

# AUTO-TUNE: FINDING AN OPTIMAL DISTANCE THRESHOLD FOR INFERRING HIV TRANSMISSION CLUSTERS

Steven Weaver<sup>1\*</sup>, Vanessa Davila-Conn<sup>2</sup>, Joel Wertheim<sup>3</sup>, and Sergei L. Kosakovsky Pond<sup>1</sup>

<sup>1</sup> Center for Viral Evolution, Temple University, Philadelphia, PA, USA

<sup>2</sup> Center for Research in Infectious Diseases, National Institute of Respiratory Diseases, Mexico City, Mexico

<sup>3</sup> Department of Medicine, University of California, San Diego, CA

Correspondence\*:

Steven Weaver

sweaver@temple.edu

## 2 ABSTRACT

3 Choosing an appropriate distance threshold is an important part of inferring a transmission network  
4 to determine the relative growth of clusters within a localized epidemic. This distance threshold  
5 determines how close two consensus sequences must be in order for a link to be created between  
6 them in the network. Using a distance threshold that is too high can result in a network with many  
7 unnecessary links, making it difficult to interpret and analyze. On the other hand, using a distance  
8 threshold that is too low can result in a network with too few links, which may not capture key  
9 insights into rapidly growing clusters among patients with shared attributes that could benefit from  
10 public health intervention measures.

11 Here, we present a heuristic scoring approach for tuning a distance threshold by associating each  
12 tested threshold against the maximal number of clusters created across all thresholds and the  
13 difference between the ratio ( $R_{12}$ ) of the largest cluster in the network to the second largest cluster  
14 at each iteration. The number of clusters is normalized between  $[0, 1]$  then gated via a Gompertz  
15 function transform. Meanwhile, the distribution of all  $R_{12}$  ratios are converted to  $Z$  scores, and  
16 normalized relative to the largest positive  $Z$  score across all candidate distances. The priority score  
17 is the sum of aforementioned two components.

18 Published research using the HIV-TRACE software package frequently use the default threshold  
19 of 1.5% for HIV pol gene sequences. We apply our scoring heuristic to outbreaks with different  
20 characteristics, such as regional or temporal variability, and demonstrate the utility of using the  
21 scoring mechanism's suggested distance threshold to identify clusters exhibiting risk factors that  
22 would have otherwise been more difficult to identify. For example, while we found that a 1.5%  
23 distance threshold is typical for US-like epidemics, recent outbreaks like the CRF07\_BC subtype  
24 among men who have sex with men (MSM) in China has been found to have a lower optimal  
25 threshold of 0.5% to better capture the transition from injected drug use (IDU) to MSM as the  
26 primary risk factor. Alternatively, in communities surrounding Lake Victoria, where there has been  
27 sustained transmission for several years, we found that a larger distance threshold is suitable to

capture a more risk factor diverse populace with sparse sampling over a longer period of time. Such identification may allow for more informed intervention action by respective public health officials.

Keywords: molecular epidemiology, HIV, network, transmission cluster, surveillance

## 1 INTRODUCTION

Choosing an appropriate distance threshold is an important part of using a transmission network to track the spread of a contagious disease. This distance threshold determines how close two individuals must be in order for a link to be created between them in the network.

Using a distance threshold that is too small can result in a network with many unnecessary links, making it difficult to interpret and analyze. On the other hand, using a distance threshold that is too large can result in a network with too few links, making it difficult to accurately track the spread of the disease.

To ensure that the transmission network is useful and informative, it is important to carefully consider the appropriate distance threshold. This may vary depending on the specific disease and the context in which it is spreading. For example, a highly contagious respiratory illness may require a smaller distance threshold than a less contagious illness that is primarily spread through direct contact.

In general, the goal is to strike a balance between having enough links to accurately track the spread of the disease, while not having so many links that the network becomes difficult to interpret. This can be achieved through careful analysis and consideration of the specific disease and context.

Overall, choosing an appropriate distance threshold is an important step in using a transmission network to track the spread of a contagious disease. It can help ensure that the network is useful and informative, and can ultimately aid in efforts to control and prevent the spread of the disease.

