

EPFL

Data Science in Practice

Covid-19 Risk Factor Prediction



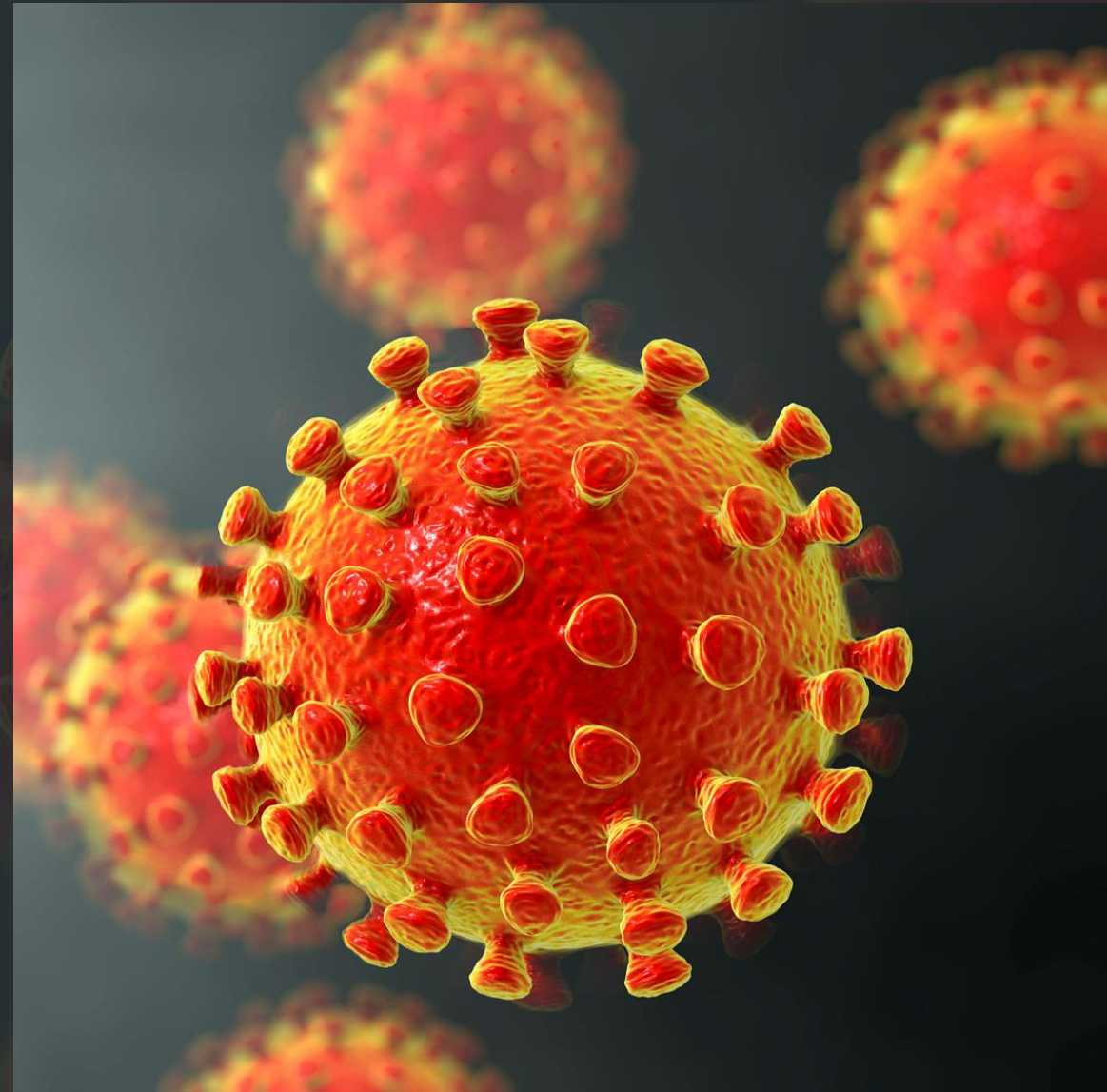
*R. Chaouche, Y. Martinson,
C. Padovani, J. Triomphe*

Context

Introduction

- I. Problem description
- II. Data Management
- III. Network analysis
- IV. Random forest regression
- V. Interpretations

