Hierarchical Data Generator based on Tree-Structured Stick Breaking Process for Benchmarking Clustering Methods

Łukasz P. Olech^{1,2}, Michał Spytkowski¹, Halina Kwaśnicka¹, and Zbigniew Michalewicz^{2,3,4}

 ¹Department of Computational Intelligence Wrocław University of Science and Technology
 wybrzeże Stanisława Wyspiańskiego 27, 50-370 Wrocław, Poland
 {lukasz.olech,michal.spytkowski,halina.kwasnicka}@pwr.edu.pl
 ²Complexica Pty Ltd.
 155 Brebner Drive, West Lakes, SA 5021, Australia
 {lo,zm}@complexica.com
 ³Institute of Computer Science, Polish Academy of Sciences
 ul. Ordona 21, 01-237 Warsaw, Poland
 ⁴Polish-Japanese Academy of Information Technology
 ul. Koszykowa 86, 02-008 Warsaw, Poland

Abstract

Object Cluster Hierarchies is a new variant of Hierarchical Cluster Analysis that gains interest in the field of Machine Learning. Being still at an early stage of development, the lack of tools for systematic analysis of Object Cluster Hierarchies inhibits its further improvement. In this paper we address this issue by proposing a generator of synthetic hierarchical data that can be used for benchmarking Object Cluster Hierarchy methods. The article presents a thorough empirical and theoretical analysis of the generator and provides guidance on how to control its parameters. Conducted experiments show the usefulness of the data generator that is capable of producing a wide range of differently structured data. Further, benchmarking datasets that mirror the most common types of hierarchies are generated and made available to the public, together with the developed generator (http://kio.pwr.edu.pl/?page_id=396).

keywords: Artificial Data, Benchmark Data, Benchmark Data Generator, Hierarchical Clustering, Object Cluster Hierarchy, Tree-Structured Stick Breaking Process, Clustering Evaluation, Cluster Analysis.

1 Introduction

A high volume of digital data drives rapid development of analysis methods providing more effective tools to obtain data insights. The key data analysis areas of research are *regression*, *classification*, and *clustering* [8, 35, 30, 17, 20, 32].

Clustering aims at generating *meaningful* groups (clusters) of points in the provided dataset. The measure of meaningfulness of groups depends on the application and is problem-specific [12]. However, the general idea behind any measure of meaningfulness is that points within any particular cluster are as similar as possible to each other, and at the same time, points that belong to different clusters are as different as possible from each other.

The clustering methods can be characterised, among others, by the type of results they produce. One of the classic types is Flat Clustering that organises points into a predefined number of clusters where the only relation between clusters is spacial. The most popular flat clustering method is k-means [37, 15], where the number of clusters is defined by one of the input parameters. Another clustering type, called *Hierarchical Clustering (HC)* or *Hierarchical Clus*ter Analysis (HCA), produces several flat clustering solutions that are organised into a tree structure called dendrogram (a diagram representing a tree). In the dendrogram every node represents a cluster, and the tree structure indicates parent-child relations between clusters which don't exist in flat clustering. To obtain a flat clustering from a dendrogram, it has to be cut. The final (flat) solution can have a different number of clusters depending on where (at which level of the dendrogram) the cut is performed. The obtained flat clustering always cover all the observations from the input set as no data points are left in parent nodes. Hierarchical Clustering methods do not need to specify the number of clusters beforehand; this number is determined by the height where the dendrogram is cut. Since dendrograms have a partial order relation (hierarchical relation) between nodes, they provide insights into how the process of clustering was performed. By comparing different cuts of the same dendrogram, it can also give a view on the clustering from the perspective of different granularities (different number of clusters in the flat solution). One of the examples of hierarchical clustering is Agglomerative Hierarchical Clustering [10, 6].

For a variety of reasons, clustering methods, in general, still have some weaknesses [4]. From the perspective of this article the interest is in the weaknesses of hierarchical methods where the primary issue is *a semantic gap* between how humans perceive hierarchies and results produced by Hierarchical Clustering methods. Human perception describes hierarchical data (e.g., [31]) as possessing the following properties:

- 1. the data can be present in any node in the hierarchy and belongs to that node without being propagated to the child clusters; and
- 2. the data in the child groups should represent equal or more precise concepts than the data in the corresponding parent that should resemble more general concepts; and

3. the data in a node should be more similar to data in parent and child nodes than to unrelated nodes located in other subtrees of the hierarchy.

Object Cluster Hierarchy (OCH) [33, 28, 34] is a prospective extension to Hierarchical Clustering (HC) paradigm [19] that aims at satisfying the above three properties, and it is further described in Section 3.

Even though the above-mentioned properties might be challenging to address, the first point has been (at least partially) incorporated in a few alreadypublished methods such as Tree-Structured Stick Breaking for Hierarchical Clustering (TSSB-HC) [2], Bayesian Rose Trees [5], Inheritance Retention Variance Hierarchical Clustering (IRV-HC) [33], or modified hierarchic Gaussian Mixture Model (Hk++) [28]. Even though OCH is at an early stage of development, there are increasingly more new Object Cluster Hierarchy clustering methods being published. For that reason there is a need to establish a systematic benchmarking approach for OCH. The existing benchmark datasets created to validate HC methods are not suitable to fully validate OCH methods as they don't consider the differences between OCH and HC structures. Because of that, there is no publicly available set of benchmarking data which would assist OCH development.

This lack of commonly accessible benchmark datasets is addressed in this paper. The main contribution of the paper is the development of *a new method* generating hierarchical structures of data with assumed, user-defined properties. The additional benefit for researchers is the establishment of *a new set of* benchmarks — hierarchical structures of data with the ground truth assignment. The implemented generator, together with the benchmarking datasets are freely available online at http://kio.pwr.edu.pl/?page_id=396 along with instructions on how to use it.

The published datasets can serve as a baseline benchmark for methods generating OCH. Additionally, researches can generate new structures of data according to their needs. Furthermore, such data and generator can help with proposing and testing new clustering quality measures. All these applications should significantly boost the research on OCH.

In the article, the term Object Cluster Hierarchy or its abbreviation (OCH) appears in a variety of different contexts which are best explained when compared to the classical Hierarchical Clustering. The term Hierarchical Clustering represents the *concept* (approach), and similarly, Object Cluster Hierarchy represents the *new concept* (paradigm). Furthermore, a dendrogram is a *result* of a Hierarchical Clustering *method* (e.g., hierarchical agglomerative clustering). In the case of Object Cluster Hierarchy, the *result* is called an Object Cluster Hierarchy as well, or more explicitly an Object Cluster Hierarchy structure/result. A *method* is referred to as Object Cluster Hierarchy method or Object Cluster Hierarchy denotes both the approach and the clustering result. The exact meaning is clear, depending on the context of use.

The paper is organised as follows. The next section provides a literature review, followed by a more in-depth introduction to the OCH in Section 3. Section 4 presents the details of the generator. In Section 5, the meaning of generator parameters and their influence on generated data are described; this section also provides insights on how to control these parameters. Conducted experiments and their results are discussed in Section 6, whereas, the published benchmarking datasets are presented in Section 7. Section 8 concludes the paper.

2 Overview of benchmarking in the analysis of clustering methods

Clustering methods can be verified on real and/or synthetic data. The first type of data offers the advantage of representing real-world cases. However, such data may not always be available in sufficient quantities nor in a form that facilitates testing. Then, the second type of data can be used. Artificially generated data should have properties that imitate real data. A generator that produces artificial data has an additional advantage of allowing for finer control over the data used in testing, i.e., specific aspects of clustering can be tested independently.

Any new clustering method should be carefully evaluated and such evaluations are often based on a comparison to other methods. Usually, such a comparison is made by running different methods on a number of commonly used benchmark datasets and then by comparing their results using a set of evaluation measures, e.g., [28, 27, 21, 7, 23, 18, 39, 14, 32]. Benchmark dataset can be made available to the public by researchers (from academia) or come from different companies (from industry). Additionally, there are many publicly available repositories which collect and provide multiple different datasets in one place, that can be used as benchmarking baseline for method comparison, e.g., UCI repository [13], KEEL dataset repository [3], and many others (e.g., [16, 38]).