## 2 METHODS

### 2.1 Scoring Heuristic Procedure

Network threshold selection procedure proceeds as follows:

1. For each candidate threshold  $d_L$ , in increasing order, ranging from the smallest genetic distance in the dataset, up to either the largest distance or a predetermined maximal threshold, we compute two network statistics:  $R_{12}$ , the ratio of the largest cluster to the second largest cluster, and  $C$  – the number of clusters in the network.
2. A priority score is assigned to each  $d_L$ . This score measures two properties of the threshold: Does  $R_{12}$  jump at  $d_L$ ? How far is the number of clusters  $C$  at  $d_L$  from the maximal number of clusters over all threshold values? Let there be  $N$  overall  $d_L$  candidate values, and assume we are examining the  $i$ th candidate,  $d_L^i$  with  $W < i \leq N - W$  ( $W$  is a positive integer defined below).
  - a. The  $R_{12}$  jump is computed by looking at the normalized ratio of the mean  $R_{12}$  values computed over the leading window  $d_L^{i+1} d_L^{i+W}$  and the trailing window  $d_L^{i-W} d_L^{i-1}$ . The width of the window,  $W$ , is defined as  $((\lceil \frac{N}{100} \rceil, 3), 30)$ . The distribution of ratios is converted to  $Z$

scores, and normalized relative to the largest positive  $Z$  score across all candidate distances, yielding the jump component of the score.

- b. The number of clusters,  $C_i$  at threshold  $d_L^i$  is first normalized to  $[0, 1]$  through  $\frac{C_{max}-C_i}{C_{max}-C_{min}}$  and next gated via a Gompertz function transform  $1 - e^{-e^{-25x+3}}$ . This function provides an ad hoc means for penalizing having too few clusters relative to the maximum over all ranges. For example, a threshold that yields 95% of the maximal number of clusters receives a score of 0.996, while a threshold that yields 85% - a score of 0.376.
- c. The priority score for  $d_L^i$  is the sum of the two components defined in (a) and (b).
3. The threshold with the highest priority score will be selected as the suggested automatic distance threshold, if the score is high enough (1.9 or more), and either of the two conditions hold.
  - a. No other thresholds have priority scores of 1.9 or higher
  - b. If other thresholds have priority scores of 1.9 or higher, then the range of thresholds represented by these options is small (no more than  $\log(N)$  times the mean step between successive  $d_L^i$ ).
4. If no single threshold can be selected in step 3, then the one with the highest priority score is suggested, and an inspection of the plot like the one on the analyze page is recommended to ensure that the threshold is sensible.

## 2.2 Assortativity

Degree-weighted homophily (DWH) is a measure of similarity between nodes in a network based on their attributes (such as demographic characteristics or behaviors) and their degree (i.e., the number of connections they have to other nodes in the network). It is used to quantify the extent to which nodes with similar attributes tend to be connected to each other more frequently than would be expected by chance. DWH is calculated as the ratio of the observed number of connections between nodes with similar attributes to the expected number of connections between such nodes, based on their degree.

In mathematical terms, it is defined as:

$$DWH = \frac{W_M + W_C - 2W_X}{\frac{d_{in}}{nodes_{in}^2} + \frac{d_{out}}{nodes_{out}^2}} \quad (1)$$

Where

- $W_M$  : Weight of in-group connections
- $W_C$  : Weight of out-group connections
- $W_X$  : Weight of cross-group connections
- $d_{in}$  : In-group degree
- $d_{out}$  : Out-group degree
- $nodes_{in}$  : number of in-group nodes
- $nodes_{out}$  : number of out-group nodes

DWH ranges from -1 to 1. A DWH value of 0 indicates that there is no more homophily than expected with chance, while a value of 1 indicates that there is perfect homophily (e.g. Birds always

link to birds, and only birds). A value of -1 is achieved for perfectly disassortative networks (e.g. Bird never linking with another bird).