Conclusion





Introduction

Introduction

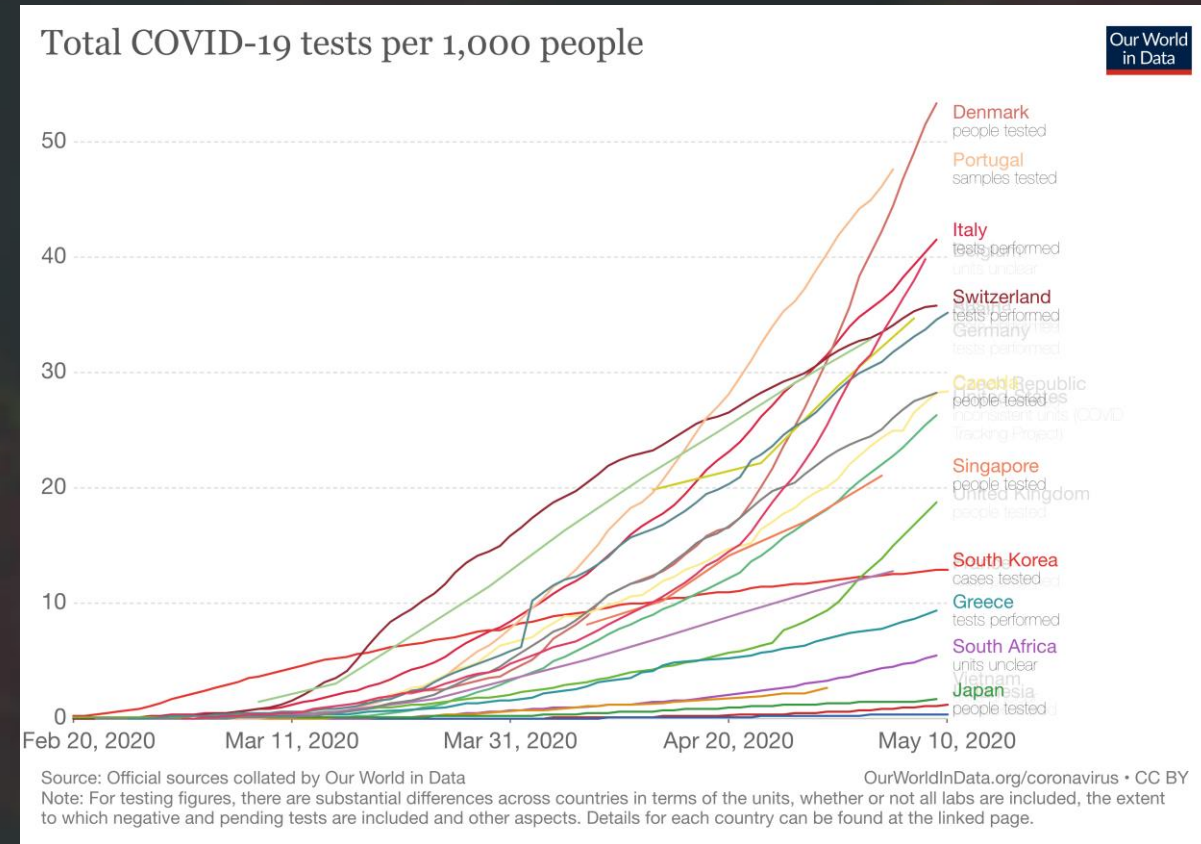
→ Virus has been spreading worldwide

- 4 million confirmed cases
- Over 265,000 deaths

→ Global Lock – down

- Shut-down of large pans of the economy
- Long-term implications

→ Tests to target SARS-CoV-2





I. Problem Description

Understanding the need



→ **Health issues**

- High transmissibility and more than 265,000 deaths (06 May)
- Risk factors diseases known but poor literature about drug targets*



→ **Expensive testing campaigns**

- \$100 /genome tests
- \$30 million for testing Switzerland (35 test/1000 habitants)



→ **Tests do not have an optimal efficiency yet**

* Specific molecule , usually a protein, to which the drug binds to produce its effect

Challenges



→ **Mapping** a network of the major risk factor diseases for COVID-19



→ **Create a model** to identify which targets of those risk factor diseases are more likely to be associated to COVID-19

Strategy



→ Existing project from Kaggle

Predicted relation between diseases and COVID-19 with a score system



→ Our project

Predicts relation between the targets of those diseases and the COVID-19

Expected outputs and outcomes



→ **Accurate network**

- Will help scientist figure out gene families.
- Expected communities could focus on the most represented genes.

→ **Accurate prediction of targets to pay attention on**

- Will help doctors to avoid additionnal infection due to COVID-19
- Will help to create more accurate tests



II. Data Set

Initial Data Set: Kaggle project (1)



→ **CORD-19 Dataset challenge**

59,000 scholarly articles, including over 47,000 with full text, about COVID-19, SARS-CoV-2, and related coronaviruses created by The White House and a coalition of leading research groups

→ **Mondo Disease DataBase**

Aggregates and merges genetic associations curated from both literature and newly-derived loci¹ from UK Biobank

→ **Open Targets DataBase**

Semi-automatically constructed ontology that merges in multiple disease resources to yield a coherent merged ontology

Initial Data Set: Kaggle project (2)



→ **Connects diseases**

- With other diseases
- With drug targets

→ **Association scores between diseases and COVID-19**

Data Reconstruction



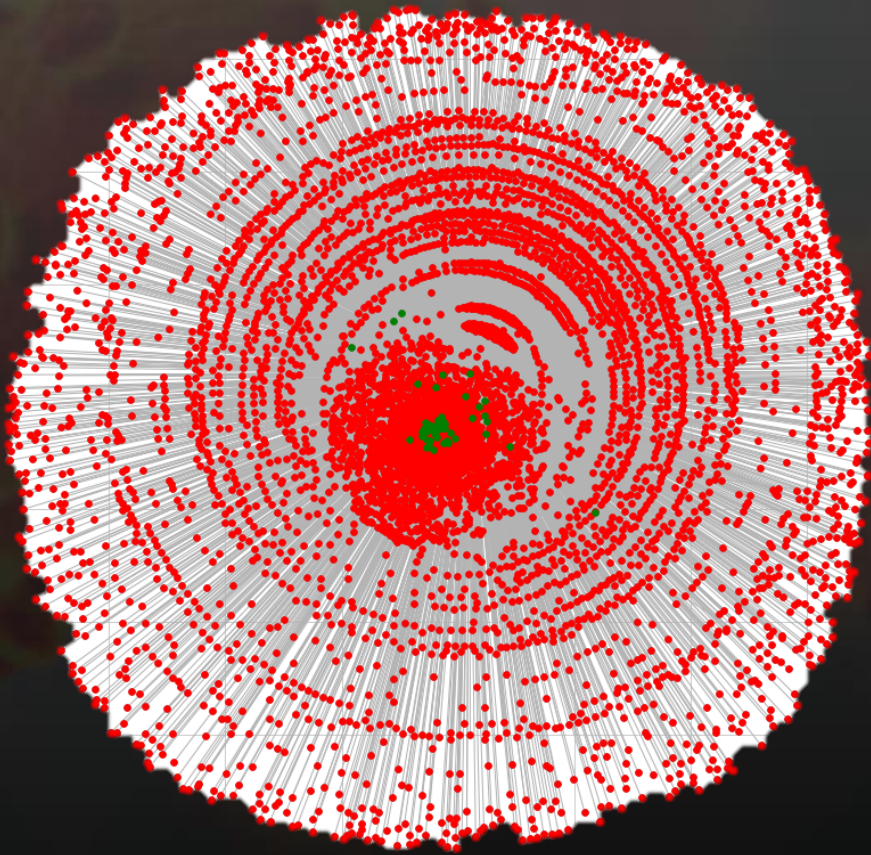
→ **From the output of the previous study**

- Combine disease-target association score with disease-COVID-19 association score
- Leverage the Open Targets Platform



III. Network Model

Model



→ Dataset I

- Nodes: **Diseases**, **Genes**
- Edges: weighted score of gene on disease
- Goal: Centrality ranking
- Output: Overall Ranking

