

LDA2Net: Digging under the surface of COVID-19 topics in literature

Topic 5 companion sheet

G. Minello

C.R.M.A. Santagiustina

M. Warglien

This file contains the following supplementary information for Topic 5 of the manuscript “*LDA2Net*: Digging under the surface of COVID-19 topics in scientific literature”:

- Human label and automatic n-gram label proposals (Table 1)
- Summary measures (Table 2)
- Network of top 25 bigrams (Figure 1)
- Wordclouds of top 25 words by node relevance measure (Figure 2)
- Wordclouds of top 25 bigrams by edge relevance measure (Figure 3)
- Filtered (0.99 percentile) topic network (Figure 4)

Table 1: Human and automatic label proposals. Automatic label candidate for largest word community of the topic. In parenthesis: absolute frequency of the walk out of a sample of size 1000.

Human label	2-gram label	3-gram label	4-gram label
epidemic models	mathematical->model (10.3%)	epidemic->used->model (6.7%)	reproduction->number->R0->model (2.9%)

Here follows the set of topic-specific measures that have been used to classify the topic and to analyse its structural properties (see manuscript for details):

Table 2: Summary measures

	JSD	Mean propensity	Variance propensity	Modularity	Barrat Clustering Coeff.
value	0.637157	0.010230	0.000681	0.000000	0.599124
rank	52	116	105	27	92

Based on the aforementioned measures, Topic 5 has been classified as a SPECIALIZED topic.

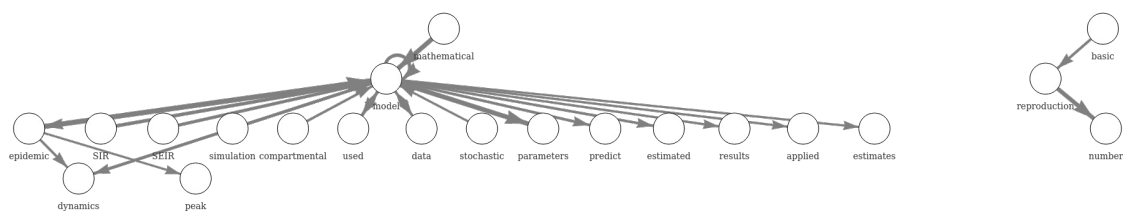


Figure 1: Network of top 25 bigrams (i.e., edges) by weight.

A word cloud of terms related to epidemiology and statistics. The words are arranged in a circular pattern, with 'model!' being the largest and most central word. Other prominent words include 'epidemic', 'modelling', 'estimated', 'rate', 'time', 'data', 'parameters', 'peak', 'simulation', 'scenarios', 'population', 'used', 'dynamics', 'reproduction', 'mathematical', 'estimates', 'spread', 'estimation', 'snow', 'modeling', 'results', 'number', and 'simulations'. The words are in various shades of grey and black, with different font sizes and orientations.

[illegible]

model
parameters
epidemic
reproduction
different
show
used
basic
number
obtained
results
data
estimates
scenarios
applied
predicted
spread
simulations
population
peak
simulation
dynamics
R0
time
estimated

Out-degree

Betweenness

PageRank

Figure 2: Top 25 unigrams (i.e., nodes) by measure.

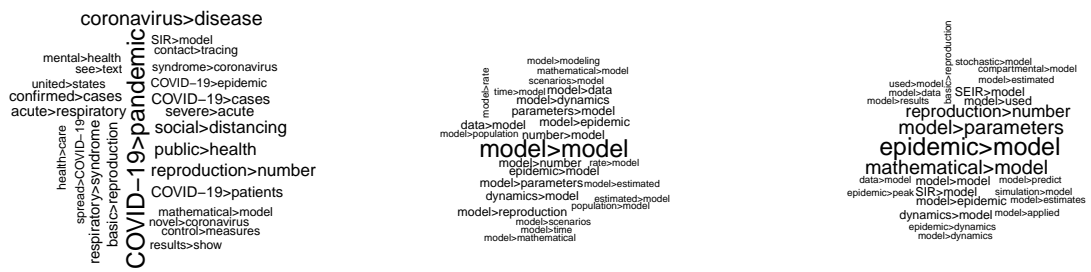


Figure 3: Top 25 bigrams (i.e., edges) by measure.