For example, Adams et al. [2] proposed a method called TSSB-HC and tested it on two datasets – the CIFAR-100 image set [22] and a sample of 1,740 documents from the NIPS 1-12¹ datasets. Still, both of them lack hierarchical structure annotations (even though a hierarchy might exist within the data) so that the ability to establish an OCH structure has not been verified. Spytkowski et al. [33] have proposed an extension of that method called IRV-HC. A comparison with the base method was made using a few synthetic benchmarks generated from a stochastic model of known parameters. Two measures were used: internal (Average Mixture Model Likelihood) and external (Class Purity), but again the class inclusion hierarchy was not considered. Blundell et al. [5] used several datasets to present performance of their hierarchical method Bayesian Rose Trees. One of them was a synthetically generated dataset in the form of binary vectors which is not the main focus of OCH as the vast majority of data are numerical. The other was Spambase Dataset from UCI repository –

¹https://cs.nyu.edu/~roweis/data.html

a subset of the CMU newsgroup dataset reduced to 4 categories. The authors also used the CEDAR Buffalo² digits dataset in two versions for testing – a subset of the full dataset and a sample of only the 0, 2 and 4 digits. All of these datasets were used in the same way as they would be for a flat clustering algorithm, ignoring the unique capabilities of the method. In [28], a GMM-based Hierarchical Clustering method called Hk++ was presented. UCI repository datasets [13] were used for its verification, including the Iris, Wine, Glass Identification and Image Segmentation datasets. However, these datasets are not annotated with the inclusion of the hierarchical relation between classes, so this aspect of OCH couldn't be verified. In that case, elementary, synthetic and hierarchical datasets were also used, but with the purpose to show the concept rather than thoroughly verify the method capabilities.

The above publications are the closest to the concept of Object Cluster Hierarchy that we were able to find in the literature. However, the use of data that is not annotated with a hierarchy of classes does not allow authors to verify results with respect to the obtained hierarchical structure. In the case of Rose Trees [5] in particular, the aspect was omitted due to the very conservative approach towards testing method results. In the case of TSSB method [2], the hierarchy was thoroughly examined empirically and presented to the reader in a visual format. The authors of IRV-HC [33] focused on highlighting how the additional properties of the proposed method impact the final result regarding statistical characteristics, not external validation. In the case of Hk++ [28], the authors attempted to find a way to verify the resulting hierarchies by generating synthetic hierarchical datasets and evaluating results by a modified F-Measure. However, the synthetically generated data used for testing was too simple to support a comprehensive evaluation.

3 A brief introduction to Object Cluster Hierarchies

Object Cluster Hierarchy [33, 28] is a recent variant of Hierarchical Clustering. The HC paradigm [1, 9, 11, 25, 26], whether agglomerative or divisive, produces a dendrogram showing all levels of aggregations. Although there is a hierarchy relationship between nodes in a dendrogram, there is no hierarchy relationship between objects. It is because all objects are assigned only to the dendrogram leaves, and clusters are generated by cutting the tree at any particular level. The level at which a dendrogram is cut determines the number of clusters in the final solution. Any node in the tree, except for leaves, does not have objects assigned to it. Thus, the structure of the generated clusters is flat.

The OCH paradigm extends HC by allowing objects to be assigned to any node in the hierarchy tree. Researchers have already developed methods with such capability, enabling the hierarchy relation between clustered objects to be obtained, e.g., [2, 5, 33, 28].

²http://www.cedar.buffalo.edu/Databases/

Within this paradigm, we have formulated three critical requirements [33] to reflect a semantic (ontological) approach to Hierarchical Clustering:

- 1. *Inheritance* every object that belongs to a given group also belongs to the parents' groups, up to the root;
- 2. *Retention* objects are not required to be located in the tree leaves;
- 3. Variance groups situated lower in the hierarchy are more specific, i.e., every child cannot have higher variation than its ancestors.

These requirements characterise human perception of hierarchy that can also be found in images [31, 2], documents [2, 5], and community structures in social networks [36, 24].

In Figure 1, a comparison between Hierarchical Clustering (a) and Object Cluster Hierarchy (b) is presented based on a simple example. In the former, the final clustering is flat, and the number of clusters depends on the level where the dendrogram is cut. By cutting the tree from Figure 1a at the bottom of the hierarchy, a set of seven clusters is formed, each of them containing one object. Regardless of where the hierarchy is cut, the resulting clustering consists of the same seven objects. In comparison, in Figure 1b, the whole OCH represents clustering — partition of all seven letters. There is no need to cut such a hierarchy. Due to hierarchical relations, objects from child clusters conceptually belong to the parent clusters. Root always contains all the objects (i.e., the whole set), whereas leaves contain only what belongs to them.



Figure 1: Examples of a dendrogram (a) and an Object Cluster Hierarchy (b). Letters represent objects and squares are groups. The arrows show the partial order relation. Note, that the diagrams represent *the final location* of objects within the structures, but conceptually, every object belongs to all its predecessor groups (including the root). In the dendrogram, the partial order relation exists only between the clusters, whereas in the Object Cluster Hierarchy the partial order relation exists between both the clusters and the objects.

Hierarchical Clustering dendrogram always contains all the data points at the bottom of the hierarchy (in leaves). A hierarchy where data points are placed only in leaves is also a valid OCH. Thus, any possible HC result is always an OCH. However, the opposite does not always hold true due to the *Retention* requirement in OCH. Hence, Object Cluster Hierarchy is an abstraction over Hierarchical Clustering.

In one of the early papers on the subject of OCH [2], the authors pointed out that many data arise from a latent hierarchy, for example, a set of text documents or images. An OCH can model such data. In that paper the authors proposed a nonparametric method allowing to discover trees of unbounded width and depth by inferring them during a learning process. That method can assign objects to any node so some nodes can be empty, i.e., without any objects assigned. However, this method is not capable of generating OCH as it does not satisfy the *Variance* requirement – objects belonging to a child node can vary more than objects assigned to the ancestor nodes of that child. Starting from this method as the prototype, the ongoing development has been carried out to propose an improved version satisfying that requirement. The results were presented in [33].

The newly developed method couldn't be comprehensively compared with others against OCH-related characteristics. The problem, to the best of our knowledge, is twofold. Firstly, there are no appropriate evaluation measures dedicated to Object Cluster Hierarchies that would account for the distinguishing characteristics. Secondly, there are no available benchmark datasets with known properties, which could be used for testing different OCH generation methods. The former problem was partially addressed in [34] where several external measures dedicated to OCH have been presented. Further, new external and internal validation techniques have been developed and are planned to be published soon, but their experimental studies require appropriate testing datasets. These two issues led to the development of a method generating synthetic data structures with known characteristics that are presented in this article. Ability to model the attributes of created datasets allows for more in-depth analysis of both methods and evaluation measures.

4 Generator model

The data generated by this model can be interpreted as coming from an infinite mixture model. Most commonly an infinite mixture model is composed of infinite, indexed distributions from which the data is drawn. In such case the mixture weights can be drawn from a Dirichlet distribution, using the Stick Breaking Process. Generally, the distributions for the mixture components are unrelated to each other. This generator uses a similar approach. However, it bases its mixture weights on the Tree-Structured Stick Breaking Process, as described in [2], which arranges the mixture components into a hierarchical structure where each component is a separate node. This structure also defines the relationship between the mixture components as children distributions are based on their parent's distribution parameters.

Hierarchies generated using the Stick Breaking Process possesses the following characteristics [2, 33]:

- 1. Every node can have an unlimited number of child nodes, thus the hierarchy depth and breadth are not limited;
- 2. The children of a node are indexed and ordered. However, this indexing is not important after the generation process finishes;
- 3. The hierarchy may contain empty nodes, that is, nodes that do not have any data directly assigned to them. However, such nodes might still have data belonging to them indirectly through the *Inheritance* requirement;
- 4. A child node distribution parameters are generated based on its parent distribution parameters and a kernel describing the transition;
- 5. The shape of the generated hierarchy depends on several control hyperparameters described in Section 5.

Throughout this paper the following symbols are used to describe the generator and the generated model:

X	- set of all data points, or objects,
x_i	- an object or data point with unique identifier i represented by
	a vector of features,
Θ	- set of all clusters,
ϵ	- specific cluster from Θ ,
ϵ_{x_i}	- cluster of object x_i ,
$\epsilon\epsilon_i$	- the <i>i</i> -th child of cluster ϵ , if the Object Cluster Hierarchy is defined,
ϵ_{\varnothing}	- the root cluster of the Object Cluster Hierarchy,
X_{ϵ}	- set of all objects in cluster ϵ ,
$X_{E_{\epsilon}}$	- set of all objects in hierarchy subtree starting with node ϵ ,
S	- number of elements in set S ,
$ \epsilon $	- a depth of a node ϵ ,
θ_{ϵ}	- distribution of a node ϵ . Specifically, $\theta_{\epsilon_{\varnothing}}$ is
	a distribution of the root node, θ_{ϵ_c} is a distribution of a node
	ϵ_c , and $\theta_{\epsilon\epsilon_i}$ is a distribution of an <i>i</i> -th child of cluster ϵ ,
$Beta(\alpha,\beta)$	- Beta distribution with shape parameters α and β ,
$Gauss(\mu,\sigma)$	- Gaussian distribution with mean μ and standard deviation $\sigma.$

Additionally, to the symbols above, several values are provided to the generator as parameters. The use of these parameters is further described in Section 5, and their influence on the final result is empirically shown in Section 6:

d - the dimensionality of the generated data points, n - the number of data points to be generated, α_0, λ - input parameters controlling the hierarchy depth, used by

	equation $\alpha(\epsilon) = \alpha_0 \lambda^{ \epsilon }$,
γ	- parameter controlling the width of a tree structure,
p,q	- parameters controlling the specificity of the generated data;
	they influence how much smaller the deviation of points'
	features in the child node should be in comparison with points
	in the parent node,
$\theta_{\epsilon \alpha}$	- the distribution of the root node.