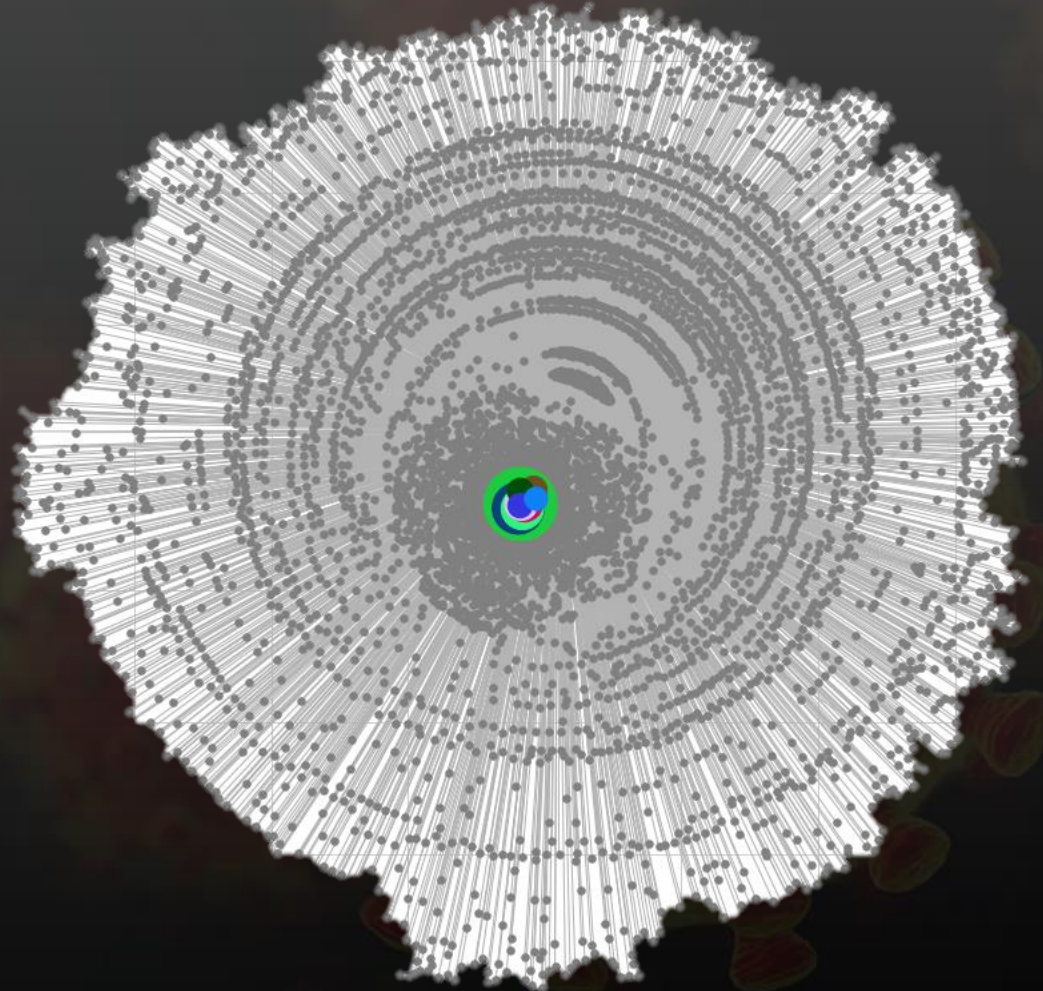
Model

- Centralities: Degree, Betweenness, Closeness, PageRank
- Scoring: Point ranking system for each node
- Positive: Finds key players in the network, robust
- Drawbacks: No Covid infos, unexpected outputs

