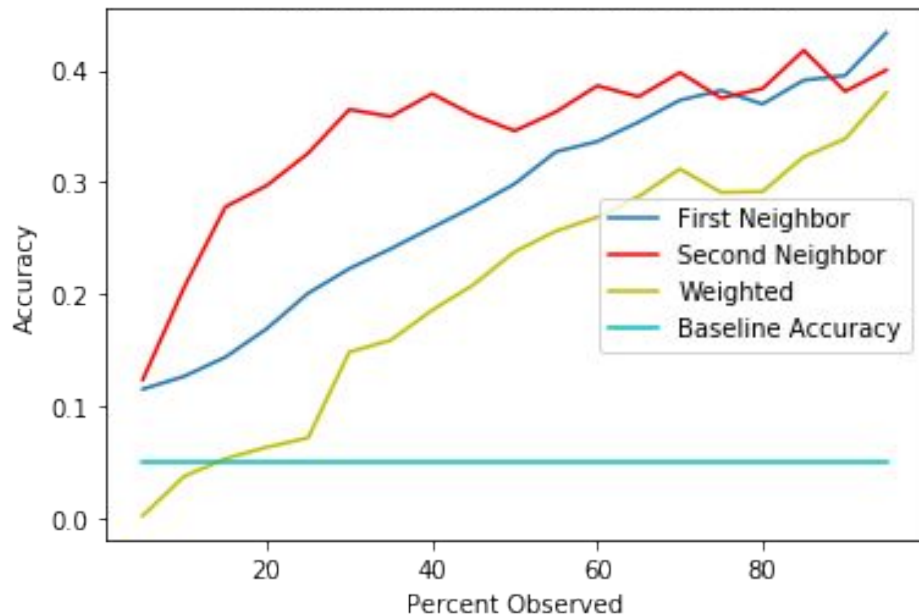
Normal Projection

Percent Correct vs Percent Observed

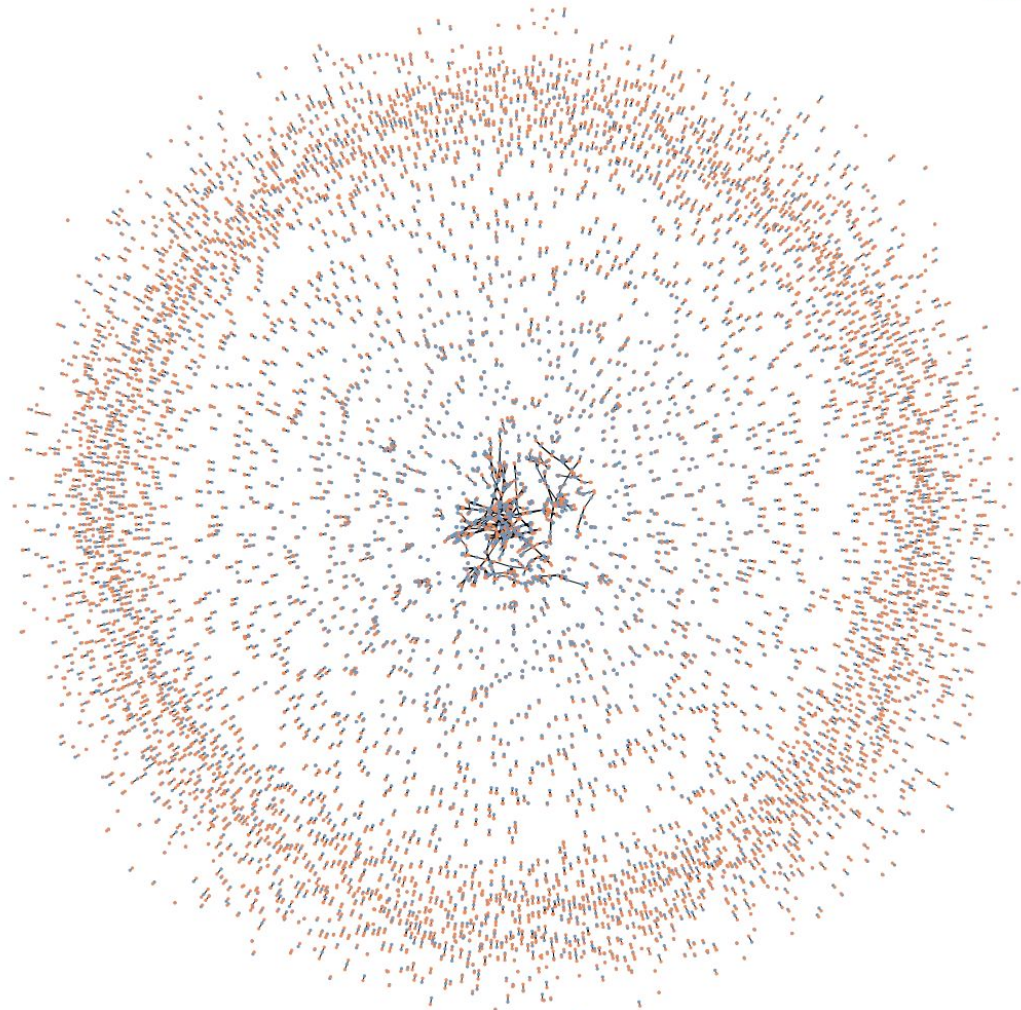


Breadth First Traversal

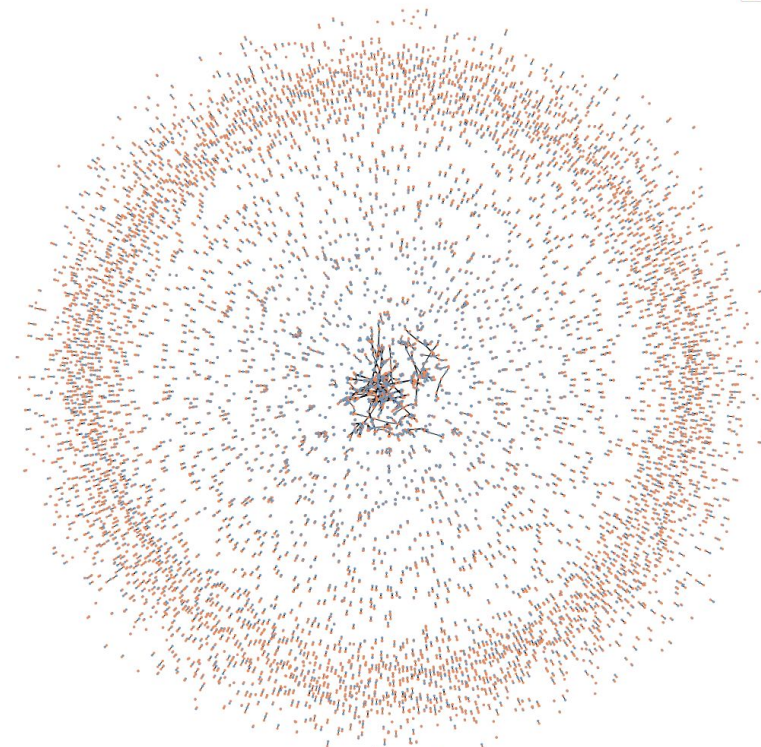
Percent Correct vs Percent Observed



OMIM Data