Two conditional probabilities are used to determine from which node data is generated. The first is the conditional probability of a datum remaining in node ϵ , at depth $|\epsilon|$, when entering the node:

$$\nu_{\epsilon} = P(x \in X_{\epsilon} | x \in X_{E_{\epsilon}}), \tag{1}$$

$$\nu_{\epsilon} \sim Beta(1, \alpha(\epsilon)), \quad \alpha(\epsilon) = \alpha_0 \lambda^{|\epsilon|}.$$
 (2)

The second is the conditional probability of a datum being transferred to the subtree $\epsilon \epsilon_i$ if it does not remain in node ϵ and hasn't been transferred to any of the previous siblings (i.e., did not travel down sibling subtrees with a lower indices $\epsilon \epsilon_j$, j < i):

$$\psi_{\epsilon\epsilon_i} = P(x \in X_{E_{\epsilon\epsilon_i}} | x \in X_{E_{\epsilon}} \land x \notin X_{\epsilon} \land \neg \exists_{j < i} x \in X_{E_{\epsilon\epsilon_i}}), \tag{3}$$

$$\psi_{\epsilon\epsilon_i} \sim Beta(1,\gamma). \tag{4}$$

Additionally, we need to define the kernel. We begin with a specified root node distribution $\theta_{\epsilon_{\sigma}}$ given as a starting parameter of the generation method. The values set at this point are the means (μ) and standard deviations (σ) for each of the Gaussian distributions in the *d* different dimensions:

$$\theta_{\epsilon_{\varnothing}} = (Gauss(\mu_{\epsilon_{\varnothing}1}, \sigma_{\epsilon_{\varnothing}1}), ..., Gauss(\mu_{\epsilon_{\varnothing}d}, \sigma_{\epsilon_{\varnothing}d})).$$
(5)

From there, for any node for which we need the distribution, we can draw the distribution based on the parent's distribution. The child's mean values are drawn directly from the parent distribution, and the child's standard deviation is based on a scaling factor $(\Delta \sigma_n)$ drawn from the Beta distribution. The values are taken separately for each dimension:

$$\theta_{\epsilon\epsilon_i} = (Gauss(\Delta\mu_1, \sigma_{\epsilon 1}\Delta\sigma_1), ..., Gauss(\Delta\mu_d, \sigma_{\epsilon d}\Delta\sigma_d)), \tag{6}$$

$$\Delta \mu_n \sim Gauss(\mu_{\epsilon n}, \sigma_{\epsilon n}), \quad n = 1, ..., d, \tag{7}$$

$$\Delta \sigma_n \sim Beta(p,q), \quad n = 1, ..., d.$$
(8)

With the kernel defined we can now generate data from the model. We begin with the hyperparameters and the probability distribution for the root node $(\theta_{\epsilon_{\varnothing}} = (Gauss(\mu_{\epsilon_{\varnothing}1}, \sigma_{\epsilon_{\varnothing}1}), ..., Gauss(\mu_{\epsilon_{\varnothing}d}, \sigma_{\epsilon_{\varnothing}d}))).$

The following process continues until n points are generated:

Step 1: If |X| < n go to Step2, else end.

- **Step 2:** Randomly draw an insertion point $i_x \sim Uni(0,1), i_x \in (0,1)$.
- **Step 3:** Set the root node as the current node $(\epsilon_c := \epsilon_{\emptyset})$, depth is $0 (|\epsilon_c| := 0)$.
- **Step 4:** If ν for the current node is not yet known, draw the value $\nu_{\epsilon_c} \sim Beta(1, \alpha_0 \lambda^{|\epsilon_c|})$.
- **Step 5:** If $i_x \leq \nu_{\epsilon_c}$ then $x \sim \theta_{\epsilon_c}$ $(x \sim \theta_{\epsilon_{\varnothing}}$ if $\epsilon_c = \epsilon_{\varnothing}$), the point belongs to the current node $(X_{\epsilon_c} := \{x\} \cup X_{\epsilon_c})$, go to **Step 1**, else move on to **Step 6**.
- **Step 6:** Adjust i_x to new value: $i_x := (i_x \nu_{\epsilon_c})/(1 \nu_{\epsilon_c})$.
- **Step 7:** Set the current child node index $(\epsilon_c \epsilon_i)$ to the first child node of the current node: i := 0.
- **Step 8:** If ψ for the current child node is not yet known, draw the value: $\psi_{\epsilon_c \epsilon_i} \sim Beta(1, \gamma)$.
- **Step 9:** If $\theta_{\epsilon_c \epsilon_i}$ for the current child node is not yet known, draw the values based on the parent of the node:

 $\begin{aligned} \theta_{\epsilon_c \epsilon_i} &= (Gauss(\Delta \mu_1, \sigma_{\epsilon_c 1} \Delta \sigma_1), ..., Gauss(\Delta \mu_d, \sigma_{\epsilon_c d} \Delta \sigma_d)), \\ \Delta \mu_1 \text{ is drawn from the first dimension of the parent node } (\Delta \mu_1 \sim Gauss(\mu_{\epsilon_c 1}, \sigma_{\epsilon_c 1})), \end{aligned}$

 $\Delta \mu_d$ is drawn from the *d*-th dimension of the parent node $(\Delta \mu_d \sim Gauss(\mu_{\epsilon_c d}, \sigma_{\epsilon_c d})), \Delta \sigma_1 \sim Beta(p, q),$

 $\Delta \sigma_d \sim Beta(p,q).$

Step 10: If $i_x \leq \psi_{\epsilon_c \epsilon_i}$ go to **Step 11**, else go to **Step 12**.

- **Step 11:** Adjust the value of i_x to new value: $i_x := i_x/\psi_{\epsilon_c \epsilon_i}$. Make the current child the current node ($\epsilon_c := \epsilon_c \epsilon_i$) and increase depth ($|\epsilon_c| := |\epsilon_c| + 1$). Go to **Step 4**.
- **Step 12:** Adjust the value of i_x to new value $i_x := (i_x \psi_{\epsilon_c \epsilon_i})/(1 \psi_{\epsilon_c \epsilon_i})$. Increment child index of currently relevant child node (i := i + 1). Go to **Step 8**.

The generation process described above is illustrated as a block diagram in Figure 2.



Figure 2: Block diagram of Object Cluster Hierarchy generator.

5 Parameter selection

Hierarchical data might follow hierarchies of different characteristics, e.g., depth, width, the average number of objects per node. Thus, to be applicable to various problems, the generator should provide high flexibility in generating a variety of hierarchies. The primary interest is in the structure of the hierarchy, that is whether the hierarchy is tall or short, wide or narrow, as well as the distribution of data across the levels of the hierarchy. Additionally, the difference between data in parent and child nodes can also be important in some cases. All of these are controlled by several parameters in the model:

- 1. hierarchy depth: α_0 , λ or in a more general sense, the $\alpha(\epsilon)$ function,
- 2. hierarchy width: γ ,
- 3. data specificity: p, q,
- 4. root node distribution: $\theta_{\epsilon_{\varnothing}}$.

The remaining part of this section is organised into five subsections. In the following four sub-sections all the above-mentioned parameter groups are discussed. In Section 5.1 the $\alpha(\epsilon)$ function and its influence on the depth of hierarchy are presented. The γ parameter and its impact on the hierarchy width are described in Section 5.2. The differences in data distributions between nodes in a hierarchy (especially parent-child nodes) are discussed in Section 5.3. As the hierarchy generation process is iterative and top-down, the initial (root) distribution parameters have an impact on the final hierarchy – this is discussed in Section 5.4. A reassignment post-processing procedure allowing for the generated hierarchies to be denoised is described in Section 5.5. Furthermore, the influence of parameters on the generated trees discussed theoretically in this section is also empirically demonstrated later in the paper.

