# Predicting Missing Links in Infection Networks: Accelerate contact tracing investigations using network theory.

**Pavlos Alexandros Dimitriou[1,2], V. Silvestros[3], C.Pitris[1,2], P. Kolios[1,4]**

[1]University of Cyprus, KIOS Research and Innovation Center of Excellence, Nicosia, Cyprus.

[2]University of Cyprus, Dept. of Electrical & Computer Engineering, Nicosia, Cyprus.

[3]Unit for Surveillance and Control of Communicable Diseases, Ministry of Health, Nicosia, Cyprus.

[4]University of Cyprus, Dept. of Computer Science, Nicosia, Cyprus.

**24-27 September, World Congress of Epidemiology**

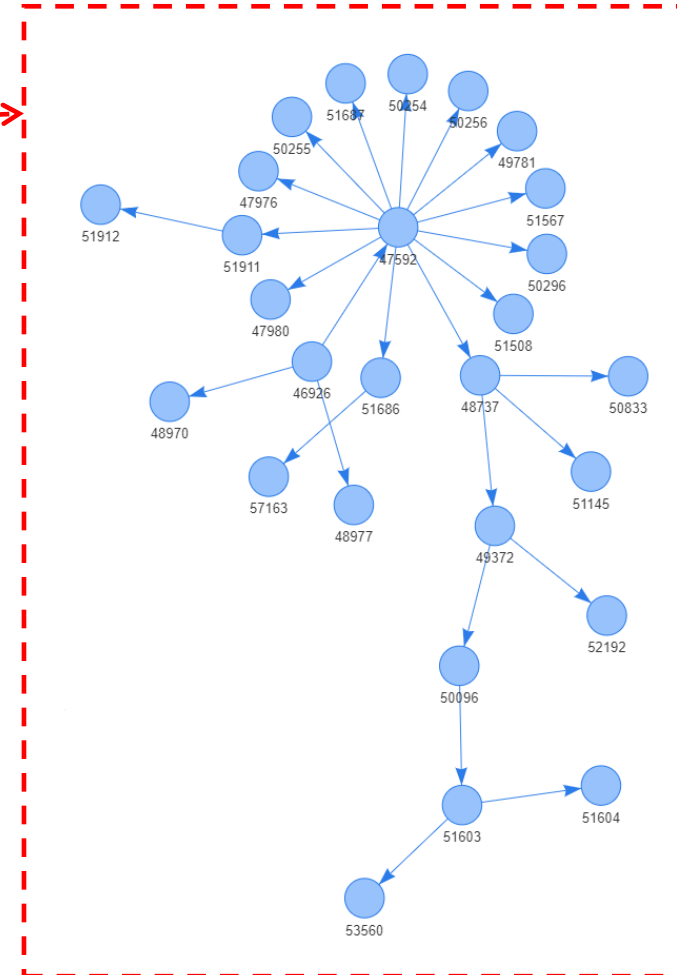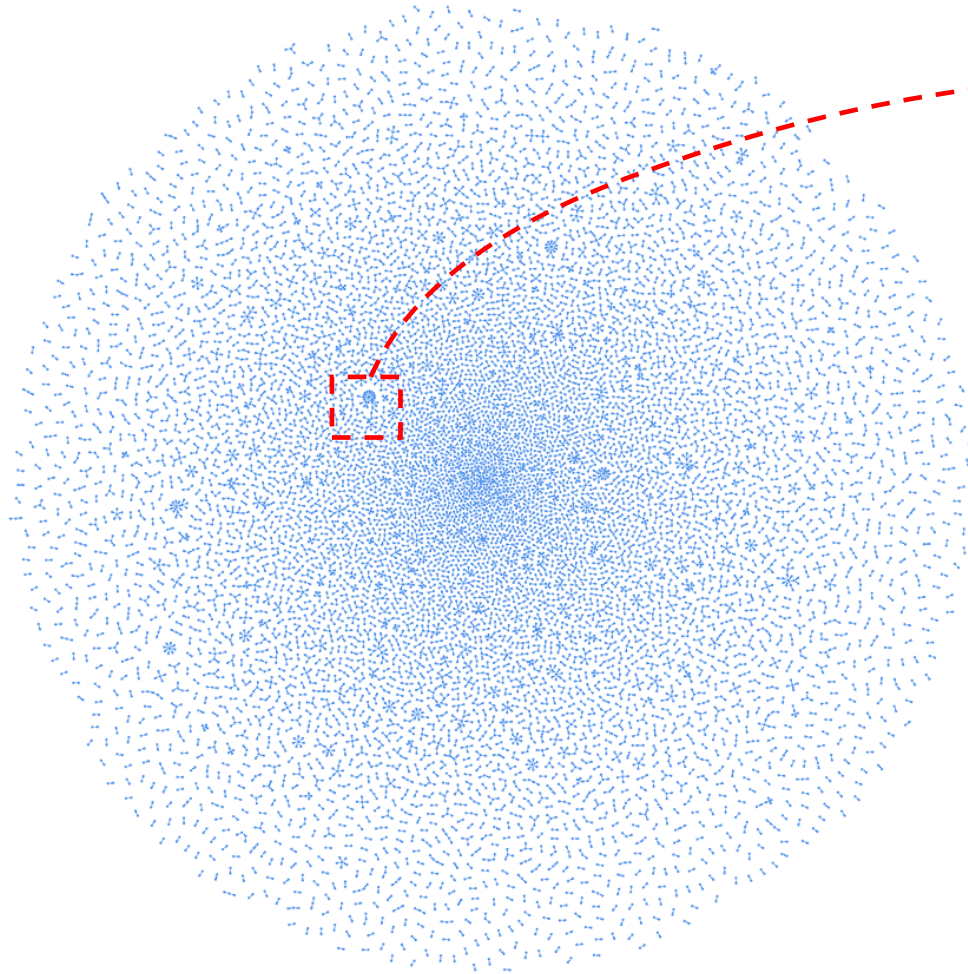# Building networks from contact tracing data