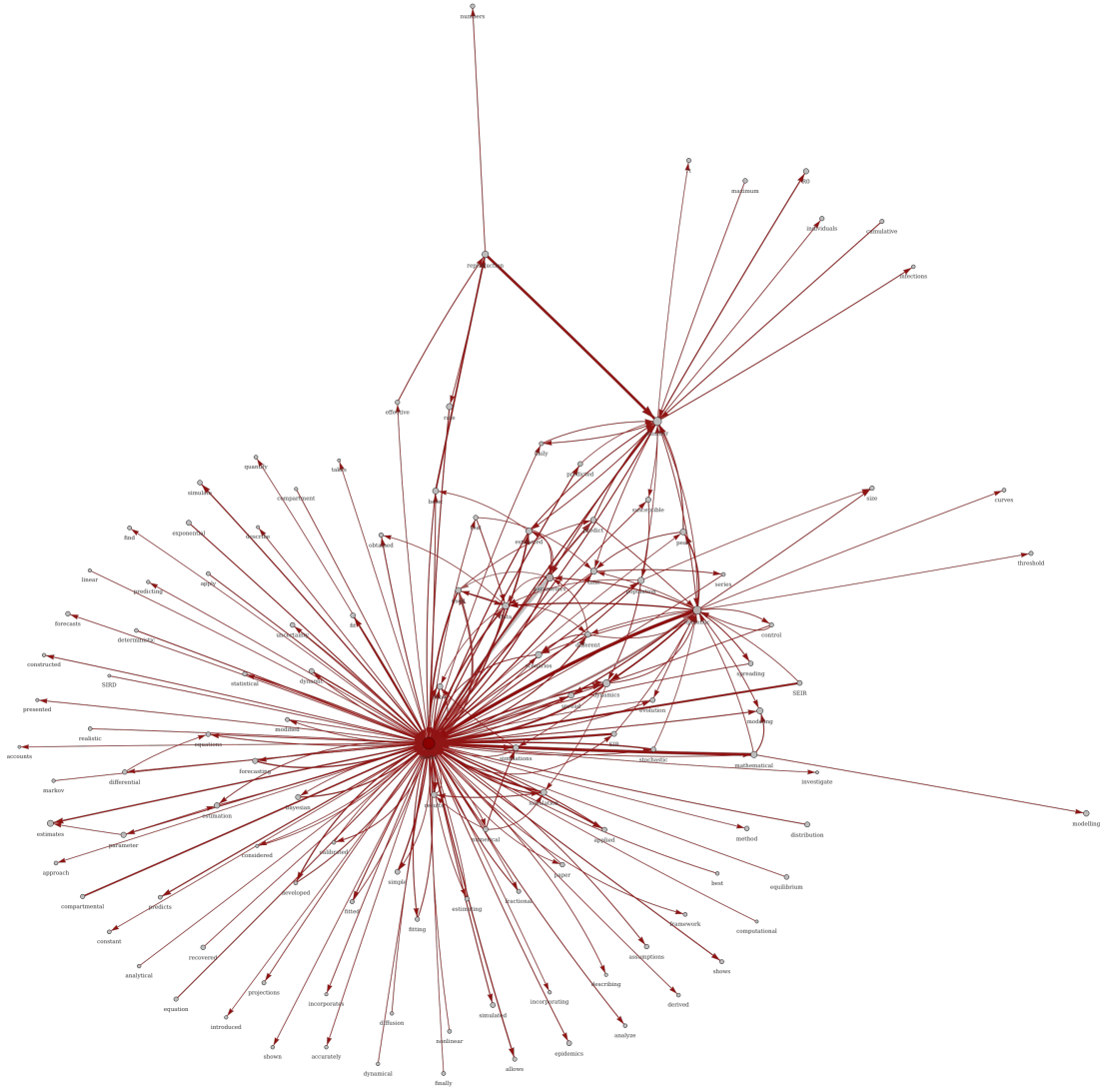


Figure 4: Filtered topic network (by weight). Layout based on Fruchterman-Reingold algorithm. Node size is proportional to topic-specific word probability provided by LDA. Edge width is proportional to topic-specific bigram weight provided by LDA2Net method. Node and edge color represent their betweenness centrality. Isolated nodes have been removed after filtration.