5.1 Controlling hierarchy depth

The depth of the hierarchy is controlled by the function $\alpha(\epsilon) = \alpha_0 \lambda^{|\epsilon|}$. The higher the probability of a datum remaining in a node, the fewer data will travel deep down the tree, and thus the tree will be shallower. On the other hand, if the probability is low, the data will, on average, travel deeper into the tree. The average probability of data remaining in a given node is based on the selected α function, which influences the structure of the tree based on the α_0 and λ parameters:

$$\mathbf{E}[x \in X_{\epsilon}, |\epsilon| = 0] = \frac{1}{1 + \alpha_0},\tag{9}$$

$$\mathbf{E}[x \in X_{\epsilon}, |\epsilon| = n] = \frac{\prod_{i=0}^{n-1} \alpha_0 \lambda^i}{\prod_{j=0}^n (1 + \alpha_0 \lambda^j)}.$$
(10)

Additionally the variance can also be calculated:

$$\mathbf{var}[x \in X_{\epsilon}, |\epsilon| = 0] = \frac{\alpha_0}{(1 + \alpha_0)^2 (2 + \alpha_0)},\tag{11}$$

$$\mathbf{var}[x \in X_{\epsilon}, |\epsilon| = n] =$$

$$= \frac{2 \prod_{i=0}^{n-1} \alpha_0 \lambda^i}{(1 + \alpha_0 \lambda^n) \prod_{j=0}^n (2 + \alpha_0 \lambda^j)}$$

$$- (\mathbf{E}[x \in X_{\epsilon}, |\epsilon| = n])^2.$$
(12)

It is possible to predict the shape of the tree and data distribution based on α_0 and λ parameters:

- $\alpha_0 = 1, \lambda = 1$: the structure of the tree is chaotic and hard to predict, the further away the parameters move from these values the more stable the tree becomes;
- $\alpha_0 < 1, \lambda \leq 1$: shallow structure, data located primarily at the top of the tree;
- α₀ ≤ 1, λ > 1: similar to the above case, the depth of the tree increases but most data is located at the top of the tree;
- $\alpha_0 \ge 1, \lambda < 1$: the structure is deep, but data is not located at the top, the bigger α_0 starts out, and smaller λ is the more data will move down the tree into the central or lower region;
- $\alpha_0 > 1, \lambda \ge 1$: deep structure, data located primarily at the top of the tree but spread out.

5.2 Controlling hierarchy width

The width of the tree is based on the value of the γ parameter at a given node depth j. Given that $x \in X_{E_{\epsilon}}$ and $x \notin X_{\epsilon}$ the average probability of data being generated from a specific subtree (based on the index) can be used to estimate the number of children a node can potentially have:

$$\mathbf{E}[x \in X_{E_{\epsilon\epsilon_i}}, i=1] = \frac{1}{1+\gamma},\tag{13}$$

$$\mathbf{E}[x \in X_{E_{\epsilon\epsilon_i}}, i=n] = \frac{\gamma^{n-1}}{(1+\alpha_0\lambda^j)^n}.$$
(14)

Variance for these values can also be calculated:

$$\mathbf{var}[x \in X_{E_{\epsilon\epsilon_i}}, i=1] = \frac{\gamma}{(1+\gamma)^2(2+\gamma)},\tag{15}$$

$$\operatorname{var}[x \in X_{E_{\epsilon\epsilon_i}}, i=n] = \frac{2\gamma^{n-1}}{(1+\gamma)(2+\gamma)^n} - \left(\mathbf{E}[x \in X_{E_{\epsilon\epsilon_i}}, i=n]\right)^2.$$
(16)

Influence of the parameter γ on the generated hierarchies is as follows:

- $\gamma = 1$: the number of children is chaotic and difficult to predict;
- $\gamma < 1$: narrower tree, fewer children per node on average;
- $\gamma > 1$: wider tree, more children per node on average.

5.3 Controlling data specificity

When a new group is considered, the parameters for that group data points distribution are drawn based on the parent distribution and the kernel parameters p and q. An important aspect of the generated model is that data becomes more specific at lower nodes following the OCH principles. However, the values that are taken for the kernel change the average proportion of standard data deviation between the parent and child. This is based on the expected standard deviation of the new node compared to the old node (taken separately in each dimension):

$$\mathbf{E}[\sigma_{\epsilon\epsilon_i d}] = \sigma_{\epsilon d} \frac{p}{p+q},\tag{17}$$

$$\mathbf{var}[\sigma_{\epsilon\epsilon_i d}] = \sigma_{\epsilon d} \frac{pq}{(p+q)^2(p+q+1)}.$$
(18)

By selecting p and q, the rate at which the nodes become more specific can be altered. The lower the mean is, the more specific every child will be (on average), the higher the variance is, the more variety there will be in how the child nodes relate to their parent.

5.4 Influence of starting distribution on results

Due to the relative nature of the model (i.e., the specific values generated from the model are calculated relative to each other, starting from the root distribution, as shown in equations 5, 6, 7, 8), the choice of initial distribution parameters is not very important. The data generated from the model can be scaled afterwards to any desired values as well as moved in any direction along any dimension. Because of this, the generator assumes a following data distribution for the root node:

$$\theta_{\epsilon_{\varnothing}} = (Gauss(\mu_{\epsilon_{\varnothing}1}, \sigma_{\epsilon_{\varnothing}1}), ..., Gauss(\mu_{\epsilon_{\varnothing}d}, \sigma_{\epsilon_{\varnothing}d})),$$
(19)

where

$$\mu_{\epsilon_{\mathscr{O}}1} = \mu_{\epsilon_{\mathscr{O}}2} = \dots = \mu_{\epsilon_{\mathscr{O}}d} = 0, \tag{20}$$

and

$$\sigma_{\epsilon_{\varnothing}1} = \sigma_{\epsilon_{\varnothing}2} = \dots = \sigma_{\epsilon_{\varnothing}d} = \sigma_{max}.$$
(21)

Every dimension of the root node is described by a normal distribution with zero mean and the standard deviation of value σ_{max} which is a method's parameter. Data generated from the model can be then post-processed to a more desirable spread of values. This is done by applying scaling and translation to all the data generated by the model as well as the parameters of each group node.

5.5 Reassignment post-processing

As it is shown in Section 4, the assignment of points to groups is conducted in a top-down manner separately for every point. It starts from the root node and moves down the hierarchy a particular path considering a different sequence of nodes for a point assignment. For every visited node the conditional probabilities ν_{ϵ} and $\psi_{\epsilon\epsilon_i}$ are calculated (Equations (4) and (6)) and their values determine which sequence of nodes is considered before a point is finally assigned. A point is assigned to the first node for which certain conditions are met (see **Step 5** in Section 4). As this process is greedy and stochastic, only a subset of nodes will be considered for point assignment which does not guarantee that a node with the highest probability will be selected. It introduces noise to the nodes and biases their data distributions.

To address such behaviour, a hierarchy can undergo one form of post-processing after being generated. This process, referred in this paper as *reassignment*, moves the data between clusters in such a way that each object belongs to the cluster it is most likely to be generated from:

$$\forall_{x \in X} \left(x \in X_{\epsilon_a} \Leftrightarrow \neg \exists_{\epsilon_b \neq \epsilon_a} L(x | \theta_{\epsilon_a}) < L(x | \theta_{\epsilon_b}) \right) \tag{22}$$

The process does not modify the number of clusters, hierarchy relations between them or their parameters in any way. It merely relocates data to reduce noise and produce cleaner clusters.

6 Experiments

The generator was tested with a number of different goals in mind. The tests serve to empirically investigate the analytical and intuitive properties of the introduced parameters provided in the previous sections. Thus, a large part of the experiments serves to verify how the different parameter values affect the generated hierarchies. Further, the experiments aim at demonstrating various properties of the generated hierarchies and generator flexibility in producing differently-structured hierarchical data. The influence of the reassignment postprocessing (Section 5.5) on the generated hierarchies has also been investigated. The goal of the reassignment process is to reduce noise in the generated data by moving objects to the node for which the likelihood of being drawn from is the highest. A comparison of post-processed hierarchies with unmodified ones was performed. All of these experiments were done with a primary objective to provide a potential user with a comprehensive understanding of the generator and to assist him/her in choosing the best parameter set for any given use case. An additional benefit from conducted experiments is the establishment of new benchmarking datasets for OCH that are ready to use for method comparison without a need to use the generator.