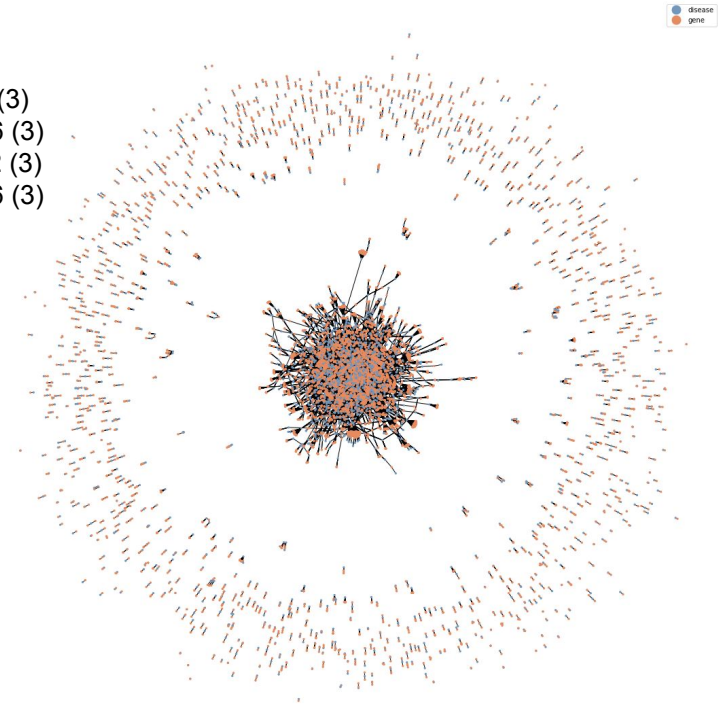
Condensing Diseases



Spinocerebellar ataxia 1, 164400 (3)
Spinocerebellar ataxia 10, 603516 (3)
Spinocerebellar ataxia 11, 604432 (3)
Spinocerebellar ataxia 12, 604326 (3)



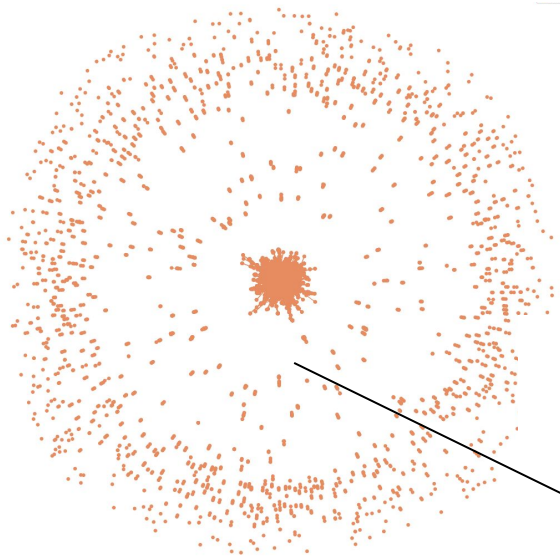
Spinocerebellar ataxia



