

TreeScan™ User Guide

for version 2.1



By Martin Kulldorff

July, 2022

<http://www.treescan.org/>

Table of Contents

Table of Contents.....	2
Introduction.....	4
TreeScan Software	4
Download and Installation	5
Test Run	5
Help System	6
Sample Data Sets	6
Statistical Methodology	8
Tree-Structured Variable.....	8
Tree Terminology	9
Tree-Based Scan Statistic.....	10
Poisson Probability Model	11
Bernoulli Probability Model	11
Purely Temporal Scan Statistic	12
Tree-Temporal Scan Statistic	13
Conditional versus Unconditional Analyses	13
Secondary Clusters.....	14
Comparison with Other Methods.....	15
Other Scan Statistics	15
Classification and Regression Trees (CART)	15
Input Data.....	16
Data Requirements	16
Tree File	16
Count File.....	16
Control File	18
Cut File.....	19
TreeScan Import Wizard.....	19
TreeScan File Format.....	20
Basic TreeScan Features	22
Analysis Tab	22
Input Tab.....	25
Output Tab	27
Advanced Features	30
Temporal Window Tab	30
Adjustments Tab	32
Inference Tab	34
Sequential Analysis Tab.....	36
Power Evaluation Tab	38
Advanced Input Tab.....	41
Additional Output Tab	42
Running TreeScan	44
Specifying Analysis and Data Options.....	44
Launching the Analysis.....	44
Status Messages	45
Warnings and Errors	45
Saving Analysis Parameters	46
Sequential Analysis.....	47
Parallel Processors	47
Batch Mode	48
Computing Time	48
Memory Requirements.....	49
Results.....	50

Standard Results File (*.txt).....	50
Mathematical Formulas.....	52
HTML Results File (*.html)	53
Comma Delimited Results File (*.csv)	53
Simulated Log Likelihood Ratios File (*_llr.csv).....	54
Miscellaneous.....	55
New Versions.....	55
Random Number Generator	55
Contact Us.....	55
Acknowledgements	56
Frequently Asked Questions.....	57
Input Data.....	57
Results.....	58
Operating Systems	58
TreeScan Bibliography.....	59
Suggested Citations.....	59
Methodology	59
Selected Applications.....	59
Other References Mentioned in this User Guide.....	60
Index	62

Introduction

TreeScan Software

Purpose

TreeScan is a free data mining software that allows users to analyze large data sets using different versions of the tree-based scan statistic.¹ The software was originally designed for pharmacovigilance, with the purpose to detect unsuspected drug and vaccine adverse reactions using large electronic health plan databases. The software may also be used for similar problems in other medical as well as non-medical fields, whenever the data can be classified into a hierarchical tree-like structure.

Tree Structure and Cluster Detection

To perform an analysis, the user must provide a pre-determined hierarchical tree structure of their data. For example, the tree may consist of (i) ICD-10 medical diagnosis codes with related diagnoses located on the same branch of the tree, (ii) pharmaceutical drugs with similar drugs on the same branch of the tree, or (iii) occupational codes with similar occupations on the same branch of the tree. The tree can have two or more hierarchical levels, representing the number of increasingly smaller sized ‘branches’ as one moves further away from the trunk towards the ‘leaves’. On each leaf and/or branch of the tree, there are observed and expected counts of some outcome. By considering cuts on different branches of the tree, closer or further away from the trunk, the tree-based scan statistic scans the tree for clusters where there are significantly more cases than expected, evaluating both very specific outcome definitions represented by a leaf as well as large groups of related outcomes represented by a big branch.

Data Granularity

A key feature of the tree-based scan statistic is that the granularity of the data does not have to be pre-specified. For example, when we are looking for potential adverse reactions to a pharmaceutical drug, it is impossible to know a priori if the drug may cause a very specific health outcome such as cardiac dysrhythmia, a more general set of related outcomes, such as a variety of heart problems, or an even wider set such as different types of cardiovascular issues. To cast the data mining net as wide as possible, the tree-based scan statistic simultaneously tests many overlapping groups of related outcomes.

Temporal Data

The tree-temporal scan statistic is an extension of the tree-based scan statistic, where scanning is also done over time, to detect temporal clusters on one or more branches of the tree. The TreeScan software can also perform a purely temporal scan statistic, in which case the tree is non-existent or ignored.

Multiple Testing

Multiple testing is by nature present in all forms of data mining. The tree-based scan statistic automatically adjusts for multiple testing. This is critical in order to conduct proper statistical analysis without generating a large number of false positives.

Developers and Funders

The TreeScan™ software was developed by Martin Kulldorff together with Information Management Services Inc. Financial support for TreeScan has been received from the National Institutes of Health and the Food and Drug Administration. Their financial support is greatly appreciated. The contents of TreeScan are the responsibility of the developer and do not necessarily reflect the official views of the funders.

Related Topics: [Statistical Methodology](#), [TreeScan Bibliography](#), [Acknowledgments](#)


Download and Installation

To install TreeScan, go to the TreeScan web site at: <http://www.treescan.org/> and select the TreeScan download link. Choose between the TreeScan versions for Linux, Mac or Windows. After downloading the TreeScan installation executable to your computer, click on its icon and install the software by following the step-wise instructions.

Related Topics: [Test Run](#), [New Versions](#).

Test Run

Before using your own data, we recommend trying one of the sample data sets provided with the software. Use these to get an idea of how to run TreeScan. To perform a test run:

1. Click on the TreeScan application icon.
2. Click on 'Open Saved Session'.
3. Select one of the parameter files, which has a *.prm ending.
4. Click on 'Open'.
5. Click on the Execute  button. A new window will open with the program running in the top section and a Warnings/Errors section below. When the program finishes running, the results will be displayed.

Note: The sample files should not produce warnings or errors.

Related Topics: [Sample Data Sets](#).

Help System

The TreeScan help system consists of three parts:

- i. The TreeScan User Guide in PDF format, which you are currently reading, located in the same folder as the TreeScan executable. It can also be obtained from the TreeScan web site (www.treescan.org/techdoc.html) or directly within the TreeScan software by selecting Help > User Guide. You may print this as a single document for easy reference.
- ii. Methodology papers describe the details about the statistical methods available in the TreeScan software¹. These papers are listed in the TreeScan bibliography, which can be found both at the end of this User Guide and on the web (<http://www.treescan.org/references.html>).
- iii. The sample data sets described below makes it easy to familiarize oneself with the software.

Related Topics: [Sample Data Sets](#), [TreeScan Bibliography](#).

Sample Data Sets

Three sample data sets are provided with the software. They are automatically downloaded to your computer together with the software itself. The numbers are completely made up, and do not reflect any real data. Their purpose is to illustrate the input file content and format and to conduct a simple test run.

Tree Scan Statistic, Poisson Model

Parameter file: Poisson.prm

Count file: Poisson.cas

Format: <nodeID> , <#cases> , <population>

Tree file: tree.tre

Format: <nodeID> , <parent nodeID>

Study period: n/a

Tree Scan Statistic, Bernoulli Model

Parameter file: Bernoulli.prm

Count file: Bernoulli.cas

Format: <nodeID> , <#cases> , <#controls>

Tree file: tree.tre

Format: <nodeID> , <parent nodeID>

Study period: n/a

Tree-Temporal Scan Statistic

Parameter file: TreeTemporal.prm

Count file: TreeTemporal.cas

Format: <nodeID> , <#cases> , <time of cases>

Tree file: tree.tre

Format: <nodeID> , <parent nodeID>

Study period: 1 to 28

Purely Temporal Scan Statistic

Parameter file: TemporalOnly.prm

Count file: TemporalOnly.cas

Format: <#cases> , <time of cases>

Tree file: none

Study period: 1 to 28

Related Topics: [Test Run](#), [Help System](#).

Statistical Methodology

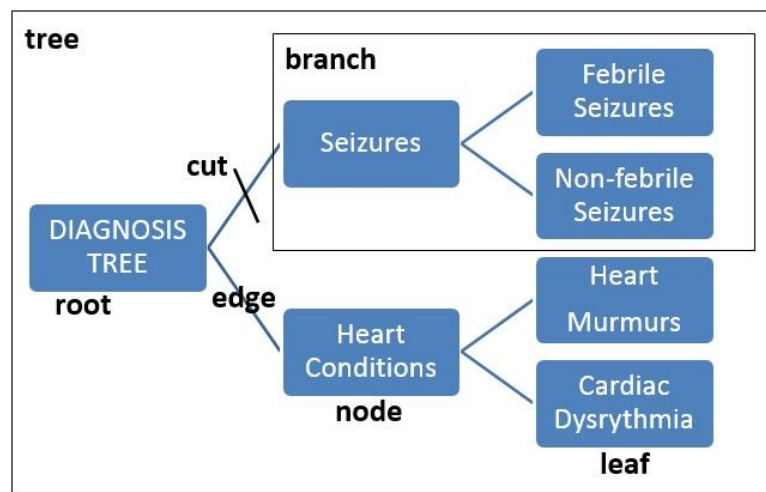
For the best description of the statistical methodology underlying the TreeScan software, we recommend reading the scientific papers describing the various versions of the tree-based scan statistic.¹ These are listed in the TreeScan bibliography at the end of the User Guide. Here we only provide a brief summary.

Related Topics: [TreeScan Software](#), [Basic TreeScan Features](#), [Analysis Tab](#), [Methodology Papers](#).

Tree-Structured Variable

The foundation of the tree-based scan statistic is a hierarchical tree structured variable. The tree must be pre-specified by the user, and it is not something that is created by the method. The tree should represent some relational/hierarchical aspect of the data. Here are a few examples of potential trees:

Disease Diagnosis Tree: In vaccine safety surveillance, the interest may be in detecting any unanticipated adverse reaction to a new vaccine. The tree will then consist of all medical diagnoses that could conceivably be caused by the vaccine. For example, non-febrile seizures and febrile seizures may be close to each other on one part of the tree while heart murmurs and cardiac dysrhythmias may be close to each other on another part of the tree.



Example of a Very Small Diagnoses Tree

Pharmaceutical Drug Tree: In pharmacovigilance, the interest may be detecting drugs or a class of related drugs that cause drug induced liver injury.⁴ The tree will then consist of all pharmaceutical drugs with drugs in the same class close to each other on the tree. For example, rofecoxib (vioxx) and celecoxib (celebrex) may be very close to each other on the tree since they are both cox-2 inhibiting nonsteroidal anti-inflammatory drugs, while ibuprofen is a little further away on the tree but on the same big branch since it is also a nonsteroidal anti-inflammatory drug but not a cox-2 inhibitor.

Occupational Tree: In occupational disease surveillance, the interest may be to determine if people with a particular occupation or group of related occupations are at higher risk for some disease.¹ The tree will then consist of all occupations of interest with closely related occupations close to each other on the tree. For example, elementary school teachers and high school teacher

may be close to each other on one branch of the tree, while coal miner and iron ore miners may be close to each other on another branch of the tree.

Related Topics: [TreeScan Software](#), [Tree Terminology](#), [Tree File](#), [Methodology Papers](#).

Tree Terminology

The terminology used for the tree-based scan statistic is derived from graph theory and computer science. For a graphical depiction, see the image, Example of a Very Small Diagnoses Tree, on the previous page.

Tree A hierarchical variable. In mathematical terminology, it is a directed acyclic simple graph consisting of nodes/vertices and edges.

Node A vertex on the tree. Each node contains a number of cases, a population at risk and/or a number of controls. The population must be non-negative and may for example represent expected counts. The cases must be non-negative integers. A node may have zero population and zero cases, and there will typically be many such nodes. If the population is zero, the number of cases must also be zero. Each node must have a unique name identifier, called the node identifier or the node ID.

Edge A line that connects two nodes on the tree.

Parent The parent of a node is the node that is immediately above it on the tree, connected through an edge.

Child The child of a node is a node that is immediately below it on the tree, connected through an edge. A node may have zero, one or multiple children. Note that a node can be a parent to one node while it is a child to another node.

Leaf A node that does not have any children. On most trees, the majority of nodes are leaves.

Root A node that does not have a parent. A typical tree has only one root, but there can be more than one.

Siblings Two or more nodes that share the same parent.

Descendants The descendants of a node are all its children, its children's children (grandchildren), its children's children's children (great grandchildren), and so on.

Ancestors The ancestors of a node are its parent, its parent's parent (grandparent), its parent's parent's parent (great grandparent), and so on

Cut A cut may be done on any one of the edges, just above a node. A simple cut defines a branch, the collection of nodes that include the node just below the cut together with all its descendants. The cut is identified by the name of the node that is just below the cut.

Branch The collection of nodes defined by a cut, including the node just below the cut and all its descendants.

Cluster A branch that has more observed cases than expected. Some clusters are statistically significant, and hence unlikely to have occurred by chance, but most clusters are not statistically significant. For the tree-temporal scan statistic, the cluster is defined by the collection of nodes together with a time interval.

Related Topics: [Tree Structured Variable](#), [Tree File](#).

Tree-Based Scan Statistic

Under the null hypothesis, a case is equally likely to occur anywhere on the tree, in any of the nodes, in such a way that the expected number of cases in any node is proportional to the population of that node. Under the alternative hypothesis, there is one or more branches on the tree where cases have a higher probability of occurring, constituting a cluster. The goal of the tree-based scan statistic is (i) to detect branches that contain a cluster of at least two cases without pre-specifying the branch a priori, and (ii) to determine whether the detected clusters are statistically significant after adjusting for the multiple testing inherent in the many overlapping branches evaluated.

The tree-based scan statistic is a likelihood ratio test. For each branch of the tree, the likelihood is calculated under both the null and alternative hypotheses, using a mathematical formula that depends on the probability model, as described in the tree scan methodology papers.¹ The branch with the maximum likelihood ratio is the most likely cluster, or the most likely cut, that is, the cluster of cases that is least likely to have occurred by chance. For computational reasons, the maximum log likelihood ratio (LLR) is used as the test statistic rather than the maximum likelihood ratio. Since one is a monotone function of the other, the two give exactly the same statistical test.