Due to the stochastic nature of the generator, the results presented in this section were obtained by averaging over 100 generated hierarchies for each of the used parameter sets (s00 - s07) shown in Table 1. Each of the used parameter

sets represents a different type of hierarchical structure. Some parameters such as n = 10,000, d = 2, p = 1, q = 5, $\sigma_{min} = 0.05$ and $\sigma_{max} = 10$ remained constant across all experiments whereas α_0 , λ , γ varied. For every generator run two hierarchies were produced: as generated from the statistical model (the initial assignment of data to nodes), and after reassignment of data (*reassigned* datasets). Datasets with the initial assignment of data are referred to as s00 - s07 depending on which parameter set-up was used (Table 1). For the reassigned hierarchies, the naming convention is the same with an additional letter 'r' (s00r - s07r). Regardless of which variant of the hierarchy (with or without the reassigned procedure) has been created, generator parameters remain the same as the reassignment procedure is applied *after* the data is created. In other words, the reassignment procedure does not influence the data generation process.

Several quantitative measures were used to investigate the aggregate properties of the generated hierarchies:

- \bar{N} the number of nodes in the hierarchy, averaged over all generated hierarchies,
- \overline{L} the number of leaves in the hierarchy (nodes with no children or with empty children only), averaged over all generated hierarchies,
- \overline{D} the depth of the hierarchy, averaged over all generated hierarchies,
- \overline{B} the breadth of the hierarchy, averaged over all levels in a hierarchy, and over all hierarchies generated,
- \bar{P} the average length of all paths in a hierarchy, averaged over all generated hierarchies.

Table 1: Generator parameters used to create experimental data sets from s00 to s07 (or s00r to s07r correspondingly if the reassigned post-processing procedure has been used; the parameters are the same since reassignment is applied after the dataset is created). The remaining parameters are shared between test sets: $n = 10,000, d = 2, p = 1, q = 5, \sigma_{min} = 0.05$ and $\sigma_{max} = 10$. Parameter set selection is based on previous research in this area [2]

Set	α_0	λ	γ
s00 and $s00r$	1	0.5	0.2
s01 and $s01r$	1	1.0	0.2
s02 and $s02r$	1	1.0	1.0
s03 and $s03r$	5	0.5	0.2
s04 and $s04r$	5	1.0	0.2
s05 and $s05r$	5	0.5	1.0
$s\theta 6$ and $s\theta 6r$	25	0.5	0.2
s07 and $s07r$	25	0.5	1.0

All of the reported average measures are accompanied by standard deviations. Since B and P are averages of averages instead of a standard deviation, an *average* of standard deviations over all generated hierarchies is provided. All defined measures are reported separately for the initially generated (Table 2) and reassigned hierarchies (Table 3). The remaining experiments results (Figures 4 to 13) are presented as histograms averaged over the 100 generated hierarchies for each parameter set. The histograms (except for Figures 12 and 13) present measures across different levels of hierarchies providing a more in-depth (structural) view on the generated data. Figures 4, 6, 8, 10 and 12 show results when the reassignment post-processing has not been executed, whereas Figures 5, 7, 9, 11 and 13 show values after the post-processing. The reported measures are:

- average width per level (Figures 4 and 5),
- average number of objects per node per level (Figures 6 and 7),
- average number of children per node per level (Figures 8 and 9),
- average number of leaves per level (Figures 10 and 11),
- average number of nodes with a given number of children (Figures 12 and 13).

Table 2: Accumulative characteristics of generated hierarchies **without** the reassignment procedure. Average \bar{X} values together with standard deviation $\sigma_{\bar{X}}$ (or an average of standard deviations $\bar{\sigma}_{\bar{X}}$) are provided.

Cat	Nodes		Leaves		Dep	oth	Bre	adth	Path length	
sei	\bar{N}	$\sigma_{ar{N}}$	\bar{L}	$\sigma_{ar{L}}$	\bar{D}	$\sigma_{ar{D}}$	\bar{B}	$\bar{\sigma}_{\bar{B}}$	\bar{P}	$\bar{\sigma}_{\bar{P}}$
s00	17.58	7.80	8.75	4.33	4.06	0.82	3.40	2.14	2.86	0.93
s01	95.23	53.06	31.00	17.85	11.73	2.80	7.14	5.10	6.30	2.50
s02	556.81	329.74	271.80	188.79	12.33	2.20	40.44	41.33	5.40	1.99
s03	58.19	21.33	25.84	10.62	6.41	0.73	7.84	5.97	4.36	1.13
s04	3090.88	944.13	483.62	187.81	52.54	6.83	58.14	59.92	19.39	7.88
s05	485.43	149.83	297.62	108.54	6.88	0.55	61.69	64.09	4.25	1.00
s06	175.71	61.50	67.57	26.01	8.83	0.64	17.84	14.97	6.20	1.32
s07	2071.07	536.50	1109.17	367.70	9.27	0.51	201.88	223.79	5.88	1.16

The results of the conducted experiments can be confronted with prior analytical estimations of the effect that parameters have on the structure of the hierarchy (Sections 5 and 5.1 to 5.3). The simplest case is the γ parameter (Section 5.2). This parameter is responsible for the formation of child nodes and as such, the breadth of the hierarchy. For datasets that differ only by the γ value (s01 and s02 or s06 and s07), the distribution of data per level is very similar (Figures 4 and 5). It is because the data distribution is controlled by the α function, which is not influenced by γ . On the other hand, there is a

Table 3: Accumulative characteristics of generated hierarchies with the reassignment procedure. Average \bar{X} values together with standard deviation $\sigma_{\bar{X}}$ (or an average of standard deviations $\bar{\sigma}_{\bar{X}}$) are provided.

Cat	Nodes		Leaves		Dep	oth	Bre	adth	Path length	
Sei	\bar{N}	$\sigma_{ar{N}}$	\bar{L}	$\sigma_{ar{L}}$	\bar{D}	$\sigma_{ar{D}}$	\bar{B}	$\bar{\sigma}_{\bar{B}}$	\bar{P}	$\bar{\sigma}_{\bar{P}}$
$s\theta\theta r$	18.11	8.17	9.30	4.68	4.05	0.82	3.50	2.23	2.86	0.92
s01r	98.56	54.99	34.76	19.50	11.70	2.82	7.40	5.31	6.30	2.48
s02r	642.02	367.21	363.20	211.42	12.23	2.19	46.88	47.62	5.43	1.96
s03r	59.88	22.11	27.64	11.33	6.40	0.74	8.07	6.15	4.35	1.14
s04r	3099.03	936.39	597.36	176.93	52.05	6.96	58.90	60.44	19.17	7.78
s05r	552.61	151.38	366.31	100.87	6.87	0.56	70.30	72.92	4.25	1.03
s06r	180.47	63.97	73.14	28.05	8.81	0.65	18.34	15.40	6.16	1.35
s07r	2310.21	524.57	1375.60	302.29	9.27	0.51	225.19	250.42	5.87	1.18

significant change in the width of the hierarchy, approximately by one order of magnitude (10 times higher for higher γ), as it was predicted by the prior analysis.

The influence of the α_0 and λ parameters on generated hierarchies is difficult to describe (Section 5.1) as the two parameters are intervoven together within the α function (Equation (2)); it also depends on the level of a hierarchy that is currently considered. However, the impact of this function is the best presented in Figures 6 and 7, especially when results for datasets s00, s01, s02 and s04are compared with the results for s03, s05, s06 and s07. The first set has a clear tendency to retain data in higher nodes (nodes that are located closer to the root). In comparison, the other set has the main mass of data located in the lower nodes (nodes located distant from the root). Especially with the $s\theta 6$ and s07 the majority of objects are located in lower nodes, close to the 5th level. For these two datasets, we can see that α on average starts out with low values (due to a high value of α_0) and because of λ parameter being smaller than 1, inclines as moving lower in the hierarchy (compare the influence of parameters on the Beta distribution³). For a low value of α , the probability of retaining data in a node is on average low (see steps 2-5 in Figure 2). Thus the nodes close to the root do not retain data, but as the value of α increases, more objects gather in the lower nodes of the hierarchy before eventually, the remaining data is passed on to the lowest nodes (leaves). From this, we come to an important conclusion about the importance of these two parameters. In cases where it is undesirable to have many generic (root level) objects, and it is important to have clearly distinct, specific (lower level) objects parameters α_0 and λ must have values similar to those present in s06 and s07 – high α_0 and λ that controls the decline of the *Beta* function value over levels to be smaller than 1. Such behaviour was earlier predicted from the analytical study of the parameters (Section 5.1), and the conducted experiments show that behaviour

³http://eurekastatistics.com/beta-distribution-pdf-grapher/

empirically. It appears that the bulk of data is retained at the level in which $\alpha(\epsilon)$ drops below 1.