Data store in excel tables



| Case Id | Contact case Id | Gender | District | Age group | NACE | Infection date |
|---|---|---|---|---|---|---|
| 1 | 11 | M | Pafos | 20-29 | P85.4 | 12/1/2021 |
| 2 | 12 | M | Pafos | 3-5 | U6 | 12/1/2021 |
| 3 | nan | M | Pafos | 30-39 | P85.4 | 12/1/2021 |
| 4 | 3 | F | Pafos | 40-49 | P85.4 | 12/1/2021 |
| 5 | 12 | M | Pafos | 3-5 | U6 | 12/1/2021 |
| 6 | 12,13 | F | Pafos | 0-2 | U6 | 12/1/2021 |
| 7 | 12 | F | Pafos | 30-39 | nan | 12/1/2021 |
| 8 | 4 | F | Pafos | 40-49 | nan | 12/1/2021 |
| 9 | 3 | F | Pafos | 60-69 | R65 | 12/1/2021 |
| 10 | 12 | M | Pafos | 20-29 | T97 | 12/1/2021 |
| 11 | 3 | F | Pafos | 50-59 | P85.4 | 12/1/2021 |
| 12 | 3 | F | Pafos | 60-69 | R65 | 12/1/2021 |
| 13 | 3 | M | Pafos | 70-79 | R65 | 12/1/2021 |

Challenging to analyse

WCE WORLD CONGRESS OF EPIDEMIOLOGY 2024

# Building networks from contact tracing data

**Infection network:**



Identify transmission patterns and super-spreaders

NODES ARE THE INFECTED
------ INDIVIDUALS

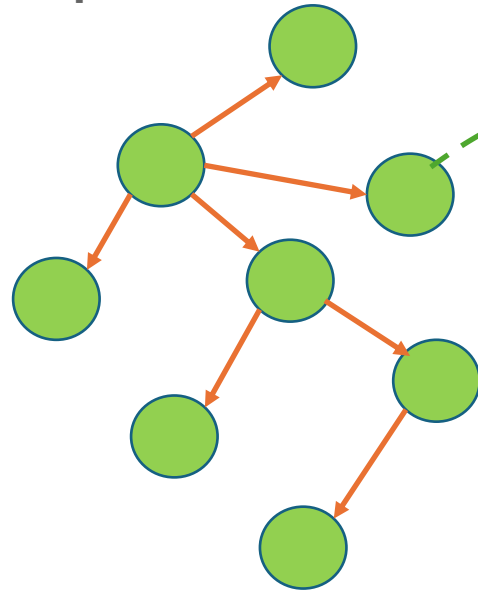------ EDGES ARE THE EPIDEMIOLOGICAL
LINKS BETWEEN THEM

# Infection networks are sparse
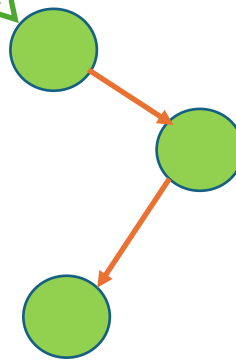
# Can we predict missing links?



**Infection Network**
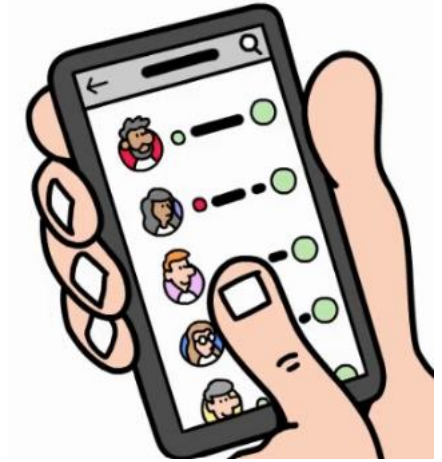
Component 1

Component 2

?