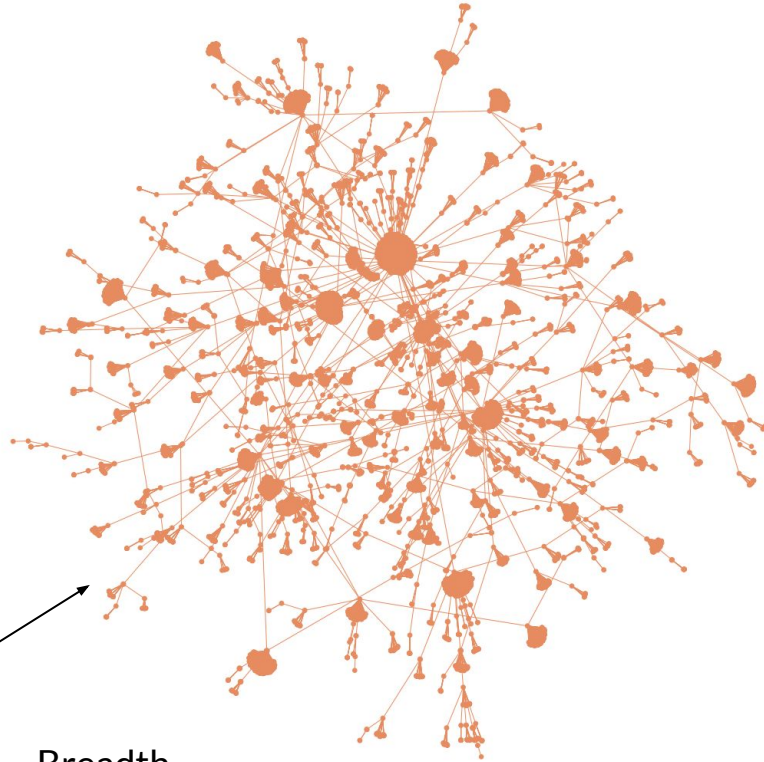
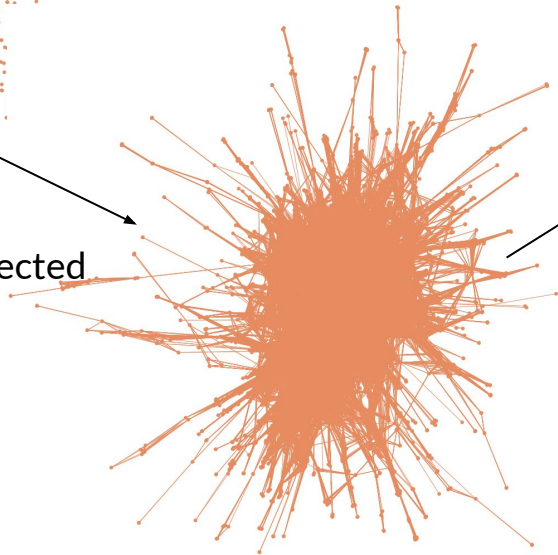
number of nodes, n = 21733
number of edges, m = 24398
mean degree, $\langle k \rangle$ = 2.25
mean geodesic distance, $\langle \ell \rangle$ = 1.67
clustering coefficient, C = 0.00