A very prominent behaviour of the generator seen in all test cases is the production of – what will be referred to from this point onward – *trailing divisions* of data. Trailing divisions occur when the generator attempts to split small remaining partitions of data. This happens both in the right (higher index) children of a populated node and lowers down the hierarchy as presented in Figure 3. In both cases, it is possible to observe large numbers of nodes with a low number of children, usually one or zero (in the latter case the node becoming a leaf node), as well as many nodes that are not populated with data. Trailing divisions reveal the fractal nature of trees generated by the procedure, which manifests itself both when producing direct children for a node (horizontal self-similarity) and going down the hierarchy (vertical self-similarity).



Figure 3: A simple schema of the location of trailing divisions.

The above phenomenon can be visualised horizontally as an ordered set of all children of any node being statistically similar to the ordered set of all children of the node except the first one. It is a direct effect of the Tree-Structured Stick Breaking process (TSSB) [2]. Similarly to the above, from a vertical point of view, in any tree with $\lambda = 1$ (the conditional probability of a datum being assigned to a particular node is constant for all the nodes, see Equation (2)) all sub-trees of a node are statistically similar to that node. In the presented experimental data, these trailing divisions are visible as the falling off "tail" in the histograms of data instances distribution per level (Figures 6 and 7) as well as the cause of the high deviation of width within the trees (Figures 4 and 5). Unfortunately, due to the nature of the TSSB distribution, it is impossible to avoid this behaviour without post-processing. No forms of such post-processing were employed for the experiments presented in this paper.

Since the γ value does not change between levels (see Equation (4)), a critical factor in considering how many children a node will have is the number of data passing through the node during the hierarchy generation process. Because of that, nodes that are located closer to the root are more likely to have more children than smaller nodes lower down the hierarchy. It has also been verified experimentally, and it is shown in Figures 8 and 9. Additionally, the more children a node has, the more of them will be small nodes, i.e., nodes through which few objects pass, resulting primarily in leaf nodes or nodes with a single

child. This behaviour of smaller nodes also transfers lower down the hierarchy where less data reaches, leading to similar behaviour.

Finally, the reassigned test cases show a tendency for data to move down hierarchy levels. It is best shown when comparing Figure 6 with Figure 7. Intuitively, the groups located lower down in the hierarchy are more specific and potentially conflicting data would be prone to moving down into the more specific child clusters during the reassignment process. However, despite this tendency, the hierarchies structure do not change. Hierarchies retain most of their initial characteristics with the mass of data being shifted down towards the lower levels. Due to the post-processing applied to these datasets, they can be better suited for initial testing of grouping methods as they contain less noise. Testing using both types of datasets (unfiltered and filtered) may be the preferred and most valuable approach in every case where the features of the objects are considered.

7 Benchmarking dataset

One of the goals of this article is to provide the research community with a systematic and comprehensive approach for benchmarking OCH methods. To achieve this goal, the benchmarking dataset of 160 hierarchies was created. These hierarchies were chosen from the 1,600 hierarchies analysed in Section 6.

The process of selecting 160 hierarchies to be included in the benchmarking dataset was as follows. For every parameter set-ups out of 16 possibilities (see Table 1), 100 hierarchies h_i^s were generated, where

- s ∈ {s00, s00r, s01, s01r, s02, s02r, s03, s03r, s04, s04r, s05, s05r, s06, s06r, s07, s07r},
- $1 \le i \le 100.$

Each of these 1,600 hierarchies has been described by a vector

$$descr(h_i^s) = (N_i^s, L_i^s, D_i^s, B_i^s, P_i^s),$$
(23)

where

- N_i^s is the number of nodes in the *i*-th hierarchy for parameter set-up s,
- L_i^s is the number of leaves in the *i*-th hierarchy for parameter set-up s,
- D_i^s is the depth of the *i*-th hierarchy for parameter set-up s,
- B_i^s is the breadth of the *i*-th hierarchy for parameter set-up s,
- P_i^s is the average length of all paths in the *i*-th hierarchy for parameter set-up s.



Figure 4: Average hierarchy width (B) on every hierarchy level (number of nodes on every level) without execution of reassignment procedure. Vertical axes show hierarchy width and horizontal axes indicate hierarchy level.



Figure 5: Average hierarchy width (B) on every hierarchy level (number of nodes on every level) with execution of reassignment procedure. Vertical axes show hierarchy width and horizontal axes indicate hierarchy level.



Figure 6: Average distribution of data instances among hierarchy levels **without** execution of reassignment procedure. Vertical axes show the number of instances and horizontal axes indicate hierarchy level.



Figure 7: Average distribution of data instances among hierarchy levels with execution of reassignment procedure. Vertical axes show the number of instances and horizontal axes indicate hierarchy level.



Figure 8: Distribution of the average number of children per node among hierarchy levels **without** execution of reassignment procedure. Vertical axes show the number of children and horizontal axes indicate hierarchy level.



Figure 9: Distribution of the average number of children per node among hierarchy levels **with** execution of reassignment procedure. Vertical axes show the number of children and horizontal axes indicate hierarchy level.



Figure 10: Average number of leaf nodes (L) on every hierarchy level **without** execution of reassignment procedure. Vertical axes show the number of children and horizontal axes indicate hierarchy level.



Figure 11: Average number of leaf nodes (L) on every hierarchy level with execution of reassignment procedure. Vertical axes show the number of children and horizontal axes indicate hierarchy level.



Figure 12: Average number of child nodes for every node in generated hierarchies **without** execution of reassignment procedure. Horizontal axes show the number of children and vertical axes show the number of occurrences (count) in the hierarchies.



Figure 13: Average number of child nodes for every node in generated hierarchies with execution of reassignment procedure. Horizontal axes show the number of children and vertical axes show the number of occurrences (count) in the hierarchies.

Additionally, for every parameters set-ups a vector of average values ${\cal H}^s_{avg}$ has been established

$$H^{s}_{avg} = (\bar{N}^{s}, \bar{L}^{s}, \bar{D}^{s}, \bar{B}^{s}, \bar{P}^{s}), \qquad (24)$$

where \bar{N}^s , \bar{L}^s , \bar{D}^s , \bar{B}^s , \bar{P}^s are values of \bar{N} , \bar{L} , \bar{D} , \bar{B} , \bar{P} from Tables 2 and 3 for a corresponding s.

For all values of parameter s separately, vectors $descr(h_i^s)$ and H_{avg}^s have been min-max scaled to range from 0 to 1. In the scaled space, Euclidean distance has been computed between an average vector and all the corresponding hierarchy vectors. Based on the distances the top 10 closest hierarchies have been identified. These hierarchies were published as benchmarking data for a particular parameter set-up. The procedure has been repeated for all parameter set-ups resulting in a collection of 160 hierarchies that are publicly available (http://kio.pwr.edu.pl/?page_id=396).

Quantitative measures for the published dataset are reported in Table 4. This table presents results in the same format as in Tables 2 and 3 with two additional measures:

- \overline{C} the average number of immediate children for every internal node (node with at least one non-empty child) in a hierarchy, averaged over all generated hierarchies,
- \overline{I} the average number of instances per node in a hierarchy, averaged over all generated hierarchies.

The results presented in Table 4 are very similar to these presented in Tables 2 and 3 with differences mainly attributed to the stochastic nature of the generator. Similarities in the results indicate that the chosen set of published 160 hierarchies is a representative sample and all the conclusions from Section 6 apply to them as well. Specifically, the trends (how the hierarchy statistics change between levels) showed in Figures 4 to 13 apply to the published sample as well.

The average number of children per internal node (\bar{C}) is usually between 1 and 2, which indicates that the hierarchies are quite narrow (Table 4). A manual inspection of hierarchies confirmed that nodes with multiple (e.g., 5) children happen but far less often. The reported average number of instances per node (\bar{I}) differs significantly between sets, e.g., 2.74 for s04 and 679.08 for s00. Furthermore, high standard deviation values $(\bar{\sigma}_{\bar{I}})$ indicate high variability in node sizes within hierarchies that might be compensated by the reassignment procedure.

The reassignment post-processing procedure has the largest impact on the instances per node measure, whereas the differences in other metrics (when hierarchies for reassigned and not reassigned parameter set-ups are compared) are not that significant. They rather stem from differences between samples themselves (separate hierarchies have been chosen for s00 and for s00r). Reassignment procedure acts as a denoising component relocating instances between nodes but leaving hierarchy structures unchanged. This is in line with earlier observations made in this article.