The distribution of the test statistic under the null hypothesis is not known, making it impossible to conduct inference using analytical methods, but we do know how to generate data under the null hypothesis. Because of that, we can do inference using computer simulations and Monte Carlo hypothesis testing.¹³ First we generate 999 (say) random replicas of the data generated under the null hypothesis. For each of these data sets, we calculate the maximum log likelihood ratio. Note that these maxima are attained for different branches for different random data sets. If the real data set is also generated under the null hypothesis, then there is a 5% probability that the maximum LLR from the real data set is among the 50 highest LLRs from the real and random data sets. If it is, then we can reject the null hypothesis at the $\alpha=0.05$ significance level. This probability is exact, so the hypothesis test is neither conservative nor liberal.

The Monte Carlo p-value is calculated as $p=R/(S+1)$, where R is the rank of the maximum LLR from the real data compared to the random data sets and S is the number of simulated Monte Carlo replications. In order for the p-value to be a 'nice looking' number, the number of simulations is restricted to 999 or some other number ending in 999 such as 1999, 9999 or 99999. That way it is always clear whether to reject or not reject the null hypothesis for typical cut-off values such as 0.05, 0.01 and 0.001. Additional Monte Carlo replications will increase statistical power, but beyond 999, the increase is marginal.

Related Topics: [TreeScan Software](#), [Basic TreeScan Features](#), [Analysis Tab](#), [Methodology Papers](#).

Poisson Probability Model

With the Poisson model,¹ the number of cases in each node is Poisson-distributed. For the unconditional version, the expected number of cases under the null hypothesis is provided in the count file. For the conditional version, the expected number of cases under the null hypothesis is proportional to its population, in such a way that the total expected across the whole tree is equal to the total number of observed cases. For the conditional Poisson model, the results will be the same if the population is multiplied by the same constant everywhere on the tree.

Example: For the Poisson model, cases may be health outcomes while exposed to a particular pharmaceutical drug, while the population is the age adjusted expected number of health outcomes.

Related Topics: [Analysis Tab](#), [Bernoulli Probability Model](#), [Methodology Papers](#).

Bernoulli Probability Model

With the Bernoulli model, there are cases and non-cases represented by a 0/1 variable. These variables may represent people with or without a disease, people being exposed or unexposed, or two different time periods of a person's life in a self-control type analysis. They may reflect cases and controls from a larger population, or they may together constitute the population as a whole. Whatever the situation may be, these variables will be denoted as cases and controls throughout the user guide, and their total will be denoted as the population. For each node, the population is fixed and non-random.

In the unconditional version, the user must specify the probability of being a case under the null hypothesis. In the conditional version, the analysis is conditioned on the total number of cases observed in the whole tree.

Example 1: In comparative drug safety cohort study, cases may be adverse events among recipients of drug A while the controls are adverse events among recipients of comparator drug B. If a 1:1 matching is used, the probability of being a case would be $\frac{1}{2}$, using the unconditional Bernoulli model.

Example 2: In a vaccine safety self-control design, cases may be health outcomes occurring in a risk window 1-14 days after vaccination while controls are health outcomes occurring in a control

window 29-56 days after vaccination. If the unconditional version is used, the probability of being a case would be 1/3, since the control window is twice as long as the risk window.

If an unconditional Bernoulli model is selected, one should indicate if data comes from a self-control design. This does not affect the statistical analysis per-se, but it uses a different self-control formula to calculate the relative risk, the excess number of cases and the attributable risk.

Related Topics: [Analysis Tab](#), [Poisson Probability Model](#), [Methodology Papers](#).

Purely Temporal Scan Statistic

While the tree and tree-temporal scan statistics are the key features of the TreeScan software, it is also possible to run a purely temporal scan statistic. In fact, the purely temporal scan statistic is a special case of the tree-temporal scan statistic when there is only one node in the 'tree'.

There are two versions of the temporal scan statistic, based on the uniform and Bernoulli probability models respectively. With the uniform model, cases under the null hypothesis are independently distributed over time according to the uniform probability distribution, with equal probability at each time. The alternative hypothesis is that there is some time interval where the cases occur with higher probability. Only case data are needed to do the analysis, and the analysis is conditioned on the total number of cases.

Temporal scan statistics are used to detect and evaluate the statistical significance of temporal clusters, without a prior specification of the risk window. The method adjusts for multiple testing of the many possible cluster times and lengths. For example, with a study time period of 56 days, it may consider all 2002 time intervals with a length of no more than 28 days, as possible temporal clusters, including [1-4], [3-6], [5], [6-23], [16-43], [31-32] and [29-56]. Alternatively, it may only consider a subset of those time intervals, as requested and specified by the user.

The maximum temporal cluster length should never be more than 50% of the total study time period. The reason for this is that a 'cluster' in a longer time interval would more accurately be thought of as two negative clusters with than expected fewer cases, before and after the time interval. For example, a 'cluster' in the [2-55] time period would more naturally be interpreted as a lack of cases on days 1 and 56 rather than an excess number of cases on days [2-55].

When the purely temporal scan statistic is used with the Bernoulli probability model, it is necessary to have both case and control data. Under the null hypothesis, every observation has the same probability of being a case, equal to the ratio of the total number of cases divided by the total number of cases plus controls. The analysis is conditioned on the total number of cases, the total number of controls, as well as the collection of the case and control times combined, but obviously not on the times for the cases nor the times for the controls.

Related Topics: [Analysis Tab](#), [Methodology Papers](#), [Tree Temporal Scan Statistic](#).

Tree-Temporal Scan Statistic

The tree-temporal scan statistic is a fusion of the standard tree-based scan statistic¹ and the temporal scan statistic.^{14,18} The temporal component in the tree-temporal scan statistic operates in the same way as described for the purely temporal scan statistic.

The tree-temporal scan statistic requires a time for each case. For the uniform probability model, only case data are needed, with no information needed about controls or a background population at risk. For the Bernoulli model, both cases and controls are needed.

In addition to the tree-structured variable, it is necessary to specify a study time period. That could be, for example, the 1 to 56 days following the initial (incident) use of a pharmaceutical drug. For each observed case, one records not only the node to which it belongs but also the time of the case, which should fall within the study time period.

With the tree-temporal scan statistic, we are in essence performing multiple temporal scan statistics, one for each of the many overlapping branches of the tree, adjusting for the multiple testing stemming both from the many branches and the many time intervals evaluated. Each time interval is evaluated on each of the branches, so if there are for example 1000 nodes on the tree and 2002 potential time intervals, there would be 2,002,000 potential clusters to evaluate and for which we need to adjust for multiple testing. If these were 2 million independent tests with independent non-overlapping data, there would be a huge loss in power when adjusting for all that multiple testing. With scan statistics, such a large loss in power does not happen, since the 2 million potential clusters are highly overlapping with each other. Hence, the penalty for adjusting for the multiple testing is fairly modest and not as much as one may originally think.

The tree-temporal scan statistic conditions the analysis on the number of cases observed in each node. This means that, unlike the pure tree scan statistic, there is no probability distribution to model the number of cases in each node. That number is now deterministic. What is probabilistic is the time of each case, which under the null hypothesis is assumed to be uniform across the study time period. Under the alternative hypothesis, there is at least one branch for which there is a temporal cluster of cases during a shorter or longer time interval.

Related Topics: [Analysis Tab](#), [Methodology Papers](#), [Purely Temporal Scan Statistic](#).

Conditional versus Unconditional Analyses

Scan statistics exist in both unconditional and conditional forms. Unconditional scan statistics need external information about the true probabilities underlying the null hypothesis. Conditional scan statistics are conditioned on the total number of cases observed, either in the whole data set, or in each of a subset of the data. The conditional Poisson and Bernoulli versions of the tree-based scan statistics both condition on the total number of cases observed in the whole tree. In this way, the inference is done as a comparison between the risk in different parts of the tree, and the total number of cases in the tree does not make it either more or less likely

to reject the null hypothesis. The tree-temporal scan statistics not only conditions on the total number of cases in the whole tree, but also on the vector of cases in each of the nodes.

Related Topics: [Bernoulli Probability Model](#), [Poisson Probability Model](#), [Tree-Temporal Scan Statistic](#), [Methodology Papers](#).

Secondary Clusters

TreeScan also identifies secondary clusters in the data set in addition to the most likely cluster, and orders them according to their likelihood ratio test statistic. There will often be a secondary cluster that overlaps with the most likely cluster, with the cuts being made higher or lower on the same branch of the tree. There may also be secondary clusters that do not overlap with the most likely cluster. The p-values for such clusters should be interpreted in terms of the ability of the secondary cluster to reject the null hypothesis on its own strength, whether or not the more likely clusters are true clusters or not. Hence, these p-values are not adjusted for the fact that there may be other clusters in the data.

Related Topics: [Analysis Results](#).

Comparison with Other Methods

Other Scan Statistics

Scan statistics were first studied in detail by Joseph Naus.^{[19](#)} Scan statistics has been developed for one dimension, such as time, as well as for spatial and spatio-temporal data. Excellent reviews of scan statistics have been provided in books by Glaz and Balakrishnan^{[15](#)}, Glaz, Naus and Wallenstein^{[16](#)} and Glaz, Pozdnyakov and Wallenstein.^{[17](#)} For spatial, temporal and spatio-temporal scan statistics, the free SaTScanTM software can be used, which is available at www.satscan.org.

Classification and Regression Trees (CART)

When using the tree-based scan statistic, the tree must be defined by the user, and as such, it is a tree-type variable, different from a categorical, ordinal, multinomial or continuous variable. The purpose of the tree-based scan statistic is not to create the tree. The method is hence very different from statistical data mining methods such as Classification and Regression Trees (CART), where the purpose is to create a tree structure.^{[10](#)}

Input Data

Data Requirements

Required Files: The input data should be provided in a number of files. A count file with the observed data is always needed, and a tree file defining the structure of the tree is needed except for a purely temporal analysis.

File Format: The data input files must be in TreeScan ASCII file format or you may use the TreeScan import wizard for dBase, Excel, comma delimited or space delimited files. Using such files, the wizard will automatically generate TreeScan file format files. Both options are described below.

Related Topics: [Input Tab](#), [Count File](#), [Tree File](#), [TreeScan Import Wizard](#), [TreeScan File Format](#).

Tree File

The tree file defines the structure of the tree. Each node on the tree must be assigned a Node ID, which can be any text string. Each line of the file represents one node. The first column contains the node ID. The second column denotes the node ID of the 'parent' of that node, that is, the closest node one level up on the hierarchical tree. If there is no parent node, the second column is left blank.

Most analyses use a single tree. If you simultaneously want to evaluate multiple trees, just include all of them in the same tree file, one after the other. If there are multiple trees, some leaves and nodes may be in only one of the trees while other may be in two or more trees. In the latter case, a node may have one parent in one of the trees and another parent in another tree.

The tree file is not needed for the purely temporal scan statistic.

Note: All nodes must be included in this file, including those that do not have a parent node.

Related Topics: [Input Tab](#), [Tree Structured Variable](#), [Tree Terminology](#), [TreeScan Import Wizard](#), [TreeScan File Format](#).

Count File

The count file provides information about observed cases. For the Poisson model, it should also contain expected cases or population numbers. The count file is used for all probability models and it should contain columns with the following information:

Node ID: Any numerical value or string of characters. Not used for a purely temporal analysis.

Number of Cases: The number of observed cases for the specified node. Must be a non-negative integer.

Population (Poisson model): The population size for the specified node. For the conditional Poisson model this could be raw population numbers, a covariate adjusted population at risk, or, the expected number of cases under the null hypothesis. If the unconditional Poisson model is used, it must be the expected counts. If the population size is zero for a particular node, it may either be left out or included in the count file with a specified value as zero. The population can be specified as integers or decimal numbers. It cannot be a negative number.

Time (Tree-temporal and purely temporal scan statistics): Specified in a generic format, which typically represent days, weeks, months or years, although it could also represent seconds, minutes, decades or centuries. Must be an integer, but could be either negative or positive.

Example: If there were 3 seizures 7 days after receiving the measles-mumps-rubella-varicella vaccine and 5 seizures 8 days after the vaccination, the following information should be provided for the tree-temporal model:

Seizure, 3, 7
Seizure, 5, 8

Note: Multiple lines may be used for different cases in the same node/leaf, having the same time attribute. TreeScan will automatically add them. For example, the above information can also be provided as:

Seizure, 1, 7
Seizure, 1, 7
Seizure, 1, 7
Seizure, 4, 8
Seizure, 1, 8

Censored Data

When the tree-temporal or purely temporal scan statistic is used, subjects are followed for a certain length study period and the timing of events are noted. If some subjects under observation does not have complete follow-up during the whole study period, those observations are censored. TreeScan can analyze such censored data. The test statistic used is no longer a likelihood ratio test statistic, but rather a pseudolikelihood. During the randomization step, when data are generated under the null hypothesis, the events of censored observations are not randomized according to a uniformly distribution over the whole study period but according to a uniform distribution from the start of the study period until the time of censoring. This assures that the inference is correct with accurate p-values, irrespectively of the test statistic used.

The censoring feature will work well for most data, but it should not be used if there is a large amount of censoring early on. For example, with a one-year study period, it should not be used if most subjects are censored at 60 days. It is then better to run an uncensored analysis with a shorter study period.

If a subject is censored early on, it does not provide much information, and TreeScan requires all subject to be observed for at least 50 time units. Subjects with a censoring time less than 50 are ignored and excluded from the analysis. TreeScan also requires that all subjects are observed for at least 10% of the total length of the study period. Hence, if the study period is days 1 to 1000, all subjects that are censored before day 100 are ignored.

The censoring time added in a forth column in the case file. The censoring time is the last time on which the subject was followed and on which an event would have been recorded if it had occurred. If no censoring time is provided, it is assumed that the observation is not censored but followed until the end of the study period. To specify a censoring time that is equal to the end of the study period is equivalent to not specifying a censoring time.

Related Topics: [Input Tab](#), [TreeScan Import Wizard](#), [TreeScan File Format](#), [ControlFile](#).

Control File

The control file is only used for the Bernoulli model. It has the same format as the count file.

Node ID: Any numerical value or string of characters. Not used for a purely temporal analysis.

Number of Controls: The number of observed controls for the specified node. Must be a non-negative integer.

Time (Tree-temporal and purely temporal scan statistics): Specified in a generic format, which typically represent days, weeks, months or years, although it could also represent seconds, minutes, decades or centuries. Must be an integer, but can be either negative or positive.

Related Topics: [Input Tab](#), [TreeScan Import Wizard](#), [Count File](#), [TreeScan File Format](#).

