

# Accelerating Epidemiological Investigation Analysis by using NLP and Knowledge Reasoning: A Case Study on COVID-19.

## Abstract:

Triple-based Multi-Task Neural Network (MTT-NN) jointly reasoning epidemiological entities and relations from case reports, an epidemiological knowledge graph, and its corresponding inference engine are built to uncover the infection modes, sources and pathways.

## OpenAI

### Covid-19

travel trajectories before the disease.

including demographic data & clinical manifestations of confirmed cases, travel trajectories before the disease.

tracing and analyzing infection source and pathway for each patients, then discussing & generating epidemiological insights with visualization.

## Related work

① Current methods for epidemiological investigation analysis.

② state of art: medical information extraction.

③ applying knowledge graph in epidemiology research

II: automatic extraction of concepts, entities, events, as well as relations & associated attributes from free texts.

general work in this paper:

design & implements a neural network based on BERT

④ an inference engine is constructed to discover infection modes, pathways and networks based on KG.

## System Description & epidemiological Ktir

An epidemiological case report for an individual patient usually includes 4 types of information:

① disease graphs (patient ID, name, gender, residence place, epidemiological contact history).

② recent active log (date, location, vehicle, manifesting persons).

③ disease related event: diagnosis and treatment information (e.g. hospitalization date, treatment date).

④ the information of close contacts (e.g. family members).

Epidemiological information extraction via a novel deep neural network specifically designed for case reports.

Nested tuples: an epidemiological case report can be structured

1 patient info tuple: (patient ID, age, gender, residence place, epidemiological contact history).

1 more event tuple: (time, location, vehicle, manifesting persons)

1 more social relation tuple: (id, name, relationship).

MTT-NN is specifically designed with a hierarchical nested tuple labeling structure to extract corresponding nested tuples and their corresponding attributes in an end-to-end way.

ensuring → tuple beyond labeling → position embedding in token → tuple attribute labeling.

Infection Network Inference & Infection Source Identification

structured epidemiological information are organized as knowledge graph, we use epidemiological inference algorithm to discover the implicit infection relations among confirmed patients and find out all possible infection sources.

## Patient Description

$V = \{1, \dots, N\}$ : N individual patients.

$G_{\text{inf}} = (V, E_{\text{inf}})$ : directed graph

$E_{\text{inf}} = \{(i,j)\}$ : directed edge, infection relation.

$G_{\text{contact}} = (V, E_{\text{c}})$ : undirected graph

$E_{\text{c}} = \{(i,j)\}$ : undirected edges, contact relation ... with a independent probability p.

Contract: the occurrence of a physical association of two individuals that could be sufficient for disease transmission.

though not actual contacts between infected & susceptible individuals are guaranteed to result in transmission.

\* the infection network is a subgraph of the contact network, shown in Figure 4.

epidemiological investigation case reports  $\Rightarrow$  infection network from contact network, with the history observation O(t).

Our goal is to infer the infection source & transmission path (pathway along which the epidemic spreads) from the contact network, with the history observation O(t).

the inferred sources, transmission paths with nodes are the directed subgraph of the undirected contact network. one out its directed infection network.  $G_{\text{infection}} = (V, E_{\text{infection}})$



$$S, g = \arg \max_{S \subseteq V, g \subseteq G_{\text{infection}}} P(O|S, g)$$

s: infection source

g: inferred infection network

$P(O|S, g)$ : probability of observation O of all node states after the disease transmits in contact network.

4 steps pipeline:

- ① construct the contact network from structured case reports stored in KGs [with the state of each node at every point of time t].
- ↑ the individuals in the epidemiological KGs are the nodes in the contact network.
- contact relation: ② causal relationship between those patients who have been interacting multiple times with the same time duration & location.
- ↑ the contact events associated with more patients who have been interacting multiple times with the same time duration & location.

Note: if graph not be connected graph.

③ Only look loops.

2 spot nodes. the observed spots in the same time t or the same duration of events can be a transmitting medium of the disease.

add spot nodes to the node set V.

$E_{\text{infection}}$ ,  $E_{\text{contact}}$

$\xrightarrow{x} \xrightarrow{y} \xrightarrow{z}$

3<sup>rd</sup> transmission time  $t_3$  for each node. After this time, the node state change from susceptible to infected.

Note: since the source node is not defined yet, we use absolute time instead of relative time from source node.

↑ individual node: each date of reported case as the time of each node.

↑ spot node: time of associated event.

④ Divide the contact network into connected subgraphs to reduce computation complexity.



⑤ For each subgraph, identify the possible infection sources and transmission paths to get infection subgraph.

⇒ identify all combined alternative of i & j

⇒  $S, g$  which maximize  $P(O|S, g)$

Ex: 1  $x \rightarrow y$

    y as infection relation

    y  $\rightarrow z$

Rule: For any  $x \in V$ , set  $P(O|x) = p$  [ $p$ : avg. possibility value of all nodes].

    if x has contact history with confirmed cases(s) before time t, set  $P(O|x) = m$ .

    if x has contact history with suspected cases(s) or possibly came from epidemic areas

    ... in addition,  $m > w_1 > w_2 > 1$

If  $x$  has contact history with suspected cases or patients) come from epidemic activity or known in epidemic activity, weight  $w_{x,y} = 1$ ,  $\dots$   
 $w_1, w_2$  are the empirical constant value.  
 Rule2: For any  $x, y \in V$ , contact edge  $E \subseteq E$ ,  $t_x$ : onset time of  $x$ ,  $t_y$ : onset time of  $y$ .  
 If the directed infection edge  $l_{xy}$ :  $x \rightarrow y$   
 If  $y \neq \infty$   
 $P(O|l_{xy}) > P(O|l_{yz})$ .  
 Rule3: For any  $x, y \in V$ ,  $E \subseteq E$ ,  $t_x$ : onset time of  $x$ ,  $t_y$ : onset time of  $y$ .  
 If  $t_x = t_y$  &  $P(O|l_{xy}) > P(O|l_{yz})$ :  $P(O|l_{xy}) > P(O|l_{yz})$ .  
 Rule4: For any  $x, y \in V$ ,  $E \subseteq E$ , the onset time  $t_x, t_y, t_z$ .  
 If  $t_x < t_y < t_z$ :  $P(O|l_{xy}) > P(O|l_{yz})$   
 $P(O|l_{xy}) = \sum_{E \subseteq E} P(O|l_{xy})$   
~~由这个点传播到另一个点概率最大和边缘接触入门槛 [一个病人只能被感染一次]~~

~~arg max  $P(O|L)$   $L$ : set of directed edge of directed infection subgraph.  $\Rightarrow$  source [the node which in-degree=0].~~  
 ~~$V$ : nodes~~

~~② Validate the inferred graphs. --- 究竟正~~  
~~If conflicts happen, we should check the results. [Conflicts between original epidemiological investigation data & inference result] Example caused by wrong reported date etc.~~  
~~manually fix~~

Dong Liu  
2021.11.26