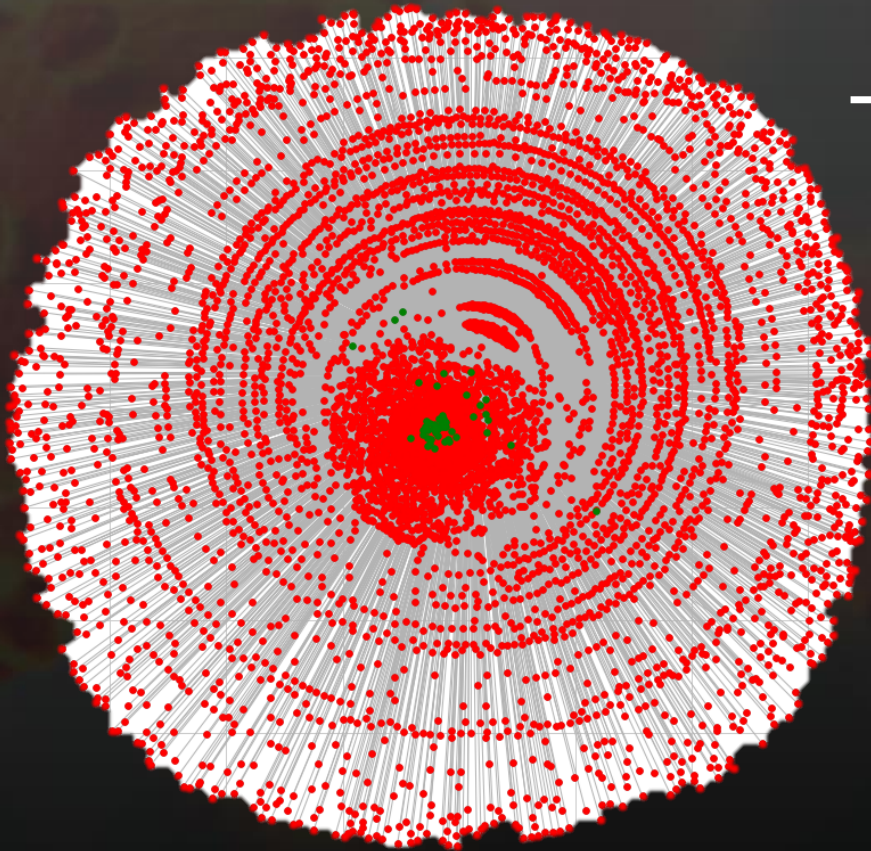
Results

→ Very centered, close to gene locations

1. Bone Disease
2. Psychiatric Disorder
3. Biological process
4. Protein measurement
5. Diabetes Mellitus



Model

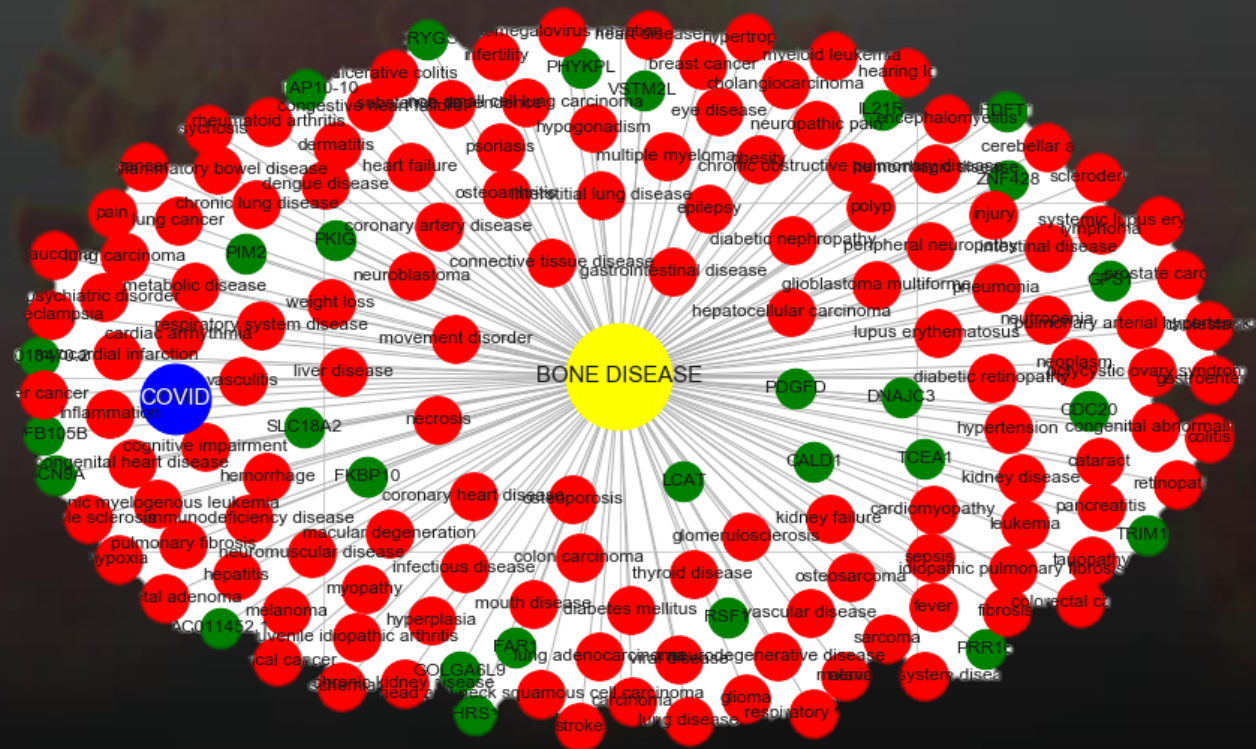


→ Dataset II

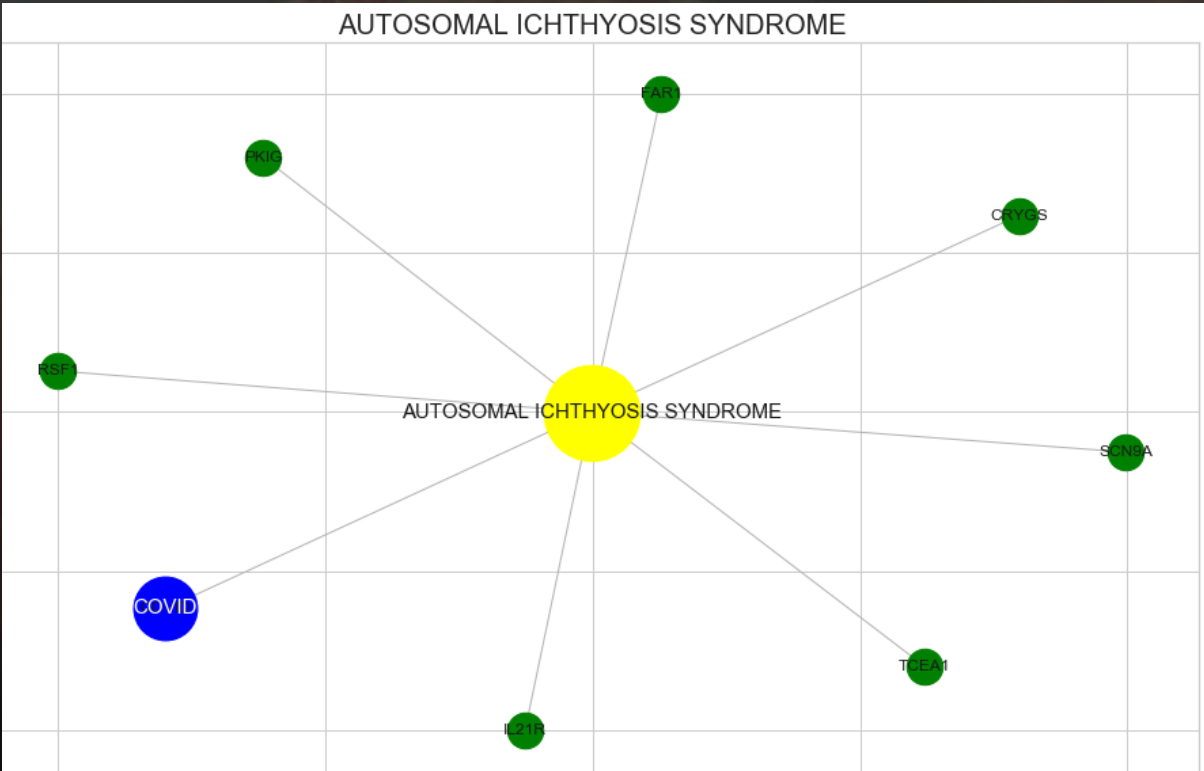
- Nodes: **Diseases**, **Genes**
- Edges: weighted score of gene on disease + disease to disease
- Goal: Central nodes own network understanding + COVID-19 scoring comparison