DWH is used in social network analysis and in the study of how different attributes are related to the formation of connections between individuals. It is used as a way to measure the similarity of attributes between individuals in a network.

## 2.3 Implementation

The software implementation involves a step-by-step process that utilizes the HIV-TRACE suite of packages. It starts with calculating pairwise distances with the tn93 tool and a supplied multiple sequence alignment. This generated pairwise distances are supplied to the hivnetworkcsv script while providing the -A keyword argument. A brief outline of the software's implementation are as follows

1. Calculate pairwise distances: The user first calculates the pairwise distances using the tn93 fast pairwise distance calculator, providing the necessary threshold value and the input FASTA file. The command for this step is

```
tn93 -t 0.030 pol.fasta > pairwise_distances.15.tn93.csv
```

Please note that the threshold should include the maximal range one is intending to test.

2. Compute distance threshold scores: The hivnetworkcsv script is then executed with the required input file, format, and autotune option to generate a tab-separated output file, as shown below

```
hivnetworkcsv -i pairwise_distances.15.tn93.csv -f plain -A 0 > autotune_report.tsv
```

3. Visualize the report: Users can upload the generated autotune\_report.tsv file to <http://autotune.datamonkey.org/analyze> for visualization and further analysis of the data. This web-based platform provides an interactive environment to explore scores and other metrics across the range of tested outputs.

4. Run HIV-TRACE: Once AUTO-TUNED threshold(s) are settled upon after review, the user runs the HIV-TRACE command with the appropriate input FASTA file, distance threshold, and other required arguments. The output is saved as a JSON file. An example command is

```
hivtrace -i ./INPUT.FASTA -a resolve -r HXB2_prnt -t < autotune_threshold > -m 500 -g .05 > hivtrace.results.json
```

### 2.3.1 Optional : Compute Assortativity Metrics

5. Annotate results: The hivnetworkannotate script is used to annotate the results obtained from the HIV-TRACE step with attributes. The script takes the JSON results file, node attributes file, schema file, and a resolve flag as input.

```
hivnetworkannotate -n hivtrace.results.json -a node_attributes.json -g schema.json -r
```

For more information, users can refer to the hivnetworkannotate documentation.

6. Analyze the results with DWH: After the results file has been annotated, the user can proceed to the assortativity page, <http://autotune.datamonkey.org/assortativity>, for further analysis of the output.

AUTO-TUNE is readily accessible on GitHub as part of the hivclustering repository (<https://github.com/veg/hivclustering>). It is integrated into the command-line interface of the software as the -A or -auto-profile argument. hivclustering is a key component of the HIV-TRACE suite of tools, a resource for the inference, analysis and visualization of HIV transmission networks.

The Degree Weighted Homophily (DWH) calculation tool, an integral component of the assortativity step, is developed using TypeScript, a statically typed superset of JavaScript that ensures robustness and scalability. In an effort to promote accessibility and ease of integration, the DWH tool is packaged and distributed through the Node Package Manager (NPM), enabling researchers and developers to conveniently incorporate this advanced analytical tool into their own projects and workflows. DWH can be used in-browser or as a command line tool, allowing researchers and developers to employ the tool in an interactive command-line interface or integrate it into larger software applications, thus catering to a diverse array of technical needs and preferences. Instructions for usage and installation is found on Github (<https://github.com/veg/dwh>).

The described workflow offers a systematic approach to analyze potential distance thresholds for one's data with AUTO-TUNE, from calculating pairwise distances to visualizing and annotating results.

## 2.4 Visualization

Visualizations of AUTO-TUNE results are accessible at <http://autotune.datamonkey.org/analyze>. It is a dynamic and interactive web-based platform that offers visualization and analysis of results generated by AUTO-TUNE. The website provides a comprehensive view of the data by generating various plots across candidate distance thresholds. These include a score plot, allowing users to identify trends and anomalies across the full range of thresholds. Additionally, it generates a graph showing the number of clusters across candidate thresholds, one of the components that contribute to the score. The site also includes an R1/R2 plot that displays the ratio of the largest cluster to the second largest cluster across candidate thresholds, which is the other metric that contributes to the scoring heuristic.