number of nodes, n = 16844
number of edges, m = 22804
mean degree, $\langle k \rangle$ = 2.71
mean geodesic distance, $\langle \ell \rangle$ = 6.39
clustering coefficient, C = 0.00

OMIM Gene Network

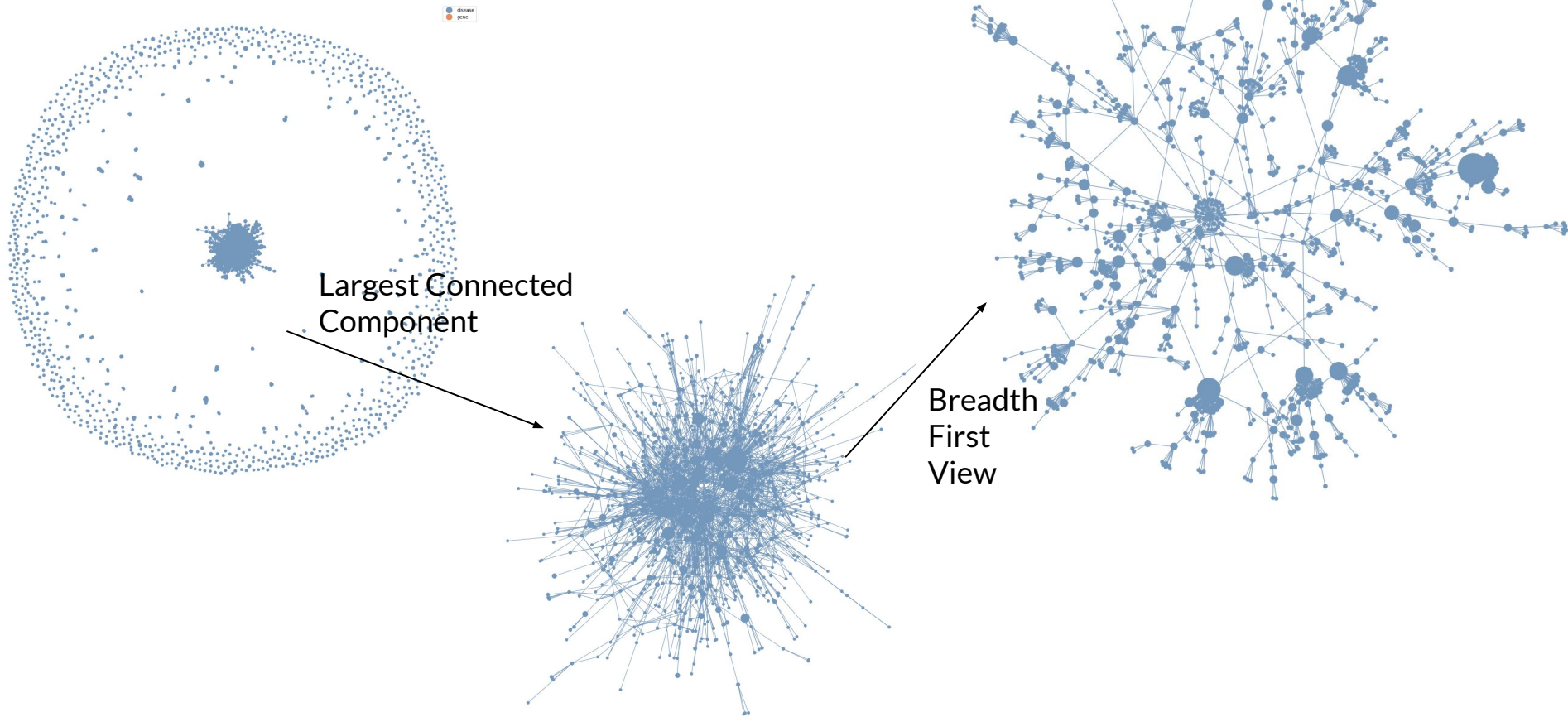


Largest Connected
Component



Breadth
First
View

OMIM Disease Network



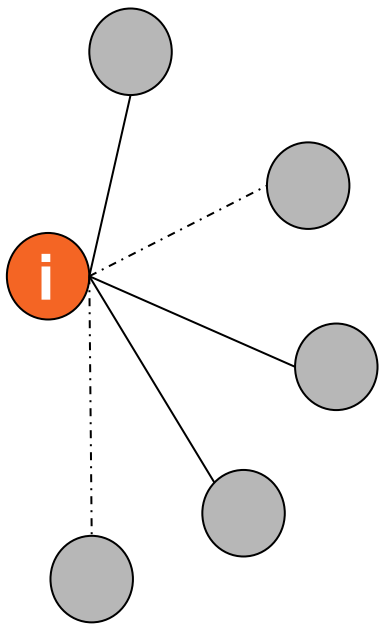
Link Prediction

Value of Link Prediction

Link prediction between disease and gene would help researchers narrow their research to relevant genes.

N-Fold Validation

1. Hide proportion of edges from i
2. Calculate similarity measures for every possible unobserved edge s.t. $(i, j) \nexists j$ in G
3. Rank the unobserved edges by measure, pick top P (hyper-parameter) edges to predict as edges
4. Calculate accuracy metrics
5. Average over all i in G



Unobserved Edges

Similarity Measure	i	j	Pred
10	1	2	1
3	1	4	1
1	1	6	0
...			
0	1	45	0

P

Similarity Measures

$\Gamma(x)$ - set of x 's neighbors