	Max	Madaa		Loonoo		Douth		Drog dth		Child per		Instances	
Set	1000	ies	Lea	Leaves		Depin		Dreaum		Int. Node		per Node	
	\bar{N}	$\sigma_{ar{N}}$	\bar{L}	$\sigma_{ar{L}}$	\bar{D}	$\sigma_{ar{D}}$	\bar{B}	$\bar{\sigma}_{\bar{B}}$	\bar{C}	$\bar{\sigma}_{ar{C}}$	\bar{I}	$ar{\sigma}_{ar{I}}$	
s00	15.10	2.60	7.00	1.56	3.90	0.32	3.09	1.69	1.75	0.65	679.08	1705.60	
s00r	13.90	2.23	6.70	1.16	3.70	0.48	2.98	1.58	1.81	0.63	736.32	1602.21	
s01	109.20	27.29	32.10	8.09	13.10	1.20	7.68	4.81	1.41	0.53	98.32	522.55	
s01r	85.40	23.68	28.50	7.53	11.60	1.65	6.72	4.36	1.49	0.56	125.29	532.46	
s02	592.70	116.10	225.60	36.93	12.60	0.97	43.43	43.04	1.62	1.18	17.41	215.91	
s02r	547.50	132.65	308.50	83.20	11.60	0.84	43.70	41.21	2.28	1.58	19.24	180.24	
s03	50.90	8.40	21.70	4.16	6.20	0.42	7.08	5.20	1.71	0.67	201.36	707.08	
s03r	50.80	7.86	23.20	5.20	6.10	0.32	7.15	5.27	1.80	0.71	201.38	441.08	
s04	3711.80	502.46	445.50	58.94	50.80	2.04	71.93	70.96	1.14	0.30	2.74	45.35	
s04r	2839.10	242.45	548.70	47.38	49.50	3.14	56.21	55.77	1.24	0.40	3.55	37.76	
s05	468.90	47.42	260.80	25.08	6.90	0.32	59.34	60.97	2.25	1.75	21.53	127.33	
s05r	512.50	87.08	336.90	55.25	6.50	0.53	68.32	67.81	2.93	2.15	19.99	64.46	
s06	167.20	26.73	60.70	9.87	8.50	0.53	17.64	14.40	1.56	0.70	61.38	252.74	
s06r	173.60	20.49	69.50	7.46	8.30	0.48	18.73	15.11	1.66	0.71	58.37	96.44	
s07	2221.70	242.26	1032.20	106.40	9.10	0.32	220.53	243.58	1.87	1.36	4.56	23.28	
s07r	2285.50	210.33	1360.20	123.66	9.40	0.52	220.00	247.80	2.47	1.81	4.41	8.71	

Table 4: Accumulative characteristics of hierarchies published as benchmarking dataset. Average \bar{X} values together with standard deviation $\sigma_{\bar{X}}$ (or an average of standard deviations $\bar{\sigma}_{\bar{X}}$) are provided.

The key summary characteristics of the published hierarchies are presented in Table 5. They are combined for the reassigned and not reassigned variants for the reasons provided in the paragraphs above. The table provides a quick reference for a potential user assisting in a decision of which hierarchies to use in a particular use-case. The decision whether to use the reassigned variant of not should be based on the fact that reassignment denoise the hierarchy, so it is expected to be easier to cluster the points.

One use-case for the benchmarking dataset might be a development of a new OCH method. If a prototype is being developed, and authors want to evaluate the performance of the method on very high and narrow hierarchies, s04 or s04r should be used. On the other hand, if performance on not high but wider hierarchies is to be validated, s05 or set05r are a better choice. An additional benefit of the benchmarking dataset is that the complexity of a particular hierarchy aspect (width, depth, number of nodes, number of instances per node) can be changed gradually in the series of experiments validating method scalability. For example, to test the method's performance on hierarchies of similar breadth and depth, the initial experiments might be conducted on s00 or s00rwhich, on average, are 3.9 deep (\overline{D}) , 3.09 nodes wide (\overline{B}) , and consist of 15.1 nodes (\bar{N}) (for the reassignment variant $\bar{D} = 3.7$; $\bar{B} = 2.98$ and $\bar{N} = 13.9$). The second round of tests can be conducted on s03 or s03r hierarchies that are still

characterised by a similar breadth and depth, but in this case, the hierarchies are higher, wider ($\bar{D} = 6.2$; $\bar{B} = 7.08$ for s03 and $\bar{D} = 6.1$; $\bar{B} = 7.15$ for s03r) and have an increased number of nodes (50.9 for s03 and 50.8 for s03r).

The established benchmarking dataset provides research communities with a standardised approach to compare their methods. The published hierarchies cover a variety of different hierarchical structures and point distributions allowing for a comprehensive method evaluation but they don't provide all the possible hierarchies. In that case, the generator described in this article can be used in order to generate hierarchies with the desired characteristics.

Table 5: Summary description of hierarchies published as benchmarking dataset.

Set	Summary Description
<u></u>	– the smallest hierarchical structures
000 000	– similar depth and breadth
SUUT	– the highest number of instances per node
	- the most longitudinal hierarchies (high and narrow)
s01	– high number of instances per node
s01r	– low number of nodes
	– low average number of children per node
	– medium number of nodes
s02	– wide hierarchies
s02r	– low average number of instances per node
	– medium number of leaves
s03	- similar to e00 and e00r but a little larger structures
s03r	- similar to sob and sobr but a intile larger structures
	– the highest hierarchies
s04	– the largest number of nodes
s04r	– very narrow
	– the lowest average number of children per node
s05	– very wide hierarchies
s00 e05r	– medium number of nodes and leaves
3001	– the largest number of children per node
s06	$-$ similar to $e^{\Omega t}$ and $e^{\Omega tr}$ but a bit larger structures
s06r	similar to sol and soll but a bit larger structures
	– very large number of nodes
e07	– the widest hierarchies
307 c077r	– the largest number of leaves
3011	– very low average number of instances per node
	– very high average number of children per node

8 Conclusions

In this paper a novel generator of synthetic hierarchical data and a set of benchmarking datasets have been proposed. They aim at providing the research community with a systematic way to benchmark OCH clustering methods and to assist with further OCH development (e.g., development of clustering validation measures). This research has presented a thorough (theoretical and empirical) analysis of the generator that should provide the reader with a comprehensive understanding of how to use it.

The experiments presented in the previous sections highlight both strengths and weaknesses of the proposed generator. A prominent strength is a high range of different tree structures that can be generated and the ability to fine-control these structures using the introduced parameters. Because of that, a wide range of different hierarchy types, often seen in the real-world problems [2], has been generated and made publicly available⁴.

Published benchmarks and the ability to create more hierarchies using the generator is laying a solid foundation for further development of the concept of Object Cluster Hierarchies. Generated hierarchies can assist not only in clustering methods development and comparison but also in research of OCH quality measures. The latter is important since existing internal and external measures do not fully recognise the differences between OCH and Hierarchical Clustering results.

From a practical point of view, it is also beneficial that the generator's parameters can be separated into groups, each controlling a different aspect of the hierarchy. Vertical distribution of data is controlled by α_0 and λ , hierarchy width depends on the value of γ , and p, q controls the data specificity. As shown in this article, every parameter set has an interpretation, and its effect on the generated hierarchy is straightforward. This allows for a further fine-tuning towards desired test data. The generation process scales with the number of points to generate, expanding the hierarchy as more elements are generated.

One of the conclusions that emerged from the experiments was that generated hierarchies would display a degree of self-similarity replicating the same general form both vertically and horizontally. We called it as a generation of trailing divisions. Because of that, a few specific areas of the generated hierarchies are not fully controlled. A remedy to this issue is to use a post-processing procedure similar to the reassignment process described in Section 5.5. This should result in cleaner hierarchies and give the user more control over their overall structure.

In its current form, the generator is limited to a generation of normallydistributed, multidimensional, and uncorrelated real value data. It can be extended to use different kernels leading to different structures of generated hierarchies or generators operating on different types of data.

Furthermore, the self-similar (fractal) nature of the hierarchies suggests a potential for the generator to be described using the language of fractals and

⁴http://kio.pwr.edu.pl/?page_id=396

especially L-Systems [29]. Describing the generation process in that form may provide a different (more granular) view over the details of the hierarchical structure.