Link prediction

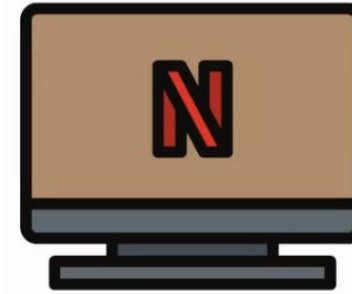RECONSTRUCTION OF THE TRANSMISSION CHAINS

# Link Prediction and Its Application:

- Widely used in different domains
  - Recommender system
    - friendships in social networks
    - e-commerce websites

**WE TREAT LINK PREDICTION AS A SUPERVISED CLASSIFICATION TASK**

Feature vector $(x)$

$$x = [f_{0,0}, \ldots f_{i,j}, \ldots, f_{n,n}]$$

Class label (y)

$$y = [y_{0,0}, \ldots y_{i,j}, \ldots, y_{n,n}]$$

# Creating Feature Vectors:

1. **Node2vec**:

   - Features are calculated solely on network characteristics

**FEATURE VECTORS BASED ON THE CONNECTIONS AND TOPOLOGY OF THE INFECTION NETWORK**

2. **Shallow embeddings with handcrafted features:**

**EPIDEMIOLOGICAL FEATURES + NETWORK STRUCTURE**

- THE TIME DIFFERENCE BETWEEN INFECTIONS
- THE PHYSICAL DISTANCE BETWEEN INDIVIDUALS BASED ON THEIR RESIDENCE
- THE AGE DIFFERENCE BETWEEN INFECTED INDIVIDUALS
- THE OVERLAP OF OCCUPATIONS (NACE CODES) AMONG CHAINS OF INFECTED INDIVIDUALS
- THE OVERLAP OF POSTCODES AMONG INFECTION CHAINS

# Splitting the Dataset

- Sampling positive and negative edges to create training and test set



Full network

Training network

Test network

# Machine Learning: Classification algorithms



GB (Gradient Boosting)

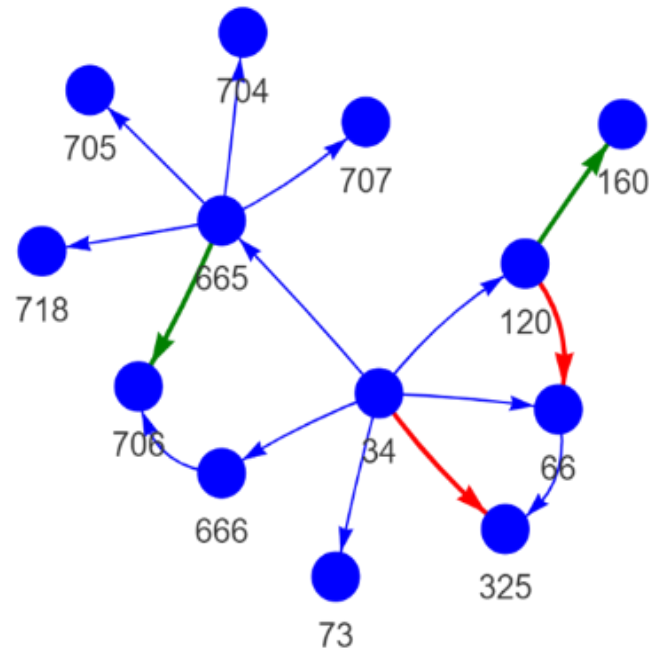LogReg (Logistic Regression)

KNN (k-Nearest Neighbors)

Machine learning classifiers

DT (Decision Tree)

LDA (Linear Discriminant Analysis)

SVM (Support Vector Machine)

RF (Random Forest)