1. Common Neighbors $O(Nk^2)$

$$\# \text{ Common Neighbors} = |\Gamma(x) \cap \Gamma(y)|$$

2. Jaccard Coefficient $O(Nk^2)$

$$\frac{\# \text{ Common Neighbors}}{\# \text{ Neighbors}} \quad \left| \frac{\Gamma(x) \cap \Gamma(y)}{\Gamma(x) \cup \Gamma(y)} \right|$$

3. Preferential Attachment

$$\text{Degree of } x * \text{Degree of } y \quad |\Gamma(x)| * |\Gamma(y)|$$

4. Adamic/Adar Index

$$\# \text{ Common Neighbors with degree penalty} \quad \sum_{z \in \Gamma(x) \cap \Gamma(y)} \frac{1}{\log |\Gamma(z)|}$$

5. Leicht-Holme-Newman Index

$$\frac{\# \text{ Common Neighbors}}{\text{Degree Product}} \quad \frac{|\Gamma(x) \cap \Gamma(y)|}{k_x * k_y}$$

6. Katz_B Index

Path lengths
between x and y
weighted by B

$$\sum_{l=1}^{\infty} \beta^l \cdot |\text{paths}_{xy}^{(l)}|$$

Disclaimers

- OMIM network is far from fully observed
- Link prediction measures were created to be on unipartite networks, because of common neighbors