Cut File

The cut file is optional, as specified on the advanced Input Tab. If used, it should contain the following information:

Node ID: Any numerical value or string of characters.

Cut type: The type of cuts used below this particular node. The options are simple cuts, pairs cuts, triplets cuts and ordinal cuts, which are denoted as follows:


simple cut:	s	or	simple
pairs cut:	p	or	pairs
triplets cut:	t	or	triplets
ordinal cut:	o	or	ordinal

The definitions of these cuts are provided under Complex Cuts in the Advanced Input Tab section.

Note: Not all node IDs must be included in the Cut File. The default is simple cuts, so that will be used if a node is not listed in this file, or if this file is not provided.

Related Topics: Advanced Input Tab, Complex Cuts, [Tree File](#), [TreeScan Import Wizard](#), [TreeScan File Format](#).

TreeScan Import Wizard

The TreeScan Import Wizard can be used to import dBase, comma delimited, or space delimited files. It works for all import files. Launch the Import Wizard by clicking on the File Import  button furthest to the right of the text field for the file that you want to import. Follow the steps below to import files. Use the **Next** and **Previous** buttons to navigate between the dialogs.

Step 1 – Selecting the Source File

1. At the bottom of the Select Source File dialog, select the file type extension you are looking for. If you are unsure, select the All Files option. Supported file formats are: dBase III/IV, CSV, Excel, Text (*.txt) and TreeScan file formats.
2. Browse the folders and highlight the file you want to open. It will appear in the File Name text field.
3. Click on Open. The TreeScan Import Wizard will now appear.

Step 2: Specifying the File Structure

If you are importing a dBase or an Excel file, this step is automatically skipped. For all other source files, you need to specify the file structure using the File Format dialog box.

1. First specify the delimiting character and grouping indicator of the file.
2. If there are extraneous lines in the beginning of the file, type the number of lines that you would like to ignore in the text field below data sample area.
3. Click on **Next** to proceed to the next dialog box.

Step 3: Matching Source File Variables with TreeScan Variables

The top grid in this dialog box links the TreeScan variables with the input file variables from the source file. The bottom grid displays sample data from the chosen input file.

1. To match the variables, click on one of the places where it says `unassigned`.
2. Select the appropriate variable from the input file to go with the chosen TreeScan variable.
3. When all the required and optional variables that you selected have been matched, click on the Execute button to import the file. This will create a temporary file in TreeScan ASCII file format.

Step 4: Saving the Imported File

The imported file, which is in TreeScan ASCII file format, must be saved at least temporarily. The default is to save it to the TEMP directory and after the analysis is completed you may erase the file. You can also save it to some other directory of your choice and use it for future analyses without having to recreate it by using the Import Wizard again.

Related Topics: [Input Tab](#), [Count File](#), [Tree File](#), [Cut File](#).

TreeScan File Format

As an alternative to using the TreeScan Import Wizard, it is possible to directly write the name of the input files in the text fields provided on the Input Tab, or to browse the file directories for the desired input files using the button to the right of that box. The files must then be in TreeScan comma delimited file format, which are ASCII files with one row for each node/time combination with columns as defined below. Such files can be created using any text editor and most spreadsheets. The order of the columns in the file is very important, but the rows can be in any order. The optional variables, defined above, are optional columns in the TreeScan file format.

Each input file will have a multiple rows, and the following columns:

Tree File Format (*.tre):

<node ID> , <parent node ID>

Count File Format (*.cas):

<node ID> , <#cases> , <population> (Poisson model)

<node ID> , <#cases> (Bernoulli model)

<node ID> , <#cases> , <time> (tree-temporal model)

<#cases> , <time> (purely temporal model)

Control File Format (*.cas):

<node ID> , <#controls> (Bernoulli model)

<node ID> , <#controls> , <time> (tree-temporal model)

<#controls> , <time> (purely temporal model)

Cut File Format (*.cut):

<node ID> , <cut type>

Alternative Hypothesis File Format (*.alt):

<node ID> , <relative risk> (Poisson model)

<node ID> , <probability> (Bernoulli model)

<node ID> , <relative risk> , <start time> , <end time> (tree-temporal model)

<relative risk> , <start time> , <end time> (purely temporal model)

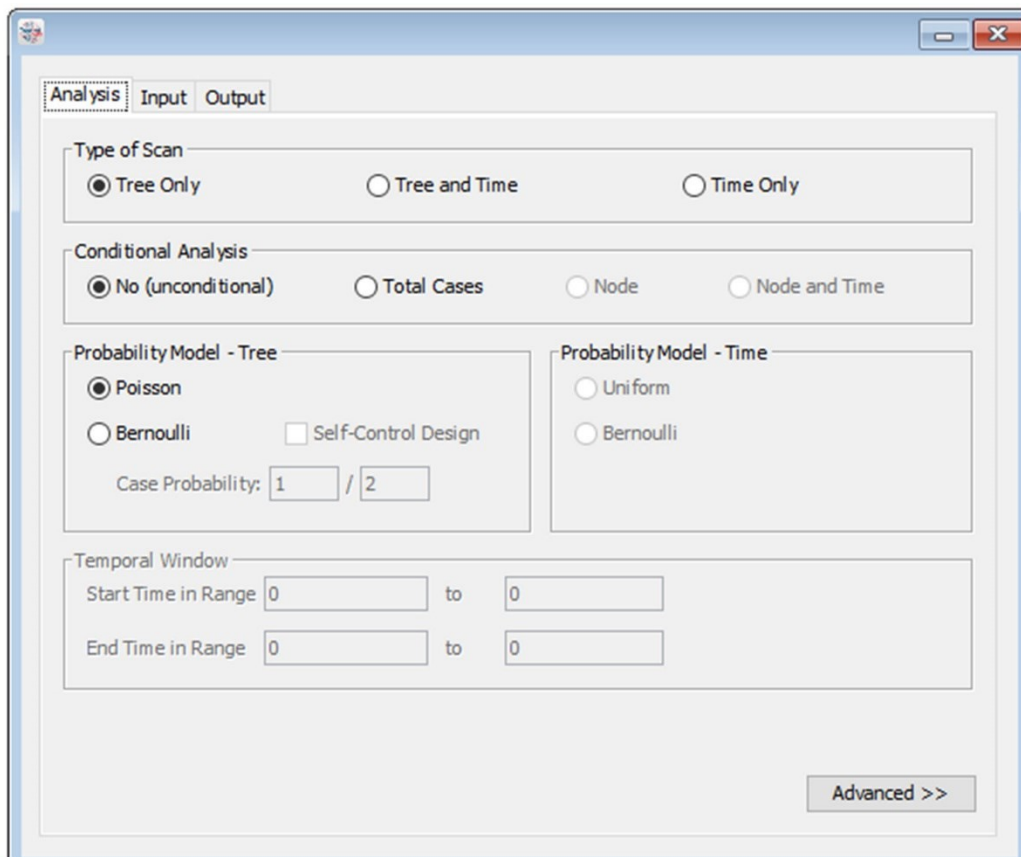
Related Topics: [Input Tab](#), [Count File](#), [Tree File](#), [Cut File](#), [TreeScan Import Wizard](#).

Basic TreeScan Features

Most TreeScan analyses can be performed using a set of basic analysis and data features. The users specify these on three different window tabs for analysis, input and output options respectively. These contain all required specifications for a TreeScan analysis as well as a few optional ones. Additional features, all optional, can be specified on the advanced features tabs.

Related Topics: [Statistical Methodology](#), [Input Tab](#), [Analysis Tab](#), [Output Tab](#), [Advanced Features](#).

Analysis Tab



The screenshot shows the 'Analysis Tab' dialog box with three tabs: 'Analysis', 'Input', and 'Output'. The 'Analysis' tab is selected. It contains several sections: 'Type of Scan' with radio buttons for 'Tree Only' (selected), 'Tree and Time', and 'Time Only'; 'Conditional Analysis' with radio buttons for 'No (unconditional)' (selected), 'Total Cases', 'Node', and 'Node and Time'; 'Probability Model - Tree' with radio buttons for 'Poisson' (selected) and 'Bernoulli', a checkbox for 'Self-Control Design', and a 'Case Probability' field set to '1 / 2'; 'Probability Model - Time' with radio buttons for 'Uniform' and 'Bernoulli'; and 'Temporal Window' with 'Start Time in Range' and 'End Time in Range' fields, each set to '0' with a 'to' separator. An 'Advanced >>' button is in the bottom right corner.

Analysis Tab Dialog Box

The Analysis Tab is used to specify the fundamental analysis options. Additional optional features are available by clicking on the Advanced button in the lower right corner.

Type of Scan

TreeScan may be used for a purely tree-based analysis, scanning only the tree, as well as for a tree-temporal scan statistic which simultaneously scans both the tree and a temporal time period. It is also possible to run a purely temporal scan statistic, not using a tree at all.

Conditional Analysis

With the TreeScan software, it is possible to conduct either unconditional or conditional analyses, where the analysis is conditioned on the total number of cases in the whole tree or in parts of the tree. In an unconditional analysis, random data is generated either from expected counts (Poisson model) or by using a pre-specified event probability (Bernoulli model). The expected counts and the event probability must be provided by the user.

With a conditional scan statistic, the analysis is conditioned on some aspect of the observed data. The analysis can be conditioned on the total number of cases found in the tree as a whole. This means that each random data set has exactly the same number of total cases as the real data set. For the Poisson model this means that under the null hypothesis, the expected number of cases in a particular node is no longer identical to the expected count provided by the user. Instead, the expected number of cases in that node is the total number of cases times the proportion of the expected count that is in that node. For the Bernoulli model, it means that under the null hypothesis, the probability of being a case is equal the total number of cases in the whole tree divided by the total number of observations (cases and controls).

For the tree-temporal scan statistic, the analysis is always conditioned on the number of cases in each node. That is, a node has exactly the same number of cases in each of the random data sets as in the real data set. This means that the methods does not evaluate whether there are branches on the tree with more cases than expected. Instead, it evaluates if there are branches on the tree that has a temporal cluster of cases. In addition to conditioning on the number of cases in each node, a tree-temporal analysis may also condition on the total number of cases at each time point. The purpose of this is to adjust for any purely temporal variation in the data, that is common to the tree as a whole.

For a purely temporal scan statistic, the analysis is always conditioned on the total number of cases observed.

Probability Model - Tree

There are two different probability models that can be used for the number of cases observed in each node: Poisson and Bernoulli.

Poisson Model: The Poisson model should be used when the background population reflects a certain risk density such as total person years, or, some covariate adjusted expected counts.

Bernoulli Model: The Bernoulli model should be used for 0/1 type data, such as individuals who may or may not have a disease. Those who have the disease are 'cases' and those who don't have the disease are 'controls'. For an unconditional analysis, it is necessary to specify the probability of being a case, which must be greater than 0 and smaller than 1. If you are doing a self-control type analysis, that should be indicated. This does not affect the statistical analysis, but ensures that the self-control formulas are used when calculating the relative risk, the excess number of cases and the attributable risk.

When the tree-temporal scan statistic is selected, the analysis is conditioned on the number of cases in each node, so there is no probability model for the number of cases on the nodes. That is, the number of cases in each node is deterministic rather than probabilistic.

Probability Model – Time

For the tree-temporal and purely temporal scan statistics, the two models available for the time dimension are uniform and Bernoulli. With the uniform model, a case is equally likely to occur during any of the days (or other unit) in the data time range, when the null hypothesis is true. Under the alternative hypothesis, there is some time interval where cases are more likely to occur.

With the Bernoulli model, each observation has the same probability of being a case when the null hypothesis is true. Under the alternative hypothesis, there is some time interval where observations have a higher probability of being a case compared to outside the time interval.

Note: The tree only Bernoulli model is a special case of the tree-temporal scan statistic with the uniform probability model, when there is only two time periods in the data time range, such as [1,2], and when the start and end ranges of the temporal window are [1,1] and [1,1].

Temporal Window Range

The maximum temporal window size is always set to be at most 50% of the data time period. That is, the time period inside the cluster can never be larger than the sum of the time periods before and after the cluster. This is to ensure that we are evaluating an excess number of cases inside the window, rather than a deficit of cases during a very small time period at the very beginning and/or end of the data time range.

As an addition to this firm restriction set by the software, it is necessary to specify the collection of windows to be evaluated in terms of a range of start times for the temporal clusters as well as a range of end times. For example, if [1,2] is selected as the start range and [4,5] as the end range, only the following temporal clusters are evaluated: [1,4], [1,5], [2,4] and [2,5]. Most commonly, a much larger collection of windows are used. By specifying both the start and end range to be identical to the data time range on the Input Tab, all possible windows with a length less than 50% are evaluated.



This option is only relevant and available for the tree-temporal scan statistic.

Related Topics: [Basic TreeScan Features](#), [Statistical Methodology](#), [Poisson Model](#), [Bernoulli Model](#), [Tree-Temporal Model](#), [Conditional Versus Unconditional Analyses](#).

Input Tab

The screenshot shows the 'Input' tab of a dialog box. It has three tabs: 'Input', 'Analysis', and 'Output'. The 'Input' tab is active. It contains three text input fields for file names: 'Tree File (not used for Time Only scan):', 'Count File:', and 'Control File (Bernoulli Only):'. Each field has a small 'Open File' icon (a folder with a plus sign) to its right. Below these fields is a 'Time Precision' section with five radio buttons: 'None' (selected), 'Generic', 'Year', 'Month', and 'Day'. At the bottom is a 'Data Time Range' section with two sets of date pickers. The 'Start Date' is set to Year: 2000, Month: 1, Day: 1. The 'End Date' is set to Year: 2000, Month: 12, Day: 31. An 'Advanced >>' button is located at the bottom right of the dialog box.

Input Tab Dialog Box

The Input Tab is used to specify the names of the input data files as well as the nature of the data in these files. If the files are in TreeScan comma delimited file format, they may be specified either by writing the name in the text box or by using the Open File  button. If they are not in TreeScan comma delimited file format, they must be specified using the TreeScan import wizard, by clicking on the File Import  button. Both the TreeScan comma delimited file format and the TreeScan import wizard are described in the Input Data section.

Tree File Name

Specify the name of the input file that defines the tree structure. This file is required for all analyses except the time only scan statistic, irrespectively of the probability model used.

Count File Name

Specify the name of the input file with the observed case data. This file is required for all scan statistics, irrespectively of the probability model used.