Results

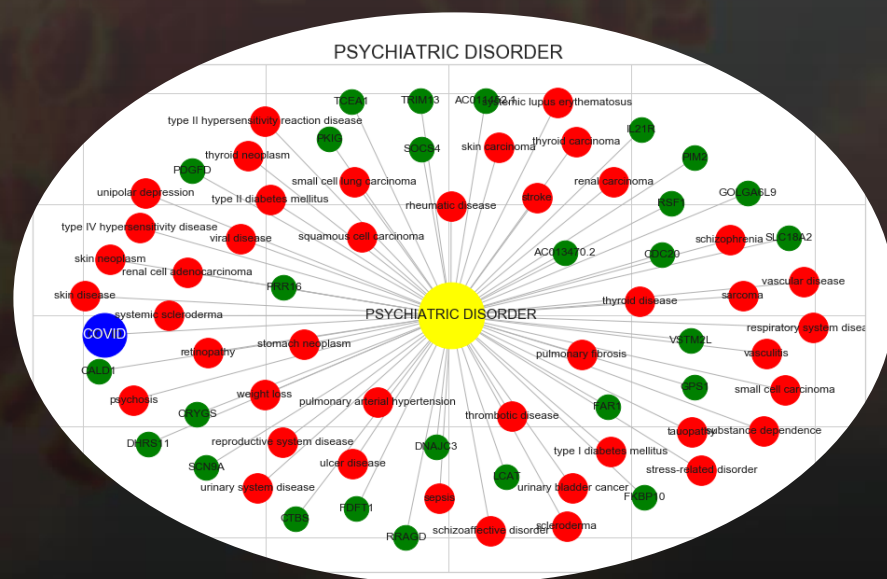
1. Overall Ranking



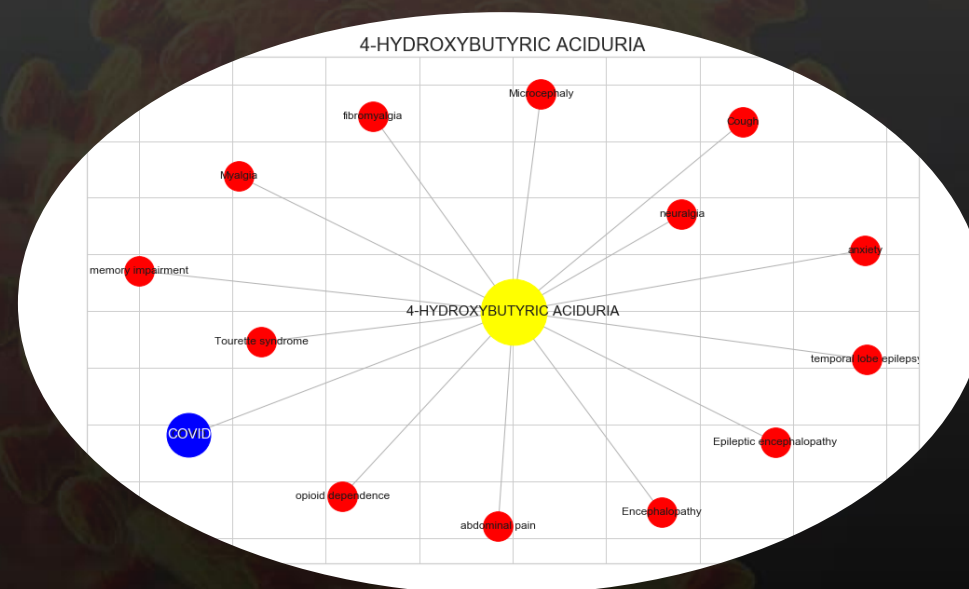
1. Covid Ranking



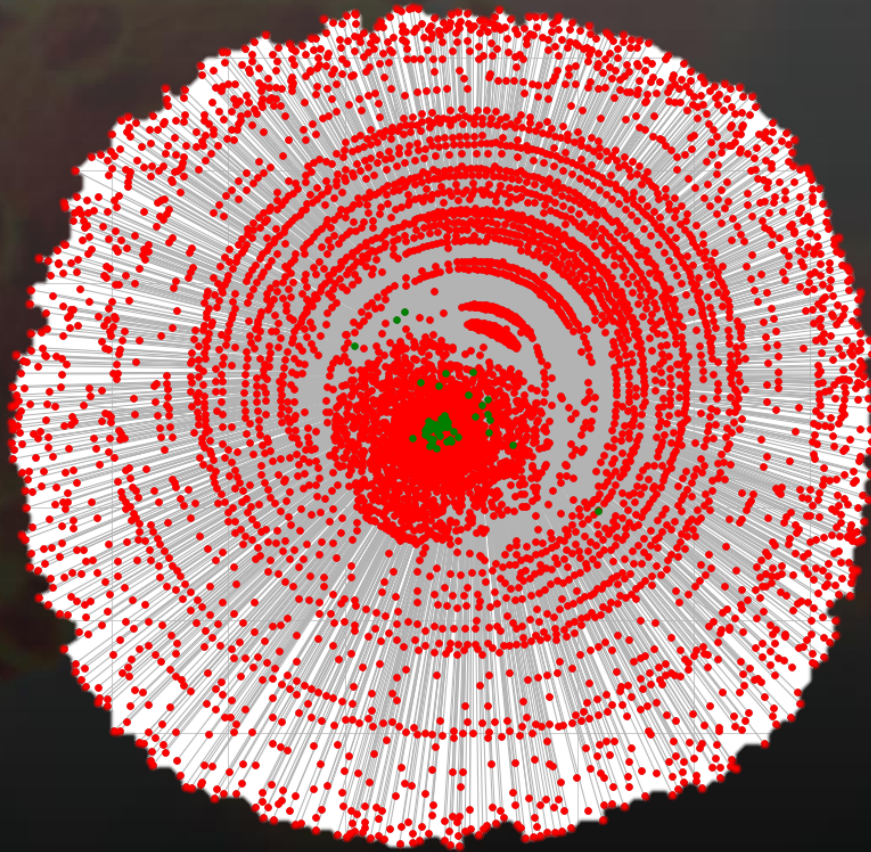
2. Overall Ranking



2. Covid Ranking

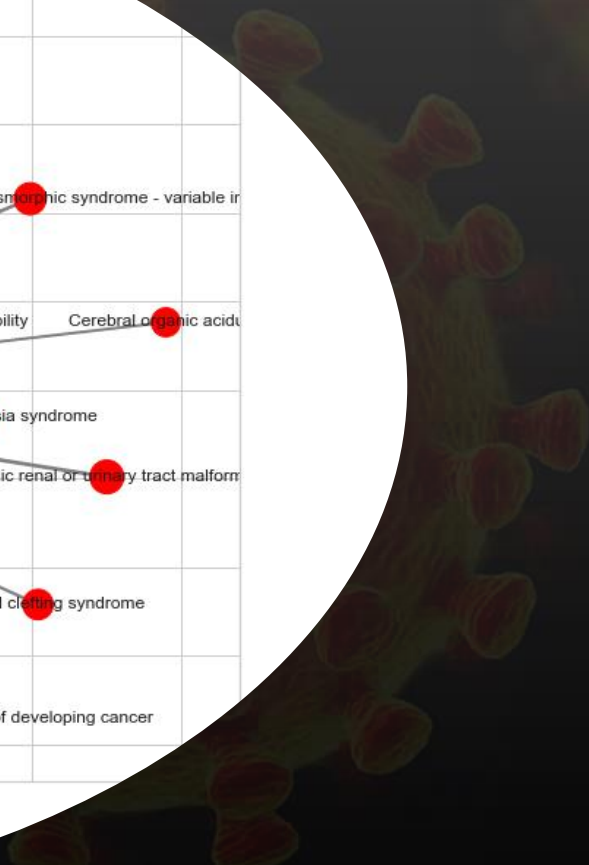


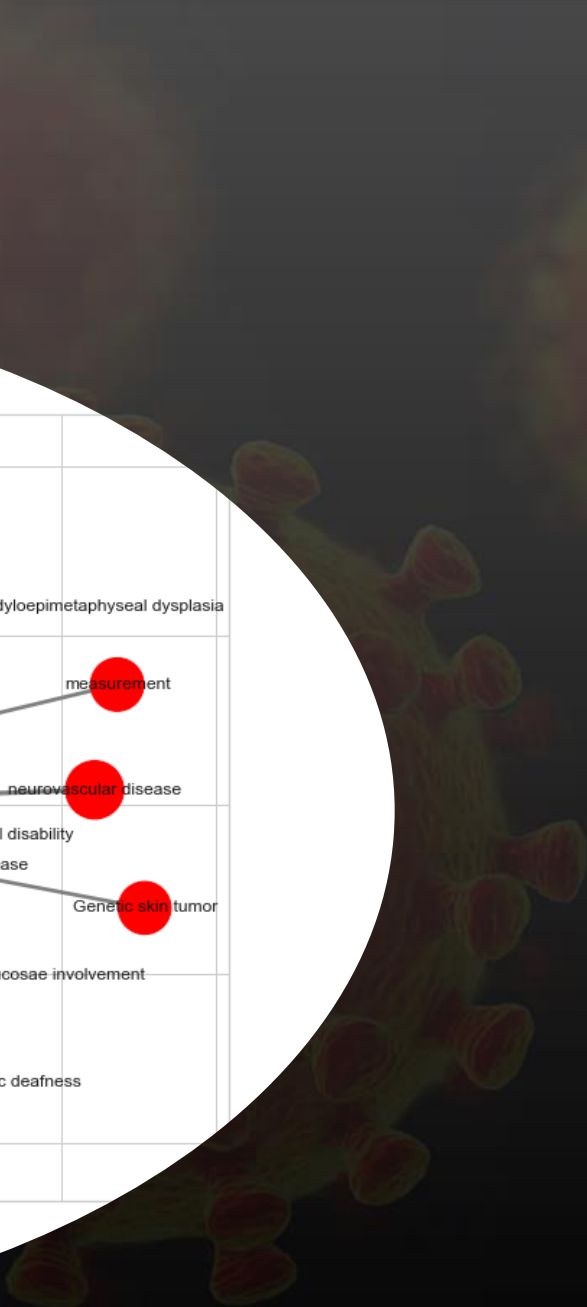
Model



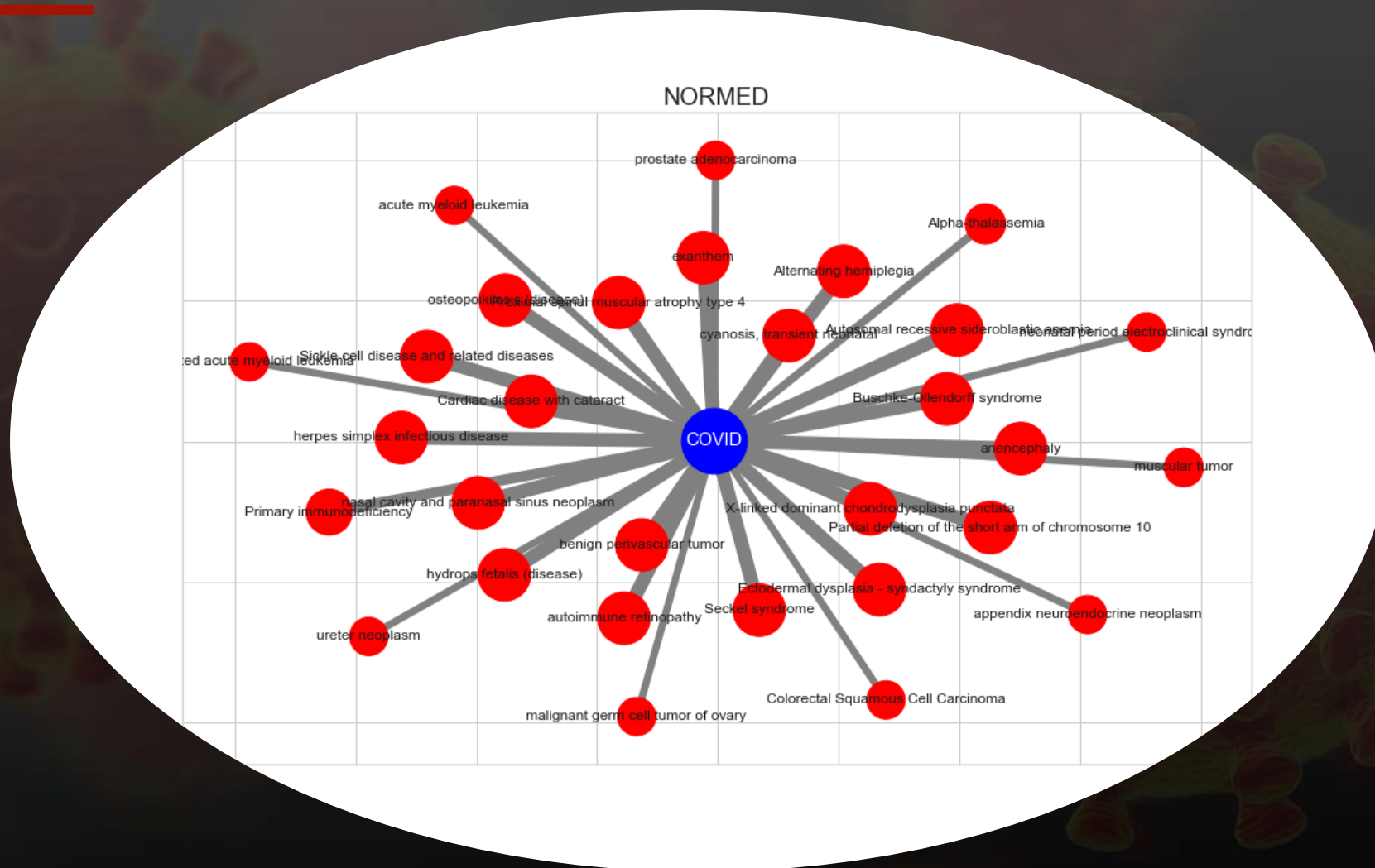
→ Dataset II

- **Nodes: Communités**
- **Edges: Covid score**
- **Goal: Community influence**





Results





IV. Random Forest Model

Pre-processing



→ Data set

- Self generated
- No pre-processing needed
- 70/30 split for train/test data

Model



→ Grid search model

- **Parameters**
 - Maximum depth
 - Number of estimators
 - Minimum samples per leaf
- **385 fits**