Important Metrics

Recall - $TP/(TP+FN)$

Proportion of actual edges that are detected

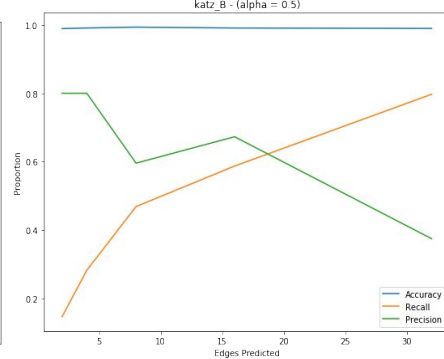
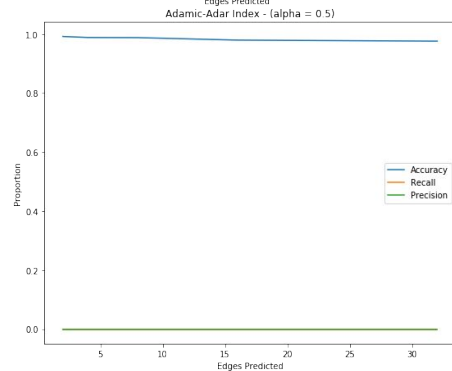
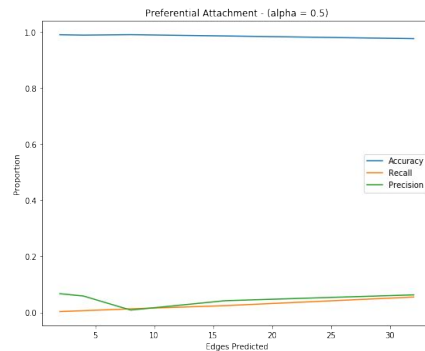
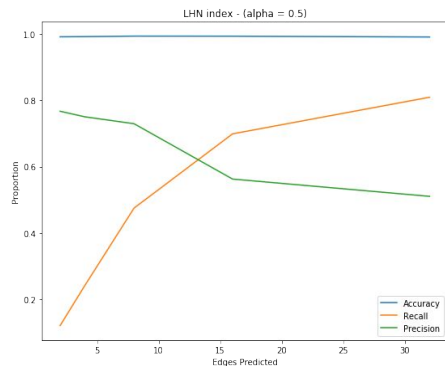
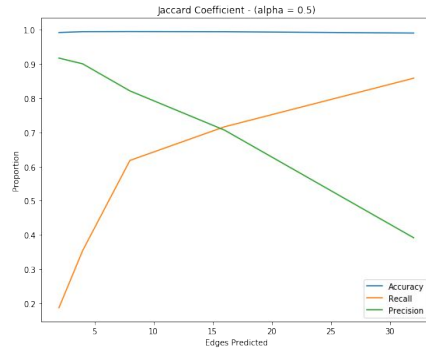
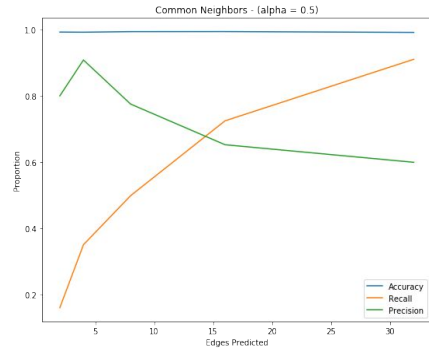
Precision - $TP/(FP+TP)$