References

- M. Abdolali and M. Rahmati. Neither global nor local: A hierarchical robust subspace clustering for image data. <u>Information Sciences</u>, 514:333 – 353, 2020.
- [2] R. P. Adams, Z. Ghahramani, and M. I. Jordan. Tree-structured stick breaking for hierarchical data. In Proceedings of the 23rd International <u>Conference on Neural Information Processing Systems - Volume 1</u>, NIPS'10, pages 19–27, USA, 2010. Curran Associates Inc.
- [3] J. Alcalá-Fdez, A. Fernández, J. Luengo, J. Derrac, S. García, L. Sánchez, and F. Herrera. Keel data-mining software tool: data set repository, integration of algorithms and experimental analysis framework. <u>Journal of</u> Multiple-Valued Logic & Soft Computing, 17, 2011.
- [4] C. Blundell, Y. W. Teh, and K. Heller. Discovering non-binary hierarchical structures with bayesian rose trees. In <u>Mixtures: Estimation and</u> Applications, pages 161–185. John Wiley & Sons, 05 2011.
- [5] C. Blundell, Y. W. Teh, and K. A. Heller. Bayesian rose trees. In P. Grünwald and P. Spirtes, editors, <u>UAI 2010</u>, Proceedings of the Twenty-Sixth <u>Conference on Uncertainty in Artificial Intelligence</u>, Catalina Island, CA, <u>USA</u>, July 8-11, 2010, pages 65–72. AUAI Press, 2010.
- [6] Z. Cai, X. Yang, T. Huang, and W. Zhu. A new similarity combining reconstruction coefficient with pairwise distance for agglomerative clustering. Information Sciences, 508:173 – 182, 2020.
- [7] H. Cao, S. Bernard, R. Sabourin, and L. Heutte. Random forest dissimilarity based multi-view learning for radiomics application. <u>Pattern</u> Recognition, 88:185 – 197, 2019.
- [8] A. Cena and M. Gagolewski. Genie+owa: Robustifying hierarchical clustering with owa-based linkages. Information Sciences, 520:324 336, 2020.
- [9] V. Cohen-Addad, V. Kanade, F. Mallmann-Trenn, and C. Mathieu. <u>Hierarchical Clustering: Objective Functions and Algorithms</u>, pages 378– 397. Society for Industrial and Applied Mathematics, 2018.
- [10] V. Cohen-addad, V. Kanade, F. Mallmann-trenn, and C. Mathieu. Hierarchical clustering: Objective functions and algorithms. <u>J. ACM</u>, 66(4):26:1– 26:42, June 2019.

- [11] G. Costa, G. Manco, R. Ortale, and E. Ritacco. Hierarchical clustering of {XML} documents focused on structural components. <u>Data and Knowledge</u> Engineering, 84:26 – 46, 2013.
- [12] N. J. de Vries, Ł. P. Olech, and P. Moscato. <u>Introducing Clustering with</u> <u>a Focus in Marketing and Consumer Analysis</u>, pages 165–212. Springer International Publishing, Cham, 2019.
- [13] D. Dheeru and E. Karra Taniskidou. UCI machine learning repository, 2017.
- [14] P. Douglas, S. Harris, A. Yuille, and M. S. Cohen. Performance comparison of machine learning algorithms and number of independent components used in fmri decoding of belief vs. disbelief. <u>NeuroImage</u>, 56(2):544 – 553, 2011. Multivariate Decoding and Brain Reading.
- [15] A. K. Dubey, U. Gupta, and S. Jain. Comparative study of k-means and fuzzy c-means algorithms on the breast cancer data. <u>International Journal</u> on Advanced Science, Engineering and Information Technology, 8(1):18–29, 2018.
- [16] P. Fränti and S. Sieranoja. K-means properties on six clustering benchmark datasets, 2018.
- [17] L. d. l. Fuente-Tomas, B. Arranz, G. Safont, P. Sierra, M. Sanchez-Autet, A. Garcia-Blanco, and M. P. Garcia-Portilla. Classification of patients with bipolar disorder using k-means clustering. PLOS ONE, 14(1):1–15, 01 2019.
- [18] J. GarcAŋa, B. Crawford, R. Soto, and G. Astorga. A clustering algorithm applied to the binarization of swarm intelligence continuous metaheuristics. Swarm and Evolutionary Computation, 44:646 – 664, 2019.
- [19] C. Hennig, M. Meila, F. Murtagh, and R. Rocci. <u>Handbook of cluster</u> analysis. CRC Press, 2015.
- [20] S. Jiang, G. Chen, X. Song, and L. Liu. Deep patch representations with shared codebook for scene classification. <u>ACM Trans. Multimedia Comput.</u> Commun. Appl., 15(1s):5:1–5:17, Jan. 2019.
- [21] T. Jo. <u>Text Clustering: Evaluation</u>, pages 249–268. Springer International Publishing, Cham, 2019.
- [22] A. Krizhevsky and G. Hinton. Learning multiple layers of features from tiny images. University of Toronto, 1(4), 2009.
- [23] S. K. Lakshmanaprabu, K. Shankar, M. Ilayaraja, A. W. Nasir, V. Vijayakumar, and N. Chilamkurti. Random forest for big data classification in the internet of things using optimal features. <u>International Journal of</u> Machine Learning and Cybernetics, Jan 2019.

- [24] E. Massaro and F. Bagnoli. Hierarchical community structure in complex (social) networks. <u>Acta Physica Polonica B Proceedings Supplement</u>, 7:379, 02 2014.
- [25] B. Moseley and J. Wang. Approximation bounds for hierarchical clustering: Average linkage, bisecting k-means, and local search. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, <u>Advances in Neural Information Processing Systems</u> 30, pages 3094–3103. Curran Associates, Inc., 2017.
- [26] F. Murtagh and P. Contreras. Algorithms for hierarchical clustering: an overview. Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery, 2(1):86–97, 2012.
- [27] P. B. Myszkowski, M. E. Skowroński, Ł. P. Olech, and K. Oślizło. Hybrid ant colony optimization in solving multi-skill resource-constrained project scheduling problem. Soft Computing, 19(12):3599–3619, Dec 2015.
- [28] Ł. P. Olech and M. Paradowski. <u>Hierarchical Gaussian Mixture Model with</u> <u>Objects Attached to Terminal and Non-terminal Dendrogram Nodes</u>, pages 191–201. Springer International Publishing, Cham, 2016.
- [29] P. Prusinkiewicz and J. Hanan. <u>Lindenmayer systems</u>, fractals, and plants, volume 79. Springer Science & Business Media, 2013.
- [30] M. Z. Rodriguez, C. H. Comin, D. Casanova, O. M. Bruno, D. R. Amancio, L. d. F. Costa, and F. A. Rodrigues. Clustering algorithms: A comparative approach. PLOS ONE, 14(1):1–34, 01 2019.
- [31] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, and L. Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. <u>International Journal of</u> Computer Vision (IJCV), 115(3):211–252, 2015.
- [32] N. Sharma, A. Bajpai, and R. Litoriya. Comparison the various clustering algorithms of weka tools. <u>International Journal of Emerging Technology</u> and Advanced Engineering, 2:73–80, 2012.
- [33] M. Spytkowski and H. Kwasnicka. Hierarchical clustering through bayesian inference. In <u>ICCCI (1)</u>, volume 7653 of <u>Lecture Notes in Computer Science</u>, pages 515–524. Springer, 2012.
- [34] M. Spytkowski, L. P. Olech, and H. Kwaśnicka. <u>Hierarchy of Groups</u> <u>Evaluation Using Different F-Score Variants</u>, pages 654–664. Springer Berlin Heidelberg, Berlin, Heidelberg, 2016.
- [35] K. Sri Dhivya Krishnan and P. T. V. Bhuvaneswari. Multiple linear regression-based prediction model to detect hexavalent chromium in drinking water. In H. S. Behera, J. Nayak, B. Naik, and A. Abraham, editors, <u>Computational Intelligence in Data Mining</u>, pages 493–504, Singapore, 2019. Springer Singapore.

- [36] T. M. G. Tennakoon and R. Nayak. Discovering influence hierarchy based on frequent social interactions. In <u>2018 IEEE/ACM</u> <u>International Conference on Advances in Social Networks Analysis and</u> <u>Mining (ASONAM)</u>, pages 575–576, Aug 2018.
- [37] C. TÃőrnÄČucÄČ, D. GÃşmez-PÃľrez, J. L. BalcÃązar, and J. L. MontaÃśa. Global optimality in k-means clustering. <u>Information Sciences</u>, 439-440:79 – 94, 2018.
- [38] J. Vanschoren, J. N. van Rijn, B. Bischl, and L. Torgo. Openml: Networked science in machine learning. SIGKDD Explorations, 15(2):49–60, 2013.
- [39] X. Wang, Z. Lei, X. Guo, C. Zhang, H. Shi, and S. Z. Li. Multi-view subspace clustering with intactness-aware similarity. <u>Pattern Recognition</u>, 88:50 – 63, 2019.