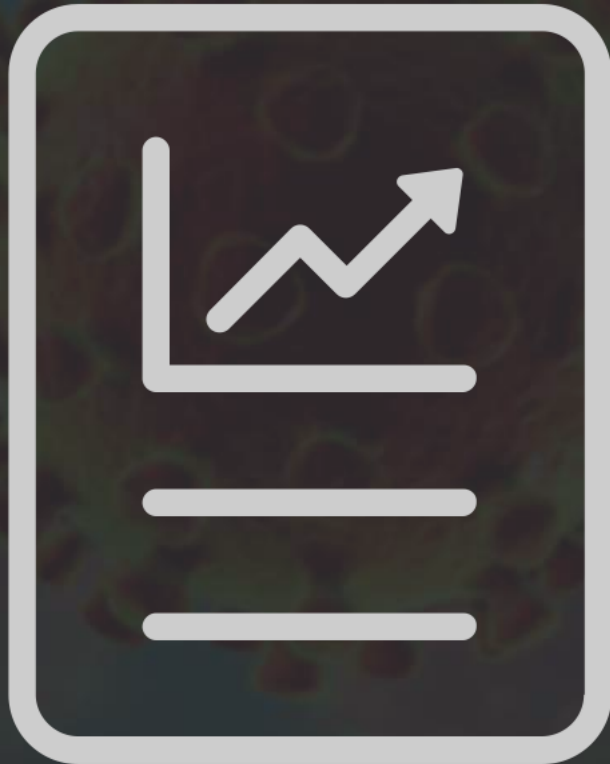
Model



→ Random search model

- **Parameters**
 - Number of estimators
 - Maximum depth
 - Minimum sample split
 - Minimum samples per leaf
 - Number of features considered when looking for the best split
- **25 fits due to simulation length**