Proportion of edges detected that are actual edges

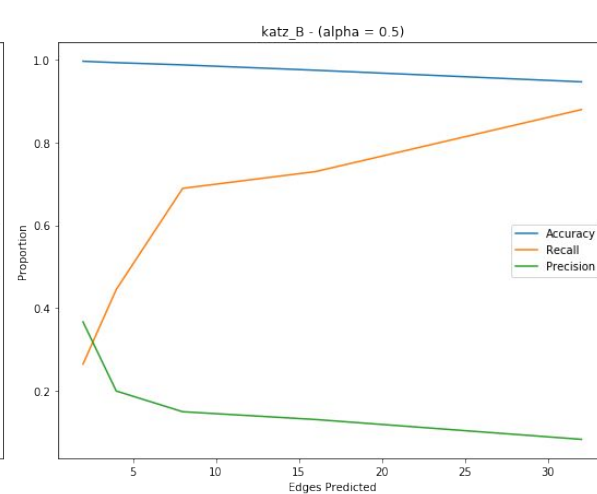
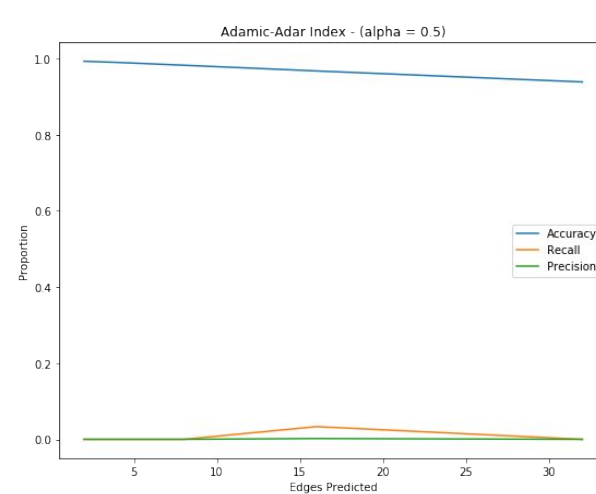
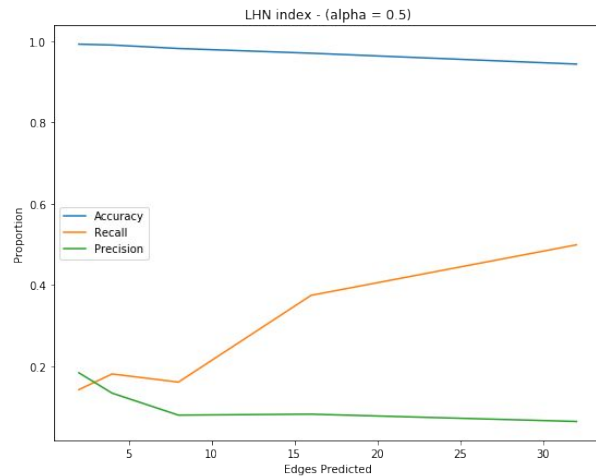
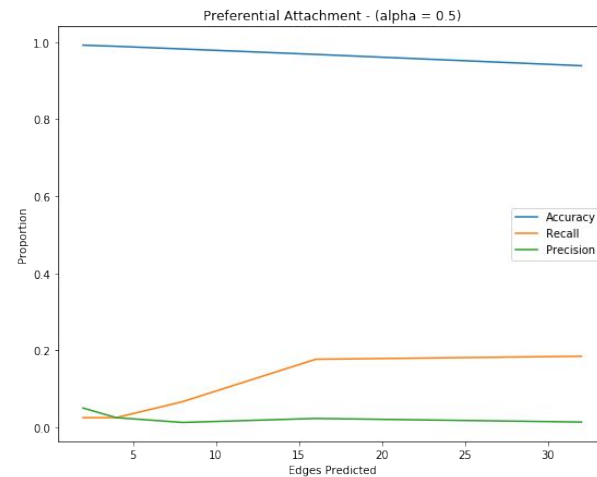
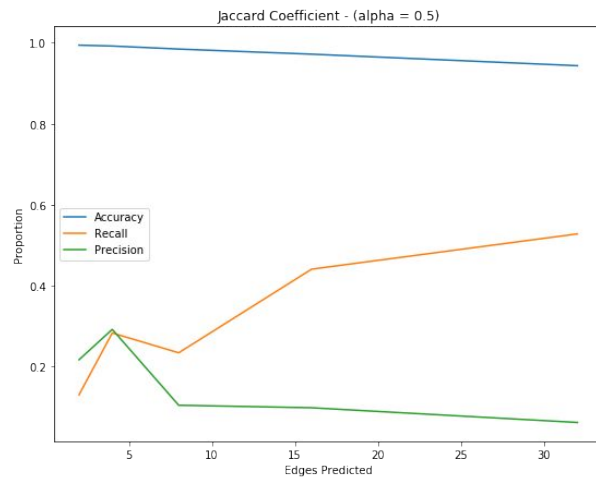
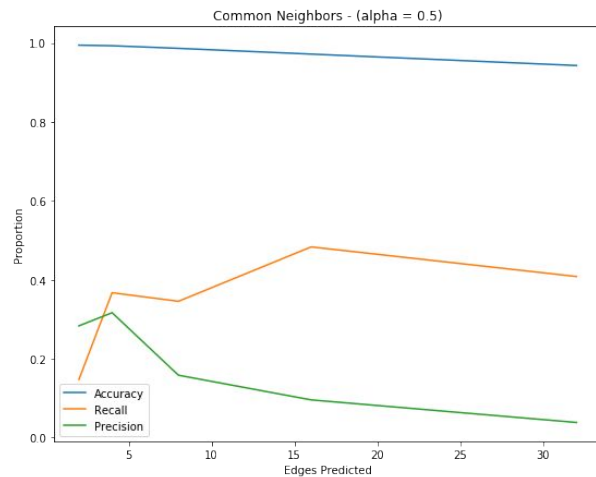
Disclaimers

- OMIM network is far from fully observed
- Link prediction measures were created to be on unipartite networks, because of common neighbors

Gene Network



Disease Network



Future Directions

- Estimate class for OMIM
- Create link prediction for bipartite network