Control File Name

Specify the name of the input file with the observed control data. This file is only used for Bernoulli data.

Time Precision

Indicates whether the count file and the control file (when applicable) contain information about the time of each case (and control). If they do, this section also indicates whether the time precision should be read as generic, days, months or years. If the time precision is specified to be days but the precision in the count or control file is in month or year, then there will be an error. If the time precision is specified as years, but the count or control file includes some dates specified in terms of the month or day, then the month or day will be ignored.

Data Time Range

Only used for the tree-temporal and purely temporal scan statistics. Specify the start and end of the time range for which temporal data was collected. All times in the count file should fall on or between the start and end times.

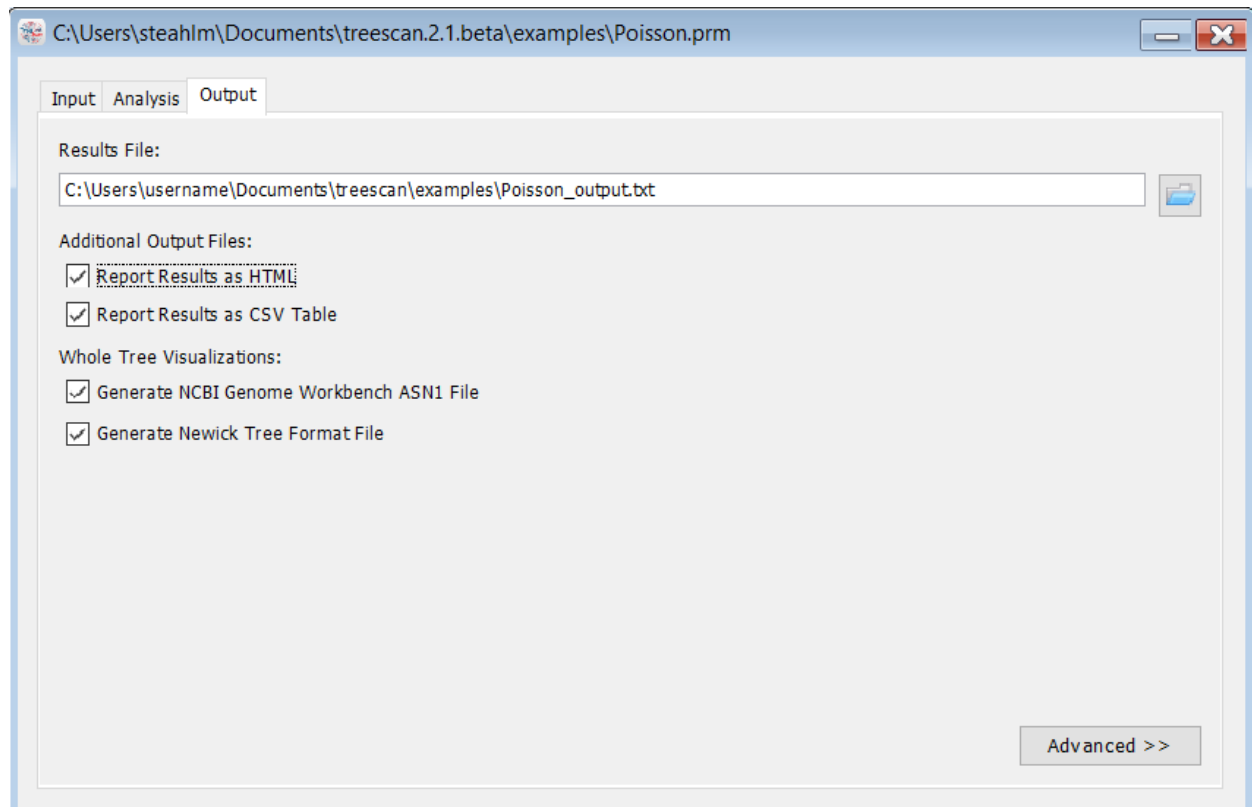
Data Time Range Start: The earliest time included in the study period.

Data Time Range End: The latest time included in the study period.

The times may be negative or positive. Provided in a generic format, the times can represent seconds, minutes, hours, days, weeks, months, years or any other unit of time.

Related Topics: [Basic TreeScan Features](#), [Count File](#), [Tree File](#).

Output Tab



Output Tab Dialog Box

The Output Tab is used to set parameters defining the output information provided by TreeScan. These parameters do not influence the actual analysis; just the way that the results are reported.

Text Output Format

A standard text based results file is automatically shown after the completion of the calculations. It contains information about the clusters detected, summary information about the data, computing time and the analysis parameters chosen. Specify the name of this file. Other optional output files will have the same name, but with different filename extensions.

Warning: If you specify the name of a file that already exists, the old file will be overwritten and lost.

HTML Output Format

If requested, TreeScan will create a HTML file that shows the results in a table, using the web browser. It will automatically launch after the analysis is complete. The name of this output file

is the same as the name of the text output format file, but with a different filename extension (*.html).

CVS Table Output Format

If requested, TreeScan will create a comma delimited (CSV) output file, that will show the detected clusters in table format. This file can easily be imported into Excel, SAS, R or other software for further formatting, depiction or analyses. The name of this output file is the same as the name of the text output format file, but with a different filename extension (*.csv). You must manually open this file after the TreeScan run is completed.

Generate NCBI Genome Workbench ASN1 File

If requested, TreeScan will create an ASN1 file when the scan is run. The name of this output file is the same as the name of the text output format file, but with '_ncbi' appended to the end and a different filename extension (*.asn).

Generate Newick Tree Format File

If requested, TreeScan will create a Newick file when the scan is run. The name of this output file is the same as the name of the text output format file, but with a different filename extension (*.nwk).

Viewing Additional Output Files

In order to view additional output files created during the analysis, a separate program is required. There are several such programs available online. The following tools are recommended:

NCBI Genome Workbench. This is a full-featured desktop application that is the intended use for the .asn output file (NCBI ASN1) and can be used to view and filter on tree nodes and reported cuts, as well as many other features. The program can be downloaded here: <https://www.ncbi.nlm.nih.gov/tools/gbench/>. The NCBI Genome Workbench application has many capabilities to assist researchers but for the purposes of TreeScan, we're only interested in Tree View feature. More information on that feature can be found here: <https://www.ncbi.nlm.nih.gov/tools/gbench/tutorial3/#treeviewfeatures>.

NCBI Tree Viewer. This is an online version of the tree viewing, which can be considered a lesser version of the Workbench program described above with a limited set of functionality. It can view either an NCBI ASN1 file or a Newick file (.nwk). It can be found here: <https://www.ncbi.nlm.nih.gov/tools/treeviewer/>.

IcyTree. This is an online Newick file viewer with the ability to attach metadata from csv – which works well with our csv table file. When using this tool, you first upload the Newick file. Afterward, you can attach the metadata csv file from the File menu and assign the label. This resource can be found here: <https://icytree.org/>.

Related Topics: [*Basic TreeScan Features*](#), [*Analysis Results*](#).

Advanced Features

While most TreeScan analyses can be performed using the features on the three basic tabs for analysis, input and output, a few additional options are available as advanced features. These features are reached through the Advanced button on the lower right corner of each of the three main tabs. 'Advanced' should be interpreted as 'additional' or 'uncommon' rather than 'complex', 'difficult' or 'better'.

Since many of the advanced options depend on the selections made on the Analysis and Input Tabs, it is recommended that those two tabs be filled in first.

Related Topics: [Basic TreeScan Features](#), [Inference Tab](#), [Additional Output Tab](#).

Temporal Window Tab

Advanced Analysis Options

Temporal Window Adjustments Inference Sequential Power Evaluation

Maximum Temporal Size

☒ is 50.0 percent of the data time range ($\leq 50\%$)

☐ is 1 data time units

Minimum Temporal Window

Minimum temporal size is 2 data time units

☒ Ensure that temporal window length is at least 20.0 percent of the between time zero and the end of the temporal window

☐ Prospective Evaluation

Temporal Window

☐ Include only windows with:

Start time in range: 0 to 0

End time in range: 0 to 0

Set Defaults Close

Temporal Window Tab Dialog Box

This tab is reached by clicking the Advanced button in the lower right corner of the Analysis Tab. It is only relevant when using the tree-temporal or tree-only scan statistic.

Maximum Temporal Cluster Size

The maximum temporal cluster size can be specified in terms of a percentage of the study period as a whole or as a certain number of time units. If specified as a percent, it can be at most 50 percent. If specified in time units, it can be set to at most half the length of the data time range specified on the Input Tab.

Minimum Temporal Cluster Size

A minimum temporal cluster size can be specified in terms of a certain number of time units. The default value is two time units. When one is specified, there is no minimum restriction on the temporal cluster size.

In addition to a fixed minimum temporal length on the cluster, it is also possible to define a minimum as a percentage so that small clusters are allowed soon after exposure but not later on. For example, it may be meaningful to evaluate a two-day cluster 3-4 days after exposure but not 103-104 days after exposure. With this option, it is possible to set a minimum so that the length of the risk window is at least X percent of the distance from zero to the end of the risk window. Suppose the risk window is [a,b]. For a chosen value of X in the range [0,100], the mathematical formula is then the requirement that $(b-a+1)/b \geq X$. The default is $X=0.20=20\%$. If no such minimum restriction is wanted this feature can be deselected, making it equivalent to $X=0$.

Prospective Evaluation

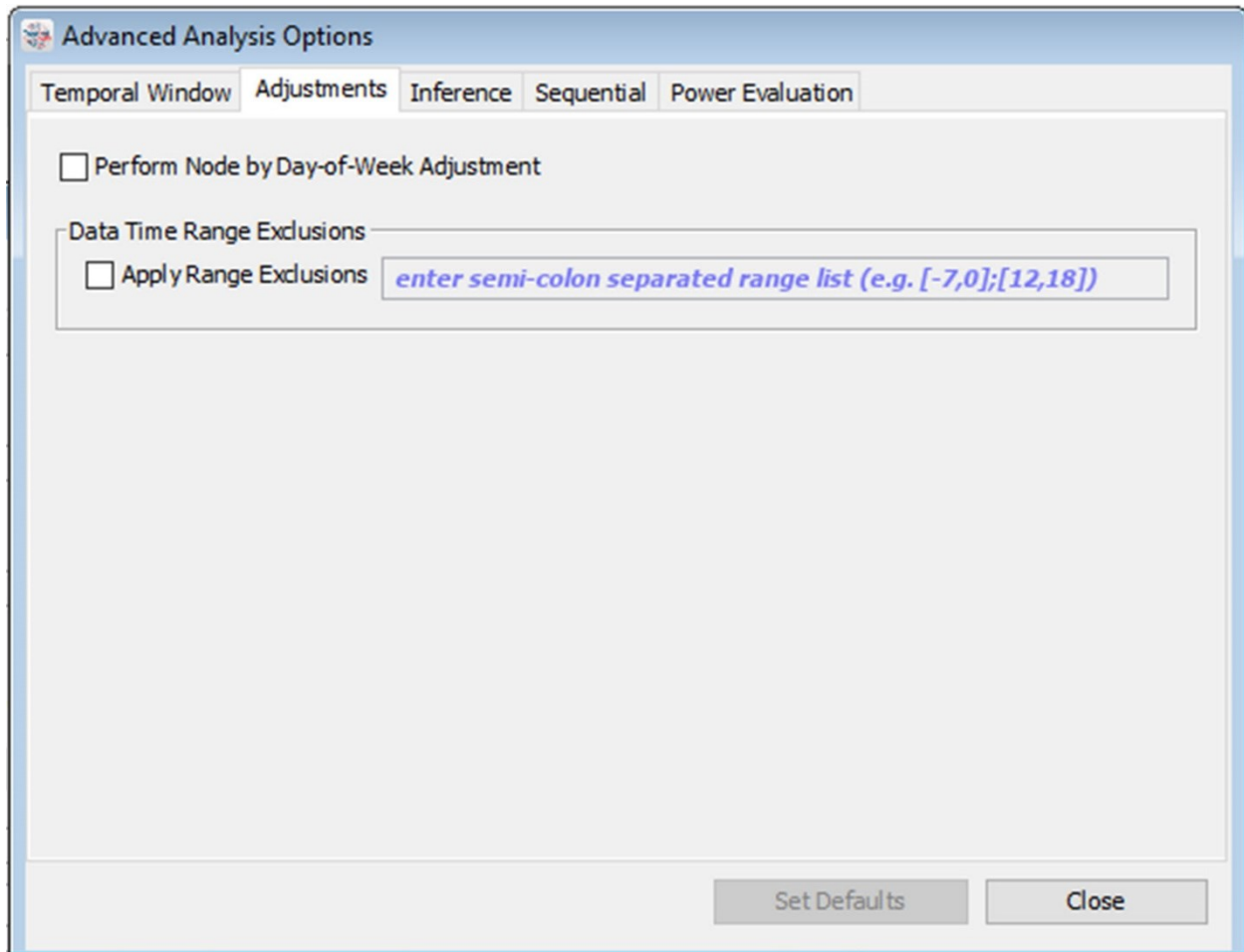
The prospective option is used for the early detection of disease outbreaks, when scans are repeated every day, week, month or year. Only alive clusters, clusters that reach all the way to current time as defined by the study period end date, are then searched for.

Flexible Temporal Window Definition

TreeScan will evaluate all temporal windows less than the specified maximum, and for prospective scans the same is true with the added restriction that the end of the window is identical to the study period end date. When needed, TreeScan can be more flexible than that, and it is possible via these settings to define the scanning window as any time period that starts within a predefined 'start range' and ends within a predefined 'end range'.

Related Topics: [Analysis Tab](#), [Input Tab](#), [Tree-Temporal Scan Statistic](#).

Adjustments Tab



Adjustments Tab Dialog Box

This tab is reached by clicking the Advanced button in the lower right corner of the Analysis Tab.

Day-of-Week Adjustment

By checking this box, TreeScan will automatically perform day-of-week adjustment. This option is only relevant for the tree-temporal and purely temporal scan statistics.

Suppose that we are looking for adverse events after vaccination, and the time is the number of days between the vaccination and the adverse event. If *both* the vaccinations and the outcome have a day-of-week pattern, such as being more common on weekdays than weekends, the assumption of a uniform probability over time is no longer valid. The day-of-week adjustment takes care of this, using stratified randomization. If only the exposure, or only the outcome, has a day-of-week effect, the uniform assumption holds, and there is no need to do a day-of-week adjustment.