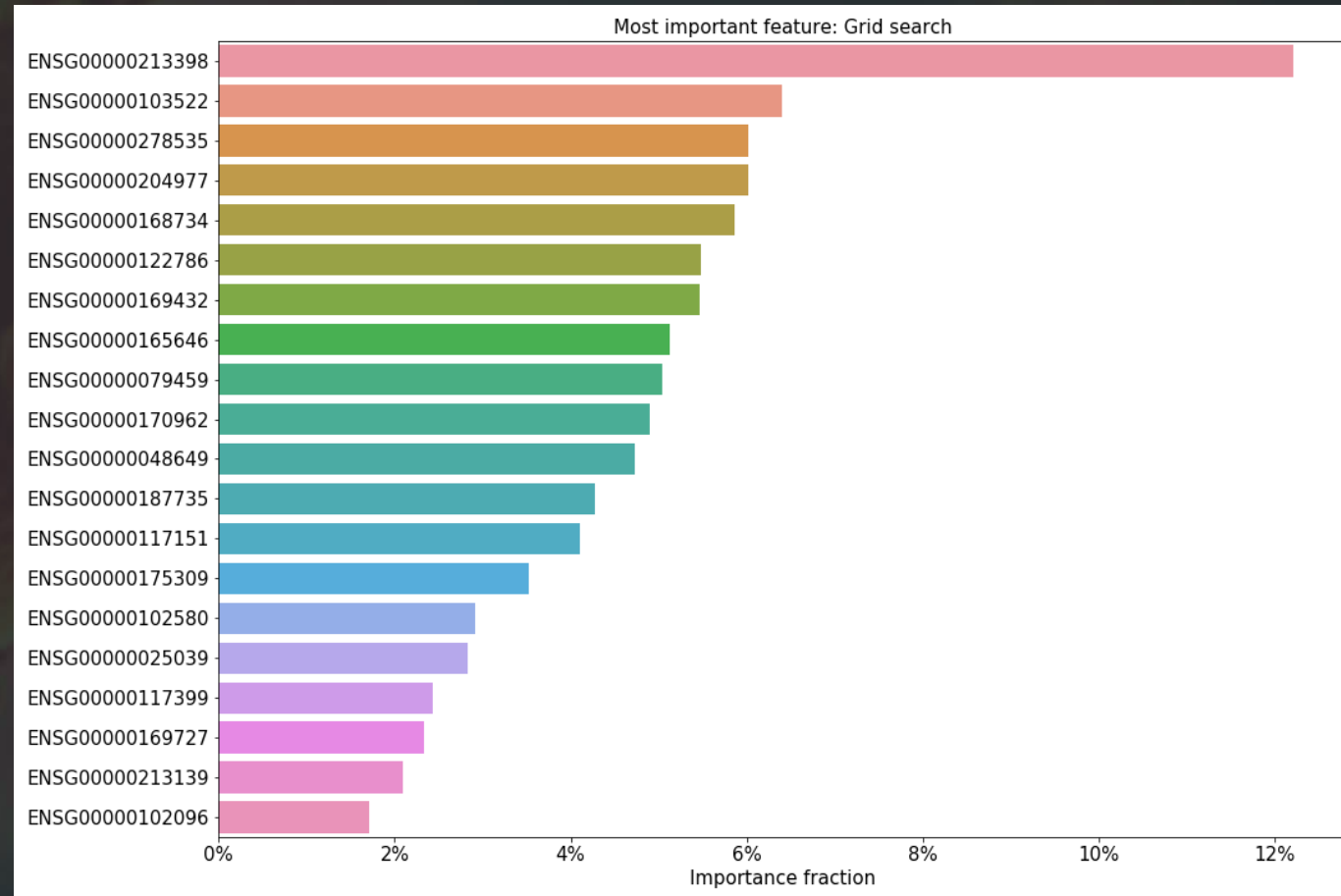
Results



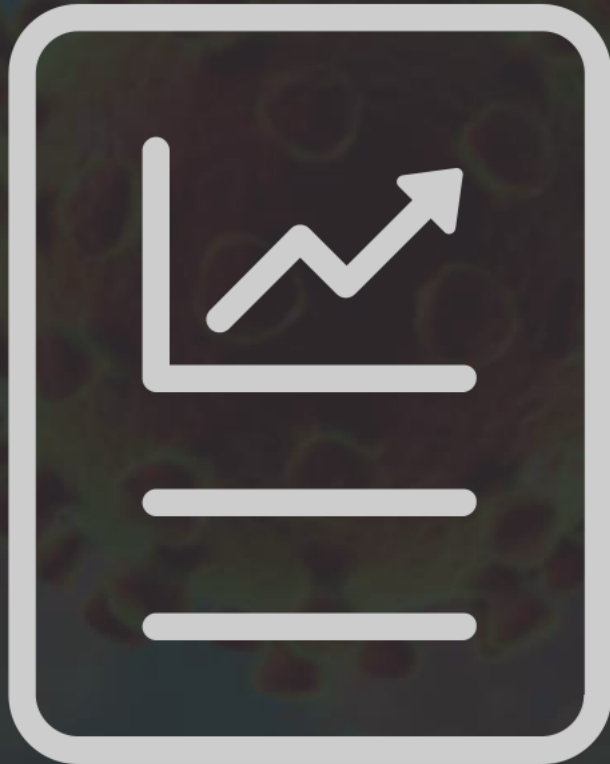
→ Grid search

- 22 targets out of 100 compose 95% of the importance
- Lecithin-cholesterol acyltransferase at 12%
- Next 18 between 2 and 8%
- $r^2 = 0.001369837276703811$
 - 35% less than baseline

Results



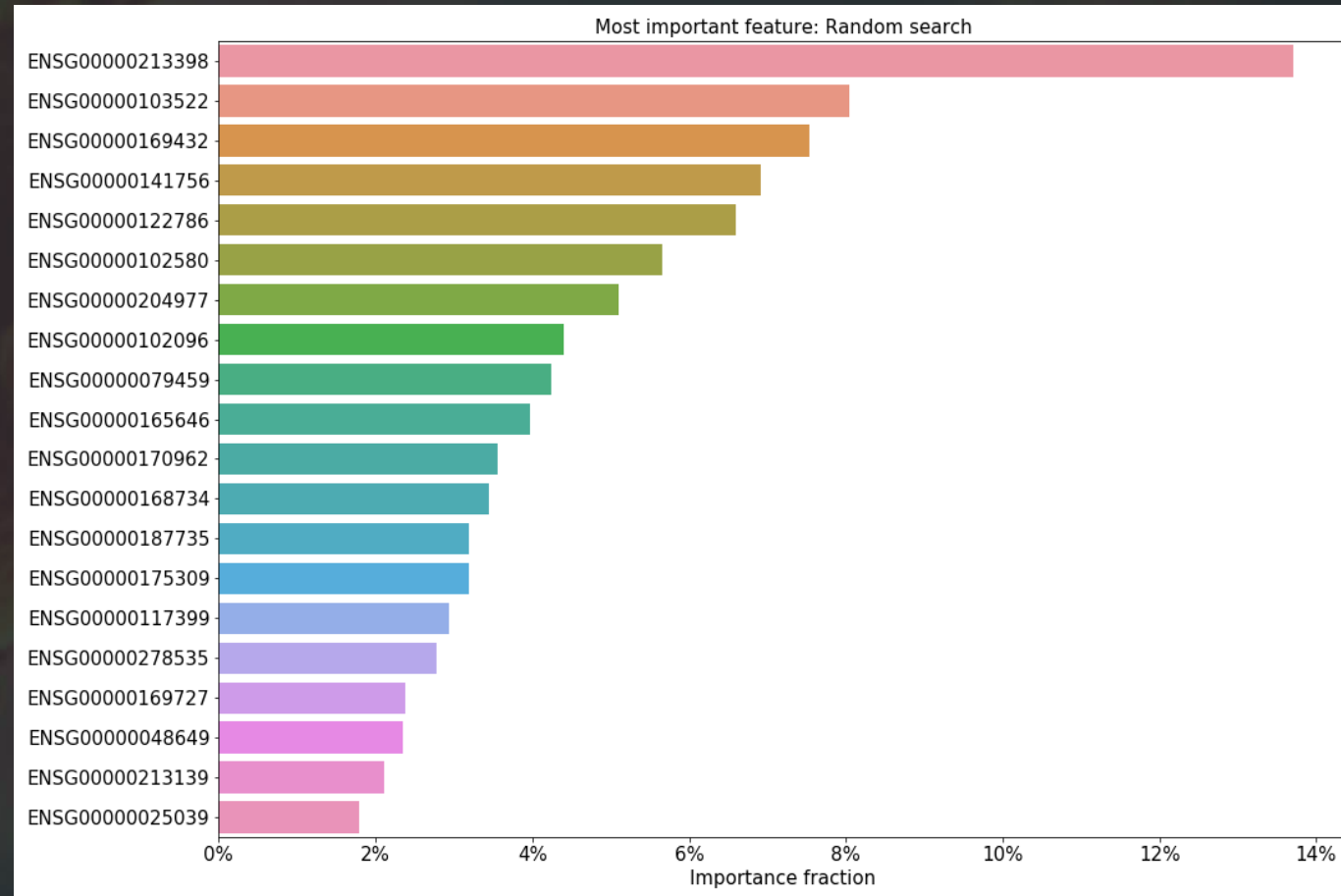
Results



→ Random search

- 21 targets out of 100 compose 95% of the importance
- Lecithin-cholesterol acyltransferase at 14%
- Next 18 between 2 and 8%
- $r^2 = 0.00048797298004477074$
 - 47% less than baseline

Results





V. Interpretations

Benefits for the research on the virus



→ **Network analysis**

- **Disease ranking:**

- Bone disease, psychiatric disorder, biological process, protein measurement and diabetes mellitus

- **Target ranking:**

- LCAT, SCN9A, SLC18A2, PDGFD and CALD1

→ **Random forest regression**

- **Target identification helps reduce research time and efforts**
- **Reduce costs**
- **Speed up the vaccine development**

Improving and reusing the model

→ Improvements and limitations



- Based on the hypothesis that genetics and genomics play a role in the impact on COVID-19
- Association score between diseases and COVID-19 built on a partial data-set
- Main limitation relying in the amount of data used: targets and diseases
- Can be run with a larger amount of data to expect better results



Conclusion

Conclusions

- Closest diseases to COVID-19 identified
- Most important target risk-factors identified
- Reduction of the time spent on vaccine development
- Help governments make proportionate decisions



*Thank you for
your attention*

*R. Chaouche, Y. Martinson,
C. Padovani, J. Triomphe*