# Performance Metrics: Node2vec Link prediction

| | Wave 1 | | Wave 2 | | Wave 3 | | Wave 4 | |
|---|---|---|---|---|---|---|---|---|
| | AUC | F1-Score | AUC | F1-Score | AUC | F1-Score | AUC | F1-Score |
| RF | 0.54 | 0.49 | 0.54 | 0.46 | 0.57 | 0.53 | 0.67 | 0.65 |
| LDA | 0.50 | 0.51 | 0.48 | 0.47 | 0.55 | 0.57 | 0.65 | 0.71 |
| SVM | 0.54 | 0.53 | 0.50 | 0.51 | 0.61 | 0.64 | 0.70 | 0.75 |
| LOGREG | 0.53 | 0.52 | 0.49 | 0.49 | 0.58 | 0.60 | 0.68 | 0.75 |
| KNN | 0.56 | 0.61 | 0.49 | 0.50 | 0.58 | 0.65 | 0.64 | 0.62 |
| NB | 0.52 | 0.42 | 0.51 | 0.47 | 0.58 | 0.55 | 0.62 | 0.71 |
| GB | 0.57 | 0.57 | 0.52 | 0.52 | 0.58 | 0.61 | 0.69 | 0.71 |
| DT | 0.53 | 0.52 | 0.50 | 0.51 | 0.54 | 0.54 | 0.62 | 0.62 |

# Performance Metrics:
# Shallow embeddings with handcrafted features

|  | Wave 1 | | Wave 2 | | Wave 3 | | Wave 4 | |
|---|---|---|---|---|---|---|---|---|
|  | AUC | F1-Score | AUC | F1-Score | AUC | F1-Score | AUC | F1-Score |
| RF | 0.87 | 0.85 | 0.82 | 0.82 | 0.93 | 0.93 | 0.94 | 0.94 |
| LDA | 0.85 | 0.86 | **0.82** | **0.83** | 0.90 | 0.90 | 0.91 | 0.91 |
| SVM | 0.56 | 0.57 | 0.55 | 0.61 | 0.54 | 0.59 | 0.54 | 0.53 |
| LOGREG | 0.54 | 0.17 | 0.53 | 0.24 | 0.53 | 0.10 | 0.53 | 0.11 |
| KNN | 0.78 | 0.75 | 0.74 | 0.75 | 0.67 | 0.67 | 0.73 | 0.73 |
| NB | 0.56 | 0.57 | 0.55 | 0.58 | 0.54 | 0.61 | 0.53 | 0.63 |
| GB | **0.90** | **0.88** | 0.82 | 0.82 | **0.94** | **0.94** | **0.95** | **0.95** |
| DT | 0.81 | 0.85 | 0.79 | 0.80 | 0.91 | 0.90 | 0.92 | 0.92 |

WCE
WORLD CONGRESS OF EPIDEMIOLOGY 2024

Koios

# Link prediction: Complete transmissions chains



Chain 1    Chain 2    Chain 3

PREDICTING MISSING LINKS IN INFECTION NETWORKS: ACCELERATE CONTACT TRACING INVESTIGATIONS USING NETWORK THEORY.

24-27 September, World Congress of Epidemiology

COVID-19

University of Cyprus

Nodes are the infected individuals

Edges are the epidemiological links between them

WHILE ALSO HELPING EPIDEMIOLOGISTS ACCELERATE CONTACT TRACING INVESTIGATIONS

WHICH CREATES GAPS IN UNDERSTANDING THE FULL TRANSMISSION PICTURE

RECONSTRUCTION OF THE TRANSMISSION CHAINS

LINK PREDICTION AND ITS APPLICATION

SUCH AS RECOMMENDER SYSTEMS

SPLITTING THE DATASET

Once the feature vectors are created

We start by sampling positive and negative edges

FULL NETWORK

THE FIRST APPROACH UTILIZES NODE2VEC

FEATURE VECTORS BASED ON THE CONNECTIONS AND TOPOLOGY OF THE INFECTION NETWORK

IN THE SECOND APPROACH

EPIDEMIOLOGICAL FEATURES + NETWORK STRUCTURE

- The time difference between infections
- The physical distance between individuals based on their residence
- The age difference between infected individuals
- The overlap of occupations once codes among chains of infected individuals
- The overlap of postcodes among infection chains

CREATING FEATURE VECTORS FOR LINK PREDICTION

2

WHAT WE CONSIDER TO BE "SIMILAR" IN THE INFECTION NETWORK

USING THESE FEATURE VECTORS AND LABELS

WE TREAT LINK PREDICTION AS A SUPERVISED CLASSIFICATION TASK

(Y)

(X)

LIKE AGE, LOCATION AND OCCUPATION

POSITIVE EDGES

( TRAINING SET )