An assortativity tool is available at <http://autotune.datamonkey.org/assortativity>, and is an advanced analytical tool engineered to facilitate the calculation of Degree Weighted Homophily (DWH) values. It utilizes the DWH NPM package to generate a tabular representation of DWH values corresponding to each value for a selected attribute annotation, providing an exhaustive examination of the interrelationships for the field. A notable feature is the computation of the panmictic range, which involves a label permutation test to generate the null distribution of DWH values. This feature establishes a comparative baseline that aids in determining the significance of homophily versus what would be expected by chance. Lastly, the site also provides a plot of the fraction of pairwise connections, normalized by degree, for each value pertinent to the selected field. This visual depiction facilitates an intuitive comprehension of the distribution and interconnections within the dataset.

The site aims to offer a user-friendly interface for data visualization, playing a crucial role in interpreting and understanding AUTO-TUNE's output data. The visualization code is available on Github (<https://github.com/stevenweaver/autotune-app/>).

## 2.5 Data curation and analysis

# 3 RESULTS

## 3.1 Case Study 1: Middle Tennessee, Seattle, and Alberta

## 4 DISCUSSION

### CONFLICT OF INTEREST STATEMENT

178 The authors declare that the research was conducted in the absence of any commercial or financial  
179 relationships that could be construed as a potential conflict of interest.

### AUTHOR CONTRIBUTIONS

180 The Author Contributions section is mandatory for all articles, including articles by sole authors. If  
181 an appropriate statement is not provided on submission, a standard one will be inserted during  
182 the production process. The Author Contributions statement must describe the contributions of  
183 individual authors referred to by their initials and, in doing so, all authors agree to be accountable  
184 for the content of the work. Please see here for full authorship criteria.

### FUNDING

185 Details of all funding sources should be provided, including grant numbers if applicable. Please  
186 ensure to add all necessary funding information, as after publication this is no longer possible.

### ACKNOWLEDGMENTS

187 This is a short text to acknowledge the contributions of specific colleagues, institutions, or agencies  
188 that aided the efforts of the authors.

### SUPPLEMENTAL DATA

189 Supplementary Material should be uploaded separately on submission, if there are Supplementary  
190 Figures, please include the caption in the same file as the figure. LaTeX Supplementary Material  
191 templates can be found in the Frontiers LaTeX folder.

### DATA AVAILABILITY STATEMENT

192 The datasets [GENERATED/ANALYZED] for this study can be found in the [NAME OF  
193 REPOSITORY] [LINK].

### FIGURE CAPTIONS

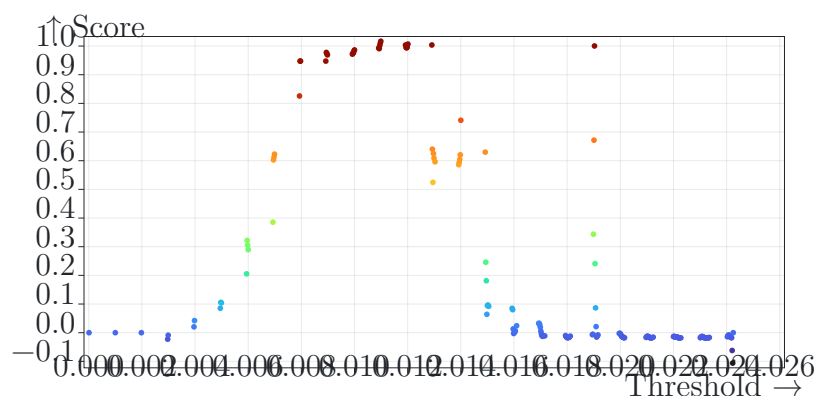


Figure 1. Enter the caption for your figure here. Repeat as necessary for each of your figures



Figure 2a. This is Subfigure 1.



Figure 2b. This is Subfigure 2.

Figure 2. Enter the caption for your subfigure here. (A) This is the caption for Subfigure 1. (B) This is the caption for Subfigure 2.