Multiple Data Time Ranges

When using the tree-temporal scan statistic, it is most common to have a single time range as the study period, e.g. days 1 to 56. Sometimes it is of interest to use two or more non-consecutive time periods though. For example, if the study period is from 60 days before exposure to 60 days after exposure, one may wish to exclude the 14 days prior to exposure as well as the day of exposure in order to minimize the risk of bias due to confounding by indication or contra-indication. As another example, if we are interested in studying the risk of seizures after the DTaP vaccine, we may wish to exclude days 7-10 after vaccination since there is a known risk of seizures 7-10 days after measles containing vaccines and such vaccines are often given on the same day as the DTaP vaccine.

After defining the data time range on the Input Tab, as the first time in the first time range, to the last time in the last time range, then specify the intermediate time ranges that should be removed. This is done on the Adjustments Tab, one of the advanced analysis tabs. The time periods to exclude are specified in the following format: [-14,0];[7,10] If the input file includes events in any of these time periods, they are ignored in the analysis.

This feature is only available for the tree-temporal scan statistic that is conditioned on both time and node.

Related Topics: [Analysis Tab](#), [Tree-Temporal Scan Statistic](#).

Inference Tab

The screenshot shows a software window titled "Advanced Analysis Options" with a blue header bar. Below the header is a tabbed interface with four tabs: "Temporal Window", "Adjustments", "Inference" (which is selected and highlighted), and "Sequential", followed by "Power Evaluation". The "Inference" tab contains four distinct sections, each with a title and a text input field:

- Monte Carlo Replications:** A text input field containing the value "999". The label reads "Number of replications (0, 9, 999, or value ending in 999):".
- Tree Levels:** A text input field containing the placeholder text "enter comma separated list of integers - root is 1". To the left of the field is a checkbox labeled "Do not evaluate tree levels:" which is currently unchecked.
- Prospective Analyses:** A dropdown menu with the text "Daily" selected. The label reads "How frequently are analyses performed?".
- Minimum Number of Cases:** A text input field containing the value "2", followed by the text "cases". The label reads "Restrict cuts to have at least:". Below this section is a large, empty rectangular area.

At the bottom right of the dialog box are two buttons: "Set Defaults" and "Close".

Inference Tab Dialog Box

This tab is reached by clicking the Advanced button in the lower right corner of the Analysis Tab.

Monte Carlo Replications

To calculate p-values for detected clusters, TreeScan uses computer simulations to generate a number of random replications of the data set under the null hypothesis. If the maximum likelihood ratio calculated for the most likely cluster in the real data set is high compared to the maximum likelihood ratios calculated for the most likely clusters in the random data sets, that is evidence against the null hypothesis and for the existence of clusters.

The test statistic is calculated for each random replication as well as for the real data set, and if the latter is among the 5 percent highest, then the test is significant at the 0.05 level. If it is among the 1 percent highest, the test is significant at the 0.01 level, and so on. This is called Monte Carlo hypothesis testing, and was first proposed by Dwass.¹⁴ Irrespective of the number of Monte Carlo replications chosen, the hypothesis test is unbiased, generating a correct significance level that is neither conservative nor liberal nor an estimate. The number of replications does affect the power of the test though, with more replications giving slightly higher power. A higher number also increases the computing time.

In TreeScan, the number of replications must be at least 999 to ensure excellent power for all types of data sets. For small to medium size data sets, where computing time is not an issue, at least 9999 replications are recommended.

Note: The number of Monte Carlo replications can be specified to be 0 or 9 in order to do a test run, but that will not provide p-values for the analysis.

Tree Levels

The level on the tree is defined as its distance from a node to the root node, in terms of the number of intermediate node levels. The root nodes are defined as level 1; children of a root node are defined as level 2; grandchildren of a root node are defined as level 3; and so on. If a node is a descendant of a root node in more than one way, the level is defined as the distance from its closest root node.

By default, TreeScan will evaluate cuts on all levels of the tree. If there are some levels on the tree that should not be evaluated, it is possible to exclude those levels. This is done on the advanced inference tab, by selecting 'Do not evaluate tree levels', and then specifying those tree levels that should not be evaluated, with commas in between. For example, by specifying [2,3,7], TreeScan will not evaluate nodes 2, 3 or 7, but only 1,4,5,6,8 and beyond.

Frequency of Prospective Analyses

When doing a prospective scan, the frequency of the analysis is normally the same as the precision of time data. For example, daily data feeds are analyzed daily, weekly data feeds are analyzed weekly and monthly data feeds are analyzed monthly. If the prospective analyses are

conducted less frequently than the time units, that must be specified on this tab. For example, daily data may be available, but the daily data only arrives once a week so that analyses can then only be conducted on a weekly basis. Prospective Evaluation must be enabled on the Temporal Windows tab before analyses frequency can be set here.

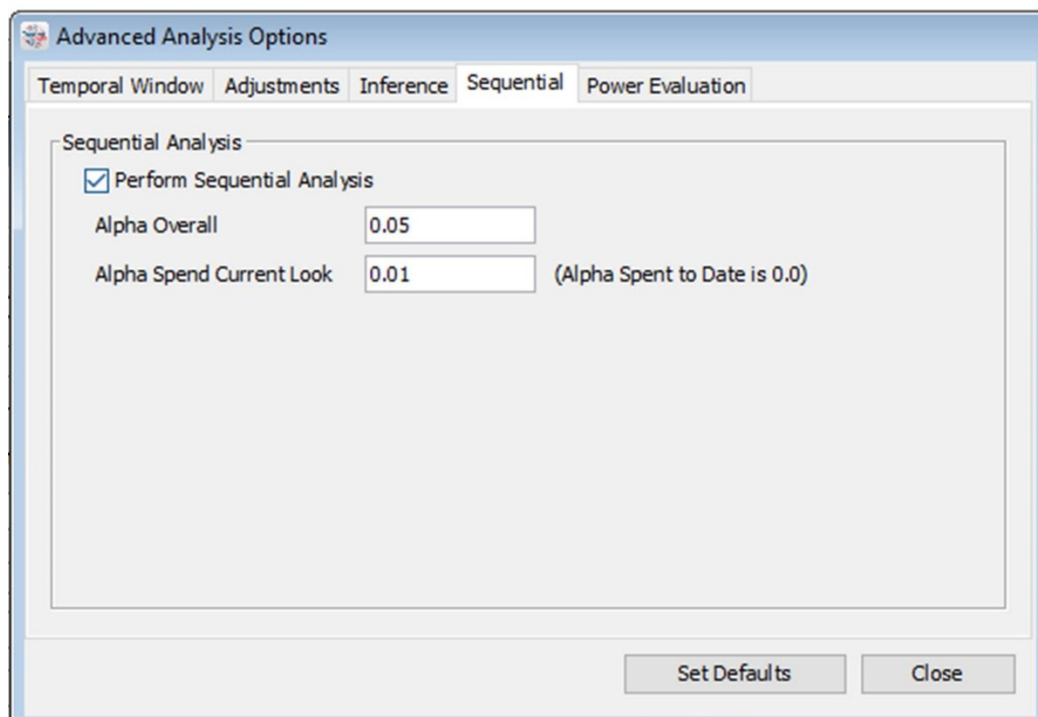
Minimum Number of Cases

It is possible to require that the detected cuts have at least a minimum number of cases. The default value is 2. By putting a higher value, such as 5, there will be no cuts with less than 5 cases. This will slightly increase the power to find clusters with 5 or more cases. Note that this feature does not restrict the collection of potential cuts that are evaluated.

The minimum number of cases should be selected before the analysis is done as part of the study specifications. It is not appropriate to try different values on the minimum number of cases, and then select the results with the lowest p-value. That invalidates the statistical inference and creates biased p-values.

Related Topics: [Analysis Tab](#), [Tree-Based Scan Statistic](#), [Random Number Generator](#), [Analysis Results](#).

Sequential Analysis Tab



Sequential Tab Dialog Box

This tab is reached by clicking the Advanced button on the lower right corner of the Analysis Tab.

Sequential statistical analysis is used when the data is analyzed repeatedly as new data arrives, gradually increasing the sample size. This option is only available for the tree only Bernoulli model.

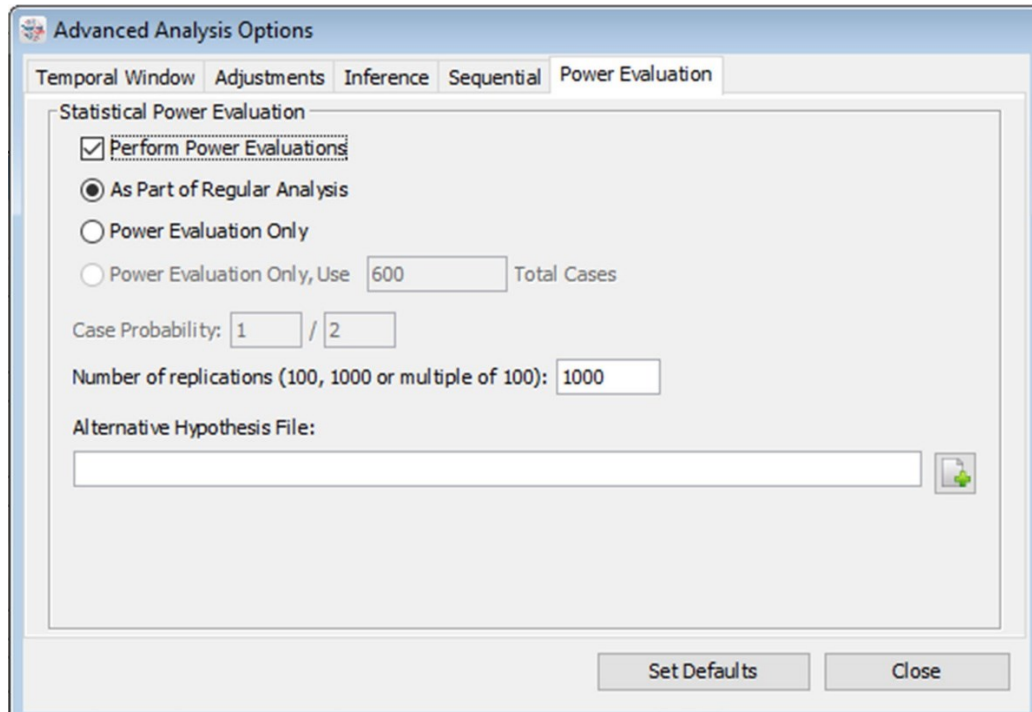
Overall Alpha Level: The overall alpha level determines the combined probability of rejecting the null hypothesis when it is true, at any time during the repeated analyses. This must be set before the first analysis and cannot subsequently change. The default value is 0.05.

Alpha to Spend at the Current Look: At each of the repeated analyses, one must specify how much to spend at that look. This must be less or equal than the difference of the overall alpha level minus the amount that has already been spent. It does not have to be the same at each look. For example, if there is very little additional data, one may want to spend less as compared to when there is a lot of additional data. It may only depend on the sample size of each new batch of data; but not on the observations in the data. The alpha to spend at each look can be specified as one goes, just before each look. While it can, it does not have to be defined apriori for all future looks.

TreeScan will automatically remember how much alpha has already been spent, and this is indicated on the Sequential Tab.

Related Topics: [Analysis Tab](#), [Bernoulli Probability Model](#), [Sequential Analysis](#).

Power Evaluation Tab



Power Evaluation Tab Dialog Box

This tab is reached by clicking the Advanced button in the lower right corner of the Analysis Tab.

Statistical Power Evaluation

The power of the tree-based scan statistic varies for different branches and leaves on the tree. The higher the expected counts in a Poisson analysis, the higher the power is for a specific relative risk. Likewise for a fixed relative risk, for the Bernoulli model, the power is higher if the number of cases and controls is higher for a branch or leaf. Of course, power also increases with increased relative risk, and it also depends on the size and nature of the tree. For a fixed attributable risk, that is, a fixed expected excess number of cases, the power is higher for a node or leaf with fewer expected cases in a Poisson model and fewer controls in a Bernoulli model.

With the TreeScan software, it is possible to estimate the power either as part of a regular analysis or as a stand-alone exercise. The feature is available for the conditional and unconditional Poisson models, for the unconditional and conditional Bernoulli models, for the purely temporal scan statistic, and for the tree-temporal scan statistic that is conditioned on the node but not on time. This feature is not available for the tree-temporal scan statistic that is conditioned on both the node and time.

If the power evaluation is done as part of a regular analysis, the estimated power will be added to the results file. If only a power evaluation is requested, the software will not evaluate any cuts on the tree and not produce any likelihood ratios, relative risks, p-values, etc. If only a power

evaluation is performed for the conditional Poisson model, the user has the option to either condition the analysis on the total number of cases in the count file or to specify a different total number of cases on the Power Evaluation Tab.

The statistical power is estimated using data that are simulated under both the null and the alternative hypothesis. The number of simulated replications of the data needs to be specified. It is recommended to use at least 999 replications under the null hypothesis and at least 1000 replications under the alternative hypothesis, but more is always better in order to increase the precision of the power estimates.

Alternative Hypothesis: The alternative hypothesis is defined in the Alternative Hypothesis File, which is a plain text file that consists of two columns when doing a Tree Only analysis. In the first column, specify the node ID for which there should be an excess risk. In the second column, specify the excess risk. For the Poisson model, specify it in terms of the relative risk. For the Bernoulli model, specify it in terms of a probability. For the alternative hypothesis, it is possible to specify an excess risk for multiple nodes. All of them will then have an excess risk. This is done by having multiple rows in the Alternative Hypothesis File, without an empty row in between. Note that the excess risk can be different for different nodes within the same alternative hypothesis, except for the conditional Bernoulli model, where the excess risk must be the same for all nodes with an excess risk.

For a tree-temporal analysis conditioned on the node only, the alternative hypothesis file consists of four columns. The first column is the node and the second column is the relative risk. The third and fourth columns are the start and end times of the time interval that has that excess relative risk compared to the times without excess risk for that same node. Since the analyses are conditioned on the node, there is no need to specify relative risks at one node versus the other nodes. The same node can have two lines with, for example, a relative risk of 2 for times 10 to 15 and a relative risk of 3 for times 16 to 20. If a relative risk of 2 is specified for times 10 to 16 and a relative risk of 3 for times 14 to 20, then the times 14 to 16 will have a relative risk of $2 \times 3 = 6$ compared to the times without any excess risk.