NEGATIVE EDGES

POSITIVE

NEGATIVE EDGES

WHILE RESPECTING THE NETWORK STRUCTURE

( TEST SET )

MACHINE LEARNING CLASSIFICATION ALGORITHMS

NAIVE BAYES / LOGISTIC REGRESSION / Types of Classification Algorithms / K-NEAREST NEIGHBOR / DECISION TREE / SUPPORT VECTOR MACHINE (SVM) / RANDOM FOREST

PERFORMANCE ON PREDICTING MISSING LINKS

APPROACH I:

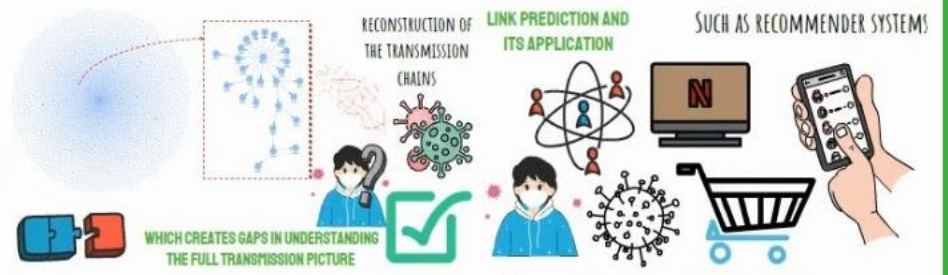| | Wave 1 | | Wave 2 | | Wave 3 | | Wave 4 | |
|---|---|---|---|---|---|---|---|---|
| | AUC | F1-Score | AUC | F1-Score | AUC | F1-Score | AUC | F1-Score |
| RF | 0.54 | 0.49 | 0.54 | 0.46 | 0.57 | 0.53 | 0.67 | 0.65 |
| LDA | 0.50 | 0.51 | 0.48 | 0.47 | 0.55 | 0.57 | 0.65 | 0.71 |
| SVM | 0.54 | 0.53 | 0.50 | 0.51 | 0.61 | 0.64 | 0.70 | 0.75 |
| LOGREG | 0.53 | 0.52 | 0.49 | 0.49 | 0.58 | 0.60 | 0.66 | 0.75 |
| KNN | 0.56 | 0.61 | 0.49 | 0.50 | 0.58 | 0.65 | 0.64 | 0.62 |
| NB | 0.52 | 0.42 | 0.51 | 0.47 | 0.58 | 0.55 | 0.62 | 0.71 |
| GB | 0.57 | 0.57 | 0.52 | 0.52 | 0.58 | 0.61 | 0.69 | 0.71 |
| DT | 0.53 | 0.52 | 0.50 | 0.51 | 0.54 | 0.54 | 0.62 | 0.62 |

WITH THE NODE2VEC ALGORITHM

IDENTIFYING MISSING LINKS WITHIN THE INFECTION NETWORKS

APPROACH 2:

| | Wave 1 | | Wave 2 | | Wave 3 | | Wave 4 | |
|---|---|---|---|---|---|---|---|---|
| | AUC | F1-Score | AUC | F1-Score | AUC | F1-Score | AUC | F1-Score |
| RF | 0.97 | 0.85 | 0.82 | 0.82 | 0.93 | 0.93 | 0.90 | 0.94 |
| LDA | 0.85 | 0.86 | 0.82 | 0.83 | 0.90 | 0.90 | 0.90 | 0.93 |
| SVM | 0.90 | 0.90 | 0.83 | 0.83 | 0.91 | 0.90 | 0.91 | 0.94 |
| LOGREG | 0.54 | 0.17 | 0.53 | 0.34 | 0.53 | 0.10 | 0.10 | 0.03 |
| KNN | 0.79 | 0.75 | 0.74 | 0.75 | 0.67 | 0.67 | 0.73 | 0.79 |
| NB | 0.78 | 0.57 | 0.55 | 0.34 | 0.58 | 0.54 | 0.54 | 0.64 |
| DT | 0.81 | 0.81 | 0.77 | 0.76 | 0.82 | 0.84 | 0.85 | 0.88 |

ACHIEVING AN AUC AND F1 SCORE OF 95%

OCCUPATION OVERLAPS, AND POSTCODE OVERLAPS

EPIDEMIOLOGICAL & NETWORK FEATURES

Infection Network

RECONSTRUCT TRANSMISSION CHAINS

IDENTIFYING MISSING LINKS AND ACCELERATING CONTACT TRACING INVESTIGATIONS

THANK YOU!

KIOS

Alexandros Dimitriou, Cyprus

WCE WORLD CONGRESS OF EPIDEMIOLOGY 2024