For a purely temporal power evaluation, the alternative hypothesis file has one or more lines and three columns. The first column is the relative risk. The second and third columns are the start and end times of the time interval that has that excess relative risk compared to the times without any excess risk. If a time is included in multiple rows, the total relative risk is the product of the specified relative risks.

Multiple alternative hypotheses can be evaluated within the same TreeScan run. This is done by adding additional rows to the Alternative Hypothesis File, keeping an empty line between the different alternative hypotheses. It is computationally much faster to do it this way than to conduct a separate run for each alternative, since the different alternative can use the same set of data that is generated under the null hypothesis.

Excess Risk versus Less Risk: For the Poisson model, it is typical to specify an alternative hypothesis with a relative risk greater than one, and for the Bernoulli model, a probability that is greater than the probability specified on the main Analysis Tab. This ensures that the alternative

hypothesis contain an excess risk somewhere on the tree. Technically, it is also possible to specify smaller numbers, but that does not make sense when one is looking for clusters of excess risk.

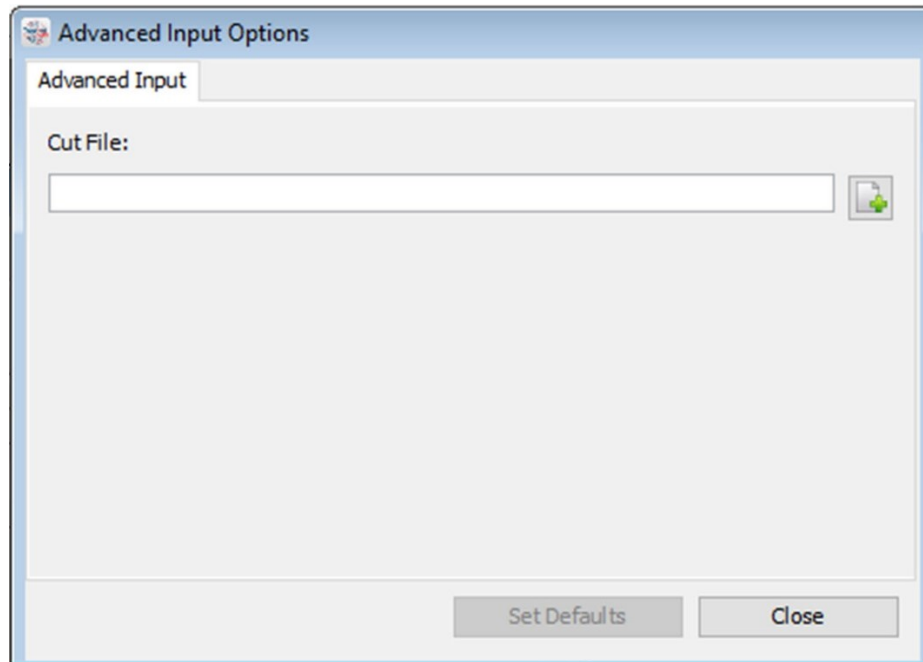
Conditional Poisson Model: The power estimation for the conditional Poisson model can theoretically be done in either of two ways, by either conditioning or by not conditioning the power estimation on the total number of cases that were actually observed. For example, suppose that we have a tree with 100 leaves each of which under the null hypothesis has 1 expected case. Moreover, suppose that we want to know the power if one node with 10 leaves has a relative risk of 2, so that the expected number of cases in that node is 20. Now, under the alternative hypothesis, we would sometimes observe 110 total cases, sometime 105 total cases and sometimes 115 total cases, etc. Suppose we observe 105 total cases. The question is then, do we want to know what the power is considering that there were 105 total cases observed, or do we want to know what the power is ignoring the total number of cases that were actually observed. The true power will be slightly different for the two different approaches, since the power is slightly higher when a few more total cases are observed. TreeScan uses the former approach. If the actual total number of cases (=sample size) is already known when the power evaluation is performed, use that number. If the actual number is not already known, use the approximate sample size that you expect to have.

One advantage of the approach that TreeScan uses is that it is not necessary to specify the actual magnitude of the expected counts for each node under the null hypothesis, but only the relative size of those expected counts. It is also computationally faster to calculate.

Conditional Bernoulli Model: As for the conditional Poisson model, the power for the conditional Bernoulli model is also calculated conditioned on the total number of cases that were actually observed in the data. For the Bernoulli model power calculation though, it is necessary to specify the exact probabilities under the null and alternative hypotheses, something which is not needed when doing the actual analysis using this probability model. That is, it is necessary to specify the probability of a case under the null hypothesis and the probability of a case under the alternative hypothesis. This can be a little tricky, since those probabilities under the alternative hypothesis must be such that they can realistically generate the total number of cases in the data. If they are off by too much, TreeScan will generate an error message. The problem can easily be solved by either changing the total number of cases or the underlying probabilities.

Related Topics: [Analysis Tab](#).

Advanced Input Tab



Advanced Input Tab Dialog Box

This tab is reached by clicking the Advanced button in the lower right corner of the Input Tab.

Complex Cuts

As the default, TreeScan uses simple cuts on the tree, where the cut is on a single edge, defining the branch that is a potential cluster. When a node has more than two children, it is also possible to perform pair cuts, triplet cuts or ordinal cuts, as described below. The type of cuts requested is specified in the Cut File, and the name of that file is specified on the Advanced Input Tab.

Suppose a node has five children. With simple cuts below the node, each child and all of its descendants is a cut defining a branch with potential cluster. With a simple cut above the node, all five children and all their descendants is also a branch defining a potential cluster. Two of the children together, but without the other three siblings, cannot define a potential cluster though, when using simple cuts. When using pair cuts, pairs of two siblings with their descendants will also define a potential cluster, consisting of two of the five branches emanating from that node. A group of three siblings with their descendants will not be a candidate for a cluster, when a pair cut is used, but such trios will be considered if triplet cuts are used. When a triplet cut is specified, pair and simple cuts are also utilized. When a pair cut is specified, simple cuts are also used.

With an ordinal cut, there needs to be an order on the children, and we only consider pairs, triplets or larger groups of children that are next to each other in the ordering. For example, suppose the five children are A, B, C, D and E. As potential clusters, we would then consider for example [A,B], [B,C], [C,D,E], and [A,B,C,D], but not [A,C] or [B,E]. The order used for the children is alphabetical. If you want to order the children in some other way, the NodeID's need to be

renamed. For example, if the original names were VeryLow, Low, Medium, High and Very High, that can be renamed as 1VeryLow, 2Low, 3Medium, 4High and 5Very High. Whether the order is on one direction or the opposite does not matter, generating the same analysis and the same results.

Related Topics: [Advanced Features](#), [Input Tab](#), [Cut File](#).

Additional Output Tab

Advanced Output Options

Additional Output

Attributable Risk

☐ Report attributable risk based on exposed.

Log Likelihood Ratios

☐ Report Simulated Log Likelihood Ratios

Critical Values

☐ Report critical values for an observed cut to be significant

Temporal Graphs

☐ Produce Temporal Graphs

☒ Most likely cluster only

☐ 1 most likely clusters, one graph for each

☐ All significant clusters, one graph for each, with p-value less than:

Set Defaults Close

Additional Output Tab Dialog Box

This tab is reached by clicking the Advanced button in the lower right corner of the Output Tab.

Attributable Risk

For some analysis, it is of interest to calculate the attributable risk in addition to the relative risk. For this to work, it is necessary to specify the total number of exposed individuals, including those that may not have contributed any cases to the tree.

Related Topics: [Output Tab](#).

Simulated Log Likelihood Ratios

TreeScan can print all the simulated log likelihood ratios to a comma delimited output file. This output is not needed for running or interpreting a TreeScan analysis. The option is available for statistical researchers who want to study the distributional properties of the tree-based scan statistic. The name of this output file is the same as the name of the text output format file, but with a *.llr.csv filename extension.

Related Topics: [Output Tab.](#)

Critical Values

If requested, TreeScan will provide the critical values for rejecting the null hypothesis at the $\alpha = 0.05, 0.01$ and 0.001 levels. This output is not needed for running or interpreting a TreeScan analysis.


Related Topics: [Output Tab.](#)

Temporal Graphs

When a tree only or tree time scan is conducted, selecting this option allows TreeScan to produce temporal graphs depicting the observed and expected counts over time, both inside and outside the cluster. It will also show the ratio of observed over expected counts. You can select whether to only produce a graph for the most likely cut, or if multiple graphs for a fixed number of cuts or all cuts with a p-value less than some specified value should be produced. The graphs are generated as HTML files that can be opened in any web browser. The name and location of the file is listed in the parameter section of the standard results file. Once you have opened the HTML file, you can edit it and also generate the temporal graphs in PNG, JPEG, PDF and SVG formats. These other formats do not need to be pre-specified.

Running TreeScan

Specifying Analysis and Data Options

The TreeScan program requires that you specify parameters defining analysis, input and output options for the analysis you wish to conduct. A tabbed dialog is provided for this purpose. To access the parameter tab dialog, either press the  button or select the File/New menu item. Once the tabbed dialog has been opened, specify the parameters for your session on the following three main tabs:

- Analysis Tab
- Input Tab
- Output Tab


See the section on Basic TreeScan Features for instructions on how to fill in these tabs.

Most analyses can be performed using only these three tabs. For each tab, there are additional features that can be selected by first clicking on the Advanced button in the lower right corner of the tab. These additional features may be useful in special circumstances.

The available choices for some features may depend on what was selected in other places, and they will then be deactivated accordingly.

Related Topics: [Basic TreeScan Features](#), [Sample Data Sets](#), [Test Run](#), [Input Tab](#), [Analysis Tab](#), [Output Tab](#), [Advanced Features](#), [Launching the Analysis](#).

Launching the Analysis

Once the data input files have been created, and the parameters defining the analysis input and output options have been specified, select the Start Analysis  button to launch the analysis and produce the results file. Once the analysis has been completed, the standard results file will appear in the job status window.

Multiple parameter session windows may be opened simultaneously for data entry, and multiple analyses may be run concurrently. If you are running multiple analyses concurrently, please make sure that the output files have different names, or they will overwrite each other.

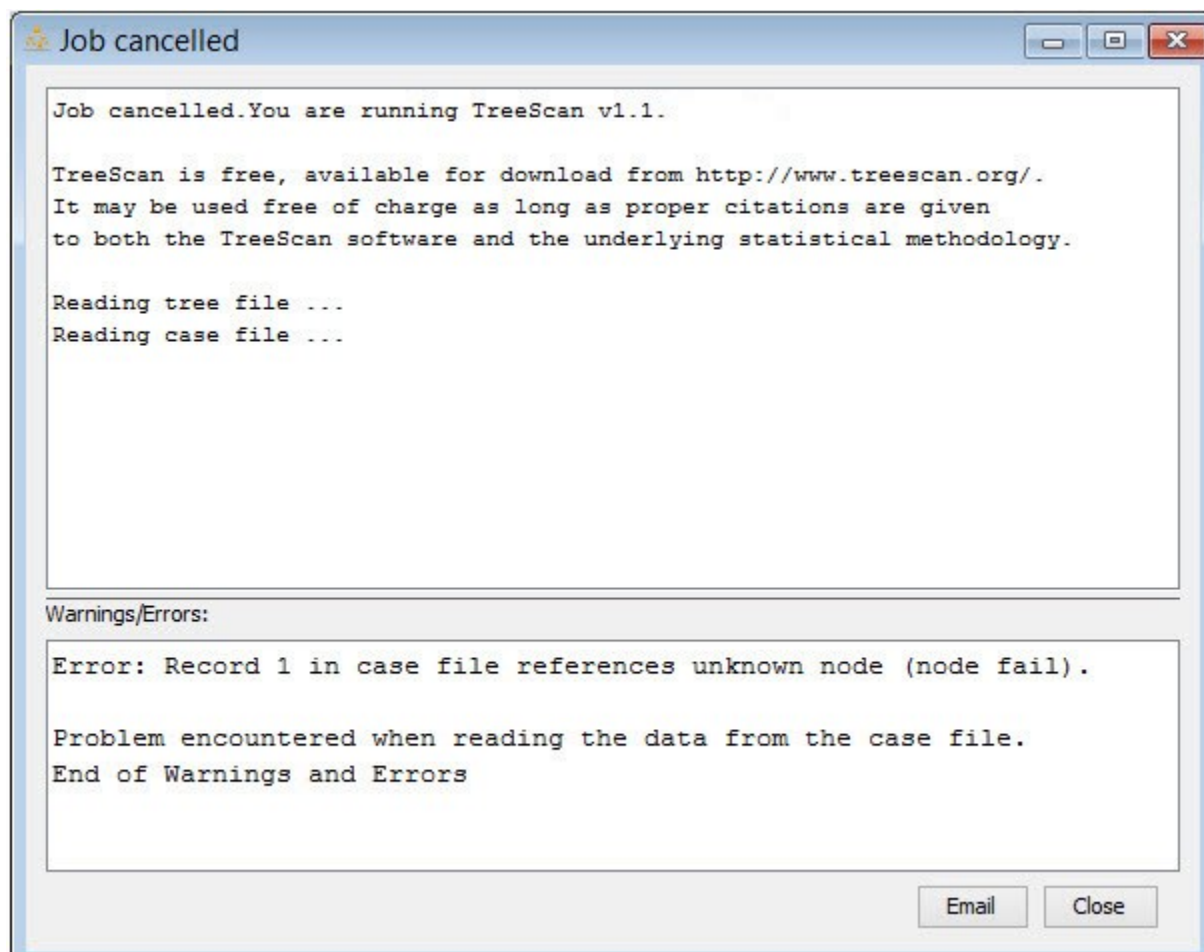
Related Topics: [Input Data](#), [Test Run](#), [Specifying Analysis and Data Options](#), [Status Messages, Warnings and Errors](#), [Computing Time](#), [Batch Mode](#).

Status Messages

Status messages are displayed as the program executes the analysis, as the data is read, and at each step of the analysis. Normal status messages are displayed in the top box of the job status window. Warnings and error messages are displayed in the bottom box of the job status window. Upon successful completion of the calculations, the standard results file will be shown in the job status window.

Related Topics: [Launching the Analysis](#), [Warnings and Errors](#).

Warnings and Errors



TreeScan Status Messages and Warnings/Errors Dialog Box

Warning Messages

TreeScan may produce warnings as the job is executing. A message is then displayed in the Warnings/Errors box on the bottom of the job status window. A warning will not stop the execution of the analysis. If a warning occurs, please review the message and access the help system if further information is required.

Error Messages

If a serious problem occurs during the run, an error message will be displayed in the Warnings/Errors box on the bottom of the job status window and the job will be terminated. The user may resolve most errors by reviewing the message and using the help system.

One of the most common errors is that the input files are not in the required format, or that the file contents are incompatible with each other. When this occurs, an error message will be shown specifying the nature and location of the problem. Such error messages are designed to help with data cleaning.

If the error message cannot be resolved, you may press the email button on the job status window. This will generate an automatic email message to TreeScan technical support. The contents of the “Warnings/Errors” box will be automatically placed in the e-mail message. All a user needs to do is press their e-mail Send key. Users may also print the contents of the Warnings/Errors box and even select, copy (ctrl c) and paste (ctrl v) the contents if necessary.

Related Topics: [Input Data](#), [Data Requirements](#), [TreeScan Support](#).

Saving Analysis Parameters


Analysis parameters, specified on the Parameter tab dialog, can be saved and reused for future analyses. It is recommended that you save the parameters with a “.prm” file extension. The parameter file is stored in an ASCII text file format.

Save Analysis Parameters

1. If the parameters have not previously been saved, select Save As from the File menu. A ‘Save Parameter File As’ dialog will open.
2. Select a directory location from the ‘Look In’ drop-down menu at the top of the dialog box.
3. Enter a name for your parameter file in the ‘File Name’ text box. It is recommended that the ‘Save As Type’ selection remain as Parameter Files (*.prm).
4. Press the Save button.

Once the parameter file is initially saved, save changes to the file by selecting ‘Save’ on the File menu. The file will save without opening the ‘Save Parameter File As’ dialog.

Open a Saved Parameter File

1. Select 'Open' from the File menu or click on the  button in the toolbar. A Select Parameter File dialog will open.
2. Locate the desired file using the Look in drop-down menu.
3. Once the file is located, highlight the file name by clicking on it.
4. Press the Open button.

A Parameter tab dialog will open containing the saved parameter settings. The location and name of the parameter file is listed in the title bar of this dialog.

Related Topics: [Specifying Analysis and Data Options](#), [Basic TreeScan Features](#), [Advanced Features](#), [Batch Mode](#).

Sequential Analysis

For sequential analyses, TreeScan is run multiple times, as the sample size grows with additional data. Each run, for each of the looks, is run just as a regular TreeScan analysis, but the analysis parameters have to stay the same, except for the amount of alpha to spend at the current looks, which can be changed over time.

It is important that the name of the results file stays the same. If not, TreeScan will think that it is a new sequential analysis that is strating. This also means that if one runs multiple sequential analyses in tandem, for different data sets or outcomes, then each one must have a different output file name.


Related Topics: [Sequential Analysis Tab](#).

Parallel Processors

If you have parallel processors on your computer, TreeScan will take advantage of this by running different Monte Carlo simulations using different processors, thereby increasing the speed of the calculations. The default is that TreeScan will use all processors that the computer has. If you want to restrict the number, you can do that by clicking on Session > Execute Options, and selecting the maximum number of processors that TreeScan is allowed to use.

Related Topics: [Analysis Results](#).

Batch Mode

TreeScan is most easily run by clicking the Start Analysis  button at the top of the TreeScan window, after filling out the various parameter fields in the Windows interface.

An alternative approach is to skip the windows interface and launch the TreeScan calculation engine directly by either:

1. Dragging a parameter file onto the 'treescan32.exe' executable.
2. Writing 'treescan32.exe *.prm' in a batch file or at the command prompt, where *.prm is the name of the parameter file.

By using the batch mode version, it is possible to write special software that incorporates the TreeScan calculation engine with other applications. To use TreeScan in this manner requires a reasonable amount of computer skill and sophistication.

When running TreeScan in batch mode, it is easiest to first create the parameter file using the TreeScan windows interface, and then save that file. Changes can be made in the same way, but it is also possible to change the parameter manually using any text editor or automatically by using some other software product. If only a few parameters should change compared to what is in an existing parameter file, name the parameter file on the command prompt together with instruction on which parameters should change. The command line parameter values will then over-ride the parameter values specified in the parameter file.

When the batch mode version of TreeScan is run, the standard results file does not automatically pop up on the screen, but must be opened manually using any available text editor such as Notepad.

Related Topics: [Launching the Analysis](#), [Basic TreeScan Features](#), [Advanced Features](#), [Saving Analysis Parameters](#).

Computing Time

The tree-based scan statistic can be computer intensive to calculate. The computing time depends on a wide variety of variables, and depending on the data set and the analytical options chosen; it could range from a few seconds to several days or weeks. The main variables that increase the computing time are the number of nodes and the total number of observed cases. For the tree-temporal scan statistic, the number of time periods and the range of the window start and end times will also influence the computing time. Unconditional analyses are typically faster than conditional analyses.

Related Topics: [Memory Requirements](#).

Memory Requirements

TreeScan uses dynamic memory allocation. If there is insufficient memory available on the computer to run the analysis, there are several options available for working around the limitation:

- Close other applications.
- Aggregate the data into fewer nodes.
- Aggregate the temporal data into fewer time periods.
- Run the program on a computer with more memory.

It is highly desirable that there is sufficient RAM to cover all the memory needs, as TreeScan runs considerably slower when the swap file is used, so these techniques may also be used to avoid the swap file.

Related Topics: [Computing Time](#), [Warnings and Errors](#).

Results

As output, TreeScan always creates a standard plain text based results file in ASCII format. This file is automatically shown after the analysis is complete. It can also be opened using any text editor, such as NotePad or Word. Two additional optional output files with the same information can be generated in HTML and comma delimited CSV formats respectively. If the former is requested, it is automatically launched when the analysis is done.

Related Topics: [Output Tab](#).

Standard Results File (*.txt)

The standard results file is automatically shown after the calculations are completed. It is fairly self-explanatory, but for proper interpretation it is recommended to read the section on statistical methodology, or even better, one of the methodology papers listed in the bibliography.

SUMMARY OF DATA

Provides data concerning the tree as a whole. For example, it includes information about the size and depth of the tree and the total number of cases. Use this to check that the input data files contain the data that you expect it to contain. For a sequential analysis, the summary will include the amount of alpha spent to date, out of the total alpha available.

MOST LIKELY CUTS

Summary information about the most likely cuts, that is, the tree branches with the cluster of cases that are least likely to be due to chance. They are presented in order by their likelihood, from larger to smaller; that is, from stronger to weaker clusters.

Node Identifier: This is the node below where the cut was made.

Node Cases: The total number of cases in the node regardless of whether they are inside or outside the time window of the cluster. Only reported for the tree-temporal scan statistic, as it is otherwise identical to the number of cases.

Time Window: The temporal window of the cluster. Only reported for the tree-temporal and purely temporal scan statistics.

Cases: The number of observed cases in the cluster. For a purely tree-based scan statistic, this is the number of observed cases in the node where the cut was made. For the tree-temporal scan statistic, it is the number of observed cases in the node that are also within the time window of the cluster.

Expected: The number of expected cases in the cluster, under the assumption that the null hypothesis is true, and conditioned on whatever the analysis is conditioned on. The mathematical formulas are provided below.

Relative Risk: The risk in the cluster divided by the risk outside the cluster. For the Poisson and Bernoulli models, outside means for other parts of the tree. For the tree-temporal and purely temporal scan statistics, outside means outside the cluster time window but inside the same node. The formula is different for different probability models. See below.

Excess Cases: The excess number of cases. The formula is different for different probability models. See below.

Attributable Risk (optional): The attributable risk is the excess number of cases divided by the total number of people exposed that were specified on the Additional Output Tab.

Log Likelihood Ratio: The natural logarithm of the likelihood ratio for the cut. This is the test statistic, and a larger value is evidence for a cluster and against the null hypothesis.

P-value: The p-values are adjusted for the multiple testing stemming from the multitude of cuts evaluated. This means that under the null-hypothesis of complete randomness there is a 5% chance that the p-value for the most likely cut will be smaller than 0.05 and a 95% chance that it will be bigger. Under the null hypothesis there will always be some area with a rate higher than expected just by chance alone. Hence, even though the most likely cut always has an excess rate, the p-value may actually be very close or identical to one.

Signalled: For a sequential analysis, no p-values are provided. Instead, it is indicated if a node has signalled, which means that the null hypothesis has been rejected because of that node. If the node has signaled, it will indicate at which look it did signal. Once a node has signaled, it can never 'unsignal' at future, even if the relative risk goes below one or the excess cases goes below zero.

PARAMETER SETTINGS

A list of the parameter settings used for the analysis.

COMPUTATIONAL INFORMATION

Information about when the program was run, how long it took to execute, and the number of processors used.

Related Topics: [Output Tab](#), [Additional Output Tab](#), [Mathematical Formulas](#), [HTML Results File](#), [Comma Delimited Results File](#).

Mathematical Formulas

Here we give the mathematical formulas for various result variables. First we define the following notation:

c = Number of cases in the cluster

C = Total number of cases, in the whole tree

n = Population in the node

N = Total population in the whole tree

p = Case probability (for the unconditional Bernoulli model)

w = Length of the temporal cluster window (for tree-temporal and temporal scan statistics)

T = Length of the data time range (for tree-temporal and temporal scan statistics)

z = Number of cases in the temporal cluster window, summed over the whole tree

Depending on the probability model used, the mathematical formulas are:

Type of Scan	Model	Conditional	Expected	Relative Risk	Excess Cases
Tree Only	Poisson	No	n	$\frac{c}{n}$	$c - n$
Tree Only	Poisson	Total cases	$n \frac{C}{N}$	$\frac{c/n}{(C - c)/(N - n)}$	$c - n \frac{(C - c)}{(N - n)}$
Tree Only	Bernoulli	No, standard	np	$\frac{c}{np}$	$c - np$
Tree Only	Bernoulli	No, self-control	np	$\frac{c/p}{(n - c)/(1 - p)}$	$c - p \frac{(n - c)}{(1 - p)}$
Tree Only	Bernoulli	Total cases	$n \frac{C}{N}$	$\frac{c/n}{(C - c)/(N - n)}$	$c - n \frac{(C - c)}{(N - n)}$
Tree-Temporal	Uniform	Node	$\frac{wn}{T}$	$\frac{c/w}{(n - c)/(T - w)}$	$c - w \frac{(n - c)}{(T - w)}$
Tree-Temporal	Uniform	Node and time	$\frac{nz}{C}$	$\frac{c(C - n - z + c)}{(n - c)(z - c)}$	$c - \frac{(n - c)(z - c)}{(C - n - z + c)}$
Purely Temporal	Uniform	Total cases	$\frac{wC}{T}$	$\frac{c/w}{(C - c)/(T - w)}$	$c - w \frac{(C - c)}{(T - w)}$

Related Topics: [Output Tab](#), [Analysis Results](#), [Standard Results File](#).

HTML Results File (*.html)

The HTML results file has the same information as the standard ASCII results file, but it is presented as an HTML table. If selected, it is automatically launched when the analysis is completed. The name of the file is the same as the standard results file, except for the *.html extension.

Tree Visualization

In addition to the text and table output formats, TreeScan can also show the results as a tree, reflecting the tree-based structure of the data. This tree visualization is done after the analysis is complete, by clicking the Tree Visualization button at the bottom of the HTML output file. For large trees with many leaves and nodes, it takes a little bit of time and patience to generate the Tree Visualization. Once the tree is up, one can choose to show or hide different branches of the tree. The leaves and node of the tree are colored according to p-values, or if preferred, according to the relative risks. When clicking on a node, detailed information about that node is shown.

Related Topics: [Output Tab](#), [Analysis Results](#), [Standard Results File](#), [Comma Delimited Results File](#).

Comma Delimited Results File (*.csv)

The Comma Delimited results file takes the information from each of the most likely cuts, presented as a comma delimited CSV table. This file is useful for importing the TreeScan output into other software. It needs to be opened manually. The columns are:

Cut: The index number of the most likely cuts, in decreasing order with respect to log likelihood ratios.

Node Identifier: The node identifier.

Node Cases: The total number of cases in the node. This is only reported for the tree-temporal scan statistic as it is identical to Cases for the other scan types.

Time Window Start: The start time of the detected cluster window. Only reported for the tree-temporal and purely temporal scan statistics.

Time Window End: The end time of the detected cluster window. Only reported for the tree-temporal and purely temporal scan statistics.

Cases: The observed number of cases in the cluster.

Expected Cases: The expected number of cases in the cluster, when the null hypothesis is true, and conditioned on whatever the analysis was conditioned on.

Relative Risk: Relative risk.

Excess Cases: The excess number of cases.

Attributable Risk: Attributable risk.

Log Likelihood Ratio: The log likelihood ratio.

P-value: The p-value for the cluster.

The name of the file is the same as the standard results file, except for the *.csv extension.

Related Topics: [Output Tab](#), [Analysis Results](#), [Standard Results File](#), [HTML Results File](#).


Simulated Log Likelihood Ratios File (*_llr.csv)

The log likelihood ratio test statistics that were calculated from each of the random data sets are not provided as part of the standard output. As an advanced option, they can be printed to a special file. There is no need for this file when doing a regular analysis, but it can be interesting for statistical researchers who are studying the distributional properties of the tree-based scan statistic under various scenarios. By default it has the same name as the output file but with the extension *_llr.csv

Related Topics: [Output Tab](#), [Analysis Results](#).

Miscellaneous

New Versions

To check whether there is a later version of TreeScan than the one you are currently using, simply click on the update button  on the tool bar. If a newer version exists, you will be asked whether you want to download and install it. You can request that TreeScan automatically checks for new versions once a week, once a month or every time TreeScan is used. Alternatively, you can set TreeScan to only check for new versions manually, when you decide to do so.

At any given time, it is also possible to download the latest version of TreeScan at 'www.treescan.org'.

Related Topics: [Download and Installation](#).

Random Number Generator

The choice of random number generator is critical for any software creating simulated data. TreeScan uses a Lehmer random number generator¹² with modulus $2^{31}-1 = 2147483647$ and multiplier 48271, which is known to perform well.¹⁴

Related Topics: [Monte Carlo Replications](#).

Contact Us

Please direct technical questions about installation and running the program, as well as the web site, to:

techsupport@treescan.org

Please direct substantive questions about the statistical methods and suggestions about new features to:

Martin Kulldorff

Department of Medicine, Division of Pharmacoepidemiology and Pharmacoeconomics

Harvard Medical School and Brigham and Women's Hospital

1620 Tremont Street, suite 3030, Boston, MA 02120, USA

Email: kulldorff@treescan.org

Acknowledgements

Financial Support

National Institutes of Health, National Library of Medicine, through grant #RC1LM010371 [TreeScan v1.0]

Food and Drug Administration, Center for Biologics Evaluation and Research, Division of Epidemiology and Biostatistics, Mini-Sentinel Post-licensure Rapid Immunization Safety Monitoring (PRISM) Program [TreeScan v1.1, 1.2, 1.3]

Food and Drug Administration, Center for Drug Evaluation and Research, Sentinel Program [TreeScan v2.0]

Their financial support is greatly appreciated. The contents of TreeScan are the responsibility of the developer and do not necessarily reflect the official views of funders.

Comments and Suggestions

Feedback from users is greatly appreciated. Very valuable suggestions concerning the TreeScan software have been received from various individuals, including:

Judith Maro, Harvard Medical School and Harvard Pilgrim Health Care Institute

Michael Nguyen, Food and Drug Administration

Shirley Wang, Harvard Medical School and Brigham and Women's Hospital

Katherine Yih, Harvard Medical School and Harvard Pilgrim Health Care Institute

Logo

The TreeScan logo was created by Bulkhead Design (www.bulkheaddesign.com).



Frequently Asked Questions

Input Data

1. I tried running TreeScan using one of the sample data sets and all went well, but when I try it on my own data there is an error. What should I do?

TreeScan makes sure that the input data is compatible with each other, and with the options specified on the windows interface. For example, it complains if a node ID in the count file is not present in the tree file, because it must know where to assign those cases. For most data sets there is some need for data cleaning and TreeScan is designed to help with this process by spotting and pointing out any inconsistencies found.

2. I have constructed the ASCII input files exactly according to the description in the TreeScan User Guide, but TreeScan complains that they are not in the correct format. What is wrong?

The most likely explanation is that the files are in UNICODE rather than ASCII format. Just convert to ASCII and it should work.

3. In my data, there is zero or only one case in many nodes. Can I use TreeScan for such sparse data?

Yes, you certainly can. One of the main reasons for using TreeScan is to avoid arbitrary aggregation of the data, letting the scan statistic consider different smaller or larger aggregations by considering different cuts on the tree. With finer resolution of the input data, TreeScan can evaluate more different cluster locations and sizes without restrictions imposed by higher level groupings.

4. What is the minimum size of the tree needed to run TreeScan?

Technically, the tree-based scan statistic can be run using only two nodes, providing correct inference. Unless there is a time component, there is no point using scan statistic for such data though, for which a regular chi-square statistic can be used instead, as there is no multiple testing to adjust for. With three nodes or more, the fundamental scan statistic concept of including different combinations of locations into the potential clusters is being utilized. In most practical applications though, the tree-based scan statistic is used for data sets with hundreds or thousands of nodes.

Regarding the tree-temporal scan statistic, it can make sense to run it with only two nodes on the tree, in which case it will look for temporal clusters in either or both of the nodes.

5. When should I use the Bernoulli versus the Poisson model?
Use the Bernoulli model when you have binary data, such as cases and controls in a matched design or exposed and unexposed time periods in a self-control design. Use the Poisson model when you have cases and expected counts from a historical or concurrent comparison population at risk.
6. I have memory problems when running TreeScan. What should I do?
Make sure you are running TreeScan in 64-bit mode. For this you must (i) have a 64 bit computer, and (ii) have 64-bit Java installed on your computer.

Results

7. I get an error stating that the output file could not be created. Why?
In Windows, permission to write to the "Program Files" folder is given only to administrators and power users of that machine. If the output file path includes the "Program Files" folder and you do not have administrative or power user privileges on your computer, Windows prevents TreeScan from creating the output file in the designated location. The solution is to specify a different directory for the output file.
8. Since the TreeScan results are based on Monte Carlo simulated random data, why are the p-values the same when I run the analysis twice?
All computer-based simulations are based on pseudo-random number generators. When the same seed is used, exactly the same sequence of pseudo-random numbers will be generated. Since TreeScan uses the same seed for every run, you obtain the same result for two runs when the input data is the same.
9. In the results for a Poisson analysis, the population is smaller than the number of cases. How can that be?

When using the Poisson model, the 'population' could be a variety of things. If it is the total number of people exposed to a drug, then the population must be greater or equal to the number of cases. Population could also be expected counts, calculated using some other statistical software; it could be person-years; or it could be the number of million people exposed. In these three scenarios, the observed cases can be less than the population. Note that, with the conditional Poisson model, results will be the same if you multiply all population counts with the same constant.

Operating Systems

10. Is TreeScan available for Windows/Mac/Linux?

The TreeScan software is available for Linux, Mac and Windows, and all three versions can be downloaded from the www.treescan.org web site.

TreeScan Bibliography

Suggested Citations

The TreeScan software may be used freely, with the requirement that proper references are provided to the scientific papers describing the statistical methods. Depending on the application, the suggested citations can be found among the methodological papers below.

Methodology

General Statistical Theory

1. Kulldorff M, Fang Z, Walsh S. A tree-based scan statistic for database disease surveillance. *Biometrics*, 2003,59:323-331.
2. Kulldorff M, TreeScan User Guide, v2.0, 2020.

Statistical Power Evaluation

3. Maro JC, Nguyen MD, Dashevsky I, Baker MA, Kulldorff M. Statistical Power for Postlicensure Medical Product Safety Data Mining. *eGEMs: The Journal for Electronic Health Data and Records*, 5:6, 2017.

Selected Applications

Drug Safety Surveillance

4. Kulldorff M, Dashevsky I, Avery TR, Chan KA, Davis RL, Graham D, Platt R, Andrade SE, Boudreau D, Gunter MJ, Herrinton LJ, Pawloski P, Raebel MA, Roblin D, Brown JS. Drug safety data mining with a tree-based scan statistic. *Pharmacoepidemiology and Drug Safety*, 2013, 22:517-523.
5. Brown JS, Petronis KR, Bate A, Zhang F, Dashevsky I, Kulldorff M, Avery TR, Davis RL, Chan KA, Andrade SE, Boudreau D, Gunter MJ, Herrinton L, Pawllowski PA, Raebel MA, Roblin D, Smith D, Reynolds R. Drug adverse event detection in health plan data using the gamma Poisson shrinker and comparison to the tree-based scan statistic. *Pharmaceutics*, 2013,5:179-200.
6. Schachterle SE, Hurley S, Liu Q, Petronis KR, Bate A. An Implementation and Visualization of the Tree-Based Scan Statistic for Safety Event Monitoring in Longitudinal Electronic Health Data. *Drug Safety*, 42:727-741, 2019

Vaccine Safety Surveillance

7. Yih WK, Maro JC, Nguyen M, Baker MA, Balsbaugh C, Cole DV, Dashevsky I, Mba-Jonas A, Kulldorff M. Assessment of quadrivalent human papillomavirus vaccine safety using the self-controlled tree-temporal scan statistic signal-detection method in the Sentinel system. *American Journal of Epidemiology*, 187:1269-1276, 2018.
8. Li R, Weintraub E, McNeil MM, Kulldorff M, Lewis EM, Nelson J, Xu S, Qian L, Klein NP, Destefano F. Meningococcal conjugate vaccine safety surveillance in the Vaccine Safety Datalink using a tree-temporal scan data mining method. *Pharmacoepidemiology and Drug Safety*, 27:391-397, 2018
9. Yih WK, Kulldorff M, Dashevsky I, Maro JC. Using the Self-Controlled Tree-Temporal Scan Statistic to Assess the Safety of Live Attenuated Herpes Zoster Vaccine. *American Journal of Epidemiology*, 188:1383-1388, 2019.

Occupational Disease Surveillance

10. Kulldorff M, Fang Z, Walsh S. A tree-based scan statistic for database disease surveillance. *Biometrics*, 2003,59:323-331.

Manufacturing

11. Mahaux O, Bauchau V, Zeinoun Z, Van Holle L. Tree-based scan statistic–Application in manufacturing-related safety signal detection. *Vaccine*, 37:49-55, 2019.

History

12. Alfani G, Gráda C. The timing and causes of famines in Europe. *Nature Sustainability*, 1:283-288, 2018

Other References Mentioned in this User Guide

13. Breiman L, Friedman J, Stone CJ, Olshen RA. *Classification and regression trees*. CRC press, 1984.
14. Dwass M. Modified randomization tests for nonparametric hypotheses. *Annals of Mathematical Statistics*, 1957; 28:181-187
15. Glaz J, Balakrishnan N (editors). *Scan Statistics and Applications*. Birkhäuser: Boston, 1999.
16. Glaz J, Naus JI, Wallenstein S. *Scan Statistics*. Springer Verlag: New York, 2001.
17. Glaz J, Pozdnyakov V, Wallenstein S (editors). *Scan Statistics: Theory and Applications*. Birkhäuser: Boston, 2009.

18. Lehmer DH. Mathematical methods in large-scale computing units. In Proceedings of the second symposium on large scale digital computing machinery. Cambridge, USA: Harvard Univ. Press, 1951.
19. Naus J. The distribution of the size of maximum cluster of points on the line. Journal of the American Statistical Association, 60:532-538, 1965.
20. Park SK, Miller KW. Random number generators: Good ones are hard to find. Communications of the ACM, 31:1192-1201, 1988.

Index

- Additional Output Files, 29, 30
- Advanced Features, 31
- Analysis
 - Launching, 51
 - Saving Parameters, 53
 - Specifying Options, 51
- Analysis Tab, 22
 - Probability Model, 23, 24
- Ancestors, 9
- Basic TreeScan Features, 22
- Batch Mode, 55
- Bernoulli Distribution, 23
- Bernoulli Model, 7, 11
- Bernoulli versus Poisson model, 65
- Branch, 10
- Case File, 16, 26
 - Format, 16
- Cases, 57
- Child, 9
- Citations
 - Suggested, 66
- Classification and Regression Trees (CART), 15
- Cluster, 10
- Clusters, 58
 - Maximum Temporal Size, 34, 35
 - Most Likely, 14
 - Secondary**, 58
- Comma Delimited Results File (*.csv), 60
- Computational Information, 58
- Computing Time, 55
- Contact Us, 62
- Control File, 18, 19, 26
 - Format, 18, 19
- Coordinates File, 16
- Cut**, 10
- cut type, 19
- Data
 - Control, 26
 - Requirements, 16
- Data Granularity, 4
- Data Requirements
 - Data Requirements, 16
- Descendants, 9
- Drug Safety Surveillance, 66
- Edge, 9
- End Date, 26
- Excess Cases, 58
- Expected, 57
- Frequently Asked Questions, 64
 - Analysis, 64
 - Input Data, 64
 - Results, 65

- Help System, 6
- HTML Results File (*.html), 60
- Import File
 - SaTScan Import Wizard, 19
- Inference Tab, 32, 36, 38, 44
- Input Data, 16, 64
 - Case File, 16
 - Control File, 18, 19
 - Coordinates File, 16
- Input File
 - SaTScan ASCII File Format, 20
- Input Files Tab, 51
- Input Tab, 25
 - Case File, 26
 - Study Period, 26
- Launching an Analysis, 51
- Leaf, 9
- Log Likelihood Ratio, 58
- Mathematical Formulas, 59
- Maximum Temporal Cluster Size, 34, 35
- memory problems, 65
- Memory Requirements, 56
- Methodology Papers, 66
- minimum tree size, 64
- Monte Carlo Replications, 36, 41, 44
- Monte Carlo simulated random data, 65
- Most Likely Cuts, 57
- Multiple Testing, 5
- New Versions, 62
- Node, 9
- Node Cases, 57
- node ID**, 18, 19
- Node Identifier, 57
- Occupational Disease Surveillance, 67
- Output
 - Results File, 29
 - Simulated Log Likelihood Ratios File, 60, 61
- Output Files Tab, 29, 51
- Output Tab, 28
 - Additional Output Files, 29, 30
 - Results File Name, 29
- Parallel Processors, 54
- Parameter Settings, 58
- Parent, 9
- Poisson Distribution, 23
- Poisson Model, 6, 11
- population, 17
- Probability Model, 23, 24
 - Bernoulli**, 23
 - Poisson**, 23
- Probability Model Comparison, 14
- P-value, 58

- Random Number Generator, 62
- Relative Risk, 58
- Results File, 29
- Results File Name, 29
- Results of Analysis, 57
- Root, 9
- Running SaTScan, 51
- Sample Data Sets, 6
- SaTScan ASCII File Format, 20
- SaTScan Import Wizard, 19
- Saving Analysis Parameters, 53
- Secondary Clusters, 14
- Siblings, 9
- Simulated Log Likelihood Ratios, 50
- Simulated Log Likelihood Ratios File, 60, 61
- Simulated Log Likelihood Ratios File (*_llr.csv), 61
- Space-Time Permutation Model, 12, 13
- Space-Time Scan Statistic, 55
- Standard Results File, 57
- Start Date, 26
- Statistical Methodology
 - Bernoulli Model, 11
 - Poisson Model, 11
 - Probability Model Comparison, 14
 - Space-Time Permutation Model, 12, 13
- Status Message, 52
- Study Period, 26
- Summary of Data, 57
- Technical Support, 53, 62
- Temporal Window Tab
 - Maximum Temporal Cluster Size, 34, 35
- Test Run, 5
- Time Window, 57
- Tree, 9
- Tree Structure and Cluster Detection, 4
- Tree Terminology
 - Tree, 9
- Tree-Based Scan Statistic, 10
- Tree-Structured Variable, 8
- Tree-Temporal Scan Model, 7
- Updates and Revisions, 62
- Vaccine Safety Surveillance, 67
- Warnings and Errors, 52